

IBM SPSS Modeler  
14.2-Modellierungsknoten



*Hinweis:* Lesen Sie vor der Verwendung dieser Informationen und des zugehörigen Produkts die allgemeinen Informationen unter Hinweise auf S. .

Dieses Dokument enthält eigentumsrechtlich geschützte Informationen von SPSS Inc., an IBM Company. Sie werden im Rahmen einer Lizenzvereinbarung bereitgestellt und sind durch Copyright-Gesetze geschützt. Die in dieser Publikation enthaltenen Informationen umfassen keinerlei Produktgewährleistungen und keine der Aussagen in diesem Handbuch darf als solche ausgelegt werden.

Wenn Sie Informationen an IBM bzw. SPSS senden, räumen Sie IBM und SPSS das nicht ausschließliche Recht ein, die Informationen in jeglicher Form zu verwenden bzw. weiterzugeben, die dem Unternehmen geeignet erscheint, ohne dass ihm daraus Verbindlichkeiten Ihnen gegenüber entstehen.

© **Copyright IBM Corporation 1994, 2011..**

---

# Vorwort

IBM® SPSS® Modeler ist die auf Unternehmensebene einsetzbare Data-Mining-Workbench von IBM Corp.. Mit SPSS Modeler können Unternehmen und Organisationen die Beziehungen zu ihren Kunden bzw. zu den Bürgern durch ein tief greifendes Verständnis der Daten verbessern. Organisationen benutzen die mithilfe von SPSS Modeler gewonnenen Erkenntnisse zur Bindung profitabler Kunden, zur Ermittlung von Cross-Selling-Möglichkeiten, zur Gewinnung neuer Kunden, zur Ermittlung von Betrugsfällen, zur Reduzierung von Risiken und zur Verbesserung der Verfügbarkeit öffentlicher Dienstleistungen.

Die visuelle Benutzeroberfläche von SPSS Modeler erleichtert die Anwendung des spezifischen Geschäftswissens der Benutzer, was zu leistungsstärkeren Vorhersagemodellen führt und die Zeit bis zur Lösungserstellung verkürzt. SPSS Modeler bietet zahlreiche Modellierungsverfahren, beispielsweise Algorithmen für Vorhersage, Klassifizierung, Segmentierung und Assoziationserkennung. Nach der Modellerstellung ermöglicht IBM® SPSS® Modeler Solution Publisher die unternehmensweite Bereitstellung für Entscheidungsträger oder in einer Datenbank.

## Über IBM Business Analytics

IBM Business Analytics-Software bietet vollständige, einheitliche und genaue Informationen, auf die Entscheidungsträger vertrauen, um die Unternehmensleistung zu steigern. Ein umfassendes Portfolio von Anwendungen für [Unternehmensinformationen](#), [Vorhersageanalysen](#), [Verwaltung der Finanzleistung und Strategie](#) sowie [Analysen](#) bietet sofort klare und umsetzbare Einblicke in die aktuelle Leistung und ermöglicht die Vorhersage zukünftiger Ergebnisse. In Kombination mit umfassenden Branchenlösungen, bewährten Vorgehensweisen und professionellen Dienstleistungen können Unternehmen jeder Größe optimale Produktivität erreichen, die Entscheidungsfindung zuverlässig automatisieren und bessere Ergebnisse erzielen.

Als Teil dieses Portfolios unterstützt die IBM SPSS Predictive Analytics-Software Unternehmen dabei, zukünftige Ereignisse vorherzusagen und aktiv auf diese Erkenntnisse zu reagieren, um bessere Geschäftsergebnisse zu erzielen. Kunden aus den Bereichen Wirtschaft, Behörden und Bildung aus aller Welt verlassen sich auf die IBM SPSS-Technologie. Sie bringt Ihnen beim Gewinnen, Halten und Ausbauen neuer Kundenbeziehungen einen Wettbewerbsvorteil und verringert gleichzeitig das Betrugs- sowie andere Risiken. Durch Integration der IBM SPSS-Software in den täglichen Betrieb können diese Unternehmen qualifizierte Vorhersagen treffen und dadurch die Entscheidungsfindung so ausrichten und automatisieren, dass Geschäftsziele erreicht werden und ein messbarer Wettbewerbsvorteil entsteht. Wenn Sie weitere Informationen wünschen oder einen Mitarbeiter kontaktieren möchten, ist dies unter <http://www.ibm.com/spss> möglich.

## Technischer Support

Kunden mit Wartungsvertrag können den technischen Support in Anspruch nehmen. Kunden können sich an den technischen Support wenden, wenn sie Hilfe bei der Arbeit mit IBM Corp.-Produkten oder bei der Installation in einer der unterstützten Hardware-Umgebungen benötigen. Die Kontaktdaten des Technischen Supports finden Sie auf der IBM Corp.-Website

unter <http://www.ibm.com/support>. Sie müssen bei der Kontaktaufnahme Ihren Namen, Ihre Organisation und Ihre Supportvereinbarung angeben.

---

# Inhalt

<b>1</b>	<b>Informationen zu IBM SPSS Modeler</b>	<b>1</b>
	IBM SPSS Modeler Server . . . . .	1
	IBM SPSS Modeler-Optionen . . . . .	2
	IBM SPSS Text Analytics . . . . .	2
	IBM SPSS Modeler-Dokumentation . . . . .	3
	Anwendungsbeispiele . . . . .	4
	Ordner "Demos" . . . . .	4
<b>2</b>	<b>Einführung in die Modellierung</b>	<b>6</b>
	Erstellen des Streams . . . . .	8
	Durchsuchen des Modells . . . . .	13
	Bewertung des Modells . . . . .	18
	Scoren von Datensätzen . . . . .	22
	Zusammenfassung . . . . .	23
<b>3</b>	<b>Übersicht über die Modellbildung</b>	<b>24</b>
	Überblick über Modellierungsknoten . . . . .	24
	Erstellung von aufgeteilten Modellen . . . . .	31
	Aufteilung und Partitionierung . . . . .	33
	Modellierungsknoten zur Unterstützung aufgeteilter Modelle. . . . .	34
	Von der Aufteilung betroffene Merkmale . . . . .	35
	Feldoptionen der Modellierungsknoten . . . . .	36
	Verwenden von Häufigkeits- und Gewichtungsfeldern . . . . .	39
	Analyseoptionen bei Modellierungsknoten . . . . .	41
	Neigungsbewertungen . . . . .	44
	Modell-Nuggets . . . . .	45
	Modellverknüpfungen . . . . .	46
	Ersetzen eines Modells . . . . .	49
	Die Modellpalette . . . . .	50
	Durchsuchen von Modell-Nuggets . . . . .	52
	Modell-Nugget-Übersicht/-Information . . . . .	54
	Bedeutsamkeit des Prädiktors . . . . .	54
	Modelle für Ensembles . . . . .	57
	Modell-Nuggets für aufgeteilte Modelle. . . . .	66
	Verwendung von Modell-Nuggets in Streams . . . . .	68

Erneutes Erzeugen eines Modellierungsknotens . . . . .	70
Importieren und Exportieren von Modellen als PMML . . . . .	70
Nicht verfeinerte Modelle . . . . .	74

## **4 Screening von Modellen 75**

Screening von Feldern und Datensätzen . . . . .	75
Merkmalsauswahlknoten . . . . .	75
Einstellungen für das Merkmalsauswahlmodell . . . . .	76
Merkmalsauswahloption . . . . .	78
Modell-Nuggets vom Typ "Merkmalsauswahl" . . . . .	80
Ergebnisse des Merkmalsauswahlmodells. . . . .	80
Auswählen der Felder nach Wichtigkeit. . . . .	82
Generieren eines Filters aus einem Merkmalsauswahlmodell . . . . .	82
Anomalieerkennungsknoten . . . . .	83
Anomalieerkennung – Modelloptionen. . . . .	85
Anomalieerkennung– Expertenoptionen. . . . .	86
Modell-Nuggets vom Typ "Anomalieerkennung" . . . . .	88
Anomalieerkennungsmodelle – Details . . . . .	89
Anomalieerkennungsmodell – Übersicht . . . . .	90
Anomalieerkennungsmodell – Einstellungen . . . . .	91

## **5 Knoten für die automatisierte Modellierung 93**

Knoten für die automatisierte Modellierung – Algorithmuseinstellungen . . . . .	94
Knoten für die automatisierte Modellierung – Stoppregeln . . . . .	95
Knoten "Automatischer Klassifizierer" . . . . .	96
Knoten "Automatischer Klassifizierer" - Modelloptionen . . . . .	98
Knoten "Automatischer Klassifizierer" – Expertenoptionen . . . . .	100
Fehlklassifizierungskosten . . . . .	103
Knoten "Automatischer Klassifizierer" - Optionen für Verwerfen . . . . .	104
Knoten "Automatischer Klassifizierer" - Einstellungsoptionen . . . . .	105
Knoten "Auto-Numerisch" . . . . .	107
Knoten "Auto-Numerisch" - Modelloptionen . . . . .	108
Knoten "Auto-Numerisch" - Expertenoptionen. . . . .	111
Knoten "Auto-Numerisch" - Einstellungsoptionen . . . . .	113
Knoten "Autom. Cluster" . . . . .	114
Knoten "Autom. Cluster" - Modelloptionen . . . . .	115

Knoten "Autom. Cluster" - Expertenoptionen . . . . .	117
Knoten "Autom. Cluster" - Optionen für Verwerfen . . . . .	119
Nugget für automatisierte Modellierung . . . . .	120
Generieren von Knoten und Modellen . . . . .	123
Generieren von Evaluationsdiagrammen . . . . .	123
Evaluationsdiagramme . . . . .	124

## **6 Decision Trees (Entscheidungsbäume) 126**

Entscheidungsbaum-Modelle . . . . .	126
Interactive Tree Builder . . . . .	129
Erweitern und Reduzieren des Baums . . . . .	131
Definieren benutzerdefinierter Aufteilungen. . . . .	132
Aufteilungsdetails und Ersatztrenner . . . . .	134
Anpassen der Baumansicht. . . . .	136
Gewinne . . . . .	137
Risiken . . . . .	146
Speichern der Baummodelle und Ergebnisse . . . . .	146
Generieren von Filter- und Auswahlknoten. . . . .	151
Generieren einer Regelmenge aus einem Entscheidungsbaum . . . . .	151
Direktes Erstellen eines Baummodells . . . . .	152
Entscheidungsbaumknoten. . . . .	154
C&R-Baumknoten . . . . .	155
CHAID-Knoten . . . . .	156
QUEST-Knoten . . . . .	157
Entscheidungsbaumknoten - Feldoptionen. . . . .	158
Entscheidungsbaumknoten - Erstellungsoptionen . . . . .	159
Modelloptionen für Entscheidungsbaumknoten . . . . .	174
C5.0-Knoten . . . . .	177
Modelloptionen für C5.0-Knoten. . . . .	179
Entscheidungsbaummodell-Nuggets. . . . .	181
Modell-Nuggets bei einzelnen Bäumen . . . . .	183
Modell-Nuggets zur Verbesserung/Verstärkung und für sehr große Daten-Sets. . . . .	193
Regelmengen-Modell-Nuggets . . . . .	194
Regelmenge – Registerkarte "Modell". . . . .	196
Projekte aus AnswerTree 3.0 importieren . . . . .	197

## **7 Bayes-Netzwerk-Modelle** **198**

Bayes-Netzwerk-Knoten . . . . .	198
Modelloptionen für Bayes-Netzwerk-Knoten . . . . .	200
Expertenoptionen für Bayes-Netzwerk-Knoten . . . . .	203
Modell-Nuggets vom Typ "Bayes-Netzwerk" . . . . .	205
Einstellungen im Bayes-Netzwerk-Modell . . . . .	206
Bayes-Netzwerk-Modell – Übersicht . . . . .	208

## **8 Neural Networks (Neuronale Netze)** **210**

Das neuronale Netzwerkmodell . . . . .	211
Verwenden von neuronalen Netzwerken mit Legacy-Streams . . . . .	212
Ziele . . . . .	212
Grundeinstellungen . . . . .	214
Stoppregeln . . . . .	215
Ensembles . . . . .	216
Erweitert . . . . .	217
Modelloptionen . . . . .	218
Modellübersicht . . . . .	219
Bedeutsamkeit des Prädiktors . . . . .	220
Vorhersage nach Beobachtung . . . . .	221
Klassifizierung . . . . .	222
Netzwerk . . . . .	223
Einstellungen . . . . .	224

## **9 Entscheidungsliste** **226**

Entscheidungslistenmodell – Optionen . . . . .	232
Entscheidungslistenknoten – Expertenoptionen . . . . .	234
Modell-Nugget vom Typ "Entscheidungsliste" . . . . .	235
Entscheidungslisten-Modell-Nugget – Einstellungen . . . . .	237
Decision List Viewer . . . . .	238
Arbeitsmodellbereich . . . . .	238
Registerkarte "Alternativen" . . . . .	240
Registerkarte "Snapshots" . . . . .	242
Arbeiten mit Decision List Viewer . . . . .	244

Linearknoten . . . . .	266
Lineare Modelle. . . . .	267
Ziele . . . . .	268
Grundeinstellungen . . . . .	270
Modellauswahl . . . . .	271
Ensembles . . . . .	273
Erweitert . . . . .	274
Modelloptionen . . . . .	274
Modellübersicht . . . . .	275
Automatische Datenaufbereitung . . . . .	276
Bedeutsamkeit des Prädiktors . . . . .	277
Vorhersage nach Beobachtung . . . . .	278
Residuen . . . . .	279
Ausreißer . . . . .	280
Effekte . . . . .	281
Koeffizienten . . . . .	283
Geschätzte Mittel . . . . .	285
Modellerstellungsübersicht . . . . .	286
Einstellungen . . . . .	287
Logistikknoten . . . . .	287
Logistikknoten – Modelloptionen . . . . .	288
Hinzufügen von Termen zu einem logistischen Regressionsmodell . . . . .	294
Expertenoptionen für Logistikknoten . . . . .	296
Logistische Regression – Konvergenzoptionen . . . . .	297
Logistische Regression – Erweiterte Ausgabe . . . . .	298
Logistische Regression – Optionen für die Schrittkriterien . . . . .	301
Logistik-Modell-Nugget . . . . .	302
Logistik-Modell-Nugget – Details. . . . .	303
Logistik-Modell-Nugget – Übersicht. . . . .	305
Logistik-Modell-Nugget – Einstellungen. . . . .	306
Logistik-Modell-Nugget – Erweiterte Ausgabe . . . . .	309
Faktor/PCA-Knoten . . . . .	311
Faktor/PCA-Knoten – Modelloptionen . . . . .	312
Faktor/PCA-Knoten – Expertenoptionen . . . . .	313
Faktor/PCA-Knoten – Rotationsoptionen . . . . .	314
Modell-Nugget vom Typ “Faktor/PCA” . . . . .	315
Modell-Nugget vom Typ “Faktor/PCA” – Gleichungen . . . . .	316
Modell-Nugget vom Typ “Faktor/PCA” – Übersicht. . . . .	316
Modell-Nuggets vom Typ “Faktor/PCA” – Erweiterte Ausgabe . . . . .	318
Diskriminanzknoten . . . . .	319
Diskriminanzknoten – Modelloptionen . . . . .	320

Diskriminanzknoten – Expertenoptionen . . . . .	321
Diskriminanzknoten – Ausgabeoptionen . . . . .	322
Diskriminanzknoten – Schrittoptionen . . . . .	324
Diskriminanz-Modell-Nugget . . . . .	325
Diskriminanz-Modell-Nugget – Erweiterte Ausgabe . . . . .	326
Diskriminanz-Modell-Nugget – Einstellungen . . . . .	326
Diskriminanz-Modell-Nugget – Übersicht . . . . .	327
GenLin-Knoten . . . . .	328
Feldoptionen für den GenLin-Knoten . . . . .	330
Modelloptionen für den GenLin-Knoten . . . . .	331
Expertenoptionen für den GenLin-Knoten . . . . .	332
Verallgemeinerte lineare Modelle – Iterationen . . . . .	336
Verallgemeinerte lineare Modelle – Erweiterte Ausgabe . . . . .	337
GenLin-Modell-Nugget . . . . .	339
GenLin-Modell-Nugget – Erweiterte Ausgabe . . . . .	340
GenLin-Modell-Nugget – Einstellungen . . . . .	340
GenLin-Modell-Nugget – Übersicht . . . . .	341
Cox-Knoten . . . . .	342
Feldoptionen für Cox-Knoten . . . . .	343
Modelloptionen für Cox-Knoten . . . . .	344
Expertenoptionen für Cox-Knoten . . . . .	347
Einstellungsoptionen für Cox-Knoten . . . . .	351
Cox-Modell-Nugget . . . . .	352
Cox-Regression – Ausgabeeinstellungen . . . . .	353
Cox-Regression – Erweiterte Ausgabe . . . . .	353

## **11 Cluster-Modelle**

**354**

Kohonen-Knoten . . . . .	355
Optionen des Kohonen-Knotenmodells . . . . .	357
Expertenoptionen für den Kohonen-Knoten . . . . .	359
Modell-Nuggets vom Typ "Kohonen" . . . . .	361
Übersicht über das Kohonen-Modell . . . . .	361
K-Means-Knoten . . . . .	362
Optionen für K-Means-Knotenmodelle . . . . .	363
Expertenoptionen für K-Means-Knoten . . . . .	364
Modell-Nuggets vom Typ "K-Means" . . . . .	365
Übersicht über das K-Means-Modell . . . . .	366
TwoStep-Cluster-Knoten . . . . .	367
Optionen für TwoStep-Cluster-Knotenmodelle . . . . .	368

Modell-Nuggets vom Typ "TwoStep-Cluster" . . . . .	370
Übersicht über das TwoStep-Modell . . . . .	370
Der Cluster-Viewer . . . . .	371
Cluster-Viewer -Registerkarte Modell . . . . .	372
Navigieren in der Clusteranzeige . . . . .	382
Erzeugen von Diagrammen aus Clustermodellen . . . . .	385

## **12 Assoziationsregeln**

**388**

Tabellendaten im Vergleich zu Transaktionsdaten . . . . .	389
A Priori-Knoten . . . . .	390
Modelloptionen für den A Priori-Knoten . . . . .	391
Expertenoptionen für den A Priori-Knoten . . . . .	392
CARMA-Knoten . . . . .	394
Feldoptionen für den CARMA-Knoten . . . . .	395
Modelloptionen für den CARMA-Knoten . . . . .	398
Expertenoptionen für den CARMA-Knoten . . . . .	399
Assoziationsregelmodell-Nuggets . . . . .	400
Nähere Informationen zum Assoziationsregelmodell-Nugget . . . . .	401
Einstellungen beim Assoziationsregelmodell-Nugget . . . . .	407
Übersicht über das Assoziationsregelmodell-Nugget . . . . .	409
Generieren einer Regelmenge aus einem Assoziationsmodell-Nugget . . . . .	410
Erstellen eines gefilterten Modells . . . . .	411
Scoring von Assoziationsregeln . . . . .	412
Bereitstellung von Assoziationsmodellen . . . . .	413
Sequenzknoten . . . . .	416
Feldoptionen für den Sequenzknoten . . . . .	417
Modelloptionen für den Sequenzknoten . . . . .	419
Expertenoptionen für den Sequenzknoten . . . . .	420
Sequenzmodell-Nuggets . . . . .	422
Nähere Informationen zum Sequenzmodell-Nugget . . . . .	424
Sequenzmodell-Nugget - Einstellungen . . . . .	427
Sequenzmodell-Nugget-Übersicht . . . . .	427
Generieren eines Regel-Superknotens aus einem Sequenzmodell-Nugget . . . . .	428

## **13 Zeitreihenmodelle**

**431**

Wozu dienen Vorhersagen? . . . . .	431
Zeitreihendaten . . . . .	431
Merkmale von Zeitreihen . . . . .	432

Autokorrelation und partielle Autokorrelationsfunktionen . . . . .	436
Reihentransformationen . . . . .	437
Prädiktorreihen. . . . .	438
Zeitreihen – Modellierungsknoten . . . . .	438
Voraussetzungen . . . . .	440
Zeitreihenmodelle – Optionen . . . . .	442
Zeitreihen – Expert Modeler-Kriterien . . . . .	444
Zeitreihen – Kriterien für exponentielles Glätten. . . . .	446
Zeitreihen – ARIMA-Kriterien. . . . .	447
Transferfunktionen. . . . .	450
Umgang mit Ausreißern . . . . .	452
Generieren von Zeitreihenmodellen . . . . .	453
Generieren mehrerer Modelle . . . . .	453
Verwenden von Zeitreihenmodellen bei der Prognoseerstellung . . . . .	453
Erneute Schätzung und Vorhersage . . . . .	454
Zeitreihen-Modell-Nugget . . . . .	454
Zeitreihen – Modellparameter . . . . .	458
Zeitreihen – Modellresiduen . . . . .	459
Zeitreihen – Modellübersicht . . . . .	460
Zeitreihen – Modelleinstellungen. . . . .	461

## **14 Selbstlern-Antwortknotenmodelle 463**

SLRM-Knoten . . . . .	463
Feldoptionen für den SLRM-Knoten . . . . .	464
Modelloptionen für den SLRM-Knoten . . . . .	465
Einstellungsoptionen für den SLRM-Knoten . . . . .	467
SLRM-Modell-Nuggets . . . . .	469
SLRM-Modell – Einstellungen . . . . .	470

## **15 Modelle vom Typ “Support Vector Machine” 473**

Informationen zu SVM . . . . .	473
So funktioniert SVM: . . . . .	473
Feinabstimmung von SVM-Modellen. . . . .	475
SVM-Knoten. . . . .	476
Modelloptionen für SVM-Knoten . . . . .	477
Expertenoptionen für SVM-Knoten. . . . .	477

SVM-Modell-Nugget .....	479
Einstellungen beim SVM-Modell .....	480
<b>16 Nächste-Nachbarn-Modelle</b>	<b>482</b>
KNN-Knoten .....	482
Zieloptionen für KNN-Knoten .....	483
KNN-Knoten - Einstellungen .....	484
KNN-Modell-Nugget .....	494
Modellansicht .....	495
KNN-Modell-Einstellungen .....	502
 <b>Anhang</b>	
 <b>A Hinweise</b>	 <b>504</b>
 <b>Index</b>	 <b>507</b>



# **Informationen zu IBM SPSS Modeler**

IBM® SPSS® Modeler ist ein Set von Data-Mining-Tools, mit dem Sie auf der Grundlage Ihres Geschäftswissens schnell und einfach Vorhersagemodelle erstellen und zur Erleichterung der Entscheidungsfindung in die Betriebsabläufe einbinden können. SPSS Modeler, das auf der Grundlage des den Industrienormen entsprechenden Modells CRISP-DM entwickelt wurde, unterstützt den gesamten Data-Mining-Prozess, von den Daten bis hin zu besseren Geschäftsergebnissen.

SPSS Modeler bietet eine Vielzahl von Modellbildungsmethoden, die aus dem maschinellen Lernen, der künstlichen Intelligenz und der Statistik stammen. Mit den in der Modellierungspalette verfügbaren Methoden können Sie aus Ihren Daten neue Informationen ableiten und Vorhersagemodelle erstellen. Jede Methode besitzt ihre Stärken und eignet sich besonders für bestimmte Problemtypen.

SPSS Modeler kann als Standalone-Produkt oder in Verbindung mit SPSS Modeler Server erworben werden. Außerdem ist eine Reihe von Zusatzoptionen verfügbar, die in den folgenden Abschnitten kurz dargelegt werden. Weitere Informationen finden Sie unter <http://www.ibm.com/software/analytics/spss/products/modeler/>.

## **IBM SPSS Modeler Server**

SPSS Modeler verwendet eine Client/Server-Architektur zur Verteilung von Anforderungen für ressourcenintensive Vorgänge an leistungsstarke Serversoftware, wodurch bei größeren Daten-Sets eine schnellere Leistung erzielt werden kann. Neben den hier aufgeführten Produkten können auch weitere Produkte bzw. Aktualisierungen verfügbar sein. Weitere Informationen finden Sie unter <http://www.ibm.com/software/analytics/spss/products/modeler/>.

**SPSS Modeler.** SPSS Modeler Clementine Client ist eine im Funktionsumfang vollständige Version des installierten Produkts und kann auf dem Desktop-Computer des Benutzers ausgeführt werden. Es kann im lokalen Modus als Standalone-Produkt oder im verteilten Modus zusammen mit IBM® SPSS® Modeler Server verwendet werden, um im Falle von großen Daten-Sets die Leistung zu verbessern.

**SPSS Modeler Server.** SPSS Modeler Server wird ständig im verteilten Analysemodus zusammen mit einer oder mehreren IBM® SPSS® Modeler-Installationen ausgeführt, wodurch eine herausragende Leistung bei großen Daten-Sets erzielt werden kann, da speicherintensive Operationen auf dem Server durchgeführt werden können, ohne Daten auf den Client-Computer herunterzuladen. SPSS Modeler Server bietet außerdem Unterstützung für SQL-Optimierung sowie Funktionen zur Modellierung innerhalb der Datenbank, wodurch Leistungsfähigkeit und Automatisierung weiter verbessert werden. Es muss mindestens eine SPSS Modeler-Installation vorhanden sein, um eine Analyse durchzuführen.

## **IBM SPSS Modeler-Optionen**

Die folgenden Komponenten und Funktionen können separat erworben und für die Verwendung mit SPSS Modeler lizenziert werden. Beachten Sie, dass zu einem späteren Zeitpunkt möglicherweise noch weitere Produkte und Updates erhältlich sind. Weitere Informationen finden Sie unter <http://www.ibm.com/software/analytics/spss/products/modeler/>.

- SPSS Modeler Server-Zugriff, der eine bessere Skalierbarkeit und bessere Leistungsfähigkeit bei großen Daten-Sets bietet sowie Unterstützung für SQL-Optimierung und Funktionen zur Modellierung innerhalb der Datenbank.
- SPSS Modeler Solution Publisher, für Scoring in Echtzeit oder automatisiertes Scoring außerhalb der SPSS Modeler-Umgebung. Für weitere Informationen siehe Thema [IBM SPSS Modeler Solution Publisher in Kapitel 2 in IBM SPSS Modeler 14.2 Solution Publisher](#).
- Adapter zum Deployment für IBM SPSS Collaboration and Deployment Services oder die Thin-Client-Anwendung IBM SPSS Modeler Advantage. Für weitere Informationen siehe Thema [Speichern und Bereitstellen von IBM SPSS Collaboration and Deployment Services Repository-Objekten in Kapitel 9 in IBM SPSS Modeler 14.2- Benutzerhandbuch](#).

## **IBM SPSS Text Analytics**

IBM® SPSS® Text Analytics ist ein vollständig integriertes Zusatzprodukt für SPSS Modeler, das hoch entwickelte linguistische Technologien und die Verarbeitung natürlicher Sprache (Natural Language Processing, NLP) benutzt, um eine schnelle Verarbeitung einer großen Vielfalt an unstrukturierten Textdaten zu ermöglichen, um die Schlüsselkonzepte zu extrahieren und zu ordnen und um diese Konzepte in Kategorien zusammenzufassen. Extrahierte Konzepte und Kategorien können mit bestehenden strukturierten Daten, beispielsweise demografischen Informationen, kombiniert und mithilfe der vollständigen Suite der Data-Mining-Tools von IBM® SPSS® Modeler auf die Modellierung angewendet werden, um bessere und fokussiertere Entscheidungen zu ermöglichen.

- Der Text-Mining-Knoten bietet die Modellierung von Konzepten und Kategorien sowie eine interaktive Workbench, in der Sie eine erweiterte Untersuchung von Textlinks und Clustern durchführen, Ihre eigenen Kategorien erstellen und die Vorlagen für linguistische Ressourcen verfeinern können.
- Eine Reihe von Importformaten wird unterstützt, darunter Blogs und andere webbasierte Quellen.
- Benutzerdefinierte Vorlagen, Bibliotheken und Wörterbücher für bestimmte Domänen, wie CRM und Genomforschung, sind ebenfalls eingeschlossen.

*Anmerkung:* Für den Zugriff auf diese Komponente ist eine separate Lizenz erforderlich. Weitere Informationen finden Sie unter <http://www.ibm.com/software/analytics/spss/products/modeler/>.

## **IBM SPSS Modeler-Dokumentation**

Die vollständige Dokumentation im Online-Hilfe-Format finden Sie im Hilfe-Menü von SPSS Modeler. Dazu gehören die Dokumentation für SPSS Modeler, SPSS Modeler Server und SPSS Modeler Solution Publisher sowie das Anwendungshandbuch und weiteres Material zur Unterstützung.

Die vollständige Dokumentation für die einzelnen Produkte im PDF-Format finden Sie im Ordner *Documentation* auf der jeweiligen Produkt-DVD.

- **IBM SPSS Modeler-Benutzerhandbuch.** Allgemeine Einführung in die Verwendung von SPSS Modeler, in der u. a. die Erstellung von Daten-Streams, der Umgang mit fehlenden Werten, die Erstellung von CLEM-Ausdrücken, die Arbeit mit Projekten und Berichten sowie das Packen von Streams für das Deployment in IBM SPSS Collaboration and Deployment Services, Predictive Applications (Prognoseanwendungen) oder IBM SPSS Modeler Advantage beschrieben werden.
- **Quellen-, Prozess- und Ausgabeknoten in IBM SPSS Modeler.** Beschreibung aller Knoten, die zum Lesen, zum Verarbeiten und zur Ausgabe von Daten in verschiedenen Formaten verwendet werden. Im Grunde sind sie alle Knoten, mit Ausnahme der Modellierungsknoten.
- **IBM SPSS Modeler Modellierungsknoten.** Beschreibungen sämtlicher für die Erstellung von Data Mining-Modellen verwendeter Knoten. IBM® SPSS® Modeler bietet eine Vielzahl von Modellbildungsmethoden, die aus dem maschinellen Lernen, der künstlichen Intelligenz und der Statistik stammen. [Für weitere Informationen siehe Thema Überblick über Modellierungsknoten in Kapitel 3 auf S. 24.](#)
- **IBM SPSS Modeler-Algorithmushandbuch.** Beschreibung der mathematischen Grundlagen der in SPSS Modeler verwendeten Modellierungsmethoden.
- **IBM SPSS Modeler-Anwendungshandbuch.** Die Beispiele in diesem Handbuch bieten eine kurze, gezielte Einführung in bestimmte Modellierungsmethoden und -verfahren. Eine Online-Version dieses Handbuchs kann auch über das Hilfe-Menü aufgerufen werden. [Für weitere Informationen siehe Thema Anwendungsbeispiele in IBM SPSS Modeler 14.2-Benutzerhandbuch.](#)
- **Skripterstellung und Automatisierung in IBM SPSS Modeler.** Informationen zur Automatisierung des Systems über Skripterstellung, einschließlich der Eigenschaften, die zur Bearbeitung von Knoten und Streams verwendet werden können.
- **IBM SPSS Modeler Deployment-Handbuch.** Informationen zum Ausführen von SPSS Modeler-Streams und -Szenarien als Schritte bei der Verarbeitung von Jobs im IBM® SPSS® Collaboration and Deployment Services Deployment Manager.
- **IBM SPSS Modeler CLEF-Entwicklerhandbuch.** CLEF bietet die Möglichkeit, Drittanbieterprogramme, wie Datenverarbeitungsroutinen oder Modellierungsalgorithmen, als Knoten in SPSS Modeler zu integrieren.
- **In-Database Mining-Handbuch für IBM SPSS Modeler.** Informationen darüber, wie Sie Ihre Datenbank dazu einsetzen, die Leistung zu verbessern, und wie Sie die Palette der Analysefunktionen über Drittanbieteralgorithmen erweitern.
- **IBM SPSS Modeler Server- und -Leistungshandbuch.** Informationen zur Konfiguration und Verwaltung von IBM® SPSS® Modeler Server.

- **IBM SPSS Modeler Administration Console – Benutzerhandbuch.** Informationen zur Installation und Nutzung der Konsolen-Benutzeroberfläche zur Überwachung und Konfiguration von SPSS Modeler Server. Die Konsole ist als Plugin für die Deployment Manager-Anwendung implementiert.
- **IBM SPSS Modeler Solution Publisher-Handbuch.** SPSS Modeler Solution Publisher ist eine Zusatzkomponente, mit der Unternehmen Streams zur Verwendung außerhalb der SPSS Modeler-Standardumgebung veröffentlichen können.
- **IBM SPSS Modeler-Handbuch zu CRISP-DM.** Schritt-für-Schritt-Anleitung für das Data-Mining mit SPSS Modeler unter Verwendung der CRISP-DM-Methode.

## **Anwendungsbeispiele**

Mit den Data-Mining-Tools in SPSS Modeler kann eine große Bandbreite an geschäfts- und unternehmensbezogenen Problemen gelöst werden; die Anwendungsbeispiele dagegen bieten jeweils eine kurze, gezielte Einführung in spezielle Modellierungsmethoden und -verfahren. Die hier verwendeten Daten-Sets sind viel kleiner als die riesigen Datenbestände, die von einigen Data-Mining-Experten verwaltet werden müssen, die zugrunde liegenden Konzepte und Methoden sollten sich jedoch auch auf reale Anwendungen übertragen lassen.

Sie können auf die Beispiele zugreifen, indem Sie im Menü “Hilfe” in SPSS Modeler auf die Option Anwendungsbeispiele klicken. Die Datendateien und Beispiel-Streams wurden im Ordner *Demos*, einem Unterordner des Produktinstallationsverzeichnis, installiert. [Für weitere Informationen siehe Thema Ordner “Demos” in IBM SPSS Modeler 14.2- Benutzerhandbuch.](#)

**Beispiele für die Datenbank-Modellierung.** Die Beispiele finden Sie im *IBM SPSS Modeler In-Database Mining-Handbuch*.

**Skriptbeispiele.** Die Beispiele finden Sie im *IBM SPSS Modeler Handbuch für die Skripterstellung und Automatisierung*.

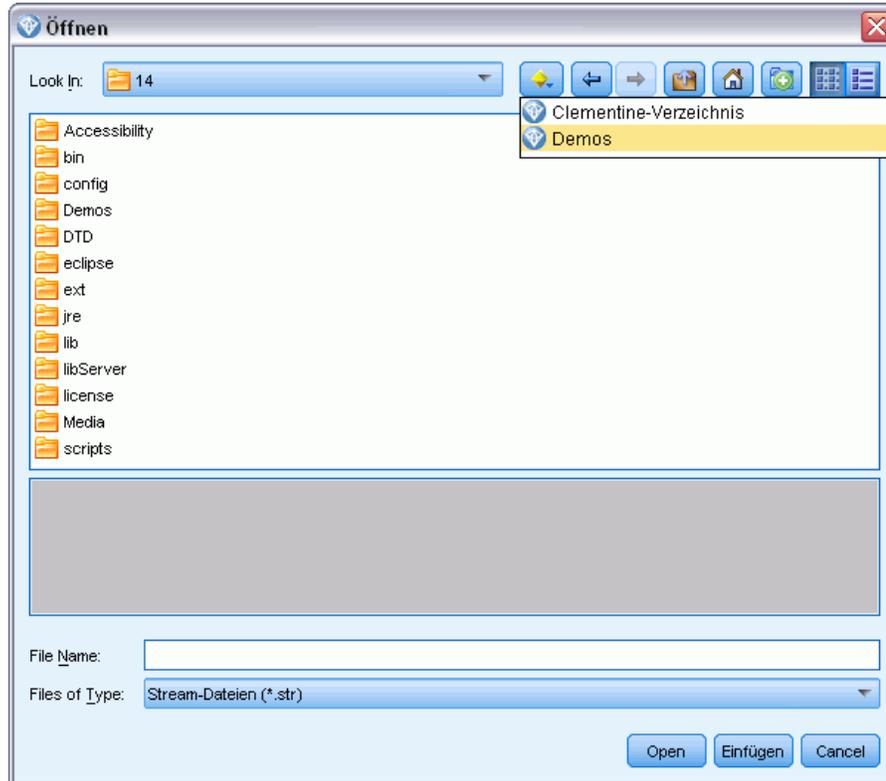
## **Ordner “Demos”**

Die in den Anwendungsbeispielen verwendeten Datendateien und Beispiel-Streams wurden im Ordner *Demos*, einem Unterordner des Produktinstallationsverzeichnis, installiert. Auf diesen Ordner können Sie auch über die Programmgruppe IBM SPSS Modeler 14.2 im

Windows-Startmenü oder durch Klicken auf *Demos* in der Liste der zuletzt angezeigten Verzeichnisse im Dialogfeld "Datei öffnen" zugreifen.

**Abbildung 1-1**

*Auswahl des Ordners "Demos" in der Liste der zuletzt angezeigten Verzeichnisse*

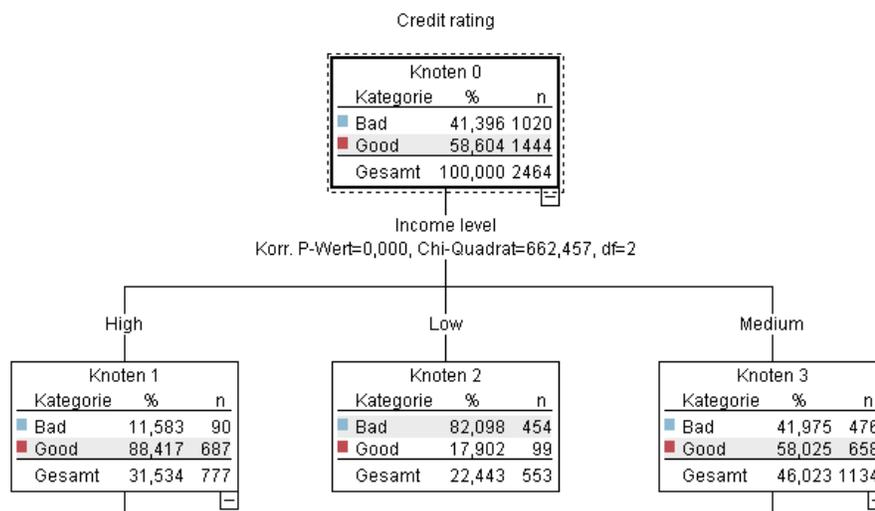


# Einführung in die Modellierung

Ein Modell ist eine Menge von Regeln, Formeln bzw. Gleichungen, mit der ein Ergebnis auf der Grundlage einer Menge von Eingabefeldern bzw. -variablen vorhergesagt werden kann. Eine Finanzinstitut verwendet z. B. möglicherweise ein Modell, um auf Basis von bekannten Informationen über vorherige Kreditantragsteller vorherzusagen, ob ein Kreditantragsteller ein geringes oder hohes Risiko darstellt.

Die Möglichkeit zur Vorhersage eines Ergebnisses ist das zentrale Ziel von Prognoseanalysen und ein Verständnis des Modellierungsprozesses ist der Schlüssel für die Verwendung von IBM® SPSS® Modeler.

Abbildung 2-1  
Ein einfaches Entscheidungsbaum-Modell



In diesem Beispiel wird ein **Entscheidungsbaum**-Modell verwendet, das Datensätze aufzeichnet (und eine Reaktion vorhersagt), wobei eine Reihe von Entscheidungsregeln verwendet wird wie beispielsweise:

IF income = Medium  
AND cards <5  
THEN -> 'Good'

In diesem Beispiel wird zwar ein Modell vom Typ "CHAID" (Chi-squared Automatic Interaction Detection) verwendet, es ist jedoch als allgemeine Einführung gedacht und die meisten Konzepte gelten im Wesentlichen auch für andere Modellierungstypen in SPSS Modeler.

Um ein Modell zu verstehen, müssen Sie zunächst ein Verständnis für die darin verwendeten Daten entwickeln. Die Daten in diesem Beispiel enthalten Informationen über die Kunden einer Bank. Es werden folgende Felder verwendet:

Feldname	Beschreibung
Credit_rating	Kreditrating: 0 = Schlecht, 1 = Gut, 9 = fehlende Werte
Alter	Alter in Jahren
Einkommen	Einkommenstufe: 1 = Niedrig, 2 = Mittel, 3 = Hoch
Credit_cards	Anzahl der Kreditkarten: 1 = Weniger als fünf, 2 = Fünf oder mehr
Bildung	Bildungsniveau: 1 = Hauptschulabschluss, 2 = Hochschulabschluss
Car_loans	Anzahl der Autokredite: 1 = Keine oder einen, 2 = Mehr als zwei

Die Bank führt eine Datenbank historischer Informationen über Kunden, die bei der Bank Kredite in Anspruch genommen haben, in der auch festgehalten wird, ob ein Kredit zurückgezahlt wurde (Bonität = Gut) oder nicht (Bonität = Schlecht). Mithilfe dieser vorhandenen Daten will die Bank ein Modell erstellen, das vorhersagen kann, mit welcher Wahrscheinlichkeit zukünftige Kreditantragsteller ihren Kreditverpflichtungen nicht nachkommen.

Anhand eines Entscheidungsbaum-Modells können Sie die Charakteristiken der beiden Kundengruppen analysieren und die Wahrscheinlichkeit von Kreditausfällen vorhersagen.

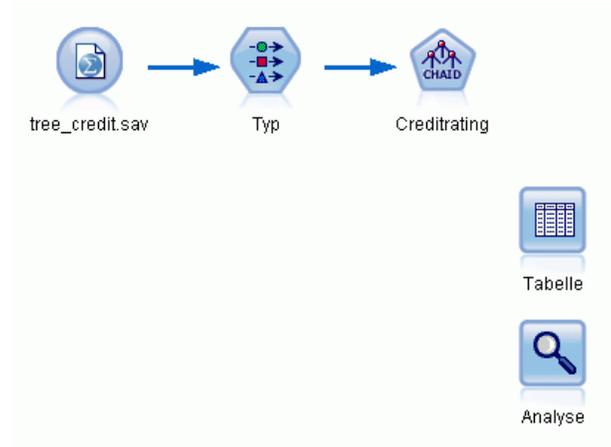
Für dieses Beispiel wird der Stream *modelingintro.str* verwendet, der im Ordner *Demos* unter dem Unterordner *streams* verfügbar ist. Die Datendatei ist *tree\_credit.sav*. [Für weitere Informationen siehe Thema Ordner "Demos" in Kapitel 1 in IBM SPSS Modeler 14.2- Benutzerhandbuch.](#)

Werfen wir nun einen Blick auf den Stream.

- ▶ Wählen Sie im Hauptmenü folgende Menüoption:  
Datei > Stream öffnen
- ▶ Klicken Sie auf der Symbolleiste des Dialogfelds "Öffnen" auf das Gold-Nugget-Symbol und wählen Sie den Ordner "Demos" aus.
- ▶ Doppelklicken Sie auf den Ordner *streams*.
- ▶ Doppelklicken Sie auf die Datei *modelingintro.str*.

## Erstellen des Streams

Abbildung 2-2  
Modellierungs-Stream



Um einen Stream zum Erzeugen eines Modells zu erstellen, benötigen Sie mindestens drei Elemente:

- Ein Quellenknoten, der Daten aus einer externen Quelle einliest, in diesem Fall eine IBM® SPSS® Statistics-Datendatei.
- Ein Quellen- oder Typknoten, der Feldeigenschaften wie beispielsweise das Messniveau (die Daten, die das Feld enthält) und die Rolle der einzelnen Felder als Ziel oder Eingabe in der Modellierung angibt.
- Ein Modellierungsknoten, der bei Ausführung des Streams ein Modell-Nugget erstellt.

In diesem Beispiel verwenden wir einen CHAID-Modellierungsknoten. CHAID (Chi-squared Automatic Interaction Detection) ist eine Klassifizierungsmethode für die Erstellung von Entscheidungsbäumen mit bestimmten Statistiktypen namens Chi-Quadrat-Statistiken zur Identifizierung der optimalen Splits.

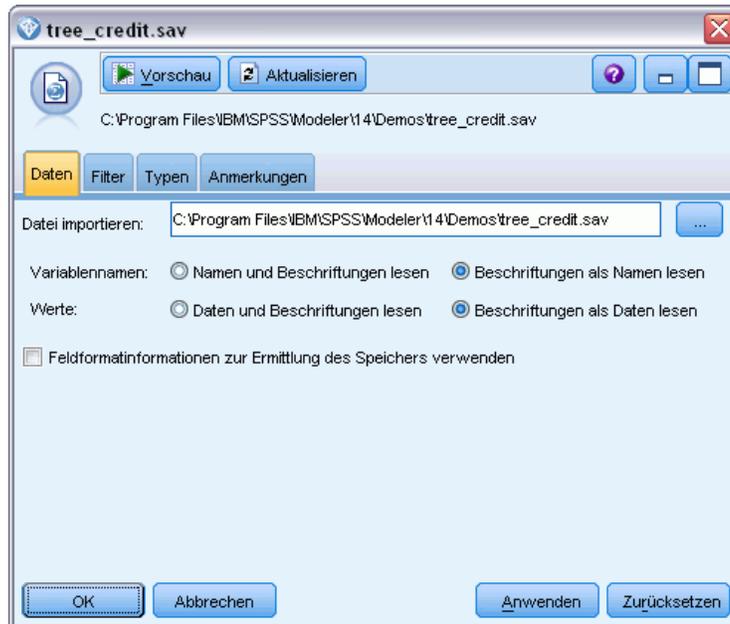
Wenn Messniveaus im Quellenknoten angegeben sind, kann auf den separaten Typknoten verzichtet werden. Hinsichtlich der Funktion ist das Ergebnis dasselbe.

Dieser Stream weist außerdem Tabellen- und Analyseknöten auf, mit denen die Scoring-Ergebnisse angezeigt werden, nachdem das Modell-Nugget erstellt und in den Stream aufgenommen wurde.

Der Statistikdatei-Quellenknoten liest Daten in SPSS Statistics-Format aus der Datendatei *tree\_credit.sav* ein, die im Ordner *Demos* installiert wurde. (Eine spezielle Variable mit der Bezeichnung *\$CLEO\_DEMOS* dient zur Referenzierung diese Ordners in der aktuellen IBM®

SPSS® Modeler-Installation. Dadurch wird sichergestellt, dass der Pfad gültig ist, unabhängig vom aktuellen Installationsordner bzw. der jeweiligen Version.)

Abbildung 2-3  
Einlesen von Daten mit einem Statistikdateiquellenknoten



Der Typknoten gibt das **Messniveau** für die einzelnen Felder an. Das Messniveau ist eine Kategorie, die den Datentyp für das Feld anzeigt. Unsere Quelldatendatei verwendet drei verschiedene Messniveaus.

Ein Feld des Typs **Stetig** (z. B. das Feld *Alter*) enthält stetige numerische Werte, während ein Feld des Typs **Nominal** (z. B. das Feld *Kreditrating*) zwei oder mehr bestimmte Werte enthält, z. B. *Schlecht*, *Gut* oder *Keine früheren Schulden*. Ein Feld des Typs **Ordinal** (z. B. *Einkommen*

in Kategorien) beschreibt Daten mit mehreren unterschiedlichen Werten, die eine natürliche Reihenfolge aufweisen — in diesem Fall *Niedrig*, *Mittel* und *Hoch*.

Abbildung 2-4  
Festlegen des Ziels und der Eingabefelder mit dem Typknoten



Der Typknoten legt für jedes Feld außerdem die **Rolle** fest, die jedes Feld bei der Modellierung spielt. Für das Feld *Kreditrating*, das angibt, ob ein bestimmter Kunde seinen Kreditverpflichtungen nicht nachgekommen ist, ist die Rolle als *Ziel* festgelegt. Hierbei handelt es sich also um das **Ziel** oder das Feld, für das wir den Wert vorhersagen möchten.

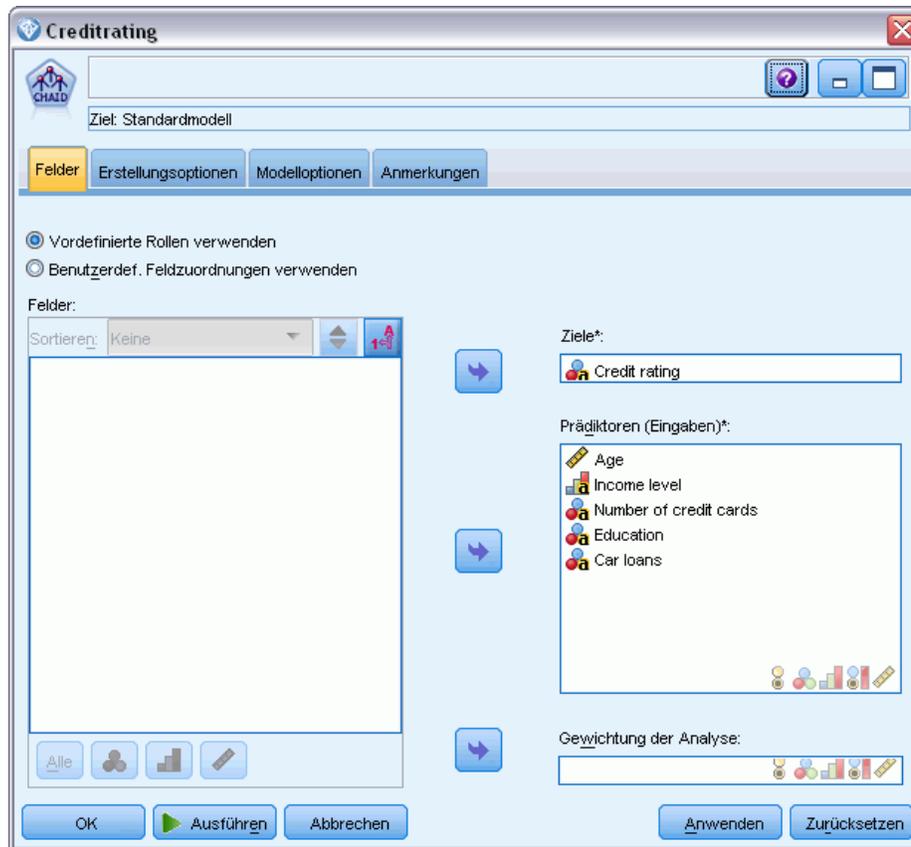
Für die anderen Felder ist die Rolle auf *Eingabe* eingestellt. Eingabefelder werden manchmal auch als **Prädiktoren** bezeichnet oder als Felder, mit deren Werten der Modellierungsalgorithmus den Wert des Zielfelds vorhersagt.

Der CHAID-Modellierungsknoten generiert das Modell.

Auf der Registerkarte "Felder" im Modellierungsknoten wird die Option Vordefinierte Rollen verwenden ausgewählt. Dies bedeutet, dass die im Typknoten angegebenen Ziele und Eingaben verwendet werden sollen. Wir können die Feldrollen hier ändern, doch in diesem Beispiel belassen wir sie unverändert.

- Klicken Sie auf die Registerkarte “Erstellungsoptionen”.

Abbildung 2-5  
CHAID-Modellierungsknoten, Registerkarte “Felder”



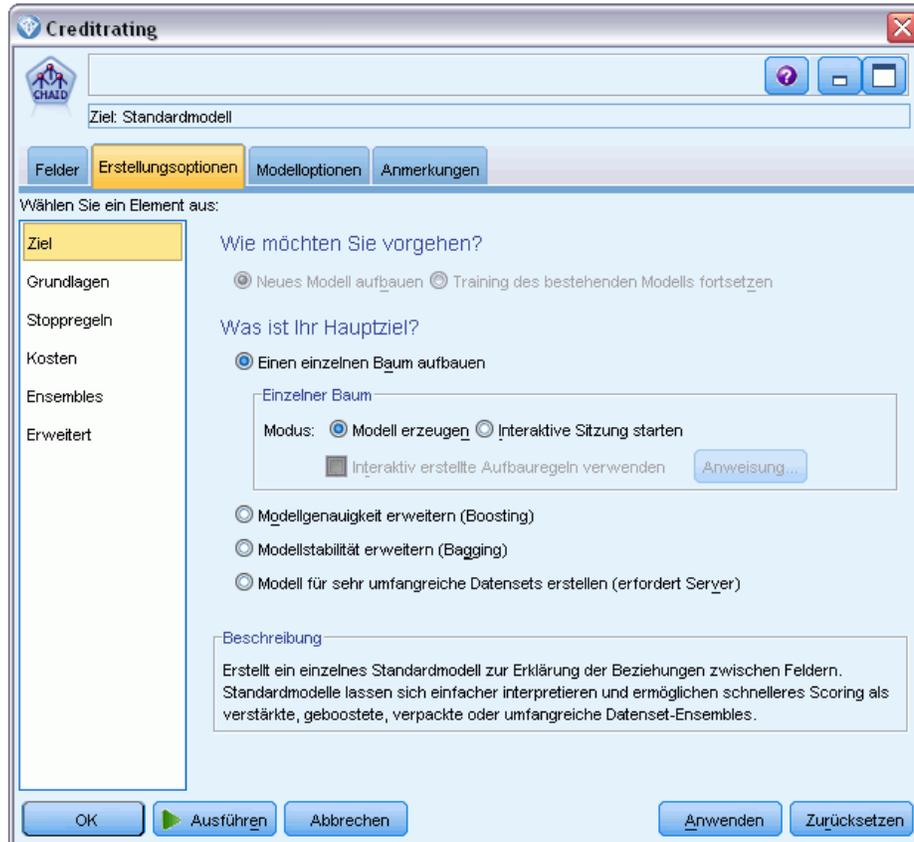
Hier finden Sie einige Optionen, in denen Sie die Art des aufzubauenden Modells festlegen können.

Da wir ein komplett neues Modell möchten, verwenden wir die Standardoption Neues Modell aufbauen.

Außerdem möchten wir nur ein einzelnes Standardentscheidungsbaummodell ohne Erweiterungen, weshalb wir auf die Standardzieloption Einzelnen Baum aufbauen zurückgreifen.

Sie können optional eine interaktive Modellierungssitzung starten, mit der Sie eine Feinabstimmung des Modells vornehmen können. Im vorliegenden Beispiel wird jedoch einfach ein Modell mit der Standardmuseinstellung Modell erzeugen generiert.

Abbildung 2-6  
CHAID-Modellierungsknoten, Registerkarte "Erstellungsoptionen"



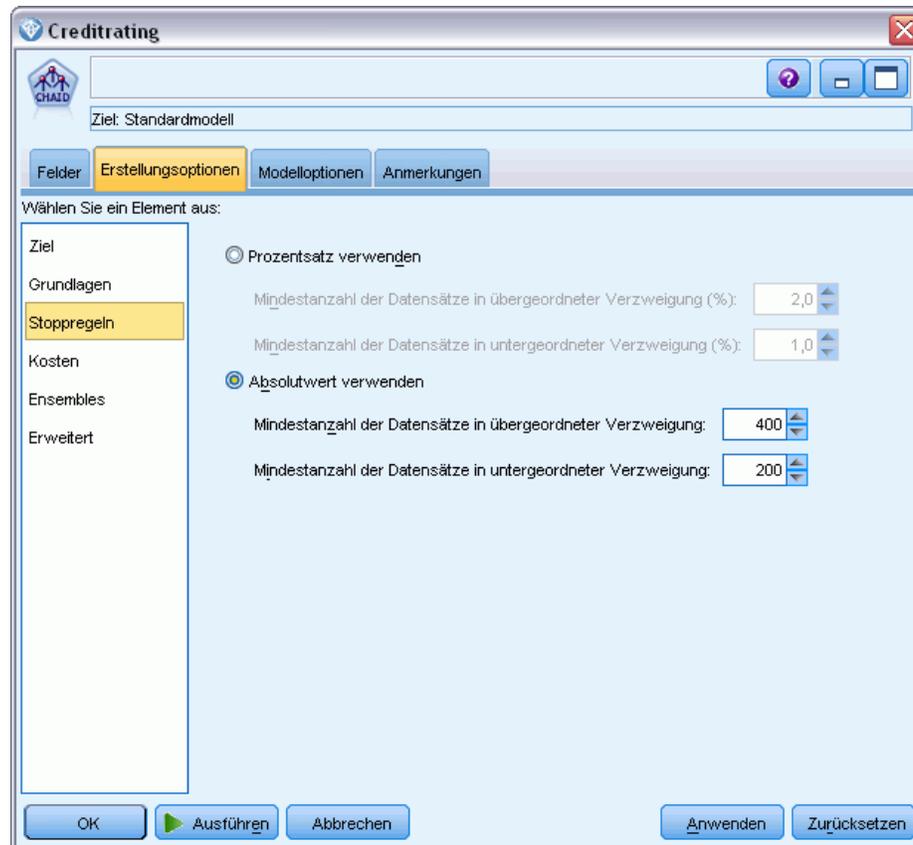
Für dieses Beispiel möchten wir einen einfach strukturierten Baum verwenden und begrenzen deshalb die Baumerweiterung, indem wir die minimale Anzahl der Fälle für über- und untergeordnete Knoten erhöhen.

- ▶ Wählen Sie auf der Registerkarte "Erstellungsoptionen" im linken Navigationsbereich Stoppregelein aus.
- ▶ Wählen Sie die Option Absolutwert verwenden aus.
- ▶ Legen Sie für Mindestanzahl der Datensätze in übergeordneter Verzweigung 400 fest.

- Legen Sie für Mindestanzahl der Datensätze in untergeordneter Verzweigung 200 fest.

Abbildung 2-7

Festlegen der Grenzkriterien beim Erstellen von Entscheidungsbäumen



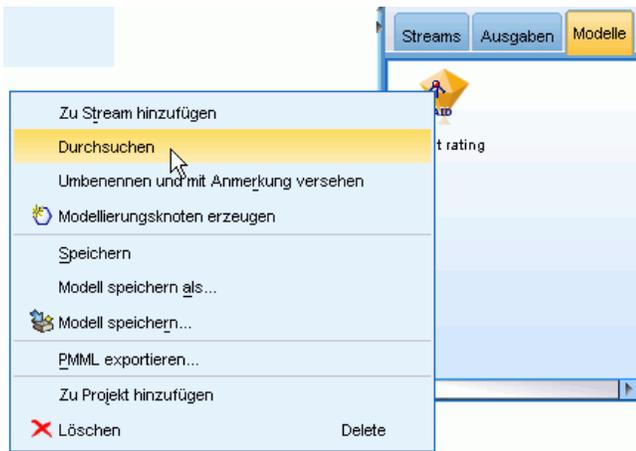
Wir können in diesem Beispiel alle anderen Standardoptionen verwenden, klicken Sie also auf Ausführen, um das Modell zu erstellen. (Alternativ können Sie mit der rechten Maustaste auf den Knoten klicken und im Kontextmenü Ausführen auswählen oder den Knoten auswählen und Ausführen im Menü "Extras" wählen.)

## Durchsuchen des Modells

Nach Abschluss der Ausführung wird das Modell-Nugget der Modellpalette rechts oben im Anwendungsfenster hinzugefügt. Zusätzlich wird es in der Stream-Zeichenfläche mit einer Verknüpfung zum Modellierungsknoten gezeigt, von dem aus es erstellt wurde. Um die

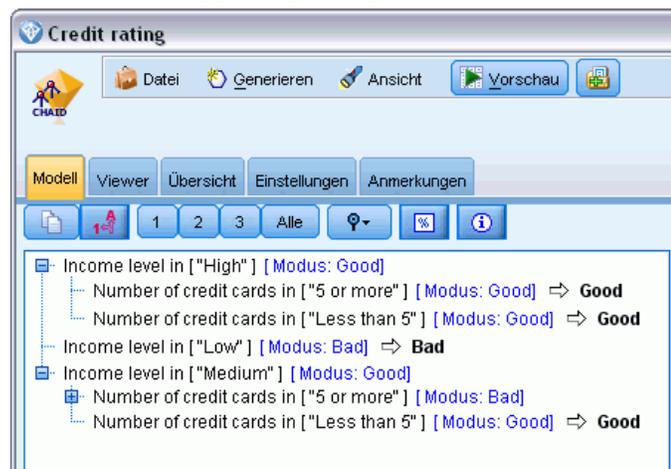
Modelldetails anzuzeigen, klicken Sie mit der rechten Maustaste auf den generierten Modellknoten und wählen Durchsuchen (in der Modellpalette) oder Bearbeiten (in der Zeichenfläche).

Abbildung 2-8  
Modellpalette



Im Fall des CHAID-Nuggets zeigt die Registerkarte "Modell" die Details in Form einer Regelmenge. Im Wesentlichen handelt es sich hierbei um eine Reihe von Regeln, die dazu verwendet werden können, einzelne Datensätze untergeordneten Knoten basierend auf den Werten verschiedener Eingabefelder zuzuweisen.

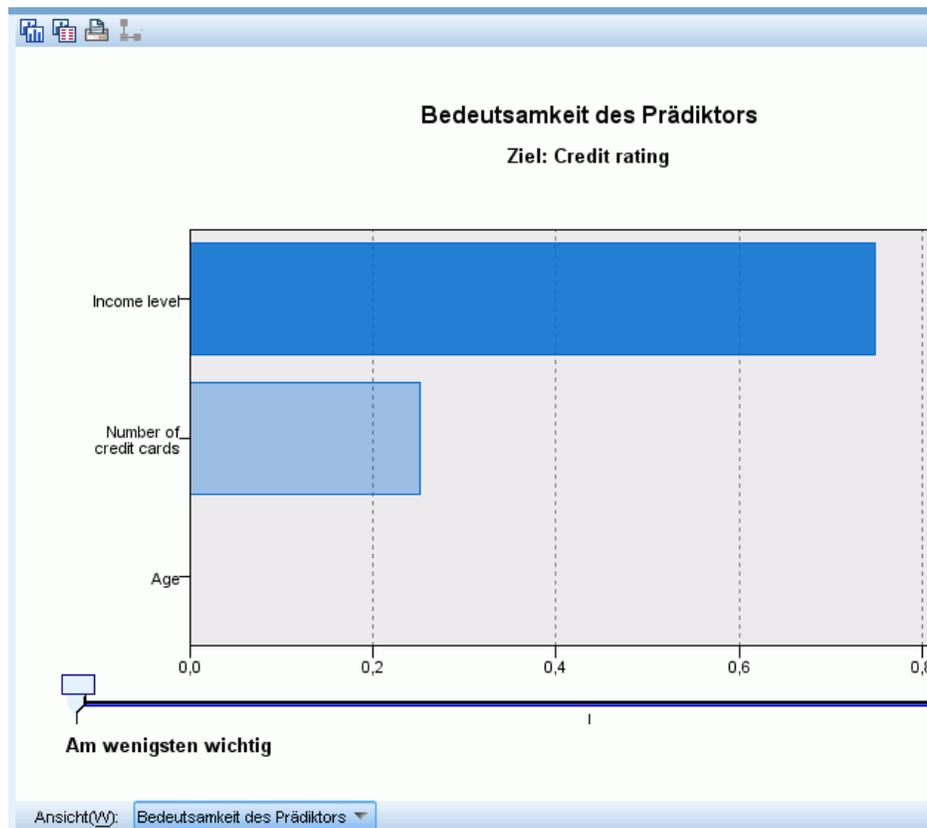
Abbildung 2-9  
CHAID-Modell-Nugget, Regelmenge



Für jeden Entscheidungsbaum-Endknoten (also die Baumknoten, die nicht weiter aufgeteilt werden) wird die Vorgersage *Gut* oder *Schlecht* getroffen. In jedem Fall wird die Vorhersage für Datensätze, die unter diesen Knoten fallen, durch den **Modus** bestimmt, also durch die häufigste Antwort.

Rechts neben der Regelmenge zeigt die Registerkarte "Modell" das Diagramm "Bedeutsamkeit der Prädiktoren", das die relative Wichtigkeit jedes Prädiktors beim Schätzen des Modells zeigt. Das zeigt uns, dass die *Einkommen in Kategorien* in diesem Fall eindeutig die größte Bedeutung hat, und dass der einzige andere bedeutsame Faktor die *Anzahl der Kreditkarten* ist.

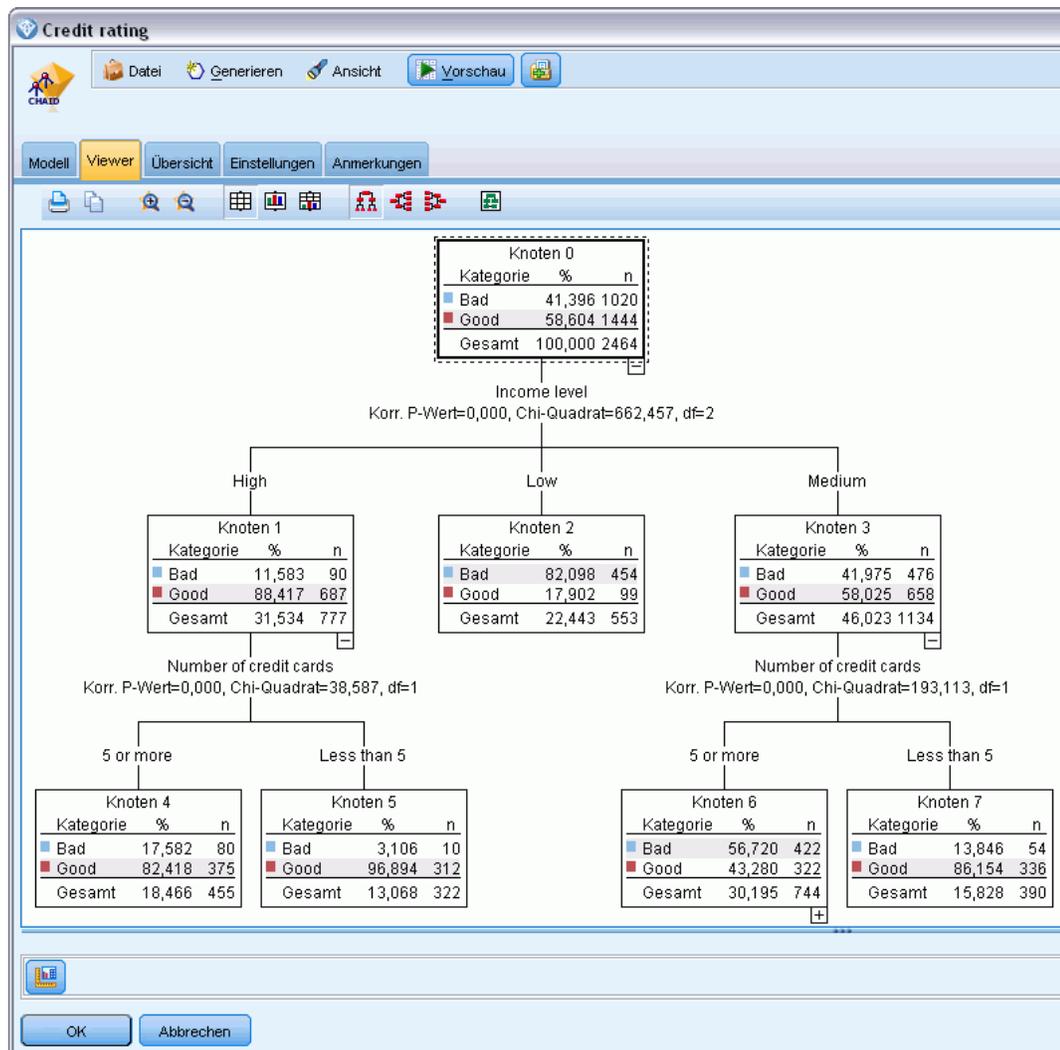
Abbildung 2-10  
Bedeutsamkeit der Prädiktoren – Diagramm



Auf der Registerkarte "Viewer" im Modell-Nugget wird dasselbe Modell in Form eines Baums angezeigt, mit einem Knoten bei jedem Entscheidungspunkt. Mit den Zoom-Steuerelementen auf der Symbolleiste können Sie die Ansicht eines bestimmten Knotens vergrößern bzw. die Ansicht verkleinern, um einen größeren Ausschnitt aus dem Baum zu sehen.

Abbildung 2-11

Registerkarte "Viewer" im Modell-Nugget, "Verkleinern" ausgewählt



Im oberen Teil des Baums fasst der erste Knoten (Knoten 0) alle Datensätze im Daten-Set zusammen. Knapp über 40 % der Fälle im Daten-Set sind als hochriskant eingestuft. Da dieser Anteil ziemlich hoch ist, interessiert es uns, ob der Baum Informationen darüber enthält, welche Faktoren dafür verantwortlich sind.

Wie wir sehen, findet die erste Aufteilung bei der *Einkommen in Kategorien* statt Datensätze, bei denen die Einkommensstufe in der Kategorie *Niedrig* liegt, werden Knoten 2 zugewiesen. Entsprechend enthält diese Kategorie den höchsten Prozentsatz an Kreditausfällen. Die Kreditvergabe an Kunden in dieser Kategorie bringt offensichtlich ein hohes Risiko mit sich.

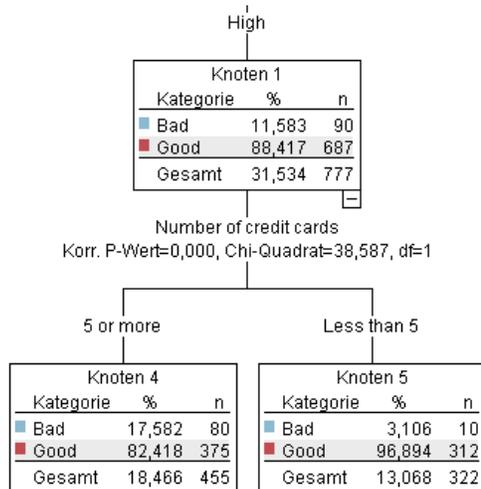
Bei 16 % der Kunden in dieser Kategorie kam es allerdings *nicht* zum Kreditausfall, die Prognose ist stimmt also nicht in jedem Fall. Kein Modell kann jede Antwort korrekt vorhersagen, aber ein gutes Modell sollte es ermöglichen, die auf der Grundlage der verfügbaren Daten *wahrscheinlichste* Antwort für die einzelnen Datensätze vorherzusagen.

Wenn wir die Kunden mit hohem Einkommen betrachten (Knoten 1), ist das Risiko bei der überwiegenden Mehrheit (89 %) entsprechend gering. Aber mehr als 1 aus 10 dieser Kunden ist ebenfalls seinen Kreditverpflichtungen nicht nachgekommen. Ist es möglich, die Kreditvergabekriterien zu verfeinern, um das Risiko zu minimieren?

Wie Sie sehen, hat das Modell diese Kunden auf Basis der Anzahl ihrer Kreditkarten in zwei Unterkategorien (Knoten 4 und 5) aufgeteilt. Wenn wir Kredite nur an Kunden mit hohem Einkommen vergeben, die weniger als 5 Kreditkarten besitzen, können wir unsere Erfolgsquote von 89 % auf 97 % erhöhen und somit ein noch zufriedenstellenderes Ergebnis erzielen.

Abbildung 2-12

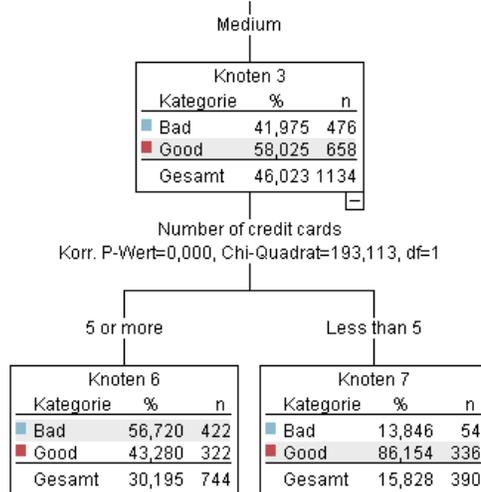
Baumansicht der Kunden mit hohem Einkommen



Aber was ist mit den Kunden in der Kategorie mit mittlerem Einkommen (Knoten 3)? Die Verteilung auf gute und schlechte Bonität fällt bei ihnen viel gleichmäßiger aus.

Auch hier sind wieder die Unterkategorien (in diesem Fall Knoten 6 und 7) sehr hilfreich. Wenn wir Kredite nur an Kunden mit mittlerem Einkommen vergeben, die weniger als 5 Kreditkarten besitzen, können wir unsere Erfolgsquote wieder von 58 % auf 85 % erhöhen und somit ein noch zufriedenstellenderes Ergebnis erzielen.

Abbildung 2-13  
Baumansicht der Kunden mit mittlerem Einkommen



Wir haben gesehen, dass jeder Datensatz, der in diesem Modell verarbeitet wird, einem spezifischen Knoten und der Prognose *Gut* oder *Schlecht* zugewiesen wird, je nachdem, welche die häufigste Antwort für den jeweiligen Knoten ist.

Dieser Vorgang der Zuweisung von Vorhersagen zu einzelnen Datensätzen wird als **Scoring** bezeichnet. Indem wir die Datensätze scoren, die auch zur Schätzung des Modells verwendet wurden, können wir evaluieren, mit welcher Genauigkeit das Modell für die Trainingsdaten (die Daten, für die das Ergebnis berechnet werden soll) funktioniert. Sehen wir uns an, wie das funktioniert.

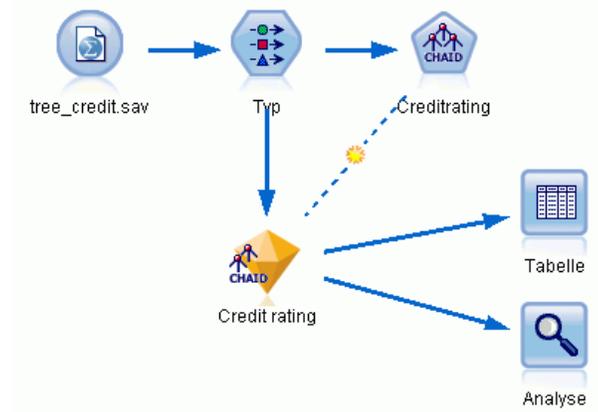
## Bewertung des Modells

Wir haben das Modell durchsucht, um zu verstehen, wie das Scoring funktioniert. Aber um zu evaluieren, *mit welcher Genauigkeit* es funktioniert, müssen wir einige Datensätze scoren und die vom Modell vorhergesagten Ergebnisse mit den tatsächlichen Ergebnissen vergleichen. Nun

werden wir dieselben Datensätze bewerten, die zum Schätzen des Modells verwendet wurden, und können damit die beobachteten und vorhergesagten Antworten vergleichen.

Abbildung 2-14

Anhängen des Modell-Nuggets an Ausgabeknoten zur Modellevaluation



- Fügen Sie zur Anzeige der Scores bzw. Vorhersagen den Tabellenknoten zum Modell-Nugget hinzu, doppelklicken Sie auf den Tabellenknoten und klicken Sie auf Ausführen.

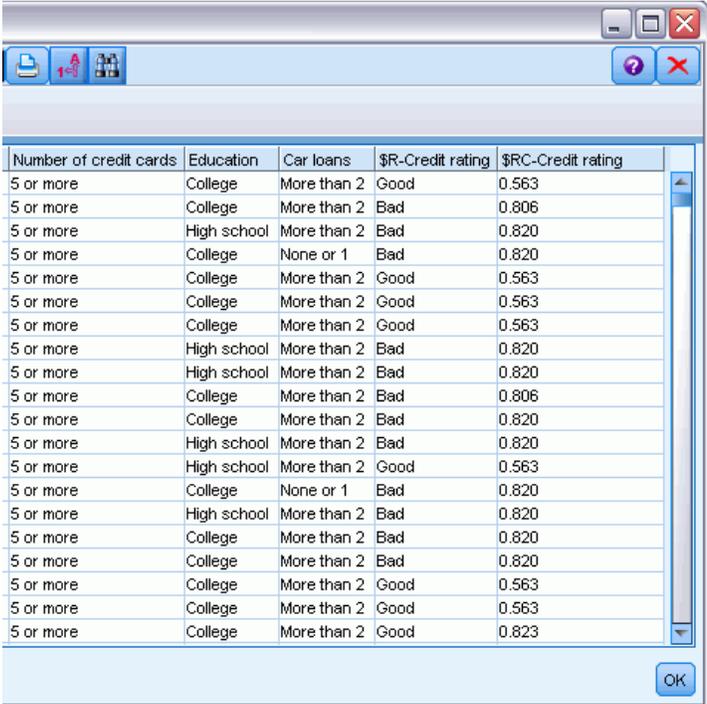
Die Tabelle zeigt die vorhergesagten Scores unter einem Feldnamen (*\$R-Credit rating*) an, der vom Modell erstellt wurde. Wir können diese Werte mit dem ursprünglichen Feld *Kreditrating* vergleichen, das die tatsächlichen Antworten enthält.

Gemäß der Konvention beruhen die Namen der während des Scorens generierten Felder auf dem Zielfeld, tragen jedoch ein Standard-Präfix, wie beispielsweise *\$R-* für Vorhersagen oder *\$RC-* für Konfidenzwerte. Verschiedene Modelltypen verwenden verschiedene Präfix-Sets. Ein

**Konfidenzwert** ist die Schätzung des Modells (auf einer Skala von 0,0 bis 1,0) bezüglich der Genauigkeit der einzelnen vorhergesagten Werte.

Abbildung 2-15

Table mit generierten Scores und Konfidenzwerten



Number of credit cards	Education	Car loans	\$R-Credit rating	\$RC-Credit rating
5 or more	College	More than 2	Good	0.563
5 or more	College	More than 2	Bad	0.806
5 or more	High school	More than 2	Bad	0.820
5 or more	College	None or 1	Bad	0.820
5 or more	College	More than 2	Good	0.563
5 or more	College	More than 2	Good	0.563
5 or more	College	More than 2	Good	0.563
5 or more	High school	More than 2	Bad	0.820
5 or more	High school	More than 2	Bad	0.820
5 or more	College	More than 2	Bad	0.806
5 or more	College	More than 2	Bad	0.820
5 or more	High school	More than 2	Bad	0.820
5 or more	High school	More than 2	Good	0.563
5 or more	College	None or 1	Bad	0.820
5 or more	High school	More than 2	Bad	0.820
5 or more	College	More than 2	Bad	0.820
5 or more	College	More than 2	Bad	0.820
5 or more	College	More than 2	Good	0.563
5 or more	College	More than 2	Good	0.563
5 or more	College	More than 2	Good	0.823

Erwartungsgemäß stimmt der vorhergesagte Wert bei vielen – nicht jedoch bei allen – Datensätzen mit dem tatsächlichen Ergebnis überein. Der Grund hierfür besteht darin, dass jeder CHAID-Endknoten eine Mischung von Ergebnissen aufweist. Die Vorhersage stimmt mit dem *häufigsten* überein, ist jedoch für alle anderen im Knoten falsch. (Wir erinnern uns an die Minderheit von 16 % der Kunden mit niedrigem Einkommen, die Ihren Kredit zurückgezahlt haben.)

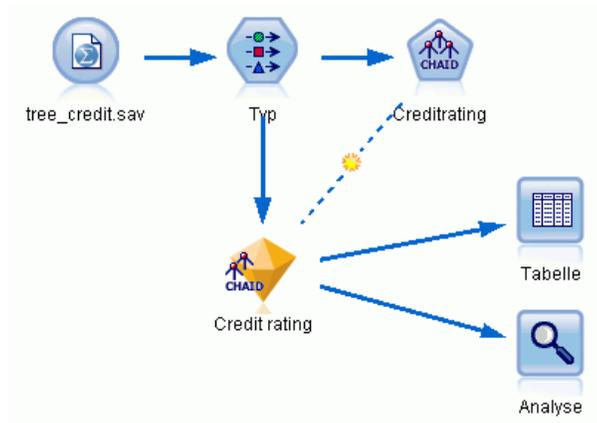
Um dies zu vermeiden, könnten wir damit fortfahren, den Baum in immer kleinere Verzweigungen aufzuspalten, bis jeder Knoten 100%ig einheitlich wäre – nur *Gut* oder nur *Schlecht*, ohne gemischte Antworten. Ein derartiges Modell wäre jedoch extrem kompliziert und ließe sich vermutlich nicht gut auf andere Daten-Sets verallgemeinern.

Um herauszufinden, wie viele der Vorhersagen genau zutreffen, könnten wir die Tabelle durchlesen und die Datensätze zählen, bei denen der Wert im vorhergesagten Feld *\$R-Credit rating* dem Wert im Feld *Credit rating* entspricht. Zum Glück gibt es eine viel einfachere Methode: Wir können einen Analyseknoten verwenden, der dies automatisch erledigt.

- Verbinden Sie das Modell-Nugget mit dem Analyseknoten.

- Doppelklicken Sie auf den Analyseknotten und klicken Sie auf Ausführen.

Abbildung 2-16  
Einfügen eines Analyseknottens



Die Analyse zeigt, dass für 1899 von 2464 Datensätzen (über 77%) der vom Modell vorhergesagte Wert mit der tatsächlichen Antwort übereinstimmte.

Abbildung 2-17  
Analyseergebnisse für den Vergleich zwischen den beobachteten und vorhergesagten Ergebnissen

Das Bild zeigt ein Fenster mit dem Titel 'Analyse von [Credit rating]'. In der Mitte ist eine Tabelle mit den Analyseergebnissen für den Vergleich zwischen den beobachteten und vorhergesagten Ergebnissen zu sehen. Die Tabelle enthält folgende Daten:

Ergebnisse für Zielfeld Credit rating		
Vergleichen von \$R-Credit rating mit Credit rating		
<b>Korrekt</b>	1.960	79,55%
<b>Falsch</b>	504	20,45%
<b>Gesamt</b>	2.464	

Das Fenster enthält auch eine Menüleiste mit 'Datei' und 'Bearbeiten', sowie Schaltflächen für 'Analyse', 'Anmerkungen', 'Alles ausblenden' und 'Alles anzeigen'. Ein 'OK'-Knopf befindet sich unten rechts.

Das Ergebnis wird durch die Tatsache eingeschränkt, dass die gescorten Datensätze dieselben sind, die zur Schätzung des Modells verwendet werden. In einer realen Situationen könnten Sie einen Partitionsknoten verwenden, um die Daten in separate Stichproben für Training und Evaluierung aufzuteilen.

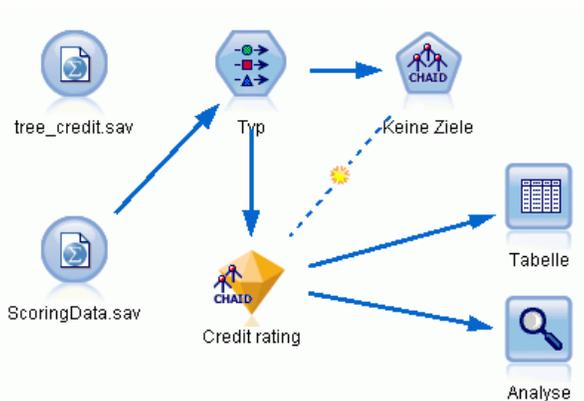
Durch Verwendung einer Stichprobenpartition zur Generierung des Modells und einer weiteren Stichprobenpartition zum Testen des Modells können Sie einen wesentlich besseren Anhaltspunkt dafür erhalten, wie gut sich das Modell auf andere Daten-Sets generalisieren lässt.

Mit dem Analyseknoden können wir das Modell an Datensätzen testen, bei denen wir das tatsächliche Ergebnis bereits kennen. Im nächsten Schritt wird gezeigt, wie wir mit dem Modell Datensätze scoren können, deren Ergebnis wir noch nicht kennen. Es könnten z. B. Personen miteinbezogen werden, die noch keine Kunden der Bank sind, die aber potenzielle Ziele für Werberundschreiben sind.

## Scoren von Datensätzen

Zuvor haben wir dieselben Datensätze gescort, die zur Schätzung des Modells verwendet wurden, um zu evaluieren, wie genau das Modell war. Jetzt werden wir sehen, wie wir einen anderen Datensatz verwenden als den zur Erstellung des Modells. Dies ist das Ziel der Modellierung mit einem Zielfeld: Untersuchung von Datensätzen, bei denen das Ergebnis bekannt ist, um Muster zu ermitteln, mit denen sich Ergebnisse vorhersagen lassen, die noch nicht bekannt sind.

Abbildung 2-18  
Angliedern neuer Daten zum Scoring



Sie können den Statistikdatei-Quellenknoten so aktualisieren, dass er auf eine andere Datendatei verweist, oder Sie können einen neuen Quellenknoten hinzufügen, der die zu scorenden Daten einliest. In beiden Fällen muss das neue Daten-Set dieselben Eingabefelder enthalten wie das Modell (*Age* (Alter), *Income level* (Einkommenskategorie), *Education* (Bildung) usw.), nicht jedoch das Zielfeld *Credit Rating* (Kreditrating).

Alternativ können Sie das Modell-Nugget zu einem beliebigen Stream hinzufügen, der die erwarteten Eingabefelder enthält. Es ist egal, ob die Daten aus einer Datei oder einer Datenbank eingelesen wurden; der Quellentyp ist unerheblich, solange die Feldnamen und -typen mit den im Modell verwendeten übereinstimmen.

Sie können das Modell-Nugget auch als separate Datei speichern, das Modell im PMML-Format exportieren, um sie in anderen Anwendungen zu nutzen, die dieses Format unterstützen, oder das Modell in IBM® SPSS® Collaboration and Deployment Services speichern, was unternehmensweites Deployment, Scoring und unternehmensweite Verwaltung der Modelle ermöglicht.

Unabhängig von der verwendeten Infrastruktur funktioniert das Modell auf dieselbe Weise.

## **Zusammenfassung**

In diesem Beispiel werden die grundlegenden Schritte für Erstellung, Evaluation und Scoring eines Modells erläutert.

- Der Modellierungsknoten schätzt das Modell durch Untersuchung von Datensätzen, deren Ergebnis bekannt ist, und erstellt ein Modell-Nugget. Dieser Vorgang wird auch als Trainieren des Modells bezeichnet.
- Das Modell-Nugget kann zu jedem Stream mit den erwarteten Feldern hinzugefügt werden, um Datensätze zu scoren. Durch Scoren der Datensätze, deren Ergebnis Sie bereits kennen (z. B. bestehende Kunden), können Sie die Leistung des Modells evaluieren.
- Sobald Sie mit der Leistungsfähigkeit des Modells zufrieden sind, können Sie neue Daten (beispielsweise potenzielle Kunden) scoren, um vorherzusagen, wie diese reagieren.
- Die zum Trainieren bzw. Schätzen des Modells verwendeten Daten können auch als analytische oder historische Daten bezeichnet werden; die Scoring-Daten können auch als operationale Daten bezeichnet werden.

# Übersicht über die Modellbildung

## Überblick über Modellierungsknoten

IBM® SPSS® Modeler bietet eine Vielzahl von Modellbildungsmethoden, die aus dem maschinellen Lernen, der künstlichen Intelligenz und der Statistik stammen. Mit den in der Modellierungspalette verfügbaren Methoden können Sie aus Ihren Daten neue Informationen ableiten und Vorhersagemodelle erstellen. Jede Methode besitzt ihre Stärken und eignet sich besonders für bestimmte Problemtypen.

Im *SPSS Modeler-Anwendungshandbuch* finden Sie Beispiele für viele dieser Methoden sowie eine allgemeine Einführung in den Modellierungsprozess. Dieses Handbuch ist als Online-Lernprogramm und im PDF-Format verfügbar. [Für weitere Informationen siehe Thema Anwendungsbeispiele in Kapitel 1 in IBM SPSS Modeler 14.2- Benutzerhandbuch.](#)

Modellierungsmethoden werden in drei Kategorien unterteilt:

- Klassifikation
- Assoziation
- Segmentierung

### **Klassifizierungsmodelle**

Bei *Klassifizierungsmodellen* werden die Werte von einem oder mehreren Feldern **Eingabe** verwendet, um den Wert von einem oder mehreren Feldern “Ausgabe” oder **Ziel** vorherzusagen. Zu diesen Techniken zählen zum Beispiel: Entscheidungsbäume (C&R-Baum-, QUEST-, CHAID- and C5.0-Algorithmen), Regression (lineare, logistische, verallgemeinerte lineare und Cox-Regressionsalgorithmen), neuronale Netze, Support Vector Machines und Bayes-Netzwerke.

Klassifizierungsmodelle können Unternehmen ein bekanntes Ergebnis vorhersagen. Beispielsweise, ob ein Kunde kaufen wird oder nicht, oder ob eine Transaktion mit einem bekannten Betrugsmuster übereinstimmt. Zu den Modellierungstechniken gehören Maschinelles Lernen, Regelinduktion, Identifikation von Untergruppen, statistische Methoden und die Erzeugung mehrerer Modelle.

### *Klassifikationsknoten*



Mit dem Knoten “Autom. Klassifizierer” können Sie eine Reihe verschiedener Modelle für binäre Ergebnisse (“Ja” oder “Nein”, “Abwanderung” oder “Keine Abwanderung” usw.) erstellen und vergleichen, um den besten Ansatz für die jeweilige Analyse auszuwählen. Es wird eine Reihe von Modellierungsalgorithmen unterstützt, sodass Sie die gewünschten Methoden, die spezifischen Optionen für die jeweilige Methode und die Kriterien zum Vergleich der Ergebnisse auswählen können. Der Knoten generiert eine Gruppe von Modellen, die auf den angegebenen Optionen beruhen, und erstellt anhand der von Ihnen angegebenen Kriterien eine Rangordnung der besten Kandidaten. [Für weitere Informationen siehe Thema Knoten “Automatischer Klassifizierer” in Kapitel 5 auf S. 96.](#)



Der Knoten “Auto-Numerisch” schätzt und vergleicht mit einer Reihe verschiedener Methoden Modelle für die Ergebnisse stetiger numerischer Bereiche. Der Knoten arbeitet auf dieselbe Weise wie der Knoten “Automatischer Klassifizierer”: Sie können die zu verwendenden Algorithmen auswählen und in einem Modellierungsdurchlauf mit mehreren Optionskombinationen experimentieren. Folgende Algorithmen werden unterstützt: C&RT-Baum, CHAID, lineare Regression, verallgemeinerte lineare Regression und Support Vector Machines (SVM). Modelle können anhand von Korrelation, relativem Fehler bzw. Anzahl der verwendeten Variablen verglichen werden. [Für weitere Informationen siehe Thema Knoten “Auto-Numerisch” in Kapitel 5 auf S. 107.](#)



Der Knoten für Klassifizierungs- und Regressions-Bäume (C&RT-Bäume) erstellt einen Entscheidungsbaum, mit dem Sie zukünftige Beobachtungen vorhersagen oder klassifizieren können. Bei dieser Methode wird eine rekursive Partitionierung verwendet, um die Trainingsdatensätze in Segmente aufzuteilen. Dabei wird bei jedem Schritt die Unreinheit verringert und ein Knoten im Baum wird als “rein” betrachtet, wenn 100 % der Fälle in eine bestimmte Kategorie des Zielfelds fallen. Ziel- und Eingabefelder können numerische Bereiche oder kategorial (nominal, ordinal oder Flags) sein. Alle Aufteilungen sind binär (nur zwei Untergruppen). [Für weitere Informationen siehe Thema C&R-Baumknoten in Kapitel 6 auf S. 155.](#)



Der QUEST-Knoten bietet eine binäre Klassifizierungsmethode zum Erstellen von Entscheidungsbäumen, die dazu dient, die für große C&R-Baum-Analysen erforderliche Verarbeitungszeit zu verkürzen. Gleichzeitig soll die in den Klassifizierungsbaummodellen festgestellte Tendenz verringert werden, die darin besteht, dass Eingaben bevorzugt werden, die mehr Aufteilungen erlauben. Eingabefelder können stetig (numerische Bereiche) sein, das Zielfeld muss aber kategorial sein. Alle Aufteilungen sind binär. [Für weitere Informationen siehe Thema QUEST-Knoten in Kapitel 6 auf S. 157.](#)



Der CHAID-Knoten erzeugt Entscheidungsbäume unter Verwendung von Chi-Quadrat-Statistiken zur Ermittlung optimaler Aufteilungen. Im Gegensatz zu den Knoten vom Typ “C&RT-Baum” und “QUEST” kann CHAID nichtbinäre Bäume generieren, d. h. Bäume mit Aufteilungen mit mehr als zwei Verzweigungen. Ziel- und Eingabefelder können in einem numerischen Bereich (stetig) oder kategorial sein. Exhaustive CHAID ist eine Änderung von CHAID, die noch gründlicher vorgeht, indem sie alle möglichen Aufteilungen untersucht, allerdings mehr Rechenzeit beansprucht. [Für weitere Informationen siehe Thema CHAID-Knoten in Kapitel 6 auf S. 156.](#)



Der C5.0-Knoten erstellt entweder einen Entscheidungsbaum oder ein Regel-Set. Das Modell teilt die Stichprobe auf der Basis des Felds auf, das auf der jeweiligen Ebene den maximalen Informationsgewinn liefert. Das Zielfeld muss kategorial sein. Es sind mehrere Aufteilungen in mehr als zwei Untergruppen zulässig. [Für weitere Informationen siehe Thema C5.0-Knoten in Kapitel 6 auf S. 177.](#)



Der Knoten “Entscheidungsliste” kennzeichnet Untergruppen bzw. Segmente, die eine höhere oder geringere Wahrscheinlichkeit für ein bestimmtes binäres Ergebnis aufweisen als die Gesamtpopulation. Sie könnten beispielsweise nach Kunden suchen, deren Abwanderung unwahrscheinlich ist oder die mit großer Wahrscheinlichkeit positiv auf eine Kampagne reagieren. Sie können Ihr Geschäftswissen in das Modell integrieren, indem Sie eigene, benutzerdefinierte Segmente hinzufügen und eine Vorschau anzeigen, in der alternative Modelle nebeneinander angezeigt werden, um die Ergebnisse zu vergleichen. Entscheidungslistenmodelle bestehen aus einer Liste von Regeln, bei denen jede Regel eine Bedingung und ein Ergebnis aufweist. Regeln werden in der vorgegebenen Reihenfolge angewendet und die erste Regel, die zutrifft, bestimmt das Ergebnis. [Für weitere Informationen siehe Thema Entscheidungsliste in Kapitel 9 auf S. 226.](#)



Bei linearen Regressionsmodellen wird ein stetiges Ziel auf der Basis linearer Beziehungen zwischen dem Ziel und einem oder mehreren Prädiktoren vorhergesagt. [Für weitere Informationen siehe Thema Lineare Modelle in Kapitel 10 auf S. 267.](#)



Die lineare Regression ist ein statistisches Verfahren zur Zusammenfassung von Daten und die Erstellung von Prognosen durch Anpassung einer geraden Linie oder Fläche, mit der die Diskrepanzen zwischen den vorhergesagten und den tatsächlichen Ausgabewerten minimiert werden.



Der Faktor/PCA-Knoten bietet leistungsstarke Datenreduktionsverfahren zur Verringerung der Komplexität der Daten. Die Hauptkomponentenanalyse (PCA) findet lineare Kombinationen der Eingabefelder, die die Varianz im gesamten Set der Felder am besten erfassen, wenn die Komponenten orthogonal (senkrecht) zueinander sind. Mit der Faktorenanalyse wird versucht, die zugrunde liegenden Faktoren zu bestimmen, die die Korrelationsmuster innerhalb eines Sets beobachteter Felder erklären. Bei beiden Ansätzen besteht das Ziel darin, eine kleinere Zahl abgeleiteter Felder zu finden, mit denen die Informationen in der ursprünglichen Menge der Felder effektiv zusammengefasst werden können. [Für weitere Informationen siehe Thema Faktor/PCA-Knoten in Kapitel 10 auf S. 311.](#)



Der Merkmalsauswahlknoten sichtet die Eingabefelder, um auf der Grundlage einer Reihe von Kriterien (z. B. dem Prozentsatz der fehlenden Werte) zu entscheiden, ob diese entfernt werden sollen. Anschließend erstellt er eine Wichtigkeitsrangfolge der verbleibenden Eingaben in Bezug auf ein angegebenes Ziel. Beispiel: Angenommen, Sie haben ein Daten-Set mit Hunderten potenzieller Eingaben. Welche davon sind voraussichtlich für die Modellierung von medizinischen Behandlungsergebnissen von Bedeutung? [Für weitere Informationen siehe Thema Merkmalsauswahlknoten in Kapitel 4 auf S. 75.](#)



Bei der Diskriminanzanalyse werden strengere Annahmen als bei der logistischen Regression verwendet, sie kann jedoch eine wertvolle Alternative oder Ergänzung zu einer logistischen Regressionsanalyse sein, wenn diese Annahmen erfüllt sind. [Für weitere Informationen siehe Thema Diskriminanzknoten in Kapitel 10 auf S. 319.](#)



Die logistische Regression ist ein statistisches Verfahren zur Klassifizierung von Datensätzen auf der Grundlage der Werte von Eingabefeldern. Sie ist analog zur linearen Regression, außer dass statt eines numerischen Bereichs ein kategoriales Zielfeld verwendet wird. [Für weitere Informationen siehe Thema Logistikknoten in Kapitel 10 auf S. 287.](#)



Das verallgemeinerte lineare Modell erweitert das allgemeine lineare Modell so, dass die abhängige Variable über eine angegebene Verknüpfungsfunktion in linearem Zusammenhang zu den Faktoren und Kovariaten steht. Außerdem ist es mit diesem Modell möglich, dass die abhängige Variable eine von der Normalverteilung abweichende Verteilung aufweist. Es deckt die Funktionen einer großen Bandbreite an Statistikmodellen ab, darunter lineare Regression, logistische Regression, loglineare Modelle für Häufigkeitsdaten und Überlebensmodelle mit Intervallzensurierung. [Für weitere Informationen siehe Thema GenLin-Knoten in Kapitel 10 auf S. 328.](#)



Der Knoten vom Typ “Cox-Regression” ermöglicht Ihnen auch bei zensierten Datensätzen die Erstellung eines Überlebensmodells für Daten über die Zeit bis zum Eintreten des Ereignisses. Das Modell erstellt eine Überlebensfunktion, die die Wahrscheinlichkeit vorhersagt, dass das untersuchte Ereignis für bestimmte Werte der Eingabevariablen zu einem bestimmten Zeitpunkt ( $t$ ) eingetreten ist. [Für weitere Informationen siehe Thema Cox-Knoten in Kapitel 10 auf S. 342.](#)



Der Knoten “Support Vector Machine” (SVM) ermöglicht die Klassifizierung von Daten in eine von zwei Gruppen ohne Überanpassung. SVM eignet sich gut für umfangreiche Daten-Sets, beispielsweise solche mit einer großen Anzahl an Eingabefeldern. [Für weitere Informationen siehe Thema SVM-Knoten in Kapitel 15 auf S. 476.](#)



Mithilfe des Bayes-Netzwerk-Knotens können Sie ein Wahrscheinlichkeitsmodell erstellen, indem Sie beobachtete und aufgezeichnete Hinweise mit Weltwissen kombinieren, um die Wahrscheinlichkeit ihres Vorkommens zu ermitteln. Der Knoten ist speziell für Netzwerke vom Typ “Tree Augmented Naïve Bayes” (TAN) und “Markov-Decke” gedacht, die in erster Linie zur Klassifizierung verwendet werden. [Für weitere Informationen siehe Thema Bayes-Netzwerk-Knoten in Kapitel 7 auf S. 198.](#)



Mithilfe des Knotens für das Selbstlern-Antwortmodell (Self-Learning Response Model, SLRM) können Sie ein Modell erstellen, in dem das Modell anhand eines einzelnen neuen Falls oder einer kleinen Anzahl neuer Fälle neu eingeschätzt werden kann, ohne dass das Modell mit allen Daten neu trainiert werden muss. [Für weitere Informationen siehe Thema SLRM-Knoten in Kapitel 14 auf S. 463.](#)



Der Zeitreihenknoten berechnet Schätzungen für die exponentielle Glättung sowie univariate und multivariate ARIMA-Modelle (ARIMA steht für Autoregressive Integrated Moving Average (autoregressiver integrierter gleitender Durchschnitt)) für Zeitreihendaten und erstellt Vorhersagen über die zukünftige Leistung. Einem Zeitreihenknoten muss stets ein Zeitintervallknoten vorangehen. [Für weitere Informationen siehe Thema Zeitreihen – Modellierungsknoten in Kapitel 13 auf S. 438.](#)



Der Knoten “ $k$ -Nächste Nachbarn” (KNN) verknüpft einen neuen Fall mit der Kategorie oder dem Wert der  $k$  Objekte, die ihm im Prädiktorraum am nächsten liegen, wobei  $k$  eine Ganzzahl ist. Ähnliche Fälle liegen nah beieinander und Fälle mit geringer Ähnlichkeit sind weit voneinander entfernt. [Für weitere Informationen siehe Thema KNN-Knoten in Kapitel 16 auf S. 482.](#)

### **Assoziationsmodelle**

*Assoziationsmodelle* finden Muster in Ihren Daten, wenn ein oder mehrere Elemente (wie zum Beispiel Ereignisse, Einkäufe oder Attribute) mit einem oder mehreren Elementen verknüpft werden. Die Modelle erstellen Regel-Sets, die diese Beziehungen definieren. Dabei können die Felder in den Daten sowohl als Eingaben als auch als Ziele dienen. Sie können solche Assoziationen natürlich auch manuell finden, doch erledigen Algorithmen für Assoziationsregeln das sehr viel schneller. Außerdem können sie komplexere Muster untersuchen. A Priori- und Carma-Modelle sind Beispiele für den Einsatz solcher Algorithmen. Ein weiterer Typ von Assoziationsmodell ist ein Sequenzerkennungsmodell, das sequenzielle Muster in zeitlich strukturierten Daten erkennt.

Assoziationsmodelle sind am nützlichsten bei der Vorhersage mehrerer Ergebnisse – wenn beispielsweise Kunden, die Produkt X gekauft haben, auch Y und Z gekauft haben. Assoziationsmodelle verknüpfen eine bestimmte Schlussfolgerung (wie die Entscheidung, etwas zu kaufen) mit einem Set von Bedingungen. Der Vorteil von Algorithmen für Assoziationsregeln im Vergleich zu Algorithmen für Standardentscheidungsbäume (C5.0 und C&RT) besteht darin, dass zwischen beliebigen Attributen Verbindungen bestehen können. Ein Entscheidungsbaumalgorithmus erstellt Regeln mit nur einer Schlussfolgerung, während Assoziationsalgorithmen viele Regeln zu finden versuchen, von denen jede zu einer anderen Schlussfolgerung kommen kann.

### *Assoziationsknoten*



Der A-Priori-Knoten extrahiert eine Regelmenge aus den Daten und daraus die Regeln mit dem höchsten Informationsgehalt. A Priori bietet fünf verschiedene Methoden zur Auswahl von Regeln und verwendet ein ausgereiftes Indizierungsschema zur effizienten Verarbeitung großer Daten-Sets. Bei großen Problemen ist A Priori in der Regel schneller zu trainieren, es gibt keine willkürliche Begrenzung für die Anzahl der Regeln, die beibehalten werden können, und es können Regeln mit bis zu 32 Vorbedingungen verarbeitet werden. Bei A Priori müssen alle Ein- und Ausgabefelder kategorial sein; dafür bietet es jedoch eine bessere Leistung, da es für diesen Datentyp optimiert ist. [Für weitere Informationen siehe Thema A Priori-Knoten in Kapitel 12 auf S. 390.](#)



Beim CARMA-Modell wird eine Regelmenge aus den Daten extrahiert, ohne dass Sie Eingabe- oder Ziel-Felder angeben müssen. Im Gegensatz zu A Priori bietet der CARMA-Knoten Erstellungseinstellungen für die Regelunterstützung (Unterstützung für Antezedens und Sukzedens) und nicht nur für die Antezedens-Unterstützung. Die erstellten Regeln können somit für eine größere Palette an Anwendungen verwendet werden, beispielsweise um eine Liste mit Produkten und Dienstleistungen (Antezedenzen) zu finden, deren Nachfolger (Sukzedens) das Element darstellt, das Sie in der Ferienzeit desselben Jahres bewerben möchten. [Für weitere Informationen siehe Thema CARMA-Knoten in Kapitel 12 auf S. 394.](#)



Der Sequenzknoten erkennt Assoziationsregeln in sequenziellen oder zeitorientierten Daten. Eine Sequenz ist eine Liste mit Element-Sets, die in einer vorhersagbaren Reihenfolge auftreten. Beispiel: Ein Kunde, der einen Rasierer und After-Shave-Lotion kauft, kauft möglicherweise beim nächsten Einkauf Rasiercreme. Der Sequenzknoten basiert auf dem CARMA-Assoziationsregelalgorithmus, der eine effiziente bidirektionale Methode zum Suchen von Sequenzen verwendet. [Für weitere Informationen siehe Thema Sequenzknoten in Kapitel 12 auf S. 416.](#)

### **Segmentierungsmodelle**

Bei *Segmentierungsmodellen* werden die Daten in Segmente oder Cluster von Datensätzen unterteilt, die ähnliche Muster der Eingabefelder aufweisen. Da sie ausschließlich auf die Eingabefelder konzentriert sind, gibt es bei Segmentierungsmodellen kein Konzept für Ausgabe- oder Zielfelder. Beispiele für Segmentierungsmodelle sind Kohonen-Netzwerke, K-Means-Clustering, Two-Step-Clustering und die Anomalieerkennung.

Segmentierungsmodelle (auch „Clustermodelle“ genannt) sind nützlich in Fällen, in denen das genaue Ergebnis unbekannt ist (beispielsweise, wenn neue Betrugsmuster ermittelt werden oder wenn in der Kundenbasis bestimmte Gruppen identifiziert werden sollen). Clustermodelle konzentrieren sich auf die Ermittlung ähnlicher Datensätze und Beschriftung der Datensätze anhand der Gruppe, in die sie gehören. Dies erfolgt ohne den Vorteil bereits zuvor vorhandener Kenntnisse der Gruppen und der zugehörigen Merkmale. Dies unterscheidet Clustermodelle von anderen Modellierungsverfahren: Es gibt kein zuvor definiertes Ausgabe- oder Zielfeld für das vorherzusagende Modell. Für diese Modelle gibt es keine richtigen oder falschen Antworten. Ihr Wert wird durch die Möglichkeit bestimmt, interessante Gruppierungen in den Daten zu erfassen und sinnvolle Beschreibungen dieser Gruppierungen zu liefern. Cluster-Modelle werden häufig verwendet, um Cluster oder Segmente zu erstellen, die dann als Eingaben in nachfolgenden Analysen verwendet werden (z. B. die Segmentierung potenzieller Kunden in homogene Untergruppen).

### *Segmentierungsknoten*



Mit dem Knoten „Autom. Cluster“ können Sie Clustering-Modelle, die Gruppen und Datensätze mit ähnlichen Merkmalen identifizieren, schätzen und vergleichen. Die Funktionsweise des Knotens gleicht der von anderen Knoten für automatisierte Modellierung, und Sie können in einem einzigen Modellierungsdurchgang mit mehreren Optionskombinationen experimentieren. Modelle können mithilfe grundlegender Messwerte für Filterung und Rangfolge der Nützlichkeit von Cluster-Modellen verglichen werden, um ein Maß auf der Basis der Wichtigkeit von bestimmten Feldern zu liefern. [Für weitere Informationen siehe Thema Knoten „Autom. Cluster“ in Kapitel 5 auf S. 114.](#)



Der K-Means-Knoten teilt das Daten-Set in unterschiedliche Gruppen (oder Cluster) auf. Bei diesem Verfahren wird eine festgelegte Anzahl von Clustern definiert, den Clustern werden iterativ Datensätze zugewiesen und die Cluster-Zentren werden angepasst, bis eine weitere Verfeinerung keine wesentliche Verbesserung des Modells mehr darstellen würde. Statt zu versuchen, ein Ergebnis vorherzusagen, versucht K-Means mithilfe eines als „nicht überwachtetes Lernen“ bezeichneten Verfahrens Muster im Set der Eingabefelder zu entdecken. [Für weitere Informationen siehe Thema K-Means-Knoten in Kapitel 11 auf S. 362.](#)



Der Kohonen-Knoten erstellt eine Art von neuronalem Netzwerk, das verwendet werden kann, um ein Clustering der Datenmenge in einzelne Gruppen vorzunehmen. Wenn das Netz voll trainiert ist, sollten ähnliche Datensätze auf der Ausgabekarte eng nebeneinander stehen, während Datensätze, die sich unterscheiden, weit voneinander entfernt sein sollten. Die Zahl der von jeder Einheit im Modell-Nugget erfassten Beobachtungen gibt Aufschluss über die starken Einheiten. Dadurch wird ein Eindruck von der ungefähren Zahl der Cluster vermittelt. [Für weitere Informationen siehe Thema Kohonen-Knoten in Kapitel 11 auf S. 355.](#)



Der TwoStep-Knoten verwendet eine aus zwei Schritten bestehende Clusterbildungsmethode. Im ersten Schritt wird ein einzelner Durchlauf durch die Daten vorgenommen, bei dem die Eingangsrohdaten zu einem verwaltbaren Set von Unterclustern komprimiert werden. Im zweiten Schritt werden die Untercluster mithilfe einer hierarchischen Methode zur Clusterbildung nach und nach in immer größere Cluster zusammengeführt. TwoStep hat den Vorteil, dass die optimale Anzahl an Clustern für die Trainingsdaten automatisch geschätzt wird. Mit dem Verfahren können gemischte Feldtypen und große Daten-Sets effizient verarbeitet werden. [Für weitere Informationen siehe Thema TwoStep-Cluster-Knoten in Kapitel 11 auf S. 367.](#)



Der Anomalieerkennungsknoten ermittelt ungewöhnliche Fälle bzw. “Ausreißer”, die nicht den Mustern der “normalen” Daten entsprechen. Mit diesem Knoten können Ausreißer ermittelt werden, selbst wenn sie keinem bereits bekannten Muster entsprechen und selbst wenn Sie nicht genau wissen, wonach Sie suchen. [Für weitere Informationen siehe Thema Anomalieerkennungsknoten in Kapitel 4 auf S. 83.](#)

### **Modelle für In-Database Mining**

SPSS Modeler unterstützt die Integration mit Data-Mining-Tools und Daten-Modellierungstools von Datenbankherstellern wie Oracle Data Miner, IBM DB2 InfoSphere Warehouse und Microsoft Analysis Services. Sie können Modelle in der Datenbank erstellen, bewerten und speichern – ohne dazu die SPSS Modeler-Anwendung verlassen zu müssen. Vollständige Informationen finden Sie im *SPSS Modeler In-Database Mining Handbuch*, das sich auf dem Produkt-DVD befindet.

### **IBM SPSS Statistics Modelle**

Wenn auf Ihrem Computer eine Kopie von IBM® SPSS® Statistics installiert und lizenziert ist, können Sie auf bestimmte SPSS Statistics-Routinen in SPSS Modeler zugreifen und diese ausführen, um Modelle zu erstellen und zu scoden. [Für weitere Informationen siehe Thema IBM SPSS Statistics-Knoten – Überblick in Kapitel 8 in IBM SPSS Modeler 14.2- Quellen-, Prozess- und Ausgabeknoten.](#)

### **Weitere Informationen**

Es steht auch eine detaillierte Dokumentation zu den Modellierungsalgorithmen zur Verfügung. Weitere Informationen finden Sie im *SPSS Modeler-Algorithmushandbuch*, das sich auf dem Produkt-DVD befindet.

## Erstellung von aufgeteilten Modellen

Die Erstellung aufgeteilter Modelle ermöglicht es Ihnen, einen einzelnen Stream für die Erstellung getrennter Modelle für jeden möglichen Wert eines Flag-, nominalen oder stetigen Eingabefelds zu verwenden, wobei auf alle daraus resultierenden Modelle von einem einzelnen Modell-Nugget aus zugegriffen werden kann. Die möglichen Werte für die Eingabefelder könnten sehr unterschiedliche Effekte auf das Modell haben. Durch aufgeteilte Modellierung können Sie ganz einfach das am besten geeignete Modell für jeden möglichen Feldwert mit einer einzigen Ausführung des Streams erstellen.

Bitte beachten Sie, dass die Aufteilungsfunktion in interaktiven Modellierungssitzungen nicht verwendet werden kann. Bei der interaktiven Modellierung geben Sie jedes Modell einzeln an, weswegen die Verwendung der Aufteilungsfunktion, über die mehrere Modelle automatisch erstellt werden, nicht von Vorteil wäre.

Die aufgeteilte Modellierung wird angewendet, indem man ein bestimmtes Eingabefeld als Aufteilungsfeld angibt. Dies ist möglich, indem die Feldrolle in der Typspezifikation auf Aufteilen eingestellt wird:

Abbildung 3-1  
Festlegung eines Eingabefelds als Aufteilungsfeld

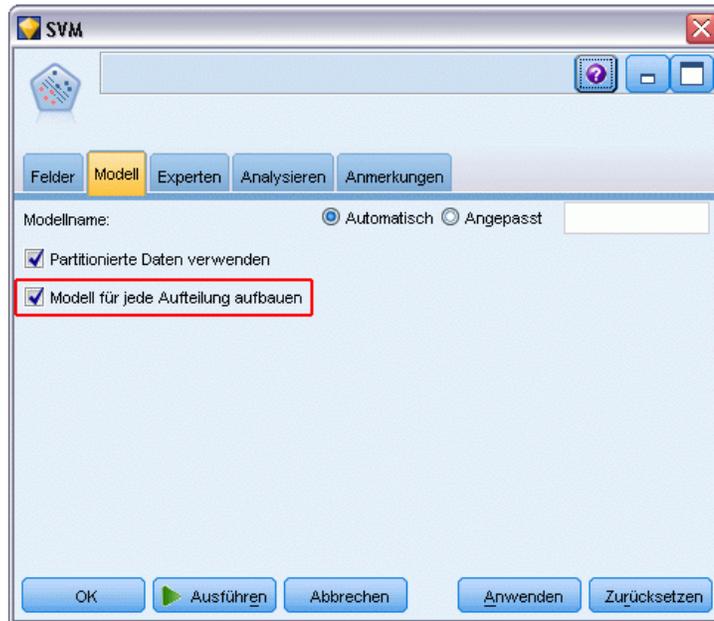


Sie können nur Felder mit einem Messniveau Flag, Nominal, Ordinal oder Stetig als Aufteilungsfelder festlegen.

Sie können mehr als ein Eingabefeld als Aufteilungsfeld festlegen. Dadurch kann jedoch die Anzahl erstellter Modelle erheblich gesteigert werden. Für jede mögliche Kombination der Werte der ausgewählten Aufteilungsfelder wird ein Modell erstellt. Wenn beispielsweise drei Eingabefelder mit je drei möglichen Werten als Aufteilungsfelder festgelegt werden, führt dies zu einer Erstellung von 27 unterschiedlichen Modellen.

Selbst nachdem Sie ein Feld oder mehrere Felder als Aufteilungsfelder festgelegt haben, können Sie mit Hilfe eines Kontrollkästchens im Dialogfeld des Modellierungsknotens wählen, ob aufgeteilte Modelle oder ein einzelnes Modell erstellt werden sollen.

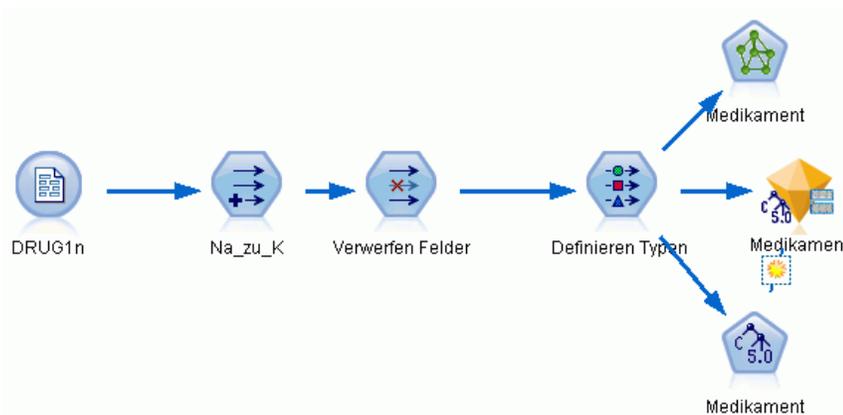
**Abbildung 3-2**  
Option zur Erstellung von aufgeteilten Modellen



Wenn Aufteilungsfelder definiert sind, das Kontrollkästchen aber nicht aktiviert ist, wird nur ein einzelnes Modell generiert. Ebenso wird, wenn das Kontrollkästchen aktiviert, aber kein Aufteilungsfeld definiert ist, die Aufteilung ignoriert und ein einzelnes Modell generiert.

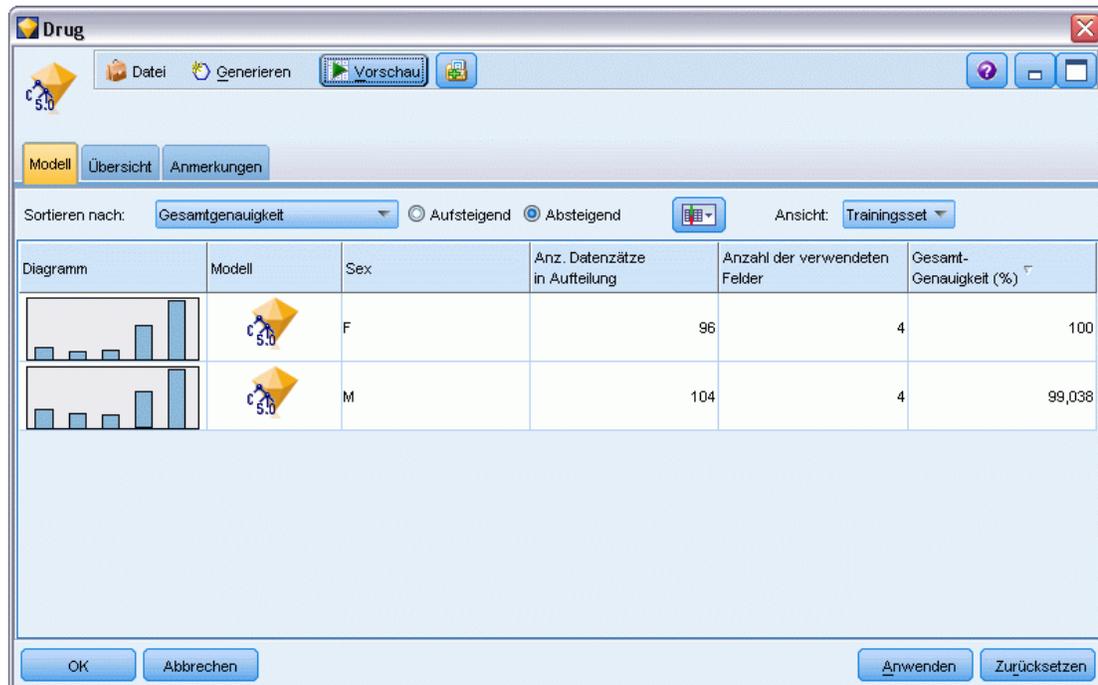
Wenn Sie einen Stream ausführen, werden für jeden möglichen Wert der Aufteilungsfelder im Hintergrund separate Modelle generiert, aber es wird nur ein einzelnes Modell-Nugget in der Modellpalette und der Stream-Zeichenfläche erstellt. Ein Nugget für aufgeteilte Modelle wird durch das Aufteilungs-Symbol gekennzeichnet:

**Abbildung 3-3**  
Nugget für aufgeteilte Modelle in einem Stream



Wenn Sie das Nugget für aufgeteilte Modelle durchsuchen, sehen Sie eine Liste der separaten Modelle, die erstellt wurden.

Abbildung 3-4  
Viewer für aufgeteilte Modelle



Sie können ein individuelles Modell aus einer Liste untersuchen, indem Sie im Viewer auf sein Nugget-Symbol doppelklicken. Damit wird ein Standard-Browserfenster für das individuelle Modell geöffnet. Wenn sich das Nugget im Zeichenbereich befindet, wird durch Doppelklicken auf ein Miniaturdiagramm ein Diagramm in voller Größe geöffnet. [Für weitere Informationen siehe Thema Viewer für aufgeteilte Modelle auf S. 66.](#)

Nachdem ein Modell als aufgeteiltes Modell erstellt wurde, können Sie den Aufteilungsprozess nicht mehr rückgängig machen; auch weiter abwärts vorgenommene Aufteilungen können von einem Aufteilungsmodellierungsknoten oder -Nugget aus nicht rückgängig gemacht werden.

**Beispiel.** Ein national operierender Einzelhändler möchte Schätzungen der Verkäufe nach Produktkategorie für jedes seiner Geschäfte im ganzen Land vornehmen. Unter Verwendung von Aufteilungsmodellierung legt er das Speicherfeld für seine Eingabedaten als Aufteilungsfeld fest und kann so für jede Kategorie in jedem Geschäft mit Hilfe eines einzigen Vorgangs separate Modelle erstellen. Durch die so gewonnenen Informationen kann er die Lagerbestände viel genauer kontrollieren als es anhand eines einzelnen Modells möglich wäre.

## Aufteilung und Partitionierung

Aufteilung und Partitionierung haben einige gemeinsame Eigenschaften, werden aber auf vollkommen unterschiedliche Arten verwendet.

Die **Partitionierung** unterteilt das Daten-Set zufällig in entweder zwei oder drei Teile – Training, Testen und (optional) Validierung – und wird verwendet, um die Leistung eines einzelnen Modells zu testen.

Die **Aufteilung** unterteilt das Daten-Set in so viele Teile, wie es mögliche Werte für ein Aufteilungsfeld gibt, und wird verwendet, um mehrere Modelle zu erstellen.

Partitionierung und Aufteilung sind Vorgänge, die vollkommen unabhängig voneinander sind. In einem Modellierungsknoten können Sie einen von ihnen, beide oder keinen auswählen.

### **Modellierungsknoten zur Unterstützung aufgeteilter Modelle**

Eine Reihe von Modellierungsknoten können aufgeteilte Modelle erstellen. Die Ausnahmen sind “Autom. Cluster”, “Zeitreihe”, “Faktor/PCA”, “Merkmalsauswahl”, SLRM, die Assoziationsmodelle (A Priori, Carma und Sequenz), die Clustermodelle (K-Means, Kohonen, Two Step und Anomaly), Statistics-Modell sowie die Knoten zur Modellierung innerhalb der Datenbank.

Die folgenden Modellierungsknoten unterstützen aufgeteilte Modellierung:



C&R-Baum



Bayes-Netz



QUEST



GenLin



CHAID



KNN



C5,0



Cox



Netzwerk



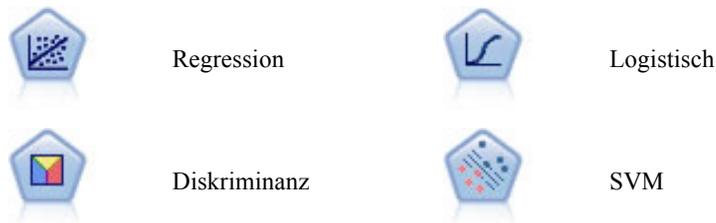
Automatischer Klassifizierer



Entscheidungsliste



Auto-Numerisch



## Von der Aufteilung betroffene Merkmale

Die Verwendung von aufgeteilten Modellen beeinflusst eine Reihe von IBM® SPSS® Modeler-Merkmalen auf mehrere Arten. Dieser Abschnitt bietet Richtlinien zur Nutzung von aufgeteilten Modellen im Zusammenhang mit anderen Knoten in einem Stream.

### **Knoten "Datensatzoperationen"**

Beim Verwenden von aufgeteilten Modellen in einem Stream, der einen **Stichproben**-Knoten enthält, schichten Sie Datensätze nach dem Aufteilungsfeld, um gleichmäßige Stichproben von Datensätzen zu erhalten. Diese Option ist verfügbar, wenn Sie "Komplex" als Stichprobenmethode wählen. [Für weitere Informationen siehe Thema Einstellungen unter "Klumpen und Schichtung" in Kapitel 3 in IBM SPSS Modeler 14.2- Quellen-, Prozess- und Ausgabeknoten.](#)

Wenn der Stream einen **Balancierung**knoten enthält, beachten Sie, dass die Balancierung für das vollständige Set der Eingabedatensätze gilt, nicht für eine Untergruppe von Datensätzen innerhalb einer Aufteilung.

Beim Aggregieren von Datensätzen mithilfe eines **Aggregat**knotens legen Sie die aufgeteilten Felder als Schlüsselfelder fest, wenn Sie Aggregate für jede Aufteilung berechnen möchten. [Für weitere Informationen siehe Thema Aggregatknoten. in Kapitel 3 in IBM SPSS Modeler 14.2- Quellen-, Prozess- und Ausgabeknoten.](#)

### **Feldoperationsknoten**

Im Knoten **Typ** geben Sie an, welches Feld bzw. welche Felder als aufgeteilte Felder dienen sollen. [Für weitere Informationen siehe Thema Typknoten in Kapitel 4 in IBM SPSS Modeler 14.2- Quellen-, Prozess- und Ausgabeknoten.](#)

Beachten Sie: Der Knoten **Ensemble** wird zur Kombination von zwei oder mehr Modell-Nuggets verwendet, lässt sich jedoch nicht benutzen, um die Aktion des Aufteilens umzukehren, da sich die aufgeteilten Modelle in einem einzigen Modell-Nugget befinden.

### **Modellierungsknoten**

Aufgeteilte Modelle unterstützen keine Berechnung der Bedeutsamkeit von Prädiktoren (die relative Bedeutsamkeit der Prädiktor-Eingabefelder bei der Modellschätzung). Die Einstellungen der Bedeutsamkeit von Prädiktoren werden bei der Erstellung von aufgeteilten Modellen ignoriert.

Der Knoten **KNN** (Nächster Nachbar) unterstützt nur dann aufgeteilte Modelle, wenn er auf die Vorhersage eines Zielfelds eingestellt ist. Die alternative Einstellung (nur nächste Nachbarn identifizieren) erzeugt kein Modell. Wenn die Option “Automatisch k wählen” aktiviert ist, kann jedes aufgeteilte Modell über eine unterschiedliche Anzahl nächster Nachbarn verfügen. Das Gesamtmodell verfügt so über eine Reihe von generierten Spalten gleich der größten Anzahl nächster Nachbarn, die in allen aufgeteilten Modellen gefunden wird. Für die aufgeteilten Modelle, deren Anzahl nächster Nachbarn dieses Maximum unterschreitet, gibt es eine entsprechende Anzahl von Spalten, die mit \$null-Werten gefüllt sind. [Für weitere Informationen siehe Thema KNN-Knoten in Kapitel 16 auf S. 482.](#)

### ***Datenbank-Modellierungsknoten***

Die Knoten für Modellierung innerhalb der Datenbank unterstützen keine aufgeteilten Modelle.

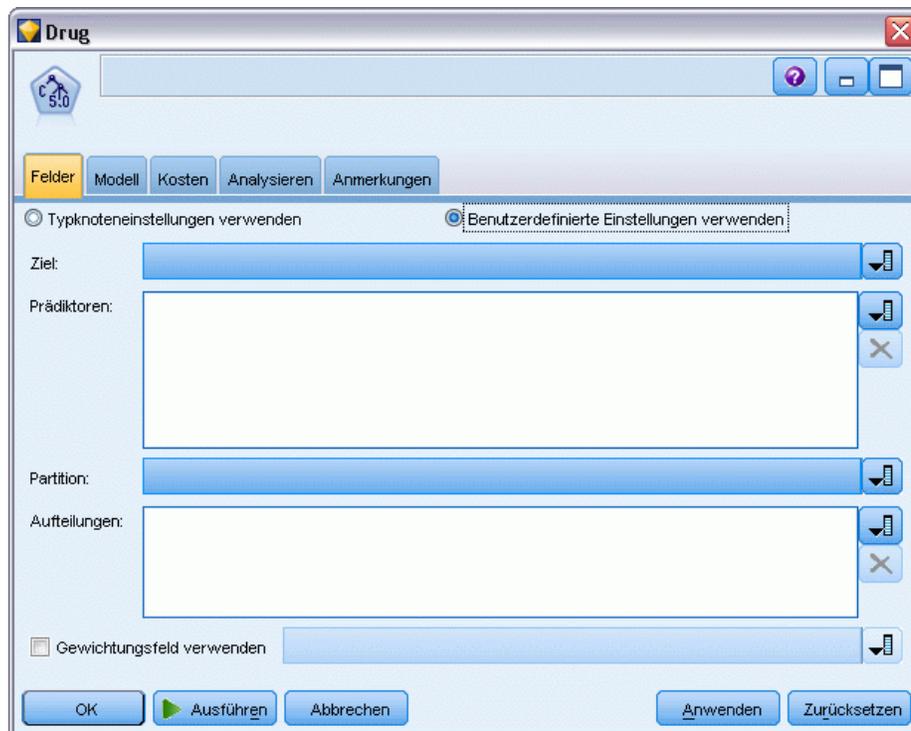
### ***Modell-Nuggets***

**PMML exportieren** ist aus einem aufgeteilten Modell-Nugget nicht möglich, da das Nugget mehrere Modelle enthält und PMML eine solche Zusammenfassung nicht unterstützt. Das Exportieren von Text oder HTML ist jedoch möglich.

## ***Feldoptionen der Modellierungsknoten***

Alle Modellierungsknoten besitzen die Registerkarte “Felder”, auf der Sie die Felder festlegen können, die beim Erstellen des Modells verwendet werden.

Abbildung 3-5  
C5.0-Modellierungsknoten, Registerkarte "Felder"



Bevor Sie ein Modell erstellen können, müssen Sie festlegen, welche Felder als Ziele und als Eingaben verwendet werden sollen. Von wenigen Ausnahmen abgesehen, verwenden alle Modellierungsknoten die Feldinformationen des oberhalb liegenden Typknotens. Wenn Sie einen Typknoten benutzen, um Eingabe- und Zielfelder auszuwählen, brauchen Sie auf dieser Registerkarte keine Änderungen vorzunehmen. (Eine Ausnahme bilden der Sequenzknoten und der Textextraktionsknoten, deren Feldeinstellungen im Modellierungsknoten angegeben sein müssen.)

**Typknoteneinstellungen verwenden.** Diese Option weist den Knoten an, die Feldinformationen von einem weiter oben liegenden Typknoten zu verwenden. Dies ist die Standardeinstellung.

**Benutzerdefinierte Einstellungen verwenden.** Diese Option weist den Knoten an, die hier angegebenen Feldinformationen anstelle der in einem weiter oben liegenden Typknoten angegebenen zu verwenden. Geben Sie nach Auswahl dieser Option wie erforderlich die unten stehenden Felder an.

*Anmerkung: Nicht alle Felder werden für alle Knoten angezeigt.*

- **Als Aktualisierungsmodell verwenden (nur A Priori-, CARMA-, MS-Assoziationsregeln-, ISWAssociation- und Oracle A Priori-Knoten).** Aktivieren Sie dieses Kontrollkästchen, wenn die Quelldaten in **Transaktionsformat** vorliegen. Datensätze in diesem Format enthalten zwei Felder, eines für eine ID und eines für den Inhalt. Jeder Datensatz steht für ein einzelnes Element. Zugeordnete Elemente werden verknüpft, indem sie dieselbe ID erhalten. Deaktivieren Sie dieses Feld, wenn die Daten in **Tabellenformat** vorliegen, in dem Elemente durch separate Flags repräsentiert werden, wobei jedes Flag-Feld für das Vorhandensein oder

die Abwesenheit eines bestimmten Elements steht und jeder Datensatz ein vollständiges Set an zugehörigen Elementen repräsentiert. [Für weitere Informationen siehe Thema Tabellendaten im Vergleich zu Transaktionsdaten in Kapitel 12 auf S. 389.](#)

- **ID.** Wählen Sie für Transaktionsdaten ein ID-Feld aus der Liste aus. Als ID-Feld können numerische oder symbolische Felder verwendet werden. Jeder eindeutige Wert in diesem Feld sollte eine bestimmte Analyseeinheit darstellen. Bei einer Warenkorbanwendung könnte z. B. jede ID einen einzelnen Kunden darstellen. Für eine Webprotokoll-Analyseanwendung könnte jede ID einen Computer (nach IP-Adresse) oder einen Benutzer (nach Anmeldedaten) darstellen.
- **IDs sind zusammenhängend.** (Nur A Priori- und CARMA-Knoten) Wenn Ihre Daten vorsortiert sind, sodass alle Datensätze mit derselben ID im Daten-Stream zusammengefasst sind, wählen Sie diese Option, um die Verarbeitung zu beschleunigen. Wenn Ihre Daten nicht vorsortiert sind (oder Sie nicht sicher sind), lassen Sie diese Option deaktiviert. Die Daten werden dann vom Knoten automatisch sortiert.

*Hinweis:* Wenn Ihre Daten nicht sortiert werden und Sie diese Option auswählen, könnten Sie ungünstige Ergebnisse in Ihrem Modell erhalten.

- **Inhalt.** Geben Sie die Inhaltsfelder für das Modell an. Diese Felder enthalten die Elemente, die für die Assoziationsmodellierung interessant sind. Sie können mehrere Flag-Felder angeben (falls die Daten in tabellarischer Form vorliegen) oder ein einzelnes nominales Feld (falls die Daten im Transaktionsformat vorliegen).
- **Ziel.** Wählen Sie die Zielfelder für Modelle aus, die eines oder mehrere Zielfelder benötigen. Dies ist so, als würden Sie in einem Typknoten für die Rolle eines Felds den Wert *Ziel* festlegen.
- **Evaluierung.** (Nur für “Autom. Cluster”-Modelle.) Für Cluster-Modelle ist kein Ziel angegeben. Sie können jedoch ein Evaluierungsfeld auswählen, um das Wichtigkeitsniveau zu ermitteln. Darüber hinaus können Sie evaluieren, wie gut die Cluster Werte dieses Felds differenzieren. Dies wiederum gibt an, ob die Cluster zur Vorhersage dieses Felds verwendet werden können.
- **Eingaben.** Wählen Sie das Eingabefeld bzw. die Eingabefelder aus. Dies ist so, als würden Sie in einem Typknoten für die Rolle eines Felds den Wert *Eingabe* festlegen.
- **Partition.** In diesem Feld können Sie ein Feld angeben, das verwendet wird, um die Daten in getrennte Stichproben für die Trainings-, Test- und Validierungsphase der Modellbildung aufzuteilen. Indem Sie mit einer Stichprobe das Modell erstellen und es mit einer anderen Stichprobe testen, erhalten Sie einen guten Hinweis dafür, wie gut das Modell sich für größere Datenmengen generalisieren lässt, die den aktuellen Daten ähneln. Wenn mehrere Partitionsfelder mithilfe von Typ- oder Partitionsknoten erstellt wurden, muss in jedem Modellierungsknoten, der die Partitionierung verwendet, auf der Registerkarte “Felder” ein einzelnes Partitionsfeld ausgewählt werden. (Wenn nur eine einzige Partition vorhanden ist, wird diese immer automatisch verwendet, wenn die Partitionierung aktiviert ist.) [Für weitere Informationen siehe Thema Partitionsknoten in Kapitel 4 in IBM SPSS Modeler 14.2-Quellen-, Prozess- und Ausgabeknoten.](#) Beachten Sie außerdem, dass die Partitionierung auch auf der Registerkarte “Modelloptionen” des Knotens aktiviert sein muss, wenn die ausgewählte Partitionierung in Ihrer Analyse angewendet werden soll. (Wenn die Auswahl der Option aufgehoben ist, kann die Partitionierung ohne Änderung der Feldeinstellungen deaktiviert werden.)

- **Aufteilen.** Wählen Sie für Aufteilungsmodelle das Aufteilungsfeld bzw. die Aufteilungsfelder. Dies ist so, als würden Sie in einem Typknoten für die Rolle eines Felds den Wert *Aufteilung* festlegen. Sie können nur Felder mit einem Messniveau *Flag*, *Nominal*, *Ordinal* oder *Stetig* als Aufteilungsfelder festlegen. Als Aufteilungsfelder gewählte Felder können nicht als Ziel-, Prädiktor-, Partitions-, Häufigkeits- oder Gewichtungsfelder verwendet werden. [Für weitere Informationen siehe Thema Erstellung von aufgeteilten Modellen auf S. 31.](#)
- **Häufigkeitsfeld verwenden.** Mit dieser Option können Sie ein Feld als Häufigkeitsgewichtung angeben. Dies sollten Sie tun, wenn die in Ihren Trainingsdaten enthaltenen Datensätze jeweils mehr als eine Einheit darstellen — wenn Sie beispielsweise aggregierte Daten verwenden. Die Feldwerte sollten die Anzahl der Einheiten sein, die von jedem Datensatz repräsentiert werden. [Für weitere Informationen siehe Thema Verwenden von Häufigkeits- und Gewichtungsfeldern auf S. 39.](#)

*Hinweis:* Wenn Sie die Meldung Metadaten (in Eingabe-/Ausgabefeldern) nicht gültig sehen, stellen Sie sicher, dass Sie alle geforderten Felder angegeben haben, z. B. das Häufigkeitsfeld.

- **“Feld gewichten” verwenden.** Mit dieser Option können Sie ein Feld als Fallgewichtung festlegen. Fallgewichtungen werden verwendet, um Differenzen in der Varianz zwischen den Ebenen des Ausgabefelds zu berücksichtigen. [Für weitere Informationen siehe Thema Verwenden von Häufigkeits- und Gewichtungsfeldern auf S. 39.](#)
- **Sukzedenzen.** Wählen Sie die für Regelinduktionsknoten (A Priori) in der resultierenden Regelmenge als Sukzedenzen zu verwendenden Felder aus. (Dies entspricht den in einem Typknoten vorhandenen Feldern mit der Rolle *Ziel* oder *Beides*.)
- **Antezedenzen.** Wählen Sie die für Regelinduktionsknoten (A Priori) in der resultierenden Regelmenge als Antezedenzen zu verwendenden Felder aus. (Dies entspricht den in einem Typknoten vorhandenen Feldern mit der Rolle *Eingabe* oder *Beides*.)

Bei einigen Modellen weicht die Registerkarte “Felder” von den in diesem Abschnitt beschriebenen ab.

- [Für weitere Informationen siehe Thema Feldoptionen für den Sequenzknoten in Kapitel 12 auf S. 417.](#)
- [Für weitere Informationen siehe Thema Feldoptionen für den CARMA-Knoten in Kapitel 12 auf S. 395.](#)

## **Verwenden von Häufigkeits- und Gewichtungsfeldern**

Häufigkeits- und Gewichtungsfelder dienen dazu, einigen Datensätzen eine größere Bedeutsamkeit zu verleihen als anderen, beispielsweise, weil Sie wissen, dass ein Bevölkerungsteil in den Trainingsdaten unterrepräsentiert ist (Gewicht) oder weil ein Datensatz für eine Reihe identischer Fälle steht (Häufigkeit).

- Häufigkeitsfeldwerte müssen als positive ganze Zahlen angegeben werden. Datensätze mit einer Häufigkeitsgewichtung kleiner oder gleich 0 werden von der Analyse ausgeschlossen. Nicht als ganze Zahlen angegebene Häufigkeitsgewichtungen werden auf die nächstliegende ganze Zahl gerundet.
- Fallgewichtungswerte müssen als positive Zahlen angegeben werden, müssen aber keine ganzen Zahlen sein. Datensätze mit einer Fallgewichtung kleiner oder gleich 0 werden von der Analyse ausgeschlossen.

### **Scoren von Häufigkeits- und Gewichtungsfeldern**

Häufigkeits- und Gewichtungsfelder werden beim Trainieren von Modellen verwendet, nicht jedoch beim Scoren, da der Score für die einzelnen Datensätze auf den jeweiligen Merkmalen des Datensatzes beruht, unabhängig davon, wie viele Fälle er umfasst. Hier ein Beispiel: Angenommen, es liegen folgende Daten vor:

Verheiratet	Antwort
Ja	Ja
Ja	Ja
Ja	Ja
Ja	Nein
Nein	Ja
Nein	Nein
Nein	Nein

Auf dieser Grundlage schließen Sie, dass drei von vier verheirateten Personen auf die Werbeaktion antworteten und zwei von drei unverheirateten Personen nicht antworteten. Daher scoren Sie etwaige neue Datensätze entsprechend:

Verheiratet	S-Antwort	SRP-Antwort
Ja	Ja	0,75 (drei/vier)
Nein	Nein	0,67 (zwei/drei)

Alternativ können Sie die Trainingsdaten mithilfe eines Häufigkeitsfelds in kompakterer Form speichern:

Verheiratet	Antwort	Frequency
Ja	Ja	3
Ja	Nein	1
Nein	Ja	1
Nein	Nein	2

Da dies für genau dasselbe Datenset steht, erstellen Sie damit dasselbe Modell und sagen die Antworten ausschließlich auf der Grundlage des Ehestandes voraus. Wenn Ihre Scoring-Daten zehn verheiratete Personen enthalten, sagen Sie für alle jeweils *Ja* voraus, unabhängig davon, ob sie als zehn separate Datensätze vorgelegt werden oder als ein einziger Datensatz mit der Häufigkeit 10. Das Gewicht ist zwar im Allgemeinen keine ganze Zahl, aber zeigt dennoch

in ähnlicher Weise die Bedeutsamkeit eines Datensatzes an. Daher werden Häufigkeits- und Gewichtungsfelder beim Scoring von Datensätzen nicht verwendet.

### ***Evaluation und Vergleich von Modellen***

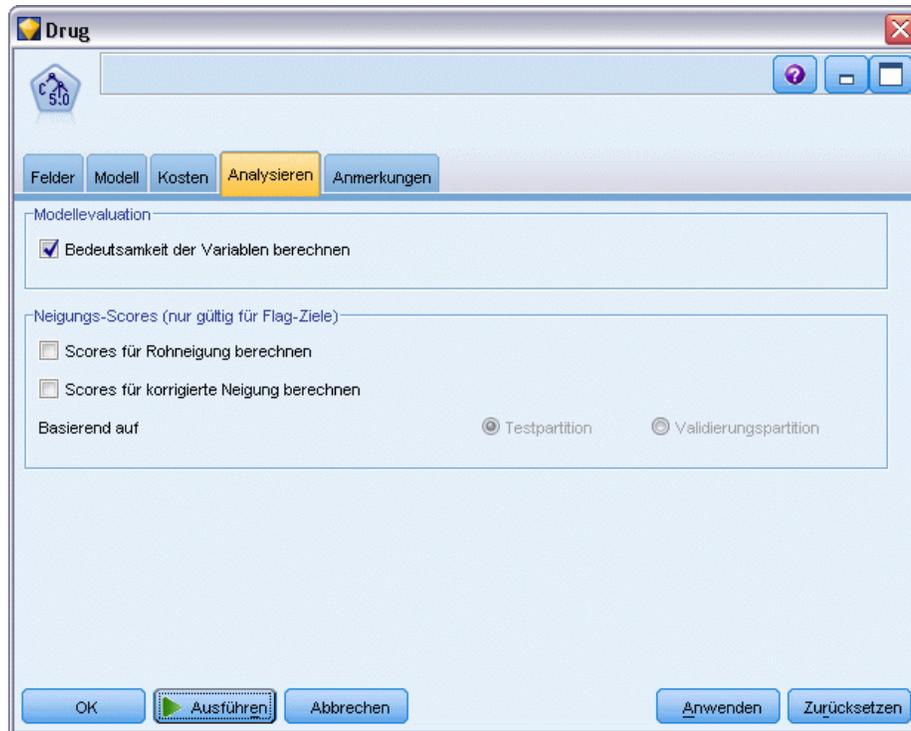
Einige Modelltypen unterstützen Häufigkeitsfelder, einige Gewichtsfelder und einige beide Arten von Feldern. In allen Fällen, in denen sie zulässig sind, werden sie jedoch ausschließlich für die Modellerstellung verwendet und bei der Evaluation von Modellen mithilfe eines Evaluations- oder Analyseknötens nicht verwendet. Ebenso wenig verwendet werden sie bei der Rangeinteilung von Modellen mithilfe der meisten von den Knöten vom Typ "Automatischer Klassifizierer" und "Auto-Numerisch" unterstützten Methoden.

- So werden Häufigkeits- und Gewichtswerte (beispielsweise beim Vergleich von Modellen) mithilfe von Evaluationsdiagrammen ignoriert. Dadurch wird zwar ein Vergleich zwischen Modellen, die diese Felder verwenden, und Modellen ohne diese Felder möglich, es bedeutet jedoch auch, dass für eine genaue Evaluation ein Datensatz verwendet werden muss, das die Grundgesamtheit genau darstellt, ohne dass dafür auf ein Häufigkeits- oder Gewichtsfeld zurückgegriffen wird. In der Praxis können Sie dies tun, indem Sie sicherstellen, dass die Modelle mithilfe einer Teststichprobe evaluiert werden, in der der Wert des Häufigkeitsfelds immer null oder 1 ist. (Diese Einschränkung gilt nur bei der Evaluation von Modellen; wenn die Häufigkeits- bzw. Gewichtswerte sowohl für die Trainings- als auch für die Teststichprobe stets 1 wären, gäbe es keinen Grund, diese Felder überhaupt zu verwenden.)
- Bei Verwendung von "Automatischer Klassifizierer" kann die Häufigkeit berücksichtigt werden, wenn die Modelle auf der Grundlage des Profits in Ränge eingeteilt werden. In diesem Fall wird also diese Methode empfohlen.
- Falls erforderlich, können Sie die Daten mithilfe eines Partitionsknötens in Trainings- und Teststichproben aufspalten. [Für weitere Informationen siehe Thema Partitionsknöten in Kapitel 4 in IBM SPSS Modeler 14.2- Quellen-, Prozess- und Ausgabeknötens.](#)

## ***Analyseoptionen bei Modellierungsknötens***

Zahlreiche Modellierungsknöten enthalten die Registerkarte "Analysieren", der Sie Informationen zur Bedeutsamkeit der Prädiktoren sowie Scores für Rohneigung und angepasste Neigung entnehmen können.

Abbildung 3-6  
Registerkarte "Analysieren" im Modellierungsknoten



### **Modellevaluation**

**Bedeutsamkeit der Prädiktoren berechnen.** Bei Modellen, die zu einem angemessenen Maß an Bedeutsamkeit führen, können Sie ein Diagramm anzeigen, in dem die relative Wichtigkeit der einzelnen Prädiktoren bei der Modellschätzung angegeben wird. Normalerweise ist es sinnvoll, die Modellierungsbemühungen auf die wichtigsten Prädiktoren zu konzentrieren und diejenigen Prädiktoren zu verwerfen bzw. zu ignorieren, die die geringste Bedeutung haben. Beachten Sie, dass die Berechnung der Bedeutsamkeit der Prädiktoren bei einigen Modellen längere Zeit in Anspruch nehmen kann, insbesondere bei der Arbeit mit großen Daten-Sets, und daher bei einigen Modellen standardmäßig deaktiviert ist. Die Bedeutsamkeit der Prädiktoren ist für Entscheidungslistenmodelle nicht verfügbar. [Für weitere Informationen siehe Thema Bedeutsamkeit des Prädiktors auf S. 54.](#)

### **Neigungs-Scores**

Neigungs-Scores können im Modellierungsknoten oder auf der Registerkarte "Einstellungen" im Modell-Nugget aktiviert werden. Diese Funktionalität ist nur verfügbar, wenn das ausgewählte Ziel ein Flag-Feld ist. [Für weitere Informationen siehe Thema Neigungsbewertungen auf S. 44.](#)

**Scores für Rohneigung berechnen.** Rohneigungs-Scores werden ausschließlich auf der Grundlage der Trainingsdaten aus dem Modell abgeleitet. Wenn das Modell den Wert *wahr* (wird antworten) vorhersagt, ist die Neigung mit  $P$  identisch. Dabei ist  $P$  die Wahrscheinlichkeit der Vorhersage. Wenn das Modell den Wert "falsch" vorhersagt, wird die Neigung als  $(1 - P)$  berechnet.

- Wenn Sie bei der Modellerstellung diese Option auswählen, werden standardmäßig Neigungs-Scores im Modell-Nugget aktiviert. Sie können jedoch immer festlegen, dass Rohneigungs-Scores im Modell-Nugget aktiviert werden sollen, unabhängig davon, ob Sie sie im Modellierungsknoten auswählen.
- Beim Scoring des Modells werden Rohneigungs-Scores in einem Feld hinzugefügt, bei dem die Buchstaben *RP* an das Standard-Präfix angehängt sind. Wenn sich die Vorhersagen beispielsweise im Feld *\$R-Abwanderung* befinden, lautet der Name des Felds für den Neigungs-Score *\$RRP-Abwanderung*.

**Scores für korrigierte Neigung berechnen.** Rohneigungen basieren ausschließlich auf vom Modell angegebenen Schätzern. Beim Modell kann jedoch eine Überanpassung vorliegen, was zu übermäßig optimistischen Schätzern für die Neigung führt. Korrigierte Neigungen versuchen, dies zu kompensieren, indem untersucht wird, wie leistungsfähig das Modell bei den Test- bzw. Validierungspartitionen ist, und die Neigungen entsprechend angepasst werden, um einen besseren Schätzer zu erzeugen.

- Diese Einstellung ist nur möglich, wenn ein gültiges Partitionsfeld im Stream vorhanden ist. [Für weitere Informationen siehe Thema Partitionsknoten in Kapitel 4 in IBM SPSS Modeler 14.2- Quellen-, Prozess- und Ausgabeknoten.](#)
- Anders als rohe Konfidenz-Scores müssen Scores für die korrigierte Neigung bei der Erstellung des Modells berechnet werden; anderenfalls stehen Sie beim Scoring des Modell-Nuggets nicht zur Verfügung.
- Beim Scoring des Modells werden Scores für die korrigierte Neigung in einem Feld hinzugefügt, bei dem die Buchstaben *AP* an das Standard-Präfix angehängt sind. Wenn sich die Vorhersagen beispielsweise im Feld *\$R-Abwanderung* befinden, lautet der Name des Felds für den Neigungs-Score *\$RAP-Abwanderung*. Scores für die korrigierte Neigung stehen bei logistischen Regressionsmodellen nicht zur Verfügung.
- Bei der Berechnung der Scores für die korrigierte Neigung darf die für die Berechnung verwendete Test- bzw. Validierungspartition nicht ausbalanciert worden sein. Um dies zu vermeiden, müssen Sie darauf achten, dass in etwaigen weiter oben im Stream befindlichen Balancierungsknoten die Option Balancierung nur für Trainingsdaten durchführen ausgewählt wurde. [Für weitere Informationen siehe Thema Festlegen der Optionen für den Balancierungsknoten in Kapitel 3 in IBM SPSS Modeler 14.2- Quellen-, Prozess- und Ausgabeknoten.](#) Außerdem gilt: Wenn weiter oben im Stream eine komplexe Stichprobe gezogen wurde, werden dadurch die Scores für die korrigierte Neigung ungültig.
- Scores für die korrigierte Neigung stehen bei verstärkten Baum- und Regelmengenmodellen nicht zur Verfügung. [Für weitere Informationen siehe Thema Verbesserte C5.0-Modelle in Kapitel 6 auf S. 190.](#)

**Basierend auf.** Um Scores für die angepasste Neigung berechnen zu können, muss im Stream ein Partitionsfeld vorhanden sein. Sie können angeben, ob die Test- bzw. Validierungspartition für diese Berechnung verwendet werden soll. Um bestmögliche Ergebnisse zu erzielen, sollte die Test- bzw. Validierungspartition mindestens so viele Datensätze enthalten wie die

Partition, die zum Trainieren des ursprünglichen Modells verwendet wurde. [Für weitere Informationen siehe Thema Partitionsknoten in Kapitel 4 in IBM SPSS Modeler 14.2- Quellen-, Prozess- und Ausgabeknoten.](#)

## Neigungsbewertungen

Bei Modellen, die eine Vorhersage mit den Werten *Ja* und *Nein* ergeben, können Sie neben den standardmäßigen Vorhersage- und Konfidenzwerten auch Neigungs-Scores anfordern. Neigungs-Scores geben die Wahrscheinlichkeit eines bestimmten Ergebnisses bzw. einer bestimmten Antwort an. Beispiel:

Tabelle 3-1  
Neigungs-Scores

Kunde	Neigung, zu antworten
Karl Schmidt	35%
Marion Schmidt	15%

Neigungs-Scores stehen nur bei Modellen mit Flag-Zielen zur Verfügung und geben die Wahrscheinlichkeit des für das Feld definierten *Wahr*-Werts an, der in einem Quellen- oder Typknoten angegeben wurde. [Für weitere Informationen siehe Thema Angabe von Werten für ein Flag in Kapitel 4 in IBM SPSS Modeler 14.2- Quellen-, Prozess- und Ausgabeknoten.](#)

### Neigungs-Scores im Vergleich zu Konfidenz-Scores

Neigungs-Scores unterscheiden sich von Konfidenz-Scores, die für die aktuelle Vorhersage (*Ja* oder *Nein*) gelten. Wenn die Vorhersage beispielsweise *Nein* lautet, bedeutet eine hohe Konfidenz in Wirklichkeit eine hohe Wahrscheinlichkeit, dass es zu *keiner* Antwort kommt. Neigungs-Scores vermeiden diese Einschränkung, um einen leichteren Vergleich zwischen allen Datensätzen zu ermöglichen. So wird eine *Nein*-Vorhersage mit der Konfidenz *0,85* zu einer Rohneigung von *0,15* (d. h. *1 minus 0,85*).

Tabelle 3-2  
Konfidenz-Scores

Kunde	Prognose	Confidence
Karl Schmidt	Wird antworten	.35
Marion Schmidt	Wird nicht antworten	.85

### Neigungs-Scores

- Neigungs-Scores können auf der Registerkarte “Analysieren” im Modellierungsknoten oder auf der Registerkarte “Einstellungen” im Modell-Nugget aktiviert werden. Diese Funktionalität ist nur verfügbar, wenn das ausgewählte Ziel ein Flag-Feld ist. [Für weitere Informationen siehe Thema Analyseoptionen bei Modellierungsknoten auf S. 41.](#)
- Neigungs-Scores können je nach der verwendeten Ensemble-Methode auch vom Ensemble-Knoten berechnet werden. [Für weitere Informationen siehe Thema Ensemble-Knoten in Kapitel 4 in IBM SPSS Modeler 14.2- Quellen-, Prozess- und Ausgabeknoten.](#)

### **Berechnen der Scores für korrigierte Neigung**

Scores für korrigierte Neigung werden im Rahmen der Modellerstellung berechnet und stehen ansonsten nicht zur Verfügung. Nach der Modellerstellung wird das Modell mithilfe von Daten aus der Test- bzw. Validierungspartition gescort und es wird ein neues Modell erstellt, das die Scores für die korrigierte Neigung ausgeben soll. Dazu wird die Leistung des ursprünglichen Modells bei dieser Partition analysiert. Je nach Modelltyp kann eine von zwei Methoden zur Berechnung der Scores für die korrigierte Neigung verwendet werden.

- Bei Regelmengen- und Baummodellen werden Scores für die korrigierte Neigung durch erneute Berechnung der Häufigkeit der einzelnen Kategorien an jedem Baumknoten (bei Baummodellen) bzw. der Unterstützung und Konfidenz der einzelnen Regeln (bei Regelmengenmodellen) generiert. Die führt zu einem neuen Regelmengen- bzw. Baummodell, das zusammen mit dem ursprünglichen Modell gespeichert wird, um immer dann verwendet zu werden, wenn Scores für die korrigierte Neigung angefordert werden. Jedes Mal, wenn das ursprüngliche Modell auf neue Daten angewendet wird, kann das neue Modell anschließend auf die Rohneigungs-Scores angewendet werden, um die korrigierten Scores zu generieren.
- Bei anderen Modellen werden die durch Scores des ursprünglichen Modells anhand der Test- bzw. Validierungspartition erstellten Datensätze anschließend nach ihrem Rohneigungs-Score klassiert. Als Nächstes wird ein neuronales Netzwerkmodell trainiert, das eine nichtlineare Funktion definiert, die eine Zuordnung zwischen der Rohneigung in den einzelnen Klassen und der mittleren beobachteten Neigung in derselben Klasse erstellt. Wie zuvor für Baummodelle angemerkt, wird das resultierende neuronale Netzwerkmodell zusammen mit dem ursprünglichen Modell gespeichert und kann jedes Mal auf die Rohneigungs-Scores angewendet werden, wenn Scores für die korrigierte Neigung angefordert werden.

**Vorsicht bei fehlenden Werten in der Testpartition.** Die Verarbeitung fehlender Eingabewerte in der Test-/Validierungspartition variiert nach Modelltyp (mehr dazu unter "Einzelne Modell-Scoring-Algorithmen"). Das C5-Modell kann bei fehlenden Eingaben angepasste Neigungen nicht berechnen.

## **Modell-Nuggets**

Abbildung 3-7  
Modell-Nugget



Ein Model-Nugget ist ein Behälter für ein Modell, d. h. das Set von Regeln, Formeln oder Gleichungen, die das Ergebnis Ihrer Operationen zur Modellerstellung in IBM® SPSS® Modeler repräsentieren. Der Hauptzweck eines Nuggets ist das Scoring von Daten, um Vorhersagen zu generieren oder eine weitere Analyse der Modelleigenschaften zu erlauben. Durch Öffnen eines Modell-Nuggets am Bildschirm können Sie verschiedene Details zum Modell wie z. B. die relative Wichtigkeit der Eingabefelder beim Erstellen des Modells sehen. Zur Anzeige der Vorhersagen müssen Sie einen weiteren Prozess oder Ausgabeknoten anfügen und ausführen. [Für weitere Informationen siehe Thema Verwendung von Modell-Nuggets in Streams auf S. 68.](#)

Abbildung 3-8  
Modellverknüpfung vom Modellierungsknoten zum Modell-Nugget



Wenn Sie einen Modellierungsknoten erfolgreich ausführen, wird ein entsprechendes Modell-Nugget auf den Stream-Zeichenbereich platziert, wo es durch ein goldfarbenedes Rautensymbol repräsentiert wird (daher der Name "Nugget"). Auf der Stream-Zeichenfläche wird das Nugget mit einer Verbindung (einer durchgehenden Linie) zum nächsten passenden Knoten vor dem Modellierungsknoten sowie einer Verknüpfung (einer gepunkteten Linie) zum Modellierungsknoten angezeigt.

Das Nugget wird auch in die Modellpalette in der rechten oberen Ecke des SPSS Modeler-Fensters platziert. An beiden Positionen können Nuggets ausgewählt und durchsucht werden, um Details des Modells anzuzeigen.

Nuggets werden stets in der Modellpalette platziert, nachdem ein Modellierungsknoten erfolgreich ausgeführt wurde. Sie können eine Benutzeroption festlegen, die steuert, ob das Nugget zusätzlich auf der Stream-Zeichenfläche platziert wird. [Für weitere Informationen siehe Thema Festlegen von Benachrichtigungsoptionen in Kapitel 12 in IBM SPSS Modeler 14.2- Benutzerhandbuch.](#)

In den folgenden Themenabschnitten finden Sie Informationen zur Verwendung von Modell-Nuggets in SPSS Modeler. Wenn Sie ein tieferes Verständnis der verwendeten Algorithmen wünschen, lesen Sie im *SPSS Modeler-Algorithmushandbuch* nach, das Sie im Ordner *\Documentation* auf den DVD für IBM® SPSS® Modeler finden.

## Modellverknüpfungen

Standardmäßig wird ein Nugget auf der Zeichenfläche mit einer Verknüpfung zu dem Modellierungsknoten angezeigt, der das Nugget erstellt hat. Dies ist vor allem in komplexen Streams mit mehreren Nuggets nützlich und ermöglicht Ihnen das Nugget zu identifizieren, das von jedem Modellierungsknoten aktualisiert wird. Jede Verknüpfung enthält ein Symbol, um anzuzeigen, ob das Modell beim Ausführen des Modellierungsknotens ersetzt wird. [Für weitere Informationen siehe Thema Ersetzen eines Modells auf S. 49.](#)

### Definieren und Entfernen von Modellverknüpfungen

Sie können Verknüpfungen auf der Zeichenfläche manuell definieren und entfernen. Wenn Sie eine neue Verknüpfung definieren, ändert der Cursor seine Form zum Verknüpfungs-Cursor.

Abbildung 3-9  
Verknüpfungs-Cursor



**Definieren einer neuen Verknüpfung (Kontextmenü)**

- ▶ Klicken Sie mit der rechten Maustaste auf den Modellierungsknoten, von dem die Verknüpfung ausgehen soll.
- ▶ Wählen Sie Modellverknüpfung definieren aus dem Kontextmenü.
- ▶ Klicken Sie auf das Nugget, der das Ende der Verknüpfung darstellen soll.

**Definieren einer neuen Verknüpfung (Hauptmenü)**

- ▶ Klicken Sie auf den Modellierungsknoten, von dem die Verknüpfung ausgehen soll.
- ▶ Wählen Sie im Hauptmenü Folgendes:  
Bearbeiten > Knoten > Modellverknüpfung definieren
- ▶ Klicken Sie auf das Nugget, der das Ende der Verknüpfung darstellen soll.

**Entfernen einer vorhandenen Verknüpfung (Kontextmenü)**

- ▶ Klicken Sie mit der rechten Maustaste auf das Nugget am Ende der Verknüpfung.
- ▶ Wählen Sie Modellverknüpfung entfernen aus dem Kontextmenü.  
Alternative:
  - ▶ Klicken Sie mit der rechten Maustaste auf das Symbol in der Mitte der Verknüpfung.
  - ▶ Wählen Sie Verknüpfung entfernen aus dem Kontextmenü.

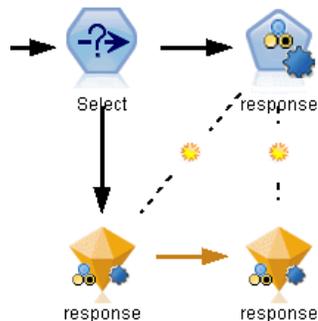
**Entfernen einer vorhandenen Verknüpfung (Hauptmenü)**

- ▶ Klicken Sie auf den Modellierungsknoten oder das Nugget, von dem Sie die Verknüpfung entfernen möchten.
- ▶ Wählen Sie im Hauptmenü Folgendes:  
Bearbeiten > Knoten > Modellverknüpfung entfernen

### Kopieren und Einfügen von Modellverknüpfungen

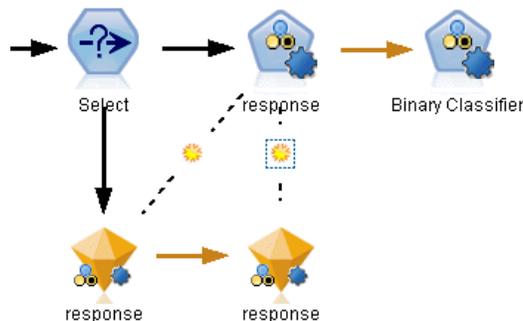
Wenn Sie ein verknüpftes Nugget ohne seinen Modellierungsknoten kopieren und in denselben Stream einfügen, wird das Nugget mit einer Verknüpfung zum Modellierungsknoten eingefügt. Die neue Verknüpfung hat denselben Modellersetzungsstatus (siehe Ersetzen eines Modells auf S. 49) wie die ursprüngliche Verknüpfung:

Abbildung 3-10  
Kopieren und Einfügen eines verknüpften Nuggets



Wenn Sie ein Nugget zusammen mit seinem verknüpften Modellierungsknoten kopieren und einfügen, wird die Verknüpfung beibehalten, wenn die Objekte in denselben oder in einen neuen Stream eingefügt werden:

Abbildung 3-11  
Kopieren und Einfügen eines verknüpften Nuggets



*Hinweis:* Wenn Sie ein verknüpftes Nugget ohne seinen Modellierungsknoten kopieren und in einen neuen Stream einfügen (oder in einen Superknoten, der den Modellierungsknoten nicht enthält), wird die Verknüpfung getrennt und nur das Nugget eingefügt.

### Modellverknüpfungen und Superknoten

Wenn Sie einen Superknoten so definieren, dass er entweder den Modellierungsknoten oder das Modell-Nugget eines verknüpften Modells (jedoch nicht beide) enthält, wird die Verknüpfung unterbrochen. Ein Erweitern des Superknotens stellt die Verknüpfung nicht wieder her. Dies ist nur möglich, indem Sie die Erstellung des Superknotens widerrufen.

## Ersetzen eines Modells

Sie können wählen, ob ein bestehendes Nugget bei der Neuausführung des Modellierungsknotens, der das Nugget erstellt hat, ersetzt (d. h. aktualisiert) wird. Wenn Sie die Ersetzungsoption deaktivieren, wird bei der erneuten Ausführung des Modellierungsknotens ein neues Nugget angelegt.

*Hinweis:* Das Ersetzen eines Modells unterscheidet sich vom Aktualisieren eines Modells, das sich auf das Aktualisieren eines Modells in einem Szenario bezieht. [Für weitere Informationen siehe Thema Modellaktualisierung in Kapitel 9 in IBM SPSS Modeler 14.2- Benutzerhandbuch.](#)

Jede Verknüpfung vom Modellierungsknoten zum Nugget enthält ein Symbol, um anzuzeigen, ob das Modell beim erneuten Ausführen des Modellierungsknotens ersetzt wird.

Abbildung 3-12  
Modellverknüpfung mit aktivierter Modellersetzung



Die Verknüpfung wird anfangs mit aktivierter Modellersetzung gezeigt, dargestellt durch das kleine Sonnensymbol in der Verknüpfung. In diesem Status wird bei erneuter Ausführung des Modellierungsknotens an einem Ende der Verknüpfung einfach das Nugget am anderen Ende aktualisiert.

Abbildung 3-13  
Modellverknüpfung mit deaktivierter Modellersetzung



Wenn die Modellersetzung deaktiviert ist, wird das Verknüpfungssymbol durch einen grauen Punkt ersetzt. In diesem Status wird bei erneuter Ausführung des Modellierungsknotens an einem Ende der Verknüpfung der Zeichenfläche eine neue, aktualisierte Version des Nuggets hinzugefügt.

In beiden Fällen wird das vorhandene Nugget in der Modellpalette aktualisiert oder ein neues Nugget hinzugefügt, abhängig von der Systemoption Vorheriges Modell ersetzen. [Für weitere Informationen siehe Thema Festlegen von Benachrichtigungsoptionen in Kapitel 12 in IBM SPSS Modeler 14.2- Benutzerhandbuch.](#)

### Reihenfolge der Ausführung

Wenn Sie einen Stream mit mehreren Verzweigungen ausführen, die Modell-Nuggets enthalten, wird der Stream zunächst evaluiert, um sicherzustellen, dass eine Verzweigung mit aktivierter Modellersetzung vor einer Verzweigung ausgeführt wird, die das resultierende Modell-Nugget verwendet.

Bei komplexeren Anforderungen können Sie die Ausführungsreihenfolge manuell durch Erstellung eines Skripts festlegen.

### **Ändern der Modellersetzungseinstellung**

So ändern Sie die Einstellung für Modellersetzung:

- ▶ Klicken Sie mit der rechten Maustaste auf das Symbol in der Verknüpfung.
- ▶ Wählen Sie wie gewünscht Modellersetzung aktivieren (deaktivieren).

*Hinweis:* Die Modellersetzungseinstellung einer Modellverknüpfung überschreibt die Einstellung auf der Registerkarte “Benachrichtigungen” des Dialogfelds “Benutzeroptionen” (Extras > Optionen > Benutzeroptionen).

### **Die Modellpalette**

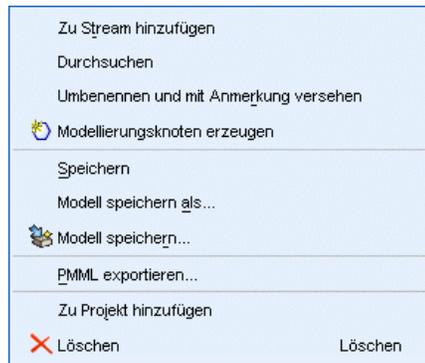
Mit der Modellpalette (auf der Registerkarte “Modelle” im Manager-Fenster) können Sie Modell-Nuggets auf verschiedene Weise verwenden, untersuchen und bearbeiten.

Abbildung 3-14  
Modellpalette



Wenn Sie mit der rechten Maustaste auf ein Modell-Nugget in der Modellpalette klicken, wird ein Kontextmenü mit folgenden Optionen geöffnet:

Abbildung 3-15  
Kontextmenü für Modell-Nuggets

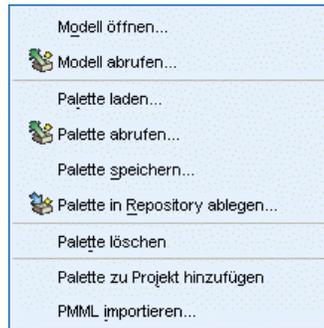


- **Zu Stream hinzufügen.** Fügt das generierte Modell-Nugget zum derzeit aktiven Stream hinzu. Wenn im Stream ein Knoten ausgewählt wurde, wird das Modell-Nugget mit dem ausgewählten Knoten verbunden, wenn eine solche Verbindung möglich ist. Anderenfalls erfolgt die Verbindung zum nächsten möglichen Knoten. Das Nugget wird mit einer Verknüpfung zum Modellierungsknoten gezeigt, von dem aus das Modell erstellt wurde, sofern sich der Knoten noch im Stream befindet.
- **Durchsuchen.** Öffnet den Modell-Browser für das Nugget.
- **Umbenennen und mit Anmerkung versehen.** Ermöglicht das Umbenennen des Modell-Nuggets und/oder die Bearbeitung der Anmerkung für das Nugget.
- **Erzeugen eines Modellierungsknotens.** Wenn Sie ein Modell-Nugget ändern oder aktualisieren möchten und der zum Erstellen des Modells verwendete Stream nicht verfügbar ist, können Sie mithilfe dieser Option erneut einen Modellierungsknoten mit denselben Optionen erzeugen, die zum Erstellen des ursprünglichen Modells verwendet wurden.
- **Modell speichern, Modell speichern unter.** Speichert das Modell-Nugget in einer externen Binärdatei für generierte Modelle (.gm).
- **Modell speichern.** Speichert das Modell-Nugget in IBM® SPSS® Collaboration and Deployment Services Repository. [Für weitere Informationen siehe Thema Informationen zu IBM SPSS Collaboration and Deployment Services Repository in Kapitel 9 in IBM SPSS Modeler 14.2- Benutzerhandbuch.](#)
- **PMML exportieren.** Exportiert das Modell-Nugget als Predictive Model Markup Language (PMML), die zum Scoring neuer Daten außerhalb von IBM® SPSS® Modeler verwendet werden kann. PMML exportieren ist für alle generierten Modellknoten verfügbar. *Hinweis:* Für die Verwendung dieser Funktion ist eine IBM® SPSS® Modeler Server-Lizenz erforderlich. [Für weitere Informationen siehe Thema Festlegen von PMML-Exportoptionen in Kapitel 12 in IBM SPSS Modeler 14.2- Benutzerhandbuch.](#)
- **Zu Projekt hinzufügen.** Speichert das generierte Modell-Nugget und fügt es zum aktuellen Projekt hinzu. Auf der Registerkarte “Klassen” wird das Nugget zum Ordner “Generierte Modelle” hinzugefügt. Auf der Registerkarte “CRISP-DM” wird er zur Standard-Projektphase

hinzugefügt. (Unter [Einrichten der Standard-Projektphase](#) finden Sie Informationen zum Ändern der Standard-Projektphase.)

- **Löschen.** Löscht das Modell-Nugget aus der Palette.

Abbildung 3-16  
Kontextmenü für Modellpalette



Wenn Sie mit der rechten Maustaste auf einen nicht belegten Bereich in der Modellpalette klicken, wird ein Kontextmenü mit folgenden Optionen geöffnet:

- **Modell öffnen.** Lädt ein Modell-Nugget, das zuvor in SPSS Modeler erstellt wurde.
- **Modell abrufen.** Ruft ein gespeichertes Modell aus einem IBM SPSS Collaboration and Deployment Services-Repository ab.
- **Palette laden.** Lädt eine gespeicherte Modellpalette aus einer externen Datei.
- **Palette abrufen.** Ruft eine gespeicherte Modellpalette aus einem IBM SPSS Collaboration and Deployment Services-Repository ab.
- **Palette speichern.** Speichert den gesamten Inhalt der Modellpalette in einer externen Datei für generierte Modellpaletten (.gen).
- **Palette speichern.** Speichert den gesamten Inhalt der Modellpalette in einem IBM SPSS Collaboration and Deployment Services-Repository.
- **Palette löschen.** Löscht alle Nuggets aus der Palette.
- **Palette zu Projekt hinzufügen.** Speichert die Modellpalette und fügt sie dem aktuellen Projekt hinzu. Auf der Registerkarte “Klassen” wird das Nugget zum Ordner “Generierte Modelle” hinzugefügt. Auf der Registerkarte “CRISP-DM” wird er zur Standard-Projektphase hinzugefügt.
- **PMML importieren.** Lädt ein Modell aus einer externen Datei. Sie können PMML-Modelle öffnen, durchsuchen und scoren, die von IBM® SPSS® Statistics oder anderen Anwendungen, die dieses Format unterstützen, erstellt wurden. [Für weitere Informationen siehe Thema Importieren und Exportieren von Modellen als PMML in IBM SPSS Modeler 14.2-Benutzerhandbuch.](#)

## **Durchsuchen von Modell-Nuggets**

Mit den Browsern für Modell-Nuggets können Sie die Ergebnisse Ihrer Modelle überprüfen und verwenden. Über den Browser können Sie das generierte Modell speichern, drucken oder exportieren, die Modellzusammenfassung überprüfen und Anmerkungen für das Modell

anzeigen oder bearbeiten. Bei einigen Typen von Modell-Nuggets können Sie auch neue Knoten generieren, beispielsweise Filterknoten oder Regelmengenknoten. Bei einigen Modellen können Sie außerdem Modellparameter, wie Regeln oder Cluster-Zentren, anzeigen. Bei einigen Modelltypen (baumbasierten Modellen und Cluster-Modellen) können Sie eine grafische Darstellung der Struktur des Modells anzeigen. Die Steuerelemente für die Verwendung der Modell-Nugget-Browser sind unten beschrieben.

### **Menüs**

**Menü "Datei"**. Alle Modell-Nuggets weisen ein Dateimenü auf, das eine Untergruppe der folgenden Optionen enthält:

- **Speicherknoten**. Speichert das Modell-Nugget in einer Knotendatei (.nod).
- **Knoten speichern**. Speichert das Modell-Nugget in einem IBM SPSS Collaboration and Deployment Services-Repository.
- **Kopf-/Fußzeile**. Erlaubt die Bearbeitung der Kopf- und Fußzeile der Seite zum Drucken über das Nugget.
- **Seite einrichten**. Erlaubt die Bearbeitung der Seiteneinrichtung zum Drucken über das Nugget.
- **Druckvorschau**. Zeigt als Vorschau an, wie das Nugget beim Ausdruck aussieht. Wählen Sie im Untermenü die Informationen aus, die in der Vorschau angezeigt werden sollen.
- **Drucken**. Druckt den Inhalt des Nuggets. Wählen Sie im Untermenü die Informationen aus, die gedruckt werden sollen.
- **Druckansicht**. Druckt die aktuelle Ansicht oder alle Ansichten.
- **Text exportieren**. Exportiert den Inhalt des Nuggets in eine Textdatei. Wählen Sie im Untermenü die Informationen aus, die exportiert werden sollen.
- **HTML exportieren**. Exportiert den Inhalt des Nuggets in eine HTML-Datei. Wählen Sie im Untermenü die Informationen aus, die exportiert werden sollen.
- **PMML exportieren**. Exportiert das Modell als Predictive Model Markup Language (PMML), die mit anderen PMML-kompatiblen Softwareprodukten verwendet werden kann. *Hinweis:* Für die Verwendung dieser Funktion ist eine IBM® SPSS® Modeler Server-Lizenz erforderlich. [Für weitere Informationen siehe Thema Festlegen von PMML-Exportoptionen in Kapitel 12 in IBM SPSS Modeler 14.2- Benutzerhandbuch.](#)
- **SQL exportieren**. Exportiert das Modell als SQL (Structured Query Language), die mit anderen Datenbanken bearbeitet und verwendet werden kann.

*Hinweis:* Der SQL-Export ist nur aus folgenden Modellen möglich: C5, C&RT, CHAID, QUEST, Lineare Regression, Logistische Regression, Netzwerk, Faktor und Entscheidungsliste. [Für weitere Informationen siehe Thema Knoten zur Unterstützung der SQL-Erzeugung in Kapitel 6 in IBM SPSS Modeler Server 14.2-Verwaltungs- und Leistungshandbuch.](#)

**Menü "Generieren"**. Die meisten Modell-Nuggets weisen auch das Menü "Generieren" auf, mit dem Sie auf der Grundlage des Modell-Nuggets neue Knoten erstellen können. Die über dieses Menü verfügbaren Optionen hängen vom durchsuchten Modelltyp ab. Einzelheiten zu den Elementen, die aus einem bestimmten Modell generiert werden können, finden Sie in den Informationen zum jeweils generierten Modell-Nugget-Typ.

**Menü "Ansicht".** Auf der Registerkarte "Modell" eines Nugget ermöglicht dieses Menü das Ein- bzw. Ausblenden der verschiedenen Visualisierungs-Symbolleisten, die im aktuellen Modus verfügbar sind. Um alle Symbolleisten verfügbar zu machen, wählen Sie in der Symbolleiste "Allgemein" die Option "Bearbeitungsmodus" (das Pinselsymbol).

**Schaltfläche "Vorschau"** Einige Modell-Nuggets verfügen über die Schaltfläche "Vorschau". Sie ermöglicht die Anzeige eines Auszugs der Modelldaten, inklusive der vom Modellierungsprozess erstellten Zusatzfelder. Die Standardanzahl der angezeigten Zeilen ist 10. Sie können diese Einstellung jedoch in den Stream-Eigenschaften ändern. [Für weitere Informationen siehe Thema Festlegen von Optionen für Streams in Kapitel 5 in IBM SPSS Modeler 14.2- Benutzerhandbuch.](#)

**Schaltfläche "Zu aktuellem Projekt hinzufügen".** Speichert das generierte Modell-Nugget und fügt es zum aktuellen Projekt hinzu. Auf der Registerkarte "Klassen" wird das Nugget zum Ordner "Generierte Modelle" hinzugefügt. Auf der Registerkarte "CRISP-DM" wird er zur Standard-Projektphase hinzugefügt. (Unter [Einrichten der Standard-Projektphase](#) finden Sie Informationen zum Ändern der Standard-Projektphase.)

## **Modell-Nugget-Übersicht/-Information**

Auf der Registerkarte "Übersicht" oder der Informationsansicht für ein Modell-Nugget werden Informationen über die Felder, die Aufbaueinstellungen und die Modellschätzung angezeigt. Die Ergebnisse werden in einer Baumansicht dargestellt, die durch Klicken auf bestimmte Elemente erweitert bzw. reduziert werden kann.

**Analyse.** Zeigt Informationen zum Modell an. Die konkreten Details variieren nach Modelltyp und werden jeweils im Abschnitt zu den einzelnen Modell-Nuggets behandelt. Wenn Sie einen Analyseknotten ausgeführt haben, der an diesen Modellierungsknoten angehängt ist, werden die Informationen aus dieser Analyse ebenfalls in diesem Abschnitt angezeigt. [Für weitere Informationen siehe Thema Analyseknotten in Kapitel 6 in IBM SPSS Modeler 14.2- Quellen-, Prozess- und Ausgabeknoten.](#)

**Felder.** Listet die Felder auf, die bei der Erstellung des Modells als Ziel bzw. als Eingaben verwendet werden. Listet bei aufgeteilten Modellen außerdem die Felder auf, welche die Aufteilung bestimmen.

**Aufbaueinstellungen/-optionen.** Enthält Informationen zu den beim Aufbau des Modells verwendeten Einstellungen.

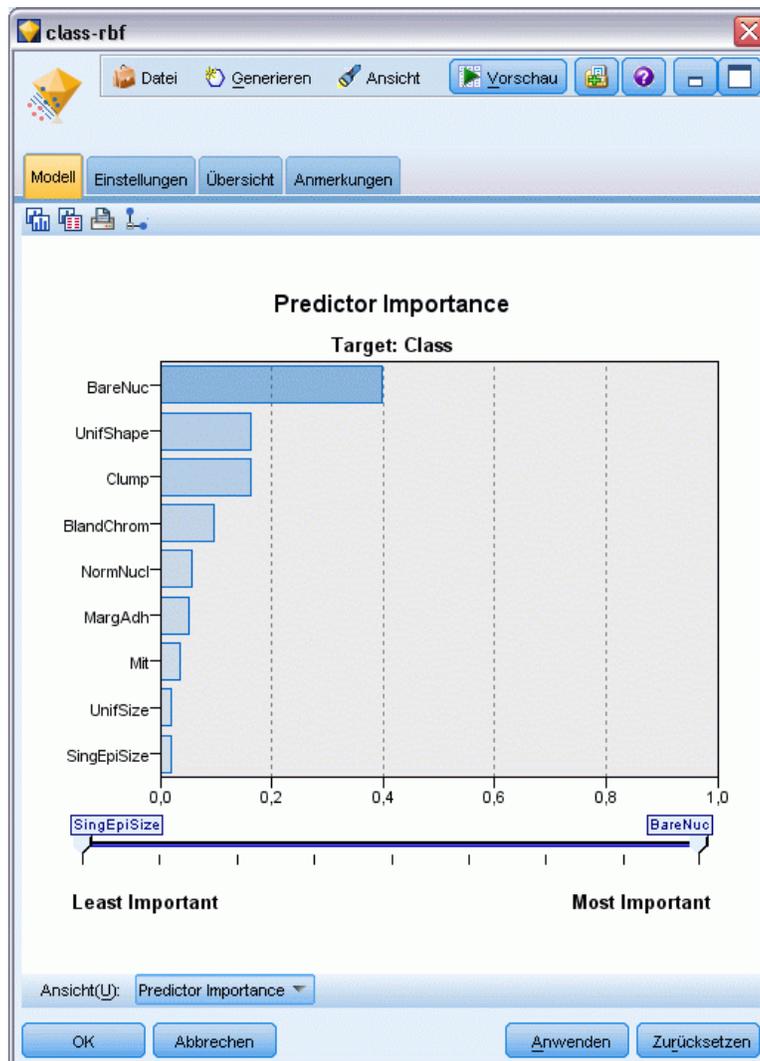
**Trainingsübersicht.** In diesem Abschnitt werden folgende Informationen angezeigt: der Typ des Modells, der für seine Erstellung verwendete Stream, der Benutzer, der ihn erstellt hat, der Erstellungszeitpunkt und die für den Aufbau des Modells benötigte Zeit.

## **Bedeutsamkeit des Prädiktors**

Normalerweise ist es sinnvoll, die Modellierungsbemühungen auf die wichtigsten Prädiktorfelder zu konzentrieren und diejenigen Prädiktoren zu verwerfen bzw. zu ignorieren, die die geringste Bedeutung haben. Dabei unterstützt Sie das Diagramm für die Bedeutsamkeit der Prädiktoren, da es die relative Bedeutsamkeit der einzelnen Prädiktoren für das Modell angibt. Da die Werte

relativ sind, beträgt die Summe der Werte aller Prädiktoren im Diagramm 1,0. Die Bedeutsamkeit der Prädiktoren bezieht sich nicht auf die Genauigkeit des Modells. Sie bezieht sich lediglich auf die Bedeutsamkeit der einzelnen Prädiktoren für eine Vorhersage und nicht auf die Genauigkeit der Vorhersage.

Abbildung 3-17  
Bedeutsamkeit der Prädiktoren - Diagramm



Die Bedeutsamkeit von Prädiktoren steht für Modelle zur Verfügung, die ein angemessenes statistisches Maß für Bedeutsamkeit erstellen. Dazu gehören neuronale Netzwerke, Entscheidungsbäume (C&RT-Baum, C5.0, CHAID und QUEST), Bayes-Netzwerke, Diskriminanz, SVM- und SLRM-Modelle, lineare und logistische Regression, verallgemeinerte lineare Modelle und Nächste-Nachbarn-(KNN)Modelle. Für die meisten dieser Modelle kann die Bedeutsamkeit der Prädiktoren auf der Registerkarte "Analysieren" im Modellierungsknoten aktiviert werden. [Für weitere Informationen siehe Thema Analyseoptionen](#)

bei [Modellierungsknoten auf S. 41](#). Informationen zu KNN-Modellen finden Sie unter Nachbarn auf S. 487.

*Anmerkung:* Die Bedeutsamkeit von Prädiktoren wird von aufgeteilten Modellen nicht unterstützt. Die Einstellungen der Bedeutsamkeit von Prädiktoren werden bei der Erstellung von aufgeteilten Modellen ignoriert. [Für weitere Informationen siehe Thema Erstellung von aufgeteilten Modellen auf S. 31](#).

Die Berechnung der Bedeutsamkeit der Prädiktoren kann erheblich länger dauern als die Modellerstellung, insbesondere bei großen Daten-Sets. Die Berechnung dauert bei SVM und logistischer Regression länger als bei anderen Modellen und ist daher für diese Modelle standardmäßig deaktiviert. Bei Verwendung eines Daten-Sets mit einer großen Anzahl an Prädiktoren kann ein anfängliches Screening mit einem Merkmalsauswahlknoten zu schnelleren Ergebnissen führen (siehe unten).

- Die Bedeutsamkeit der Prädiktoren wird aus der Testpartition berechnet, sofern verfügbar. Anderenfalls werden die Trainingsdaten verwendet. [Für weitere Informationen siehe Thema Partitionsknoten in Kapitel 4 in IBM SPSS Modeler 14.2- Quellen-, Prozess- und Ausgabeknoten](#).
- Bei SLRM-Modellen ist die Bedeutsamkeit der Prädiktoren verfügbar, sie wird jedoch durch den SLRM-Algorithmus berechnet. [Für weitere Informationen siehe Thema SLRM-Modell-Nuggets in Kapitel 14 auf S. 469](#).
- Mit den Diagrammtools von IBM® SPSS® Modeler können Sie in das Diagramm eingreifen, es bearbeiten und speichern.
- Optional können Sie anhand der Informationen im Diagramm für die Bedeutsamkeit der Prädiktoren einen Filterknoten generieren. [Für weitere Informationen siehe Thema Filtern von Variablen auf der Grundlage der Bedeutsamkeit auf S. 57](#).

### ***Bedeutsamkeit der Prädiktoren und Merkmalsauswahl***

In einigen Fällen kann es so aussehen, als ob das in einem Modell-Nugget angezeigte Diagramm für die Bedeutsamkeit der Prädiktoren zu ähnlichen Ergebnissen führt wie der Merkmalsauswahlknoten. Während bei der Merkmalsauswahl die einzelnen Eingabefelder hinsichtlich der Stärke ihrer Beziehung zum angegebenen Ziel unabhängig von anderen Eingaben in Ränge eingeteilt werden, gibt das Diagramm für die Bedeutsamkeit der Prädiktoren die relative Wichtigkeit der einzelnen Eingaben für *dieses* konkrete Modell an. Daher führt die Merkmalsauswahl beim Screening der Eingaben zu einem konservativeren Ergebnis. Wenn beispielsweise sowohl *Berufsbezeichnung* als auch *Berufskategorie* in einem starken Zusammenhang zum Gehalt stehen, zeigt die Merkmalsauswahl an, dass beide von Bedeutung sind. Bei der Modellierung werden jedoch auch Interaktionen (Wechselwirkungen) und Korrelationen berücksichtigt. Dies kann dazu führen, dass nur eine von zwei Eingaben verwendet wird, wenn beide größtenteils dieselben Informationen bieten. In der Praxis ist die Merkmalsauswahl am nützlichsten für ein erstes Screening, insbesondere beim Umgang mit großen Datensets mit zahlreichen Variablen, während die Bedeutsamkeit der Prädiktoren bei der Feinabstimmung des Modells nützlicher ist.

### **Filtern von Variablen auf der Grundlage der Bedeutsamkeit**

Optional können Sie anhand der Informationen im Diagramm für die Bedeutsamkeit der Prädiktoren einen Filterknoten generieren.

Markieren Sie ggf. die Prädiktoren, die Sie in das Diagramm aufnehmen möchten, und wählen Sie folgende Optionen aus den Menüs aus:

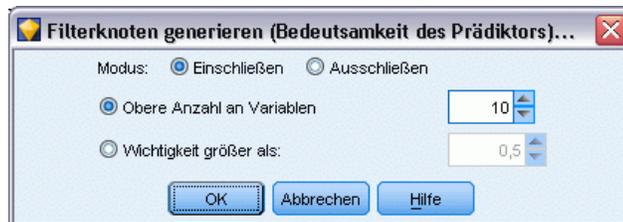
Erzeugen > Filterknoten (Bedeutsamkeit des Prädiktors)

ODER

> Felddauswahl (Bedeutsamkeit der Prädiktoren)

Abbildung 3-18

*Filtern von Prädiktoren auf der Grundlage der Bedeutsamkeit*



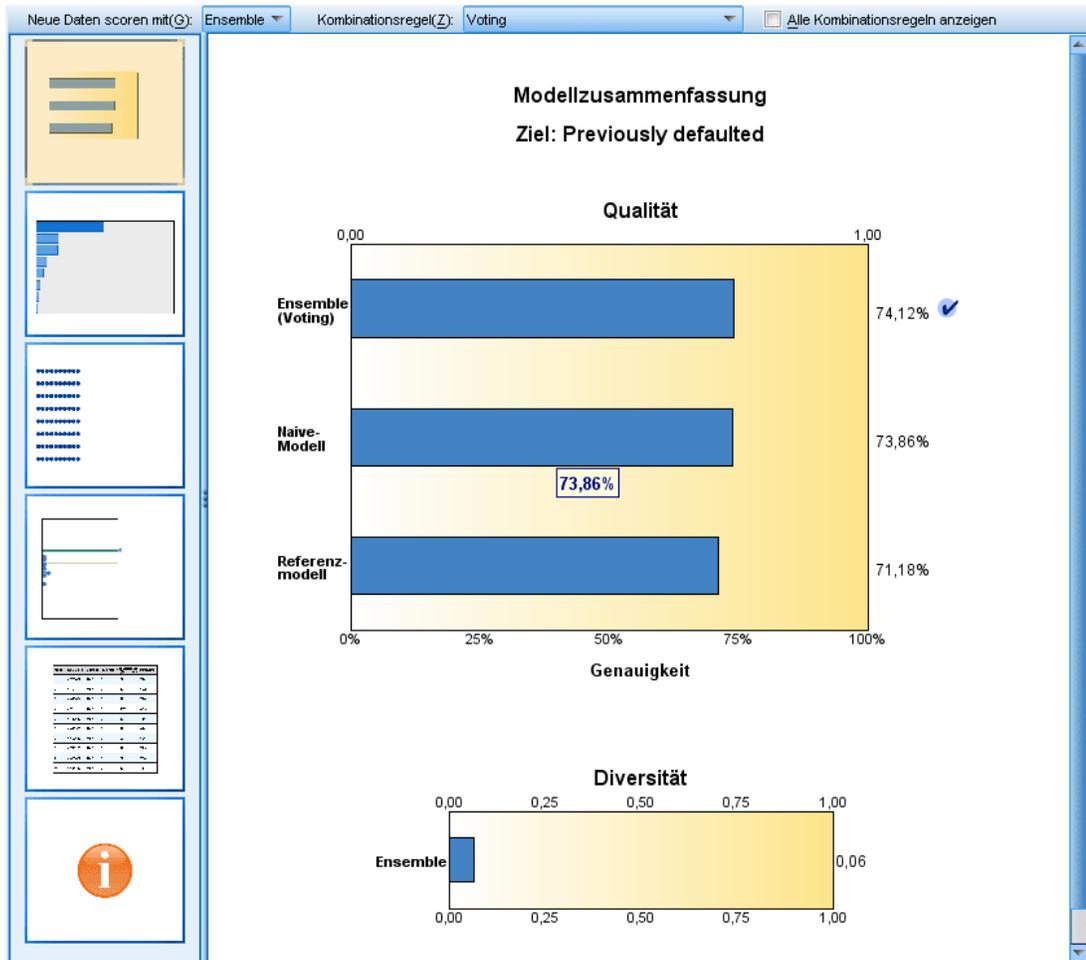
**Obere Anzahl an Variablen.** Schließt die wichtigsten Prädiktoren bis zur angegebenen Anzahl ein bzw. aus.

**Wichtigkeit größer als.** Schließt alle Prädiktoren ein bzw. aus, deren relative Wichtigkeit den angegebenen Wert übersteigt.

### **Modelle für Ensembles**

Das Modell für ein Ensemble bietet Informationen zu den Komponentenmodellen im Ensemble und zur Leistung des Ensembles als Ganzes.

Abbildung 3-19  
Ansicht "Modellzusammenfassung"



In der (von der Ansicht unabhängigen) Hauptsymboleiste können Sie auswählen, ob das Ensemble oder ein Referenzmodell für die Bewertung (Scoring) verwendet werden soll. Wenn das Ensemble für die Bewertung verwendet wird, können Sie auch die Kombinationsregel auswählen. Für diese Änderungen ist keine erneute Ausführung des Modells erforderlich. Die getroffene Auswahl wird jedoch zur Bewertung (Scoring) und/oder zur nachfolgenden Modellauswertung im Modell (Nugget) gespeichert. Außerdem wirkt sie sich auf die aus dem Ensemble-Viewer exportierte PMML aus.

**Kombinationsregel.** Beim Scoring eines Ensembles wird diese Regel angewendet, um die vorhergesagten Werte aus den Basismodellen für die Berechnung des Score-Werts für das Ensemble zu kombinieren.

- Ensemble-Vorhersagewerte für **kategoriale** Ziele können unter Verwendung von Voting, der höchsten Wahrscheinlichkeit oder der höchsten mittleren Wahrscheinlichkeit kombiniert werden. Mit **Voting** wird die Kategorie ausgewählt, die in der Menge aller Basismodelle am häufigsten die höchste Wahrscheinlichkeit aufweist. Mit **Höchste Wahrscheinlichkeit** wird

die Kategorie ausgewählt, die in der Menge aller Basismodelle die höchste Wahrscheinlichkeit überhaupt aufweist. Mit **Höchste mittlere Wahrscheinlichkeit** wird die Kategorie mit dem höchsten Wert ausgewählt, wenn der Mittelwert der Kategoriewahrscheinlichkeiten aus der Menge aller Basismodelle berechnet wird.

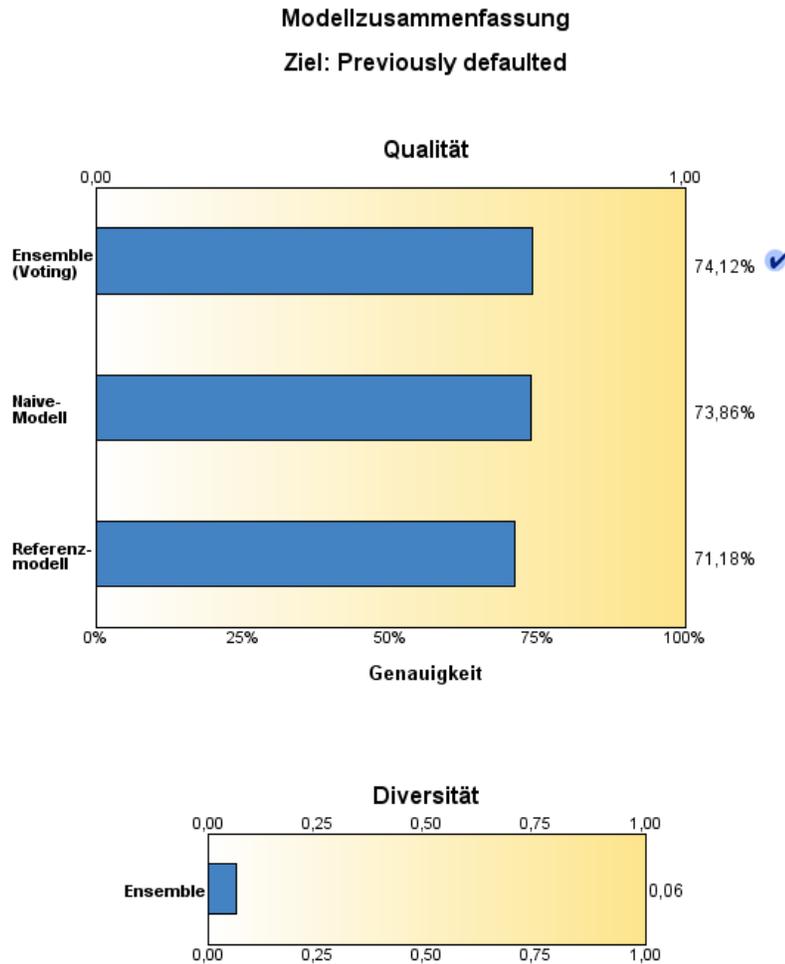
- Ensemble-Vorhersagewerte für **stetige** Ziele können unter Verwendung des Mittelwerts oder Medians der Vorhersagewerte aus den Basismodellen kombiniert werden.

Die Standardvorgabe beruht auf den während der Modellerstellung angegebenen Spezifikationen. Durch Ändern der Kombinationsregel wird die Modellgenauigkeit neu berechnet und alle Ansichten der Modellgenauigkeit werden aktualisiert. Das Diagramm für die Bedeutsamkeit des Prädiktors wird ebenfalls aktualisiert. Dieses Steuerelement ist deaktiviert, wenn das Referenzmodell für die Bewertung (Scoring) ausgewählt wurde.

**Alle Kombinationsregeln anzeigen.** Bei Auswahl dieser Option werden Ergebnisse für alle verfügbaren Kombinationsregeln im Diagramm zur Modellqualität angezeigt. Das Diagramm "Komponentenmodellgenauigkeit" wird ebenfalls aktualisiert und zeigt nun Bezugslinien für die einzelnen Voting-Methoden.

## Modellzusammenfassung

Abbildung 3-20  
Ansicht "Modellzusammenfassung"



Mit der Ansicht "Modellzusammenfassung" erhalten Sie eine momentane, übersichtliche Zusammenfassung über Ensemble-Qualität und -Vielfalt.

**Qualität.** Dieses Diagramm zeigt die Genauigkeit des endgültigen Modells im Vergleich mit einem Referenzmodell und einem naiven Modell. Die Genauigkeit wird nach dem Prinzip "größer ist besser" dargestellt. Das "beste" Modell hat die höchste Genauigkeit. Bei kategorialen Zielen ist die Genauigkeit einfach der Prozentsatz der Datensätze, für den der vorhergesagte Wert mit dem beobachteten Wert übereinstimmt. Bei stetigen Zielen ist die Genauigkeit 1 minus dem Verhältnis zwischen dem mittleren absoluten Fehler bei der Vorhersage (Durchschnitt der Absolutwerte der vorhergesagten Werte minus beobachtete Werte) und dem Bereich der vorhergesagten Werte (größter vorhergesagter Wert minus kleinster vorhergesagter Wert).

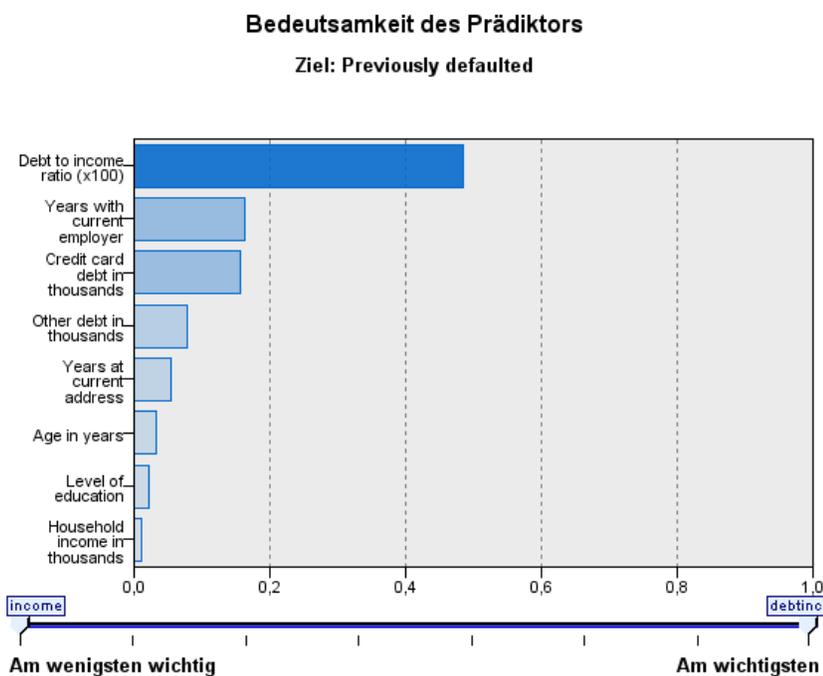
Bei Ensembles mit Bootstrap-Aggregation ist das Referenzmodell ein Standardmodell, das auf der gesamten Trainingspartition beruht. Bei verstärkten Ensembles ist das Referenzmodell das Modell der ersten Komponente.

Das naive Modell stellt die Genauigkeit dar, die bestünde, wenn kein Modell erstellt würde, und weist alle Datensätze der Modalkategorie zu. Für stetige Ziele wird das naive Modell nicht berechnet.

**Diversität.** Das Diagramm zeigt die “Meinungsdiversität” unter den zum Erstellen des Ensembles verwendeten Komponentenmodellen an, dargestellt nach dem Prinzip “größer ist besser”. Es handelt sich hierbei um die Messung der Variation der Vorhersagen zwischen den verschiedenen Basismodellen. Für verstärkte Ensemble-Modelle steht die Option “Diversität” nicht zur Verfügung und sie wird auch nicht für stetige Ziele angezeigt.

### Bedeutsamkeit des Prädiktors

Abbildung 3-21  
Ansicht “Bedeutsamkeit des Prädiktors”

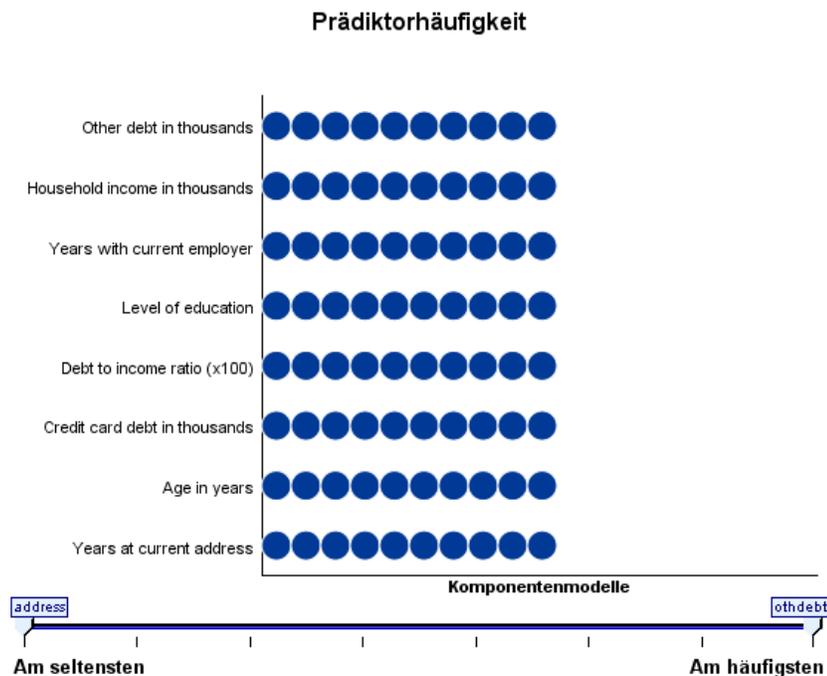


Normalerweise ist es sinnvoll, die Modellierungsbemühungen auf die wichtigsten Prädiktorfelder zu konzentrieren und diejenigen Prädiktoren zu verwerfen bzw. zu ignorieren, die die geringste Bedeutung haben. Dabei unterstützt Sie das Diagramm für die Bedeutsamkeit der Prädiktoren, da es die relative Bedeutsamkeit der einzelnen Prädiktoren für das Modell angibt. Da die Werte relativ sind, beträgt die Summe der Werte aller Prädiktoren im Diagramm 1,0. Die Bedeutsamkeit der Prädiktoren bezieht sich nicht auf die Genauigkeit des Modells. Sie bezieht sich lediglich auf die Bedeutsamkeit der einzelnen Prädiktoren für eine Vorhersage und nicht auf die Genauigkeit der Vorhersage.

Die Bedeutsamkeit des Prädiktors steht nicht für alle Ensemble-Modelle zur Verfügung. Die Menge der Prädiktoren (Einflussvariablen) kann zwischen verschiedenen Komponentenmodellen variieren, die Wichtigkeit kann jedoch für Prädiktoren berechnet werden, die in mindestens einem Komponentenmodell verwendet werden.

### Prädiktorhäufigkeit

Abbildung 3-22  
Ansicht "Prädiktorhäufigkeit"

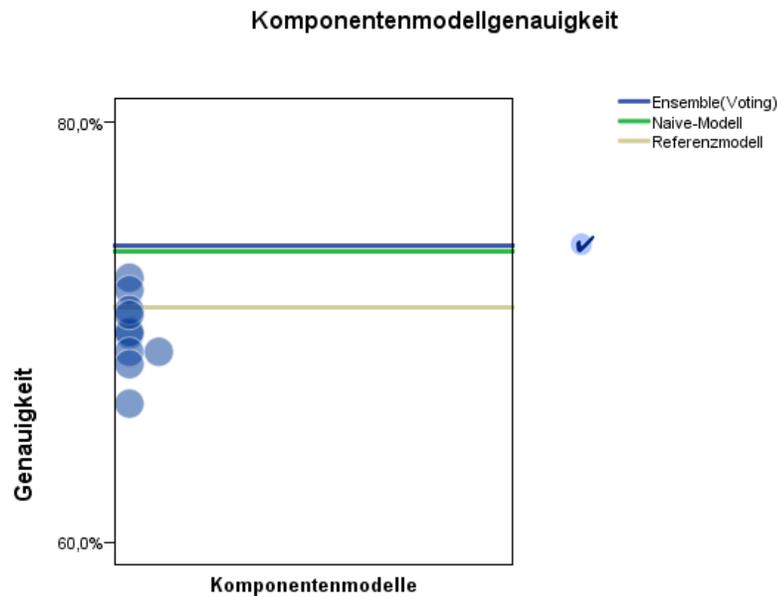


Die Menge der Prädiktoren (Einflussvariablen) kann aufgrund der Auswahl der verwendeten Modellierungsmethode bzw. der Auswahl der Prädiktoren zwischen verschiedenen Komponentenmodellen variieren. Das Diagramm "Prädiktorhäufigkeit" ist ein Punktdiagramm, das die Verteilung der Prädiktoren in den verschiedenen Komponentenmodellen im Ensemble zeigt. Jeder Punkt steht für eine oder mehrere Komponenten, die den Prädiktor enthält/enhalten. Prädiktoren werden auf der y-Achse dargestellt und in absteigender Reihenfolge ihrer Häufigkeit sortiert. Somit ist der oberste Prädiktor derjenige, der in der größten Anzahl an Komponentenmodellen verwendet wurde, und der unterste derjenige, der in den wenigsten Modellen verwendet wurde. Die ersten 10 Prädiktoren werden angezeigt.

Die am häufigsten vorkommenden Prädiktoren (Einflussvariablen) sind normalerweise auch die wichtigsten. Dieses Diagramm ist nicht brauchbar bei Methoden, bei denen die Menge der Prädiktoren zwischen den verschiedenen Komponentenmodellen variieren kann.

### Komponentenmodellgenauigkeit

Abbildung 3-23  
Ansicht "Komponentenmodellgenauigkeit"



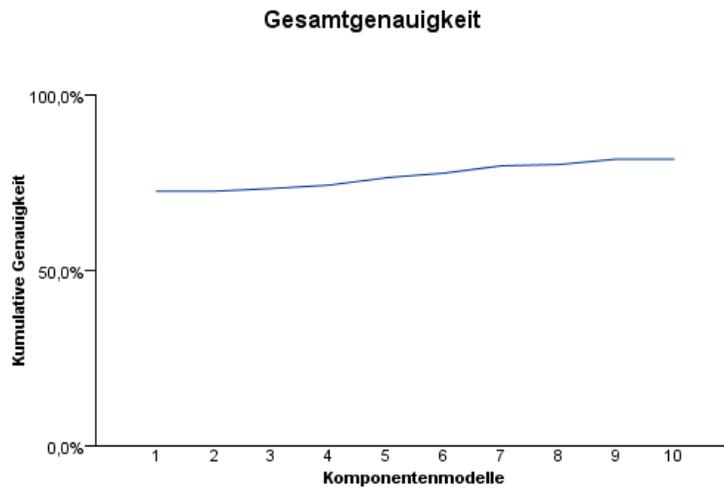
Bei diesem Diagramm handelt es sich um ein Punktdiagramm der Vorhersagegenauigkeit für Komponentenmodelle. Jeder Punkt steht für ein oder mehrere Komponentenmodelle, wobei der Genauigkeitsgrad auf der y-Achse dargestellt wird. Fahren Sie mit der Maus über einen Punkt, um Informationen zum zugehörigen Einzelkomponentenmodell abzurufen.

**Bezugslinien.** In diesem Diagramm werden farbkodierte Linien für das Ensemble sowie für das Referenzmodell und die naiven Modelle angezeigt. Neben der Linie, die zu dem für die Bewertung (Scoring) verwendeten Modell gehört, wird ein Häkchen angezeigt.

**Interaktivität.** Das Diagramm wird aktualisiert, wenn Sie die Kombinationsregel ändern.

**Verstärkte Ensembles.** Für verstärkte Ensembles wird ein Liniendiagramm angezeigt.

Abbildung 3-24  
Ansicht "Gesamtgenauigkeit," verstärktes Ensemble



### Komponentenmodelldetails

Abbildung 3-25  
Ansicht "Komponentenmodelldetails"

**Komponentenmodelldetails**

Modell	Genauigkeit	Methode	Prädiktoren	Modellgröße (Synapsen)	Datensätze
1	72,6%		8	77	700
2	72,0%		8	107	692
3	69,9%		8	92	708
4	70,0%		8	107	685
5	71,1%		8	107	706
6	69,1%		8	92	690
7	69,1%		8	92	696
8	70,8%		8	122	703
9	66,6%		8	62	726
10	68,5%		8	107	701

In der Tabelle werden Informationen zu Komponentenmodellen, nach Zeile aufgelistet, angezeigt. Standardmäßig werden die Komponentenmodelle in aufsteigender Reihenfolge nach der Modellnummer sortiert. Sie können die Zeilen in aufsteigender oder absteigender Reihenfolge nach den Werten jeder beliebigen Spalte sortieren.

**Modell.** Eine Nummer, die die Reihenfolge angibt, in der das Komponentenmodell erstellt wurde.

**Genauigkeit.** Als Prozentwert angegebene Gesamtgenauigkeit.

**Methode.** Die Modellierungsmethode.

**Prädiktoren.** Die Anzahl der im Komponentenmodell verwendeten Prädiktoren (Einflussvariablen).

**Modellgröße.** Die Modellgröße hängt von der Modellierungsmethode ab: Bei Bäumen handelt es sich um die Anzahl der Knoten im Baum, bei linearen Modellen um die Anzahl der Koeffizienten, bei neuronalen Netzen um die Anzahl der Synapsen.

**Datensätze.** Die gewichtete Anzahl an Eingabedatensätzen in der Trainingsstichprobe.

### Automatische Datenaufbereitung

Abbildung 3-26  
Ansicht "Automatische Datenaufbereitung"

#### Automatische Datenvorbereitung

Ziel: Total sales

Feld	Rolle	Durchgeführte Aktionen
Age category	Prädiktor	Zerstreute Kategorien für maximale Zuordnung mit Ziel zusammenführen
Primary keyword set	Prädiktor	Zerstreute Kategorien für maximale Zuordnung mit Ziel zusammenführen
Promotion	Prädiktor	Messniveau von kontinuierlich zu ordinal ändern
Secondary keyword set	Prädiktor	Zerstreute Kategorien für maximale Zuordnung mit Ziel zusammenführen

Wenn der Name des ursprünglichen Felds X ist, lautet der Name des transformierten Felds "X\_transformed". Das Originalfeld wird aus der Analyse ausgeschlossen und das transformierte Feld wird stattdessen eingeschlossen.

Diese Ansicht zeigt Informationen darüber an, welche Felder ausgeschlossen wurden und wie transformierte Felder im Schritt "automatische Datenaufbereitung" (ADP) abgeleitet wurden. Für jedes transformierte oder ausgeschlossene Feld listet die Tabelle den Feldnamen, die Rolle in der Analyse und die im ADP-Schritt vorgenommene Aktion auf. Die Felder werden in aufsteigender alphabetischer Reihenfolge der Feldnamen sortiert.

Wenn die Aktion Ausreißer trimmen angezeigt wird, bedeutet dies, dass Werte stetiger Prädiktoren, die über einem Cutoff-Wert liegen (drei Standardabweichungen vom Mittelwert), auf den Cutoff-Wert gesetzt wurden.

## Modell-Nuggets für aufgeteilte Modelle

Das Nugget für ein aufgeteiltes Modell bietet Zugriff auf alle einzelnen Modelle, die durch die Aufteilungen entstanden sind.

Ein Nugget für aufgeteilte Modelle enthält:

- eine Liste der erstellten aufgeteilten Modelle zusammen mit einem Set von Statistikdaten zu jedem Modell
- Informationen zum Gesamtmodell

In der Liste der aufgeteilten Modelle können Sie einzelne Modelle öffnen, um sie weiter zu untersuchen.

## Viewer für aufgeteilte Modelle

Auf der Registerkarte “Modell” sind alle Modelle aufgelistet, die in einem Nugget enthalten sind, und sie enthält zahlreiche Statistikdaten über die aufgeteilten Modelle. Je nach Modellierungsknoten hat sie zwei allgemeine Formen.

Abbildung 3-27  
Viewer für aufgeteilte Modelle

Diagramm	Modell	Sex	Anz. Datensätze in Aufteilung	Anzahl der verwendeten Felder	Gesamt-Genauigkeit (%)
		F	96	4	100
		M	104	4	99,038

**Sortieren nach.** Verwenden Sie diese Liste, um die Reihenfolge der aufgeführten Modelle zu wählen. Sie können die Liste auf der Basis der Werte in einer der angezeigten Spalten aufsteigend oder absteigend sortieren. Alternativ können Sie auf die Überschrift einer Spalte klicken, um eine Sortierung nach dieser Spalte vorzunehmen. Standard ist absteigende Sortierung mit Gesamtgenauigkeit.

**Menü "Spalten anzeigen/ausblenden".** Klicken Sie auf diese Schaltfläche, um ein Menü zu öffnen, in dem Sie einzelne Spalten zum Anzeigen oder Ausblenden wählen können.

**Ansicht.** Wenn Sie Partitionierung verwenden, können Sie die Ergebnisse für die Trainings- bzw. Testdaten anzeigen lassen. [Für weitere Informationen siehe Thema Partitionsknoten in Kapitel 4 in IBM SPSS Modeler 14.2- Quellen-, Prozess- und Ausgabeknoten.](#)

Für jede Aufteilung werden die folgenden Daten angezeigt:

**Diagramm.** Ein Miniaturdiagramm, das die Datenverteilung für dieses Modell anzeigt. Wenn sich das Nugget im Zeichenbereich befindet, wird das Diagramm durch Doppelklicken auf das Miniaturdiagramm in voller Größe geöffnet.

**Modell.** Ein Symbol des Modelltyps. Doppelklicken Sie auf das Symbol, um das Modell-Nugget für diese bestimmte Aufteilung zu öffnen.

**Aufteilungsfelder.** Die Felder, die im Modellierungsknoten als Aufteilungsfelder festgelegt sind, sowie Ihre möglichen Werte.

**Anzahl der Datensätze in der Aufteilung.** Die Anzahl von Datensätzen, die an dieser bestimmten Aufteilung beteiligt sind.

**Anzahl verwendeter Felder.** Teilt aufgeteilte Modelle auf der Grundlage der verwendeten Eingabefelder in Ränge ein.

**Gesamtgenauigkeit (%).** Der Prozentsatz der Datensätze, der korrekt vom aufgeteilten Modell vorhergesagt wird, im Verhältnis zur Gesamtzahl der Datensätze in dieser Aufteilung.

Abbildung 3-28  
Modellanzeige "Aufteilen"

**Aufteilungsgruppen**

ed	Genauigkeit	Modellgröße (Synapsen)	Datensätze
Did not complete high school	68,3%	62	372
High school degree	66,8%	42	198
Some college	72,3%	42	87
College degree	59,2%	12	38
Post-undergraduate degree	.	.	.

Modelle konnten nicht für eine oder mehrere Aufteilungsgruppen erstellt werden.

**Aufteilen.** Die Spaltenüberschrift zeigt das Feld bzw. die Felder, die zum Erstellen der Aufteilungen verwendet wurden, und die Zellen enthalten die Aufteilungswerte. Doppelklicken Sie auf eine Aufteilung, um eine Modellanzeige für das für die betreffende Aufteilung erstellte Modell anzuzeigen.

**Genauigkeit.** Als Prozentwert angegebene Gesamtgenauigkeit.

**Modellgröße.** Die Modellgröße hängt von der Modellierungsmethode ab: Bei Bäumen handelt es sich um die Anzahl der Knoten im Baum, bei linearen Modellen um die Anzahl der Koeffizienten, bei neuronalen Netzen um die Anzahl der Synapsen.

**Datensätze.** Die gewichtete Anzahl an Eingabedatensätzen in der Trainingsstichprobe.

## Verwendung von Modell-Nuggets in Streams

Modell-Nuggets werden in Streams platziert, damit Sie neue Daten scoren und neue Knoten generieren können. Das **Scoring** der Daten ermöglicht die Verwendung der aus der Modellerstellung gewonnenen Informationen zur Erstellung von Vorhersagen für neue Datensätze. Zur Anzeige der Scoring-Ergebnisse müssen Sie dem Nugget einen Endknoten hinzufügen (d. h. einen Verarbeitungs- oder Ausgabeknoten) und den Endknoten ausführen.

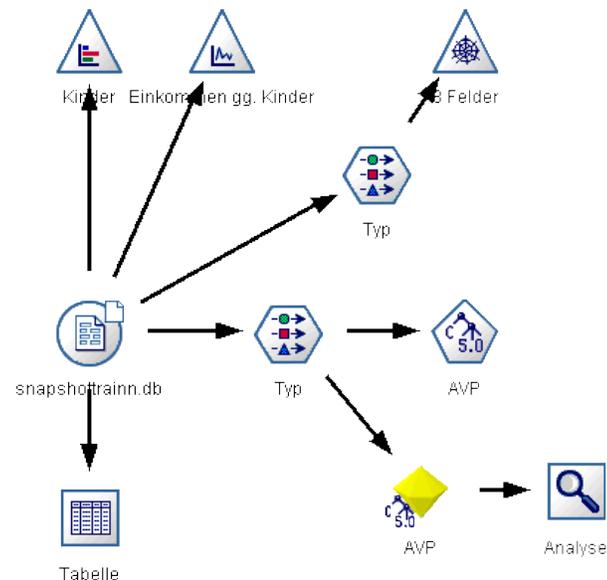
Bei einigen Modellen bieten Modell-Nuggets auch weitere Informationen über die Qualität der Vorhersage, beispielsweise Konfidenzwerte oder Entfernungen von den Cluster-Zentren. Durch das Generieren neuer Knoten können Sie ganz einfach neue Knoten auf der Grundlage der Struktur des generierten Modells erstellen. So ermöglichen beispielsweise die meisten Modelle, die eine Eingabefeldauswahl durchführen, die Erstellung von Filterknoten, die nur diejenigen Eingabefelder übergeben, die das Modell als wichtig identifiziert hat.

**So verwenden Sie ein Modell-Nugget zum Scoring von Daten:**

- Verbinden Sie das Modell-Nugget mit einer Datenquelle oder einem Stream, der ihm Daten übergeben soll.

Abbildung 3-29

Verwenden eines Modell-Nuggets für das Scoring



- Fügen Sie einen oder mehrere Verarbeitungs- oder Ausgabeknoten (beispielsweise einen Tabellen- oder Analyseknoden) zum Modell-Nugget hinzu oder verbinden Sie sie damit.
- Führen Sie einen der Knoten unterhalb des Modell-Nuggets aus.

*Hinweis:* Sie können nicht den Knoten für nicht verfeinerte Regeln für das Scoring von Daten verwenden. Um Daten auf der Grundlage eines -Assoziationsmodells zu scoren, erstellen Sie mit dem Knoten für nicht verfeinerte Regeln ein Regelmengen-Nugget und führen damit das Scoring durch. [Für weitere Informationen siehe Thema Generieren einer Regelmenge aus einem Assoziationsmodell-Nugget in Kapitel 12 auf S. 410.](#)

**So verwenden Sie ein Modell-Nugget zum Generieren von Verarbeitungsknoten:**

- Durchsuchen Sie auf der Palette das Modell oder bearbeiten Sie es im Stream-Zeichenbereich.
- Wählen Sie den gewünschten Knotentyp im Generierungsmenü des Modell-Nugget-Browsers aus. Die verfügbaren Optionen hängen vom Typ des Modell-Nuggets ab. Einzelheiten zu den Elementen, die aus einem bestimmten Modell generiert werden können, finden Sie in den Informationen zum jeweils generierten Modell-Nugget-Typ.

## **Erneutes Erzeugen eines Modellierungsknotens**

Wenn Sie ein Modell-Nugget ändern oder aktualisieren möchten und der zum Erstellen des Modells verwendete Stream nicht verfügbar ist, können Sie erneut einen Modellierungsknoten mit denselben Optionen erzeugen, die zum Erstellen des ursprünglichen Modells verwendet wurden.

- ▶ Um ein Modell erneut zu erstellen, klicken Sie in der Modellalette mit der rechten Maustaste auf das gewünschte Modell und wählen Sie die Option Modellierungsknoten erzeugen.
- ▶ Alternativ können Sie beim Durchsuchen eines Modells im Menü “Generieren” die Option Modellierungsknoten erzeugen auswählen.

Der erneut erzeugte Modellierungsknoten sollte in den meisten Fällen von der Funktionsweise her identisch mit dem Knoten sein, mit dem das ursprüngliche Modell erstellt wurde.

- Bei Entscheidungsbaummodellen können zusammen mit dem Knoten auch weitere Einstellungen gespeichert werden, die während der interaktiven Sitzung angegeben wurden, und die Option Interaktiv erstellte Aufbauregeln verwenden ist in dem neu erzeugten Modellierungsknoten aktiviert.
- Bei Entscheidungslistenmodellen ist die Option Gespeicherte Informationen aus interaktiver Sitzung verwenden aktiviert. [Für weitere Informationen siehe Thema Entscheidungslistenmodell – Optionen in Kapitel 9 auf S. 232.](#)
- Bei Zeitreihenmodellen ist die Option Schätzung unter Verwendung bestehender Modelle fortsetzen aktiviert. Diese Einstellung gestattet die erneute Erzeugung des vorherigen Modells mit aktuellen Daten. [Für weitere Informationen siehe Thema Zeitreihenmodelle – Optionen in Kapitel 13 auf S. 442.](#)

## **Importieren und Exportieren von Modellen als PMML**

PMML (Predictive Model Markup Language) ist ein XML-Format zur Beschreibung von Data-Mining-Modellen und statistischen Modellen, einschließlich der Eingaben zu den Modellen, der zur Vorbereitung der Daten für das Data-Mining verwendeten Transformationen sowie der Parameter, die die Modelle selbst definieren. IBM® SPSS® Modeler kann PMML importieren und exportieren, wodurch es möglich ist, Modelle mit anderen Anwendungen auszutauschen, die dieses Format unterstützen, zum Beispiel IBM® SPSS® Statistics.

*Hinweis:* Für den Export von PMML ist eine IBM® SPSS® Modeler Server-Lizenz erforderlich.

Weitere Informationen zu PMML finden Sie auf der Website der Data Mining Group (<http://www.dmg.org>).

### **So exportieren Sie ein Modell:**

PMML-Export wird für die meisten der in SPSS Modeler erstellten Modelltypen unterstützt. [Für weitere Informationen siehe Thema Modelltypen, die PMML unterstützen in IBM SPSS Modeler 14.2- Benutzerhandbuch.](#)

- ▶ Klicken Sie mit der rechten Maustaste auf ein Modell-Nugget in der Modellalette. (Alternativ können Sie auf ein Modell-Nugget im Zeichenbereich klicken und das Menü “Datei” auswählen.)

- Klicken Sie im Menü auf PMML exportieren.

Abbildung 3-30

Exportieren von Modellen im PMML-Format



- Geben Sie im Dialogfeld “Exportieren” (oder “Speichern”) ein Zielverzeichnis und einen eindeutigen Namen für das Modell an.

*Hinweis:* Im Dialogfeld “Benutzeroptionen” können Sie die Optionen für den PMML-Export ändern. Klicken Sie im Hauptmenü auf:

Werkzeuge > Optionen > Benutzeroptionen

und wählen Sie die Registerkarte “PMML”.

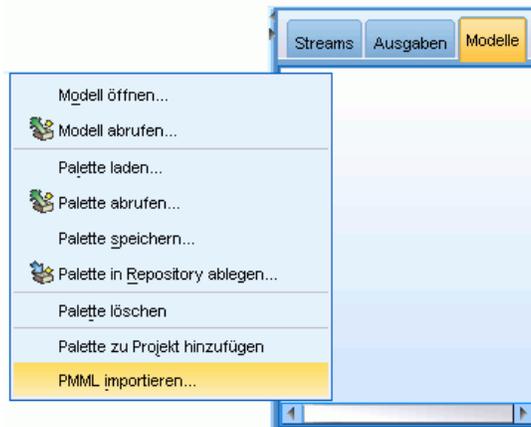
Für weitere Informationen siehe Thema Festlegen von PMML-Exportoptionen in Kapitel 12 in *IBM SPSS Modeler 14.2- Benutzerhandbuch*.

**So importieren Sie ein als PMML gespeichertes Modell:**

Modelle, die aus SPSS Modeler oder einer anderen Anwendung als PMML exportiert wurden, können in die Modellpalette importiert werden. Für weitere Informationen siehe Thema [Modelltypen, die PMML unterstützen](#) in *IBM SPSS Modeler 14.2- Benutzerhandbuch*.

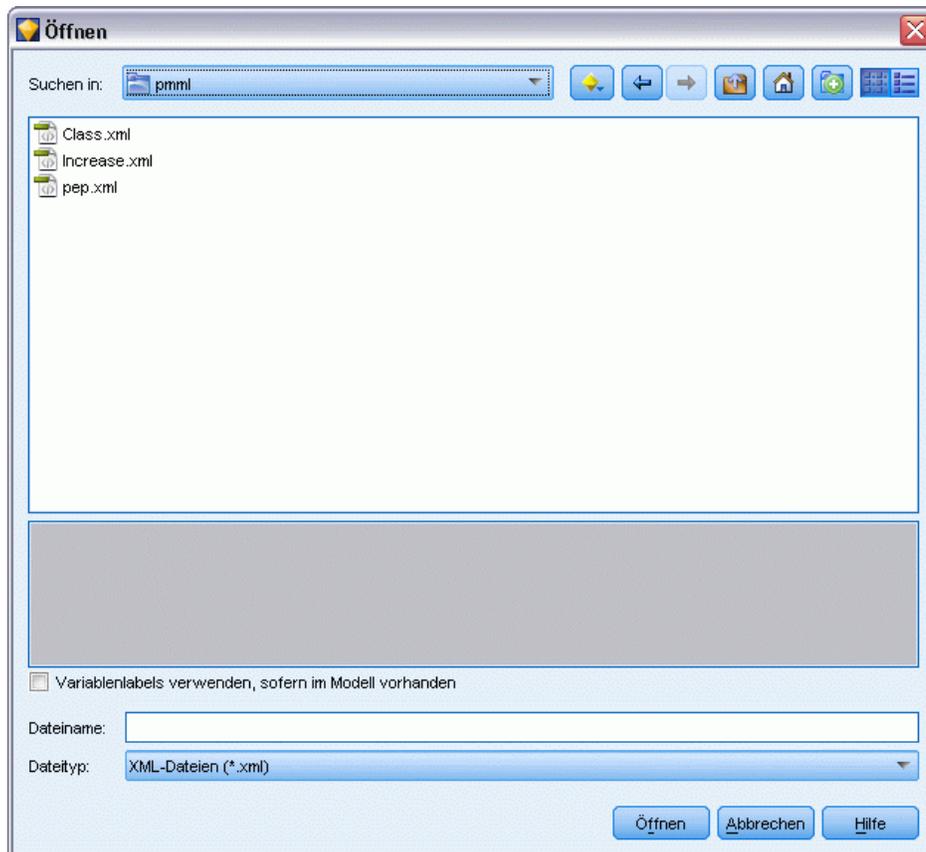
- Klicken Sie in der Modellpalette mit der rechten Maustaste auf die Palette und wählen Sie aus dem Menü die Option PMML importieren.

Abbildung 3-31  
Importieren von Modellen im PMML-Format



- ▶ Wählen Sie die zu importierende Datei aus und geben Sie nach Bedarf Optionen für Variablenlabels an.
- ▶ Klicken Sie auf Öffnen.

Abbildung 3-32  
Auswählen der XML-Datei für ein mit PMML gespeichertes Modell



**Verwenden Sie Variablenlabels, falls im Modell vorhanden.** Im PMML-Code können sowohl die Variablennamen als auch die Variablenlabels (beispielsweise “Referrer ID” für *RefID*) für Variablen im Datenlexikon gefunden. Wählen Sie diese Option aus, um Variablenlabels zu verwenden, wenn diese im ursprünglich exportierten PMML-Code vorhanden sind.

Wenn Sie die Optionen zur Variablenbeschriftung ausgewählt haben, im PMML-Code jedoch keine Variablenlabels vorhanden sind, werden die Variablennamen wie üblich verwendet.

## **Modelltypen, die PMML unterstützen**

### **PMML-Export**

**SPSS Modeler Modelle.** Die folgenden in IBM® SPSS® Modeler erstellten Modelle können als PMML 4,0 exportiert werden:

- C&R-Baum
- QUEST
- CHAID
- Lineare Regression
- Netzwerk
- C5,0
- Logistische Regression
- Genlin
- SVM
- Bayes-Netz
- A Priori
- Carma
- Sequenz
- Cluster-Zentren
- Kohonen
- Two Step
- KNN
- Statistics Modell

Das folgende in SPSS Modeler erstellte Modell kann als PMML 3.2 exportiert werden:

- Entscheidungsliste

**Datenbankeigene Modelle.** Bei Modellen, die mithilfe von datenbankeigenen Algorithmen generiert wurden, steht der PMML-Export nur bei IBM InfoSphere Warehouse-Modellen zur Verfügung. Modelle, die mithilfe von Analysis Services von Microsoft oder mit Oracle Data Miner erstellt wurden, können nicht exportiert werden. Beachten Sie außerdem, dass als PMML exportierte IBM-Modelle danach nicht wieder in SPSS Modeler importiert werden können. [Für](#)

weitere Informationen siehe Thema [Übersicht über die Datenbank-Modellierung in Kapitel 2 in IBM SPSS Modeler 14.2 – In-Database Mining-Handbuch](#).

### **PMML-Import**

SPSS Modeler kann PMML-Modelle importieren und scoren, die von aktuellen Versionen aller IBM® SPSS® Statistics-Produkte erstellt wurden, darunter Modelle, die aus SPSS Modeler exportiert wurden, sowie Modell- bzw. Transformations-PMML, die von SPSS Statistics 17.0 oder höher generiert wurde. Dies gilt also im Grunde für jegliche PMML, die die Scoring-Engine scoren kann – mit folgenden Ausnahmen:

- A Priori-, CARMA- und Anomaly Detection-(Anomalieerkennung-)Modelle können nicht importiert werden.
- PMML-Modelle können nach dem Import in SPSS Modeler nicht durchsucht werden, obwohl sie für das Scoring verwendet werden können. (Dies gilt auch für Modelle, die ursprünglich aus SPSS Modeler exportiert wurden. Um diese Einschränkung zu vermeiden, sollten Sie das betreffende Modell als generierte Modelldatei (\*.gm) und nicht als PMML exportieren.
- Als PMML exportierte IBM InfoSphere Warehouse-Modelle können nicht importiert werden.
- Eine eingeschränkte Validierung findet beim Import statt, aber die vollständige Validierung erfolgt beim Versuch, das Modell zu scoren. Daher kann der Import erfolgreich durchgeführt werden, das Scoring aber fehlschlagen oder falsche Ergebnisse erzeugen.

### **Nicht verfeinerte Modelle**

Ein nicht verfeinertes -Modell enthält Informationen, die aus den Daten extrahiert wurden, jedoch nicht zur direkten Generierung von Vorhersagen gedacht sind. Daher kann es nicht zu Streams hinzugefügt werden. Nicht verfeinerte Modelle werden auf der Palette der generierten Modelle als “Rohdiamanten” angezeigt.

Abbildung 3-33  
Symbol für nicht verfeinerte Modelle



Informationen zum nicht verfeinerten Regelmodell erhalten Sie, wenn Sie mit der rechten Maustaste auf das Modell klicken und im Kontextmenü die Option Durchsuchen auswählen. Wie bei anderen in IBM® SPSS® Modeler generierten Modellen bieten die verschiedenen Registerkarten Übersichts- und Regelinformationen zum erstellten Modell.

**Generieren von Knoten.** Im Menü “Generieren” können Sie anhand der Regeln neue Knoten erstellen.

- **Auswahlknoten.** Generiert einen Auswahlknoten zur Auswahl von Datensätzen, für die die ausgewählte Regel gilt. Diese Option ist deaktiviert, wenn keine Regel ausgewählt wurde.
- **Regelmenge.** Generiert einen Regelmengenknoten zur Vorhersage der Werte für ein einzelnes Zielfeld. [Für weitere Informationen siehe Thema Generieren einer Regelmenge aus einem Assoziationsmodell-Nugget in Kapitel 12 auf S. 410.](#)

# Screening von Modellen

## Screening von Feldern und Datensätzen

In den vorgelagerten Phasen einer Analyse können mehrere Modellierungsknoten verwendet werden, um Felder und Datensätze zu finden, die voraussichtlich bei der Modellierung relevant sind. Sie können den Merkmalsauswahlknoten verwenden, um Felder per Screening zu untersuchen und nach Wichtigkeit zu ordnen, und den Anomalieerkennungsknoten, um ungewöhnliche Datensätze zu finden, die nicht den bekannten Mustern “normaler” Daten entsprechen.



Der Merkmalsauswahlknoten sichtet die Eingabefelder, um auf der Grundlage einer Reihe von Kriterien (z. B. dem Prozentsatz der fehlenden Werte) zu entscheiden, ob diese entfernt werden sollen. Anschließend erstellt er eine Wichtigkeitsrangfolge der verbleibenden Eingaben in Bezug auf ein angegebenes Ziel. Beispiel: Angenommen, Sie haben ein Daten-Set mit Hunderten potenzieller Eingaben. Welche davon sind voraussichtlich für die Modellierung von medizinischen Behandlungsergebnissen von Bedeutung? [Für weitere Informationen siehe Thema Merkmalsauswahlknoten auf S. 75.](#)



Der Anomalieerkennungsknoten ermittelt ungewöhnliche Fälle bzw. “Ausreißer”, die nicht den Mustern der “normalen” Daten entsprechen. Mit diesem Knoten können Ausreißer ermittelt werden, selbst wenn sie keinem bereits bekannten Muster entsprechen und selbst wenn Sie nicht genau wissen, wonach Sie suchen. [Für weitere Informationen siehe Thema Anomalieerkennungsknoten auf S. 83.](#)

Beachten Sie, dass bei der Anomalieerkennung ungewöhnliche Datensätze oder Fälle mithilfe einer Cluster-Analyse ermittelt werden, die auf der im Modell ausgewählten Menge an Feldern beruht – ohne Berücksichtigung eines speziellen Zielfelds (abhängigen Felds) und unabhängig davon, ob diese Felder für das Muster relevant sind, das Sie vorherzusagen versuchen. Aus diesem Grund sollten Sie die Anomalieerkennung in Kombination mit der Merkmalsauswahl oder einem anderen Verfahren für Screening und Rangordnung von Feldern verwenden. Beispielsweise können Sie mithilfe der Merkmalsauswahl die wichtigsten Felder in Bezug auf ein bestimmtes Ziel ermitteln und anschließend mit der Anomalieerkennung die Datensätze finden, die in Bezug auf diese Felder besonders ungewöhnlich sind. (Eine alternative Vorgehensweise besteht darin, ein Entscheidungsbaummodell zu erstellen und anschließend alle falsch klassifizierte Datensätze als potenzielle Anomalien zu untersuchen. Diese Methode lässt sich jedoch nicht so leicht reproduzieren bzw. in größerem Maßstab automatisieren.)

## Merkmalsauswahlknoten

Ein Problem beim Data-Mining kann darin bestehen, dass Hunderte oder sogar Tausende Felder potenziell als Eingaben in Frage kommen. Als Folge davon muss aufwendig untersucht werden, welche Felder bzw. Variablen in das Modell aufgenommen werden sollen. Um die Auswahlmöglichkeiten einzugrenzen, können mithilfe des Merkmalsauswahlalgorithmus die Felder ermittelt werden, die für eine bestimmte Analyse am wichtigsten sind. Wenn Sie

beispielsweise versuchen, die Ergebnisse medizinischer Behandlungen anhand einer Reihe von Faktoren vorherzusagen, welche Faktoren sind dann vermutlich am wichtigsten?

Die Merkmalsauswahl besteht aus drei Schritten:

- **Screening.** Eliminiert unwichtige und problematische Eingaben und Datensätze bzw. Fälle, beispielsweise Eingabefelder mit zu vielen fehlenden Werten oder Eingaben, die eine so starke oder geringe Variation aufweisen, dass sie nicht brauchbar sind.
- **Einstufen.** Sortiert die verbleibenden Eingaben und weist ihnen Ränge je nach Wichtigkeit zu.
- **Auswahl.** Ermittelt die Untermenge an Merkmalen, die in den nachfolgenden Modellen verwendet werden sollen, beispielsweise indem nur die wichtigsten Eingaben beibehalten und alle anderen gefiltert bzw. ausgeschlossen werden.

In einer Zeit, in der viele Unternehmen mit einer gewaltigen Datenflut umgehen müssen, können die Vorteile, die die Merkmalsauswahl für die Vereinfachung und Beschleunigung des Modellierungsprozesses bietet, erheblich sein. Indem die Aufmerksamkeit schnell auf die wichtigsten Felder gelenkt wird, lässt sich der Berechnungsaufwand verringern, schwache, aber wichtige Beziehungen, die ansonsten leicht übersehen werden, können einfacher aufgespürt werden und schließlich erhalten Sie einfachere, genauere und leichter erklärable Modelle. Wenn Sie die Anzahl der im Modell verwendeten Felder verringern, stellen Sie möglicherweise fest, dass Sie die Scoring-Zeiten verkürzen sowie die bei zukünftigen Wiederholungen zu sammelnde Datenmenge reduzieren können.

Die Verringerung der Anzahl der Felder kann insbesondere für Modelle wie die logistische Regression nützlich sein, für die eine Obergrenze von 350 Feldern gilt.

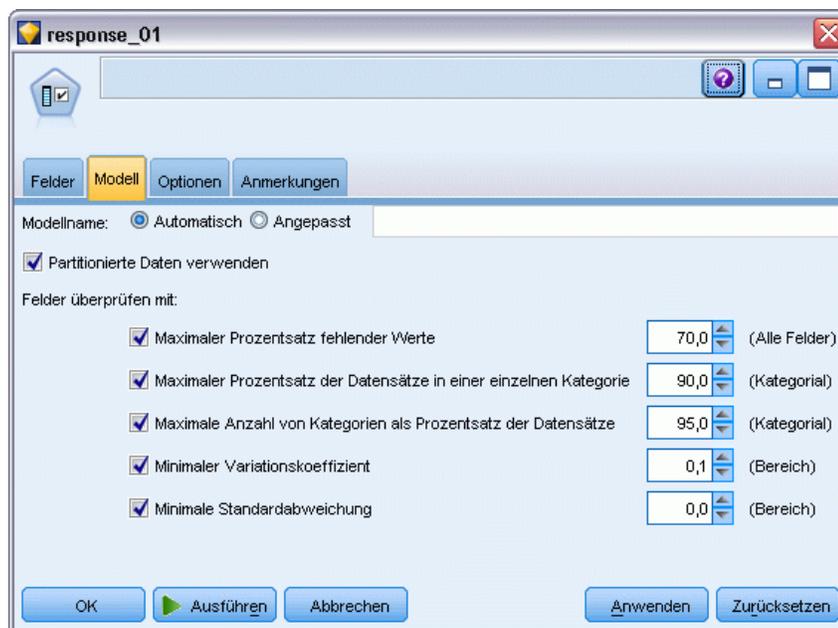
**Beispiel.** Eine Telefongesellschaft verfügt über ein Data Warehouse, das Informationen zu Reaktionen auf eine spezielle Werbeaktion enthält, die an 5.000 Kunden des Unternehmens gerichtet war. Die Daten enthalten eine Vielzahl von Feldern, darunter das Alter der Kunden, Ihr Beschäftigungsverhältnis, ihr Einkommen und statistische Daten zu ihrer Telefonnutzung. Drei Zielfelder zeigen jeweils an, ob der Kunde auf die drei Angebote reagierte oder nicht. Das Unternehmen möchte anhand dieser Daten vorhersagen, welche Kunden mit der größten Wahrscheinlichkeit auf ähnliche Angebote in Zukunft reagieren. [Für weitere Informationen siehe Thema Screening von Prädiktoren \(Merkmalsauswahl\) in Kapitel 10 in IBM SPSS Modeler 14.2- Anwendungshandbuch.](#)

**Anforderungen.** Ein einzelnes Zielfeld (mit der Rolle *Ziel*) sowie mehrere Eingabefelder, die in Bezug auf das Ziel untersucht bzw. nach Rang geordnet werden sollen. Sowohl Ziel- als auch Eingabefelder können das Messniveau *Stetig* (numerischer Bereich) oder *Kategorial* aufweisen.

## ***Einstellungen für das Merkmalsauswahlmodell***

Die Einstellungen auf der Registerkarte "Modell" beinhalten Standardmodelloptionen sowie Einstellungen, mit denen Sie eine Feineinstellung der für das Screening von Eingabefeldern verwendeten Kriterien vornehmen können.

Abbildung 4-1  
Registerkarte "Merkmalsauswahlmodell"



**Modellname.** Sie können den Modellnamen automatisch basierend auf den Ziel- oder ID-Feldnamen (oder dem Modelltyp in Fällen, in denen kein solches Feld angegeben ist) generieren oder einen benutzerdefinierten Namen eingeben.

### Screening von Eingabefeldern

Zum Screening gehört das Entfernen von Eingaben bzw. Fällen, die hinsichtlich der Beziehung zwischen Eingabe und Ziel keine nützlichen Informationen hinzufügen. Die Screeningoptionen beruhen auf Attributen des betreffenden Felds ohne Berücksichtigung der Vorhersagekraft in Bezug auf das ausgewählte Zielfeld. Die untersuchten Felder werden aus den Berechnungen für die Rangordnung der Eingaben ausgenommen und können optional gefiltert bzw. aus den bei der Modellierung verwendeten Daten entfernt werden.

Ein Screening der Felder kann auf folgenden Kriterien beruhen:

- **Maximaler Prozentsatz fehlender Werte.** Untersucht Felder mit zu vielen fehlenden Werten, ausgedrückt als Prozentsatz der Gesamtzahl an Datensätzen. Felder mit einem großen Prozentsatz an fehlenden Werten bieten wenig Informationen für die Vorhersage.
- **Maximaler Prozentsatz der Datensätze in einer einzelnen Kategorie.** Untersucht Felder, bei denen zu viele Datensätze (im Verhältnis zur Gesamtzahl der Datensätze) in dieselbe Kategorie fallen. Wenn beispielsweise 95 % der Kunden in der Datenbank denselben Autotyp fahren, ist die Aufnahme dieser Information für die Unterscheidung der einzelnen Kunden untereinander nicht hilfreich. Alle Felder, die das angegebene Maximum überschreiten, werden im Screening untersucht. Diese Option gilt nur für kategoriale Felder.
- **Maximale Anzahl von Kategorien als Prozentsatz der Datensätze.** Untersucht Felder mit zu vielen Kategorien im Verhältnis zur Gesamtzahl der Datensätze. Wenn ein hoher Prozentsatz der Kategorien nur einen einzelnen Fall enthält, ist das Feld voraussichtlich von begrenztem

Nutzen. Beispiel: Wenn jeder Kunde einen anderen Hut trägt, ist diese Information für die Modellierung von Verhaltensmustern mit großer Wahrscheinlichkeit unbrauchbar. Diese Option gilt nur für kategoriale Felder.

- **Minimaler Variationskoeffizient.** Untersucht Felder mit einem Varianzkoeffizienten, der kleiner oder gleich dem angegebenen Mindestwert ist. Dieses Maß ist der Quotient aus der Standardabweichung des Eingabefelds und dem Mittelwert des Eingabefelds. Wenn dieser Wert nahe bei null liegt, liegt nur eine geringe Variabilität in den Werten für die betreffende Variable vor. Diese Option gilt nur für stetige Felder (numerischer Bereich).
- **Minimale Standardabweichung.** Untersucht Felder mit einer Standardabweichung, die kleiner oder gleich dem angegebenen Mindestwert ist. Diese Option gilt nur für stetige Felder (numerischer Bereich).

**Datensätze mit fehlenden Daten.** Datensätze oder Fälle mit fehlenden Werten für das Zielfeld oder fehlenden Werten für alle Eingaben werden automatisch aus allen Berechnungen für die Rangfolge ausgeschlossen.

### ***Merkmalsauswahloption***

Auf der Registerkarte “Optionen” können Sie die Standardeinstellungen für die Auswahl bzw. den Ausschluss von Eingabefeldern im Modell-Nugget angeben. Anschließend können Sie das Modell zu einem Stream hinzufügen, um die Untermenge der Felder auszuwählen, die in nachfolgenden Modellerstellungsvorgängen verwendet werden sollen. Alternativ können Sie diese Einstellungen nach der Modellgeneration durch die Auswahl bzw. das Aufheben der Auswahl weiterer Felder im Modellbrowser überschreiben. Die Standardeinstellungen ermöglichen es jedoch, das Modell-Nugget ohne weitere Änderungen anzuwenden, was insbesondere für die Skripterstellung nützlich sein kann.

[Für weitere Informationen siehe Thema Ergebnisse des Merkmalsauswahlmodells auf S. 80.](#)

Abbildung 4-2  
Merkmalsauswahl – Registerkarte "Optionen"



Die folgenden Optionen sind verfügbar:

**Alle Felder mit Rangzahl.** Wählt die Felder auf der Grundlage ihres Ranges (*bedeutsam*, *marginal* oder *unbedeutend*) aus. Sie können die Beschriftung für jeden Rang bearbeiten sowie die Cutoff-Werte ändern, die verwendet werden um Datensätze einem bestimmten Rang zuzuweisen.

**Obere Anzahl an Feldern.** Wählt die obersten  $n$  Felder nach Wichtigkeit aus.

**Wichtigkeit größer als.** Wählt alle Felder aus, deren Wichtigkeit den angegebenen Wert übersteigt.

Das Zielfeld bleibt unabhängig von der Auswahl immer erhalten.

#### **Optionen für die Rangeinteilung nach Wichtigkeit**

**Ausschließlich kategoriale Werte.** Wenn alle Eingaben und das Ziel kategorial sind, stehen für die Rangordnung nach Wichtigkeit vier verschiedene Maße zur Auswahl:

- **Pearson'sches Chi-Quadrat.** Testet auf Unabhängigkeit von Ziel und Eingabe ohne Angabe der Stärke oder Verwendung (Richtung) einer bestehenden Beziehung.
- **Likelihood-Quotienten-Chi-Quadrat.** Ähnlich dem Pearson'schen Chi-Quadrat; testet jedoch außerdem auf Ziel-Eingabe-Unabhängigkeit.
- **Cramer-V.** Ein Assoziationsmaß auf der Grundlage der Chi-Quadrat-Statistik nach Pearson. Die Werte reichen von 0 (keine Assoziation) bis 1 (vollkommene Assoziation).
- **Lambda.** Ein Assoziationsmaß, das die proportionale Fehlerrückgang angibt, die sich ergibt, wenn die Variable zur Vorhersage des Zielwerts verwendet wird. Der Wert 1 gibt an, dass das Eingabefeld das Ziel perfekt vorhersagt, wohingegen der Wert 0 bedeutet, dass die Eingabe keine nützlichen Informationen über das Ziel bietet.

**Teilweise kategoriale Werte.** Wenn einige – jedoch nicht alle – Eingaben kategorial sind und das Ziel ebenfalls kategorial ist, lässt sich die Rangordnung nach Wichtigkeit entweder auf der Grundlage des Pearson'schen Chi-Quadrats oder des Likelihood-Quotienten-Chi-Quadrats ermitteln. (Cramer- $V$  und Lambda sind nur verfügbar, wenn alle Eingaben kategorial sind.)

**“Kategorial” und “stetig” im Vergleich.** Bei der Rangbewertung einer kategorialen Eingabe anhand eines stetigen Ziels oder umgekehrt (eines der beiden Elemente ist kategorial, nicht jedoch beide) wird die  $F$ -Statistik verwendet.

**Beides stetig.** Bei der Rangbewertung einer stetigen Eingabe anhand eines stetigen Ziels wird die  $t$ -Statistik auf der Grundlage des Korrelationskoeffizienten verwendet.

## **Modell-Nuggets vom Typ “Merkmalsauswahl”**

Modell-Nuggets vom Typ “Merkmalsauswahl” zeigen die Bedeutsamkeit der einzelnen Eingaben in Bezug auf ein ausgewähltes Ziel an (gemäß der Rangeinstufung durch den Merkmalsauswahlknoten). Alle Felder, die vor der Rangeinstufung per Screening ausgeschlossen wurden, werden ebenfalls aufgeführt. [Für weitere Informationen siehe Thema Merkmalsauswahlknoten auf S. 75.](#)

Bei der Ausführung eines Streams, der ein Modell-Nugget vom Typ “Merkmalsauswahl” enthält, fungiert das Modell als Filter, mit dem nur ausgewählte Eingaben (in der aktuellen Auswahl auf der Registerkarte “Modell” angezeigt) beibehalten werden. Sie können beispielsweise alle als bedeutsam eingestuft Felder auswählen (eine der Standardoptionen) oder manuell eine Teilmenge der Felder auf der Registerkarte “Modell” auswählen. Das Zielfeld bleibt unabhängig von der Auswahl ebenfalls erhalten. Alle anderen Felder werden ausgeschlossen.

Die Filterung beruht ausschließlich auf dem Feldnamen; wenn Sie beispielsweise *Alter* und *Einkommen* auswählen, wird jedes Feld, das mit einem dieser Namen übereinstimmt, beibehalten. Das Modell aktualisiert die Rangwerte für die Felder nicht anhand neuer Daten; es filtert einfach nur die Felder anhand der ausgewählten Namen. Aus diesem Grund sollten Sie bei der Anwendung des Modells auf neue oder aktualisierte Daten Vorsicht walten lassen. Im Zweifelsfall wird eine erneute Generierung des Modells empfohlen.

## **Ergebnisse des Merkmalsauswahlmodells**

Auf der Registerkarte “Modell” für ein Modell-Nugget vom Typ “Merkmalsauswahl” werden die Rangwertung und die Bedeutsamkeit für alle Eingaben im oberen Fensterbereich angezeigt. Außerdem haben Sie die Möglichkeit, mithilfe der Kontrollkästchen in der Spalte auf der linken Seite Felder für die Filterung auszuwählen. Bei der Ausführung des Streams werden nur die aktivierten Felder beibehalten. Die anderen Felder werden verworfen. Die Standardauswahl beruht auf den im Modellerstellungsknoten angegebenen Optionen, Sie können jedoch nach Bedarf weitere Felder auswählen bzw. deren Auswahl aufheben.

Im unteren Fensterbereich werden Eingaben aufgelistet, die gemäß des Prozentsatzes an fehlenden Werten oder aufgrund anderer im Modellierungsknoten angegebener Kriterien aus der Rangwertung ausgenommen wurden. Wie bei den in die Rangwertung aufgenommenen Feldern, können Sie mithilfe der Kontrollkästchen in der Spalte auf der linken Seite auswählen, ob diese Felder eingeschlossen oder verworfen werden sollen. [Für weitere Informationen siehe Thema Einstellungen für das Merkmalsauswahlmodell auf S. 76.](#)

Abbildung 4-3  
Ergebnisse des Merkmalsauswahlmodells

The screenshot shows the 'response\_01' application window. The main table displays the following data:

	Rang	Feld	Messung	Wichtigkeit	Wert
<input checked="" type="checkbox"/>	1	ed	Stetig	★ Important	1,0
<input checked="" type="checkbox"/>	2	ownpc	Nominal	★ Important	1,0
<input checked="" type="checkbox"/>	3	edcat	Ordinal	★ Important	1,0
<input checked="" type="checkbox"/>	4	internet	Nominal	★ Important	1,0
<input checked="" type="checkbox"/>	5	equip	Nominal	★ Important	1,0
<input checked="" type="checkbox"/>	6	owngame	Nominal	★ Important	1,0
<input checked="" type="checkbox"/>	7	equipmon	Stetig	★ Important	1,0
<input checked="" type="checkbox"/>	8	confer	Nominal	★ Important	1,0
<input checked="" type="checkbox"/>	9	ebill	Nominal	★ Important	1,0
<input checked="" type="checkbox"/>	10	callwait	Nominal	★ Important	1,0
<input type="checkbox"/>	11	forward	Nominal	★ Important	1,0
<input type="checkbox"/>	12	tollmon	Stetig	★ Important	1,0
<input type="checkbox"/>	13	multiline	Nominal	★ Important	1,0
<input type="checkbox"/>	14	ownipod	Nominal	★ Important	1,0
<input type="checkbox"/>	15	callid	Nominal	★ Important	1,0
<input type="checkbox"/>	16	equipten	Stetig	★ Important	1,0
<input type="checkbox"/>	17	tollfree	Nominal	★ Important	1,0
<input type="checkbox"/>	18	tollten	Stetig	★ Important	1,0
<input type="checkbox"/>	19	churn	Nominal	★ Important	1,0
<input type="checkbox"/>	20	spousedcat	Ordinal	★ Important	1,0

Ausgewählte Felder: 19    Gesamtzahl der verfügbaren Felder: 128

Filter:  > 0,95    ≤ 0,95    < 0,9

9 Ausgefilterte Felder:

	Feld	Messung	Grund
<input checked="" type="checkbox"/>	ownvcr	Nominal	Einzelkategorie zu groß
<input checked="" type="checkbox"/>	owntv	Nominal	Einzelkategorie zu groß
<input checked="" type="checkbox"/>	owndvd	Nominal	Einzelkategorie zu groß
<input checked="" type="checkbox"/>	owncd	Nominal	Einzelkategorie zu groß
<input checked="" type="checkbox"/>	Inwireten	Stetig	Zu viele fehlende Werte
<input checked="" type="checkbox"/>	Inwirem...	Stetig	Zu viele fehlende Werte
<input checked="" type="checkbox"/>	Inequip...	Stetig	Variationskoeffizient unter Schwellenwert
<input checked="" type="checkbox"/>	commut...	Nominal	Einzelkategorie zu groß

Buttons: OK, Abbrechen, Anwenden, Zurücksetzen

- Um die Liste nach Rang, Feldname, Wichtigkeit oder einer anderen der angezeigten Spalten zu sortieren, klicken Sie auf die Spaltenüberschrift. Wenn Sie lieber die Symbolleiste verwenden, wählen Sie das gewünschte Element in der Liste "Sortieren nach" aus. Mit den nach unten bzw. oben zeigenden Pfeilen können Sie die Sortierrichtung ändern.
- Sie können mithilfe der Symbolleiste alle Felder aktivieren bzw. deaktivieren und auf das Dialogfeld "Felder markieren" zugreifen, in dem Sie Felder nach Rangordnung oder Wichtigkeit auswählen können. Außerdem können Sie beim Klicken auf die Felder die Umschalt- oder Strg-Taste gedrückt halten, um die Auswahl zu erweitern, und mithilfe

der Leertaste eine Gruppe ausgewählter Felder aktivieren bzw. deaktivieren. [Für weitere Informationen siehe Thema Auswählen der Felder nach Wichtigkeit auf S. 82.](#)

- Die Schwellenwerte für die Einordnung von Eingaben als “bedeutsam”, “marginal” bzw. “unbedeutend” werden in der Legende unterhalb der Tabelle angezeigt. Diese Werte werden im Modellierungsknoten angegeben. [Für weitere Informationen siehe Thema Merkmalsauswahloption auf S. 78.](#)

### **Auswählen der Felder nach Wichtigkeit**

Beim Scoring von Daten mithilfe eines Modell-Nuggets vom Typ “Merkmalsauswahl” bleiben alle Felder, die (mithilfe der Kontrollkästchen in der Spalte auf der linken Seite) aus der Liste der in Ränge eingeteilten bzw. per Screening untersuchten Felder ausgewählt wurden, erhalten. Die anderen Felder werden verworfen. Um die Auswahl zu ändern, können Sie mithilfe der Symbolleiste das Dialogfeld “Felder markieren” aufrufen, in dem Sie Felder nach Rang oder Wichtigkeit auswählen können.

Abbildung 4-4  
Dialogfeld “Felder markieren”



**Alle Felder, die markiert sind als.** Wählt alle als bedeutsam, marginal oder unbedeutend markierten Felder aus.

**Obere Anzahl an Feldern.** Ermöglicht die Auswahl der obersten  $n$  Felder nach Wichtigkeit.

**Wichtigkeit größer als.** Wählt alle Felder aus, deren Wichtigkeit den angegebenen Schwellenwert übersteigt.

### **Generieren eines Filters aus einem Merkmalsauswahlmodell**

Auf der Grundlage der Ergebnisse eines Merkmalsauswahlmodells können Sie einen oder mehrere Filterknoten generieren, die Untergruppen von Feldern auf der Grundlage ihrer Wichtigkeit in Bezug auf das angegebene Ziel ein- bzw. ausschließen. Das Modell-Nugget kann zwar auch als Filter verwendet werden, hiermit jedoch erhalten Sie die Flexibilität, mit verschiedenen Untergruppen von Feldern zu experimentieren, ohne das Modell kopieren oder bearbeiten zu müssen. Das Zielfeld wird stets vom Filter beibehalten, unabhängig davon ob “Einschließen” oder “Ausschließen” ausgewählt wurde.

Abbildung 4-5  
Generieren eines Filterknotens



**Einschließen/Ausschließen.** Sie können auswählen, welche Felder ein- bzw. ausgeschlossen werden sollen – Sie können beispielsweise die obersten 10 Felder einschließen oder alle Felder, die als unbedeutend markiert sind, ausschließen.

**Ausgewählte Felder.** Schließt alle aktuell in der Tabelle ausgewählten Felder ein bzw. aus.

**Alle Felder, die markiert sind als.** Wählt alle als bedeutsam, marginal oder unbedeutend markierten Felder aus.

**Obere Anzahl an Feldern.** Ermöglicht die Auswahl der obersten  $n$  Felder nach Wichtigkeit.

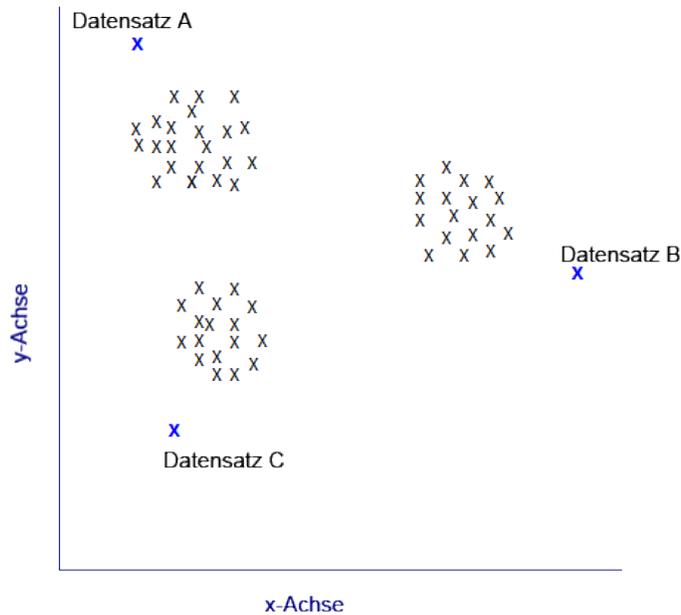
**Wichtigkeit größer als.** Wählt alle Felder aus, deren Wichtigkeit den angegebenen Schwellenwert übersteigt.

## Anomalieerkennungsknoten

Anomalieerkennungsmodelle ermitteln ungewöhnliche Fälle bzw. Ausreißer in den Daten. Im Gegensatz zu anderen Modellierungsmethoden, bei denen Regeln zu ungewöhnlichen Fällen gespeichert sind, speichern Anomalieerkennungsmodelle ausführliche Informationen darüber, wie das "normale" Verhalten aussieht. Auf diese Weise können Ausreißer selbst dann erkannt werden, wenn sie keinem bekannten Muster entsprechen. Dies ist insbesondere in Einsatzgebieten wie der Betrugserkennung von Nutzen, bei denen ständig neue Muster auftreten. Die Anomalieerkennung ist eine nicht überwachte Methode; dies bedeutet, dass kein Trainings-Daten-Set mit bekannten Betrugsfällen als Grundlage erforderlich ist.

Während herkömmliche Methoden zur Erkennung von Ausreißern in der Regel nur ein bis zwei Variablen gleichzeitig betrachten, kann die Anomalieerkennung zahlreiche Felder untersuchen und somit Cluster oder Vergleichsgruppen bilden, in die ähnliche Datensätze fallen. Die einzelnen Datensätze können dann jeweils mit den anderen Datensätzen in der Vergleichsgruppe verglichen werden, um so mögliche Anomalien zu erkennen. Je weiter ein Fall vom normalen Zentrum entfernt ist, desto größer ist die Wahrscheinlichkeit, dass dieser Fall ungewöhnlich ist. Der Algorithmus kann beispielsweise die Datensätze zu drei unterschiedlichen Clustern zusammenfassen und dann die Datensätze mit einem Flag versehen, die weit vom Zentrum des jeweiligen Clusters entfernt sind.

Abbildung 4-6  
Erkennen potenzieller Anomalien über die Clusterbildung



Jeder Datensatz wird einem Anomalieindex zugewiesen, der dem Verhältnis des Gruppenabweichungsindex zum Durchschnitt des Clusters darstellt, zu dem der Fall gehört. Je größer der Wert dieses Index ist, desto stärker ist die Abweichung des Falls vom Durchschnitt. Unter normalen Umständen würden Fälle mit einem Anomalieindex kleiner als 1 oder ggf. auch 1,5 nicht als Anomalien betrachtet, weil die Abweichung nahezu mit dem Durchschnitt übereinstimmt oder nur wenig höher liegt. Fälle mit einem Indexwert größer 2 sind dagegen vielversprechende Anomaliekandidaten, weil die Abweichung hierbei mindestens das Zweifache des Durchschnitts beträgt.

Die Anomalieerkennung ist eine Untersuchungsmethode, mit der ungewöhnliche Fälle oder Datensätze rasch erkannt werden, die als Kandidaten für die weitere Analyse infrage kommen. Diese Kandidaten gelten als *vermutete* Anomalien, die sich bei näherer Untersuchung als tatsächliche Anomalien herausstellen können (oder auch nicht). Unter Umständen stufen Sie einen Datensatz als völlig normal ein, den Sie jedoch beim Aufbauen eines Modells von den Daten abschirmen möchten. Umgekehrt gilt: Wenn der Algorithmus wiederholt falsche Anomalien zurückliefert, kann dies auf einen Fehler oder ein Artefakt bei der Datensammlung hinweisen.

Beachten Sie, dass bei der Anomalieerkennung ungewöhnliche Datensätze oder Fälle mithilfe einer Cluster-Analyse ermittelt werden, die auf der im Modell ausgewählten Menge an Feldern beruht – ohne Berücksichtigung eines speziellen Zielfelds (abhängigen Felds) und unabhängig davon, ob diese Felder für das Muster relevant sind, das Sie vorherzusagen versuchen. Aus diesem Grund sollten Sie die Anomalieerkennung in Kombination mit der Merkmalsauswahl oder einem anderen Verfahren für Screening und Rangordnung von Feldern verwenden. Beispielsweise können Sie mithilfe der Merkmalsauswahl die wichtigsten Felder in Bezug auf ein bestimmtes Ziel ermitteln und anschließend mit der Anomalieerkennung die Datensätze finden, die in Bezug auf diese Felder besonders ungewöhnlich sind. (Eine alternative Vorgehensweise besteht darin, ein Entscheidungsbaummodell zu erstellen und anschließend alle falsch klassifizierten Datensätze

als potenzielle Anomalien zu untersuchen. Diese Methode lässt sich jedoch nicht so leicht reproduzieren bzw. in größerem Maßstab automatisieren.)

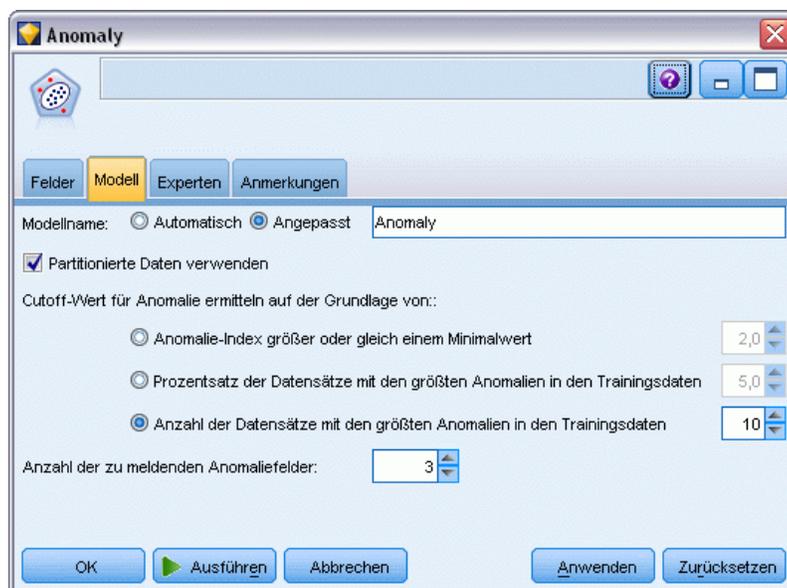
**Beispiel.** Bei der Untersuchung landwirtschaftlicher Subventionen auf mögliche Fälle von Betrug kann die Anomalieerkennung verwendet werden, um Abweichungen von der Norm aufzudecken, indem diejenigen Datensätze gekennzeichnet werden, die Unregelmäßigkeiten aufweisen und weiter untersucht werden müssen. Sie sind in erster Linie an Subventionsanträgen interessiert, die für den Typ und die Größe des landwirtschaftlichen Betriebs offenbar zu viel (oder zu wenig) Geld beantragen.

**Anforderungen.** Ein oder mehrere Eingabefelder. Beachten Sie, dass nur bei Feldern, bei denen eine Rolle auf Eingabe gesetzt ist, Quellen- oder Typknoten als Eingabe verwendet werden können. Zielfelder (Rolle auf Ziel oder Beides gesetzt) werden ignoriert.

**Stärken.** Durch die Kennzeichnung von Fällen, die einem bekannten Regel-Set *nicht* entsprechen (anstatt diejenigen Fälle zu kennzeichnen, die den Regeln entsprechen) können Anomalieerkennungsmodelle ungewöhnliche Fälle ermitteln, selbst wenn diese keinem zuvor bekannten Muster folgen. Bei Verwendung in Kombination mit der Merkmalsauswahl kann mithilfe der Anomalieerkennung ein Screening großer Datenmengen durchgeführt werden, um die relevantesten Datensätze relativ schnell zu ermitteln.

## Anomalieerkennung – Modelloptionen

Abbildung 4-7  
Anomalieerkennung – Registerkarte "Modell"



**Modellname.** Sie können den Modellnamen automatisch basierend auf den Ziel- oder ID-Feldnamen (oder dem Modelltyp in Fällen, in denen kein solches Feld angegeben ist) generieren oder einen benutzerdefinierten Namen eingeben.

**Cutoff-Wert für Anomalie ermitteln auf der Grundlage von.** Gibt die Methode für die Bestimmung des Cutoff-Werts zur Kennzeichnung von Anomalien an. Die folgenden Optionen sind verfügbar:

- **Minimale Anomalie-Indexebene.** Gibt den minimalen Cutoff-Wert für die Kennzeichnung von Anomalien an. Datensätze, die diesen Schwellenwert erreichen oder überschreiten, werden gekennzeichnet.
- **Prozentsatz der Datensätze mit den größten Anomalien in den Trainingsdaten.** Legt den Schwellenwert automatisch auf einem Niveau fest, bei dem der angegebene Prozentsatz an Datensätzen in den Trainingsdaten gekennzeichnet wird. Der resultierende Cutoff wird als Parameter in das Modell aufgenommen. Beachten Sie, dass mit dieser Option bestimmt wird, wie der Cutoff-Wert festgelegt wird, *nicht* jedoch der tatsächliche Prozentsatz der beim Scoring zu kennzeichnenden Datensätze. Die tatsächlichen Scoring-Ergebnisse können je nach den Daten abweichen.
- **Anzahl der Datensätze mit den größten Anomalien in den Trainingsdaten.** Legt den Schwellenwert automatisch auf einem Niveau fest, bei dem die angegebene Anzahl an Datensätzen in den Trainingsdaten gekennzeichnet wird. Der resultierende Schwellenwert wird als Parameter in das Modell aufgenommen. Beachten Sie, dass mit dieser Option bestimmt wird, wie der Cutoff-Wert festgelegt wird, *nicht* jedoch die konkrete Anzahl der beim Scoring zu kennzeichnenden Datensätze. Die tatsächlichen Scoring-Ergebnisse können je nach den Daten abweichen.

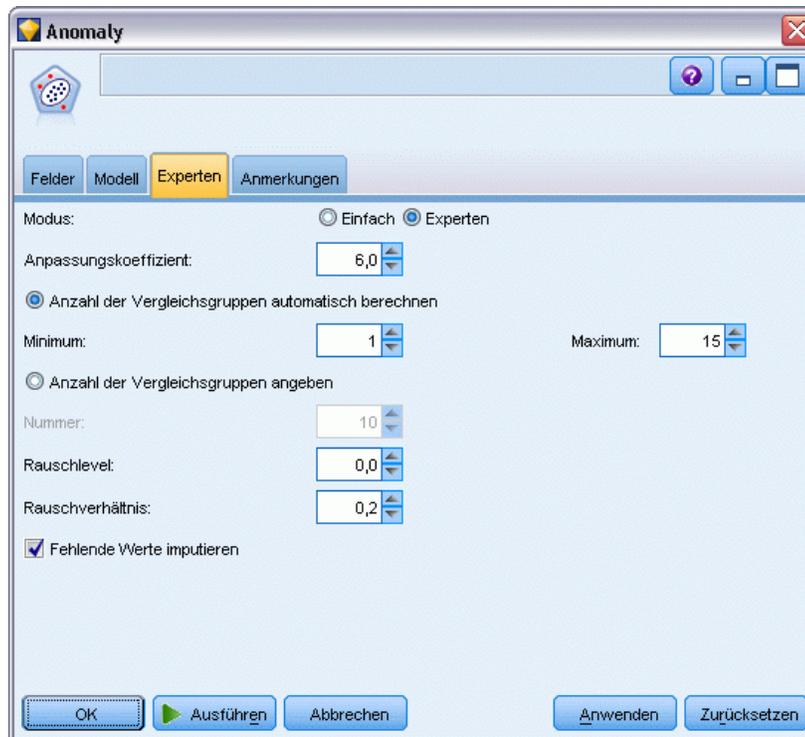
*Hinweis:* Unabhängig davon, wie der Cutoff-Wert bestimmt wird, hat er keine Auswirkungen auf den zugrunde liegenden Anomalieindexwert, der für die einzelnen Datensätze gemeldet wird. Er legt einfach den Schwellenwert fest, ab dem die Datensätze beim Schätzen oder Scoring des Modells als anomal gekennzeichnet werden sollen. Wenn Sie später eine größere oder kleinere Anzahl von Datensätzen untersuchen möchten, können Sie einen Auswahlknoten verwenden, um eine Teilmenge der Datensätze auf der Grundlage des Anomalieindexwertes ( $\$0\text{-AnomalyIndex} > X$ ) zu identifizieren.

**Anzahl der zu meldenden Anomaliefelder.** Gibt an, wie viele Felder gemeldet werden sollen, um anzugeben, warum ein bestimmter Datensatz als Anomalie gekennzeichnet wird. Die Felder mit den größten Anomalien werden gemeldet. Diese Felder sind definiert als diejenigen, die die größte Abweichung von der Feldnorm für den Cluster aufweisen, dem der Datensatz zugeordnet ist.

## ***Anomalieerkennung– Expertenoptionen***

Um Optionen für fehlende Werte und andere Einstellungen anzugeben, setzen Sie den Modus auf der Registerkarte “Experten” auf Experten.

Abbildung 4-8  
Anomalieerkennung – Registerkarte "Experten"



**Anpassungskoeffizient.** Wert, der zum Balancieren des relativen Gewichts verwendet wird, das den stetigen Feldern (numerischer Bereich) und den kategorialen Feldern bei der Berechnung der Distanz zugewiesen wird. Größere Werte erhöhen den Einfluss der stetigen Felder. Dieser Wert darf kein Nullwert sein.

**Anzahl der Vergleichsgruppen automatisch berechnen.** Mit der Anomalieerkennung lässt sich eine schnelle Analyse einer großen Anzahl von möglichen Lösungen durchführen, um die optimale Anzahl an Vergleichsgruppen für die Trainingsdaten auszuwählen. Sie können den Bereich erweitern oder einengen, indem Sie die minimale bzw. maximale Anzahl an Vergleichsgruppen festlegen. Mit größeren Werten kann das System einen breiteren Bereich möglicher Lösungen erkunden, allerdings erhöht sich dadurch die Verarbeitungszeit.

**Anzahl der Vergleichsgruppen angeben.** Wenn Sie wissen, wie viele Cluster in Ihr Modell aufgenommen werden sollen, wählen Sie diese Option und geben Sie die Anzahl der Vergleichsgruppen ein. Die Auswahl dieser Option führt im Allgemeinen zu einer besseren Leistung.

**Rauschlevel und Rauschverhältnis.** Diese Einstellungen bestimmen, wie Ausreißer bei der zweistufigen Clusterbildung behandelt werden. In der ersten Stufe wird ein Cluster-Merkmalbaum (CF-Baum) verwendet, um die Daten aus einer sehr großen Anzahl einzelner Datensätze auf eine überschaubare Anzahl von Clustern zu verdichten. Der Baum wird anhand von Ähnlichkeitsmaßen erstellt. Wenn ein Knoten des Baums zu viele Datensätze enthalten würde, wird er in untergeordnete Knoten aufgespalten. In der zweiten Phase beginnt die hierarchische Clusterbildung an den Endknoten des CD-Baums. Die Rauschverarbeitung

ist bei der ersten Datenübergabe aktiviert und bei der zweiten Datenübergabe deaktiviert. Die Fälle im Rausch-Cluster aus der ersten Datenübergabe werden in der zweiten Datenübergabe den regulären Clustern zugewiesen.

- **Ebene des Rauschens.** Geben Sie einen Wert zwischen 0 und 0,5 an. Diese Einstellung ist nur relevant, sofern der CF-Baum während der Wachstumsphase gefüllt wird, wenn er also keine weiteren Fälle in einem Blattknoten annehmen kann und kein Blattknoten aufgeteilt werden kann.

Wenn der CF-Baum gefüllt wird und die Rausch-Ebene auf 0 gesetzt ist, wird der Schwellenwert erhöht und der CF-Baum wird mit allen Fällen neu erstellt. Nach der abschließenden Clusteranalyse werden die Werte, die keinem Cluster zugewiesen werden konnten, als Ausreißer bezeichnet. Der Ausreißer-Cluster erhält die Identifikationsnummer -1 und wird nicht in die Auszählung der Anzahl der Cluster aufgenommen, d. h., wenn Sie  $n$  Cluster und Rauschverarbeitung angeben, gibt der Algorithmus  $n$  Cluster und einen Rausch-Cluster aus. In der Praxis bedeutet dies, dass bei einer Erhöhung dieses Werts der Algorithmus mehr Spielraum hat, ungewöhnliche Datensätze in den Baum einzupassen. Er muss sie also nicht einem gesonderten Ausreißer-Cluster zuweisen.

Wenn der CF-Baum gefüllt wird und die Rauschebene größer als 0 ist, wird der CF-Baum neu gebildet, nachdem alle Daten in "dünn besetzten" Blättern in einem eigenen Rausch-Blatt abgelegt wurden. Ein Blatt wird als "dünn besetzt" betrachtet, wenn das Verhältnis der Anzahl der Fälle im dünn besetzten Blatt zu der Anzahl der Fälle im größten Blatt kleiner ist als die Rauschebene. Nach der Neubildung des Baums werden die Ausreißer nach Möglichkeit im CF-Baum positioniert. Andernfalls werden die Ausreißer für die zweite Phase der Clusterbildung verworfen.

- **Rauschverhältnis.** Gibt an, welcher Anteil des der Komponente zugeordneten Arbeitsspeichers für die Rausch-Pufferung verwendet werden soll. Dieser Wert liegt im Bereich von 0,0 bis 0,5. Wenn das Hinzufügen eines gegebenen Falls zu einem Blatt des Baums zu einer Dichte unterhalb dieses Schwellenwerts führen würde, wird das Blatt nicht geteilt. Wenn die Dichte den Schwellenwert überschreitet, wird das Blatt geteilt und ein weiterer kleiner Cluster wird zum CF-Baum hinzugefügt. Durch die Erhöhung dieser Einstellung strebt der Algorithmus also möglicherweise schneller hin zu einem einfacheren Baum.

**Fehlende Werte vorschreiben.** Setzt bei stetigen Feldern für fehlende Werte den Feldmittelwert ein. Bei kategorialen Feldern werden fehlende Kategorien kombiniert und als gültige Kategorie behandelt. Wenn die Auswahl dieser Option aufgehoben wird, werden alle Datensätze mit fehlenden Werten aus der Analyse ausgeschlossen.

## ***Modell-Nuggets vom Typ "Anomalieerkennung"***

Modell-Nuggets vom Typ "Anomalieerkennung" enthalten alle Informationen, die vom Anomalieerkennungsmodell erfasst wurden, sowie Informationen zu den Trainingsdaten und dem Schätzworgang.

Bei der Ausführung eines Streams, der ein Modell-Nugget vom Typ "Anomalieerkennung" enthält, wird eine Reihe neuer Felder zum Stream hinzugefügt. Diese werden durch die Auswahl auf der Registerkarte "Einstellungen" im Modell-Nugget festgelegt. [Für weitere Informationen](#)

siehe Thema [Anomalieerkennungsmodell – Einstellungen auf S. 91](#). Neue Feldnamen beruhen auf dem Modellnamen und tragen das Präfix  $\$O$ , wie in der folgenden Tabelle zusammengefasst:

$\$O$ -Anomaly	Flag-Feld das angibt, ob der Datensatz anomal ist oder nicht.
$\$O$ -AnomalyIndex	Der Anomalieindexwert für den Datensatz.
$\$O$ -PeerGroup	Gibt die Vergleichsgruppe an, der der Datensatz zugewiesen ist.
$\$O$ -Field- $n$	Name des Felds, das den $n$ . Rang in der Reihenfolge der anomalsten Felder einnimmt (hinsichtlich der Abweichung von der Cluster-Norm).
$\$O$ -FieldImpact- $n$	Variabler Abweichungsindex für das Feld. Dieser Wert misst die Abweichung von der Feldnorm für den Cluster, dem der Datensatz zugewiesen ist.

Optional können Sie Scores für nicht anomale Datensätze unterdrücken, um die Lesbarkeit der Ergebnisse zu verbessern.

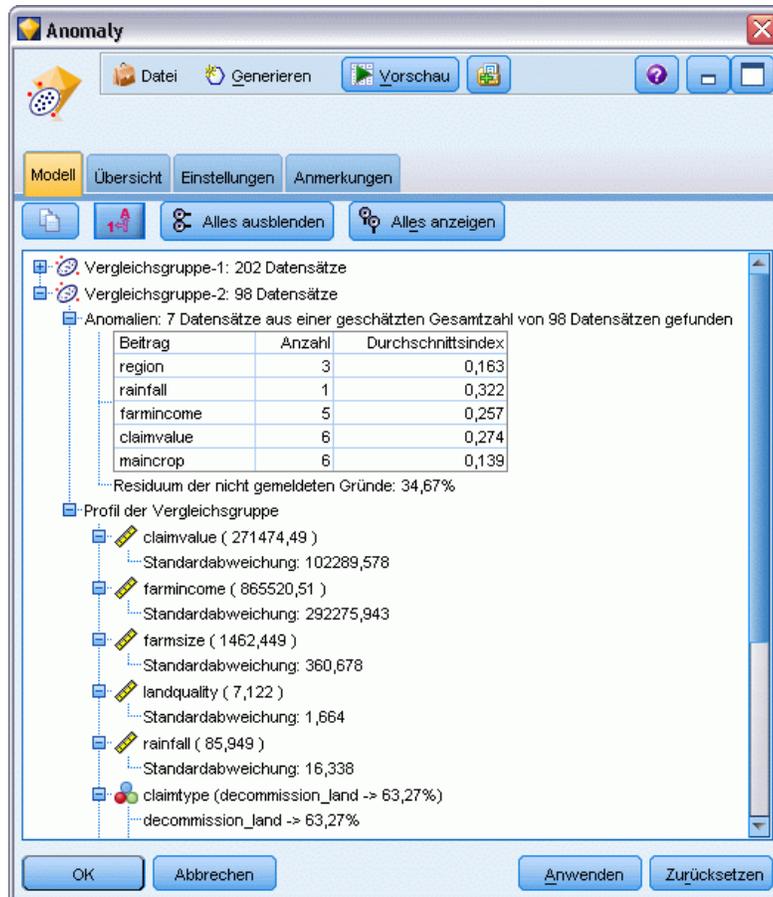
Abbildung 4-9  
Scoring von Ergebnissen mit Unterdrückung nicht anomaler Datensätze

	id	$\$O$ -Anomaly	$\$O$ -AnomalyIndex	$\$O$ -PeerGroup	$\$O$ -Field-1	$\$O$ -FieldImpact-1	$\$O$ -Field-2	$\$O$ -FieldImpact-2
1	id633	T	1.600	2	claimvalue	0.358	farmincome	0.275
2	id647	T	1.403	2	farminco...	0.334	claimvalue	0.161
3	id654	T	1.495	2	rainfall	0.322	maincrop	0.181
4	id703	T	1.358	1	rainfall	0.230	region	0.219
5	id704	T	1.427	2	farminco...	0.287	maincrop	0.190
6	id739	T	1.684	2	claimvalue	0.404	farmincome	0.233
7	id752	T	1.770	2	claimvalue	0.391	farmincome	0.155
8	id791	T	1.386	1	maincrop	0.236	rainfall	0.163
9	id813	T	1.641	1	region	0.181	landquality	0.160
10	id883	T	1.350	2	region	0.187	maincrop	0.169

## Anomalieerkennungsmodelle – Details

Auf der Registerkarte “Modell” für ein generiertes Anomalieerkennungsmodell werden Informationen zu den Vergleichsgruppen im Modell angezeigt.

Abbildung 4-10  
 Modell-Nuggets vom Typ "Anomalieerkennung" – Details

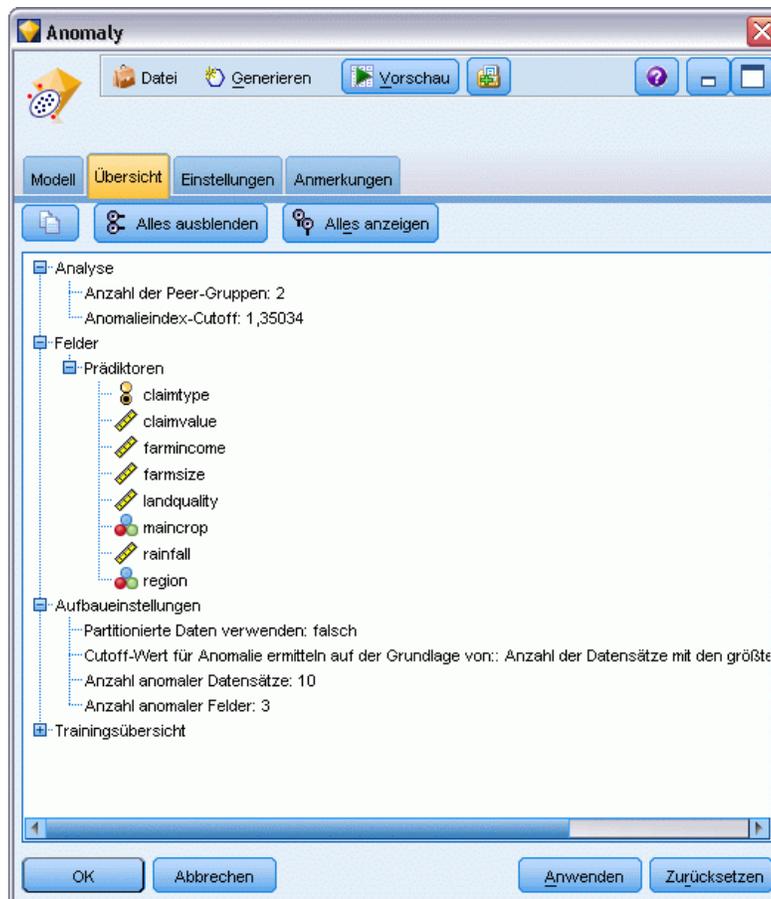


Beachten Sie, dass die gemeldeten Größen und Statistiken für die Vergleichsgruppe auf den Trainingsdaten beruhen und etwas vom tatsächlichen Scoring-Ergebnis abweichen können, selbst wenn dieselben Daten verwendet werden.

### **Anomalieerkennungsmodell – Übersicht**

Auf der Registerkarte "Übersicht" für ein Modell-Nugget vom Typ "Anomalieerkennung" werden Informationen über die Felder, die Aufbaueinstellungen und den Schätzvorgang angezeigt. Außerdem werden die Anzahl der Vergleichsgruppen sowie die Anzahl des Cutoff-Werts angezeigt, der verwendet wird, um Datensätze als anomal zu kennzeichnen.

Abbildung 4-11  
Modell-Nuggets vom Typ "Anomalieerkennung" – Übersicht



### **Anomalieerkennungsmodell – Einstellungen**

Auf der Registerkarte "Einstellungen" können Sie Optionen für das Scoring des Modell-Nuggets angeben.

Abbildung 4-12  
Scoring-Optionen für ein Modell-Nugget vom Typ "Anomalieerkennung"



**Anomale Datensätze kennzeichnen mit.** Gibt an, wie anomale Datensätze in der Ausgabe behandelt werden sollen.

- **Flag und Index.** Erstellt ein Flag-Feld, das für alle Datensätze, die den im Modell enthaltenen Cutoff-Wert überschreiten, auf *True* (Wahr) gesetzt ist. Der Anomalieindex wird außerdem für jeden Datensatz in einem separaten Feld angegeben. [Für weitere Informationen siehe Thema Anomalieerkennung – Modelloptionen auf S. 85.](#)
- **Nur Flag.** Erstellt ein Flag-Feld, jedoch ohne den Anomalieindex für die einzelnen Datensätze zu melden.
- **Nur Index.** Meldet den Anomalieindex, ohne ein Flag-Feld zu erstellen.

**Anzahl der zu meldenden Anomaliefelder.** Gibt an, wie viele Felder gemeldet werden sollen, um anzugeben, warum ein bestimmter Datensatz als Anomalie gekennzeichnet wird. Die Felder mit den größten Anomalien werden gemeldet. Diese Felder sind definiert als diejenigen, die die größte Abweichung von der Feldnorm für den Cluster aufweisen, dem der Datensatz zugeordnet ist.

**Datensätze verwerfen.** Wählen Sie diese Option, um alle nicht anomalen Datensätze aus dem Stream zu verwerfen. Dadurch können Sie sich leichter auf potenzielle Anomalien in abwärtsgelegenen Knoten konzentrieren. Alternativ können Sie festlegen, dass alle anomalen Datensätze verworfen werden sollen, um die nachfolgende Analyse auf diejenigen Datensätze zu begrenzen, die nicht auf der Grundlage des Modells als potenzielle Anomalien gekennzeichnet wurden.

*Hinweis:* Aufgrund kleiner Unterschiede beim Runden stimmt die tatsächliche Anzahl beim Scoring gekennzeichneten Datensätze eventuell nicht mit der Anzahl der beim Trainieren des Modells gekennzeichneten Datensätze überein, selbst wenn in beiden Fällen dieselbe Datengrundlage verwendet wurde.

# ***Knoten für die automatisierte Modellierung***

Die automatisierten Modellierungsknoten schätzen und vergleichen eine Reihe von verschiedenen Modellierungsmethoden und ermöglichen Ihnen, eine Vielzahl von Methoden in einem einzigen Modellierungsdurchgang auszuprobieren. Sie können die zu verwendenden Modellierungsalgorithmen und die jeweils spezifischen Optionen auswählen, inklusive Kombinationen, die sich anderenfalls gegenseitig ausschließen würden. Beispielsweise müssen Sie sich nicht zwischen den Methoden “Schnell”, “Dynamisch” und “Reduzieren” für ein neuronales Netzwerk entscheiden, sondern können alle drei Methoden ausprobieren. Der Knoten untersucht jede mögliche Kombination von Optionen, stuft jedes in Frage kommende Modell auf der Basis des angegebenen Werts und speichert die geeignetsten Kombinationen in Scoring oder weiterer Analyse.

Sie können je nach Anforderungen Ihrer Analyse zwischen drei Knoten für automatisierte Modellierung wählen:



Mit dem Knoten “Autom. Klassifizierer” können Sie eine Reihe verschiedener Modelle für binäre Ergebnisse (“Ja” oder “Nein”, “Abwanderung” oder “Keine Abwanderung” usw.) erstellen und vergleichen, um den besten Ansatz für die jeweilige Analyse auszuwählen. Es wird eine Reihe von Modellierungsalgorithmen unterstützt, sodass Sie die gewünschten Methoden, die spezifischen Optionen für die jeweilige Methode und die Kriterien zum Vergleich der Ergebnisse auswählen können. Der Knoten generiert eine Gruppe von Modellen, die auf den angegebenen Optionen beruhen, und erstellt anhand der von Ihnen angegebenen Kriterien eine Rangordnung der besten Kandidaten. [Für weitere Informationen siehe Thema Knoten “Automatischer Klassifizierer” auf S. 96.](#)



Der Knoten “Auto-Numerisch” schätzt und vergleicht mit einer Reihe verschiedener Methoden Modelle für die Ergebnisse stetiger numerischer Bereiche. Der Knoten arbeitet auf dieselbe Weise wie der Knoten “Automatischer Klassifizierer”: Sie können die zu verwendenden Algorithmen auswählen und in einem Modellierungsdurchlauf mit mehreren Optionskombinationen experimentieren. Folgende Algorithmen werden unterstützt: C&RT-Baum, CHAID, lineare Regression, verallgemeinerte lineare Regression und Support Vector Machines (SVM). Modelle können anhand von Korrelation, relativem Fehler bzw. Anzahl der verwendeten Variablen verglichen werden. [Für weitere Informationen siehe Thema Knoten “Auto-Numerisch” auf S. 107.](#)



Mit dem Knoten “Autom. Cluster” können Sie Clustering-Modelle, die Gruppen und Datensätze mit ähnlichen Merkmalen identifizieren, schätzen und vergleichen. Die Funktionsweise des Knotens gleicht der von anderen Knoten für automatisierte Modellierung, und Sie können in einem einzigen Modellierungsdurchgang mit mehreren Optionskombinationen experimentieren. Modelle können mithilfe grundlegender Messwerte für Filterung und Rangfolge der Nützlichkeit von Cluster-Modellen verglichen werden, um ein Maß auf der Basis der Wichtigkeit von bestimmten Feldern zu liefern. [Für weitere Informationen siehe Thema Knoten “Autom. Cluster” auf S. 114.](#)

Die besten Modelle werden in einem einzigen kombinierten Modell-Nugget gespeichert, damit Sie in ihnen navigieren, sie vergleichen und die gewünschten Modelle für das Scoring auswählen können.

- Nur für binäre, nominale und numerische Ziele können Sie mehrere Scoring-Modelle auswählen und die Scores in einem einzigen Modell-Ensemble kombinieren. Durch die Kombination der Vorhersagen aus mehreren Modellen lassen sich Begrenzungen, die einzelne Modelle aufweisen, vermeiden. Dadurch kann häufig eine höhere Gesamtgenauigkeit erreicht werden als mit einem der Modelle allein.
- Optional können Sie einen Drill-Down für die Ergebnisse durchführen und Modellierungsknoten oder Modell-Nuggets für jedes einzelne Modell generieren lassen, das Sie weiterverwenden oder eingehender untersuchen möchten.

### ***Modelle und Ausführungszeit***

Abhängig von dem Daten-Set und der Anzahl an Modellen kann die Ausführung von Knoten für automatisierte Modellierung Stunden oder noch länger dauern. Achten Sie bei der Auswahl von Optionen auf die Anzahl der erstellten Modelle. Nach Möglichkeit sollten Sie die Modellierungsdurchläufe für die Nacht oder das Wochenende planen, wenn die Systemressourcen weniger ausgelastet sind.

- Falls erforderlich kann die Anzahl der Datensätze im ursprünglichen Trainingsdurchlauf mithilfe eines Partitions- oder Stichprobenknotens reduziert werden. Nachdem Sie die Auswahl auf einige wenige infrage kommende Modelle eingeschränkt haben, kann das vollständige Daten-Set wiederhergestellt werden. Weitere Informationen finden Sie unter [Stichprobenknoten](#) oder [Partitionsknoten](#).
- Mit der Merkmalsauswahl können Sie die Anzahl der Eingabefelder verringern. [Für weitere Informationen siehe Thema Merkmalsauswahlknoten in Kapitel 4 auf S. 75](#). Alternativ können Sie durch die anfänglichen Modellierungsdurchgänge Felder und Optionen ermitteln, deren eingehendere Untersuchung lohnenswert sein könnte. Wenn beispielsweise die besten Modelle jeweils dieselben drei Felder zu verwenden scheinen, ist dies ein deutlicher Hinweis darauf, dass diese Felder beibehalten werden sollten.
- Optional können Sie die Zeit, die für die Schätzung eines Modells aufgewendet werden soll, begrenzen und die für Screening und Rangordnung der Modelle zu verwendenden Evaluationsmaße angeben.

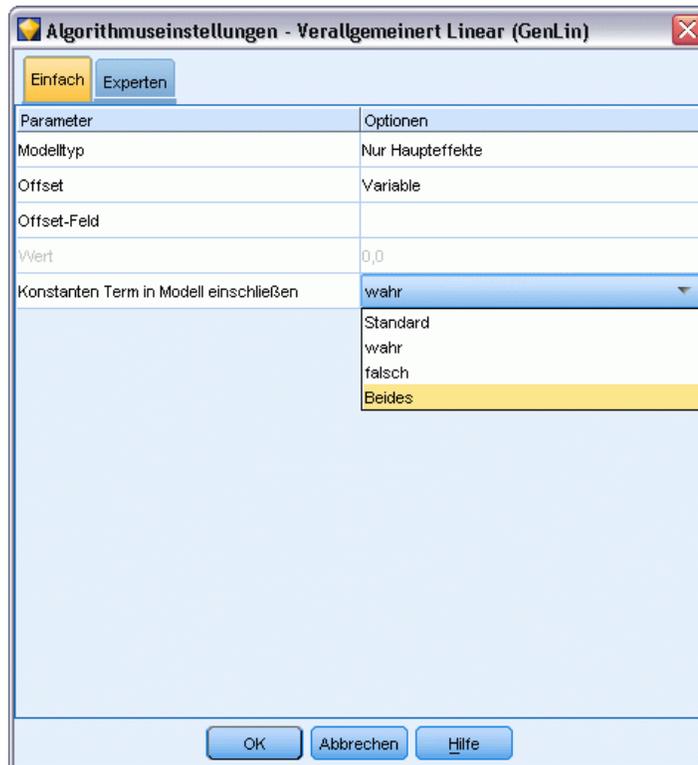
## ***Knoten für die automatisierte Modellierung – Algorithmuseinstellungen***

Sie können für die einzelnen Modelltypen entweder die Standardeinstellungen verwenden oder Optionen für den jeweiligen Modelltyp auswählen. Die einzelnen Optionen ähneln den Optionen, die in den gesonderten Modellierungsknoten verfügbar sind, mit dem Unterschied, dass nicht eine bestimmte Einstellung ausgewählt werden muss, sondern in den meisten Fällen beliebig viele Einstellungen verwendet werden können. So können Sie beispielsweise beim Vergleich von Netzwerkmodellen mehrere verschiedene Trainingsmethoden auswählen und jede Methode mit und ohne Zufallsstartwert ausprobieren. Es werden alle möglichen Kombinationen der ausgewählten Optionen verwendet. Dadurch wird es sehr einfach, viele verschiedene Modelle in

einem einzelnen Durchgang zu generieren. Seien Sie jedoch vorsichtig, da die Auswahl mehrerer Einstellungen die Anzahl der Modelle sehr schnell vervielfachen kann.

Abbildung 5-1

Auswahl der Algorithmeinstellungen für die automatisierte Modellierung



### ***Zum Auswählen der Optionen für den jeweiligen Modelltyp***

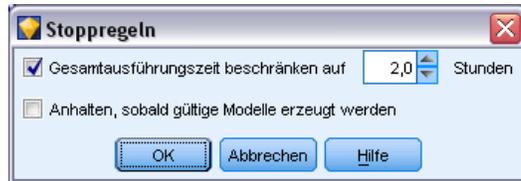
- ▶ am automatisierten Modellierungsknoten wählen Sie die Registerkarte Experten aus.
- ▶ Klicken Sie auf die Spalte Modellparameter für den Modelltyp.
- ▶ Wählen Sie im Dropdown-Menü die Option Angeben.
- ▶ Wählen Sie im Dialogfeld Algorithmeinstellungen die Optionen in der Spalte Optionen aus.

*Hinweis:* Weitere Optionen sind auf der Registerkarte “Experten” im Dialogfeld Algorithmusoptionen verfügbar.

## ***Knoten für die automatisierte Modellierung – Stoppregeln***

Für die Knoten “Automatisierte Modellierung” angegebene Stoppregeln beziehen sich auf die gesamte Knotenausführung, nicht auf das Stoppen einzelner, vom Knoten erstellter Modelle.

Abbildung 5-2  
Stoppregeln



**Gesamtausführungszeit beschränken auf.** (Nur für neuronale Netzwerk-, K-Means-, Kohonen-, TwoStep-, SVM-, KNN-, Bayes Net- und C&R-Baum-Modelle) Stoppt die Ausführung nach einer angegebenen Anzahl an Stunden. Alle bis zu diesem Zeitpunkt generierten Modelle werden in das Modell-Nugget aufgenommen, es werden jedoch keine weiteren Modelle erstellt.

**Anhalten, sobald gültige Modelle erzeugt werden.** Stoppt die Ausführung, wenn ein Modell alle auf der Registerkarte "Verwerfen" (für den Knoten "Automatischer Klassifizierer" oder "Autom. Cluster") oder der Registerkarte "Modell" (für den Knoten "Auto-Numerisch") angegebenen Kriterien erfüllt. [Für weitere Informationen siehe Thema Knoten "Automatischer Klassifizierer" - Optionen für Verwerfen auf S. 104.](#) [Für weitere Informationen siehe Thema Knoten "Autom. Cluster" - Optionen für Verwerfen auf S. 119.](#)

## ***Knoten "Automatischer Klassifizierer"***

Mit dem Knoten "Automatischer Klassifizierer" können Sie mit verschiedenen Methoden mehrere Modelle für nominale (Set) oder binäre Ziele (Ja/Nein) schätzen und vergleichen, wodurch Sie eine Vielzahl von Methoden in einem einzigen Modellierungsdurchgang ausprobieren können. Sie können die gewünschten Algorithmen auswählen und mit mehreren Kombinationen von Optionen experimentieren. Beispielsweise müssen Sie sich nicht zwischen den Methoden "Schnell", "Dynamisch" und "Reduzieren" für ein neuronales Netzwerk entscheiden, sondern können alle drei Methoden ausprobieren. Der Knoten untersucht jede mögliche Kombination von Optionen, stuft jedes in Frage kommende Modell auf der Basis des angegebenen Werts und speichert die geeignetsten Modelle für Scoring oder weitere Analyse. [Für weitere Informationen siehe Thema Knoten für die automatisierte Modellierung auf S. 93.](#)

Abbildung 5-3  
Modellierungsergebnisse mit automatischem Klassifizierer

Verwenden?	Diagramm	Modell	Erstellungszeit (Min.)	Maximaler Profit	Maximaler Profit	Lift{Oberste ...	Gesamt-Genauigkeit	Anzahl der verwendeten	Fläche unter Kurve
<input checked="" type="checkbox"/>		C5 1 < 1		4.906,667	8	2,203	92,861	10	0,777
<input checked="" type="checkbox"/>		C&...3		4.602,692	9	2,778	92,365	8	0,924
<input checked="" type="checkbox"/>		CH...3		4.145,668	8	2,851	91,706	4	0,927

**Beispiel.** Ein Einzelhandelsunternehmen verfügt über historische Daten, die die Angebote verfolgen, die bestimmten Kunden in früheren Werbeaktionen unterbreitet wurden. Das Unternehmen möchte nun profitablere Ergebnisse erzielen, indem es für jeden Kunden das richtige Angebot ermittelt. [So geht's](#)

**Anforderungen.** Ein Zielfeld mit einem Messniveau des Typs *Nominal* oder *Flag* (mit der Rolle Ziel) und mindestens ein Eingabefeld (mit der Rolle Eingabe). Für ein Flag-Feld gilt der für das Ziel definierte *Wahr*-Wert bei der Berechnung von Profiten, Lifts und verwandten Statistiken als Treffer. Eingabefelder können ein Messniveau von *Stetig* oder *Kategorial* aufweisen, mit der Einschränkung, dass einige Eingaben für bestimmte Modelltypen möglicherweise nicht geeignet sind. Ordinale Felder beispielsweise, die als Eingaben in Modellen vom Typ “C&R-Baum”, “CHAID” und “QUEST” verwendet werden sollen, müssen einen numerischen Speichertyp (nicht “Zeichenkette”) aufweisen und werden anderenfalls von diesen Modellen ignoriert. Ebenso können stetige Eingabefelder in einigen Fällen klassiert werden. Die Anforderungen sind dieselben wie bei Verwendung der einzelnen Modellierungsknoten; so funktioniert ein Bayes-Netz-Modell immer auf dieselbe Weise, unabhängig davon, ob es über den Knoten “Bayes-Netz” oder den Knoten “Automatischer Klassifizierer” generiert wurde.

**Häufigkeits- und Gewichtungsfelder.** Häufigkeit und Gewicht dienen dazu, einigen Datensätzen eine größere Bedeutsamkeit zu verleihen als anderen, beispielsweise, weil der Benutzer weiß, dass ein Teil der übergeordneten Grundgesamtheit im erstellten Daten-Set unterrepräsentiert ist (Gewicht) oder weil ein Datensatz für eine Reihe identischer Fälle steht (Häufigkeit). Häufigkeitsfelder können, sofern angegeben, von Modellen vom Typ “C&RT-Baum”, “CHAID”, “QUEST”, “Entscheidungsliste” und “Bayes-Netz” verwendet werden. Gewichtungsfelder können von Modellen vom Typ “C&RT”, “CHAID” und “C5.0” verwendet werden. Andere

Modelltypen ignorieren diese Felder und erstellen die Modelle in jedem Fall. Häufigkeits- und Gewichtungsfelder werden nur für die Modellerstellung verwendet. Bei der Evaluation bzw. beim Scoring von Modellen werden sie nicht berücksichtigt. [Für weitere Informationen siehe Thema Verwenden von Häufigkeits- und Gewichtungsfeldern in Kapitel 3 auf S. 39.](#)

### **Unterstützte Modelltypen**

Folgende Modelltypen werden unterstützt: “Netzwerk”, “C&RT-Baum”, “QUEST”, “CHAID”, “C5.0”, “Logistische Regression”, “Entscheidungsliste”, “Bayes-Netz”, “Diskriminanz”, “Nächster Nachbar” und “SVM”. [Für weitere Informationen siehe Thema Knoten “Automatischer Klassifizierer” – Expertenoptionen auf S. 100.](#)

## **Knoten “Automatischer Klassifizierer” - Modelloptionen**

Auf der Registerkarte “Modell” des Knotens “Automatischer Klassifizierer” können Sie die Anzahl der zu erstellenden Modelle sowie die zum Vergleich der Modelle verwendeten Kriterien angeben.

Abbildung 5-4

Knoten “Automatischer Klassifizierer” Registerkarte “Modell”

The screenshot shows the 'response' dialog box with the 'Modell' tab selected. The window title is 'response'. At the top, it says 'Geschätzte Anzahl der auszuführenden Modelle: 9'. Below this are tabs for 'Felder', 'Modell', 'Experten', 'Verwerfen', 'Einstellungen', and 'Anmerkungen'. The 'Modell' tab is active. The 'Modellname:' field has two radio buttons: 'Automatisch' (selected) and 'Angepasst'. There are two checked checkboxes: 'Partitionierte Daten verwenden' and 'Modell für jede Aufteilung aufbauen'. The 'Modelle in Ränge einteilen nach:' dropdown is set to 'Gesamtgenauigkeit'. The 'Modelle in Ränge einteilen mithilfe von:' has two radio buttons: 'Trainingspartition' and 'Testpartition' (selected). The 'Anzahl der zu verwendenden Modelle:' is set to 3. There is a checked checkbox for 'Einflussvariablenwichtigkeit berechnen'. Below this are two sections: 'Profitkriterien (nur gültig für Flag-Ziele)' and 'Lift-Kriterien (nur gültig für Flag-Ziele)'. The 'Profitkriterien' section has three rows: 'Kosten' (radio 'Fest' selected, value 5,0), 'Einkünfte' (radio 'Fest' selected, value 10,0), and 'Gewicht' (radio 'Fest' selected, value 1,0). Each row has a 'Variabel' option with a dropdown arrow. The 'Lift-Kriterien' section has a 'Für Lift-Berechnung zu verwendendes Perzentil:' set to 30. At the bottom are buttons for 'OK', 'Ausführen', 'Abbrechen', 'Anwenden', and 'Zurücksetzen'.

**Modellname.** Sie können den Modellnamen automatisch basierend auf den Ziel- oder ID-Feldnamen (oder dem Modelltyp in Fällen, in denen kein solches Feld angegeben ist) generieren oder einen benutzerdefinierten Namen eingeben.

**Partitionierte Daten verwenden.** Wenn ein Partitionsfeld definiert ist, gewährleistet diese Option, dass nur Daten aus der Trainingspartition für die Modellerstellung verwendet werden. [Für weitere Informationen siehe Thema Partitionsknoten in Kapitel 4 in IBM SPSS Modeler 14.2- Quellen-, Prozess- und Ausgabeknoten.](#)

**Aufteilungsmodelle erstellen.** Erstellt ein separates Modell für jeden möglichen Wert von Eingabefeldern, die als Aufteilungsfelder angegeben werden. [Für weitere Informationen siehe Thema Erstellung von aufgeteilten Modellen in Kapitel 3 auf S. 31.](#)

**Modelle in Ränge einteilen nach.** Legt fest, welche Kriterien für Vergleich und Rangordnung von Modellen verwendet werden sollen. Zu den Optionen gehören die Gesamtgenauigkeit, die Fläche unter der ROC-Kurve, Profit, Lift und die Anzahl der Felder. Beachten Sie, dass alle diese Maße im Zusammenfassungsbericht angegeben werden, unabhängig davon, welches davon an dieser Stelle ausgewählt wird.

*Hinweis:* Für ein nominales (Set-) Ziel ist die Rangbildung auf entweder **Gesamtgenauigkeit** oder **Anzahl der Felder** eingeschränkt.

Bei der Berechnung von Profiten, Lifts und verwandten Statistiken gilt der für das Zielfeld definierte *Wahr*-Wert als Treffer.

- 
- 
- 
- 
- 

**Modelle in Ränge einteilen mithilfe von.** Wenn eine Partition verwendet wird, können Sie angeben, ob die Ränge auf dem Trainingsdaten-Set oder auf dem Testdaten-Set beruhen sollen. Bei großen Daten-Sets lässt sich die Leistungsfähigkeit durch die Verwendung einer Partition für ein erstes Screening der Modelle u. U. erheblich verbessern.

**Anzahl der zu verwendenden Modelle.** Legt die maximale Anzahl der Modelle fest, die in dem vom Knoten erstellten Modell-Nugget aufgeführt werden sollen. Die Modelle mit dem höchsten Rang werden gemäß dem angegebenen Rangordnungskriterium aufgeführt. Beachten Sie, dass eine Erhöhung dieses Grenzwerts die Leistungsgeschwindigkeit verringern kann. Der höchste zulässige Wert ist 100.

**Bedeutsamkeit der Prädiktoren berechnen.** Bei Modellen, die zu einem angemessenen Maß an Bedeutsamkeit führen, können Sie ein Diagramm anzeigen, in dem die relative Wichtigkeit der einzelnen Prädiktoren bei der Modellschätzung angegeben wird. Normalerweise ist es sinnvoll, die Modellierungsbemühungen auf die wichtigsten Prädiktoren zu konzentrieren und diejenigen Prädiktoren zu verwerfen bzw. zu ignorieren, die die geringste Bedeutung haben. Beachten Sie, dass sich bei einigen Modellen der Zeitaufwand für die Berechnung durch die Bedeutsamkeit der Prädiktoren erhöhen kann. Außerdem wird diese Option nicht empfohlen, wenn Sie einfach einen allgemeinen Vergleich zwischen vielen verschiedenen Modellen wünschen. Diese Option ist von

größerem Nutzen, wenn die Analyse bereits auf eine Handvoll Modelle eingeeignet wurde, die detaillierter untersucht werden sollen. [Für weitere Informationen siehe Thema Bedeutsamkeit des Prädiktors in Kapitel 3 auf S. 54.](#)

**Profitkriterien.***Anmerkung.* Nur für Flag-Ziele. Der Profit entspricht dem Umsatz für jeden Datensatz abzüglich der Kosten für den betreffenden Datensatz. Die Profite für ein Quantil entsprechen einfach der Summe der Profite für alle Datensätze im Quantil. Profite gelten definitionsgemäß nur für Treffer, Kosten dagegen für alle Datensätze.

- **Kosten.** Geben Sie die Kosten für die einzelnen Datensätze an. Wählen Sie die Option Fest oder Variabel für den Umsatz. Bei festen Kosten geben Sie den Wert der Kosten ein. Bei variablen Kosten klicken Sie auf die Feldauswahl-Schaltfläche und bestimmen Sie ein Feld als Kostenfeld.
- **Umsatz.** Geben Sie den Umsatz für die einzelnen Datensätze ein, die als Treffer gelten. Wählen Sie die Option Fest oder Variabel für den Umsatz. Bei einem festen Umsatz geben Sie den Wert des Umsatzes ein. Bei einem variablen Umsatz klicken Sie auf die Feldauswahl-Schaltfläche und bestimmen Sie ein Feld als Umsatzfeld.
- **Gewicht.** Wenn die Datensätze in den Daten für mehrere Einheiten stehen, können Sie die Ergebnisse mithilfe der Häufigkeitsgewichtungen anpassen. Geben Sie die Gewichtung für die einzelnen Datensätze im Feld Fest oder Variabel an. Bei einer festen Gewichtung geben Sie den Wert für das Gewicht an (die Anzahl der Einheiten pro Datensatz). Bei variablen Gewichtungen klicken Sie auf die Feldauswahl—Schaltfläche und bestimmen Sie ein Feld als Gewichtsfeld.

**Lift-Kriterien.***Anmerkung.* Nur für Flag-Ziele. Gibt das für Lift-Berechnungen zu verwendende Perzentil an. Diesen Wert können Sie auch beim Vergleichen der Ergebnisse ändern. [Für weitere Informationen siehe Thema Nugget für automatisierte Modellierung auf S. 120.](#)

### ***Knoten “Automatischer Klassifizierer” – Expertenoptionen***

Auf der Registerkarte “Experten” des Knotens “Automatischer Klassifizierer” können Sie eine Partition anwenden (sofern verfügbar), die zu verwendenden Algorithmen auswählen und Stoppregeln angeben.

Abbildung 5-5  
Knoten "Automatischer Klassifizierer" Registerkarte "Experten"



**Verwendete Modelle.** Wählen Sie anhand der Kontrollkästchen in der Spalte auf der linken Seite die Modelltypen (Algorithmen) aus, die in den Vergleich aufgenommen werden sollen. Je mehr Typen Sie auswählen, desto mehr Modelle werden erstellt und desto länger dauert die Verarbeitung.

**Modelltyp.** Listet die verfügbaren Algorithmen auf (siehe unten).

**Modellparameter.** Sie können für die einzelnen Modelltypen entweder die Standardeinstellungen verwenden oder mithilfe von Angeben Optionen für den jeweiligen Modelltyp auswählen. Die einzelnen Optionen ähneln den in den separaten Modellierungskonten verfügbaren Optionen, mit dem Unterschied, dass mehrere Optionen bzw. Kombinationen ausgewählt werden können. Beispiel: Beim Vergleich von neuronalen Netzwerkmodellen können Sie, anstatt eine der sechs Trainingsmethoden auszuwählen, alle sechs auswählen, um sechs Modelle in einem einzigen Durchgang zu trainieren.

**Anzahl der Modelle.** Listet die Anzahl der Modelle auf, die auf der Grundlage der aktuellen Einstellungen für die einzelnen Algorithmen erstellt wurden. Bei einer Kombination von Optionen kann die Anzahl der Modelle schnell recht groß werden. Daher wird dringend empfohlen, auf diesen Wert zu achten, insbesondere bei Verwendung großer Daten-Sets.

**Maximale Zeit für Erstellung eines einzelnen Modells beschränken.** (Nur für K-Means-, Kohonen-, TwoStep-, SVM-, KNN-, Bayes Net- und Entscheidungslisten-Modelle) Legt ein maximales Zeitlimit für jedes beliebige Modell fest. Wenn beispielsweise das Training für ein bestimmtes Modell aufgrund einer komplexen Interaktion unerwartet viel Zeit in Anspruch nimmt, wird durch diese Option vermieden, dass das Modell den gesamten Modellierungsdurchlauf aufhält.

*Hinweis:* Falls es sich um ein nominales (Set-) Ziel handelt, steht die Option “Entscheidungsliste” nicht zur Verfügung.

### Unterstützte Algorithmen



Der Netzwerkknoten verwendet ein vereinfachtes Modell der Art und Weise, wie ein menschliches Gehirn Informationen verarbeitet. Es funktioniert, indem eine große Anzahl miteinander verbundener einfacher Verarbeitungseinheiten simuliert wird, die abstrakten Versionen von Neuronen ähnlich sind. Neuronale Netze sind leistungsstarke Mehrzweck-Schätzer, für deren Training und Anwendung nur sehr geringe statistische oder mathematische Kenntnisse erforderlich sind.



Der C5.0-Knoten erstellt entweder einen Entscheidungsbaum oder ein Regel-Set. Das Modell teilt die Stichprobe auf der Basis des Felds auf, das auf der jeweiligen Ebene den maximalen Informationsgewinn liefert. Das Zielfeld muss kategorial sein. Es sind mehrere Aufteilungen in mehr als zwei Untergruppen zulässig. [Für weitere Informationen siehe Thema C5.0-Knoten in Kapitel 6 auf S. 177.](#)



Der Knoten für Klassifizierungs- und Regressions-Bäume (C&RT-Bäume) erstellt einen Entscheidungsbaum, mit dem Sie zukünftige Beobachtungen vorhersagen oder klassifizieren können. Bei dieser Methode wird eine rekursive Partitionierung verwendet, um die Trainingsdatensätze in Segmente aufzuteilen. Dabei wird bei jedem Schritt die Unreinheit verringert und ein Knoten im Baum wird als “rein” betrachtet, wenn 100 % der Fälle in eine bestimmte Kategorie des Zielfelds fallen. Ziel- und Eingabefelder können numerische Bereiche oder kategorial (nominal, ordinal oder Flags) sein. Alle Aufteilungen sind binär (nur zwei Untergruppen). [Für weitere Informationen siehe Thema C&R-Baumknoten in Kapitel 6 auf S. 155.](#)



Der QUEST-Knoten bietet eine binäre Klassifizierungsmethode zum Erstellen von Entscheidungsbaummodellen, die dazu dient, die für große C&R-Baum-Analysen erforderliche Verarbeitungszeit zu verkürzen. Gleichzeitig soll die in den Klassifizierungsbaummodellen festgestellte Tendenz verringert werden, die darin besteht, dass Eingaben bevorzugt werden, die mehr Aufteilungen erlauben. Eingabefelder können stetig (numerische Bereiche) sein, das Zielfeld muss aber kategorial sein. Alle Aufteilungen sind binär. [Für weitere Informationen siehe Thema QUEST-Knoten in Kapitel 6 auf S. 157.](#)



Der CHAID-Knoten erzeugt Entscheidungsbäume unter Verwendung von Chi-Quadrat-Statistiken zur Ermittlung optimaler Aufteilungen. Im Gegensatz zu den Knoten vom Typ “C&RT-Baum” und “QUEST” kann CHAID nichtbinäre Bäume generieren, d. h. Bäume mit Aufteilungen mit mehr als zwei Verzweigungen. Ziel- und Eingabefelder können in einem numerischen Bereich (stetig) oder kategorial sein. Exhaustive CHAID ist eine Änderung von CHAID, die noch gründlicher vorgeht, indem sie alle möglichen Aufteilungen untersucht, allerdings mehr Rechenzeit beansprucht. [Für weitere Informationen siehe Thema CHAID-Knoten in Kapitel 6 auf S. 156.](#)



Die logistische Regression ist ein statistisches Verfahren zur Klassifizierung von Datensätzen auf der Grundlage der Werte von Eingabefeldern. Sie ist analog zur linearen Regression, außer dass statt eines numerischen Bereichs ein kategoriales Zielfeld verwendet wird. [Für weitere Informationen siehe Thema Logistikknoten in Kapitel 10 auf S. 287.](#)



Der Knoten "Entscheidungsliste" kennzeichnet Untergruppen bzw. Segmente, die eine höhere oder geringere Wahrscheinlichkeit für ein bestimmtes binäres Ergebnis aufweisen als die Gesamtpopulation. Sie könnten beispielsweise nach Kunden suchen, deren Abwanderung unwahrscheinlich ist oder die mit großer Wahrscheinlichkeit positiv auf eine Kampagne reagieren. Sie können Ihr Geschäftswissen in das Modell integrieren, indem Sie eigene, benutzerdefinierte Segmente hinzufügen und eine Vorschau anzeigen, in der alternative Modelle nebeneinander angezeigt werden, um die Ergebnisse zu vergleichen. Entscheidungslistenmodelle bestehen aus einer Liste von Regeln, bei denen jede Regel eine Bedingung und ein Ergebnis aufweist. Regeln werden in der vorgegebenen Reihenfolge angewendet und die erste Regel, die zutrifft, bestimmt das Ergebnis. [Für weitere Informationen siehe Thema Entscheidungsliste in Kapitel 9 auf S. 226.](#)



Mithilfe des Bayes-Netzwerk-Knotens können Sie ein Wahrscheinlichkeitsmodell erstellen, indem Sie beobachtete und aufgezeichnete Hinweise mit Weltwissen kombinieren, um die Wahrscheinlichkeit ihres Vorkommens zu ermitteln. Der Knoten ist speziell für Netzwerke vom Typ "Tree Augmented Naive Bayes" (TAN) und "Markov-Decke" gedacht, die in erster Linie zur Klassifizierung verwendet werden. [Für weitere Informationen siehe Thema Bayes-Netzwerk-Knoten in Kapitel 7 auf S. 198.](#)



Bei der Diskriminanzanalyse werden strengere Annahmen als bei der logistischen Regression verwendet, sie kann jedoch eine wertvolle Alternative oder Ergänzung zu einer logistischen Regressionsanalyse sein, wenn diese Annahmen erfüllt sind. [Für weitere Informationen siehe Thema Diskriminanzknoten in Kapitel 10 auf S. 319.](#)



Der Knoten " $k$ -Nächste Nachbarn" (KNN) verknüpft einen neuen Fall mit der Kategorie oder dem Wert der  $k$  Objekte, die ihm im Prädiktorraum am nächsten liegen, wobei  $k$  eine Ganzzahl ist. Ähnliche Fälle liegen nah beieinander und Fälle mit geringer Ähnlichkeit sind weit voneinander entfernt. [Für weitere Informationen siehe Thema KNN-Knoten in Kapitel 16 auf S. 482.](#)



Der Knoten "Support Vector Machine" (SVM) ermöglicht die Klassifizierung von Daten in eine von zwei Gruppen ohne Überanpassung. SVM eignet sich gut für umfangreiche Daten-Sets, beispielsweise solche mit einer großen Anzahl an Eingabefeldern. [Für weitere Informationen siehe Thema SVM-Knoten in Kapitel 15 auf S. 476.](#)

## Fehlklassifizierungskosten

In einigen Zusammenhängen sind bestimmte Fehler kostspieliger als andere. Zum Beispiel kann es kostspieliger sein, einen Kreditantragsteller mit hohem Risiko als niedriges Risiko zu klassifizieren (eine Art von Fehler) als einen Kreditantragsteller mit niedrigem Risiko als hohes Risiko (eine andere Art von Fehler). Anhand von Fehlklassifizierungskosten können Sie die relative Bedeutung verschiedener Arten von Vorhersagefehlern angeben.

Fehlklassifizierungskosten sind im Grunde Gewichtungen, die auf bestimmte Ergebnisse angewendet werden. Diese Gewichtungen werden in das Modell mit einbezogen und können die Vorhersage (als Schutz gegen kostspielige Fehler) ändern.

Mit Ausnahme von C5.0-Modellen werden Fehlklassifizierungskosten beim Scoring von Modellen nicht angewendet und bei der Rangbildung bzw. beim Vergleich von Modellen mithilfe eines Knotens vom Typ "Automatischer Klassifizierer", eines Evaluationsdiagramms oder eines Analyseknosens nicht berücksichtigt. Ein Modell, das Kosten beinhaltet, führt nicht unbedingt zu weniger Fehlern als ein Modell, das keine Kosten beinhaltet, und es weist auch nicht unbedingt eine höhere Gesamtgenauigkeit auf, aber es ist möglicherweise in praktischer Hinsicht leistungsfähiger, da es eine integrierte Verzerrung zugunsten von *weniger teuren* Fehlern aufweist.

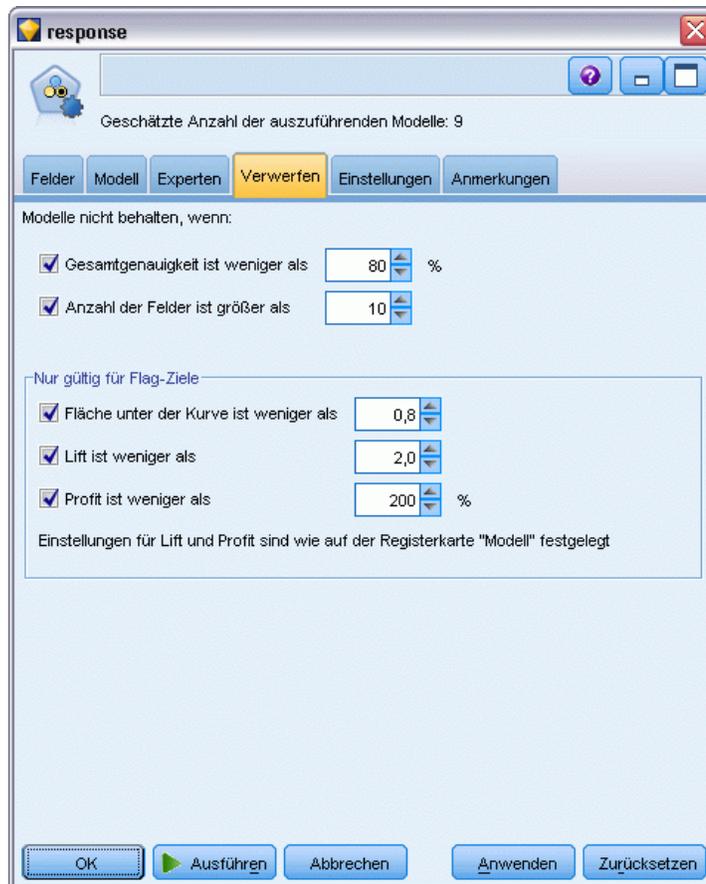
Die Kostenmatrix zeigt die Kosten für jede mögliche Kombination aus prognostizierter Kategorie und tatsächlicher Kategorie. Standardmäßig werden alle Fehlklassifizierungskosten auf 1,0 gesetzt. Um benutzerdefinierte Kostenwerte einzugeben, wählen Sie Fehlklassifizierungskosten verwenden und geben Ihre benutzerdefinierten Werte in die Kostenmatrix ein.

Um die Fehlklassifizierungskosten zu ändern, wählen Sie die Zelle aus, die der gewünschten Kombination aus vorhergesagten und tatsächlichen Werten entspricht, löschen den Inhalt der Zelle und geben die gewünschten Kosten für die Zelle ein. Die Kosten sind nicht automatisch symmetrisch. Wenn Sie z. B. die Kosten einer Fehlklassifizierung von *A* als *B* auf 2,0 setzen, weisen die Kosten für eine Fehlklassifizierung von *B* als *A* weiterhin den Standardwert 1,0 auf, es sei denn, Sie ändern diesen Wert ausdrücklich.

### ***Knoten "Automatischer Klassifizierer" - Optionen für Verwerfen***

Auf der Registerkarte "Verwerfen" des Knotens "Automatischer Klassifizierer" können Sie automatisch Modelle verwerfen, die bestimmte Kriterien nicht erfüllen. Diese Modelle werden nicht im Zusammenfassungsbericht aufgeführt.

Abbildung 5-6  
Knoten "Automatischer Klassifizierer" Registerkarte "Verwerfen"



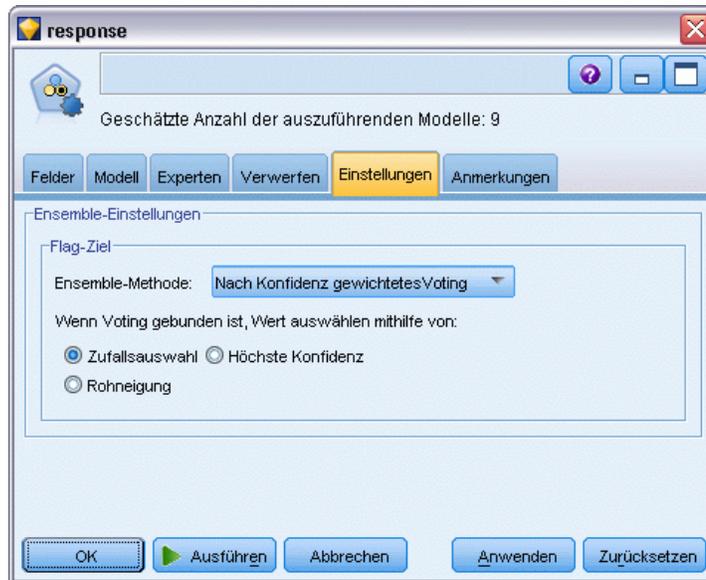
Sie können eine Untergrenze für "Gesamtgenauigkeit" sowie eine Obergrenze für die Anzahl der im Modell verwendeten Variablen festlegen. Für Flag-Ziele können Sie zusätzlich eine Untergrenze für "Lift", "Profit" und "Fläche unter der Kurve" angeben. Hierbei werden Lift und Profit gemäß den Festlegungen im Modellierungsknoten bestimmt. [Für weitere Informationen siehe Thema Knoten "Automatischer Klassifizierer" - Modelloptionen auf S. 98.](#)

Optional können Sie den Knoten so konfigurieren, dass die Ausführung gestoppt wird, sobald erstmals ein Modell generiert wurde, das alle angegebenen Kriterien erfüllt. [Für weitere Informationen siehe Thema Knoten für die automatisierte Modellierung – Stopregeln auf S. 95.](#)

### **Knoten "Automatischer Klassifizierer" - Einstellungsoptionen**

In der Registerkarte "Einstellungen" des Knotens "Automatischer Klassifizierer" können Sie die Score-Zeit-Optionen konfigurieren, die für das Nugget verfügbar sind.

Abbildung 5-7  
Knoten "Automatischer Klassifizierer" Registerkarte "Einstellungen"



**Ensemble-Methode.** Für Ziele können Sie aus den folgenden Ensemble-Methoden wählen:

- Voting
- Nach Konfidenz gewichtetes Voting
- Nach Rohneigung gewichtetes Voting (nur bei Flag-Zielen).
- Höchste Konfidenz hat Vorrang
- Durchschnittliche Rohneigung (nur bei Flag-Zielen).

Für weitere Informationen siehe Thema Ensemble-Knoten – Einstellungen in Kapitel 4 in *IBM SPSS Modeler 14.2- Quellen-, Prozess- und Ausgabeknoten*.

**Wenn Voting gebunden ist, Wert auswählen mithilfe von.** Bei Voting-Methoden können Sie auswählen, wie Gleichstände aufgelöst werden sollen:

- **Zufallsauswahl.** Einer der gebundenen Werte (Werte mit Gleichstand) wird nach dem Zufallsprinzip ausgewählt.
- **Höchste Konfidenz.** Der gebundene Wert, der mit der höchsten Konfidenz vorhergesagt wurde, gewinnt. Beachten Sie, dass es sich hierbei nicht unbedingt um die höchste Konfidenz aller vorhergesagten Werte handelt.
- **Rohneigung.** (nur bei Flag-Zielen) Der gebundene Wert, der mit der höchsten absoluten Neigung vorhergesagt wurde. Dabei berechnet sich die absolute Rohneigung wie folgt:

$$\text{abs}(0.5 - \text{propensity}) * 2$$

## Knoten "Auto-Numerisch"

Mit dem Knoten "Auto-Numerisch" können Sie mit verschiedenen Methoden mehrere Modelle für Ergebnisse stetiger numerischer Bereiche schätzen und vergleichen, wodurch Sie eine Vielzahl von Methoden in einem einzigen Modellierungsdurchgang ausprobieren können. Sie können die gewünschten Algorithmen auswählen und mit mehreren Kombinationen von Optionen experimentieren. Beispielsweise könnten Sie Immobilienwerte mithilfe von Modellen vom Typ "Netzwerk", "Lineare Regression", "C&RT" und "CHAID" vorhersagen, um zu ermitteln, welches Modell die beste Leistung bringt, und Sie könnten verschiedene Kombinationen der Regressionsmethoden "Schrittweise", "Vorwärts" und "Rückwärts" ausprobieren. Der Knoten untersucht jede mögliche Kombination von Optionen, stuft jedes in Frage kommende Modell auf der Basis des angegebenen Werts und speichert die geeignetsten Kombinationen in Scoring oder weiterer Analyse. [Für weitere Informationen siehe Thema Knoten für die automatisierte Modellierung auf S. 93.](#)

Abbildung 5-8  
Auto-Numerisch - Ergebnisse

Verwenden?	Diagramm	Modell	Erstellungszeit (Min.)	Korrelation	Anzahl der verwendeten	Relativer Fehler
<input checked="" type="checkbox"/>		Generaliz...	< 1	0,915	7	0,162
<input checked="" type="checkbox"/>		Regressi...	< 1	0,9	5	0,19
<input checked="" type="checkbox"/>		CHAID Tr...	< 1	0,892	5	0,204

**Beispiel.** Eine Gemeinde möchte die Immobiliensteuern mit größerer Genauigkeit schätzen und Werte für bestimmte Immobilien nach Bedarf anpassen, ohne jedes einzelne Anwesen besichtigen zu müssen. Mithilfe des Knotens "Auto-Numerisch" kann der Analyst eine Reihe von Modellen generieren und vergleichen, die Immobilienwerte auf der Grundlage von Gebäudetyp, Lage, Größe und anderen bekannten Faktoren vorhersagen. [Für weitere Informationen siehe Thema Eigenschaftswerte \(Auto-Numerisch\) in Kapitel 6 in IBM SPSS Modeler 14.2-Anwendungshandbuch.](#)

**Anforderungen.** Ein einzelnes Zielfeld (mit der Rolle Ziel) und mindestens ein Eingabefeld (mit der Rolle Eingabe). Beim Ziel muss es sich um ein stetiges Feld (numerischer Bereich) handeln, beispielsweise *Alter* oder *Einkommen*. Eingabefelder können stetig oder kategorial sein, mit der Einschränkung, dass einige Eingaben für bestimmte Modelltypen nicht geeignet sind. So können beispielsweise Modelle vom Typ “C&RT-Baum” kategoriale Zeichenkettenfelder als Eingaben verwenden, während lineare Regressionsmodelle diese Felder nicht verwenden können und sie ignorieren, wenn sie angegeben sind. Die Anforderungen sind dieselben wie bei Verwendung der einzelnen Modellierungsknoten. So funktioniert beispielsweise ein CHAID-Modell immer auf dieselbe Weise, unabhängig davon, ob es aus einem CHAID-Knoten oder aus einem Knoten vom Typ “Auto-Numerisch” generiert wurde.

**Häufigkeits- und Gewichtungsfelder.** Häufigkeit und Gewicht dienen dazu, einigen Datensätzen eine größere Bedeutsamkeit zu verleihen als anderen, beispielsweise, weil der Benutzer weiß, dass ein Teil der übergeordneten Grundgesamtheit im erstellten Daten-Set unterrepräsentiert ist (Gewicht) oder weil ein Datensatz für eine Reihe identischer Fälle steht (Häufigkeit). Häufigkeitsfelder können, sofern angegeben, von Algorithmen vom Typ “C&RT-Baum” und “CHAID” verwendet werden. Gewichtungsfelder können von Algorithmen vom Typ “C&RT”, “CHAID” und “GenLin” verwendet werden. Andere Modelltypen ignorieren diese Felder und erstellen die Modelle in jedem Fall. Häufigkeits- und Gewichtungsfelder werden nur für die Modellerstellung verwendet. Bei der Evaluation bzw. beim Scoring von Modellen werden sie nicht berücksichtigt. [Für weitere Informationen siehe Thema Verwenden von Häufigkeits- und Gewichtungsfeldern in Kapitel 3 auf S. 39.](#)

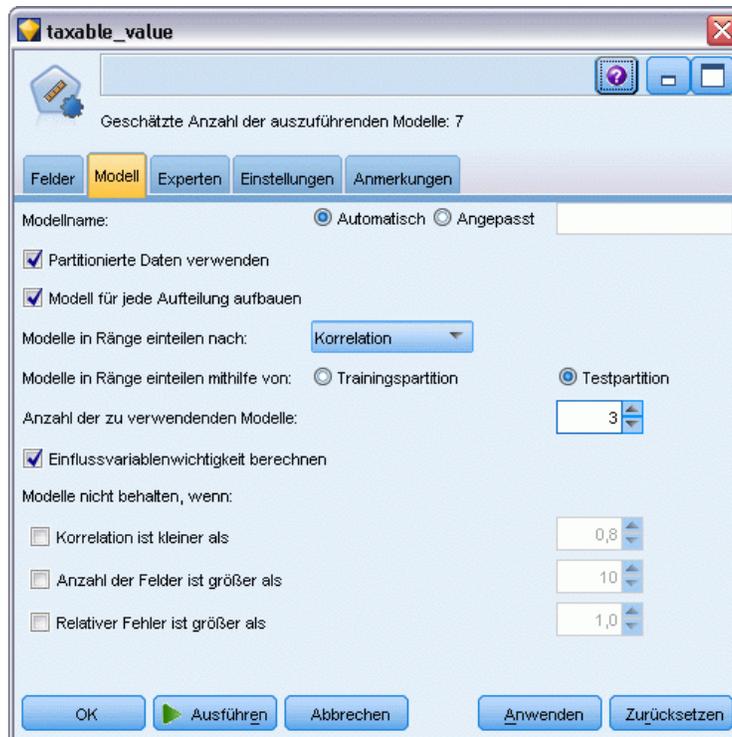
#### **Unterstützte Modelltypen**

Folgende Modelltypen werden unterstützt: “Netzwerk”, “C&RT-Baum”, “CHAID”, “Regression”, “GenLin”, “Nächster Nachbar” und SVM. [Für weitere Informationen siehe Thema Knoten “Auto-Numerisch” - Expertenoptionen auf S. 111.](#)

### **Knoten “Auto-Numerisch” - Modelloptionen**

Auf der Registerkarte “Modell” des Knotens “Auto-Numerisch” können Sie die Anzahl der zu speichernden Modelle sowie die zum Vergleich der Modelle verwendeten Kriterien angeben.

Abbildung 5-9  
Knoten "Auto-Numerisch": Registerkarte "Modell"



**Modellname.** Sie können den Modellnamen automatisch basierend auf den Ziel- oder ID-Feldnamen (oder dem Modelltyp in Fällen, in denen kein solches Feld angegeben ist) generieren oder einen benutzerdefinierten Namen eingeben.

**Partitionierte Daten verwenden.** Wenn ein Partitionsfeld definiert ist, gewährleistet diese Option, dass nur Daten aus der Trainingspartition für die Modellerstellung verwendet werden. [Für weitere Informationen siehe Thema Partitionsknoten in Kapitel 4 in IBM SPSS Modeler 14.2- Quellen-, Prozess- und Ausgabeknoten.](#)

**Aufteilungsmodelle erstellen.** Erstellt ein separates Modell für jeden möglichen Wert von Eingabefeldern, die als Aufteilungsfelder angegeben werden. [Für weitere Informationen siehe Thema Erstellung von aufgeteilten Modellen in Kapitel 3 auf S. 31.](#)

**Modelle in Ränge einteilen nach.** Legt fest, welche Kriterien zum Vergleichen von Modellen verwendet werden sollen.

- **Korrelation.** Die Pearson-Korrelation zwischen dem beobachteten Wert für die einzelnen Datensätze und dem vom Modell vorhergesagten Wert. Die Korrelation ist ein Maß für den linearen Zusammenhang zwischen zwei Variablen. Dabei deuten Werte nahe bei 1 auf eine stärkere Beziehung hin. (Die Korrelationswerte liegen im Bereich zwischen  $-1$  für eine exakte gegenläufige (negative) Beziehung) und  $+1$  für eine exakte gleichgerichtete (positive) Beziehung. Der Wert 0 bedeutet, dass keine lineare Beziehung besteht. Ein Modell mit einer negativen Korrelation weist den niedrigsten Rang auf.)

- **Anzahl der Felder.** Die Anzahl der Felder, die als Prädiktoren im Modell verwendet werden. Durch die Auswahl von Modellen mit weniger Feldern lässt sich in einigen Fällen die Datenvorbereitung rationalisieren und die Leistung verbessern.
- **Relativer Fehler.** Der relative Fehler ist der Quotient aus der Varianz der beobachteten Werte von den vom Modell vorhergesagten Werten und der Varianz der beobachteten Werte vom Mittelwert. Es wird also verglichen, wie gut die Leistungsfähigkeit des Modells in Bezug auf ein **Null-Modell** (leeres Modell) oder ein Modell nur mit dem **konstanten Term** ist, das einfach den Mittelwert des Zielfelds als Vorhersage ergibt. Bei einem guten Modell sollte dieser Wert kleiner als 1 sein, was darauf hinweist, dass das Modell genauer als das Null-Modell ist. Ein Modell mit einem relativen Fehler von mehr als 1 ist weniger genau als das Null-Modell und daher nicht brauchbar. Bei linearen Regressionsmodellen ist der relative Fehler gleich dem Quadrat der Korrelation und fügt ebenfalls keine neuen Informationen hinzu. Bei nichtlinearen Modellen steht der relative Fehler in keinem Zusammenhang zur Korrelation und bietet ein zusätzliches Maß für die Bewertung der Leistungsfähigkeit des Modells.

**Modelle in Ränge einteilen mithilfe von.** Wenn eine Partition verwendet wird, können Sie angeben, ob die Ränge auf der Trainingspartition oder auf der Testpartition beruhen sollen. Bei großen Daten-Sets lässt sich die Leistungsfähigkeit durch die Verwendung einer Partition für ein erstes Screening der Modelle u. U. erheblich verbessern. [Für weitere Informationen siehe Thema Partitionsknoten in Kapitel 4 in IBM SPSS Modeler 14.2- Quellen-, Prozess- und Ausgabeknoten.](#)

**Anzahl der zu verwendenden Modelle.** Legt die maximale Anzahl der Modelle fest, die in dem vom Knoten erstellten Modell-Nugget aufgeführt werden sollen. Die Modelle mit dem höchsten Rang werden gemäß dem angegebenen Rangordnungskriterium aufgeführt. Durch Erhöhen dieser Obergrenze können Sie die Ergebnisse für einen größeren Anteil an Modellen vergleichen, allerdings kann dadurch auch die Verarbeitungsgeschwindigkeit sinken. Der höchste zulässige Wert ist 100.

**Bedeutsamkeit der Prädiktoren berechnen.** Bei Modellen, die zu einem angemessenen Maß an Bedeutsamkeit führen, können Sie ein Diagramm anzeigen, in dem die relative Wichtigkeit der einzelnen Prädiktoren bei der Modellschätzung angegeben wird. Normalerweise ist es sinnvoll, die Modellierungsbemühungen auf die wichtigsten Prädiktoren zu konzentrieren und diejenigen Prädiktoren zu verwerfen bzw. zu ignorieren, die die geringste Bedeutung haben. Beachten Sie, dass sich bei einigen Modellen der Zeitaufwand für die Berechnung durch die Bedeutsamkeit der Prädiktoren erhöhen kann. Außerdem wird diese Option nicht empfohlen, wenn Sie einfach einen allgemeinen Vergleich zwischen vielen verschiedenen Modellen wünschen. Diese Option ist von größerem Nutzen, wenn die Analyse bereits auf eine Handvoll Modelle eingeengt wurde, die detaillierter untersucht werden sollen. [Für weitere Informationen siehe Thema Bedeutsamkeit des Prädiktors in Kapitel 3 auf S. 54.](#)

**Modelle nicht behalten, wenn.** Dient zur Angabe von Schwellenwerten für Korrelation, relativen Fehler und Anzahl der verwendeten Felder. Modelle, die eines dieser Kriterien nicht erfüllen, werden verworfen und nicht im Zusammenfassungsbericht aufgeführt.

- **Korrelation ist kleiner als.** Die minimale Korrelation (als absoluter Wert), die gegeben sein muss, damit ein Modell in den Zusammenfassungsbericht aufgenommen wird.

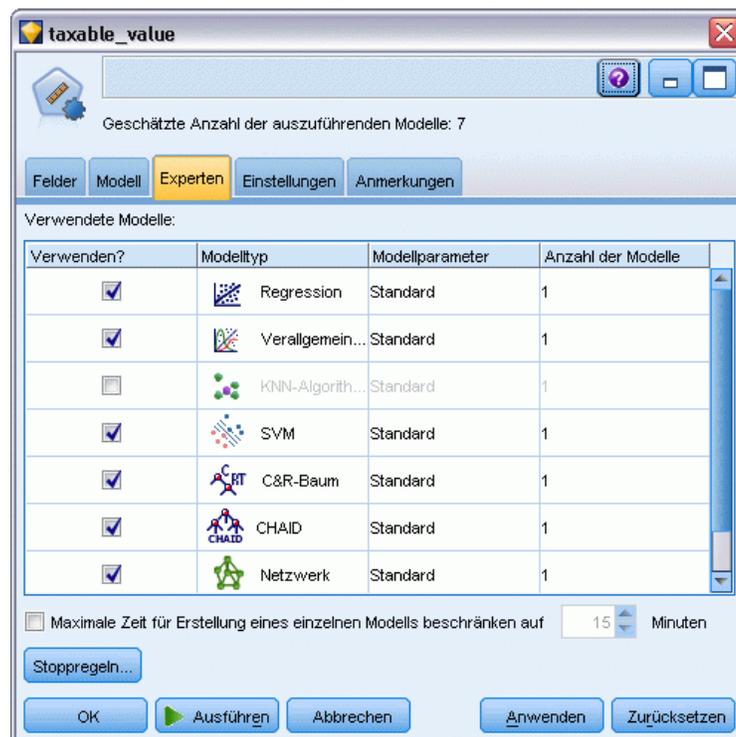
- **Anzahl der Felder ist größer als.** Die maximale Anzahl an Feldern, die von einem aufzunehmenden Modell verwendet werden können.
- **Relativer Fehler ist größer als.** Der maximale relative Fehler, der für ein aufzunehmendes Modell zulässig ist.

Optional können Sie den Knoten so konfigurieren, dass die Ausführung gestoppt wird, sobald erstmals ein Modell generiert wurde, das alle angegebenen Kriterien erfüllt. [Für weitere Informationen siehe Thema Knoten für die automatisierte Modellierung – Stoppregeln auf S. 95.](#)

### Knoten "Auto-Numerisch" - Expertenoptionen

Auf der Registerkarte "Experten" des Knotens "Auto-Numerisch" können Sie die zu verwendenden Algorithmen und Optionen auswählen und Stoppregeln angeben.

Abbildung 5-10  
Knoten "Auto-Numerisch": Registerkarte "Experten"



**Verwendete Modelle.** Wählen Sie anhand der Kontrollkästchen in der Spalte auf der linken Seite die Modelltypen (Algorithmen) aus, die in den Vergleich aufgenommen werden sollen. Je mehr Typen Sie auswählen, desto mehr Modelle werden erstellt und desto länger dauert die Verarbeitung.

**Modelltyp.** Listet die verfügbaren Algorithmen auf (siehe unten).

**Modellparameter.** Sie können für die einzelnen Modelltypen entweder die Standardeinstellungen verwenden oder mithilfe von Angebenen Optionen für den jeweiligen Modelltyp auswählen. Die einzelnen Optionen ähneln den in den separaten Modellierungskonten verfügbaren Optionen, mit

dem Unterschied, dass mehrere Optionen bzw. Kombinationen ausgewählt werden können. Beispiel: Beim Vergleich von neuronalen Netzwerkmodellen können Sie, anstatt eine der sechs Trainingsmethoden auszuwählen, alle sechs auswählen, um sechs Modelle in einem einzigen Durchgang zu trainieren.

**Anzahl der Modelle.** Listet die Anzahl der Modelle auf, die auf der Grundlage der aktuellen Einstellungen für die einzelnen Algorithmen erstellt wurden. Bei einer Kombination von Optionen kann die Anzahl der Modelle schnell recht groß werden. Daher wird dringend empfohlen, auf diesen Wert zu achten, insbesondere bei Verwendung großer Daten-Sets.

**Maximale Zeit für Erstellung eines einzelnen Modells beschränken.** (Nur für K-Means-, Kohonen-, TwoStep-, SVM-, KNN-, Bayes Net- und Entscheidungslisten-Modelle) Legt ein maximales Zeitlimit für jedes beliebige Modell fest. Wenn beispielsweise das Training für ein bestimmtes Modell aufgrund einer komplexen Interaktion unerwartet viel Zeit in Anspruch nimmt, wird durch diese Option vermieden, dass das Modell den gesamten Modellierungsdurchlauf aufhält.

### Unterstützte Algorithmen



Der Netzwerkknoten verwendet ein vereinfachtes Modell der Art und Weise, wie ein menschliches Gehirn Informationen verarbeitet. Es funktioniert, indem eine große Anzahl miteinander verbundener einfacher Verarbeitungseinheiten simuliert wird, die abstrakten Versionen von Neuronen ähnlich sind. Neuronale Netze sind leistungsstarke Mehrzweck-Schätzer, für deren Training und Anwendung nur sehr geringe statistische oder mathematische Kenntnisse erforderlich sind.



Der Knoten für Klassifizierungs- und Regressions-Bäume (C&RT-Bäume) erstellt einen Entscheidungsbaum, mit dem Sie zukünftige Beobachtungen vorhersagen oder klassifizieren können. Bei dieser Methode wird eine rekursive Partitionierung verwendet, um die Trainingsdatensätze in Segmente aufzuteilen. Dabei wird bei jedem Schritt die Unreinheit verringert und ein Knoten im Baum wird als "rein" betrachtet, wenn 100 % der Fälle in eine bestimmte Kategorie des Zielfelds fallen. Ziel- und Eingabefelder können numerische Bereiche oder kategorial (nominal, ordinal oder Flags) sein. Alle Aufteilungen sind binär (nur zwei Untergruppen). [Für weitere Informationen siehe Thema C&R-Baumknoten in Kapitel 6 auf S. 155.](#)



Der CHAID-Knoten erzeugt Entscheidungsbäume unter Verwendung von Chi-Quadrat-Statistiken zur Ermittlung optimaler Aufteilungen. Im Gegensatz zu den Knoten vom Typ "C&RT-Baum" und "QUEST" kann CHAID nichtbinäre Bäume generieren, d. h. Bäume mit Aufteilungen mit mehr als zwei Verzweigungen. Ziel- und Eingabefelder können in einem numerischen Bereich (stetig) oder kategorial sein. Exhaustive CHAID ist eine Änderung von CHAID, die noch gründlicher vorgeht, indem sie alle möglichen Aufteilungen untersucht, allerdings mehr Rechenzeit beansprucht. [Für weitere Informationen siehe Thema CHAID-Knoten in Kapitel 6 auf S. 156.](#)



Die lineare Regression ist ein statistisches Verfahren zur Zusammenfassung von Daten und die Erstellung von Prognosen durch Anpassung einer geraden Linie oder Fläche, mit der die Diskrepanzen zwischen den vorhergesagten und den tatsächlichen Ausgabewerten minimiert werden.



Das verallgemeinerte lineare Modell erweitert das allgemeine lineare Modell so, dass die abhängige Variable über eine angegebene Verknüpfungsfunktion in linearem Zusammenhang zu den Faktoren und Kovariaten steht. Außerdem ist es mit diesem Modell möglich, dass die abhängige Variable eine von der Normalverteilung abweichende Verteilung aufweist. Es deckt die Funktionen einer großen Bandbreite an Statistikmodellen ab, darunter lineare Regression, logistische Regression, loglineare Modelle für Häufigkeitsdaten und Überlebensmodelle mit Intervallzensurierung. [Für weitere Informationen siehe Thema GenLin-Knoten in Kapitel 10 auf S. 328.](#)



Der Knoten “ $k$ -Nächste Nachbarn” (KNN) verknüpft einen neuen Fall mit der Kategorie oder dem Wert der  $k$  Objekte, die ihm im Prädiktorraum am nächsten liegen, wobei  $k$  eine Ganzzahl ist. Ähnliche Fälle liegen nah beieinander und Fälle mit geringer Ähnlichkeit sind weit voneinander entfernt. [Für weitere Informationen siehe Thema KNN-Knoten in Kapitel 16 auf S. 482.](#)



Der Knoten “Support Vector Machine” (SVM) ermöglicht die Klassifizierung von Daten in eine von zwei Gruppen ohne Überanpassung. SVM eignet sich gut für umfangreiche Daten-Sets, beispielsweise solche mit einer großen Anzahl an Eingabefeldern. [Für weitere Informationen siehe Thema SVM-Knoten in Kapitel 15 auf S. 476.](#)

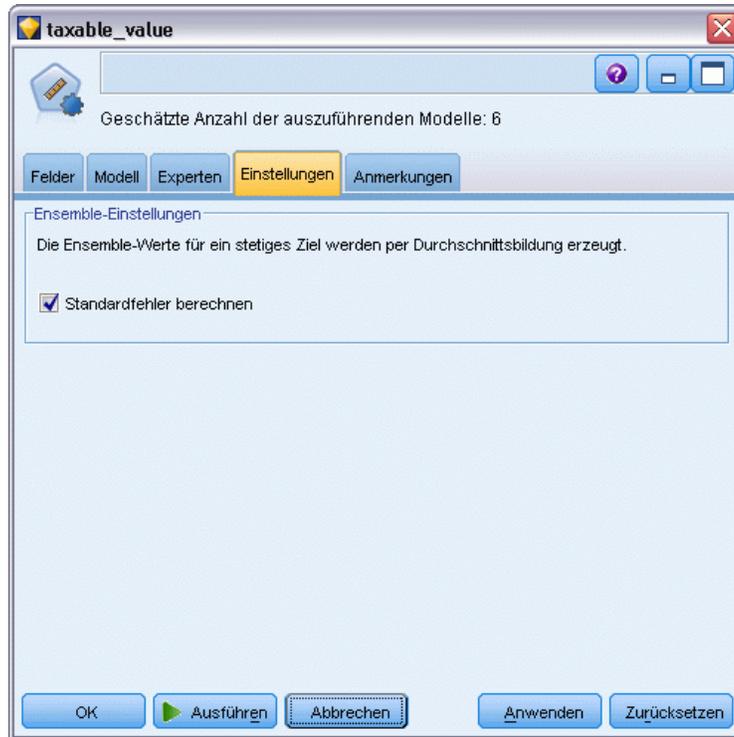


Bei linearen Regressionsmodellen wird ein stetiges Ziel auf der Basis linearer Beziehungen zwischen dem Ziel und einem oder mehreren Prädiktoren vorhergesagt. [Für weitere Informationen siehe Thema Lineare Modelle in Kapitel 10 auf S. 267.](#)

### **Knoten “Auto-Numerisch” - Einstellungsoptionen**

In der Registerkarte “Einstellungen” des Knotens “Auto-Numerisch” können Sie die Score-Zeit-Optionen konfigurieren, die für das Nugget verfügbar sind.

Abbildung 5-11  
Knoten "Auto-Numerisch": Registerkarte "Einstellungen"



**Standardfehler berechnen.** Für ein stetiges Ziel (numerischer Bereich) wird standardmäßig eine Standardfehlerberechnung durchgeführt, um den Unterschied zwischen den gemessenen oder geschätzten Werten und den wahren Werten zu berechnen sowie um zu zeigen, wie hoch die Übereinstimmung dieser Schätzungen war.

## ***Knoten "Autom. Cluster"***

Mit dem Knoten "Autom. Cluster" können Sie Clustering-Modelle, die Gruppen und Datensätze mit ähnlichen Merkmalen identifizieren, schätzen und vergleichen. Die Funktionsweise des Knotens gleicht der von anderen Knoten für automatisierte Modellierung, und Sie können in einem einzigen Modellierungsdurchgang mit mehreren Optionskombinationen experimentieren. Modelle können mithilfe grundlegender Messwerte für Filterung und Rangfolge der Nützlichkeit von Cluster-Modellen verglichen werden, um ein Maß auf der Basis der Wichtigkeit von bestimmten Feldern zu liefern.

Clustering-Modelle werden häufig verwendet, um Gruppen zu identifizieren, die als Eingabe für nachfolgende Analysen dienen können. Beispielsweise können Sie Kundengruppen auf der Basis von demografischen Merkmalen wie Einkommen oder von den Dienstleistungen, die sie in der Vergangenheit erworben haben, als Ziel nehmen. Dies ist ohne vorherige Kenntnis über die Gruppen und deren Merkmale möglich - Sie wissen u. U. gar nicht, wie viele Gruppen gesucht sind oder welche Funktionen für ihre Definition verwendet werden sollen. Clustering-Modelle werden häufig als nicht überwachte Lernmodelle bezeichnet, da sie kein Zielfeld benutzen und

keine bestimmte Vorhersage liefern, die sich als “wahr” oder “falsch” bewerten lässt. Der Wert eines Clustering-Modells wird durch die Möglichkeit bestimmt, interessante Gruppierungen in den Daten zu erfassen und sinnvolle Beschreibungen dieser Gruppierungen zu liefern. [Für weitere Informationen siehe Thema Cluster-Modelle in Kapitel 11 auf S. 354.](#)

Abbildung 5-12  
Autom. Cluster - Ergebnisse

Verwenden?	Diagramm	Modell	Erstellungszeit (Min.)	Silhouette	Anzahl der Cluster	Kleinsten Cluster (N)	Kleinsten Cluster (%)	Größten Cluster (N)	Größten Cluster (%)	Kleinsten/Größten	Wichtigkeit
<input checked="" type="checkbox"/>		K-...	< 1	0,229	5	137	12	372	32	0,368	1
<input type="checkbox"/>		Tw...	< 1	0,229	7	62	5	271	23	0,229	1
<input type="checkbox"/>		Ko...	< 1	0,206	9	6	0	285	25	0,021	1

**Anforderungen.** Eines oder mehrere Felder, die relevante Merkmale definieren. Cluster-Modelle verwenden Zielfelder nicht auf die gleiche Weise wie andere Modelle, da die keine spezifischen Vorhersagen treffen, die sich als wahr oder falsch bewerten lassen. Stattdessen werden sie verwendet, um Gruppen von Fällen zu identifizieren, die möglicherweise zusammenhängen. Beispielsweise können Sie anhand eines Cluster-Modells nicht vorhersagen, ob ein bestimmter Kunde positiv oder negativ auf ein Angebot reagiert. Sie können jedoch ein Cluster-Modell verwenden, um Kunden Gruppen basierend auf ihrer Tendenz zu einer bestimmten Reaktion zuzuweisen. Gewichts- und Häufigkeitsfelder werden nicht verwendet.

**Evaluierungsfelder.** Wenn kein Ziel verwendet wird, können Sie optional eines oder mehrere Evaluierungsfelder für den Vergleich von Modellen angeben. Der Nutzen eines Cluster-Modells kann dadurch bewertet werden, dass gemessen wird, wie gut (oder schlecht) die Cluster diese Felder differenzieren.

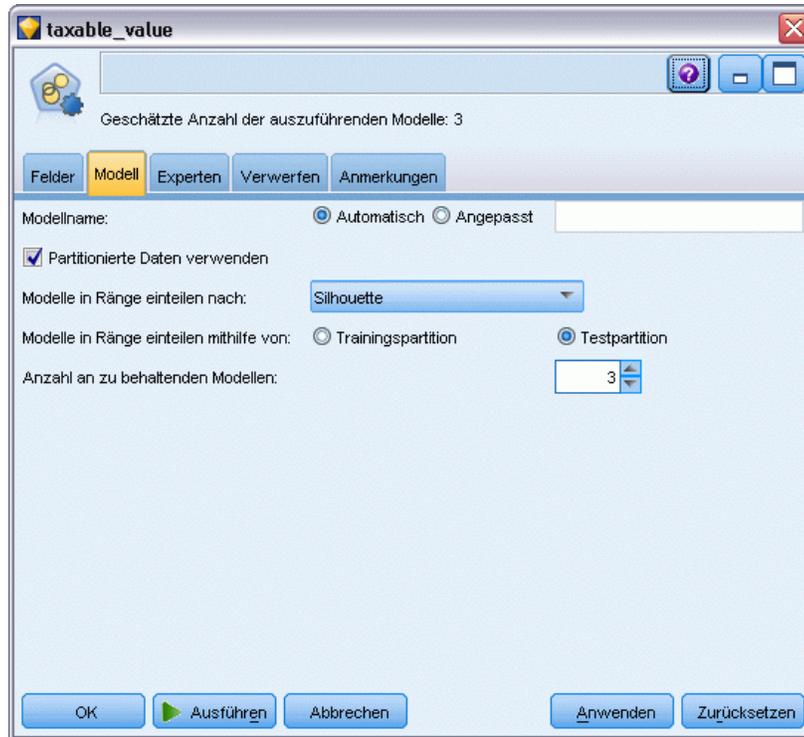
### Unterstützte Modelltypen

Unterstützte Modelltypen sind “TwoStep”, “K-Means” und “Kohonen”.

## Knoten “Autom. Cluster” - Modelloptionen

Auf der Registerkarte “Modell” des Knotens “Autom. Cluster” können Sie die Anzahl der zu speichernden Modelle sowie die zum Vergleich der Modelle verwendeten Kriterien angeben.

Abbildung 5-13  
Knoten "Autom. Cluster": Registerkarte "Modell"



**Modellname.** Sie können den Modellnamen automatisch basierend auf den Ziel- oder ID-Feldnamen (oder dem Modelltyp in Fällen, in denen kein solches Feld angegeben ist) generieren oder einen benutzerdefinierten Namen eingeben.

**Partitionierte Daten verwenden.** Wenn ein Partitionsfeld definiert ist, gewährleistet diese Option, dass nur Daten aus der Trainingspartition für die Modellerstellung verwendet werden. [Für weitere Informationen siehe Thema Partitionsknoten in Kapitel 4 in IBM SPSS Modeler 14.2- Quellen-, Prozess- und Ausgabeknoten.](#)

**Modelle in Ränge einteilen nach.** Legt fest, welche Kriterien für Vergleich und Rangordnung von Modellen verwendet werden sollen.

- **Silhouette.** Ein Index zur Messung von Cluster-Zusammenhalt und -Abgrenzung. Siehe *Rangenteilung mit Silhouetten* für weitere Informationen.
- **Anzahl der Cluster.** Die Anzahl der Cluster im Modell.
- **Größe des kleinsten Clusters.** Die kleinste Clustergröße.
- **Größe des größten Clusters.** Die größte Clustergröße.
- **Kleinstes/größtes Cluster.** Das Größenverhältnis zwischen dem kleinsten und dem größten Cluster.
- **Wichtigkeit.** Die Bedeutung des Felds Evaluierung in der Registerkarte Felder. Beachten Sie, dass dies nur berechnet werden kann, wenn das Feld Evaluierung angegeben wurde.

**Modelle in Ränge einteilen mithilfe von.** Wenn eine Partition verwendet wird, können Sie angeben, ob die Ränge auf dem Trainingsdaten-Set oder auf dem Testdaten-Set beruhen sollen. Bei großen Daten-Sets lässt sich die Leistungsfähigkeit durch die Verwendung einer Partition für ein erstes Screening der Modelle u. U. erheblich verbessern.

**Anzahl der zu behaltenden Modelle.** Legt die maximale Anzahl der Modelle fest, die in dem vom Knoten erstellten Nugget aufgeführt werden sollen. Die Modelle mit dem höchsten Rang werden gemäß dem angegebenen Rangordnungskriterium aufgeführt. Beachten Sie, dass eine Erhöhung dieses Grenzwerts die Leistungsgeschwindigkeit verringern kann. Der höchste zulässige Wert ist 100.

### ***Rangeinteilung mit Silhouetten***

Die standardmäßige Rangeinteilung mit Silhouetten verwendet einen Standardwert von 0, da ein Wert kleiner 0 (also ein negativer Wert) angibt, dass der durchschnittliche Abstand zwischen einem Fall und Punkten in seinem zugeordneten Cluster größer ist als der minimale durchschnittliche Abstand zu Punkten in einem anderen Cluster. Daher können Modelle mit einer negativen Silhouette einfach verworfen werden.

Die Rangeinteilung ist eigentlich ein modifizierter Silhouetten-Koeffizient, der die Konzepte von Cluster-Zusammenhalt (Favorisierung von Modellen mit eng zusammengehörenden Clustern) und Cluster-Abgrenzung (Favorisierung von Modellen mit stark separierten Clustern) kombiniert. Der durchschnittliche Silhouetten-Koeffizient ist einfach der Durchschnitt für alle Fälle der folgenden Berechnung für jeden Einzelfall:

$$(B - A) / \max.(A, B)$$

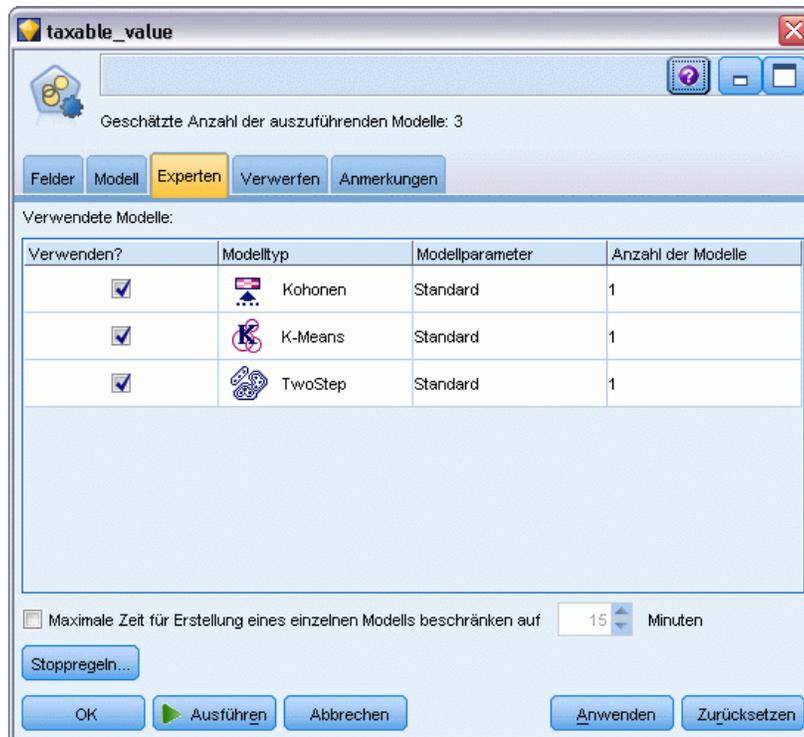
$A$  ist die Entfernung zwischen Fall und Zentroid des Clusters, zu dem der Fall gehört und  $B$  ist der minimale Abstand zwischen Fall und Zentroid jedes anderen Clusters.

Der Silhouetten-Koeffizient (und sein Durchschnittswert) liegen zwischen -1 (stellvertretend für ein sehr schlechtes Modell) und 1 (stellvertretend für ein ausgezeichnetes Modell). Der Durchschnitt kann auf der Ebene aller Fälle (führt zu einer Gesamt-Silhouette) oder auf Cluster-Ebene (führt zu einer Cluster-Silhouette) durchgeführt werden. Entfernungen können mithilfe euklidischer Entfernungen berechnet werden.

## ***Knoten "Autom. Cluster" - Expertenoptionen***

Auf der Registerkarte "Experten" des Knotens "Autom. Cluster" können Sie eine Partition anwenden (sofern verfügbar), die zu verwendenden Algorithmen auswählen und Stoppregeln angeben.

Abbildung 5-14  
Knoten "Autom. Cluster": Registerkarte "Experten"



**Verwendete Modelle.** Wählen Sie anhand der Kontrollkästchen in der Spalte auf der linken Seite die Modelltypen (Algorithmen) aus, die in den Vergleich aufgenommen werden sollen. Je mehr Typen Sie auswählen, desto mehr Modelle werden erstellt und desto länger dauert die Verarbeitung.

**Modelltyp.** Listet die verfügbaren Algorithmen auf (siehe unten).

**Modellparameter.** Sie können für die einzelnen Modelltypen entweder die Standardeinstellungen verwenden oder mithilfe von Angebenen Optionen für den jeweiligen Modelltyp auswählen. Die einzelnen Optionen ähneln den in den separaten Modellierungskonten verfügbaren Optionen, mit dem Unterschied, dass mehrere Optionen bzw. Kombinationen ausgewählt werden können. Beispiel: Beim Vergleich von neuronalen Netzwerkmodellen können Sie, anstatt eine der sechs Trainingsmethoden auszuwählen, alle sechs auswählen, um sechs Modelle in einem einzigen Durchgang zu trainieren.

**Anzahl der Modelle.** Listet die Anzahl der Modelle auf, die auf der Grundlage der aktuellen Einstellungen für die einzelnen Algorithmen erstellt wurden. Bei einer Kombination von Optionen kann die Anzahl der Modelle schnell recht groß werden. Daher wird dringend empfohlen, auf diesen Wert zu achten, insbesondere bei Verwendung großer Daten-Sets.

**Maximale Zeit für Erstellung eines einzelnen Modells beschränken.** (Nur für K-Means-, Kohonen-, TwoStep-, SVM-, KNN-, Bayes Net- und Entscheidungslisten-Modelle) Legt ein maximales Zeitlimit für jedes beliebige Modell fest. Wenn beispielsweise das Training für ein bestimmtes

Modell aufgrund einer komplexen Interaktion unerwartet viel Zeit in Anspruch nimmt, wird durch diese Option vermieden, dass das Modell den gesamten Modellierungsdurchlauf aufhält.

### **Unterstützte Algorithmen**



Der K-Means-Knoten teilt das Daten-Set in unterschiedliche Gruppen (oder Cluster) auf. Bei diesem Verfahren wird eine festgelegte Anzahl von Clustern definiert, den Clustern werden iterativ Datensätze zugewiesen und die Cluster-Zentren werden angepasst, bis eine weitere Verfeinerung keine wesentliche Verbesserung des Modells mehr darstellen würde. Statt zu versuchen, ein Ergebnis vorherzusagen, versucht K-Means mithilfe eines als "nicht überwachtes Lernen" bezeichneten Verfahrens Muster im Set der Eingabefelder zu entdecken. [Für weitere Informationen siehe Thema K-Means-Knoten in Kapitel 11 auf S. 362.](#)



Der Kohonen-Knoten erstellt eine Art von neuronalem Netzwerk, das verwendet werden kann, um ein Clustering der Datenmenge in einzelne Gruppen vorzunehmen. Wenn das Netz voll trainiert ist, sollten ähnliche Datensätze auf der Ausgabekarte eng nebeneinander stehen, während Datensätze, die sich unterscheiden, weit voneinander entfernt sein sollten. Die Zahl der von jeder Einheit im Modell-Nugget erfassten Beobachtungen gibt Aufschluss über die starken Einheiten. Dadurch wird ein Eindruck von der ungefähren Zahl der Cluster vermittelt. [Für weitere Informationen siehe Thema Kohonen-Knoten in Kapitel 11 auf S. 355.](#)

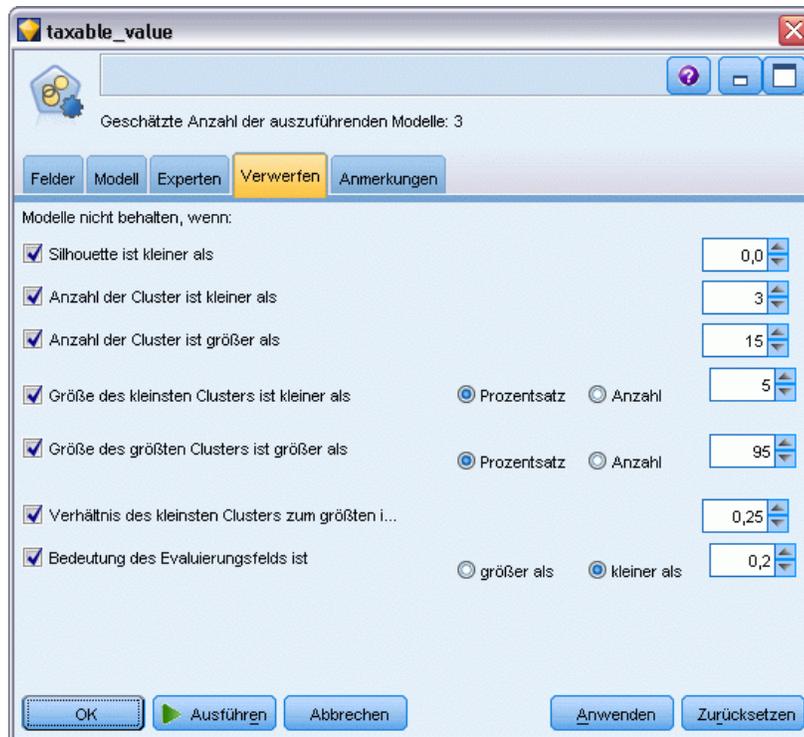


Der TwoStep-Knoten verwendet eine aus zwei Schritten bestehende Clusterbildungsmethode. Im ersten Schritt wird ein einzelner Durchlauf durch die Daten vorgenommen, bei dem die Eingangsrohdaten zu einem verwaltbaren Set von Unterclustern komprimiert werden. Im zweiten Schritt werden die Untercluster mithilfe einer hierarchischen Methode zur Clusterbildung nach und nach in immer größere Cluster zusammengeführt. TwoStep hat den Vorteil, dass die optimale Anzahl an Clustern für die Trainingsdaten automatisch geschätzt wird. Mit dem Verfahren können gemischte Feldtypen und große Daten-Sets effizient verarbeitet werden. [Für weitere Informationen siehe Thema TwoStep-Cluster-Knoten in Kapitel 11 auf S. 367.](#)

### **Knoten "Autom. Cluster" - Optionen für Verwerfen**

Auf der Registerkarte "Verwerfen" des Knotens "Autom. Cluster" können Sie automatisch Modelle verwerfen, die bestimmte Kriterien nicht erfüllen. Diese Modelle werden nicht auf dem Modell-Nugget aufgeführt.

Abbildung 5-15  
Knoten "Autom. Cluster": Registerkarte "Verwerfen"



Sie können den Silhouetten-Mindestwert, Clusterzahlen, Clustergrößen und Wichtigkeit des Evaluierungsfelds für das Modell angeben. Silhouette sowie Anzahl und Größe von Clustern richten sich nach den Angaben im Modellierungsknoten. [Für weitere Informationen siehe Thema Knoten "Autom. Cluster" - Modelloptionen auf S. 115.](#)

Optional können Sie den Knoten so konfigurieren, dass die Ausführung gestoppt wird, sobald erstmals ein Modell generiert wurde, das alle angegebenen Kriterien erfüllt. [Für weitere Informationen siehe Thema Knoten für die automatisierte Modellierung – StoppregeIn auf S. 95.](#)

## ***Nugget für automatisierte Modellierung***

Beim Ausführen eines Knotens für automatisierte Modellierung schätzt der Knoten in Frage kommende Modelle für alle möglichen Einstellungskombinationen, stuft jedes in Frage kommende Modell auf der Basis des angegebenen Werts und speichert die geeignetsten Modelle in einem zusammengesetzten Nugget für automatisierte Modellierung. Dieses Modell-Nugget enthält ein Set aus einem oder mehreren vom Knoten generierten Modellen, die individuell durchgesehen oder zur Verwendung beim Scoring ausgewählt werden können. Für jedes Modell werden Modelltyp und Erstellungszeit zusammen mit der Anzahl anderer Maße passend zum Modelltyp aufgeführt. Sie können die Tabelle nach einer beliebigen Spalte sortieren, um rasch die interessantesten Modelle ermitteln zu können.

Abbildung 5-16  
Auto-Numerisch - Ergebnisse

Verwenden?	Diagramm	Modell	Erstellungszeit (Min.)	Korrelation	Anzahl der verwendeten	Relativer Fehler
<input checked="" type="checkbox"/>		Generaliz...	< 1	0,915	7	0,162
<input checked="" type="checkbox"/>		Regressi...	< 1	0,9	5	0,19
<input checked="" type="checkbox"/>		CHAID Tr...	< 1	0,892	5	0,204

- Zum Durchsehen von einem der einzelnen Modell-Nuggets doppelklicken Sie auf das Nugget-Symbol. Sie können dann einen Modellknoten für dieses Modell im Stream-Zeichenbereich generieren oder eine Kopie des Modell-Nuggets in der Modellalette.
- Miniaturansichten von Diagrammen erlauben eine rasche visuelle Einschätzung für jeden Modelltyp, wie unten zusammengefasst. Durch Doppelklicken auf eine Miniaturansicht können Sie ein Diagramm in voller Größe generieren. Das Diagramm in voller Größe zeigt bis zu 1.000 Punkten und beruht auf einer Stichprobe, wenn das Daten-Set mehr Punkte enthält. (Bei Streudiagrammen wird das Diagramm jeweils, wenn es angezeigt wird, neu generiert, sodass etwaige Änderungen an den Daten oberhalb des Knotens – wie beispielsweise die Aktualisierung einer Zufallsstichprobe oder einer Partition, wenn Startwert für Zufallsgenerator festlegen nicht ausgewählt ist – stets berücksichtigt werden, wenn das Streudiagramm erneut gezeichnet wird.)
- Mithilfe der Symbolleiste können Sie bestimmte Spalten auf der Registerkarte “Modell” ein- bzw. ausblenden oder die Spalte ändern, die für die Sortierung der Tabelle verwendet wird. (Außerdem können Sie die Sortierung durch Klicken auf die Spaltenüberschriften ändern.)
- Mit der Schaltfläche “Löschen” können Sie unbenutzte Modelle permanent entfernen.
- Um die Anordnung der Spalten zu ändern, klicken Sie auf eine Spaltenüberschrift und ziehen Sie die Spalte an die gewünschte Position.
- Wenn eine Partition gerade verwendet wird, können Sie die Ergebnisse für die Trainings- bzw. Testpartition anzeigen lassen. [Für weitere Informationen siehe Thema Partitionsknoten in Kapitel 4 in IBM SPSS Modeler 14.2- Quellen-, Prozess- und Ausgabeknoten.](#)

Die jeweils verwendeten Spalten hängen vom Typ der zu vergleichenden Modelle ab, wie unten angegeben.

### ***Binäre Ziele***

- Bei binären Modellen zeigt das Miniaturdiagramm die Verteilung der tatsächlichen Werte, überlagert mit den vorhergesagten Werten, und bietet so einen schnellen Überblick darüber, wie viele Datensätze in den einzelnen Kategorien korrekt vorhergesagt wurden.
- Rangordnungskriterien entsprechen den Optionen im Knoten “Automatischer Klassifizierer”. [Für weitere Informationen siehe Thema Knoten “Automatischer Klassifizierer” - Modelloptionen auf S. 98.](#)
- Für den maximalen Profit wird außerdem das Perzentil angegeben, in dem dieses Maximum auftritt.
- Beim kumulativen Lift können Sie das ausgewählte Perzentil mithilfe der Symbolleiste ändern.

### ***Nominale Ziele***

- Bei nominale (Set-) Modellen zeigt das Miniaturdiagramm die Verteilung der tatsächlichen Werte, überlagert mit den vorhergesagten Werten, und bietet so einen schnellen Überblick darüber, wie viele Datensätze in den einzelnen Kategorien korrekt vorhergesagt wurden.
- Rangordnungskriterien entsprechen den Optionen im Knoten “Automatischer Klassifizierer”. [Für weitere Informationen siehe Thema Knoten “Automatischer Klassifizierer” - Modelloptionen auf S. 98.](#)

### ***Stetige Ziele***

- Bei stetigen Modelle (numerischer Bereich) bietet das Diagramm für jedes Modell eine grafische Darstellung der vorhergesagten gegenüber den beobachteten Werten und ermöglicht einen schnellen Überblick über die Korrelation zwischen diesen Werten. Bei einem guten Modell sollten die Punkte tendenziell entlang der Diagonalen gruppiert und nicht zufällig im ganzen Diagramm verstreut sein.
- Rangordnungskriterien entsprechen den Optionen im Knoten “Auto-Numerisch”. [Für weitere Informationen siehe Thema Knoten “Auto-Numerisch” - Modelloptionen auf S. 108.](#)

### ***Cluster-Ziele.***

- Bei Cluster-Modellen bietet das Diagramm für jedes Modell eine grafische Darstellung von Zahlen gegenüber Clustern und ermöglicht einen schnellen Überblick über die Cluster-Verteilung.
- Rangordnungskriterien entsprechen den Optionen im Knoten “Autom. Cluster”. [Für weitere Informationen siehe Thema Knoten “Autom. Cluster” - Modelloptionen auf S. 115.](#)

### ***Auswählen von Modellen zum Scoring***

In der Spalte Verwenden? können Sie Modelle für das Scoring auswählen.

- Nur für binäre, nominale und numerische Ziele können Sie mehrere Scoring-Modelle auswählen und die Scores in einem einzigen Modell-Nugget kombinieren. Durch die Kombination der Vorhersagen aus mehreren Modellen lassen sich Begrenzungen, die einzelne

Modelle aufweisen, vermeiden. Dadurch kann häufig eine höhere Gesamtgenauigkeit erreicht werden als mit einem der Modelle allein.

- Für Cluster-Modelle kann nur jeweils ein Scoring-Modell ausgewählt werden. Standardmäßig wird das Modell mit dem höchsten Rang zuerst ausgewählt.

## **Generieren von Knoten und Modellen**

Sie können eine Kopie des zusammengesetzten Nuggets für automatisierte Modellierung oder des Knotens für automatisierte Modellierung, aus dem das Nugget erstellt wurde, generieren. Dies kann beispielsweise dann hilfreich sein, wenn Sie nicht über den Original-Stream verfügen, aus dem das Nugget für automatisierte Modellierung erstellt wurde. Für alle im Nugget für automatisierte Modellierung aufgeführten Modelle kann ein Modell-Nugget oder ein Modellierungsknoten generiert werden.

### ***Nugget für automatisierte Modellierung***

- ▶ Wählen Sie im Modell “Generieren” die Option Modelle in Palette, um der Modellpalette das Nugget für automatisierte Modellierung hinzuzufügen. Das generierte Modell kann gespeichert oder in der vorliegenden Form verwendet werden, ohne dass der Stream erneut ausgeführt werden muss.
- ▶ Alternativ können Sie im Menü “Generieren” die Option Modellierungsknoten generieren auswählen, um dem Stream-Zeichenbereich den Modellierungsknoten hinzuzufügen. Anhand dieses Knotens können Sie die ausgewählten Modelle erneut schätzen, ohne den gesamten Modellierungsdurchlauf wiederholen zu müssen.

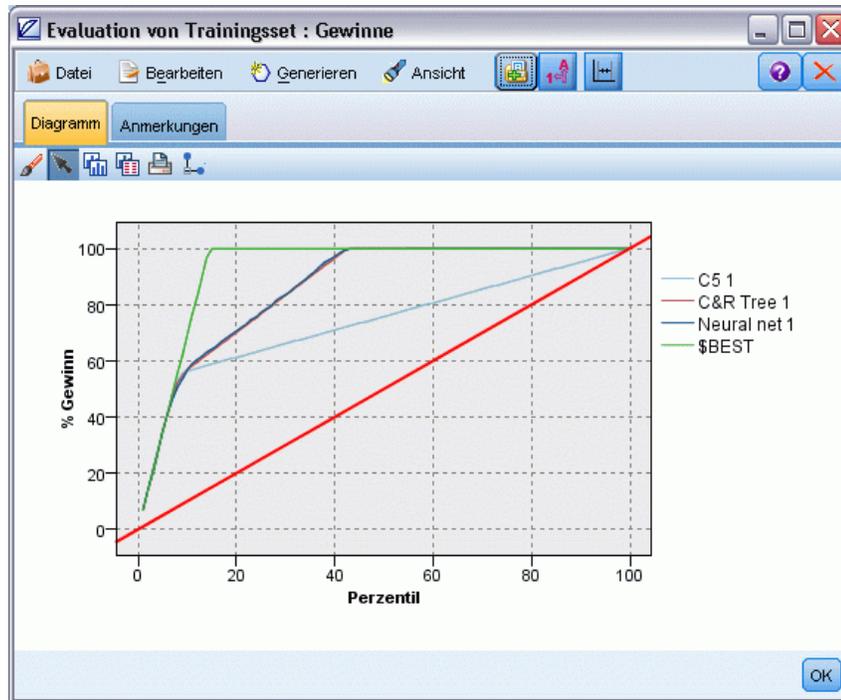
### ***Einzelnes Modellierungs-Nugget***

- ▶ Doppelklicken Sie im Menü Modell auf das benötigte Nugget. In einem neuen Dialogfeld wird eine Kopie dieses Nuggets erstellt.
- ▶ Wählen Sie im Menü “Generieren” des neuen Dialogfelds die Option Modelle in Palette, um der Modellpalette das einzelne Nugget hinzuzufügen.
- ▶ Alternativ können Sie im Menü “Generieren” des neuen Dialogfelds die Option Modellierungsknoten generieren auswählen, um dem Stream-Zeichenbereich den einzelnen Modellierungsknoten hinzuzufügen.

## **Generieren von Evaluationsdiagrammen**

Bei binären Modellen können Sie Evaluationsdiagramme generieren, die eine visuelle Möglichkeit zur Bewertung und zum Vergleich der einzelnen Modelle bieten. Evaluationsdiagramme stehen nicht für Modelle zur Verfügung, die mithilfe der Knoten vom Typ “Auto-Numerisch” oder “Autom. Cluster” generiert wurden. [Für weitere Informationen siehe Thema Evaluationsknoten in Kapitel 5 in IBM SPSS Modeler 14.2- Quellen-, Prozess- und Ausgabeknoten.](#)

Abbildung 5-17  
Trefferdiagramm (kumulativ) mit bester Linie und Basis



- ▶ Wählen Sie im Nugget für automatisierte Modellierung des Knotens “Automatischer Klassifizierer” in der Spalte *Verwenden?* die Modelle aus, die Sie auswerten möchten.
- ▶ Wählen Sie im Menü “Generieren” die Option Evaluationsdiagramm(e).

Abbildung 5-18  
Generieren von Evaluationsdiagrammen



- ▶ Wählen Sie den gewünschten Diagrammtyp und die anderen gewünschten Optionen. Für weitere Informationen siehe Thema Evaluation – Registerkarte “Plot” in Kapitel 5 in *IBM SPSS Modeler 14.2- Quellen-, Prozess- und Ausgabeknoten*.

## Evaluationsdiagramme

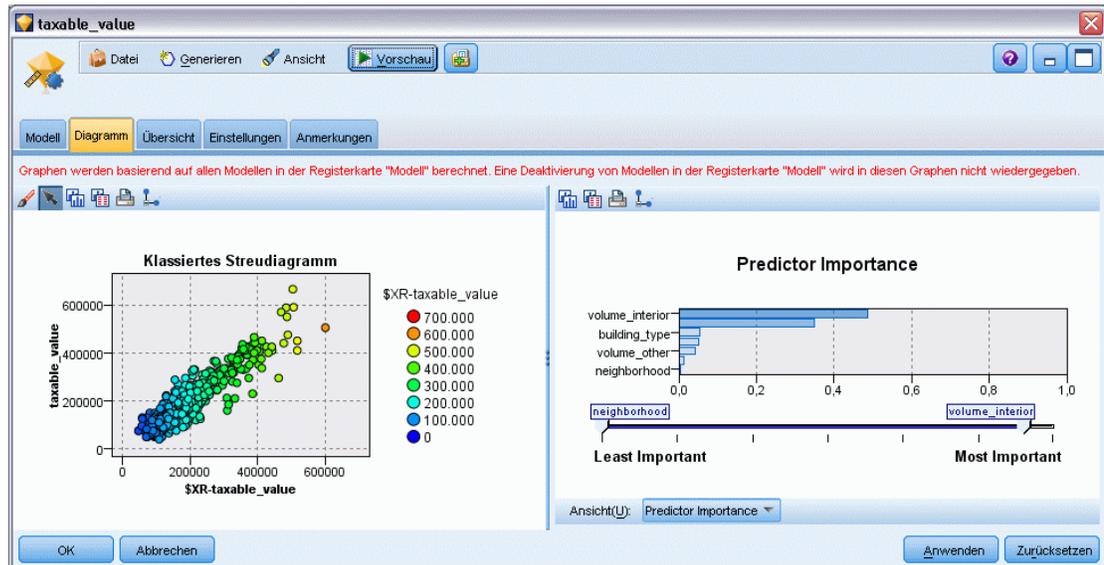
Auf der Registerkarte “Modell” des Nuggets für automatisierte Modellierung können Sie einen Drill-Down durchführen, um einzelne Diagramme für jedes der angezeigten Modelle anzuzeigen. Für die Nuggets “Automatischer Klassifizierer” und “Auto-Numerisch” werden auf

der Registerkarte “Diagramm” sowohl ein Diagramm als auch die Bedeutsamkeit der Prädiktoren angezeigt, die die Ergebnisse aller Modelle zusammen widerspiegeln. [Für weitere Informationen siehe Thema Bedeutsamkeit des Prädiktors in Kapitel 3 auf S. 54.](#)

Für “Automatischer Klassifizierer” wird ein Verteilungsdiagramm angezeigt, wohingegen für “Auto-Numerisch” ein Multidiagramm (auch Streudiagramm genannt) angezeigt wird. [Für weitere Informationen siehe Thema Häufig verwendete Funktionen von Diagrammknoten in Kapitel 5 in IBM SPSS Modeler 14.2- Quellen-, Prozess- und Ausgabeknoten.](#)

Abbildung 5-19

Auto-Numerisch - Multidiagramm für die Ensemble-Modelle im Nugget für automatisierte Modellierung

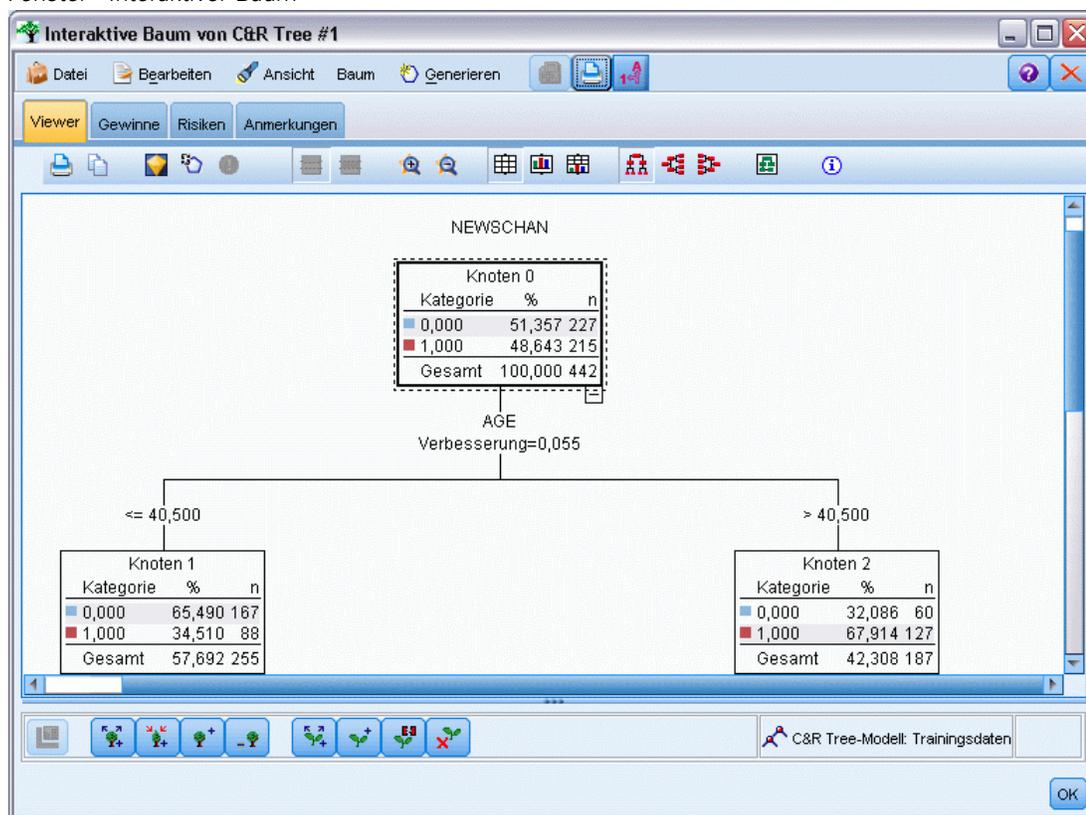


# Decision Trees (Entscheidungsbäume)

## Entscheidungsbaum-Modelle

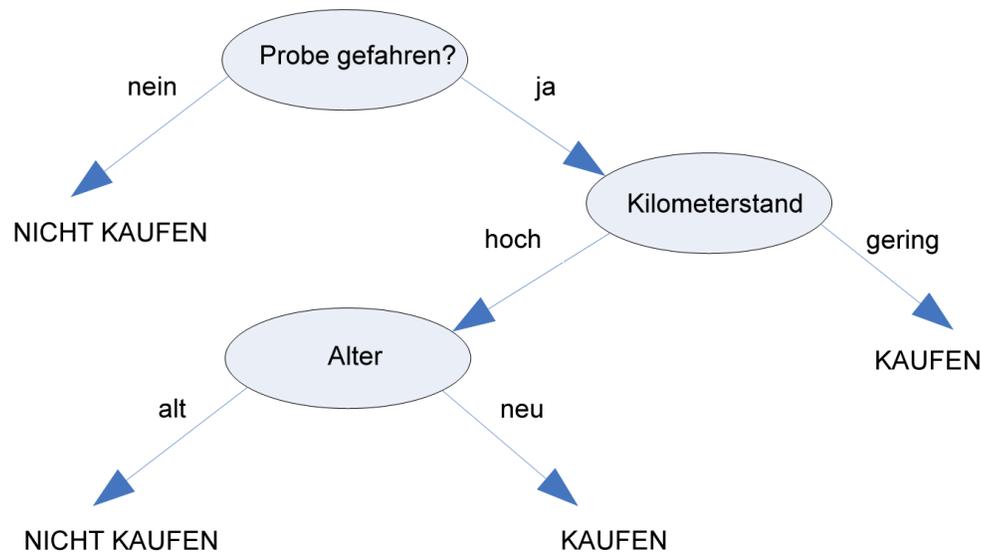
Mithilfe von Entscheidungsbaummodellen werden Klassifizierungssysteme entwickelt, die zukünftige Beobachtungen auf der Grundlage eines Satzes von Entscheidungsregeln vorhersagen oder klassifizieren. Wenn die Daten in Klassen aufgeteilt sind, die Sie interessieren (z. B. Darlehen mit hohem Risiko im Gegensatz zu Darlehen mit niedrigem Risiko, Abonnenten gegenüber Personen ohne Abonnement Wähler im Gegensatz zu Nichtwählern oder Bakterienarten), können Sie mit diesen Daten Regeln erstellen, die Sie zur Klassifizierung alter oder neuer Fälle mit maximaler Genauigkeit verwenden können. So können Sie z. B. einen Baum erstellen, der das Kreditrisiko oder die Kaufabsicht basierend auf Alter und anderen Faktoren klassifiziert.

Abbildung 6-1  
Fenster "Interaktiver Baum"



Dieser Ansatz, manchmal bekannt als **Regelinduktion**, hat mehrere Vorteile. Zunächst wird die Argumentationskette hinter dem Modell deutlich, wenn Sie durch die Struktur blättern. Dies steht im Gegensatz zu anderen "Blackbox"-Modellierungstechniken, bei denen die interne Logik nicht so leicht zu durchschauen ist.

Abbildung 6-2  
Einfacher Entscheidungsbaum für einen Autokauf



Zweitens: Der Prozess berücksichtigt in seiner Regel automatisch nur die Attribute, die beim Fällen einer Entscheidung tatsächlich eine Rolle spielen. Attribute, die nicht zur Genauigkeit des Baums beitragen, werden ignoriert. Dies kann zu sehr hilfreichen Informationen über die Daten führen und kann dazu verwendet werden, die Daten auf die relevanten Felder zu reduzieren, bevor ein anderes Verfahren zum Maschinenlernen trainiert wird, z. B. ein neuronales Netz.

Entscheidungsbaummodell-Nuggets können in eine Sammlung von Wenn-dann-Regeln (eine **Regelmenge**) umgewandelt werden, die die Informationen in vielen Fällen in einer verständlicheren Form darstellen. Die Entscheidungsbaumdarstellung ist nützlich, wenn Sie sehen möchten, wie die Attribute in den Daten die Gesamtheit in Untergruppen **teilen** oder **aufteilen**, die für das Problem relevant sind. Die Regelmengendarstellung ist dann nützlich, wenn Sie sehen möchten, in welchem Zusammenhang bestimmte Elementgruppen mit einer bestimmten Schlussfolgerung stehen. Die folgende Regel präsentiert uns z. B. ein **Profil** für eine Gruppe von Fahrzeugen, die einen Kauf wert sind:

```

IF tested = 'yes'
AND mileage = 'low'
THEN -> 'BUY'.
  
```

### **Baumerstellungsalgorithmen**

Vier Algorithmen stehen für die Durchführung der Klassifizierungs- und Segmentierungsanalyse zur Verfügung. Diese Algorithmen führen im Grunde alle dieselben Operationen aus. Sie prüfen alle Felder Ihres Daten-Sets, um diejenigen herauszufinden, die die beste Klassifizierung oder Vorhersage liefern, indem sie die Daten in Untergruppen aufteilen. Der Vorgang wird rekursiv angewendet, wobei die Untergruppen in immer kleinere Einheiten aufgeteilt werden, bis der Baum erstellt ist (wie von bestimmten Grenzkriterien definiert). Die bei der Baumerstellung verwendeten Ziel- und Eingabefelder können je nach verwendetem Algorithmus stetig (numerischer Bereich)

oder kategorial sein. Wenn ein stetiges Ziel verwendet wird, wird ein Regressionsbaum generiert; wenn ein kategoriales Ziel verwendet wird, wird ein Klassifizierungsbaum generiert.



Der Knoten für Klassifizierungs- und Regressions-Bäume (C&RT-Bäume) erstellt einen Entscheidungsbaum, mit dem Sie zukünftige Beobachtungen vorherhersagen oder klassifizieren können. Bei dieser Methode wird eine rekursive Partitionierung verwendet, um die Trainingsdatensätze in Segmente aufzuteilen. Dabei wird bei jedem Schritt die Unreinheit verringert und ein Knoten im Baum wird als “rein” betrachtet, wenn 100 % der Fälle in eine bestimmte Kategorie des Zielfelds fallen. Ziel- und Eingabefelder können numerische Bereiche oder kategorial (nominal, ordinal oder Flags) sein. Alle Aufteilungen sind binär (nur zwei Untergruppen). [Für weitere Informationen siehe Thema C&R-Baumknoten auf S. 155.](#)



Der CHAID-Knoten erzeugt Entscheidungsbäume unter Verwendung von Chi-Quadrat-Statistiken zur Ermittlung optimaler Aufteilungen. Im Gegensatz zu den Knoten vom Typ “C&RT-Baum” und “QUEST” kann CHAID nichtbinäre Bäume generieren, d. h. Bäume mit Aufteilungen mit mehr als zwei Verzweigungen. Ziel- und Eingabefelder können in einem numerischen Bereich (stetig) oder kategorial sein. Exhaustive CHAID ist eine Änderung von CHAID, die noch gründlicher vorgeht, indem sie alle möglichen Aufteilungen untersucht, allerdings mehr Rechenzeit beansprucht. [Für weitere Informationen siehe Thema CHAID-Knoten auf S. 156.](#)



Der QUEST-Knoten bietet eine binäre Klassifizierungsmethode zum Erstellen von Entscheidungsbaummodellen, die dazu dient, die für große C&R-Baum-Analysen erforderliche Verarbeitungszeit zu verkürzen. Gleichzeitig soll die in den Klassifizierungsbaummodellen festgestellte Tendenz verringert werden, die darin besteht, dass Eingaben bevorzugt werden, die mehr Aufteilungen erlauben. Eingabefelder können stetig (numerische Bereiche) sein, das Zielfeld muss aber kategorial sein. Alle Aufteilungen sind binär. [Für weitere Informationen siehe Thema QUEST-Knoten auf S. 157.](#)



Der C5.0-Knoten erstellt entweder einen Entscheidungsbaum oder ein Regel-Set. Das Modell teilt die Stichprobe auf der Basis des Felds auf, das auf der jeweiligen Ebene den maximalen Informationsgewinn liefert. Das Zielfeld muss kategorial sein. Es sind mehrere Aufteilungen in mehr als zwei Untergruppen zulässig. [Für weitere Informationen siehe Thema C5.0-Knoten auf S. 177.](#)

### **Allgemeine Verwendung der baumbasierten Analyse**

Im Folgenden werden einige allgemeine Anwendungsbereiche der baumbasierten Analyse erläutert:

**Segmentierung.** Ermitteln Sie die Personen, die wahrscheinlich Mitglieder einer bestimmten Klasse sind.

**Schichtung.** Weisen Sie die Fälle einer von mehreren Kategorien zu, beispielsweise Gruppen mit hohem, mittlerem und niedrigem Risiko.

**Prognose.** Erstellen Sie Regeln und verwenden Sie diese, um zukünftige Ereignisse vorherzusagen. Vorhersage kann auch den Versuch bezeichnen, Vorhersageattribute Werten einer kontinuierlichen Variablen zuzuordnen.

**Datenreduktion und Variablenüberwachung.** Wählen Sie eine sinnvolle Untergruppe von Prädiktoren aus einer großen Menge Variablen aus, um ein formales parametrisches Modell zu erstellen.

**Interaktionsidentifizierung.** Ermitteln Sie Beziehungen, die nur bestimmten Untergruppen angehören und geben Sie diese in einem formalen parametrischen Modell an.

**Kategoriezusammenführung und Einteilen von kontinuierlichen Variablen.** Kodieren Sie die Gruppenprädiktorkategorien und kontinuierlichen Variablen mit minimalem Informationsverlust um.

## ***Interactive Tree Builder***

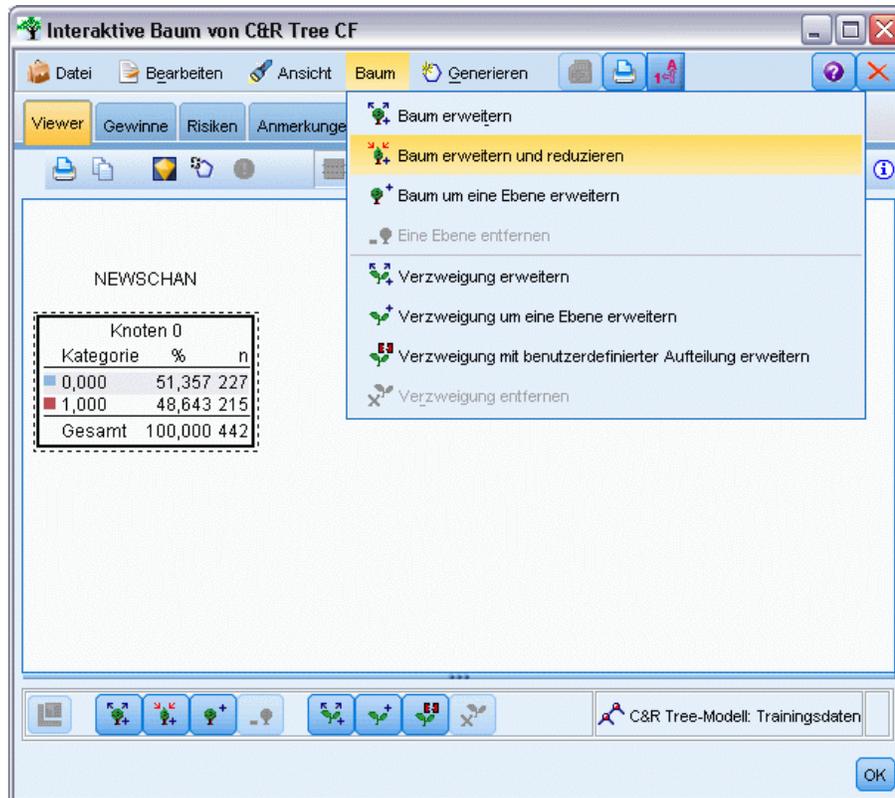
Sie können ein Baummodell entweder automatisch erstellen, indem Sie den Algorithmus auf jeder Ebene den besten Split auswählen lassen, oder Sie können den interaktiven Tree Builder nutzen und den Baum vor dem Speichern des Modell-Nuggets auf der Grundlage Ihres Geschäftswissens verfeinern oder vereinfachen.

- ▶ Erstellen Sie einen Stream und fügen Sie einen der drei Entscheidungsbaumknoten “C&R-Baum”, “CHAID” oder “QUEST” hinzu.

*Hinweis:* Für C5.0-Bäume wird keine interaktive Baumerstellung unterstützt.

- ▶ Öffnen Sie den Knoten, wählen Sie auf der Registerkarte “Felder” die Ziel- und Prädiktorfelder und geben Sie bei Bedarf zusätzliche Modelloptionen an. Spezifische Anleitungen finden Sie in der Dokumentation zu den einzelnen Baumerstellungsknoten.
- ▶ Wählen Sie auf dem Panel “Ziel” der Registerkarte “Erstellungsoptionen” den Befehl interaktive Sitzung starten.
- ▶ Klicken Sie auf Ausführen, um den Tree Builder zu starten.

Abbildung 6-3  
Interaktiver Tree Builder – Fenster



Der aktuelle Baum wird ab dem Stammknoten angezeigt. Sie können den Baum Ebene für Ebene bearbeiten und reduzieren und, bevor Sie eines oder mehrere Modelle erstellen, auf Gewinne, Risiken und zugehörige Informationen zugreifen.

### Kommentare

- Für C&RT-, CHAID- und QUEST-Knoten müssen alle im Modell verwendeten ordinalen Felder numerisch (und nicht als Zeichenkette) gespeichert sein. Im Bedarfsfall können Sie die Felder mit dem Umkodierungsknoten konvertieren. [Für weitere Informationen siehe Thema Umkodierungsknoten in Kapitel 4 in IBM SPSS Modeler 14.2- Quellen-, Prozess- und Ausgabeknoten.](#)
- Optional können Sie mit einem Partitionsfeld die Daten in Trainings- und Test-Stichproben trennen. [Für weitere Informationen siehe Thema Partitionsknoten in Kapitel 4 in IBM SPSS Modeler 14.2- Quellen-, Prozess- und Ausgabeknoten.](#)
- Sie können ein Modell mit Tree Builder oder, wie andere IBM® SPSS® Modeler-Modelle, direkt aus dem Modellierungsknoten generieren. [Für weitere Informationen siehe Thema Direktes Erstellen eines Baummodells auf S. 152.](#)

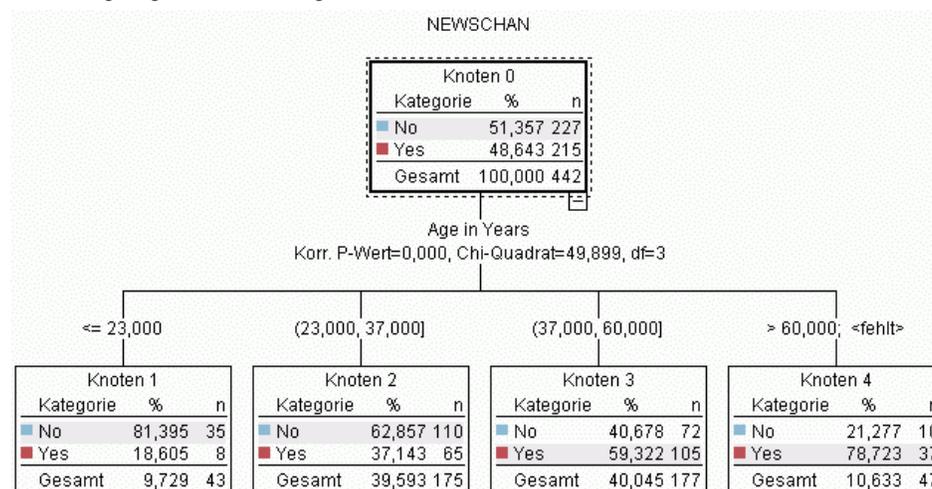
## Erweitern und Reduzieren des Baums

Auf der Registerkarte "Viewer" im Tree Builder können Sie den aktuellen Baum ab dem Stammknoten anzeigen.

- ▶ Um den Baum zu erweitern, wählen Sie in den Menüs folgende Optionsfolge:  
Baum > Baum erweitern
- Das System erstellt den Baum, indem jeder Zweig so lange rekursiv aufgeteilt wird, bis ein oder mehrere Grenzkriterien erfüllt sind. Auf der Grundlage der verwendeten Modellbildungsmethode wird bei jeder Aufteilung automatisch der beste Prädiktor ausgewählt.
- ▶ Sie können auch nur eine Ebene hinzufügen, indem Sie Baum um eine Ebene erweitern auswählen.
- ▶ Um eine Verzweigung unter einem bestimmten Knoten hinzuzufügen, wählen Sie den Knoten aus und klicken Sie auf Verzweigung erweitern.
- ▶ Um den für die Aufteilung verwendeten Prädiktor festzulegen, wählen Sie den gewünschten Knoten aus und klicken Sie auf Verzweigung mit benutzerdefinierter Aufteilung erweitern. [Für weitere Informationen siehe Thema Definieren benutzerdefinierter Aufteilungen auf S. 132.](#)
- ▶ Um eine Verzweigung zu reduzieren, wählen Sie einen Knoten aus und klicken Sie auf Verzweigung entfernen. Der ausgewählte Knoten wird gelöscht.
- ▶ Um die untere Ebene des Baums zu entfernen, wählen Sie Eine Ebene entfernen.
- ▶ Ausschließlich für C&R- und QUEST-Bäume können Sie Struktur erweitern und reduzieren wählen, um die Reduktion auf der Grundlage eines Kostenkomplexitätsalgorithmus durchzuführen, der die Risikoschätzung auf der Grundlage der Anzahl der Terminalknoten anpasst, was in der Regel zu einer einfacheren Struktur führt. [Für weitere Informationen siehe Thema C&R-Baumknoten auf S. 155.](#)

### Lesen von Aufteilungsregeln auf der Registerkarte "Viewer"

Abbildung 6-4  
Aufteilungsregeln auf der Registerkarte "Viewer"



Bei der Anzeige von Aufteilungsregeln auf der Registerkarte “Viewer” bedeuten eckige Klammern, dass der angegebene Wert im Bereich enthalten ist, während runde Klammern anzeigen, dass der Wert aus dem Bereich ausgeschlossen ist. Der Ausdruck (23,37] bedeutet somit von 23 (ausgeschlossen) bis einschließlich 37, also von etwas über 23 bis 37. Auf der Registerkarte “Modell” wird dieselbe Bedingung wie folgt angezeigt:

Age > 23 and Age <= 37

**Unterbrechen der Baumerweiterung.** Um die Baumerweiterung zu unterbrechen (wenn er beispielsweise länger als erwartet dauert), klicken Sie in der Symbolleiste auf die Schaltfläche “Ausführung anhalten”.

Abbildung 6-5  
Schaltfläche “Ausführung anhalten”



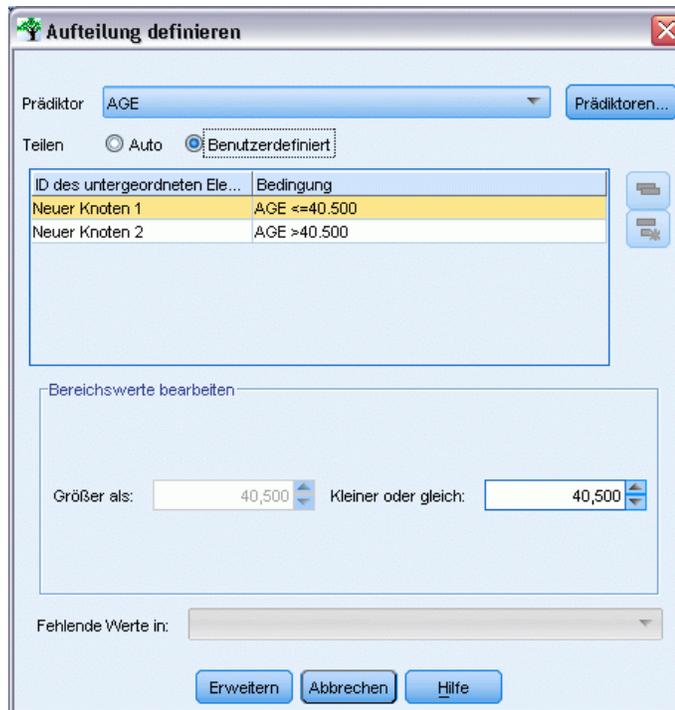
Diese Schaltfläche ist nur während der Strukturierung aktiviert. Sie stoppt den aktuellen Aufbauvorgang am gerade erreichten Punkt, wobei alle bereits hinzugefügten Knoten erhalten bleiben, ohne dass Änderungen gespeichert werden oder das Fenster geschlossen wird. Tree Builder bleibt geöffnet. Sie können nun ein Modell generieren, Richtlinien aktualisieren oder die Ausgabe im benötigten Format exportieren.

## ***Definieren benutzerdefinierter Aufteilungen***

Im Dialogfeld “Aufteilung definieren” können Sie den Prädiktor auswählen und für alle Aufteilungen Bedingungen festlegen.

- ▶ Wählen Sie in Tree Builder auf der Registerkarte “Viewer” einen Knoten aus und treffen Sie in den Menüs folgende Auswahl:  
Baum > Verzweigung mit benutzerdefinierter Aufteilung erweitern

Abbildung 6-6  
Dialogfeld "Aufteilung definieren"



- ▶ Wählen Sie in der Dropdown-Liste den gewünschten Prädiktor aus oder klicken Sie auf die Schaltfläche Prädiktoren, damit zu jedem Prädiktor Einzelheiten angezeigt werden. [Für weitere Informationen siehe Thema Anzeigen der Prädiktordetails auf S. 134.](#)
- ▶ Für jede Aufteilung können Sie die Standardbedingungen übernehmen oder über Angepasst die für die Aufteilung gewünschten Bedingungen festlegen.
  - Für stetige Prädiktoren (numerischer Bereich), können Sie die Felder Bereichswerte bearbeiten verwenden, um den Wertebereich anzugeben, der in jeden neuen Knoten fällt.
  - Für kategoriale Prädiktoren können Sie die Felder Set-Werte bearbeiten oder Ordinale Werte bearbeiten verwenden, um die spezifischen Werte (oder Wertebereiche, wenn es sich um einen ordinalen Prädiktor handelt) anzugeben, die jedem neuen Knoten zugeordnet sind.
- ▶ Wählen Sie Erweitern, damit die Verzweigung mit dem ausgewählten Prädiktor noch einmal erweitert wird.

Unabhängig von Stoppregeln kann der Baum in der Regel ohne Prädiktor aufgeteilt werden. Ausgenommen sind lediglich die Fälle, wo es sich um einen reinen Knoten handelt (d. h. 100 % der Fälle fallen in dieselbe Zielklasse, wodurch nichts mehr aufgeteilt werden kann) oder wenn der ausgewählte Prädiktor konstant ist (also nichts gegen ihn aufgeteilt werden kann).

**Fehlende Werte in.** Nur für CHAID-Knoten gibt es bei der Festlegung einer benutzerdefinierten Aufteilung die Option, fehlende Werte eines bestimmten Prädiktors einem bestimmten untergeordneten Knoten zuzuordnen. (Bei "C&R-Baum" und "QUEST" werden fehlende Werte mit im Algorithmus definierten Ersatztrennern bearbeitet. [Für weitere Informationen siehe Thema Aufteilungsdetails und Ersatztrenner auf S. 134.](#))

### Anzeigen der Prädiktordetails

Das Dialogfeld “Prädiktor auswählen” zeigt Statistiken der für den aktuellen Split verfügbaren Prädiktoren (oder “Konkurrenten”, wie sie zuweilen auch genannt werden).

Abbildung 6-7  
Dialogfeld “Prädiktor auswählen”



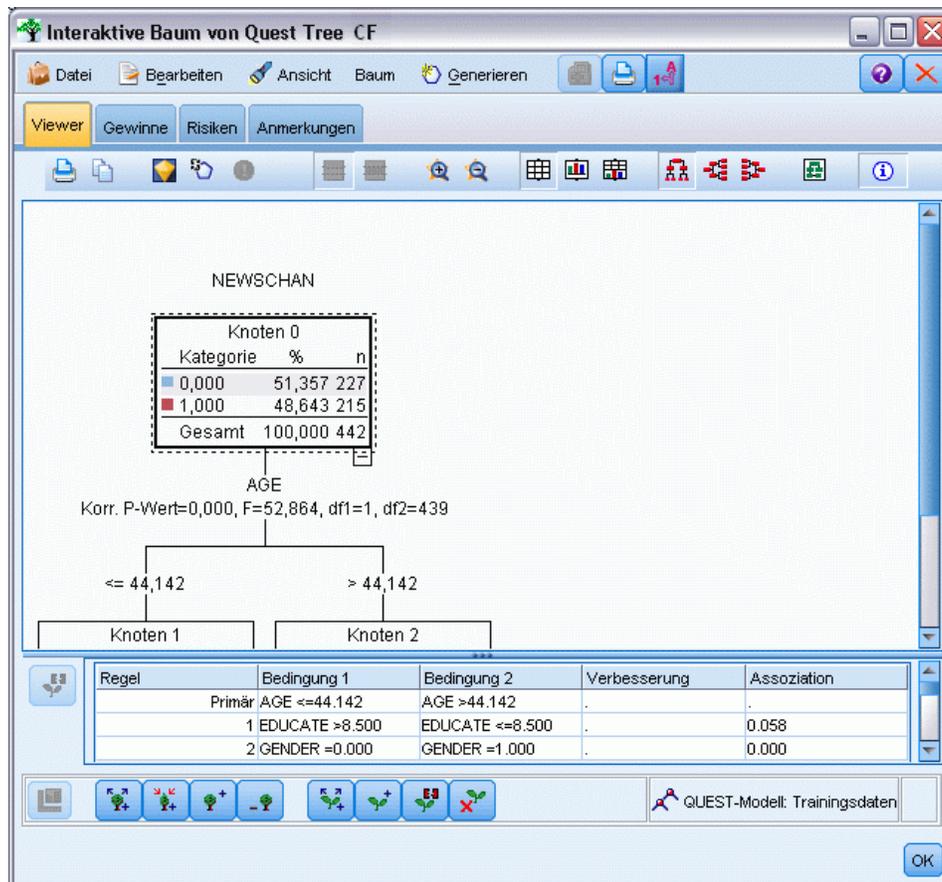
Prädiktor	Knoten	Verbesserung
AGE	2	0.055
GENDER	2	0.011
INC	2	0.010
ORGS	2	0.004
EDUCATE	2	0.004
TVDAY	2	0.004
CHILDS		

- Bei CHAID und Exhaustive CHAID wird für jeden kategorialen Prädiktor die Chi-Quadrat-Statistik aufgeführt. Wenn es sich bei einem Prädiktor um einen numerischen Bereich handelt, wird die  $F$ -Statistik angezeigt. Die Chi-Quadrat-Statistik ist ein Maß dafür, wie unabhängig das Zielfeld vom Aufteilungsfeld ist. Eine hoch ausfallende Chi-Quadrat-Statistik weist in der Regel auf eine geringere Wahrscheinlichkeit hin, d. h. die Chance, dass zwei Felder unabhängig sind, ist niedriger, was wiederum bedeutet, dass der Split gut ist. Freiheitsgrade werden aufgenommen, weil diese die Tatsache berücksichtigen, dass es bei einer dreifachen Aufteilung einfacher ist, eine hoch ausfallende Statistik und eine geringe Wahrscheinlichkeit zu erhalten, als dies bei einer zweifachen Aufteilung der Fall ist.
- Bei “C&R-Baum” und “QUEST” wird für jeden Prädiktor die Verbesserung angezeigt. Je größer die Verbesserung ist, desto stärker reduziert der Einsatz des Prädiktors die zwischen den über- und den untergeordneten Knoten entstehende Unreinheit. (Ein reiner Knoten liegt dann vor, wenn alle Fälle in eine einzige Zielkategorie fallen. Je geringer die Unreinheit des gesamten Baums ist, desto besser passt das Modell zu den Daten.) Anders ausgedrückt, weist ein hoher Verbesserungswert in der Regel auf eine brauchbare Aufteilung für diesen Baumtyp hin. Das verwendete Unreinheitsmaß wird im Baumerstellungsknoten festgelegt.

### Aufteilungsdetails und Ersatztrenner

Sie können jeden Knoten auf der Registerkarte “Viewer” auswählen und über die Schaltfläche für die Aufteilungsinformationen (rechts in der Symbolleiste) die Details über die Aufteilung des Knotens anzeigen. Zusammen mit der verwendeten Aufteilungsregel werden die relevanten Statistiken angezeigt. Für kategoriale C&RT-Bäume werden die Verbesserung und die Assoziation angezeigt. Die Assoziation ist ein Maß für die Entsprechung zwischen einem Ersatztrenner und dem primären Aufteilungsfeld, wobei der “beste” Ersatztrenner in der Regel derjenige ist, der sich dem Aufteilungsfeld am weitesten angleicht. Bei “C&R-Baum” und “QUEST” werden auch alle anstelle des primären Prädiktors verwendeten Ersatztrenner aufgeführt.

Abbildung 6-8  
Fenster des interaktiven Tree Builder mit Aufteilungsinformationen



- Um die Aufteilung des ausgewählten Knotens zu bearbeiten, klicken Sie auf der linken Seite des Ersatztrennerbereichs auf das Symbol, damit das Dialogfeld “Aufteilung definieren” geöffnet wird. (Sie können das Verfahren abkürzen, indem Sie in der Liste einen Ersatztrenner auswählen und diesen durch Anklicken des Symbols zum primären Aufteilungsfeld machen.)

**Ersatzfelder.** Wo zutreffend, werden für den ausgewählten Knoten alle Ersatzfelder für das primäre Aufteilungsfeld angezeigt. Ersatzfelder sind alternative Felder, die verwendet werden, wenn der primäre Prädiktorwert für einen bestimmten Datensatz fehlt. Die maximale Anzahl an Ersatzfeldern, die für eine bestimmte Aufteilung erlaubt ist, ist im Baumerstellungsknoten angegeben, die tatsächliche Anzahl richtet sich jedoch nach den Trainingsdaten. Im Allgemeinen gilt: Je mehr fehlende Daten, desto mehr Ersatzfelder werden wahrscheinlich verwendet. Für andere Entscheidungsbaummodelle bleibt diese Registerkarte leer.

*Anmerkung:* Um im Modell berücksichtigt zu werden, müssen die Ersatzfelder während der Trainingsphase ermittelt werden. Wenn die Trainings-Stichprobe keine fehlenden Werte enthält, werden keine Ersatzfelder ermittelt. Alle Datensätze mit fehlenden Werten, die während des Testens oder der Bewertung gefunden werden, fallen automatisch in den untergeordneten Knoten mit der größten Anzahl an Datensätzen. Wenn während des Testens oder Bewertens fehlende

Werte erwartet werden, können Sie sicher sein, dass auch in den Trainings-Stichproben Werte fehlen. Für CHAID-Bäume stehen keine Ersatzfelder zur Verfügung.

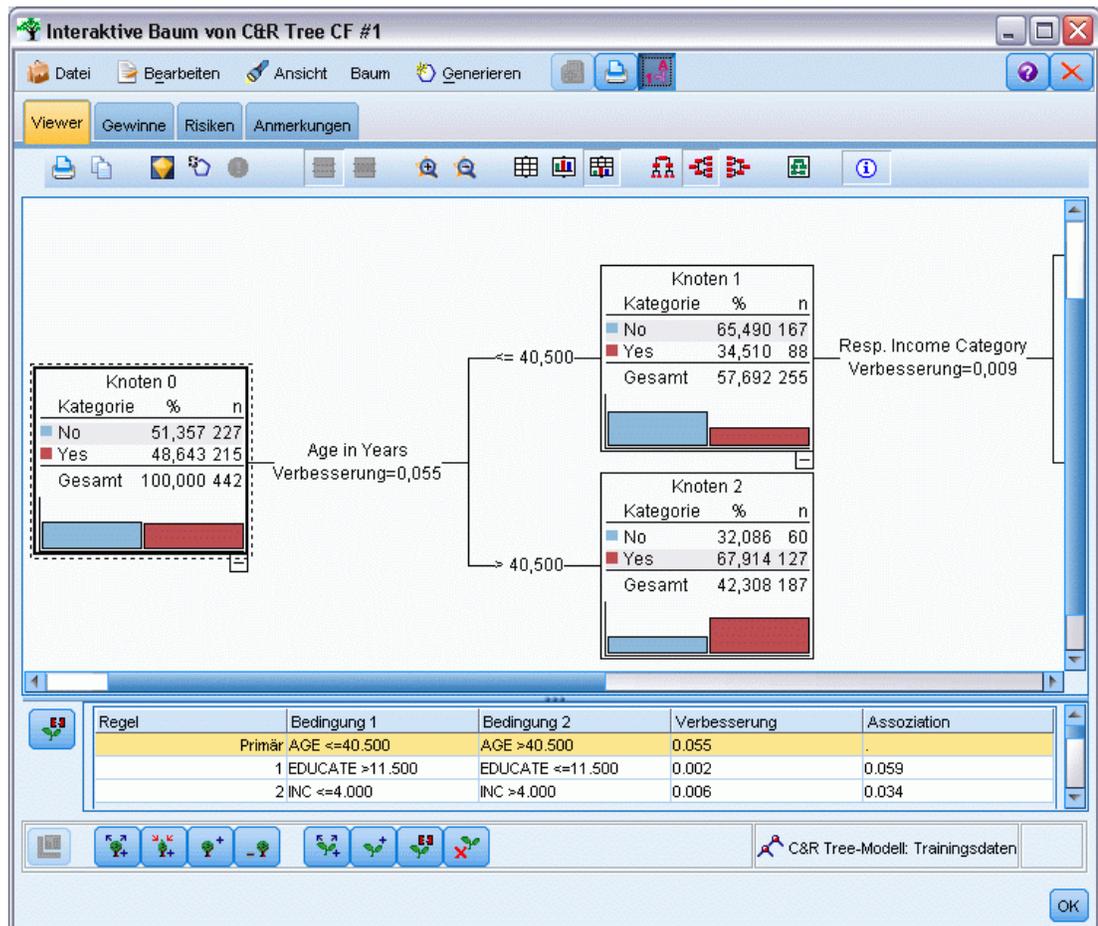
Obwohl bei CHAID-Knoten keine Ersatztrenner benutzt werden, haben Sie in einer benutzerdefinierten Aufteilung die Option, diese einem bestimmten untergeordneten Knoten zuzuordnen. [Für weitere Informationen siehe Thema Definieren benutzerdefinierter Aufteilungen auf S. 132.](#)

## Anpassen der Baumansicht

Auf der Registerkarte “Viewer” von Tree Builder wird der aktuelle Baum angezeigt. Standardmäßig sind alle Verzweigungen des Baums erweitert. Sie können Verzweigungen allerdings ausblenden und erweitern oder weitere Einstellungen anpassen.

Abbildung 6-9

Von links nach rechts: Aufteilungsregeldetails, Knotendiagramme und eingblendete Beschriftungen



- Klicken Sie auf das Minuszeichen (–) in der rechten unteren Ecke eines übergeordneten Knotens, um alle untergeordneten Knoten auszublenden. Klicken Sie auf das Pluszeichen (+) in der rechten unteren Ecke eines übergeordneten Knotens, um dessen untergeordnete Knoten anzuzeigen.

- Die Ausrichtung des Baums (von oben nach unten, von links nach rechts, von rechts nach links) können Sie im Menü “Ansicht” oder über die Symbolleiste ändern.
- Zum Anzeigen oder Ausblenden von Feld- und Wertelabels klicken Sie auf die Schaltfläche “Feld- und Wertelabels anzeigen” in der Hauptsymbolleiste.
- Die Schaltflächen mit den Vergrößerungsgläsern vergrößern oder verkleinern die Anzeige. Wenn Sie rechts in der Symbolleiste auf die Baumübersichtsschaltfläche klicken, wird ein Diagramm des gesamten Baums angezeigt.
- Wenn ein Partitionsfeld verwendet wird, können Sie die Baumansicht zwischen der Trainings- und der Testpartition umschalten (Ansicht > Partition). Wenn die Teststichprobe angezeigt wird, kann der Baum zwar angezeigt, aber nicht bearbeitet werden. (Die aktuelle Partition wird in der unteren linken Fensterecke in der Statusleiste angegeben.)
- Zum Anzeigen von Details über die aktuelle Aufteilung klicken Sie auf die Schaltfläche “Aufteilungsinformationen” (die Schaltfläche “i” ganz rechts in der Symbolleiste). [Für weitere Informationen siehe Thema Aufteilungsdetails und Ersatztrenner auf S. 134.](#)
- Für jeden Knoten Statistiken, Diagramme oder beides anzeigen (siehe unten).

### **Anzeigen von Statistiken und Diagrammen**

**Knotenstatistiken.** Bei einem kategorialen Zielfeld zeigt die Tabelle in jedem Knoten die Anzahl und den Prozentsatz der in jeder Kategorie enthaltenen Datensätze sowie den Prozentsatz der gesamten vom Knoten dargestellten Stichprobe. Bei einem stetigen Zielfeld (numerischer Bereich) zeigt die Tabelle die durchschnittliche und die Standardabweichung, die Anzahl der Datensätze und den vorhergesagten Wert des Zielfelds.

**Knotendiagramme.** Bei einem kategorialen Zielfeld wird das Diagramm als Balkendiagramm der in jeder Kategorie des Zielfelds enthaltenen Prozentsätze ausgegeben. Vor jeder Tabellenzeile befindet sich ein Farbmuster, das die Farbe angibt, die für die Zielfeldkategorie in den Diagrammen des Knotens verwendet wird. Bei einem stetigen Zielfeld (numerischer Bereich) wird das Diagramm als Histogramm des Zielfelds der im Knoten enthaltenen Datensätze angezeigt.

## **Gewinne**

Die Registerkarte “Gewinne” zeigt Statistiken aller im Baum vorhandenen Endknoten. Gewinne bieten ein Maß dafür, wie weit der Mittelwert oder der Anteil eines bestimmten Knotens vom Gesamtmittelwert abweicht. In der Regel gilt, je größer der Unterschied, desto besser kann der Baum als Werkzeug zur Entscheidungsfindung eingesetzt werden. Ein Index- oder “Lift”-Wert für einen Knoten von 148 % zeigt beispielsweise an, dass die Wahrscheinlichkeit, dass in diesem Knoten enthaltene Datensätze unter die Zielkategorie fallen, fast anderthalb Mal so hoch ist, wie dies für das gesamte Daten-Set der Fall ist.

Für CRT-Baum- und QUEST-Knoten, für die ein Set zur Verhinderung übermäßiger Anpassung angegeben ist, werden zwei Sets von Statistikdaten angezeigt:

- Baumaufbau-Set – Trainingsstichprobe ohne Set zur Verhinderung übermäßiger Anpassung
- Set zur Verhinderung der übermäßigen Anpassung

Für andere interaktive CRT-Bäume und QUEST-Bäume und für alle interaktiven CHAID-Bäume werden nur die Baumaufbau-Set-Statistiken angezeigt.

Abbildung 6-10  
Registerkarte "Gewinne"

Baumerweiterungs-Set							Set zur Vermeidung einer Überanpassung						
Knoten	Knoten: n	Knoten (...)	Gewinn: n	Gewinn (...)	Treffer (%)	Index (%)	Knoten	Knoten: n	Knoten (...)	Gewinn: n	Gewinn (...)	Treffer (%)	Index (%)
2	73,00	48,34	45,00	64,29	61,64	132,97	2	37,00	52,86	23,00	69,70	62,16	131,86
1	78,00	51,66	25,00	35,71	32,05	69,14	1	33,00	47,14	10,00	30,30	30,30	64,28

Die Registerkarte "Gewinne" bietet folgende Optionen:

- Anzeige knotenweiser, kumulativer oder quantiler Statistiken.
- Anzeigen von Gewinnen oder Profiten.
- Umschalten der Anzeige zwischen Tabellen und Diagrammen.
- Auswählen der Zielkategorie (nur für kategoriale Ziele).
- Sortieren der Tabelle in auf- oder absteigender Reihenfolge, basierend auf dem Indexprozentsatz. Wenn Statistiken für mehrere Partitionen angezeigt werden, wird die Sortierung immer in der Trainings- und nicht in der Teststichprobe vorgenommen.

Eine in der Tabelle der Gewinne getroffene Auswahl wird in der Regel in der Baumansicht aktualisiert und umgekehrt. Wenn Sie beispielsweise in der Tabelle eine Zeile auswählen, dann wird im Baum der entsprechende Knoten ausgewählt.

### **Klassifizierungsgewinne**

Bei Klassifizierungsbäumen (mit einer kategorialen Zielvariablen) zeigt der Indexprozentsatz des Gewinns, wie stark der Anteil einer bestimmten Zielkategorie jedes Knotens von dem Gesamtanteil abweicht.

#### **Knotenweise Statistiken**

In dieser Ansicht enthält die Tabelle in jeder Zeile einen Endknoten. Beispiel: Wenn Ihre Direktmarketingkampagne insgesamt eine Trefferquote von 10 % erzielt hat, dabei aber 20 % der unter Knoten X fallenden Datensätze positiv waren, dann liegt der Indexprozentsatz dieses Knotens bei 200 %. Dieser Wert drückt aus, dass die Wahrscheinlichkeit, dass die in dieser Gruppe enthaltenen Teilnehmer kaufen, doppelt so hoch ist wie in der Gesamtpopulation.

Für CRT-Baum- und QUEST-Knoten, für die ein Set zur Verhinderung übermäßiger Anpassung angegeben ist, werden zwei Sets von Statistikdaten angezeigt:

- Baumaufbau-Set – Trainingsstichprobe ohne Set zur Verhinderung übermäßiger Anpassung
- Set zur Verhinderung der übermäßigen Anpassung

Für andere interaktive CRT-Bäume und QUEST-Bäume und für alle interaktiven CHAID-Bäume werden nur die Baumaufbau-Set-Statistiken angezeigt.

Abbildung 6-11  
Knotenweise Gewinnstatistiken

Baumerweiterungs-Set							Set zur Vermeidung einer Überanpassung						
Knoten	Knoten: n	Knoten (...)	Gewinn: n	Gewinn (...)	Treffer (%)	Index (%)	Knoten	Knoten: n	Knoten (...)	Gewinn: n	Gewinn (...)	Treffer (%)	Index (%)
2	73,00	48,34	45,00	64,29	61,64	132,97	2	37,00	52,86	23,00	69,70	62,16	131,86
1	78,00	51,66	25,00	35,71	32,05	69,14	1	33,00	47,14	10,00	30,30	30,30	64,28

**Knoten.** Die ID des aktuellen Knotens (die auf der Registerkarte “Viewer” ausgegeben wird).

**Knoten: n.** Die Gesamtzahl der in diesem Knoten enthaltenen Datensätze.

**Knoten (%).** Der Prozentsatz aller im Daten-Set enthaltenen Datensätze, die unter diesen Knoten fallen.

**Gewinn: n.** Die Anzahl der Datensätze mit der ausgewählten Zielkategorie, die unter diesen Knoten fallen. Dieser Wert gibt an, wie viele der insgesamt im Daten-Set enthaltenen Datensätze, die unter die Zielkategorie fallen, sich in diesem Knoten befinden.

**Gewinn (%).** Der Prozentsatz aller in der Zielkategorie enthaltenen Datensätze des gesamten Daten-Sets, die unter diesen Knoten fallen.

**Treffer (%).** Der Prozentsatz der Datensätze des aktuellen Knotens, die unter die Zielkategorie fallen. Treffer werden in diesem Kontext manchmal auch als “Hits” bezeichnet.

**Index (%).** Der Trefferprozentsatz des aktuellen Knotens, ausgedrückt als Prozentsatz des Trefferprozentsatzes des gesamten Daten-Sets. Beispiel: Ein Indexwert von 300 % zeigt an, dass die Wahrscheinlichkeit, dass in diesem Knoten enthaltene Datensätze unter die Zielkategorie fallen, dreimal so hoch ist, wie dies im gesamten Daten-Set der Fall ist.

### **Kumulative Statistiken**

In der kumulativen Ansicht enthält die Tabelle einen Knoten pro Zeile, wobei die Statistiken aber kumulativ und auf- oder absteigend nach Indexprozentsatz sortiert sind. Wenn beispielsweise eine absteigende Sortierung vorliegt, wird in der ersten Zeile der Knoten mit dem höchsten Indexprozentsatz ausgegeben und in den darauffolgenden Zeilen erscheinen die kumulierten Werte der jeweiligen Zeile mit den darüberliegenden Zeilen.

Abbildung 6-12  
Kumulative Gewinne, absteigend nach Indexprozensatz sortiert

Baumerweiterungs-Set							Set zur Vermeidung einer Überanpassung						
Knoten	Knoten: n	Knoten (...)	Gewinn: n	Gewinn (...)	Treffer (%)	Index (%)	Knoten	Knoten: n	Knoten (...)	Gewinn: n	Gewinn (...)	Treffer (%)	Index (%)
2	73,00	48,34	45,00	64,29	61,64	132,97	2	37,00	52,86	23,00	69,70	62,16	131,86
1	151,00	100,00	70,00	100,00	46,36	100,00	1	70,00	100,00	33,00	100,00	47,14	100,00

Der kumulierte Indexprozensatz sinkt mit jeder Zeile, da Knoten mit immer niedrigeren Trefferprozenten hinzugefügt werden. Der kumulative Index der letzten Zeile liegt immer bei 100 %, da an diesem Punkt das gesamte Daten-Set enthalten ist.

### Quantile

In dieser Ansicht wird in jeder Zeile anstelle eines Knotens ein Quantil angezeigt. Bei den Quantilen handelt es sich entweder um Quartile, Quintile (Fünftel), Dezile (Zehntel), Vingtile (Zwanzigstel) oder Perzentile (Hundertstel). Sofern mehrere Knoten benötigt werden, um den Prozentsatz zu erreichen, können in einem Quantil mehrere Knoten aufgeführt sein (wenn z. B. Quartile angezeigt werden, die beiden höchsten Knoten jedoch weniger als 50 % aller Fälle enthalten). Die übrige Tabelle ist kumulativ und genau wie die kumulative Ansicht zu interpretieren.

Abbildung 6-13  
Gewinne pro Quartil, absteigend nach Indexprozensatz sortiert

Baumerweiterungs-Set							Set zur Vermeidung einer Überanpassung						
Knoten	Perzentil	Perzentil: n	Gewinn: n	Gewinn (...)	Treffer (%)	Index (%)	Knoten	Perzentil	Perzentil: n	Gewinn: n	Gewinn (...)	Treffer (%)	Index (%)
2	25,00	38,00	23,00	33,46	61,64	132,97	2	25,00	18,00	11,00	33,91	62,16	131,86
2,1	50,00	76,00	46,00	65,66	60,48	130,45	2	50,00	35,00	22,00	65,93	62,16	131,86
1	75,00	113,00	58,00	82,60	51,17	110,38	2,1	75,00	53,00	28,00	84,39	52,54	111,46
1	100,00	151,00	70,00	100,00	46,36	100,00	1	100,00	70,00	33,00	100,00	47,14	100,00

### Klassifizierung der Profite und ROI

Die Gewinnstatistiken für Klassifizierungsbäume können auch mit Profit und ROI (Return on Investment) ausgegeben werden. Für jede Kategorie können Sie im Dialogfeld "Profite definieren" Einnahmen und Ausgaben angeben.

- Öffnen Sie die Registerkarte "Gewinne" (mit dem Symbol \$/\$) und klicken Sie in der Symbolleiste auf die Schaltfläche "Profit", um das Dialogfeld zu öffnen.

Abbildung 6-14  
Dialogfeld "Profite definieren"



- ▶ Geben Sie für jede Kategorie des Zielfelds Einnahmen- und Ausgabenwerte ein.

Beispiel: Wenn es Sie 0,48 \$ kostet, ein Angebot an einen Kunden zu schicken und die Einnahme bei einer positiven Antwort für ein dreimonatiges Abonnement 9,95 \$ beträgt, dann kostet Sie jede *Nein*-Antwort 0,48 \$ und jede *Ja*-Antwort bringt Ihnen 9,47 \$ (9,95 – 0,48).

In der Tabelle "Gewinne" wird der **Profit** als die Summe der Einnahmen abzüglich der Ausgaben für alle im Endknoten enthaltenen Datensätze berechnet. **ROI** ist der Gesamtprofit geteilt durch die Gesamtausgaben in einem Knoten.

#### **Kommentare**

- Profitwerte wirken sich nur auf die in der Tabelle "Gewinne" angezeigten durchschnittlichen Profit- und ROI-Werte aus und bieten eine für Ihr Endergebnis brauchbarere statistische Anzeige. Die grundlegende Baummodellstruktur bleibt unverändert. Profite dürfen nicht mit Fehlklassifizierungskosten verwechselt werden, die im Baumerstellungsknoten angegeben sind und zum Schutz gegen teure Fehler als Faktor in das Modell eingehen.
- Profitangaben sind zwischen interaktiven Baumerstellungssitzungen nicht persistent.

#### **Regressionsgewinne**

Für Regressionsbäume kann eine knotenweise, kumulativ knotenweise und eine quantile Ansicht gewählt werden. Durchschnittswerte werden in der Tabelle ausgegeben. Diagramme sind nur für Quantile verfügbar.

#### **Gewinndiagramme**

Alternativ zu Tabellen können auf der Registerkarte "Gewinne" Diagramme angezeigt werden.

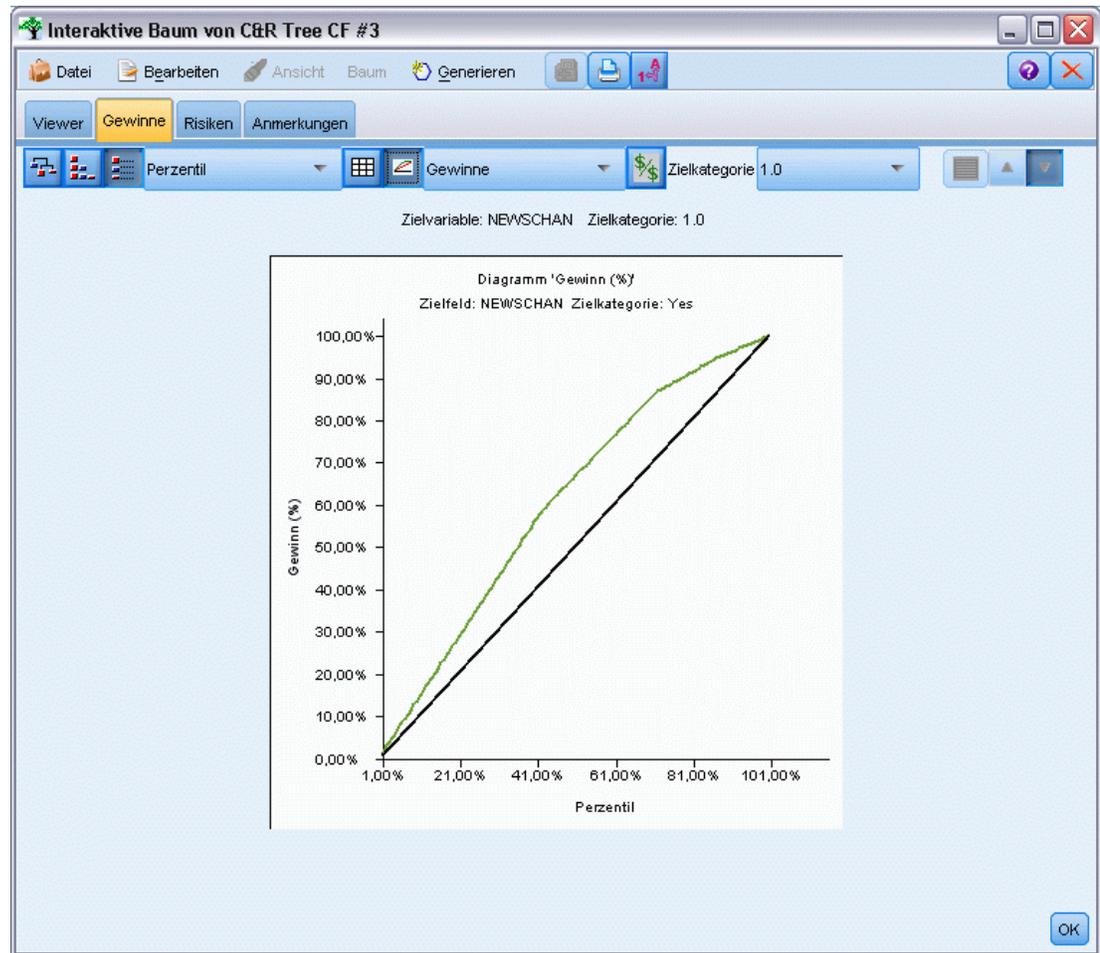
- ▶ Wählen Sie auf der Registerkarte "Gewinne" das Symbol für Quantile (das dritte von links auf der Symbolleiste). (Für knotenweise oder kumulative Statistiken sind keine Diagramme verfügbar.)
- ▶ Klicken Sie auf das Diagrammsymbol.
- ▶ Wählen Sie in der Dropdown-Liste die Einheit aus (Perzentile, Dezile etc.), die angezeigt werden soll.
- ▶ Wählen Sie Gewinne, Treffer oder Lift, um die angezeigte Messung zu ändern.

### Gewinndiagramm

Das Gewinndiagramm bildet die Werte der Tabellenspalte *Gewinn (%)* ab. Gewinne sind als der Anteil der in jedem Inkrement enthaltenen Treffer im Verhältnis zur Gesamtzahl der im Baum enthaltenen Treffer definiert. Dabei kommt folgende Gleichung zum Einsatz:

$$(\text{Treffer im Inkrement} / \text{Gesamtzahl Treffer}) \times 100 \%$$

Abbildung 6-15  
Gewinndiagramm



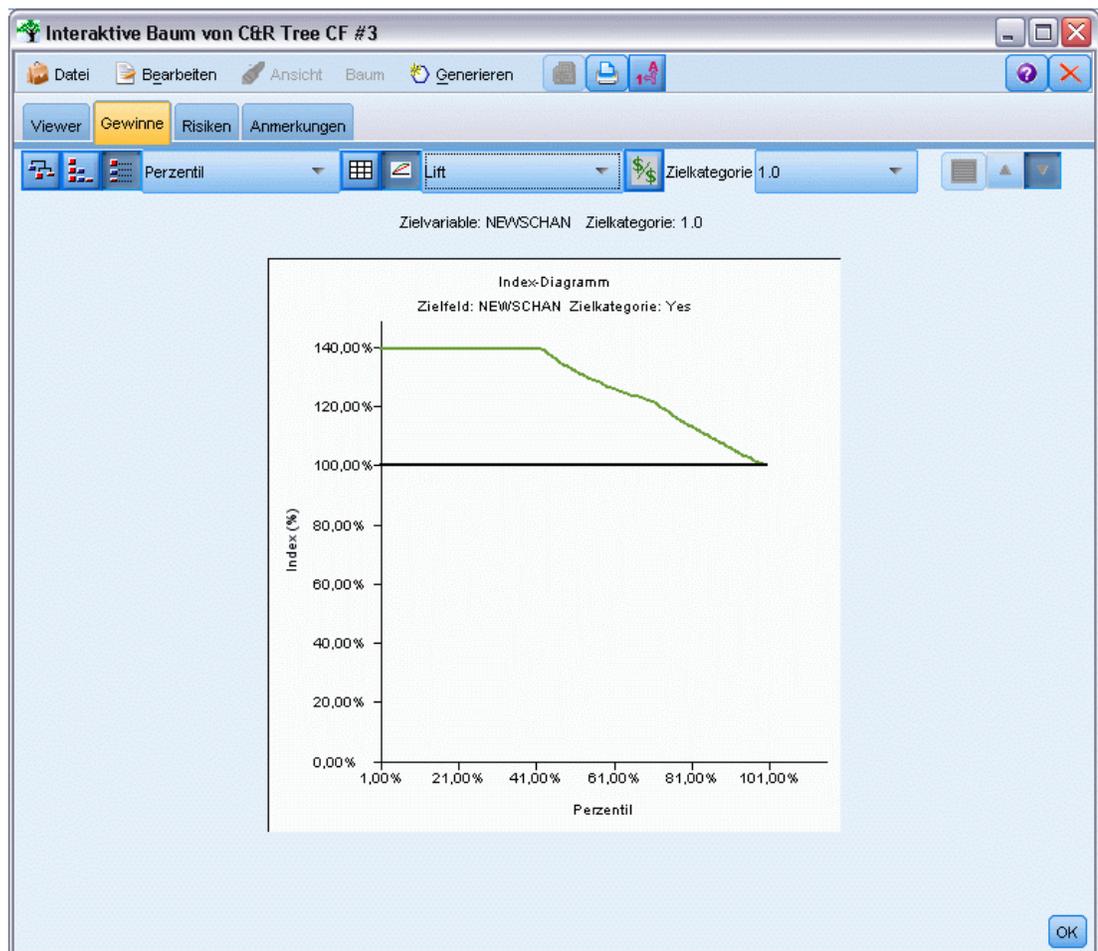
Das Diagramm illustriert, wie weit Sie das Netz auswerfen müssen, um einen bestimmten Prozentsatz aller im Baum enthaltenen Treffer zu erzielen. Die diagonale Linie bildet die für die gesamte Stichprobe erwarteten Treffer ab, wenn das Modell nicht verwendet wird. In diesem Fall ist die Trefferrate konstant, da die Wahrscheinlichkeit eines Treffers für alle Personen gleich ist. Um das Ergebnis zu verdoppeln, müssen Sie doppelt so viele Personen ansprechen. Die gekrümmte Linie zeigt an, wie weit Sie Ihre Treffer verbessern können, wenn Sie nur die einschließen, deren Prozentsatz hinsichtlich des Gewinns höher ausfällt. Wenn Sie beispielsweise die obersten 50 % einschließen, erhalten Sie über 70 % der positiven Treffer. Je steiler die Kurve, desto höher ist der Gewinn.

### Lift Chart

Das Lift Chart bildet die Werte der Tabellenspalte *Index (%)* ab. Dieses Diagramm vergleicht den Prozentsatz der in jedem Inkrement enthaltenen Datensätze, bei denen es sich um Treffer handelt, mit dem Prozentsatz aller im Trainings-Daten-Set enthaltenen Treffer. Dabei wird folgende Gleichung zugrunde gelegt:

$$(\text{Treffer im Inkrement} / \text{Datensätze im Inkrement}) / (\text{Gesamtzahl Treffer} / \text{Gesamtzahl Datensätze})$$

Abbildung 6-16  
Lift Chart (Index)

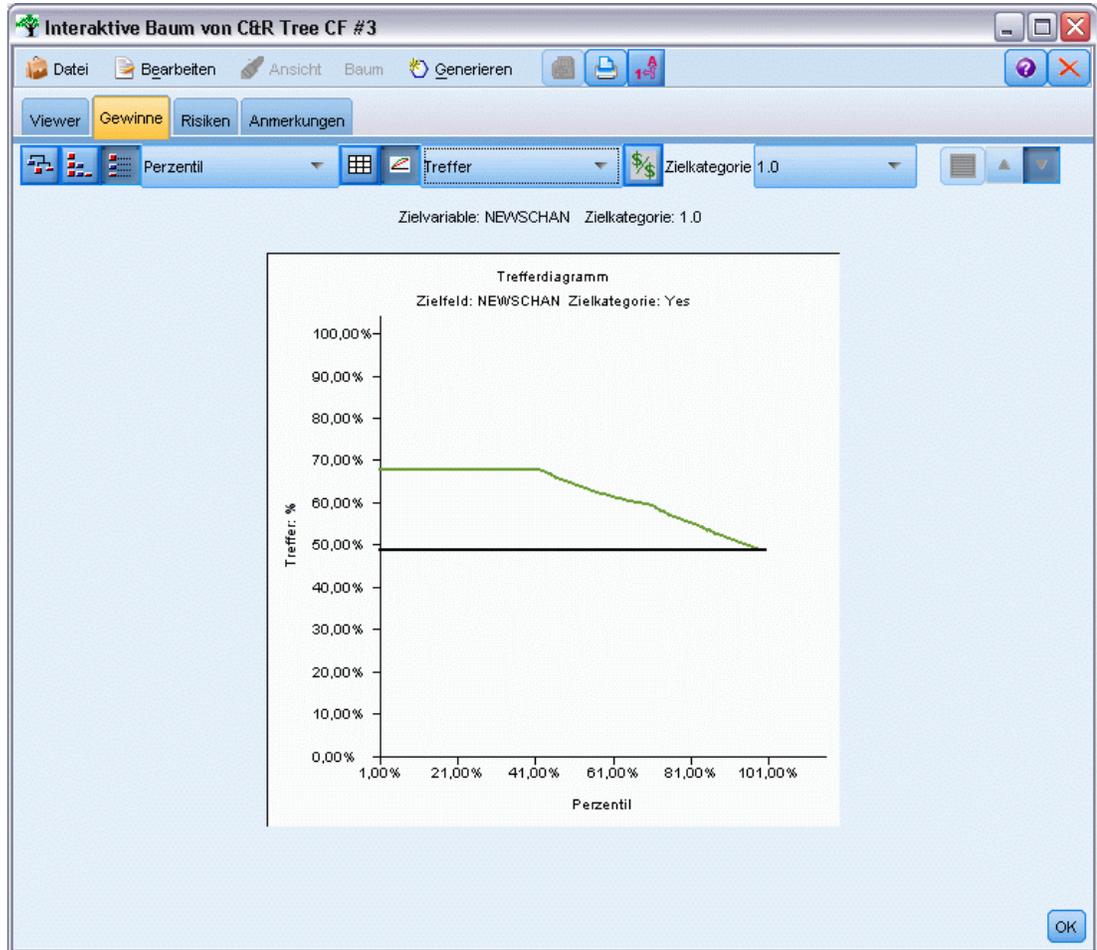


### Trefferdiagramm

Das Trefferdiagramm bildet die Werte der Tabellenspalte *Treffer (%)* ab. Ergebnis ist der Prozentsatz der im Inkrement enthaltenen Datensätze, bei denen es sich um Treffer handelt, wobei folgende Gleichung angewendet wird:

$$(\text{Treffer im Inkrement} / \text{Datensätze im Inkrement}) \times 100 \%$$

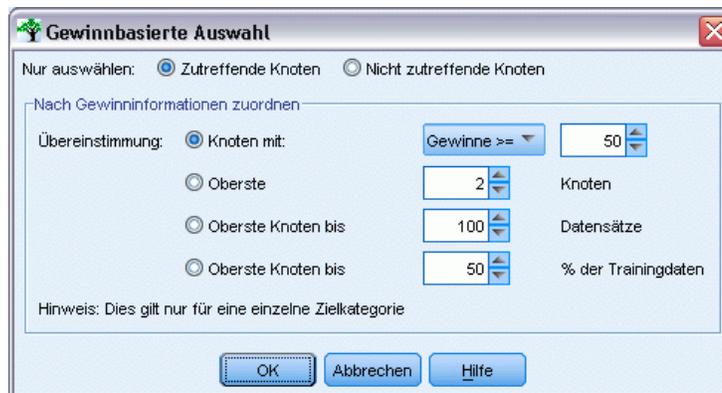
Abbildung 6-17  
Trefferdiagramm



### **Gewinnbasierte Auswahl**

Im Dialogfeld "Gewinnbasierte Auswahl" können Sie Endknoten mit den höchsten (oder niedrigsten) Gewinnen automatisch auf der Grundlage einer vorgegebenen Regel oder eines Grenzwerts auswählen. Auf der Grundlage dieser Auswahl können Sie dann einen Auswahlknoten generieren.

Abbildung 6-18  
Dialogfeld "Gewinnbasierte Auswahl"



- ▶ Wählen Sie auf der Registerkarte "Gewinne" die knotenweise oder die kumulative Anzeige sowie die Zielkategorie aus, die die Grundlage für die Auswahl bilden soll. (Eine Auswahl basiert auf der aktuellen Tabellenanzeige und ist für Quantile nicht verfügbar.)
- ▶ Wählen Sie auf der Registerkarte "Gewinne" folgende Menüoptionen:  
Bearbeiten > Endknoten auswählen > Gewinnbasierte Auswahl

**Nur auswählen.** Sie können zutreffende Knoten *oder* nicht zutreffende Knoten auswählen, um beispielsweise *alle außer* den 100 höchsten Knoten auszuwählen.

**Nach Gewinninformationen zuordnen.** Ordnet Knoten auf der Grundlage der Gewinnstatistiken der aktuellen Zielkategorie zu. Dazu gehören:

- Knoten, deren Gewinn-, Treffer- oder Lift-Wert (Index) einem bestimmten Grenzwert entspricht, wie beispielsweise einem Trefferwert größer oder gleich 50 %.
- Die  $n$  höchsten Knoten, basierend auf dem Gewinn für die Zielkategorie.
- Die höchsten Knoten, bis zu einer vorgegebenen Anzahl Knoten.
- Die höchsten Knoten, bis zu einem vorgegebenen Prozentsatz der Trainings-Daten.

- ▶ Klicken Sie auf OK, um die Auswahl auf der Registerkarte "Viewer" zu aktualisieren.
- ▶ Um auf der Grundlage der aktuellen Auswahl auf der Registerkarte "Viewer" einen neuen Auswahlknoten zu erstellen, wählen Sie im Menü "Generieren" die Option Auswahlknoten. [Für weitere Informationen siehe Thema Generieren von Filter- und Auswahlknoten auf S. 151.](#)

*Hinweis:* Da Sie zurzeit Knoten anstelle von Datensätzen oder Prozentzahlen ausgewählt haben, kann nicht immer eine perfekte Übereinstimmung mit dem Auswahlkriterium erreicht werden. Das System wählt komplette Knoten *bis zu* der vorgegebenen Ebene aus. Wenn Sie beispielsweise die höchsten 12 Fälle auswählen, von denen 10 im ersten und 2 im zweiten Knoten liegen, dann wird nur der erste Knoten ausgewählt.

## Risiken

Risiken geben an, wie groß die Gefahr einer Fehlklassifizierung auf jeder Ebene ist. Die Registerkarte "Risiken" enthält eine punktuelle Risikoschätzung und (für kategoriale Ausgaben) eine Fehlklassifizierungstabelle.

Abbildung 6-19  
Fehlklassifizierungstabelle für ein kategoriales Ziel

The screenshot shows the 'Risiken' tab in the 'Interaktive Baum von NEWSCHAN' software. It displays two risk assessment sections, each with a risk estimate, standard error, and a confusion matrix.

**Baumerweiterungs-Set**

Risikoschätzung: **0,351**  
Standardfehler: **0,039**

Vorhergesagt		Fehlklassifizierungsmatrix		
		0.0	1.0	Gesamt
0.0	53	28	81	
1.0	25	45	70	
Gesamt	78	73	151	

**Set zur Vermeidung einer Überanpassung**

Risikoschätzung: **0,343**  
Standardfehler: **0,057**

Vorhergesagt		Fehlklassifizierungsmatrix		
		0.0	1.0	Gesamt
0.0	23	14	37	
1.0	10	23	33	
Gesamt	33	37	70	

- Bei numerischen Vorhersagen bildet das Risiko eine Gesamtschätzung der in jedem Endknoten vorhandenen Varianz.
- Bei kategorialen Vorhersagen bildet das Risiko den Anteil der falsch klassifizierten Fälle, angepasst um alle A Priori- oder Fehlklassifizierungskosten.

## Speichern der Baummodelle und Ergebnisse

Die Ergebnisse Ihrer interaktiven Baumerstellung können Sie auf mehrere Arten speichern oder exportieren:

- Generieren eines Modells auf der Grundlage des aktuellen Baums (Generieren > Modell erzeugen).
- Speichern der für die Erweiterung des aktuellen Baums verwendeten Richtlinien. Wenn Sie den Baumerstellungsknoten das nächste Mal ausführen, wird der aktuelle Baum automatisch wieder mit allen von Ihnen festgelegten benutzerdefinierten Aufteilungen aufgebaut.
- Exportieren der Modell-, Gewinn- und Risikoinformationen. [Für weitere Informationen siehe Thema Exportieren der Modell-, Gewinn- und Risikoinformationen auf S. 150.](#)

Im Tree Builder oder in einem Baummodell-Nugget können Sie:

- Einen Filter generieren oder einen Knoten basierend auf dem aktuellen Baum auswählen. [Für weitere Informationen siehe Thema Generieren von Filter- und Auswahlknoten auf S. 151.](#)
- Ein Regelmengen-Nugget generieren, das die Baumstruktur als Set von Regeln darstellt, das die Endverzweigungen des Baums definiert. [Für weitere Informationen siehe Thema Generieren einer Regelmenge aus einem Entscheidungsbaum auf S. 151.](#)
- Bei Baummodell-Nuggets können Sie außerdem das Modell im PMML-Format exportieren. [Für weitere Informationen siehe Thema Die Modellpalette in Kapitel 3 auf S. 50.](#) Wenn das Modell benutzerdefinierte Aufteilungen enthält, werden diese Informationen in der exportierten PMML nicht beibehalten. (Die Aufteilung wird beibehalten, die Tatsache, dass sie benutzerdefiniert und nicht vom Algorithmus gewählt ist, jedoch nicht.)
- Generieren Sie ein Diagramm auf der Basis des ausgewählten Teils des aktuellen Baums. *Anmerkung:* Dies ist nur für ein Nugget möglich, wenn es mit anderen Knoten in einem Stream verbunden ist. [Für weitere Informationen siehe Thema Erzeugen von Diagrammen auf S. 191.](#)

*Hinweis:* Der interaktive Baum als solcher kann nicht gespeichert werden. Damit Ihre Arbeit nicht verloren geht, müssen Sie ein Modell generieren und/oder die Aufbauregeln aktualisieren, bevor Sie das Tree Builder-Fenster schließen.

### Generieren eines Modells mit Tree Builder

Um ein Modell auf der Grundlage des aktuellen Baums zu generieren, wählen Sie in den Tree Builder-Menüs folgende Optionen:

Erzeugen > Modell

Abbildung 6-20  
Generieren eines Entscheidungsbaummodells



Sie können folgende Optionen auswählen:

**Modellname.** Sie können einen benutzerdefinierten Namen angeben oder den Namen auf der Grundlage des Namens des Modellierungsknotens automatisch generieren lassen.

**Knoten erstellen auf.** Sie können den Knoten zum Zeichenbereich, zur Palette Generierte Modelle oder zu Beiden hinzufügen.

**Interaktiv erstellte Aufbauregeln einschließen.** Zum Einschließen der Aufbauregeln vom aktuellen Baum im generierten Modell, markieren Sie dieses Kontrollkästchen. Auf diese Weise können Sie den Baum bei Bedarf erneut generieren. [Für weitere Informationen siehe Thema Aufbauregeln für Bäume auf S. 148.](#)

### **Aufbauregeln für Bäume**

Bei Modellen vom Typ “C&R-Baum”, “CHAID” und “QUEST” werden durch Aufbauregeln die Bedingungen festgelegt, nach denen der Baum um jeweils eine Ebene erweitert wird. Aufbauregeln werden jedes Mal angewendet, wenn der interaktive Tree Builder vom Knoten aus gestartet wird.

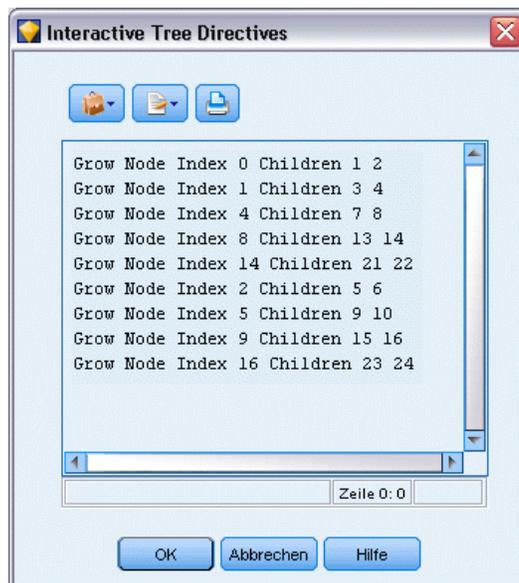
- Aufbauregeln werden meist verwendet, um einen zuvor interaktiv erstellten Baum noch einmal zu generieren. [Für weitere Informationen siehe Thema Aktualisieren der Aufbauregeln auf S. 150.](#) Sie können Aufbauregeln auch manuell bearbeiten. Dabei sollten Sie allerdings vorsichtig sein.
- Die Aufbauregeln sind sehr spezifisch für die Struktur des Baums, den sie beschreiben. Daher kann jede Änderung der zugrunde liegenden Daten oder der Modellierungsoptionen dazu führen, dass ein bislang gültiger Satz an Aufbauregeln Fehler verursacht. Wenn beispielsweise der CHAID-Algorithmus eine zweifache Aufteilung auf der Grundlage aktualisierter Daten in eine dreifache Aufteilung ändert, führen alle zuvor auf der Grundlage der zweifachen Aufteilung erstellten Richtlinien zu Fehlern.

*Hinweis:* Beim direkten Generieren eines Modells (ohne Tree Builder) werden sämtliche Aufbauregeln ignoriert.

### **Bearbeiten von Aufbauregeln**

- ▶ Um gespeicherte Aufbauregeln anzuzeigen oder zu bearbeiten, öffnen Sie den Baumerstellungsknoten und wählen Sie das Panel “Ziel” auf der Registerkarte “Erstellungsoptionen”.
- ▶ Wählen Sie Interaktive Sitzung starten, um die Steuerungen zu aktivieren, wählen Sie Interaktiv erstellte Aufbauregeln verwenden und klicken Sie auf Aufbauregeln.

Abbildung 6-21  
Aufbauregeln für Bäume



### **Syntax der Aufbauregeln**

Aufbaueregeln legen die Bedingungen für die Erweiterung des Baums fest, beginnend mit dem Stammknoten. Um den Baum beispielsweise um eine Ebene zu erweitern:

```
Grow Node Index 0 Children 1 2
```

Da kein Prädiktor angegeben ist, wählt der Algorithmus die beste Aufteilung aus.

Beachten Sie, dass die erste Aufteilung immer im Stammknoten (Index 0) erfolgen muss und dass die Indexwerte für beide untergeordneten Elemente angegeben werden müssen (in diesem Fall 1 und 2). Die Angabe von `Grow Node Index 2 Children 3 4` ist ungültig, wenn Sie nicht zunächst den Stamm erweitert haben, der Knoten 2 erstellt hat.

So erweitern Sie den Baum:

```
Baum erweitern
```

So erweitern und reduzieren Sie den Baum (nur C&R-Baum):

```
Grow_And_Prune Tree
```

So legen Sie eine benutzerdefinierte Aufteilung für einen stetigen Prädiktor fest:

```
Grow Node Index 0 Children 1 2 Spliton
( "EDUCATE", Interval ( NegativeInfinity, 12.5)
  Interval ( 12.5, Infinity ))
```

So nehmen Sie eine Aufteilung eines nominalen Prädiktors mit zwei Werten vor:

```
Grow Node Index 2 Children 3 4 Spliton
( "GENDER", Group( "0.0" )Group( "1.0" ))
```

So nehmen Sie eine Aufteilung eines nominalen Prädiktors mit mehreren Werten vor:

```
Grow Node Index 6 Children 7 8 Spliton
( "ORGS", Group( "2.0","4.0" )
  Group( "0.0","1.0","3.0","6.0" ))
```

So nehmen Sie eine Aufteilung eines ordinalen Prädiktors vor:

```
Grow Node Index 4 Children 5 6 Spliton
( "CHILDS", Interval ( NegativeInfinity, 1.0)
  Interval ( 1.0, Infinity ))
```

*Hinweis:* Bei der Angabe benutzerdefinierter Aufteilungen wird bei Feldnamen und Werten (EDUCATE, GENDER, CHILDS usw.) zwischen Groß- und Kleinschreibung unterschieden.

### **Aufbaueregeln für CHAID-Bäume**

Aufbaueregeln für CHAID-Bäume reagieren besonders empfindlich auf Änderungen der Daten oder des Modells, da sie, anders als "C&R-Baum" und "QUEST" nicht auf die Verwendung binärer Aufteilungen beschränkt sind. Die folgende Syntax sieht beispielsweise völlig korrekt

aus, führt aber zu einem Fehler, wenn der Algorithmus den Stammknoten in mehr als zwei untergeordnete Elemente aufteilt:

```
Grow Node Index 0 Children 1 2
```

```
Grow Node Index 1 Children 3 4
```

Bei CHAID ist es möglich, dass der Knoten 0 auch 3 oder 4 untergeordnete Elemente besitzt, was dazu führt, dass die zweite Zeile einen Fehler verursacht.

### **Verwenden von Aufbauregeln in Skripten**

Aufbauregeln können auch in Skripten eingebettet werden, indem sie in dreifache Anführungszeichen eingeschlossen werden. [Für weitere Informationen siehe Thema Blöcke mit Literaltext in Kapitel 3 in IBM SPSS Modeler 14.2 – Handbuch für die Skripterstellung und Automatisierung.](#)

### **Aktualisieren der Aufbauregeln**

Um Ihre Ergebnisse einer interaktiven Baumerstellung zu erhalten, können Sie die für das Generieren des aktuellen Baums aufgestellten Aufbauregeln speichern. So können Sie den Baum später neu generieren und bearbeiten, was bei einem Modell-Nugget nicht mehr möglich ist.

- ▶ Um Richtlinien zu aktualisieren, wählen Sie in den Tree Builder-Menüs folgende Optionen:  
Datei > Aufbauregeln aktualisieren

Aufbauregeln werden in dem Modellierungsknoten gespeichert, mit dem der Baum erstellt wurde (“C&R-Baum”, “QUEST” oder “CHAID”) und können dazu genutzt werden, den aktuellen Baum noch einmal zu generieren. [Für weitere Informationen siehe Thema Aufbauregeln für Bäume auf S. 148.](#)

### **Exportieren der Modell-, Gewinn- und Risikoinformationen**

Über Tree Builder können Sie Modell-, Gewinn- und Risikostatistiken als Text-, HTML- oder Bildformate exportieren.

- ▶ Wählen Sie im Tree Builder-Fenster die Registerkarte oder die Ansicht, die Sie exportieren wollen.
- ▶ Wählen Sie die folgenden Befehle aus den Menüs aus:  
Datei > Exportieren
- ▶ Wählen Sie Text, HTML oder Diagramm sowie die Elemente aus, die Sie aus dem Untermenü exportieren wollen.

Soweit zutreffend, wird der Export auf der Grundlage der aktuellen Auswahl durchgeführt.

**Exportieren von Text- oder HTML-Formaten.** Sie können Gewinn- oder Risikostatistiken für die Trainings- oder die Testpartition (sofern definiert) exportieren. Der Export erfolgt auf der Grundlage der aktuell auf der Registerkarte “Gewinne” getroffenen Auswahl. So können Sie beispielsweise knotenweise, kumulative oder quantile Statistiken auswählen.

**Exportieren von Grafiken.** Sie können den aktuellen Baum so exportieren, wie er auf der Registerkarte “Viewer” angezeigt wird. Oder Sie exportieren die Gewinn diagramme für die Trainings- oder Testpartition (sofern definiert). Zu den verfügbaren Formaten gehören *.JPEG*, *.PNG* und *.BMP*. Bei Gewinnen basiert der Export auf der aktuell auf der Registerkarte “Gewinne” (nur verfügbar, wenn ein Diagramm angezeigt wird) vorgenommenen Auswahl.

### **Generieren von Filter- und Auswahlknoten**

- ▶ Im Tree Builder-Fenster oder beim Durchsuchen eines Modell-Nuggets für ein Entscheidungsbaummodell wählen Sie in den Menüs folgende Optionen:  
Erzeugen > Filterknoten

*oder*

> Auswahlknoten

**Filterknoten.** Generiert einen Knoten, der alle vom aktuellen Baum nicht benutzten Felder herausfiltert. Mit diesem Verfahren schränken Sie das Daten-Set schnell ein, sodass es nur noch solche Felder enthält, die vom Algorithmus als wichtige Felder ausgewählt werden. Wenn oberhalb dieses Entscheidungsbaumknotens ein Typknoten liegt, leitet das Modell-Nugget des Filtermodells alle Felder mit der Rolle *Ziel* weiter.

**Auswahlknoten.** Generiert einen Knoten, der alle Datensätze auswählt, die in den aktuellen Knoten fallen. Diese Option setzt voraus, dass auf der Registerkarte “Viewer” eine oder mehrere Verzweigungen ausgewählt sind.

Das Modell-Nugget wird im Stream-Zeichenbereich abgelegt.

### **Generieren einer Regelmenge aus einem Entscheidungsbaum**

Sie können ein Modell-Nugget vom Typ “Regelmenge” generieren, das die Baumstruktur als Menge von Regeln darstellt, mit denen die Endzweigungen des Baums definiert werden. Regelmengen enthalten meist die wichtigsten Informationen eines gesamten Entscheidungsbaums, allerdings mit einem weniger komplexen Modell. Der wichtigste Unterschied besteht darin, dass es bei einer Regelmenge möglich ist, dass für einen bestimmten Datensatz mehr als eine oder aber überhaupt keine Regel gilt. Sie können beispielsweise alle Regeln nehmen, die als Ergebnis *Nein* vorhersagen, und dann alle, die *Ja* vorhersagen. Wenn mehrere Regeln gelten, dann wird jeder Regel ein gewichtetes “Votum” zugeordnet, das auf der dieser Regel zugeordneten Konfidenz basiert, und die endgültige Vorhersage ergibt sich aus der Kombination der gewichteten Voten aller für den fraglichen Datensatz geltenden Regeln. Wenn keine Regel gilt, wird dem Datensatz eine Standardvorhersage zugeordnet.

Regelmengen können nur aus Bäumen mit kategorialen Zielfeldern generiert werden (nicht aus Regressionsbäumen).

- ▶ Im Tree Builder-Fenster oder beim Durchsuchen eines Modell-Nuggets für ein Entscheidungsbaummodell wählen Sie in den Menüs folgende Optionen:  
Erzeugen > Regelsatz

Abbildung 6-22  
Dialogfeld "Regelmenge generieren"



**Name der Regelmenge.** Ermöglicht die Angabe eines Namens des neuen Modell-Nuggets vom Typ "Regelmenge".

**Knoten erstellen auf.** Steuert den Standort des neuen Modell-Nuggets vom Typ "Regelmenge". Wählen Sie Zeichenbereich, Generierte Modelle oder Beides aus.

**Mindestzahl Instanzen.** Legen Sie fest, wie viele Instanzen (Anzahl der Datensätze, für die die Regel gilt) im Modell-Nugget vom Typ "Regelmenge" mindestens beibehalten werden sollen. Regeln, deren Unterstützung unter dem angegebenen Wert liegt, werden nicht in die neue Regelmenge aufgenommen.

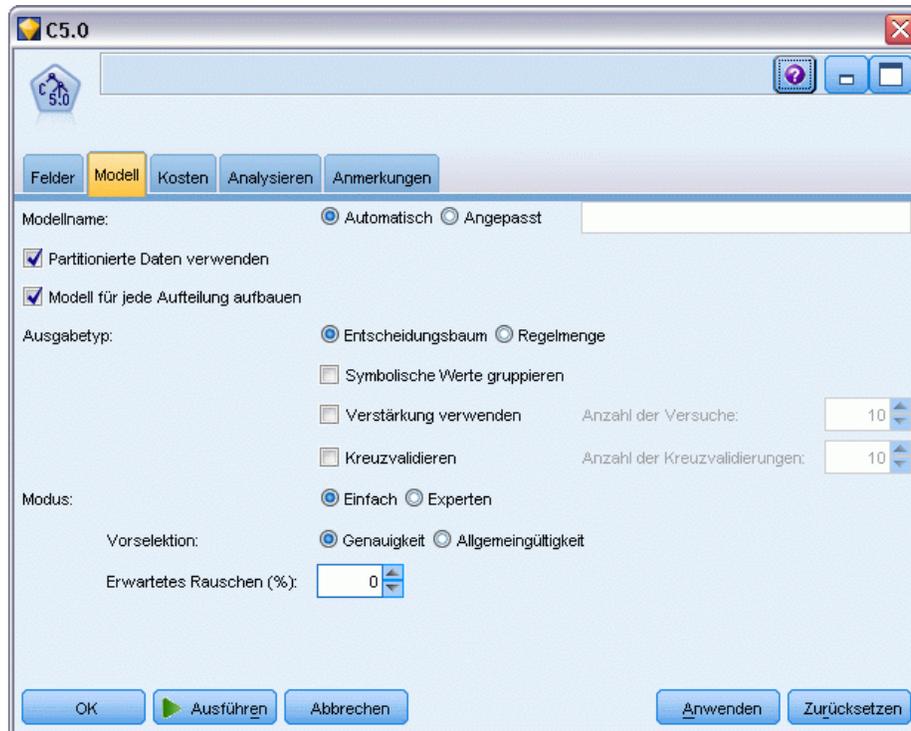
**Minimale Konfidenz.** Geben Sie die minimale Konfidenz für Regeln an, die im Modell-Nugget vom Typ "Regelmenge" erhalten bleiben sollen. Regeln, deren Konfidenz unter dem angegebenen Wert liegt, werden nicht in die neue Regelmenge aufgenommen.

## ***Direktes Erstellen eines Baummodells***

Alternativ zum interaktiven Tree Builder können Sie ein Entscheidungsbaummodell auch direkt aus dem Knoten erstellen, sobald der Stream ausgeführt wird. Dies ist konsistent mit anderen Modellerstellungsknoten. Für C5.0-Baummodelle, die nicht vom interaktiven Tree Builder unterstützt werden, kann ausschließlich diese Methode genutzt werden.

- ▶ Erstellen Sie einen Stream und fügen Sie einen der Entscheidungsbaumerstellungsknoten "C&R-Baum", "CHAID", "QUEST" bzw. "C5.0" hinzu.

Abbildung 6-23  
Direktes Erstellen eines C5.0-Baums



- ▶ Wählen Sie für “C&R-Baum”, “QUEST” oder “CHAID” auf dem Panel “Ziel” der Registerkarte “Erstellungsoptionen” eines der Hauptziele aus. Wenn Sie “Einzelnen Baum aufbauen” wählen, stellen Sie sicher, dass der Modus auf Modell erzeugen gesetzt ist.

Für C5.0 setzen Sie auf der Registerkarte “Modell” die Option Ausgabetyyp auf Entscheidungsbaum.

- ▶ Wählen Sie die Ziel- und Prädiktorfelder aus und legen Sie die zusätzlich benötigten Modelloptionen fest. Spezifische Anleitungen finden Sie in der Dokumentation zu den einzelnen Baumerstellungsknoten.
- ▶ Führen Sie den Stream aus, damit das Modell generiert wird.

### **Kommentare**

- Wenn Sie Bäume nach dieser Methode generieren, werden die Aufbauregeln für Bäume ignoriert.
- Egal ob interaktiv oder direkt, beide Methoden zur Erstellung von Entscheidungsbäumen führen zu ähnlichen Modellen. Die Frage ist lediglich, wie viel Kontrolle Sie während des Ablaufs ausüben wollen.

## Entscheidungsbaumknoten

Die Entscheidungsbaumknoten in IBM® SPSS® Modeler gewähren Zugriff auf die früher eingeführten Baumerstellungsalgorithmen:

- C&R-Baum
- QUEST
- CHAID
- C5.0

Für weitere Informationen siehe Thema Entscheidungsbaum-Modelle auf S. 126.

Die Algorithmen gleichen sich dahingehend, dass alle einen Entscheidungsbaum aufbauen können, indem sie rekursiv die Daten in immer kleinere Untergruppen aufteilen. Es gibt jedoch einige entscheidende Unterschiede.

**Eingabefelder.** Die Eingabefelder (Prädiktoren) können folgende Typen (Messniveaus) sein: stetig, kategorial, Flag, nominal oder ordinal. Für weitere Informationen siehe Thema Messniveaus in Kapitel 4 in *IBM SPSS Modeler 14.2- Quellen-, Prozess- und Ausgabeknoten*.

**Zielfelder.** Es kann nur ein Zielfeld angegeben werden. Bei “C&R-Baum” und “CHAID” kann das Ziel stetig, kategorial, ein Flag, nominal oder ordinal sein. Bei “QUEST” kann es kategorial, ein Flag oder nominal sein. Bei “C5.0” kann das Ziel ein Flag, nominal oder ordinal sein.

**Art der Aufteilung.** C&R-Baum und QUEST unterstützen nur binäre Aufteilungen (das heißt, jeder Knoten des Baums kann nur in zwei Verzweigungen aufgeteilt werden). Dagegen unterstützen CHAID und C5.0 die Aufteilung in mehr als zwei Verzweigungen gleichzeitig.

**Für die Aufteilung verwendetes Verfahren.** Die Algorithmen unterscheiden sich in den Kriterien, die für die Aufteilungsentscheidung verwendet werden. Bei Vorhersage einer kategorialen Ausgabe in C&R-Baum mit ein Streuungsmaß verwendet (standardmäßig ist das der Gini-Koeffizient, Sie können das jedoch ändern). Für stetige Ziele wird das Verfahren der kleinsten quadratischen Abweichung verwendet. CHAID verwendet den Chi-Quadrat-Test; bei QUEST wird ein Chi-Quadrat-Test für kategoriale Prädiktoren verwendet und die Varianzanalyse bei stetigen Eingaben. Bei C5.0 wird ein Informationstheorie-Maß verwendet, das Informationsgewinnverhältnis.

**Behandlung fehlender Werte.** Alle Algorithmen lassen fehlende Werte in den Prädiktorfeldern zu, wenden jedoch unterschiedliche Verfahrensweisen für deren Behandlung an. C&R-Baum und QUEST verwenden bei Bedarf Ersatz-Vorhersagefelder, um einen Datensatz mit fehlenden Werten beim Training weiterzuleiten. CHAID erstellt für die fehlenden Werte eine separate Kategorie und lässt sie zur Verwendung bei der Baumerstellung zu. Bei C5.0 wird ein Fraktionierungsverfahren eingesetzt, das einen Bruchbereich eines Datensatzes entlang jeder Verzweigung des Baums von einem Knoten weiterreicht, bei dem die Aufteilung auf einem Feld mit fehlendem Wert basiert.

**Beschneidung.** Bei C&R-Baum, QUEST und C5.0 gibt es die Option, den Baum vollständig aufzubauen und ihn anschließend durch Entfernen von Aufteilungen der untersten Ebene zu beschneiden, die keine signifikante Auswirkung auf die Genauigkeit des Baums besitzen. Bei

allen Entscheidungsbaumalgorithmen ist es jedoch möglich, die Mindestgröße der Untergruppen zu steuern, was dazu beiträgt, Verzweigungen mit wenigen Datensätzen zu vermeiden.

**Interaktive Baumerstellung.** C&R-Baum, QUEST und CHAID bieten die Option, eine interaktive Sitzung zu starten. So ist es möglich, den Baum jeweils auf einer Ebene zu erstellen, die Aufteilungen zu bearbeiten und den Baum zu reduzieren, bevor das Modell erstellt wird. C5.0 bietet keine Option zur interaktiven Bearbeitung.

**A-priori-Wahrscheinlichkeiten.** C&R-Baum und QUEST unterstützen die Spezifikation von A-priori-Wahrscheinlichkeiten für Kategorien bei der Vorhersage eines kategorialen Zielfelds. A-priori-Wahrscheinlichkeiten sind Schätzungen der gesamten relativen Häufigkeit für jede Zielkategorie in der Gesamtheit, aus der die Trainingsdaten gezogen werden. Mit anderen Worten: Es handelt sich um die Wahrscheinlichkeitsschätzungen, die Sie für jeden möglichen Zielwert vornehmen würden, bevor Sie etwas über die Prädiktorwerte wissen. CHAID und C5.0 unterstützen die Spezifizierung von A-priori-Wahrscheinlichkeiten nicht.

**Regelmengen.** Für Modelle mit kategorialen Zielfeldern bieten die Entscheidungsbaumknoten die Option, das Modell in Form einer Regelmenge zu erstellen, was bei der Interpretation teilweise einfacher ist als ein komplexer Entscheidungsbaum. Bei C&R-Baum, QUEST und CHAID lässt sich eine Regelmenge in einer interaktiven Sitzung erstellen; bei C5.0 kann diese Option im Modellierungsknoten spezifiziert werden. Außerdem ist es bei allen Entscheidungsbaummodellen möglich, eine Regelmenge im Modell-Nugget zu erstellen. [Für weitere Informationen siehe Thema Generieren einer Regelmenge aus einem Entscheidungsbaum auf S. 151.](#)

## **C&R-Baumknoten**

Der Klassifizierungs- und Regressionsbaumknoten (C&R) bildet eine baumbasierte Klassifizierungs- und Vorhersagemethode. Ähnlich wie C5.0 verwendet diese Methode eine rekursive Partitionierung, um die Trainingsdatensätze in Segmente mit ähnlichen Ausgabefeldwerten aufzuteilen. Der Knoten vom Typ "C&R-Baum" beginnt mit der Untersuchung der Eingabefelder, um die beste Aufteilung zu finden, die anhand der Reduktion in einem aus der Aufteilung resultierenden Unreinheitsindex gemessen wird. Die Aufteilung definiert zwei Untergruppen, die anschließend beide in zwei weitere Untergruppen aufgeteilt werden. Dies wird so lange fortgesetzt, bis die Grenzkriterien erreicht sind. Alle Aufteilungen erfolgen binär (nur zwei Untergruppen).

### **Beschneidung**

Bei C&RT-Bäumen haben Sie die Option, den Baum zuerst zu erweitern und dann auf der Grundlage eines Kostenkomplexitätsalgorithmus, der die Risikoschätzung auf der Grundlage der Anzahl der Endknoten anpasst, zu reduzieren. Diese Methode, die eine große Erweiterung des Baums zulässt, bevor dieser nach komplexeren Kriterien reduziert wird, führt zu kleineren Bäumen mit besseren Kreuzvalidierungseigenschaften. Wenn die Anzahl der Endknoten vergrößert wird, verringert dies in der Regel das Risiko für die aktuellen (Trainings-)Daten. Das tatsächliche Risiko kann aber größer sein, wenn das Modell auf unbekannte Daten generalisiert wird. Angenommen es liegt der Extremfall vor, dass Sie für jeden im Trainings-Daten-Set vorhandenen Datensatz einen separaten Endknoten besitzen. Das geschätzte Risiko liegt bei 0 %, da jeder Datensatz in einen

eigenen Knoten fällt, das Fehlklassifizierungsrisiko für unbekannte (Test-)Daten liegt allerdings mit großer Sicherheit über 0 %. Die Kostenkomplexitätsmessung versucht, dies zu kompensieren.

**Beispiel.** Ein Kabelfernsehunternehmen hat eine Marketingstudie in Auftrag gegeben, um zu ermitteln, welche Kunden ein Abonnement für einen interaktiven Nachrichtenservice über Kabel erwerben würden. Mithilfe der Daten aus der Studie können Sie einen Stream erstellen, in dem das Zielfeld die Absicht angibt, das Abonnement zu erwerben, und in dem als Prädiktorfelder Alter, Geschlecht, Bildung, Einkommenskategorie, wöchentlicher Fernsehkonsum und Anzahl der Kinder verwendet werden. Durch Anwendung eines Knotens vom Typ “C&R-Baum” der Kinder können Sie die Antworten vorhersagen und klassifizieren, um eine möglichst hohe Trefferquote für die Kampagne zu erzielen.

**Anforderungen.** Um ein C&R-Baummodell zu trainieren, benötigen Sie mindestens ein *Eingabe*-Feld und genau ein *Ziel*-Feld. Ziel- und Eingabefelder können stetig (in einem numerischen Bereich) oder kategorial sein. Felder, die auf *Beides* oder *Keine* gesetzt sind, werden ignoriert. Die Typen der im Modell verwendeten Felder müssen vollständig als Instanz generiert sein und alle im Modell verwendeten Ordinalfelder (sortiertes Set) müssen numerisch (und nicht als Zeichenkette) gespeichert sein. Im Bedarfsfall können Sie die Felder mit dem Umkodierungsknoten konvertieren. [Für weitere Informationen siehe Thema Umkodierungsknoten in Kapitel 4 in IBM SPSS Modeler 14.2- Quellen-, Prozess- und Ausgabeknoten.](#)

**Stärken.** C&R-Baummodelle verhalten sich bei Problemen mit fehlenden Daten und einer großen Feldzahl sehr robust. Sie benötigen für die Schätzung in der Regel keine langen Trainingsphasen. Darüber hinaus sind C&R-Baummodelle tendenziell leichter verständlich als andere Modelltypen — die aus dem Modell abgeleiteten Regeln lassen sich sehr direkt interpretieren. Im Gegensatz zu C5.0 können C&R-Bäume stetige genauso wie kategoriale Ausgabefelder verarbeiten.

## CHAID-Knoten

CHAID (Chi-squared Automatic Interaction Detection) ist eine Klassifizierungsmethode für die Erstellung von Entscheidungsbäumen mit Chi-Quadrat-Statistiken zur Identifizierung der optimalen Splits.

CHAID untersucht zuerst die zwischen allen Eingabefeldern und dem Ergebnis vorhandenen Kreuztabellen und testet die Signifikanz mit einem Chi-Quadrat-Unabhängigkeitstest. Wenn mehr als eine dieser Beziehungen statistisch signifikant ist, wählt CHAID das signifikanteste Eingabefeld aus (kleinster *P*-Wert). Wenn eine Eingabe mehr als zwei Kategorien besitzt, werden diese verglichen und solche Kategorien gegeneinander reduziert, deren Ergebnis keinen Unterschied aufweist. Dies erfolgt, indem sukzessive alle Kategorienpaare mit dem am wenigsten signifikanten Unterschied verbunden werden. Diese Kategoriezusammenführung wird gestoppt, wenn die Abweichung aller verbleibenden Kategorien das angegebene Testniveau erreicht hat. Bei nominalen Eingabefeldern können alle Kategorien zusammengeführt werden. Bei einem ordinalen Set können nur zusammenhängende Kategorien zusammengeführt werden.

Exhaustive CHAID ist eine Änderung von CHAID, die noch gründlicher vorgeht, indem sie alle für jeden Prädiktor möglichen Aufteilungen untersucht, allerdings mehr Rechenzeit beansprucht.

**Anforderungen.** Ziel- und Eingabefelder können stetig oder kategorial sein. Knoten können auf jeder Ebene in zwei oder mehr Untergruppen aufgeteilt werden. Alle im Modell verwendeten ordinalen Felder müssen numerisch (nicht als Zeichenkette) gespeichert sein. Im Bedarfsfall

können Sie die Felder mit dem Umkodierungsknoten konvertieren. [Für weitere Informationen siehe Thema Umkodierungsknoten in Kapitel 4 in IBM SPSS Modeler 14.2- Quellen-, Prozess- und Ausgabeknoten.](#)

**Stärken.** Im Gegensatz zu den Knoten vom Typ “C&R-Baum” und “QUEST” kann CHAID nichtbinäre Bäume generieren, d. h. Bäume mit Aufteilungen mit mehr als zwei Verzweigungen. CHAID erstellt daher tendenziell breitere Bäume als die binären Erweiterungsmethoden. CHAID funktioniert mit allen Eingaben und akzeptiert sowohl Fallgewichtungs- als auch Häufigkeitsvariablen.

## **QUEST-Knoten**

QUEST (Quick, Unbiased, Efficient Statistical Tree) ist eine binäre Klassifizierungsmethode zum Erstellen von Entscheidungsbäumen. Diese Methode wurde primär in der Absicht entwickelt, die Verarbeitungszeit zu verkürzen, die für große C&RT-Analysen mit vielen Variablen oder mit vielen Fällen benötigt wird. Ein zweites Ziel von QUEST war die Senkung der in den Klassifizierungsbaummodellen festgestellten Tendenz, Eingaben zu bevorzugen, die mehr Aufteilungen erlauben. Dabei handelt es sich um stetige Eingabefelder (numerischer Bereich) oder um solche mit vielen Kategorien.

- QUEST verwendet eine Folge von auf signifikanten Tests basierenden Regeln, um die im Knoten vorhandenen Eingabefelder zu bewerten. Zu Auswahlzwecken muss gegebenenfalls für jede in einem Knoten vorhandene Eingabe nur ein einziger Test durchgeführt werden. Im Gegensatz zu “C&R-Baum” werden nicht alle Aufteilungen untersucht. Und im Gegensatz zu “C&R-Baum” und “CHAID” werden beim Bewerten eines Eingabefelds für die Auswahl die Kategoriekombinationen nicht getestet. Dies beschleunigt die Analyse.
- Aufteilungen werden festgelegt, indem eine quadratische Diskriminanzanalyse durchgeführt wird, die die ausgewählte Eingabe für Gruppen verwendet, die durch die Zielkategorien gebildet werden. Diese Methode führt gegenüber einer erschöpfenden Suche (C&R-Baum) wiederum zu einer Steigerung der Geschwindigkeit bei der Bestimmung der optimalen Aufteilung.

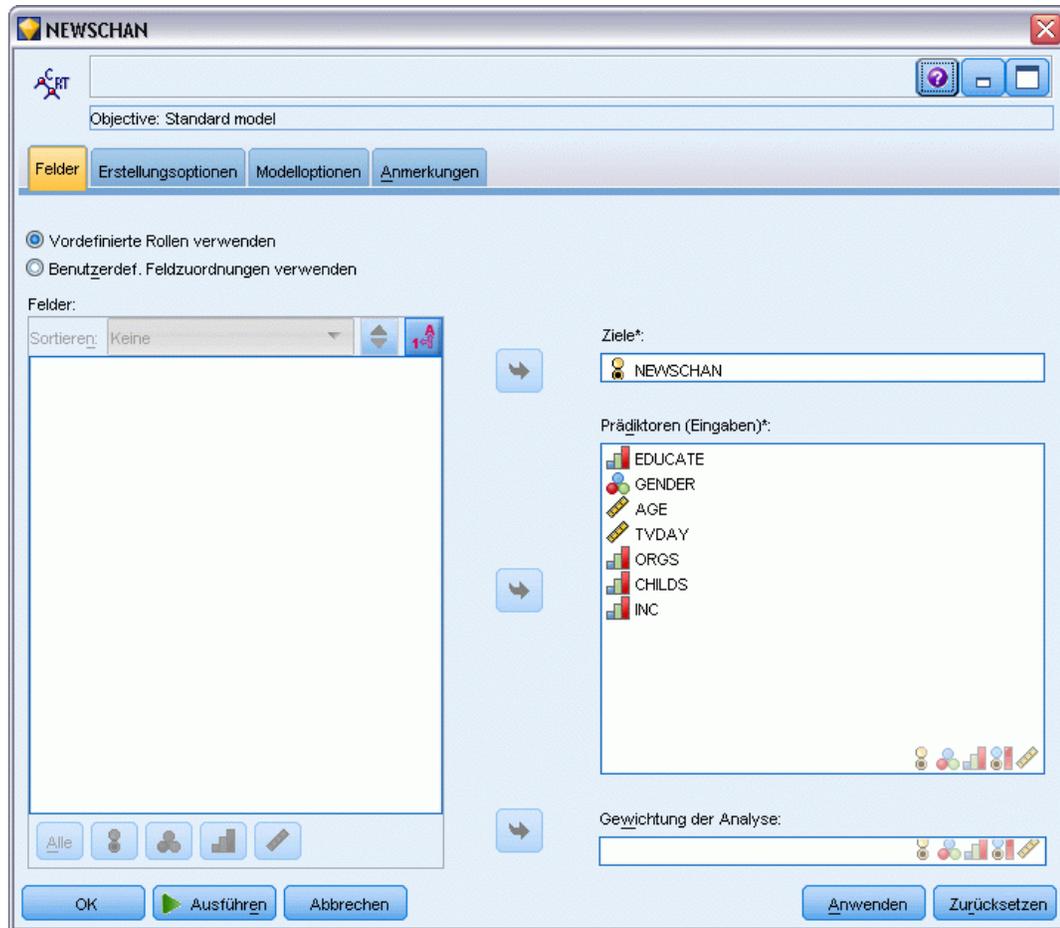
**Anforderungen.** Eingabefelder können stetig (numerische Bereiche) sein, das Zielfeld muss aber kategorial sein. Alle Aufteilungen sind binär. Gewichtungsfelder können nicht eingesetzt werden. Alle im Modell verwendeten ordinalen Felder (sortiertes Set) müssen numerisch (nicht als Zeichenkette) gespeichert sein. Im Bedarfsfall können Sie die Felder mit dem Umkodierungsknoten konvertieren. [Für weitere Informationen siehe Thema Umkodierungsknoten in Kapitel 4 in IBM SPSS Modeler 14.2- Quellen-, Prozess- und Ausgabeknoten.](#)

**Stärken.** Genau wie “CHAID”, aber im Gegensatz zu “C&R-Baum”, verwendet “QUEST” statistische Tests, um zu entscheiden, ob ein Eingabefeld benutzt wird. Das Verfahren trennt auch die Eingabeauswahl von der Aufteilung und verwendet jeweils unterschiedliche Kriterien. Dies stellt einen Unterschied zu CHAID dar, wo das statistische Testergebnis, das die Variablenauswahl bestimmt, auch die Aufteilung erzeugt. “C&R-Baum” verfährt ähnlich, indem die Messung der Unreinheitsänderung sowohl die Auswahl des Eingabefelds als auch die Aufteilung bestimmt.

## Entscheidungsbaumknoten - Feldoptionen

Auf der Registerkarte “Felder” geben Sie an, ob Sie die Feldrolleneinstellungen verwenden möchten, die bereits in den Knoten oberhalb definiert wurden, oder ob Sie die Feldzuweisungen manuell vornehmen möchten.

Abbildung 6-24  
Knoten “C&R”, Registerkarte “Felder”



**Vordefinierte Rollen verwenden** Bei dieser Option werden die Rolleneinstellungen (Ziele, Prädiktoren usw.) aus einem Typknoten oberhalb verwendet (oder die Registerkarte “Typen” eines weiter oben im Stream gelegenen Quellenknotens). [Für weitere Informationen siehe Thema Festlegen der Feldrolle in Kapitel 4 in IBM SPSS Modeler 14.2- Quellen-, Prozess- und Ausgabeknoten.](#)

**Benutzerdefinierte Feldzuweisungen verwenden** Wählen Sie diese Option, wenn Sie die Ziele, Prädiktoren und andere Rollen in diesem Bildschirm manuell zuweisen möchten.

**Felder.** Verwenden Sie die Pfeiltasten, um den verschiedenen Rollenfeldern rechts im Bildschirm Objekte aus dieser Liste manuell zuzuweisen. Die Symbole zeigen das gültige Messniveau für jedes Rollenfeld an.

Klicken Sie auf die Schaltfläche **Alle**, damit alle Felder in der Liste ausgewählt werden, oder klicken Sie auf die Schaltfläche für ein individuelles Messniveau, um alle Felder mit diesem Messniveau auszuwählen.

**Ziel.** Wählen Sie ein Feld als Ziel für die Vorhersage aus.

**Prädiktoren (Eingaben).** Wählen Sie ein oder mehrere Felder als Eingaben für die Vorhersage aus.

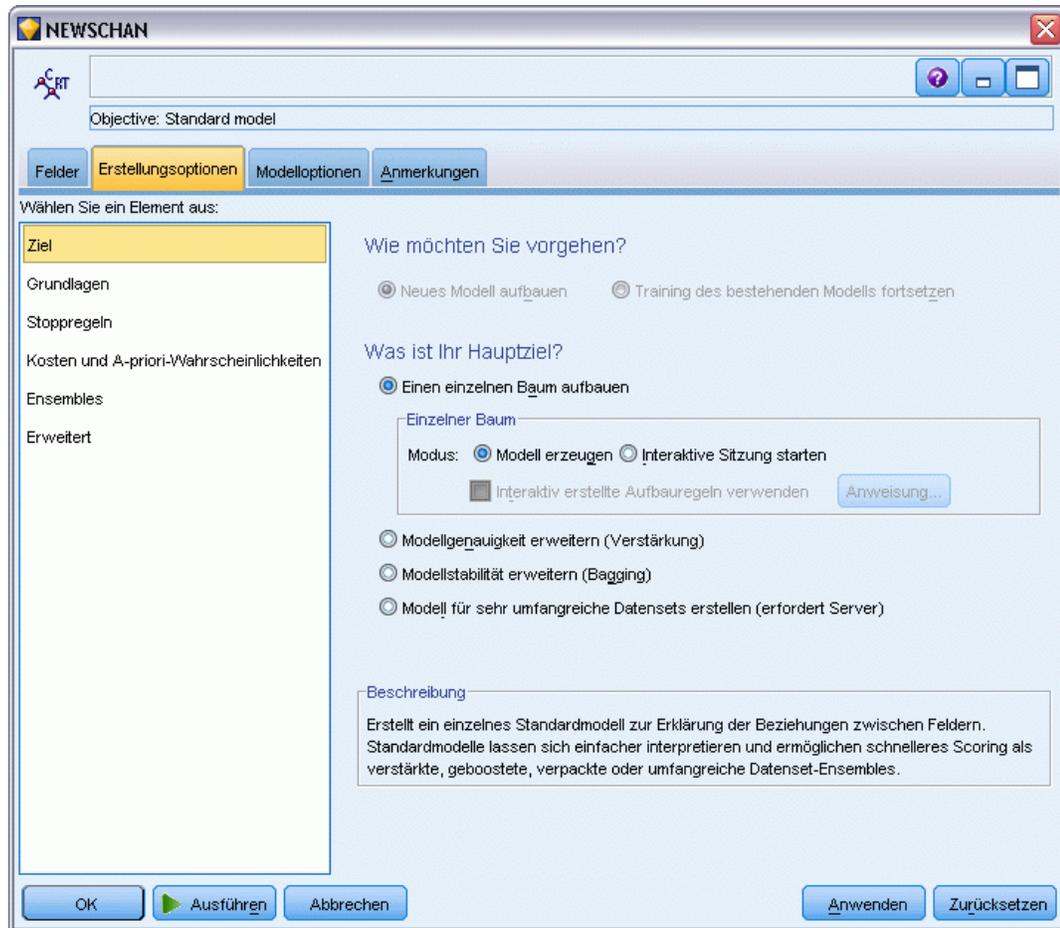
**Analysegewichtung.** (nur CHAID und C&R-Baum) Spezifizieren Sie das Feld hier, um es als Fallgewichtung zu verwenden. Fallgewichtungen werden verwendet, um Differenzen in der Varianz zwischen den Ebenen des Ausgabefelds zu berücksichtigen. [Für weitere Informationen siehe Thema Verwenden von Häufigkeits- und Gewichtungsfeldern in Kapitel 3 auf S. 39.](#)

### ***Entscheidungsbaumknoten - Erstellungsoptionen***

Auf der Registerkarte "Erstellungsoptionen" legen Sie sämtliche Optionen zum Erstellen des Modells fest. Selbstverständlich können Sie auch einfach auf die Schaltfläche **Ausführen** klicken, um ein Modell mit allen Standardoptionen zu erstellen, doch in der Regel werden Sie Einstellungen zur Anpassung an Ihre eigenen Zwecke vornehmen wollen.

Hier können Sie wählen, ob Sie ein neues Modell erstellen oder ein vorhandenes aktualisieren möchten. Außerdem legen Sie das Hauptziel des Knotens fest: Erstellung eines Standardmodells, Erstellung eines Modells mit verbesserter Genauigkeit oder Stabilität oder Erstellung eines Modells für sehr große Daten-Sets.

Abbildung 6-25  
Knoten "C&R, Registerkarte "Erstellungsoptionen"



### Was möchten Sie tun?

**Neues Modell aufbauen.** (Standard) Erstellt jedes Mal ein vollständig neues Modell, wenn Sie einen Stream mit diesem Modellierungsknoten ausführen.

**Training des bestehenden Modells fortsetzen.** In der Standardeinstellung wird bei jeder Ausführung eines Modellierungsknotens ein völlig neues Modell erstellt. Bei Auswahl dieser Option wird das Training mit dem letzten, vom Knoten erfolgreich aufgebauten Modell fortgesetzt. Dies ermöglicht die Aktualisierung eines bestehenden Modells, ohne dass ein Zugriff auf die ursprünglichen Daten erforderlich ist, und kann zu einer wesentlich schnelleren Leistung führen, da *ausschließlich* die neuen bzw. aktualisierten Datensätze in den Stream eingespeist werden. Details zum vorherigen Modell werden zusammen mit dem Modellierungsknoten gespeichert, wodurch diese Option auch dann verwendet werden kann, wenn das vorherige Modell-Nugget nicht mehr im Stream oder in der Modellpalette verfügbar ist.

*Hinweis:* Diese Option wird nur aktiviert, wenn Sie Modell für sehr große Daten-Sets erstellen als Ziel wählen.

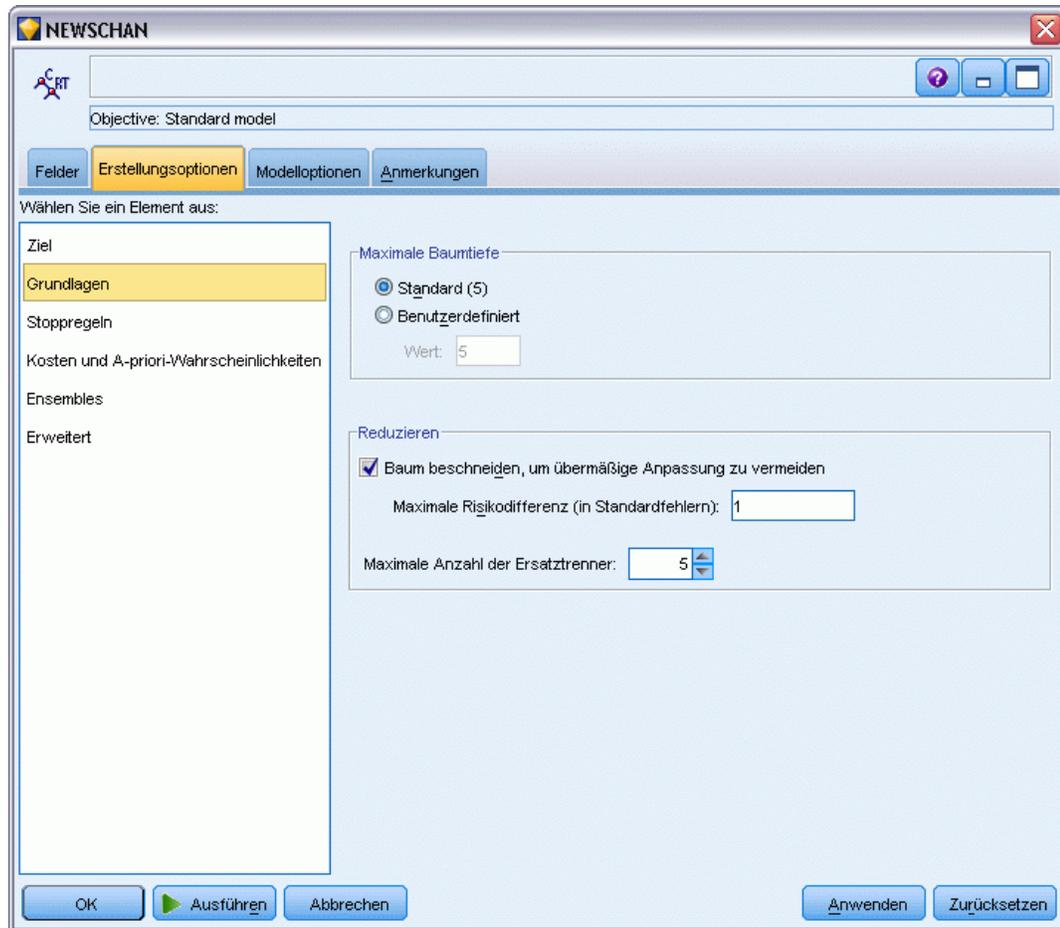
**Wie lautet Ihr Hauptziel?**

- **Einzelnen Baum aufbauen.** Erstellt ein Standardmodell mit individuellem Entscheidungsbaum. Standardmodelle können grundsätzlich einfacher interpretiert und schneller gescort werden, als Modelle, die unter Verwendung der anderen Zieloptionen erstellt werden.  
**Modalwert.** Legt fest, welche Methode für die Modellbildung verwendet wird. Modell erzeugen erstellt ein Modell automatisch, sobald der Stream ausgeführt wird. Interaktive Sitzung starten öffnet den Tree Builder, mit dem Sie Ihr Modell Ebene für Ebene erstellen, Aufteilungen bearbeiten und nach Wunsch reduzieren können, bevor Sie das Modell-Nugget erstellen.  
**Interaktiv erstellte Aufbauregeln verwenden.** Mit dieser Option legen Sie Richtlinien fest, die angewendet werden, wenn aus dem Knoten ein interaktiver Baum generiert wird. Sie können beispielsweise die Aufteilungen der ersten und zweiten Ebene bestimmen, die dann beim Tree Builder-Start automatisch angewendet werden. Richtlinien einer interaktiven Baumerstellung können Sie auch speichern, um den Baum zu einem späteren Zeitpunkt noch einmal zu erstellen. [Für weitere Informationen siehe Thema Aktualisieren der Aufbauregeln auf S. 150.](#)
- **Modellgenauigkeit erhöhen (verbessern).** Mit dieser Option wählen Sie ein spezielles Verfahren, das als **Verbesserung** bekannt ist, um die Modellgenauigkeitsquote zu erhöhen. Die Verbesserung funktioniert so, dass mehrere Modelle in einer Folge erstellt werden. Das erste Modell wird auf die übliche Weise erstellt. Anschließend wird ein zweites Modell erstellt, bei dem besonders die Datensätze berücksichtigt werden, bei denen es im ersten Modell zu Fehlklassifizierungen kam. Das dritte Modell wird in Bezug auf die im zweiten Modell enthaltenen Fehler erstellt usw. Zum Schluss werden die Fälle klassifiziert, indem der gesamte Modellsatz auf ihnen angewendet wird, wobei ein gewichtetes Voting-Verfahren genutzt wird, um die einzelnen Vorhersagen zu einer Gesamtvorhersage zu kombinieren. Eine Verbesserung kann die Genauigkeit eines Entscheidungsbaummodells signifikant verbessern, macht aber auch ein längeres Training notwendig.
- **Modellstabilität steigern (verstärken).** Mit dieser Option wählen Sie ein spezielles Verfahren, das als **Verstärkung** (Bootstrap-Aggregation) bekannt ist, um die Stabilität des Modells zu verbessern und eine Überanpassung zu vermeiden. Mit dieser Option werden mehrere Modelle erstellt und kombiniert, um zuverlässigere Vorhersagen zu erhalten. Mithilfe dieser Option gebildete Modelle benötigen mehr Zeit bei der Erstellung und beim Scoring als Standardmodelle.
- **Modell für sehr große Daten-Sets erstellen.** Wählen Sie diese Option, wenn Sie mit Daten-Sets arbeiten, die zu groß für die Erstellung eines Modells mithilfe der anderen Zieloptionen sind. Diese Option unterteilt die Daten in kleinere Datenblöcke und erstellt auf jedem Block ein Modell. Die genauesten Modelle werden anschließend automatisch ausgewählt und zu einem einzigen Modell-Nugget zusammengefasst. Durch Auswahl der Option Training des bestehenden Modells fortsetzen auf diesem Bildschirm können Sie eine inkrementelle Modellaktualisierung durchführen. *Hinweis:* Diese Option ist für sehr große Daten-Sets gedacht und benötigt eine Verbindung mit IBM® SPSS® Modeler Server. [Für weitere Informationen siehe Thema Verbindung mit IBM SPSS Modeler Server in Kapitel 3 in IBM SPSS Modeler 14.2- Benutzerhandbuch.](#)

**Entscheidungsbaumknoten - Grundeinstellungen**

Hier nehmen Sie die Grundeinstellungen für die Erstellung des Entscheidungsbaums vor.

Abbildung 6-26  
Grundeinstellungen für Entscheidungsbaum



**Algorithmus zur Baumerweiterung.** (nur CHAID) Wählen Sie den CHAID-Algorithmus aus, den Sie verwenden möchten. Exhaustive CHAID ist eine Änderung von CHAID, die noch gründlicher vorgeht, indem sie alle für jeden Prädiktor möglichen Aufteilungen untersucht, allerdings mehr Rechenzeit beansprucht.

**Maximale Baumtiefe.** Legen Sie die maximale Anzahl der Ebenen unter dem Stammknoten fest (wie oft die Stichprobe rekursiv aufgeteilt wird). Die Standardeinstellung ist 5; wählen Sie Benutzerdefiniert und geben Sie einen Wert ein, um eine andere Anzahl von Ebenen festzulegen.

**Reduzieren (nur C&RT und QUEST)**

**Baum reduzieren, um zu große Anpassung zu vermeiden.** Die Reduktion besteht im Entfernen von Aufteilungen der unteren Ebene, die keine signifikante Auswirkung auf die Genauigkeit des Baums besitzen. Durch Reduktion kann der Baum vereinfacht werden, wodurch er leichter zu

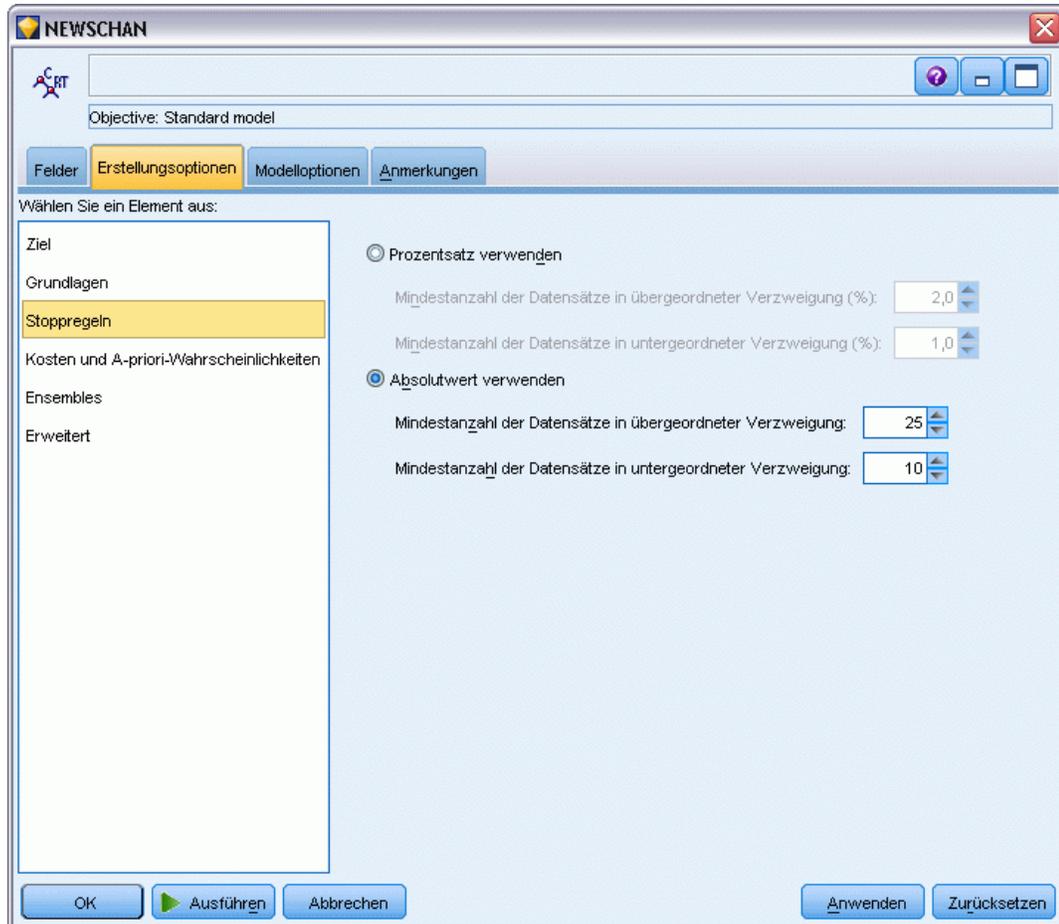
interpretieren ist. In einigen Fällen kann er auch besser generalisiert werden. Lassen Sie diese Option deaktiviert, wenn Sie den vollständigen Baum ohne Reduktion erhalten wollen.

- **Maximale Risikodifferenz (in Standardfehlern).** Damit können Sie eine etwas liberalere Reduktionsregel vorgeben. Die Standardfehlerregel erlaubt dem Algorithmus die Auswahl des einfachsten Baums, dessen Risikoschätzung nahe bei (oder möglichst über) der des untergeordneten Baums mit dem geringsten Risiko liegt. Der Wert gibt die Größe der zulässigen Differenz zwischen der Risikoschätzung für den reduzierten Baum und der des Baums mit dem kleinsten Risiko an. Wenn Sie beispielsweise 2 angeben, kann ein Baum ausgewählt werden, dessen Risikoschätzung ( $2 \times$  Standardfehler) größer ist als die des gesamten Baums.

**Maximale Anzahl Ersatztrenner.** Ersatztrenner bieten ein Verfahren für die Behandlung fehlender Werte. Für jede Aufteilung ermittelt der Algorithmus die Eingabefelder, die dem ausgewählten Aufteilungsfeld am ähnlichsten sind. Diese Felder sind die **Ersatztrenner** für diese Aufteilung. Wenn ein Datensatz klassifiziert werden muss, aber in einem Aufteilungsfeld ein Wert fehlt, kann für die Aufteilung der entsprechende Wert eines Ersatztrennerfelds benutzt werden. Eine höhere Einstellung sorgt für mehr Flexibilität bei der Behandlung fehlender Werte, verursacht allerdings auch eine stärkere Arbeitsspeicherverwendung und längere Trainingszeiten.

## Entscheidungsbaumknoten - Stoppregeln

Abbildung 6-27  
Optionen für Stoppregeln

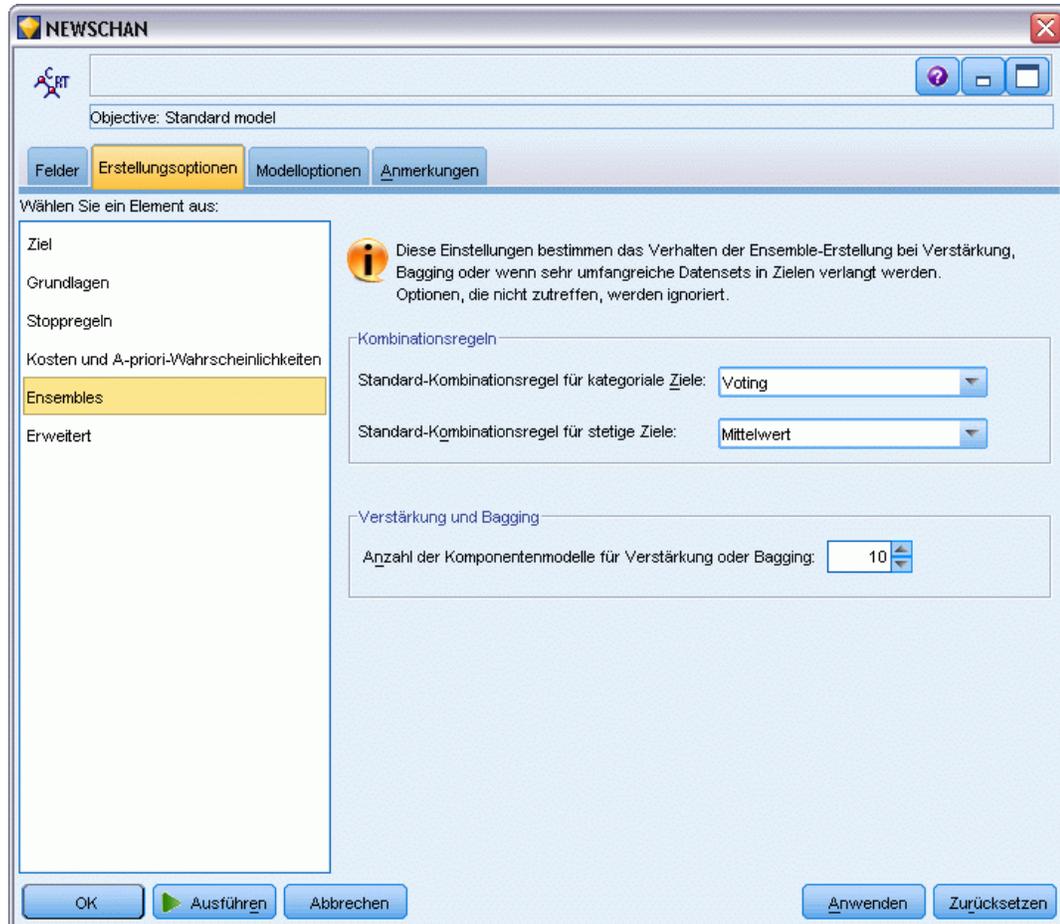


Diese Optionen steuern, wie der Baum konstruiert wird. Grenzregeln legen fest, wann mit dem Aufteilen bestimmter Verzweigungen des Baums aufgehört wird. Damit legen Sie die minimale Verzweigungsgröße fest. Dies verhindert Aufteilungen, die zu sehr kleinen Untergruppen führen. Die Option Mindestanzahl der Datensätze in übergeordneter Verzweigung verhindert eine Aufteilung, wenn die Anzahl der im aufzuteilenden Knoten (die **übergeordnete Verzweigung**) enthaltenen Datensätze geringer ist, als der hier angegebene Wert. Die Option Mindestanzahl der Datensätze in untergeordneter Verzweigung verhindert eine Aufteilung, wenn die Anzahl der in einer durch die Aufteilung erzeugten Verzweigung (die **untergeordnete Verzweigung**) enthaltenen Datensätze geringer ist, als der hier angegebene Wert.

- **Prozentsatz verwenden.** Hiermit können Sie Größen als Prozentsätze der gesamten Trainingsdaten angeben.
- **Absolutwert verwenden.** Hiermit können Sie Größen als absolute Anzahl von Datensätzen angeben.

## Entscheidungsbaumknoten - Ensembles

Abbildung 6-28  
Optionen für Ensembles



Diese Einstellungen legen das Verhalten der Ensemblebildung fest, die erfolgt, wenn auf der Registerkarte “Ziele” die Option “Verbesserung”, “Verstärkung” oder “Sehr große Daten-Sets” ausgewählt ist. Optionen, die für das ausgewählte Ziel nicht gelten, werden ignoriert.

**Bagging und sehr umfangreiche Daten-Sets.** Beim Scoren eines Ensembles wird diese Regel angewendet, um die vorhergesagten Werte aus den Basismodellen für die Berechnung des Score-Werts für das Ensemble zu kombinieren.

- Standard-Kombinierungsregel für kategoriale Ziele.** Ensemble-Vorhersagewerte für kategoriale Ziele können mithilfe von “Voting”, “höchster Wahrscheinlichkeit” oder “höchste mittlere Wahrscheinlichkeit” kombiniert werden. Bei **Voting** wird die Kategorie gewählt, die in allen Basismodellen am häufigsten die höchste Wahrscheinlichkeit erreicht. Bei **Höchste Wahrscheinlichkeit** wird die Kategorie gewählt, die in allen Basismodellen den höchsten Einzelwert bei der höchsten Wahrscheinlichkeit erzielt. Bei **Höchste mittlere**

**Wahrscheinlichkeit** wird die Kategorie gewählt, die den höchsten Wert erreicht, wenn der Mittelwert der Wahrscheinlichkeiten für alle Kategorien in den Basismodellen berechnet wird.

- **Standard-Kombinierungsregel für stetige Ziele.** Ensemble-Vorhersagewerte für stetige Ziele können unter Verwendung des Mittelwerts oder Medians der Vorhersagewerte aus den Basismodellen kombiniert werden.

Hinweis: Wenn als Ziel die Verbesserung der Modellgenauigkeit ausgewählt wurde, wird die Auswahl zum Kombinieren der Regeln ignoriert. Bei der Verbesserung wird für das Scoring der kategorialen Ziele stets eine gewichtete Mehrheit verwendet und für das Scoring stetiger Ziele ein gewichteter Median.

**Verbesserung und Verstärkung.** Geben Sie die Anzahl der zu erstellenden Basismodelle an, wenn als Ziel die Verbesserung der Modellgenauigkeit oder -stabilität angegeben ist. Im Falle der Verstärkung ist das die Anzahl der Bootstrap-Stichproben. Muss eine positive ganze Zahl sein.

### C&R-Baum- und QUEST-Knoten - Kosten & A Priori

Abbildung 6-29  
Einstellung der Fehlklassifizierungskosten und A-priori-Wahrscheinlichkeiten

NEWSCHAN

Objective: Standard model

Felder Erstellungsoptionen Modelloptionen Anmerkungen

Wählen Sie ein Element aus:

- Ziel
- Grundlagen
- Stoppregeln
- Kosten und A-priori-Wahrscheinlichkeiten**
- Ensembles
- Erweitert

Fehlklassifizierungskosten

Fehlklassifizierungskosten verwenden

Vorhergesagt

	0	1
0	0.0	1.0
1	1.0	0.0

Ist

Vorgänger

Auf Trainingsdaten basierend  Für alle Klassen gleich  Benutzerdefiniert

Wert	Wahrscheinlich...
0	0,5
1	0,5

A-prioris mit Fehlklassifizierungskosten anpassen

Normalisieren Gleichsetzen

Anwenden Zurücksetzen

OK **Ausführen** Abbrechen

### **Fehlklassifizierungskosten**

In einigen Zusammenhängen sind bestimmte Fehler kostspieliger als andere. Zum Beispiel kann es kostspieliger sein, einen Kreditantragsteller mit hohem Risiko als niedriges Risiko zu klassifizieren (eine Art von Fehler) als einen Kreditantragsteller mit niedrigem Risiko als hohes Risiko (eine andere Art von Fehler). Anhand von Fehlklassifizierungskosten können Sie die relative Bedeutung verschiedener Arten von Vorhersagefehlern angeben.

Fehlklassifizierungskosten sind im Grunde Gewichtungen, die auf bestimmte Ergebnisse angewendet werden. Diese Gewichtungen werden in das Modell mit einbezogen und können die Vorhersage (als Schutz gegen kostspielige Fehler) ändern.

Mit Ausnahme von C5.0-Modellen werden Fehlklassifizierungskosten beim Scoring von Modellen nicht angewendet und bei der Rangbildung bzw. beim Vergleich von Modellen mithilfe eines Knotens vom Typ "Automatischer Klassifizierer", eines Evaluationsdiagramms oder eines Analyseknosens nicht berücksichtigt. Ein Modell, das Kosten beinhaltet, führt nicht unbedingt zu weniger Fehlern als ein Modell, das keine Kosten beinhaltet, und es weist auch nicht unbedingt eine höhere Gesamtgenauigkeit auf, aber es ist möglicherweise in praktischer Hinsicht leistungsfähiger, da es eine integrierte Verzerrung zugunsten von *weniger teuren* Fehlern aufweist.

Die Kostenmatrix zeigt die Kosten für jede mögliche Kombination aus prognostizierter Kategorie und tatsächlicher Kategorie. Standardmäßig werden alle Fehlklassifizierungskosten auf 1,0 gesetzt. Um benutzerdefinierte Kostenwerte einzugeben, wählen Sie Fehlklassifizierungskosten verwenden und geben Ihre benutzerdefinierten Werte in die Kostenmatrix ein.

Um die Fehlklassifizierungskosten zu ändern, wählen Sie die Zelle aus, die der gewünschten Kombination aus vorhergesagten und tatsächlichen Werten entspricht, löschen den Inhalt der Zelle und geben die gewünschten Kosten für die Zelle ein. Die Kosten sind nicht automatisch symmetrisch. Wenn Sie z. B. die Kosten einer Fehlklassifizierung von *A* als *B* auf 2,0 setzen, weisen die Kosten für eine Fehlklassifizierung von *B* als *A* weiterhin den Standardwert 1,0 auf, es sei denn, Sie ändern diesen Wert ausdrücklich.

### **A-priori-Wahrscheinlichkeiten**

Mit diesen Optionen können Sie die A-priori-Wahrscheinlichkeiten für Kategorien bei einer Vorhersage eines kategorialen Zielfelds angeben. **A-priori-Wahrscheinlichkeiten** sind Schätzungen der gesamten relativen Häufigkeit für jede Zielkategorie in der Gesamtheit, aus der die Trainingsdaten gezogen werden. Mit anderen Worten: Es handelt sich um die Wahrscheinlichkeitsschätzungen, die Sie für jeden möglichen Zielwert vornehmen würden, *bevor* Sie etwas über die Prädiktorwerte wissen. Es gibt drei Methoden, A-priori-Werte festzulegen:

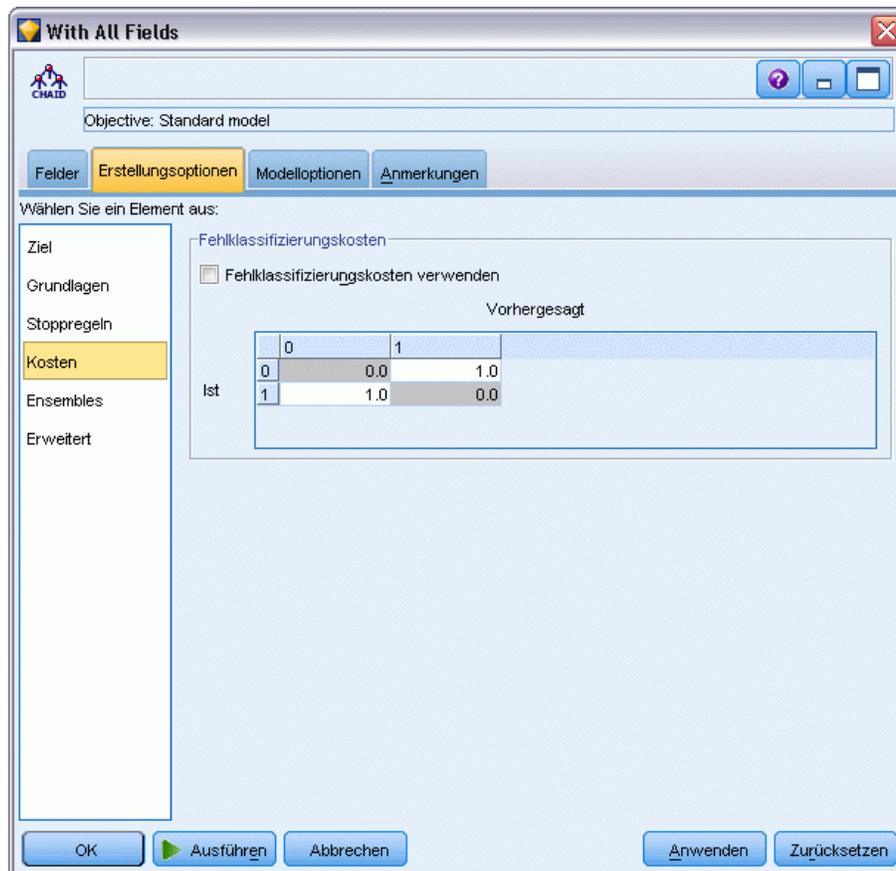
- **Auf Trainingsdaten basierend.** Dies ist die Standardeinstellung. A-priori-Wahrscheinlichkeiten basieren auf den relativen Häufigkeiten der Kategorien in den Trainingsdaten.
- **Für alle Klassen gleich.** A-priori-Wahrscheinlichkeiten für alle Kategorien werden als  $1/k$  definiert, wobei  $k$  die Zahl der Zielkategorien darstellt.
- **Angepasst.** Sie können eigene A-priori-Wahrscheinlichkeiten angeben. Die Startwerte für A-priori-Wahrscheinlichkeiten werden für alle Klassen gleich gesetzt. Sie können die Wahrscheinlichkeiten für einzelne Kategorien auf benutzerdefinierte Werte einstellen. Um die Wahrscheinlichkeit einer bestimmten Kategorie anzupassen, wählen Sie die Wahrscheinlichkeitszelle in der Tabelle aus, die der gewünschten Kategorie entspricht, löschen den Inhalt der Zelle und geben den gewünschten Wert ein.

Die A-priori-Wahrscheinlichkeiten für alle Kategorien sollten sich auf 1,0 summieren (die **Wahrscheinlichkeitsbeschränkung**). Wenn sie keine Summe von 1,0 bilden, wird eine Warnmeldung ausgegeben und es besteht die Möglichkeit, die Werte automatisch normalisieren zu lassen. Diese automatische Anpassung behält die Anteile über die Kategorien hinweg bei, während die Wahrscheinlichkeitsbeschränkung erzwungen wird. Sie können diese Anpassung jederzeit durchführen, indem Sie auf die Schaltfläche Normalisieren klicken. Um die Tabelle auf gleiche Werte für alle Kategorien zurückzusetzen, klicken Sie auf die Schaltfläche Gleichsetzen.

**A-priori-Wahrscheinlichkeiten anhand der Fehlklassifizierungskosten korrigieren.** Mit dieser Option können Sie die A Priori auf der Grundlage der (auf der Registerkarte "Kosten" angegebenen) Fehlklassifizierungskosten anpassen. Dadurch können Sie Kosteninformationen für Bäume, die mit dem Twoing-Unreinheitsmaß arbeiten, direkt in den Vorgang zur Baumerweiterung aufnehmen. (Wenn diese Option nicht ausgewählt ist, werden Kosteninformationen nur beim Klassifizieren der Datensätze und bei der Berechnung der Risikoschätzungen für Bäume auf der Grundlage des Twoing-Maßes verwendet.)

### CHAID-Knoten - Kosten

Abbildung 6-30  
Fehlklassifizierungskosten im CHAID-Knoten



In einigen Zusammenhängen sind bestimmte Fehler kostspieliger als andere. Zum Beispiel kann es kostspieliger sein, einen Kreditantragsteller mit hohem Risiko als niedriges Risiko zu klassifizieren (eine Art von Fehler) als einen Kreditantragsteller mit niedrigem Risiko als hohes Risiko (eine andere Art von Fehler). Anhand von Fehlklassifizierungskosten können Sie die relative Bedeutung verschiedener Arten von Vorhersagefehlern angeben.

Fehlklassifizierungskosten sind im Grunde Gewichtungen, die auf bestimmte Ergebnisse angewendet werden. Diese Gewichtungen werden in das Modell mit einbezogen und können die Vorhersage (als Schutz gegen kostspielige Fehler) ändern.

Mit Ausnahme von C5.0-Modellen werden Fehlklassifizierungskosten beim Scoring von Modellen nicht angewendet und bei der Rangbildung bzw. beim Vergleich von Modellen mithilfe eines Knotens vom Typ "Automatischer Klassifizierer", eines Evaluationsdiagramms oder eines Analyseknosens nicht berücksichtigt. Ein Modell, das Kosten beinhaltet, führt nicht unbedingt zu weniger Fehlern als ein Modell, das keine Kosten beinhaltet, und es weist auch nicht unbedingt eine höhere Gesamtgenauigkeit auf, aber es ist möglicherweise in praktischer Hinsicht leistungsfähiger, da es eine integrierte Verzerrung zugunsten von *weniger teuren* Fehlern aufweist.

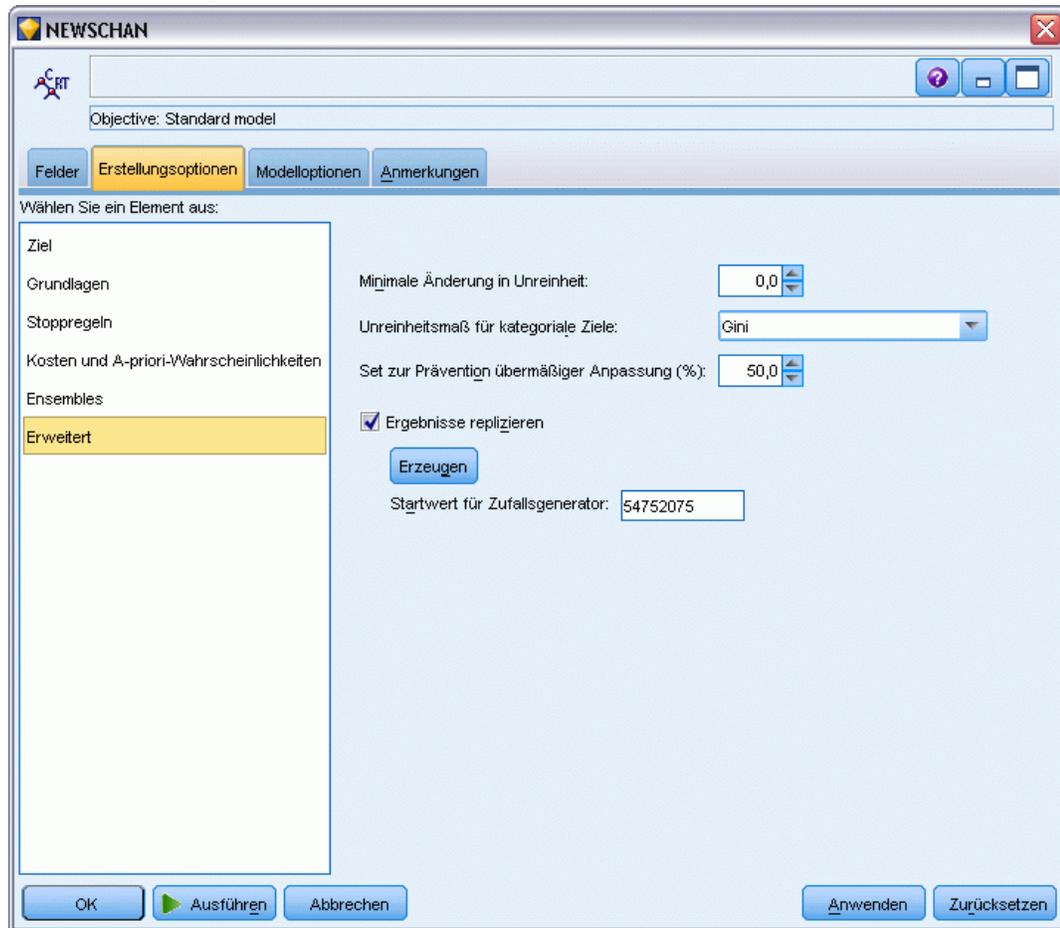
Die Kostenmatrix zeigt die Kosten für jede mögliche Kombination aus prognostizierter Kategorie und tatsächlicher Kategorie. Standardmäßig werden alle Fehlklassifizierungskosten auf 1,0 gesetzt. Um benutzerdefinierte Kostenwerte einzugeben, wählen Sie Fehlklassifizierungskosten verwenden und geben Ihre benutzerdefinierten Werte in die Kostenmatrix ein.

Um die Fehlklassifizierungskosten zu ändern, wählen Sie die Zelle aus, die der gewünschten Kombination aus vorhergesagten und tatsächlichen Werten entspricht, löschen den Inhalt der Zelle und geben die gewünschten Kosten für die Zelle ein. Die Kosten sind nicht automatisch symmetrisch. Wenn Sie z. B. die Kosten einer Fehlklassifizierung von *A* als *B* auf 2,0 setzen, weisen die Kosten für eine Fehlklassifizierung von *B* als *A* weiterhin den Standardwert 1,0 auf, es sei denn, Sie ändern diesen Wert ausdrücklich.

### **Knoten "C&R" – Erweitert**

Mit den erweiterten Optionen können Sie die Feinabstimmung der Baumerstellung vornehmen.

Abbildung 6-31  
Einstellen erweiterter Optionen für den C&R-Baumknoten



**Minimale Änderung in Unreinheit.** Legt die minimale Änderung in der Unreinheit fest, damit im Baum eine neue Aufteilung erstellt wird. **Unreinheit** bezieht sich auf das Ausmaß, in dem durch den Baum definierte Untergruppen in jeder Gruppe eine große Reihe von Ausgabefeldwerten besitzen. Bei kategorialen Zielen wird ein Knoten als “rein” betrachtet, wenn 100 % der im Knoten vorhandenen Fälle in eine bestimmte Kategorie des Zielfelds fallen. Das Ziel der Baumerstellung besteht darin, Untergruppen mit ähnlichen Ausgabewerte zu erstellen, also die Unreinheit in jedem Knoten zu minimieren. Wenn die beste für eine Verzweigung mögliche Aufteilung die Unreinheit um weniger als den vorgegebenen Betrag reduziert, wird die Aufteilung nicht durchgeführt.

**Unreinheitsmaß für kategoriale Ziele.** Geben Sie für kategoriale Zielfelder die für die Messung der im Baum vorhandenen Unreinheit verwendete Methode an. (Für kontinuierliche Ziele wird diese Option ignoriert und als Unreinheitsmaß immer die **kleinste quadratische Abweichung** verwendet.)

- Gini ist ein allgemeines Unreinheitsmaß, das auf Wahrscheinlichkeiten der Zugehörigkeit zu einer Kategorie einer Verzweigung basiert.

- Twoing ist ein Unreinheitsmaß, das die binäre Aufteilung betont und eher zu einer Aufteilung in annähernd gleichgroße Verzweigungen führt.
- Ordinal fügt eine weitere Einschränkung hinzu, indem nur zusammenhängende Zielklassen zu Gruppen zusammengefasst werden können, was nur bei Ordinalzielen möglich ist. Wenn diese Option für ein nominales Ziel ausgewählt ist, wird standardmäßig das Standard-Twoing-Maß verwendet.

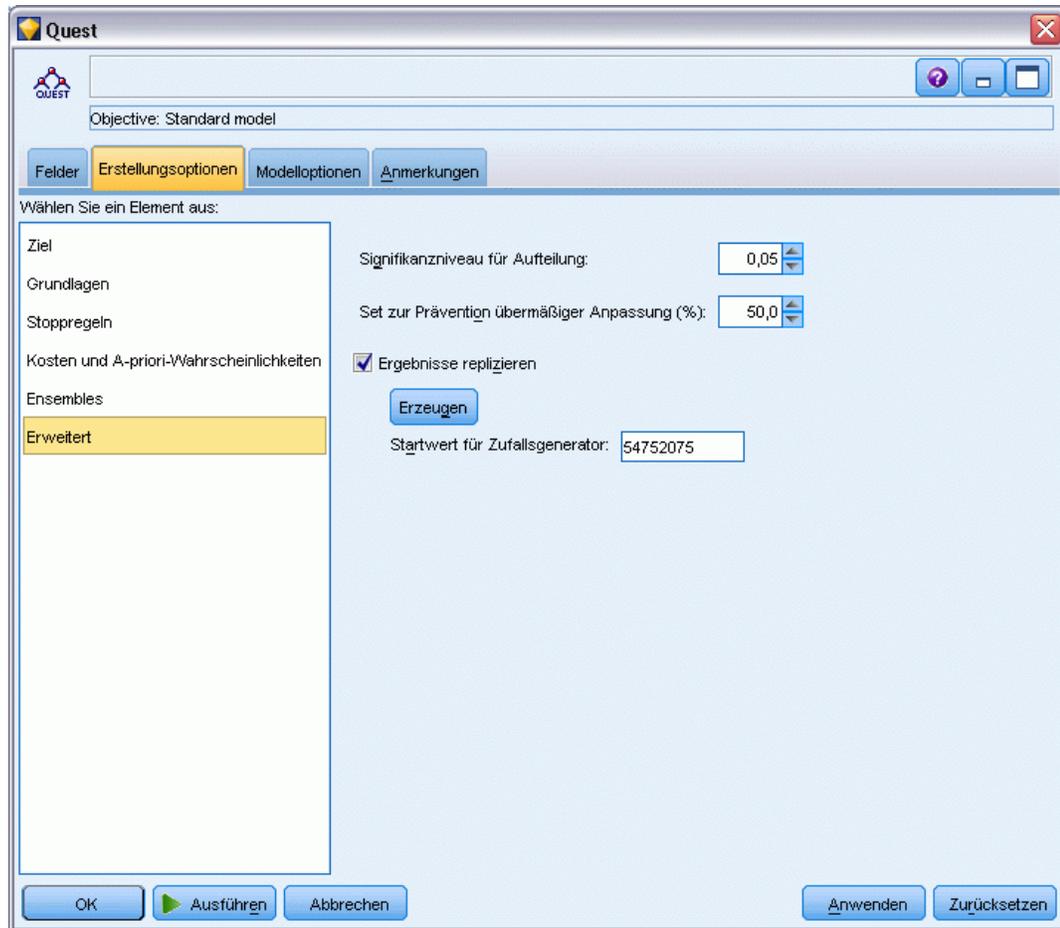
**Verhinderung der übermäßigen Anpassung eingestellt.** Bei dem Algorithmus werden Datensätze intern in ein Modellerstellungs-Set und ein Set zur Verhinderung der übermäßigen Anpassung aufgeteilt. Zweites ist ein unabhängiges Set an Datensätzen, das dazu verwendet wird, Fehler während des Trainings zu erfassen. So kann verhindert werden, dass die Methode zufällige Variationen in den Daten modelliert. Geben Sie einen Prozentsatz an Datensätzen an. Der Standardwert ist 30.

**Ergebnisse reproduzieren.** Durch Einstellen eines Startwerts für den Zufallsgenerator können Analysen reproduziert werden. Geben Sie eine ganze Zahl ein oder klicken Sie auf Generieren. Dadurch wird eine pseudozufällige Ganzzahl zwischen 1 und 2147483647 (einschließlich) erzeugt.

### ***QUEST-Knoten - Erweitert***

Mit den erweiterten Optionen können Sie die Feinabstimmung der Baumerstellung vornehmen.

Abbildung 6-32  
Einstellen erweiterter Optionen für den QUEST-Knoten



**Signifikanzniveau für Aufteilung.** Legt das Signifikanzniveau (Alpha) für das Aufteilen von Knoten fest. Der Wert muss zwischen 0 und 1 liegen. Niedrigere Werte führen in der Regel zu Bäumen mit weniger Knoten.

**Verhinderung der übermäßigen Anpassung eingestellt.** Bei dem Algorithmus werden Datensätze intern in ein Modellerstellungs-Set und ein Set zur Verhinderung der übermäßigen Anpassung aufgeteilt. Zweiteres ist ein unabhängiges Set an Datensätzen, das dazu verwendet wird, Fehler während des Trainings zu erfassen. So kann verhindert werden, dass die Methode zufällige Variationen in den Daten modelliert. Geben Sie einen Prozentsatz an Datensätzen an. Der Standardwert ist 30.

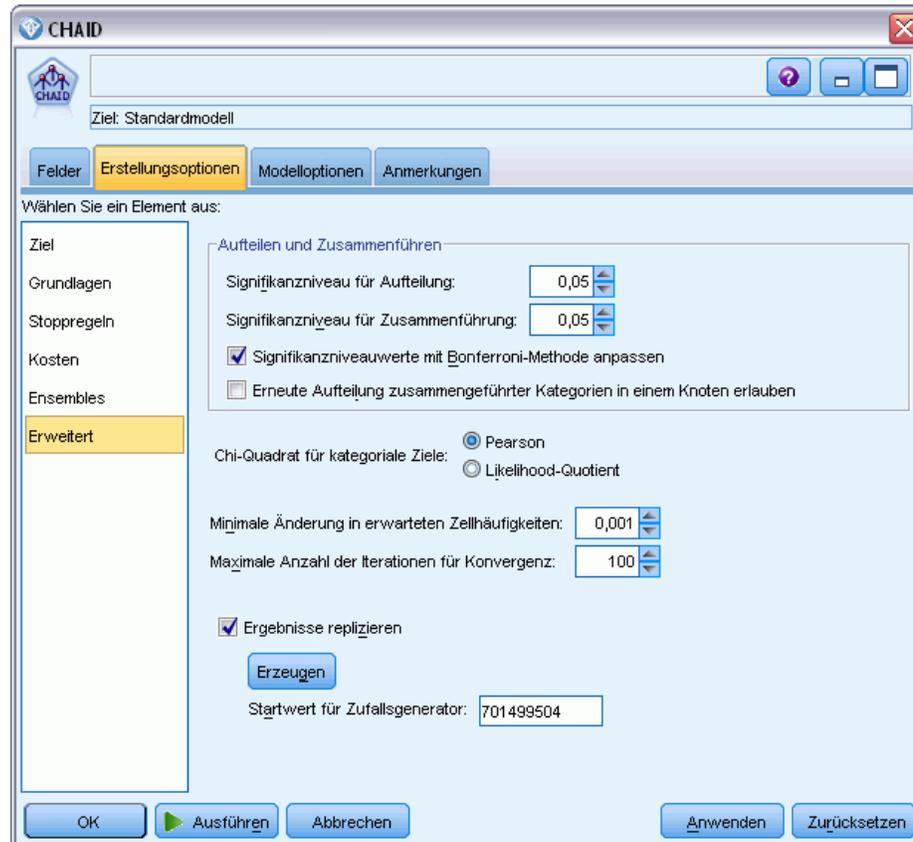
**Ergebnisse reproduzieren.** Durch Einstellen eines Startwerts für den Zufallsgenerator können Analysen reproduziert werden. Geben Sie eine ganze Zahl ein oder klicken Sie auf Generieren. Dadurch wird eine pseudozufällige Ganzzahl zwischen 1 und 2147483647 (einschließlich) erzeugt.

### CHAID-Knoten- Erweitert

Mit den erweiterten Optionen können Sie die Feinabstimmung der Baumerstellung vornehmen.

Abbildung 6-33

Einstellen erweiterter Optionen für den CHAID-Knoten



**Signifikanzniveau für Aufteilung.** Legt das Signifikanzniveau (Alpha) für das Aufteilen von Knoten fest. Der Wert muss zwischen 0 und 1 liegen. Niedrigere Werte führen in der Regel zu Bäumen mit weniger Knoten.

**Signifikanzschwelle für Zusammenführung.** Legt das Signifikanzniveau (Alpha) für das Zusammenführen von Kategorien fest. Der Wert muss größer als 0 und kleiner oder gleich 1 sein. Um zu vermeiden, dass Kategorien zusammengeführt werden, geben Sie den Wert 1 an. Bei stetigen Zielen bedeutet dies, dass die Anzahl der Kategorien für die Variable im Endbaum mit der angegebenen Anzahl von Intervallen übereinstimmt. Diese Option ist für Exhaustive CHAID nicht verfügbar.

**Signifikanzwerte mit der Bonferroni-Methode anpassen.** Passt beim Testen der verschiedenen Kategoriekombinationen eines Prädiktors die Signifikanzwerte an. Die Werte werden auf der Grundlage der Anzahl der Tests angepasst, die sich direkt auf die Anzahl der Kategorien und das Messniveau eines Prädiktors bezieht. Dies ist in der Regel wünschenswert, da eine bessere Kontrolle der falsch positiven Fehlerrate stattfindet. Das Deaktivieren dieser Option erhöht die Leistung Ihrer Analyse beim Auffinden tatsächlicher Differenzen, führt aber zu einer höheren

falsch positiven Rate. Insbesondere für kleine Stichproben kann das Deaktivieren dieser Option ratsam sein.

**Erneutes Aufteilen zusammengeführter Kategorien innerhalb eines Knotens zulassen.** Der CHAID-Algorithmus versucht Kategorien zusammenzuführen, um den einfachsten, das Modell beschreibenden Baum zu erzeugen. Wenn diese Option ausgewählt ist, ist es zulässig, dass zusammengeführte Kategorien noch einmal aufgeteilt werden, wenn dies zu einer besseren Lösung führt.

**Chi-Quadrat für kategoriale Ziele.** Für kategoriale Ziele können Sie die für die Berechnung der Chi-Quadrat-Statistik verwendete Methode angeben.

- **Pearson.** Diese Methode liefert schnellere Berechnungen, sollte bei kleineren Stichproben jedoch nur nach sorgfältiger Erwägung verwendet werden.
- **Likelihood-Quotient.** Diese Methode ist robuster als Pearson, benötigt aber mehr Rechenzeit. Diese Methode eignet sich ideal für kleine Stichproben. Bei stetigen Zielen wird immer diese Methode verwendet.

**Minimale Änderung in der erwarteten Zellhäufigkeit.** Beim Schätzen der Zellhäufigkeiten (des nominalen Modells und des ordinalen Modells der Zeileneffekte) kommt eine iterative Prozedur (Epsilon) zum Einsatz, um ein Konvergieren gegen die optimale Schätzung zu erreichen, die im Chi-Quadrat-Test für eine bestimmte Aufteilung verwendet wird. Epsilon bestimmt, wie groß die Änderung sein muss, damit Iterationen fortgesetzt werden. Wenn die aus der letzten Iteration resultierende Änderung unter dem festgelegten Wert liegt, wird die Iteration beendet. Wenn Sie Probleme damit haben, dass der Algorithmus nicht konvergiert, können Sie diesen Wert erhöhen oder die maximale Anzahl der Iterationen so lange reduzieren, bis die Konvergenz stattfindet.

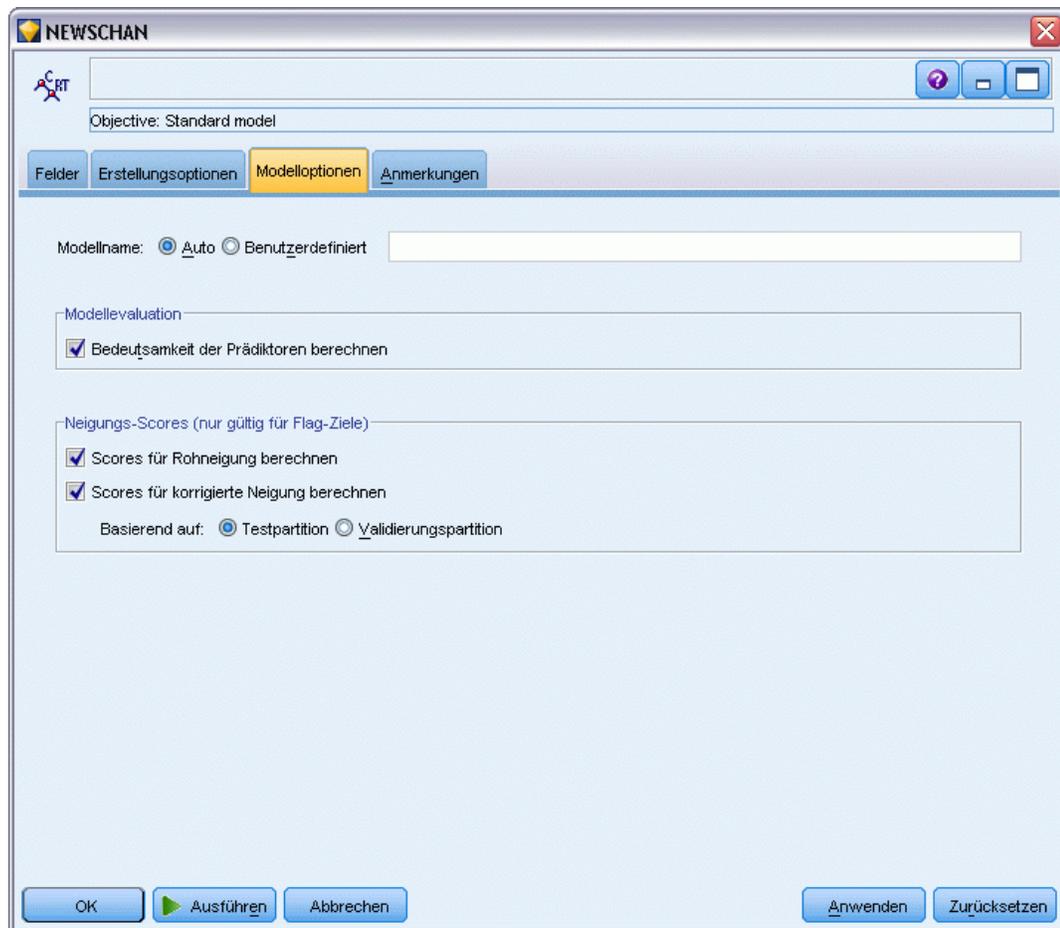
**Maximale Anzahl der Iterationen für Konvergenz.** Legt die maximale Anzahl der Iterationen fest, nach der aufgehört wird — egal ob eine Konvergenz stattgefunden hat oder nicht.

**Ergebnisse reproduzieren.** Durch Einstellen eines Startwerts für den Zufallsgenerator können Analysen reproduziert werden. Geben Sie eine ganze Zahl ein oder klicken Sie auf Generieren. Dadurch wird eine pseudozufällige Ganzzahl zwischen 1 und 2147483647 (einschließlich) erzeugt.

## ***Modelloptionen für Entscheidungsbaumknoten***

Auf der Registerkarte “Modelloptionen” können Sie bestimmen, ob Sie einen Namen für das Modell angeben oder automatisch erstellen lassen möchten. Außerdem können Sie auswählen, ob Sie Informationen zur Bedeutsamkeit der Prädiktoren sowie Scores für Rohneigung und angepasste Neigung für Flag-Ziele erhalten möchten.

Abbildung 6-34  
Einstellen der Modelloptionen für einen Entscheidungsbaumknoten



**Modellname.** Sie können den Modellnamen automatisch basierend auf den Ziel- oder ID-Feldnamen (oder dem Modelltyp in Fällen, in denen kein solches Feld angegeben ist) generieren oder einen benutzerdefinierten Namen eingeben.

### **Modellevaluation**

**Bedeutsamkeit der Prädiktoren berechnen.** Bei Modellen, die zu einem angemessenen Maß an Bedeutsamkeit führen, können Sie ein Diagramm anzeigen, in dem die relative Wichtigkeit der einzelnen Prädiktoren bei der Modellschätzung angegeben wird. Normalerweise ist es sinnvoll, die Modellierungsbemühungen auf die wichtigsten Prädiktoren zu konzentrieren und diejenigen Prädiktoren zu verwerfen bzw. zu ignorieren, die die geringste Bedeutung haben. Beachten Sie, dass die Berechnung der Bedeutsamkeit der Prädiktoren bei einigen Modellen längere Zeit in Anspruch nehmen kann, insbesondere bei der Arbeit mit großen Daten-Sets, und daher bei einigen Modellen standardmäßig deaktiviert ist. Die Bedeutsamkeit der Prädiktoren ist für Entscheidungslistenmodelle nicht verfügbar. [Für weitere Informationen siehe Thema Bedeutsamkeit des Prädiktors in Kapitel 3 auf S. 54.](#)

### **Neigungs-Scores**

Neigungs-Scores können im Modellierungsknoten oder auf der Registerkarte “Einstellungen” im Modell-Nugget aktiviert werden. Diese Funktionalität ist nur verfügbar, wenn das ausgewählte Ziel ein Flag-Feld ist. [Für weitere Informationen siehe Thema Neigungsbewertungen in Kapitel 3 auf S. 44.](#)

**Scores für Rohneigung berechnen.** Rohneigungs-Scores werden ausschließlich auf der Grundlage der Trainingsdaten aus dem Modell abgeleitet. Wenn das Modell den Wert *wahr* (wird antworten) vorhersagt, ist die Neigung mit P identisch. Dabei ist P die Wahrscheinlichkeit der Vorhersage. Wenn das Modell den Wert “falsch” vorhersagt, wird die Neigung als  $(1 - P)$  berechnet.

- Wenn Sie bei der Modellerstellung diese Option auswählen, werden standardmäßig Neigungs-Scores im Modell-Nugget aktiviert. Sie können jedoch immer festlegen, dass Rohneigungs-Scores im Modell-Nugget aktiviert werden sollen, unabhängig davon, ob Sie sie im Modellierungsknoten auswählen.
- Beim Scoring des Modells werden Rohneigungs-Scores in einem Feld hinzugefügt, bei dem die Buchstaben *RP* an das Standard-Präfix angehängt sind. Wenn sich die Vorhersagen beispielsweise im Feld *\$R-Abwanderung* befinden, lautet der Name des Felds für den Neigungs-Score *\$RRP-Abwanderung*.

**Scores für korrigierte Neigung berechnen.** Rohneigungen basieren ausschließlich auf vom Modell angegebenen Schätzern. Beim Modell kann jedoch eine Überanpassung vorliegen, was zu übermäßig optimistischen Schätzern für die Neigung führt. Korrigierte Neigungen versuchen, dies zu kompensieren, indem untersucht wird, wie leistungsfähig das Modell bei den Test- bzw. Validierungspartitionen ist, und die Neigungen entsprechend angepasst werden, um einen besseren Schätzer zu erzeugen.

- Diese Einstellung ist nur möglich, wenn ein gültiges Partitionsfeld im Stream vorhanden ist. [Für weitere Informationen siehe Thema Partitionsknoten in Kapitel 4 in IBM SPSS Modeler 14.2- Quellen-, Prozess- und Ausgabeknoten.](#)
- Anders als rohe Konfidenz-Scores müssen Scores für die korrigierte Neigung bei der Erstellung des Modells berechnet werden; anderenfalls stehen Sie beim Scoring des Modell-Nuggets nicht zur Verfügung.
- Beim Scoring des Modells werden Scores für die korrigierte Neigung in einem Feld hinzugefügt, bei dem die Buchstaben *AP* an das Standard-Präfix angehängt sind. Wenn sich die Vorhersagen beispielsweise im Feld *\$R-Abwanderung* befinden, lautet der Name des Felds für den Neigungs-Score *\$RAP-Abwanderung*. Scores für die korrigierte Neigung stehen bei logistischen Regressionsmodellen nicht zur Verfügung.
- Bei der Berechnung der Scores für die korrigierte Neigung darf die für die Berechnung verwendete Test- bzw. Validierungspartition nicht ausbalanciert worden sein. Um dies zu vermeiden, müssen Sie darauf achten, dass in etwaigen weiter oben im Stream befindlichen Balancierungsknoten die Option Balancierung nur für Trainingsdaten durchführen ausgewählt wurde. [Für weitere Informationen siehe Thema Festlegen der Optionen für den Balancierungsknoten in Kapitel 3 in IBM SPSS Modeler 14.2- Quellen-, Prozess- und](#)

*Ausgabeknoten.* Außerdem gilt: Wenn weiter oben im Stream eine komplexe Stichprobe gezogen wurde, werden dadurch die Scores für die korrigierte Neigung ungültig.

- Scores für die korrigierte Neigung stehen bei verstärkten Baum- und Regelmengenmodellen nicht zur Verfügung. [Für weitere Informationen siehe Thema Verbesserte C5.0-Modelle auf S. 190.](#)

**Basierend auf.** Um Scores für die angepasste Neigung berechnen zu können, muss im Stream ein Partitionsfeld vorhanden sein. Sie können angeben, ob die Test- bzw. Validierungspartition für diese Berechnung verwendet werden soll. Um bestmögliche Ergebnisse zu erzielen, sollte die Test- bzw. Validierungspartition mindestens so viele Datensätze enthalten wie die Partition, die zum Trainieren des ursprünglichen Modells verwendet wurde. [Für weitere Informationen siehe Thema Partitionsknoten in Kapitel 4 in IBM SPSS Modeler 14.2- Quellen-, Prozess- und Ausgabeknoten.](#)

## C5.0-Knoten

Dieser Knoten benutzt den C5.0-Algorithmus, um entweder einen **Entscheidungsbaum** oder eine **Regelmenge** zu erstellen. Ein C5.0-Modell teilt die Stichprobe auf der Basis des Felds auf, das den maximalen **Informationsgewinn** liefert. Jede durch die erste Aufteilung definierte Unterstichprobe wird dann wieder aufgeteilt — üblicherweise auf der Grundlage eines anderen Felds. Das Verfahren wird so lange fortgesetzt, bis die Unterstichproben nicht weiter aufgeteilt werden können. Zum Schluss werden die Aufteilungen der untersten Ebene noch einmal untersucht, wobei solche entfernt oder **reduziert** werden, die nicht wesentlich zum Wert des Modells beitragen.

*Hinweis:* Der Knoten “C5.0” kann nur ein kategoriales Ziel vorhersagen. Bei der Analyse von Daten mit kategorialen Feldern (nominal oder ordinal) fasst der Knoten mit größerer Wahrscheinlichkeit Kategorien zu einer Gruppe zusammen als C5.0-Versionen vor Version 11.0.

C5.0 kann zwei Arten von Modellen erstellen. Ein **Entscheidungsbaum** ist eine einfache Beschreibung der vom Algorithmus gefundenen Aufteilungen. Jeder Endknoten (oder Blattknoten) beschreibt eine bestimmte Untermenge der Trainingsdaten. Und jeder in den Trainingsdaten vorhandene Fall gehört zu genau einem im Baum vorhandenen Endknoten. Somit ist für jeden in einem Entscheidungsbaum vorhandenen Datensatz genau eine Vorhersage möglich.

Eine **Regelmenge** ist dagegen eine Menge von Regeln, mit der versucht wird, Vorhersagen für einzelne Datensätze zu erstellen. Regelmengen werden aus Entscheidungsbaum abgeleitet und stellen eine vereinfachte oder konzentrierte Version der im Entscheidungsbaum gefundenen Informationen dar. Regelmengen enthalten meist die wichtigsten Informationen eines gesamten Entscheidungsbaums, allerdings mit einem weniger komplexen Modell. Regelmengen arbeiten anders als Entscheidungsbaum und besitzen daher nicht dieselben Eigenschaften. Der wichtigste Unterschied besteht darin, dass es bei einer Regelmenge möglich ist, dass für einen bestimmten Datensatz mehr als eine oder aber überhaupt keine Regel gilt. Wenn mehrere Regeln gelten, dann wird jeder Regel ein gewichtetes “Votum” zugeordnet, das auf der dieser Regel zugeordneten Konfidenz basiert, und die endgültige Vorhersage ergibt sich aus der Kombination der gewichteten Voten aller für den fraglichen Datensatz geltenden Regeln. Wenn keine Regel gilt, wird dem Datensatz eine Standardvorhersage zugeordnet.

**Beispiel.** Ein Medizinforscher hat Daten über eine Gruppe von Patienten zusammengetragen, die alle an der gleichen Krankheit leiden. Im Behandlungsverlauf sprach jeder Patient auf eines von fünf Medikamenten an. Sie können ein C5.0-Modell in Verbindung mit anderen Knoten verwenden, um herauszufinden, welches Medikament für einen Patienten geeignet sein kann, der später an derselben Krankheit leidet. [Für weitere Informationen siehe Thema Medikamentöse Behandlung \(Explorative Diagramme/C5.0\) in Kapitel 9 in IBM SPSS Modeler 14.2- Anwendungshandbuch.](#)

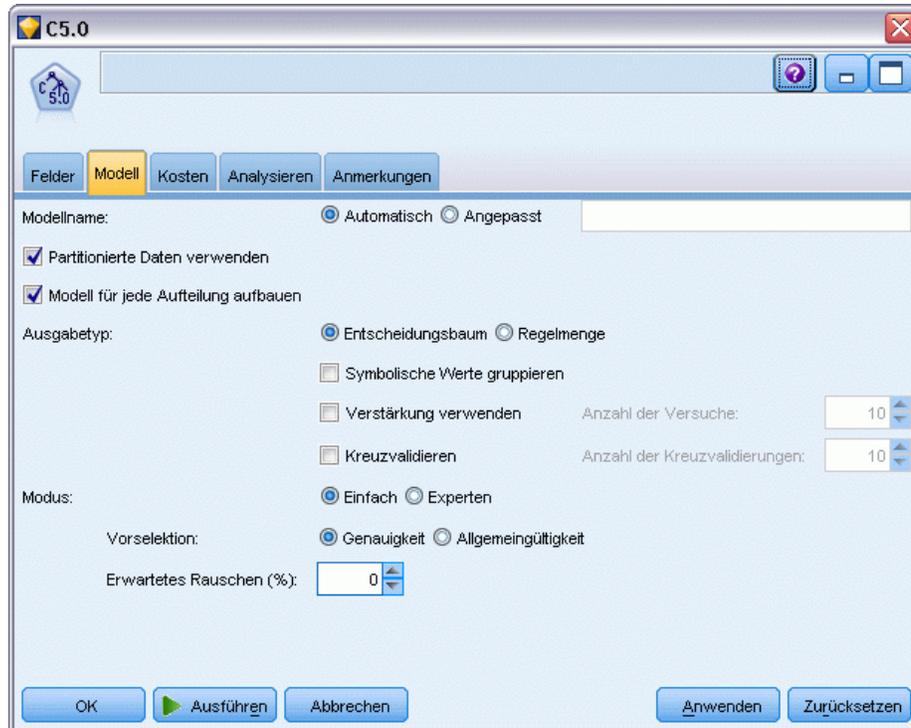
**Anforderungen.** Um ein C5.0-Modell zu trainieren, muss genau ein kategoriales (d. h. vom Typ “Nominal” oder “Ordinal”) *Ziel*-Feld und mindestens ein *Eingabe*-Feld beliebigen Typs vorliegen. Felder, die auf *Beides* oder *Keine* gesetzt sind, werden ignoriert. Bei den im Modell verwendeten Feldern müssen die Typen vollständig instanziiert sein. Außerdem kann ein Gewichtsfeld angegeben werden.

**Stärken.** C5.0-Modelle verhalten sich bei Problemen mit fehlenden Daten und einer großen Anzahl von Eingabefeldern sehr robust. Sie benötigen für die Schätzung in der Regel keine langen Trainingsphasen. Darüber hinaus sind C5.0-Modelle tendenziell leichter verständlich als andere Modelltypen, da sich die aus dem Modell abgeleiteten Regeln sehr direkt interpretieren lassen. C5.0 bietet außerdem die leistungsstarke Methode der **Verstärkung**, mit der die Genauigkeit der Klassifizierung gesteigert wird.

*Anmerkung:* Die Geschwindigkeit für die Erstellung von C5.0-Modellen kann eventuell durch die Aktivierung der parallelen Verarbeitung erhöht werden. [Für weitere Informationen siehe Thema Festlegen von Optimierungsoptionen in Kapitel 12 in IBM SPSS Modeler 14.2- Benutzerhandbuch.](#)

## Modelloptionen für C5.0-Knoten

Abbildung 6-35  
Modelloptionen für C5.0-Knoten



**Modellname.** Geben Sie den Namen des zu erstellenden Modells an.

- **Auto.** Wenn diese Option ausgewählt ist, wird der Modellname automatisch auf der Grundlage der Namen der Zielfelder generiert. Dies ist die Standardeinstellung.
- **Benutzerdefiniert.** Wählen Sie diese Option, um für das durch diesen Knoten erstellte Modell-Nugget einen eigenen Namen anzugeben.

**Partitionierte Daten verwenden.** Wenn ein Partitionsfeld definiert ist, gewährleistet diese Option, dass nur Daten aus der Trainingspartition für die Modellerstellung verwendet werden. [Für weitere Informationen siehe Thema Partitionsknoten in Kapitel 4 in IBM SPSS Modeler 14.2- Quellen-, Prozess- und Ausgabeknoten.](#)

**Aufteilungsmodelle erstellen.** Erstellt ein separates Modell für jeden möglichen Wert von Eingabefeldern, die als Aufteilungsfelder angegeben werden. [Für weitere Informationen siehe Thema Erstellung von aufgeteilten Modellen in Kapitel 3 auf S. 31.](#)

**Ausgabetypp.** Legen Sie hier fest, ob das resultierende Modell-Nugget ein Entscheidungsbaum oder eine Regelmenge sein soll.

**Symbolische Werte gruppieren.** Wenn diese Option ausgewählt ist, versucht C5.0 symbolische Werte zu gruppieren, die in Bezug auf das Ausgabefeld ähnliche Muster aufweisen. Wenn diese Option nicht ausgewählt ist, erzeugt C5.0 für jeden Wert des symbolischen Felds einen untergeordneten Knoten, der zum Aufteilen des übergeordneten Knotens verwendet wird. C5.0

teilt beispielsweise das Feld *FARBE* (mit den Werten *ROT*, *GRÜN* und *BLAU*) standardmäßig in drei Teile. Wenn diese Option jedoch ausgewählt ist und die Datensätze mit *FARBE = ROT* sehr ähnlich aussehen wie Datensätze mit *FARBE = BLAU*, dann wird eine Zweiteilung durchgeführt, bei der alle Datensätze mit *GRÜN* eine Gruppe bilden und alle mit *BLAU* und *ROT* gemeinsam eine zweite.

**Verstärkung verwenden.** Der C5.0-Algorithmus verwendet eine spezielle Methode zur Verbesserung der Genauigkeitsrate, die als **Verstärkung** bezeichnet wird. Sie arbeitet so, dass sie mehrere Modelle in einer Folge erstellt. Das erste Modell wird auf die übliche Weise erstellt. Anschließend wird ein zweites Modell erstellt, bei dem besonders die Datensätze berücksichtigt werden, bei denen es im ersten Modell zu Fehlklassifizierungen kam. Das dritte Modell wird in Bezug auf die im zweiten Modell enthaltenen Fehler erstellt usw. Zum Schluss werden die Fälle klassifiziert, indem der gesamte Modellsatz auf ihnen angewendet wird, wobei ein gewichtetes Voting-Verfahren genutzt wird, um die einzelnen Vorhersagen zu einer Gesamtvorhersage zu kombinieren. Eine Verstärkung kann die Genauigkeit eines C5.0-Modells signifikant verbessern, macht aber auch ein längeres Training notwendig. Mit der Option Anzahl der Versuche können Sie steuern, wie viele Modelle für das verstärkte Modell verwendet werden. Diese Funktion basiert auf der Forschung von Freund & Schapire, mit einigen eigentumsrechtlich geschützten Verbesserungen für die Behandlung verrauschter Daten.

**Kreuzvalidieren.** Mit dieser Option verwendet C5.0 einen Satz von Modellen, die aus einer Untermenge der Trainingsdaten erstellt werden, um die Genauigkeit eines aus dem gesamten Daten-Set erstellten Modells zu schätzen. Dies ist nützlich, wenn das Daten-Set zu klein ist, um in herkömmliche Trainings- und Test-Sets aufgeteilt zu werden. Die Kreuzvalidierungsmodelle werden nach der Berechnung der Genauigkeitsschätzung entfernt. Sie können auch die **Anzahl der Kreuzvalidierungen** oder die Anzahl der für die Kreuzvalidierung verwendeten Modelle festlegen. Beachten Sie, dass die Modellbildung und die Kreuzvalidierung in früheren IBM® SPSS® Modeler-Versionen zwei separate Vorgänge darstellten. In der aktuellen Version ist kein gesonderter Modellbildungsschritt erforderlich. Modellbildung und Kreuzvalidierung werden gleichzeitig durchgeführt.

**Modalwert.** Für ein Training vom Typ Einfach werden die meisten C5.0-Parameter automatisch eingestellt. Ein Training vom Typ Experten bietet Ihnen eine direktere Kontrolle über die Trainingsparameter.

### **Optionen des Modus "Einfach"**

**Vorselektion.** Standardmäßig versucht C5.0 den genauest möglichen Baum zu erstellen. In einigen Fällen kann dies zu einer Überanpassung führen, die eine schwache Leistung bedingt, sobald das Modell auf neue Daten angewendet wird. Wählen Sie die Option Allgemeingültigkeit, um Algorithmeinstellungen zu verwenden, mit denen dieses Problem weniger häufig auftritt.

*Hinweis:* Es gibt keine Garantie dafür, dass mit der Option Allgemeingültigkeit erstellte Modelle besser generalisiert werden können als andere. Wenn die Allgemeingültigkeit ein wichtiger Punkt ist, müssen Sie Ihr Modell immer gegen eine zurückgehaltene Teststichprobe validieren.

**Erwartetes Rauschen (%).** Geben Sie den erwarteten Anteil der im Trainings-Set enthaltenen verrauschten oder fehlerhaften Daten an.

### **Optionen des Expertenmodus**

**Reduktionsgrad.** Legt fest, in welchem Umfang der Entscheidungsbaum bzw. die Regelmenge reduziert werden. Wenn Sie diesen Wert erhöhen, erhalten Sie einen kleineren und prägnanteren Baum. Geben Sie einen niedrigeren Wert an, wenn Sie einen genaueren Baum erhalten wollen. Diese Einstellung wirkt sich ausschließlich auf die lokale Reduktion aus (siehe unten “Globale Reduktion verwenden”).

**Minimale Anzahl der Datensätze pro untergeordneter Verzweigung.** Anhand der Größe der Untergruppen kann die Anzahl der in jeder Verzweigung des Baums durchgeführten Aufteilungen eingeschränkt werden. Eine Verzweigung wird nur dann aufgeteilt, wenn zwei oder mehr der daraus entstehenden Unterverzweigungen mindestens so viele Datensätze des Trainings-Sets enthalten. Der Standardwert ist 2. Erhöhen Sie den Wert, um ein **Übertrainieren** mit verrauschten Daten zu verhindern.

**Globale Reduktion verwenden.** Bäume werden in zwei Phasen reduziert: Zuerst in einer lokalen Reduktionsphase, in der die untergeordneten Bäume untersucht und Verzweigungen reduziert werden, um die Genauigkeit des Modells zu steigern. Die zweite, globale Reduktionsphase berücksichtigt den Baum als Ganzes und reduziert schwache untergeordnete Bäume. Die globale Reduktion wird standardmäßig durchgeführt. Wenn Sie die globale Reduktionsphase auslassen wollen, müssen Sie diese Option deaktivieren.

**Vorselektion.** Mit dieser Option untersucht C5.0 die Nützlichkeit der Prädiktoren, bevor die Modellbildung gestartet wird. Als irrelevant eingestufte Prädiktoren werden dann aus dem Modellbildungsvorgang ausgeschlossen. Diese Option ist oft bei Modellen mit vielen Prädiktorfeldern hilfreich und kann eine Überanpassung verhindern.

*Anmerkung:* Die Geschwindigkeit für die Erstellung von C5.0-Modellen kann eventuell durch die Aktivierung der parallelen Verarbeitung erhöht werden. [Für weitere Informationen siehe Thema Festlegen von Optimierungsoptionen in Kapitel 12 in IBM SPSS Modeler 14.2- Benutzerhandbuch.](#)

## **Entscheidungsbaummodell-Nuggets**

Modell-Nuggets vom Typ “Entscheidungsbaum” stellen die Baumstrukturen für die Vorhersage eines bestimmten Ausgabefelds dar, das von einem der Knoten für die Entscheidungsbaum-Modellierung (“C&R-Baum”, “CHAID”, “QUEST”, “C5.0”) entdeckt wurde. Die Baummodelle können direkt aus dem Baumerstellungsknoten oder indirekt aus dem interaktiven Tree Builder generiert werden. [Für weitere Informationen siehe Thema Interactive Tree Builder auf S. 129.](#)

### **Scoring von Baummodellen**

Wenn Sie einen Stream ausführen, der Baum-Modell-Nugget enthält, hängt das jeweils erzielte Ergebnis vom Baumtyp ab.

- Bei Klassifizierungsbäumen (kategoriales Ziel) werden zwei neue Felder, die den vorhergesagten Wert und die Konfidenz für die einzelnen Datensätze enthalten, zu den Daten hinzugefügt. Die Vorhersage beruht auf der häufigsten Kategorie für den Endknoten, dem der Datensatz zugewiesen ist; wenn die Mehrzahl der Antworten in einem bestimmten Knoten *ja* lautet, ist die Vorhersage für alle diesem Knoten zugewiesenen Datensätzen “Ja”.

- Bei Regressionsbäumen werden lediglich vorhergesagte Werte generiert, es werden keine Konfidenzen zugewiesen.
- Optional kann für Modelle vom Typ “CHAID”, “QUEST”, und “C&R-Baum” ein weiteres Feld hinzugefügt werden, das die ID für den Knoten angibt, dem die einzelnen Datensätze zugewiesen werden.

Die neuen Feldnamen werden durch Hinzufügen von Präfixen aus dem Modellnamen abgeleitet. Bei “C&R-Baum”, “CHAID” und “QUEST” lautet das Präfix *\$R-* für das Vorhersagefeld *\$RC-* für das Konfidenzfeld und *\$RI-* für das Knoten-ID-Feld. Bei C5.0-Bäumen lautet das Präfix *\$C-* für das Vorhersagefeld und *\$CC-* für das Konfidenzfeld. Wenn mehrere Entscheidungsbaumknoten vorliegen, enthalten die neuen Feldnamen Zahlen im *Präfix*, um sie gegebenenfalls unterscheiden zu können, beispielsweise durch *\$RI-* und *\$RC1-*, *\$R2-*.

### **Arbeiten mit Modell-Nuggets vom Typ “Entscheidungsbaum”**

Sie können zum Modell gehörige Informationen auf verschiedene Weise speichern bzw. exportieren.

*Hinweis:* Viele dieser Optionen stehen auch im Tree Builder-Fenster zur Verfügung.

Im Tree Builder oder in einem Baummodell-Nugget können Sie:

- Einen Filter generieren oder einen Knoten basierend auf dem aktuellen Baum auswählen. [Für weitere Informationen siehe Thema Generieren von Filter- und Auswahlknoten auf S. 151.](#)
- Ein Regelmengen-Nugget generieren, das die Baumstruktur als Set von Regeln darstellt, das die Endverzweigungen des Baums definiert. [Für weitere Informationen siehe Thema Generieren einer Regelmenge aus einem Entscheidungsbaum auf S. 151.](#)
- Bei Baummodell-Nuggets können Sie außerdem das Modell im PMML-Format exportieren. [Für weitere Informationen siehe Thema Die Modellpalette in Kapitel 3 auf S. 50.](#) Wenn das Modell benutzerdefinierte Aufteilungen enthält, werden diese Informationen in der exportierten PMML nicht beibehalten. (Die Aufteilung wird beibehalten, die Tatsache, dass sie benutzerdefiniert und nicht vom Algorithmus gewählt ist, jedoch nicht.)
- Generieren Sie ein Diagramm auf der Basis des ausgewählten Teils des aktuellen Baums. *Anmerkung:* Dies ist nur für ein Nugget möglich, wenn es mit anderen Knoten in einem Stream verbunden ist. [Für weitere Informationen siehe Thema Erzeugen von Diagrammen auf S. 191.](#)
- Nur bei verstärkten C5.0-Modellen: Auswählen von Einzelner Entscheidungsbaum (Zeichenbereich) oder Einzelner Entscheidungsbaum (Palette der generierten Modelle), um eine neue einzelne Regelmenge zu erstellen, die aus der aktuell ausgewählten Regelmenge abgeleitet ist. [Für weitere Informationen siehe Thema Verbesserte C5.0-Modelle auf S. 190.](#)

*Hinweis:* Der Regelerstellungsknoten wurde zwar durch den Knoten “C&R-Baum” ersetzt, die ursprünglich mit einem Regelerstellungsknoten erstellten Entscheidungsbaumknoten in bestehenden Streams funktionieren jedoch weiterhin ordnungsgemäß.

## Modell-Nuggets bei einzelnen Bäumen

Wenn Sie am Modellierungsknoten Einzelnen Baum aufbauen als Hauptziel auswählen, enthält das daraus resultierende Modell-Nugget die folgenden Registerkarten:

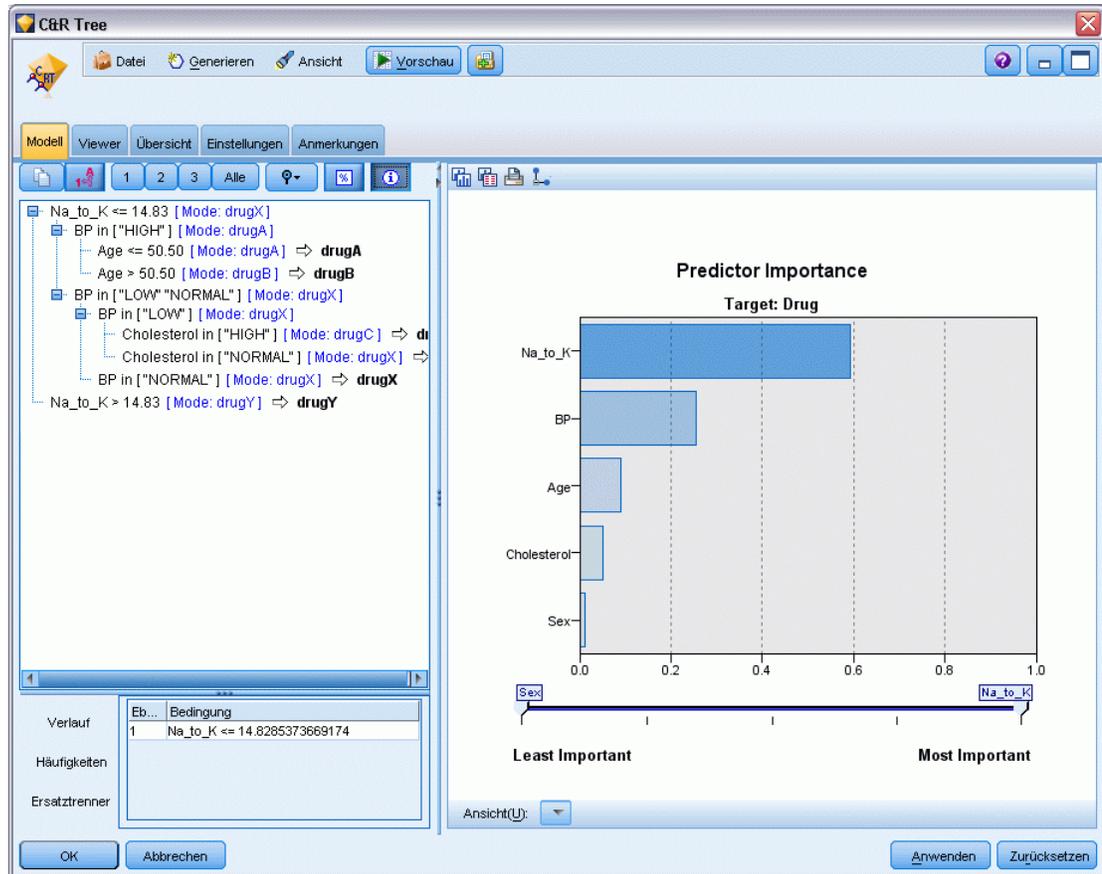
Tabulator	Beschreibung	Weitere Informationen
Modell	Zeigt die Regeln an, die das Modell definieren.	Für weitere Informationen siehe Thema Regeln für Entscheidungsbaummodelle auf S. 183.
Viewer	Zeigt die Baumansicht des Modells an.	Für weitere Informationen siehe Thema Viewer für Entscheidungsbaummodelle auf S. 187.
Zusammenfassung	Zeigt Informationen über die Felder, die Aufbaueinstellungen und die Modellschätzung an.	Für weitere Informationen siehe Thema Modell-Nugget-Übersicht/-Information in Kapitel 3 auf S. 54.
Einstellungen	Hier können Sie Optionen für Konfidenzen und für die SQL-Erzeugung während des Modell-Scorings angeben.	Für weitere Informationen siehe Thema Einstellungen für Modell-Nuggets vom Typ "Entscheidungsbaum"/"Regelmenge" auf S. 188.
Anmerkung	Hier können Sie beschreibende Anmerkungen hinzufügen, einen benutzerdefinierten Namen angeben, QuickInfo-Text hinzufügen und Suchwörter für das Modell angeben.	Für weitere Informationen siehe Thema Anmerkungen in Kapitel 5 in <i>IBM SPSS Modeler 14.2- Benutzerhandbuch</i> .

## Regeln für Entscheidungsbaummodelle

Die Registerkarte "Modell" für ein Entscheidungsbaum-Nugget zeigt die Regeln an, die das Modell definieren. Optional können auch ein Diagramm für die Bedeutsamkeit des Prädiktors und ein drittes Fenster mit Informationen zu Verlauf, Häufigkeiten und Ersatztrennern angezeigt werden.

*Hinweis:* Wenn Sie die Option Modell für sehr große Daten-Sets erstellen in der Registerkarte "Erstellungsoptionen" des CHAID-Knotens (Panel "Ziel") wählen, zeigt die Registerkarte "Modell" nur Details zu den Baumregeln an.

Abbildung 6-36  
Entscheidungsbaummodell-Nugget



### Baumregeln

Im linken Fenster wird eine Liste mit Bedingungen angezeigt, die die Aufteilung der vom Algorithmus ermittelten Daten definieren. Im Wesentlichen handelt es sich hierbei um eine Reihe von Regeln, die dazu verwendet werden können, einzelne Datensätze untergeordneten Knoten basierend auf den Werten verschiedener Prädiktoren zuzuweisen.

Entscheidungsbäume funktionieren durch die rekursive Aufteilung der Daten basierend auf den Werten der Eingabefelder. Die Datenaufteilungen werden als **Verzweigungen** bezeichnet. Die erste Verzweigung (manchmal als **Stamm** bezeichnet) umfasst alle Datensätze. Der Stamm ist in Untergruppen bzw. **untergeordnete Verzweigungen** aufgeteilt, basierend auf dem Wert eines bestimmten Eingabefelds. Jede untergeordnete Verzweigung kann in weitere Verzweigungen aufgeteilt werden, die wiederum aufgeteilt werden können usw. Auf der niedrigsten Ebene des Baums befinden sich Verzweigungen ohne weitere Aufteilungen. Diese Verzweigungen heißen **Endverzweigungen** (oder **Blätter**).

Der Regelbrowser zeigt die Eingabewerte, die jede Partition oder Verzweigung definieren, sowie eine Übersicht über die Werte der Ausgabefelder für die Datensätze dieser Aufteilung. Allgemeine Informationen zur Verwendung des Modellbrowsers finden Sie unter [Durchsuchen von Modell-Nuggets](#).

Bei Aufteilungen, die auf numerischen Feldern basieren, wird die Verzweigung durch eine Zeile ähnlich der folgenden dargestellt:

Feldname Beziehungswert [Übersicht]

wobei *Beziehung* eine numerische Beziehung darstellt. Eine Verzweigung, die z. B. durch Werte größer als 100 für das Feld *Einkünfte* definiert ist, würde wie folgt dargestellt:

Einkünfte > 100 [Übersicht].

Bei Aufteilungen, die auf symbolischen Feldern basieren, wird die Verzweigung durch eine Zeile ähnlich der folgenden dargestellt:

Feldname = Wert [Übersicht] oder Feldname in [Werte] [Übersicht]

wobei *Werte* die Feldwerte darstellt, die von der Verzweigung definiert werden. Eine Verzweigung, die z. B. Datensätze einschließt, bei denen der Wert *RegionNorden*, *Westen* oder *Süden* sein kann, würde wie folgt dargestellt:

Region in ["Norden" "Westen" "Süden"] [Übersicht]

Für Endverzweigungen wird ebenfalls eine Vorhersage ausgegeben, indem am Ende der Regelbedingung ein Pfeil und der vorhergesagte Wert hinzugefügt werden. Ein Blatt, das z. B. als *Einkünfte* > 100 definiert ist und den Wert *hoch* für das Ausgabefeld vorhersagt, würde wie folgt dargestellt:

Einkünfte > 100 [Modus: hoch] • hoch

Die **Übersicht** für die Verzweigung ist für symbolische und numerische Ausgabefelder unterschiedlich definiert. Für Bäume mit numerischen Ausgabefeldern stellt die Übersicht den **durchschnittlichen** Wert für die Verzweigung dar und die **Standardabweichung** der Verzweigung ist die Abweichung zwischen dem Durchschnitt für die Verzweigung und dem Durchschnitt der übergeordneten Verzweigung. Im Falle von Bäumen mit symbolischen Ausgabefeldern ist die Übersicht der **Modus** bzw. der häufigste Wert für Datensätze in der Verzweigung.

Um eine Verzweigung vollständig zu beschreiben, müssen Sie die Bedingung einschließen, die die Verzweigung definiert sowie die Bedingungen, die die Aufteilungen weiter oben im Baum definieren. Beispiel: Im Baum

```
Einkünfte > 100
Region = "Norden"
Region in ["Süden" "Osten" "Westen"]
Einkünfte <= 200
```

wird die durch die zweite Zeile dargestellte Verzweigung definiert durch die Bedingungen *Einkünfte* > 100 und *Region* = "Norden".

Wenn Sie auf die Schaltfläche Instanzen/Konfidenz anzeigen in der Symbolleiste klicken, werden bei jeder Regel auch die Informationen darüber angezeigt, für wie viele Datensätze die Regel gilt (**Instanzen**), sowie der Anteil der Datensätze, für die die gesamte Regel wahr ist (**Konfidenz**).

### **Bedeutsamkeit des Prädiktors**

Optional kann auf der Registerkarte “Modell” auch ein Diagramm, das die relative Bedeutsamkeit der einzelnen Prädiktoren für die Schätzung des Modells angibt, angezeigt werden. Normalerweise ist es sinnvoll, die Modellierungsbemühungen auf die wichtigsten Prädiktoren zu konzentrieren und diejenigen Prädiktoren zu verwerfen bzw. zu ignorieren, die die geringste Bedeutung haben. Beachten Sie, dass dieses Diagramm nur verfügbar ist, wenn vor dem Generieren des Modells Bedeutsamkeit der Prädiktoren berechnet auf der Registerkarte “Analysieren” ausgewählt wurde. [Für weitere Informationen siehe Thema Bedeutsamkeit des Prädiktors in Kapitel 3 auf S. 54.](#)

### **Zusätzliche Modellinformationen**

Wenn Sie in der Symbolleiste auf Fenster mit weiteren Informationen klicken, wird unten im Fenster ein Teilfenster mit detaillierten Informationen über die ausgewählte Regel angezeigt. Das Informationsfenster enthält drei Registerkarten.

Abbildung 6-37  
Im Informationsfenster angezeigte Ersatztrenner

Verlauf	Regel	
	Primär	Cholesterol in [ "NORMAL" ]
Häufigkeiten	1	Age > 30
	2	Na_to_K > 10.6886346572169
Ersatztrenner		

**Verlauf.** Auf dieser Registerkarte werden die Aufteilungsbedingungen vom Stammknoten abwärts zum ausgewählten Knoten verfolgt. Auf diese Weise erhalten Sie eine Liste mit Bedingungen, die festlegen, wann ein Datensatz dem ausgewählten Knoten zugewiesen wird. Datensätze, für die alle Bedingungen wahr sind, werden diesem Knoten zugewiesen.

**Häufigkeiten.** Für Modelle mit symbolischen Zielfeldern zeigt diese Registerkarte für jeden möglichen Zielwert die Anzahl der Datensätze an, die diesem Knoten (in den Trainingsdaten) zugewiesen sind und diesen Zielwert aufweisen. Die Zahl für die Häufigkeit, ausgedrückt als Prozentwert (bis maximal drei Dezimalstellen) wird ebenfalls angezeigt. Für Modelle mit numerischen Zielwerten bleibt diese Registerkarte leer.

**Ersatzfelder.** Wo zutreffend, werden für den ausgewählten Knoten alle Ersatzfelder für das primäre Aufteilungsfeld angezeigt. Ersatzfelder sind alternative Felder, die verwendet werden, wenn der primäre Prädiktorwert für einen bestimmten Datensatz fehlt. Die maximale Anzahl an Ersatzfeldern, die für eine bestimmte Aufteilung erlaubt ist, ist im Baumerstellungsknoten angegeben, die tatsächliche Anzahl richtet sich jedoch nach den Trainingsdaten. Im Allgemeinen gilt: Je mehr fehlende Daten, desto mehr Ersatzfelder werden wahrscheinlich verwendet. Für andere Entscheidungsbaummodelle bleibt diese Registerkarte leer.

*Anmerkung:* Um im Modell berücksichtigt zu werden, müssen die Ersatzfelder während der Trainingsphase ermittelt werden. Wenn die Trainings-Stichprobe keine fehlenden Werte enthält, werden keine Ersatzfelder ermittelt. Alle Datensätze mit fehlenden Werten, die während des Testens oder der Bewertung gefunden werden, fallen automatisch in den untergeordneten Knoten mit der größten Anzahl an Datensätzen. Wenn während des Testens oder Bewertens fehlende

Werte erwartet werden, können Sie sicher sein, dass auch in den Trainings-Stichproben Werte fehlen. Für CHAID-Bäume stehen keine Ersatzfelder zur Verfügung.

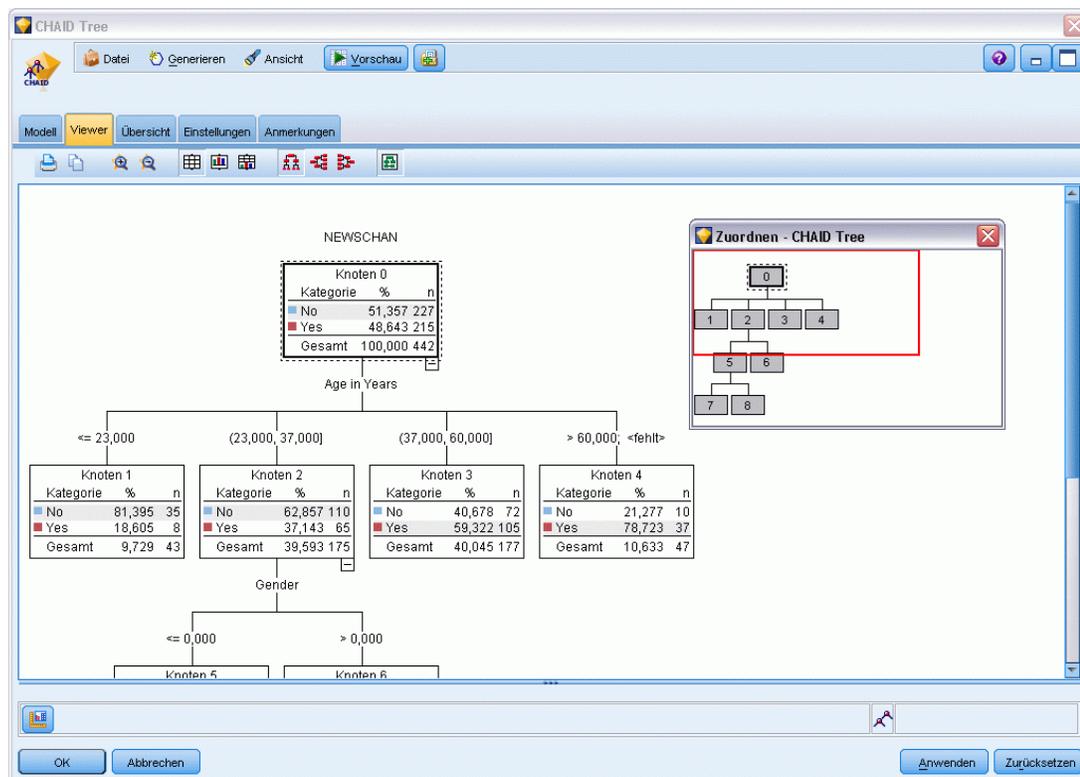
### Viewer für Entscheidungsbaummodelle

Die Registerkarte “Viewer” eines Modell-Nuggets vom Typ “Entscheidungsbaum” ähnelt der Anzeige im Tree Builder. Der Hauptunterschied besteht darin, dass der Baum beim Modell-Nugget nicht erweitert oder bearbeitet werden kann. Andere Optionen für die Betrachtung und Anpassung der Anzeige sind zwischen den beiden Komponenten ähnlich. [Für weitere Informationen siehe Thema Anpassen der Baumansicht auf S. 136.](#)

*Hinweis:* Die Registerkarte “Viewer” wird für CHAID-Modell-Nuggets nicht angezeigt, wenn Sie die Option Modell für sehr große Daten-Sets erstellen auf der Registerkarte “Erstellungsoptionen” des Bereichs “Ziel” auswählen.

Abbildung 6-38

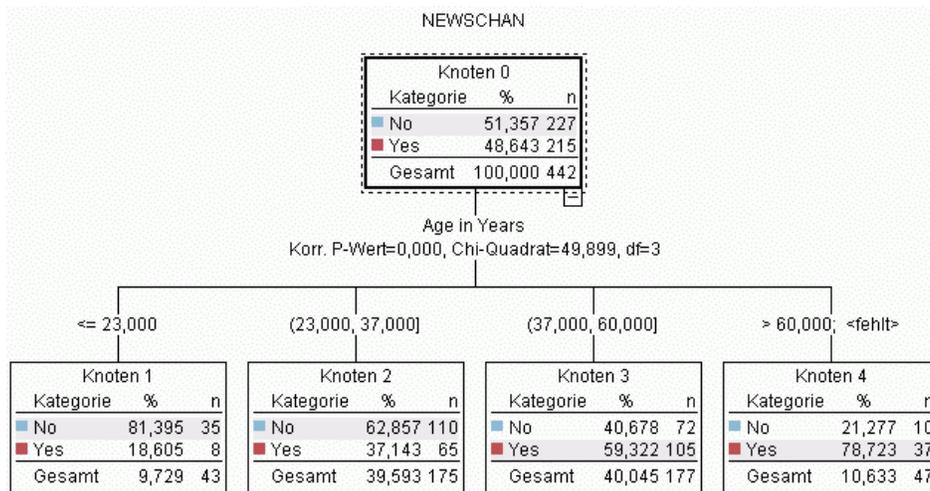
Registerkarte “Viewer” eines Entscheidungsbaums mit Baumübersichtsfenster



Bei der Anzeige von Aufteilungsregeln auf der Registerkarte “Viewer” bedeuten eckige Klammern, dass der angegebene Wert im Bereich enthalten ist, während runde Klammern anzeigen, dass der Wert aus dem Bereich ausgeschlossen ist. Der Ausdruck  $(23,37]$  bedeutet somit von 23 (ausgeschlossen) bis einschließlich 37, also von etwas über 23 bis 37. Auf der Registerkarte “Modell” wird dieselbe Bedingung wie folgt angezeigt:

Age > 23 and Age <= 37

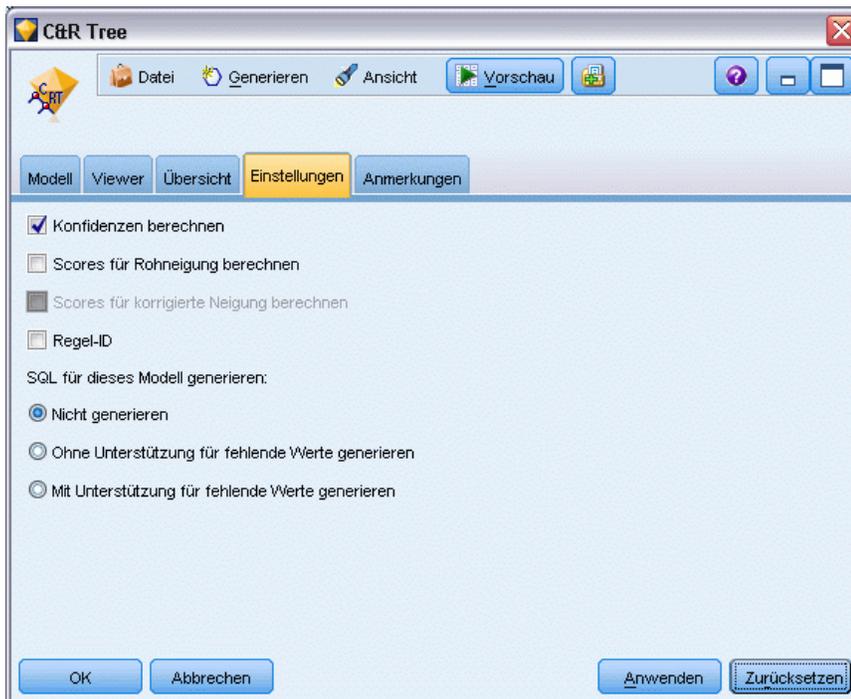
Abbildung 6-39  
Aufteilungsregeln auf der Registerkarte "Viewer"



### Einstellungen für Modell-Nuggets vom Typ "Entscheidungsbaum"/"Regelmenge"

Auf der Registerkarte "Einstellungen" für ein Modell-Nugget vom Typ "Entscheidungsbaum" oder "Regelmenge" können Sie Optionen für Konfidenzen und für die SQL-Erzeugung während des Modell-Scorings angeben. Diese Registerkarte ist erst verfügbar, nachdem das Modell-Nugget einem Stream hinzugefügt wurde.

Abbildung 6-40  
Einstellungen für Entscheidungsbaummodell-Nugget



**Konfidenzen berechnen.** Wählen Sie diese Option aus, um Konfidenzen in die Scoring-Operationen aufzunehmen. Beim Scoring von Modellen in der Datenbank können Sie durch Ausschluss der Konfidenzen effizienter eine SQL erzeugen. Bei Regressionsbäumen werden keine Konfidenzen zugewiesen.

*Hinweis:* Wenn Sie die Option Modell für sehr große Daten-Sets erstellen auf der Registerkarte “Erstellungsoptionen” des CHAID-Modells (Bereich “Methode”) wählen, ist dieses Kontrollkästchen nur in den Modell-Nuggets für kategoriale Ziele von “Nominal” oder “Flag” verfügbar.

**Scores für Rohneigung berechnen.** Bei Modellen mit einem Flag-Ziel (das als Vorhersage “Ja” bzw. “Nein” ausgibt) können Sie Neigungs-Scores anfordern, die die Wahrscheinlichkeit des für das Zielfeld angegebenen wahren Ergebnisses angeben. Diese Werte werden zusätzlich zu den anderen Vorhersage- und Konfidenzwerten angegeben, die ggf. während des Scorings erstellt werden.

*Hinweis:* Wenn Sie die Option Modell für sehr große Daten-Sets erstellen auf der Registerkarte “Erstellungsoptionen” des CHAID-Modells (Bereich “Methode”) wählen, ist dieses Kontrollkästchen nur in den Modell-Nuggets für kategoriale Ziele von “Flag” verfügbar.

**Scores für korrigierte Neigung berechnen.** Rohneigungs-Scores beruhen ausschließlich auf den Trainingsdaten und sind aufgrund der Neigung vieler Modelle zur übermäßigen Anpassung an die Trainingsdaten möglicherweise zu optimistisch. Bei korrigierten Neigungen wird versucht, dies durch Evaluation der Modellleistung anhand einer Test- bzw. Validierungspartition zu kompensieren. Bei dieser Option muss im Stream ein Partitionsfeld definiert sein und die Scores für die korrigierte Neigung müssen im Modellierungsknoten aktiviert werden, bevor das Modell generiert wird.

*Hinweis:* Scores für die korrigierte Neigung stehen bei verstärkten Baum- und Regelmengenmodellen nicht zur Verfügung. [Für weitere Informationen siehe Thema Verbesserte C5.0-Modelle auf S. 190.](#)

**Regel-ID.** Bei Modellen vom Typ “CHAID”, “QUEST” und “C&R-Baum” fügt diese Option ein Feld in der Scoring-Ausgabe hinzu, das die ID für den Endknoten angibt, dem der jeweilige Datensatz zugewiesen ist.

*Hinweis:* Wenn diese Option ausgewählt wurde, ist die SQL-Erzeugung nicht verfügbar.

**SQL für dieses Modell generieren.** Bei Verwendung von Daten aus einer Datenbank kann SQL-Code per Pushback zur Ausführung an die Datenbank zurückübertragen werden. Dadurch kann bei vielen Operationen eine bessere Leistung erzielt werden. [Für weitere Informationen siehe Thema SQL-Optimierung in Kapitel 6 in IBM SPSS Modeler Server 14.2-Verwaltungs- und Leistungshandbuch.](#)

Wählen Sie eine der folgenden Optionen aus, um die SQL-Erzeugung zu aktivieren bzw. deaktivieren.

- **Nicht generieren.** Wählen Sie diese Option aus, um die SQL-Erzeugung für das Modell zu deaktivieren.
- **Ohne Unterstützung für fehlende Werte generieren.** Wählen Sie diese Option, um die SQL-Erzeugung ohne den Aufwand für den Umgang mit fehlenden Werten zu aktivieren. Bei dieser Option wird die Vorhersage einfach auf null gesetzt (\$null\$), wenn beim Scoring eines Falles ein fehlender Wert gefunden wird.

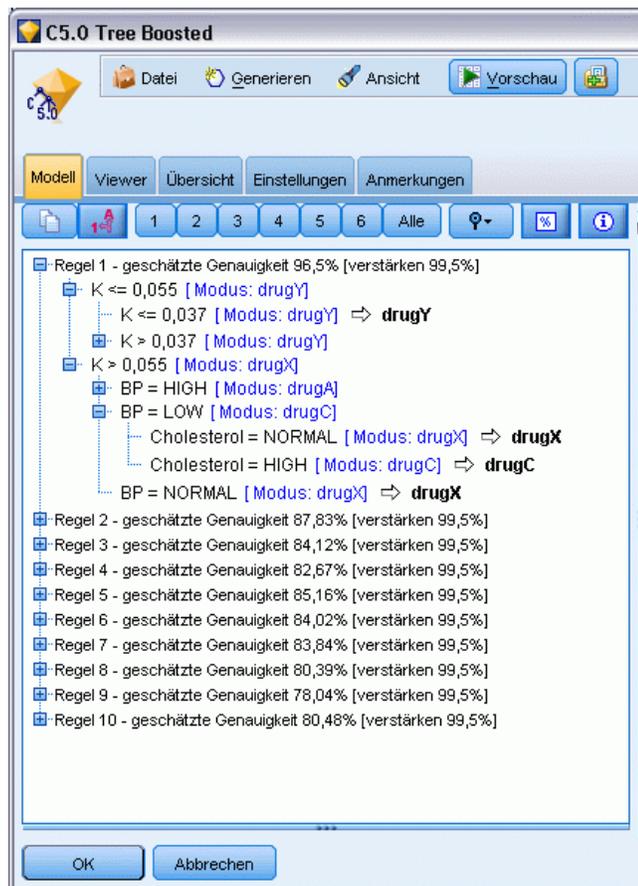
*Hinweis:* Diese Option ist für CHAID-Modelle nicht verfügbar. Bei anderen Modelltypen steht sie nur für Entscheidungsbäume (nicht für Regelmengen) zur Verfügung.

- **Mit Unterstützung für fehlende Werte generieren.** Bei Modellen vom Typ “CHAID”, “QUEST” und “C&R-Baum” können sie die SQL-Erzeugung mit vollständiger Unterstützung für fehlende Werte aktivieren. Dabei wird SQL so generiert, dass die fehlenden Werte als im Modell angegeben behandelt werden. Bei C&RT-Bäumen beispielsweise werden Ersetzungsregeln und ein Rückgriff auf den größten Unterknoten verwendet.

*Hinweis:* Bei C5.0-Modellen steht diese Option nur für Regelmengen (nicht für Entscheidungsbäume) zur Verfügung.

### Verbesserte C5.0-Modelle

Abbildung 6-41  
Verbessertes C5.0 Modell-Nugget, Registerkarte “Modell”



Beim Erstellen eines verstärkten C5.0-Modells (entweder Regelmenge oder Entscheidungsbaum) wird eigentlich ein Set aus verwandten Modellen erstellt. Der Modellregel-Browser für ein verstärktes C5.0-Modell zeigt die Liste der Modelle auf der obersten Ebene der Hierarchie an. Außerdem wird die geschätzte Genauigkeit der einzelnen Modelle und die Gesamtgenauigkeit des Ensembles der verstärkten Modelle angezeigt. Um die Regeln oder Aufteilungen für ein

bestimmtes Modell zu untersuchen, wählen Sie das betreffende Modell aus und erweitern Sie es, wie bei einer Regel oder Verzweigung in einem einzelnen Modell.

Außerdem können Sie ein bestimmtes Modell aus dem Set der verstärkten Modelle extrahieren und ein Modell-Nugget vom Typ "Regelmenge" erstellen, das nur dieses eine Modell enthält. Um eine neue Regelmenge aus einem verstärkten C5.0-Modell zu erstellen, wählen Sie die gewünschte Regelmenge bzw. den gewünschten Baum aus und wählen Sie im Menü "Generieren" entweder die Option Einzelner Entscheidungsbaum (Palette der generierten Modelle) oder Einzelner Entscheidungsbaum (Zeichenbereich).

### **Erzeugen von Diagrammen**

Die Baumknoten bieten eine Menge an Informationen, jedoch sind diese nicht unbedingt immer in einem leicht zugänglichen Format für Unternehmensanwender. Zur Bereitstellung der Daten auf eine Art, die problemlos in Geschäftsberichte, Präsentationen u.s.w. integriert werden kann, können aus ausgewählten Daten Diagramme erstellt werden. Beispielsweise können Sie von der Registerkarte "Modell" oder "Viewer" eines Modell-Nuggets oder der Registerkarte "Viewer" eines interaktiven Baums ein Diagramm für den ausgewählten Teil eines Baums generieren und damit nur ein Diagramm für die Fälle im ausgewählten Baum oder Verzweigungsknoten erstellen.

*Hinweis:* Sie können aus einem Nugget nur dann ein Diagramm erstellen, wenn das Nugget mit anderen Knoten in einem Stream verbunden ist.

#### **Ein Diagramm erstellen**

Wählen Sie als ersten Schritt die Informationen aus, die im Diagramm gezeigt werden sollen:

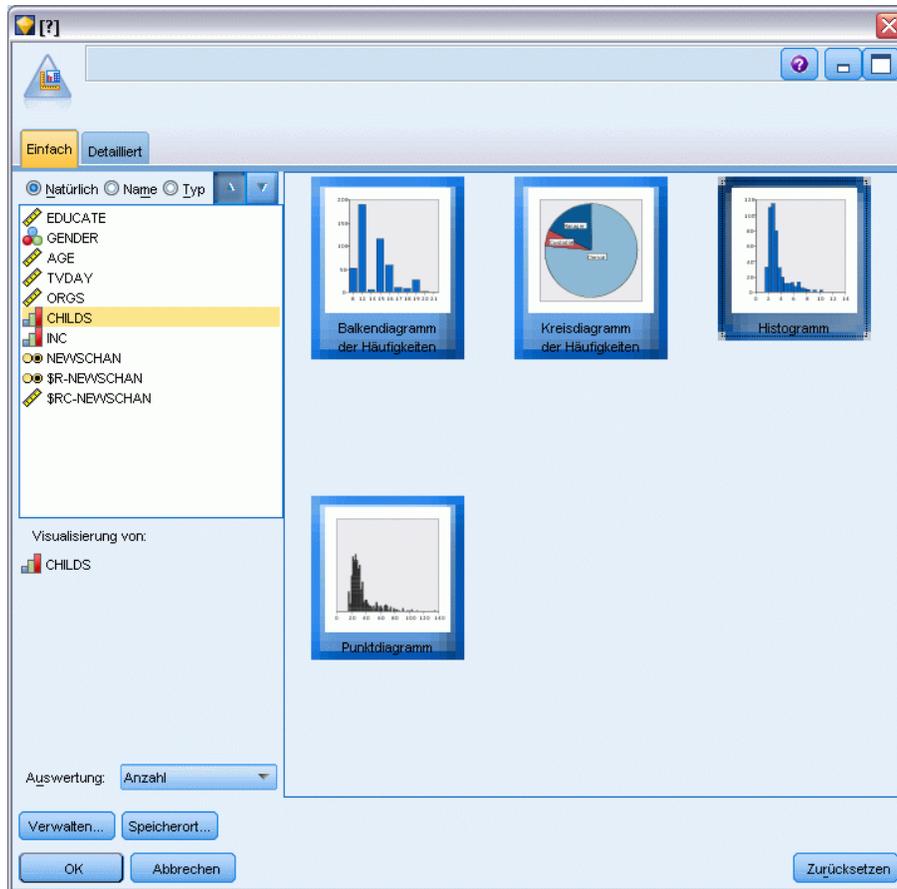
- Erweitern Sie auf der Registerkarte "Modell" eines Nuggets die Liste der Bedingungen und Regeln im rechten Fensterbereich und wählen Sie die gewünschte aus.
- Erweitern Sie auf der Registerkarte "Viewer" eines Nuggets die Liste der Verzweigungen und wählen Sie den gewünschten Knoten aus.
- Erweitern Sie auf der Registerkarte "Viewer" eines interaktiven Baums die Liste der Verzweigungen und wählen Sie den gewünschten Knoten aus.

*Hinweis:* Der erste Knoten in der Hierarchie lässt sich auf keiner Viewer-Registerkarte auswählen.

Unabhängig von den anzuzeigenden Daten erstellen Sie ein Diagramm immer auf die gleiche Weise:

- ▶ Wählen Sie aus dem Menü "Generieren" die Option Diagramm (von Auswahl). Oder klicken Sie alternativ in der Registerkarte "Viewer" auf die Schaltfläche Diagramm (von Auswahl) in der unteren linken Ecke. Die Registerkarte "Einfach" der Diagrammtafel wird angezeigt.

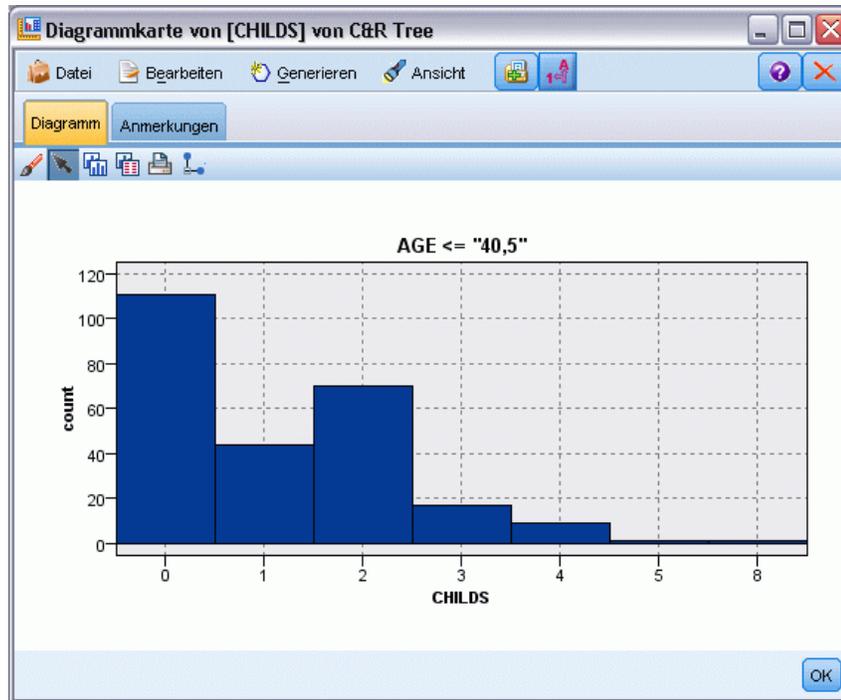
Abbildung 6-42  
Dialogfeld Grafiktafelknoten, Registerkarte "Basis"



*Hinweis:* Wenn Sie die Grafiktafel auf diese Art anzeigen, sind nur die Registerkarten "Basis" und "Details" verfügbar. Für weitere Informationen siehe Thema Diagrammtafelknoten in Kapitel 5 in *IBM SPSS Modeler 14.2- Quellen-, Prozess- und Ausgabeknoten*.

- ▶ Mithilfe der Einstellungen auf den Registerkarten "Basis" oder "Details" können Sie die Details angeben, die auf dem Diagramm angezeigt werden sollen.
- ▶ Klicken Sie auf "OK", um das Diagramm zu erstellen.

Abbildung 6-43  
Histogramm erzeugen mit der Grafiktafelregisterkarte "Basis"



Die Überschrift des Diagramms identifiziert die Knoten oder Regeln, die im Diagramm berücksichtigt werden.

### **Modell-Nuggets zur Verbesserung/Verstärkung und für sehr große Daten-Sets**

Wenn Sie am Modellierungsknoten Modellgenauigkeit erhöhen (verbessern), Modellstabilität steigern (verstärken) oder Modell für sehr große Daten-Sets erstellen als Hauptziel auswählen, erstellt IBM® SPSS® Modeler ein Ensemble aus mehreren Modellen. [Für weitere Informationen siehe Thema Modelle für Ensembles in Kapitel 3 auf S. 57.](#)

Das daraus resultierende Modell-Nugget enthält die folgenden Registerkarten. Die Registerkarte "Modell" bietet mehrere unterschiedliche Modellansichten.

Tabulator	Ansicht	Beschreibung	Weitere Informationen
Modell	Modellzusammenfassung	Zeigt eine Übersicht der Ensemble-Qualität und (mit Ausnahme von verbesserten Modellen und stetigen Zielen) -Vielfältigkeit an, ein Maß dafür, wie stark die Vorhersagen unter den verschiedenen Modellen abweichen.	<a href="#">Für weitere Informationen siehe Thema Modellzusammenfassung in Kapitel 3 auf S. 60.</a>
	Bedeutsamkeit des Prädiktors	Zeigt eine Tabelle an, die die relative Wichtigkeit der einzelnen Prädiktoren	<a href="#">Für weitere Informationen siehe Thema Bedeutsamkeit</a>

Tabulator	Ansicht	Beschreibung	Weitere Informationen
		(Eingabefeld) für die Schätzung des Modells angibt.	des Prädiktors in Kapitel 3 auf S. 61.
	Prädiktorhäufigkeit	Zeigt eine Tabelle an, die die relative Häufigkeit angibt, mit der jeder Prädiktor in dem Modell-Set verwendet wird.	Für weitere Informationen siehe Thema Prädiktorhäufigkeit in Kapitel 3 auf S. 62.
	Komponenten-Modell-Genauigkeit	Zeigt ein Diagramm der Vorhersagegenauigkeit für jedes Modell im Ensemble.	
	Komponenten-Modell-Details	Zeigt Informationen über jedes Modell im Ensemble an.	Für weitere Informationen siehe Thema Komponentenmodelldetails in Kapitel 3 auf S. 64.
	Information	Zeigt Informationen über die Felder, die Aufbaueinstellungen und die Modellschätzung an.	Für weitere Informationen siehe Thema Modell-Nugget-Übersicht/-Information in Kapitel 3 auf S. 54.
Einstellungen		Hier können Sie Konfidenzen in die Scoring-Operation aufnehmen.	Für weitere Informationen siehe Thema Einstellungen für Modell-Nuggets vom Typ "Entscheidungsbaum"/"Regelmenge" auf S. 188.
Anmerkung		Hier können Sie beschreibende Anmerkungen hinzufügen, einen benutzerdefinierten Namen angeben, QuickInfo-Text hinzufügen und Suchwörter für das Modell angeben.	Für weitere Informationen siehe Thema Anmerkungen in Kapitel 5 in <i>IBM SPSS Modeler 14.2-Benutzerhandbuch</i> .

## Regelmengen-Modell-Nuggets

Modell-Nuggets vom Typ "Regelmenge" stellen die Regeln für die Vorhersage eines bestimmten Ausgabefelds dar, das von Modellierungsknoten für Assoziationsregeln ("A Priori") oder von einem der Baumerstellungsknoten (C&R-Baum, CHAID, QUEST oder C5.0) entdeckt wurde. Bei Assoziationsregeln muss die Regelmenge aus einem Nugget für nicht verfeinerte Regeln generiert werden. Bei Bäumen können Regelmengen über den Tree Builder, aus einem C.50-Modellierungsknoten oder aus einem beliebigen aum-Modell-Nugget generiert werden. Im Gegensatz zu Nuggets vom Typ "Nicht verfeinerte Regel" können Nuggets vom Typ "Regelmenge" in Streams platziert werden, um Prognosen zu generieren.

Bei der Ausführung eines Streams, der einen Regelmengen-Nugget enthält, werden zwei neue Felder, die den vorhergesagten Wert und die Konfidenz für die einzelnen Datensätze enthalten, zum Stream hinzugefügt. Die neuen Feldnamen werden durch Hinzufügen von Präfixen aus dem Modellnamen abgeleitet. Bei Assoziationsregelmengen lautet das Präfix \$A- für das Vorhersagefeld und \$AC- für das Konfidenzfeld. Bei C5.0-Regelmengen lautet das Präfix \$C- für das Vorhersagefeld und \$CC- für das Konfidenzfeld. Bei C&RT-Regelmengen lautet das Präfix \$R- für das Vorhersagefeld und \$RC- für das Konfidenzfeld. In einem Stream mit mehreren

Regelmengen-Nuggets in einer Reihe, die dieselben Ausgabefelder vorhersagen, enthalten die neuen Feldnamen Zahlen im *Präfix*, damit sie auseinander gehalten werden können. Beim ersten Assoziations-Regelmengenknoten im Stream werden die üblichen Namen verwendet, beim zweiten Knoten Namen, die mit *\$AI-* und *\$ACI-* beginnen, beim dritten Knoten Namen mit *\$A2-* und *\$AC2-* usw.

**Anwendung der Regeln.** Aus Assoziationsregeln erstellte Regelmengen sind anders als andere Modell-Nuggets, da für jeden Datensatz mehrere Vorhersagen generiert werden können und diese Vorhersagen nicht unbedingt alle übereinstimmen. Es gibt zwei Methoden zum Generieren von Vorhersagen aus Regelmengen:

*Hinweis:* (Regelmengen, die aus Entscheidungsbäumen generiert wurden, geben unabhängig von der verwendeten Methode dieselben Ergebnisse aus, da die aus einem Entscheidungsbaum abgeleiteten Regeln sich gegenseitig ausschließen.)

- **Voting.** Bei dieser Methode wird versucht, die Vorhersagen aller Regeln, die für den Datensatz gelten, zu kombinieren. Bei jedem Datensatz werden alle Regeln untersucht, und jede Regel, die für den Datensatz gilt, wird verwendet, um eine Prognose und eine zugehörige Konfidenz zu generieren. Für jeden Ausgabewert werden die Zahlen für die Summe der Konfidenz berechnet und der Wert mit der größten Konfidenzsumme wird als endgültige Prognose ausgewählt. Die Konfidenz für die endgültige Vorhersage ist die Konfidenzsumme für diesen Wert dividiert durch die Anzahl der Regeln, die für diesen Datensatz ausgelöst wurden.
- **Erster Treffer.** Diese Methode testet einfach die Regeln in der Reihenfolge und die erste Regel, die für den Datensatz gilt, wird zur Generierung der Prognose verwendet.

Die verwendete Methode kann in den Stream-Optionen festgelegt werden. [Für weitere Informationen siehe Thema Festlegen von Optionen für Streams in Kapitel 5 in IBM SPSS Modeler 14.2- Benutzerhandbuch.](#)

**Generieren von Knoten.** Im Menü “Generieren” können Sie anhand der Regelmenge neue Knoten erstellen.

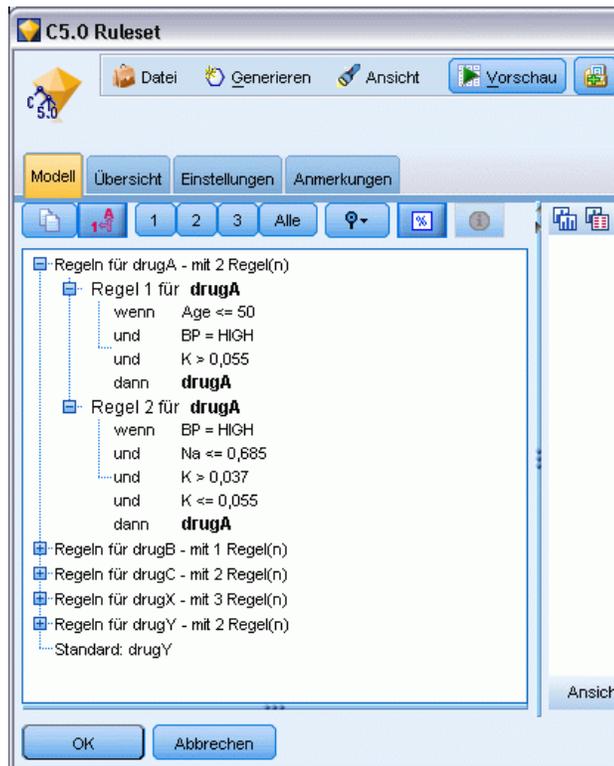
- **Filterknoten.** Erstellt einen neuen Filterknoten zum Filtern von Feldern, die nicht von den Regeln in der Regelmenge verwendet werden.
- **Auswahlknoten.** Erstellt einen neuen Auswahlknoten zur Auswahl von Datensätzen, für die die ausgewählte Regel gilt. Der generierte Knoten wählt Datensätze aus, für die alle Antezedenzen der Regel wahr sind. Für diese Option muss eine Regel ausgewählt sein.
- **Regel-Verfolgungsknoten.** Erstellt einen neuen Superknoten und berechnet ein Feld, das angibt, welche Regel zur Erstellung der Prognose für die einzelnen Datensätze verwendet wurde. Bei der Auswertung einer Regelmenge mithilfe der Methode “Erster Treffer” ist dies einfach ein Symbol, das die erste Regel angibt, die ausgelöst werden würde. Bei der Auswertung der Regelmenge mithilfe der Methode “Voting” ist dies eine komplexere Zeichenkette, die die Eingabe in den Voting-Mechanismus anzeigt.
- **Einzelner Entscheidungsbaum (Zeichenbereich)/Einzelner Entscheidungsbaum (Palette der generierten Modelle).** Erstellt eine einzelne Regelmenge, die aus der aktuell ausgewählten Regel abgeleitet wurde. Nur für **verstärkte** C5.0-Modelle verfügbar. [Für weitere Informationen siehe Thema Verbesserte C5.0-Modelle auf S. 190.](#)
- **Modell zur Palette hinzufügen.** Gibt das Modell an die Modellpalette zurück. Das ist nützlich, wenn Sie von einem Kollegen einen Datenstrom, der das Modell enthält, jedoch nicht das Modell selbst erhalten.

*Hinweis:* Die Registerkarten “Einstellungen” und “Übersicht” im Regelmengen-Nugget stimmen mit den von den Entscheidungsbaummodellen verwendeten überein.

## Regelmenge – Registerkarte “Modell”

Auf der Registerkarte “Modell” für ein Regelmengen-Nugget wird eine Liste der Regeln angezeigt, die der Algorithmus aus den Daten extrahiert hat.

Abbildung 6-44  
Regelmengen-Modell-Nugget, Registerkarte “Modell”



Die Regeln werden nach Sukzedens (vorhergesagte Kategorie) aufgeschlüsselt und in folgendem Format angezeigt:

```
wenn Antezedens_1
und Antezedens_2
...
und Antezedens_n
dann Sukzedens
```

wobei consequent und antecedent\_1 bis antecedent\_n alle Bedingungen sind. Die Regel wird als “für Datensätze, bei denen antecedent\_1 bis antecedent\_n alle wahr sind, ist consequent wahrscheinlich ebenfalls wahr” interpretiert. Wenn Sie auf die Schallfläche Instanzen/Konfidenz anzeigen in der Symbolleiste klicken, werden bei jeder Regel auch die Informationen darüber angezeigt, für wie viele Datensätze die Regel gilt, also für wie viele Datensätze die Antezedenzen

wahr sind (**Instanzen**), sowie der Anteil der Datensätze, für die die gesamte Regel wahr ist (**Konfidenz**).

Beachten Sie, dass die Berechnungsmethode für die Konfidenz bei C5.0-Regelmengen ein wenig abweicht. Bei C5.0 wird folgende Formel zur Berechnung der Konfidenz einer Regel verwendet:

$$(1 + \text{Anzahl der Datensätze, bei denen die Regel richtig ist}) / (2 + \text{Anzahl der Datensätze, bei denen die Antezedenzen der Regel wahr sind})$$

Diese Berechnung des Konfidenzschätzers wird für den Prozess der Generalisierung von Regeln aus einem Entscheidungsbaum (das Verfahren, mit dem bei C5.0 Regelmengen erstellt werden) angepasst.

## Projekte aus AnswerTree 3.0 importieren

IBM® SPSS® Modeler kann in AnswerTree 3.0 oder 3.1 gespeicherte Modelle über das Standarddialogfeld "Datei" > "Öffnen" importieren:

- ▶ Wählen Sie in den SPSS Modeler-Menüs folgende Befehlsfolge:  
Datei > Stream öffnen
- ▶ Wählen Sie in der Dropdown-Liste "Dateityp" die Option AT Project Files (\*.atp, \*.ats).  
Jedes importierte Projekt wird in einen SPSS Modeler-Stream konvertiert und enthält folgende Knoten:
  - einen Quellenknoten, der die verwendete Datenquelle definiert (z. B. eine IBM® SPSS® Statistics-Datendatei oder eine Datenbankquelle).
  - Für jedem im Projekt vorhandenen Baum (es kann mehrere enthalten) wird ein Typknoten erstellt, der die Eigenschaften aller Felder (Variablen) definiert, mit Typ, Rolle (Eingabe oder Prädiktorfeld <> Ausgabe oder vorhergesagtes Feld), fehlenden Werten und anderen Optionen.
  - Für jeden im Projekt enthaltenen Baum wird ein Partitionsknoten erstellt, der die Daten für eine Trainings- und eine Teststichprobe partitioniert. Ein Baumerstellungsknoten wird erstellt, der die Parameter definiert, mit denen der Baum generiert wird (Knoten vom Typ "C&R-Baum", "QUEST" oder "CHAID").
- ▶ Führen Sie den Stream aus, um den oder die generierten Bäume anzuzeigen.

### Kommentare

- In SPSS Modeler generierte Entscheidungsbäume können nicht in AnswerTree exportiert werden. Der Import von AnswerTree in SPSS Modeler ist nicht umkehrbar.
- In AnswerTree definierte Profite gehen beim Importieren des Projekts in SPSS Modeler verloren.

# Bayes-Netzwerk-Modelle

## Bayes-Netzwerk-Knoten

Mithilfe des Knotens **Bayes-Netzwerk** können Sie ein Wahrscheinlichkeitsmodell erstellen, indem Sie beobachtete und aufgezeichnete Hinweise mit Weltwissen (“gesundem Menschenverstand”) kombinieren, um die Wahrscheinlichkeit des Vorkommens unter Verwendung scheinbar nicht miteinander verknüpfter Attribute zu ermitteln. Der Knoten ist speziell für Netzwerke vom Typ “Tree Augmented Naïve Bayes” (TAN) und “Markov-Decke” gedacht, die in erster Linie zur Klassifizierung verwendet werden.

Bayes-Netzwerke dienen zur Erstellung von Vorhersagen in vielen verschiedenen Situationen. Hier einige Beispiele:

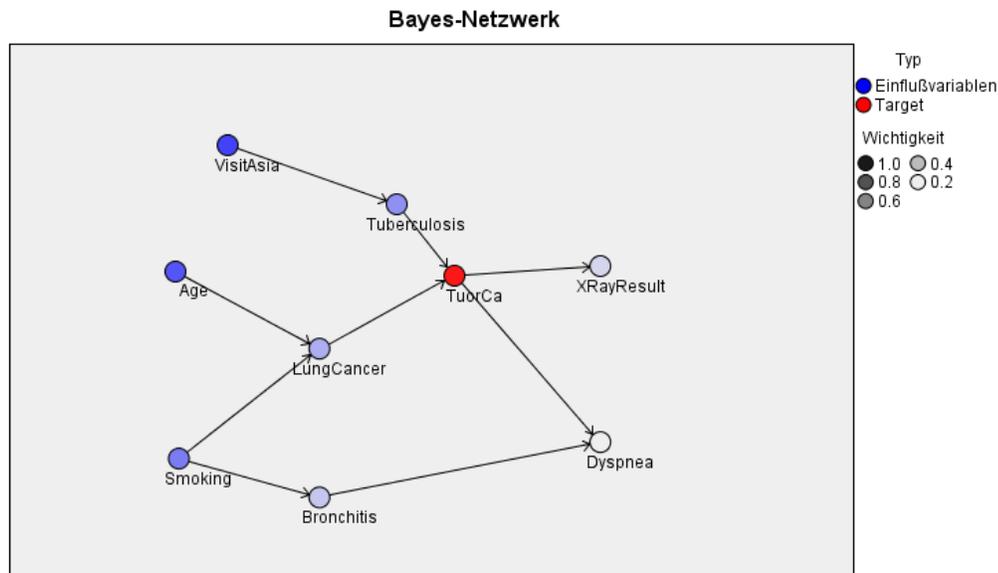
- Ermitteln von Antragstellern für Kredite mit geringem Risiko der Zahlungsunfähigkeit.
- Abschätzung, zu welchem Zeitpunkt für Geräte Wartungsarbeiten, Ersatzteile oder ein Austausch erforderlich ist, basierend auf Sensoreingaben und bestehenden Datensätzen.
- Lösen von Kundenproblemen mithilfe von Online-Tools zur Fehlerbehebung.
- Diagnostizierung und Fehlerbehebung bei Mobiltelefonnetzwerken in Echtzeit.
- Abschätzen des potenziellen Risikos und Nutzens von Forschungs- und Entwicklungsprojekten mit dem Ziel, die Ressourcen auf die aussichtsreichsten Möglichkeiten zu konzentrieren.

Ein Bayes-Netzwerk ist ein grafisches Modell, das Variablen (häufig als **Knoten** bezeichnet) in einem Daten-Set und die probabilistischen bzw. bedingten Unabhängigkeiten zwischen den Variablen anzeigt. Durch Bayes-Netzwerke können kausale Beziehungen zwischen Knoten angezeigt werden; die Verbindungen in den Netzwerken (auch als **arcs** (Bögen) bezeichnet) stehen jedoch nicht unbedingt für ein direktes Verhältnis von Ursache und Wirkung. Mithilfe eines Bayes-Netzwerks kann beispielsweise die Wahrscheinlichkeit berechnet werden, dass ein Patient unter einer bestimmten Krankheit leidet. Diese Berechnung beruht auf dem Vorliegen bzw. Fehlen bestimmter Symptome sowie anderen relevanten Daten und setzt voraus, dass sich die probabilistischen Unabhängigkeiten zwischen Symptomen und Krankheiten, die im Diagramm angezeigt werden, als wahr erweisen. Netzwerke sind bei fehlenden Informationen sehr robust und führen zu der bestmöglichen Vorhersage unter Nutzung aller vorhandenen Informationen.

Ein typisches, einfaches Beispiel eines Bayes-Netzwerks wurde von Lauritzen und Spiegelhalter (1988) erstellt. Es wird häufig als das “Asien-Modell” bezeichnet und ist eine vereinfachte Version eines Netzwerks, das zur Diagnostizierung neuer Patienten eines Arztes verwendet werden kann. Die Richtung der Verbindungen entspricht in etwa der Kausalität. Jeder Knoten steht für eine Facette, die mit dem Zustand des Patienten in Beziehung stehen könnte. So zeigt “Smoking” (Rauchen) an, dass sie Raucher sind und “VisitAsia” (BesuchAsien) zeigt an, ob sie sich in letzter Zeit in Asien aufgehalten haben. Wahrscheinlichkeitsbeziehungen werden durch die Verbindungen zwischen Knoten angezeigt. So erhöht das Rauchen beispielsweise sowohl die Wahrscheinlichkeit, dass ein Patient Bronchitis entwickelt, als auch die Wahrscheinlichkeit, dass er Lungenkrebs entwickelt, wohingegen das Alter nur mit der Möglichkeit der Entwicklung von

Lungenkrebs verknüpft zu sein scheint. Ebenso gilt, dass Anomalien auf einem Röntgenbild der Lungen entweder durch Tuberkulose oder durch Lungenkrebs hervorgerufen sein können, und die Wahrscheinlichkeit, dass ein Patient an Atemnot (Dyspnoe) leidet, erhöht ist, wenn der Patient auch an Bronchitis oder Lungenkrebs leidet.

Abbildung 7-1  
Asien-Netzwerkbeispiel von Lauritzen und Spiegelhalter



Es gibt mehrere Gründe, die für die Verwendung eines Bayes-Netzwerks sprechen können:

- Es bietet Einblicke in Kausalbeziehungen. Ausgehend davon trägt es zum Verständnis eines Problembereichs bei und ermöglicht es, die Folgen von Eingriffen vorauszusagen.
- Das Netzwerk bietet einen effizienten Ansatz zur Vermeidung einer Überanpassung an die Daten.
- Es bietet eine deutliche Visualisierung der beteiligten Beziehungen.

**Voraussetzungen.** Die Zielfelder müssen kategorial sein und können das Messniveau *Nominal*, *Ordinal* oder *Flag* aufweisen. Als Eingaben kommen Felder jedes Typs infrage. Stetige Eingabefelder (numerischer Bereich) werden automatisch klassiert. Bei verzerrten Verteilungen erzielen Sie jedoch möglicherweise bessere Ergebnisse, wenn Sie die Felder vor dem Bayes-Netzwerk-Knoten manuell mithilfe eines Klassierknotens klassieren. Verwenden Sie beispielsweise "Optimales Klassieren", wenn das Supervisor-Feld mit dem Feld Ziel des Bayes-Netzwerk-Knotens übereinstimmt. [Für weitere Informationen siehe Thema Klassierknoten in Kapitel 4 in IBM SPSS Modeler 14.2- Quellen-, Prozess- und Ausgabeknoten.](#)

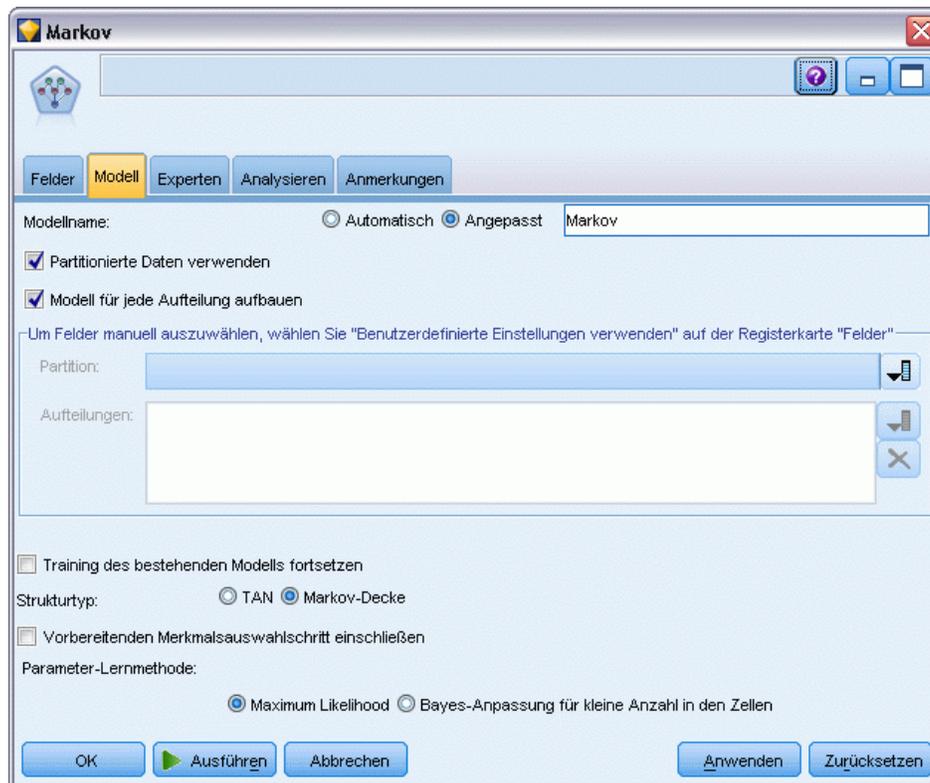
**Beispiel.** Ein Analyst einer Bank möchte in der Lage sein vorherzusagen, welche Kunden bzw. potenzielle Kunden mit hoher Wahrscheinlichkeit mit ihren Kreditrückzahlungen in Verzug geraten. Mithilfe eines Bayes-Netzwerk-Modells können Sie die Eigenschaften der Kunden ermitteln, die mit der höchsten Wahrscheinlichkeit zahlungsunfähig werden, und mehrere verschiedene Arten von Modellen erstellen, um herauszufinden, welches davon die potenziell zahlungsunfähigen Personen am besten vorhersagt. [Für weitere Informationen siehe Thema](#)

Vorhersage von Kreditausfällen (Bayes-Netzwerk) in Kapitel 18 in *IBM SPSS Modeler 14.2-Anwendungshandbuch*.

**Beispiel.** Ein Telekommunikationsbetreiber möchte die Anzahl der Kunden, die abwandern, verringern und das Modell monatlich mit den Daten des jeweiligen Vormonats aktualisieren. Mithilfe eines Bayes-Netzwerk-Modells können Sie die Eigenschaften der Kunden ermitteln, die mit der höchsten Wahrscheinlichkeit abwandern und das Training des Modells jeden Monat mit den neuen Daten fortsetzen. Für weitere Informationen siehe Thema [Erneutes Trainieren eines Modells auf monatlicher Basis \(Bayes-Netzwerk\)](#) in Kapitel 19 in *IBM SPSS Modeler 14.2-Anwendungshandbuch*.

## Modelloptionen für Bayes-Netzwerk-Knoten

Abbildung 7-2  
Bayes-Netzwerk-Knoten: Registerkarte "Modell"



**Modellname.** Sie können den Modellnamen automatisch basierend auf den Ziel- oder ID-Feldnamen (oder dem Modelltyp in Fällen, in denen kein solches Feld angegeben ist) generieren oder einen benutzerdefinierten Namen eingeben.

**Partitionierte Daten verwenden.** Wenn ein Partitionsfeld definiert ist, gewährleistet diese Option, dass nur Daten aus der Trainingspartition für die Modellerstellung verwendet werden. Für weitere Informationen siehe Thema [Partitionsknoten in Kapitel 4 in IBM SPSS Modeler 14.2- Quellen-, Prozess- und Ausgabeknoten](#).

**Modell für jede Aufteilung erstellen.** Erstellt ein separates Modell für jeden möglichen Wert von Eingabefeldern, die als Aufteilungsfelder angegeben werden. [Für weitere Informationen siehe Thema Erstellung von aufgeteilten Modellen in Kapitel 3 auf S. 31.](#)

**Partition.** In diesem Feld können Sie ein Feld angeben, das verwendet wird, um die Daten in getrennte Stichproben für die Trainings-, Test- und Validierungsphase der Modellbildung aufzuteilen. Indem Sie mit einer Stichprobe das Modell erstellen und es mit einer anderen Stichprobe testen, erhalten Sie einen guten Hinweis dafür, wie gut das Modell sich für größere Datenmengen generalisieren lässt, die den aktuellen Daten ähneln. Wenn mehrere Partitionsfelder mithilfe von Typ- oder Partitionsknoten erstellt wurden, muss in jedem Modellierungsknoten, der die Partitionierung verwendet, auf der Registerkarte "Felder" ein einzelnes Partitionsfeld ausgewählt werden. (Wenn nur eine einzige Partition vorhanden ist, wird diese immer automatisch verwendet, wenn die Partitionierung aktiviert ist.) [Für weitere Informationen siehe Thema Partitionsknoten in Kapitel 4 in IBM SPSS Modeler 14.2- Quellen-, Prozess- und Ausgabeknoten.](#) Beachten Sie außerdem, dass die Partitionierung auch auf der Registerkarte "Modelloptionen" des Knotens aktiviert sein muss, wenn die ausgewählte Partitionierung in Ihrer Analyse angewendet werden soll. (Wenn die Auswahl der Option aufgehoben ist, kann die Partitionierung ohne Änderung der Feldeinstellungen deaktiviert werden.)

**Aufteilen.** Wählen Sie für Aufteilungsmodelle das Aufteilungsfeld bzw. die Aufteilungsfelder. Dies ist so, als würden Sie in einem Typknoten für die Rolle eines Felds den Wert *Aufteilung* festlegen. Sie können nur Felder mit einem Messniveau Flag, Nominal, Ordinal oder Stetig als Aufteilungsfelder festlegen. Als Aufteilungsfelder gewählte Felder können nicht als Ziel-, Prädiktor-, Partitions-, Häufigkeits- oder Gewichtungsfelder verwendet werden. [Für weitere Informationen siehe Thema Erstellung von aufgeteilten Modellen in Kapitel 3 auf S. 31.](#)

**Training des bestehenden Modells fortsetzen.** Wenn Sie diese Option auswählen, werden die auf der Registerkarte "Modell" des Modell-Nuggets angezeigten Ergebnisse bei jeder Modellausführung erneut generiert und aktualisiert. Diese Vorgehensweise wird beispielsweise verwendet, wenn eine neue oder aktualisierte Datenquelle zu einem bestehenden Modell hinzugefügt wurde.

*Hinweis:* Damit ist lediglich eine Aktualisierung des bestehenden Netzwerks möglich. Es können keine Knoten oder Verbindungen hinzugefügt bzw. entfernt werden. Bei jedem erneuten Trainieren des Modells behält das Netzwerk dieselbe Form bei. Es ändern sich lediglich die bedingten Wahrscheinlichkeiten und die Bedeutsamkeit des Prädiktors. Wenn die neuen Daten im Großen und Ganzen den alten Daten ähneln, ist dies nicht von Bedeutung, da davon auszugehen ist, dass dieselben Elemente von Bedeutung sind. Wenn Sie jedoch überprüfen bzw. aktualisieren möchten, *was* von Bedeutung ist (im Gegensatz dazu, *wie* bedeutsam es ist), müssen Sie ein neues Modell, also ein neues Netzwerk, erstellen.

**Strukturtyp.** Dient zur Auswahl der beim Erstellen des Bayes-Netzwerks zu verwendenden Struktur:

- **TAN.** Das TAN-Modell (Tree Augmented Naïve Bayes) erstellt ein einfaches Bayes-Netzwerk-Modell, das eine Verbesserung gegenüber dem standardmäßigen Naïve Bayes-Modell darstellt. Dies liegt daran, dass dabei jeder Prädiktor neben der

Zielvariablen auch von einem anderen Prädiktor abhängig sein kann, wodurch die Klassifizierungsgenauigkeit erhöht wird.

- **Markov-Decke.** Dient zur Auswahl der Gruppe von Knoten im Daten-Set, die die übergeordneten Variablen der Zielvariablen, ihre untergeordneten Variablen sowie die Variablen auswählt, die ihren untergeordneten Variablen übergeordnet sind. Im Wesentlichen identifiziert eine Markov-Decke alle Variablen im Netzwerk, die zur Vorhersage der Zielvariablen erforderlich sind. Die Methode zur Erstellung eines Netzwerks gilt als genauer. Bei großen Daten-Sets kann sie jedoch aufgrund der großen Anzahl an Variablen den Nachteil einer längeren Verarbeitungsdauer mit sich bringen. Um den Verarbeitungsaufwand zu verringern, können Sie mithilfe der Optionen Merkmalsauswahl auf der Registerkarte “Experten” die Variablen auswählen, die in einer signifikanten Beziehung zur Zielvariablen stehen.

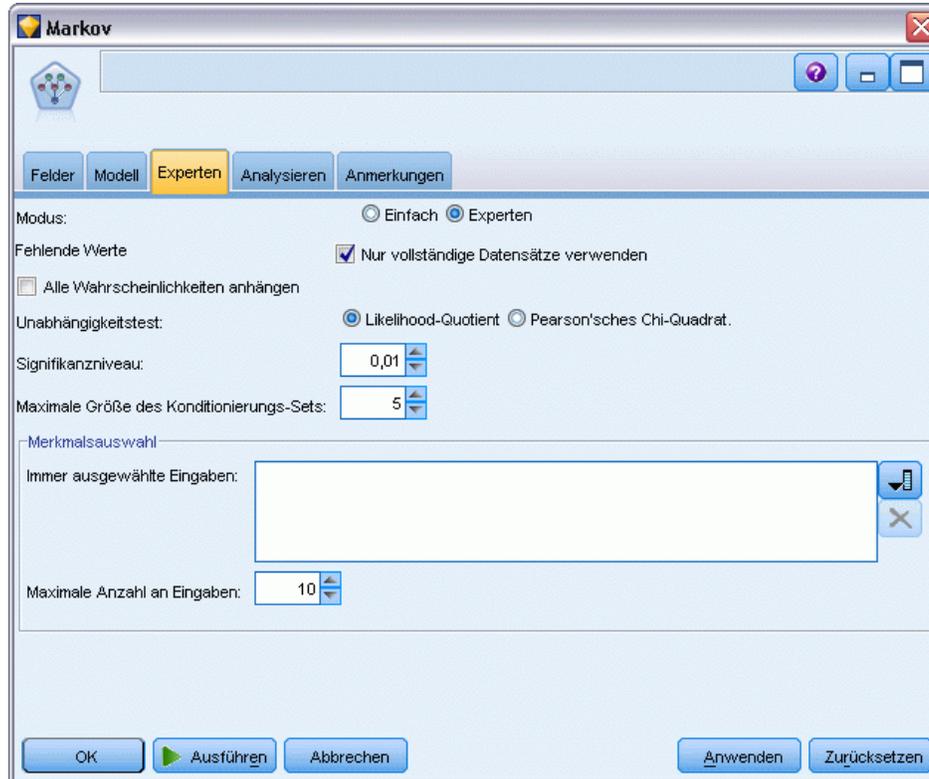
**Vorbereitenden Merkmalsauswahlschritt einschließen.** Durch Auswahl dieses Kontrollkästchens können Sie die Optionen zur Merkmalsauswahl auf der Registerkarte “Experten” nutzen.

**Parameter-Lernmethode.** Die Parameter von Bayes-Netzwerken beziehen sich auf die bedingten Wahrscheinlichkeiten für die einzelnen Knoten in Abhängigkeit von den Werten seiner jeweiligen übergeordneten Elemente. Es gibt zwei verschiedene Auswahlmöglichkeiten, mit denen Sie die Aufgabe zur Schätzung der Tabellen zur bedingten Wahrscheinlichkeit zwischen Knoten beeinflussen können, wenn die Werte der übergeordneten Elemente bekannt sind:

- **Maximum Likelihood.** Aktivieren Sie dieses Kontrollkästchen, wenn Sie ein großes Daten-Set verwenden. Dies ist die Standardauswahl.
- **Bayes-Anpassung für kleine Anzahl in den Zellen.** Bei kleineren Daten-Sets besteht die Gefahr einer Überanpassung des Modells sowie die Möglichkeit einer hohen Anzahl von Zellen mit einer Zellenhäufigkeit von 0. Mit dieser Option können Sie diese Probleme verringern, indem Sie Glättung zur Reduzierung des Effekts von Zellen mit einer Häufigkeit von 0 und etwaiger unzuverlässiger Schätzer anwenden.

## Expertenoptionen für Bayes-Netzwerk-Knoten

Abbildung 7-3  
Bayes-Netzwerk-Knoten: Registerkarte "Experten"



Mit dem Knoten "Expertenoptionen" können Sie die Feinabstimmung der Modellerstellung vornehmen. Für den Zugriff auf die Expertenoptionen setzen Sie den Modus auf der Registerkarte "Experten" auf Experten.

**Fehlende Werte.** Standardmäßig verwendet IBM® SPSS® Modeler nur Datensätze mit gültigen Werten für alle im Modell verwendeten Felder. (Dies wird zuweilen als **listenweiser Ausschluss** fehlender Werte bezeichnet.) Wenn sehr viele fehlende Daten vorliegen, werden mit diesem Ansatz möglicherweise zu viele Datensätze entfernt, sodass nicht mehr genügend Daten zu Erstellung eines guten Modells vorhanden sind. In solchen Fällen können Sie die Auswahl der Option **Nur vollständige Datensätze verwenden** aufheben. SPSS Modeler versucht dann, so viele Informationen wie möglich zur Schätzung des Modells zu verwenden, auch Datensätze, bei denen bei einigen Feldern fehlende Werte vorliegen. (Dies wird zuweilen als **paarweiser Ausschluss** fehlender Werte bezeichnet.) In einigen Situationen jedoch kann eine derartige Verwendung unvollständiger Datensätze zu Berechnungsproblemen bei der Schätzung des Modells führen.

**Alle Wahrscheinlichkeiten ausgeben.** Gibt an, ob die Wahrscheinlichkeiten für die einzelnen Kategorien des Ausgabefelds zu den einzelnen vom Knoten verarbeiteten Datensätzen hinzugefügt werden. Wenn diese Option nicht ausgewählt wurde, wird nur die Wahrscheinlichkeit der vorhergesagten Kategorie hinzugefügt.

**Unabhängigkeitstest.** Ein Test auf Unabhängigkeit dient zur Einschätzung, ob gepaarte Beobachtungen bei zwei Variablen voneinander unabhängig sind. Wählen Sie den Typ des zu verwendenden Tests aus. Folgende Optionen stehen zur Verfügung:

- **Likelihood-Quotient.** Testet auf Unabhängigkeit zwischen Ziel und Prädiktor durch Berechnung des Verhältnisses zwischen der maximalen Wahrscheinlichkeit eines Ergebnisses unter zwei verschiedenen Hypothesen.
- **Pearson'sches Chi-Quadrat.** Testet auf Unabhängigkeit zwischen Ziel und Prädiktor unter Verwendung der Nullhypothese, dass die relativen Häufigkeiten des Eintretens beobachteter Ereignisse einer angegebenen Häufigkeitsverteilung folgen.

Bayes-Netzwerk-Modelle führen bedingte Tests auf Unabhängigkeit durch, bei denen über die getesteten Paare hinaus zusätzliche Variablen verwendet werden. Außerdem untersuchen die Modelle nicht nur die Beziehungen zwischen Ziel und Prädiktoren, sondern auch die Beziehungen zwischen den Prädiktoren selbst.

*Hinweis:* Die Optionen für den Unabhängigkeitstest stehen nur zur Verfügung, wenn auf der Registerkarte "Modell" entweder Vorbereitenden Merkmalsauswahlschritt einschließen oder ein Strukturtyp der Markov-Decke ausgewählt wurde.

**Signifikanzniveau.** In Verbindung mit den Einstellungen für den Unabhängigkeitstest können Sie mit dieser Einstellung einen Trennwert festlegen, der bei der Durchführung der Tests verwendet werden soll. Je niedriger der Wert, desto weniger Verbindungen verbleiben im Netzwerk; das Standardniveau ist 0,01.

*Hinweis:* Die Option steht nur zur Verfügung, wenn auf der Registerkarte "Modell" entweder Vorbereitenden Merkmalsauswahlschritt einschließen oder ein Strukturtyp der Markov-Decke ausgewählt wurde.

**Maximale Größe des Konditionierungs-Sets.** Der Algorithmus zum Erstellen einer Struktur vom Typ "Markov-Decke" verwendet Konditionierungs-Sets mit zunehmender Größe, um Unabhängigkeitstests durchzuführen und unnötige Verbindungen aus dem Netzwerk zu entfernen. Da die Tests mit einer höheren Anzahl von Konditionierungsvariablen einen größeren Zeit- und Speicherbedarf für die Verarbeitung aufweisen, können Sie die Anzahl der aufzunehmenden Variablen begrenzen. Dies kann insbesondere bei der Verarbeitung von Daten mit starken Abhängigkeiten zwischen vielen Variablen nützlich sein. Beachten Sie jedoch, dass das so entstehende Netzwerk einige überflüssige Verbindungen enthalten kann.

Geben Sie die Maximalzahl der für die Unabhängigkeitstests zu verwendenden Konditionierungsvariablen an. Die Standardeinstellung lautet 5.

*Hinweis:* Die Option steht nur zur Verfügung, wenn auf der Registerkarte "Modell" entweder Vorbereitenden Merkmalsauswahlschritt einschließen oder ein Strukturtyp der Markov-Decke ausgewählt wurde.

**Merkmalswahl.** Mit dieser Option können Sie die Anzahl der bei der Verarbeitung des Modells verwendeten Eingaben beschränken, um die Modellerstellung zu beschleunigen. Dies ist aufgrund der möglicherweise großen Anzahl potenzieller Eingaben besonders nützlich bei Erstellung einer Struktur vom Typ "Markov-Decke", da Sie auf diese Weise die Eingaben auswählen können, die in einer signifikanten Beziehung zur Zielvariablen stehen.

*Hinweis:* Die Optionen zur Merkmalsauswahl sind nur verfügbar, wenn Sie auf der Registerkarte "Modell" die Option Vorbereitenden Merkmalsauswahlschritt einschließen auswählen.

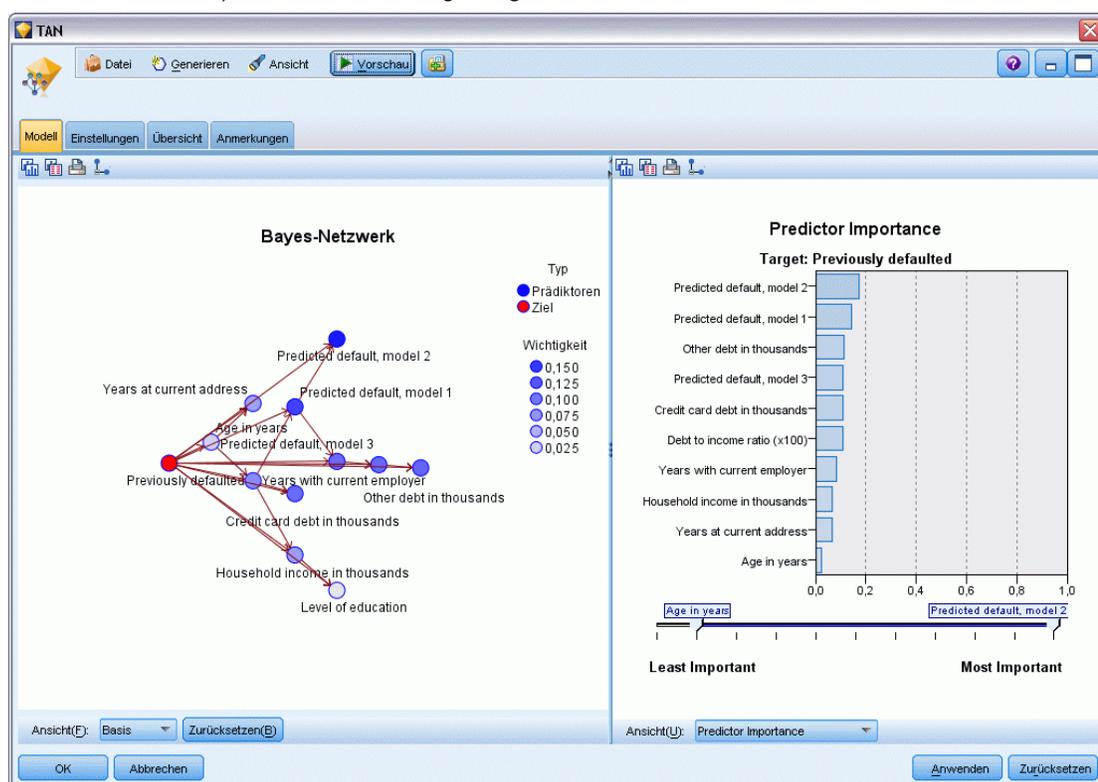
- **Eingaben immer ausgewählt** Mit der Feldauswahl-Schaltfläche rechts neben dem Textfeld können Sie Felder aus dem Daten-Set auswählen, die immer bei Erstellung des Bayes-Netzwerk-Modells verwendet werden sollen. Beachten Sie, dass das Zielfeld immer ausgewählt ist.
- **Maximale Anzahl an Eingaben.** Geben Sie die Gesamtzahl an Eingaben aus dem Daten-Set an, die beim Erstellen des Bayes-Netzwerk-Modells verwendet werden sollen. Als Höchstwert kann die Gesamtzahl der Eingaben im Daten-Set eingegeben werden.

*Hinweis:* Wenn die Anzahl der in Eingaben immer ausgewählt ausgewählten Felder den Wert von Maximale Anzahl an Eingaben überschreitet, wird eine Fehlermeldung angezeigt.

## Modell-Nuggets vom Typ "Bayes-Netzwerk"

Abbildung 7-4

Modelldetails für Bayes-Netzwerk und zugehörige Bedeutsamkeit der Prädiktoren



*Hinweis:* Wenn Sie auf der Registerkarte "Modell" des Modellknotens die Option Training der bestehenden Parameter fortsetzen ausgewählt haben, werden die auf der Registerkarte "Modell" des Modell-Nuggets angezeigten Informationen bei jeder erneuten Generierung des Modells aktualisiert.

Die Registerkarte “Modell” des Modell-Nuggets gliedert sich in zwei Bereiche:

### **Linker Bereich**

**Einfach.** Diese Ansicht enthält ein Netzwerkdiagramm mit Knoten, das die Beziehung zwischen dem Ziel und seinen wichtigsten Prädiktoren sowie die Beziehung zwischen den Prädiktoren anzeigt. Die Bedeutsamkeit der einzelnen Prädiktoren wird durch die Farbdichte angezeigt. Eine vollere Farbe zeigt einen bedeutsamen Prädiktor an und umgekehrt.

Die Klassenwerte für Knoten, die für einen Bereich stehen, werden in einer Popup-QuickInfo angezeigt, wenn Sie mit dem Mauszeiger über den Knoten fahren.

Mit den Diagrammtools von IBM® SPSS® Modeler können Sie in das Diagramm eingreifen, es bearbeiten und speichern. Beispielsweise zur Verwendung in anderen Programmen wie MS Word.

*Tip:* Wenn das Netzwerk sehr viele Knoten enthält, können Sie auf einen Knoten klicken und ihn an eine andere Stelle ziehen, um die Lesbarkeit des Diagramms zu erhöhen.

**Verteilung.** In dieser Ansicht werden die bedingten Wahrscheinlichkeiten für die einzelnen Knoten im Netzwerk als Minidiagramme angezeigt. Bewegen Sie den Mauszeiger über ein Diagramm, um dessen Werte in einer Popup-QuickInfo anzuzeigen.

### **Rechter Bereich**

**Bedeutsamkeit des Prädiktors.** Damit wird eine Tabelle angezeigt, die die relative Wichtigkeit der einzelnen Prädiktoren für die Schätzung des Modells angibt. [Für weitere Informationen siehe Thema Bedeutsamkeit des Prädiktors in Kapitel 3 auf S. 54.](#)

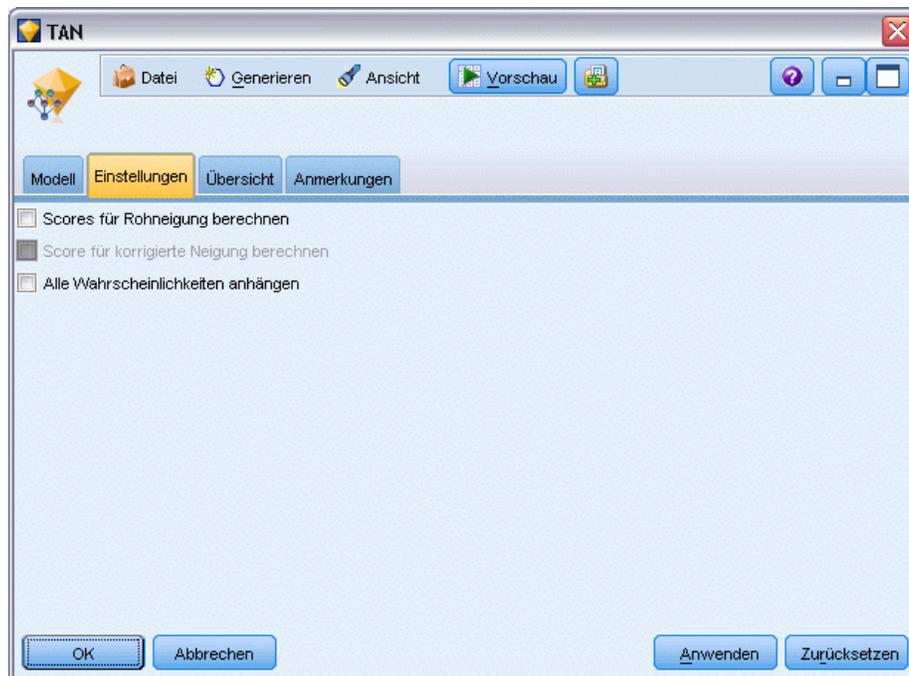
**Bedingte Wahrscheinlichkeiten.** Wenn Sie im linken Bereich einen Knoten oder ein Miniatur-Verteilungsdiagramm auswählen, wird im rechten Bereich die zugehörige Tabelle mit bedingten Wahrscheinlichkeiten angezeigt. Diese Tabelle enthält den Wert der bedingten Wahrscheinlichkeit für die einzelnen Knotenwerte und die einzelnen Kombinationen von Werten in ihren übergeordneten Knoten. Außerdem beinhaltet sie die Anzahl der für die einzelnen Datensatzwerte beobachteten Datensätze und die einzelnen Kombinationen von Werten in den übergeordneten Knoten.

## **Einstellungen im Bayes-Netzwerk-Modell**

Auf der Registerkarte “Einstellungen” für ein Modell-Nugget vom Typ “Bayes-Netzwerk” werden Optionen zum Ändern des erstellten Modells angegeben. Beispielsweise können Sie mit dem Bayes-Netzwerk-Knoten unter Verwendung derselben Daten und Einstellungen mehrere verschiedene Modelle erstellen und anschließend diese Registerkarte bei jedem Modell verwenden, um die Einstellungen leicht abzuändern und zu ermitteln, welche Auswirkungen dies auf die Ergebnisse hat.

*Hinweis:* Diese Registerkarte ist erst verfügbar, nachdem das Modell-Nugget einem Stream hinzugefügt wurde.

Abbildung 7-5  
Registerkarte "Einstellungen" für ein Bayes-Netzwerk-Modell



**Scores für Rohneigung berechnen.** Bei Modellen mit einem Flag-Ziel (das als Vorhersage "Ja" bzw. "Nein" ausgibt) können Sie Neigungs-Scores anfordern, die die Wahrscheinlichkeit des für das Zielfeld angegebenen wahren Ergebnisses angeben. Diese Werte werden zusätzlich zu den anderen Vorhersage- und Konfidenzwerten angegeben, die ggf. während des Scorings erstellt werden.

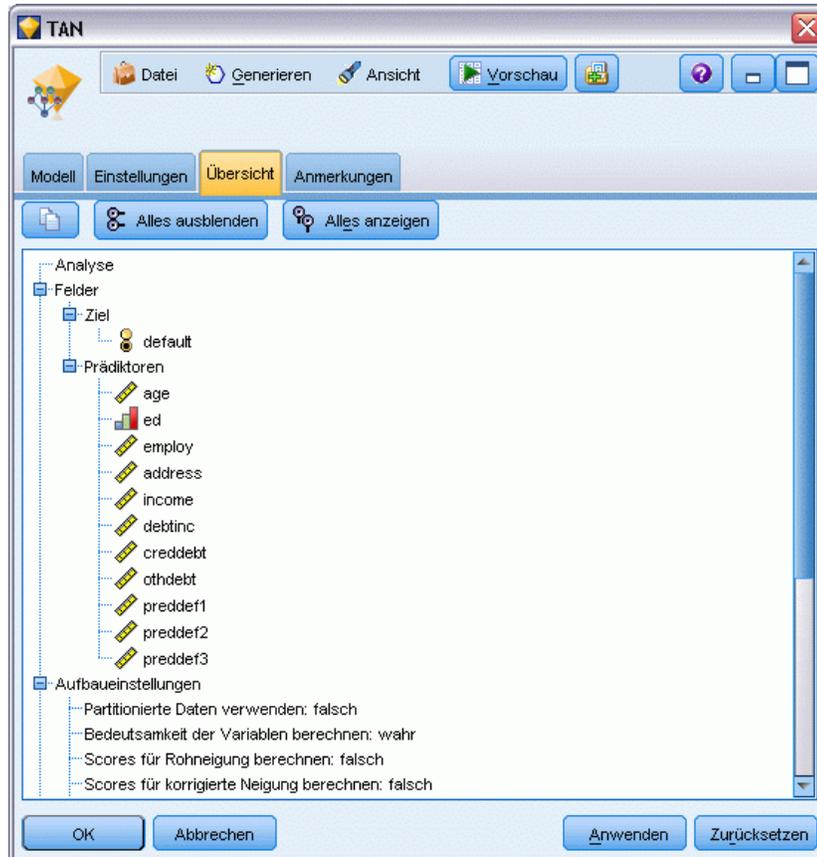
**Scores für korrigierte Neigung berechnen.** Rohneigungs-Scores beruhen ausschließlich auf den Trainingsdaten und sind aufgrund der Neigung vieler Modelle zur übermäßigen Anpassung an die Trainingsdaten möglicherweise zu optimistisch. Bei korrigierten Neigungen wird versucht, dies durch Evaluation der Modellleistung anhand einer Test- bzw. Validierungspartition zu kompensieren. Bei dieser Option muss im Stream ein Partitionsfeld definiert sein und die Scores für die korrigierte Neigung müssen im Modellierungsknoten aktiviert werden, bevor das Modell generiert wird.

**Alle Wahrscheinlichkeiten ausgeben.** Gibt an, ob die Wahrscheinlichkeiten für die einzelnen Kategorien des Ausgabefelds zu den einzelnen vom Knoten verarbeiteten Datensätzen hinzugefügt werden sollen. Wenn diese Option nicht ausgewählt wurde, wird nur die Wahrscheinlichkeit der vorhergesagten Kategorie hinzugefügt.

Die Standardeinstellung dieses Kontrollkästchens richtet sich nach dem entsprechenden Kontrollkästchen auf der Registerkarte "Experten" des Modellierungsknotens. [Für weitere Informationen siehe Thema Expertenoptionen für Bayes-Netzwerk-Knoten auf S. 203.](#)

## Bayes-Netzwerk-Modell – Übersicht

Abbildung 7-6  
Registerkarte “Übersicht” für ein Bayes-Netzwerk-Modell



Auf der Registerkarte “Übersicht” eines Modell-Nuggets finden Sie Informationen zum Modell selbst (*Analyse*), den im Modell verwendeten Feldern (*Felder*), den beim Erstellen des Modells verwendeten Einstellungen (*Aufbaueinstellungen*) und zum Modelltraining (*Trainingsübersicht*).

Beim ersten Durchsuchen des Knotens sind die Ergebnisse der Registerkarte “Übersicht” reduziert. Um die für Sie relevanten Ergebnisse anzuzeigen, vergrößern Sie sie mithilfe des Erweiterungssteuerelements links neben dem betreffenden Element oder klicken Sie auf die Schaltfläche *Alles anzeigen*, um alle Ergebnisse anzuzeigen. Um die Ergebnisse nach der Betrachtung wieder auszublenden, können Sie mit dem Erweiterungssteuerelement die gewünschten Ergebnisse reduzieren. Alternativ können Sie mit der Schaltfläche *Alles ausblenden* alle Ergebnisse ausblenden.

**Analyse.** Zeigt Informationen zum jeweiligen Modell an.

**Felder.** Listet die Felder auf, die bei der Erstellung des Modells als Ziel bzw. als Eingaben verwendet werden.

**Aufbaueinstellungen.** Enthält Informationen zu den beim Aufbau des Modells verwendeten Einstellungen.

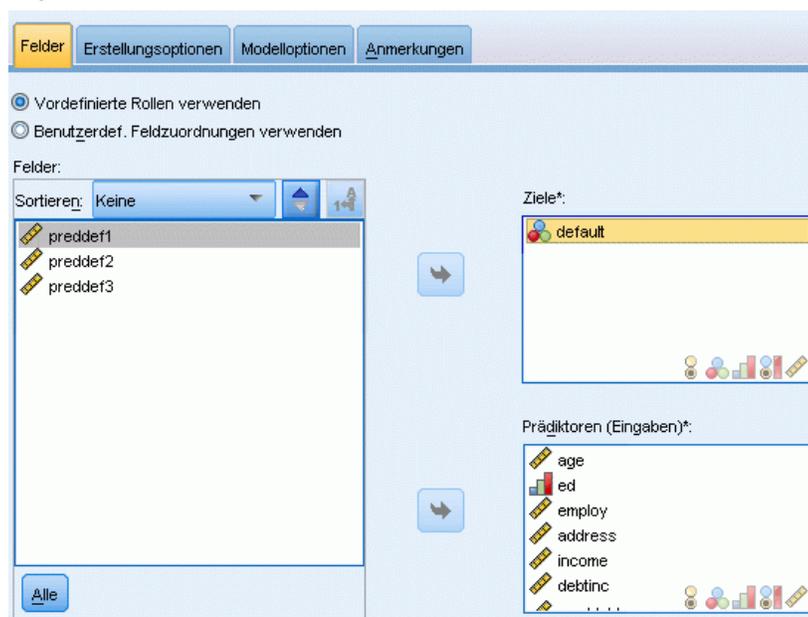
**Trainingsübersicht.** In diesem Abschnitt werden folgende Informationen angezeigt: der Typ des Modells, der für seine Erstellung verwendete Stream, der Benutzer, der ihn erstellt hat, der Erstellungszeitpunkt und die für den Aufbau des Modells benötigte Zeit.

# Neural Networks (Neuronale Netze)

Mit einem **neuronalen Netzwerk** können Näherungswerte für eine große Bandbreite an Vorhersagemodellen mit minimalen Anforderungen an die Modellstruktur und minimalen Annahmen für das Modell berechnet werden. Die Form der Beziehungen wird während des Lernprozesses bestimmt. Wenn sich eine lineare Beziehung zwischen Ziel und Prädiktoren eignet, sollten sich die Ergebnisse des neuronalen Netzwerks denen eines klassischen linearen Modells stark annähern. Ist eine nichtlineare Beziehung besser geeignet, erstellt das neuronale Netzwerk automatisch eine Näherung für die "korrekte" Modellstruktur.

Der Nachteil dieser Flexibilität besteht darin, dass das neuronale Netzwerk schwer zu interpretieren ist. Sollten Sie versuchen, einen Prozess zu erklären, der den Beziehungen zwischen Ziel und Prädiktoren zugrunde liegt, ist ein klassisches statistisches Modell besser geeignet. Wenn die Interpretierbarkeit des Modells jedoch keine große Rolle spielt, können Sie mit einem neuronalen Netzwerk gute Vorhersagen erzielen.

Abbildung 8-1  
Registerkarte "Felder"

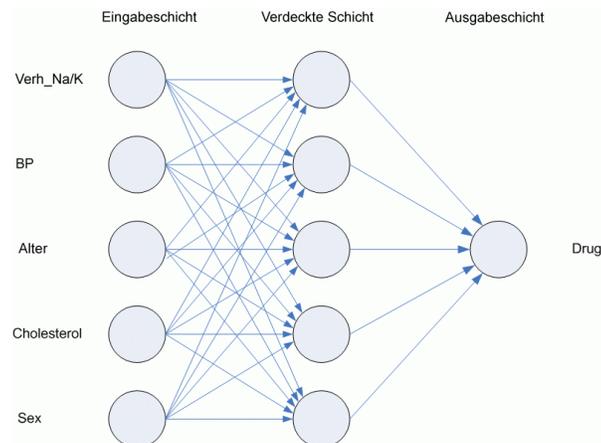


**Feldanforderungen.** Es muss mindestens ein Ziel und eine Eingabe vorhanden sein. Felder, die auf Beides oder Keine gesetzt sind, werden ignoriert. Es gibt keine Messniveaubeschränkungen bei Zielen oder Prädiktoren (Eingaben). [Für weitere Informationen siehe Thema Feldoptionen der Modellierungsknoten in Kapitel 3 auf S. 36.](#)

## Das neuronale Netzwerkmodell

Neuronale Netze sind einfache Modelle der Funktionsweise des Nervensystems. Die Grundeinheiten sind **Neuronen**, die in der Regel in **Schichten** organisiert sind, wie in der folgenden Abbildung dargestellt.

Abbildung 8-2  
Struktur eines neuronalen Netzes



Ein **neuronales Netz** ist ein vereinfachtes Modell der Art und Weise, wie ein menschliches Gehirn Informationen verarbeitet. Es funktioniert, indem eine große Anzahl miteinander verbundener Verarbeitungseinheiten simuliert wird, die abstrakten Versionen von Neuronen ähnlich sind.

Die Verarbeitungseinheiten sind in Schichten angeordnet. Diese bestehen in der Regel in einem neuronalen Netz aus drei Teilen: einer **Eingabeschicht** mit Einheiten, die die Eingabefelder darstellen; einer oder mehreren **verdeckten Schichten**; und einer **Ausgabeschicht**, mit einer Einheit bzw. Einheiten, die die Zielfelder darstellen. Die Einheiten sind mit verschiedenen Verbindungsstärken (**Gewichtungen**) verbunden. Die Eingangsdaten werden der ersten Schicht präsentiert und die Werte von jedem Neuron an jedes Neuron in der nächsten Schicht weitergeleitet. Schließlich gibt die Ausgabeschicht ein Ergebnis aus.

Das Netz lernt durch Prüfen einzelner Datensätze und generiert eine Vorhersage für jeden Datensatz. Außerdem nimmt es Änderungen der Gewichtungen vor, sobald eine falsche Vorhersage erfolgt. Dieser Vorgang wird viele Male wiederholt. Das Netz verbessert seine Vorhersagen so lange, bis mindestens eines der Grenzkriterien erfüllt ist.

Ursprünglich sind alle Gewichtungen zufällig und die Antworten, die vom Netz stammen, sind wahrscheinlich unsinnig. Das Netz lernt durch **Training**. Beispiele, für die die Ausgabe bekannt ist, werden dem Netz wiederholt präsentiert und die Antworten werden mit den bekannten Ausgaben verglichen. Die Informationen aus diesem Vergleich werden durch das Netz geleitet und die Gewichtungen schrittweise geändert. Je weiter das Training fortschreitet, desto genauer wird das Netz bei der Replizierung der bekannten Ergebnisse. Sobald das Netz trainiert ist, kann es auf zukünftige Fälle angewendet werden, bei denen das Ergebnis unbekannt ist.

## Verwenden von neuronalen Netzwerken mit Legacy-Streams

In Version 14 von IBM® SPSS® Modeler ist jetzt ein neuer Netzwerkknoten verfügbar, der Techniken zur Verbesserung/Verstärkung und Optimierung für sehr große Datensätze unterstützt. Vorhandene Streams, die den alten Knoten enthalten, werden in dieser Ausgabe weiterhin Modelle erstellen und scoren. Diese Unterstützung wird jedoch in zukünftigen Versionen wegfallen, so dass wir von nun an die Verwendung der neuen Version empfehlen.

Ab Version 13 werden Felder mit unbekanntem Wert (also Werte, die nicht in den Trainingsdaten vorhanden sind) nicht mehr automatisch als fehlende Werte behandelt und werden mit dem Wert \$null\$ gescort. Wenn Sie also Felder mit unbekanntem Wert mit einer älteren Version des neuronalen Netzwerkmodells (vor Version 13) in Version 13 oder höher als Nicht-Null-Wert scoren möchten, sollten Sie unbekannte Werte als fehlende Werte markieren (beispielsweise mithilfe des Typknotens).

## Ziele

Abbildung 8-3  
Ziele - Einstellungen

Wählen Sie ein Element aus:

Ziele  
Grundlagen  
Stoppregeln  
Ensembles  
Erweitert

Wie möchten Sie vorgehen?

Neues Modell aufbauen  
 Training des bestehenden Modells fortsetzen

Was ist Ihr Hauptziel?

Standardmodell erstellen  
 Modellgenauigkeit erweitern (Verstärkung)  
 Modellstabilität erweitern (Bagging)  
 Für sehr umfangreiche Datensätze optimieren (erfordert Server)

Beschreibung

Erstellt ein einzelnes Standardmodell zur Erklärung der Beziehungen zwischen Feldern. Standardmodelle lassen sich einfacher interpretieren und ermöglichen schnelleres Scoring als verstärkte, geboostete, verpackte oder umfangreiche Datensatz-Ensembles. Für mehrere Ziele wird immer ein Standardmodell verwendet.

### Was möchten Sie tun?

- **Neues Modell aufbauen.** Ein vollständig neues Modell aufbauen. Dies ist die übliche Wirkungsweise des Knotens.
- **Training eines bestehenden Modells fortsetzen.** Das Training wird mit dem letzten vom Knoten erfolgreich erstellten Modell fortgesetzt. Dies ermöglicht die Aktualisierung eines bestehenden Modells, ohne dass ein Zugriff auf die ursprünglichen Daten erforderlich ist,

und kann zu einer wesentlich schnelleren Leistung führen, da ausschließlich die neuen bzw. aktualisierten Datensätze in den Stream eingespeist werden. Details zum vorherigen Modell werden zusammen mit dem Modellierungsknoten gespeichert, wodurch diese Option auch dann verwendet werden kann, wenn das vorherige Modell-Nugget nicht mehr im Stream oder in der Modellpalette verfügbar ist.

*Anmerkung:* Wenn diese Option aktiviert ist, werden alle anderen Steuerungsfunktionen auf den Registerkarten “Felder” und “Erstellungsoptionen” deaktiviert.

### Wie lautet Ihr Hauptziel?

- **Erstellen eines Standardmodells.** Bei dem Verfahren wird ein einziges Modell erstellt, um das Ziel unter Verwendung der Prädiktoren vorherzusagen. In der Regel gilt, dass Standardmodelle einfacher interpretiert und schneller gescort werden können, als verbesserte, verstärkte oder große Daten-Set-Ensembles.
- **Modellgenauigkeit erhöhen (verbessern).** Bei dem Verfahren wird unter Einsatz der Verbesserung ein Ensemble-Modell erstellt. Dabei wird eine Modellsequenz erzeugt, um genauere Vorhersagen zu erhalten. Möglicherweise nimmt das Erstellen und Scoring bei Ensembles mehr Zeit in Anspruch als bei Standardmodellen.

Durch Verbesserung wird eine Reihe von “Komponentenmodellen” erstellt, von denen jede einzelne Komponente auf dem gesamten Daten-Set beruht. Vor dem Erstellen jedes aufeinander folgenden Komponentenmodells werden die Datensätze basierend auf den Residuen des vorangegangenen Komponentenmodells gewichtet. Größere Residuen erhalten eine höhere Analysegewichtung, sodass beim nächsten Komponentenmodell das Augenmerk auf einer hochwertigen Vorhersage dieser Datensätze liegt. Zusammen bilden diese Komponentenmodelle ein Ensemble-Modell. Das Ensemble-Modell scort basierend auf einer Kombinationsregel neue Datensätze. Die verfügbaren Regeln hängen vom Messniveau des Ziels ab.

- **Modellstabilität steigern (verstärken).** Bei dem Verfahren wird unter Einsatz der Verstärkung (Bootstrap-Aggregation) ein Ensemble-Modell erstellt. Dabei werden mehrere Modelle erzeugt, um zuverlässigere Vorhersagen zu erhalten. Möglicherweise nimmt das Erstellen und Scoring bei Ensembles mehr Zeit in Anspruch als bei Standardmodellen.

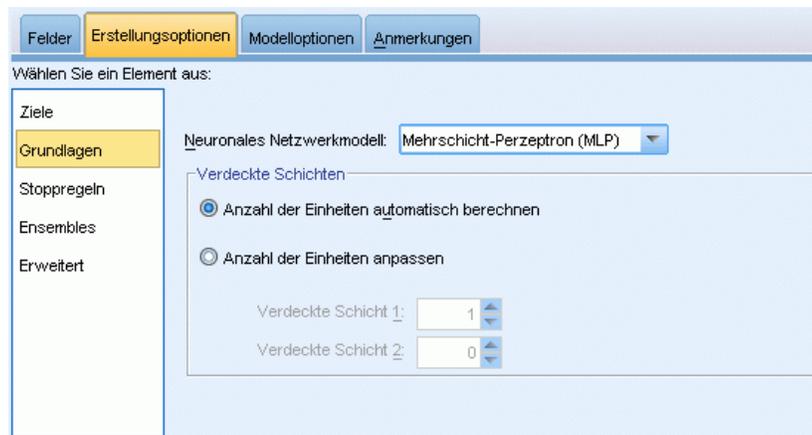
Bei der Bootstrap-Aggregation (Verstärkung) werden Reproduktionen des Trainings-Daten-Set erstellt, indem von der Ersetzung aus dem ursprünglichen Daten-Set Stichproben genommen werden. Dadurch werden Bootstrap-Stichproben mit der gleichen Größe wie beim ursprünglichen Daten-Set erstellt. Dann wird von jeder Reproduktion ein “Komponentenmodell” erstellt. Zusammen bilden diese Komponentenmodelle ein Ensemble-Modell. Das Ensemble-Modell scort basierend auf einer Kombinationsregel neue Datensätze. Die verfügbaren Regeln hängen vom Messniveau des Ziels ab.

- **Modell für sehr große Daten-Sets erstellen (IBM® SPSS® Modeler Server erforderlich).** Bei dieser Methode wird ein Ensemble-Modell durch Aufteilen des Daten-Sets in separate Datenblöcke erstellt. Diese Option ist empfehlenswert, wenn Ihr Daten-Set zu groß für die Erstellung eines der oben erwähnten Modelle oder die inkrementelle Modellerstellung ist. Unter Umständen kann das Modell mit dieser Option schneller als ein Standardmodell erstellt werden, das Scoring dauert jedoch evtl. länger als bei einem Standardmodell. Diese Option erfordert eine SPSS Modeler Serververbindung.

Wenn es mehrere Ziele gibt, wird bei dieser Methode nur ein Standardmodell erstellt, unabhängig vom ausgewählten Ziel.

## Grundeinstellungen

Abbildung 8-4  
Grundeinstellungen



**Neuronales Netzwerkmodell.** Der Modelltyp bestimmt, wie das Netzwerk die Prädiktoren durch die verdeckte(n) Schicht(en) mit den Zielen verbindet. Das **Mehrschicht-Perzeptron (MLP)** ermöglicht komplexere Beziehungen, was sich jedoch unter Umständen auf längere Trainings- und Scoringzeiten auswirkt. Die **radiale Basisfunktion (RBF)** kann die Trainings- und Scoringdauer verringern, was jedoch unter Umständen Einbußen bei der Vorhersagekraft im Vergleich zu MLP mit sich bringt.

**Verdeckte Schichten.** Die verdeckte(n) Schicht(en) eines neuronalen Netzwerks enthält nicht sichtbare Einheiten. Der Wert jeder verdeckten Einheit ist eine Funktion der Prädiktoren; die exakte Form der Funktion hängt teilweise vom Netzwerktyp ab. Ein Mehrschicht-Perzeptron kann eine oder zwei verdeckte Schichten besitzen; ein radiales Basisfunktionsnetzwerk kann nur eine verdeckte Schicht besitzen.

- **Anzahl der Einheiten automatisch berechnen.** Mit dieser Option wird ein Netzwerk mit einer verdeckten Schicht erstellt und die "beste" Anzahl an Einheiten in der verdeckten Schicht berechnet.
- **Anzahl der Einheiten anpassen.** Mit dieser Option können Sie die Anzahl der Einheiten in jeder verdeckten Schicht eingeben. Die erste verdeckte Schicht muss mindestens eine Einheit aufweisen. Durch Eingabe von 0 Einheiten in der zweiten verdeckten Schicht wird ein Mehrschicht-Perzeptron mit einer einzigen verdeckten Schicht erstellt.

*Anmerkung:* Sie sollten die Werte so wählen, dass die Anzahl an Knoten die Anzahl der stetigen Prädiktoren plus die Gesamtzahl der Kategorien in allen kategorialen Prädiktoren (Flag, nominal und ordinal) nicht übersteigt.

## Stoppregeln

Abbildung 8-5  
Stoppregeln- Einstellungen

Mit diesen Regeln wird festgelegt, wann das Training von Mehrschicht-Perzeptron-Netzwerken gestoppt werden soll; diese Einstellungen werden ignoriert, wenn der radiale Basisfunktionsalgorithmus verwendet wird. Das Training durchläuft mindestens einen Zyklus (Datenübergabe) und kann dann entsprechend den folgenden Kriterien gestoppt werden:

**Maximale Trainingszeit verwenden (pro Komponentenmodell).** Sie können wählen, ob sie eine maximale Minutenanzahl für die Ausführung des Algorithmus angeben wollen. Geben Sie eine Zahl größer 0 ein. Wenn ein Ensemble-Modell erstellt wird, ist das die zulässige Trainingszeit für jedes Komponentenmodell des Ensembles. Das Training kann ein wenig länger dauern, als die angegebene Zeit, da der aktuelle Zyklus abgeschlossen wird.

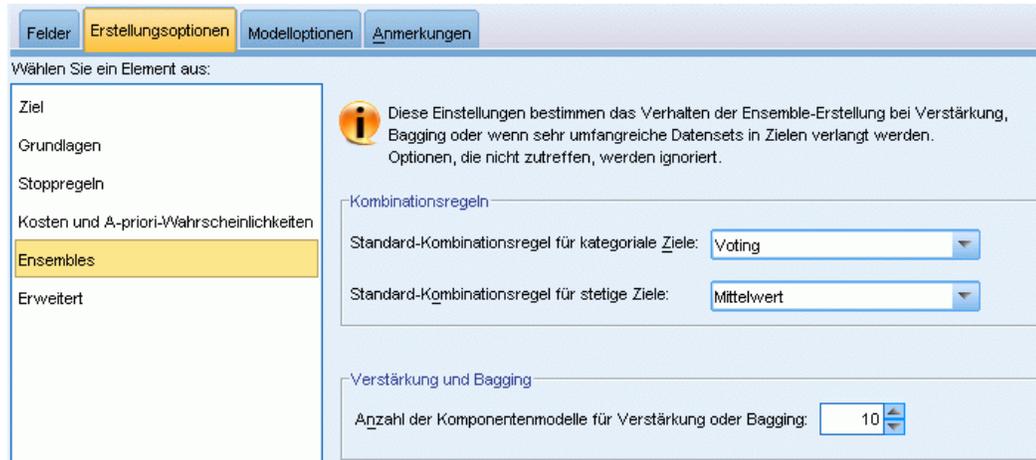
**Angeben der maximalen Anzahl an Trainingszyklen.** Die maximale Anzahl der zulässigen Trainingszyklen. Wenn die maximale Anzahl an Zyklen überschritten wird, stoppt das Training. Geben Sie eine ganze Zahl größer 0 ein.

**Mindestgenauigkeit verwenden.** Bei dieser Option wird das Training fortgesetzt, bis die angegebene Genauigkeit erreicht ist. Es kann vorkommen, dass dieser Zustand überhaupt nicht erreicht wird; sie können jedoch das Training jederzeit unterbrechen und das Netz mit der besten Genauigkeit speichern, die bisher erzielt wurde.

Der Trainingsalgorithmus stoppt auch dann, wenn der Fehler im Set zur Verhinderung übermäßiger Anpassung nach den einzelnen Zyklen nicht abnimmt, wenn die relative Änderung im Trainingsfehler gering ist, oder wenn die Quote des aktuellen Trainingsfehlers gering im Vergleich zum Anfangsfehler ist.

## Ensembles

Abbildung 8-6  
Ensemble-Einstellungen



Diese Einstellungen legen das Verhalten der Ensemblebildung fest, die erfolgt, wenn auf der Registerkarte “Ziele” die Option “Verbesserung”, “Verstärkung” oder “Sehr große Daten-Sets” ausgewählt ist. Optionen, die für das ausgewählte Ziel nicht gelten, werden ignoriert.

**Bagging und sehr umfangreiche Daten-Sets.** Beim Scoring eines Ensembles wird diese Regel angewendet, um die vorhergesagten Werte aus den Basismodellen für die Berechnung des Score-Werts für das Ensemble zu kombinieren.

- **Standard-Kombinierungsregel für kategoriale Ziele.** Ensemble-Vorhersagewerte für kategoriale Ziele können mithilfe von “Voting”, “höchster Wahrscheinlichkeit” oder “höchste mittlere Wahrscheinlichkeit” kombiniert werden. Bei **Voting** wird die Kategorie gewählt, die in allen Basismodellen am häufigsten die höchste Wahrscheinlichkeit erreicht. Bei **Höchste Wahrscheinlichkeit** wird die Kategorie gewählt, die in allen Basismodellen den höchsten Einzelwert bei der höchsten Wahrscheinlichkeit erzielt. Bei **Höchste mittlere Wahrscheinlichkeit** wird die Kategorie gewählt, die den höchsten Wert erreicht, wenn der Mittelwert der Wahrscheinlichkeiten für alle Kategorien in den Basismodellen berechnet wird.
- **Standard-Kombinierungsregel für stetige Ziele.** Ensemble-Vorhersagewerte für stetige Ziele können unter Verwendung des Mittelwerts oder Medians der Vorhersagewerte aus den Basismodellen kombiniert werden.

Hinweis: Wenn als Ziel die Verbesserung der Modellgenauigkeit ausgewählt wurde, wird die Auswahl zum Kombinieren der Regeln ignoriert. Bei der Verbesserung wird für das Scoring der kategorialen Ziele stets eine gewichtete Mehrheit verwendet und für das Scoring stetiger Ziele ein gewichteter Median.

**Verbesserung und Verstärkung.** Geben Sie die Anzahl der zu erstellenden Basismodelle an, wenn als Ziel die Verbesserung der Modellgenauigkeit oder -stabilität angegeben ist. Im Falle der Verstärkung ist das die Anzahl der Bootstrap-Stichproben. Muss eine positive ganze Zahl sein.

## Erweitert

Abbildung 8-7  
Erweiterte Einstellungen

Wählen Sie ein Element aus:

- Ziele
- Grundlagen
- Stoppregeln
- Ensembles
- Erweitert**

**i** Die Modellierung neuronaler Netzwerke zerlegt Datensätze in ein Modellerstellungsset und ein Set zur Prävention übermäßiger Anpassung. Geben Sie einen Prozentsatz der Datensätze für das Set zur Prävention übermäßiger Anpassung an.

Set zur Prävention übermäßiger Anpassung (%):

Ergebnisse replizieren

Startwert für Zufallsgenerator:

Fehlende Werte in Prädiktoren:

- Listenweise löschen
- Fehlende Werte imputieren

Mit den erweiterten Einstellungen lassen sich Optionen steuern, die nicht wirklich in andere Einstellungsgruppen passen.

**Verhinderung der übermäßigen Anpassung eingestellt.** Bei der neuronalen Netzwerkmethod werden Datensätze intern in ein Modellerstellungs-Set und ein Set zur Verhinderung der übermäßigen Anpassung aufgeteilt. Zweiteres ist ein unabhängiges Set an Datensätzen, das dazu verwendet wird, Fehler während des Trainings zu erfassen. So kann verhindert werden, dass die Methode zufällige Variationen in den Daten modelliert. Geben Sie einen Prozentsatz an Datensätzen an. Der Standardwert ist 30.

**Ergebnisse reproduzieren.** Durch Einstellen eines Startwerts für den Zufallsgenerator können Analysen reproduziert werden. Geben Sie eine ganze Zahl ein oder klicken Sie auf Generieren. Dadurch wird eine pseudozufällige Ganzzahl zwischen 1 und 2147483647 (einschließlich) erzeugt. Standardmäßig werden Analysen mit dem Startwert 229176228 reproduziert.

**Fehlende Werte in Prädiktoren.** Hier wird festgelegt, wie fehlende Werte behandelt werden sollen. **Listenweise löschen** entfernt Datensätze mit fehlenden Werten bei Prädiktoren aus der Modellerstellung. Mit **Fehlende Werte imputieren** werden fehlende Werte in Prädiktoren ersetzt und diese Datensätze in der Analyse verwendet. Stetige Felder imputieren den Mittelwert der minimalen und maximalen beobachteten Werte; kategoriale Felder imputieren die am häufigsten auftretende Kategorie. Hinweis: Datensätze mit fehlenden Werten in einem anderen auf der Registerkarte "Felder" angegebenen Feld werden stets aus der Modellerstellung ausgeschlossen.

## Modelloptionen

Abbildung 8-8  
Registerkarte "Modelloptionen"

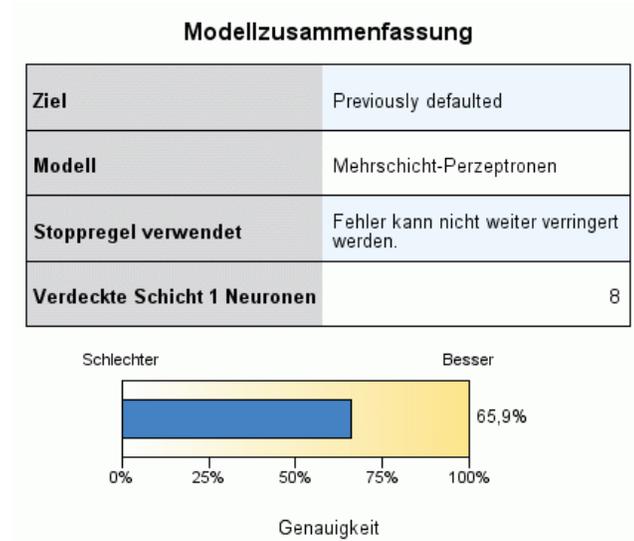
**Modellname.** Sie können den Modellnamen automatisch basierend auf den Zielfeldern generieren, oder einen benutzerdefinierten Namen eingeben. Der automatisch generierte Name ist der Zielfeldname. Bei mehreren Zielen besteht der Modellname aus den Feldnamen, die der Reihe nach durch Und-Zeichen verbunden werden. Wenn zum Beispiel *Feld1 Feld2 Feld3* Ziele sind, lautet der Modellname: *Feld1 & Feld2 & Feld3*.

**Für Scoring zur Verfügung stellen.** Beim Scoring des Modells sollten die ausgewählten Objekte in dieser Gruppe produziert werden. Wenn das Modell gescort wird, werden stets der vorhergesagte Wert (für alle Ziele) und die Konfidenz (für stetige Ziele) berechnet. Die berechnete Konfidenz kann auf der Wahrscheinlichkeit des vorhergesagten Werts (die höchste vorhergesagte Wahrscheinlichkeit) oder der Differenz zwischen der höchsten vorhergesagten Wahrscheinlichkeit und der zweithöchsten vorhergesagten Wahrscheinlichkeit basieren.

- **Vorhergesagte Wahrscheinlichkeit für kategoriale Ziele.** Mit dieser Option werden die vorhergesagten Wahrscheinlichkeiten für kategoriale Ziele produziert. Für jede Kategorie wird ein Feld erstellt.
- **Neigungs-Scores für Flag-Ziele.** Bei Modellen mit einem Flag-Ziel (das als Vorhersage "Ja" bzw. "Nein" ausgibt) können Sie Neigungs-Scores anfordern, die die Wahrscheinlichkeit des für das Zielfeld angegebenen wahren Ergebnisses angeben. Das Modell erzeugt Scores für Rohneigung. Bei aktiven Partitionen erzeugt das Modell außerdem anhand der Testpartition Scores für angepasste Neigung. [Für weitere Informationen siehe Thema Neigungsbewertungen in Kapitel 3 auf S. 44.](#)

## Modellübersicht

Abbildung 8-9  
Ansicht "Neuronales Netzwerkmodell – Übersicht"



Die Modellübersicht ist ein "Schnappschuss", eine Übersicht auf einen Blick über die Vorhersage- oder Klassifizierungsgenauigkeit des neuronalen Netzwerks.

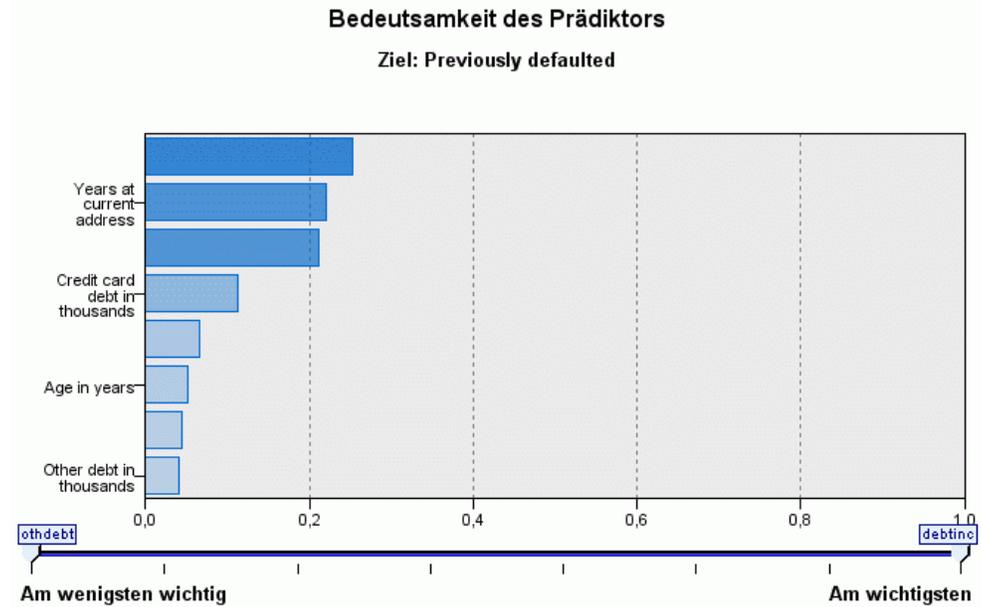
**Modellübersicht.** Die Tabelle identifiziert das Ziel, den Typ des trainierten neuronalen Netzwerks, die Stoppregel, die das Training beendet hat (wird angezeigt, wenn ein Mehrschicht-Perzeptron-Netzwerk trainiert wurde), und die Anzahl der Neuronen in jeder verdeckten Schicht des Netzwerks.

**Neuronales Netzwerk - Qualität.** Das Diagramm zeigt die Genauigkeit des endgültigen Modells an, das in einem größeren und besseren Format dargestellt wird. Bei einem kategorialen Ziel ist das einfach der Prozentsatz an Datensätzen, zu dem der vorhergesagte Wert dem beobachteten Wert entspricht. Bei stetigen Zielen ist das 1 minus das Verhältnis des mittleren absoluten Fehlers bei der Vorhersage (der Mittelwert der absoluten Werte aus vorhergesagten Werten minus beobachteten Werten) zum Bereich der vorhergesagten Werte (der höchste vorhergesagte Werte minus dem niedrigsten vorhergesagten Wert).

**Mehrere Ziele.** Wenn mehrere Ziele vorhanden sind, wird jedes Ziel in der Reihe Ziele der Tabelle angezeigt. Die im Diagramm angegebene Genauigkeit ist der Mittelwert der einzelnen Zielgenauigkeiten.

## Bedeutsamkeit des Prädiktors

Abbildung 8-10  
Ansicht "Bedeutsamkeit des Prädiktors"

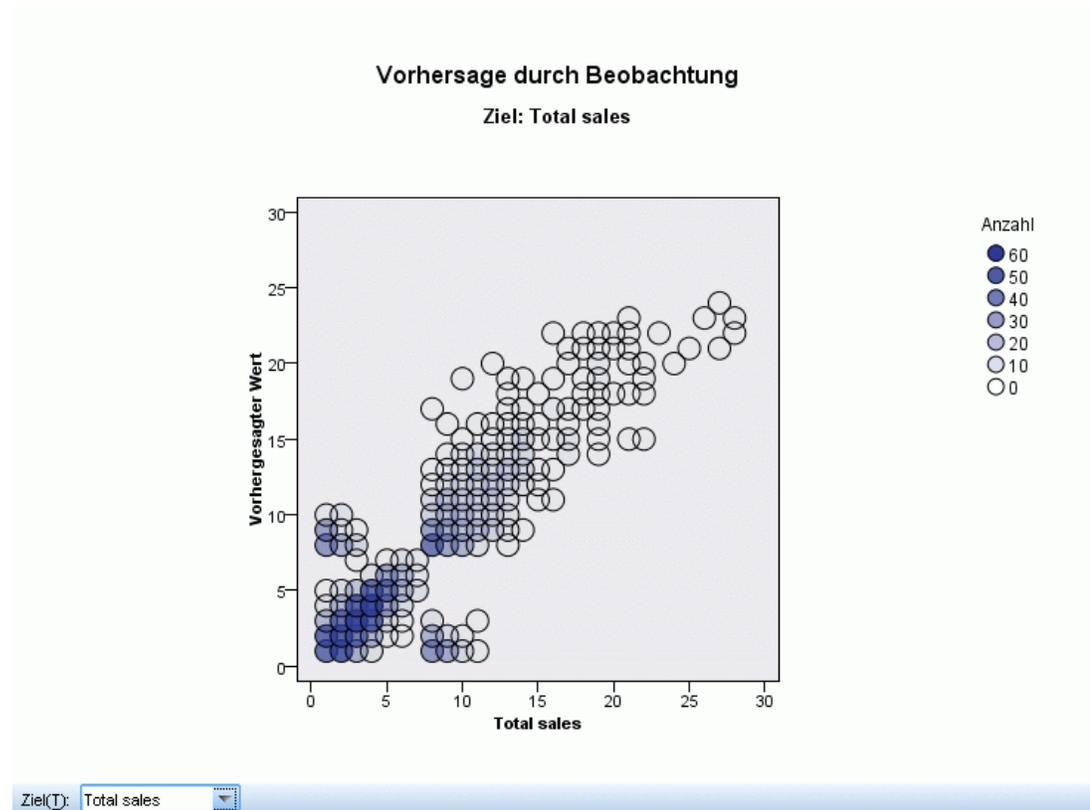


Normalerweise ist es sinnvoll, die Modellierungsbemühungen auf die wichtigsten Prädiktorfelder zu konzentrieren und diejenigen Prädiktoren zu verwerfen bzw. zu ignorieren, die die geringste Bedeutung haben. Dabei unterstützt Sie das Diagramm für die Bedeutsamkeit der Prädiktoren, da es die relative Bedeutsamkeit der einzelnen Prädiktoren für das Modell angibt. Da die Werte relativ sind, beträgt die Summe der Werte aller Prädiktoren im Diagramm 1,0. Die Bedeutsamkeit der Prädiktoren bezieht sich nicht auf die Genauigkeit des Modells. Sie bezieht sich lediglich auf die Bedeutsamkeit der einzelnen Prädiktoren für eine Vorhersage und nicht auf die Genauigkeit der Vorhersage.

**Mehrere Ziele.** Bei mehreren Zielen wird jedes Ziel in einem separaten Diagramm dargestellt. Dazu gibt es eine Dropdown-Liste Ziel, aus der das anzuzeigende Ziel ausgewählt werden kann.

## Vorhersage nach Beobachtung

Abbildung 8-11  
Ansicht "Vorhersage nach Beobachtung"

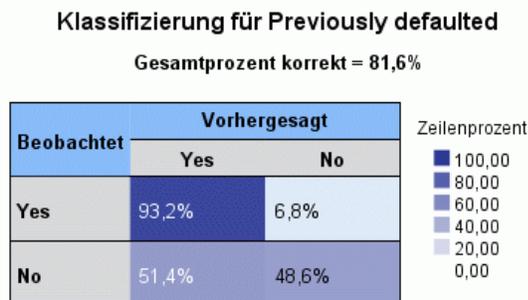


Für stetige Ziele zeigt diese Ansicht ein Bin-Streudiagramm der vorhergesagten Werte auf der vertikalen Achse durch die beobachteten Werte auf der horizontalen Achse.

**Mehrere Ziele.** Bei mehreren stetigen Zielen wird jedes Ziel in einem separaten Diagramm dargestellt. Dazu gibt es eine Dropdown-Liste Ziel, aus der das anzuzeigende Ziel ausgewählt werden kann.

## Klassifizierung

Abbildung 8-12  
Ansicht "Klassifizierung," Stil "Reihenprozent"



Bei kategorialen Zielen wird hier die Kreuzklassifikation der beobachteten versus vorhergesagten Werte in einer Hitze Karte angezeigt, plus die "Gesamtprozent korrekt"-Werte.

**Tabellenstile.** Für die Darstellung sind mehrere verschiedene Anzeigestile verfügbar, auf die über die Dropdown-Liste Stil zugegriffen werden kann.

- **Reihenprozent.** Diese Darstellung zeigt die Reihenprozentensätze (die Zellenhäufigkeit als Prozentsatz der Reihengesamtsumme) in den Zellen an. Dies ist die Standardeinstellung.
- **Zellenhäufigkeit.** Diese Option zeigt die Zellenhäufigkeit in den Zellen an. Die Schattierung der Hitze Karte basiert weiterhin auf den Reihenprozentensätzen.
- **Hitze Karte.** Diese Option zeigt keine Werte in der Zellen an, nur die Schattierung.
- **Komprimiert.** Hier werden keine Reihen- oder Spaltenüberschriften oder Werte in den Zellen angezeigt. Diese Ansicht kann nützlich sein, wenn das Ziel sehr viele Kategorien besitzt.

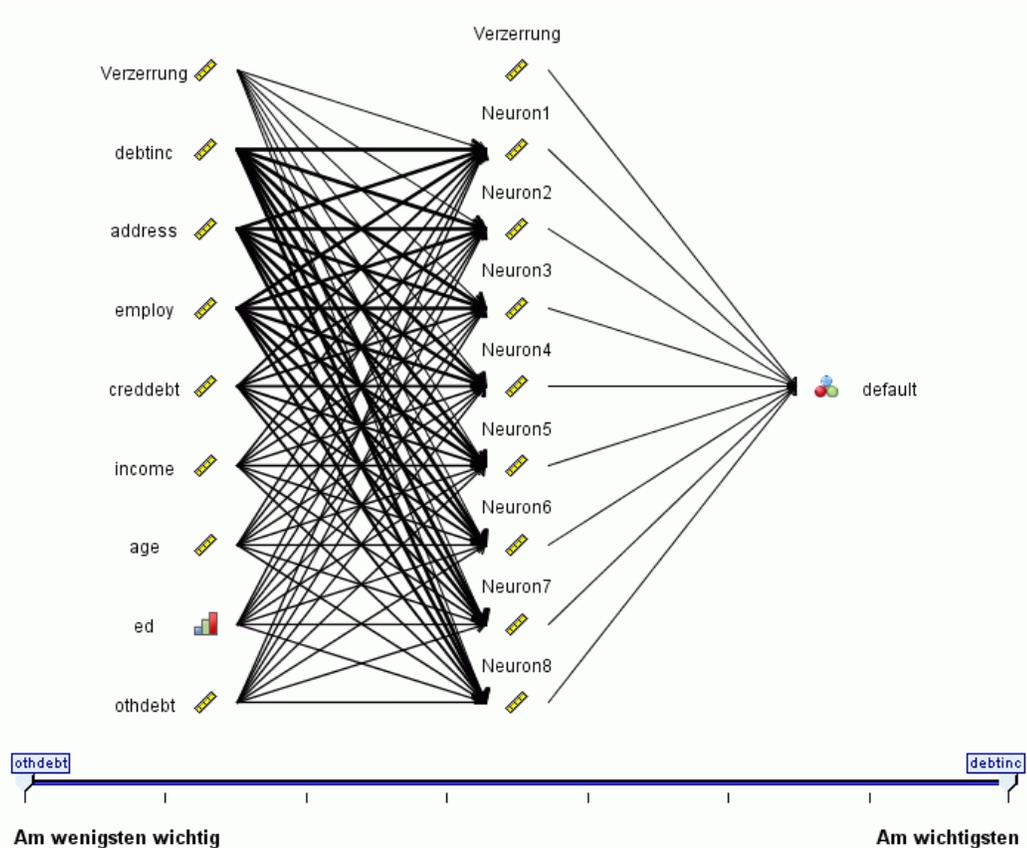
**Fehlend.** Wenn Datensätze fehlende Werte im Ziel aufweisen, werden sie in einer Reihe (Fehlende) unter allen gültigen Reihen angezeigt. Datensätze mit fehlenden Werten tragen nicht zu "Gesamtprozent korrekt" bei.

**Mehrere Ziele.** Bei mehreren kategorialen Zielen wird jedes Ziel in einer separaten Tabelle dargestellt. Dazu gibt es eine Dropdown-Liste Ziel, aus der das anzuzeigende Ziel ausgewählt werden kann.

**Große Tabellen.** Wenn das angezeigte Ziel mehr als 100 Kategorien besitzt, wird keine Tabelle angezeigt.

## Netzwerk

Abbildung 8-13  
Ansicht "Netzwerk," Eingaben links, Stil "Effekte"



Hier wird eine grafische Darstellung des neuronalen Netzwerks angezeigt.

**Diagrammstile.** Es sind zwei verschiedene Anzeigestile verfügbar, auf die über die Dropdown-Liste Stil zugegriffen werden kann.

- **Effekte.** Hier werden jeweils Prädiktor und Ziel als ein Knoten im Diagramm angezeigt, unabhängig davon, ob die Messskala stetig oder kategorial ist. Dies ist die Standardeinstellung.
- **Koeffizienten.** Hier werden mehrere Indikator-knoten für kategoriale Prädiktoren und Ziele angezeigt. Die Verbindungslinien im Diagramm im Koeffizienten-Stil sind farbig dargestellt, basierend auf dem geschätzten Wert der synaptischen Gewichtung.

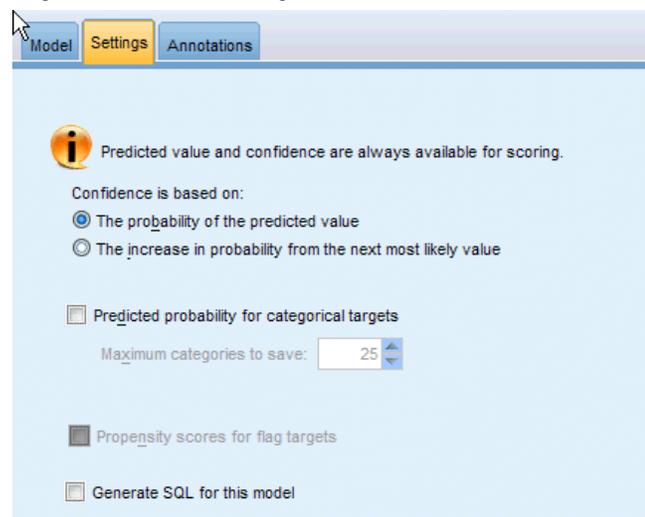
**Diagrammausrichtung.** Standardmäßig ist das Netzwerkdiagramm so angeordnet, dass die Eingaben links und die Ziele rechts dargestellt sind. Mithilfe der Steuerelemente in der Symbolleiste können Sie die Ausrichtung so ändern, dass die Eingaben oben und die Ziele unten oder die Eingaben unten und die Ziele oben dargestellt werden.

**Bedeutsamkeit des Prädiktors.** Die Verbindungslinien im Diagramm sind basierend auf der Bedeutsamkeit des Prädiktors gewichtet, wobei eine größere Linienbreite einer größeren Bedeutsamkeit entspricht. Für die Bedeutsamkeit des Prädiktors gibt es einen Schieberegler in der Symbolleiste, mit dem eingestellt wird, welche Prädiktoren im Netzwerkdiagramm gezeigt werden. Dadurch wird das Modell nicht verändert, doch Sie können sich ganz problemlos auf die wichtigsten Prädiktoren konzentrieren.

**Mehrere Ziele.** Wenn mehrere Ziele vorhanden sind, werden alle Ziele im Diagramm angezeigt.

## Einstellungen

Abbildung 8-14  
Registerkarte "Einstellungen"



Beim Scoring des Modells sollten die ausgewählten Objekte in dieser Registerkarte produziert werden. Wenn das Modell gescort wird, werden stets der vorhergesagte Wert (für alle Ziele) und die Konfidenz (für stetige Ziele) berechnet. Die berechnete Konfidenz kann auf der Wahrscheinlichkeit des vorhergesagten Werts (die höchste vorhergesagte Wahrscheinlichkeit) oder der Differenz zwischen der höchsten vorhergesagten Wahrscheinlichkeit und der zweithöchsten vorhergesagten Wahrscheinlichkeit basieren.

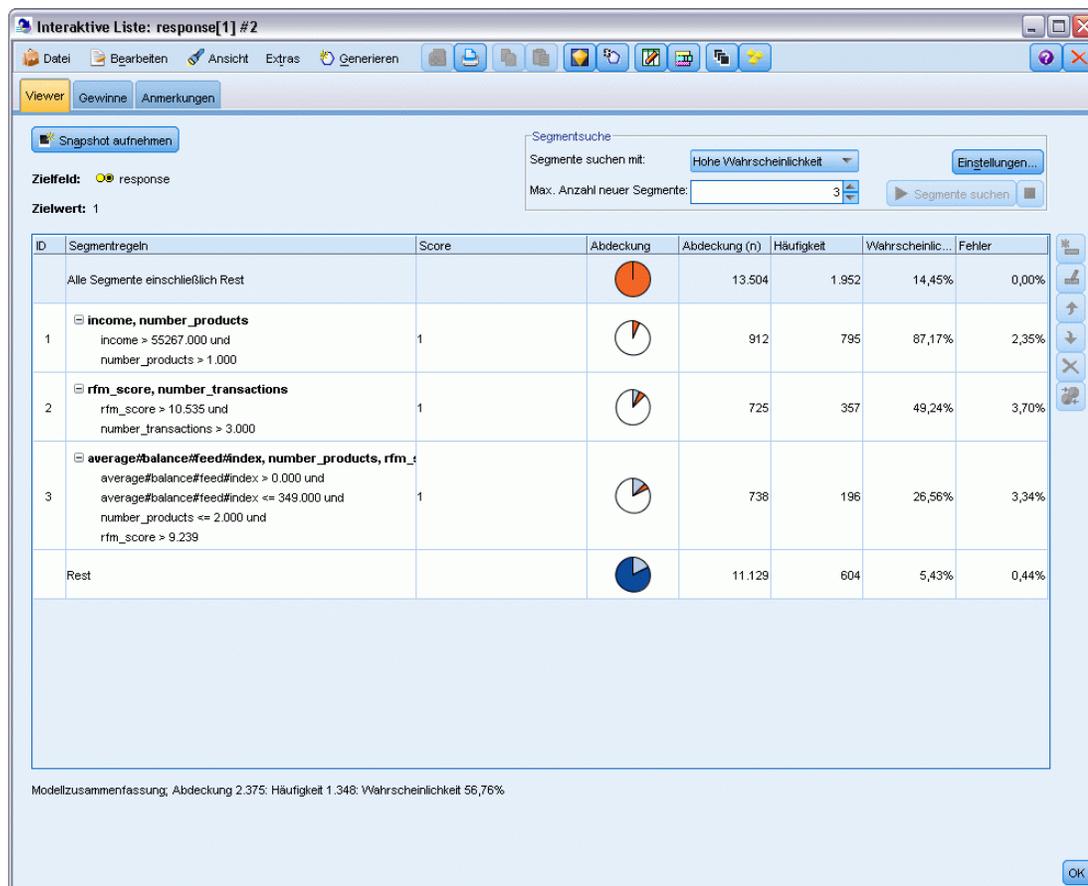
- **Vorhergesagte Wahrscheinlichkeit für kategoriale Ziele.** Mit dieser Option werden die vorhergesagten Wahrscheinlichkeiten für kategoriale Ziele produziert. Für jede Kategorie wird ein Feld erstellt.
- **Neigungs-Scores für Flag-Ziele.** Bei Modellen mit einem Flag-Ziel (das als Vorhersage “Ja” bzw. “Nein” ausgibt) können Sie Neigungs-Scores anfordern, die die Wahrscheinlichkeit des für das Zielfeld angegebenen wahren Ergebnisses angeben. Das Modell erzeugt Scores für Rohneigung. Bei aktiven Partitionen erzeugt das Modell außerdem anhand der Testpartition Scores für angepasste Neigung. [Für weitere Informationen siehe Thema Neigungsbewertungen in Kapitel 3 auf S. 44.](#)

**SQL für dieses Modell generieren.** Bei Verwendung von Daten aus einer Datenbank kann SQL-Code per Pushback zur Ausführung an die Datenbank zurückübertragen werden. Dadurch kann bei vielen Operationen eine bessere Leistung erzielt werden. [Für weitere Informationen siehe Thema SQL-Optimierung in Kapitel 6 in IBM SPSS Modeler Server 14.2-Verwaltungs- und Leistungshandbuch.](#)

## ***Entscheidungsliste***

Decision List-Modelle kennzeichnen Untergruppen bzw. **Segmente**, die eine höhere oder geringere Wahrscheinlichkeit für ein binäres Ergebnis (“Ja” bzw. “Nein”) aufweisen als die Gesamtstichprobe. Sie könnten beispielsweise nach Kunden suchen, deren Abwanderung besonders unwahrscheinlich ist oder die mit größter Wahrscheinlichkeit auf ein Angebot oder eine Kampagne ansprechen. Der Decision List Viewer bietet Ihnen vollständige Kontrolle über das Modell: Sie können Segmente bearbeiten, Ihre eigenen Geschäftsregeln hinzufügen, angeben, wie die einzelnen Segmente gescort werden, und das Modell auf zahlreiche andere Weisen anpassen, um den Anteil an Treffern über alle Segmente hinweg zu optimieren. Er eignet sich besonders gut für die Erstellung von Mailinglisten bzw. die anderweitige Ermittlung der Datensätze, die für eine bestimmte Kampagne gezielt betrachtet werden sollen. Außerdem können Sie mehrere **Mining-Aufgaben** verwenden, um verschiedene Modellierungsansätze zu kombinieren, beispielsweise durch Ermittlung der Segmente mit hohen oder niedrigen Trefferquoten innerhalb desselben Modells oder durch Aufnahme bzw. Ausschluss der einzelnen Segmente in der Scoring-Phase je nach Bedarf.

Abbildung 9-1  
Entscheidungslistenmodell



### Segmente, Regeln und Bedingungen

Ein Modell besteht aus einer Liste von Segmenten, von denen jedes durch eine Regel definiert ist, die übereinstimmende Datensätze auswählt. Eine Regel kann jeweils mehrere Bedingungen aufweisen; Beispiel:

RFM\_SCORE > 10 and  
MONTHS\_CURRENT <= 9

Die Regeln werden in der angegebenen Reihenfolge angewendet. Dabei legt die erste zutreffende Regel das Ergebnis für einen bestimmten Datensatz fest. Für sich genommen können sich Regeln bzw. Bedingungen überschneiden, durch die Reihenfolge der Regeln wird die Mehrdeutigkeit jedoch aufgelöst. Wenn keine Regel zutrifft, wird der Datensatz der Restregel zugewiesen.

### Vollständige Kontrolle über das Scoring

Im Decision List Viewer können Sie Segmente anzeigen, bearbeiten und neu anordnen und auswählen, welche Elemente für das Scoring ein- bzw. ausgeschlossen werden sollen. Sie können beispielsweise eine Gruppe von Kunden von zukünftigen Angeboten aus- und andere einschließen

und sofort ablesen, wie dies die Gesamttrefferquote beeinflusst. Decision List-Modelle geben den Score *Ja* für eingeschlossene Segmente aus und *\$null* für alles andere, einschließlich des Rests. Durch diese unmittelbare Steuerung des Scorings sind Decision List-Modelle ideal für das Erstellen von Mailinglisten und sie sind – unter anderem in Callcenter- und Marketinganwendungen – weit verbreitet im Customer Relationship Management.

Abbildung 9-2  
Entscheidungslistenmodell

### **Mining-Aufgaben, Maße und Auswahlmöglichkeiten**

Der Modellierungsvorgang wird durch **Mining-Aufgaben** gesteuert. Jede Mining-Aufgabe initiiert effektiv einen neuen Modellierungsdurchgang und gibt ein neues Set mit alternativen Modellen aus, aus denen Sie wählen können. Die Standardaufgabe beruht auf Ihren ursprünglichen Angaben im Decision List-Knoten; Sie können jedoch jede beliebige Anzahl an benutzerdefinierten Aufgaben definieren. Außerdem können Sie Aufgaben iterativ anwenden. Beispielsweise können Sie eine Suche des Typs “Hohe Wahrscheinlichkeit” für das gesamte Trainings-Set durchführen und anschließend eine Suche des Typs “Geringe Wahrscheinlichkeit” für den Rest, um Segmente mit niedriger Leistung auszusondern.

Abbildung 9-3  
Erstellen einer Mining-Aufgabe

**Mining-Aufgabe erstellen/bearbeiten: response[1]**

Einstellungen laden:

Ziel

Zielfeld:  response Zielwert: 1

Einfache Einstellungen

Segmente suchen mit:

Maximale Anzahl an neuen Segmenten:

Minimale Segmentgröße

Als Prozentsatz des vorherigen Segments (%):

Als absoluter Wert (N):

Maximale Anzahl an Alternativen:

Maximale Anzahl an Attributen pro Segment:

Wiederverwendung von Attributen innerhalb des Segments zulassen

Konfidenzintervall für neue Bedingungen (%):

Experteneinstellungen

Klassiermethode:	Gleiche Anzahl	Anzahl der Klassen:	10
Modellsuchbreite:	5	Regelsuchbreite:	5
Klassenzusammenführungsfaktor:	2.00		
Fehlende Werte in Bedingungen zulassen:	Wahr	Zwischenergebnisse verwerfen:	Wahr

Daten

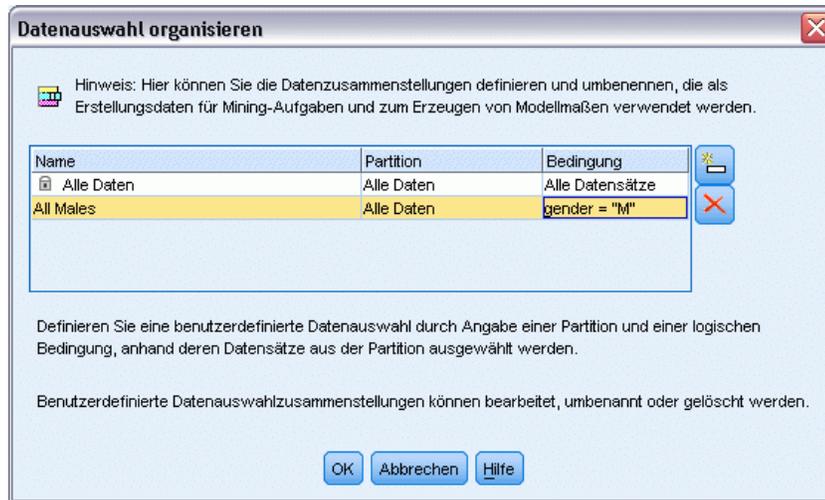
Erstellungsauswahl:

Verfügbare Felder:  Alle Felder  Benutzerdefiniert

### ***Datenauswahl***

Sie können Datenauswahlmöglichkeiten und benutzerdefinierte Modellmaße für Modellerstellung und -evaluation definieren. Sie können beispielsweise eine Datenauswahl in einer Mining-Aufgabe angeben, um das Modell auf eine bestimmte Region zuzuschneiden, und ein benutzerdefiniertes Maß erstellen, um zu evaluieren, wie leistungsfähig das betreffende Modell für das gesamte Land ist. Im Gegensatz zu Mining-Aufgaben ändern Maße nicht das zugrunde liegende Modell, sondern bieten vielmehr einen anderen Fokus zur Einschätzung, wie leistungsfähig es ist.

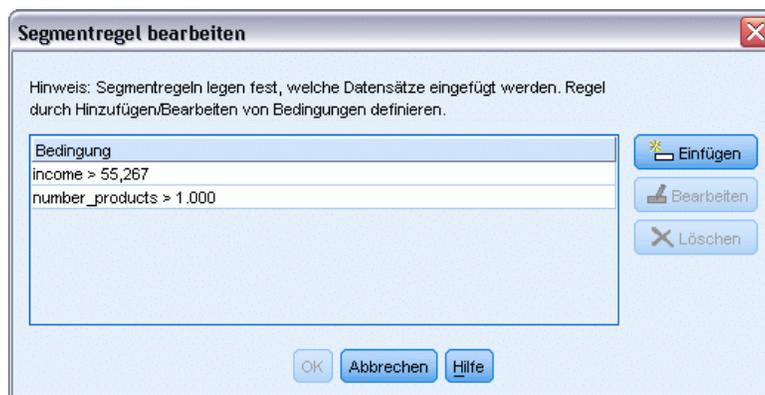
Abbildung 9-4  
Erstellen einer Datenauswahl



### Einbringen Ihres Geschäftswissens

Durch die Feinabstimmung oder Erweiterung der vom Algorithmus ermittelten Segmente können Sie mithilfe des Decision List Viewer Ihr Geschäftswissen unmittelbar in das Modell integrieren. Sie können die vom Modell generierten Segmente bearbeiten oder auf der Grundlage der von Ihnen angegebenen Regeln weitere Segmente hinzufügen. Anschließend können Sie die Änderungen übernehmen und eine Vorschau der Ergebnisse anzeigen.

Abbildung 9-5  
Angabe einer Regel



Für weitere Einblicke ermöglicht eine dynamische Verknüpfung mit Excel den Export Ihrer Daten in Excel, wo damit Präsentationsdiagramme erstellt und benutzerdefinierte Maße berechnet werden können wie beispielsweise komplexe Profit- und ROI-Werte, die während der Modellerstellung im Decision List Viewer angezeigt werden können.

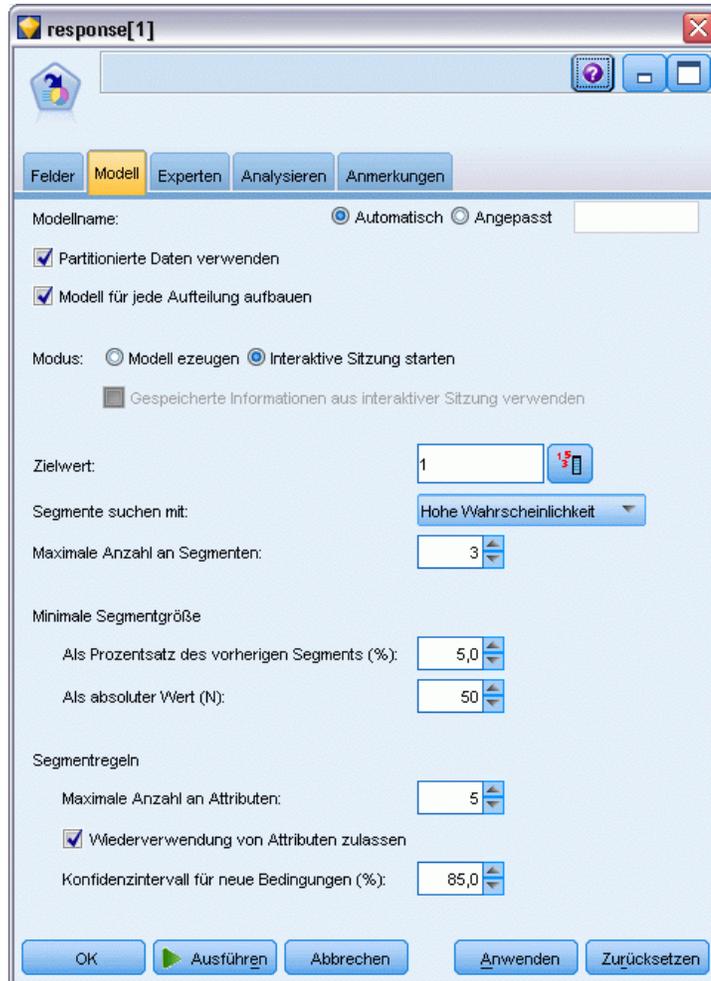
**Beispiel.** Die Marketingabteilung eines Finanzinstituts möchte in zukünftigen Kampagnen profitablere Ergebnisse erzielen, indem jedem Kunden ein speziell für ihn geeignetes Angebot unterbreitet wird. Mit dem Entscheidungslistenmodell können Sie auf der Grundlage früherer

Werbeaktionen die Eigenschaften der Kunden ermitteln, die mit der größten Wahrscheinlichkeit positiv reagieren werden, und auf der Grundlage der Ergebnisse eine Mailingliste generieren. [Für weitere Informationen siehe Thema Modellieren der Kundenreaktion \(Entscheidungsliste\) in Kapitel 12 in IBM SPSS Modeler 14.2- Anwendungshandbuch.](#)

**Anforderungen.** Ein einzelnes kategoriales Zielfeld mit einem Messniveau des Typs *Flag* oder *Nominal*, das das binäre Ergebnis angibt, das Sie vorhersagen möchten (Ja/Nein), sowie mindestens ein Eingabefeld. Wenn das Zielfeld den Typ *Nominal* aufweist, müssen Sie manuell einen einzelnen Wert auswählen, der als **Treffer** oder **Antwort** behandelt werden soll; alle anderen Werte werden als **kein Treffer** zusammengefasst. Außerdem kann ein optionales Häufigkeitsfeld angegeben werden. Kontinuierliche Datums-/Uhrzeitfelder werden ignoriert. Eingaben mit kontinuierlichem numerischen Bereich werden automatisch vom Algorithmus klassiert, wie auf der Registerkarte “Experten” des Modellierungsknotens angegeben. Eine detailliertere Kontrolle über das Klassieren erhalten Sie, wenn Sie weiter oben im Stream einen Klassierknoten einfügen und das klassierte Feld als Eingabe mit dem Messniveau *Ordinal* verwenden.

## Entscheidungslistenmodell – Optionen

Abbildung 9-6  
Entscheidungslistenknoten: Registerkarte "Modell"



**Modellname.** Sie können den Modellnamen automatisch basierend auf den Ziel- oder ID-Feldnamen (oder dem Modelltyp in Fällen, in denen kein solches Feld angegeben ist) generieren oder einen benutzerdefinierten Namen eingeben.

**Partitionierte Daten verwenden.** Wenn ein Partitionsfeld definiert ist, gewährleistet diese Option, dass nur Daten aus der Trainingspartition für die Modellerstellung verwendet werden. [Für weitere Informationen siehe Thema Partitionsknoten in Kapitel 4 in IBM SPSS Modeler 14.2- Quellen-, Prozess- und Ausgabeknoten.](#)

**Aufteilungsmodelle erstellen.** Erstellt ein separates Modell für jeden möglichen Wert von Eingabefeldern, die als Aufteilungsfelder angegeben werden. [Für weitere Informationen siehe Thema Erstellung von aufgeteilten Modellen in Kapitel 3 auf S. 31.](#)

**Modalwert.** Legt fest, welche Methode für die Modellbildung verwendet wird.

- **Modell erzeugen.** Generiert beim Ausführen des Knotens automatisch ein Modell in der Modellpalette. Das so entstandene Modell kann zum Zwecke des Scorens zu Streams hinzugefügt werden; eine weitere Bearbeitung des Modells ist jedoch nicht möglich.
- **Interaktive Sitzung starten.** Öffnet das Fenster für die interaktive Modellierung (Ausgabe) des Decision List Viewer, in dem Sie eine Auswahl aus mehreren Alternativen treffen und den Algorithmus wiederholt mit verschiedenen Einstellungen anwenden können, um das Modell progressiv zu erweitern bzw. zu ändern. [Für weitere Informationen siehe Thema Decision List Viewer auf S. 238.](#)
- **Gespeicherte Informationen aus interaktiver Sitzung verwenden.** Startet eine interaktive Sitzung unter Verwendung von zuvor gespeicherten Einstellungen. Interaktive Sitzungen können im Decision List Viewer mithilfe des Menüs “Generieren” (zur Erstellung eines Modells oder Modellierungsknotens) oder des Menüs “Datei” (zur Aktualisierung des Knotens, von dem aus die Sitzung gestartet wurde) gespeichert werden.

**Zielwert.** Gibt den Wert des Zielfelds an, in dem das zu modellierende Ergebnis angegeben wird. Wenn beispielsweise das Zielfeld “churn” (Abwanderung) mit 0 = no und 1 = yes kodiert ist, geben Sie 1 an, um Regeln festzulegen, mit denen die Datensätze für wahrscheinlich abwandernde Kunden identifiziert werden.

**Segmente suchen mit.** Gibt an, ob die Suche nach der Zielvariablen auf eine Hohe Wahrscheinlichkeit oder eine Geringe Wahrscheinlichkeit des Auftretens achten soll. Das Suchen und Ausschließen dieser Segmente kann eine geeignete Methode zum Verbessern des Modells darstellen und ist oft besonders nützlich, wenn die verbleibende Teilmenge eine geringe Wahrscheinlichkeit aufweist.

**Maximale Anzahl an Segmenten.** Legt die maximale Anzahl von zurückzugebenden Segmenten fest. Die obersten  $N$  Segmente werden erstellt, wobei das beste Segment die höchste Wahrscheinlichkeit bzw., bei mehreren Modellen mit der gleichen Wahrscheinlichkeit, die höchste Abdeckung aufweist. Die zulässige Minimaleinstellung liegt bei 1, eine Maximaleinstellung ist nicht vorhanden.

**Mindestsegmentgröße.** Die beiden Einstellungen unten geben die minimale Segmentgröße vor. Der größere der beiden Werte wird ausgewählt. Wenn z. B. der Prozentwert einer höheren Zahl entspricht als der absolute Wert, wird die prozentuale Einstellung ausgewählt.

- **Als Prozentsatz des vorherigen Segments (%).** Legt die minimale Gruppengröße als Prozentsatz der Datensätze fest. Die zulässige Minimaleinstellung liegt bei 0, die zulässige Maximaleinstellung bei 99,9.
- **Als absoluter Wert (N).** Legt die minimale Gruppengröße als absolute Anzahl der Datensätze fest. Die zulässige Minimaleinstellung liegt bei 1, eine Maximaleinstellung ist nicht vorhanden.

#### Segmentregeln.

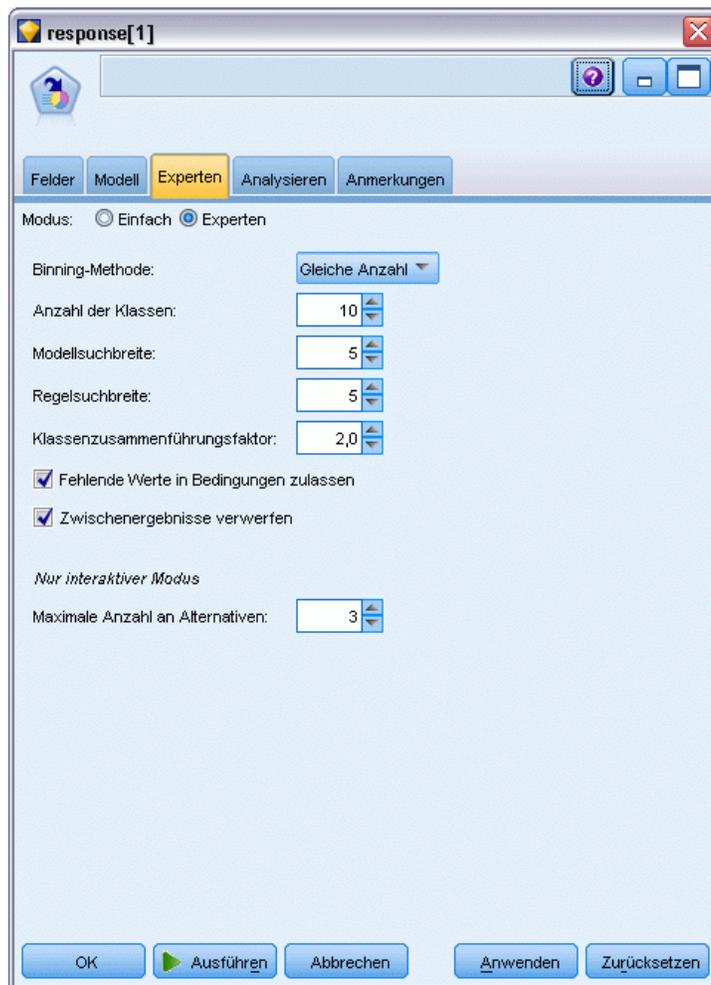
**Maximale Anzahl an Attributen.** Legt die maximale Anzahl von Bedingungen pro Segmentregel fest. Die zulässige Minimaleinstellung liegt bei 1, eine Maximaleinstellung ist nicht vorhanden.

- **Wiederverwendung von Attributen zulassen.** Bei Aktivierung dieser Option können in jedem Zyklus alle Attribute berücksichtigt werden, auch solche, die in vorherigen Zyklen verwendet wurden. Die Bedingungen für ein Segment sind in Zyklen aufgebaut, wobei in jedem Zyklus eine neue Bedingung hinzugefügt wird. Die Anzahl der Zyklen wird über die Einstellung Maximale Anzahl an Attributen festgelegt.

**Konfidenzintervall für neue Bedingungen (%).** Legt das Konfidenzintervall für das Testen der Signifikanz des Segments fest. Diese Einstellung spielt eine wichtige Rolle für die Anzahl der zurückgegebenen Segmente (sofern zutreffend) sowie für die Anzahl von Bedingungen pro Segmentregel. Je höher der Wert, desto kleiner das zurückgegebene Ergebnis-Set. Die zulässige Minimaleinstellung liegt bei 50, die zulässige Maximaleinstellung bei 99,9.

## Entscheidungslistenknoten – Expertenoptionen

Abbildung 9-7  
Entscheidungslistenknoten: Registerkarte "Experten"



Mit Expertenoptionen können Sie die Feinabstimmung des Modellbildungsvorgangs vornehmen.

**Klassiermethode.** Die Methode zum Klassieren kontinuierlicher Felder (gleiche Anzahl oder gleiche Breite).

**Anzahl der Klassen.** Die Anzahl der für kontinuierliche Felder zu erstellenden Klassen. Die zulässige Minimaleinstellung liegt bei 2, eine Maximaleinstellung ist nicht vorhanden.

**Modellsuchbreite.** Die maximale Anzahl von Modellergebnissen pro Zyklus, die für den nächsten Zyklus verwendet werden können. Die zulässige Minimaleinstellung liegt bei 1, eine Maximaleinstellung ist nicht vorhanden.

**Regelsuchbreite.** Die maximale Anzahl von Regelergebnissen pro Zyklus, die für den nächsten Zyklus verwendet werden können. Die zulässige Minimaleinstellung liegt bei 1, eine Maximaleinstellung ist nicht vorhanden.

**Klassenzusammenführungsfaktor.** Der minimale Betrag, um den ein Segment beim Zusammenführen mit dem benachbarten Segment wachsen muss. Die zulässige Minimaleinstellung liegt bei 1,01, eine Maximaleinstellung ist nicht vorhanden.

- **Fehlende Werte in Bedingungen zulassen.** Wenn True, wird in Regeln der Test IS MISSING zugelassen.
- **Zwischenergebnisse verwerfen.** Wenn True, werden nur die Endergebnisse des Suchvorgangs zurückgegeben. Ein Endergebnis ist ein Ergebnis, das im Suchvorgang nicht weiter verfeinert wird. Wenn False, werden auch Zwischenergebnisse zurückgegeben.

**Maximale Anzahl an Alternativen.** Gibt die maximale Anzahl von Alternativen an, die beim Ausführen der Mining-Aufgabe zurückgegeben werden. Die zulässige Minimaleinstellung liegt bei 1, eine Maximaleinstellung ist nicht vorhanden.

Beachten Sie, dass die Mining-Aufgabe nur die tatsächliche Anzahl an Alternativen bis zum angegebenen Maximum zurückgibt. Falls als Maximum beispielsweise 100 angegeben wurde und nur 3 Alternativen gefunden werden, werden nur diese 3 angezeigt.

## ***Modell-Nugget vom Typ "Entscheidungsliste"***

Ein Modell besteht aus einer Liste von **Segmenten**, von denen jedes durch eine **Regel** definiert ist, die übereinstimmende Datensätze auswählt. Sie können die Segmente problemlos anzeigen oder ändern, bevor Sie das Modell erzeugen und auswählen, welche Segmente ein- bzw. ausgeschlossen werden sollen. Beim Scoring ergeben Entscheidungslistenmodelle den Wert *Ja* für eingeschlossene Segmente und *\$null\$* für alles andere, einschließlich des Rests. Durch diese unmittelbare Steuerung des Scorings sind Entscheidungslistenmodelle ideal für das Erstellen von Mailinglisten und sie sind weit verbreitet im Customer Relationship Management, unter anderen in Callcenter- und Marketinganwendungen.

Abbildung 9-8  
Modell-Nugget vom Typ "Entscheidungsliste"

ID	Segmentregeln	Score	Abdeck...	Häufigkeit	Wahrsch...
	Alle Segmente einschließlich R...		13.504	1.952	14,45%
1	months_customer months_customer = "0"	Ausges...	1.747	0	0,00%
2	rfm_score rfm_score <= 0.000	Ausges...	6.003	0	0,00%
3	rfm_score, income rfm_score > 12.333 und income > 52213.000	1	555	456	82,16%
4	income income > 55267.000	1	643	551	85,69%
5	number_transactions, rfm number_transactions > 2 ur1 rfm_score > 12.333		533	206	38,65%

Wenn Sie einen Stream ausführen, der ein Entscheidungslistenmodell enthält, fügt der Knoten drei neue Felder hinzu; diese enthalten den Score, entweder *1* (d. h. *Ja*) für eingeschlossene Felder oder *\$null* für ausgeschlossene Felder, die Wahrscheinlichkeit (Trefferquote) für das Segment, in das der Datensatz fällt, sowie die ID-Nummer des Segments. Die Namen der neuen Felder werden aus dem Namen des prognostizierten Ausgabefelds abgeleitet, dem *\$D-* für den Score, *\$DP-* für die Wahrscheinlichkeit und *\$DI-* für die Segment-ID vorangestellt ist.

Das Modell wird auf der Grundlage des Zielwerts gescort, der zum Zeitpunkt der Modellerstellung angegeben wurde. Sie können Segmente manuell ausschließen, sodass sie beim Scoring den Wert *\$null* erhalten. Wenn Sie beispielsweise eine Suche des Typs "Geringe Wahrscheinlichkeit" durchführen, um Segmente mit unterdurchschnittlichen Trefferquoten zu finden, erhalten diese "niedrigen" Segmente den Score *Ja*, es sei denn, sie schließen Sie manuell aus. Falls erforderlich, können Nullen mithilfe eines Ableitungs- oder Füllerknotens zu *Nein* umkodiert werden.

## PMML

Entscheidungslistenmodelle können als PMML RuleSetModel mit dem Auswahlkriterium "Erster Treffer" gespeichert werden. Es wird jedoch erwartet, dass alle Regeln denselben Score aufweisen. Um Änderungen im Zielfeld oder Zielwert zu berücksichtigen, können mehrere Regelmengenmodelle in einer einzigen Datei gespeichert und nacheinander angewendet werden: Fälle die vom ersten Modell nicht abgedeckt werden, werden an das zweite weitergeleitet usw. Der Algorithmusname *DecisionList* gibt dieses nicht dem Standard entsprechende Verhalten an

und nur Regelmengenmodelle mit diesem Namen werden als Entscheidungslistenmodelle erkannt und als solche gescort.

## Entscheidungslisten-Modell-Nugget – Einstellungen

Auf der Registerkarte “Einstellungen” für ein Modell-Nugget vom Typ “Entscheidungsliste” können Sie Neigungs-Scores ermitteln bzw. die SQL-Optimierung deaktivieren. Diese Registerkarte ist erst verfügbar, nachdem das Modell-Nugget zu einem Stream hinzugefügt wurde.

Abbildung 9-9  
Entscheidungslisten-Modell-Nugget – Einstellungen



**Scores für Rohneigung berechnen.** Bei Modellen mit einem Flag-Ziel (das als Vorhersage “Ja” bzw. “Nein” ausgibt) können Sie Neigungs-Scores anfordern, die die Wahrscheinlichkeit des für das Zielfeld angegebenen wahren Ergebnisses angeben. Diese Werte werden zusätzlich zu den anderen Vorhersage- und Konfidenzwerten angegeben, die ggf. während des Scorings erstellt werden.

**Scores für korrigierte Neigung berechnen.** Rohneigungs-Scores beruhen ausschließlich auf den Trainingsdaten und sind aufgrund der Neigung vieler Modelle zur übermäßigen Anpassung an die Trainingsdaten möglicherweise zu optimistisch. Bei korrigierten Neigungen wird versucht, dies durch Evaluation der Modellleistung anhand einer Test- bzw. Validierungspartition zu kompensieren. Bei dieser Option muss im Stream ein Partitionsfeld definiert sein und die Scores für die korrigierte Neigung müssen im Modellierungsknoten aktiviert werden, bevor das Modell generiert wird.

**SQL für dieses Modell generieren.** Bei Verwendung von Daten aus einer Datenbank kann SQL-Code per Pushback zur Ausführung an die Datenbank zurückübertragen werden. Dadurch kann bei vielen Operationen eine bessere Leistung erzielt werden. [Für weitere Informationen](#)

siehe Thema SQL-Optimierung in Kapitel 6 in *IBM SPSS Modeler Server 14.2-Verwaltungs- und Leistungshandbuch*.

## Decision List Viewer

Mit der einfach bedienbaren, aufgabenbasierten grafischen Benutzeroberfläche von Decision List Viewer entfällt die Komplexität des Modellbildungsvorgangs, da Sie sich nicht mit Details der unteren Ebene der Data Mining-Techniken befassen müssen. Sie können Ihre gesamte Aufmerksamkeit auf die Teile der Analyse konzentrieren, die eine Benutzerintervention erforderlich machen, wie das Festlegen von Zielen, die Auswahl von Zielgruppen, die Analyse der Ergebnisse und die Auswahl des optimalen Modells.

Abbildung 9-10  
Interaktiver Entscheidungslisten-Viewer

ID	Segmentregeln	Score	Abdeckung	Abdeckung (n)	Häufigkeit	Wahrscheinlic...	Fehler
	Alle Segmente einschließlich Rest			13.504	1.952	14,45%	0,00%
1	income > 55267.000 und number_products > 1.000	1		912	795	87,17%	2,35%
2	rfm_score > 10.535 und number_transactions > 3.000	1		725	357	49,24%	3,70%
3	average#balance#feed#index, number_products, rfm_score > 0.000 und average#balance#feed#index <= 349.000 und number_products <= 2.000 und rfm_score > 9.239	1		738	196	26,56%	3,34%
	Rest			11.129	604	5,43%	0,44%

Modellzusammenfassung; Abdeckung 2.375; Häufigkeit 1.348; Wahrscheinlichkeit 56,76%

## Arbeitsmodellbereich

Der Arbeitsmodellbereich zeigt das aktuelle Modell an, einschließlich der Mining-Aufgaben und anderer Aktionen für das Arbeitsmodell.

Abbildung 9-11  
Arbeitsmodellbereich

ID	Segmentregeln	Score	Abdeckung	Abdeckung (n)	Häufigkeit	Wahrscheinlic...	Fehler
	Alle Segmente einschließlich Rest			13.504	1.952	14,45%	0,00%
1	<b>income, number_products</b> income > 55267.000 und number_products > 1.000	1		912	795	87,17%	2,35%
2	<b>rfm_score, number_transactions</b> rfm_score > 10.535 und number_transactions > 3.000	1		725	357	49,24%	3,70%
3	<b>average#balance#feed#index, number_products, rfm_s</b> average#balance#feed#index > 0.000 und average#balance#feed#index <= 349.000 und number_products <= 2.000 und rfm_score > 9.239	1		738	196	26,56%	3,34%
	Rest			11.129	604	5,43%	0,44%

**ID.** Legt die Reihenfolge der Segmente fest. Modellsegmente werden entsprechend der Reihenfolge ihrer ID-Nummer berechnet.

**Segmentregeln.** Gibt den Segmentnamen und die definierten Segmentbedingungen an. Beim Segmentnamen handelt es sich standardmäßig um den Feldnamen oder um aneinandergereihte Feldnamen, die in Bedingungen verwendet werden und durch Kommas getrennt sind.

**Score.** Steht für das vorherzusagende Feld, dessen Werte vermutlich mit den Werten anderer Felder (den Prädiktoren) in Beziehung stehen.

*Hinweis:* Folgende Optionen können für die Anzeige im Dialogfeld [Organisieren von Modellmaßen](#) ausgewählt werden.

**Abdeckung.** Das Kreisdiagramm stellt die Abdeckung der einzelnen Segmente in Bezug zur gesamten Abdeckung visuell dar.

**Abdeckung (n).** Liste der Abdeckung der einzelnen Segmente in Bezug zur gesamten Abdeckung.

**Häufigkeit.** Liste der Anzahl der Treffer in Bezug zur Abdeckung. Beispiel: Wenn die Abdeckung bei 79 liegt und die Häufigkeit bei 50, dann bedeutet dies, dass für das ausgewählte Segment 50 von 79 geantwortet haben.

**Wahrscheinlichkeit.** Zeigt die Wahrscheinlichkeit des Segments an. Beispiel: Wenn die Abdeckung bei 79 liegt und die Häufigkeit bei 50, dann bedeutet dies, dass die Wahrscheinlichkeit für das Segment 63,29 % ist (50 geteilt durch 79).

**Fehler.** Zeigt den Segmentfehler an.

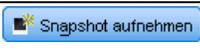
Die am unteren Rand des Bereichs eingeblendeten Informationen zeigen die Abdeckung, die Häufigkeit und die Wahrscheinlichkeit des gesamten Modells an.

### **Symboleiste des Arbeitsmodells**

Der Arbeitsmodellbereich bietet über die Symboleiste die folgenden Funktionen.

*Hinweis:* Einige Funktionen sind auch verfügbar, wenn Sie mit der rechten Maustaste auf ein Modellsegment klicken.

Tabelle 9-1  
Symbolleistenschaltflächen für das Arbeitsmodell

	Startet das Dialogfeld <a href="#">Neues Modell erzeugen</a> , das Optionen für das Erstellen eines neuen Modell-Nuggets enthält.
	Speichert den aktuellen Status der interativen Sitzung. Der Modellierungsknoten der Entscheidungsliste mit den aktuellen Einstellungen wird aktualisiert, darunter Mining-Aufgaben, Modell-Snapshots, Datenauswahl und benutzerdefinierte Maße. Zum Wiederherstellen einer Sitzung in diesem Status aktivieren Sie das Kästchen Gespeicherte Sitzungsinformationen verwenden auf der Registerkarte "Modell" des Modellierungsknotens und klicken Sie auf Ausführen.
	Zeigt das Dialogfeld "Modellmaße organisieren" an. <a href="#">Für weitere Informationen siehe Thema Organisieren von Modellmaßen auf S. 256.</a>
	Zeigt das Dialogfeld "Datenauswahl organisieren" an. <a href="#">Für weitere Informationen siehe Thema Datenauswahl organisieren auf S. 249.</a>
	Zeigt die Registerkarte "Snapshot" an. <a href="#">Für weitere Informationen siehe Thema Registerkarte "Snapshots" auf S. 242.</a>
	Zeigt die Registerkarte "Alternativen" an. <a href="#">Für weitere Informationen siehe Thema Registerkarte "Alternativen" auf S. 240.</a>
	Erstellt einen Snapshot der aktuellen Modellstruktur. Snapshots werden auf der Registerkarte "Snapshots" angezeigt und in der Regel zum Zweck des Modellvergleichs verwendet.
	Startet das Dialogfeld <a href="#">Einfügen von Segmenten</a> , das Optionen für das Erstellen neuer Modellsegmente enthält.
	Startet das Dialogfeld Bearbeiten von Segmentregeln, das Optionen zum Hinzufügen von Bedingungen zu Modellsegmenten oder zum Ändern zuvor definierter Modellsegmentbedingungen enthält.
	Verschiebt das ausgewählte Segment in der Modellhierarchie nach oben.
	Verschiebt das ausgewählte Segment in der Modellhierarchie nach unten.
	Löscht das ausgewählte Segment.
	Schließt das ausgewählte Segment entweder in das Modell ein oder davon aus. Wenn es ausgeschlossen wird, werden die Ergebnisse zum Rest hinzugefügt. Dies unterscheidet sich dahin gehend vom Löschen eines Segments, dass Sie die Option haben, es wieder zu aktivieren.

## **Registerkarte "Alternativen"**

Die Registerkarte "Alternativen" wird generiert, wenn Sie auf Segmente suchen klicken. Auf ihr werden alle alternativen Mining-Ergebnisse für das im Arbeitsmodellfenster ausgewählte Modell oder Segment aufgeführt.

- ▶ Wenn eine Alternative in ein Arbeitsmodell umgewandelt werden soll, markieren Sie die entsprechende Alternative und klicken Sie auf Laden; das alternative Modell wird im Arbeitsmodellfenster angezeigt.

*Hinweis:* Die Registerkarte “Alternativen” ist nur sichtbar, wenn Sie Maximale Anzahl an Alternativen auf der Registerkarte “Experten” im Modellierungsknoten der Entscheidungsliste auf die Erstellung mehrerer Alternativen eingestellt haben.

Abbildung 9-12  
Alternativen, Registerkarte

The screenshot shows a software window titled "Alben modellieren". It contains a table with the following data:

Name	Ziel	Anzahl der Segmente	Abdeckung	Häufigkeit	Wahrschei...
Alternative 1	1	3	2.375	1.348	56,76%
Alternative 2	1	3	2.368	1.326	56,00%
Alternative 3	1	3	2.380	1.329	55,84%

Below the table is a section titled "Alternative Vorschau" showing a detailed view of Alternative 3:

ID	Segmentregeln	Score	Abdeckun...	Häufigkeit	Wahrschei...
	Alle Segmente einschließlich Rest		13.504	1.952	14,45%
1	<b>income, number_products</b> income > 55267.000 und number_products > 1.000	1	912	795	87,17%
2	<b>rfm_score, number_transactions</b> rfm_score > 12.333 und number_transactions > 2.000	1	737	360	48,85%
3	<b>number_transactions, income</b> number_transactions > 0.000 und number_transactions <= 1.000 und income > 46072.000	1	731	174	23,80%

At the bottom of the window, there is a "Laden" button with an upward arrow, a tabbed interface with "Alternativen" selected and "Snapshots" visible, and "OK", "Abbrechen", and "Hilfe" buttons.

Alle generierten Modellalternativen zeigen bestimmte Modellinformationen an:

**Name.** Jede Alternative wird fortlaufend nummeriert. Die erste Alternative enthält in der Regel die besten Ergebnisse.

**Ziel.** Gibt den Zielwert an. Beispiel: 1, was dem Wert “true” entspricht.

**Anzahl der Segmente.** Die Anzahl der Segmentregeln, die im alternativen Modell verwendet werden.

**Abdeckung.** Die Abdeckung des alternativen Modells.

**Häufigkeit** Anzahl der Treffer in Bezug zur Abdeckung.

**Prob.** Gibt die prozentuale Wahrscheinlichkeit des alternativen Modells an.

*Hinweis:* Alternative Ergebnisse werden nicht mit dem Modell gespeichert. Ergebnisse sind nur innerhalb einer aktiven Sitzung gültig.

## **Registerkarte “Snapshots”**

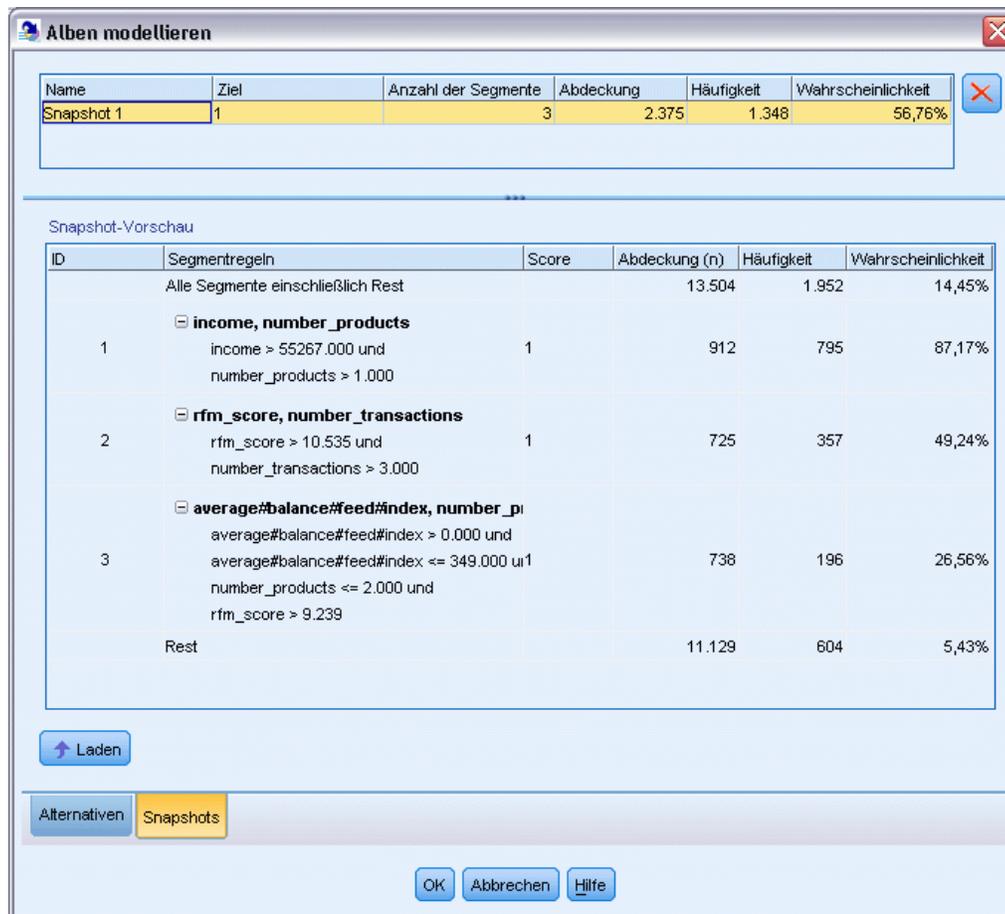
Ein Snapshot ist eine Ansicht eines Modells zu einem bestimmten Zeitpunkt. Sie können beispielsweise einen Modell-Snapshot erstellen, wenn Sie ein anderes alternatives Modell in das Arbeitsmodellfenster laden möchten und die am aktuellen Modell durchgeführten Arbeiten nicht verlieren möchten. Auf der Registerkarte “Snapshots” werden alle manuell erstellten Snapshots für eine beliebige Anzahl von Modellzuständen aufgeführt.

*Hinweis:* Snapshots werden zusammen mit dem Modell gespeichert. Sie sollten einen Snapshot erzeugen, wenn Sie das erste Modell laden. Dieser Snapshot nimmt die Originalstruktur des Modells auf und bietet Ihnen die Möglichkeit, jederzeit zum ursprünglichen Modellzustand zurückzukehren. Der Name des generierten Snapshots wird als Zeitstempel angezeigt, der dem Generierungszeitpunkt entspricht.

### **Erstellen eines Modell-Snapshots**

- ▶ Wählen Sie für die Anzeige im Arbeitsmodellfenster ein geeignetes Modell oder eine Alternative aus.
- ▶ Nehmen Sie am Arbeitsmodell alle erforderlichen Änderungen vor.
- ▶ Klicken Sie auf Snapshot aufnehmen. Auf der Registerkarte “Snapshots” wird ein neuer Snapshot angezeigt.

Abbildung 9-13  
Registerkarte "Snapshots"



**Name.** Der Name des Snapshots. Sie können den Namen eines Snapshots ändern, indem Sie auf den Namen doppelklicken.

**Ziel.** Gibt den Zielwert an. Beispiel: 1, was dem Wert "true" entspricht.

**Anzahl der Segmente.** Die Anzahl der Segmentregeln, die im Modell verwendet werden.

**Abdeckung.** Die Abdeckung des Modells.

**Häufigkeit** Anzahl der Treffer in Bezug zur Abdeckung.

**Prob.** Gibt die prozentuale Wahrscheinlichkeit des Modells an.

- ▶ Wenn ein Snapshot in ein Arbeitsmodell umgewandelt werden soll, markieren Sie den entsprechenden Snapshot und klicken Sie auf Laden; der Snapshot wird im Arbeitsmodellfenster angezeigt.
- ▶ Um einen Snapshot zu löschen, klicken Sie auf Löschen oder klicken Sie mit der rechten Maustaste auf den Snapshot und wählen Sie im Menü Löschen.

## Arbeiten mit Decision List Viewer

Ein Modell, das die Reaktionen und das Verhalten von Kunden am besten vorhersagt, wird in mehreren Phasen erstellt. Wenn Decision List Viewer gestartet wird, wird das Arbeitsmodell mit den definierten Modellsegmenten und Maßen gefüllt. Anschließend können Mining-Aufgaben ausgeführt, die Segmente/Maße geändert und ein neues Modell oder ein Modellbildungsknoten generiert werden.

Sie können eine oder mehrere Segmentregeln hinzufügen, um ein zufrieden stellendes Modell zu entwickeln. Sie können Segmentregeln zum Modell hinzufügen, indem Sie Mining-Aufgaben ausführen oder indem Sie die Funktion Segmentregel bearbeiten anwenden.

Während des Modellbildungsprozesses können Sie die Leistung des Modells bewerten, indem Sie das Modell mit den Maßdaten vergleichen, das Modell in einem Diagramm visuell darstellen oder benutzerdefinierte Excel-Maße erstellen.

Wenn Sie sich hinsichtlich der Qualität des Modells sicher sind, können Sie ein neues Modell erzeugen und es in den IBM® SPSS® Modeler-Zeichenbereich oder in die Modellpalette einfügen.

### Mining-Aufgaben

Eine **Mining-Aufgabe** ist eine Sammlung von Parametern, die festlegen, wie neue Regeln generiert werden. Einige dieser Parameter können ausgewählt werden, um Ihnen die Möglichkeit zu bieten, das Modell an neue Situationen anzupassen. Eine Aufgabe besteht aus einer Aufgabenvorlage (Typ), einer Zielgröße und einer Erstellungsauswahl (Mining-Daten-Set).

Im folgenden Abschnitt werden die verschiedenen Vorgänge von Mining-Aufgaben beschrieben:

- [Ausführen von Mining-Aufgaben](#)
- [Erstellen und Bearbeiten einer Mining-Aufgabe](#)
- [Datenauswahl organisieren](#)

### Ausführen von Mining-Aufgaben

Decision List Viewer bietet Ihnen die Möglichkeit, einem Modell Segmentregeln von Hand hinzuzufügen, indem Sie Mining-Aufgaben ausführen oder indem Sie Segmentregeln zwischen Modellen kopieren und einfügen. Eine Mining-Aufgabe enthält Informationen für das Erstellen neuer Segmentregeln (die Data Mining-Parametereinstellungen, wie die Suchstrategie, Quellenattribute, Breite der Suche, Konfidenzniveau usw.), das vorherzusagende Kundenverhalten und die zu untersuchenden Daten. Das Ziel einer Mining-Aufgabe besteht darin, die bestmöglichen Segmentregeln zu finden.

#### **So erstellen Sie eine Modellsegmentregel, indem Sie eine Mining-Aufgabe ausführen:**

- ▶ Klicken Sie auf die Zeile Rest. Wenn im Arbeitsmodellfenster bereits Modelle angezeigt werden, können Sie auch eines der Segmente auswählen, um auf der Grundlage des ausgewählten Segments nach neuen Regeln zu suchen. Verwenden Sie nach der Auswahl des Rests oder eines Segments eine der folgenden Methoden, um das Modell oder alternative Modelle zu generieren:
  - Wählen Sie im Menü "Extras" die Option Segmente suchen aus.

- Klicken Sie mit der rechten Maustaste auf die Zeile Rest/das Segment und wählen Sie Segmente suchen.
- Klicken Sie im Arbeitsmodellfenster auf die Schaltfläche Segmente suchen.

Während die Aufgabe ausgeführt wird, wird der Fortschritt am unteren Rand des Arbeitsbereichs angezeigt. Dort sehen Sie, wenn die Aufgabe abgeschlossen ist. Wie lange eine Aufgabe genau dauert, hängt von der Komplexität der Mining-Aufgabe und der Größe des Daten-Sets ab. Wenn das Ergebnis nur ein einziges Modell umfasst, wird dieses im Arbeitsmodellfenster angezeigt, sobald die Aufgabe erledigt ist. Wenn das Ergebnis jedoch mehrere Modelle enthält, werden diese in der Registerkarte "Alternativen" angezeigt.

*Hinweis:* Eine Aufgabe schließt entweder mit Modellen, ohne Modelle oder mit einem Fehler ab.

Der Vorgang zum Finden neuer Segmentregeln kann so lange wiederholt werden, bis dem Modell keine neuen Regeln mehr hinzugefügt werden. Dies bedeutet, dass alle signifikanten Gruppen oder Kunden gefunden wurden.

Eine Mining-Aufgabe kann auf einem vorhandenen Modellsegment ausgeführt werden. Wenn die Aufgabe nicht das gesuchte Ergebnis liefert, können Sie auf demselben Segment eine andere Mining-Aufgabe ausführen. So können Sie auf der Grundlage des ausgewählten Segments zusätzliche Regeln finden. Segmente, die sich unterhalb des ausgewählten Segments befinden (d. h., die nach dem ausgewählten Segment zum Modell hinzugefügt wurden), werden durch die neuen Segmente ersetzt, da jedes Segment von seinen Vorgängern abhängt.

### ***Erstellen und Bearbeiten einer Mining-Aufgabe***

Eine Mining-Aufgabe ist der Mechanismus, der nach der Sammlung von Regeln sucht, die einem Datenmodell zugrunde liegen. Neben den in der ausgewählten Vorlage definierten Suchkriterien definiert eine Aufgabe außerdem das Ziel (die Frage, die der Analyse zugrunde liegt, z. B., wie viele Kunden wahrscheinlich auf ein Mailing reagieren werden) und identifiziert die zu verwendenden Daten-Sets. Das Ziel einer Mining-Aufgabe besteht darin, die bestmöglichen Modelle zu finden.

#### ***Erstellen einer Mining-Aufgabe***

So erstellen Sie eine Mining-Aufgabe:

- ▶ Wählen Sie das Segment, aus dem Sie zusätzliche Segmentbedingungen ermitteln möchten.
- ▶ Klicken Sie auf Einstellungen. Das Dialogfeld "Mining-Aufgabe erstellen/bearbeiten" wird geöffnet. Dieses Dialogfeld bietet Optionen für die Definition der Mining-Aufgabe.
- ▶ Nehmen Sie die erforderlichen Änderungen vor und klicken Sie auf OK, um in das Arbeitsmodellfenster zurückzukehren. Decision List Viewer verwendet die Einstellungen als Standards, die für jede Aufgabe ausgeführt werden, bis eine alternative Aufgabe oder Einstellung gewählt wird.
- ▶ Klicken Sie auf Segmente suchen, um die Mining-Aufgabe auf dem ausgewählten Segment zu starten.

### **Mining-Aufgabe bearbeiten**

Das Dialogfeld “Mining-Aufgabe erstellen/bearbeiten” bietet Optionen zur Definition einer neuen Mining-Aufgabe oder zum Bearbeiten einer vorhandenen.

Die meisten für Mining-Aufgaben verfügbaren Parameter entsprechen denen für den Entscheidungslistenknoten. Die Ausnahmen werden unten gezeigt. [Für weitere Informationen siehe Thema Entscheidungslistenmodell – Optionen auf S. 232.](#)

Abbildung 9-14  
Dialogfeld “Mining-Aufgabe erstellen/bearbeiten”

**Mining-Aufgabe erstellen/bearbeiten: response[1]**

Einstellungen laden:

Ziel

Zielfeld:  response Zielwert: 1

Einfache Einstellungen

Segmente suchen mit:

Maximale Anzahl an neuen Segmenten:

Minimale Segmentgröße

Als Prozentsatz des vorherigen Segments (%):

Als absoluter Wert (N):

Maximale Anzahl an Alternativen:

Maximale Anzahl an Attributen pro Segment:

Wiederverwendung von Attributen innerhalb des Segments zulassen

Konfidenzintervall für neue Bedingungen (%):

Experteneinstellungen

Klassiermethode:	Gleiche Anzahl	Anzahl der Klassen:	10
Modellsuchbreite:	5	Regelsuchbreite:	5
Klassenzusammenführungsfaktor:	2.00		
Fehlende Werte in Bedingungen zulassen:	Wahr	Zwischenergebnisse verwerfen:	Wahr

Daten

Erstellungsauswahl:

Verfügbare Felder:  Alle Felder  Benutzerdefiniert

**Einstellungen laden:** Wenn Sie mehrere Mining-Aufgaben erstellt haben, wählen Sie die erforderliche Aufgabe aus.

**Neu...** Klicken Sie auf diese Option, um eine neue Mining-Aufgabe auf der Basis der Einstellungen der aktuell angezeigten Aufgabe zu erstellen.

**Target**

**Zielfeld:** Steht für das vorherzusagende Feld, dessen Werte vermutlich mit den Werten anderer Felder (den Prädiktoren) in Beziehung stehen.

**Zielwert.** Gibt den Wert des Zielfelds an, in dem das zu modellierende Ergebnis angegeben wird. Wenn beispielsweise das Zielfeld “churn” (Abwanderung) mit 0 = no und 1 = yes kodiert ist, geben Sie 1 an, um Regeln festzulegen, mit denen die Datensätze für wahrscheinlich abwandernde Kunden identifiziert werden.

**Einfache Einstellungen (SimpleSettings)**

**Maximale Anzahl an Alternativen.** Gibt die Anzahl von Alternativen an, die beim Ausführen der Mining-Aufgabe angezeigt werden. Die zulässige Minimaleinstellung liegt bei 1, eine Maximaleinstellung ist nicht vorhanden.

**Experteneinstellungen (ExpertSettings)**

**Bearbeiten...** Öffnet das Dialogfeld Erweiterte Parameter bearbeiten, in dem Sie die erweiterten Einstellungen festlegen können. [Für weitere Informationen siehe Thema Erweiterte Parameter bearbeiten auf S. 247.](#)

**Data**

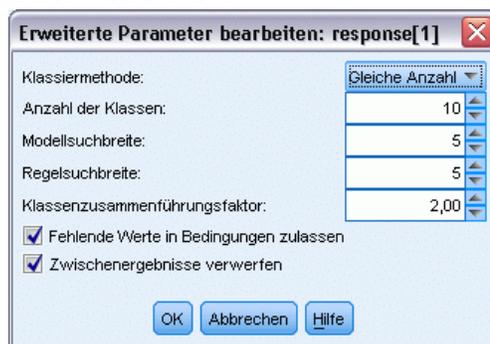
**Erstellungsauswahl.** Bietet Optionen zur Angabe des Evaluationsmaßes, das Decision List Viewer analysieren soll, um neue Regeln zu finden. Die aufgeführten Evaluationsmaße werden im Dialogfeld “Datenauswahl organisieren” erstellt/bearbeitet.

**Verfügbare Felder.** Bietet Optionen zur Anzeige aller Felder oder zur manuellen Auswahl der anzuzeigenden Felder.

**Bearbeiten...** Wenn die Option Benutzerdefiniert ausgewählt ist, wird das Dialogfeld Verfügbare Felder anpassen geöffnet, in dem Sie auswählen können, welche Felder als Segmentattribute verfügbar sind, die die Mining-Aufgabe findet. [Für weitere Informationen siehe Thema Verfügbare Felder anpassen auf S. 248.](#)

**Erweiterte Parameter bearbeiten**

Abbildung 9-15  
Erweiterte Parameter



Das Dialogfeld “Erweiterte Parameter bearbeiten” bietet die folgenden Konfigurationsoptionen.

**Klassiermethode.** Die Methode zum Klassieren kontinuierlicher Felder (gleiche Anzahl oder gleiche Breite).

**Anzahl der Klassen.** Die Anzahl der für kontinuierliche Felder zu erstellenden Klassen. Die zulässige Minimaleinstellung liegt bei 2, eine Maximaleinstellung ist nicht vorhanden.

**Modellsuchbreite.** Die maximale Anzahl von Modellergebnissen pro Zyklus, die für den nächsten Zyklus verwendet werden können. Die zulässige Minimaleinstellung liegt bei 1, eine Maximaleinstellung ist nicht vorhanden.

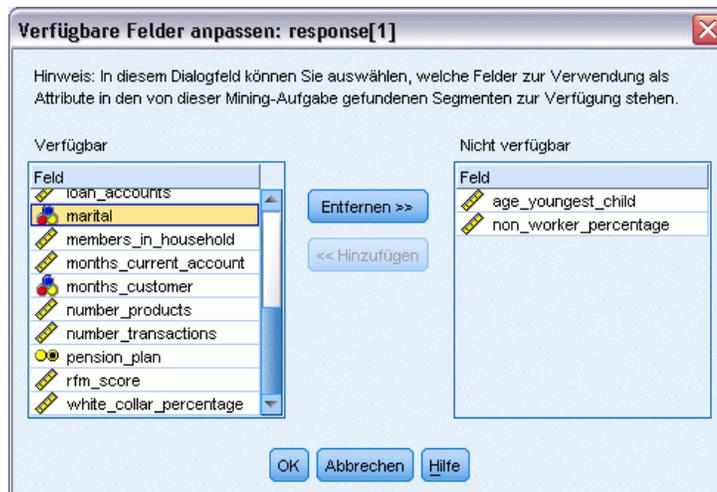
**Regelsuchbreite.** Die maximale Anzahl von Regelergebnissen pro Zyklus, die für den nächsten Zyklus verwendet werden können. Die zulässige Minimaleinstellung liegt bei 1, eine Maximaleinstellung ist nicht vorhanden.

**Klassenzusammenführungsfaktor.** Der minimale Betrag, um den ein Segment beim Zusammenführen mit dem benachbarten Segment wachsen muss. Die zulässige Minimaleinstellung liegt bei 1,01, eine Maximaleinstellung ist nicht vorhanden.

- **Fehlende Werte in Bedingungen zulassen.** Wenn True, wird in Regeln der Test IS MISSING zugelassen.
- **Zwischenergebnisse verwerfen.** Wenn True, werden nur die Endergebnisse des Suchvorgangs zurückgegeben. Ein Endergebnis ist ein Ergebnis, das im Suchvorgang nicht weiter verfeinert wird. Wenn False, werden auch Zwischenergebnisse zurückgegeben.

### Verfügbare Felder anpassen

Abbildung 9-16  
Dialogfeld “Verfügbare Felder anpassen”



Im Dialogfeld “Verfügbare Felder anpassen” haben Sie die Möglichkeit, die Felder auszuwählen, die als Segmentattribute verfügbar sind, die die Mining-Aufgabe findet.

**Verfügbar.** Führt die Felder auf, die aktuell als Segmentattribute verfügbar sind. Um Felder aus der Liste zu entfernen, wählen Sie die entsprechenden Felder aus und klicken Sie auf Entfernen>>. Die ausgewählten Felder werden aus der Liste der verfügbaren Felder in die Liste der nicht verfügbaren Felder verschoben.

**Nicht verfügbar.** Führt die Felder auf, die nicht als Segmentattribute verfügbar sind. Um diese Felder in die Liste der verfügbaren Felder aufzunehmen, wählen Sie die entsprechenden Felder aus und klicken Sie auf << Hinzufügen. Die ausgewählten Felder werden aus der Liste der nicht verfügbaren Felder in die Liste der verfügbaren Felder verschoben.

### **Datenauswahl organisieren**

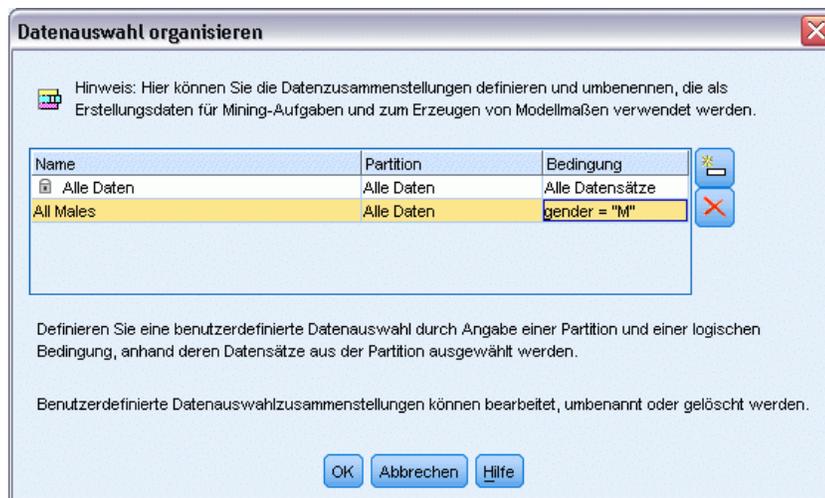
Durch das Organisieren einer Datenauswahl (einem Mining-Daten-Set) können Sie festlegen, welche Evaluationsmaße Decision List Viewer analysieren soll, um neue Regeln zu suchen, und welche Datenauswahl als Grundlage für die Maße verwendet wird.

#### **So organisieren Sie eine Datenauswahl:**

- ▶ Wählen Sie im Menü “Extras” die Option Datenauswahl organisieren oder klicken Sie mit der rechten Maustaste auf ein Segment und wählen Sie die Option. Das Dialogfeld “Datenauswahl organisieren” wird geöffnet.

Abbildung 9-17

Dialogfeld “Datenauswahl organisieren”

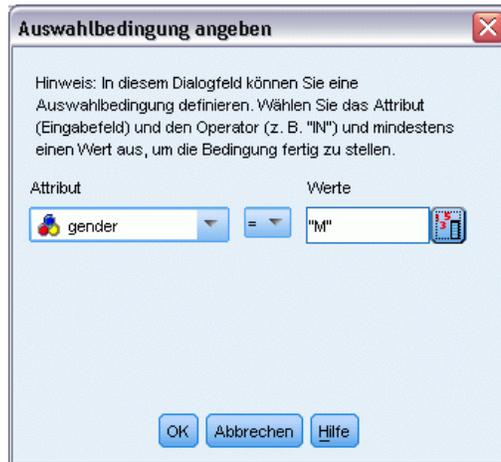


*Hinweis:* Im Dialogfeld “Datenauswahl organisieren” können Sie außerdem eine vorhandene Datenauswahl bearbeiten oder löschen.

- ▶ Klicken Sie auf die Schaltfläche Neue Datenauswahl hinzufügen. Zur vorhandenen Tabelle wird ein neuer Datenauswahleintrag hinzugefügt.
- ▶ Klicken Sie auf Name und geben Sie für die Auswahl einen geeigneten Namen ein.
- ▶ Klicken Sie auf Partition und wählen Sie einen Partitionstyp aus.

- Klicken Sie auf Bedingung und wählen Sie eine Bedingungsoption aus. Wenn Angeben ausgewählt ist, wird das Dialogfeld “Auswahlbedingung angeben” geöffnet, das Optionen zur Angabe bestimmter Feldbedingungen enthält.

Abbildung 9-18  
Dialogfeld “Auswahlbedingung angeben”



- Definieren Sie die entsprechende Bedingung und klicken Sie auf OK.

Die Datenauswahl ist im Dialogfeld “Mining-Aufgabe erstellen/bearbeiten” in der Dropdown-Liste “Erstellungsauswahl” verfügbar. In der Liste können Sie auswählen, welches Evaluationsmaß für eine bestimmte Mining-Aufgabe verwendet wird.

### **Segmentregeln**

Modellsegmentregeln finden Sie, indem Sie eine Mining-Aufgabe ausführen, die auf einer Aufgabenvorlage basiert. Sie können einem Modell manuell Regeln hinzufügen, wenn Sie die Funktionen “Segment einfügen” oder “Segmentregel bearbeiten” verwenden.

Wenn Sie neue Segmentregeln mithilfe einer Mining-Aufgabe suchen, werden die Ergebnisse (sofern vorhanden) auf der Registerkarte “Viewer” des Dialogfelds “Interaktive Liste” angezeigt. Sie können Ihr Modell schnell verfeinern, indem Sie eines der alternativen Ergebnisse aus dem Dialogfeld “Alben modellieren” auswählen und auf Laden klicken. Auf diese Art können Sie so lange mit verschiedenen Ergebnissen experimentieren, bis Sie ein Modell erstellen können, das Ihre optimale Zielgruppe genau beschreibt.

### **Einfügen von Segmenten**

Sie können einem Modell manuell Regeln hinzufügen, wenn Sie die Funktion “Segment einfügen” verwenden.

#### **So fügen Sie eine Segmentregelbedingung zu einem Modell hinzu:**

- Wählen Sie im Dialogfeld “Interaktive Liste” eine Position aus, an der Sie ein neues Segment hinzufügen möchten. Das neue Segment wird direkt über dem ausgewählten Segment eingefügt.

- ▶ Wählen Sie im Menü Bearbeiten die Option Segment einfügen aus oder greifen Sie darauf zu, indem Sie mit der rechten Maustaste auf ein Segment klicken.

Das Dialogfeld “Segment einfügen” wird geöffnet, in dem Sie neue Segmentregelbedingungen einfügen können.

- ▶ Klicken Sie auf Einfügen. Das Dialogfeld “Bedingung einfügen” wird geöffnet, in dem Sie die Attribute für die neue Regelbedingung definieren können.

- ▶ Wählen Sie in den Dropdown-Listen ein Feld und einen Operator aus.

*Hinweis:* Wenn Sie den Operator Nicht in auswählen, funktioniert die ausgewählte Bedingung als Ausschlussbedingung und wird im Dialogfeld Regel einfügen in Rot angezeigt. Wenn beispielsweise die Bedingung `region = 'TOWN'` in Rot angezeigt wird, bedeutet dies, dass TOWN vom Ergebnis-Set ausgeschlossen ist.

- ▶ Geben Sie einen oder mehrere Werte ein oder klicken Sie auf das Symbol Wert einfügen, um das Dialogfeld **Wert einfügen** zu öffnen. In diesem Dialogfeld können Sie einen Wert wählen, der für das ausgewählte Feld definiert ist. Das Eingabefeld verheiratet würde beispielsweise die Werte ja und nein liefern.

- ▶ Klicken Sie auf OK, um zum Dialogfeld “Segment einfügen” zurückzukehren. Klicken Sie ein zweites Mal auf OK, um das erstellte Segment dem Modell hinzuzufügen.

Das neue Segment wird an der angegebenen Position im Modell angezeigt.

### ***Bearbeiten von Segmentregeln***

Mithilfe der Funktion “Segmentregel bearbeiten” können Sie Segmentregelbedingungen hinzufügen, ändern oder löschen.

#### **So ändern Sie eine Segmentregelbedingung:**

- ▶ Wählen Sie das Modellsegment aus, das Sie bearbeiten möchten.
- ▶ Wählen Sie im Menü “Bearbeiten” die Option Segmentregel bearbeiten oder klicken Sie mit der rechten Maustaste auf die Regel, um auf diese Option zuzugreifen.

Das Dialogfeld “Segmentregel bearbeiten” wird geöffnet.

- ▶ Wählen Sie die entsprechende Bedingung aus und klicken Sie auf Bearbeiten.

Das Dialogfeld “Bedingung bearbeiten” wird geöffnet, in dem Sie die Attribute für die ausgewählte Regelbedingung definieren können.

- ▶ Wählen Sie in den Dropdown-Listen ein Feld und einen Operator aus.

*Hinweis:* Wenn Sie den Operator Nicht in auswählen, funktioniert die ausgewählte Bedingung als Ausschlussbedingung und wird im Dialogfeld Segmentregel bearbeiten in Rot angezeigt. Wenn

beispielsweise die Bedingung `region = 'TOWN'` in Rot angezeigt wird, bedeutet dies, dass TOWN vom Ergebnis-Set ausgeschlossen ist.

- ▶ Geben Sie einen oder mehrere Werte ein oder klicken Sie auf die Schaltfläche **Wert einfügen**, um das Dialogfeld **Wert einfügen** zu öffnen. In diesem Dialogfeld können Sie einen Wert wählen, der für das ausgewählte Feld definiert ist. Das Eingabefeld `verheiratet` würde beispielsweise die Werte ja und nein liefern.
- ▶ Klicken Sie auf OK, um zum Dialogfeld **Segmentregel bearbeiten** zurückzukehren. Klicken Sie ein zweites Mal auf OK, um zum Arbeitsmodell zurückzukehren.

Das ausgewählte Modell wird mit den aktualisierten Regelbedingungen angezeigt.

### **Löschen von Segmentregelbedingungen**

**So löschen Sie eine Segmentregelbedingung:**

- ▶ Wählen Sie das Modellsegment aus, das die Regelbedingungen enthält, die Sie löschen möchten.
- ▶ Wählen Sie im Menü "Bearbeiten" die Option **Segmentregel bearbeiten** oder klicken Sie mit der rechten Maustaste auf das Segment, um auf diese Option zuzugreifen.

Das Dialogfeld "Segmentregel bearbeiten" wird geöffnet, in dem Sie eine oder mehrere Segmentregelbedingungen löschen können.

- ▶ Wählen Sie die entsprechende Regelbedingung aus und klicken Sie auf **Löschen**.
- ▶ Klicken Sie auf OK.

Das Löschen einer oder mehrerer Segmentregelbedingungen bewirkt, dass die Maßzahlen im Arbeitsmodellbereich aktualisiert werden.

### **Kopieren von Segmenten**

Decision List Viewer bietet ein bequemes Verfahren zum Kopieren von Modellsegmenten. Wenn Sie ein Segment eines Modells auf ein anderes Modell anwenden möchten, dann kopieren Sie das Segment einfach in einem Modell (oder schneiden Sie es aus) und fügen Sie es in ein anderes Modell ein. Sie können auch ein Segment eines Modells kopieren, das im Bereich "Alternative Vorschau" angezeigt wird, und es in das Modell einfügen, das im Arbeitsmodellfenster angezeigt wird. Diese Funktionen zum Ausschneiden, Kopieren und Einfügen verwenden zum Speichern und Abrufen temporärer Daten die Zwischenablage des Systems. Dies bedeutet, dass die Bedingungen und das Ziel in die Zwischenablage kopiert werden. Die Inhalte der Zwischenablage sind nicht für die Verwendung in Decision List Viewer reserviert und können auch in andere Anwendungen eingefügt werden. Wenn der Inhalt der Zwischenablage beispielsweise in einen Texteditor eingefügt wird, werden die Bedingungen und das Ziel im XML-Format eingefügt.

So kopieren Sie Modellsegmente oder schneiden diese aus:

- ▶ Wählen Sie das Modellsegment aus, das Sie in einem anderen Modell verwenden möchten.

- ▶ Wählen Sie im Menü “Bearbeiten” die Option **Kopieren** (oder **Ausschneiden**) oder klicken Sie mit der rechten Maustaste auf das Modellsegment und wählen Sie **Kopieren** oder **Ausschneiden**.
- ▶ Öffnen Sie das entsprechende Modell (in das Sie das Modellsegment einfügen möchten).
- ▶ Wählen Sie eines der Modellsegmente und klicken Sie auf **Einfügen**.

*Hinweis:* Anstatt der Befehle **Ausschneiden**, **Kopieren** und **Einfügen** können Sie auch folgende Tastenkombinationen verwenden: **Strg+X**, **Strg+C** und **Strg+V**.

Das kopierte (oder ausgeschnittene) Segment wird oberhalb des ausgewählten Modellsegments eingefügt. Die Maße des eingefügten Segments und der darunterliegenden Segmente werden (neu) berechnet.

*Hinweis:* Beide Modelle müssen bei diesem Vorgang auf derselben Modellvorlage basieren und dasselbe Ziel enthalten. Andernfalls wird eine Fehlermeldung angezeigt.

### **Alternative Modelle**

Wenn mehrere Ergebnisse vorhanden sind, zeigt die Registerkarte “Alternativen” die Ergebnisse der einzelnen Mining-Aufgaben an. Jedes Ergebnis besteht aus den Bedingungen der ausgewählten Daten, die am stärksten mit dem Ziel übereinstimmen, sowie allen Alternativen, die als ausreichend gut eingestuft werden. Die Gesamtzahl der angezeigten Alternativen hängt von den Suchkriterien ab, die im Analysevorgang verwendet werden.

#### **So zeigen Sie alternative Modelle an:**

- ▶ Klicken Sie auf der Registerkarte “Alternativen” auf ein alternatives Modell. Im Fenster “Alternative Vorschau” werden die alternativen Modellsegmente angezeigt bzw. sie ersetzen die aktuellen Modellsegmente.
- ▶ Um im Arbeitsmodellfenster mit einem alternativen Modell zu arbeiten, wählen Sie das Modell aus und klicken Sie im Fenster “Alternative Vorschau” auf **Laden** oder klicken Sie mit der rechten Maustaste in der Registerkarte “Alternativen” auf den Namen einer Alternative und wählen Sie **Laden**.

*Hinweis:* Alternative Modelle werden nicht gespeichert, wenn ein neues Modell erzeugt wird.

### **Anpassen eines Modells**

Daten sind nicht statisch. Kunden ziehen um, heiraten und ändern ihren Arbeitsplatz. Produkte fallen aus dem Marktfokus und veralten.

Decision List Viewer bietet Unternehmensanwendern die Flexibilität, Modelle einfach und schnell an neue Situationen anzupassen. Sie können ein Modell ändern, indem Sie es bearbeiten, mit Prioritäten versehen, löschen oder bestimmte Modellsegmente deaktivieren.

### **Prioritäten für Segmente zuweisen**

Sie können Segmentregeln in beliebiger Reihenfolge eine Rangfolge zuweisen. Standardmäßig werden Modellsegmente in der Reihenfolge ihrer Priorität angezeigt, wobei das erste Segment die höchste Priorität besitzt. Wenn Sie einem oder mehreren Segmenten eine andere Priorität zuweisen, wird das Modell entsprechend geändert. Sie können das Modell an die Anforderungen anpassen, indem Sie Segmente in eine höhere oder niedrigere Prioritätsposition verschieben.

#### **So weisen Sie Modellsegmenten Prioritäten zu:**

- ▶ Wählen Sie das Modellsegment aus, dem Sie eine andere Priorität zuweisen möchten.
- ▶ Klicken Sie in der Symbolleiste des Arbeitsmodellfensters auf eine der beiden Pfeilschaltflächen, um das ausgewählte Modellsegment in der Liste nach oben oder nach unten zu verschieben.

Nach dem Zuweisen der Priorität werden alle vorherigen Bewertungsergebnisse neu berechnet und die neuen Werte werden angezeigt.

### **Löschen von Segmenten**

#### **So löschen Sie eines oder mehrere Segmente:**

- ▶ Wählen Sie ein Modellsegment aus.
- ▶ Wählen Sie im Menü "Bearbeiten" die Option Segment löschen oder klicken Sie in der Symbolleiste des Arbeitsmodellfensters auf die Schaltfläche "Löschen".

Die Maße werden für das geänderte Modell neu berechnet und das Modell entsprechend geändert.

### **Ausschließen von Segmenten**

Wenn Sie nach bestimmten Gruppen suchen, werden Sie Geschäftsaktionen wahrscheinlich auf der Grundlage einer Auswahl von Modellsegmenten entscheiden. Wenn Sie ein Modell bereitstellen, können Sie auswählen, Segmente innerhalb des Modells auszuschließen. Ausgeschlossene Segmente werden als Nullwerte gescort. Wenn ein Segment ausgeschlossen wird, bedeutet dies nicht, dass es nicht verwendet wird. Es bedeutet, dass alle Datensätze, die dieser Regel entsprechen, aus der Mailingliste ausgeschlossen werden. Die Regel wird weiterhin angewendet, aber in anderer Form.

#### **So schließen Sie bestimmte Modellsegmente aus:**

- ▶ Wählen Sie im Arbeitsmodellfenster ein Segment aus.
- ▶ Klicken Sie in der Symbolleiste des Arbeitsmodellfensters auf die Schaltfläche Zwischen Segmentausschluss umschalten. In der ausgewählten Spalte "Ziel" des ausgewählten Segments wird nun Ausgeschlossen angezeigt.

*Hinweis:* Anders als beim Löschen von Segmenten stehen ausgeschlossene Segmente weiterhin für die Verwendung im endgültigen Modell zur Verfügung. Ausgeschlossene Segmente wirken sich auf die Diagrammergebnisse aus.

### **Zielwert ändern**

Im Dialogfeld “Zielwert ändern” können Sie den Zielwert für das aktuelle Zielfeld ändern.

Snapshots und Sitzungsergebnisse mit einem anderen Zielwert als dem des Arbeitsmodells werden dahin gehend gekennzeichnet, dass der Tabellenhintergrund der entsprechenden Spalte gelb dargestellt wird. Damit wird angezeigt, dass das Snapshot-/Sitzungsergebnis veraltet ist.

Im Dialogfeld **Mining-Aufgabe erstellen/bearbeiten** wird der Zielwert für das aktuelle Arbeitsmodell angezeigt. Der Zielwert wird nicht mit der Mining-Aufgabe gespeichert. Er wird stattdessen dem Wert des Arbeitsmodells entnommen.

Wenn Sie ein gespeichertes Modell als Arbeitsmodell übernehmen, das einen anderen Zielwert besitzt als das aktuelle Arbeitsmodell (indem Sie beispielsweise ein alternatives Ergebnis oder eine Kopie eines Snapshots bearbeiten), wird der Zielwert des gespeicherten Modells dahingehend geändert, dass er mit dem des Arbeitsmodells übereinstimmt (der im Arbeitsmodellfenster angezeigte Zielwert ändert sich nicht). Die Modellmaße werden mit dem neuen Ziel erneut bewertet.

### **Neues Modell erzeugen**

Das Dialogfeld “Neues Modell erzeugen” bietet Optionen für die Benennung des Modells und zur Festlegung, wo der neue Knoten erstellt werden soll.

**Modellname.** Mit der Option Angepasst können Sie den automatisch erzeugten Namen anpassen oder einen eindeutigen Namen für den Knoten erstellen, der im Stream-Zeichenbereich angezeigt wird.

**Knoten erstellen auf.** Wenn Sie Zeichenbereich auswählen, wird das neue Modell im Arbeitszeichenbereich platziert. Wenn Sie Generierte Modelle auswählen, wird das neue Modell in der Modellpalette platziert. Wenn Sie Beide auswählen, wird das neue Modell sowohl im Zeichenbereich als auch in der Modellpalette platziert.

**Status der interaktiven Sitzung einbeziehen.** Wenn diese Option aktiviert ist, wird der Status der interaktiven Sitzung im generierten Modell erhalten. Wenn Sie später aus dem Modell einen Modellbildungsknoten generieren, wird der Status übernommen und für die Initialisierung der interaktiven Sitzung verwendet. Unabhängig davon, ob die Option aktiviert ist, scort das Modell neue Daten identisch. Wenn die Option nicht ausgewählt ist, ist das Modell immer noch in der Lage, einen Erstellungsknoten zu erstellen, dabei wird es sich aber um einen allgemeineren Erstellungsknoten handeln, der eine neue interaktive Sitzung startet, anstatt die alte Sitzung dort wieder aufzunehmen, wo sie verlassen wurde. Wenn Sie die Knoteneinstellungen ändern und dann einen gespeicherten Status ausführen, werden Ihre geänderten Einstellungen aber ignoriert und die Einstellungen des gespeicherten Status angewendet.

*Hinweis:* Die Standardmaße sind die einzigen Maße, die mit dem Modell abgelegt werden. Zusätzliche Maße werden zusammen mit dem interaktiven Status gespeichert. Das generierte Modell repräsentiert nicht den gespeicherten Status der interaktiven Mining-Aufgabe. Nach dem Starten von Decision List Viewer werden die ursprünglich im Viewer vorgenommenen Einstellungen angezeigt.

Für weitere Informationen siehe Thema Erneutes Erzeugen eines Modellierungsknotens in Kapitel 3 auf S. 70.

### ***Modellauswertung***

Eine erfolgreiche Modellbildung erfordert vor der Implementierung des Modells in einer Produktionsumgebung eine sorgfältige Auswertung des Modells. Decision List Viewer bietet eine gewisse Anzahl von statistischen und betriebswirtschaftlichen Methoden, die für die Bewertung der Auswirkung des Modells in der Realität verwendet werden können. Dazu gehören Gewinn diagramme und eine vollständige Interoperabilität mit Excel, wodurch für die Bewertung der Auswirkung der Bereitstellung Kosten-/Nutzen-Szenarios simuliert werden können.

Modelle können auf folgende Arten ausgewertet werden:

- Mithilfe der in Decision List Viewer vorhandenen statistischen und betriebswirtschaftlichen Methoden (Wahrscheinlichkeit, Häufigkeit).
- Evaluierung von aus Microsoft Excel importierten Maßen.
- Visualisierung des Modells mithilfe eines Gewinn diagramms.

### ***Organisieren von Modellmaßen***

Decision List Viewer bietet Optionen für die Definition von Maßen, die als Spalten berechnet und angezeigt werden. Jedes Segment kann die standardmäßige Abdeckung, Häufigkeit, Wahrscheinlichkeit und Fehlermaße als Spalten enthalten. Sie können außerdem neue Maße erstellen, die als Spalten angezeigt werden.

#### ***Definieren von Modellmaßen***

**So fügen Sie zu Ihrem Modell ein Maß hinzu bzw. definieren Sie ein vorhandenes Maß:**

- ▶ Wählen Sie im Menü "Extras" die Option Modellmaße organisieren oder klicken Sie mit der rechten Maustaste auf das Modell, um diese Option auszuwählen. Das Dialogfeld "Modellmaße organisieren" wird geöffnet.

Abbildung 9-19  
Dialogfeld "Modellmaße organisieren"



- ▶ Klicken Sie auf die Schaltfläche Neues Modellmaß hinzufügen (rechts von der Spalte "Anzeigen"). In der Tabelle wird ein neues Maß angezeigt.
- ▶ Geben Sie einen Namen für das Maß sowie den entsprechenden Typ, die Anzeigoption und die Auswahl ein. In der Spalte "Anzeigen" wird angegeben, ob das Maß für das Arbeitsmodell angezeigt wird. Wenn Sie ein vorhandenes Maß definieren, wählen Sie eine entsprechende Metrik und Auswahl und geben Sie an, ob das Maß für das Arbeitsmodell angezeigt wird.
- ▶ Klicken Sie auf OK, um zum Arbeitsbereich von Decision List Viewer zurückzukehren. Wenn für das neue Maß die Spalte "Anzeigen" aktiviert wurde, wird das neue Maß für das Arbeitsmodell angezeigt.

### **Benutzerdefinierte Maße in Excel**

Für weitere Informationen siehe Thema [Auswertung in Excel](#) auf S. 258.

### **Aktualisieren von Maßen**

In bestimmten Fällen kann es notwendig sein, die Modellmaße neu zu berechnen, beispielsweise wenn Sie ein vorhandenes Modell auf ein neues Kunden-Set anwenden.

#### **So werden Modellmaße neu berechnet (aktualisiert):**

- Wählen Sie im Menü "Bearbeiten" die Option Alle Maße aktualisieren.

oder

- Drücken Sie die Taste F5.

Alle Maße werden neu berechnet und die neuen Werte werden im Arbeitsmodell angezeigt.

### **Auswertung in Excel**

Decision List Viewer kann mit Microsoft Excel integriert werden, wodurch Sie Ihre eigenen Wertberechnungen und Gewinnformeln direkt im Modellbildungsablauf verwenden können, um Kosten-Nutzen-Szenarios zu simulieren. Über die Verknüpfung mit Excel können Sie Daten in Excel exportieren, wo diese verwendet werden können, um Präsentationsgrafiken zu erstellen, benutzerdefinierte Maße zu berechnen, wie komplexe Gewinn- und ROI-Maße, und diese während der Modellbildung in Decision List Viewer anzuzeigen.

Für weitere Informationen siehe [Thema Berechnen von benutzerdefinierten Maßen mithilfe von Excel in Kapitel 12 in IBM SPSS Modeler 14.2- Anwendungshandbuch](#).

*Hinweis:* Damit Sie mit einer Excel-Tabelle arbeiten können, muss der Experte für analytisches CRM die Konfigurationsdaten für die Synchronisierung von Decision List Viewer mit Microsoft Excel definieren. Die Konfiguration befindet sich in der Tabelle einer Excel-Datei. Sie gibt an, welche Informationen von Decision List Viewer an Excel und umgekehrt übertragen werden.

Die folgenden Schritte sind nur dann gültig, wenn MS Excel installiert ist. Wenn Excel nicht installiert ist, werden die Optionen für die Synchronisierung von Modellen mit Excel nicht angezeigt.

#### **So synchronisieren Sie Modelle mit MS Excel:**

- ▶ Öffnen Sie das Modell, führen Sie eine interaktive Sitzung aus und wählen Sie im Menü “Extras” die Option Modellmaße organisieren aus.
- ▶ Wählen Sie Ja für die Option Benutzerdefinierte Maße in Excel berechnen. Das Feld Arbeitsmappe wird aktiviert, in dem Sie eine vorkonfigurierte Excel-Arbeitsmappenvorlage auswählen können.
- ▶ Klicken Sie auf die Schaltfläche Verbindung mit Excel herstellen. Das Dialogfeld “Öffnen” wird angezeigt, in dem Sie in Ihrem lokalen Netzwerkdateisystem zum Speicherort der vorkonfigurierten Vorlage navigieren können.
- ▶ Wählen Sie die entsprechende Excel-Vorlage aus und klicken Sie auf Öffnen. Die ausgewählte Excel-Vorlage wird gestartet. Wechseln Sie mithilfe der Windows-Taskleiste (oder durch Drücken von Alt-Tab) zurück zum Dialogfeld “Eingaben für benutzerdefinierte Maße”.
- ▶ Wählen Sie die Zuordnungen zwischen den in der Excel-Vorlage definierten Metrikenamen und den Metrikenamen des Modells aus und klicken Sie auf OK.

Nachdem die Verknüpfung hergestellt ist, startet Excel mit der vorkonfigurierten Excel-Vorlage, die die Modellregeln in einer Tabelle anzeigt. Die in Excel berechneten Ergebnisse werden in Decision List Viewer als neue Spalten angezeigt.

*Hinweis:* Excel-Maße werden nicht mit dem Modell gespeichert. Sie gelten nur während der aktiven Sitzung. Sie können jedoch Snapshots erstellen, die Excel-Maße enthalten. Die in Snapshot-Ansichten gespeicherten Excel-Maße eignen sich für den historischen Vergleich und werden beim erneuten Öffnen nicht aktualisiert. [Für weitere Informationen siehe Thema Registerkarte “Snapshots” auf S. 242.](#) Die Excel-Maße werden erst in Snapshots angezeigt, wenn Sie die Verbindung zu der Excel-Vorlage wiederherstellen.

### **MS Excel-Integration - Setup**

Die Integration von Decision List Viewer und Microsoft Excel erfolgt über vorkonfigurierte Excel-Tabellenvorlagen. Die Vorlage besteht aus drei Arbeitsblättern:

**Modellmaße.** Zeigt die importierten Decision List Viewer-Maße, die benutzerdefinierten Excel-Maße und die Summen der Berechnungen (die im Arbeitsblatt “Einstellungen” definiert sind).

**Einstellungen.** Enthält die Variablen, mit denen Berechnungen auf der Grundlage der importierten Decision List Viewer-Maße und der benutzerdefinierten Excel-Maße erstellt werden.

**Konfiguration.** Enthält Optionen, mit denen festgelegt wird, welche Maße aus Decision List Viewer importiert werden, und mit denen benutzerdefinierte Excel-Maße definiert werden.

**ACHTUNG:** Die Struktur des Arbeitsblatts “Konfiguration” ist streng definiert. Bearbeiten Sie **KEINESFALLS** Zellen im grau schattierten Bereich.

- **Maße aus Modell.** Gibt an, welche Decision List Viewer-Maße in den Berechnungen verwendet werden.
- **Maße an Modell.** Gibt an, welche von Excel generierten Maße an Decision List Viewer zurückgegeben werden. Die von Excel generierten Maße werden in Decision List Viewer als neue Maßspalten angezeigt.

*Hinweis:* Excel-Maße werden beim Generieren eines neuen Modells nicht erhalten. Sie gelten nur während der aktiven Sitzung.

### **Ändern der Modellmaße**

In den folgenden Beispielen werden verschiedene Möglichkeiten zum Ändern von Modellmaßen erläutert:

- Ändern eines bestehenden Maßes.
- Importieren eines weiteren Standardmaßes aus dem Modell.
- Exportieren eines weiteren Standardmaßes in das Modell.

#### **Ändern eines bestehenden Maßes**

- ▶ Öffnen Sie die Vorlage und wählen Sie das Arbeitsblatt “Konfiguration” aus.
- ▶ Bearbeiten Sie einen beliebigen Wert für Name oder Beschreibung, indem Sie ihn markieren und überschreiben.

Beachten Sie: Um ein Maß zu ändern – beispielsweise, um den Benutzer statt zur Eingabe der Häufigkeit zur Eingabe der Wahrscheinlichkeit aufzufordern –, müssen Sie lediglich den Namen und die Beschreibung unter Maße aus Modell– ändern. Die Änderung wird dann im Modell angezeigt und der Benutzer kann das entsprechende Maß für die Zuordnung auswählen.

### **Importieren eines weiteren Standardmaßes aus dem Modell**

- ▶ Öffnen Sie die Vorlage und wählen Sie das Arbeitsblatt “Konfiguration” aus.
  - ▶ Wählen Sie die folgenden Befehle aus den Menüs aus:  
Werkzeuge > Schutz > Schutz des Arbeitsblatts aufheben
  - ▶ Wählen Sie Zelle A5, die gelb schattiert ist, und das Wort End enthält.
  - ▶ Wählen Sie die folgenden Befehle aus den Menüs aus:  
Einfügen > Zeilen
  - ▶ Geben Sie Name und Beschreibung des neuen Maßes ein. Beispielsweise könnten Sie Error (Fehler) und Error associated with segment (Zum Segment gehöriger Fehler) eingeben.
  - ▶ Geben Sie in Zelle C5 die Formel =COLUMN('Model Measures'!N3) ein.
  - ▶ Geben Sie in Zelle D5 die Formel =ROW('Model Measures'!N3)+1 ein.
- Durch diese Formeln wird das neue Maß in Spalte N des Arbeitsblatts “Model Measures” angezeigt, die derzeit leer ist.
- ▶ Wählen Sie die folgenden Befehle aus den Menüs aus:  
Werkzeuge > Schutz > Arbeitsblatt schützen
  - ▶ Klicken Sie auf OK.
  - ▶ Vergewissern Sie sich auf dem Arbeitsblatt “Model Measures”, dass Zelle N3 Error als Titel für die neue Spalte aufweist.
  - ▶ Wählen Sie die gesamte Spalte N aus.
  - ▶ Wählen Sie die folgenden Befehle aus den Menüs aus:  
Format > Zellen
  - ▶ Standardmäßig haben alle Zellen die Zahlenkategorie Allgemein. Klicken Sie auf Prozentsatz, um die Darstellungsart der Zahlen zu ändern. Dadurch können Sie Ihre Zahlen besser in Excel überprüfen. Außerdem können Sie dadurch die Daten auf andere Weise nutzen, beispielsweise als Ausgabe für ein Diagramm.
  - ▶ Klicken Sie auf OK.
  - ▶ Speichern Sie das Arbeitsblatt als Excel 2003-Vorlage, mit einem eindeutigen Namen und der Dateierweiterung *.xlt*. Um die neue Vorlage leichter wieder auffinden zu können, sollten Sie sie am Speicherort der vorkonfigurierten Vorlage auf Ihrem lokalen System oder Netzwerkdateisystem speichern.

**Exportieren eines weiteren Standardmaßes in das Modell**

- ▶ Öffnen Sie die Vorlage, zu der Sie die Spalte "Error" (Fehler) im vorherigen Beispiel hinzugefügt haben; wählen Sie das Arbeitsblatt "Konfiguration" aus.
- ▶ Wählen Sie die folgenden Befehle aus den Menüs aus:  
Werkzeuge > Schutz > Schutz des Arbeitsblatts aufheben
- ▶ Wählen Sie Zelle A14, die gelb schattiert ist, und das Wort End enthält.
- ▶ Wählen Sie die folgenden Befehle aus den Menüs aus:  
Einfügen > Zeilen
- ▶ Geben Sie Name und Beschreibung des neuen Maßes ein. Beispielsweise könnten Sie Scaled Error (Skalierter Fehler) und Scaling applied to error from Excel (Skalierung auf Fehler aus Excel angewendet) eingeben.
- ▶ Geben Sie in Zelle C14 die Formel =COLUMN('Model Measures'!O3) ein.
- ▶ Geben Sie in Zelle D14 die Formel =ROW('Model Measures'!O3)+1 ein.  
  
Diese Formeln geben an, dass die Spalte O das neue Maß für das Modell liefert.
- ▶ Wählen Sie das Arbeitsblatt "Einstellungen" aus.
- ▶ Geben Sie in Zelle A17 die Beschreibung ' - Scaled Error (Skalierter Fehler) ein.
- ▶ Geben Sie in Zelle B17 den Skalierungsfaktor 10 ein.
- ▶ Geben Sie auf dem Arbeitsblatt "Model Measures" (Modellmaße) die Beschreibung Scaled Error (Skalierter Fehler) in Zelle O3 als Titel für die neue Spalte ein.
- ▶ Geben Sie in Zelle O4 die Formel =N4\*Settings!\$B\$17 ein.
- ▶ Wählen Sie die Ecke von Zelle O4 aus und ziehen Sie sie nach unten auf Zelle O22, um die Formel in jede Zelle zu kopieren.
- ▶ Wählen Sie die folgenden Befehle aus den Menüs aus:  
Werkzeuge > Schutz > Arbeitsblatt schützen
- ▶ Klicken Sie auf OK.
- ▶ Speichern Sie das Arbeitsblatt als Excel 2003-Vorlage, mit einem eindeutigen Namen und der Dateierweiterung *.xlt*. Um die neue Vorlage leichter wieder auffinden zu können, sollten Sie sie am Speicherort der vorkonfigurierten Vorlage auf Ihrem lokalen System oder Netzwerkdateisystem speichern.

Wenn Sie über diese Vorlage eine Verbindung mit Excel herstellen, ist der Fehlerwert als neues benutzerdefiniertes Maß verfügbar.

## Visualisieren von Modellen

Die Auswirkung eines Modells wird am deutlichsten, wenn es visuell dargestellt wird. Mithilfe eines Gewinn diagrams erhalten Sie einen wertvollen täglichen Einblick in den betriebswirtschaftlichen und technischen Nutzen Ihres Modells, indem Sie die Effekte mehrerer Alternativen in Echtzeit untersuchen. Im Abschnitt [Gewinndiagramm](#) wird der Vorteil eines Modells im Vergleich zu einer zufälligen Entscheidungsfindung besprochen. Hier erfolgt außerdem der direkte Vergleich mehrerer Diagramme beim Vorhandensein alternativer Modelle.

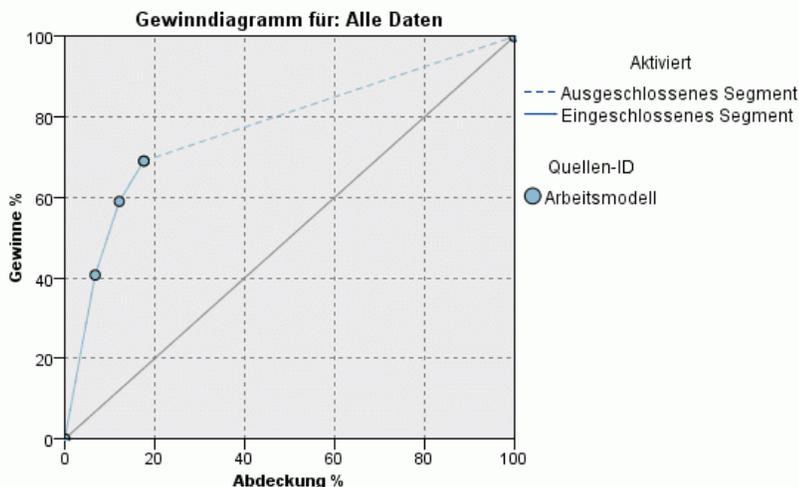
### Gewinndiagramm

Das Gewinn diagram bildet die Werte der Tabellenspalte *Gewinn %* ab. Gewinne sind als der Anteil der in jedem Inkrement enthaltenen Treffer im Verhältnis zur Gesamtzahl der im Baum enthaltenen Treffer definiert. Dabei kommt folgende Gleichung zum Einsatz:

$$(\text{Treffer im Inkrement} / \text{Gesamtzahl Treffer}) \times 100 \%$$

Gewinndiagramme illustrieren, wie weit Sie das Netz auswerfen müssen, um einen bestimmten Prozentsatz aller im Baum enthaltenen Treffer zu erzielen. Die diagonale Linie bildet die für die gesamte Stichprobe erwarteten Treffer ab, wenn das Modell nicht verwendet wird. In diesem Fall ist die Trefferrate konstant, da die Wahrscheinlichkeit eines Treffers für alle Personen gleich ist. Um das Ergebnis zu verdoppeln, müssen Sie doppelt so viele Personen ansprechen. Die gekrümmte Linie zeigt an, wie weit Sie Ihre Treffer verbessern können, wenn Sie nur die einschließen, deren Prozentsatz hinsichtlich des Gewinns höher ausfällt. Wenn Sie beispielsweise die obersten 50 % einschließen, erhalten Sie über 70 % der positiven Treffer. Je steiler die Kurve, desto höher ist der Gewinn.

Abbildung 9-20  
Registerkarte "Gewinne"



#### So zeigen Sie ein Gewinn diagram an:

- Öffnen Sie einen Stream, der einen Entscheidungslistenknoten enthält, und starten Sie von diesem Knoten aus eine interaktive Sitzung.

- Klicken Sie auf die Registerkarte Gewinne. Je nachdem, welche Partitionen angegeben sind, werden ein oder zwei Diagramme angezeigt (zwei Diagramme werden angezeigt, wenn für die Modellmaße beispielsweise die Trainings- und die Testpartition definiert sind).

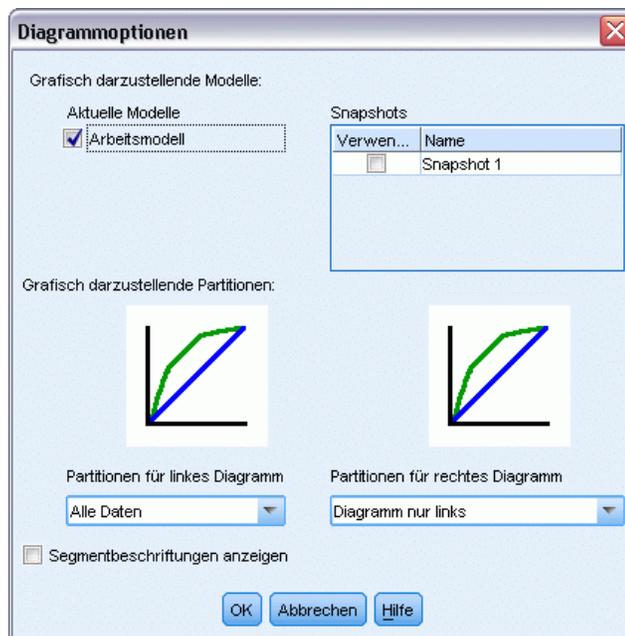
Die Diagramme werden standardmäßig als Segmente angezeigt. Sie können die Anzeige der Diagramme in Quantile umschalten, indem Sie Quantile auswählen und dann im Dropdown-Menü das entsprechende Quantilverfahren auswählen.

*Hinweis:* Informationen zur Arbeit mit Diagrammen finden Sie unter [Bearbeiten von Visualisierungen](#).

## Diagrammoptionen

Die Funktion “Diagrammoptionen” bietet Optionen zur Auswahl der Modelle und Snapshots, die als Diagramm dargestellt werden, welche Partitionen angezeigt werden und ob Segmentbezeichnungen angezeigt werden oder nicht.

Abbildung 9-21  
Dialogfeld “Diagrammoptionen”



### Modelle für das Diagramm

**Aktuelle Modelle.** Hier können Sie auswählen, welche Modelle als Diagramm dargestellt werden. Sie können das Arbeitsmodell oder eines der erstellten Snapshot-Modelle auswählen.

### Partitionen für das Diagramm

**Partitionen für linkes Diagramm.** Die Dropdown-Liste enthält Optionen für die Anzeige aller definierten Partitionen oder aller Daten.

**Partitionen für rechtes Diagramm.** Die Dropdown-Liste enthält Optionen für die Anzeige aller definierten Partitionen, aller Daten oder nur des linken Diagramms. Wenn Diagramm nur links ausgewählt ist, wird nur das linke Diagramm angezeigt.

**Segmentbeschriftungen anzeigen.** Wenn diese Option aktiviert ist, werden die Segmentbeschriftungen in den Diagrammen angezeigt.

# Statistische Modelle

Statistische Modelle verwenden mathematische Gleichungen, um Informationen zu kodieren, die aus den Daten extrahiert wurden. Mitunter können mithilfe statistischer Modellierungstechniken geeignete Modelle sehr schnell bereitgestellt werden. Selbst bei Problemen, bei denen flexiblere Techniken zum Maschinernen (z. B. neuronale Netze) letztendlich bessere Ergebnisse liefern, können Sie statistische Modelle als Basisvorhersagemodelle einsetzen, um die Leistung fortgeschrittener Techniken zu beurteilen.

Die folgenden Knoten für die statistische Modellierung stehen zur Verfügung.



Bei linearen Regressionsmodellen wird ein stetiges Ziel auf der Basis linearer Beziehungen zwischen dem Ziel und einem oder mehreren Prädiktoren vorhergesagt. [Für weitere Informationen siehe Thema Lineare Modelle auf S. 267.](#)



Die logistische Regression ist ein statistisches Verfahren zur Klassifizierung von Datensätzen auf der Grundlage der Werte von Eingabefeldern. Sie ist analog zur linearen Regression, außer dass statt eines numerischen Bereichs ein kategoriales Zielfeld verwendet wird. [Für weitere Informationen siehe Thema Logistikknoten auf S. 287.](#)



Der Faktor/PCA-Knoten bietet leistungsstarke Datenreduktionsverfahren zur Verringerung der Komplexität der Daten. Die Hauptkomponentenanalyse (PCA) findet lineare Kombinationen der Eingabefelder, die die Varianz im gesamten Set der Felder am besten erfassen, wenn die Komponenten orthogonal (senkrecht) zueinander sind. Mit der Faktorenanalyse wird versucht, die zugrunde liegenden Faktoren zu bestimmen, die die Korrelationsmuster innerhalb eines Sets beobachteter Felder erklären. Bei beiden Ansätzen besteht das Ziel darin, eine kleinere Zahl abgeleiteter Felder zu finden, mit denen die Informationen in der ursprünglichen Menge der Felder effektiv zusammengefasst werden können. [Für weitere Informationen siehe Thema Faktor/PCA-Knoten auf S. 311.](#)



Bei der Diskriminanzanalyse werden strengere Annahmen als bei der logistischen Regression verwendet, sie kann jedoch eine wertvolle Alternative oder Ergänzung zu einer logistischen Regressionsanalyse sein, wenn diese Annahmen erfüllt sind. [Für weitere Informationen siehe Thema Diskriminanzknoten auf S. 319.](#)



Das verallgemeinerte lineare Modell erweitert das allgemeine lineare Modell so, dass die abhängige Variable über eine angegebene Verknüpfungsfunktion in linearem Zusammenhang zu den Faktoren und Kovariaten steht. Außerdem ist es mit diesem Modell möglich, dass die abhängige Variable eine von der Normalverteilung abweichende Verteilung aufweist. Es deckt die Funktionen einer großen Bandbreite an Statistikmodellen ab, darunter lineare Regression, logistische Regression, loglineare Modelle für Häufigkeitsdaten und Überlebensmodelle mit Intervallzensurierung. [Für weitere Informationen siehe Thema GenLin-Knoten auf S. 328.](#)

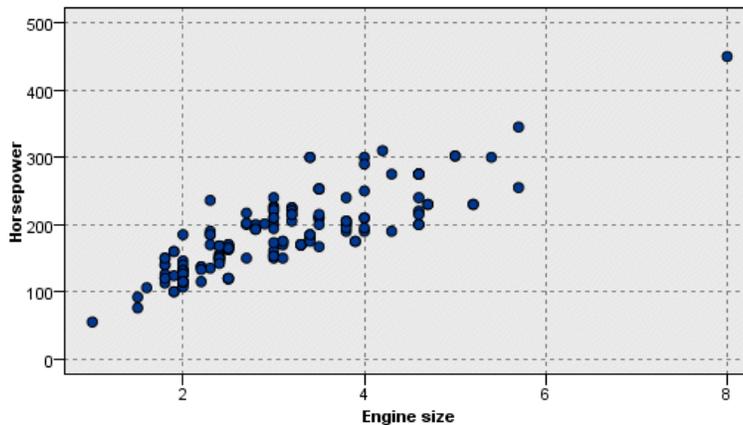


Der Knoten vom Typ “Cox-Regression” ermöglicht Ihnen auch bei zensierten Datensätzen die Erstellung eines Überlebensmodells für Daten über die Zeit bis zum Eintreten des Ereignisses. Das Modell erstellt eine Überlebensfunktion, die die Wahrscheinlichkeit vorhersagt, dass das untersuchte Ereignis für bestimmte Werte der Eingabevariablen zu einem bestimmten Zeitpunkt ( $t$ ) eingetreten ist. [Für weitere Informationen siehe Thema Cox-Knoten auf S. 342.](#)

## Linearknoten

Die lineare Regression ist ein verbreitetes statistisches Verfahren zur Klassifizierung von Datensätzen auf der Grundlage der Werte von numerische Eingabefeldern. Die lineare Regression entspricht einer geraden Linie oder Fläche, die die Diskrepanzen zwischen den vorhergesagten und den tatsächlichen Werten minimiert.

Abbildung 10-1  
Einfaches lineares Regressionsdiagramm



**Anforderungen.** In linearen Regressionsmodellen können nur numerische Felder verwendet werden. Es werden genau ein Zielfeld (mit der Rolle *Ziel*) und mindestens ein Prädiktor (mit der Rolle *Eingabe*) benötigt. Felder mit der Rolle *Beides* oder *Keine* werden, wie auch nicht numerische Felder, ignoriert. (Nicht numerische Felder können, falls erforderlich, mithilfe eines Ableitungsknotens neu kodiert werden.) [Für weitere Informationen siehe Thema Umkodieren von Werten mit dem Ableitungsknoten in Kapitel 4 in IBM SPSS Modeler 14.2- Quellen-, Prozess- und Ausgabeknoten.](#))

**Stärken.** Lineare Regressionsmodelle sind relativ einfach und bieten eine leicht zu interpretierende mathematische Formel für das Generieren von Prognosen. Da die lineare Regressionsmodellierung ein seit langem etabliertes statistisches Verfahren ist, liegen umfassende Kenntnisse über die Eigenschaften dieser Modelle vor. Lineare Regressionsmodelle lassen sich üblicherweise sehr schnell trainieren. Der Linearknoten bietet Methoden für die automatische Feldauswahl zum Entfernen nicht signifikanter Eingabefelder aus der Gleichung.

*Hinweis:* In Fällen, bei denen das Zielfeld keinen kontinuierlichen Bereich darstellt, sondern kategorial ist, wie beispielsweise *ja/nein* oder *Abwanderung/Keine Abwanderung*, kann die logistische Regression als Alternative verwendet werden. Die logistische Regression bietet

außerdem Unterstützung für nicht numerische Eingaben, sodass eine Umkodierung dieser Felder nicht mehr erforderlich ist. [Für weitere Informationen siehe Thema Logistikknoten auf S. 287.](#)

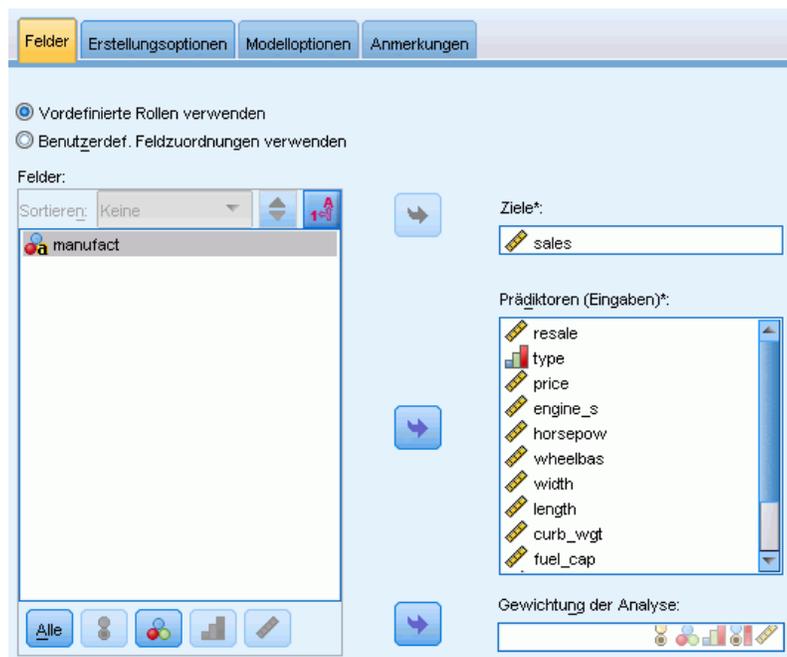
## Lineare Modelle

Bei linearen Modellen wird ein stetiges Ziel auf der Basis linearer Beziehungen zwischen dem Ziel und einem oder mehreren Prädiktoren vorhergesagt.

Lineare Modelle sind relativ einfach und bieten eine leicht zu interpretierende mathematische Formel für das Scoring. Die Eigenschaften dieser Modelle sind umfassend bekannt und sie lassen sich üblicherweise sehr schnell im Vergleich zu anderen Modelltypen (beispielsweise neuronale Netze oder Entscheidungsbäume) im selben Datenblatt (Daten-Set) erstellen.

**Beispiel.** Eine Versicherungsgesellschaft mit beschränkten Ressourcen für die Untersuchung der Versicherungsansprüche von Hauseigentümern möchte ein Modell zur Schätzung der Kosten von Ansprüchen erstellen. Durch die Bereitstellung dieses Modells in Servicecentern können die Mitarbeiter beim telefonischen Gespräch mit einem Kunden die Informationen zum Anspruch eingeben und auf Basis von Daten aus der Vergangenheit sofort die zu erwartenden Kosten des Anspruchs abrufen.

Abbildung 10-2  
Registerkarte "Felder"



**Feldanforderungen.** Es müssen ein Ziel und mindestens eine Eingabe vorhanden sein. Standardmäßig werden Felder mit den vordefinierten Rollen "Beides" oder "Keine" nicht verwendet. Das Ziel muss stetig (Skala) sein. Für die Prädiktoren (Eingaben) gibt es keine Einschränkungen bezüglich der Messniveaus. Als Faktoren im Modell werden kategoriale Felder (Flag, nominal und ordinal) und stetige Felder als Kovariaten verwendet. [Für weitere Informationen siehe Thema Feldoptionen der Modellierungsknoten in Kapitel 3 auf S. 36.](#)

## Ziele

Abbildung 10-4  
Ziele - Einstellungen

Wählen Sie ein Element aus:

- Ziele
- Grundlagen
- Modellauswahl
- Ensembles
- Erweitert

Wie möchten Sie vorgehen?

Neues Modell aufbauen  Training des bestehenden Modells fortsetzen

Was ist Ihr Hauptziel?

Standardmodell erstellen

Modellgenauigkeit erweitern (Verstärkung)

Modellstabilität erweitern (Bagging)

Modell für sehr umfangreiche Datensets erstellen (erfordert Server)

Beschreibung

Erstellt ein einzelnes Standardmodell zur Erklärung der Beziehungen zwischen Feldern. Standardmodelle lassen sich einfacher interpretieren und ermöglichen schnelleres Scoring als verstärkte, geboostete, verpackte oder umfangreiche Datenset-Ensembles.

### Was möchten Sie tun?

- **Neues Modell aufbauen.** Ein vollständig neues Modell aufbauen. Dies ist die übliche Wirkungsweise des Knotens.
- **Training eines bestehenden Modells fortsetzen.** Das Training wird mit dem letzten vom Knoten erfolgreich erstellten Modell fortgesetzt. Dies ermöglicht die Aktualisierung eines bestehenden Modells, ohne dass ein Zugriff auf die ursprünglichen Daten erforderlich ist, und kann zu einer wesentlich schnelleren Leistung führen, da ausschließlich die neuen bzw. aktualisierten Datensätze in den Stream eingespeist werden. Details zum vorherigen Modell werden zusammen mit dem Modellierungsknoten gespeichert, wodurch diese Option auch dann verwendet werden kann, wenn das vorherige Modell-Nugget nicht mehr im Stream oder in der Modellpalette verfügbar ist.

*Anmerkung:* Wenn diese Option aktiviert ist, werden alle anderen Steuerungsfunktionen auf den Registerkarten “Felder” und “Erstellungsoptionen” deaktiviert.

### Wie lautet Ihr Hauptziel?

- **Erstellen eines Standardmodells.** Bei dem Verfahren wird ein einziges Modell erstellt, um das Ziel unter Verwendung der Prädiktoren vorherzusagen. In der Regel gilt, dass Standardmodelle einfacher interpretiert und schneller gescort werden können, als verbesserte, verstärkte oder große Daten-Set-Ensembles.
- **Modellgenauigkeit erhöhen (verbessern).** Bei dem Verfahren wird unter Einsatz der Verbesserung ein Ensemble-Modell erstellt. Dabei wird eine Modellsequenz erzeugt, um genauere Vorhersagen zu erhalten. Möglicherweise nimmt das Erstellen und Scoring bei Ensembles mehr Zeit in Anspruch als bei Standardmodellen.

Durch Verbesserung wird eine Reihe von “Komponentenmodellen” erstellt, von denen jede einzelne Komponente auf dem gesamten Daten-Set beruht. Vor dem Erstellen jedes aufeinander folgenden Komponentenmodells werden die Datensätze basierend auf den Residuen des vorangegangenen Komponentenmodells gewichtet. Größere Residuen erhalten eine höhere Analysegewichtung, sodass beim nächsten Komponentenmodell das Augenmerk auf einer hochwertigen Vorhersage dieser Datensätze liegt. Zusammen bilden diese Komponentenmodelle ein Ensemble-Modell. Das Ensemble-Modell scort basierend auf einer Kombinationsregel neue Datensätze. Die verfügbaren Regeln hängen vom Messniveau des Ziels ab.

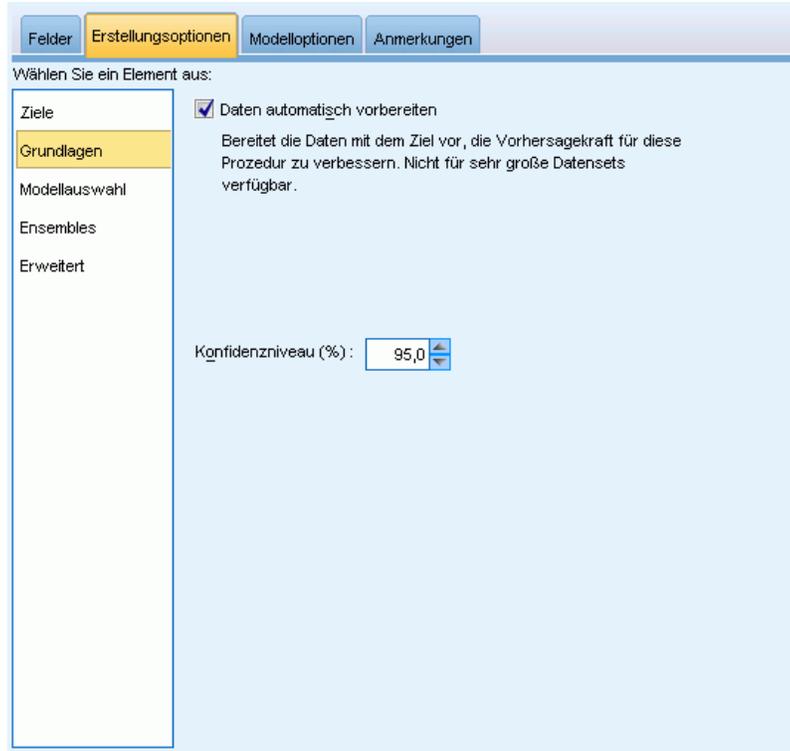
- **Modellstabilität steigern (verstärken).** Bei dem Verfahren wird unter Einsatz der Verstärkung (Bootstrap-Aggregation) ein Ensemble-Modell erstellt. Dabei werden mehrere Modelle erzeugt, um zuverlässigere Vorhersagen zu erhalten. Möglicherweise nimmt das Erstellen und Scoring bei Ensembles mehr Zeit in Anspruch als bei Standardmodellen.

Bei der Bootstrap-Aggregation (Verstärkung) werden Reproduktionen des Trainings-Daten-Set erstellt, indem von der Ersetzung aus dem ursprünglichen Daten-Set Stichproben genommen werden. Dadurch werden Bootstrap-Stichproben mit der gleichen Größe wie beim ursprünglichen Daten-Set erstellt. Dann wird von jeder Reproduktion ein “Komponentenmodell” erstellt. Zusammen bilden diese Komponentenmodelle ein Ensemble-Modell. Das Ensemble-Modell scort basierend auf einer Kombinationsregel neue Datensätze. Die verfügbaren Regeln hängen vom Messniveau des Ziels ab.

- **Modell für sehr große Daten-Sets erstellen (IBM® SPSS® Modeler Server erforderlich).** Bei dieser Methode wird ein Ensemble-Modell durch Aufteilen des Daten-Sets in separate Datenblöcke erstellt. Diese Option ist empfehlenswert, wenn Ihr Daten-Set zu groß für die Erstellung eines der oben erwähnten Modelle oder die inkrementelle Modellerstellung ist. Unter Umständen kann das Modell mit dieser Option schneller als ein Standardmodell erstellt werden, das Scoring dauert jedoch evtl. länger als bei einem Standardmodell. Diese Option erfordert eine SPSS Modeler Serververbindung.

## Grundeinstellungen

Abbildung 10-5  
Grundeinstellungen



**Automatische Datenaufbereitung.** Mit dieser Option kann die Prozedur das Ziel und die Prädiktoren intern transformieren, um die Vorhersagekraft des Modells zu maximieren. Alle Transformationen werden mit dem Modell gespeichert und auf neue Daten zum Scoring angewendet. Die Originalversionen der transformierten Felder werden vom Modell ausgeschlossen. Standardmäßig wird folgende automatische Datenaufbereitung durchgeführt.

- **Verarbeitung von Datum und Zeit.** Jeder Datumsprädiktor wird in einen neuen stetigen Prädiktor transformiert, der die Zeit enthält, die seit einem Referenzdatum (1970-01-01) vergangen ist. Jeder Zeitprädiktor wird in einen neuen stetigen Prädiktor transformiert, der die Zeit enthält, die seit einer Referenzzeit (00:00:00) vergangen ist.
- **Messniveau anpassen.** Stetige Prädiktoren mit weniger als fünf distinkten Werten werden in ordinale Felder umgewandelt. Ordinale Prädiktoren mit mehr als zehn distinkten Werten werden in stetige Prädiktoren umgewandelt.
- **Ausreißer-Behandlung.** Werte stetiger Prädiktoren, die über einem Cutoff-Wert liegen (drei Standardabweichungen vom Mittelwert), werden auf den Cutoff-Wert gesetzt.

- **Behandlung fehlender Werte.** Fehlende Werte nominaler Prädiktoren werden durch den Modus der Trainingspartition ersetzt. Fehlende Werte ordinaler Prädiktoren werden durch den Median der Trainingspartition ersetzt. Fehlende Werte stetiger Prädiktoren werden durch den Mittelwert der Trainingspartition ersetzt.
- **Überwachte Zusammenführung.** Mit dieser Option erstellen Sie ein sparsameres Modell, indem die Anzahl der zu verarbeitenden Felder in Zusammenhang mit dem Ziel reduziert wird. Ähnliche Kategorien werden anhand der Beziehung zwischen der Eingabe und dem Ziel identifiziert. Kategorien, die sich nicht signifikant unterscheiden (d. h. einen p-Wert aufweisen, der größer als 0,1 ist), werden zusammengeführt. Hinweis: Wenn alle Kategorien zu einer verschmolzen werden, werden die originalen und abgeleiteten Versionen des Felds aus dem Modell ausgeschlossen, da sie als Einflussgrößen keinen Wert haben.

**Konfidenzniveau.** Dieses Konfidenzniveau wird zur Berechnung der Intervallschätzungen der Modellkoeffizienten in der Ansicht [Koeffizienten](#) verwendet. Geben Sie einen Wert größer 0 und kleiner 100 ein. Der Standardwert ist 95.

## Modellauswahl

Abbildung 10-6  
Modellauswahl - Einstellungen

The screenshot shows the 'Modellauswahl - Einstellungen' dialog box with the following settings:

- Tab: **Erstellungsoptionen**
- Wählen Sie ein Element aus: **Modellauswahl**
- Methode zur Modellauswahl: **Schrittweise vorwärts**
- Auswahl schrittweise vorwärts:
  - Kriterien für Aufnahme bzw. Ausschluss: **Informationskriterium (AICC)**
  - Effekte einschließen mit p-Werten kleiner als: **0,05**
  - Effekte ausschließen mit p-Werten größer als: **0,1**
  - Maximale Anzahl von Effekten im endgültigen Modell anpassen  
Maximale Anzahl an Effekten.:
  - Maximale Schrittzahl anpassen  
Maximale Schrittzahl:
- Auswahl der besten Untergruppen:
  - Kriterien für Aufnahme bzw. Ausschluss: **Informationskriterium (AICC)**

**Modellauswahlmethode.** Wählen Sie eine der Modellauswahlmethoden (Details unten) oder alle Prädiktoren aus, wodurch einfach alle verfügbaren Prädiktoren als Haupteffekt-Modellterme eingegeben werden. Standardmäßig wird Schrittweise vorwärts verwendet.

**Auswahl "Schrittweise vorwärts".** Diese Option beginnt ohne Effekte im Modell und nimmt jeweils einen Effekt auf bzw. schließt ihn aus, bis entsprechend den Kriterien bei "Schrittweise vorwärts" keine weiteren Vorgänge möglich sind.

- **Kriterien für Aufnahme/Ausschluss.** Diese Statistik wird zur Bestimmung verwendet, ob ein Effekt im Modell aufgenommen oder aus diesem ausgeschlossen werden soll. Das Informationskriterium (AICC) basiert auf der Wahrscheinlichkeit des Trainings-Sets für das Modell und wird zur Penalisierung übermäßig komplexer Modelle angepasst. Die F-Statistik beruht auf einem statistischen Test der Verbesserung des Modellfehlers. Korrigiertes R-Quadrat beruht auf der Anpassungsgüte des Trainings-Sets und wird zur Penalisierung übermäßig komplexer Modelle angepasst. Das Kriterium zur Verhinderung übermäßiger Anpassung (ASE) basiert auf der Anpassungsgüte (mittlere quadratische Abweichung; Average Squared Error, ASE) des Sets zur Verhinderung übermäßiger Anpassung. Das Set zur Verhinderung übermäßiger Anpassung ist eine zufällige Teilprobe von etwa 30 % des Original-Daten-Sets, die nicht zum Trainieren des Modells verwendet wird.

Wenn ein anderes Kriterium als F-Statistik gewählt wird, wird bei jedem Schritt der Effekt im Modell aufgenommen, der dem größten positiven Zuwachs des Kriteriums entspricht. Alle Effekte, die einer Abnahme des Kriteriums entsprechen, werden aus dem Modell ausgeschlossen.

Wenn F-Statistik als Kriterium gewählt wird, wird bei jedem Schritt der Effekt mit dem geringsten  $p$ -Wert kleiner als der festgelegte Schwellenwert, Einschließen von Effekten mit  $p$ -Werten kleiner als , in das Modell aufgenommen. Der Standardwert ist 0.05. Alle Effekte im Modell mit einem  $p$ -Wert größer als der festgelegte Schwellenwert, Entfernen von Effekten mit  $p$ -Werten größer als werden ausgeschlossen. Der Standardwert ist 0.10.

- **Anpassen der maximalen Anzahl an Effekten im endgültigen Modell.** Standardmäßig können alle verfügbaren Effekte in das Modell eingegeben werden. Wenn alternativ der schrittweise Algorithmus einen Schritt bei der festgelegten maximalen Anzahl an Effekten beendet, stoppt der Algorithmus beim aktuellen Effekt-Set.
- **Anpassen der maximalen Anzahl an Schritten.** Der schrittweise Algorithmus stoppt nach einer bestimmten Anzahl von Schritten. Standardmäßig ist das dreimal die Anzahl an verfügbaren Effekten. Alternativ kann eine positive Ganzzahl als maximale Anzahl an Schritten angegeben werden.

**Auswahl "Beste Untergruppen".** Diese Option überprüft "alle möglichen" Modelle oder zumindest eine größere Untergruppe der möglichen Modelle als "Schrittweise vorwärts", um die beste Möglichkeit entsprechend dem Kriterium "Beste Untergruppen" auszuwählen. Das Informationskriterium (AICC) basiert auf der Wahrscheinlichkeit des Trainings-Sets für das Modell und wird zur Penalisierung übermäßig komplexer Modelle angepasst. Korrigiertes R-Quadrat beruht auf der Anpassungsgüte des Trainings-Sets und wird zur Penalisierung übermäßig komplexer Modelle angepasst. Das Kriterium zur Verhinderung übermäßiger Anpassung (ASE) basiert auf der Anpassungsgüte (mittlere quadratische Abweichung; Average Squared Error, ASE) des Sets zur Verhinderung übermäßiger Anpassung. Das Set zur Verhinderung übermäßiger

Anpassung ist eine zufällige Teilprobe von etwa 30 % des Original-Daten-Sets, die nicht zum Trainieren des Modells verwendet wird.

Das Modell mit dem höchsten Wert für das Kriterium wird als das beste Modell ausgewählt.

*Anmerkung:* Die Auswahl “Beste Untergruppen” ist rechenintensiver als die Auswahl “Schrittweise vorwärts”. Wenn “Beste Untergruppen” zusammen mit “Verbesserung”, “Verstärkung” oder “Sehr große Daten-Sets” verwendet wird, kann das Erstellen deutlich länger dauern als das Erstellen eines Standard-Modells mithilfe der Auswahl “Schrittweise vorwärts”.

## Ensembles

Abbildung 10-7  
Ensemble-Einstellungen

The screenshot shows a software interface for setting ensemble options. It has four tabs: 'Felder', 'Erstellungsoptionen', 'Modelloptionen', and 'Anmerkungen'. The 'Erstellungsoptionen' tab is active. On the left, a sidebar lists categories: 'Ziele', 'Grundlagen', 'Modellauswahl', 'Ensembles', and 'Erweitert'. The 'Ensembles' category is highlighted. The main area contains a heading 'Wählen Sie ein Element aus:' followed by a list of options: 'Ziele', 'Grundlagen', 'Modellauswahl', 'Ensembles', and 'Erweitert'. Below this, there is a descriptive text: 'Diese Einstellungen bestimmen das Verhalten der Ensemble-Erstellung bei Boosting, Bagging oder wenn sehr umfangreiche Datensets in Zielen verlangt werden. Optionen, die nicht zutreffen, werden ignoriert.' There are two main sections: 'Bagging und sehr umfangreiche Daten-Sets' with a dropdown menu for 'Standard-Kombinationsregel für stetige Ziele' set to 'Mittelwert', and 'Boosting und Bagging' with a spin box for 'Anzahl der Komponentenmodelle für Boosting und/oder Bagging' set to 10.

Diese Einstellungen legen das Verhalten der Ensemblebildung fest, die erfolgt, wenn auf der Registerkarte “Ziele” die Option “Verbesserung”, “Verstärkung” oder “Sehr große Daten-Sets” ausgewählt ist. Optionen, die für das ausgewählte Ziel nicht gelten, werden ignoriert.

**Bagging und sehr umfangreiche Daten-Sets.** Beim Scoring eines Ensembles wird diese Regel angewendet, um die vorhergesagten Werte aus den Basismodellen für die Berechnung des Score-Werts für das Ensemble zu kombinieren.

- **Standard-Kombinierungsregel für stetige Ziele.** Ensemble-Vorhersagewerte für stetige Ziele können unter Verwendung des Mittelwerts oder Medians der Vorhersagewerte aus den Basismodellen kombiniert werden.

Hinweis: Wenn als Ziel die Verbesserung der Modellgenauigkeit ausgewählt wurde, wird die Auswahl zum Kombinieren der Regeln ignoriert. Bei der Verbesserung wird für das Scoring der kategorialen Ziele stets eine gewichtete Mehrheit verwendet und für das Scoring stetiger Ziele ein gewichteter Median.

**Verbesserung und Verstärkung.** Geben Sie die Anzahl der zu erstellenden Basismodelle an, wenn als Ziel die Verbesserung der Modellgenauigkeit oder -stabilität angegeben ist. Im Falle der Verstärkung ist das die Anzahl der Bootstrap-Stichproben. Muss eine positive ganze Zahl sein.

## Erweitert

Abbildung 10-8  
Erweiterte Einstellungen

Felder Erstellungsoptionen **Modelloptionen** Anmerkungen

Wählen Sie ein Element aus:

- Ziele
- Grundlagen
- Modellauswahl
- Ensembles
- Erweitert**

Ergebnisse replizieren

**Erzeugen**

Startwert für Zufallsgenerator:

**Ergebnisse reproduzieren.** Durch Einstellen eines Startwerts für den Zufallsgenerator können Analysen reproduziert werden. Der Zufallszahlengenerator wird verwendet, um zu wählen, welche Datensätze sich im Set zur Verhinderung übermäßiger Anpassung befinden. Geben Sie eine ganze Zahl ein oder klicken Sie auf Generieren. Dadurch wird eine pseudozufällige Ganzzahl zwischen 1 und 2147483647 (einschließlich) erzeugt. Der Standardwert ist 54752075.

## Modelloptionen

Abbildung 10-9  
Registerkarte "Modelloptionen"

Felder Erstellungsoptionen **Modelloptionen** Anmerkungen

Modellname:  Automatisch  Benutzerdefiniert

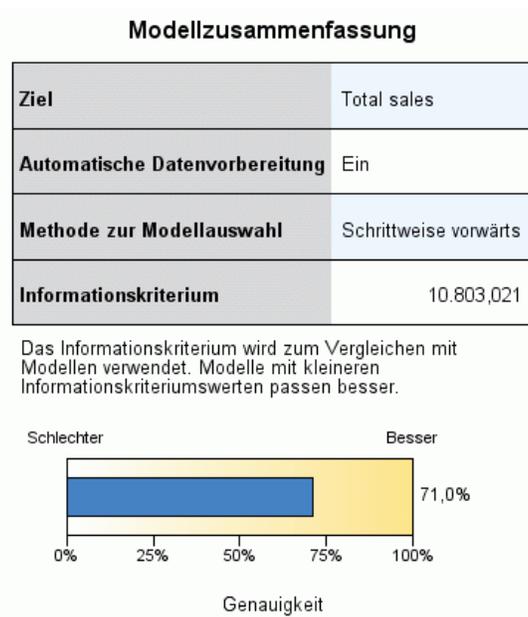
Vorhergesagter Wert ist immer für Scoring verfügbar. Keine anderen Scoring-Optionen stehen zur Verfügung.

**Modellname.** Sie können den Modellnamen automatisch basierend auf den Zielfeldern generieren, oder einen benutzerdefinierten Namen eingeben. Der automatisch generierte Name ist der Zielfeldname.

Hinweis: Der vorhergesagte Wert wird immer berechnet, wenn das Modell gescort wird. Der Name des neuen Felds ist der Name des Zielfelds mit einem vorangestellten  $\$L$ -. Bei einem Zielfeld mit der Bezeichnung *Umsatz* beispielsweise erhält das neue Feld den Namen  $\$L$ -Umsatz.

## Modellübersicht

Abbildung 10-10  
Ansicht Modellübersicht



Die Modellübersicht ist ein “Schnappschuss”, eine Übersicht auf einen Blick über das Modell und dessen Anpassungsgüte.

**Tabelle.** Mit der Tabelle werden einige Modelleinstellungen hoher Stufe wie etwa die folgenden ermittelt:

- Der Name des Ziels, das auf der Registerkarte **Felder** angegeben wurde;
- Ob die automatische Datenaufbereitung durchgeführt wurde, wie in den Einstellungen **Basis** angegeben wurde;
- Die Modellauswahlmethode und das Auswahlkriterium, die in den Einstellungen für die **Modellauswahl** angegeben wurden. Der Wert für das Auswahlkriterium des finalen Modells wird ebenfalls angezeigt und in einem Format dargestellt, bei dem kleinere Werte vorteilhafter sind.

**Diagramme.** Das Diagramm zeigt die Genauigkeit des endgültigen Modells an, das in einem größeren und besseren Format dargestellt wird. Der Wert ist  $100 \times$  der eingestellten  $R^2$  für das endgültige Modell.

## Automatische Datenaufbereitung

Abbildung 10-11  
Ansicht "Automatische Datenaufbereitung"

### Automatische Datenvorbereitung

Ziel: Total sales

Feld	Rolle	Durchgeführte Aktionen
Age category	Prädiktor	Zerstreute Kategorien für maximale Zuordnung mit Ziel zusammenführen
Primary keyword set	Prädiktor	Zerstreute Kategorien für maximale Zuordnung mit Ziel zusammenführen
Promotion	Prädiktor	Messniveau von kontinuierlich zu ordinal ändern
Secondary keyword set	Prädiktor	Zerstreute Kategorien für maximale Zuordnung mit Ziel zusammenführen

Wenn der Name des ursprünglichen Felds X ist, lautet der Name des transformierten Felds "X\_transformed". Das Originalfeld wird aus der Analyse ausgeschlossen und das transformierte Feld wird stattdessen eingeschlossen.

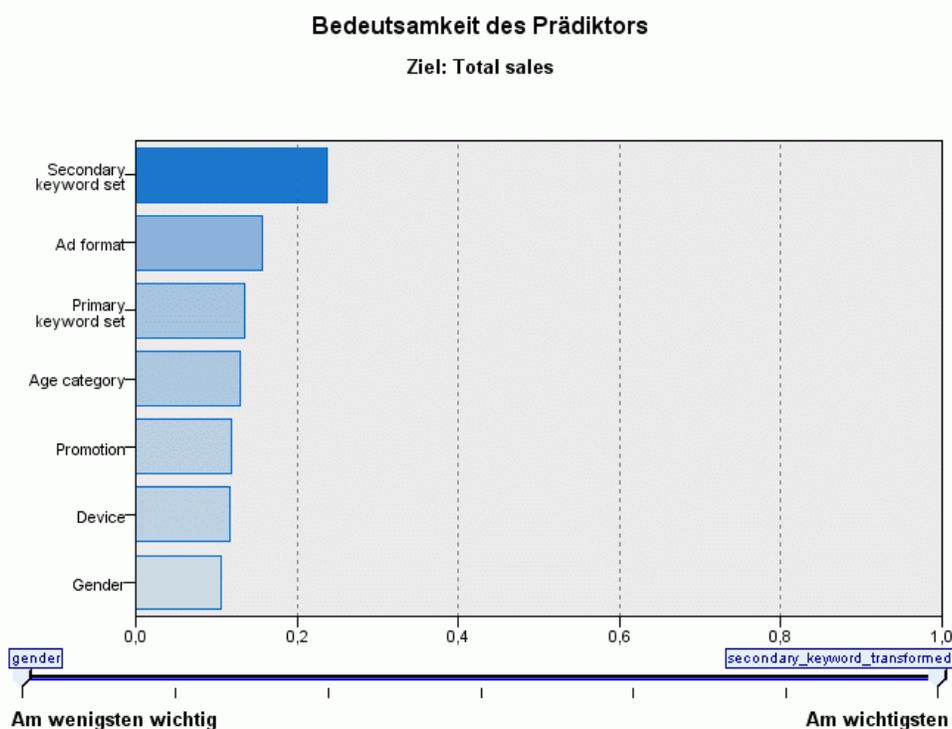
Diese Ansicht zeigt Informationen darüber an, welche Felder ausgeschlossen wurden und wie transformierte Felder im Schritt "automatische Datenaufbereitung" (ADP) abgeleitet wurden. Für jedes transformierte oder ausgeschlossene Feld listet die Tabelle den Feldnamen, die Rolle in der Analyse und die im ADP-Schritt vorgenommene Aktion auf. Die Felder werden in aufsteigender alphabetischer Reihenfolge der Feldnamen sortiert. Zu den möglichen Aktionen, die für die Felder durchgeführt wurden, zählen:

- Mit Ableitung der Dauer: Monate wird die verstrichene Zeit in Monaten aus den Werten in einem Feld mit Datumsangaben und dem aktuellen Datum des Systems berechnet.
- Mit Ableitung der Dauer: Stunden wird die verstrichene Zeit in Stunden aus den Werten in einem Feld mit Zeitangaben und der aktuellen Zeit des Systems berechnet.
- Mit Messniveau von stetig auf ordinal ändern werden stetige Felder mit weniger als fünf eindeutigen Werten in ordinale Felder umgewandelt.
- Mit Messniveau von ordinal auf stetig ändern werden ordinale Felder mit über zehn eindeutigen Werten in stetige Felder umgewandelt.
- Mit Ausreißer entfernen werden Werte stetiger Prädiktoren, die über einem Cutoff-Wert liegen (drei Standardabweichungen vom Mittelwert), auf den Cutoff-Wert gesetzt.
- Fehlende Werte ersetzen ersetzt fehlende Werte von nominalen Feldern durch den Modus, von ordinalen Feldern durch den Median und von stetigen Feldern durch den Mittelwert.

- Kategorien zur Maximierung des Zielzusammenhangs zusammenführen ermittelt “ähnliche” Prädiktorkategorien auf der Grundlage der Beziehung zwischen der Eingabe und dem Ziel. Kategorien, die sich nicht signifikant unterscheiden (d. h. einen  $p$ -Wert aufweisen, der größer als 0,05 ist), werden zusammengeführt.
- Konstanten Prädiktor ausschließen/nach Ausreißer-Behandlung/nach der Zusammenführung von Kategorien entfernt Prädiktoren, die einen einzelnen Wert aufweisen, möglicherweise nachdem andere ADP-Aktionen ausgeführt wurden.

## Bedeutsamkeit des Prädiktors

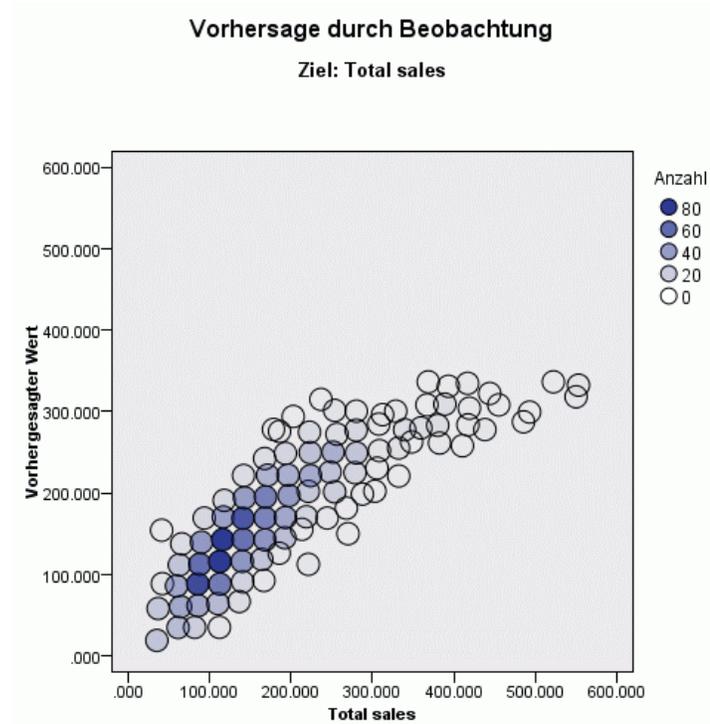
Abbildung 10-12  
Ansicht “Bedeutsamkeit des Prädiktors”



Normalerweise ist es sinnvoll, die Modellierungsbemühungen auf die wichtigsten Prädiktorfelder zu konzentrieren und diejenigen Prädiktoren zu verwerfen bzw. zu ignorieren, die die geringste Bedeutung haben. Dabei unterstützt Sie das Diagramm für die Bedeutsamkeit der Prädiktoren, da es die relative Bedeutsamkeit der einzelnen Prädiktoren für das Modell angibt. Da die Werte relativ sind, beträgt die Summe der Werte aller Prädiktoren im Diagramm 1,0. Die Bedeutsamkeit der Prädiktoren bezieht sich nicht auf die Genauigkeit des Modells. Sie bezieht sich lediglich auf die Bedeutsamkeit der einzelnen Prädiktoren für eine Vorhersage und nicht auf die Genauigkeit der Vorhersage.

## Vorhersage nach Beobachtung

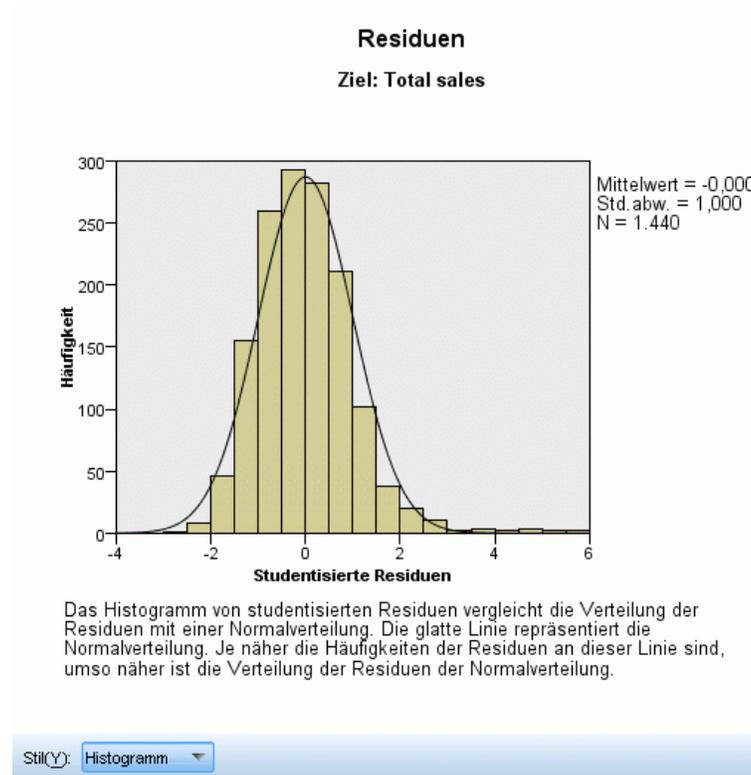
Abbildung 10-13  
Ansicht "Vorhersage nach Beobachtung"



Diese Ansicht zeigt ein Bin-Streudiagramm der vorhergesagten Werte auf der vertikalen Achse durch die beobachteten Werte auf der horizontalen Achse. Idealerweise sollten die Werte entlang einer 45-Grad-Linie liegen. In dieser Ansicht können Sie erkennen, ob bestimmte Datensätze vom Modell besonders schlecht vorhergesagt werden.

## Residuen

Abbildung 10-14  
Ansicht "Residuen," Histogrammstil



Diese Ansicht zeigt ein Diagnosediagramm von Modellresiduen.

**Diagrammstile.** Für die Diagramme sind verschiedene Anzeigestile verfügbar, auf die über die Dropdown-Liste Stil zugegriffen werden kann.

- **Histogramm.** Diese Ansicht zeigt ein Bin-Histogramm der studentisierten Residuen, das mit der normalen Verteilung überlagert ist. Lineare Modelle gehen davon aus, dass Residuen eine normale Verteilung aufweisen. Das Histogramm sollte sich also idealerweise einer nahezu glatten Linie annähern.
- **P-P-Diagramm.** Diese Ansicht zeigt ein Wahrscheinlichkeits-Wahrscheinlichkeits-Diagramm, bei dem die studentisierten Residuen mit einer normalen Verteilung verglichen werden. Wenn die Steigung der Diagrammpunkte weniger steil als die normale Linie ist, zeigen die Residuen eine größere Schwankung als eine normale Verteilung; ist die Steigung steiler, zeigen die Residuen weniger Schwankung als eine normale Verteilung. Wenn die Diagrammpunkte eine S-förmige Kurve aufweisen, ist die Verteilung der Residuen verzerrt.

## Ausreißer

Abbildung 10-15  
Ansicht "Ausreißer"

**Ausreißer**  
**Ziel: Total sales**

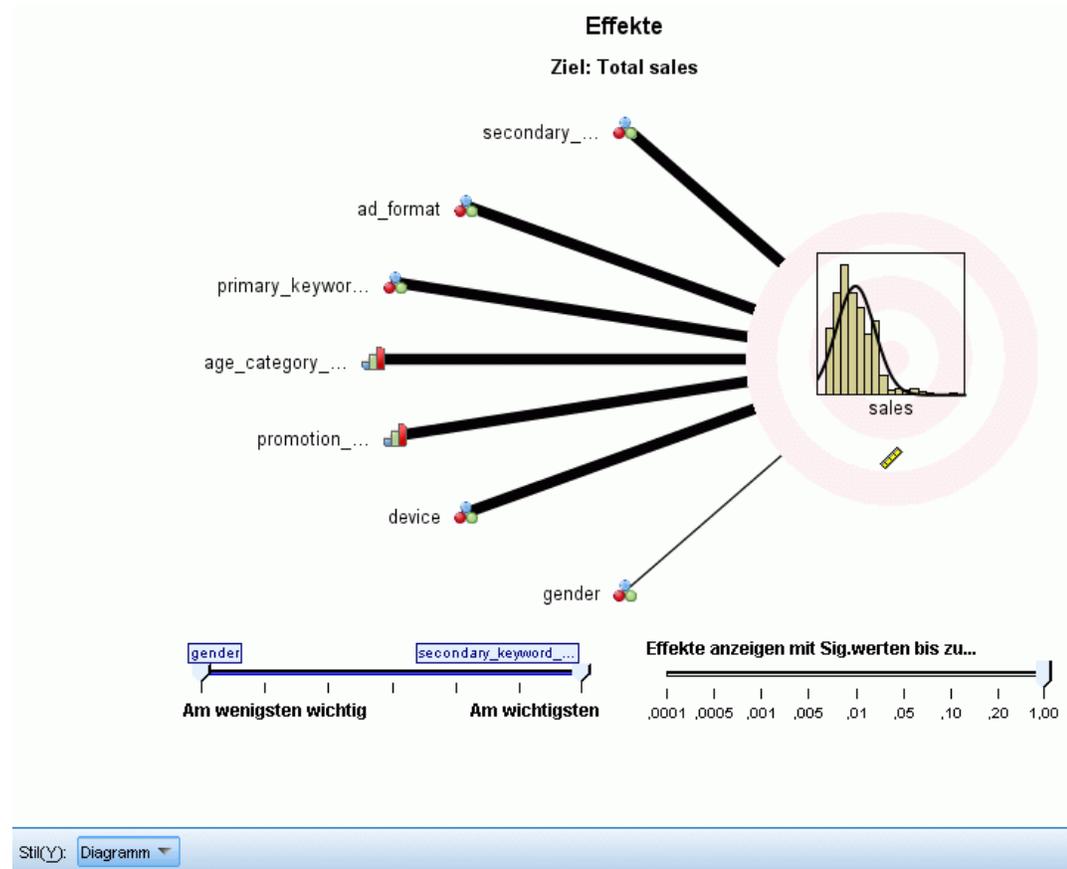
Total sales	Cook-Distanz
560.040	0,026
566.440	0,025
548.990	0,018
539.630	0,018
485.430	0,014
543.240	0,014

In dieser Tabelle sind Datensätze aufgelistet, die einen unverhältnismäßigen Einfluss auf das Modell ausüben. Außerdem werden die Datensatz-ID (sofern in der Registerkarte "Felder" angegeben), der Zielwert und die Cook-Distanz angezeigt. Die Cook-Distanz ist ein Maß dafür, wie stark sich die Residuen aller Datensätze ändern würden, wenn ein spezieller Datensatz von der Berechnung der Modellkoeffizienten ausgeschlossen würde. Ein großer Wert der Cook-Distanz zeigt an, dass der Ausschluss eines Datensatzes von der Berechnung die Koeffizienten substantiell verändert, und sollte daher als einflussreich betrachtet werden.

Einflussreiche Datensätze sollten sorgfältig untersucht werden, um zu entscheiden, ob ihnen bei der Schätzung des Modells weniger Gewicht gegeben werden kann, ob die extremen Werte auf einen akzeptablen Schwellenwert verringert werden können oder ob die einflussreichen Datensätze vollständig entfernt werden sollen.

## Effekte

Abbildung 10-16  
Ansicht "Effekte", Diagrammstil



Diese Ansicht zeigt die Größe der einzelnen Effekte im Modell.

**Stile.** Für die Diagramme sind verschiedene Anzeigestile verfügbar, auf die über die Dropdown-Liste Stil zugegriffen werden kann.

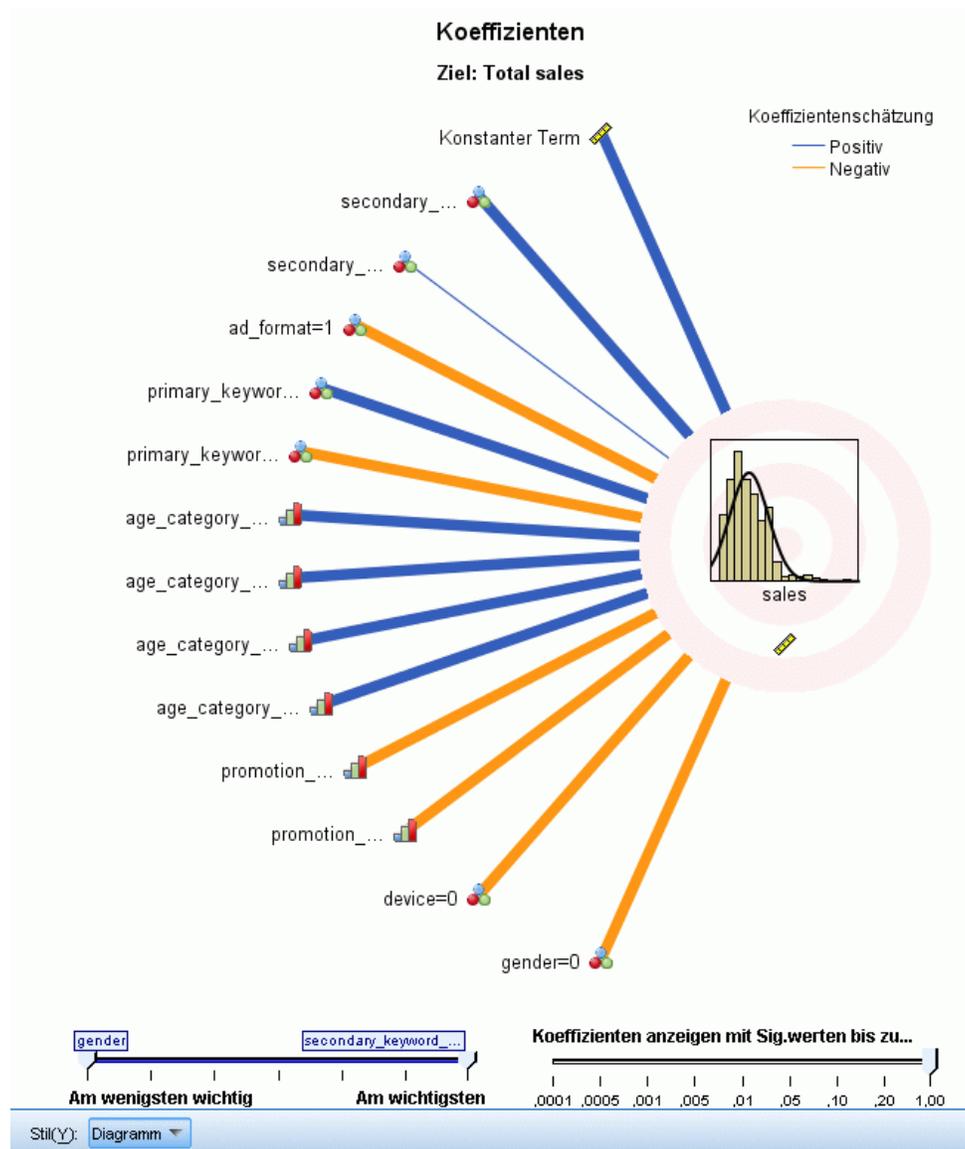
- **Diagramm.** In diesem Diagramm sind die Effekte von oben nach unten nach absteigender Bedeutsamkeit der Prädiktoren sortiert. Verbindungslinien im Diagramm sind basierend auf der Effektsignifikanz gewichtet, wobei eine größere Linienbreite signifikanteren Effekten entspricht (kleinere  $p$ -Werte). Wenn Sie den Mauszeiger über eine Verbindungslinie bewegen, wird eine QuickInfo mit dem  $p$ -Wert und der Bedeutung des Effekts angezeigt. Dies ist die Standardeinstellung.
- **Tabelle.** Diese Ansicht zeigt eine ANOVA-Tabelle für das Gesamtmodell und die einzelnen Modelleffekte. Die einzelnen Effekte sind von oben nach unten nach absteigender Bedeutsamkeit der Prädiktoren sortiert. Beachten Sie, dass die Tabelle standardmäßig minimiert ist, sodass nur die Ergebnisse des Gesamtmodells angezeigt werden. Klicken Sie in der Tabelle auf die Zelle für das korrigierte Modell, um die Ergebnisse für die einzelnen Modelleffekte anzuzeigen.

**Bedeutsamkeit des Prädiktors.** Für die Bedeutsamkeit des Prädiktors gibt es einen Schieberegler, mit dem eingestellt wird, welche Prädiktoren in der Ansicht gezeigt werden. Dadurch wird das Modell nicht verändert, doch Sie können sich ganz problemlos auf die wichtigsten Prädiktoren konzentrieren. Standardmäßig werden die zehn besten Effekte angezeigt.

**Signifikanz.** Mit dem Signifikanz-Schieberegler kann noch weiter angegeben werden, welche Effekte in der Anzeige dargestellt werden. Diese Einstellungen gehen über die Eingaben, die auf der Bedeutsamkeit der Prädiktoren beruhen, hinaus. Effekte, deren Signifikanzwerte größer als der Wert des Schiebereglers sind, werden ausgeblendet. Dadurch wird das Modell nicht verändert, doch Sie können sich ganz problemlos auf die wichtigsten Effekte konzentrieren. Standardmäßig ist der Wert 1,00 eingestellt, so dass keine Effekte basierend auf der Signifikanz herausgefiltert werden.

## Koeffizienten

Abbildung 10-17  
Ansicht "Koeffizienten", Diagrammstil



Diese Ansicht zeigt den Wert der einzelnen Koeffizienten im Modell. Hinweis: Faktoren (kategoriale Prädiktoren) sind innerhalb des Modells indikatorkodiert, sodass Faktoren, die **Effekte** enthalten, in der Regel mehrere zugehörige **Koeffizienten** aufweisen. Mit Ausnahme der Kategorie für den redundanten (Referenz-)Parameter erhält jede Kategorie einen solchen Koeffizienten.

**Stile.** Für die Diagramme sind verschiedene Anzeigestile verfügbar, auf die über die Dropdown-Liste Stil zugegriffen werden kann.

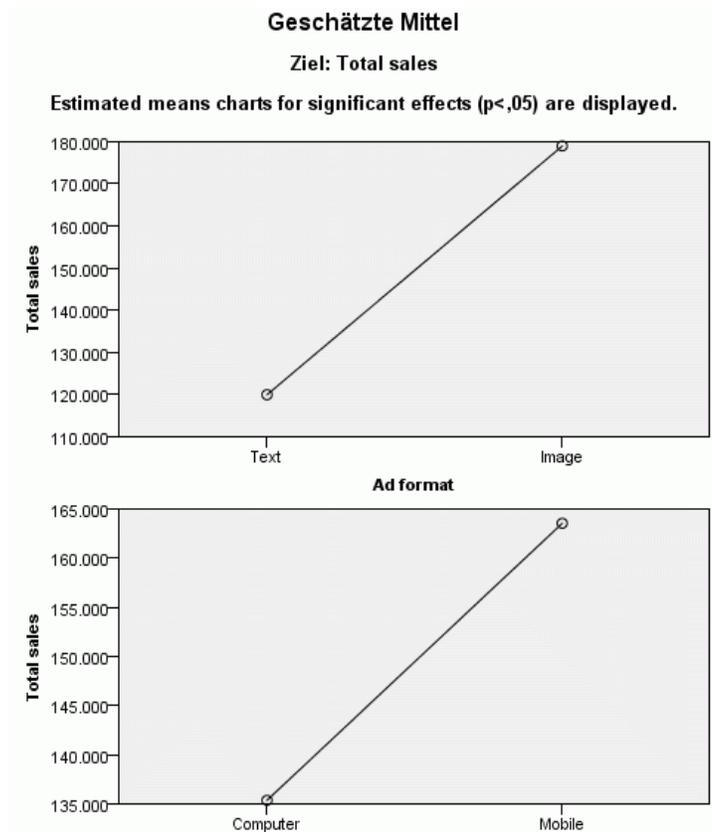
- **Diagramm.** In diesem Diagramm werden die konstanten Terme zuerst angezeigt, und dann die Effekte von oben nach unten nach absteigender Bedeutsamkeit der Prädiktoren sortiert. In Faktoren, die Effekte enthalten, werden die Koeffizienten in aufsteigender Reihenfolge der Datenwerte sortiert. Verbindungslinien im Diagramm sind basierend auf dem Vorzeichen des Koeffizienten farbig dargestellt (siehe Diagrammschlüssel) und auf der Grundlage der Koeffizientensignifikanz gewichtet, wobei eine größere Linienbreite signifikanteren Koeffizienten entspricht (kleinere  $p$ -Werte). Wenn Sie den Mauszeiger über eine Verbindungslinie bewegen, wird eine QuickInfo mit dem Wert des Koeffizienten, seinem  $p$ -Wert und der Bedeutung des Effekts angezeigt, mit dem der Parameter verbunden ist. Dies ist der Standardstil.
- **Tabelle.** Diese Tabelle zeigt die Werte, Signifikanztests und Konfidenzintervalle für die einzelnen Modellkoeffizienten. Nach dem konstanten Term sind die einzelnen Effekte von oben nach unten nach absteigender Bedeutsamkeit der Prädiktoren sortiert. In Faktoren, die Effekte enthalten, werden die Koeffizienten in aufsteigender Reihenfolge der Datenwerte sortiert. Beachten Sie, dass die Tabelle standardmäßig minimiert ist, sodass nur der Koeffizient, die Signifikanz und die Bedeutung der einzelnen Modellparameter angezeigt werden. Klicken Sie zum Anzeigen des Standardfehlers, der  $t$ -Statistik und des Konfidenzintervalls in der Tabelle auf die Zelle Koeffizient. Wenn Sie den Mauszeiger in der Tabelle über den Namen eines Modellparameters bewegen, wird eine QuickInfo mit dem Namen des Parameters, dem Effekt, mit dem der Parameter verbunden ist, und (für kategoriale Prädiktoren) den Wertelabels angezeigt, die mit dem Modellparameter verbunden sind. Dies kann besonders hilfreich sein, um die neuen Kategorien anzuzeigen, die erstellt werden, wenn bei der automatischen Datenaufbereitung ähnliche Kategorien eines kategorialen Prädiktors zusammengeführt werden.

**Bedeutsamkeit des Prädiktors.** Für die Bedeutsamkeit des Prädiktors gibt es einen Schieberegler, mit dem eingestellt wird, welche Prädiktoren in der Ansicht gezeigt werden. Dadurch wird das Modell nicht verändert, doch Sie können sich ganz problemlos auf die wichtigsten Prädiktoren konzentrieren. Standardmäßig werden die zehn besten Effekte angezeigt.

**Signifikanz.** Mit dem Signifikanz-Schieberegler kann noch weiter angegeben werden, welche Koeffizienten in der Anzeige dargestellt werden. Diese Einstellungen gehen über die Eingaben, die auf der Bedeutsamkeit der Prädiktoren beruhen, hinaus. Koeffizienten, deren Signifikanzwerte größer als der Wert des Schiebereglers sind, werden ausgeblendet. Dadurch wird das Modell nicht verändert, doch Sie können sich ganz problemlos auf die wichtigsten Koeffizienten konzentrieren. Standardmäßig ist der Wert 1,00 eingestellt, so dass keine Koeffizienten basierend auf der Signifikanz herausgefiltert werden.

## Geschätzte Mittel

Abbildung 10-18  
Ansicht "Geschätzte Mittel"



Diese Diagramme werden für signifikante Prädiktoren angezeigt. Das Diagramm zeigt den vom Modell geschätzten Zielwert auf der vertikalen Achse für jeden Prädiktorwert auf der horizontalen Achse, wobei alle anderen Prädiktoren konstant gehalten werden. Es gewährt eine nützliche Visualisierung der Effekte der einzelnen Prädiktorkoeffizienten auf dem Ziel.

*Anmerkung:* wenn keine Prädiktoren signifikant sind, werden keine geschätzten Mittel produziert.

## Modellerstellungsübersicht

Abbildung 10-19

Ansicht "Modellerstellungsübersicht"; Algorithmus "Schrittweise vorwärts"

### Übersicht über Modellerstellung

Ziel: Total sales

	Schritt						
	1	2	3	4	5	6	7
<b>Informationskriterium</b>	11.949,413	11.597,758	11.347,000	11.118,878	10.965,287	10.816,338	10.803,021
<b>secondary_keyword_transformed</b>	✓	✓	✓	✓	✓	✓	✓
<b>ad_format</b>		✓	✓	✓	✓	✓	✓
<b>primary_keyword_transformed</b>			✓	✓	✓	✓	✓
<b>Effekt age_category_transformed</b>				✓	✓	✓	✓
<b>promotion_transformed</b>					✓	✓	✓
<b>device</b>						✓	✓
<b>gender</b>							✓

Die Modellerstellungsmethode ist "Schrittweise vorwärts" mit dem "Informationskriterium".  
Ein Häkchen bedeutet, dass sich der Effekt bei diesem Schritt im Modell befindet.

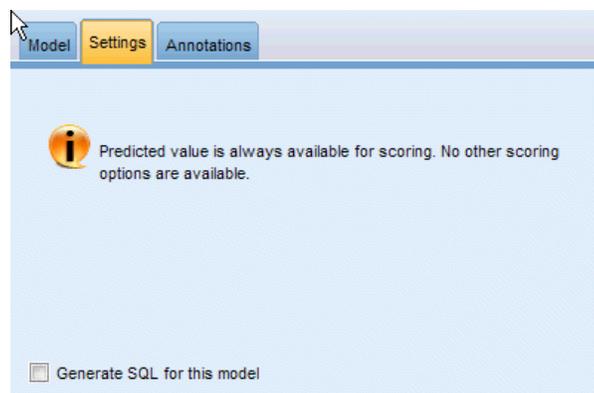
Wenn ein anderer Modellauswahlalgorithmus als Keiner in den Einstellungen "Modellauswahl" gewählt wird, werden einige Details zum Modellerstellungsprozess angegeben.

**Schrittweise vorwärts** Wenn der Auswahlalgorithmus "Schrittweise vorwärts" ist, zeigt die Tabelle die letzten zehn Schritte im schrittweisen Algorithmus an. Für jeden Schritt werden der Wert des Auswahlkriteriums und die Effekte im Modell an diesem Schritt angezeigt. Auf diese Weise bekommen Sie einen Eindruck davon, wie groß der Beitrag der einzelnen Schritte zum Modell ist. In jeder Spalte können Sie die Reihen so sortieren, dass Sie noch leichter erkennen können, welche Effekte sich bei einem bestimmten Schritt im Modell befinden.

**Beste Untergruppen.** Wenn der Auswahlalgorithmus "Beste Untergruppen" ist, zeigt die Tabelle die zehn besten Modelle an. Für jedes Modell werden der Wert des Auswahlkriteriums und die Effekte im Modell angezeigt. So erhalten Sie einen Eindruck der Stabilität der besten Modelle; wenn sie zu vielen ähnlichen Effekten mit wenigen Unterschieden neigen, können Sie sich auf das "Top"-Modell verlassen; wenn sie dagegen sehr unterschiedliche Effekte aufweisen, sind eventuell einige Effekte zu ähnlich und sollten kombiniert (oder entfernt) werden. In jeder Spalte können Sie die Reihen so sortieren, dass Sie noch leichter erkennen können, welche Effekte sich bei einem bestimmten Schritt im Modell befinden.

## Einstellungen

Abbildung 10-20  
Registerkarte "Einstellungen"



Hinweis: Der vorhergesagte Wert wird immer berechnet, wenn das Modell gescort wird. Der Name des neuen Felds ist der Name des Zielfelds mit einem vorangestellten  $SL$ -. Bei einem Zielfeld mit der Bezeichnung *Umsatz* beispielsweise erhält das neue Feld den Namen  $SL$ -Umsatz.

**SQL für dieses Modell generieren.** Bei Verwendung von Daten aus einer Datenbank kann SQL-Code per Pushback zur Ausführung an die Datenbank zurückübertragen werden. Dadurch kann bei vielen Operationen eine bessere Leistung erzielt werden. [Für weitere Informationen siehe Thema SQL-Optimierung in Kapitel 6 in IBM SPSS Modeler Server 14.2-Verwaltungs- und Leistungshandbuch.](#)

## Logistikknoten

**Logistische Regression**, auch als **nominale Regression** bekannt, ist ein statistisches Verfahren zur Klassifizierung von Datensätzen anhand der Werte der Eingabefelder. Sie ist analog zur linearen Regression, außer dass statt eines numerischen Zielfelds ein kategoriales verwendet wird. Es werden sowohl binomiale Modelle (für Ziele mit zwei diskreten Kategorien) als auch multinomiale Modelle (für Ziele mit mehr als zwei Kategorien) unterstützt.

Die logistische Regression funktioniert durch Erstellung einer Menge von Gleichungen, die die Werte der Eingabefelder mit den Wahrscheinlichkeiten in Relation setzen, die den einzelnen Ausgabefeldkategorien zugeordnet sind. Nach der Generierung des Modells kann es zur Schätzung der Wahrscheinlichkeiten für neue Daten verwendet werden. Für jeden Datensatz wird eine Wahrscheinlichkeit der Zugehörigkeit für jede mögliche Ausgabekategorie berechnet. Die Zielkategorie mit der höchsten Wahrscheinlichkeit wird als vorhergesagter Ausgabewert für den betreffenden Datensatz zugewiesen.

**Beispiel für ein binomiales Modell.** Ein Telekommunikationsanbieter ist besorgt über die Anzahl an Kunden, die er an Mitbewerber verliert. Mithilfe von Daten über die Servicenutzung können Sie ein binomiales Modell erstellen, mit dem Sie prognostizieren können, welche Kunden mit hoher Wahrscheinlichkeit zu einem anderen Anbieter wechseln, und Ihre Angebote entsprechend anpassen, um so viele Kunden wie möglich zu halten. Ein binomiales Modell wird verwendet, weil das Ziel zwei verschiedene Kategorien aufweist (hohe/geringe Wechselwahrscheinlichkeit). [Für weitere Informationen siehe Thema Kundenabwanderung bei](#)

Telekommunikationsunternehmen (binomiale logistische Regression) in Kapitel 14 in *IBM SPSS Modeler 14.2- Anwendungshandbuch*.

*Hinweis:* Bei binomialen Modellen (und nur dort) müssen Zeichenkettenfelder auf 8 Zeichen beschränkt sein. Längere Zeichenketten können, falls erforderlich, mithilfe eines Umkodierungsknotens neu kodiert werden. Für weitere Informationen siehe Thema Umkodierungsknoten in Kapitel 4 in *IBM SPSS Modeler 14.2- Quellen-, Prozess- und Ausgabeknoten*. Für weitere Informationen siehe Thema Reduzieren der Länge der Zeichenkette für die Eingabedaten (Umkodierung) in Kapitel 11 in *IBM SPSS Modeler 14.2- Anwendungshandbuch*.

**Beispiel für ein multinomiales Modell.** Ein Telekommunikationsanbieter hat seinen Kundenstamm nach Servicenutzungsmustern in vier Gruppen unterteilt hat. Wenn Sie mithilfe von demografischen Daten die Gruppenzugehörigkeit prognostizieren, können Sie ein multinomiales Modell erstellen, mit dem Sie potenzielle Kunden in Gruppen einteilen und dann die Angebote für die einzelnen Kunden anpassen können. Für weitere Informationen siehe Thema Klassifizieren von Kunden im Telekommunikationsbereich (multinomiale logistische Regression) in Kapitel 13 in *IBM SPSS Modeler 14.2- Anwendungshandbuch*.

**Anforderungen.** Es werden mindestens ein Eingabefeld und genau ein kategoriales Zielfeld mit mindestens zwei Kategorien benötigt. Bei einem binomialen Modell muss das Ziel über ein Messniveau des Typs *Flag* verfügen. Bei einem multinomialen Modell kann das Ziel ein Messniveau von *Flag* oder *Nominal* mit mindestens zwei Kategorien aufweisen. Felder, die auf *Beides* oder *Keine* gesetzt sind, werden ignoriert. Bei den im Modell verwendeten Feldern müssen die Typen vollständig instanziiert sein.

**Stärken.** Logistische Regressionsmodelle sind häufig ziemlich genau. Sie können symbolische und numerische Eingabefelder verarbeiten. Sie können die vorhergesagten Wahrscheinlichkeiten für alle Zielkategorien angeben, sodass der zweitbeste Kandidat problemlos ermittelt werden kann. Logistische Modelle sind am effektivsten, wenn es sich bei der Gruppenmitgliedschaft um ein echt kategoriales Feld handelt; wenn die Gruppenmitgliedschaft auf Werten eines kontinuierlichen Bereichsfelds (z. B. hoher IQ gegenüber niedrigem IQ) basiert, sollten Sie die Anwendung der linearen Regression in Erwägung ziehen, um die umfassenderen Informationen nutzen zu können, die der vollständige Wertebereich bietet. Logistische Modelle können auch eine automatische Feldauswahl durchführen, auch wenn andere Ansätze, wie beispielsweise Baummodelle oder die Merkmalsauswahl, diese Aufgabe bei großen Daten-Sets schneller durchführen können. Und schließlich sind viele Analysten und Data-Mining-Experten gut mit logistischen Modellen vertraut, weshalb sie als Basis verwendet werden können, mit der andere Modellierungstechniken verglichen werden können.

Bei der Verarbeitung großer Daten-Sets können Sie die Leistung deutlich verbessern, indem Sie den Likelihood-Quotienten-Test, eine erweiterte Ausgabeoption, deaktivieren. Für weitere Informationen siehe Thema Logistische Regression – Erweiterte Ausgabe auf S. 298.

## **Logistiknoten – Modelloptionen**

**Modellname.** Sie können den Modellnamen automatisch basierend auf den Ziel- oder ID-Feldnamen (oder dem Modelltyp in Fällen, in denen kein solches Feld angegeben ist) generieren oder einen benutzerdefinierten Namen eingeben.

**Partitionierte Daten verwenden.** Wenn ein Partitionsfeld definiert ist, gewährleistet diese Option, dass nur Daten aus der Trainingspartition für die Modellerstellung verwendet werden. [Für weitere Informationen siehe Thema Partitionsknoten in Kapitel 4 in IBM SPSS Modeler 14.2- Quellen-, Prozess- und Ausgabeknoten.](#)

**Aufteilungsmodelle erstellen.** Erstellt ein separates Modell für jeden möglichen Wert von Eingabefeldern, die als Aufteilungsfelder angegeben werden. [Für weitere Informationen siehe Thema Erstellung von aufgeteilten Modellen in Kapitel 3 auf S. 31.](#)

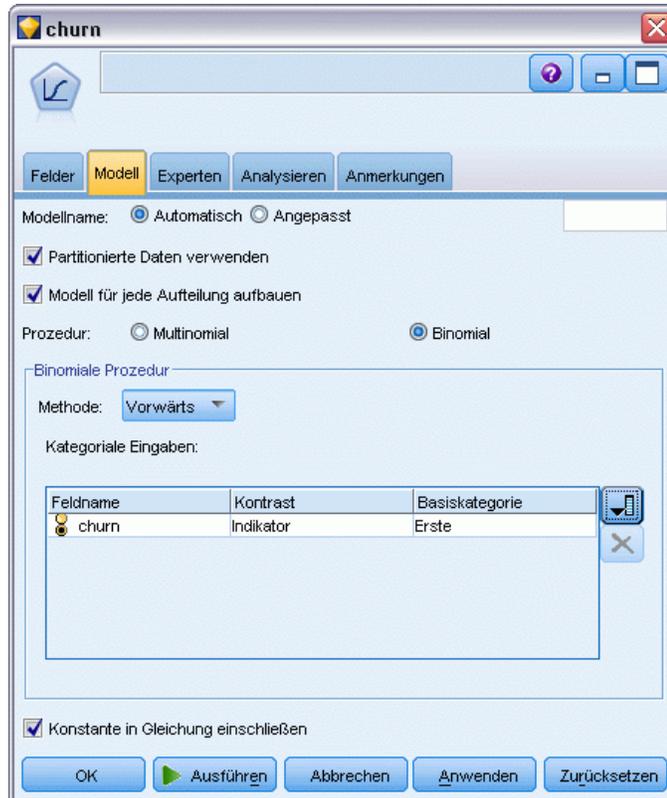
**Prozedur.** Gibt an, ob ein binomiales oder ein multinomiales Modell erstellt wird. Welche Optionen im Dialogfeld zur Verfügung stellen, hängt davon ab, welche Art von Modellierungsprozedur ausgewählt wurde.

- **Binomial.** Wird verwendet, wenn das Zielfeld ein Flag- oder nominales Feld mit zwei diskreten Werten (dichotom) ist, wie *Ja/Nein*, *Ein/Aus*, *männlich/weiblich*.
- **Multinomial.** Wird verwendet, wenn das Zielfeld ein nominales Feld mit mehr als zwei Werten ist. Sie können Haupteffekte, Gesättigt or Benutzerdefiniert auswählen.

**Konstante in Gleichung einschließen.** Mit dieser Option wird bestimmt, ob die entstehenden Gleichungen einen konstanten Term enthalten. In den meisten Fällen sollten Sie diese Option aktiviert lassen.

## Binomiale Modelle

Abbildung 10-21  
Logistikknoten – Optionen für binomiale Modelle



Für binomiale Modelle stehen folgende Methoden und Optionen zur Verfügung:

**Methode.** Dient zur Angabe der bei der Erstellung des logistischen Regressionsmodells verwendeten Methode.

- **Eingabe.** Dies ist das Standardverfahren, bei dem alle Terme direkt in die Gleichung aufgenommen werden. Beim Erstellen des Modells wird keine Feldauswahl durchgeführt.
- **Vorwärts.** Die Feldauswahlmethode “Vorwärts” erstellt das Modell, indem schrittweise nach vorn gegangen wird. Bei dieser Methode ist das ursprüngliche Modell das einfachste und nur die Konstante und die Terme können zum Modell hinzugefügt werden. Bei jedem Schritt werden die Terme, die sich noch nicht im Modell befinden, darauf getestet, wie sehr sie das Modell verbessern würden und der beste davon wird zum Modell hinzugefügt. Wenn keine Terme mehr hinzugefügt werden können oder der beste der in Frage kommenden Terme nicht zu einer hinreichend großen Verbesserung des Modells führen würde, wird das endgültige Modell generiert.
- **Rückwärts.** Die Methode “Rückwärts” ist im Grunde das Gegenteil der Methode “Vorwärts”. Bei dieser Methode enthält das ursprüngliche Modell alle Terme als Prädiktoren und es können nur Terme aus dem Modell entfernt werden. Modellterme, die kaum zum Modell beitragen, werden nach und nach entfernt, bis keine Terme mehr ohne eine signifikante Verschlechterung des Modells entfernt werden können. So entsteht das endgültige Modell.

**Kategoriale Eingaben.** Listet die Felder auf, die als kategorial gekennzeichnet sind, d. h. diejenigen mit einem Messniveau von “Flag”, “Nominal” oder “Ordinal”. Sie können den Kontrast und die Basiskategorie für jedes kategoriale Feld angeben.

- **Feldname.** Diese Spalte enthält die Feldnamen der kategorialen Eingaben und ist bereits automatisch mit allen Flag- und nominalen Werten in den Daten ausgefüllt. Um kontinuierliche oder numerische Eingaben in diese Spalte hinzuzufügen, klicken Sie auf das Symbol “Felder hinzufügen” auf der rechten Seite der Liste und wählen Sie die erforderlichen Eingaben aus.

- **Kontrast.** Die Interpretation der Regressionskoeffizienten für ein kategoriales Feld hängt von den verwendeten Kontrasten ab. Der Kontrast bestimmt, wie die Hypothesentests zum Vergleich der geschätzten Mittel eingerichtet werden. Beispiel: Wenn Sie wissen, dass ein kategoriales Feld eine implizite Reihenfolge aufweist, beispielsweise ein Muster oder eine Gruppierung, können Sie den Kontrast verwenden, um diese Reihenfolge zu modellieren. Folgende Kontraste sind verfügbar:

**Indikator.** Die Kontraste kennzeichnen das Vorhandensein oder Nichtvorhandensein einer Kategoriezugehörigkeit. Dies ist die Standardmethode.

**Einfach.** Jede Kategorie des Prädiktorfelds – mit Ausnahme der Referenzkategorie selbst – wird mit der Referenzkategorie verglichen.

**Differenz.** Jede Kategorie des Prädiktorfelds – mit Ausnahme der ersten Kategorie – wird mit dem durchschnittlichen Effekt der vorherigen Kategorien verglichen. Dies ist auch als umgekehrte Helmert-Kontraste bekannt.

**Helmert.** Jede Kategorie des Prädiktorfelds – mit Ausnahme der letzten Kategorie – wird mit dem durchschnittlichen Effekt der nachfolgenden Kategorien verglichen.

**Wiederholt.** Jede Kategorie des Prädiktorfelds – mit Ausnahme der ersten Kategorie – wird mit der Kategorie verglichen, die ihr unmittelbar vorangeht.

**Polynomial.** Orthogonale polynomiale Kontraste. Es wird angenommen, dass zwischen den Kategorien die gleichen Abstände vorliegen. Polynomiale Kontraste sind nur für numerische Felder verfügbar.

**Abweichung.** Jede Kategorie des Prädiktorfelds – mit Ausnahme der Referenzkategorie – wird mit dem Gesamteffekt verglichen.

- **Basiskategorie.** Gibt an, wie die Referenzkategorie für den ausgewählten Kontrasttyp bestimmt wird. Wählen Sie Erste, um die erste Kategorie für das Eingabefeld zu verwenden (alphabetische Sortierung), oder wählen Sie Letzte, um die letzte Kategorie zu verwenden. Der Standardwert lautet “Erste”.

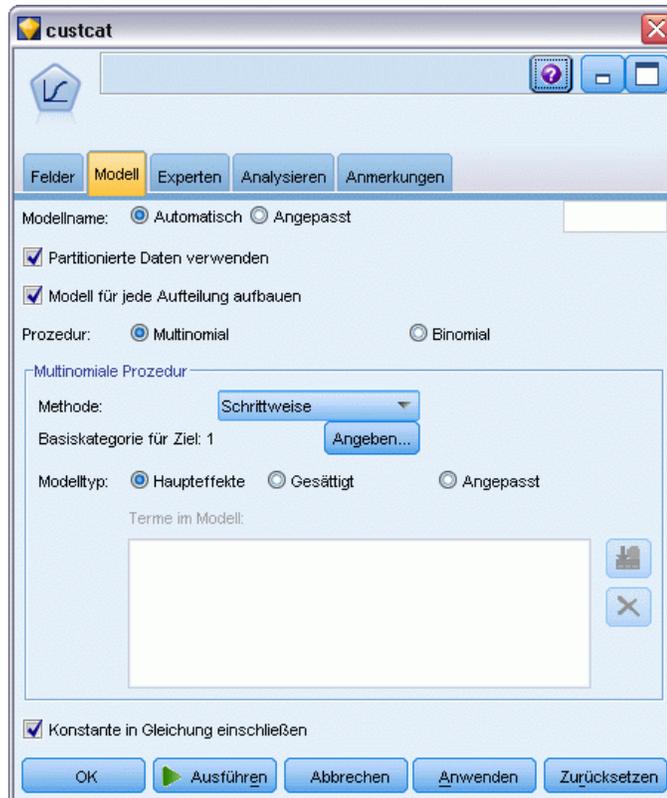
*Hinweis:* Dieses Feld steht bei den Kontrasteinstellungen “Differenz”, “Helmert”, “Wiederholt” oder “Polynomial” nicht zur Verfügung.

Der Schätzer des Effekts der einzelnen Felder auf die Gesamttrefferquote (Gesamtantwort) wird als Anstieg oder Absinken der Likelihood der einzelnen anderen Kategorien relativ zur Bezugskategorie berechnet. Dadurch können Sie besser die Felder und Werte ermitteln, die mit höherer Wahrscheinlichkeit zu einer bestimmten Antwort führen.

Die Basiskategorie wird in der Ausgabe als “0,0” angezeigt. Dies liegt daran, dass der Vergleich mit ihr selbst zu einem leeren Ergebnis führt. Alle anderen Kategorien werden als Gleichungen in Bezug auf die Basiskategorie angezeigt. [Für weitere Informationen siehe Thema Logistik-Modell-Nugget – Details auf S. 303.](#)

### **Multinomiale Modelle**

Abbildung 10-22  
Logistiknoten – Optionen für multinomiale Modelle



Für multinomiale Modelle stehen folgende Methoden und Optionen zur Verfügung:

**Methode.** Dient zur Angabe der bei der Erstellung des logistischen Regressionsmodells verwendeten Methode.

- **Eingabe.** Dies ist das Standardverfahren, bei dem alle Terme direkt in die Gleichung aufgenommen werden. Beim Erstellen des Modells wird keine Feldauswahl durchgeführt.
- **Schrittweise.** Bei der Methode “Schrittweise” der Feldauswahl wird, wie der Name andeutet, die Gleichung in Schritten erstellt. Das anfängliche Modell ist das einfachste Modell, das möglich ist. Es enthält keine Modellterme (außer der Konstanten) in der Gleichung. Bei jedem Schritt werden die Terme, die noch nicht zum Modell hinzugefügt wurden, bewertet, und wenn der beste dieser Terme in signifikanter Weise zur Vorhersagekraft des Modells beiträgt, wird er hinzugefügt. Außerdem werden die derzeit im Modell enthaltenen Terme neu bewertet, um zu ermitteln, ob einige davon ohne signifikante Beeinträchtigung des Modells entfernt werden können. Wenn dies der Fall ist, werden sie entfernt. Der Vorgang wird wiederholt und andere

Terme werden hinzugefügt und/oder entfernt. Wenn keine weiteren Terme zur Verbesserung des Modells hinzugefügt werden können und keine Terme mehr entfernt werden können, ohne das Modell zu beeinträchtigen, wird das endgültige Modell generiert.

- **Vorwärts.** Die Feldauswahlmethode “Vorwärts” ähnelt der Methode “Schrittweise” dahingehend, dass das Modell in Schritten aufgebaut wird. Allerdings ist bei dieser Methode das ursprüngliche Modell das einfachste und die Konstante und die Terme können nur zum Modell hinzugefügt werden. Bei jedem Schritt werden die Terme, die sich noch nicht im Modell befinden, darauf getestet, wie sehr sie das Modell verbessern würden und der beste davon wird zum Modell hinzugefügt. Wenn keine Terme mehr hinzugefügt werden können oder der beste der in Frage kommenden Terme nicht zu einer hinreichend großen Verbesserung des Modells führen würde, wird das endgültige Modell generiert.
- **Rückwärts.** Die Methode “Rückwärts” ist im Grunde das Gegenteil der Methode “Vorwärts”. Bei dieser Methode enthält das ursprüngliche Modell alle Terme als Prädiktoren und es können nur Terme aus dem Modell entfernt werden. Modellterme, die kaum zum Modell beitragen, werden nach und nach entfernt, bis keine Terme mehr ohne eine signifikante Verschlechterung des Modells entfernt werden können. So entsteht das endgültige Modell.
- **Schrittweise rückwärts.** Die Methode “Schrittweise rückwärts” ist im Grunde das Gegenteil der Methode “Schrittweise”. Bei dieser Methode enthält das ursprüngliche Modell alle Terme als Prädiktoren. Bei jedem Schritt werden die Terme im Modell evaluiert und alle Terme, die ohne signifikante Beeinträchtigung des Modells entfernt werden können, werden entfernt. Außerdem werden die zuvor entfernten Terme erneut evaluiert, um zu ermitteln, ob der beste dieser Terme in signifikanter Weise zur Vorhersagekraft des Modells beiträgt. Ist dies der Fall, so wird er wieder in das Modell aufgenommen. Wenn keine Terme mehr entfernt werden können, ohne das Modell wesentlich zu beeinträchtigen, und keine weiteren Terme zur Verbesserung des Modells hinzugefügt werden können, wird das endgültige Modell generiert.

*Hinweis:* Die automatischen Methoden (“Schrittweise”, “Vorwärts” und “Rückwärts”) sind sehr anpassungsfähige Lernmethoden und weisen eine starke Tendenz zur übermäßigen Anpassung an die Trainingsdaten auf. Bei der Verwendung dieser Methoden ist es ganz besonders wichtig, die Validität des entstehenden Modells zu überprüfen – entweder mit neuen Daten oder mithilfe einer zurückgehaltenen Teststichprobe, die mit dem Partitionsknoten erstellt wurde. [Für weitere Informationen siehe Thema Partitionsknoten in Kapitel 4 in IBM SPSS Modeler 14.2- Quellen-, Prozess- und Ausgabeknoten.](#)

**Basiskategorie für Ziel.** Gibt an, wie die Referenzkategorie bestimmt wird. Diese wird als Basis verwendet, anhand deren die Regressionsgleichungen für alle anderen Kategorien im Ziel geschätzt werden. Wählen Sie Erste, um die erste Kategorie für das aktuelle Zielfeld zu verwenden (alphabetische Sortierung), oder wählen Sie Letzte, um die letzte Kategorie zu verwenden. Alternativ können Sie durch Klicken auf Angeben eine bestimmte Kategorie auswählen und dann den gewünschten Wert in der Liste auswählen. Die verfügbaren Werte lassen sich für jedes Feld in einem Typknoten definieren. [Für weitere Informationen siehe Thema Verwenden des Dialogfelds “Werte” in Kapitel 4 in IBM SPSS Modeler 14.2- Quellen-, Prozess- und Ausgabeknoten.](#)

Häufig wird die Kategorie als Basiskategorie angegeben, an der das geringste Interesse besteht, beispielsweise ein Lockartikel. Die anderen Kategorien werden dann auf relative Weise in Bezug zur Basiskategorie gesetzt, um zu bestimmen, wodurch sie mit höherer Wahrscheinlichkeit in ihre eigene Kategorie fallen. Dadurch können Sie besser die Felder und Werte ermitteln, die mit höherer Wahrscheinlichkeit zu einer bestimmten Antwort führen.

Die Basiskategorie wird in der Ausgabe als “0,0” angezeigt. Dies liegt daran, dass der Vergleich mit ihr selbst zu einem leeren Ergebnis führt. Alle anderen Kategorien werden als Gleichungen in Bezug auf die Basiskategorie angezeigt. [Für weitere Informationen siehe Thema Logistik-Modell-Nugget – Details auf S. 303.](#)

**Modelltyp.** Es gibt drei Optionen zur Definition der Terme im Modell. **Haupteffekte**-Modelle beinhalten nur die einzelnen Eingabefelder und testen nicht die Interaktionen (multiplikativen Effekte) zwischen den Eingabefeldern. **Gesättigte** Modelle enthalten alle Interaktionen sowie die Haupteffekte der Eingabefelder. Gesättigte Modelle sind besser zur Erfassung komplexer Beziehungen in der Lage, sind jedoch auch wesentlich schwieriger zu interpretieren und neigen wesentlich stärker zur übermäßigen Anpassung. Aufgrund der potenziell großen Anzahl möglicher Kombinationen sind die automatischen Methoden zur Feldauswahl (alle Methoden außer “Einschluss”) für gesättigte Modelle deaktiviert. **Benutzerdefinierte Modelle** enthalten nur die von Ihnen angegebenen Terme (Haupteffekte und Interaktionen). Verwenden Sie bei der Auswahl dieser Option die Liste “Terme im Modell”, um Terme zum Modell hinzuzufügen oder daraus zu entfernen.

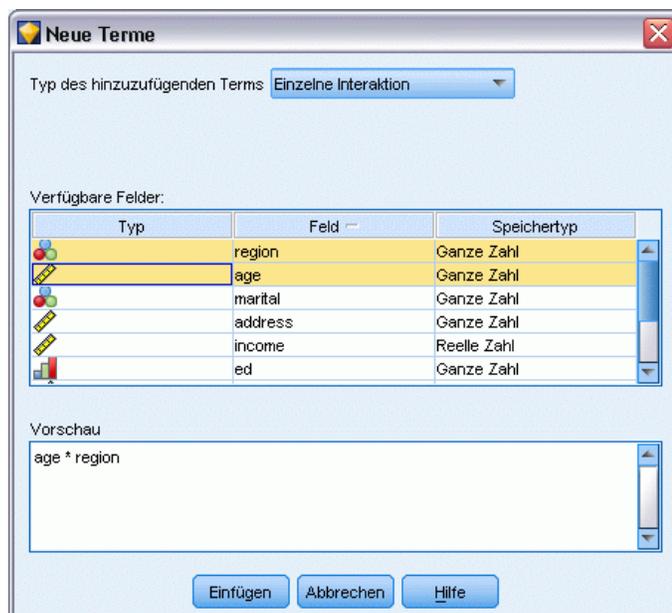
**Terme im Modell.** Beim Erstellen eines benutzerdefinierten Modells müssen Sie die Terme im Modell explizit angeben. Die Liste zeigt die aktuelle Menge an Termen für das Modell. Mit den Schaltflächen auf der rechten Seite der Liste “Terme im Modell” können Sie Modellterme hinzufügen bzw. entfernen.

- ▶ Um Terme zum Modell hinzuzufügen, klicken Sie auf die Schaltfläche *Neue Terme im Modell hinzufügen*.
- ▶ Zum Löschen von Termen wählen Sie die gewünschten Terme aus und klicken Sie auf die Schaltfläche *Ausgewählte Terme im Modell löschen*.

### ***Hinzufügen von Termen zu einem logistischen Regressionsmodell***

Beim Anfordern eines benutzerdefinierten logistischen Regressionsmodells können Sie Terme zum Modell hinzufügen, indem Sie auf der Registerkarte für das logistische Regressionsmodell auf die Schaltfläche *Neue Terme im Modell hinzufügen* klicken. Ein neues Dialogfeld wird geöffnet, in dem Sie Terme angeben können.

Abbildung 10-23  
Logistische Regression – Dialogfeld "Neue Terme"



**Typ des hinzuzufügenden Terms.** Es gibt mehrere Methoden zum Hinzufügen von Termen zum Modell, je nach der Auswahl der Eingabefelder in der Liste der verfügbaren Felder.

- **Einzelne Interaktion.** Fügt den Term ein, der für die Interaktion aller ausgewählten Felder steht.
- **Haupteffekte.** Fügt für jedes ausgewählte Eingabefeld einen Haupteffekt-Term (das Feld selbst) ein.
- **Alle zweifachen Interaktionen.** Fügt für jedes mögliche Paar ausgewählter Eingabefelder einen Zweifach-Interaktions-Term (das Produkt der Eingabefelder) ein. Beispiel: Bei Auswahl der Eingabefelder  $A$ ,  $B$  und  $C$  in der Liste der verfügbaren Felder werden bei dieser Methode die Terme  $A * B$ ,  $A * C$  und  $B * C$  eingefügt.
- **Alle dreifachen Interaktionen.** Fügt für jede mögliche Kombination aus jeweils drei ausgewählten Eingabefeldern einen Dreifach-Interaktions-Term (das Produkt der Eingabefelder) ein. Beispiel: Bei Auswahl der Eingabefelder  $A$ ,  $B$ ,  $C$  und  $D$  in der Liste der verfügbaren Felder werden bei dieser Methode die Terme  $A * B * C$ ,  $A * B * D$ ,  $A * C * D$  und  $B * C * D$  eingefügt.
- **Alle vierfachen Interaktionen.** Fügt für jede mögliche Kombination aus jeweils vier ausgewählten Eingabefeldern einen Vierfach-Interaktions-Term (das Produkt der Eingabefelder) ein. Beispiel: Bei Auswahl der Eingabefelder  $A$ ,  $B$ ,  $C$ ,  $D$  und  $E$  in der Liste der verfügbaren Felder werden bei dieser Methode die Terme  $A * B * C * D$ ,  $A * B * C * E$ ,  $A * B * D * E$ ,  $A * C * D * E$  und  $B * C * D * E$  eingefügt.

**Verfügbare Felder.** Listet die verfügbaren Eingabefelder auf, die bei der Konstruktion der Modellterme verwendet werden sollen.

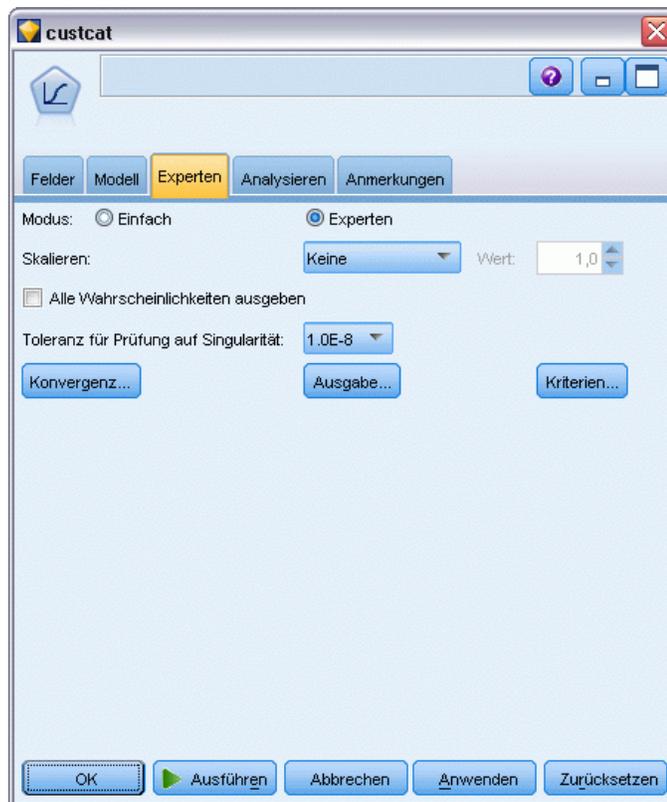
**Vorschau.** Zeigt die Terme an, die beim Klicken auf Einfügen zum Modell hinzugefügt werden. Dabei werden die ausgewählten Felder und der Termtyp zugrunde gelegt.

**Einfügen.** Fügt (auf der Grundlage der aktuellen Auswahl von Feldern und des Termtyps) Terme in das Modell ein und schließt das Dialogfeld.

### **Expertenoptionen für Logistikknoten**

Wenn Sie über umfassende Kenntnisse im Bereich der logistischen Regression verfügen, können Sie mithilfe der Expertenoptionen eine Feinabstimmung des Trainingsvorgangs vornehmen. Für den Zugriff auf die Expertenoptionen setzen Sie den Modus auf der Registerkarte "Experten" auf Experten.

Abbildung 10-24  
Logistische Regression – Registerkarte "Experten"



**Skalieren (nur bei multinomialen Modellen).** Hier können Sie den Skalierungswert für die Streuung angeben, mit dem die Schätzung der Parameter-Kovarianzmatrix korrigiert wird. Bei der Option Pearson wird der Skalierungswert unter Verwendung der Chi-Quadrat-Statistik nach Pearson geschätzt. Bei der Option Devianz wird der Skalierungswert unter Verwendung der Devianzfunktion (Likelihood-Quotienten-Chi-Quadrat) geschätzt. Außerdem können Sie auch einen eigenen, benutzerdefinierten Skalierungswert angeben. Hierbei muss es sich um einen positiven numerischen Wert handeln.

**Alle Wahrscheinlichkeiten ausgeben.** Bei Auswahl dieser Option werden die Wahrscheinlichkeiten für die einzelnen Kategorien des Ausgabefelds zu den einzelnen Datensätzen hinzugefügt, die vom Knoten verarbeitet werden. Wenn diese Option nicht ausgewählt wurde, wird nur die Wahrscheinlichkeit der vorhergesagten Kategorie hinzugefügt.

Beispielsweise kommen zu einer Tabelle, die die Ergebnisse eines multinomialen Modells mit drei Kategorien enthält, fünf neue Spalten hinzu. In einer Spalte wird die Wahrscheinlichkeit dafür angegeben, dass das Ergebnis korrekt prognostiziert wird, in der nächsten Spalte wird die Wahrscheinlichkeit angegeben, dass diese Prognose ein Treffer oder ein Fehlschlag ist, und in drei weiteren Spalten wird die Wahrscheinlichkeit angegeben, dass die Prognose für die einzelnen Kategorien ein Treffer oder ein Fehlschlag ist. [Für weitere Informationen siehe Thema Logistik-Modell-Nugget auf S. 302.](#)

*Hinweis:* Bei binomialen Modellen ist diese Option immer ausgewählt.

**Toleranz für Prüfung auf Singularität.** Dient zur Angabe der Toleranz, die bei der Prüfung auf Singularitäten verwendet wird.

**Konvergenz.** Mit diesen Optionen können Sie die Parameter für die Modellkonvergenz festlegen. Bei der Ausführung des Modells steuern die Konvergenzkriterien, wie viele Male die verschiedenen Parameter wiederholt durchlaufen werden, um zu ermitteln, wie gut sie passen. Je häufiger die Parameter durchprobiert werden, desto enger liegen die Ergebnisse beieinander (d. h. die Ergebnisse konvergieren). [Für weitere Informationen siehe Thema Logistische Regression – Konvergenzoptionen auf S. 297.](#)

**Ausgabe.** Mit diesen Optionen können Sie zusätzliche Statistiken anfordern, die in der erweiterten Ausgabe des vom Knoten erstellten Modell-Nugget erscheinen. [Für weitere Informationen siehe Thema Logistische Regression – Erweiterte Ausgabe auf S. 298.](#)

**Kriterien.** Mit diesen Optionen können Sie die Kriterien zum Hinzufügen und Entfernen von Feldern mit den Schätzmethoden “Schrittweise”, “Vorwärts”, “Rückwärts” oder “Schrittweise rückwärts” festlegen. (Die Schaltfläche ist deaktiviert, wenn die Methode “Einschluss” ausgewählt ist.) [Für weitere Informationen siehe Thema Logistische Regression – Optionen für die Schrittkriterien auf S. 301.](#)

## **Logistische Regression – Konvergenzoptionen**

Sie können die Konvergenzparameter für die Schätzung des logistischen Regressionsmodells festlegen.

Abbildung 10-25  
Logistische Regression – Konvergenzoptionen



**Maximale Anzahl der Iterationen.** Dient zur Angabe der maximalen Anzahl der Iterationen, die für die Schätzung des Modells verwendet werden.

**Maximalzahl für Schritt-Halbierung.** Die Schritthalbierung ist ein Verfahren, das von der logistischen Regression verwendet wird, um Komplexitäten im Schätzvorgang zu verarbeiten. Unter normalen Umständen sollten Sie die Standardeinstellung verwenden.

**Log-Likelihood-Konvergenz.** Iterationen werden angehalten, wenn die relative Änderung der Log-Wahrscheinlichkeit (Log-Likelihood) kleiner als dieser Wert ist. Wenn der Wert gleich 0 ist, wird dieses Kriterium nicht angewendet.

**Parameter-Konvergenz.** Die Iterationen werden angehalten, wenn die absolute oder relative Änderung in den Parameterschätzungen kleiner als dieser Wert ist. Wenn der Wert gleich 0 ist, wird dieses Kriterium nicht angewendet.

**Delta (nur bei multinomialen Modellen).** Sie können einen Wert zwischen 0 und 1 angeben, der zu jeder leeren Zelle hinzugefügt werden soll (Kombination aus Ein- und Ausgabefeldwerten). Dadurch kann der Schätzalgorithmus besser mit Daten umgehen, wenn viele mögliche Kombinationen von Feldwerten relativ zur Anzahl der Datensätze in den Daten vorhanden sind. Der Standardwert ist 0.

### **Logistische Regression – Erweiterte Ausgabe**

Wählen Sie die optionalen Ausgaben aus, die in der erweiterten Ausgabe des Regressions-Modell-Nuggets angezeigt werden sollen. Zur Anzeige der erweiterten Ausgabe durchsuchen Sie das Modell-Nugget und klicken Sie auf die Registerkarte Erweitert. [Für weitere Informationen siehe Thema Logistik-Modell-Nugget – Erweiterte Ausgabe auf S. 309.](#)

### Binomialoptionen

Abbildung 10-26

Logistische Regression – Ausgabeoptionen für binomiale Modelle



Wählen Sie die für das Modell zu generierenden Ausgabetypen aus. [Für weitere Informationen siehe Thema Logistik-Modell-Nugget – Erweiterte Ausgabe auf S. 309.](#)

**Anzeigen.** Hier können Sie auswählen, ob die Ergebnisse bei jedem Schritt angezeigt werden sollen oder ob gewartet werden soll, bis alle Schritte durchlaufen wurden.

**Konfidenzintervall für exp(B).** Dient zur Auswahl der Konfidenzintervalle für die einzelnen Koeffizienten (als Beta angezeigt) im Ausdruck. Geben Sie das Niveau des Konfidenzintervalls an (Standard: 95 %).

**Restdiagnose.** Fordert eine Tabelle mit den fallweisen Diagnosen der Residuen an.

- **Ausreißer außerhalb (Standardabweichung).** Listet nur Fälle mit Residuen auf, bei denen der absolute standardisierte Wert der aufgelisteten Variablen mindestens so groß ist wie der von Ihnen angegebene Wert. Der Standardwert lautet 2.
- **Alle Fälle.** Schließt alle Fälle in die Tabelle mit den fallweisen Diagnosen der Residuen ein.

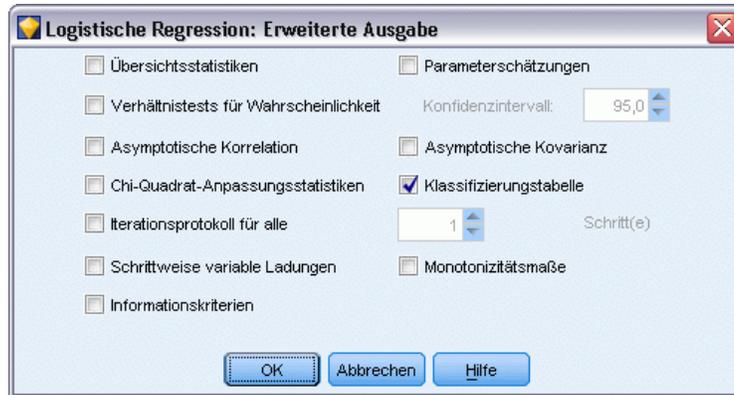
*Hinweis:* Da bei dieser Option alle Eingabedatensätze aufgeführt werden, kann dies zu einer außergewöhnlich großen Tabelle im Bericht führen, mit einer Zeile pro Datensatz.

**Klassifikationsschwellenwert.** Ermöglicht die Bestimmung der Trennwerte für die Fallklassifizierung. Fälle mit vorhergesagten Werten, die den Klassifikationsschwellenwert übersteigen, werden als positiv, vorhergesagte Werte, die unter dem Schwellenwert liegen, als negativ klassifiziert. Um die Standardeinstellung zu ändern, geben Sie einen Wert zwischen 0,01 und 0,99 ein.

### Optionen für multinomiale Modelle

Abbildung 10-27

Logistische Regression – Ausgabeoptionen für multinomiale Modelle



Wählen Sie die für das Modell zu generierenden Ausgabetypen aus. [Für weitere Informationen siehe Thema Logistik-Modell-Nugget – Erweiterte Ausgabe auf S. 309.](#)

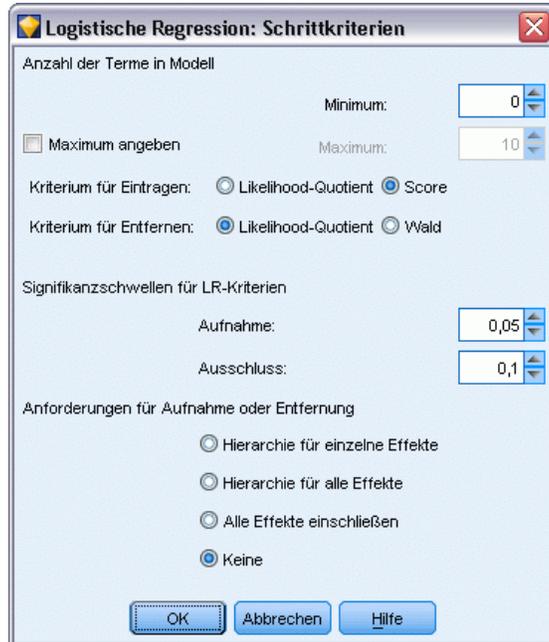
*Hinweis:* Die Auswahl der Option für Likelihood-Quotienten-Tests erhöht die erforderliche Verarbeitungszeit für den Aufbau eines logistischen Regressionsmodells beträchtlich. Wenn die Erstellung des Modells zu lange dauert, sollten Sie diese Option deaktivieren oder stattdessen die Wald- und Score-Statistiken verwenden. [Für weitere Informationen siehe Thema Logistische Regression – Optionen für die Schrittkriterien auf S. 301.](#)

**Iterationsprotokoll für alle.** Dient zur Auswahl des Schrittintervalls für das Drucken des Iterationsstatus in der erweiterten Ausgabe.

**Konfidenzintervall.** Das Konfidenzintervall für Koeffizienten in den Gleichungen. Geben Sie das Niveau des Konfidenzintervalls an (Standard: 95 %).

## Logistische Regression – Optionen für die Schrittkriterien

Abbildung 10-28  
Logistische Regression – Schrittkriterien



**Anzahl der Terme im Modell (nur bei multinomialen Modellen).** Für die Modelle vom Typ “Rückwärts” und “Schrittweise rückwärts” können Sie die Mindestzahl der Terme im Modell angeben, für “Vorwärts” und “Schrittweise vorwärts” die Höchstzahl der Terme. Wenn Sie einen Mindestwert über 0 angeben, enthält das Modell die angegebene Anzahl an Termen, selbst wenn einige davon auf der Grundlage statistischer Kriterien entfernt worden wären. Bei den Modellen “Vorwärts”, “Schrittweise” und “Einschluss” wird die Mindesteinstellung ignoriert. Wenn Sie einen Maximalwert angeben, werden einige Terme möglicherweise aus dem Modell weggelassen, auch wenn sie aufgrund der statistischen Kriterien ausgewählt worden wären. Bei den Modellen “Rückwärts”, “Schrittweise rückwärts” und “Einschluss” wird die Einstellung Maximum angeben ignoriert.

**Kriterium für Eintragen (nur bei multinomialen Modellen).** Wählen Sie Wert, um die Verarbeitungsgeschwindigkeit zu maximieren. Die Option Likelihood-Quotient kann zu robusteren Schätzungen führen, die Berechnung kann jedoch länger dauern. Standardmäßig wird die Score-Statistik verwendet.

**Kriterium für Entfernen.** Wählen Sie Likelihood-Quotient für ein robusteres Modell. Zur Verkürzung der für die Modellerstellung erforderlichen Zeit können Sie Wald auswählen. Wenn in den Daten jedoch eine vollständige oder quasi vollständige Trennung vorliegt (kann durch die Registerkarte “Erweitert” im Modell-Nugget bestimmt werden), wird die Wald-Statistik besonders unzuverlässig und sollte daher nicht verwendet werden. Standardmäßig wird der Statistiktyp “Likelihood-Quotient” verwendet. Bei binomialen Modellen gibt es die zusätzliche Option Bedingt. Diese bietet Ausschlusstests auf der Grundlage der Wahrscheinlichkeit der Likelihood-Quotienten-Statistik, die auf bedingten Parameterschätzern beruht.

**Signifikanzschwellen für LR-Kriterien.** Mit dieser Option können Sie Auswahlkriterien auf der Grundlage der statistischen Wahrscheinlichkeit ( $p$ -Wert) angeben, die den einzelnen Feldern zugeordnet ist. Felder werden nur zum Modell hinzugefügt, wenn der zugehörige  $p$ -Wert kleiner ist als der Wert für Aufnahme, und nur dann entfernt, wenn der  $p$ -Wert größer ist als der Wert für Ausschluss. Der Wert für Aufnahme muss unter dem Wert für Ausschluss liegen.

**Anforderungen für Aufnahme bzw. Ausschluss (nur bei multinomialen Modellen).** Bei einigen Anwendungen hat es, mathematisch gesehen, keinen Sinn, Interaktionsterme zum Modell hinzuzufügen, es sei denn, das Modell enthält außerdem die Terme niedrigerer Ordnung für die zum Interaktionsterm gehörenden Felder. So ist es vielleicht nicht sinnvoll,  $A * B$  in das Modell aufzunehmen, es sei denn,  $A$  und  $B$  kommen ebenfalls im Modell vor. Mit diesen Optionen können Sie festlegen, wie während der schrittweisen Termauswahl mit solchen Abhängigkeiten umgegangen werden soll.

- **Hierarchie für einzelne Effekte.** Effekte höherer Ordnung (Interaktionen, an denen mehr Felder beteiligt sind) werden nur dann in das Modell aufgenommen, wenn alle Effekte niedrigerer Ordnung (Haupteffekte oder Interaktionen mit weniger Feldern) für die betreffenden Felder bereits im Modell enthalten sind, und Effekte niedrigerer Ordnung werden nicht entfernt, wenn Effekte höherer Ordnung, die dieselben Felder betreffen, im Modell vorhanden sind. Diese Option gilt nur für kategoriale Felder. [Für weitere Informationen siehe Thema Messniveaus in Kapitel 4 in IBM SPSS Modeler 14.2- Quellen-, Prozess- und Ausgabeknoten.](#)
- **Hierarchie für alle Effekte.** Diese Option funktioniert genau wie die vorherige, außer dass sie auf alle Felder angewendet wird.
- **Alle Effekte einschließen.** Effekte können nur dann im Modell vorkommen, wenn alle in dem Effekt eingeschlossenen Effekte ebenfalls im Modell vorkommen. Diese Option ähnelt der Option Hierarchie für alle Effekte, mit der Ausnahme, dass die stetige Felder leicht abweichend behandelt werden. Damit ein Effekt einen anderen Effekt einschließt, muss der eingeschlossene Effekt (niedrigerer Ordnung) *alle* stetigen Felder enthalten, die im einschließenden Effekt (höherer Ordnung) enthalten sind, und bei den kategorialen Feldern des eingeschlossenen Effekts muss es sich um eine Untergruppe der diskreten Felder im einschließenden Effekt handeln. Beispiel: Wenn  $A$  und  $B$  kategoriale Felder sind und  $X$  ein stetiges Feld ist, dann schließt der Term  $A * B * X$  die Terme  $A * X$  und  $B * X$  ein.
- **Keine.** Es werden keine Beziehungen erzwungen; die Terme werden unabhängig zum Modell hinzugefügt und daraus entfernt.

## Logistik-Modell-Nugget

Ein Modell-Nugget vom Typ "Logistisch" steht für die Gleichung, die durch einen Logistikknoten geschätzt wurde. Diese enthält alle Informationen, die vom logistischen Regressionsmodell erfasst wurden, sowie Informationen über die Struktur und Leistung des Modells. Dieser Gleichungstyp kann auch von anderen Modellen, wie Oracle SVM, generiert werden.

Wenn Sie einen Stream ausführen, der ein Modell-Nugget vom Typ "Logistisch" enthält, fügt der Knoten zwei neue Felder hinzu, die die Prognose des Modells und die zugehörige Wahrscheinlichkeit enthalten. Die Namen der neuen Felder werden aus dem Namen des prognostizierten Ausgabefelds abgeleitet, dem  $\$L$ - für die vorhergesagte Kategorie und  $\$LP$ - für die zugehörige Wahrscheinlichkeit vorangestellt ist. Bei einem Ausgabefeld mit der Bezeichnung *Farbpräf* beispielsweise erhalten die neuen Felder die Namen  $\$L$ -*Farbpräf* und  $\$LP$ -*Farbpräf*. Wenn Sie außerdem im Logistikknoten die Option Alle Wahrscheinlichkeiten ausgeben ausgewählt

haben, wird für jede Kategorie des Ausgabefelds ein zusätzliches Feld hinzugefügt, das die Wahrscheinlichkeit enthält, die zu der entsprechenden Kategorie für die einzelnen Datensätze gehört. Diese zusätzlichen Felder werden auf der Grundlage der Werte des Ausgabefelds benannt, denen  $SLP$ - vorangestellt wurde. Wenn für *Farbpräf* die Werte *Rot*, *Grün* und *Blau* zulässig sind, werden drei neue Felder hinzugefügt:  $SLP$ -*Rot*,  $SLP$ -*Grün* und  $SLP$ -*Blau*.

**Erstellen eines Filterknotens.** Im Menü “Generieren” können Sie einen neuen Filterknoten zur Übergabe der Eingabefelder auf der Grundlage der Ergebnisse erstellen. Die Felder, die aufgrund von Multikollinearität aus dem Modell herausgenommen wurden, werden vom generierten Knoten gefiltert ebenso wie Felder, die nicht im Modell verwendet werden.

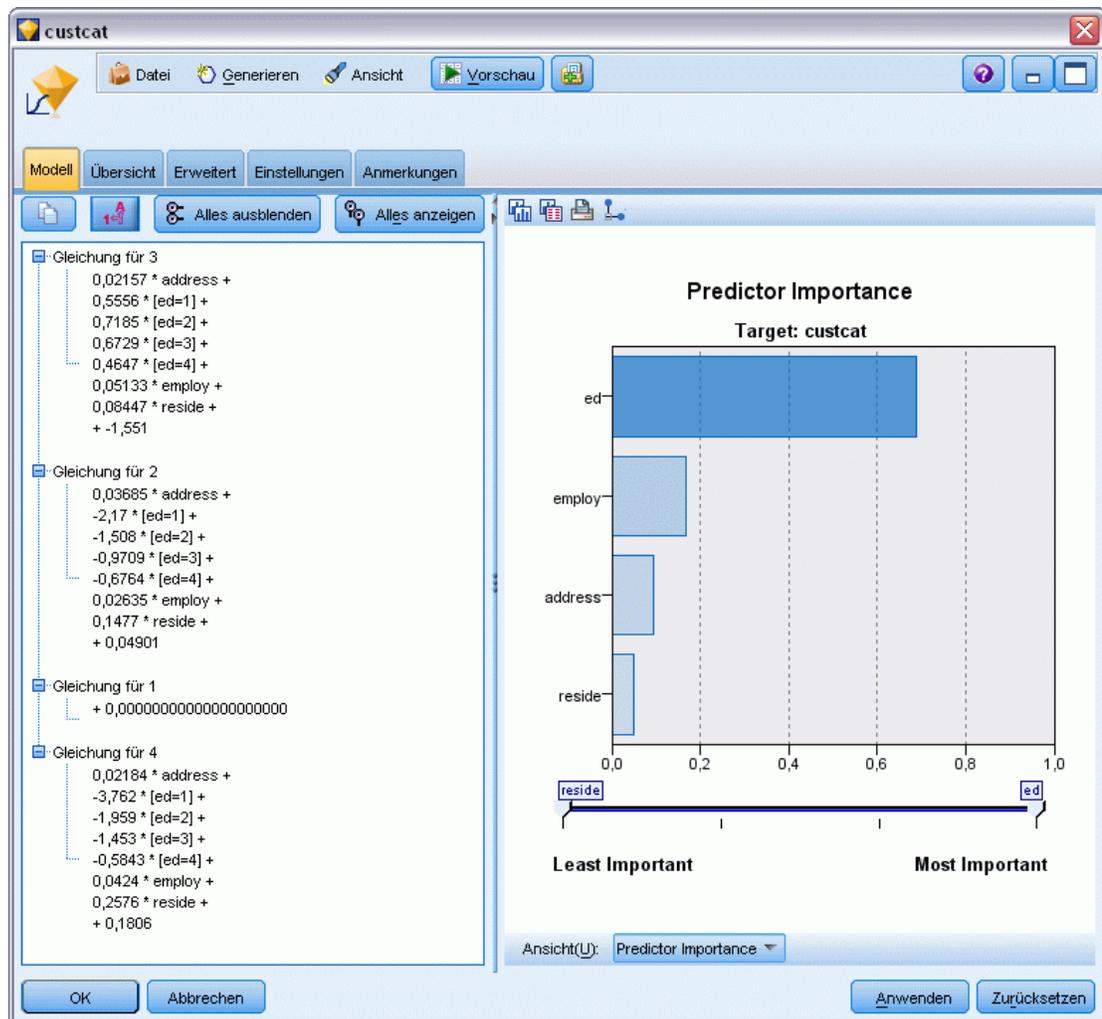
## ***Logistik-Modell-Nugget – Details***

Bei multinomialen Modellen weist die Registerkarte “Modell” in einem Modell-Nugget vom Typ “Logistisch” eine geteilte Anzeige auf. Dabei werden die Modellgleichungen im linken und die Bedeutsamkeit der Prädiktoren im rechten Fensterbereich angezeigt. Bei binomialen Modellen wird auf der Registerkarte nur die Bedeutsamkeit der Prädiktoren angezeigt. [Für weitere Informationen siehe Thema Bedeutsamkeit des Prädiktors in Kapitel 3 auf S. 54.](#)

### ***Modellgleichungen***

Bei multinomialen Modellen werden im linken Fensterbereich die tatsächlich für das logistische Regressionsmodell geschätzten Gleichungen angezeigt. Für jede Kategorie im Zielfeld (mit Ausnahme der Basiskategorie) gibt es jeweils genau eine Gleichung. Die Gleichungen werden in einem Baumformat angezeigt. Dieser Gleichungstyp kann auch von bestimmten anderen Modellen, wie beispielsweise Oracle SVM, generiert werden.

Abbildung 10-29  
 Details für das Logistik-Modell-Nugget mit Anzeige der Bedeutsamkeit der Prädiktoren



**Gleichung für.** Zeigt die Regressionsgleichungen, die bei einem Vorgegebenen Satz an Prädiktorwerten zur Ableitung der Wahrscheinlichkeiten für die Zielkategorie verwendet werden. Die letzte Kategorie des Zielfelds wird als **Basiskategorie** betrachtet; die angezeigten Gleichungen bieten für ein bestimmtes Set an Prädiktorwerten die Log-Odds für die anderen Zielkategorien relativ zur Basiskategorie. Die prognostizierte Wahrscheinlichkeit für die einzelnen Kategorien des jeweiligen Prädiktormusters wird aus diesen Log-Odds-Werten abgeleitet.

#### **Wie werden die Wahrscheinlichkeiten berechnet?**

Bei jeder Gleichung werden die Log-Odds für eine bestimmte Zielkategorie relativ zur Basiskategorie berechnet. Bei **Log-Odds**, auch als **Logit** bezeichnet, handelt es sich um den Quotienten aus der Wahrscheinlichkeit für eine angegebene Zielkategorie und der Wahrscheinlichkeit der Basiskategorie, wobei auf das Ergebnis der natürliche Logarithmus angewendet wird. Bei der Basiskategorie sind die Chancen für die Kategorie relativ zu sich selbst

1,0 und daher ist Log-Odds gleich 0. Dies kann als implizite Gleichung für die Basiskategorie betrachtet werden, bei der alle Koeffizienten gleich 0 sind.

Um die Wahrscheinlichkeit aus dem Log-Odds-Wert für eine bestimmte Zielkategorie abzuleiten, verwenden Sie den von der Gleichung für die betreffende Kategorie berechneten Logit-Wert und wenden Sie folgende Formel an:

$$P(\text{group}_i) = \exp(g_i) / \sum_k \exp(g_k)$$

Dabei ist  $g$  der berechnete Wert für Log-Odds,  $i$  der Kategorieindex und  $k$  liegt im Bereich von 1 bis zur Anzahl der Zielkategorien.

### ***Bedeutsamkeit des Prädiktors***

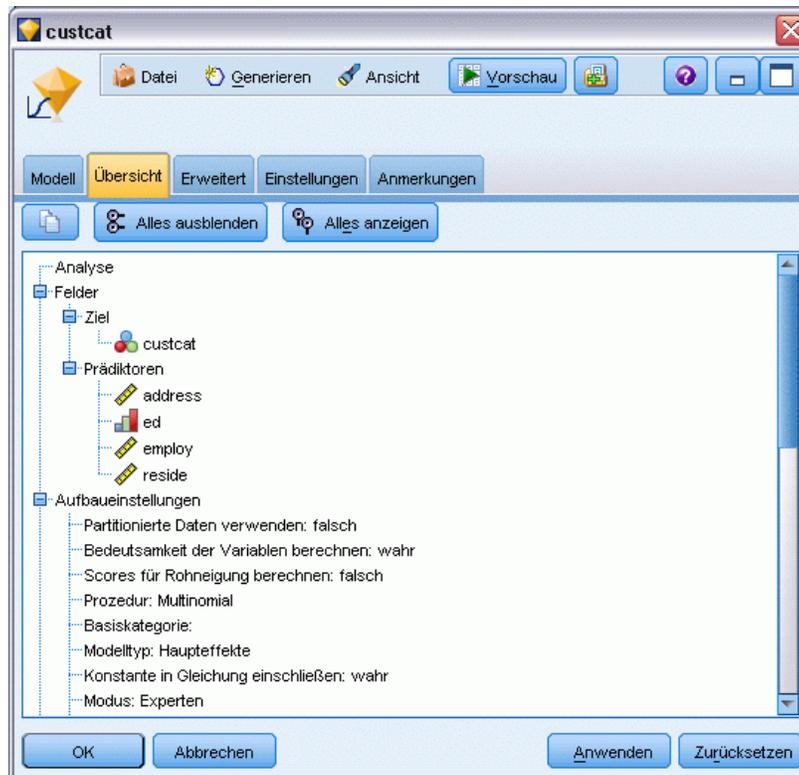
Optional kann auf der Registerkarte “Modell” auch ein Diagramm, das die relative Bedeutsamkeit der einzelnen Prädiktoren für die Schätzung des Modells angibt, angezeigt werden. Normalerweise ist es sinnvoll, die Modellierungsbemühungen auf die wichtigsten Prädiktoren zu konzentrieren und diejenigen Prädiktoren zu verwerfen bzw. zu ignorieren, die die geringste Bedeutung haben. Beachten Sie, dass dieses Diagramm nur verfügbar ist, wenn vor dem Generieren des Modells Bedeutsamkeit der Prädiktoren berechnetauf der Registerkarte “Analysieren” ausgewählt wurde. [Für weitere Informationen siehe Thema Bedeutsamkeit des Prädiktors in Kapitel 3 auf S. 54.](#)

*Anmerkung:* Die Berechnung der Bedeutsamkeit der Prädiktoren kann bei der logistischen Regression länger dauern als bei anderen Modelltypen und ist nicht standardmäßig auf der Registerkarte “Analysieren” aktiviert. Die Auswahl dieser Option kann die Leistung verlangsamen, insbesondere bei großen Daten-Sets.

## ***Logistik-Modell-Nugget – Übersicht***

In der Übersicht für ein logistisches Regressionsmodell werden die Felder und Einstellungen angezeigt, die zum Generieren des Modells verwendet wurden. Wenn Sie einen Analyseknoden ausgeführt haben, der an diesen Modellierungsknoten angehängt ist, werden die Informationen aus dieser Analyse ebenfalls in diesem Abschnitt angezeigt. [Für weitere Informationen siehe Thema Analyseknoden in Kapitel 6 in IBM SPSS Modeler 14.2- Quellen-, Prozess- und Ausgabe knoden.](#) Allgemeine Informationen zur Verwendung des Modellbrowsers finden Sie unter [Durchsuchen von Modell-Nuggets auf S. 52.](#)

Abbildung 10-30  
Logistisches Regressions-Modell-Nugget – Registerkarte “Übersicht”



### **Logistik-Modell-Nugget – Einstellungen**

Auf der Registerkarte “Einstellungen” für ein Modell-Nugget vom Typ “Logistisch” werden während des Modell-Scorings Optionen für Konfidenzen, Wahrscheinlichkeiten, Neigungs-Scores und SQL-Erzeugung angegeben. Diese Registerkarte steht erst zur Verfügung, nachdem das Modell-Nugget zu einem Stream hinzugefügt wurde, und zeigt je nach Modell- und Zieltyp verschiedene Optionen an.

Abbildung 10-31  
Einstellungen für ein multinomiales Modell mit einem nominalen Ziel



### **Multinomiale Modelle**

Für multinomiale Modelle stehen folgende Optionen zur Verfügung:

**Konfidenzen berechnen.** Gibt an, ob während des Scorings die Konfidenzen berechnet werden sollen.

**Scores für Rohneigung berechnen (nur bei Flag-Feldern).** Bei Modellen mit Flag-Zielen (und nur dort) können Sie Scores für die Rohneigung anfordern, die die Likelihood des für das Zielfeld angegebenen Ergebnisses *true* (wahr) anzeigen. Diese Werte werden zusätzlich zu den standardmäßigen Vorhersage- und Konfidenzwerten ausgegeben. Scores für die korrigierte Neigung sind nicht verfügbar. [Für weitere Informationen siehe Thema Analyseoptionen bei Modellierungsknoten in Kapitel 3 auf S. 41.](#)

**Alle Wahrscheinlichkeiten ausgeben.** Gibt an, ob die Wahrscheinlichkeiten für die einzelnen Kategorien des Ausgabefelds zu den einzelnen vom Knoten verarbeiteten Datensätzen hinzugefügt werden. Wenn diese Option nicht ausgewählt wurde, wird nur die Wahrscheinlichkeit der vorhergesagten Kategorie hinzugefügt. Bei einem nominalen Ziel mit drei Kategorien beispielsweise enthält das Scoring-Ergebnis eine Spalte für jede der drei Kategorien sowie eine vierte Spalte, die die Wahrscheinlichkeit der vorhergesagten Kategorie angibt. Beispiel: Wenn die Wahrscheinlichkeiten für die Kategorien *Rot*, *Grün* und *Blau* 0,6; 0,3 bzw. 0,1 betragen, ist die vorhergesagte Kategorie *Rot*, mit einer Wahrscheinlichkeit von 0,6.

**SQL für dieses Modell generieren.** Bei Verwendung von Daten aus einer Datenbank kann SQL-Code per Pushback zur Ausführung an die Datenbank zurückübertragen werden. Dadurch kann bei vielen Operationen eine bessere Leistung erzielt werden. [Für weitere Informationen siehe Thema SQL-Optimierung in Kapitel 6 in IBM SPSS Modeler Server 14.2-Verwaltungs- und Leistungshandbuch.](#)

*Hinweis:* Bei multinomialen Modellen steht keine SQL-Erzeugung zur Verfügung, wenn die Option Alle Wahrscheinlichkeiten anhängen ausgewählt wurde, oder – bei Modellen mit nominalen Zielen – wenn Konfidenzen berechnen ausgewählt wurde. SQL-Erzeugung mit Konfidenzberechnungen wird nur für multinomiale Modelle mit Flag-Zielen unterstützt. Für binomiale Modelle steht die SQL-Erzeugung nicht zur Verfügung.

### **Binomiale Modelle**

Bei binomialen Modellen sind Konfidenzen und Wahrscheinlichkeiten immer aktiviert und die Einstellungen, mit denen diese Optionen deaktiviert werden können, sind nicht verfügbar. Für binomiale Modelle steht die SQL-Erzeugung nicht zur Verfügung. Die einzige Einstellung, die bei binomialen Modellen geändert werden kann, ist die Möglichkeit zur Berechnung der Scores für die Rohneigung. Wie bereits für multinomiale Modelle angegeben, gilt dies nur für Modelle mit Flag-Zielen. [Für weitere Informationen siehe Thema Analyseoptionen bei Modellierungsknoten in Kapitel 3 auf S. 41.](#)

## Logistik-Modell-Nugget – Erweiterte Ausgabe

Abbildung 10-32

Beispiel-Gleichungsknoten der logistischen Regression – Registerkarte "Erweitert"

The screenshot shows the 'custcat' dialog box in SPSS. The 'Erweitert' (Advanced) tab is selected. The main area displays the 'Case Processing Summary' table for the 'Nominal Regression' model. The table lists the following variables and their corresponding counts and marginal percentages:

		N	Marginal Percentage
custcat	Basic service	266	26.6%
	E-service	217	21.7%
	Plus service	281	28.1%
	Total service	236	23.6%
region	Zone 1	322	32.2%
	Zone 2	334	33.4%
	Zone 3	344	34.4%
marital	Unmarried	505	50.5%
	Married	495	49.5%
ed	Did not complete high school	204	20.4%
	High school degree	287	28.7%
	Some college	209	20.9%
	College degree	234	23.4%
	Post-undergraduate degree	66	6.6%
retire	No	953	95.3%
	Yes	47	4.7%

At the bottom of the dialog box, there are buttons for 'OK', 'Abbrechen', 'Anwenden', and 'Zurücksetzen'.

Die erweiterte Ausgabe für die logistische Regression (auch als **nominale Regression** bekannt) bietet detaillierte Informationen zum geschätzten Modell und seiner Leistung. Die meisten der in der erweiterten Ausgabe enthaltenen Informationen weisen einen hohen Fachlichkeitsgrad auf und zur richtigen Interpretation dieser Ausgaben sind umfassende Kenntnisse der logistischen Regressionsanalyse erforderlich.

**Warnungen.** Zeigt etwaige Warnungen oder potenzielle Probleme mit den Ergebnissen an.

**Zusammenfassung der Fallverarbeitung.** Listet die Anzahl der verarbeiteten Datensätze auf, nach den einzelnen symbolischen Feldern im Modell aufgeschlüsselt.

**Schrittübersicht (optional).** Listet die Effekte auf, die bei Verwendung der automatischen Feldauswahl bei jedem Schritt der Modellerstellung hinzugefügt bzw. entfernt werden.

*Hinweis:* Wird nur für die Methoden “Schrittweise”, “Vorwärts”, “Rückwärts” bzw. “Schrittweise rückwärts” angezeigt.

**Iterationsprotokoll (optional).** Zeigt das Iterationsprotokoll von Parameterschätzern für jede  $n$ -te Iteration, ausgehend von den ursprünglichen Schätzern. Dabei ist  $n$  der Wert des Druckintervalls. Bei der Standardvorgabe wird jede Iteration gedruckt ( $n=1$ ).

**Modellanpassungsinformationen (multinomiale Modelle).** Zeigt den Likelihood-Quotienten-Test für Ihr Modell (endgültig) im Vergleich zu einem, bei dem alle Parameterkoeffizienten 0 sind (nur konstanter Term).

**Klassifizierung (optional).** Zeigt die Matrix der vorhergesagten und tatsächlichen Ausgabefeldwerte mit den zugehörigen Prozentsätzen an.

**Chi-Quadrat-Anpassungsstatistiken (optional).** Zeigt die Chi-Quadrat-Statistiken nach Pearson sowie die Likelihood-Quotienten-Chi-Quadrat-Statistiken an. Diese Statistiken testen die Gesamtanpassung des Modells für die Trainingsdaten.

**Hosmer-Lemeshow-Anpassungsgüte (optional)** Zeigt die Ergebnisse der Gruppierung von Fällen in Risikodezile und des Vergleichs der beobachteten Wahrscheinlichkeit mit der erwarteten Wahrscheinlichkeit innerhalb jedes Dezils. Diese Statistik für die Anpassungsgüte ist robuster als die herkömmliche Statistik für die Anpassungsgüte, die in multinomialen Modellen verwendet wird, insbesondere bei Modellen mit kontinuierlichen Kovariaten und bei Studien mit kleinen Stichprobenumfängen.

**Pseudo-R-Quadrat (optional).** Zeigt die  $R$ -Quadrat-Maße für die Anpassungsgüte nach Cox und Snell, Nagelkerke und McFadden an. Diese Statistiken sind in gewisser Weise analog zu der  $R$ -Quadrat-Statistik in der linearen Regression.

**Monotonizitätsmaße (optional).** Zeigt die Anzahl konkordanter Paare, diskordanter Paare und gebundener Paare in den Daten an sowie den Prozentsatz der Gesamtzahl der Paare, den die einzelnen Gruppen darstellen. In dieser Tabelle werden außerdem die Werte Somers-D, Goodman-und-Kruskal-Gamma, Kendall-Tau-a und Konkordanzindex C angezeigt.

**Informationskriterien (optional).** Zeigt das Akaike-Informationskriterium (AIC) und das Schwarz-Bayes-Informationskriterium (BIC).

**Verhältnistest für Wahrscheinlichkeit (optional).** Zeigt Statistiken, die testen, ob die Koeffizienten der Modelleffekte statistisch von 0 abweichen. Signifikante Eingabefelder sind Felder mit sehr niedrigen Signifikanzniveaus in der Ausgabe (mit *Sig.* beschriftet).

**Parameterschätzungen (optional).** Zeigt die Schätzer der Gleichungskoeffizienten, Tests für diese Koeffizienten, aus den Koeffizienten abgeleitete Quotenverhältnisse (beschriftet mit  $Exp(B)$ ) sowie Konfidenzintervalle für die Quotenverhältnisse an.

**Asymptotische Kovarianz-/Korrelationsmatrix (optional).** Zeigt die asymptotischen Kovarianzen und/oder Korrelationen der Koeffizientenschätzungen an.

**Beobachtete und vorhergesagte Häufigkeiten (optional).** Zeigt für jede Kovariaten-Struktur die beobachteten und vorhergesagten Häufigkeiten für die einzelnen Ausgabefeldwerte an. Diese Tabelle kann ziemlich groß sein, insbesondere bei Modellen mit numerischen Eingabefeldern.

Wenn die resultierende Tabelle so groß werden würde, dass sie völlig unhandlich wird, wird sie weggelassen und eine Warnung ausgegeben.

## **Faktor/PCA-Knoten**

Der Faktor/PCA-Knoten bietet leistungsstarke Datenreduktionsverfahren zur Verringerung der Komplexität der Daten. Es stehen zwei ähnliche, aber doch völlig getrennte Ansätze zur Verfügung.

- Die **Hauptkomponentenanalyse (PCA)** findet lineare Kombinationen der Eingabefelder, die die Varianz im gesamten Set der Felder am besten erfassen, wenn die Komponenten orthogonal (lotrecht) zueinander sind. PCA gilt für sämtliche Varianz, darunter sowohl gemeinsame als auch nur für bestimmte Felder geltende Varianz.
- Mit der **Faktorenanalyse** wird versucht, die zugrunde liegenden Konzepte oder **Faktoren** zu bestimmen, die die Korrelationsmuster innerhalb eines Sets beobachteter Felder erklären. Die Faktorenanalyse zielt nur auf die gemeinsame Varianz ab. Varianz, die nur für bestimmte Felder gilt, wird bei der Modellschätzung nicht berücksichtigt. Der Faktor/PCA-Knoten bietet mehrere Methoden der Faktorenanalyse.

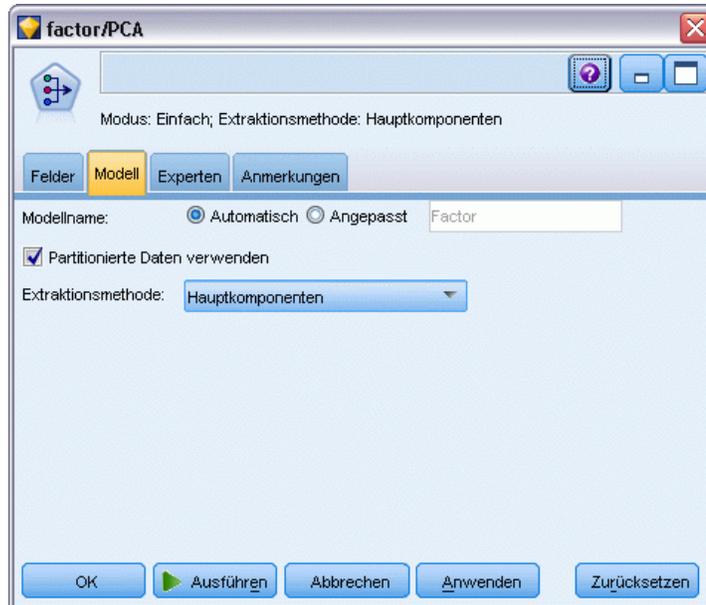
Bei beiden Ansätzen besteht das Ziel darin, eine kleinere Zahl abgeleiteter Felder zu finden, mit denen die Informationen in der ursprünglichen Menge der Felder effektiv zusammengefasst werden können.

**Anforderungen.** In PCA-Faktor-Modellen können nur numerische Felder verwendet werden. Zum Schätzen einer Faktoranalyse oder PCA ist mindestens ein Feld mit der Rolle *Eingabe* erforderlich. Felder, deren Rolle auf *Ziel*, *Beides* oder *Keine* festgelegt ist, wie nicht-numerische Felder.

**Stärken.** Die Faktorenanalyse und die Hauptkomponentenanalyse (PCA) können die Komplexität der Daten effektiv reduzieren, ohne den Informationsgehalt wesentlich zu beeinträchtigen. Mit diesen Verfahren können Sie robustere Modelle erstellen, die schneller ausgeführt werden können, als dies mit den rohen Eingabefeldern der Fall wäre.

## Faktor/PCA-Knoten – Modelloptionen

Abbildung 10-33  
PCA/Factor - Registerkarte "Modell"



**Modellname.** Sie können den Modellnamen automatisch basierend auf den Ziel- oder ID-Feldnamen (oder dem Modelltyp in Fällen, in denen kein solches Feld angegeben ist) generieren oder einen benutzerdefinierten Namen eingeben.

**Partitionierte Daten verwenden.** Wenn ein Partitionsfeld definiert ist, gewährleistet diese Option, dass nur Daten aus der Trainingspartition für die Modellerstellung verwendet werden. [Für weitere Informationen siehe Thema Partitionsknoten in Kapitel 4 in IBM SPSS Modeler 14.2- Quellen-, Prozess- und Ausgabeknoten.](#)

**Extraktionsmethode.** Dient zur Angabe der für die Datenreduktion verwendeten Methode.

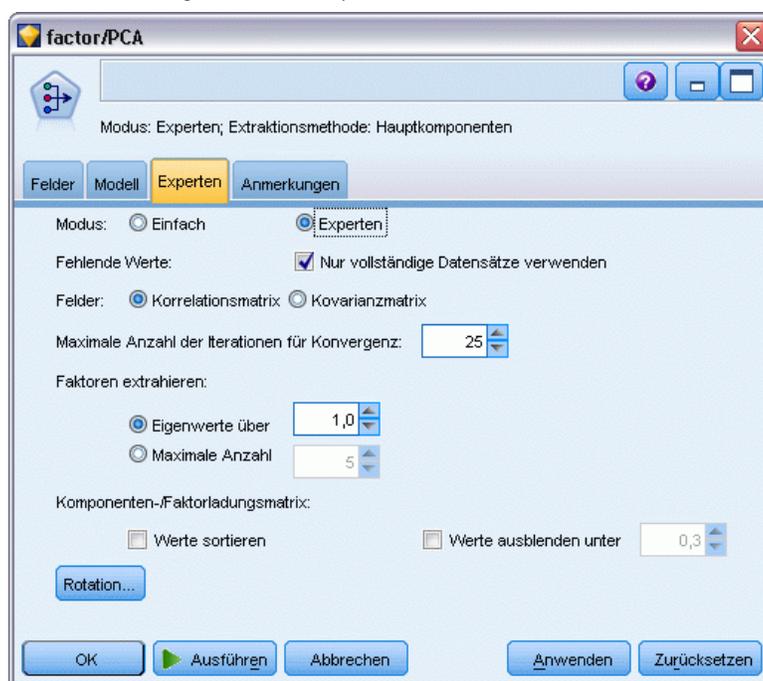
- **Hauptkomponenten.** Dies ist die Standardmethode. Dabei wird die Hauptkomponentenanalyse (PCA) verwendet, um Komponenten zu finden, die die Eingabefelder zusammenfassen.
- **Ungewichtete kleinste Quadrate.** Diese Faktorenanalyseverfahren beruht auf der Suche nach dem Faktoren-Set, das am besten das Muster der Beziehungen (Korrelationen) zwischen den Eingabefeldern reproduzieren kann.
- **Verallgemeinerte kleinste Quadrate.** Diese Faktorenanalyseverfahren ist ähnlich der Methode der ungewichteten kleinsten Quadrate, mit dem Unterschied, dass Gewichtung verwendet wird, um die verstärkte Berücksichtigung von Feldern mit einer großen Menge an spezieller (nicht gemeinsamer) Varianz aufzuheben.
- **Maximum Likelihood.** Bei dieser Faktorenanalyseverfahren werden Faktorgleichungen erstellt, die höchstwahrscheinlich zu dem beobachteten Muster der Beziehungen (Korrelationen) in den Eingabefeldern geführt haben. Hierbei werden Annahmen über die Form dieser Beziehungen zugrunde gelegt. Insbesondere geht die Methode davon aus, dass für die Trainingsdaten eine multivariate Normalverteilung gilt.

- **Hauptachsen-Faktorenanalyse.** Diese Faktorenanalysemethode ähnelt stark der Hauptkomponentenmethode, mit der Ausnahme, dass sie sich nur auf die gemeinsame Varianz konzentriert.
- **Alpha-Faktorisierung.** Diese Faktorenanalysemethode betrachtet die Felder in der Analyse als Beispiel für die Grundgesamtheit potenzieller Eingabefelder. Dadurch wird die statistische Reliabilität der Faktoren maximiert.
- **Image-Faktorisierung.** Diese Faktorenanalysemethode verwendet die Datenschätzung zur Isolation der gemeinsamen Varianz und zum Ermitteln der Faktoren, die sie beschreiben.

### Faktor/PCA-Knoten – Expertenoptionen

Wenn Sie über umfassende Kenntnisse in den Bereichen Faktoranalyse und PCA verfügen, können Sie mithilfe der Expertenoptionen eine Feinabstimmung des Trainingsvorgangs vornehmen. Für den Zugriff auf die Expertenoptionen setzen Sie den Modus auf der Registerkarte “Experten” auf Experten.

Abbildung 10-34  
PCA/Faktor – Registerkarte “Experten”



**Fehlende Werte.** Standardmäßig verwendet IBM® SPSS® Modeler nur Datensätze mit gültigen Werten für alle im Modell verwendeten Felder. (Dies wird zuweilen als **listenweiser Ausschluss** fehlender Werte bezeichnet.) Wenn sehr viele fehlende Daten vorliegen, werden mit diesem Ansatz möglicherweise zu viele Datensätze entfernt, sodass nicht mehr genügend Daten zu Erstellung eines guten Modells vorhanden sind. In solchen Fällen können Sie die Auswahl der Option **Nur vollständige Datensätze verwenden** aufheben. SPSS Modeler versucht dann, so viele Informationen wie möglich zur Schätzung des Modells zu verwenden, auch Datensätze, bei denen bei einigen Feldern fehlende Werte vorliegen. (Dies wird zuweilen als **paarweiser Ausschluss**

fehlender Werte bezeichnet.) In einigen Situationen jedoch kann eine derartige Verwendung unvollständiger Datensätze zu Berechnungsproblemen bei der Schätzung des Modells führen.

**Felder.** Dient zur Angabe, ob die Korrelationsmatrix (Standard) oder die Kovarianzmatrix der Eingabefelder für die Schätzung des Modells verwendet werden soll.

**Maximale Anzahl der Iterationen für Konvergenz.** Dient zur Angabe der maximalen Anzahl der Iterationen, die für die Schätzung des Modells verwendet werden.

**Faktoren extrahieren.** Es gibt zwei Methoden zur Auswahl der Anzahl der Faktoren, die aus den Eingabefeldern extrahiert werden sollen.

- **Eigenwerte über.** Bei dieser Option werden alle Faktoren oder Komponenten beibehalten, die Eigenwerte aufweisen, die größer sind als das angegebene Kriterium. **Eigenwerte** messen die Fähigkeit der einzelnen Faktoren oder Komponenten zur Zusammenfassung der Varianz in der Menge der Eingabefelder. Das Modell führt bei Verwendung der Korrelationsmatrix zur Beibehaltung aller Faktoren oder Komponenten mit Eigenwerten, die größer sind als der angegebene Wert. Bei Verwendung der Kovarianzmatrix wird das Kriterium als Wert mal mittlerer Eigenwert festgelegt. Bei dieser Skalierung hat diese Option eine ähnliche Bedeutung für beide Matrixtypen.
- **Maximale Anzahl.** Bei dieser Option wird die angegebene Anzahl von Faktoren bzw. Komponenten in absteigender Reihenfolge der Eigenwerte beibehalten. Die Faktoren bzw. Komponenten, die den  $n$  höchsten Eigenwerten entsprechen, werden also beibehalten. Dabei ist  $n$  das angegebene Kriterium. Das Standardextraktionskriterium liegt bei fünf Faktoren/Komponenten.

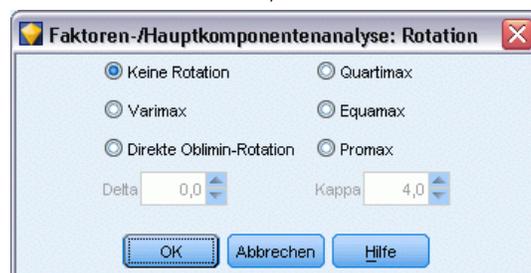
**Komponenten-/Faktorladungsmatrix.** Mit diesen Optionen wird das Format der Faktorladungsmatrix (bzw. der Komponentenladungsmatrix bei PCA-Modellen) festgelegt.

- **Werte sortieren.** Bei Auswahl dieser Option werden die Faktorladungen in der Modellausgabe numerisch sortiert.
- **Werte ausblenden unter.** Bei Auswahl dieser Option werden die Scores unterhalb des angegebenen Schwellenwerts in der Matrix ausgeblendet, damit das Muster in der Matrix besser erkannt werden kann.

**Rotation.** Mit diesen Optionen können Sie die Rotationsmethode für das Modell festlegen. [Für weitere Informationen siehe Thema Faktor/PCA-Knoten – Rotationsoptionen auf S. 314.](#)

## Faktor/PCA-Knoten – Rotationsoptionen

Abbildung 10-35  
PCA/Faktor – Rotationsoptionen



In vielen Fällen kann die mathematische Rotation des Sets der beibehaltenen Faktoren ihre Nützlichkeit und insbesondere ihre Interpretierbarkeit erhöhen. Wählen Sie eine Rotationsmethode aus:

- **Keine Rotation.** Standardoption. Es wird keine Rotation verwendet.
- **Varimax.** Eine orthogonale Rotation, bei der die Anzahl der Felder mit hohen Ladungen für die einzelnen Faktoren minimiert wird. Dadurch wird die Interpretation der Faktoren vereinfacht.
- **Direkte Oblimin-Rotation.** Eine Methode für schiefe (nichtorthogonale) Rotation. Wenn Delta gleich null 0 ist (Standard), sind die Lösungen schief. Mit zunehmendem negativem Wert von Delta werden die Faktoren weniger schiefwinklig. Um den Standardwert von 0 zu überschreiben, geben Sie eine Zahl kleiner gleich 0,8 ein.
- **Quartimax.** Eine orthogonale Methode, bei der die Anzahl der Faktoren, die für die Erklärung der einzelnen Felder erforderlich sind, minimiert wird. Dadurch wird die Interpretation der beobachteten Felder vereinfacht.
- **Equamax.** Eine Rotationsmethode, bei der es sich um eine Kombination der Varimax-Methode, die die Faktoren vereinfacht, und der Quartimax-Methode, die die Felder vereinfacht, handelt. Die Anzahl der Felder mit hoher Ladung bei einem Faktor und die Anzahl der Faktoren, die für die Erklärung eines Felds erforderlich sind, werden minimiert.
- **Promax.** Eine schiefe Rotation, bei der die Faktoren korreliert sein dürfen. Sie lässt sich schneller berechnen als eine direkte Oblimin-Rotation, sodass sie auch für große Daten-Sets verwendet werden kann. Kappa steuert die Schiefe der Lösung (den Grad, in dem die Faktoren korreliert werden können).

## **Modell-Nugget vom Typ “Faktor/PCA”**

Ein Modell-Nugget vom Typ “Faktor/PCA” stellt das Modell der Faktorenanalyse und der Hauptkomponentenanalyse (PCA) dar, das durch einen Faktor/PCA-Knoten erstellt wurde. Sie enthalten alle Informationen, die vom trainierten Modell erfasst wurden, sowie Informationen über Leistung und Merkmale des Modells.

Wenn Sie einen Stream ausführen, der ein Faktorgleichungsmodell enthält, fügt der Knoten ein neues Feld für jeden Faktor bzw. jede Komponente im Modell hinzu. Die neuen Feldnamen werden vom Modellnamen abgeleitet, mit  $\$F$ - präfigiert und mit  $-n$  suffigiert. Dabei ist  $n$  die Nummer des Faktors bzw. der Komponente. Wenn Ihr Modell beispielsweise den Namen *Faktor* aufweist und drei Faktoren enthält, werden die neuen Felder wie folgt benannt:  $\$F$ -Faktor-1,  $\$F$ -Faktor-2 und  $\$F$ -Faktor-3.

Um besser zu verstehen, was das Faktor-Modell kodiert hat, können Sie weiter unten im Stream weitere Analysen durchführen. Eine günstige Methode zur Anzeige des Ergebnisses des Faktor-Modells besteht in der Anzeige der Korrelationen zwischen den Faktoren und den Eingabefeldern mithilfe eines Statistikknötens. Dadurch wird aufgezeigt, welche Eingabefelder welche Faktoren stark belasten, und hilft bei der Ermittlung, ob den Faktoren eine Bedeutung oder Interpretation zugrunde liegt. [Für weitere Informationen siehe Thema Statistikknötens in Kapitel 6 in IBM SPSS Modeler 14.2- Quellen-, Prozess- und Ausgabeknoten.](#)

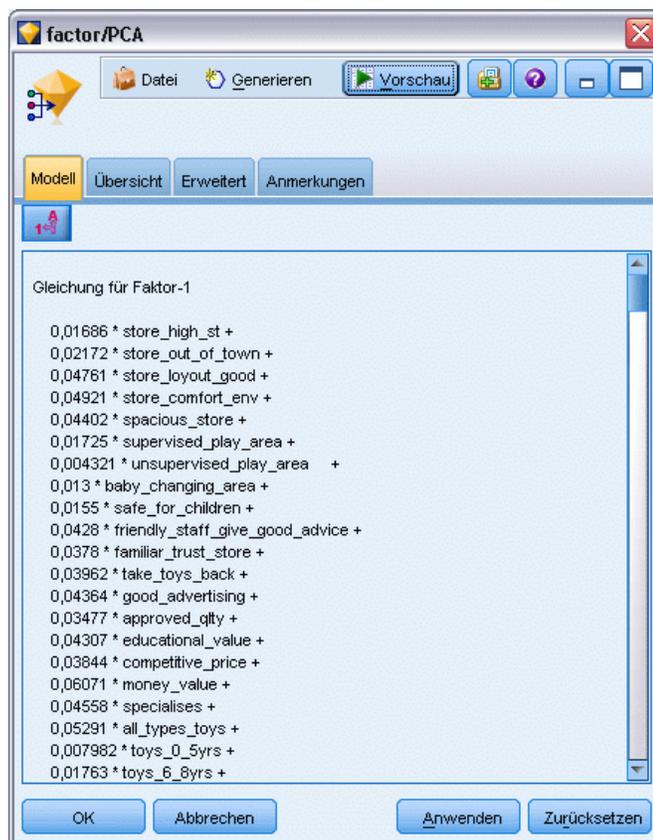
Außerdem können Sie das Faktor-Modell mit den in der erweiterten Ausgabe verfügbaren Informationen bewerten. Zur Anzeige der erweiterten Ausgabe klicken Sie im Modell-Nugget-Browser auf die Registerkarte Erweitert. Die erweiterte Ausgabe enthält zahlreiche detaillierte Informationen und ist für Benutzer mit umfassenden Kenntnissen im Bereich der

Faktoranalyse bzw. der Hauptkomponentenanalyse (PCA) gedacht. [Für weitere Informationen siehe Thema Modell-Nuggets vom Typ “Faktor/PCA” – Erweiterte Ausgabe auf S. 318.](#)

### **Modell-Nugget vom Typ “Faktor/PCA” – Gleichungen**

Auf der Registerkarte “Modell” für ein Modell-Nugget vom Typ “Faktor/PCA” wird die Faktor-Score-Gleichung für die einzelnen Faktoren angezeigt. Faktor- bzw. Komponenten-Scores werden berechnet, indem jeder Eingabefeldwert mit seinem Koeffizienten multipliziert und die Summe der Ergebnisse gebildet wird.

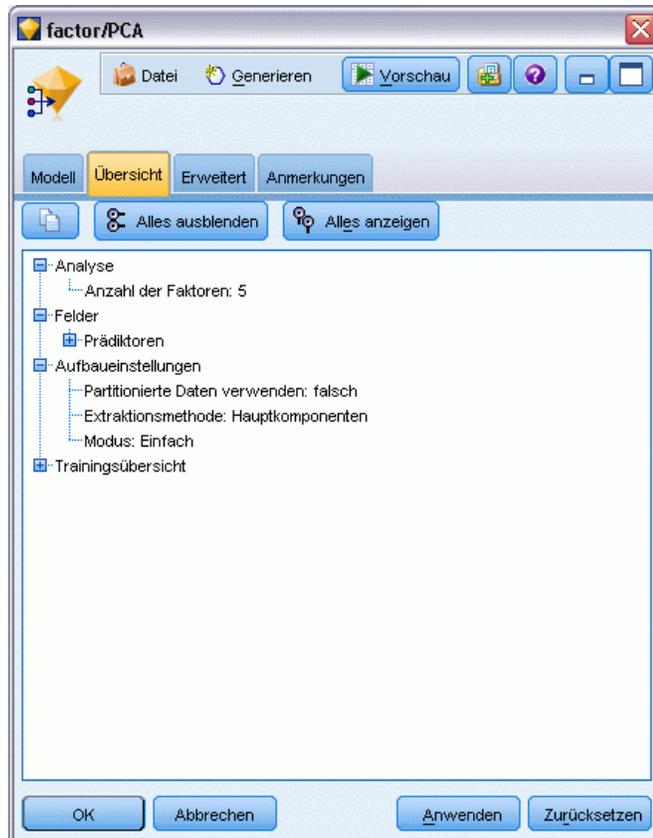
Abbildung 10-36  
PCA/Faktor-Nugget - Registerkarte “Modell”



### **Modell-Nugget vom Typ “Faktor/PCA” – Übersicht**

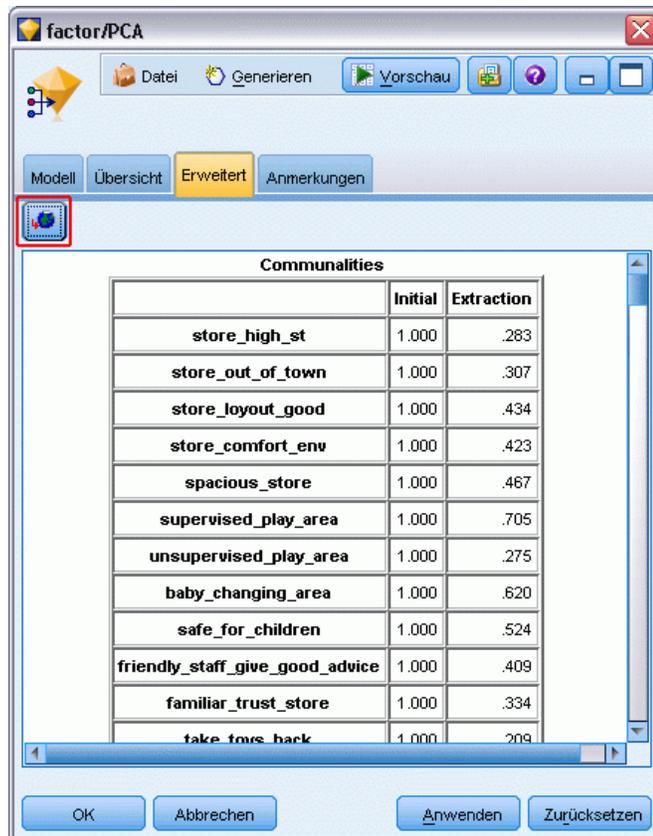
Auf der Registerkarte “Übersicht” für ein Faktor-Modell werden die Anzahl der im Faktor/PCA-Modell beibehaltenen Faktoren sowie zusätzliche Informationen zu den für die Generierung des Modells verwendeten Feldern und Einstellungen angezeigt. [Für weitere Informationen siehe Thema Durchsuchen von Modell-Nuggets in Kapitel 3 auf S. 52.](#)

Abbildung 10-37  
Beispiel-Faktorgleichungsknoten – Registerkarte "Übersicht"



## Modell-Nuggets vom Typ "Faktor/PCA" – Erweiterte Ausgabe

Abbildung 10-38  
Beispiel-Faktorgleichungsknoten – Registerkarte "Erweitert"



Die erweiterte Ausgabe für die Faktorenanalyse bietet detaillierte Informationen zum geschätzten Modell und dessen Leistung. Die meisten der in der erweiterten Ausgabe enthaltenen Informationen weisen einen hohen Fachlichkeitsgrad auf und zur richtigen Interpretation dieser Ausgaben sind umfassende Kenntnisse der Faktorenanalyse erforderlich.

**Warnungen.** Zeigt etwaige Warnungen oder potenzielle Probleme mit den Ergebnissen an.

**Kommunalitäten.** Zeigt an, welcher Anteil der Varianz der einzelnen Felder durch die Faktoren oder Komponenten erklärt wird. *Initial* gibt die Anfangskommunalitäten mit dem vollständigen Faktoren-Set aus (das Modell wird mit so vielen Faktoren gestartet, wie Eingabefelder vorhanden sind) und *Extraction* gibt die Kommunalitäten auf der Grundlage des beibehaltenen Faktoren-Sets aus.

**Erklärte Gesamtvarianz.** Zeigt die von den Faktoren im Modell erklärte Gesamtvarianz an. *Anfängliche Eigenwerte* zeigt die Varianz an, die vom vollständigen Set der Anfangsfaktoren erklärt wird. *Extrahierte Summen von quadrierten Faktorladungen* zeigt die Varianz an, die von den im Modell beibehaltenen Faktoren erklärt wird. *Rotierte Summen von quadrierten Ladungen* zeigt die Varianz an, die von den rotierten Faktoren erklärt wird. Beachten Sie, dass bei schiefen Rotationen *Rotierte Summen von quadrierten Ladungen* nur die Summen der quadrierten Ladungen und keine Varianzprozentsätze angezeigt werden.

**Faktor-(bzw. Komponenten-)Matrix** Zeigt die Korrelationen zwischen Eingabefeldern und nicht rotierten Faktoren.

**Rotierte Faktor-(bzw. Komponenten-)Matrix.** Zeigt die Korrelationen zwischen Eingabefeldern und rotierten Faktoren für orthogonale Rotationen an.

**Mustermatrix.** Zeigt die partiellen Korrelationen zwischen Eingabefeldern und rotierten Faktoren für schiefe Rotationen an.

**Strukturmatrix.** Zeigt die einfachen Korrelationen zwischen Eingabefeldern und rotierten Faktoren für schiefe Rotationen an.

**Faktorkorrelationsmatrix.** Zeigt die Korrelationen zwischen den Faktoren für schiefe Rotationen an.

## Diskriminanzknoten

Die Diskriminanzanalyse dient zur Erstellung eines Vorhersagemodells der Gruppenzugehörigkeit. Das Modell besteht aus einer Diskriminanzfunktion (oder, bei mehr als zwei Gruppen, einem Set von Diskriminanzfunktionen) auf der Grundlage derjenigen linearen Kombinationen der Einflussvariablen (Prädiktorvariablen), die die beste Diskriminanz zwischen den Gruppen ergeben. Die Funktionen werden aus einer Stichprobe der Fälle erzeugt, bei denen die Gruppenzugehörigkeit bekannt ist. Diese Funktionen können dann auf neue Fälle mit Messungen für die Prädiktorvariablen, aber unbekannter Gruppenzugehörigkeit angewandt werden.

**Beispiel.** Ein Telekommunikationsunternehmen kann mithilfe der Diskriminanzanalyse Kunden anhand der Nutzungsdaten in Gruppen einteilen. Dadurch kann das Unternehmen potenzielle Kunden scoren und sich gezielt denjenigen zuwenden, die mit der größten Wahrscheinlichkeit zu den einträglichsten Gruppen gehören. [Für weitere Informationen siehe Thema Klassifizieren von Kunden im Telekommunikationsbereich \(Diskriminanzanalyse\) in Kapitel 22 in IBM SPSS Modeler 14.2- Anwendungshandbuch.](#)

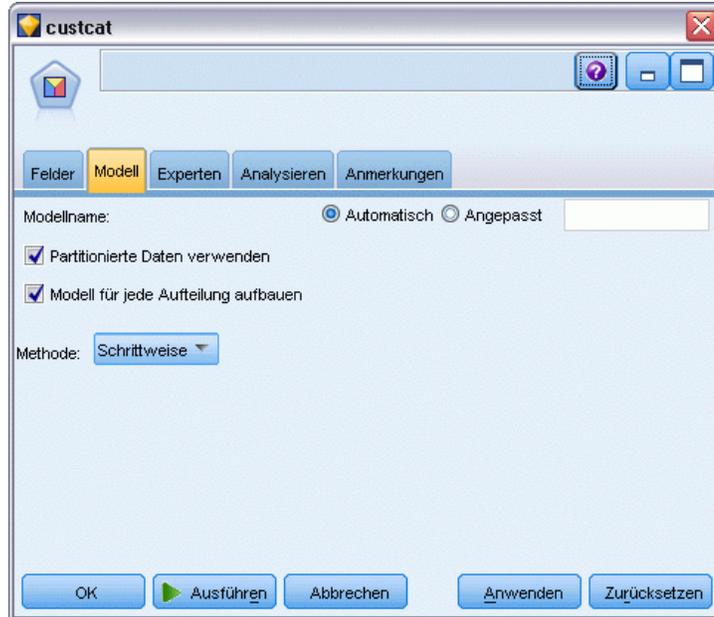
**Anforderungen.** Es werden mindestens ein Eingabefeld und genau ein Zielfeld benötigt. Bei dem Ziel muss es sich um ein kategoriales Feld (mit dem Messniveau *Flag* oder *Nominal*) mit dem Speichertyp "Zeichenkette" oder "Ganze Zahl" handeln. (Der Speichertyp kann, falls erforderlich, mithilfe eines Füller- oder Ableitungsknotens konvertiert werden. [Für weitere Informationen siehe Thema Speichertypkonvertierung mithilfe des Füllerknotens in Kapitel 4 in IBM SPSS Modeler 14.2- Quellen-, Prozess- und Ausgabeknoten.](#)) Felder, die auf *Beides* oder *Keine* gesetzt sind, werden ignoriert. Bei den im Modell verwendeten Feldern müssen die Typen vollständig instanziiert sein.

**Stärken.** Sowohl die Diskriminanzanalyse als auch die logistische Regression eignen sich jeweils als Klassifizierungsmodell. Die Diskriminanzanalyse geht jedoch von mehr Annahmen über die Eingabefelder aus, beispielsweise davon, dass sie normalverteilt sind und stetig sein sollten, und bietet bessere Ergebnisse, wenn diese Anforderungen erfüllt sind, insbesondere bei kleinem Stichprobenumfang.

## Diskriminanzknoten – Modelloptionen

Abbildung 10-39

Dialogfeld des Diskriminanzknotens, Registerkarte "Modell"



**Modellname.** Sie können den Modellnamen automatisch basierend auf den Ziel- oder ID-Feldnamen (oder dem Modelltyp in Fällen, in denen kein solches Feld angegeben ist) generieren oder einen benutzerdefinierten Namen eingeben.

**Partitionierte Daten verwenden.** Wenn ein Partitionsfeld definiert ist, gewährleistet diese Option, dass nur Daten aus der Trainingspartition für die Modellerstellung verwendet werden. [Für weitere Informationen siehe Thema Partitionsknoten in Kapitel 4 in IBM SPSS Modeler 14.2- Quellen-, Prozess- und Ausgabeknoten.](#)

**Aufteilungsmodelle erstellen.** Erstellt ein separates Modell für jeden möglichen Wert von Eingabefeldern, die als Aufteilungsfelder angegeben werden. [Für weitere Informationen siehe Thema Erstellung von aufgeteilten Modellen in Kapitel 3 auf S. 31.](#)

**Methode.** Zur Eingabe von Prädiktoren in das Modell stehen folgende Optionen zur Verfügung:

- **Eingabetaste.** Dies ist das Standardverfahren, bei dem alle Terme direkt in die Gleichung aufgenommen werden. Terme, die nicht in signifikanter Weise zur Vorhersagekraft des Modells beitragen, werden nicht hinzugefügt.
- **Schrittweise.** Das anfängliche Modell ist das einfachste Modell, das möglich ist. Es enthält keine Modellterme (außer der Konstanten) in der Gleichung. Bei jedem Schritt werden die Terme, die noch nicht zum Modell hinzugefügt wurden, bewertet, und wenn der beste dieser Terme in signifikanter Weise zur Vorhersagekraft des Modells beiträgt, wird er hinzugefügt.

*Hinweis:* Die Methode "Schrittweise" weist eine starke Tendenz zur übermäßigen Anpassung an die Trainingsdaten auf. Bei der Verwendung dieser Methoden ist es ganz besonders wichtig, die Validität des entstehenden Modells mithilfe einer zurückgehaltenen Teststichprobe oder mit neuen Daten zu überprüfen.

## Diskriminanzknoten – Expertenoptionen

Wenn Sie über umfassende Kenntnisse im Bereich der Diskriminanzanalyse verfügen, können Sie mithilfe der Expertenoptionen eine Feinabstimmung des Trainingsvorgangs vornehmen. Für den Zugriff auf die Expertenoptionen stellen Sie auf der Registerkarte “Experten” Modus auf Experten ein.

Abbildung 10-40  
Dialogfeld des Diskriminanzknotens, Registerkarte “Experten”



**A-priori-Wahrscheinlichkeiten.** Diese Option bestimmt, ob die Klassifikationskoeffizienten für A-priori-Kennntnis über Gruppenzugehörigkeiten angepasst werden.

- **Alle Gruppen gleich.** Es wird von gleichen A-priori-Wahrscheinlichkeiten für alle Gruppen ausgegangen; dies hat keine Auswirkungen auf die Koeffizienten.
- **Von Gruppengrößen berechnen.** Die beobachteten Gruppengrößen in Ihrem Beispiel bestimmen die A-priori-Wahrscheinlichkeiten der Gruppenzugehörigkeit. Falls beispielsweise 50 % der in der Analyse aufgenommenen Beobachtungen in die erste, 25 % in die zweite und 25 % in die dritte Gruppe fallen, werden die Klassifikationskoeffizienten angepasst, um die Wahrscheinlichkeit der Zugehörigkeit zur ersten Gruppe in Bezug auf die anderen beiden Gruppen zu erhöhen.

**Kovarianzmatrix verwenden.** Sie können wählen, ob zur Klassifikation der Fälle die Kovarianzmatrix innerhalb der Gruppen oder die gruppenspezifische Kovarianzmatrix verwendet werden soll.

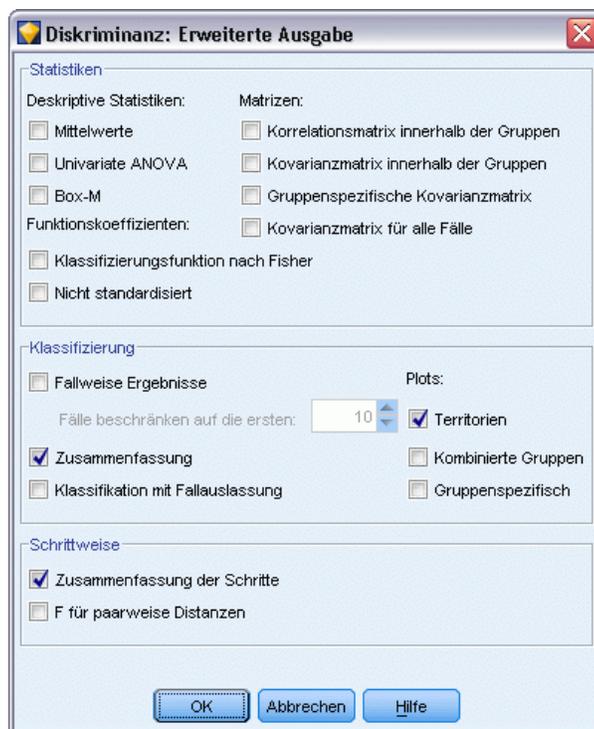
- **Innerhalb der Gruppen.** Zur Klassifizierung von Fällen wird die gemeinsame Kovarianzmatrix innerhalb der Gruppen verwendet.
- **Gruppenspezifisch.** Für die Klassifizierung werden gruppenspezifische Kovarianzmatrizen verwendet. Da die Klassifizierung auf Diskriminanzfunktionen und nicht auf ursprünglichen Variablen basiert, entspricht diese Option nicht immer der Verwendung einer quadratischen Diskriminanzfunktion.

**Ausgabe.** Mit diesen Optionen können Sie zusätzliche Statistiken anfordern, die in der erweiterten Ausgabe des vom Knoten erstellten Modell-Nugget angezeigt werden. [Für weitere Informationen siehe Thema Diskriminanzknoten – Ausgabeoptionen auf S. 322.](#)

**Kriterien.** Mit diesen Optionen können Sie die Kriterien zum Hinzufügen und Entfernen von Feldern mit der Schätzmethode “Schrittweise” festlegen. (Die Schaltfläche ist deaktiviert, wenn die Methode “Einschluss” ausgewählt ist.) [Für weitere Informationen siehe Thema Diskriminanzknoten – Schrittoptionen auf S. 324.](#)

## Diskriminanzknoten – Ausgabeoptionen

Abbildung 10-41  
Diskriminanzknoten – Erweiterte Ausgabeoptionen



Wählen Sie die optionalen Ausgaben aus, die in der erweiterten Ausgabe des Modell-Nuggets vom Typ “Logistische Regression” angezeigt werden sollen. Zur Anzeige der erweiterten Ausgabe durchsuchen Sie das Modell-Nugget und klicken Sie auf die Registerkarte Erweitert. [Für weitere Informationen siehe Thema Diskriminanz-Modell-Nugget – Erweiterte Ausgabe auf S. 326.](#)

**Deskriptive Statistiken.** Verfügbare Optionen sind Mittelwerte (einschließlich Standardabweichungen), univariate ANOVA und Box’ *M*-Test.

- **Mittelwerte.** Zeigt Gesamt- und Gruppenmittelwerte sowie Standardabweichungen für die unabhängigen Variablen an.

- **Univariate ANOVA.** Führt für jede unabhängige Variable eine einfaktorische Varianzanalyse durch, d. h. einen Test auf Gleichheit der Gruppenmittelwerte.
- **Box-M.** Ein Test auf Gleichheit der Kovarianzmatrizen der Gruppen. Bei hinreichend großen Stichproben bedeutet ein nichtsignifikanter p-Wert, dass die Anhaltspunkte für unterschiedliche Matrizen nicht ausreichend sind. Der Test ist empfindlich gegenüber Abweichungen von der multivariaten Normalverteilung.

**Funktionskoeffizienten.** Verfügbare Optionen sind Klassifikationskoeffizienten nach Fisher und nicht standardisierte Koeffizienten.

- **Fisher.** Zeigt die Koeffizienten der Klassifizierungsfunktion nach Fisher an, die direkt für die Klassifizierung verwendet werden können. Es wird ein eigenes Set von Koeffizienten der Klassifizierungsfunktion für jede Gruppe ermittelt. Ein Fall wird der Gruppe zugewiesen, für die er den größten Diskriminanzwert (Klassifizierungsfunktionswert) aufweist.
- **Nichtstandardisiert.** Zeigt die nichtstandardisierten Koeffizienten der Diskriminanzfunktion an.

**Matrizen.** Als Koeffizientenmatrizen für unabhängige Variablen stehen die Korrelationsmatrix innerhalb der Gruppen, die Kovarianzmatrix innerhalb der Gruppen, die gruppenspezifische Kovarianzmatrix und die Kovarianzmatrix für alle Fälle zur Verfügung.

- **Korrelationsmatrix innerhalb der Gruppen.** Zeigt eine gemeinsame Korrelationsmatrix innerhalb der Gruppen an, die als Mittel der separaten Kovarianzmatrizen für alle Gruppen vor der Berechnung der Korrelationen bestimmt wird.
- **Kovarianzmatrix innerhalb der Gruppen.** Zeigt eine gemeinsame Kovarianzmatrix innerhalb der Gruppen an, die sich von der Gesamt-Kovarianzmatrix unterscheiden kann. Die Matrix wird als Mittel der einzelnen Kovarianzmatrizen für alle Gruppen berechnet.
- **Gruppenspezifische Kovarianzmatrix.** Zeigt separate Kovarianzmatrizen für jede Gruppe an.
- **Kovarianzmatrix für alle Fälle.** Zeigt die Kovarianzmatrix für alle Fälle an, so als wären sie aus einer einzigen Stichprobe.

**Klassifikation.** Folgende Ausgaben gehören zu den Klassifikationsergebnissen.

- **Fallweise Ergebnisse.** Für jeden Fall werden Codes für die tatsächliche Gruppe, die vorhergesagte Gruppe, A-posteriori-Wahrscheinlichkeiten und Diskriminanzwerte angezeigt.
- **Zusammenfassende Tabelle.** Die Anzahl der Fälle, die auf Grundlage der Diskriminanzanalyse jeder der Gruppen richtig oder falsch zugeordnet werden. Zuweilen auch als Klassifikationsmatrix bezeichnet.
- **Klassifikation mit Fallauslassung.** Jeder Fall der Analyse wird durch Funktionen aus allen anderen Fällen unter Auslassung dieses Falls klassifiziert. Diese Klassifikation wird auch als "U-Methode" bezeichnet.
- **Territorien.** Ein Diagramm der Grenzen, mit denen Fälle auf der Grundlage von Funktionswerten in Gruppen klassifiziert werden. Die Zahlen entsprechen den Gruppen, in die die Fälle klassifiziert wurden. Der Mittelwert jeder Gruppe wird durch einen darin liegenden Stern (\*) angezeigt. Dieses Diagramm wird nicht angezeigt, wenn nur eine einzige Diskriminanzfunktion vorliegt.

- **Kombinierte Gruppen.** Erzeugt ein alle Gruppen umfassendes Streudiagramm der Werte für die ersten beiden Diskriminanzfunktionen. Wenn nur eine Funktion vorliegt, wird stattdessen ein Histogramm angezeigt.
- **Gruppenspezifisch.** Erzeugt gruppenspezifische Streudiagramme der Werte für die ersten beiden Diskriminanzfunktionen. Wenn nur eine Funktion vorliegt, werden stattdessen Histogramme angezeigt.

**Schrittweise.** Zusammenfassung der Schritte zeigt nach jedem Schritt Statistiken für alle Variablen an; F für paarweise Distanzen zeigt eine Matrix mit paarweisen  $F$ -Quotienten für jedes Gruppenpaar an. Die  $F$ -Quotienten können für Signifikanztests der Mahalanobis-Abstände zwischen Gruppen verwendet werden.

### Diskriminanzknoten – Schrittoptionen

Abbildung 10-42

Diskriminanzknoten – Optionen für die Methode "Schrittweise"



**Methode.** Wählen Sie die Statistiken aus, die für die Aufnahme oder den Ausschluss neuer Variablen dienen sollen. Die Optionen Wilks-Lambda, nicht erklärte Varianz, Mahalanobis-Abstand, kleinster  $F$ -Quotient und Rao- $V$  stehen zur Verfügung. Mit Rao- $V$  können Sie den Mindestanstieg von  $V$  für eine einzugebende Variable angeben.

- **Wilks-Lambda.** Eine Auswahlmethode für Variablen bei der schrittweisen Diskriminanzanalyse. Die Aufnahme von Variablen in die Gleichung erfolgt anhand der jeweiligen Verringerung von Wilks-Lambda. Bei jedem Schritt wird diejenige Variable aufgenommen, die den Gesamtwert von Wilks-Lambda am meisten vermindert.
- **Nicht erklärte Varianz.** Bei jedem Schritt wird die Variable aufgenommen, welche die Summe der nicht erklärten Variation zwischen den Gruppen minimiert.
- **Mahalanobis-Abstand.** Dieses Maß gibt an, wie weit die Werte der unabhängigen Variablen eines Falles vom Mittelwert aller Fälle abweichen. Ein großer Mahalanobis-Abstand charakterisiert einen Fall, der bei einer oder mehreren unabhängigen Variablen Extremwerte besitzt.

- **Kleinster F-Quotient.** Eine Methode für die Variablenauswahl in einer schrittweisen Analyse. Sie beruht auf der Maximierung eines F-Quotienten, der aus dem Mahalanobis-Abstand zwischen den Gruppen errechnet wird.
  - **Rao-V.** Ein Maß für die Unterschiede zwischen Gruppenmittelwerten. Auch Lawley-Hotelling-Spur genannt. Bei jedem Schritt wird die Variable aufgenommen, die den Anstieg des Rao-V maximiert. Wenn Sie diese Option ausgewählt haben, geben Sie den Minimalwert ein, den eine Variable für die Aufnahme in die Analyse aufweisen muss.
- Kriterien.** Verfügbar sind F-Wert verwenden und F-Wahrscheinlichkeit verwenden. Geben Sie Werte für die Aufnahme und den Ausschluss der Variablen an.
- **F-Wert verwenden.** Eine Variable wird in ein Modell aufgenommen, wenn ihr F-Wert größer ist als der Aufnahmewert. Sie wird ausgeschlossen, wenn der F-Wert kleiner ist als der Ausschlusswert. Der Aufnahmewert muss größer sein als der Ausschlusswert und beide Werte müssen positiv sein. Um mehr Variablen in das Modell aufzunehmen, senken Sie den Aufnahmewert. Um mehr Variablen aus dem Modell auszuschließen, erhöhen Sie den Ausschlusswert.
  - **Wahrscheinlichkeit von F verwenden.** Eine Variable wird in das Modell aufgenommen, wenn das Signifikanzniveau ihres F-Werts kleiner ist als der Aufnahmewert. Sie wird ausgeschlossen, wenn das Signifikanzniveau größer ist als der Ausschlusswert. Der Aufnahmewert muss kleiner sein als der Ausschlusswert und beide Werte müssen positiv sein. Um mehr Variablen in das Modell aufzunehmen, erhöhen Sie den Aufnahmewert. Um mehr Variablen aus dem Modell auszuschließen, senken Sie den Ausschlusswert.

## Diskriminanz-Modell-Nugget

Modell-Nuggets vom Typ “Diskriminanz” stehen für die durch Diskriminanzknoten geschätzten Gleichungen. Sie enthalten alle Informationen, die vom Diskriminanzmodell erfasst wurden, sowie Informationen über die Struktur und Leistung des Modells.

Wenn Sie einen Stream ausführen, der ein Modell-Nugget vom Typ “Diskriminanz” enthält, fügt der Knoten zwei neue Felder hinzu, die die Prognose des Modells und die zugehörige Wahrscheinlichkeit enthalten. Die Namen der neuen Felder werden aus dem Namen des prognostizierten Ausgabefelds abgeleitet, dem  $\$D$ - für die vorhergesagte Kategorie und  $\$DP$ - für die zugehörige Wahrscheinlichkeit vorangestellt ist. Bei einem Ausgabefeld mit der Bezeichnung *Farbpräf* beispielsweise erhalten die neuen Felder die Namen  $\$D$ -*Farbpräf* und  $\$DP$ -*Farbpräf*.

**Erstellen eines Filterknotens.** Im Menü “Generieren” können Sie einen neuen Filterknoten zur Übergabe der Eingabefelder auf der Grundlage der Ergebnisse erstellen.

### Bedeutsamkeit des Prädiktors

Optional kann auf der Registerkarte “Modell” auch ein Diagramm, das die relative Bedeutsamkeit der einzelnen Prädiktoren für die Schätzung des Modells angibt, angezeigt werden. Normalerweise ist es sinnvoll, die Modellierungsbemühungen auf die wichtigsten Prädiktoren zu konzentrieren und diejenigen Prädiktoren zu verwerfen bzw. zu ignorieren, die die geringste Bedeutung haben. Beachten Sie, dass dieses Diagramm nur verfügbar ist, wenn vor dem Generieren des Modells Bedeutsamkeit der Prädiktoren berechnen auf der Registerkarte “Analysieren” ausgewählt wurde. [Für weitere Informationen siehe Thema Bedeutsamkeit des Prädiktors in Kapitel 3 auf S. 54.](#)

## Diskriminanz-Modell-Nugget – Erweiterte Ausgabe

Abbildung 10-43  
Diskriminanz-Modell-Nugget – Registerkarte “Erweitert”

Analysis Case Processing Summary		
<b>Unweighted Cases</b>		
	<b>Valid</b>	<b>N</b> <b>Percent</b>
		1000 100.0
<b>Excluded</b>	<b>Missing or out-of-range group codes</b>	0 .0
	<b>At least one missing discriminating variable</b>	0 .0
	<b>Both missing or out-of-range group codes and at least one missing discriminating variable</b>	0 .0
	<b>Total</b>	0 .0
<b>Total</b>		1000 100.0
<b>Group Statistics</b>		
<b>Customer category</b>	<b>Geographic indicator</b>	<b>Valid N (listwise)</b>
		<b>Unweighted</b> <b>Weighted</b>
		266 266.000

Die erweiterte Ausgabe für die Diskriminanzanalyse bietet detaillierte Informationen zum geschätzten Modell und dessen Leistung. Die meisten der in der erweiterten Ausgabe enthaltenen Informationen weisen einen hohen Fachlichkeitsgrad auf und zur richtigen Interpretation dieser Ausgaben sind umfassende Kenntnisse der Diskriminanzanalyse erforderlich. [Für weitere Informationen siehe Thema Diskriminanzknoten – Ausgabeoptionen auf S. 322.](#)

## Diskriminanz-Modell-Nugget – Einstellungen

Auf der Registerkarte “Einstellungen” für ein Modell-Nugget vom Typ “Diskriminanz” können Sie beim Scoring des Modells Neigungs-Scores ermitteln. Diese Registerkarte ist nur für Modelle mit Flag-Zielen verfügbar und erst nachdem das Modell-Nugget einem Stream hinzugefügt wurde.

Abbildung 10-44  
Diskriminanz-Modell-Nugget, Registerkarte "Einstellungen" für Flag-Ziele



**Scores für Rohneigung berechnen.** Bei Modellen mit einem Flag-Ziel (das als Vorhersage “Ja” bzw. “Nein” ausgibt) können Sie Neigungs-Scores anfordern, die die Wahrscheinlichkeit des für das Zielfeld angegebenen wahren Ergebnisses angeben. Diese Werte werden zusätzlich zu den anderen Vorhersage- und Konfidenzwerten angegeben, die ggf. während des Scorings erstellt werden.

**Scores für korrigierte Neigung berechnen.** Rohneigungs-Scores beruhen ausschließlich auf den Trainingsdaten und sind aufgrund der Neigung vieler Modelle zur übermäßigen Anpassung an die Trainingsdaten möglicherweise zu optimistisch. Bei korrigierten Neigungen wird versucht, dies durch Evaluation der Modellleistung anhand einer Test- bzw. Validierungspartition zu kompensieren. Bei dieser Option muss im Stream ein Partitionsfeld definiert sein und die Scores für die korrigierte Neigung müssen im Modellierungsknoten aktiviert werden, bevor das Modell generiert wird.

### ***Diskriminanz-Modell-Nugget – Übersicht***

Auf der Registerkarte “Übersicht” für ein Modell-Nugget vom Typ “Diskriminanz” werden die Felder und Einstellungen angezeigt, die zum Generieren des Modells verwendet wurden. Wenn Sie einen Analyseknoden ausgeführt haben, der an diesen Modellierungsknoten angehängt ist, werden die Informationen aus dieser Analyse ebenfalls in diesem Abschnitt angezeigt. [Für weitere Informationen siehe Thema Analyseknoden in Kapitel 6 in IBM SPSS Modeler 14.2- Quellen-, Prozess- und Ausgabeknoten.](#) Allgemeine Informationen zur Verwendung des Modellbrowsers finden Sie unter [Durchsuchen von Modell-Nuggets auf S. 52.](#)

Abbildung 10-45  
Diskriminanz-Modell-Nugget – Registerkarte "Übersicht"



## GenLin-Knoten

Das verallgemeinerte lineare Modell erweitert das allgemeine lineare Modell, sodass die abhängige Variable über eine angegebene Linkfunktion in einer linearen Relation zu den Faktoren und Kovariaten steht. Darüber hinaus kann die abhängige Variable bei diesem Modell eine nichtnormale Verteilung aufweisen. Es deckt durch seine sehr allgemein gehaltene Modellformulierung häufig verwendete statistische Modelle ab, wie beispielsweise die lineare Regression für normalverteilte Antworten, logistische Modelle für binäre Daten und loglineare Modelle für Häufigkeitsdaten, Modelle vom Typ "Log-Log komplementär" für intervallzensierte Überlebensdaten sowie viele andere statistische Modelle.

**Beispiele.** Eine Reederei kann verallgemeinerte lineare Modelle verwenden, um eine Poisson-Regression auf die Anzahl der Havarien für mehrere Schiffstypen anzuwenden, die in verschiedenen Zeiträumen gebaut wurden. Anhand des so entstandenen Modells kann ermittelt werden, welche Schiffstypen am havarieanfälligsten sind. [Für weitere Informationen siehe Thema Verwenden der Poisson-Regression für die Analyse von Schiffsschadensraten \(Verallgemeinerte lineare Modelle\) in Kapitel 24 in IBM SPSS Modeler 14.2- Anwendungshandbuch.](#)

Ein KFZ-Versicherungsunternehmen kann mithilfe von verallgemeinerten linearen Modellen eine Gammaregression an die Schadensansprüche für Autos anpassen. Anhand des so entstandenen Modells können die Faktoren ermittelt werden, die am meisten zur Anspruchshöhe beitragen. [Für weitere Informationen siehe Thema Anpassen einer Gamma-Regression an Versicherungsforderungen an Kfz-Versicherungen \(Verallgemeinerte lineare Modelle\) in Kapitel 25 in IBM SPSS Modeler 14.2- Anwendungshandbuch.](#)

Medizinforscher können mithilfe von verallgemeinerten linearen Modellen eine komplementäre Log-Log-Regression für intervallzensierte Überlebensdaten anpassen, um die Dauer bis zum Wiederauftreten eines Krankheitsbilds vorherzusagen. [Für weitere Informationen siehe Thema Analysieren von intervallzensierten Überlebensdaten \(Verallgemeinerte lineare Modelle\) in Kapitel 23 in IBM SPSS Modeler 14.2- Anwendungshandbuch.](#)

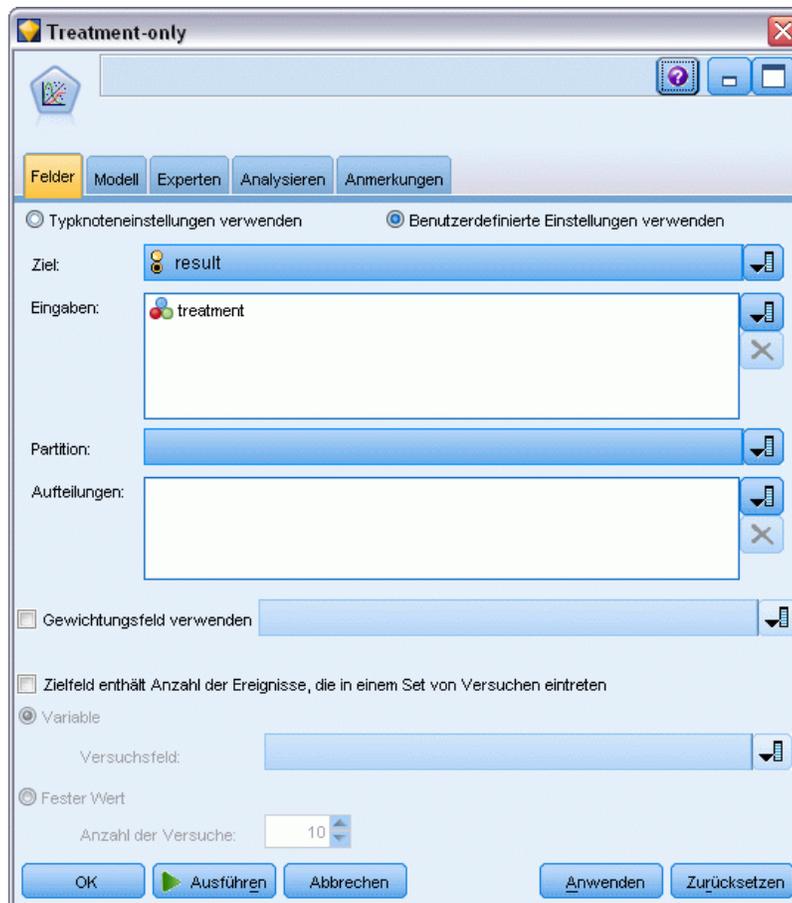
Bei verallgemeinerten linearen Modellen wird eine Gleichung erstellt, die die Werte der Eingabefelder mit den Werten der Ausgabefelder in Bezug setzt. Nach der Generierung des Modells kann es zur Schätzung der Werte für neue Daten verwendet werden. Für jeden Datensatz wird eine Wahrscheinlichkeit der Zugehörigkeit für jede mögliche Ausgabekategorie berechnet. Die Zielkategorie mit der höchsten Wahrscheinlichkeit wird als vorhergesagter Ausgabewert für den betreffenden Datensatz zugewiesen.

**Anforderungen.** Es werden mindestens ein Eingabefeld und genau ein Zielfeld (mit dem Messniveau *Stetig* oder *Flag*) mit mindestens zwei Kategorien benötigt. Bei den im Modell verwendeten Feldern müssen die Typen vollständig instanziiert sein.

**Stärken.** Das verallgemeinerte lineare Modell ist extrem flexibel, jedoch ist das Verfahren für die Auswahl der Modellstruktur nicht automatisiert, weshalb ein Grad an Vertrautheit mit den Daten nötig ist, der bei einem “Black Box”-Algorithmus nicht erforderlich ist.

## Feldoptionen für den GenLin-Knoten

Abbildung 10-46  
Dialogfeld des GenLin-Knotens, Registerkarte "Felder"



Neben den benutzerdefinierten Optionen für Ziel, Eingabe und Partition, die normalerweise auf der Registerkarte "Felder" von Modellierungsknoten verfügbar sind (siehe [Feldoptionen der Modellierungsknoten](#) auf S. 36), bietet der Knoten "GenLin" folgende Zusatzfunktionen.

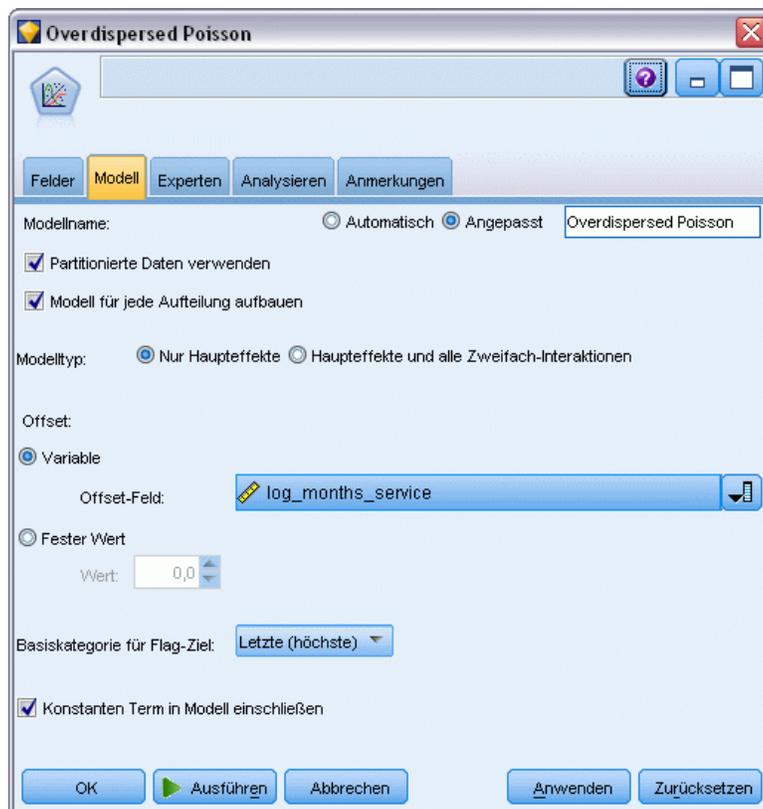
**Gewichtungsfeld verwenden.** Der Skalenparameter ist ein geschätzter Modellparameter, der mit der Varianz der Antwort zusammenhängt. Die Skalengewichte sind "bekannte" Werte, die sich zwischen den einzelnen Beobachtungen unterscheiden können. Wenn die Skalengewichtsvariable angegeben ist, wird der Skalenparameter, der mit der Varianz der Antwort zusammenhängt, für jede Beobachtung durch diese Variable geteilt. Datensätze, bei denen die Werte für das Skalengewicht kleiner oder gleich 0 sind oder fehlen, werden nicht in der Analyse verwendet.

**Zielfeld enthält Anzahl der Ereignisse, die in einem Set von Versuchen eintreten.** Wenn es sich bei der Antwort um eine Reihe von Ereignissen handelt, die in einem Set von Versuchen eintreten, enthält das Zielfeld die Anzahl der Ereignisse und Sie können eine zusätzliche Variable auswählen, die die Anzahl der Versuche enthält. Wenn die Anzahl der Versuche über alle Subjekte gleich ist, können die Versuche alternativ auch über einen festen Wert angegeben werden. Die Anzahl der Versuche sollte größer oder gleich der Anzahl der Ereignisse für die einzelnen Datensätze

sein. Bei den Ereignissen sollte es sich um nichtnegative Ganzzahlen und bei den Versuchen um positive Ganzzahlen handeln.

## Modelloptionen für den GenLin-Knoten

Abbildung 10-47  
Dialogfeld des GenLin-Knotens, Registerkarte "Modell"



**Modellname.** Sie können den Modellnamen automatisch basierend auf den Ziel- oder ID-Feldnamen (oder dem Modelltyp in Fällen, in denen kein solches Feld angegeben ist) generieren oder einen benutzerdefinierten Namen eingeben.

**Partitionierte Daten verwenden.** Wenn ein Partitionsfeld definiert ist, gewährleistet diese Option, dass nur Daten aus der Trainingspartition für die Modellerstellung verwendet werden. [Für weitere Informationen siehe Thema Partitions-knoten in Kapitel 4 in IBM SPSS Modeler 14.2- Quellen-, Prozess- und Ausgabeknoten.](#)

**Aufteilungsmodelle erstellen.** Erstellt ein separates Modell für jeden möglichen Wert von Eingabefeldern, die als Aufteilungsfelder angegeben werden. [Für weitere Informationen siehe Thema Erstellung von aufgeteilten Modellen in Kapitel 3 auf S. 31.](#)

**Modelltyp.** Es gibt zwei Optionen für den zu erstellenden Modelltyp. Haupteffekte sorgt dafür, dass das Modell nur die einzelnen Eingabefelder enthält und nicht die Interaktionen (multiplikativen Effekte) zwischen den Eingabefeldern testet. Haupteffekte und alle Zweifach-Interaktionen umfasst alle Zwei-Wege-Interaktionen sowie die Haupteffekte der Eingabefelder.

**Offset.** Der Term “Offset” ist eine “strukturelle” Einflussvariable (Prädiktor). Ihr Koeffizient wird nicht vom Modell geschätzt, sondern es wird davon ausgegangen, dass er den Wert 1 aufweist. Daher werden die Werte des Offsets einfach zur linearen Einflussvariablen der abhängigen Variablen addiert. Dies ist besonders nützlich bei Poisson-Regressionsmodellen, bei denen die verschiedenen Fälle dem relevanten Ereignis unterschiedlich stark ausgesetzt sein können. Beispielsweise gibt es bei der Modellierung der Unfallraten für einzelne Fahrer einen wichtigen Unterschied zwischen einem Fahrer, der in 3 Jahren Fahrpraxis einen Unfall verursacht hat und einem Fahrer, der in 25 Jahren einen Unfall verursacht hat. Die Anzahl der Unfälle kann als poissonverteilte Response-Variable modelliert werden, wenn die Länge der Fahrpraxis als Offset-Term einbezogen wird.

*Anmerkung:* Bei Verwendung eines variablen Offset-Felds sollte das angegebene Feld nicht auch als Eingabe verwendet werden. Setzen Sie, falls erforderlich, die Rolle des Offset-Felds in einem Quellen- oder Typknoten weiter oben im Stream auf Keine. [Für weitere Informationen siehe Thema Festlegen der Feldrolle in Kapitel 4 in IBM SPSS Modeler 14.2- Quellen-, Prozess- und Ausgabeknoten.](#)

#### **Basiskategorie für Flag-Ziel.**

Bei binären Antworten können Sie die Referenzkategorie für die abhängige Variable auswählen. Dies kann sich auf bestimmte Ausgaben, wie beispielsweise Parameterschätzer und gespeicherte Werte, auswirken, sollte jedoch nicht die Anpassungsgüte des Modells verändern. Beispiel: Angenommen, Ihre binäre Antwort nimmt die Werte 0 und 1 an:

- Standardmäßig verwendet die Prozedur die letzte Kategorie (die mit dem höchsten Wert), also 1, als Referenzkategorie. In dieser Situation wird anhand der vom Modell gespeicherten Wahrscheinlichkeiten die Wahrscheinlichkeit geschätzt, mit der ein bestimmter Fall den Wert 0 annimmt, und die Parameterschätzer sollten als in Beziehung zur Likelihood der Kategorie 0 stehend interpretiert werden.
- Wenn Sie die erste Kategorie (die mit dem niedrigsten Wert), also 0, als Referenzkategorie angeben, wird anhand der im Modell gespeicherten Wahrscheinlichkeitswerte die Wahrscheinlichkeit geschätzt, dass ein bestimmter Fall den Wert 1 annimmt.
- Wenn Sie die benutzerdefinierte Kategorie angeben und für Ihre Variable Labels definiert sind, können Sie die Referenzkategorie durch Auswahl eines Werts aus der Liste festlegen. Dies kann nützlich sein, wenn Sie bei der Festlegung eines Modells nicht mehr wissen, wie genau eine bestimmte Variable kodiert war.

**Konstanter Term in Modell einschließen.** Der konstante Term wird gewöhnlich in das Modell aufgenommen. Wenn anzunehmen ist, dass die Daten durch den Koordinatenursprung verlaufen, können Sie den konstanten Term ausschließen.

### ***Expertenoptionen für den GenLin-Knoten***

Wenn Sie über umfassende Kenntnisse im Bereich der verallgemeinerten linearen Modelle verfügen, können Sie mithilfe der Expertenoptionen eine Feinabstimmung des Trainingsvorgangs vornehmen. Für den Zugriff auf die Expertenoptionen stellen Sie auf der Registerkarte “Experten” Modus auf Experten ein.

Abbildung 10-48  
Dialogfeld des GenLin-Knotens, Registerkarte "Experten"

### **Verteilung im Zielfeld und Linkfunktion**

#### **Verteilung.**

Diese Auswahl gibt die Verteilung der abhängigen Variablen an. Die Möglichkeit einer anderen Verteilung als "Normal" und einer anderen Verknüpfungsfunktion als "Identität" ist die wichtigste Verbesserung des verallgemeinerten linearen Modells gegenüber dem allgemeinen linearen Modell. Es gibt zahlreiche mögliche Kombinationen aus Verteilung und Verknüpfungsfunktion und es können mehrere davon für das jeweils vorliegende Daten-Set geeignet sein. Daher können Sie sich in Ihrer Wahl durch theoretische Vorüberlegungen leiten lassen oder davon, welche Kombination am besten zu passen scheint.

- **Binomial.** Diese Verteilung ist nur für Variablen geeignet, die eine binäre Antwort oder eine Anzahl von Ereignissen repräsentieren.

- **Gamma.** Diese Verteilung eignet sich für Variablen mit positiven Skalenwerten, die in Richtung größerer positiver Werte verzerrt sind. Wenn ein Datenwert kleiner oder gleich 0 ist oder fehlt, wird der entsprechende Fall nicht in der Analyse verwendet.
- **Invers normal.** Diese Verteilung eignet sich für Variablen mit positiven Skalenwerten, die in Richtung größerer positiver Werte verzerrt sind. Wenn ein Datenwert kleiner oder gleich 0 ist oder fehlt, wird der entsprechende Fall nicht in der Analyse verwendet.
- **Negativ binomial.** Diese Verteilung lässt sich als Anzahl der Versuche betrachten, die erforderlich sind, um  $k$  Erfolge zu beobachten, und eignet sich für Variablen mit nichtnegativen ganzzahligen Werten. Wenn ein Datenwert keine ganze Zahl oder kleiner als 0 ist oder fehlt, wird der entsprechende Fall nicht in der Analyse verwendet. Der feste Wert des Hilfsparameters der negativen Binomialverteilung kann jede beliebige Zahl größer oder gleich 0 sein. Wenn der Hilfsparameter auf 0 gesetzt wird, entspricht die Verwendung dieser Verteilung der Verwendung der Poisson-Verteilung.
- **Normal.** Diese Option eignet sich für metrische Variablen, deren Werte eine symmetrische, glockenförmige Verteilung um einen Mittelwert aufweisen. Die abhängige Variable muss numerisch sein.
- **Poisson.** Diese Verteilung lässt sich als Anzahl der Vorkommnisse eines untersuchten Ereignisses in einem festen Zeitraum betrachten und eignet sich für Variablen mit nichtnegativen ganzzahligen Variablen mit nichtnegativen ganzzahligen Werten. Wenn ein Datenwert keine ganze Zahl oder kleiner als 0 ist oder fehlt, wird der entsprechende Fall nicht in der Analyse verwendet.
- **Tweedie.** Diese Verteilung eignet sich für Variablen, die durch Poisson-Mischungen von Gamma-Verteilungen repräsentiert werden können. Die Verteilung ist dahin gehend “gemischt”, dass sie sowohl Eigenschaften von stetigen Verteilungen (nimmt nichtnegative reelle Werte an) als auch von diskreten Verteilungen (positive Wahrscheinlichkeitsmasse an einem Einzelwert, 0) aufweist. Die abhängige Variable muss numerisch sein, mit Datenwerten größer oder gleich 0. Wenn ein Datenwert kleiner als 0 ist oder fehlt, wird der entsprechende Fall nicht in der Analyse verwendet. Der feste Wert des Parameters der Tweedie-Verteilung kann jede beliebige Zahl zwischen 1 und 2 sein.
- **Multinomial.** Diese Verteilung eignet sich für Variablen, die eine ordinale Antwort repräsentieren. Bei der abhängigen Variablen kann es sich um eine numerische Variable oder eine String-Variable handeln. Sie muss mindestens zwei verschiedene gültige Datenwerte aufweisen.

#### Linkfunktionen.

Die Verknüpfungsfunktion (Linkfunktion) ist eine Transformation der abhängigen Variablen, die eine Schätzung des Modells ermöglicht. Die folgenden Funktionen sind verfügbar:

- **Identität.**  $f(x)=x$ . Die abhängige Variable wird nicht transformiert. Diese Verknüpfung kann mit jeder beliebigen Verteilung verwendet werden.
- **Log-Log komplementär.**  $f(x)=\log(-\log(1-x))$ . Nur für die Binomialverteilung geeignet.
- **Cauchit (kumulativ).**  $f(x) = \tan(\pi (x - 0.5))$ , auf die kumulative Wahrscheinlichkeit der einzelnen Kategorien der Antwort angewendet. Nur für die Multinomialverteilung geeignet.
- **Log-Log komplementär (kumulativ),**  $f(x)=\ln(-\ln(1-x))$ , auf die kumulative Wahrscheinlichkeit der einzelnen Kategorien der Antwort angewendet. Nur für die Multinomialverteilung geeignet.

- **Logit (kumulativ)**  $f(x)=\ln(x / (1-x))$ , auf die kumulative Wahrscheinlichkeit der einzelnen Kategorien der Antwort angewendet. Nur für die Multinomialverteilung geeignet.
- **Log-Log negativ (kumulativ)**  $f(x)=-\ln(-\ln(x))$ , auf die kumulative Wahrscheinlichkeit der einzelnen Kategorien der Antwort angewendet. Nur für die Multinomialverteilung geeignet.
- **Probit (kumulativ)**  $f(x)=\Phi^{-1}(x)$ , auf die kumulative Wahrscheinlichkeit der einzelnen Kategorien der Antwort angewendet, wobei  $\Phi^{-1}$  die umgekehrte kumulative Standard-Normalverteilungsfunktion ist. Nur für die Multinomialverteilung geeignet.
- **Log**  $f(x)=\log(x)$ . Diese Verknüpfung kann mit jeder beliebigen Verteilung verwendet werden.
- **Log. Komplement**  $f(x)=\log(1-x)$ . Nur für die Binomialverteilung geeignet.
- **Logit**  $f(x)=\log(x / (1-x))$ . Nur für die Binomialverteilung geeignet.
- **Negativ binomial**  $f(x)=\log(x / (x+k^{-1}))$ , wobei  $k$  der Hilfsparameter der negativen Binomialverteilung ist. Nur für die negative Binomialverteilung geeignet.
- **Log-Log negativ**  $f(x)=-\log(-\log(x))$ . Nur für die Binomialverteilung geeignet.
- **Odds-Potenz**  $f(x)=[(x/(1-x))^\alpha-1]/\alpha$ , wenn  $\alpha \neq 0$ .  $f(x)=\log(x)$ , wenn  $\alpha=0$ .  $\alpha$  ist die erforderliche Zahlenangabe. Es muss sich dabei um eine reelle Zahl handeln. Nur für die Binomialverteilung geeignet.
- **Probit**  $f(x)=\Phi^{-1}(x)$ , wobei  $\Phi^{-1}$  die umgekehrte kumulative Standard-Normalverteilungsfunktion ist. Nur für die Binomialverteilung geeignet.
- **Exponent**  $f(x)=x^\alpha$ , wenn  $\alpha \neq 0$ .  $f(x)=\log(x)$ , wenn  $\alpha=0$ .  $\alpha$  ist die erforderliche Zahlenangabe. Es muss sich dabei um eine reelle Zahl handeln. Diese Verknüpfung kann mit jeder beliebigen Verteilung verwendet werden.

**Parameter.** Mit den Steuerelementen in dieser Gruppe können Sie Parameterwerte festlegen, wenn bestimmte Verteilungsoptionen gewählt werden.

- **Parameter für negativ binomial.** Für negative binomiale Verteilung geben Sie entweder einen Wert an oder Sie gestatten dem System, einen geschätzten Wert bereitzustellen.
- **Parameter für Tweedie.** Geben Sie als festen Wert des Parameters der Tweedie-Verteilung eine Zahl zwischen 1,0 und 2,0 an.

**Parameterschätzung.** Mit den Steuerelementen in dieser Gruppe können Sie Schätzmethoden festlegen und Anfangswerte für die Parameterschätzer angeben.

- **Methode.** Sie können eine Methode für die Parameterschätzung auswählen. Sie haben die Wahl zwischen “Newton-Raphson”, “Fisher-Bewertung” und einer Hybridmethode, bei der zuerst Iterationen des Fisher-Scorings durchgeführt werden und dann zur Methode “Newton-Raphson” gewechselt wird. Wenn während der Phase “Fisher-Bewertung” der Hybridmethode Konvergenz erreicht wird, bevor die maximale Anzahl an Fisher-Iterationen erreicht wurde, fährt der Algorithmus mit der Newton-Raphson-Methode fort.
- **Skalenparametermethode.** Sie können eine Schätzmethode für den Skalenparameter auswählen. Bei der Maximum-Likelihood-Methode wird der Skalenparameter zusammen mit den Modelleffekten geschätzt. Beachten Sie, dass diese Option ungültig ist, wenn die Antwort eine negative Binomialverteilung, eine Poisson-Verteilung oder eine Binomialverteilung aufweist. Die Optionen für die Abweichung und das Pearson-Chi-Quadrat schätzen den

Skalenparameter aus dem Wert der jeweiligen Statistik. Alternativ können Sie einen festen Wert für den Skalenparameter angeben.

- **Kovarianzmatrix.** Der modellbasierte Schätzer ist das Negative der verallgemeinerten Inversen der Hesse-Matrix. Der robuste Schätzer (auch Huber-/White-/Sandwich-Schätzer genannt) ist ein “korrigierter” modellbasierter Schätzer, der eine konsistente Schätzung der Kovarianz bietet, selbst wenn Varianz und Verknüpfungsfunktionen falsch angegeben wurden.

**Iterationen.** Mit diesen Optionen können Sie die Parameter für die Modellkonvergenz festlegen. Für weitere Informationen siehe Thema [Verallgemeinerte lineare Modelle – Iterationen](#) auf S. 336.

**Ausgabe.** Mit diesen Optionen können Sie zusätzliche Statistiken anfordern, die in der erweiterten Ausgabe des vom Knoten erstellten Modell-Nugget angezeigt werden. Für weitere Informationen siehe Thema [Verallgemeinerte lineare Modelle – Erweiterte Ausgabe](#) auf S. 337.

**Toleranz für Prüfung auf Singularität.** Singuläre (bzw. nicht invertierbare) Matrizen weisen linear abhängige Spalten auf, die zu ernststen Problemen für den Schätzalgorithmus führen können. Auch annähernd singuläre Matrizen können zu schlechten Ergebnissen führen, daher behandelt die Prozedur eine Matrix, deren Determinante unter dem Toleranzwert liegt, als singulär. Geben Sie einen positiven Wert ein.

## Verallgemeinerte lineare Modelle – Iterationen

Sie können die Konvergenzparameter für die Schätzung des verallgemeinerten linearen Modells festlegen.

Abbildung 10-49  
Verallgemeinerte lineare Modelle – Iterationsoptionen

**Verallgemeinerte lineare Modelle: Iterationen**

Maximale Iterationen: 100

Maximale Schritthalbierung: 5

Auf Trennung der Datenpunkte prüfen (nur Flag-Ziele)

Starten der Iteration: 20

**Konvergenzkriterien**

Es muss mindestens ein Konvergenzkriterium angegeben werden mit einem Minimum größer als 0.

	Minimum:	Typ
<input checked="" type="checkbox"/> Änderung beim Parameterschätzer	0,0000010000	Absolut
<input type="checkbox"/> Änderung bei Log-Likelihood	0,0000010000	Absolut
<input type="checkbox"/> Konvergenz der Hesse-Matrix	0,0000010000	Absolut

OK Abbrechen Hilfe

### Iterationen.

- **Maximale Iterationen.** Dies ist die maximale Anzahl der Iterationen, die im Algorithmus vorgenommen werden. Geben Sie eine nichtnegative ganze Zahl an.

- **Maximale Schritthalbierung.** Bei jeder Iteration wird die Schrittgröße um den Faktor 0,5 reduziert, bis die Log-Likelihood ansteigt oder die Maximalzahl für die Schritthalbierung erreicht ist. Geben Sie eine positive ganze Zahl ein.
- **Prüfen auf Datenpunkttrennung.** Mit dieser Option lassen Sie Tests durch den Algorithmus durchführen, mit denen sichergestellt wird, dass die Parameterschätzer eindeutige Werte aufweisen. Eine Trennung wird vorgenommen, sobald ein Modell erzeugt werden kann, in dem alle Fälle fehlerfrei klassifiziert werden. Diese Option ist für binomiale Antworten mit Binärformat verfügbar.

#### Konvergenzkriterien.

- **Parameterkonvergenz.** Mit dieser Option wird der Algorithmus nach einer Iteration angehalten, bei der die absolute oder relative Änderung bei den Parameterschätzern unter dem angegebenen (positiven) Wert liegt.
- **Log-Wahrscheinlichkeitskonvergenz.** Mit dieser Option wird der Algorithmus nach einer Iteration angehalten, bei der die absolute oder relative Änderung bei der Log-Likelihood-Funktion unter dem angegebenen (positiven) Wert liegt.
- **Konvergenz der Hesse-Matrix.** Für die Spezifikation "Absolut" wird angenommen, dass eine Konvergenz vorliegt, wenn eine Statistik auf der Basis der Konvergenz der Hesse-Matrix kleiner als der angegebene positive Wert ist. Für die Spezifikation "Relativ" wird angenommen, dass eine Konvergenz vorliegt, wenn die Statistik kleiner als das Produkt aus dem angegebenen positiven Wert und dem absoluten Wert der Log-Likelihood ist.

## Verallgemeinerte lineare Modelle – Erweiterte Ausgabe

Abbildung 10-50

Verallgemeinerte lineare Modelle – Erweiterte Ausgabeoptionen

The screenshot shows the 'Verallgemeinerte lineare Modelle: Erweiterte Ausgabe' dialog box. The options are as follows:

- Fallverarbeitungsübersicht
- Statistik für Anpassungsgüte
- Kovarianzmatrix für Parameterschätzer
- Korrelationsmatrix für Parameterschätzer
- Deskriptive Statistiken
- Parameterschätzungen
  - Exponentielle Parameterschätzer einschließen
- Lagrange-Multiplikator-Test fester Skalenparameter oder Hilfsparameter für negativ binomial
- Kontrastkoeffizienten der L-Matrizen
- Allgemeine schätzbare Funktionen
- Modellinformationen
- Modellübersichtsstatistiken
- Iterationsprotokoll
- Druckintervall: 1

**Modellereffekte**

- Analysetyp: Typ III
- Konfidenzintervallniveau (%): 95
- Chi-Quadrat-Statistik:
  - Wald
  - Likelihood-Quotient
- Konfidenzintervalltyp:
  - Wald
  - Profil-Likelihood
- Toleranzstufe: 0,0001
- Log-Likelihood-Funktion: Vollständig

Buttons: OK, Abbrechen, Hilfe

Wählen Sie die optionalen Ausgaben aus, die in der erweiterten Ausgabe des Modell-Nuggets für das verallgemeinerte lineare Modell angezeigt werden sollen. Zur Anzeige der erweiterten Ausgabe durchsuchen Sie das Modell-Nugget und klicken Sie auf die Registerkarte *Erweitert*. [Für weitere Informationen siehe Thema GenLin-Modell-Nugget – Erweiterte Ausgabe auf S. 340.](#)

Die folgenden Ausgaben sind verfügbar:

- **Fallverarbeitungsübersicht.** Zeigt die Anzahl und den Prozentsatz der Fälle an, die in die Analyse und in die Tabelle “Korrelierte Datenzusammenfassung” aufgenommen bzw. daraus ausgeschlossen werden.
- **Deskriptive Statistik.** Zeigt eine deskriptive Statistik und Zusammenfassungsinformationen über die abhängige Variable, die Kovariaten und die Faktoren an.
- **Modellinformationen.** Zeigt den Namen des Daten-Sets, die abhängige Variable bzw. die Ereignis- und Versuchsvariablen, die Offset-Variable, die Skalengewichtungsvariable, die Wahrscheinlichkeitsverteilung und die Verknüpfungsfunktion an.
- **Statistik für Anpassungsgüte.** Zeigt an: Abweichung und skalierte Abweichung, Pearson-Chi-Quadrat und skaliertes Pearson-Chi-Quadrat, Log-Likelihood, Akaike-Informationskriterium (AIC), AIC mit Korrektur für endliche Stichproben (AICC), Bayes-Informationskriterium (BIC) und konsistentes AIC (CAIC).
- **Modellzusammenfassungsverstatistik.** Zeigt Tests für die Anpassungsgüte des Modells an, darunter Likelihood-Quotienten-Statistiken für den Omnibus-Test für die Anpassungsgüte, sowie Statistiken für Kontraste des Typs I bzw. III für jeden Effekt.
- **Parameterschätzer.** Zeigt Parameterschätzer und entsprechende Teststatistiken und Konfidenzintervalle an. Wahlweise können Sie zusätzlich zu den rohen, unbearbeiteten Parameterschätzern auch potenzierte Parameterschätzer anzeigen.
- **Kovarianzmatrix für Parameterschätzer.** Zeigt die Kovarianzmatrix für die geschätzten Parameter an.
- **Korrelationsmatrix für Parameterschätzer.** Zeigt die Korrelationsmatrix für die geschätzten Parameter an.
- **Kontrastkoeffizienten-(L-)Matrizen.** Zeigt die Kontrastkoeffizienten für die Standardeffekte und für die geschätzten Randmittel an, sofern auf der Registerkarte “Geschätzte Randmittel” angefordert.
- **Allgemeine schätzbare Funktionen.** Zeigt die Matrizen für die Generierung der Kontrastkoeffizienten-(L-)Matrizen an.
- **Iterationsprotokoll.** Zeigt das Iterationsprotokoll für Parameterschätzer und Log-Likelihood an und druckt die letzte Auswertung des Gradientenvektors und der Hesse-Matrix. Die Tabelle mit dem Iterationsprotokoll zeigt Parameterschätzer für jede  $n$ -te Iteration an, beginnend mit der 0-ten Iteration (Anfangsschätzungen). Dabei ist  $n$  der Wert des Druckintervalls. Wenn das Iterationsprotokoll angefordert wird, wird die letzte Iteration unabhängig von  $n$  stets angezeigt.
- **Lagrange-Multiplikator-Test.** Zeigt die Statistiken für den Lagrange-Multiplikator-Test an, die zur Bewertung der Gültigkeit eines Skalenparameters dienen, der mithilfe des Pearson-Chi-Quadrats berechnet wurde oder für den bei der Normal-, Gamma- und inversen Normalverteilung ein fester Wert festgelegt wurde. Bei der negativen Binomialverteilung wird hiermit der feste Hilfsparameter getestet.

#### **Modelleffekte.**

- **Analysetyp.** Geben Sie den Typ der zu erstellenden Analyse an. Eine Analyse des Typs I ist im Allgemeinen dann angebracht, wenn Sie von vorneherein Gründe dafür haben, die Einflussvariablen (Prädiktoren) im Modell zu ordnen. Typ III dagegen ist allgemeiner anwendbar. Wald- oder Likelihood-Quotienten-Statistiken werden auf der Basis der Auswahl in der Chi-Square-Statistikgruppe berechnet.
- **Konfidenzintervalle.** Geben Sie für das Konfidenzniveau einen Wert an, der über 50 und unter 100 liegt. Wald-Intervalle beruhen auf der Annahme, dass die Parameter eine asymptotische Normalverteilung aufweisen. Profil-Likelihood-Intervalle sind präziser, können aber rechnerisch aufwendig sein. Die Toleranzstufe für Profil-Likelihood-Intervalle ist das Kriterium, anhand dessen der iterative Algorithmus zur Intervallberechnung gestoppt wird.
- **Log-Likelihood-Funktion.** Legt das Anzeigeformat der Log-Likelihood-Funktion fest. Die vollständige Funktion enthält einen zusätzlichen Term, der hinsichtlich der Parameterschätzer konstant ist. Er hat keine Auswirkungen auf die Parameterschätzung und wird bei einigen Softwareprodukten nicht angezeigt.

## GenLin-Modell-Nugget

Ein GenLin-Modell-Nugget steht für die Gleichungen, die durch einen GenLin-Knoten geschätzt wurden. Sie enthalten alle Informationen, die vom Modell erfasst wurden, sowie Informationen über die Struktur und Leistung des Modells.

Bei der Ausführung eines Streams, der ein GenLin-Modell-Nugget enthält, fügt der Knoten neue Felder hinzu, deren Inhalt von der Art des Zielfelds abhängt:

- **Flag-Ziel.** Fügt Felder hinzu, die die vorhergesagte Kategorie und die zugehörige Wahrscheinlichkeit enthalten sowie die Wahrscheinlichkeiten für die einzelnen Kategorien. Die Namen der ersten beiden neuen Felder werden aus dem Namen des prognostizierten Ausgabefelds abgeleitet, dem *\$G-* für die vorhergesagte Kategorie und *\$GP-* für die zugehörige Wahrscheinlichkeit vorangestellt ist. Bei einem Ausgabefeld mit der Bezeichnung *Standard* beispielsweise erhalten die neuen Felder die Namen *\$G-Standard* und *\$GP-Standard*. Diese letzteren beiden zusätzlichen Felder werden auf der Grundlage der Werte des Ausgabefelds benannt, denen *\$GP-* vorangestellt ist. Wenn für *Standard* die Werte *Ja* und *Nein* zulässig sind, lauten die Namen der neuen Felder *\$GP-Ja* und *\$GP-Nein*.
- **Stetiges Ziel.** Fügt Felder hinzu, die den vorhergesagten Mittelwert und Standardfehler enthalten.
- **Stetiges Ziel, enthält die Anzahl der Ereignisse in einem Set von Versuchen.** Fügt Felder hinzu, die den vorhergesagten Mittelwert und Standardfehler enthalten.

**Erstellen eines Filterknotens.** Im Menü “Generieren” können Sie einen neuen Filterknoten zur Übergabe der Eingabefelder auf der Grundlage der Ergebnisse erstellen.

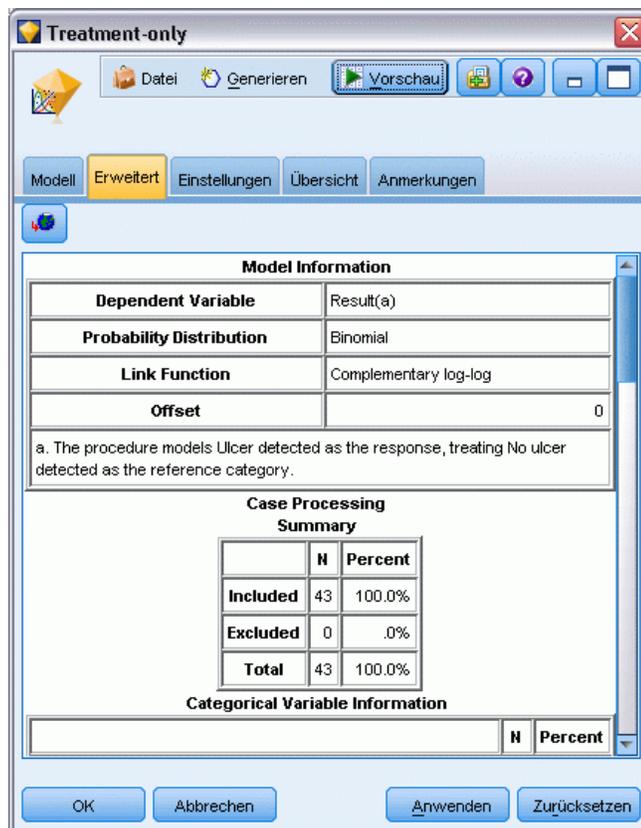
### Bedeutsamkeit des Prädiktors

Optional kann auf der Registerkarte “Modell” auch ein Diagramm, das die relative Bedeutsamkeit der einzelnen Prädiktoren für die Schätzung des Modells angibt, angezeigt werden. Normalerweise ist es sinnvoll, die Modellierungsbemühungen auf die wichtigsten Prädiktoren zu konzentrieren und diejenigen Prädiktoren zu verwerfen bzw. zu ignorieren, die die geringste Bedeutung haben. Beachten Sie, dass dieses Diagramm nur verfügbar ist, wenn vor dem Generieren des Modells

Bedeutsamkeit der Prädiktoren berechnen auf der Registerkarte “Analysieren” ausgewählt wurde. Für weitere Informationen siehe Thema Bedeutsamkeit des Prädiktors in Kapitel 3 auf S. 54.

### GenLin-Modell-Nugget – Erweiterte Ausgabe

Abbildung 10-51  
GenLin-Modell-Nugget – Registerkarte “Erweitert”



Die erweiterte Ausgabe für verallgemeinerte lineare Modelle bietet detaillierte Informationen zum geschätzten Modell und dessen Leistung. Die meisten der in der erweiterten Ausgabe enthaltenen Informationen weisen einen hohen Fachlichkeitsgrad auf und zur richtigen Interpretation dieser Ausgaben sind umfassende Kenntnisse dieser Art von Analysen erforderlich. Für weitere Informationen siehe Thema Verallgemeinerte lineare Modelle – Erweiterte Ausgabe auf S. 337.

### GenLin-Modell-Nugget – Einstellungen

Auf der Registerkarte “Einstellungen” für ein GenLin-Modell-Nugget können Sie beim Scoring des Modells Neigungs-Scores ermitteln. Diese Registerkarte ist nur für Modelle mit Flag-Zielen verfügbar und erst nachdem das Modell-Nugget einem Stream hinzugefügt wurde.

Abbildung 10-52  
GenLin-Modell-Nugget, Registerkarte "Einstellungen" für Flag-Ziele



**Scores für Rohneigung berechnen.** Bei Modellen mit einem Flag-Ziel (das als Vorhersage "Ja" bzw. "Nein" ausgibt) können Sie Neigungs-Scores anfordern, die die Wahrscheinlichkeit des für das Zielfeld angegebenen wahren Ergebnisses angeben. Diese Werte werden zusätzlich zu den anderen Vorhersage- und Konfidenzwerten angegeben, die ggf. während des Scorings erstellt werden.

**Scores für korrigierte Neigung berechnen.** Rohneigungs-Scores beruhen ausschließlich auf den Trainingsdaten und sind aufgrund der Neigung vieler Modelle zur übermäßigen Anpassung an die Trainingsdaten möglicherweise zu optimistisch. Bei korrigierten Neigungen wird versucht, dies durch Evaluation der Modellleistung anhand einer Test- bzw. Validierungspartition zu kompensieren. Bei dieser Option muss im Stream ein Partitionsfeld definiert sein und die Scores für die korrigierte Neigung müssen im Modellierungsknoten aktiviert werden, bevor das Modell generiert wird.

## GenLin-Modell-Nugget – Übersicht

Auf der Registerkarte "Übersicht" für ein GenLin-Modell-Nugget werden die Felder und Einstellungen angezeigt, die zum Generieren des Modells verwendet wurden. Wenn Sie einen Analyseknoden ausgeführt haben, der an diesen Modellierungsknoten angehängt ist, werden die Informationen aus dieser Analyse ebenfalls in diesem Abschnitt angezeigt. [Für weitere Informationen siehe Thema Analyseknoden in Kapitel 6 in IBM SPSS Modeler 14.2- Quellen-, Prozess- und Ausgabeknoten.](#) Allgemeine Informationen zur Verwendung des Modellbrowsers finden Sie unter [Durchsuchen von Modell-Nuggets auf S. 52.](#)

Abbildung 10-53  
GenLin-Modell-Nugget – Registerkarte "Übersicht"



## Cox-Knoten

Die Cox-Regression erstellt ein Vorhersagemodell für Daten, die die Zeit bis zum Eintreten des Ereignisses angeben. Das Modell erstellt eine Überlebensfunktion, die die Wahrscheinlichkeit vorhersagt, dass das relevante Ereignis zum festgelegten Zeitpunkt  $t$  für bestimmte Werte der Einflussvariablen (Prädiktorvariablen) eingetreten ist. Die Form der Überlebensfunktion und die Regressionskoeffizienten für die Einflussvariablen werden aus beobachteten Subjekten geschätzt. Anschließend kann das Modell auf neue Fälle angewendet werden, die Messungen für die Einflussvariablen enthalten. Beachten Sie, dass Informationen aus zensierten Subjekten, also Subjekten, bei denen das relevante Ereignis während der Beobachtungszeit nicht eintritt, einen nützlichen Beitrag zur Schätzung des Modells leisten.

**Beispiel.** Im Rahmen seiner Bemühungen zur Reduzierung der Kundenabwanderung ist ein Telekommunikationsunternehmen daran interessiert, die "Zeit bis zur Abwanderung" zu modellieren, um die Faktoren zu ermitteln, die für Kunden gelten, die rasch zu einem anderen Dienst wechseln. Dazu wird eine Zufallsstichprobe der Kunden ausgewählt und ihre Zeit als Kunden (unabhängig davon, ob sie noch immer aktive Kunden sind) und verschiedene demografische Felder werden aus der Datenbank gezogen. [Für weitere Informationen siehe Thema Verwenden der Cox-Regression zur Modellierung der Zeit bis zur Kundenabwanderung in Kapitel 27 in IBM SPSS Modeler 14.2- Anwendungshandbuch.](#)

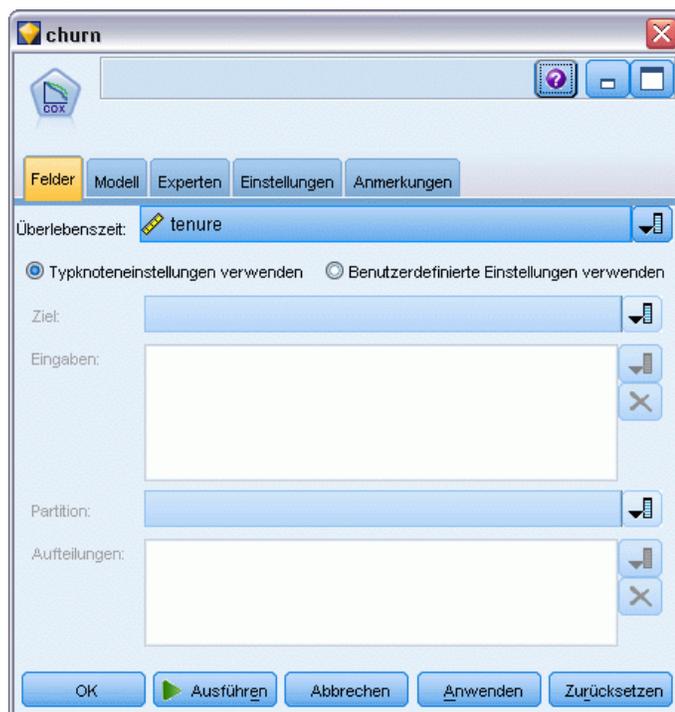
**Anforderungen.** Es werden mindestens ein Eingabefeld und genau ein Zielfeld benötigt und Sie müssen innerhalb des Cox-Knotens ein Feld für die Überlebenszeit angeben. Das Zielfeld sollte so kodiert sein, dass der Wert “falsch” Überleben anzeigt und der Wert “wahr” anzeigt, dass das relevante Ereignis eingetreten ist. Das Feld muss ein Messniveau des Typs *Flag* mit dem Speichertyp “Zeichenkette” oder “Ganze Zahl” aufweisen. (Der Speichertyp kann, falls erforderlich, mithilfe eines Füller- oder Ableitungsknotens konvertiert werden. [Für weitere Informationen siehe Thema Speichertypkonvertierung mithilfe des Füllerknotens in Kapitel 4 in IBM SPSS Modeler 14.2- Quellen-, Prozess- und Ausgabeknoten.](#)) Felder, die auf *Beides* oder *Keine* gesetzt sind, werden ignoriert. Bei den im Modell verwendeten Feldern müssen die Typen vollständig instanziiert sein. Bei der Überlebenszeit kann es sich um ein beliebiges numerisches Feld handeln.

**Datums- & Zeitangaben.** Felder vom Typ “Datum & Uhrzeit” können nicht unmittelbar zur Definition der Überlebenszeit verwendet werden; wenn Felder vom Typ “Datum & Uhrzeit” vorliegen, sollten Sie sie verwenden, um ein Feld mit Überlebenszeiten zu erstellen, das auf der Differenz zwischen dem Datum des Eintritts in die Studie und dem Datum der Beobachtung beruht. [Für weitere Informationen siehe Thema Arbeiten mit Zeit- und Datumsangaben in Kapitel 7 in IBM SPSS Modeler 14.2- Benutzerhandbuch.](#)

**Kaplan-Meier-Analyse.** Die Cox-Regression kann ohne Eingabefelder durchgeführt werden. Dies entspricht einer Kaplan-Meier-Analyse.

## Feldoptionen für Cox-Knoten

Abbildung 10-54  
Dialogfeld des Cox-Knotens, Registerkarte “Felder”



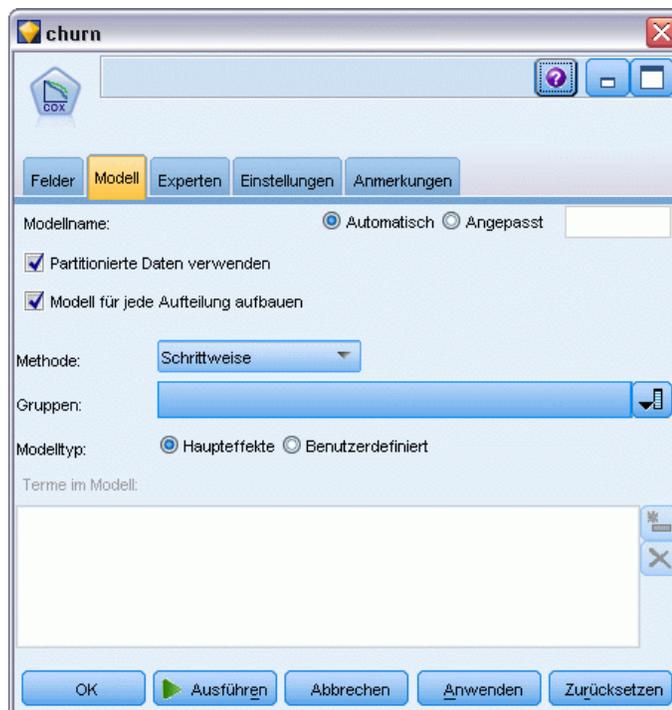
**Überlebenszeit.** Wählen Sie ein numerisches Feld (mit dem Messniveau *Stetig*), um den Knoten ausführen zu können. Die Überlebenszeit gibt die Lebensdauer des vorherzusagenden Datensatzes an. Wenn zum Beispiel die Zeit bis zur Abwanderung von Kunden modelliert wird, wäre dies das Feld, das festhält, wie lange der Kunde schon beim Unternehmen ist. Das Beitritts- oder Abwanderungsdatum des Kunden hätte keine Auswirkungen auf das Modell; nur die Beschäftigungsdauer des Kunden wäre relevant.

Die Überlebenszeit wird als Dauer ohne Einheiten angenommen. Sie müssen sicherstellen, dass die Eingabefelder mit der Überlebenszeit übereinstimmen. Beispielsweise würden Sie in einer Studie zur Messung des Kundenverlusts pro Monat als Eingabe den Monatsumsatz statt des Jahresumsatzes verwenden. Wenn Ihre Daten Start- und Endzeiten statt einer Dauer aufweisen, müssen Sie diese Zeiten für eine Dauer neu kodieren, die vor dem Cox-Knoten liegt.

Die restlichen Felder in diesem Dialogfeld sind Standardfelder, die überall in IBM® SPSS® Modeler verwendet werden. [Für weitere Informationen siehe Thema Feldoptionen der Modellierungsknoten in Kapitel 3 auf S. 36.](#)

## Modelloptionen für Cox-Knoten

Abbildung 10-55  
Dialogfeld des Cox-Knotens, Registerkarte "Modell"



**Modellname.** Sie können den Modellnamen automatisch basierend auf den Ziel- oder ID-Feldnamen (oder dem Modelltyp in Fällen, in denen kein solches Feld angegeben ist) generieren oder einen benutzerdefinierten Namen eingeben.

**Partitionierte Daten verwenden.** Wenn ein Partitionsfeld definiert ist, gewährleistet diese Option, dass nur Daten aus der Trainingspartition für die Modellerstellung verwendet werden. [Für weitere Informationen siehe Thema Partitionsknoten in Kapitel 4 in IBM SPSS Modeler 14.2- Quellen-, Prozess- und Ausgabeknoten.](#)

**Aufteilungsmodelle erstellen.** Erstellt ein separates Modell für jeden möglichen Wert von Eingabefeldern, die als Aufteilungsfelder angegeben werden. [Für weitere Informationen siehe Thema Erstellung von aufgeteilten Modellen in Kapitel 3 auf S. 31.](#)

**Methode.** Zur Eingabe von Prädiktoren in das Modell stehen folgende Optionen zur Verfügung:

- **Eingabetaste.** Dies ist das Standardverfahren, bei dem alle Terme direkt in das Modell aufgenommen werden. Beim Erstellen des Modells wird keine Feldauswahl durchgeführt.
- **Schrittweise.** Bei der Methode “Schrittweise” der Feldauswahl wird, wie der Name andeutet, das Modell in Schritten erstellt. Das anfängliche Modell ist das einfachste Modell, das möglich ist. Es enthält keine Modellterme (außer der Konstanten) im Modell. Bei jedem Schritt werden die Terme, die noch nicht zum Modell hinzugefügt wurden, bewertet, und wenn der beste dieser Terme in signifikanter Weise zur Vorhersagekraft des Modells beiträgt, wird er hinzugefügt. Außerdem werden die derzeit im Modell enthaltenen Terme neu bewertet, um zu ermitteln, ob einige davon ohne signifikante Beeinträchtigung des Modells entfernt werden können. Wenn dies der Fall ist, werden sie entfernt. Der Vorgang wird wiederholt und andere Terme werden hinzugefügt und/oder entfernt. Wenn keine weiteren Terme zur Verbesserung des Modells hinzugefügt werden können und keine Terme mehr entfernt werden können, ohne das Modell zu beeinträchtigen, wird das endgültige Modell generiert.
- **Schrittweise rückwärts.** Die Methode “Schrittweise rückwärts” ist im Grunde das Gegenteil der Methode “Schrittweise”. Bei dieser Methode enthält das ursprüngliche Modell alle Terme als Prädiktoren. Bei jedem Schritt werden die Terme im Modell evaluiert und alle Terme, die ohne signifikante Beeinträchtigung des Modells entfernt werden können, werden entfernt. Außerdem werden die zuvor entfernten Terme erneut evaluiert, um zu ermitteln, ob der beste dieser Terme in signifikanter Weise zur Vorhersagekraft des Modells beiträgt. Ist dies der Fall, so wird er wieder in das Modell aufgenommen. Wenn keine Terme mehr entfernt werden können, ohne das Modell wesentlich zu beeinträchtigen, und keine weiteren Terme zur Verbesserung des Modells hinzugefügt werden können, wird das endgültige Modell generiert.

*Hinweis:* Die automatischen Methoden (z. B. “Schrittweise” und “Schrittweise rückwärts”) sind sehr anpassungsfähige Lernmethoden und weisen eine starke Tendenz zur übermäßigen Anpassung an die Trainingsdaten auf. Bei der Verwendung dieser Methoden ist es ganz besonders wichtig, die Validität des entstehenden Modells zu überprüfen – entweder mit neuen Daten oder mithilfe einer zurückgehaltenen Teststichprobe, die mit dem Partitionsknoten erstellt wurde. [Für weitere Informationen siehe Thema Partitionsknoten in Kapitel 4 in IBM SPSS Modeler 14.2- Quellen-, Prozess- und Ausgabeknoten.](#)

**Gruppen.** Die Angabe eines Gruppenfelds führt dazu, dass der Knoten separate Modelle für die einzelnen Kategorien des Felds berechnet. Es kann sich dabei um ein beliebiges kategoriales Feld (Flag oder nominal) mit dem Speichertyp “Zeichenkette” oder “Ganze Zahl” handeln.

**Modelltyp.** Es gibt zwei Optionen zur Definition der Terme im Modell. **Haupteffekte-**Modelle beinhalten nur die einzelnen Eingabefelder und testen nicht die Interaktionen (multiplikativen Effekte) zwischen den Eingabefeldern. **Benutzerdefinierte Modelle** enthalten nur die von Ihnen angegebenen Terme (Haupteffekte und Interaktionen). Verwenden Sie bei der Auswahl dieser

Option die Liste “Terme im Modell”, um Terme zum Modell hinzuzufügen oder daraus zu entfernen.

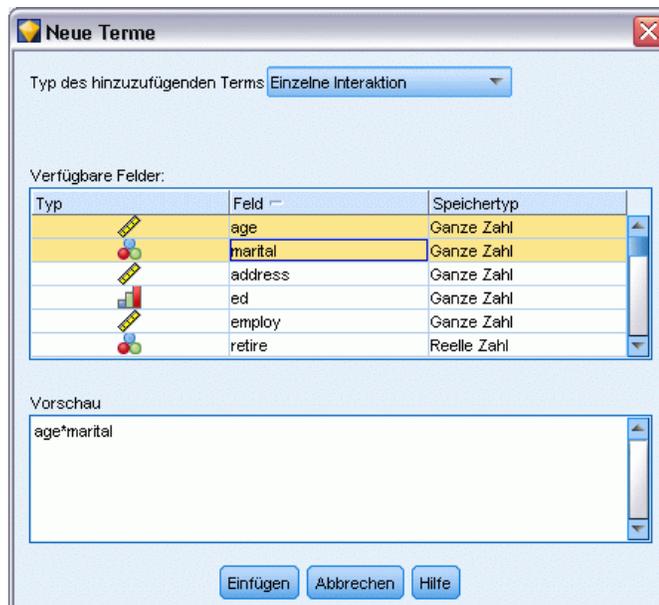
**Terme im Modell.** Beim Erstellen eines benutzerdefinierten Modells müssen Sie die Terme im Modell explizit angeben. Die Liste zeigt die aktuelle Menge an Termen für das Modell. Mit den Schaltflächen auf der rechten Seite der Liste “Terme im Modell” können Sie Modellterme hinzufügen bzw. entfernen.

- ▶ Um Terme zum Modell hinzuzufügen, klicken Sie auf die Schaltfläche *Neue Terme im Modell hinzufügen*.
- ▶ Zum Löschen von Termen wählen Sie die gewünschten Terme aus und klicken Sie auf die Schaltfläche *Ausgewählte Terme im Modell löschen*.

### Hinzufügen von Termen zu einem Cox-Regressionsmodell

Beim Anfordern eines benutzerdefinierten Modells können Sie Terme zum Modell hinzufügen, indem Sie auf der Registerkarte für das Modell auf die Schaltfläche *Neue Terme im Modell hinzufügen* klicken. Ein neues Dialogfeld wird geöffnet, in dem Sie Terme angeben können.

Abbildung 10-56  
Dialogfeld “Neue Terme”



**Typ des hinzuzufügenden Terms.** Es gibt mehrere Methoden zum Hinzufügen von Termen zum Modell, je nach der Auswahl der Eingabefelder in der Liste der verfügbaren Felder.

- **Einzelne Interaktion.** Fügt den Term ein, der für die Interaktion aller ausgewählten Felder steht.
- **Haupteffekte.** Fügt für jedes ausgewählte Eingabefeld einen Haupteffekt-Term (das Feld selbst) ein.

- **Alle zweifachen Interaktionen.** Fügt für jedes mögliche Paar ausgewählter Eingabefelder einen Zweifach-Interaktions-Term (das Produkt der Eingabefelder) ein. Beispiel: Bei Auswahl der Eingabefelder  $A$ ,  $B$  und  $C$  in der Liste der verfügbaren Felder werden bei dieser Methode die Terme  $A * B$ ,  $A * C$  und  $B * C$  eingefügt.
- **Alle dreifachen Interaktionen.** Fügt für jede mögliche Kombination aus jeweils drei ausgewählten Eingabefeldern einen Dreifach-Interaktions-Term (das Produkt der Eingabefelder) ein. Beispiel: Bei Auswahl der Eingabefelder  $A$ ,  $B$ ,  $C$  und  $D$  in der Liste der verfügbaren Felder werden bei dieser Methode die Terme  $A * B * C$ ,  $A * B * D$ ,  $A * C * D$  und  $B * C * D$  eingefügt.
- **Alle vierfachen Interaktionen.** Fügt für jede mögliche Kombination aus jeweils vier ausgewählten Eingabefeldern einen Vierfach-Interaktions-Term (das Produkt der Eingabefelder) ein. Beispiel: Bei Auswahl der Eingabefelder  $A$ ,  $B$ ,  $C$ ,  $D$  und  $E$  in der Liste der verfügbaren Felder werden bei dieser Methode die Terme  $A * B * C * D$ ,  $A * B * C * E$ ,  $A * B * D * E$ ,  $A * C * D * E$  und  $B * C * D * E$  eingefügt.

**Verfügbare Felder.** Listet die verfügbaren Eingabefelder auf, die bei der Konstruktion der Modellterme verwendet werden sollen. Beachten Sie, dass die Liste möglicherweise Felder enthält, bei denen es sich nicht um zulässige Eingabefelder handelt. Achten Sie daher sorgfältig darauf, dass alle Modellterme nur Eingabefelder enthalten.

**Vorschau.** Zeigt die Terme an, die beim Klicken auf Einfügen zum Modell hinzugefügt werden. Dabei werden die ausgewählten Felder und der oben ausgewählte Termtyp zugrunde gelegt.

**Einfügen.** Fügt (auf der Grundlage der aktuellen Auswahl von Feldern und des Termtyps) Terme in das Modell ein und schließt das Dialogfeld.

## Expertenoptionen für Cox-Knoten

Abbildung 10-57  
Dialogfeld des Cox-Knotens, Registerkarte "Experten"



**Konvergenz.** Mit diesen Optionen können Sie die Parameter für die Modellkonvergenz festlegen. Bei der Ausführung des Modells steuern die Konvergenzkriterien, wie viele Male die verschiedenen Parameter wiederholt durchlaufen werden, um zu ermitteln, wie gut sie passen. Je häufiger die Parameter durchprobiert werden, desto enger liegen die Ergebnisse beieinander (d. h. die Ergebnisse konvergieren). [Für weitere Informationen siehe Thema Konvergenzkriterien für Cox-Knoten auf S. 348.](#)

**Ausgabe.** Mit diesen Optionen können Sie zusätzliche Statistiken und Plots anfordern (einschließlich der Überlebenskurve), die in der erweiterten Ausgabe des vom Knoten erstellten generierten Modells angezeigt werden. [Für weitere Informationen siehe Thema Cox-Knoten – Erweiterte Ausgabeoptionen auf S. 349.](#)

**Kriterien.** Mit diesen Optionen können Sie die Kriterien zum Hinzufügen und Entfernen von Feldern mit der Schätzmethode “Schrittweise” festlegen. (Die Schaltfläche ist deaktiviert, wenn die Methode “Einschluss” ausgewählt ist.) [Für weitere Informationen siehe Thema Schrittkriterien für Cox-Knoten auf S. 350.](#)

### Konvergenzkriterien für Cox-Knoten

Abbildung 10-58  
Cox-Regression – Dialogfeld “Konvergenzkriterien”



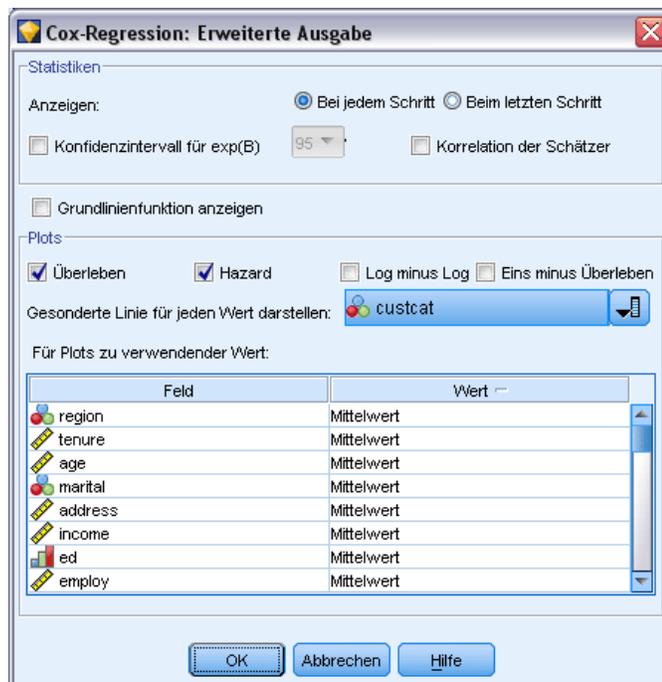
**Maximale Anzahl der Iterationen.** Hiermit können Sie die Maximalzahl der Iterationen für das Modell festlegen, die bestimmt, wie lange die Prozedur nach einer Lösung sucht.

**Log-Likelihood-Konvergenz.** Iterationen werden angehalten, wenn die relative Änderung der Log-Wahrscheinlichkeit (Log-Likelihood) kleiner als dieser Wert ist. Wenn der Wert gleich 0 ist, wird dieses Kriterium nicht angewendet.

**Parameter-Konvergenz.** Die Iterationen werden angehalten, wenn die absolute oder relative Änderung in den Parameterschätzungen kleiner als dieser Wert ist. Wenn der Wert gleich 0 ist, wird dieses Kriterium nicht angewendet.

### Cox-Knoten – Erweiterte Ausgabeoptionen

Abbildung 10-59  
Cox-Regression – Dialogfeld "Erweiterte Ausgabe"



**Statistiken.** Für die Modellparameter stehen Statistiken wie Konfidenzintervalle für  $\text{Exp}(B)$  und Korrelation der Schätzer zur Verfügung. Diese Statistiken können für jeden Schritt oder nur für den letzten Schritt angefordert werden.

**Grundlinienfunktion anzeigen.** Hiermit können Sie die Hazard-Grundlinienfunktion und die kumulative Überlebensverteilung beim Mittelwert der Kovariaten anzeigen lassen.

#### Diagramme

Diagramme können ein Hilfsmittel zur Bewertung des geschätzten Modells und zur Interpretation der Ergebnisse sein. Die Überlebens-, Hazard- und Log-Minus-Log-Funktionen sowie Eins minus Überleben können grafisch dargestellt werden.

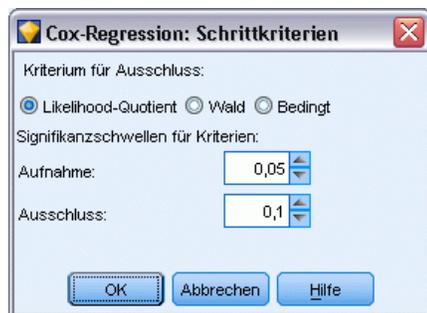
- **Überleben.** Zeigt die kumulative Überlebensfunktion auf einer linearen Skala an.
- **Hazard.** Zeigt die kumulative Hazard-Funktion auf einer linearen Skala an.
- **Log-minus-Log.** Zeigt die kumulierte Überlebensschätzung nach Anwendung der  $\ln(-\ln)$ -Transformation auf die Schätzung an.
- **Eins minus Überleben.** Erzeugt ein Diagramm der Werte "1 - Überlebensfunktion" auf einer linearen Skala.

**Gesonderte Linie für jeden Wert darstellen.** Diese Option ist nur für kategoriale Felder verfügbar.

**Für Plots zu verwendender Wert.** Da diese Funktionen von den Werten der Prädiktoren abhängen, müssen Sie konstante Werte für die Prädiktoren verwenden, um die Funktionen in Abhängigkeit von der Zeit grafisch darzustellen. In der Standardvorgabe wird der Mittelwert der einzelnen Prädiktoren als konstanter Wert verwendet. Sie können jedoch mithilfe des Gitters Ihre eigenen Werte für den Plot eingeben. Bei kategorialen Eingaben wird Indikatorkodierung verwendet, sodass ein Regressionskoeffizient für jede Kategorie (mit Ausnahme der letzten) vorhanden ist. Kategoriale Eingaben weisen also einen Mittelwert für jeden Indikatorkontrast auf, der gleich dem Anteil an Fällen in der Kategorie ist, die zum Indikatorkontrast gehört.

### Schrittkriterien für Cox-Knoten

Abbildung 10-60  
Cox-Regression – Dialogfeld “Schrittkriterien”



**Kriterium für Entfernen.** Wählen Sie Likelihood-Quotient für ein robusteres Modell. Zur Verkürzung der für die Modellerstellung erforderlichen Zeit können Sie Wald auswählen. Es gibt die zusätzliche Option Bedingt, die Ausschluss tests auf der Grundlage der Wahrscheinlichkeit der Likelihood-Quotienten-Statistik ermöglicht, welche auf bedingten Parameterschätzern beruht.

**Signifikanzschwellen für LR-Kriterien.** Mit dieser Option können Sie Auswahlkriterien auf der Grundlage der statistischen Wahrscheinlichkeit ( $p$ -Wert) angeben, die den einzelnen Feldern zugeordnet ist. Felder werden nur zum Modell hinzugefügt, wenn der zugehörige  $p$ -Wert kleiner ist als der Wert für Aufnahme, und nur dann entfernt, wenn der  $p$ -Wert größer ist als der Wert für Ausschluss. Der Wert für Aufnahme muss unter dem Wert für Ausschluss liegen.

## Einstelloptionen für Cox-Knoten

Abbildung 10-61  
Dialogfeld des Cox-Knotens, Registerkarte "Einstellungen"



**Überleben zu zukünftigen Zeitpunkten voraussagen.** Wählen Sie einen oder mehrere zukünftige Zeitpunkte aus. Das Überleben, also ob der jeweilige Fall mindestens die angegebene Zeitdauer (vom aktuellen Zeitpunkt gerechnet) überlebt, ohne dass das terminale Ereignis eintritt, wird für jeden Datensatz bei jedem Zeitwert vorhergesagt. Dabei wird jeweils eine Vorhersage pro Zeitwert erstellt. Beachten Sie, dass das Überleben der Wert "falsch" des Zielfelds ist.

- **Regelmäßige Intervalle.** Werte für die Überlebenszeit werden aus den angegebenen Werten für Zeitintervall und Anzahl der zu scorenden Zeitperioden erstellt. Wenn beispielsweise drei Zeitperioden mit dem Intervall 2 zwischen den einzelnen Zeitpunkten angefordert werden, wird das Überleben für die zukünftigen Zeitpunkte 2, 4, 6 vorhergesagt. Jeder Datensatz wird bei denselben Zeitwerten evaluiert.
- **Zeitfelder.** Für jeden Datensatz im ausgewählten Zeitfeld werden Überlebenszeiten angegeben (es wird genau ein Vorhersagefeld generiert), sodass die einzelnen Datensätze zu verschiedenen Zeitpunkten evaluiert werden können.

**Vergangene Überlebenszeit.** Dient zur Angabe der bisherigen Überlebenszeit des Datensatzes. Beispielsweise wird die bisherige Dauer des Geschäftsverhältnisses mit einem bestehenden Kunden als Feld angegeben. Das Scoring der Überlebenswahrscheinlichkeit zu einem zukünftigen Zeitpunkt hängt von der vergangenen Überlebenswahrscheinlichkeit ab.

*Anmerkung:* Die Werte für zukünftige und vergangene Überlebenszeiten müssen innerhalb des Bereichs für Überlebenszeiten in den zum Trainieren des Modells verwendeten Daten liegen. Datensätze, bei denen die Zeiten außerhalb dieses Bereichs liegen, werden als "null" gescort.

**Alle Wahrscheinlichkeiten ausgeben.** Gibt an, ob die Wahrscheinlichkeiten für die einzelnen Kategorien des Ausgabefelds zu den einzelnen vom Knoten verarbeiteten Datensätzen hinzugefügt werden. Wenn diese Option nicht ausgewählt wurde, wird nur die Wahrscheinlichkeit der vorhergesagten Kategorie hinzugefügt. Die Wahrscheinlichkeiten werden für jeden zukünftigen Zeitpunkt berechnet.

**Kumulative Hazardfunktion berechnen.** Gibt an, ob die Werte der kumulativen Hazardrate in jeden Datensatz aufgenommen werden sollen. Die kumulative Hazardrate wird für jeden zukünftigen Zeitpunkt berechnet.

## Cox-Modell-Nugget

Cox-Regressionsmodelle stehen für die Gleichungen, die durch die Cox-Knoten geschätzt wurden. Sie enthalten alle Informationen, die vom Modell erfasst wurden, sowie Informationen über die Struktur und Leistung des Modells.

Wenn Sie einen Stream ausführen, der ein generiertes Cox-Regressionsmodell enthält, fügt der Knoten zwei neue Felder hinzu, die die Prognose des Modells und die zugehörige Wahrscheinlichkeit enthalten. Die Namen der neuen Felder werden aus dem Namen des prognostizierten Ausgabefelds abgeleitet, dem  $\$C$ - für die vorhergesagte Kategorie und  $\$CP$ - für die zugehörige Wahrscheinlichkeit vorangestellt ist. Die Nummer des zukünftigen Zeitintervalls bzw. der Name des Zeitfelds, das das Zeitintervall definiert, ist als Suffix angegeben. Beispiel: Für das Ausgabefeld *churn* und zwei zukünftige Zeitintervalle, die in regelmäßigen Intervallen definiert sind, erhalten die neuen Felder die Namen  $\$C$ -*churn-1*,  $\$CP$ -*churn-1*,  $\$C$ -*churn-2* und  $\$CP$ -*churn-2*. Wenn zukünftige Zeitpunkte mit dem Zeitfeld *tenure* definiert sind, erhalten die neuen Felder die Namen  $\$C$ -*churn\_tenure* und  $\$CP$ -*churn\_tenure*.

Wenn Sie im Cox-Knoten die Einstellungsoption Alle Wahrscheinlichkeiten anhängen ausgewählt haben, werden für jeden zukünftigen Zeitpunkt zwei zusätzliche Felder hinzugefügt, die für jeden Datensatz die Wahrscheinlichkeiten für Überleben und Versagen enthalten. Diese zusätzlichen Felder werden anhand des Namens des Ausgabefelds benannt. Diesem wird  $\$CP$ -<Falsch-Wert>- für die Überlebenswahrscheinlichkeit und  $\$CP$ -<Wahr-Wert>- für die Wahrscheinlichkeit, dass das Ereignis eingetreten ist, vorangestellt und die Nummer des zukünftigen Zeitintervalls nachgestellt. Beispiel: Bei einem Ausgabefeld, bei dem der Wert "falsch" 0 und der Wert "wahr" 1 ist und zwei zukünftige Zeitintervalle in regelmäßigen Intervallen definiert sind, erhalten die neuen Felder die Namen  $\$CP$ -0-1,  $\$CP$ -1-1,  $\$CP$ -0-2 und  $\$CP$ -1-2. Wenn zukünftige Zeiten mit einem einzigen Zeitfeld *tenure* definiert sind, erhalten die neuen Felder die Namen  $\$CP$ -0-1 und  $\$CP$ -1-1, da nur ein einziges zukünftiges Intervall vorhanden ist.

Wenn Sie im Cox-Knoten die Einstellungsoption Kumulative Hazardfunktion berechnen ausgewählt haben, wird für jeden zukünftigen Zeitpunkt ein zusätzliches Feld hinzugefügt, das für jeden Datensatz die kumulative Hazard-Funktion enthält. Diese zusätzlichen Felder werden anhand des Namens des Ausgabefelds benannt. Diesem wird  $\$CH$ - vorangestellt und die Nummer des zukünftigen Zeitintervalls bzw. der Name, der das Zeitintervall definiert, nachgestellt. Beispiel: Für das Ausgabefeld *churn* und zwei zukünftige Zeitintervalle, die in regelmäßigen Intervallen definiert sind, erhalten die neuen Felder die Namen  $\$CH$ -*churn-1* und  $\$CH$ -*churn-2*. Wenn zukünftige Zeitpunkte mit dem Zeitfeld *tenure* definiert sind, erhält das neue Feld den Namen  $\$CH$ -*churn-1*.

## Cox-Regression – Ausgabeeinstellungen

Die Registerkarte “Einstellungen” des Nuggets enthält dieselben Steuerelemente wie die Registerkarte “Einstellungen” des Modellknotens. Die Standardwerte der Nugget-Steuerelemente richten sich nach den im Modellknoten festgelegten Werten. [Für weitere Informationen siehe Thema Einstellungsoptionen für Cox-Knoten auf S. 351.](#)

## Cox-Regression – Erweiterte Ausgabe

Die erweiterte Ausgabe für die Cox-Regression bietet detaillierte Informationen zum geschätzten Modell und dessen Leistung, einschließlich der Überlebenskurve. Die meisten der in der erweiterten Ausgabe enthaltenen Informationen weisen einen hohen Fachlichkeitsgrad auf und zur richtigen Interpretation dieser Ausgaben sind umfassende Kenntnisse der Cox-Regression erforderlich.

Abbildung 10-62  
Cox-Modell-Nugget – Registerkarte “Erweitert”

The screenshot shows the 'Scoring' dialog box with the 'Erweitert' tab selected. The 'Case Processing Summary' table is as follows:

		N	Percent
Cases available in analysis	Event(a)	274	27.4%
	Censored	726	72.6%
	Total	1000	100.0%
Cases dropped	Cases with missing values	0	.0%
	Cases with negative time	0	.0%
	Censored cases before the earliest event in a stratum	0	.0%
	Total	0	.0%
Total		1000	100.0%

a. Dependent Variable: Months with service

The 'Categorical Variable Codings' table for 'marital(t)' is as follows:

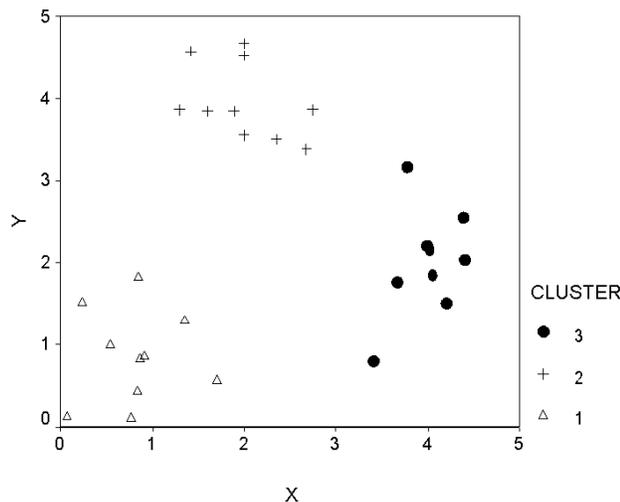
		Frequency	(1)(s)	(2)	(3)	(4)
marital(t)	0=Unmarried	505	1			
	1=Married	495	0			

# Cluster-Modelle

Clustermodelle konzentrieren sich auf die Ermittlung ähnlicher Datensätze und Beschriftung der Datensätze anhand der Gruppe, in die sie gehören. Dies geschieht ohne Vorkenntnisse zu Gruppen und ihren Eigenschaften. Vielleicht wissen Sie nicht einmal, nach wie vielen Gruppen Sie suchen sollen. Hierin liegt der Unterschied der Clustermodelle zu anderen Techniken des Maschinlernens—: Es gibt keine zuvor definierte Ausgabe und kein Zielfeld für das vorherzusagende Modell. Diese Modelle werden häufig als **nicht überwachte Lernmodelle** bezeichnet, da es keinen externen Standard gibt, mit dem die Klassifizierungsleistung des Modells beurteilt werden könnte. Für diese Modelle gibt es keine *richtigen* oder *falschen* Antworten. Ihr Wert wird durch die Möglichkeit bestimmt, interessante Gruppierungen in den Daten zu erfassen und sinnvolle Beschreibungen dieser Gruppierungen zu liefern.

Methoden zur Clusterbildung basieren auf dem Messen der Entfernungen zwischen Datensätzen und Clustern. Die Datensätze werden den Clustern auf eine Weise zugewiesen, die die Entfernung zwischen den Datensätzen minimiert, die demselben Cluster angehören.

Abbildung 11-1  
Einfaches Clusteringmodell



Es stehen drei Methoden zur Clusterbildung zur Verfügung:



Der K-Means-Knoten teilt das Daten-Set in unterschiedliche Gruppen (oder Cluster) auf. Bei diesem Verfahren wird eine festgelegte Anzahl von Clustern definiert, den Clustern werden iterativ Datensätze zugewiesen und die Cluster-Zentren werden angepasst, bis eine weitere Verfeinerung keine wesentliche Verbesserung des Modells mehr darstellen würde. Statt zu versuchen, ein Ergebnis vorherzusagen, versucht K-Means mithilfe eines als "nicht überwacht Lernen" bezeichneten Verfahrens Muster im Set der Eingabefelder zu entdecken. [Für weitere Informationen siehe Thema K-Means-Knoten auf S. 362.](#)



Der TwoStep-Knoten verwendet eine aus zwei Schritten bestehende Clusterbildungsmethode. Im ersten Schritt wird ein einzelner Durchlauf durch die Daten vorgenommen, bei dem die Eingangsrohdaten zu einem verwaltbaren Set von Unterclustern komprimiert werden. Im zweiten Schritt werden die Untercluster mithilfe einer hierarchischen Methode zur Clusterbildung nach und nach in immer größere Cluster zusammengeführt. TwoStep hat den Vorteil, dass die optimale Anzahl an Clustern für die Trainingsdaten automatisch geschätzt wird. Mit dem Verfahren können gemischte Feldtypen und große Daten-Sets effizient verarbeitet werden. [Für weitere Informationen siehe Thema TwoStep-Cluster-Knoten auf S. 367.](#)



Der Kohonen-Knoten erstellt eine Art von neuronalem Netzwerk, das verwendet werden kann, um ein Clustering der Datenmenge in einzelne Gruppen vorzunehmen. Wenn das Netz voll trainiert ist, sollten ähnliche Datensätze auf der Ausgabekarte eng nebeneinander stehen, während Datensätze, die sich unterscheiden, weit voneinander entfernt sein sollten. Die Zahl der von jeder Einheit im Modell-Nugget erfassten Beobachtungen gibt Aufschluss über die starken Einheiten. Dadurch wird ein Eindruck von der ungefähren Zahl der Cluster vermittelt. [Für weitere Informationen siehe Thema Kohonen-Knoten auf S. 355.](#)

Cluster-Modelle werden häufig verwendet, um Cluster oder Segmente zu erstellen, die dann als Eingaben in nachfolgenden Analysen verwendet werden. Ein häufiges Beispiel dafür sind die von Marktforschern verwendeten Marktsegmente, mit denen der Gesamtmarkt in homogene Untergruppen aufgeteilt wird. Jedes Segment weist besondere Eigenschaften auf, die sich auf den Erfolg der Marktforschung auswirken. Wenn Sie Data-Mining zur Optimierung Ihrer Marketingstrategie verwenden, können Sie Ihr Modell in der Regel erheblich verbessern, indem Sie die entsprechenden Segmente ermitteln und diese Segmentinformationen für Ihre Vorhersagemodelle verwenden.

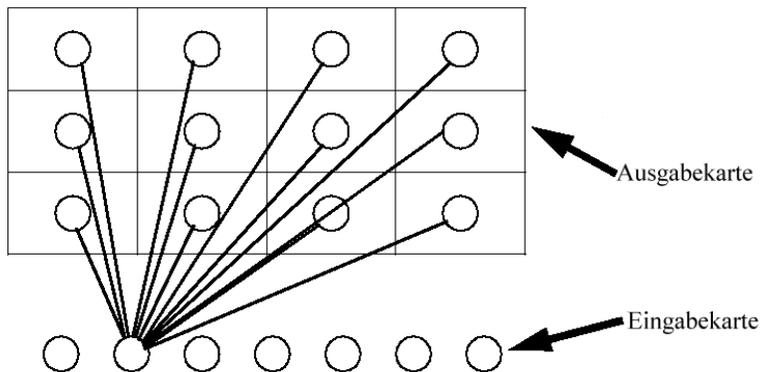
## Kohonen-Knoten

Kohonen-Netze stellen eine Form von neuronalen Netzen zur Clusterbildung dar. Sie sind auch bekannt unter dem Begriff **K-Netz (knet)** oder **SOM (selbstorganisierende Karte)**. Mit dieser Methode können Sie ein Clustering des Daten-Sets in einzelne Gruppen vornehmen, wenn Sie nicht wissen, wie diese Gruppen am Anfang aussehen. Datensätze werden in Gruppen zusammengefasst, wobei Datensätze innerhalb einer Gruppe oder eines Clusters ähnlich und Datensätze in verschiedenen Gruppen unterschiedlich sind.

Die Grundeinheiten sind **Neuronen**, die in zwei Schichten unterteilt sind: die **Eingabeschicht** und die **Ausgabeschicht** (auch als **Ausgabekarte** bezeichnet). Alle Eingabeneuronen sind mit allen Ausgabeneuronen verbunden. Mit diesen Verbindungen sind **Stärken** oder **Gewichtungen** verknüpft. Während des Trainings wetteifert jede Einheit mit allen anderen, um einen Datensatz zu "gewinnen".

Die Ausgabekarte ist ein zweidimensionales Neuronengitter ohne Verbindungen zwischen den Einheiten. Nachfolgend wird eine Karte im Format  $3 \times 4$  dargestellt, obwohl die Karten in der Regel größer sind.

Abbildung 11-2  
Struktur eines Kohonen-Netzes



Alle Eingabeneuronen sind mit allen Neuronen in der Ausgabekarte verbunden. (Zur Vereinfachung nicht angezeigt.)

Die Eingabedaten werden der Eingabeschicht präsentiert und die Werte an die Ausgabeschicht weitergeleitet. Das Ausgabeneuron mit der stärksten Reaktion soll der **Gewinner** sein und ist die Antwort für diese Eingabe.

Anfänglich sind alle Gewichtungen zufällig. Wenn eine Einheit einen Datensatz gewinnt, werden die Gewichtungen (zusammen mit denen anderer Nachbareinheiten, die kollektiv als **Nachbarschaft** bezeichnet werden) so angepasst, dass sie dem Muster der Prädiktorwerte für diesen Datensatz besser entsprechen. Alle Eingabedatensätze werden angezeigt und die Gewichtungen entsprechend aktualisiert. Dieser Vorgang wird viele Male wiederholt, bis die Änderungen nur noch gering sind. Während des Trainings werden die Gewichtungen an den Gittereinheiten so angepasst, dass sie eine zweidimensionale "Karte" der Cluster bilden (deshalb der Begriff **selbstorganisierende Karte**).

Wenn das Netz voll trainiert ist, sollten ähnliche Datensätze auf der Ausgabekarte eng nebeneinander stehen, während Datensätze, die sich stark unterscheiden, weit voneinander entfernt sein sollten.

Im Gegensatz zu den meisten Lernmethoden in IBM® SPSS® Modeler verwenden Kohonen-Netzwerke *kein* Zielfeld. Diese Art des Lernens, d. h. ohne Zielfeld, wird als **nicht überwacht Lernen** bezeichnet. Statt zu versuchen, ein Ergebnis vorherzusagen, versuchen Kohonen-Netzwerke, Muster im Set der Eingabefelder zu entdecken. In der Regel weist ein Kohonen-Netz schließlich einige Einheiten auf, die viele Beobachtungen zusammenfassen (**starke** Einheiten), und mehrere Einheiten, die keiner Beobachtung wirklich entsprechen (**schwache** Einheiten). Die starken Einheiten (und mitunter benachbarte Einheiten im Raster) repräsentieren mögliche Cluster-Zentren.

Eine weitere Einsatzmöglichkeit von Kohonen-Netzwerken findet sich bei der **Dimensionsreduzierung**. Das räumliche Merkmal des zweidimensionalen Rasters bietet eine Zuordnung der ursprünglichen  $k$ -Prädiktoren zu zwei abgeleiteten Funktionen, die die Ähnlichkeitsbeziehung in den ursprünglichen Prädiktoren bewahren. In einigen Fällen kann dies ebenso vorteilhaft sein wie die Faktoranalyse oder PCA.

Beachten Sie, dass die Methode zur Berechnung der Standardgröße des Ausgaberrasters sich im Vergleich zu früheren Versionen von SPSS Modeler geändert hat. Mit der neuen Methode werden im Allgemeinen kleinere Ausgabeschichten erzielt, die schneller zu trainieren sind und

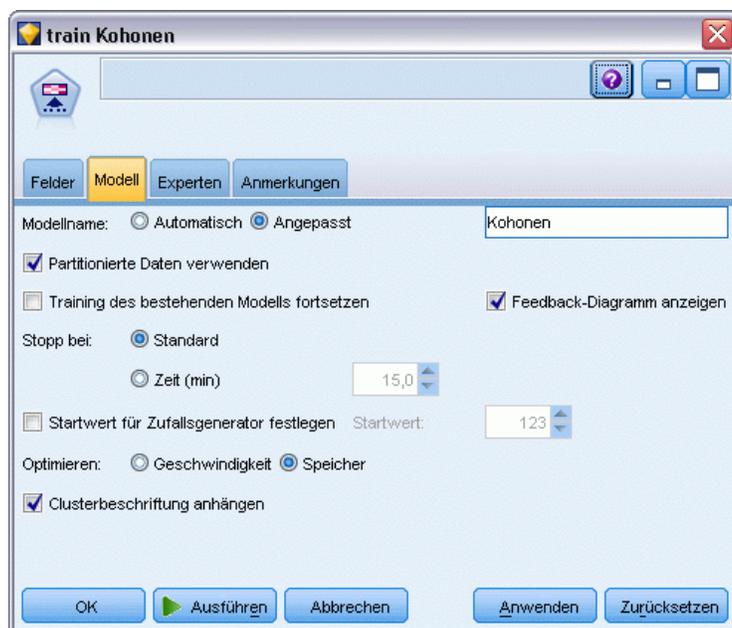
besser generalisieren. Wenn Sie mit der Standardgröße schlechte Ergebnisse erzielen, erhöhen Sie den Wert für die Größe des Ausgabereasters auf der Registerkarte "Experten". [Für weitere Informationen siehe Thema Expertenoptionen für den Kohonen-Knoten auf S. 359.](#)

**Anforderungen.** Zum Trainieren eines Kohonen-Netzes ist mindestens ein Feld mit der Rolle *Eingabe* erforderlich. Felder, deren Rolle auf *Ziel*, *Beides* oder *Keine* festgelegt ist, werden ignoriert.

**Stärken.** Um ein Kohonen-Netzwerkmodell zu erstellen, sind keine Daten über die Gruppenzugehörigkeit erforderlich. Auch die Anzahl Gruppen muss für die Suche nicht bekannt sein. Kohonen-Netzwerke beginnen mit einer großen Anzahl von Einheiten, und mit Fortschreiten des Trainings gravitieren die Einheiten zu natürlichen Clustern in den Daten. Die Zahl der von jeder Einheit erfassten Beobachtungen im Modell-Nugget gibt Aufschluss über die starken Einheiten, die einen Eindruck von der ungefähren Zahl der Cluster vermitteln.

## Optionen des Kohonen-Knotenmodells

Abbildung 11-3  
Optionen des Kohonen-Knotenmodells



**Modellname.** Sie können den Modellnamen automatisch basierend auf den Ziel- oder ID-Feldnamen (oder dem Modelltyp in Fällen, in denen kein solches Feld angegeben ist) generieren oder einen benutzerdefinierten Namen eingeben.

**Partitionierte Daten verwenden.** Wenn ein Partitionsfeld definiert ist, gewährleistet diese Option, dass nur Daten aus der Trainingspartition für die Modellerstellung verwendet werden. [Für weitere Informationen siehe Thema Partitionsknoten in Kapitel 4 in IBM SPSS Modeler 14.2- Quellen-, Prozess- und Ausgabeknoten.](#)

**Training des bestehenden Modells fortsetzen.** Standardmäßig wird bei jeder Ausführung eines Kohonen-Knotens ein komplett neues Netzwerk aufgebaut. Bei Auswahl dieser Option wird das Training mit dem letzten, vom Knoten erfolgreich aufgebauten Netz fortgesetzt.

**Feedback-Diagramm anzeigen.** Bei Auswahl dieser Option wird während des Trainings eine visuelle Darstellung des zweidimensionalen Arrays angezeigt. Die Stärke der einzelnen Knoten wird farblich dargestellt. Rot kennzeichnet eine Einheit mit vielen Datensätzen (eine **starke** Einheit), während Weiß auf eine Einheit hinweist, die wenige oder keine Datensätze enthält (eine **schwache** Einheit). Feedback wird möglicherweise nicht angezeigt, wenn die Zeit für die Modellerstellung relativ kurz ist. Diese Funktion kann die Trainingszeit verlangsamen. Deaktivieren Sie diese Option, wenn Sie die Trainingszeit beschleunigen möchten.

Abbildung 11-4  
Kohonen-Feedback-Diagramm



**Stopp bei.** Das Standardgrenzkriterium stoppt das Training basierend auf internen Parametern. Sie können auch eine Zeit als Grenzkriterium festlegen. Geben Sie die Zeit (in Minuten) für das Training des Netzwerks ein.

**Startwert für Zufallsgenerator festlegen.** Wenn kein Startwert für den Zufallsgenerator festgelegt wurde, ist die Sequenz der Zufallswerte, mit denen die Netzwerkgewichtungen initialisiert werden, bei jeder Ausführung des Knotens unterschiedlich. Dadurch ist es möglich, dass der Knoten bei verschiedenen Ausführungen unterschiedliche Modelle erstellt, auch wenn die Knoteneinstellungen und Datenwerte vollkommen identisch sind. Wenn Sie diese Option auswählen, können Sie den Startwert für den Zufallsgenerator auf einen bestimmten Wert festlegen, sodass das entstehende Modell genau reproduziert werden kann. Ein bestimmter Startwert für den Zufallsgenerator erzeugt immer dieselbe Sequenz der zufälligen Werte. In diesem Fall führt die Ausführung des Knotens immer zu demselben generierten Modell.

*Anmerkung:* Bei Verwendung der Option Startwert für Zufallsgenerator festlegen mit Datensätzen, die aus einer Datenbank eingelesen wurden, ist möglicherweise vor der Stichprobenziehung ein Sortierknoten erforderlich, um zu gewährleisten, dass bei jeder Ausführung des Knotens dasselbe Ergebnis erzielt wird. Dies liegt daran, dass der Startwert für den Zufallsgenerator von der Reihenfolge der Datensätze abhängt, die in relationalen Datenbanken nicht unbedingt gleich bleibt. Für weitere Informationen siehe Thema Sortierknoten in Kapitel 3 in *IBM SPSS Modeler 14.2- Quellen-, Prozess- und Ausgabeknoten*.

*Hinweis:* Wenn Sie nominale (Set-) Felder in Ihr Modell integrieren möchten, jedoch Speicherprobleme beim Erstellen des Modells haben oder der Aufbau des Modells zu viel Zeit in Anspruch nimmt, kodieren Sie große Set-Felder um, um die Zahl der Werte zu verringern, oder verwenden Sie ein anderes Feld mit weniger Werten als Proxy für das große Set. Wenn Sie beispielsweise Probleme mit einem *product\_id*-Feld haben, das Werte für einzelne Produkte enthält, können Sie es aus dem Modell entfernen und stattdessen ein weniger detailliertes *product\_category*-Feld hinzufügen.

**Optimieren.** Wählen Sie Optionen, die die Leistung während der Modellerstellung basierend auf Ihren persönlichen Anforderungen erhöhen.

- Wählen Sie Geschwindigkeit, um den Algorithmus anzuweisen, zur Verbesserung der Leistung keinen Datenträgerüberlauf zu verwenden.
- Wählen Sie Speicher, um den Algorithmus anzuweisen, ggf. einen Datenträgerüberlauf zu verwenden und dafür Geschwindigkeitseinbußen hinzunehmen. Diese Option ist standardmäßig aktiviert.

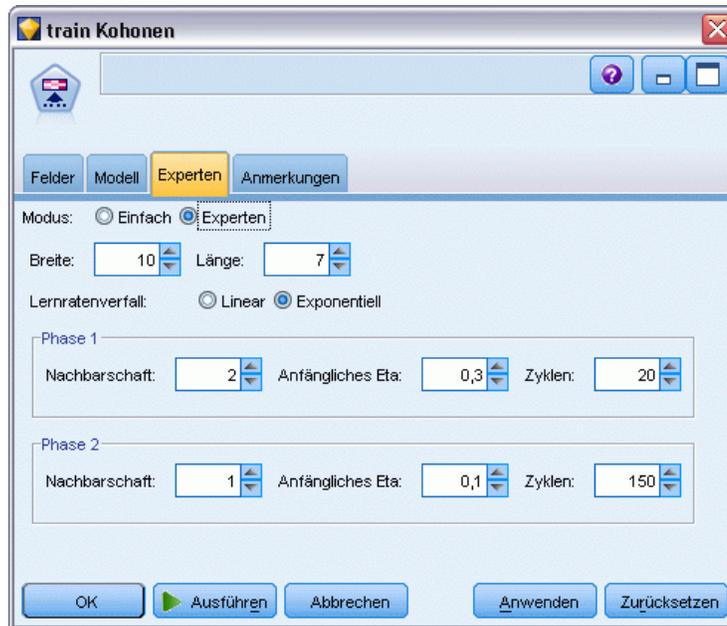
*Anmerkung:* Bei der Ausführung im verteilten Modus kann diese Einstellung durch die in der Datei *options.cfg* angegebenen Administratoroptionen überschrieben werden. [Für weitere Informationen siehe Thema Verwenden der Datei options.cfg in Kapitel 4 in IBM SPSS Modeler Server 14.2-Verwaltungs- und Leistungshandbuch.](#)

**Cluster-Beschriftung anhängen.** Diese standardmäßig für neue Modelle ausgewählte Option, die aber für Modelle aus früheren IBM® SPSS® Modeler-Versionen deaktiviert ist, erstellt ein einzelnes kategoriales Score-Feld vom selben Typ, der sowohl vom Knoten “K-Means” als auch vom Knoten “TwoStep” erstellt wird. Dieses Zeichenkettenfeld wird im Knoten “Autom. Cluster” verwendet, wenn Rangmessungen für die unterschiedlichen Modelltypen berechnet werden. [Für weitere Informationen siehe Thema Knoten “Autom. Cluster” in Kapitel 5 auf S. 114.](#)

## ***Expertenoptionen für den Kohonen-Knoten***

Für Anwender mit detaillierten Kenntnissen über Kohonen-Netzwerke stehen Expertenoptionen zur Verfügung, mit denen der Trainingsprozess verfeinert wird. Für den Zugriff auf die Expertenoptionen stellen Sie den Modus auf der Registerkarte “Experten” auf Experten ein.

Abbildung 11-5  
Kohonen-Expertenoptionen



**Breite und Länge.** Geben Sie die Größe (Breite und Länge) der zweidimensionalen Ausgabekarte als Anzahl der Ausgabeeinheiten für jede Dimension an.

**Lernratenverfall** Wählen Sie entweder den linearen oder exponentiellen Lernratenverfall. Die **Lernrate** ist ein Gewichtungsfaktor, der im Laufe der Zeit abnimmt, sodass das Netzwerk mit der Verschlüsselung weiträumiger Funktionen der Daten beginnt und sich allmählich auf feinere Details konzentriert.

**Phase 1 und Phase 2.** Das Kohonen-Netztraining gliedert sich in zwei Phasen. Phase 1 ist eine grobe Schätzphase, in der die groben Muster der Daten erfasst werden. Phase 2 ist eine Abstimmungsphase, in der die Zuordnung angepasst wird, um die feineren Merkmale der Daten zu modellieren. Für jede Phase gibt es drei Parameter:

- **Nachbarschaft.** Legt die Startgröße (Radius) der Nachbarschaft fest. Dieser Wert bestimmt die Zahl der “benachbarten” Einheiten, die beim Training zusammen mit der gewonnenen Einheit aktualisiert werden. In Phase 1 beginnt die Größe der Nachbarschaft mit *Phase 1 Nachbarschaft* und sinkt auf (*Phase 2 Nachbarschaft* + 1) ab. In Phase 2 liegt der Startwert für die Nachbarschaft bei *Phase 2 Nachbarschaft* und nimmt bis auf 1.0 ab. *Phase 1 Nachbarschaft* sollte größer sein als *Phase 2 Nachbarschaft*.
- **Anfängliches Eta.** Legt den Startwert für die Lernrate **Eta** fest. In Phase 1 beginnt der Eta-Wert bei *Phase 1 Anfängliches Eta* und sinkt auf *Phase 2 Anfängliches Eta* ab. In Phase 2 liegt der Eta-Anfangswert bei *Phase 2 Anfängliches Eta* und nimmt bis auf 0 ab. *Phase 1 Anfängliches Eta* sollte größer sein als *Phase 2 Anfängliches Eta*.
- **Zyklen.** Legt die Anzahl der Zyklen für jede Trainingsphase fest. Jede Phase hält für die Dauer der angegebenen Zahl von Durchläufen durch die Daten an.

## **Modell-Nuggets vom Typ "Kohonen"**

Modell-Nuggets von Kohonen-Modellen enthalten alle Informationen, die vom trainierten Kohonen-Netzwerk erfasst wurden, sowie Informationen über die Architektur des Netzwerks.

Wenn Sie einen Stream mit einem Kohonen-Modell-Nugget auswählen, fügt der Knoten zwei neue Felder mit den  $X$ - und  $Y$ -Koordinaten der Einheit im Kohonen-Ausgabegeraster hinzu, die am stärksten auf diesen Datensatz reagiert hat. Die neuen Feldnamen werden durch Präfigierung von  $\$KX$ - und  $\$KY$ - aus dem Modellnamen abgeleitet. Beispiel: Wenn das Modell den Namen *Kohonen* trägt, erhalten die neuen Felder die Namen  $\$KX$ -*Kohonen* und  $\$KY$ -*Kohonen*.

Für ein besseres Verständnis davon, was im Kohonen-Netz kodiert wurde, klicken Sie auf die Registerkarte "Modell" im Modell-Nugget-Browser. Dadurch wird der Cluster-Viewer angezeigt, der eine grafische Darstellung der Cluster, Felder und Wichtigkeitsniveaus bietet. [Für weitere Informationen siehe Thema Cluster-Viewer -Registerkarte Modell auf S. 372.](#)

Wenn Sie es vorziehen, die Cluster als Gitter zu visualisieren, können Sie die Ergebnisse des Kohonen-Netzes anzeigen, indem Sie mithilfe eines Plot-Knotens einen Plot der  $\$KX$ - und  $\$KY$ -Felder erstellen. (Sie sollten  $X$ -Bewegung und  $Y$ -Bewegung im Plotknoten auswählen, um zu verhindern, dass die Datensätze der einzelnen Einheiten alle übereinander abgebildet werden.) Im Plot können Sie außerdem ein symbolisches Feld überlagern, um festzustellen, wie das Kohonen-Netz die Daten in Cluster eingeteilt hat.

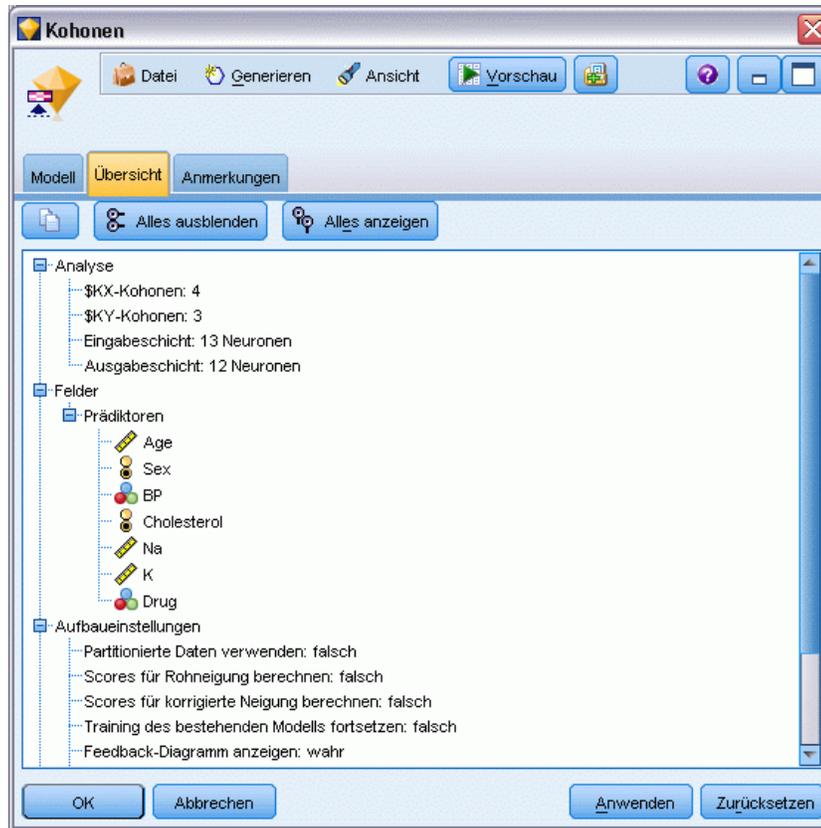
Ein weiteres leistungsstarkes Verfahren zur Gewinnung eines Einblicks in das Kohonen-Netzwerk besteht in der Verwendung der Regelinduktion zur Ermittlung der Merkmale, die die vom Netzwerk gefundenen Cluster unterscheiden. [Für weitere Informationen siehe Thema C5.0-Knoten in Kapitel 6 auf S. 177.](#)

Allgemeine Informationen zur Verwendung des Modellbrowsers finden Sie hier: [Durchsuchen von Modell-Nuggets](#)

## **Übersicht über das Kohonen-Modell**

Auf der Registerkarte "Übersicht" für ein Kohonen-Modell-Nugget werden Informationen über die Architektur bzw. Topologie des Netzwerks angezeigt. Die Länge und Breite der zweidimensionalen Kohonen-Funktionskarte (Ausgabeschicht) werden als  $\$KX$ -*Modellname* und  $\$KY$ -*Modellname* angezeigt. Für Ein- und Ausgabeschicht wird jeweils die Anzahl der Einheiten in der Schicht aufgeführt.

Abbildung 11-6  
Modell-Nugget vom Typ "Kohonen", Registerkarte "Übersicht"



## K-Means-Knoten

Der K-Means-Knoten bietet eine Methode der **Cluster-Analyse**. Mit dieser Methode können Sie ein Clustering der Daten-Sets in einzelne Gruppen vornehmen, wenn Sie nicht wissen, wie diese Gruppen am Anfang aussehen. Im Gegensatz zu den meisten Lernmethoden in IBM® SPSS® Modeler verwenden K-Means-Modelle *kein* Zielfeld. Diese Art des Lernens, d. h. ohne Zielfeld, wird als **nicht überwachtetes Lernen** bezeichnet. Statt zu versuchen, ein Ergebnis vorherzusagen, versuchen K-Means-Knoten, Muster im Set der Eingabefelder zu entdecken. Datensätze werden in Gruppen zusammengefasst, wobei Datensätze innerhalb einer Gruppe oder eines Clusters ähnlich und Datensätze in verschiedenen Gruppen unterschiedlich sind.

K-Means definiert einen Set von Cluster-Startzentren, die von Daten abgeleitet werden. Anschließend werden die einzelnen Datensätze basierend auf ihren Eingabefeldwerten dem Cluster zugewiesen, dem sie am meisten ähneln. Nachdem alle Datensätze zugewiesen wurden, werden die Cluster-Zentren aktualisiert, um die neuen Datensatz-Sets, die den einzelnen Clustern zugewiesen wurden, wiederzugeben. Die Datensätze werden nun erneut überprüft, um festzustellen, ob sie einem anderen Cluster zugewiesen werden sollten. Der Prozess der Datensatzzuweisung bzw. Cluster-Iteration wird so lange fortgesetzt, bis die maximale Anzahl an Iterationen erreicht ist oder die Änderung von einer Iteration auf die nächste einen bestimmten Schwellenwert nicht überschreitet.

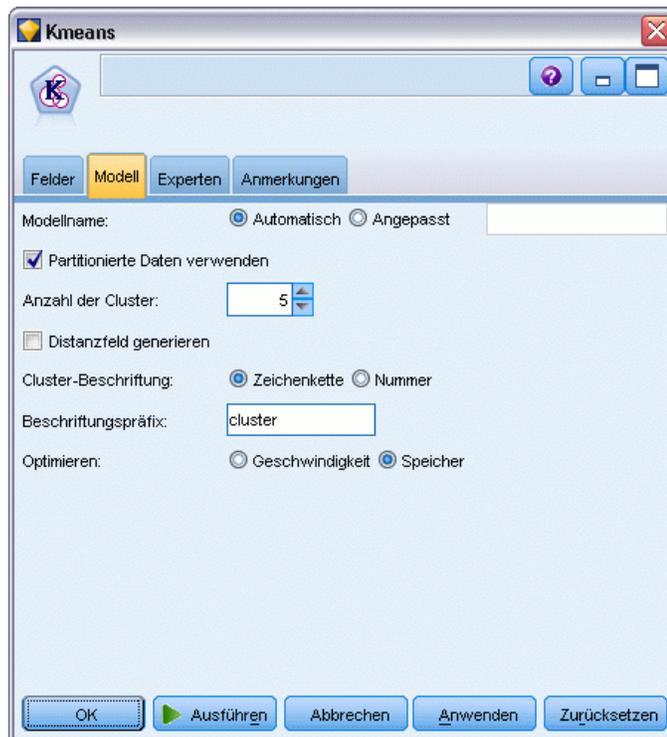
*Hinweis:* Das entstehende Modell hängt bis zu einem gewissen Grad von der Reihenfolge der Trainingsdaten ab. Eine Änderung der Datenreihenfolge und ein erneutes Erstellen des Modells kann zu einem anderen endgültigen Cluster-Modell führen.

**Anforderungen.** Zum Trainieren eines K-Means-Modells ist mindestens ein Feld mit der Rolle *Eingabe* erforderlich. Felder, deren Rolle auf *Ausgabe*, *Beides* oder *Keine* festgelegt ist, werden ignoriert.

**Stärken.** Um ein K-Means-Modell zu erstellen, sind keine Daten über die Gruppenzugehörigkeit erforderlich. Das K-Means-Modell stellt häufig die schnellste Methode zur Clusterbildung von großen Daten-Sets dar.

## Optionen für K-Means-Knotenmodelle

Abbildung 11-7  
Optionen für K-Means-Knotenmodelle



**Modellname.** Sie können den Modellnamen automatisch basierend auf den Ziel- oder ID-Feldnamen (oder dem Modelltyp in Fällen, in denen kein solches Feld angegeben ist) generieren oder einen benutzerdefinierten Namen eingeben.

**Partitionierte Daten verwenden.** Wenn ein Partitionsfeld definiert ist, gewährleistet diese Option, dass nur Daten aus der Trainingspartition für die Modellerstellung verwendet werden. [Für weitere Informationen siehe Thema Partitionsknoten in Kapitel 4 in IBM SPSS Modeler 14.2- Quellen-, Prozess- und Ausgabeknoten.](#)

**Angegebene Anzahl der Cluster.** Geben Sie die Anzahl der zu generierenden Cluster an. Der Standardwert ist 5.

**Distanzfeld generieren.** Wenn diese Option ausgewählt ist, enthält das Modell-Nugget ein Feld mit der Distanz jedes einzelnen Datensatzes vom Zentrum des entsprechend zugewiesenen Clusters.

**Clusterbeschriftung.** Geben Sie das Format für die Werte im generierten Feld "Cluster-Zugehörigkeit" an. Die Cluster-Zugehörigkeit kann als Zeichenkette mit dem festgelegten Beschriftungspräfix (z. B. "Cluster 1", "Cluster 2" usw.) oder als Zahl angegeben werden.

*Hinweis:* Wenn Sie nominale (Set-) Felder in Ihr Modell integrieren möchten, jedoch Speicherprobleme beim Erstellen des Modells haben oder der Aufbau des Modells zu viel Zeit in Anspruch nimmt, kodieren Sie große Set-Felder um, um die Zahl der Werte zu verringern, oder verwenden Sie ein anderes Feld mit weniger Werten als Proxy für das große Set. Wenn Sie beispielsweise Probleme mit einem *product\_id*-Feld haben, das Werte für einzelne Produkte enthält, können Sie es aus dem Modell entfernen und stattdessen ein weniger detailliertes *product\_category*-Feld hinzufügen.

**Optimieren.** Wählen Sie Optionen, die die Leistung während der Modellerstellung basierend auf Ihren persönlichen Anforderungen erhöhen.

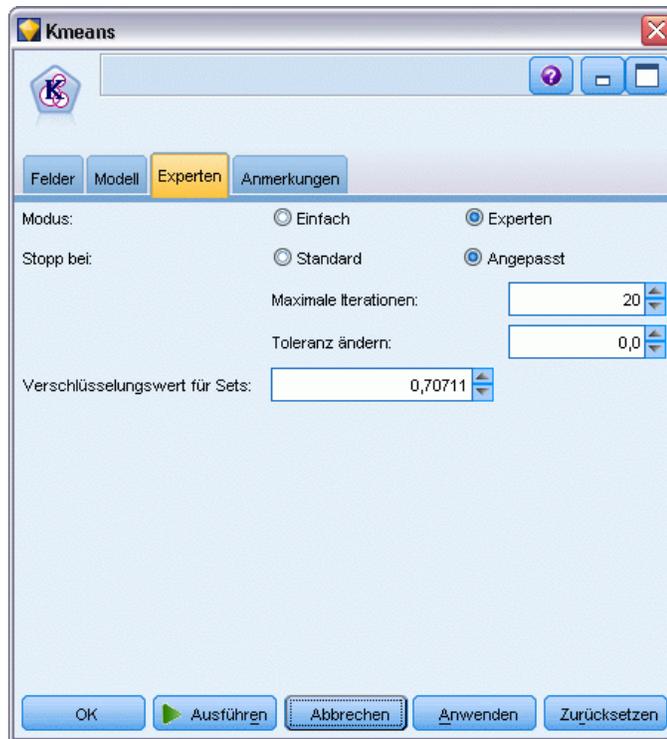
- Wählen Sie Geschwindigkeit, um den Algorithmus anzuweisen, zur Verbesserung der Leistung keinen Datenträgerüberlauf zu verwenden.
- Wählen Sie Speicher, um den Algorithmus anzuweisen, ggf. einen Datenträgerüberlauf zu verwenden und dafür Geschwindigkeitseinbußen hinzunehmen. Diese Option ist standardmäßig aktiviert.

*Anmerkung:* Bei der Ausführung im verteilten Modus kann diese Einstellung durch die in der Datei *options.cfg* angegebenen Administratoroptionen überschrieben werden. [Für weitere Informationen siehe Thema Verwenden der Datei options.cfg in Kapitel 4 in IBM SPSS Modeler Server 14.2-Verwaltungs- und Leistungshandbuch.](#)

## **Expertenoptionen für K-Means-Knoten**

Für Anwender mit detaillierten Kenntnissen über K-Means-Clusterbildung stehen Expertenoptionen zur Verfügung, mit denen der Trainingsprozess verfeinert wird. Für den Zugriff auf die Expertenoptionen stellen Sie den Modus auf der Registerkarte "Experten" auf Experten ein.

Abbildung 11-8  
K-Means-Expertenoptionen



**Stopp bei.** Geben Sie das Grenzkriterium für das Training des Modells an. Das Grenzkriterium Standard beträgt 20 Iterationen oder eine Änderung  $< 0,000001$ , je nachdem, was zuerst eintritt. Wenn Sie eigene Grenzkriterien angeben möchten, wählen Sie Benutzerdefiniert.

- **Maximalzahl der Iterationen.** Mit dieser Option können Sie das Modelltraining nach der angegebenen Anzahl von Iterationen beenden.
- **Toleranz ändern.** Mit dieser Option können Sie das Modelltraining beenden, sobald die größte Änderung der Cluster-Zentren für eine Iteration kleiner ist als das angegebene Niveau.

**Verschlüsselungswert für Sets.** Geben Sie einen Wert zwischen 0 und 1,0 für die Umkodierung von Set-Feldern als Gruppe von numerischen Feldern an. Der Standardwert ist die Wurzel aus 0,5 (rund 0,707107). Dieser Wert bietet die richtige Gewichtung für umkodierte Flag-Felder. Werte, die näher an 1,0 liegen, gewichten Set-Felder stärker als numerische Felder.

## Modell-Nuggets vom Typ "K-Means"

Modell-Nuggets für K-Means-Modelle enthalten alle Informationen, die vom Cluster-Modell erfasst wurden, sowie Informationen zu den Trainingsdaten und dem Schätzvorgang.

Wenn Sie einen Stream ausführen, der ein Modell-Nugget vom Typ "K-Means" enthält, fügt dieser Knoten zwei neue Felder hinzu, die die Cluster-Mitgliedschaft und die Entfernung vom zugewiesenen Cluster-Zentrum für den betreffenden Datensatz enthalten. Die neuen Feldnamen werden durch Präfigierung von  $\$KM-$  für die Cluster-Mitgliedschaft und  $\$KMD-$  für die

Entfernung vom Cluster-Zentrum aus dem Modellnamen abgeleitet. Beispiel: Wenn das Modell den Namen *Kmeans* trägt, erhalten die neuen Felder die Namen *\$KM-Kmeans* und *\$KMD-Kmeans*.

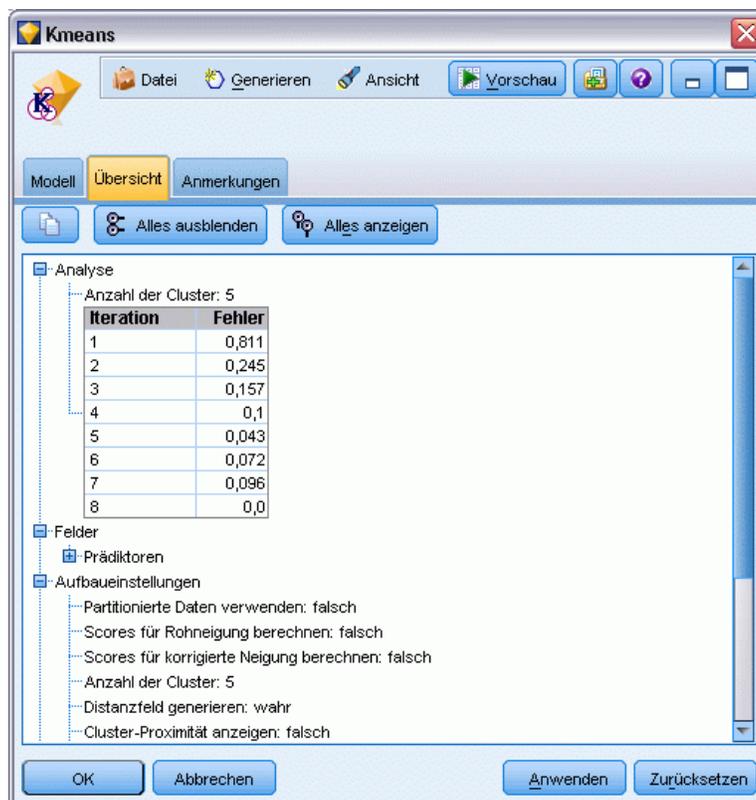
Ein leistungsstarkes Verfahren zur Gewinnung eines Einblicks in das K-Means-Netzwerk besteht in der Verwendung der Regelinduktion zur Ermittlung der Merkmale, die die vom Modell gefundenen Cluster unterscheiden. [Für weitere Informationen siehe Thema C5.0-Knoten in Kapitel 6 auf S. 177.](#) Außerdem können Sie im Modell-Nugget-Browser auf die Registerkarte “Modell” klicken, um den Cluster-Viewer anzuzeigen, der eine grafische Darstellung der Cluster, Felder und Wichtigkeitsniveaus bietet. [Für weitere Informationen siehe Thema Cluster-Viewer -Registerkarte Modell auf S. 372.](#)

Allgemeine Informationen zur Verwendung des Modellbrowsers finden Sie hier: [Durchsuchen von Modell-Nuggets](#)

### **Übersicht über das K-Means-Modell**

Auf der Registerkarte “Übersicht” eines Modell-Nuggets vom Typ “K-Means” finden Sie Informationen zu den Trainingsdaten, dem Schätzvorgang und den durch das Modell definierten Clustern. Die Anzahl der Cluster sowie das Iterationsprotokoll werden angezeigt. Wenn Sie einen Analyseknoden ausgeführt haben, der an diesen Modellierungsknoten angehängt ist, werden die Informationen aus dieser Analyse ebenfalls in diesem Abschnitt angezeigt. [Für weitere Informationen siehe Thema Analyseknoden in Kapitel 6 in IBM SPSS Modeler 14.2- Quellen-, Prozess- und Ausgabeknoten.](#)

Abbildung 11-9  
Modell-Nugget vom Typ "K-Means", Registerkarte "Übersicht"



## TwoStep-Cluster-Knoten

Der TwoStep-Cluster-Knoten bietet eine Form der **Cluster-Analyse**. Mit dieser Methode können Sie ein Clustering der Daten-Sets in einzelne Gruppen vornehmen, wenn Sie nicht wissen, wie diese Gruppen am Anfang aussehen. Ebenso wie Kohonen-Knoten und K-Means-Knoten verwenden auch TwoStep-Cluster-Modelle *kein* Zielfeld. Statt zu versuchen, ein Ergebnis vorherzusagen, versuchen TwoStep-Cluster, Muster im Set der Eingabefelder zu entdecken. Datensätze werden in Gruppen zusammengefasst, wobei Datensätze innerhalb einer Gruppe oder eines Clusters ähnlich und Datensätze in verschiedenen Gruppen unterschiedlich sind.

Beim TwoStep-Cluster handelt es sich um eine Clusterbildungsmethode in zwei Schritten. Im ersten Schritt wird ein einzelner Durchlauf durch die Daten vorgenommen, bei dem die Eingaberohdaten zu einem verwaltbaren Set von Unterclustern komprimiert werden. Im zweiten Schritt wird eine hierarchische Clusterbildungsmethode verwendet, mit der die Untercluster zu immer größeren Clustern zusammengeführt werden. Dabei ist kein erneuter Durchlauf durch die Daten erforderlich. Die hierarchische Clusterbildung bietet den Vorteil, dass vorab keine Clusteranzahl ausgewählt werden muss. Bei vielen hierarchischen Methoden zur Clusterbildung werden einzelne Datensätze als Startcluster verwendet, die dann rekursiv zu noch größeren Clustern zusammengeführt werden. Diese Methoden versagen häufig bei großen Datenmengen. Durch die anfängliche Vorclusterbildung von TwoStep wird die hierarchische Clusterbildung hingegen auch für große Daten-Sets zu einem schnellen Verfahren.

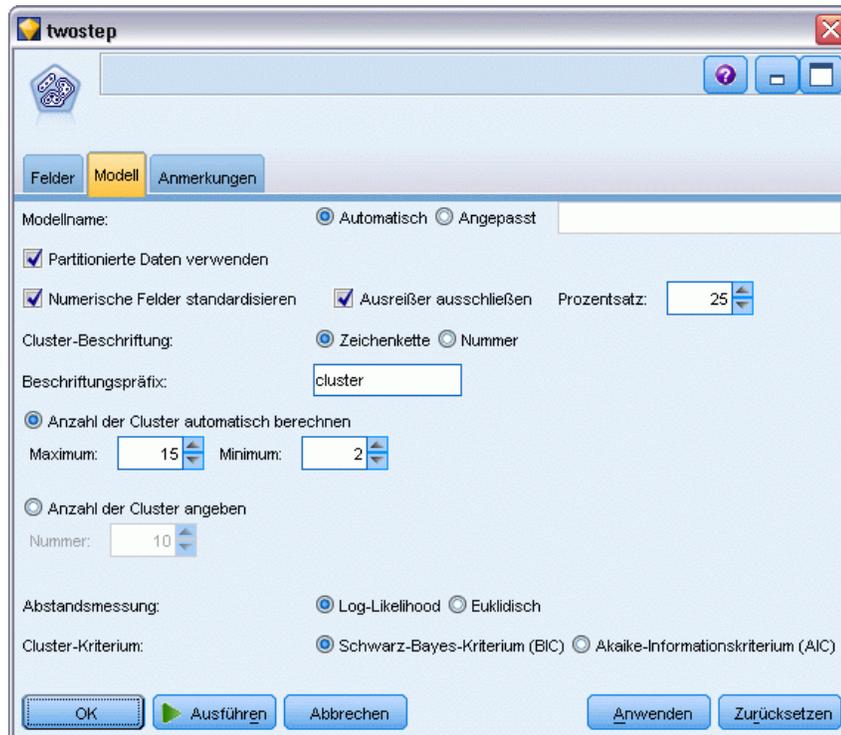
*Hinweis:* Das entstehende Modell hängt bis zu einem gewissen Grad von der Reihenfolge der Trainingsdaten ab. Eine Änderung der Datenreihenfolge und ein erneutes Erstellen des Modells kann zu einem anderen endgültigen Cluster-Modell führen.

**Anforderungen.** Zum Trainieren eines TwoStep-Cluster-Modells ist mindestens ein Feld mit der Rolle *Eingabe* erforderlich. Felder, deren Rolle auf *Ziel*, *Beides* oder *Keine* festgelegt ist, werden ignoriert. Der TwoStep-Cluster-Algorithmus verarbeitet keine fehlenden Werte. Bei der Modellerstellung werden Datensätze, die in einem der Eingabefelder Leerzeichen enthalten, ignoriert.

**Stärken.** TwoStep-Clustering kann gemischte Feldtypen verarbeiten und ist in der Lage, große Daten-Sets effizient zu verarbeiten. Es verfügt außerdem über die Fähigkeit, mehrere Clusterlösungen zu testen und die beste auszuwählen, sodass Sie nicht wissen müssen, wie viele Cluster Sie am Anfang abrufen müssen. TwoStep-Cluster können so eingestellt werden, dass **Ausreißer** oder äußerst unwahrscheinliche Fälle, die Ihre Ergebnisse verfälschen könnten, automatisch ausgeschlossen werden.

## Optionen für TwoStep-Cluster-Knotenmodelle

Abbildung 11-10  
Optionen für TwoStep-Cluster-Knotenmodelle



**Modellname.** Sie können den Modellnamen automatisch basierend auf den Ziel- oder ID-Feldnamen (oder dem Modelltyp in Fällen, in denen kein solches Feld angegeben ist) generieren oder einen benutzerdefinierten Namen eingeben.

**Partitionierte Daten verwenden.** Wenn ein Partitionsfeld definiert ist, gewährleistet diese Option, dass nur Daten aus der Trainingspartition für die Modellerstellung verwendet werden. [Für weitere Informationen siehe Thema Partitionsknoten in Kapitel 4 in IBM SPSS Modeler 14.2- Quellen-, Prozess- und Ausgabeknoten.](#)

**Numerische Felder standardisieren.** Standardmäßig nimmt TwoStep eine Standardisierung aller numerischen Eingabefelder auf dieselbe Skalierung vor, d. h. ein Mittelwert von 0 und eine Abweichung von 1. Wenn Sie die ursprüngliche Skalierung der numerischen Felder beibehalten möchten, deaktivieren Sie diese Option. Symbolische Felder sind davon nicht betroffen.

**Ausreißer ausschließen.** Wenn Sie diese Option wählen, werden Datensätze, die nicht in ein betrachtetes Cluster zu passen scheinen, automatisch von der Analyse ausgeschlossen. So wird eine Verfälschung des Ergebnisses durch derartige Fälle verhindert.

Die Erkennung von Ausreißern erfolgt während des Schritts der Vorclusterbildung. Bei Auswahl dieser Option werden Untercluster, die im Vergleich zu anderen Unterclustern wenige Datensätze enthalten, als potenzielle Ausreißer eingestuft und der Baum der Untercluster wird unter Ausschluss dieser Datensätze neu aufgebaut. Die Größe, unter der angenommen wird, dass Untercluster potenzielle Ausreißer enthalten, wird von der Option Prozent kontrolliert. Einige dieser potenziellen Ausreißerdatensätze können den neu erstellten Unterclustern hinzugefügt werden, wenn sie den neuen Untercluster-Profilen ausreichend ähneln. Die übrigen potenziellen Ausreißer, die nicht zusammengeführt werden können, werden als Ausreißer eingestuft, zu einem Cluster "Rauschen" hinzugefügt und vom Schritt der hierarchischen Clusterbildung ausgeschlossen.

Beim *Scoring* von Daten mit einem TwoStep-Modell, das Ausreißer verarbeitet, werden neue Fälle, die mehr als eine bestimmte Grenzdistanz (basierend auf der Log-Wahrscheinlichkeit) vom nächsten betrachteten Cluster entfernt sind, als Ausreißer eingestuft und dem Cluster "Rauschen" mit dem Namen -1 zugewiesen.

**Clusterbeschriftung.** Geben Sie das Format für das generierte Feld "Cluster-Zugehörigkeit" an. Die Cluster-Zugehörigkeit kann als Zeichenkette mit dem festgelegten Beschriftungspräfix (z. B. "Cluster 1", "Cluster 2" usw.) oder als Zahl angegeben werden.

**Anzahl der Cluster automatisch berechnen.** TwoStep-Cluster können eine große Anzahl von Cluster-Lösungen sehr rasch analysieren, um die optimale Anzahl von Clustern für die Trainingsdaten auszuwählen. Geben Sie einen Bereich der auszuprobierenden Lösungen an, indem Sie die maximale und minimale Anzahl der Cluster festlegen. TwoStep ermittelt die optimale Anzahl von Clustern in einem zweistufigen Prozess. In der ersten Stufe wird eine Obergrenze für die Anzahl der Cluster in dem Modell basierend darauf, wie sich das BIC (Bayes Information Criterion) beim Hinzufügen weiterer Cluster ändert, ausgewählt. In der zweiten Stufe wird die Änderung in der Mindestdistanz zwischen Clustern für alle Modelle mit weniger Clustern gesucht als die Mindest-BIC-Lösung. Das endgültige Cluster-Modell wird anhand der größten Distanzänderung ermittelt.

**Anzahl der Cluster angeben.** Wenn Sie wissen, wie viele Cluster in Ihr Modell einzubeziehen sind, wählen Sie diese Option und geben Sie die Anzahl der Cluster ein.

**Distanzmaß.** Mit dieser Auswahl legen Sie fest, wie Ähnlichkeiten zwischen zwei Clustern verarbeitet werden.

- **Log-Likelihood.** Mit dem Likelihood-Maß wird eine Wahrscheinlichkeitsverteilung für die Variablen vorgenommen. Bei stetigen Variablen wird von einer Normalverteilung, bei kategorialen Variablen von einer multinomialen Verteilung ausgegangen. Bei allen Variablen wird davon ausgegangen, dass sie unabhängig sind.
- **Euklidisch.** Das Euklidische Maß bezeichnet die “gerade” Distanz zwischen zwei Clustern. Es kann nur dann verwendet werden, wenn es sich bei sämtlichen Variablen um stetige Variablen handelt.

**Cluster-Kriterium.** Mit dieser Auswahl legen Sie fest, wie die Anzahl der Cluster vom automatischen Cluster-Algorithmus bestimmt wird. Angegeben werden kann entweder das Bayes-Informationskriterium (BIC) oder das Akaikes-Informationskriterium (AIC).

## **Modell-Nuggets vom Typ “TwoStep-Cluster”**

Modell-Nuggets für TwoStep-Cluster-Modelle enthalten alle Informationen, die vom Cluster-Modell erfasst wurden, sowie Informationen zu den Trainingsdaten und dem Schätzvorgang.

Wenn Sie einen Stream ausführen, der ein Modell-Nugget vom Typ “TwoStep-Cluster” enthält, fügt der Knoten ein neues Feld hinzu, das die Cluster-Mitgliedschaft für den betreffenden Datensatz enthält. Der neue Feldname wird durch Präfigierung von  $\$T$ - aus dem Modellnamen abgeleitet. Beispiel: Wenn das Modell den Namen *TwoStep* trägt, erhält das neue Feld den Namen  $\$T$ -*TwoStep*.

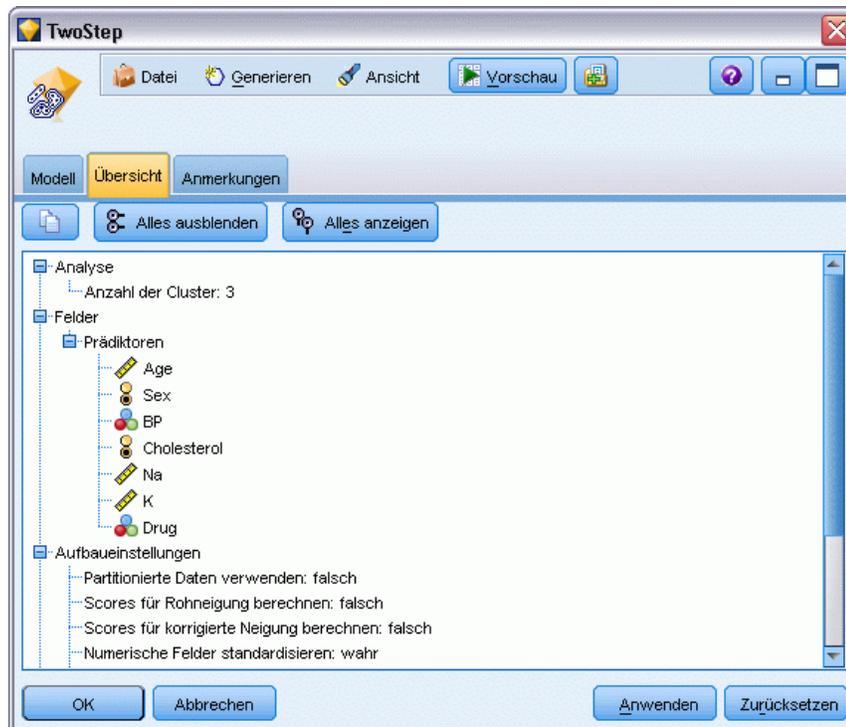
Ein leistungsstarkes Verfahren zur Gewinnung eines Einblicks in das TwoStep-Modell besteht in der Verwendung der Regelinduktion zur Ermittlung der Merkmale, die die vom Modell gefundenen Cluster unterscheiden. [Für weitere Informationen siehe Thema C5.0-Knoten in Kapitel 6 auf S. 177.](#) Außerdem können Sie im Modell-Nugget-Browser auf die Registerkarte “Modell” klicken, um den Cluster-Viewer anzuzeigen, der eine grafische Darstellung der Cluster, Felder und Wichtigkeitsniveaus bietet. [Für weitere Informationen siehe Thema Cluster-Viewer-Registerkarte Modell auf S. 372.](#)

Allgemeine Informationen zur Verwendung des Modellbrowsers finden Sie hier: [Durchsuchen von Modell-Nuggets](#)

## **Übersicht über das TwoStep-Modell**

Auf der Registerkarte “Übersicht” für ein Modell-Nugget vom Typ “TwoStep-Cluster” werden die Anzahl der gefundenen Cluster sowie Informationen zu den Trainingsdaten, dem Schätzvorgang und den verwendeten Aufbaueinstellungen angezeigt.

Abbildung 11-11  
 Beispiel-Modell-Nugget vom Typ "TwoStep-Cluster", Registerkarte "Übersicht"



Für weitere Informationen siehe Thema Durchsuchen von Modell-Nuggets in Kapitel 3 auf S. 52.

## Der Cluster-Viewer

Clustermodelle werden üblicherweise verwendet, um Gruppen (oder Cluster) ähnlicher Datensätze zu finden, die auf den untersuchten Variablen basieren, wobei die Ähnlichkeit zwischen Elementen derselben Gruppe hoch und die Ähnlichkeit zwischen Elementen verschiedener Gruppen niedrig ist. Die Ergebnisse können zur Identifizierung von Zusammenhängen verwendet werden, die ansonsten nicht offensichtlich wären. So kann es zum Beispiel die Clusteranalyse von Kundenpräferenzen, Einkommensniveau und Kaufgewohnheiten ermöglichen, die Kundentypen zu identifizieren, die mit größerer Wahrscheinlichkeit auf eine bestimmte Marketingkampagne ansprechen.

Es gibt zwei Ansätze bei der Interpretierung der Ergebnisse in einer Cluster-Darstellung:

- Untersuchen der Cluster, um die Merkmale zu bestimmen, die in einem Cluster eindeutig sind. *Enthält ein Cluster sämtliche Käufer mit hohem Einkommen? Enthält dieser Cluster mehr Datensätze als die anderen?*
- Untersuchen von Feldern in allen Clustern, um zu bestimmen, wie die Werte in den Clustern verteilt sind. *Ist der Bildungsstand entscheidend für die Zugehörigkeit zu einem Cluster? Spielt ein hoher Kreditrahmen eine Rolle bei der Zugehörigkeit zu einem Cluster oder einem anderen?*

Wenn Sie die Hauptansicht und die zahlreichen verknüpften Ansichten in der Clusteranzeige nutzen, lassen sich diese Fragen beantworten.

Die folgenden Clustermodell-Nuggets können erstellt werden in: IBM® SPSS® Modeler

- Kohonennetz-Modell-Nugget
- K-Means-Modell-Nugget
- TwoStep-Clustermodell-Nugget

Klicken Sie für weitere Informationen über die Clustermodell-Nuggets mit der rechten Maustaste auf den Modellknoten und wählen Sie Durchsuchen aus dem Kontextmenü (oder Bearbeiten für Knoten in einem Strom). Wenn Sie den Modellierungsknoten Auto Cluster verwenden, doppelklicken Sie auf den erforderlichen Cluster-Nugget in dem Auto Cluster-Modell-Nugget.

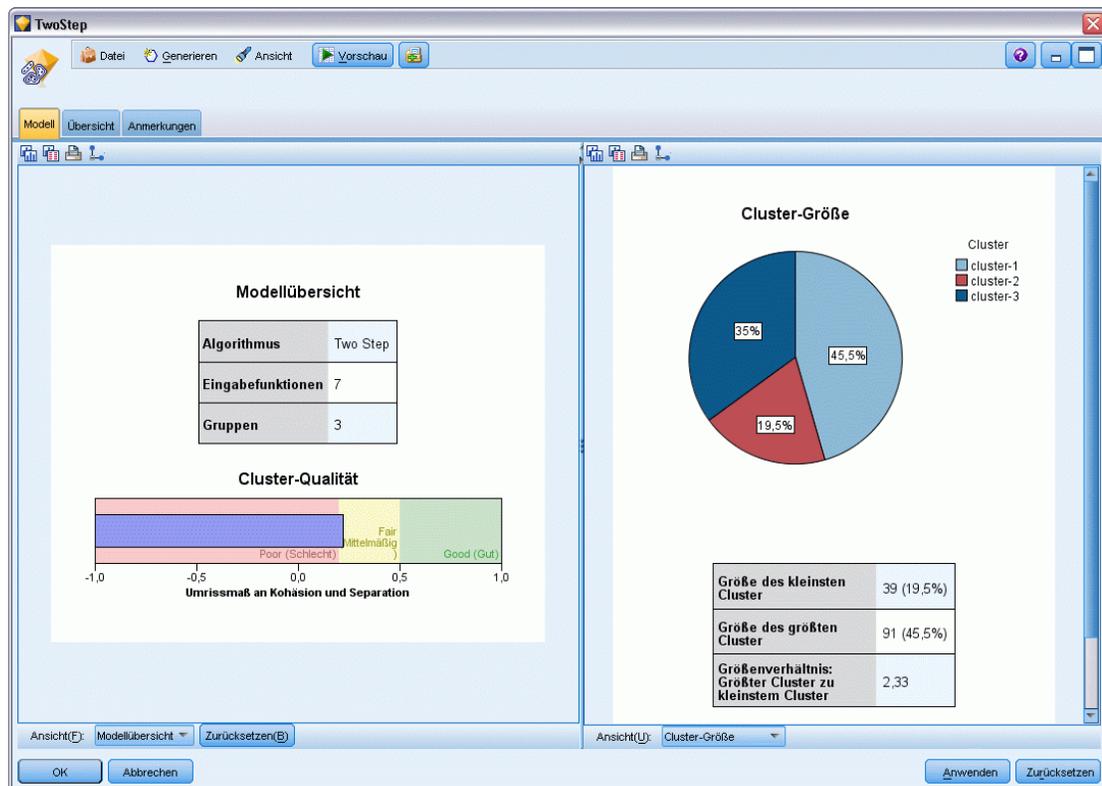
[Für weitere Informationen siehe Thema Knoten "Autom. Cluster" in Kapitel 5 auf S. 114.](#)

### ***Cluster-Viewer -Registerkarte Modell***

Die Registerkarte "Modell" für Clustermodelle zeigt eine grafische Darstellung von Übersichtsstatistiken und die Verteilung von Feldern zwischen Clustern; bekannt als **Cluster-Viewer**.

*Hinweis:* Die Registerkarte Modell ist nicht verfügbar für Modelle, die in Versionen von IBM® SPSS® Modeler vor 13 erstellt wurden.

Abbildung 11-12  
Clusteranzeige mit Standardanzeige



Die Clusteranzeige besteht aus zwei Bereichen, der Hauptansicht im linken Bereich und der verknüpften oder Hilfsansicht im rechten Bereich. Es gibt zwei Hauptansichten:

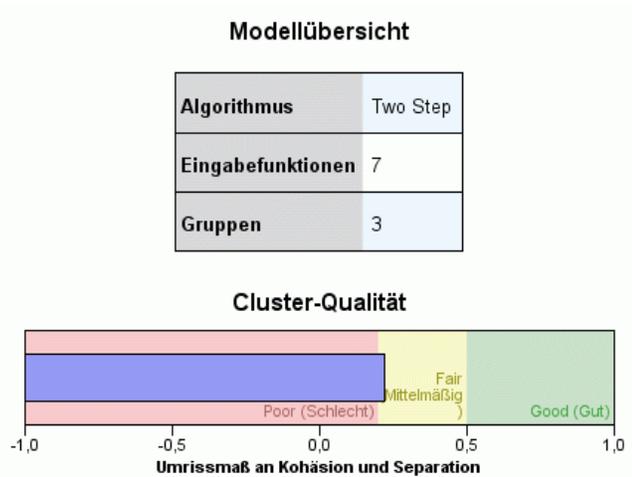
- Modellübersicht (Standard). Für weitere Informationen siehe Thema Ansicht Modellübersicht auf S. 374.
- Cluster. Für weitere Informationen siehe Thema Clusteransicht auf S. 375.

Es gibt vier verknüpfte/Hilfsansichten:

- Bedeutsamkeit des Prädiktors. Für weitere Informationen siehe Thema Ansicht Bedeutsamkeit des Prädiktors im Cluster auf S. 378.
- Clustergrößen (Standard). Für weitere Informationen siehe Thema Clustergrößenansicht auf S. 379.
- Zellenverteilung. Für weitere Informationen siehe Thema Ansicht Zellverteilung auf S. 380.
- Cluster-Vergleich. Für weitere Informationen siehe Thema Ansicht Clustervergleich auf S. 381.

### Ansicht Modellübersicht

Abbildung 11-13  
Ansicht "Modellübersicht" im Hauptpanel



Die Ansicht "Modellübersicht" zeigt eine Momentaufnahme oder eine Übersicht des Clustermodells einschließlich eines schattierten Umrissmaßes der Cluster-Kohäsion und Cluster-Separation, um schlechte, mittelmäßige und gute Ergebnisse anzuzeigen. Anhand dieser Momentaufnahme erkennen Sie schnell, ob die Qualität schlecht ist, so dass Sie dann gegebenenfalls zum Modellierungsknoten zurückkehren und die Clustermodell-Einstellungen ändern können, um ein besseres Ergebnis zu erzielen.

Die Ergebnisse "schlecht", "mittelmäßig" oder "gut" basieren auf der Arbeit von Kaufman und Rousseeuw (1990) zur Interpretation von Clusterstrukturen. In der Ansicht "Modellübersicht" entspricht ein gutes Ergebnis Daten, die von Kaufman und Rousseeuw als annehmbarer oder starker Hinweis auf eine Clusterstruktur eingestuft werden, "mittelmäßig" entspricht ihrer Einstufung als schwacher Hinweis und "schlecht" entspricht ihrer Einstufung als kein signifikanter Hinweis.

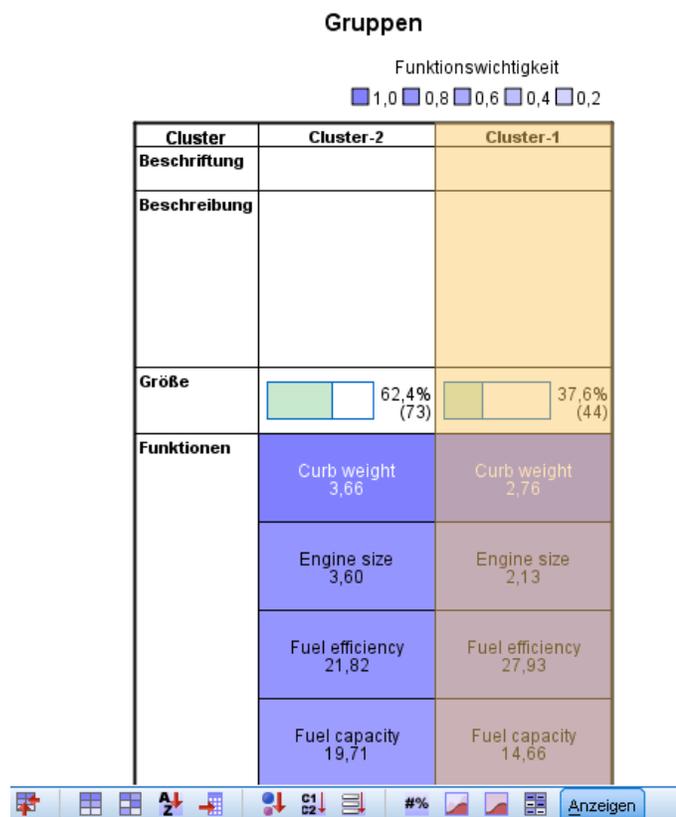
Das Umrissmaß ist ein Durchschnitt aller Datensätze  $(B-A) / \max(A,B)$ , wobei A der Abstand des Datensatzes zu seinem Clusterzentrum und B der Abstand des Datensatzes zu dem am nächsten liegenden, nicht zugehörigen Clusterzentrum ist. Ein Umrisskoeffizient von 1 würde bedeuten, dass alle Fälle direkt in ihren Clusterzentren liegen. Ein Wert von  $-1$  würde bedeuten, dass alle Fälle in den Clusterzentren anderer Cluster liegen. Ein Wert von 0 bedeutet, dass die Fälle im Durchschnitt gleich weit entfernt von ihrem eigenen Clusterzentrum und dem nächsten benachbarten Cluster liegen.

Die Übersicht beinhaltet eine Tabelle, die folgende Daten enthält:

- **Algorithmus.** Der verwendete Clustering-Algorithmus, zum Beispiel "TwoStep".
- **Eingabefunktionen.** Die Anzahl der Felder, auch bekannt als **Eingaben** oder **Einflussgrößen**.
- **Cluster.** Die Anzahl der Cluster in der Lösung.

## Clusteransicht

Abbildung 11-14  
Ansicht "Clusterzentrum" im Hauptpanel



Die Clusteransicht enthält ein Cluster-nach-Funktionen-Raster mit Clusternamen, -größen und -profilen für jeden Cluster.

Die Spalten in der Tabelle enthalten die folgenden Informationen:

- **Cluster.** Die Clusternummern werden von dem Algorithmus erstellt.
- **Bezeichnung.** Bezeichnungen für jeden Cluster (ist standardmäßig leer). Doppelklicken Sie in die Zelle, um eine Bezeichnung einzugeben, die den Clusterinhalt beschreibt; zum Beispiel "Käufer von Luxusautos".
- **Beschreibung.** Beschreibung des Clusterinhalts (ist standardmäßig leer). Doppelklicken Sie in die Zelle, um eine Beschreibung des Clusters einzugeben, zum Beispiel "Alter 55+, Berufstätige, Einkommen über \$100.000".

- **Größe.** Die Größe jedes Clusters als Prozentsatz der gesamten Cluster-Stichprobe. Jede Größenzelle in der Tabelle zeigt einen vertikalen Balken, der den Größenprozentsatz innerhalb des Clusters, einen Größenprozentsatz in numerischem Format und die Cluster-Fallzahl anzeigt.
- **Merkmale.** Die einzelnen Eingaben oder Einflussgrößen, standardmäßig nach Gesamtwichtigkeit sortiert. Wenn Spalten die gleiche Größe aufweisen, werden sie in aufsteigender Sortierfolge ihrer Clusternummern angezeigt.

Die Gesamtwichtigkeit des Merkmals wird von der Farbe der Zellenhintergrundschiattierung angezeigt; das wichtigste Merkmal ist am dunkelsten, das am wenigsten wichtige Merkmal ist ungeschattiert. Ein Hinweis oberhalb der Tabelle erläutert die Wichtigkeit, die jeder Merkmalszelle zugewiesen ist.

Wenn Sie mit der Maus über eine Zelle fahren, wird der volle Name/die Bezeichnung des Merkmals und der Wichtigkeitswert der Zelle angezeigt. Je nach Anzeige- und Merkmalstyp können auch weitere Informationen angezeigt werden. In der Ansicht "Clusterzentrum" zählen die Zellenstatistik und der Zellenwert dazu; zum Beispiel: "Mittelwert: 4.32". Bei kategorischen Merkmalen zeigt die Zelle den Namen der häufigsten (typischen) Kategorie und deren Prozentsatz.

In der Ansicht "Cluster" können Sie verschiedene Anzeigearten für die Clusterinformationen auswählen:

- Cluster und Funktionen transponieren. [Für weitere Informationen siehe Thema Cluster und Merkmale transponieren auf S. 376.](#)
- Merkmale sortieren. [Für weitere Informationen siehe Thema Merkmale sortieren auf S. 377.](#)
- Cluster sortieren. [Für weitere Informationen siehe Thema Cluster sortieren. auf S. 377.](#)
- Zelleninhalte auswählen. [Für weitere Informationen siehe Thema Zelleninhalt auf S. 377.](#)

### **Cluster und Merkmale transponieren**

Standardmäßig werden Cluster als Spalten und Merkmale als Zeilen angezeigt. Um die Anzeige umzudrehen, klicken Sie auf die Schaltfläche Cluster und Merkmale transponieren links von der Schaltfläche Merkmale sortieren nach. Dies kann zum Beispiel wünschenswert sein, wenn zahlreiche Cluster angezeigt werden, um den horizontalen Bildlauf bei der Datenansicht zu verringern.

Abbildung 11-15

Transponierte Cluster im Hauptpanel

Cluster	Beschriftung	Beschreibung	Größe	
cluster-1			45,0% (91)	BP HIGH (41,8%)
cluster-3			35,0% (70)	BP NORMAL (51,4%)
cluster-2			19,0% (39)	BP HIGH (100,0%)

### ***Merkmale sortieren***

Die Schaltflächen Merkmale sortieren nach ermöglichen Ihnen die Auswahl, wie Merkmalzellen angezeigt werden:

- **Gesamtwichtigkeit.** Das ist die standardmäßige Sortierfolge. Die Merkmale werden in absteigender Sortierfolge der Gesamtwichtigkeit sortiert, und die Sortierfolge ist dieselbe bei allen Clustern. Wenn Merkmale gebundene Wichtigkeitswerte aufweisen, sind die gebundenen Merkmale in aufsteigender Sortierfolge der Merkmalnamen aufgelistet.
- **Wichtigkeit innerhalb der Cluster.** Die Merkmale werden hinsichtlich ihrer Wichtigkeit für jeden Cluster sortiert. Wenn Merkmale gebundene Wichtigkeitswerte aufweisen, sind die gebundenen Merkmale in aufsteigender Sortierfolge der Merkmalnamen aufgelistet. Wenn diese Option ausgewählt wird, variiert üblicherweise die Sortierfolge in den Clustern.
- **Name.** Die Merkmale werden nach Namen in alphabetischer Reihenfolge sortiert.
- **Datenfolge.** Die Merkmale werden nach ihrer Reihenfolge im Datensatz sortiert.

### ***Cluster sortieren.***

Standardmäßig werden Cluster ihrer Größe nach absteigend sortiert. Mit den Schaltflächen Cluster sortieren nach können Sie die Cluster nach Namen in alphabetischer Reihenfolge sortieren, oder, wenn Sie eindeutige Bezeichnungen erstellt haben, stattdessen auch in alphanumerischer Bezeichnungsreihenfolge.

Merkmale mit derselben Bezeichnung werden nach Clustername sortiert. Wenn die Cluster nach Bezeichnung sortiert sind und Sie die Bezeichnung eines Clusters bearbeiten, wird die Sortierfolge automatisch aktualisiert.

### ***Zelleninhalt***

Mit den Schaltflächen Zellen können Sie die Anzeige der Zelleninhalte für Merkmale- und Evaluationsfelder ändern.

- **Clusterzentren.** Standardmäßig zeigen Zellen Namen/Bezeichnungen und das Maß der Zentraltendenz für jede Cluster/Merkmal-Kombination an. Für kontinuierliche Felder wird der Mittelwert angezeigt und für kategorische Felder der Modus (die am häufigsten auftretende Kategorie) mit Kategorieprozentsatz.
- **Absolute Verteilungen.** Zeigt die Merkmalnamen/-bezeichnungen und die absoluten Verteilungen der Merkmale in jedem Cluster. Bei kategorischen Merkmalen werden Balkendiagramme angezeigt, mit überlagerter Anzeige der Kategorien, die nach ihren Datenwerten aufsteigend geordnet sind. Bei kontinuierlichen Merkmalen stellt die Anzeige ein gleichmäßiges Dichtediagramm dar, bei dem die gleichen Endpunkte und Intervalle für jeden Cluster verwendet werden.

Die intensiv rote Anzeige stellt die Clusterverteilung dar, wogegen die blassere Anzeige die Gesamtdaten repräsentiert.

- **Relative Verteilungen.** Zeigt die Merkmalnamen/-bezeichnungen und die relativen Verteilungen in den Zellen. Im Allgemeinen sind die Anzeigen vergleichbar mit denen für absolute Verteilungen, nur dass stattdessen die relativen Verteilungen dargestellt sind.

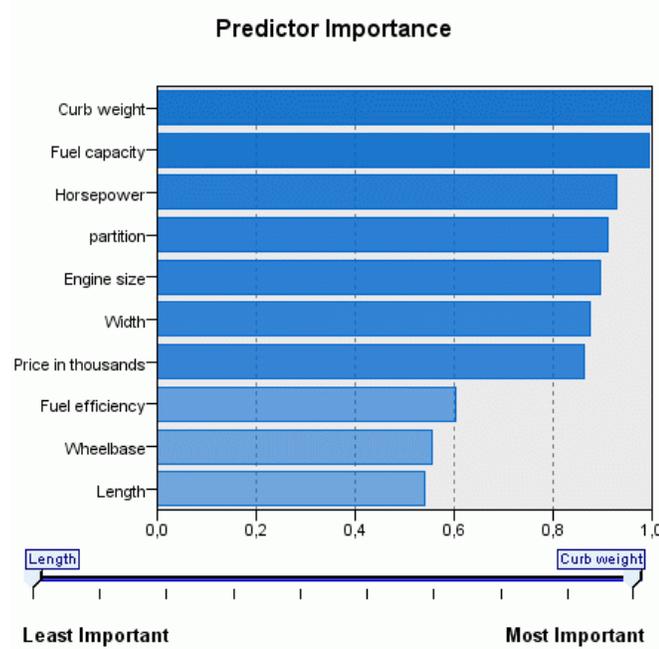
Die intensiv rote Anzeige stellt die Clusterverteilung dar, wogegen die blassere Anzeige die Gesamtdaten repräsentiert.

- **Basisansicht.** Bei sehr vielen Clustern kann es schwierig sein, sämtliche Details ohne Bildlauf zu sehen. Wählen Sie diese Ansicht, um den Bildlauf einzuschränken und die Anzeige auf eine kompaktere Version der Tabelle zu ändern.

### ***Ansicht Bedeutsamkeit des Prädiktors im Cluster***

Abbildung 11-16

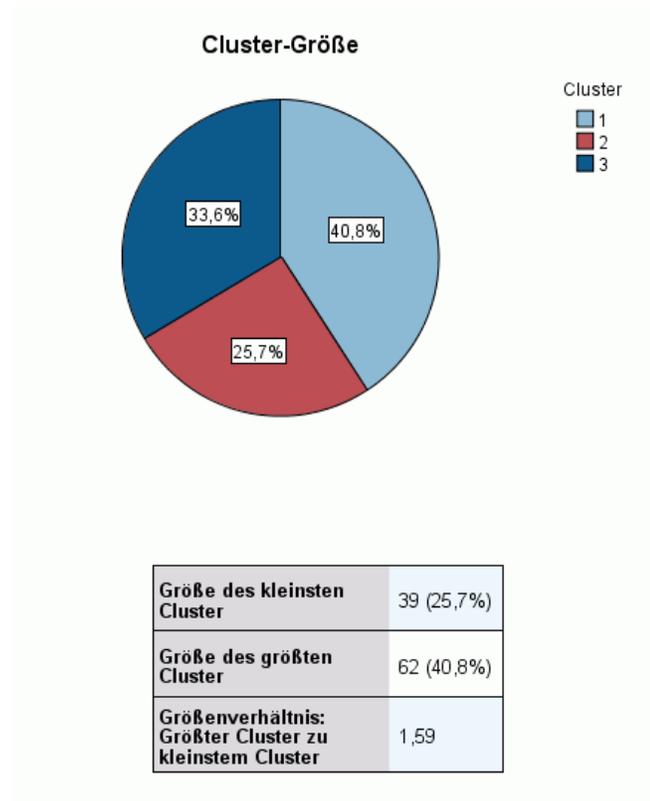
Ansicht "Bedeutsamkeit des Prädiktors im Cluster" im Verknüpfungspanel



Die Ansicht "Bedeutsamkeit des Prädiktors" zeigt die relative Wichtigkeit jedes Felds bei Schätzung des Modells. Für weitere Informationen siehe [Thema Bedeutsamkeit des Prädiktors in Kapitel 3 auf S. 54](#).

### Clustergrößenansicht

Abbildung 11-17  
Ansicht "Clustergrößen" im Verknüpfungspanel



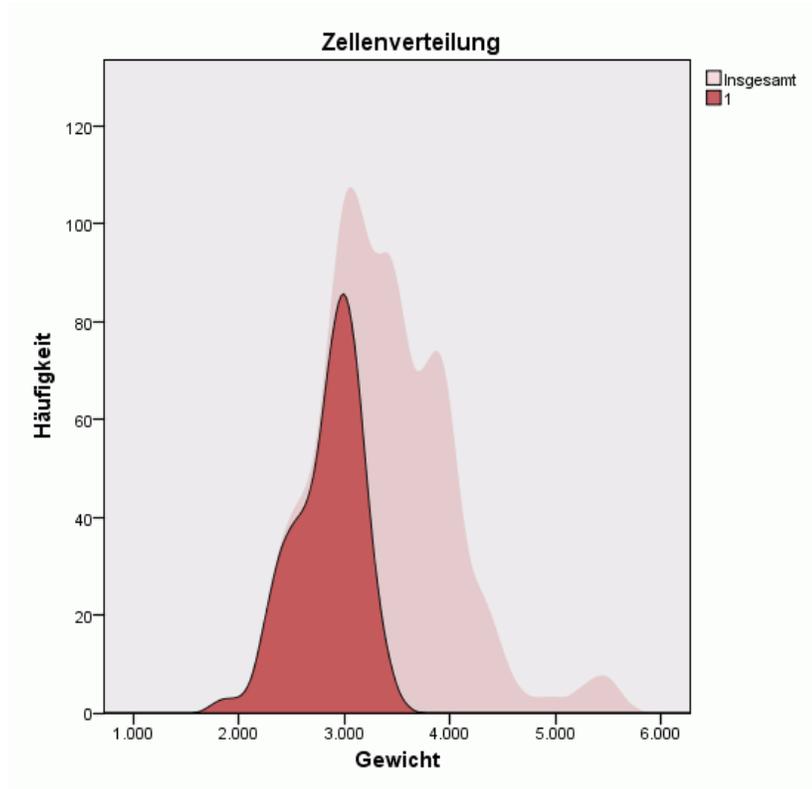
Die Ansicht "Clustergrößen" zeigt ein Tortendiagramm, das sämtliche Cluster enthält. In jedem Stückchen wird die prozentuale Größe des Clusters angezeigt; fahren Sie mit der Maus über ein Stückchen, um den Zahlwert in diesem Stück anzuzeigen.

Unterhalb des Diagramms sind in einer Tabelle die folgenden Informationen aufgelistet:

- Größe des kleinsten Clusters (als Zahlwert und Prozentsatz des Ganzen).
- Größe des größten Clusters (als Zahlwert und Prozentsatz des Ganzen).
- Verhältnis der Größe des größten Clusters zum kleinsten Cluster.

**Ansicht Zellverteilung**

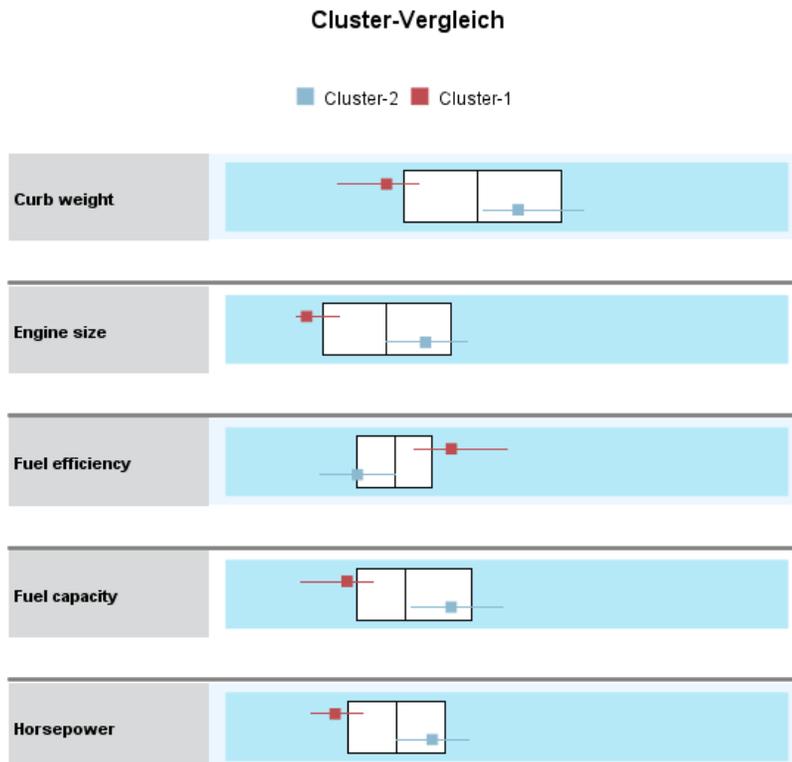
Abbildung 11-18  
Ansicht "Zellverteilung" im Verknüpfungspanel



Die Ansicht "Zellverteilung" zeigt ein erweitertes, detaillierteres Diagramm der Datenverteilung für jede Merkmalszelle, die Sie in der Tabelle im Cluster-Hauptpanel auswählen.

## Ansicht Clustervergleich

Abbildung 11-19  
Ansicht "Clustervergleich" im Verknüpfungspanel



Die Ansicht "Clustervergleich" ist eine tabellarische Grafik, bei der die Merkmale in den Zeilen und die ausgewählten Cluster in den Spalten dargestellt werden. Mit dieser Ansicht lassen sich die Faktoren besser verstehen, die die Cluster ausmachen; außerdem hilft sie dabei, die Unterschiede zwischen den Clustern zu erkennen – nicht nur im Vergleich zum Gesamtdatensatz, sondern auch untereinander.

Zum Auswählen der Cluster für die Ansicht klicken Sie oben auf die Clusterspalte im Cluster-Hauptpanel. Wenn Sie die Strg-Taste oder die Umschalttaste beim Klicken gedrückt halten, können Sie mehrere Cluster zum Vergleich auswählen oder wieder aus der Auswahl entfernen.

*Hinweis:* Sie können bis zu fünf Cluster zur Anzeige auswählen.

Die Cluster werden in der Reihenfolge ihrer Auswahl angezeigt, während die Reihenfolge der Felder mit der Option Merkmale sortieren nach festgelegt wird. Wenn Sie Wichtigkeit innerhalb der Cluster auswählen, werden die Felder immer nach ihrer Gesamtwichtigkeit sortiert.

Die Hintergrunddiagramme zeigen die Gesamtverteilungen der Merkmale:

- Kategorische Merkmale sind als Punktdiagramme dargestellt, wobei die Größe des Punktes die häufigste/typische Kategorie für jeden Cluster (nach Merkmal) anzeigt.
- Kontinuierliche Merkmale sind als Boxplots angezeigt, der die Gesamtmediane und die Interquartilbereiche anzeigt.

Vor diesen Hintergrundansichten sind Boxplots für ausgewählte Cluster dargestellt:

- Für stetige Merkmale geben quadratische Punktmarkierungen und horizontale Linien den Median- und Interquartilbereich für jedes einzelne Cluster an.
- Jeder Cluster ist mit einer anderen Farbe gekennzeichnet, die oben an der Ansicht angezeigt wird.

## Navigieren in der Clusteranzeige

Bei der Clusteranzeige handelt es sich um eine interaktive Anzeige. Sie verfügen über folgende Möglichkeiten:

- Auswählen eines Felds oder eines Clusters für weitere Details
- Vergleichen von Clustern, um die Objekte von Interesse auszuwählen
- Verändern der Anzeige
- Transponieren von Achsen
- Erzeugen von Knoten zum Ableiten, Filtern und Auswählen unter Verwendung des Menüs "Erzeugen".

### Verwendung der Symbolleisten

Sie können die Informationen, die in den Panels links und rechts angezeigt werden, mithilfe der Symbolleisteoptionen steuern. Mit der Symbolleistesteuerung können Sie die Ausrichtung der Anzeige ändern (oben-unten, links-rechts oder rechts-links). Außerdem können Sie die Clusteranzeige auf die Standardeinstellungen zurücksetzen und ein Dialogfeld öffnen, um den Inhalt der Clusteransicht im Hauptpanel zu spezifizieren.

Abbildung 11-20

Symbolleiste zum Steuern der in der Clusteranzeige angezeigten Daten



Die Optionen Merkmale sortieren nach, Cluster sortieren nach, Zellen und Anzeige sind nur verfügbar, wenn Sie die Ansicht Cluster im Hauptpanel auswählen. [Für weitere Informationen siehe Thema Clusteransicht auf S. 375.](#)

	Siehe <a href="#">Cluster und Merkmale transponieren</a> auf S. 376
	Siehe <a href="#">Merkmale sortieren nach</a> auf S. 377

	Siehe <a href="#">Cluster sortieren nach</a> auf S. 377
	Siehe <a href="#">Zellen</a> auf S. 377

### **Erzeugen von Knoten aus Clustermodellen**

Mit dem Menü “Erzeugen” können Sie auf der Basis des Clustermodells neue Knoten erstellen. Diese Option ist auf der Registerkarte “Modell” des erzeugten Modells verfügbar und ermöglicht es, Knoten auf der Basis der aktuellen Anzeige oder einer Auswahl zu erzeugen (das heißt, alle sichtbaren oder alle ausgewählten Cluster). Sie können zum Beispiel ein einzelnes Merkmal auswählen und anschließend einen Filterknoten erstellen, um alle anderen (nicht sichtbaren) Merkmale zu verwerfen. Die erzeugten Knoten werden unzusammenhängend auf der Zeichenfläche platziert. Sie können außerdem eine Kopie des Modellnuggets erstellen und zur Modellpalette hinzufügen. Denken Sie daran, vor der Ausführung die Knoten zu verknüpfen und alle erwünschten Änderungen vorzunehmen.

- **Erzeugen eines Modellierungsknotens.** Erzeugt einen Modellierungsknoten auf der Datenstrom-Zeichenfläche. Das könnte zum Beispiel nützlich sein, wenn Sie bei einem Datenstrom diese Modelleinstellungen verwenden möchten, aber nicht mehr über den Modellierungsknoten verfügen, um sie zu erzeugen.
- **Modell zur Palette hinzufügen.** Erstellt ein Nugget auf der Modellpalette. Das ist nützlich, wenn Sie von einem Kollegen einen Datenstrom, der das Modell enthält, jedoch nicht das Modell selbst erhalten.
- **Filterknoten.** Erstellt einen neuen Filterknoten, um Felder zu filtern, die von dem Clustermodell nicht verwendet werden und/oder in der aktuellen Ansicht der Clusteranzeige nicht sichtbar sind. Wenn es von diesem Clusterknoten einen vorgelagerten Typenknoten gibt, werden Felder mit der Rolle *Ziel* von dem erzeugten Filterknoten verworfen.
- **Filterknoten (aus der Auswahl).** Erzeugt einen neuen Filterknoten, um die Felder auf der Basis der Auswahl der Clusteranzeige zu filtern. Mehrere Felder können Sie durch Gedrückhalten der Strg-Taste beim Klicken auswählen. Die in der Clusteranzeige ausgewählten Felder werden nachgelagert verworfen, doch Sie können diese Einstellung ändern, indem Sie den Filterknoten vor dem Ausführen bearbeiten.
- **Auswahlknoten.** Erstellt einen neuen Auswahlknoten, um Datensätze basierend auf ihrer Zugehörigkeit zu den Clustern, die in der aktuellen Ansicht der Clusteranzeige sichtbar sind, auszuwählen. Eine Auswahlbedingung wird automatisch generiert.
- **Auswahlknoten (aus der Auswahl).** Erstellt einen neuen Auswahlknoten, um Datensätze basierend auf ihrer Zugehörigkeit zu Clustern, die in der Clusteranzeige ausgewählt wurden, auszuwählen. Wählen Sie mehrere Cluster aus, indem Sie beim Klicken die Strg-Taste gedrückt halten.

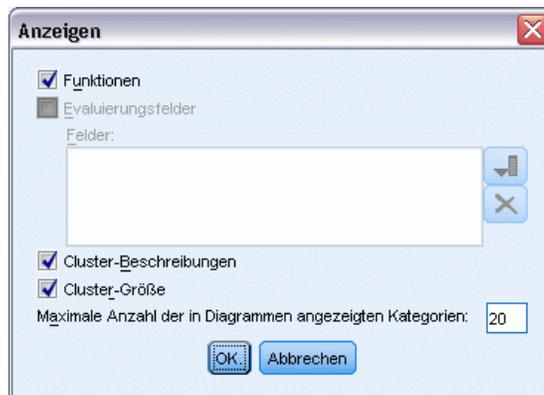
- **Ableitungsknoten.** Erzeugt einen neuen Ableitungsknoten, der ein Flag-Feld erstellt, das Datensätzen einen Wert *Wahr* oder *Falsch* zuweist, basierend auf der Zugehörigkeit zu allen in der Clusteranzeige sichtbaren Clustern. Eine Ableitungsbedingung wird automatisch generiert.
- **Ableitungsknoten (aus der Auswahl).** Erstellt einen neuen Ableitungsknoten, der ein Flag-Feld basierend auf der Zugehörigkeit zu Clustern, die in der Clusteranzeige ausgewählt wurden, ableitet. Wählen Sie mehrere Cluster aus, indem Sie beim Klicken die Strg-Taste gedrückt halten.

Im Menü “Erzeugen” können Sie nicht nur Knoten, sondern auch Diagramme erstellen. [Für weitere Informationen siehe Thema Erzeugen von Diagrammen aus Clustermodellen auf S. 385.](#)

### Anzeige “Clusteransicht steuern”

Um zu steuern, was in der Clusteransicht im Hauptpanel angezeigt wird, klicken Sie auf die Schaltfläche Anzeige. Der Anzeige-Dialog wird geöffnet.

Abbildung 11-21  
Clusteranzeige- Anzeigeoptionen



**Merkmale.** Standardmäßig ausgewählt. Deaktivieren Sie das Kästchen, um alle Eingabemerkmale auszublenden.

**Evaluierungsfelder.** Wählen Sie die anzuzeigenden Evaluierungsfelder aus (Felder, die nicht für die Erstellung des Clustermodells verwendet, sondern an die Modellanzeige zur Evaluierung der Cluster gesendet werden); standardmäßig werden keine angezeigt. *Hinweis:* Dieses Kontrollkästchen ist nicht verfügbar, wenn keine Evaluierungsfelder verfügbar sind.

**Clusterbeschreibungen.** Standardmäßig ausgewählt. Deaktivieren Sie das Kontrollkästchen, um alle Clusterbeschreibungszellen auszublenden.

**Clustergröße.** Standardmäßig ausgewählt. Deaktivieren Sie das Kontrollkästchen, um alle Clustergrößenzellen auszublenden.

**Maximale Anzahl an Kategorien.** Geben Sie die maximale Anzahl an Kategorien an, die in den Diagrammen der kategorischen Merkmale angezeigt werden sollen; der Standard ist 20.

## **Erzeugen von Diagrammen aus Clustermodellen**

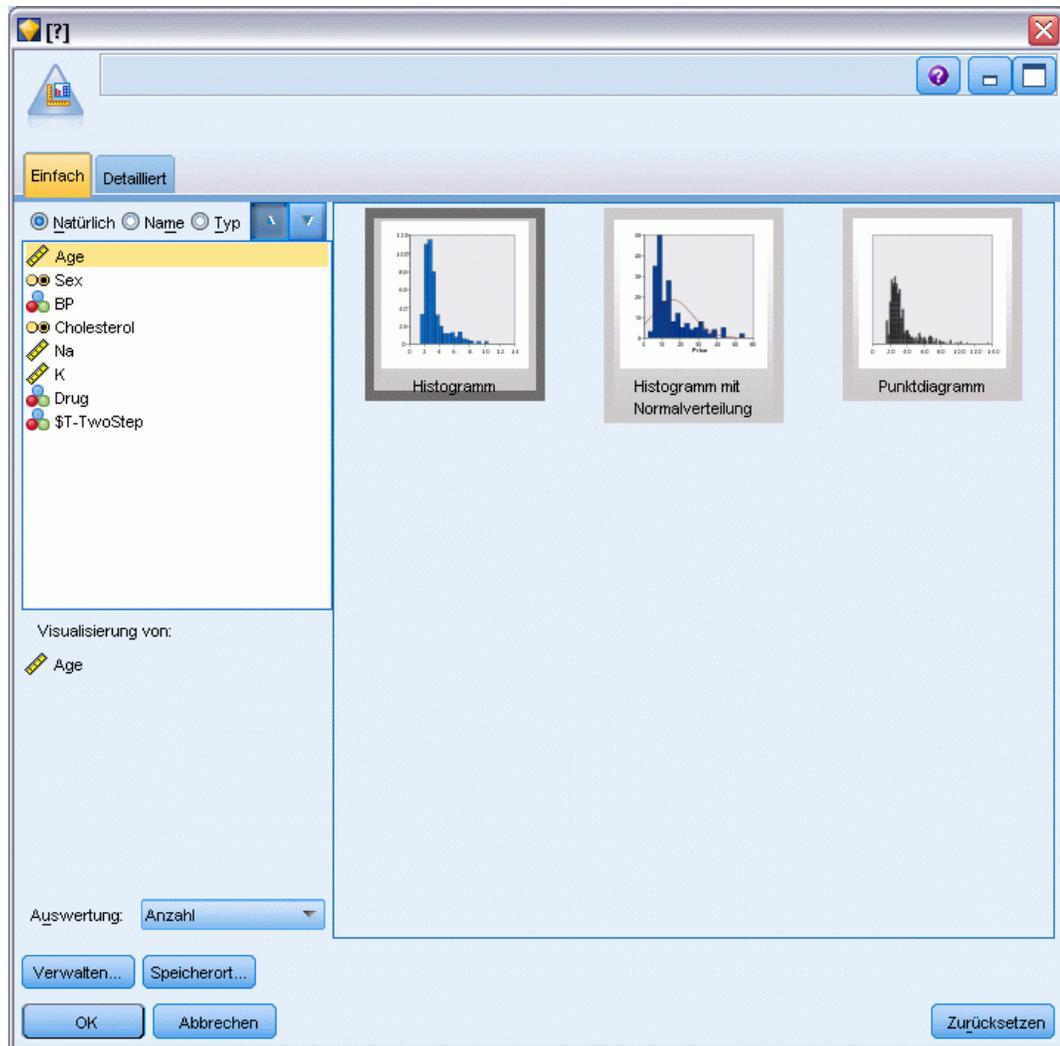
Clustermodelle beinhalten viele Informationen; sie sind jedoch häufig in einem Format, auf das Unternehmensanwender nicht so einfach zugreifen können. Zur Bereitstellung der Daten auf eine Art, die problemlos in Geschäftsberichte, Präsentationen u.s.w. integriert werden kann, können aus ausgewählten Daten Diagramme erstellt werden. In der Clusteranzeige können Sie zum Beispiel ein Diagramm für einen ausgewählten Cluster erstellen, wobei das Diagramm nur für die Fälle in diesem Cluster erzeugt wird.

*Hinweis:* Sie können in der Clusteranzeige nur dann ein Diagramm erstellen, wenn das Modellnugget mit anderen Knoten in einem Datenstrom verbunden ist.

### **Ein Diagramm erstellen**

- ▶ Öffnen Sie das Modellnugget, das die Clusteranzeige enthält.
- ▶ Wählen Sie auf der Registerkarte “Modell” in der Dropdown-Liste *Ansicht* die Option Cluster aus.
- ▶ Wählen Sie in der Hauptansicht den oder die Cluster, für die Sie ein Diagramm erstellen möchten, aus.
- ▶ Wählen Sie im Menü “Erzeugen” den Menüpunkt Diagramm (aus der Auswahl) aus; die Registerkarte “Basis” auf der Grafiktabelle wird angezeigt.

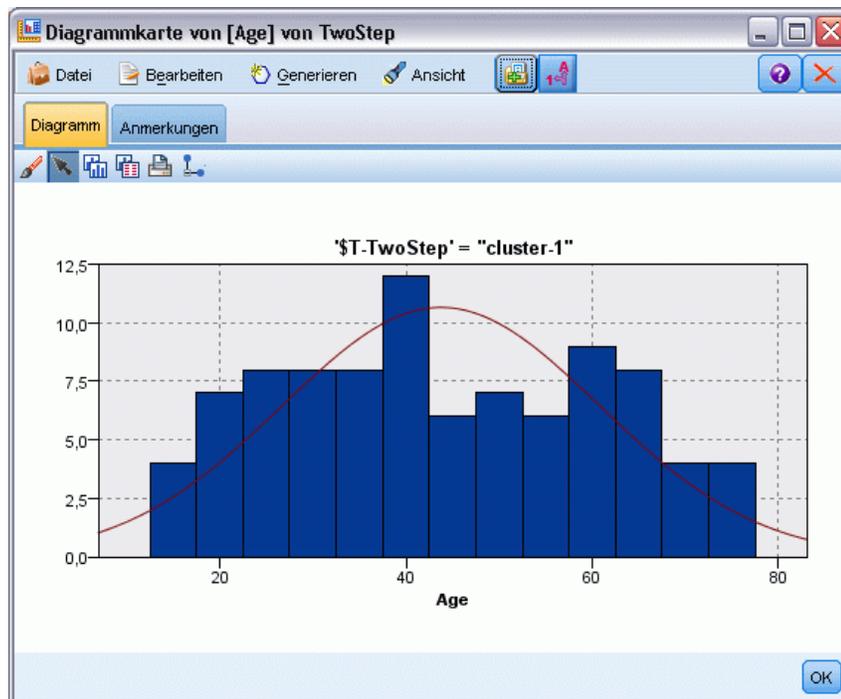
Abbildung 11-22  
Dialogfeld Grafiktafelknoten, Registerkarte "Basis"



*Hinweis:* Wenn Sie die Grafiktafel auf diese Art anzeigen, sind nur die Registerkarten "Basis" und "Details" verfügbar. Für weitere Informationen siehe Thema Diagrammtafelknoten in Kapitel 5 in *IBM SPSS Modeler 14.2- Quellen-, Prozess- und Ausgabeknoten*.

- ▶ Mithilfe der Einstellungen auf den Registerkarten "Basis" oder "Details" können Sie die Details angeben, die auf dem Diagramm angezeigt werden sollen.
- ▶ Klicken Sie auf "OK", um das Diagramm zu erstellen.

Abbildung 11-23  
Histogramm erzeugen mit der Grafiktafelregisterkarte "Basis"



Die Diagramm-Überschrift identifiziert den Modelltyp und den oder die Cluster, die eingeschlossen wurden.

# Assoziationsregeln

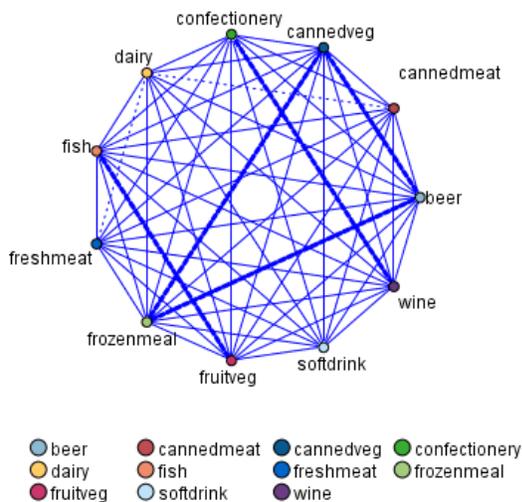
**Assoziationsregeln** ordnen eine bestimmte Schlussfolgerung (beispielsweise den Kauf eines bestimmten Produkts) einer Menge von Bedingungen (dem Kauf mehrerer anderer Produkte) zu. Beispiel: Die Regel

Bier  $\Leftarrow$  -Dosengemüse & TK-Fertiggericht (173, 17,0 %, 0,84)

besagt, dass *Bier* häufig vorkommt, wenn *Dosengemüse* und *TK-Fertiggericht* zusammen vorkommen. Die Regel ist zu 84 % zuverlässig und auf 17 % der Daten, also 173 Datensätze, anwendbar. Algorithmen für Assoziationsregeln finden automatisch die Zuweisungen, die Sie manuell finden könnten, wenn Sie Visualisierungstechniken, wie den Webknoten, anwenden.

Abbildung 12-1

Webknoten, der Verbindungen zwischen den Elementen des Warenkorbs anzeigt



Der Vorteil von Algorithmen für Assoziationsregeln im Vergleich zu Algorithmen für Standardentscheidungsbäume (C5.0 und C&R-Bäume) besteht darin, dass zwischen *beliebigen* Attributen Verbindungen bestehen können. Ein Entscheidungsbaumalgorithmus erstellt Regeln mit nur einer Schlussfolgerung, während Assoziationsalgorithmen viele Regeln zu finden versuchen, von denen jede zu einer anderen Schlussfolgerung kommen kann.

Der Nachteil von Assoziationsregeln besteht darin, dass sie versuchen, Muster innerhalb eines potenziell sehr großen Suchbereichs zu finden, also mehr Zeit für die Ausführung in Anspruch nehmen können als ein Entscheidungsbaumalgorithmus. Die Algorithmen verwenden eine Methode vom Typ **Generieren und Testen** zum Auffinden von Regeln, bei der einfache Regeln erstellt und mit dem Daten-Set verglichen werden. Die "guten" Regeln werden gespeichert und alle Regeln – nachdem sie verschiedenen Beschränkungen unterworfen wurden – werden dann weiter verfeinert. **Spezialisierung** ist der Prozess, bei dem einer Regel Bedingungen

hinzugefügt werden. Diese neuen Regeln werden dann mit den Daten verglichen und validiert. Die “besten” oder interessantesten Regeln werden dann gespeichert. Der Benutzer legt normalerweise einen Grenzwert für die mögliche Anzahl Antezedenzen in einer Regel fest. Es werden außerdem verschiedene Techniken basierend auf der Informationstheorie oder effiziente Indizierungsschemata verwendet, um den potenziell großen Suchbereich zu reduzieren.

Am Ende der Verarbeitung wird eine Tabelle mit den besten Regeln ausgegeben. Im Gegensatz zu einem Entscheidungsbaum kann dieser Satz mit Assoziationsregeln nicht direkt dazu verwendet werden, Vorhersagen auf eine Weise zu machen, wie dies mit einem Standardmodell (z. B. einem Entscheidungsbaum oder neuronalen Netz) möglich ist. Dies ist auf die vielen möglichen Schlussfolgerungen für die Regeln zurückzuführen. Es ist eine weitere Stufe der Transformation erforderlich, um die Assoziationsregeln in eine Klassifizierungsregelmenge umzuwandeln. Deshalb sind die von Assoziationsalgorithmen erstellten Assoziationsregeln bekannt als **nicht verfeinerte Modelle**. Obwohl der Benutzer diese nicht verfeinerten Modelle durchsuchen kann, können Sie nicht ausdrücklich als Klassifizierungsmodelle verwendet werden, es sei denn, der Benutzer weist das System an, aus dem nicht verfeinerten Modell ein Klassifizierungsmodell zu generieren. Dies geschieht mit dem Browser über die Menüoption “Generieren”.

Es werden zwei Algorithmen für Assoziationsregeln unterstützt:



Der A-Priori-Knoten extrahiert eine Regelmenge aus den Daten und daraus die Regeln mit dem höchsten Informationsgehalt. A Priori bietet fünf verschiedene Methoden zur Auswahl von Regeln und verwendet ein ausgereiftes Indizierungsschema zur effizienten Verarbeitung großer Daten-Sets. Bei großen Problemen ist A Priori in der Regel schneller zu trainieren, es gibt keine willkürliche Begrenzung für die Anzahl der Regeln, die beibehalten werden können, und es können Regeln mit bis zu 32 Vorbedingungen verarbeitet werden. Bei A Priori müssen alle Ein- und Ausgabefelder kategorial sein; dafür bietet es jedoch eine bessere Leistung, da es für diesen Datentyp optimiert ist. [Für weitere Informationen siehe Thema A Priori-Knoten auf S. 390.](#)



Der Sequenzknoten erkennt Assoziationsregeln in sequenziellen oder zeitorientierten Daten. Eine Sequenz ist eine Liste mit Element-Sets, die in einer vorhersagbaren Reihenfolge auftreten. Beispiel: Ein Kunde, der einen Rasierer und After-Shave-Lotion kauft, kauft möglicherweise beim nächsten Einkauf Rasiercreme. Der Sequenzknoten basiert auf dem CARMA-Assoziationsregelalgorithmus, der eine effiziente bidirektionale Methode zum Suchen von Sequenzen verwendet. [Für weitere Informationen siehe Thema Sequenzknoten auf S. 416.](#)

## ***Tabellendaten im Vergleich zu Transaktionsdaten***

Die von Assoziationsregelmodellen verwendeten Daten können im Transaktionsformat oder in Tabellenform (siehe unten) vorliegen. Hierbei handelt es sich um allgemeine Beschreibungen; die konkreten Anforderungen können abweichen, wie in der Dokumentation für die einzelnen Modelltypen erörtert. Beachten Sie: Beim Scoring von Modellen müssen die zu scorenden Daten dasselbe Format aufweisen wie die bei der Modellerstellung verwendeten Daten. Mit Tabellendaten erstellte Modelle können ausschließlich zum Scoring von Tabellendaten verwendet werden; mit Transaktionsdaten erstellte Modelle dienen ausschließlich zum Scoring von Transaktionsdaten.

### **Transaktionsformat**

Transaktionsdaten weisen einen eigenen Datensatz für jede Transaktion bzw. jedes Element auf. Wenn ein Kunde beispielsweise mehrere Einkäufe tätigt, handelt es sich bei jedem um einen eigenen Datensatz, wobei die zugehörigen Elemente durch eine Kunden-ID verknüpft sind. Dies wird auch manchmal als **Kassenrollen**-Format bezeichnet.

<b>Kunde</b>	<b>Kauf</b>
1	Marmelade
2	Milch
3	Marmelade
3	Brot
4	Marmelade
4	Brot
4	Milch

A Priori-, CARMA- und Sequenzknoten können jeweils Transaktionsdaten verwenden.

### **Tabellendaten**

In Tabellendaten (auch als **Warenkorb**- oder **Wahrheitstabellen**-Daten bezeichnet) werden die Elemente durch gesonderte Flags dargestellt, wobei jedes Flag-Feld für das Vorliegen bzw. die Abwesenheit eines bestimmten Elements steht. Jeder Datensatz steht für ein komplettes Set zugehöriger Elemente. Prinzipiell können Flag-Felder kategorial oder numerisch sein; bei manchen Modellen können jedoch genauere Anforderungen gelten.

<b>Kunde</b>	<b>Marmelade</b>	<b>Brot</b>	<b>Milch</b>
1	B	U	U
2	U	U	B
3	B	B	U
4	B	B	B

A Priori-, CARMA-, und Sequenzknoten können jeweils Tabellendaten verwenden.

## **A Priori-Knoten**

Der A Priori-Knoten erkennt Assoziationsregeln in den Daten. A Priori bietet fünf verschiedene Methoden zur Auswahl von Regeln und verwendet ein ausgereiftes Indizierungsschema zur effizienten Verarbeitung großer Daten-Sets.

**Anforderungen.** Um eine A Priori-Regelmenge zu erstellen, benötigen Sie mindestens ein *Eingabe*- und ein *Ziel*-Feld. Ein- und Ausgabefelder (mit der Rolle *Eingabe*, *Ziel* oder *Beides*) müssen symbolisch sein. Felder mit der Rolle *Keine* werden ignoriert. Feldtypen müssen vollständig instanziiert werden, bevor der Knoten ausgeführt wird. Die Daten können als Tabellen- oder Transaktionsdaten vorliegen. [Für weitere Informationen siehe Thema Tabellendaten im Vergleich zu Transaktionsdaten auf S. 389.](#)

**Stärken.** Bei großen Problemen ist A Priori in der Regel schneller zu trainieren. Außerdem gibt es keine willkürliche Begrenzung für die Anzahl der Regeln, die beibehalten werden können, und es können Regeln mit max. 32 Vorbedingungen verarbeitet werden. A Priori bietet fünf verschiedene Trainingsmethoden und somit mehr Flexibilität bei der Anpassung der Data-Mining-Methode an das aktuelle Problem.

## Modelloptionen für den A Priori-Knoten

Abbildung 12-2  
Modelloptionen für den A Priori-Knoten



**Modellname.** Sie können den Modellnamen automatisch basierend auf den Ziel- oder ID-Feldnamen (oder dem Modelltyp in Fällen, in denen kein solches Feld angegeben ist) generieren oder einen benutzerdefinierten Namen eingeben.

**Minimale Antezedens-Unterstützung.** Sie können ein Stützkriterium angeben, um die Regeln in der Regelmenge beizubehalten. **Unterstützung** bezieht sich auf den Prozentwert der Datensätze in den Trainingsdaten, für die die Antezedenzien (der “Wenn”-Teil der Regel) wahr sind. (Beachten Sie, dass diese Definition von Unterstützung von der für CARMA- und Sequenzknoten verwendeten abweicht. [Für weitere Informationen siehe Thema Modelloptionen für den Sequenzknoten auf S. 419.](#)) Wenn Sie Regeln erhalten, die auf sehr kleine Teilmengen der Daten angewendet werden, sollten Sie diese Einstellung erhöhen.

*Hinweis:* Die Definition von Unterstützung für A Priori basiert auf der Anzahl der Datensätzen mit den Antezedenzien. Dies steht im Gegensatz zu den CARMA- und Sequenzalgorithmen, für die die Definition von Unterstützung auf der Anzahl der Datensätze mit allen Elementen in einer Regel basiert (d. h. sowohl Antezedenzien als auch Sukzedenzien). Die Ergebnisse der Assoziationsmodelle zeigen sowohl die Maße für die Antezedens- als auch die Regelunterstützung.

**Minimale Regelkonfidenz** Sie können auch ein Konfidenzkriterium angeben. **Konfidenz** basiert auf den Datensätzen, für die die Antezedenzen der Regel wahr sind, und stellt den Prozentwert der Datensätze dar, für die auch die Sukzedenzen wahr sind. Mit anderen Worten, es ist der prozentuale Anteil der Vorhersagen, die, basierend auf der Regel, richtig sind. Regeln mit einer niedrigeren Konfidenz als das angegebene Kriterium werden verworfen. Wenn Sie zu viele Regeln oder uninteressante Regeln erhalten, sollten Sie diese Einstellung erhöhen. Wenn Sie zu wenige (oder überhaupt keine) Regeln erhalten, sollten Sie diese Einstellung reduzieren.

**Maximale Anzahl an Antezedenzen.** Sie können die maximale Anzahl an Vorbedingungen für jede beliebige Regel festlegen. Auf diese Weise können Sie die Komplexität der Regeln begrenzen. Wenn die Regeln zu komplex oder zu spezifisch sind, sollten Sie diese Einstellung reduzieren. Diese Einstellung hat auch einen großen Einfluss auf die Trainingszeit. Wenn das Training der Regelmenge zu viel Zeit in Anspruch nimmt, sollten Sie diese Einstellung reduzieren.

**Nur wahre Werte für Flags.** Wenn diese Option für Daten in tabellarischer Form (Wahrheitstabelle) ausgewählt ist, enthalten die resultierenden Regeln lediglich wahre Werte. Auf diese Weise können Sie die Regeln einfacher verstehen. Die Option gilt nicht für Daten im Transaktionsformat. [Für weitere Informationen siehe Thema Tabellendaten im Vergleich zu Transaktionsdaten auf S. 389.](#)

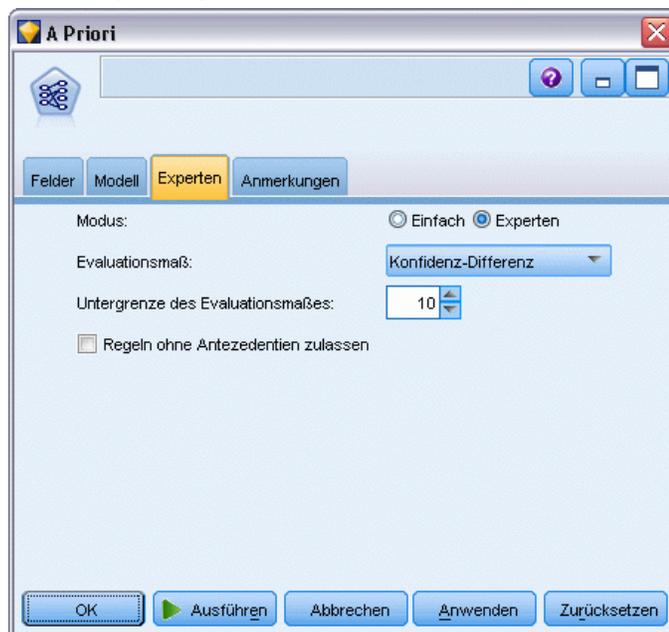
**Optimieren.** Wählen Sie Optionen, die die Leistung während der Modellerstellung basierend auf Ihren persönlichen Anforderungen erhöhen.

- Wählen Sie Geschwindigkeit, um den Algorithmus anzuweisen, zur Verbesserung der Leistung keinen Datenträgerüberlauf zu verwenden.
- Wählen Sie Speicher, um den Algorithmus anzuweisen, ggf. einen Datenträgerüberlauf zu verwenden und dafür Geschwindigkeitseinbußen hinzunehmen. Diese Option ist standardmäßig aktiviert. *Hinweis:* Bei der Ausführung im verteilten Modus kann diese Einstellung durch die in der Datei *options.cfg* angegebenen Administratoroptionen überschrieben werden. Weitere Informationen finden Sie im *IBM® SPSS® Modeler Server-Administratorhandbuch*.

### **Expertenoptionen für den A Priori-Knoten**

Für Personen mit umfassenden Kenntnissen von A Priori ermöglichen die folgenden Expertenoptionen die Feinabstimmung des Induktionsvorgangs. Für den Zugriff auf die Expertenoptionen stellen Sie den Modus auf der Registerkarte "Experten" auf Experten ein.

Abbildung 12-3  
A Priori-Expertenoptionen



**Evaluierungsmaß.** A Priori unterstützt fünf Methoden zur Auswertung möglicher Regeln.

- **Regelkonfidenz.** Die Standardmethode verwendet die Konfidenz (oder Genauigkeit) der Regel zur Regelevaluierung. Für dieses Maß ist die Option Untergrenze der Auswertungsmaßnahme deaktiviert, da sie aufgrund der Option Minimale Regelkonfidenz auf der Registerkarte “Modell” redundant ist. [Für weitere Informationen siehe Thema Modelloptionen für den A Priori-Knoten auf S. 391.](#)
- **Konfidenz-Differenz.** (Auch als **absolute Konfidenz-Differenz zum Vorgänger** bezeichnet.) Dieses Evaluierungsmaß stellt die absolute Differenz zwischen der Regelkonfidenz und der vorherigen Konfidenz dar. Diese Option vermeidet dort einen Bias, wo die Ergebnisse nicht gleichverteilt sind. Dadurch kann verhindert werden, dass “offensichtliche” Regeln beibehalten werden. Zum Beispiel kann es der Fall sein, dass 80 % der Kunden Ihr beliebtestes Produkt kaufen. Eine Regel, die den Kauf dieses beliebten Produkts mit einer 85%igen Konfidenz vorhersagt, hilft Ihnen nicht viel weiter, auch wenn 85%ige Konfidenz in absoluten Werten ziemlich gut erscheint. Setzen Sie die Untergrenze des Evaluierungsmaßes auf die Mindest-Konfidenz-Differenz, für die Regeln beibehalten werden sollen.
- **Konfidenz-Quotient.** (Auch als **Differenz des Konfidenz-Quotienten zur 1** bezeichnet.) Bei diesem Evaluierungsmaß handelt es sich um das Verhältnis der Regelkonfidenz zur vorherigen Konfidenz (oder, wenn das Verhältnis größer als 1 ist, der reziproke Wert) subtrahiert von 1. Wie bei der Konfidenz-Differenz berücksichtigt diese Methode andere Verteilungen als die Gleichverteilung. Sie ist besonders geeignet, Regeln zu finden, die seltene Ereignisse vorhersagen. Angenommen, es gibt eine seltene Krankheit, die nur bei 1 % der Patienten vorkommt. Eine Regel, mit der diese Krankheit in 10 % der Fälle vorhergesehen werden kann, ist im Gegensatz zur groben Schätzung eine erhebliche Verbesserung, obwohl eine Genauigkeit von 10 % in absoluten Werten nicht sehr beeindruckend erscheint. Setzen Sie die

Untergrenze des Evaluierungsmaßes auf die Differenz, für die Regeln eingehalten werden sollen.

- **Informationsdifferenz.** (Auch als **Informationsdifferenz zum Vorgänger** bezeichnet.) Dieses Maß basiert auf dem Maß **Informationsgewinn**. Wenn die Wahrscheinlichkeit eines bestimmten Antezedens als logischer Wert betrachtet wird (ein **Bit**), ist der Informationsgewinn der Anteil dieses Bits, der basierend auf den Antezedenzen ermittelt werden kann. Die Informationsdifferenz ist die Differenz zwischen dem Informationsgewinn, gegeben die Antezedenzen, und dem Informationsgewinn, gegeben lediglich die vorher bestehende Konfidenz des Sukzedens. Ein wichtiges Merkmal dieser Methode ist, dass die Unterstützung berücksichtigt wird, sodass die Regeln, die für mehr Datensätze gelten, für ein bestimmtes Konfidenzniveau bevorzugt werden. Setzen Sie die Untergrenze des Evaluierungsmaßes auf die Informationsdifferenz, für die Regeln eingehalten werden sollen.

*Hinweis:* Da die Skala für dieses Maß etwas weniger intuitiv ist, als dies bei anderen Skalen der Fall ist, müssen Sie u. U. mit verschiedenen Untergrenzen experimentieren, um eine zufrieden stellende Regelmenge zu erhalten.

- **Normalisiertes Chi-Quadrat.** (Auch als **normalisiertes Chi-Quadrat-Maß** bezeichnet.) Dieses Maß ist ein statistischer Assoziationsindex zwischen Antezedenzen und Sukzedenzen. Dieses Maß wird auf Werte zwischen 0 und 1 normiert. Dieses Maß ist noch stärker abhängig von der Unterstützung als das Informationsdifferenzmaß. Setzen Sie die Untergrenze des Evaluierungsmaßes auf die Informationsdifferenz, für die Regeln eingehalten werden sollen.

*Hinweis:* Da die Skala für dieses Maß wie beim Informationsdifferenzmaß etwas weniger intuitiv ist, als dies bei anderen Skalen der Fall ist, müssen Sie u. U. mit verschiedenen Untergrenzen experimentieren, um eine zufrieden stellende Regelmenge zu erhalten.

**Regeln ohne Antezedenzen zulassen.** Wählen Sie diese Option, um Regeln zuzulassen, die lediglich das Sukzedens (Element oder Element-Set) enthalten. Dies ist sinnvoll, wenn Sie an der Ermittlung häufig verwendeter Elemente oder Element-Sets interessiert sind. Zum Beispiel ist *cannedveg* eine Regel mit nur einem Element ohne Antezedens, die angibt, dass der Kauf von *Konservengemüse* in den Daten häufig vorkommt. In manchen Fällen können Sie derartige Regeln aufnehmen, wenn Sie nur an den wahrscheinlichsten Vorhersagen interessiert sind. Diese Option ist standardmäßig deaktiviert. In der Regel wird die Antezedens-Unterstützung für Regeln ohne Antezedens als 100 % ausgedrückt. Die Regelunterstützung ist gleich der Konfidenz.

## CARMA-Knoten

Der CARMA-Knoten verwendet einen Erkennungsalgorithmus für Regeln, um Assoziationsregeln in den Daten zu erkennen. Assoziationsregeln sind Anweisungen in der Form

**wenn** *Antezedens* **dann** *Sukzedens*

Wenn ein Webkunde z. B. eine drahtlose Karte und einen drahtlosen High-end-Router kauft, wird er wahrscheinlich auch einen drahtlosen Musikserver kaufen, falls ihm einer angeboten wird. Beim CARMA-Modell wird eine Regelmenge aus den Daten extrahiert, ohne dass Sie Eingabe- oder Ziel-Felder angeben müssen. Dies bedeutet, dass die generierten Regeln für einen breiteren Anwendungsbereich verwendet werden können. So können Sie z. B. von diesem Knoten generierte Regeln verwenden, um eine Liste mit Produkten und Dienstleistungen (Antezedenzen) zu finden, deren Sukzedens das Element darstellt, das Sie in der Ferienzeit desselben Jahres

bewerben möchten. Mit IBM® SPSS® Modeler können Sie ermitteln, welche Kunden die Vorgängerprodukte gekauft haben, und eine Marketing-Kampagne für das Nachfolger-Produkt ins Leben rufen.

**Anforderungen.** Im Gegensatz zu A Priori sind für den CARMA-Knoten die Felder *Eingabe* oder *Ziel* nicht erforderlich. Dies ist ein integraler Bestandteil der Funktionsweise des Algorithmus und entspricht dem Erstellen eines A Priori-Modells, wobei alle Felder auf *Beides* gesetzt sind. Sie können beschränken, welche Elemente nur als Antezedenzen oder Sukzedenzen aufgelistet werden, indem Sie das Modell filtern, nachdem es erstellt wurde. So können Sie z. B. den Modell-Browser verwenden, um eine Liste mit Produkten und Dienstleistungen (Antezedenzen) zu finden, deren Nachfolger (Sukzedens) das Element darstellt, das Sie in der Ferienzeit desselben Jahres bewerben möchten.

Um eine CARMA-Regelmenge zu erstellen, müssen Sie ein ID-Feld und mindestens ein Inhaltsfeld angeben. Das ID-Feld kann eine beliebige Rolle oder ein beliebiges Messniveau aufweisen. Felder mit der Rolle *Keine* werden ignoriert. Feldtypen müssen vollständig instanziiert werden, bevor der Knoten ausgeführt wird. Wie bei A Priori können die Daten als Tabellen- oder als Transaktionsdaten vorliegen. [Für weitere Informationen siehe Thema Tabellendaten im Vergleich zu Transaktionsdaten auf S. 389.](#)

**Stärken.** Der CARMA-Knoten basiert auf dem CARMA-Assoziationsregelalgorithmus. Im Gegensatz zu A Priori bietet der CARMA-Knoten Erstellungseinstellungen für die Regelunterstützung (Unterstützung für Antezedens und Sukzedens) und nicht für die Antezedens-Unterstützung. CARMA erlaubt auch Regeln mit mehreren Sukzedenzen. Wie bei A Priori können von einem CARMA-Knoten generierte Modelle zum Erstellen von Vorhersagen in einen Daten-Stream eingefügt werden. [Für weitere Informationen siehe Thema Modell-Nuggets in Kapitel 3 auf S. 45.](#)

### ***Feldoptionen für den CARMA-Knoten***

Vor der Ausführung eines CARMA-Knotens müssen Sie Eingabefelder auf der Registerkarte “Felder” des CARMA-Knotens angeben. Während die meisten Modellierungsknoten dieselben Optionen auf der Registerkarte “Felder” aufweisen, enthält der CARMA-Knoten mehrere einzigartige Optionen. Sämtliche Optionen werden nachfolgend beschrieben.

Abbildung 12-4  
Feldoptionen für den CARMA-Knoten



**Typknoteneinstellungen verwenden.** Diese Option weist den Knoten an, die Feldinformationen von einem weiter oben liegenden Typknoten zu verwenden. Dies ist die Standardeinstellung.

**Benutzerdefinierte Einstellungen verwenden.** Diese Option weist den Knoten an, die hier angegebenen Feldinformationen anstelle der in einem weiter oben liegenden Typknoten angegebenen zu verwenden. Geben Sie nach Auswahl dieser Option die Felder unten danach an, ob die Daten im Transaktions- oder im Tabellenformat gelesen werden.

**Transaktionsformat verwenden.** Mit dieser Option werden die Feldsteuerelemente im übrigen Dialogfeld danach geändert, ob die Daten im Transaktions- oder im Tabellenformat vorliegen. Wenn Sie mehrere Felder für Transaktionsdaten verwenden, stellen die in diesen Feldern für einen bestimmten Datensatz angegebenen Elemente solche Elemente dar, die in einer Einzeltransaktion mit einem einzelnen Zeitstempel gefunden wurden. [Für weitere Informationen siehe Thema Tabellendaten im Vergleich zu Transaktionsdaten auf S. 389.](#)

### **Tabellendaten**

Wenn Transaktionsformat verwenden nicht ausgewählt ist, werden die folgenden Felder angezeigt.

- **Eingaben.** Wählen Sie das Eingabefeld bzw. die Eingabefelder aus. Dies ist so, als würden Sie in einem Typknoten für die Rolle eines Felds den Wert *Eingabe* festlegen.
- **Partition.** In diesem Feld können Sie ein Feld angeben, das verwendet wird, um die Daten in getrennte Stichproben für die Trainings-, Test- und Validierungsphase der Modellbildung aufzuteilen. Indem Sie mit einer Stichprobe das Modell erstellen und es mit einer anderen

Stichprobe testen, erhalten Sie einen guten Hinweis dafür, wie gut das Modell sich für größere Datenmengen generalisieren lässt, die den aktuellen Daten ähneln. Wenn mehrere Partitionsfelder mithilfe von Typ- oder Partitionsknoten erstellt wurden, muss in jedem Modellierungsknoten, der die Partitionierung verwendet, auf der Registerkarte "Felder" ein einzelnes Partitionsfeld ausgewählt werden. (Wenn nur eine einzige Partition vorhanden ist, wird diese immer automatisch verwendet, wenn die Partitionierung aktiviert ist.) Für weitere Informationen siehe Thema Partitionsknoten in Kapitel 4 in *IBM SPSS Modeler 14.2-Quellen-, Prozess- und Ausgabeknoten*. Beachten Sie außerdem, dass die Partitionierung auch auf der Registerkarte "Modelloptionen" des Knotens aktiviert sein muss, wenn die ausgewählte Partitionierung in Ihrer Analyse angewendet werden soll. (Wenn die Auswahl der Option aufgehoben ist, kann die Partitionierung ohne Änderung der Feldeinstellungen deaktiviert werden.)

### **Transaktionsdaten**

Wenn Sie Transaktionsformat verwenden auswählen, werden die folgenden Felder angezeigt.

- **ID.** Wählen Sie für Transaktionsdaten ein ID-Feld aus der Liste aus. Als ID-Feld können numerische oder symbolische Felder verwendet werden. Jeder eindeutige Wert in diesem Feld sollte eine bestimmte Analyseeinheit darstellen. Bei einer Warenkorbanwendung könnte z. B. jede ID einen einzelnen Kunden darstellen. Für eine Webprotokoll-Analyseanwendung könnte jede ID einen Computer (nach IP-Adresse) oder einen Benutzer (nach Anmeldedaten) darstellen.
- **IDs sind zusammenhängend.** (Nur A Priori- und CARMA-Knoten) Wenn Ihre Daten vorsortiert sind, sodass alle Datensätze mit derselben ID im Daten-Stream zusammengefasst sind, wählen Sie diese Option, um die Verarbeitung zu beschleunigen. Wenn Ihre Daten nicht vorsortiert sind (oder Sie nicht sicher sind), lassen Sie diese Option deaktiviert. Die Daten werden dann vom Knoten automatisch sortiert.

*Hinweis:* Wenn Ihre Daten nicht sortiert werden und Sie diese Option auswählen, könnten Sie ungültige Ergebnisse in Ihrem Modell erhalten.

- **Inhalt.** Geben Sie die Inhaltsfelder für das Modell an. Diese Felder enthalten die Elemente, die für die Assoziationsmodellierung interessant sind. Sie können mehrere Flag-Felder angeben (falls die Daten in tabellarischer Form vorliegen) oder ein einzelnes nominales Feld (falls die Daten im Transaktionsformat vorliegen).

## Modelloptionen für den CARMA-Knoten

Abbildung 12-5  
Modelloptionen für den CARMA-Knoten



**Modellname.** Sie können den Modellnamen automatisch basierend auf den Ziel- oder ID-Feldnamen (oder dem Modelltyp in Fällen, in denen kein solches Feld angegeben ist) generieren oder einen benutzerdefinierten Namen eingeben.

**Minimale Regelunterstützung (%).** Sie können auch ein Stützkriterium angeben.

**Regelunterstützung** bezieht sich auf den Anteil von IDs in den Trainingsdaten, die die gesamte Regel enthalten. (Beachten Sie, dass diese Definition von der für A Priori-Knoten verwendeten Antezedens-Unterstützung abweicht.) Wenn Sie weitere gemeinsame Regeln wünschen, erhöhen Sie diese Einstellung.

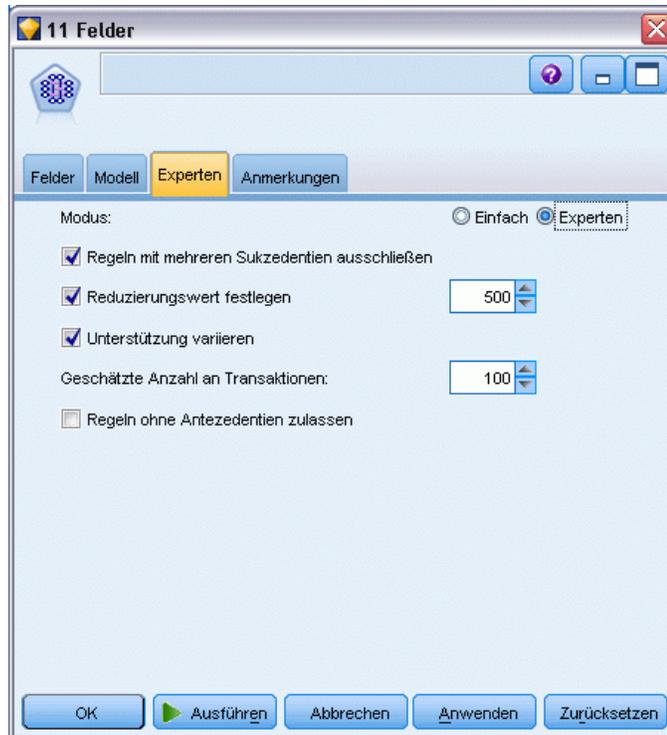
**Minimale Regelkonfidenz (%).** Sie können ein Konfidenzkriterium angeben, um die Regeln in der Regelmenge beizubehalten. **Konfidenz** bezieht sich auf den prozentualen Anteil von IDs, für die eine richtige Vorhersage gemacht wird (aus allen IDs, für die die Regel eine Vorhersage macht). Sie wird aus der Anzahl von IDs berechnet, für die die gesamte Regel gefunden wird, dividiert durch die Anzahl der IDs, für die die Vorgänger gefunden werden, basierend auf den Trainingsdaten. Regeln mit einer niedrigeren Konfidenz als das angegebene Kriterium werden verworfen. Wenn Sie uninteressante oder zu viele Regeln erhalten, sollten Sie diese Einstellung erhöhen. Wenn Sie zu wenige Regeln erhalten, sollten Sie diese Einstellung reduzieren.

**Maximale Regelgröße.** Sie können die maximale Anzahl unterschiedlicher *Element-Sets* (im Gegensatz zu *Elementen*) in einer Regel festlegen. Wenn die gewünschten Regeln relativ kurz sind, können Sie diese Einstellung reduzieren, um die Erstellung der Regelmenge zu beschleunigen.

## Expertenoptionen für den CARMA-Knoten

Für Personen mit umfassenden Kenntnissen über die Operation des CARMA-Knotens ermöglichen die folgenden Expertenoptionen die Feinabstimmung des Modellierungsvorgangs. Um auf die Expertenoptionen zuzugreifen, legen Sie auf der Registerkarte “Experten” den Modus Experten fest.

Abbildung 12-6  
Expertenoptionen für den CARMA-Knoten



**Regeln mit mehreren Sukzedenzien ausschließen.** Wählen Sie diese Option, um “doppelköpfige” Sukzedenzien – d. h., Sukzedenzien, die zwei Elemente enthalten – auszuschließen. Die Regel bread & cheese & fish -> wine&fruit beispielsweise enthält ein doppelköpfiges Sukzedens, wine&fruit. Standardmäßig sind solche Regeln enthalten.

**Reduzierungswert festlegen.** Um Speicherplatz zu sparen, entfernt (**reduziert**) der verwendete CARMA-Algorithmus in regelmäßigen Abständen während der Verarbeitung seltene Element-Sets aus seiner Liste möglicher Elemente. Wählen Sie diese Option, um die Reduktionshäufigkeit anzupassen. Die von Ihnen angegebene Zahl ermittelt die Reduktionshäufigkeit. Geben Sie einen kleinen Wert ein, um die Speicheranforderungen des Algorithmus zu reduzieren (möglicherweise aber die erforderliche Trainingszeit zu erhöhen), oder geben Sie einen hohen Wert ein, um die Trainingsgeschwindigkeit zu erhöhen (möglicherweise aber die Speicheranforderungen zu erhöhen). Der Standardwert lautet 500.

**Unterstützung variieren.** Wählen Sie diese Option, um die Effizienz zu erhöhen, indem Sie seltene Element-Sets ausführen, die den Eindruck erwecken, als wären sie häufig, wenn Sie ungleich verteilt vorkommen. Dies erreichen Sie durch ein höheres Unterstützungsniveau und durch die

Reduktion auf den auf der Registerkarte “Modell” angegebenen Wert. Geben Sie einen Wert für Geschätzte Anzahl an Transaktionen ein, um anzugeben, wie schnell das Unterstützungsniveau reduziert werden soll.

**Regeln ohne Antezedenzen zulassen.** Wählen Sie diese Option, um Regeln zuzulassen, die lediglich das Sukzedens (Element oder Element-Set) enthalten. Dies ist sinnvoll, wenn Sie an der Ermittlung häufig verwendeter Elemente oder Element-Sets interessiert sind. Zum Beispiel ist *cannedveg* eine Regel mit nur einem Element ohne Antezedens, die angibt, dass der Kauf von *Konservengemüse* in den Daten häufig vorkommt. In manchen Fällen können Sie derartige Regeln aufnehmen, wenn Sie nur an den wahrscheinlichsten Vorhersagen interessiert sind. Diese Option ist standardmäßig deaktiviert.

## ***Assoziationsregelmodell-Nuggets***

Assoziationsregelmodell-Nuggets enthalten die Regeln, die von einem der folgenden Modellierungsknoten für Assoziationsregeln entdeckt wurden:

- A Priori
- CARMA

Die Modell-Nuggets enthalten Informationen zu den Regeln, die bei der Modellerstellung aus den Daten extrahiert wurden.

### ***Anzeigen von Ergebnissen***

Sie können die von den Assoziationsmodellen (A Priori und CARMA) und den Sequenzmodellen generierten Regeln mithilfe der Registerkarte “Modell” im Dialogfeld durchsuchen. Beim Durchsuchen eines Modell-Nuggets erhalten Sie Informationen zu den Regeln sowie Optionen zum Filtern und Sortieren der Ergebnisse vor der Erstellung neuer Knoten oder der Bewertung des Modells.

### ***Bewerten des Modells***

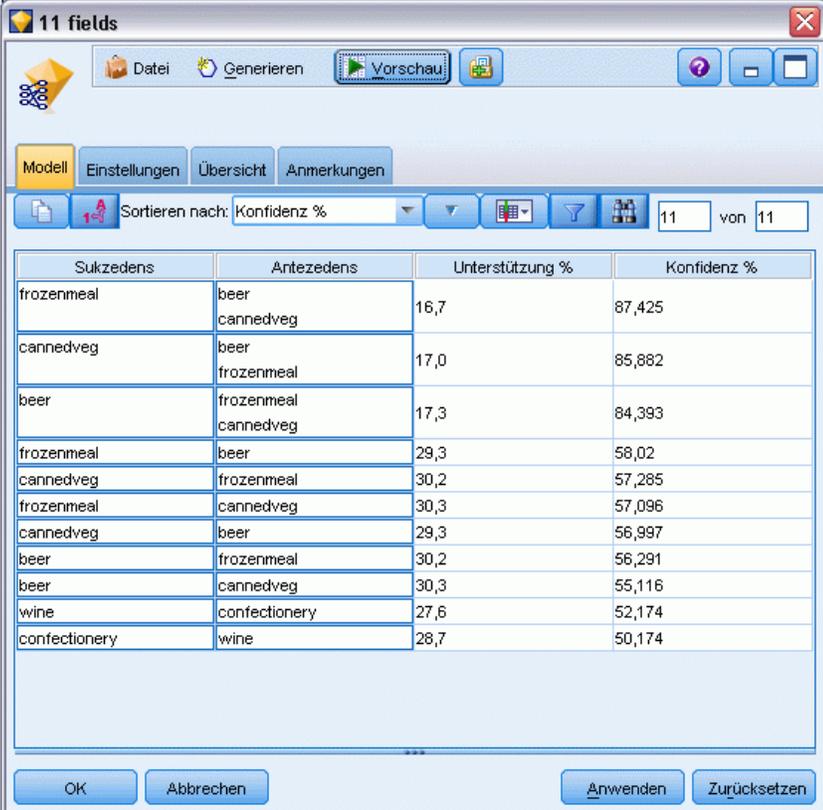
Verfeinerte Modell-Nuggets (A Priori, CARMA und Sequenz) können zum Stream hinzugefügt und für die Bewertung verwendet werden. [Für weitere Informationen siehe Thema Verwendung von Modell-Nuggets in Streams in Kapitel 3 auf S. 68.](#) Zu den für die Bewertung verwendeten Modell-Nuggets gehört eine zusätzliche Registerkarte “Einstellungen” in den entsprechenden Dialogfeldern. [Für weitere Informationen siehe Thema Einstellungen beim Assoziationsregelmodell-Nugget auf S. 407.](#)

Ein nicht verfeinertes Modell-Nugget kann nicht im Rohformat für die Bewertung verwendet werden. Stattdessen können Sie eine Regelmenge erstellen und für die Bewertung verwenden. [Für weitere Informationen siehe Thema Generieren einer Regelmenge aus einem Assoziationsmodell-Nugget auf S. 410.](#)

## Nähere Informationen zum Assoziationsregelmodell-Nugget

Auf der Registerkarte “Modell” eines Assoziationsregelmodell-Nuggets wird eine Tabelle angezeigt, die die vom Algorithmus extrahierten Regeln enthält. Jede Zeile in der Tabelle steht für eine Regel. Die erste Spalte enthält die Sukzedenzen (den “dann”-Teil der Regel), während die nächste Spalte die Antezedenzen (den “wenn”-Teil der Regel) enthält. Die weiteren Spalten enthalten Informationen zur Regel, wie Konfidenz, Unterstützung und Lift.

Abbildung 12-7  
Assoziationsregel-Nugget Registerkarte “Modell”



Sukzedens	Antezedens	Unterstützung %	Konfidenz %
frozenmeal	beer cannedveg	16,7	87,425
cannedveg	beer frozenmeal	17,0	85,882
beer	frozenmeal cannedveg	17,3	84,393
frozenmeal	beer	29,3	58,02
cannedveg	frozenmeal	30,2	57,285
frozenmeal	cannedveg	30,3	57,096
cannedveg	beer	29,3	56,997
beer	frozenmeal	30,2	56,291
beer	cannedveg	30,3	55,116
wine	confectionery	27,6	52,174
confectionery	wine	28,7	50,174

Assoziationsregeln werden häufig in folgendem Format angezeigt:

<b>Sukzedens</b>	<b>Antezedens</b>
Drug = drugY	Sex = F BP = HIGH

Die Beispielregel wird wie folgt interpretiert: *Wenn Geschlecht = “W” und BD = “HOCH”, ist das Medikament wahrscheinlich MedikamentY*, oder anders ausgedrückt: *Bei Datensätzen mit Geschlecht = “W” und BD = “HOCH” ist das Medikament wahrscheinlich MedikamentY*. Mit der Symbolleiste des Dialogfelds können Sie weitere Informationen anzeigen, beispielsweise Konfidenz, Unterstützung und Instanzen.

**Menü "Sortieren".** Mit der Schaltfläche des Menüs "Sortieren" in der Symbolleiste wird die Sortierung der Regeln gesteuert. Die Sortierrichtung (aufsteigend oder absteigend) lässt sich mit der Schaltfläche für die Sortierrichtung (nach unten bzw. oben zeigender Pfeil) ändern.

Abbildung 12-8  
Symbolleistenoptionen für das Sortieren



Regeln können nach folgenden Faktoren sortiert werden:

- Support
- Confidence
- Regelunterstützung
- Sukzedens
- Lift
- Deployability

**Menü "Anzeigen/Ausblenden"** Das Menü "Anzeigen/Ausblenden" (Symbolleistenschaltfläche "Kriterien") steuert die Optionen für die Anzeige der Regeln.

Abbildung 12-9  
Schaltfläche "Anzeigen/Ausblenden"



Die folgenden Anzeigeeoptionen stehen zur Verfügung:

- **Regel-ID** zeigt die während der Modellerstellung zugewiesene Regel-ID an. Mit einer Regel-ID können Sie identifizieren, welche Regeln für eine bestimmte Vorhersage angewendet werden. Mit Regel-IDs können Sie außerdem zusätzliche Regelinformationen zu einem späteren Zeitpunkt zusammenführen, beispielsweise Verwendbarkeit, Produktinformationen oder Antezedenzen.
- **Instanzen** zeigt Informationen über die Anzahl der eindeutigen IDs an, auf die die Regel zutrifft, d. h., für die die Antezedenzen wahr sind. Bei der Regel bread -> cheese beispielsweise wird die Anzahl der Datensätze in den Trainingsdaten, die das Antezedens Brot enthalten, als **Instanzen** bezeichnet.
- **Antezedens-Unterstützung** zeigt die Antezedens-Unterstützung an, also den Teil der IDs, für die die Antezedenzen wahr sind. Beispiel: Wenn 50 % der Trainingsdaten den Kauf von Brot beinhalten, hat die Regel bread -> cheese eine Antezedens-Unterstützung von 50 %. *Hinweis:* "Unterstützung" in diesem Kontext ist dasselbe wie Instanzen, nur als Prozentsatz ausgedrückt.
- **Konfidenz** zeigt das Verhältnis von Regelunterstützung zu Antezedens-Unterstützung an. Es wird also der Anteil der IDs mit den angegebenen Antezedenzen angegeben, bei denen auch das Sukzedens (bzw. die Sukzedenzen) wahr ist. Beispiel: Wenn 50 % der Trainingsdaten "Brot" enthalten (Antezedens-Unterstützung), jedoch nur 20 % sowohl "Brot" als auch "Käse" (Regelunterstützung), ist die Konfidenz für die Regel bread -> cheeseRule Support / Antecedent Support, in diesem Fall 40 %.

- **Regelunterstützung** zeigt den Teil der IDs an, für die die gesamte Regel, die Antezedenzen und das Sukzedens (bzw. die Sukzedenzen) wahr sind. Beispiel: Wenn 20 % der Trainingsdaten sowohl “Brot” als auch “Käse” enthalten, beträgt die Regelunterstützung für die Regel bread -> cheese 20 %.
- **Lift** zeigt das Verhältnis zwischen der Konfidenz für die Regel und der A Priori-Wahrscheinlichkeit für das Sukzedens an. Beispiel: Wenn 10 % der Gesamtbevölkerung Brot kauft, weist eine Regel, die mit einer Konfidenz von 20 % vorhersagt, ob die Leute Brot kaufen, einen Lift von  $20/10 = 2$  auf. Wenn eine andere Regel mit einer Konfidenz von 11 % besagt, dass die Leute Brot kaufen, weist diese Regel einen Lift von annähernd 1 auf, was bedeutet, dass die Antezedenzen keinen großen Unterschied für die Wahrscheinlichkeit des Sukzedens aufweisen. Im Allgemeinen sind Regeln mit einem Lift, der sich von 1 unterscheidet, interessanter als Regeln mit einem Lift von annähernd 1.
- **Verwendbarkeit** ist ein Maß dafür, welcher Prozentsatz der Trainingsdaten die Bedingungen des Antezedens erfüllen, nicht jedoch die Bedingungen des Sukzedens. Beim Einkauf von Produkten bedeutet dies im Grunde, welcher Prozentsatz des Kundenstamms das Produkt aus dem Antezedens besitzt (bzw. erworben hat), jedoch noch nicht das im Sukzedens festgelegte Produkt gekauft hat. Die Verwendbarkeitsstatistik ist definiert als  $((\text{Antecedent Support in \# of Records} - \text{Rule Support in \# of Records}) / \text{Number of Records}) * 100$ , wobei *Antezedens-Unterstützung* für die Anzahl der Datensätze steht, bei denen die Antezedenzen wahr sind, und *Regelunterstützung* für die Anzahl der Datensätze, bei denen sowohl die Antezedenzen als auch das Sukzedens wahr sind.

**Schaltfläche “Filter”.** Mit der Filterschaltfläche (Trichtersymbol) im Menü wird der untere Teil des Dialogfelds erweitert und ein Fenster mit aktiven Regelfiltern wird angezeigt. Filter werden verwendet, um die Anzahl der auf der Registerkarte “Modelle” angezeigten Regeln einzugrenzen.

Abbildung 12-10  
Schaltfläche “Filter”



Zum Erstellen von Filtern klicken Sie auf das Filtersymbol rechts neben dem erweiterten Fenster. Dadurch wird ein separates Dialogfeld geöffnet, in dem Sie Bedingungen für die Anzeige von Regeln eingeben können. Beachten Sie, dass die Filterschaltfläche häufig in Verbindung mit dem Menü “Generieren” verwendet wird, um zunächst die Regeln zu filtern und anschließend ein Modell zu erstellen, das die betreffende Untergruppe der Regeln enthält. Weitere Informationen finden Sie unter [Angeben von Filtern für Regeln](#).

**Schaltfläche “Regel suchen”.** Die Schaltfläche “Regel finden” (Fernglassymbol) ermöglicht das Durchsuchen der angezeigten Regeln nach einer angegebenen Regel-ID. Im angrenzenden Dialogfeld wird angegeben, wie viele der verfügbaren Regeln derzeit angezeigt werden. Regel-IDs werden vom Modell in der Reihenfolge ihrer Entdeckung zugewiesen und werden während der Bewertung zu den Daten hinzugefügt.

Abbildung 12-11  
Schaltfläche “Regel suchen”



So können Sie Regel-IDs neu ordnen:

- ▶ Sie können die Regel-IDs in IBM® SPSS® Modeler neu anordnen, indem Sie zuerst die Regelanzeigtabelle gemäß dem gewünschten Maß sortieren, beispielsweise “Konfidenz” oder “Lift”.
- ▶ Anschließend erstellen Sie mit den Optionen aus dem Menü “Generieren” ein gefiltertes Modell.
- ▶ Wählen Sie im Dialogfeld “Gefiltertes Modell” die Option Regeln neu nummerieren, beginnend mit und geben Sie eine Startnummer an.

Für weitere Informationen siehe Thema Erstellen eines gefilterten Modells auf S. 411.

### Angeben von Filtern für Regeln

Standardmäßig können Regelalgorithmen wie “A Priori”, “CARMA” und “Sequenz” sehr große und umständlich zu handhabende Mengen von Regeln generieren. Zugunsten einer größeren Klarheit beim Durchsuchen bzw. zur Rationalisierung der Regelbewertung sollten Sie in Erwägung ziehen, die Regeln zu filtern, sodass die für Sie relevanten Sukzedenzen und Antezedenzen deutlicher zu sehen sind. Mithilfe der Filteroptionen auf der Registerkarte “Modell” eines Regel-Browsers können Sie ein Dialogfeld zur Angabe der Filterbedingungen öffnen.

Abbildung 12-12  
Regelbrowser – Filterdialogfeld

The screenshot shows the 'Filter bearbeiten' dialog box with the following settings:

- Sukzedenzen:** Filter aktivieren (unchecked), Mindestens eines einschließen aus, Werte: beer, confectionery, wine.
- Antezedenzen:** Filter aktivieren (unchecked), Mindestens eines einschließen aus, Werte: beer, cannedveg, frozenmeal, fruitveg.
- Konfidenz:** Filter aktivieren (unchecked), Oberhalb, Min. 0, Max. 100.
- Antezedens-Unterstützung:** Filter aktivieren (checked), Oberhalb, Min. 60,0, Max. 100.
- Lift:** Filter aktivieren (unchecked), Größer als, Min. 1, Max. 1.

Buttons at the bottom: OK, Abbrechen, Hilfe.

**Sukzedenzien.** Mit der Option Filter aktivieren können Sie Optionen für das Filtern von Regeln definieren, die auf der Aufnahme bzw. dem Ausschluss angegebener Sukzedenzien beruhen. Wählen Sie Mindestens eines einschließen aus, um einen Filter zu erstellen, bei dem die Regeln mindestens einen der angegebenen Sukzedenzien enthalten. Wählen Sie alternativ Ausschließen, um einen Filter zu erstellen, der die angegebenen Sukzedenzien ausschließt. Sie können die Sukzedenzien mithilfe des Auswahlsymbols rechts neben dem Listenfeld auswählen. Dadurch wird ein Dialogfeld geöffnet, das alle Sukzedenzien auflistet, die in den generierten Regeln vorliegen.

*Hinweis:* Sukzedenzien können mehrere Elemente enthalten. Mit den Filtern wird nur überprüft, ob ein Sukzedens eines der angegebenen Elemente enthält.

**Antezedenzen.** Mit der Option Filter aktivieren können Sie Optionen für das Filtern von Regeln definieren, die auf der Aufnahme bzw. dem Ausschluss angegebener Antezedenzen beruhen. Sie können die gewünschten Elemente mithilfe des Auswahlsymbols rechts neben dem Listenfeld auswählen. Dadurch wird ein Dialogfeld geöffnet, das alle Antezedenzen auflistet, die in den generierten Regeln vorliegen.

- Wählen Sie die Option Alle einschließen aus, um den Filter als Einschlussfilter festzulegen, bei dem alle angegebenen Antezedenzen in einer Regel enthalten sein müssen.
- Wählen Sie Mindestens eines einschließen aus, um einen Filter zu erstellen, bei dem die Regeln mindestens einen der angegebenen Antezedenzen enthalten.
- Wählen Sie Ausschließen aus, um einen Filter zu erstellen, der Regeln ausschließt, die ein angegebenes Antezedens enthalten.

**Konfidenz.** Mit der Option Filter aktivieren können Sie Optionen für das Filtern von Regeln definieren, die auf dem Konfidenzniveau der jeweiligen Regel beruhen. Mit den Min- und Max-Steurelementen können Sie einen Konfidenzbereich angeben. Beim Durchsuchen der generierten Modelle wird die Konfidenz als Prozentsatz angegeben. Bei der Bewertung von Ausgaben wird die Konfidenz als Zahl zwischen 0 und 1 angegeben.

**Antezedens-Unterstützung.** Mit der Option Filter aktivieren können Sie Optionen für das Filtern von Regeln definieren, die auf dem Niveau der Antezedens-Unterstützung für die jeweilige Regel beruhen. Die Antezedens-Unterstützung gibt den Anteil der Trainingsdaten an, die dieselben Antezedenzen wie die aktuelle Regel beinhalten, ähnlich einem Popularitätsindex. Mit den Min- und Max-Steurelementen können Sie einen Bereich angeben, der zum Filtern der Regel anhand des Unterstützungsniveaus verwendet wird.

**Lift.** Mit der Option Filter aktivieren können Sie Optionen für das Filtern von Regeln definieren, die auf dem Lift-Maß für die jeweilige Regel beruhen. *Hinweis:* Die Lift-Filterung ist nur für Assoziationsmodelle verfügbar, die nach Version 8.5 erstellt wurden, und für frühere Modelle, die ein Lift-Maß enthalten. Bei Sequenzmodellen steht diese Option nicht zur Verfügung.

Klicken Sie auf OK, um alle Filter anzuwenden, die in diesem Dialogfeld aktiviert wurden.

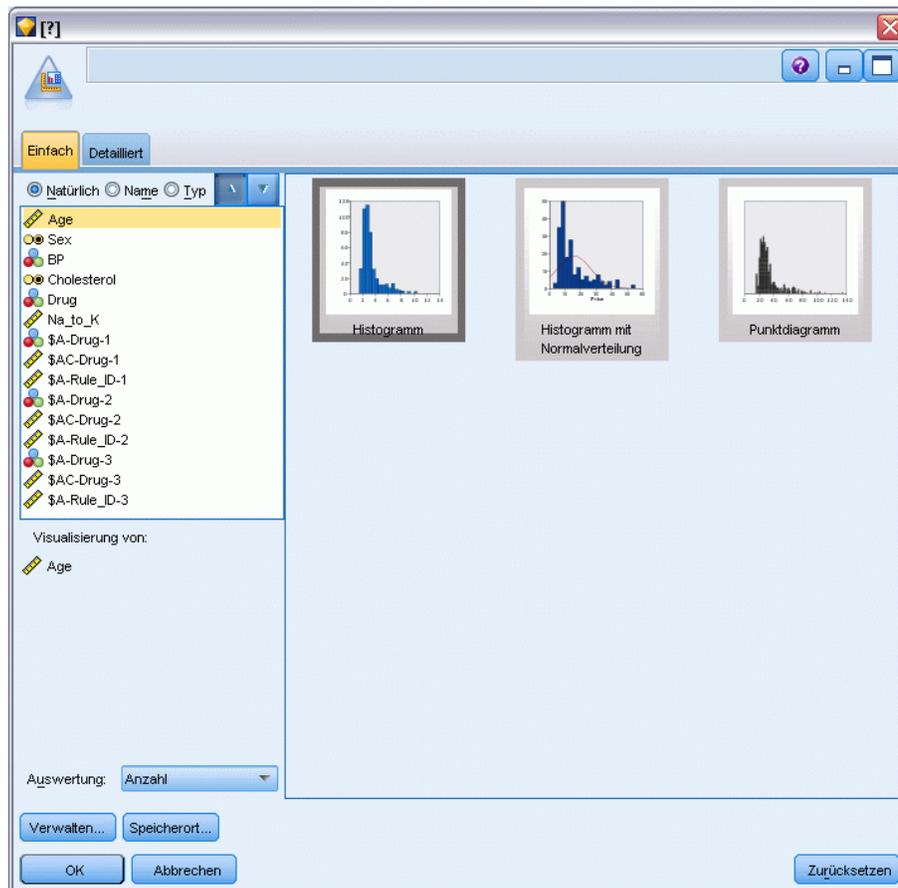
### ***Erzeugen von Diagrammen für Regeln***

Die Assoziationsknoten bieten eine Menge an Informationen, jedoch sind diese nicht unbedingt immer in einem leicht zugänglichen Format für Unternehmensanwender. Zur Bereitstellung der Daten auf eine Art, die problemlos in Geschäftsberichte, Präsentationen u.s.w. integriert werden

kann, können aus ausgewählten Daten Diagramme erstellt werden. In der Registerkarte “Modell” können Sie ein Diagramm für eine ausgewählte Regel erzeugen und damit ein Diagramm nur für die Fälle dieser Regel erstellen.

- ▶ Wählen Sie in der Registerkarte “Modell” die Regel aus, die Sie interessiert.
- ▶ Wählen Sie im Menü “Generieren” den Befehl Diagramm (von Auswahl). Die Registerkarte “Einfach” der Diagrammtafel wird angezeigt.

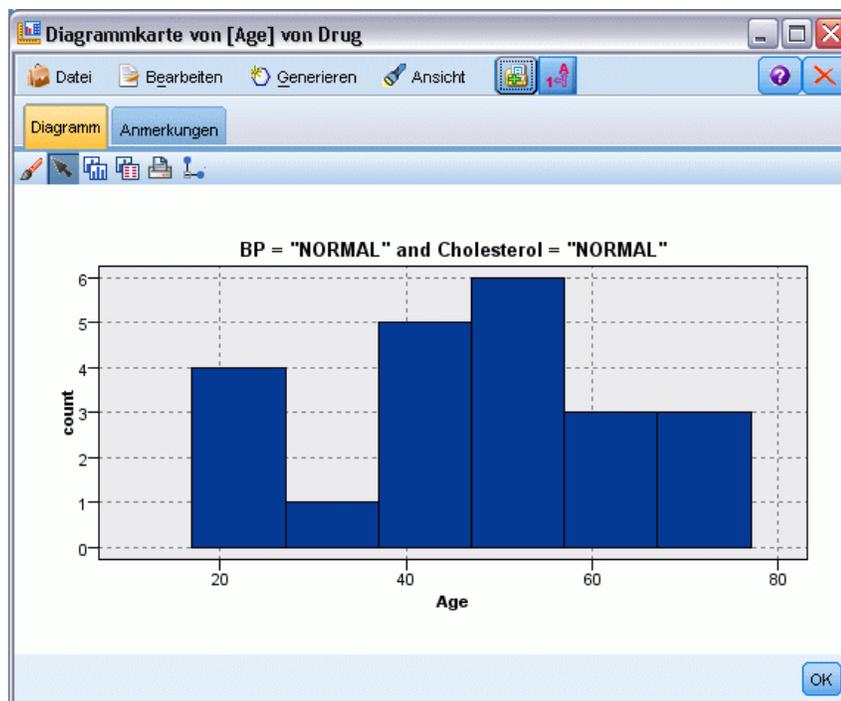
Abbildung 12-13  
Dialogfeld Grafiktafelknoten, Registerkarte “Basis”



*Hinweis:* Wenn Sie die Grafiktafel auf diese Art anzeigen, sind nur die Registerkarten “Basis” und “Details” verfügbar. Für weitere Informationen siehe Thema Diagrammtafelknoten in Kapitel 5 in *IBM SPSS Modeler 14.2- Quellen-, Prozess- und Ausgabeknoten*.

- ▶ Mithilfe der Einstellungen auf den Registerkarten “Basis” oder “Details” können Sie die Details angeben, die auf dem Diagramm angezeigt werden sollen.
- ▶ Klicken Sie auf “OK”, um das Diagramm zu erstellen.

Abbildung 12-14  
Dialogfeld Grafiktafelknoten, Registerkarte "Basis"



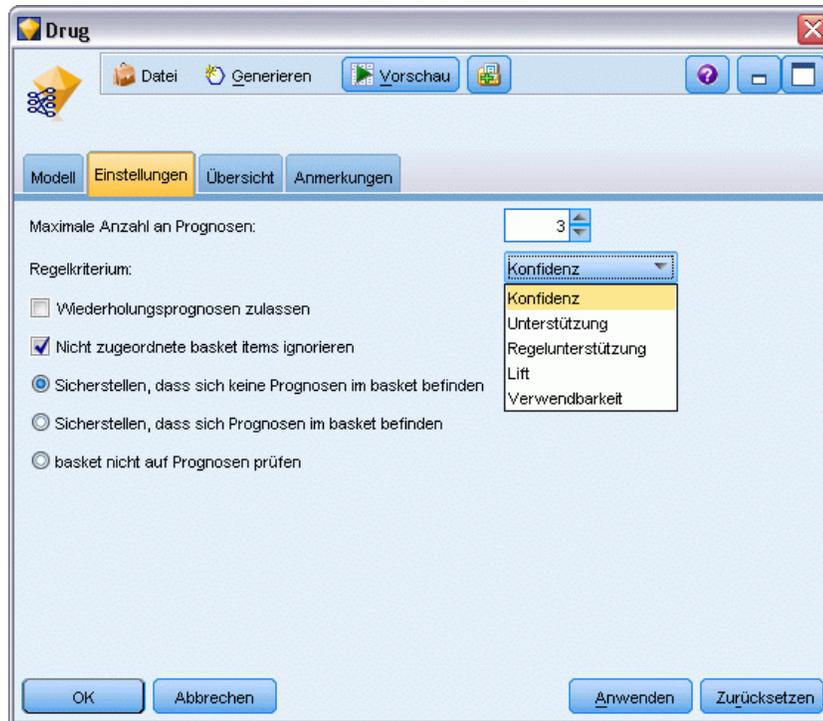
Die Überschrift des Diagramms identifiziert die Regel und Antezedens-Details, die im Diagramm berücksichtigt werden.

### **Einstellungen beim Assoziationsregelmodell-Nugget**

Die Registerkarte "Einstellungen" wird zur Angabe der Scoring-Optionen für Zuordnungsmodelle ("A Priori" und "CARMA") verwendet. Diese Registerkarte ist erst dann verfügbar, wenn das Modell-Nugget zum Zwecke des Scorings zu einem Stream hinzugefügt wurde.

*Hinweis:* Das Dialogfeld zum Durchsuchen eines nicht verfeinerten -Modells enthält nicht die Registerkarte "Einstellungen", da es nicht gescort werden kann. Zum Scoring des "nicht verfeinerten" GRI-Modells müssen Sie zunächst eine Regelmenge erstellen. [Für weitere Informationen siehe Thema Generieren einer Regelmenge aus einem Assoziationsmodell-Nugget auf S. 410.](#)

Abbildung 12-15  
Einstellungen beim Assoziationsregelmodell-Nugget



**Maximale Anzahl an Prognosen.** Dient zur Angabe der maximalen Anzahl an Vorhersagen, die für jedes Set von Warenkorbelementen (basket items) aufgenommen werden. Diese Option wird in Verbindung mit “Regelkriterium” (unten) verwendet, um “Top”-Prognosen zu erstellen. Dabei steht *Top* für das höchste Niveau an Konfidenz, Unterstützung, Lift usw., wie unten angegeben.

**Regelkriterium.** Dient zur Auswahl des Maßes, das zur Ermittlung der Stärke der Regeln verwendet wird. Die Regeln werden nach der Stärke der hier ausgewählten Kriterien sortiert, um die Top-Prognosen für ein Element-Set auszugeben. Folgende Kriterien stehen zur Verfügung:

- Confidence
- Support
- Regelunterstützung (Unterstützung \* Konfidenz)
- Lift
- Deployability

**Wiederholungsprognosen zulassen.** Wählen Sie diese Option aus, um mehrere Regeln mit demselben Sukzedens beim Scoring aufzunehmen. Bei Auswahl dieser Option können beispielsweise folgende Regeln gescort werden:

Brot & Käse -> Wein  
Käse & Obst -> Wein

Deaktivieren Sie diese Option, um beim Scoring Wiederholungsprognosen auszuschließen.

*Hinweis:* Regeln mit mehreren Sukzedenzen (bread & cheese & fruit -> wine & pate) werden nur dann als Wiederholungsprognosen angesehen, wenn alle Sukzedenzen (wine & pate) bereits vorhergesagt wurden.

**Nicht zugeordnete basket items ignorieren.** Wählen Sie diese Option, um das Vorliegen zusätzlicher Elemente im Element-Set zu ignorieren. Wenn diese Option beispielsweise für einen Warenkorb ausgewählt wird, der [tent & sleeping bag & kettle] enthält, gilt die Regel tent & sleeping bag -> gas\_stove trotz des zusätzlichen Elements (kettle) im Warenkorb.

Es kann Umstände geben, unter denen zusätzliche Elemente ausgeschlossen werden sollten. So ist es beispielsweise wahrscheinlich, dass jemand, der ein Zelt, einen Schlafsack und einen Wasserkessel kauft, bereits einen Gaskocher besitzt, worauf der Wasserkessel hindeutet. Anders ausgedrückt: Ein Gaskocher ist möglicherweise nicht die beste Prognose. In solchen Fällen sollten Sie die Auswahl von Nicht zugeordnete basket items ignorieren aufheben, um sicherzustellen, dass die Antezedenzen der Regel genau mit dem Inhalt eines Warenkorbs übereinstimmen. Standardmäßig werden nicht übereinstimmende Elemente ignoriert.

**Sicherstellen, dass sich keine Prognosen im basket befinden.** Wählen Sie diese Option aus, um sicherzustellen, dass die Elemente aus den Sukzedenzen nicht ebenfalls im Warenkorb vorhanden sind. Beispiel: Wenn das Ziel des Scorings darin besteht, eine Produktempfehlung für Möbel abzugeben, ist es unwahrscheinlich, dass bei einem Warenkorb, der bereits einen Esszimmertisch enthält, ein weiterer erworben wird. In solchen Fällen, sollten Sie diese Option auswählen. Andererseits können bei verderblichen Produkten und Einwegartikeln (wie Käse, Säuglingsnahrung oder Papiertaschentüchern) Regeln, bei denen das Sukzedens bereits im Warenkorb vorhanden ist, sinnvoll sein. Im letzteren Fall ist die nützlichste Option möglicherweise basket nicht auf Prognosen prüfen (siehe unten).

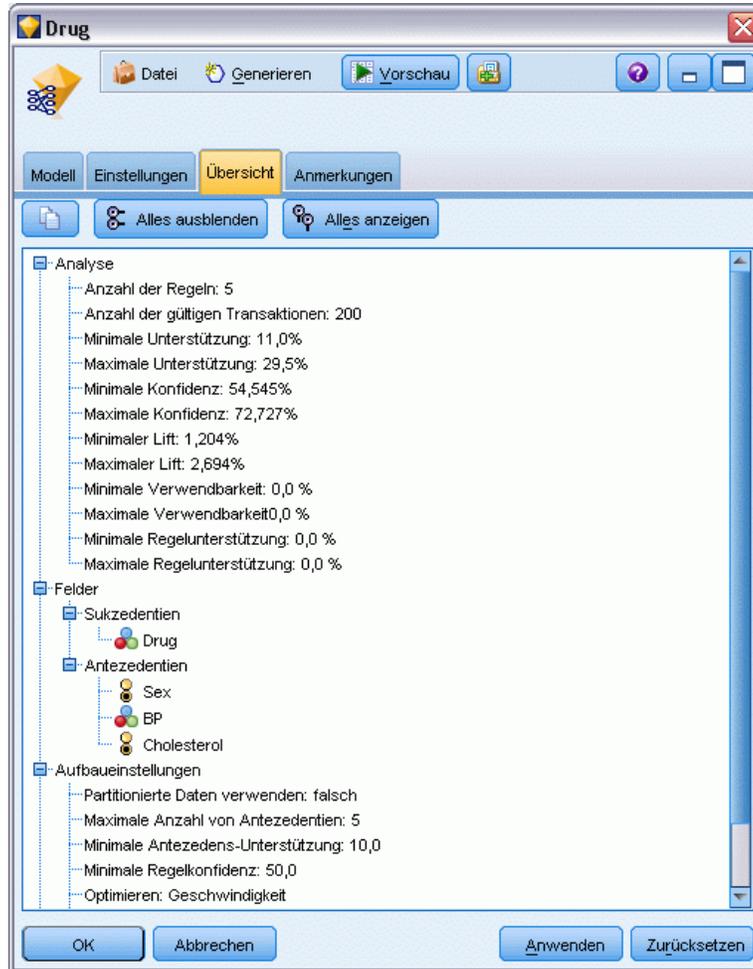
**Sicherstellen, dass sich Prognosen im basket befinden.** Wählen Sie diese Option aus, um sicherzustellen, dass die Elemente aus den Sukzedenzen auch im Warenkorb vorhanden sind. Dieser Ansatz ist sinnvoll, wenn Sie versuchen, einen Einblick in bestehende Kunden oder Transaktionen zu gewinnen. Sie könnten beispielsweise die Regeln mit dem höchsten Lift ermitteln und dann untersuchen, auf welche Kunden diese Regeln zutreffen.

**basket nicht auf Prognosen prüfen** Wählen Sie diese Option aus, um beim Scoring alle Regeln einzuschließen, unabhängig vom Vorhandensein oder Nichtvorhandensein der Sukzedenzen im Warenkorb.

## ***Übersicht über das Assoziationsregelmodell-Nugget***

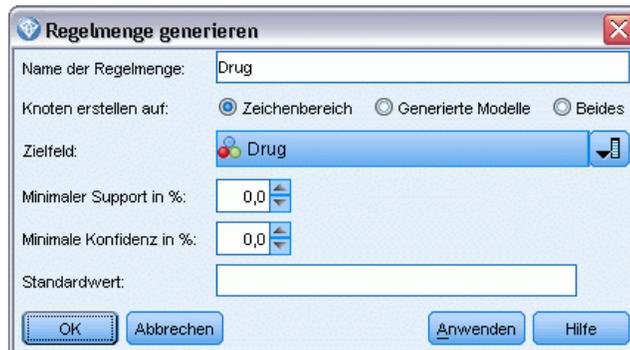
Auf der Registerkarte “Übersicht” für ein Sequenzregelmodell-Nugget werden die Anzahl der ermittelten Regeln sowie die Mindest- und Höchstwerte für Unterstützung, Lift, Konfidenz und Bereitstellbarkeit für die Regeln angegeben.

Abbildung 12-16  
Einstellungen beim Assoziationsregelmodell-Nugget



## Generieren einer Regelmenge aus einem Assoziationsmodell-Nugget

Abbildung 12-17  
Dialogfeld "Regelmenge generieren"



Assoziationsmodell-Nuggets wie “A Priori” oder “CARMA” können zur direkten Speicherung von Daten verwendet werden. Alternativ können Sie zunächst eine Untergruppe von Regeln erstellen, eine sogenannte **Regelmenge**. Regelmengen sind besonders nützlich bei der Arbeit mit einem nicht verfeinerten -Modell, das nicht unmittelbar für das Scoring verwendet werden kann. [Für weitere Informationen siehe Thema Nicht verfeinerte Modelle in Kapitel 3 auf S. 74.](#)

Um eine Regelmenge zu generieren, wählen Sie die Option Regelmenge aus dem Menü “Generieren” im Browser für Modell-Nuggets aus. Folgende Optionen für die Übersetzung der Regeln in eine Regelmenge können angegeben werden:

**Name der Regelmenge.** Ermöglicht die Angabe eines Namens für den neu generierten Regelmengenknoten.

**Knoten erstellen auf.** Steuert den Standort des neu generierten Regelmengenknotens. Wählen Sie Zeichenbereich, Generierte Modelle oder Beides aus.

**Zielfeld.** Legt fest, welches Ausgabefeld für den generierten Regelmengenknoten verwendet wird. Wählen Sie ein einzelnes Ausgabefeld in der Liste aus.

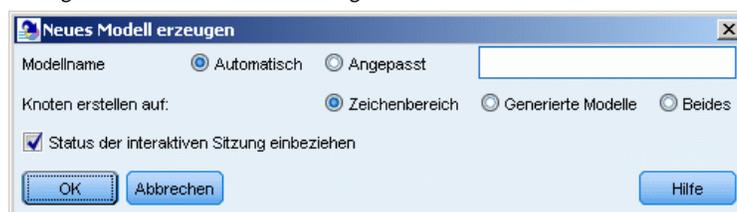
**Minimaler Support.** Geben Sie die minimale Unterstützung für Regeln an, die in der generierten Regelmenge erhalten bleiben sollen. Regeln, deren Unterstützung unter dem angegebenen Wert liegt, werden nicht in die neue Regelmenge aufgenommen.

**Minimale Konfidenz.** Geben Sie die minimale Konfidenz für Regeln an, die in der generierten Regelmenge erhalten bleiben sollen. Regeln, deren Konfidenz unter dem angegebenen Wert liegt, werden nicht in die neue Regelmenge aufgenommen.

**Standardwert.** Ermöglicht die Angabe eines Standardwerts für das Zielfeld, der gescorten Datensätzen zugewiesen wird, auf die keine Regel zutrifft.

## Erstellen eines gefilterten Modells

Abbildung 12-18  
Dialogfeld “Neues Modell erzeugen”



Um ein gefiltertes Modell aus einem Assoziationsmodell-Nugget wie einem Regelmengenknoten vom Typ “A Priori”, “CARMA” oder “Sequenz” zu erstellen, wählen Sie im Browser für Modell-Nuggets im Menü “Generieren” die Option Gefiltertes Modell. Dadurch wird ein Untergruppenmodell erstellt, das nur die Regeln enthält, die derzeit im Browser angezeigt werden. *Hinweis:* Sie können keine gefilterten Modelle für nicht verfeinerte Modelle erstellen.

Folgende Optionen für das Filtern von Regeln stehen zur Verfügung:

**Name für neues Modell.** Ermöglicht die Angabe eines Namens des neuen Knotens vom Typ “Gefiltertes Modell”.

**Knoten erstellen auf.** Steuert den Standort des neuen Knotens vom Typ “Gefiltertes Modell”. Wählen Sie Zeichenbereich, Generierte Modelle oder Beides aus.

**Regelnummerierung.** Dient zur Angabe, wie die Regel-IDs in der Regel-Untermenge, die in das gefilterte Modell eingeschlossen ist, nummeriert werden sollen.

- **Ursprüngliche Regel-ID-Nummern beibehalten.** Wählen Sie diese Option aus, um die ursprüngliche Nummerierung der Regeln beizubehalten. Standardmäßig erhalten die Regeln eine ID, die der Reihenfolge Ihrer Entdeckung durch den Algorithmus entspricht. Diese Reihenfolge kann je nach dem verwendeten Algorithmus variieren.
- **Regeln neu nummerieren, beginnend mit.** Wählen Sie diese Option, um den gefilterten Regeln neue Regel-IDs zuzuweisen. Neue IDs werden auf der Grundlage der in der Regel-Browser-Tabelle auf der Registerkarte “Modell” angezeigten Sortierreihenfolge zugewiesen. Die Nummerierung beginnt mit der Zahl, die Sie hier angeben. Sie können die Startnummer für IDs mithilfe der Pfeile auf der rechten Seite angeben.

## Scoring von Assoziationsregeln

Die Scores, die erzielt werden, wenn neue Daten ein Assoziationsregelmodell-Nugget durchlaufen, werden in separaten Feldern ausgegeben. Für jede Prognose werden drei neue Felder hinzugefügt. Dabei steht *P* für die Prognose, *C* für die Konfidenz und *I* für die Regel-ID. Die Organisation dieser Ausgabefelder ist abhängig davon, ob die Eingabedaten im Transaktions- oder im Tabellenformat vorliegen. Einen Überblick über diese Formate finden Sie unter Tabellendaten im Vergleich zu Transaktionsdaten auf S. 389.

Angenommen, Sie nehmen das Scoring von Warenkorbdaten mithilfe eines Modells vor, bei dem Prognosen auf der Grundlage der folgenden drei Regeln erzeugt werden:

Rule\_15 Brot&Wein -> Fleisch (Konfidenz 54 %)  
 Rule\_22 Käse -> Obst (Konfidenz 43 %)  
 Rule\_5 Brot&Käse -> TK-Gemüse (Konfidenz 24 %)

**Tabellendaten.** Bei Tabellendaten werden die drei Prognosen (3 ist der Standardwert) in einem einzelnen Datensatz ausgegeben.

Tabelle 12-1  
 Scores im Tabellenformat

ID	Brot	Wein	Käse	P1	C1	I1	P2	C2	I2	P3	C3	I3
Fred	1	1	1	Fleisch	0.54	15	Obst	0.43	22	TK-Gemüse	0.24	5

**Transaktionsdaten.** Bei Transaktionsdaten wird je ein separater Datensatz für die einzelnen Prognosen erzeugt. Die Prognosen werden ebenfalls in separaten Spalten hinzugefügt, die Scores werden jedoch so ausgegeben, wie sie berechnet werden. Dies führt zu Datensätzen mit unvollständigen Prognosen (vgl. das nachstehende Ausgabebeispiel). Die zweite und dritte Prognose (P2 und P3) im ersten Datensatz sind leer, ebenso wie die zugehörigen Konfidenzen und Regel-IDs. Wenn die Scores ausgegeben werden, enthält der endgültige Datensatz jedoch alle drei Prognosen.

Tabelle 12-2  
Scores im Transaktionsformat

ID	Item	P1	C1	I1	P2	C2	I2	P3	C3	I3
Fred	Brot	Fleisch	0.54	14	\$null\$	\$null\$	\$null\$	\$null\$	\$null\$	\$null\$
Fred	cheese	Fleisch	0.54	14	Obst	0.43	22	\$null\$	\$null\$	\$null\$
Fred	Wein	Fleisch	0.54	14	Obst	0.43	22	TK-Gemüse	0.24	5

Um für die Berichterstellung oder Bereitstellung nur vollständige Prognosen aufzunehmen, wählen Sie die vollständigen Datensätze mithilfe eines Auswahlknotens aus.

*Hinweis:* Die in diesen Beispielen verwendeten Feldnamen sind zugunsten größerer Klarheit abgekürzt. Während der tatsächlichen Verwendung werden die Ergebnisfelder für Assoziationsmodelle wie folgt benannt:

Neues Feld	Beispiel für Feldname
Prognose	$\$A-TRANSAKTIONSNUMMER-1$
Konfidenz (oder anderes Kriterium)	$\$AC-TRANSAKTIONSNUMMER-1$
Regel-ID	$\$A-Regel\_ID-1$

### Regeln mit mehreren Sukzedenzen

Der CARMA-Algorithmus lässt Regeln mit mehreren Sukzedenzen zu. Beispiel:

Brot -> Wein&Käse

Bei der Bewertung derartiger "doppelköpfiger" Regeln werden die Vorhersagen in dem in der folgenden Tabelle angezeigten Format zurückgegeben:

Tabelle 12-3  
Scoring-Ergebnisse einschließlich einer Prognose mit mehreren Sukzedenzen

ID	Brot	Wein	Käse	P1	C1	I1	P2	C2	I2	P3	C3	I3
Fred	1	1	1	Fleisch&Gemüse	0.54	16	Obst	0.43	22	TK-Gemüse	0.24	5

In einigen Fällen müssen Sie derartige Scores vor der Bereitstellung aufteilen. Zur Aufteilung einer Prognose mit mehreren Sukzedenzen müssen Sie das Feld mithilfe der CLEM-String-Funktionen analysieren. [Für weitere Informationen siehe Thema String-Funktionen in Kapitel 8 in IBM SPSS Modeler 14.2- Benutzerhandbuch.](#)

## Bereitstellung von Assoziationsmodellen

Beim Scoring von Assoziationsmodellen werden die Prognosen und die Konfidenzen in separaten Spalten ausgegeben. (*P* steht dabei für die Prognose, *K* für die Konfidenz und *I* für die Regel-ID.) Dies gilt für Eingabedaten sowohl im Tabellenformat als auch im Transaktionsformat. [Für weitere Informationen siehe Thema Scoring von Assoziationsregeln auf S. 412.](#)

Abbildung 12-19  
Scores im Tabellenformat mit Prognosen in Spalten

	ID	A	B	C	P1	C1	I1	P2	C2	I2	P3	C3	I3
1	Tom	1	1	1	D	9	1	E	5	23	F	3	9
2	Bob	0	1	1	F	3	9	E	2	15	D	1	4

Bei der Vorbereitung von Scores für die Bereitstellung stellt sich ggf. heraus, dass die Ausgabedaten in ein Format transponiert werden müssen, bei dem die Prognosen nicht in Spalten ausgegeben werden, sondern in Zeilen (je eine Prognose pro Zeile, auch als "Kassenrollenformat" bezeichnet).

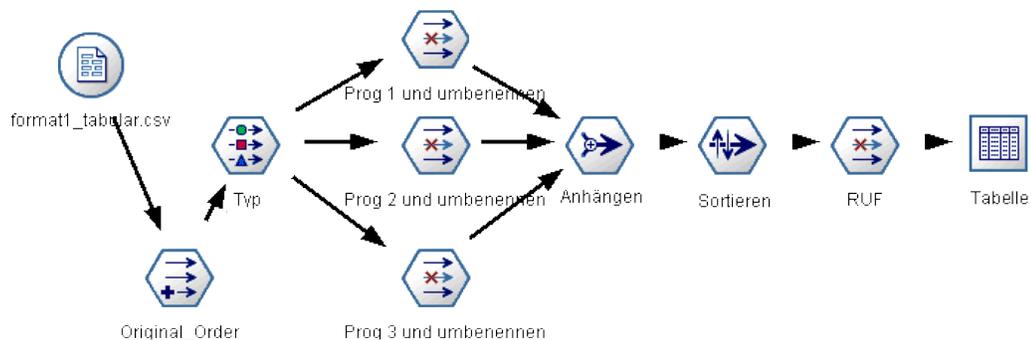
Abbildung 12-20  
Transponierte Scores mit Prognosen in Zeilen

	ID	A	B	C	Prog	Krit	Regel	ID
1	Tom	1	1	1	D	9	1	1
2	Tom	1	1	1	E	5	23	3
3	Tom	1	1	1	F	3	9	9
4	Bob	0	1	1	F	3	9	9
5	Bob	0	1	1	.	\$null\$	\$null\$	
6	Bob	0	1	1	.	\$null\$	\$null\$	

### Transponieren von Scores im Tabellenformat

Sie können Scores im Tabellenformat mithilfe einer Reihe von Schritten in IBM® SPSS® Modeler von Spalten in Zeilen transponieren (siehe nachfolgende Schritte).

Abbildung 12-21  
Beispiel-Stream für die Transposition von Tabellendaten in das Kassenrollenformat



- Überprüfen Sie die gegenwärtige Reihenfolge der Prognosen mithilfe der Funktion @INDEX in einem Ableitungsknoten und speichern Sie diesen Indikator in einem neuen Feld, wie beispielsweise *Original\_order*.

- ▶ Fügen Sie einen Typknoten hinzu, um sicherzustellen, dass alle Felder instanziiert sind.
- ▶ Mit einem Filterknoten können Sie die Standardfelder für Prognose, Konfidenz und ID (*P1*, *C1*, *I1*) in gewöhnliche Felder umbenennen, wie beispielsweise *Prog*, *Krit* und *Regel-ID*, die später an die Datensätze angehängt werden. Für jede generierte Prognose wird jeweils ein Filterknoten benötigt.

Abbildung 12-22

Filtern der Felder für die Prognosen 1 und 3 bei Umbenennung der Felder für Prognose 2.



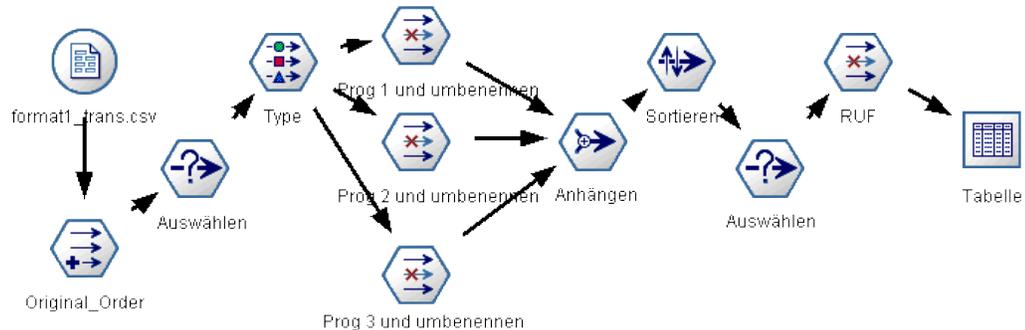
- ▶ Verwenden Sie einen Anhangknoten, um Scores für die freigegebenen Elemente *Prog*, *Krit* und *Regel-ID* anzuhängen.
- ▶ Hängen Sie einen Sortierknoten an, um die Datensätze für das Feld *Original\_order* in aufsteigender Reihenfolge und für *Crit* in absteigender Reihenfolge zu sortieren. Bei "Crit" handelt es sich um das Feld, das zum Sortieren der Vorhersagen nach Kriterien wie Konfidenz, Lift und Unterstützung verwendet wird.
- ▶ Filtern Sie mithilfe eines weiteren Filterknotens das Feld *Original\_order* aus der Ausgabe. Nun können die Daten bereitgestellt werden.

### **Transponieren von Scores im Transaktionsformat**

Für die Transposition von Transaktions-Scores wird ein ähnlicher Vorgang verwendet. Beispiel: Der unten dargestellte Stream transponiert Scores in ein Format, bei dem jede Zeile eine einzelne Prognose erhält, wie für die Bereitstellung erforderlich.

Abbildung 12-23

Beispiel-Stream für die Transposition von Transaktionsdaten in das Kassenrollenformat



Abgesehen von zwei weiteren Auswahlknoten ist der Vorgang mit dem zuvor für Tabellendaten beschriebenen identisch.

- Der erste Auswahlknoten wird verwendet, um Regel-IDs in angrenzenden Datensätzen zu vergleichen und nur eindeutige nicht definierte Datensätze aufzunehmen. Bei diesem Auswahlknoten werden CLEM-Ausdrücke zur Auswahl der Datensätze verwendet:  $ID \neq @OFFSET(ID,-1)$  or  $@OFFSET(ID,-1) = undef$ .
- Der zweite Auswahlknoten wird verwendet, um überflüssige Regeln oder Regeln, bei denen Rule\_ID einen Nullwert aufweist, zu verwerfen. Bei diesem Auswahlknoten werden folgende CLEM-Ausdrücke zum Verwerfen von Datensätzen verwendet:  $not(@NULL(Rule\_ID))$ .

Weitere Informationen zum Transponieren von Scores für die Bereitstellung erhalten Sie beim Technischen Support von .

## Sequenzknoten

Der Sequenzknoten erkennt Muster in sequenziellen oder zeitorientierten Daten, und zwar im Format **Brot -> Käse**. Die Elemente einer Sequenz sind **Element-Sets**, die eine einzelne Transaktion ausmachen. Beispiel: Wenn eine Person in den Supermarkt geht und Brot und Milch kauft und dann ein paar Tage später zurückkehrt und Käse kauft, kann das Kaufverhalten dieser Person als zwei Element-Sets dargestellt werden. Der erste Element-Set enthält Brot und Milch, der zweite Käse. Eine **Sequenz** ist eine Liste mit Element-Sets, die in einer vorhersagbaren Reihenfolge auftreten. Der Sequenzknoten erkennt häufige Sequenzen und erstellt einen generierten Modellknoten, der für Vorhersagen verwendet werden kann.

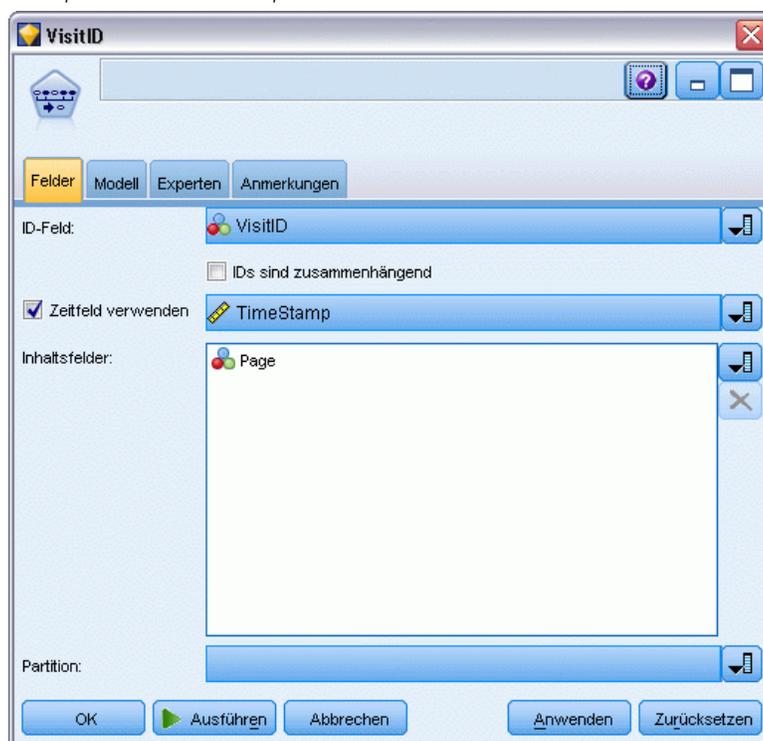
**Anforderungen.** Um eine Sequenz-Regelmengung zu erstellen, müssen Sie ein ID-Feld, ein optionales Zeitfeld und mindestens ein Inhaltsfeld angeben. Beachten Sie, dass diese Einstellungen auf der Registerkarte "Felder" des Modellierungsknotens vorgenommen werden müssen. Sie können nicht aus einem aufwärts liegenden Typknoten gelesen werden. Das ID-Feld kann eine beliebige Rolle oder ein beliebiges Messniveau aufweisen. Wenn Sie ein Zeitfeld angeben, kann es jede beliebige Rolle aufweisen, muss jedoch in numerischem, Datums-, Uhrzeit- oder Zeitstempelformat gespeichert werden. Wenn Sie kein Zeitfeld angeben, verwendet der Sequenzknoten einen implizierten Zeitstempel, wobei als Zeitwerte Zeilennummern verwendet werden. Inhaltsfelder können ein beliebiges Messniveau und eine beliebige Rolle aufweisen, sämtliche Inhaltsfelder

müssen jedoch vom selben Typ sein. Wenn es sich um numerische Felder handelt, müssen es Bereiche ganzer Zahlen (keine reellen Bereiche) sein.

**Stärken.** Der Sequenzknoten basiert auf dem CARMA-Assoziationsregelalgorithmus, der eine effiziente bidirektionale Methode zum Suchen von Sequenzen verwendet. Außerdem kann der von einem Sequenzknoten generierte Modellknoten in einen Daten-Stream eingefügt werden, um Vorhersagen zu erstellen. Der generierte Modellknoten kann auch Superknoten zum Erkennen und Zählen spezifischer Sequenzen und zum Erstellen von Vorhersagen basierend auf bestimmten Sequenzen generieren.

### **Feldoptionen für den Sequenzknoten**

Abbildung 12-24  
Feldoptionen für den Sequenzknoten



Vor der Ausführung eines Sequenzknotens müssen Sie ID- und Inhaltsfelder auf der Registerkarte “Felder” des Sequenzknotens angeben. Wenn Sie ein Zeitfeld verwenden möchten, müssen Sie auch dieses hier angeben.

**ID-Feld.** Wählen Sie ein ID-Feld aus der Liste aus. Als ID-Feld können numerische oder symbolische Felder verwendet werden. Jeder eindeutige Wert in diesem Feld sollte eine bestimmte Analyseeinheit darstellen. Bei einer Warenkorbanwendung könnte z. B. jede ID einen einzelnen

Kunden darstellen. Für eine Webprotokoll-Analyseanwendung könnte jede ID einen Computer (nach IP-Adresse) oder einen Benutzer (nach Anmelddaten) darstellen.

- **IDs sind zusammenhängend.** Wenn Ihre Daten vorsortiert sind, sodass alle Datensätze mit derselben ID im Daten-Stream zusammengefasst sind, wählen Sie diese Option, um die Verarbeitung zu beschleunigen. Wenn Ihre Daten nicht vorsortiert sind (oder Sie nicht sicher sind), lassen Sie diese Option deaktiviert. Die Daten werden dann vom Sequenzknoten automatisch sortiert.

*Hinweis:* Wenn Ihre Daten nicht sortiert werden und Sie diese Option auswählen, könnten Sie ungültige Ergebnisse in Ihrem Sequenzmodell erhalten.

**Zeitfeld.** Wenn Sie ein Feld in den Daten verwenden möchten, um die Uhrzeiten von Ereignissen anzugeben, wählen Sie Zeitfeld verwenden und geben Sie das zu verwendende Feld an. Das Zeitfeld muss numerische, Datums-, Uhrzeit- oder Zeitstempelwerte enthalten. Wenn kein Zeitfeld angegeben ist, wird davon ausgegangen, dass die Datensätze in sequenzieller Reihenfolge aus der Datenquelle ankommen und Datensatznummern als Zeitwerte verwendet werden (der erste Datensatz kommt zur Uhrzeit "1" vor; der zweite zur Uhrzeit "2" usw.).

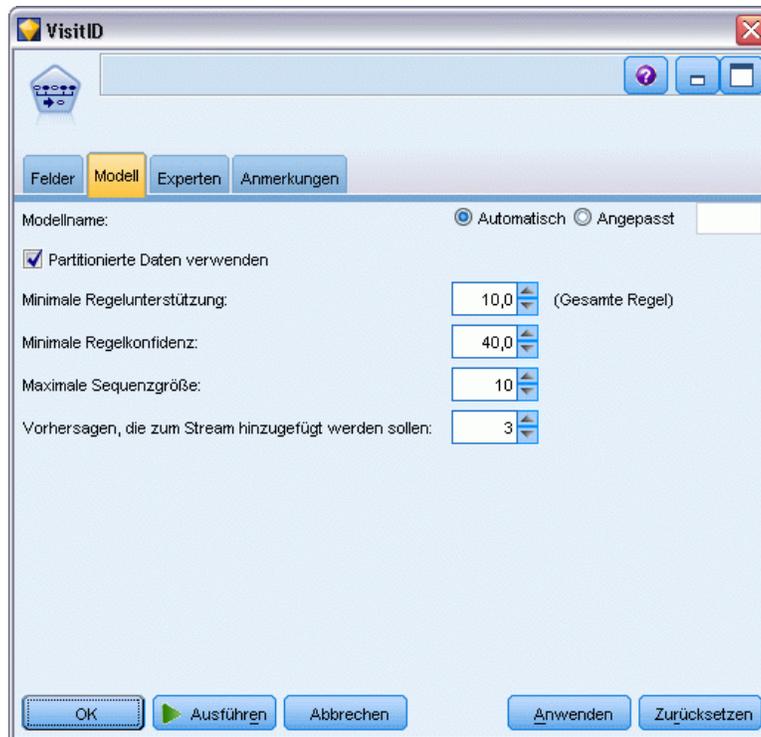
**Inhaltsfelder.** Geben Sie die Inhaltsfelder für das Modell an. Diese Felder enthalten die in der Sequenzmodellierung interessanten Ereignisse.

Der Sequenzknoten kann Daten im Tabellenformat und im Transaktionsformat verarbeiten. Wenn Sie mehrere Felder für Transaktionsdaten verwenden, stellen die in diesen Feldern für einen bestimmten Datensatz angegebenen Elemente solche Elemente dar, die in einer Einzeltransaktion mit einem einzelnen Zeitstempel gefunden wurden. [Für weitere Informationen siehe Thema Tabellendaten im Vergleich zu Transaktionsdaten auf S. 389.](#)

**Partition.** In diesem Feld können Sie ein Feld angeben, das verwendet wird, um die Daten in getrennte Stichproben für die Trainings-, Test- und Validierungsphase der Modellbildung aufzuteilen. Indem Sie mit einer Stichprobe das Modell erstellen und es mit einer anderen Stichprobe testen, erhalten Sie einen guten Hinweis dafür, wie gut das Modell sich für größere Datenmengen generalisieren lässt, die den aktuellen Daten ähneln. Wenn mehrere Partitionsfelder mithilfe von Typ- oder Partitionsknoten erstellt wurden, muss in jedem Modellierungsknoten, der die Partitionierung verwendet, auf der Registerkarte "Felder" ein einzelnes Partitionsfeld ausgewählt werden. (Wenn nur eine einzige Partition vorhanden ist, wird diese immer automatisch verwendet, wenn die Partitionierung aktiviert ist.) [Für weitere Informationen siehe Thema Partitionsknoten in Kapitel 4 in IBM SPSS Modeler 14.2- Quellen-, Prozess- und Ausgabeknoten.](#) Beachten Sie außerdem, dass die Partitionierung auch auf der Registerkarte "Modelloptionen" des Knotens aktiviert sein muss, wenn die ausgewählte Partitionierung in Ihrer Analyse angewendet werden soll. (Wenn die Auswahl der Option aufgehoben ist, kann die Partitionierung ohne Änderung der Feldeinstellungen deaktiviert werden.)

## Modelloptionen für den Sequenzknoten

Abbildung 12-25  
Modelloptionen für den Sequenzknoten



**Modellname.** Sie können den Modellnamen automatisch basierend auf den Ziel- oder ID-Feldnamen (oder dem Modelltyp in Fällen, in denen kein solches Feld angegeben ist) generieren oder einen benutzerdefinierten Namen eingeben.

**Partitionierte Daten verwenden.** Wenn ein Partitionsfeld definiert ist, gewährleistet diese Option, dass nur Daten aus der Trainingspartition für die Modellerstellung verwendet werden. [Für weitere Informationen siehe Thema Partitionsknoten in Kapitel 4 in IBM SPSS Modeler 14.2- Quellen-, Prozess- und Ausgabeknoten.](#)

**Minimale Regelunterstützung (%).** Sie können auch ein Stützkriterium angeben.

**Regelunterstützung** bezieht sich auf den Anteil von IDs in den Trainingsdaten, die die gesamte Sequenz enthalten. Wenn Sie weitere gemeinsame Sequenzen wünschen, erhöhen Sie diese Einstellung.

**Minimale Regelkonfidenz (%).** Sie können ein Konfidenzkriterium angeben, um die Sequenzen im Sequenz-Set beizubehalten. **Konfidenz** bezieht sich auf den prozentualen Anteil von IDs, für die eine richtige Vorhersage gemacht wird, aus allen IDs, für die die Regel eine Vorhersage macht. Sie wird aus der Anzahl von IDs berechnet, für die die gesamte Sequenz gefunden wird, dividiert durch die Anzahl der IDs, für die die Antezedenzen gefunden werden, basierend auf den Trainingsdaten. Sequenzen mit einer niedrigeren Konfidenz als das angegebene Kriterium werden verworfen. Wenn Sie zu viele Sequenzen oder uninteressante Sequenzen erhalten, sollten Sie diese Einstellung erhöhen. Wenn Sie zu wenige Sequenzen erhalten, sollten Sie diese Einstellung reduzieren.

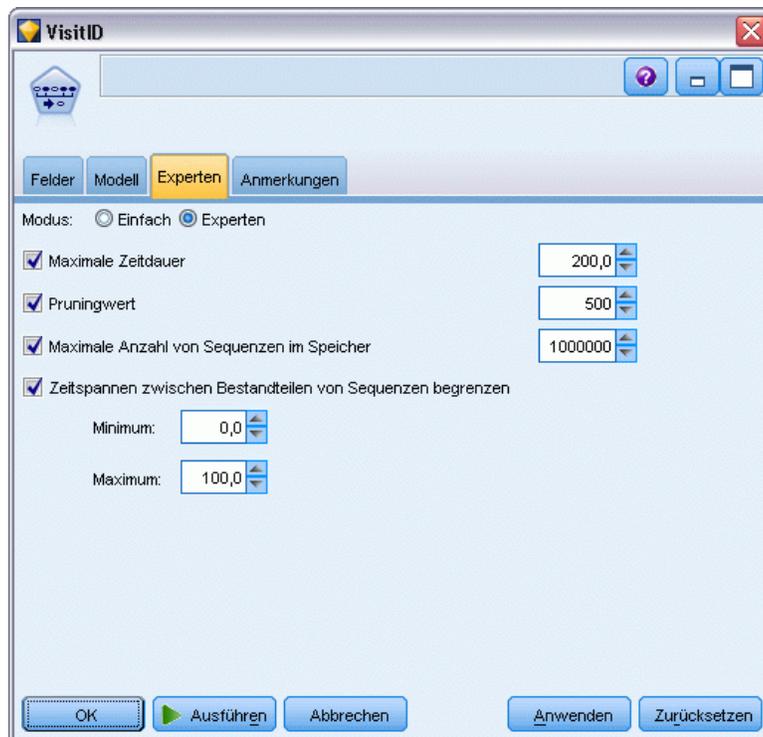
**Maximale Sequenzgröße.** Sie können die maximale Anzahl unterschiedlicher *Element-Sets* (im Gegensatz zu *Elementen*) in einer Sequenz festlegen. Wenn die gewünschten Sequenzen relativ kurz sind, können Sie diese Einstellung reduzieren, um die Erstellung des Sequenz-Sets zu beschleunigen.

**Vorhersagen, die zum Stream hinzugefügt werden sollen.** Geben Sie die Anzahl der Vorhersagen an, die dem Stream vom resultierenden generierten Modellknoten hinzugefügt werden sollen. [Für weitere Informationen siehe Thema Sequenzmodell-Nuggets auf S. 422.](#)

### Expertenoptionen für den Sequenzknoten

Für Personen mit umfassenden Kenntnissen über die Funktionsweise des Sequenzknotens ermöglichen die folgenden Expertenoptionen die Feinabstimmung des Modellerstellungsvorgangs. Für den Zugriff auf die Expertenoptionen stellen Sie den Modus auf der Registerkarte "Experten" auf Experten ein.

Abbildung 12-26  
Expertenoptionen für den Sequenzknoten



**Maximale Zeitdauer.** Falls diese Option ausgewählt ist, werden die Sequenzen auf solche mit einer Dauer (die Zeit zwischen dem ersten und letzten Element-Set) beschränkt, die dem angegebenen Wert entspricht oder darunter liegt. Wenn Sie kein Zeitfeld angegeben haben, wird die Zeitdauer in den Rohdaten in Zeilen (Datensätzen) ausgedrückt. Wenn es sich bei dem verwendeten Zeitfeld um ein Uhrzeit-, Datums- oder Zeitstempel-Feld handelt, wird die Zeitdauer in Sekunden ausgedrückt. Im Falle von numerischen Feldern wird die Zeitdauer in denselben Einheiten ausgedrückt wie das Feld selbst.

**Reduzierungswert festlegen.** Der im Sequenzknoten verwendete CARMA-Algorithmus entfernt (**reduziert**) während der Verarbeitung seltene Element-Sets aus seiner Liste potenzieller Element-Sets. Wählen Sie diese Option, um die Reduktionshäufigkeit anzupassen. Die angegebene Zahl legt die Reduktionshäufigkeit fest. Geben Sie einen kleinen Wert ein, um die Speicheranforderungen des Algorithmus zu reduzieren (möglicherweise aber die erforderliche Trainingszeit zu erhöhen), oder geben Sie einen hohen Wert ein, um die Trainingsgeschwindigkeit zu erhöhen (möglicherweise aber die Speicheranforderungen zu erhöhen).

**Maximale Anzahl von Sequenzen im Speicher.** Falls diese Option ausgewählt ist, beschränkt der CARMA-Algorithmus seinen Speicher mit möglichen Sequenzen auf die Anzahl der angegebenen Sequenzen. Wählen Sie diese Option, wenn IBM® SPSS® Modeler während der Erstellung von Sequenzmodellen zu viel Speicher belegt. Beachten Sie, dass der von Ihnen hier angegebene maximale Sequenzwert der Anzahl der möglichen Sequenzen entspricht, die intern bei der Modellerstellung aufgezeichnet wird. Diese Zahl sollte viel größer sein als die Zahl der Sequenzen, die Sie im endgültigen Modell erwarten.

**Zeitspannen zwischen Bestandteilen von Sequenzen begrenzen.** Mit dieser Option können Sie Beschränkungen der Zeitspannen festlegen, die zwischen verschiedenen Element-Sets liegen. Falls ausgewählt, werden die Element-Sets mit Zeitspannen, die unter der angegebenen Minimum gap oder über der angegebenen Maximalen Zeitspanne liegen, nicht als Teil einer Sequenz betrachtet. Verwenden Sie diese Option, um zu vermeiden, dass Sequenzen gezählt werden, die lange Zeitintervalle enthalten oder die in einer sehr kurzen Zeitspanne auftreten.

*Hinweis:* Wenn es sich bei dem verwendeten Zeitfeld um ein Uhrzeit-, Datums- oder Zeitstempel-Feld handelt, wird die Zeitdauer in Sekunden ausgedrückt. Im Falle von numerischen Feldern wird die Zeitdauer in denselben Einheiten ausgedrückt wie das Zeitfeld.

Betrachten wir z. B. diese Liste mit Transaktionen:

ID	Time	Inhalt
1001	1	apples
1001	2	Brot
1001	5	cheese
1001	6	dressing

Wenn Sie mit diesen Daten ein Modell bilden, wobei der minimale Abstand auf 2 gesetzt ist, erhalten Sie die folgenden Sequenzen:

Äpfel -> Käse

Äpfel -> Dressing

Brot -> Käse

Brot -> Dressing

Sequenzen, wie apples -> bread erscheinen nicht, weil der Abstand zwischen apples und bread kleiner als der Mindestabstand ist. Ähnlich verhält es sich, wenn die Daten stattdessen wie folgt lauten:

ID	Time	Inhalt
1001	1	apples
1001	2	Brot
1001	5	cheese
1001	20	dressing

und der maximale Abstand auf 10 gesetzt ist; dann erhalten Sie keine Sequenzen mit dressing, weil der Abstand zwischen cheese und dressing zu groß ist, um als Teil derselben Sequenz betrachtet zu werden.

## Sequenzmodell-Nuggets

Sequenzmodell-Nuggets stellen die Sequenzen dar, die in einem bestimmten, vom Sequenzknoten ermittelten Ausgabefeld gefunden wurden und zum Erstellen von Prognosen einem Stream hinzugefügt werden können.

Bei der Ausführung eines Streams, der einen Sequenzknoten enthält, fügt der Knoten ein Felderpaar hinzu, das die Prognosen und die zugeordneten Konfidenzwerte für die einzelnen Prognosen aus dem Sequenzmodell den Daten hinzufügt. Standardmäßig werden drei Felderpaare mit den drei Top-Prognosen (und den zugehörigen Konfidenzwerten) hinzugefügt. Sie können die Anzahl der generierten Prognosen ändern, wenn Sie das Modell durch Festlegung der Sequenzknoten-Modelloptionen zum Erstellungszeitpunkt erstellen, oder auch auf der Registerkarte "Einstellungen", nachdem das Modell-Nugget einem Stream hinzugefügt wurde. [Für weitere Informationen siehe Thema Sequenzmodell-Nugget - Einstellungen auf S. 427.](#)

Die neuen Feldnamen werden aus dem Modellnamen abgeleitet. Die Feldnamen lauten  $SS$ -sequence- $n$  für das Prognosefeld (dabei gibt  $n$  die  $n$ -te Prognose an) und  $SC$ -sequence- $n$  für das Konfidenzfeld. In einem Stream mit mehreren Sequenzregelknoten in einer Reihe enthalten die neuen Feldnamen Zahlen im Präfix, damit sie auseinander gehalten werden können. Beim ersten Sequenz-Set-Knoten im Stream werden die üblichen Namen verwendet, beim zweiten Knoten Namen, die mit  $SS1$ - und  $SC1$ - beginnen, beim dritten Knoten Namen mit  $SS2$ - und  $SC2$ - usw. Die Vorhersagen werden nach Konfidenz geordnet angezeigt, sodass  $SS$ -sequence-1 die Prognose mit der höchsten Konfidenz enthält,  $SS$ -sequence-2 die Prognose mit der zweithöchsten Konfidenz usw. Bei Datensätzen, bei denen die Anzahl der verfügbaren Prognosen kleiner ist als die Anzahl der angeforderten Prognosen, enthalten die restlichen Prognosen den Wert  $\$null$ . Beispiel: Wenn nur zwei Vorhersagen für einen bestimmten Datensatz vorgenommen werden können, weisen  $SS$ -sequence-3 und  $SC$ -sequence-3 den Wert  $\$null$  auf.

Bei jedem Datensatz werden die Regeln im Modell mit der Menge der Transaktionen verglichen, die bisher für die aktuelle ID verarbeitet wurden, einschließlich des aktuellen Datensatzes und aller vorangegangenen Datensätze mit derselben ID und früheren Zeitstempeln. Die  $k$  Regeln mit den höchsten Konfidenzwerten, die für dieses Set von Transaktionen gelten, werden verwendet, um die  $k$  Prognosen für den Datensatz zu generieren. Dabei ist  $k$  die Anzahl der Prognosen, die nach dem Hinzufügen des Modells zum Stream auf der Registerkarte "Einstellungen" angegeben wurden. (Wenn mehrere Regeln dasselbe Ergebnis für das

Transaktions-Set vorhersagen, wird nur die Regel mit der höchsten Konfidenz verwendet.) [Für weitere Informationen siehe Thema Sequenzmodell-Nugget - Einstellungen auf S. 427.](#)

Wie bei anderen Arten von Assoziationsregelmodellen muss das Datenformat mit dem Format übereinstimmen, das beim Aufbau des Sequenzmodells verwendet wurde. Mit Modellen, die mithilfe von Tabellendaten erstellt wurden, können entsprechend nur Tabellendaten gesortiert werden. [Für weitere Informationen siehe Thema Scoring von Assoziationsregeln auf S. 412.](#)

*Hinweis:* Beim Scoring von Daten mithilfe eines generierten Sequenz-Set-Knotens in einem Stream werden alle Toleranz- oder Lückeneinstellungen, die beim Erstellen des Modells ausgewählt wurden, beim Scoring ignoriert.

### **Vorhersagen aus Sequenzregeln**

Der Knoten bearbeitet die Datensätze in zeitabhängiger Weise (bzw. in Abhängigkeit von der Reihenfolge, wenn beim Erstellen des Modells kein Zeitstempelfeld verwendet wurde). Die Datensätze sollten nach dem ID-Feld und dem Zeitstempelfeld (sofern vorhanden) sortiert werden. Die Prognosen sind jedoch nicht an den Zeitstempel des Datensatzes gebunden, dem sie hinzugefügt werden. Sie beziehen sich einfach auf die Elemente, die unter Berücksichtigung des Transaktionsverlaufs für die aktuelle ID bis zum aktuellen Datensatz mit der größten Wahrscheinlichkeit *irgendwann in der Zukunft* auftreten.

Beachten Sie, dass die Vorhersagen für die einzelnen Datensätze nicht unbedingt von den Transaktionen des betreffenden Datensatzes abhängen. Wenn die Transaktionen des aktuellen Datensatzes keine spezifische Regel auslösen, werden die Regeln anhand der vorangegangenen Transaktionen für die aktuelle ID ausgewählt. Anders ausgedrückt: Wenn der aktuelle Datensatz keine verwertbaren Prognoseinformationen zur Sequenz hinzufügt, wird die Prognose aus der letzten nützlichen Transaktion für diese ID auf den aktuellen Datensatz übertragen.

Beispiel: Angenommen Sie haben ein Sequenzmodell mit nur einer einzigen Regel:

Marmelade -> Brot (0,66)

und Sie leiten folgende Datensätze an diese Regel weiter:

<b>ID</b>	<b>Kauf</b>	<b>Prognose</b>
001	jam	bread
001	milk	bread

Der erste Datensatz generiert, wie zu erwarten, eine Vorhersage für *Brot*. Der zweite Datensatz enthält ebenfalls eine Vorhersage für *Brot*, da keine Regel für *Marmelade* gefolgt von *Milch* vorliegt; daher fügt die Transaktion *Milch* keine verwertbaren Informationen hinzu und die Regel *Jam -> Bread* gilt weiterhin.

### **Erzeugen neuer Knoten**

Im Menü "Generieren" können Sie anhand des Sequenzmodells neue Superknoten erstellen.

- **Regel-Superknoten.** Erstellt einen Superknoten, der die Vorkommnisse von Sequenzen in den gescorten Daten ermitteln und zählen kann. Diese Option ist deaktiviert, wenn keine Regel ausgewählt wurde. [Für weitere Informationen siehe Thema Generieren eines Regel-Superknotens aus einem Sequenzmodell-Nugget auf S. 428.](#)
- **Modell zur Palette hinzufügen.** Gibt das Modell an die Modellpalette zurück. Das ist nützlich, wenn Sie von einem Kollegen einen Datenstrom, der das Modell enthält, jedoch nicht das Modell selbst erhalten.

### Nähere Informationen zum Sequenzmodell-Nugget

Auf der Registerkarte “Modell” eines Sequenzmodell-Nuggets werden die Regeln angezeigt, die durch den Algorithmus extrahiert wurden. Jede Zeile in der Tabelle steht für eine Regel, bei der das Antezedens (der “Wenn”-Teil der Regel) in der ersten Spalte jeweils vom Sukzedens (dem “Dann”-Teil der Regel) in der zweiten Spalte gefolgt wird.

Abbildung 12-27  
Sequenz-Nugget – Registerkarte “Modell”

Antezedens	Sukzedens	Unterstützung %	Konfidenz %
login.asp?	personal.asp	21,591	100,0
login.asp?	login.asp?	21,591	100,0
login.asp?	personal.asp	21,591	100,0
login.asp?	personal.asp	21,591	100,0
splash.htm	main.htm	70,455	83,871
splash.htm	login.asp	59,091	51,923
main.htm	login.asp	70,455	43,548
splash.htm	login.asp	76,136	41,791
main.htm	main.htm	76,136	37,313
splash.htm	main.htm	59,091	36,538
main.htm	main.htm	59,091	36,538

Die einzelnen Regeln werden in folgendem Format angezeigt:

Antezedens	Sukzedens
beer and cannedveg	beer
fish fish	fish

Die erste Beispielregel wird wie folgt interpretiert: *Bei IDs, bei denen „Bier“ und „Dosengemüse“ in derselben Transaktion vorkam, wird mit hoher Wahrscheinlichkeit ein Vorkommen von „Bier“ folgen.* Die zweite Beispielregel kann wie folgt interpretiert werden: *Bei IDs, bei denen „Fisch“ in einer Transaktion vorkam und ebenfalls „Fisch“ in einer anderen, wird mit hoher Wahrscheinlichkeit „Fisch“ ein weiteres Mal vorkommen.* In der ersten Regel werden *Bier* und *Dosengemüse* gleichzeitig eingekauft; in der zweiten Regel wird *Fisch* in zwei separaten Transaktionen erworben.

**Menü „Sortieren“.** Mit der Schaltfläche des Menüs „Sortieren“ in der Symbolleiste wird die Sortierung der Regeln gesteuert. Die Sortierrichtung (aufsteigend oder absteigend) lässt sich mit der Schaltfläche für die Sortierrichtung (nach unten bzw. oben zeigender Pfeil) ändern.

Abbildung 12-28  
Symbolleistenoptionen für das Sortieren

Sortieren nach:

Regeln können nach folgenden Faktoren sortiert werden:

- Unterstützung %
- Konfidenz %
- Regelunterstützung %
- Sukzedens
- Erstes Antezedens
- Letztes Antezedens
- Anzahl der Elemente (Antezedenzen)

Beispiel: Die nachstehende Tabelle wird in absteigender Reihenfolge nach der Anzahl der Elemente sortiert. Regeln mit mehreren Elementen im Antezedens-Satz haben Vorrang vor Regeln mit weniger Elementen.

Antezedens	Sukzedens
beer and cannedveg and frozenmeal	frozenmeal
beer and cannedveg	beer
fish fish	fish
softdrink	softdrink

**Menü “Anzeigen/Ausblenden” für Kriterien.** Das Menü “Anzeigen/Ausblenden” der Schaltfläche “Kriterien” (Rastersymbol) steuert die Optionen für die Anzeige der Regeln. Die folgenden Anzeigoptionen stehen zur Verfügung:

- **Instanzen** zeigt Informationen über die Anzahl der eindeutigen IDs an, für die die *vollständige Sequenz* vorkommt (also sowohl Antezedenzen als auch Sukzedenzen). (Hier besteht ein Unterschied zu den Assoziationsmodellen, bei denen die Anzahl der Instanzen die Anzahl der IDs bezeichnet, bei denen *ausschließlich* die Antezedenzen gelten.) Bei der Regel *bread -> cheese* beispielsweise wird die Anzahl der IDs in den Trainingsdaten, die sowohl *Brot* als auch *Käse* enthalten, als **Instanzen** bezeichnet.
- **Unterstützung** zeigt den Anteil von IDs in den Trainingsdaten, für den die Antezedenzen wahr sind. Beispiel: Wenn 50 % der Trainingsdaten das Antezedens *Brot* enthalten, hat die Regel *bread -> cheese* eine Antezedens-Unterstützung von 50 %. (Im Gegensatz zu den Assoziationsmodellen beruht die Unterstützung *nicht* auf der Anzahl der Instanzen, wie zuvor beschrieben.)
- **Konfidenz** zeigt den prozentualen Anteil von IDs, für die eine richtige Vorhersage gemacht wird, aus allen IDs, für die die Regel eine Vorhersage macht. Sie wird aus der Anzahl von IDs berechnet, für die die gesamte Sequenz gefunden wird, dividiert durch die Anzahl der IDs, für die die Antezedenzen gefunden werden, basierend auf den Trainingsdaten. Beispiel: Wenn 50 % der Trainingsdaten *cannedveg* enthalten (Antezedens-Unterstützung), jedoch nur 20 % sowohl *cannedveg* als auch *frozenmeal*, ist die Konfidenz für die Regel *cannedveg -> frozenmeal*  $\text{Rule Support} / \text{Antecedent Support}$ , in diesem Fall 40 %.
- **Regelunterstützung** für Sequenzmodelle beruht auf den Instanzen und zeigt den Teil der Trainings-Datensätze an, für die die gesamte Regel, die Antezedenzen und das Sukzedens (bzw. die Sukzedenzen) wahr sind. Beispiel: Wenn 20 % der Trainingsdaten sowohl *Brot* als auch *Käse* enthalten, beträgt die Regelunterstützung für die Regel *bread -> cheese* 20 %.

Beachten Sie, dass die Anteile auf gültigen Transaktionen beruhen (Transaktionen mit mindestens einem beobachteten Element oder Wahr-Wert) und nicht auf der Gesamtzahl der Transaktionen. Ungültige Transaktionen, d. h. Transaktionen ohne Elemente oder Wahr-Werte, werden für diese Berechnungen verworfen.

**Schaltfläche “Filter”.** Mit der Filterschaltfläche (Trichtersymbol) im Menü wird der untere Teil des Dialogfelds erweitert und ein Fenster mit aktiven Regelfiltern wird angezeigt. Filter werden verwendet, um die Anzahl der auf der Registerkarte “Modelle” angezeigten Regeln einzugrenzen.

Abbildung 12-29  
Schaltfläche “Filter”



Zum Erstellen von Filtern klicken Sie auf das Filtersymbol rechts neben dem erweiterten Fenster. Dadurch wird ein separates Dialogfeld geöffnet, in dem Sie Bedingungen für die Anzeige von Regeln eingeben können. Beachten Sie, dass die Filterschaltfläche häufig in Verbindung mit dem Menü “Generieren” verwendet wird, um zunächst die Regeln zu filtern und anschließend ein Modell zu erstellen, das die betreffende Untergruppe der Regeln enthält. Weitere Informationen finden Sie unter [Angeben von Filtern für Regeln](#).

## Sequenzmodell-Nugget - Einstellungen

Auf der Registerkarte “Einstellungen” eines Sequenzmodell-Nuggets werden die Scoring-Optionen für das Modell angezeigt. Diese Registerkarte ist erst dann verfügbar, nachdem das Modell zum Zwecke des Scoring zum Stream-Zeichenbereich hinzugefügt wurde.

Abbildung 12-30  
Sequenz-Nugget – Registerkarte “Einstellungen”

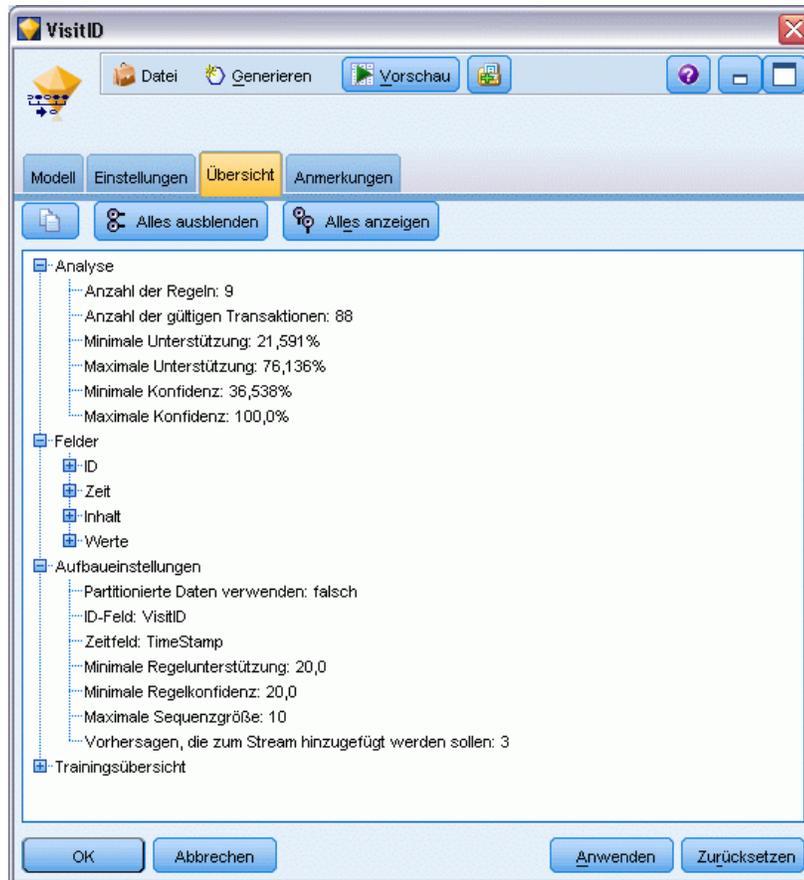


**Maximale Anzahl an Prognosen.** Dient zur Angabe der maximalen Anzahl an Vorhersagen, die für jedes Set von Warenkorbelementen (basket items) aufgenommen werden. Anhand der Regeln mit den höchsten Konfidenzwerten, die für dieses Set von Transaktionen gelten, werden die Prognosen für den Datensatz bis zum angegebenen Grenzwert generiert.

## Sequenzmodell-Nugget-Übersicht

Auf der Registerkarte “Übersicht” für ein Sequenzregelmodell-Nugget werden die Anzahl der ermittelten Regeln sowie die Mindest- und Höchstwerte für Unterstützung und Konfidenz für die Regeln angegeben. Wenn Sie einen Analyseknoden ausgeführt haben, der an diesen Modellierungsknoten angehängt ist, werden die Informationen aus dieser Analyse ebenfalls in diesem Abschnitt angezeigt. [Für weitere Informationen siehe Thema Analyseknoden in Kapitel 6 in IBM SPSS Modeler 14.2- Quellen-, Prozess- und Ausgabeknoten.](#)

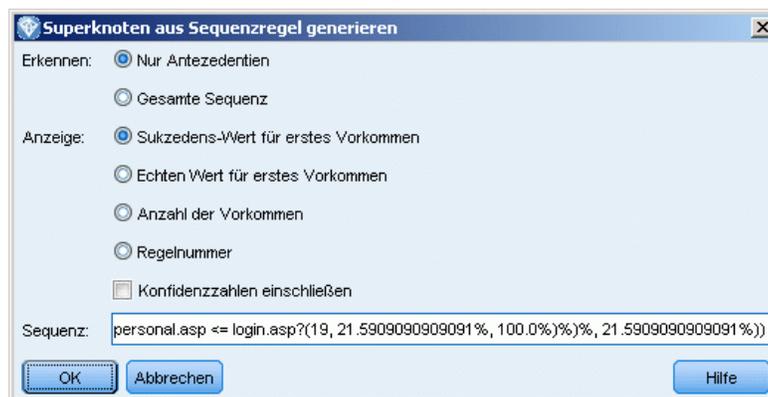
Abbildung 12-31  
 Sequenz-Nugget – Registerkarte "Übersicht"



Für weitere Informationen siehe Thema Durchsuchen von Modell-Nuggets in Kapitel 3 auf S. 52.

## Generieren eines Regel-Superknotens aus einem Sequenzmodell-Nugget

Abbildung 12-32  
 Dialogfeld zum Generieren eines Regel-Superknotens



So generieren Sie einen Regel-Superknoten auf der Grundlage einer Sequenzregel:

- ▶ Klicken Sie auf der Registerkarte “Modell” für das Sequenzregelmodell-Nugget auf die Zeile in der Tabelle, in der die gewünschte Regel verzeichnet ist.
- ▶ Wählen Sie im Regelbrowser in den Menüs die folgende Befehlsfolge:  
Erzeugen > Regel-Superknoten

*Wichtiger Hinweis:* Zur Verwendung des generierten Superknotens müssen Sie die Daten nach ID-Feld (und Zeitfeld (sofern vorhanden)) sortieren, bevor Sie sie an den Superknoten übergeben. In nicht sortierten Daten kann der Superknoten die Sequenzen nicht ordnungsgemäß erkennen.

Folgende Optionen stehen für das Generieren von Regel-Superknoten zur Verfügung:

**Erkennen.** Gibt an, wie für die an den Superknoten übergebenen Daten Übereinstimmungen definiert sind.

- **Nur Antezedenzen.** Der Superknoten ermittelt eine Übereinstimmung jedes Mal, wenn er die Antezedenzen für die ausgewählte Regel in der richtigen Reihenfolge in einem Set von Datensätzen mit derselben ID findet, unabhängig davon, ob das Sukzedens ebenfalls gefunden wird. Beachten Sie, dass hierbei nicht die Beschränkungseinstellungen für die Zeitstempeltoleranz oder die Elementlücken aus dem ursprünglichen Sequenzmodellierungsknoten berücksichtigt werden. Wenn das letzte Antezedens-Element-Set im Stream ermittelt wird (und alle anderen Antezedenzen in der richtigen Reihenfolge gefunden wurden), enthalten alle nachfolgenden Datensätze mit der aktuellen ID die unten ausgewählte Übersicht.
- **Gesamte Sequenz.** Der Superknoten ermittelt eine Übereinstimmung jedes Mal, wenn er die Antezedenzen für die ausgewählte Regel in der richtigen Reihenfolge in einem Set von Datensätzen mit derselben ID findet, unabhängig davon, ob das Sukzedens ebenfalls gefunden wird. Hierbei werden nicht die Beschränkungseinstellungen für die Zeitstempeltoleranz oder die Elementlücken aus dem ursprünglichen Sequenzmodellierungsknoten berücksichtigt. Wenn das letzte Sukzedens im Stream ermittelt wird (und alle Antezedenzen ebenfalls in der richtigen Reihenfolge gefunden wurden), enthalten der aktuelle Datensatz und alle nachfolgenden Datensätze mit der aktuellen ID die unten ausgewählte Übersicht.

**Anzeigen.** Steuert, wie Übereinstimmungsübersichten zu den Daten in der Ausgabe des Regel-Superknotens hinzugefügt werden.

- **Sukzedens-Wert für erstes Vorkommen** Bei dem zu den Daten hinzugefügten Wert handelt es sich um den Wert des Sukzedens, der auf der Grundlage des erstens Vorkommnisses der Übereinstimmung vorhergesagt wurde. Die Werte werden als neues Feld mit der Bezeichnung *rule\_n\_consequent* hinzugefügt. Dabei ist *n* die Regelnummer (gemäß der Erstellungsreihenfolge der Regel-Superknoten im Stream).
- **Echten Wert für erste Vorkommen** Der zu den Daten hinzugefügte Wert ist wahr, wenn mindestens eine Übereinstimmung für die ID vorliegt, und falsch, wenn keine Übereinstimmung vorliegt. Die Werte werden als neues Feld mit der Bezeichnung *rule\_n\_flag* hinzugefügt.
- **Anzahl der Vorkommen.** Bei dem zu den Daten hinzugefügten Wert handelt es sich um die Anzahl der Übereinstimmungen für die ID. Die Werte werden als neues Feld mit der Bezeichnung *rule\_n\_count* hinzugefügt.

- **Regelnummer.** Der hinzugefügte Wert ist die Regelnummer für die ausgewählte Regel. **Regelnummern** werden auf der Grundlage der Reihenfolge zugewiesen, in der der Superknoten zum Stream hinzugefügt wurde. Der erste Regel-Superknoten beispielsweise wird als *rule 1* betrachtet, der zweite als *rule 2* usw. Diese Option ist sinnvoll, wenn Sie mehrere Superknoten in Ihrem Stream berücksichtigen möchten. Die Werte werden als neues Feld mit der Bezeichnung *rule\_n\_number* hinzugefügt.
- **Konfidenzzahlen einschließen.** Bei Auswahl dieser Option wird die Regelkonfidenz zum Datenstrom hinzugefügt, ebenso wie die ausgewählten Übersichtsdaten. Die Werte werden als neues Feld mit der Bezeichnung *rule\_n\_confidence* hinzugefügt.

# Zeitreihenmodelle

## Wozu dienen Vorhersagen?

Bei einer Vorhersage werden die Werte für eine oder mehrere Reihen im zeitlichen Verlauf vorhergesagt. Beispiel: Sie möchten vorhersagen, wie die erwartete Nachfrage für eine Produktlinie oder eine Dienstleistung aussehen wird, um Ressourcen für die Fertigung oder Distribution zuzuordnen. Da die Implementierung der Planung von Entscheidungen zeitaufwändig ist, bilden Vorhersagen bei vielen Planungsprozessen ein wichtiges Werkzeug.

Die Methoden zur Modellierung von Zeitreihen gehen davon aus, dass sich die Geschichte wiederholt — wenn nicht genau, dann doch genau genug, dass eine Untersuchung der Vergangenheit in der Zukunft bessere Entscheidungen ermöglicht. Um z. B. die Verkaufszahlen für das nächste Jahr vorherzusagen, würden Sie wahrscheinlich damit beginnen, die Verkaufszahlen dieses und früherer Jahre zu untersuchen und eventuell vorhandene Trends oder Muster herauszuarbeiten, die sich in früheren Jahren entwickelt haben. Es kann sich aber als schwierig erweisen, Muster zu beurteilen. Wenn Ihre Verkaufszahlen beispielsweise über mehrere Wochen in Folge ansteigen, handelt es sich dann um den Beginn eines saisonbedingten Zyklus oder um den Anfang eines langfristigen Trends?

Durch statistische Modellierungsverfahren können Sie die in Ihren Vergangenheitsdaten vorhandenen Muster analysieren und projizieren, um einen Bereich zu ermitteln, in dem zukünftige Werte der Reihen wahrscheinlich liegen werden. Als Ergebnis erhalten Sie genauere Vorhersagen, auf deren Grundlage Sie Entscheidungen treffen können.

## Zeitreihendaten

Bei einer **Zeitreihe** handelt es sich um eine sortierte Sammlung von Messdaten, die in regelmäßigen Zeitabständen ermittelt wurden – beispielsweise tägliche Lagerkosten oder wöchentliche Verkaufsdaten. Die Messungen können sich auf alles beziehen, was für Sie von Interesse ist. Alle Reihen können wie folgt klassifiziert werden:

- **Abhängig.** Eine Reihe, die Sie vorhersagen möchten.
- **Prädiktor.** Eine Reihe, die bei der Erklärung des Ziel hilfreich sein kann – wenn z. B. ein Anzeigebudget verwendet wird, um den Absatz vorherzusagen. Prädiktoren können nur mit ARIMA-Modellen verwendet werden.
- **Ereignis.** Eine spezielle Prädiktorenreihe, die verwendet wird, um vorhersagbare wiederkehrende Vorfälle zu berücksichtigen — wie beispielsweise Werbeaktionen.
- **Intervention.** Eine spezielle Prädiktorenreihe, die verwendet wird, um einmalige Vorfälle in der Vergangenheit zu berücksichtigen — wie beispielsweise ein Stromausfall oder ein Streik.

Bei den Intervallen kann es sich um beliebige Zeiträume handeln. Das Intervall muss jedoch für alle Messungen identisch sein. Darüber hinaus müssen alle Intervalle, für die keine Messungen vorliegen, als fehlende Werte festgelegt werden. Demnach definiert die Anzahl der Intervalle mit

Messungen (einschließlich derer mit fehlenden Werten) den Umfang der historischen Spannweite der Daten.

## ***Merkmale von Zeitreihen***

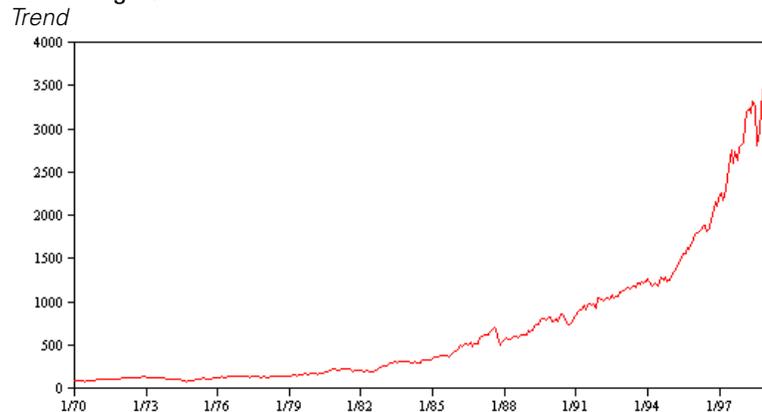
Die Untersuchung der früheren Ergebnisse einer Reihe erleichtert das Auffinden von Mustern und sorgt für bessere Vorhersagen. Bei der Darstellung als Diagramm zeigen viele Zeitreihen eines oder mehrere der folgenden Merkmale:

- Trends
- Saisonale und nichtsaisonale Zyklen
- Impulse und Schritte
- Ausreißer

### ***Trends***

Ein **Trend** ist eine allmähliche Aufwärts- oder Abwärtsveränderung der Ebene von Reihen oder der Tendenz von Reihenwerten, im Laufe der Zeit zu steigen oder zu sinken.

Abbildung 13-1



Trends sind entweder **lokal** oder **global**, eine einzelne Reihe kann jedoch beide Typen aufweisen. Reihendiagramme des Aktienindexes zeigen historisch einen globalen Aufwärtstrend. In Zeiten der Rezession tauchen lokale Abwärtstrends und in Zeiten der Hochkonjunktur lokale Aufwärtstrends auf.

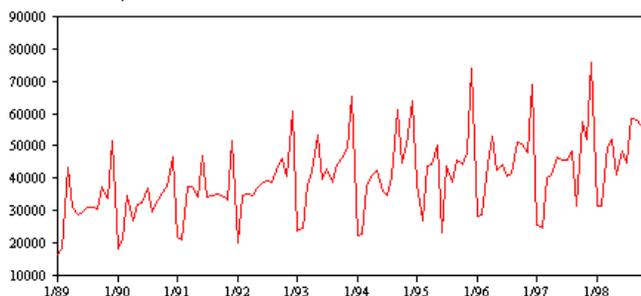
Trends können außerdem entweder **linear** oder **nichtlinear** sein. Lineare Trends sind positive oder negative additive Erhöhungen der Ebene von Reihen, vergleichbar mit dem Effekt der einfachen Kapitalverzinsung. Nichtlineare Trends sind häufig multiplikativ, mit Inkrementen, die sich proportional zu den vorherigen Reihenwerten verhalten.

Globale lineare Trends werden angepasst und liefern gute Vorhersagen sowohl bei Modellen für das exponentielle Glätten als auch mit ARIMA-Modellen. Beim Erstellen von ARIMA-Modellen werden Reihen, die Trends aufweisen, in der Regel unterschieden, um die Auswirkung des Trends zu beseitigen.

## Saisonale Zyklen

Ein **saisonaler Zyklus** ist ein sich wiederholendes, vorhersagbares Muster in den Reihenwerten.

Abbildung 13-2  
Saisonaler Zyklus



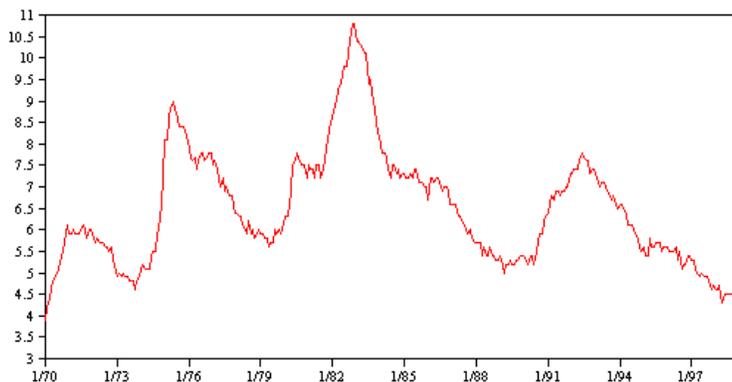
Saisonale Zyklen sind mit dem Intervall ihrer Reihen verbunden. Monatsdaten weisen z. B. in der Regel Zyklen über Quartale und Jahre auf. Eine Monatsreihe kann einen signifikanten vierteljährlichen Zyklus aufweisen, der im ersten Quartal gering ausfällt, oder einen jährlichen Zyklus mit einem Spitzenwert im Dezember. Reihen, die einen saisonalen Zyklus enthalten, zeigen eine **Saisonabhängigkeit**.

Saisonale Muster sind hilfreich, um brauchbare Anpassungen und Vorhersagen zu erhalten. Die Saisonabhängigkeit wird durch Modelle für das exponentielle Glätten sowie durch ARIMA-Modelle erfasst.

## Nichtseasonale Zyklen

Ein **nichtseasonaler Zyklus** ist ein sich wiederholendes, möglicherweise vorhersagbares Muster in den Reihenwerten.

Abbildung 13-3  
Nichtseasonaler Zyklus



Manche Reihen, wie die Arbeitslosenquote, zeigen einen deutlich zyklischen Verlauf. Die Periodizität des Zyklus verändert sich jedoch im Laufe der Zeit, weshalb es schwierig vorhersagbar ist, wann ein Höchst- oder Tiefstwert eintreten wird. Andere Reihen besitzen möglicherweise vorhersagbare Zyklen, passen aber nicht wirklich in den Gregorianischen Kalender oder haben

Zyklen, die länger als ein Jahr sind. Die Gezeiten folgen z. B. dem Mondkalender, internationale Reise- und Handelsaktivitäten im Zusammenhang mit den Olympischen Spielen steigen alle vier Jahre an und es gibt viele religiöse Feiertage, die jedes Jahr auf ein anderes Datum fallen.

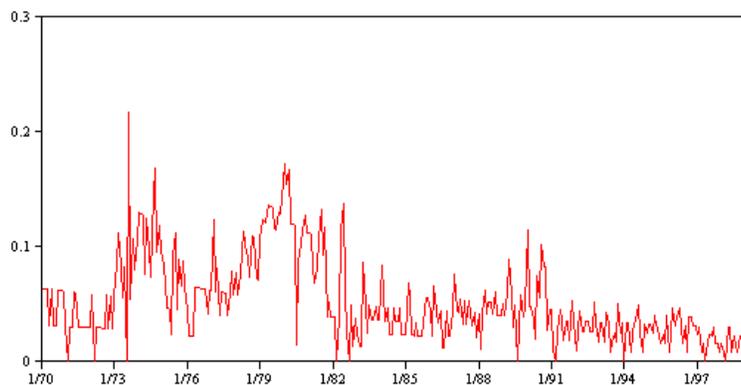
Nichtseasonale zyklische Muster sind schwierig zu modellieren und erhöhen in der Regel die Unsicherheit der Vorhersagen. Der Aktienmarkt bietet beispielsweise eine Vielzahl von Beispielen für Reihen, die sich den Anstrengungen beim Erstellen von Vorhersagen widersetzen. Trotzdem müssen nichtseasonale Muster berücksichtigt werden, wenn sie vorhanden sind. In vielen Fällen können Sie immer noch ein Modell identifizieren, das halbwegs gut zu den historischen Daten passt und Ihnen die beste Chance bietet, die Unsicherheit der Vorhersage möglichst zu minimieren.

### **Impulse und Schritte**

In vielen Reihen treten plötzliche Ebenenänderungen auf. Diese gehören in der Regel zu zwei Typen:

- Eine plötzliche *vorübergehende* Änderung oder ein **Impuls** in der Ebene der Reihe
- Eine plötzliche *permanente* Änderung oder ein **Impuls** in der Ebene der Reihe

Abbildung 13-4  
Reihen mit einem Impuls



Wenn Schritte oder Impulse beobachtet werden, ist es wichtig, eine plausible Erklärung zu finden. Zeitreihenmodelle sind so ausgelegt, dass sie allmähliche Änderungen berücksichtigen, keine plötzlichen. Demzufolge tendieren sie dazu, Impulse unterzubewerten und durch Schritte ruiniert zu werden, was zu einer schlechten Modellanpassungsgüte und unsicheren Vorhersagen führt. (Bei einigen Instanzen der Saisonabhängigkeit kann eine scheinbare plötzliche Änderung der Ebene vorliegen, während die Ebene von einem saisonalen Zeitraum zum nächsten aber konstant ist.)

Wenn eine Störung erklärt werden kann, kann sie mithilfe einer **Intervention** oder eines **Ereignisses** modelliert werden. Im August 1973 hat beispielsweise ein Erdölembargo der Organisation der Erdöl exportierenden Länder (OPEC) eine drastische Veränderung der Inflationsrate ausgelöst, die dann in den darauffolgenden Monaten wieder auf einen normalen Stand zurückkehrte. Durch die Angabe einer **Punkt-Intervention** für den Monat des Embargos können Sie die Anpassungsgüte Ihres Modells verbessern und Ihre Vorhersagen so indirekt verbessern. Ein Einzelhandelsgeschäft stellt z. B. fest, dass der Absatz an einem Tag, an dem alle Artikel mit 50 % Rabatt gekennzeichnet waren, viel höher als gewöhnlich ist. Indem die 50

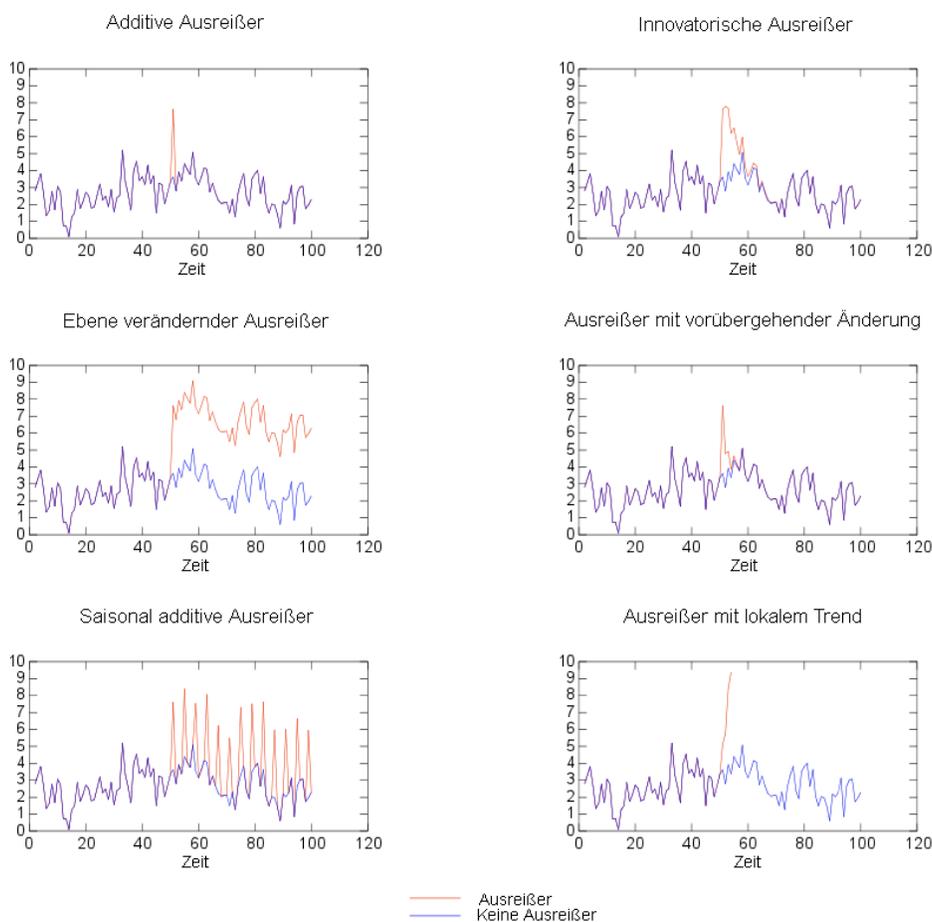
%-Rabatt-Aktion als wiederkehrendes **Ereignis** festgelegt wird, können Sie die Anpassungsgüte Ihres Modells verbessern und den Effekt der Wiederholung der Aktion in der Zukunft schätzen.

## Ausreißer

Veränderungen der Ebene einer Zeitreihe, die nicht erklärt werden können, werden als **Ausreißer** bezeichnet. Diese Beobachtungen sind nicht mit dem Rest der Reihen konsistent und können sich drastisch auf die Analyse auswirken und dementsprechend die Vorhersagefähigkeit des Modells beeinflussen.

Die folgende Abbildung zeigt verschiedene Typen von Ausreißern, die häufig in Zeitreihen auftreten. Die blaue Linie stellt eine Reihe ohne Ausreißer dar. Die roten Linien zeigen ein Muster, das vorliegen kann, wenn die Reihen Ausreißer enthalten. Diese Ausreißer werden als **deterministisch** klassifiziert, da sie sich nur auf die durchschnittliche Ebene der Reihen auswirken.

Abbildung 13-5  
Ausreißer Typen



- **Additive Ausreißer.** Ein additiver Ausreißer tritt als überraschend großer oder kleiner Wert auf, der für eine einzelne Beobachtung erscheint. Nachfolgende Beobachtungen sind nicht durch den additiven Ausreißer beeinflusst. Aufeinanderfolgende additive Ausreißer werden in der Regel als **additive Ausreißer-Patches** bezeichnet.
- **Innovatorische Ausreißer.** Ein innovatorischer Ausreißer ist gekennzeichnet durch eine anfängliche Auswirkung, deren Effekt bei nachfolgenden Beobachtungen fortbesteht. Der Einfluss des Ausreißers kann im zeitlichen Verlauf ansteigen.
- **Ebene verändernder Ausreißer.** Bei einer Ebenenänderung verschieben sich alle nach dem Ausreißer liegenden Beobachtungen auf eine neue Ebene. Im Gegensatz zu additiven Ausreißern wirkt sich ein Ebenen ändernder Ausreißer auf viele Beobachtungen aus und besitzt einen permanenten Effekt.
- **Ausreißer mit vorübergehender Änderung.** Ausreißer mit vorübergehender Änderung sind ähnlich wie Ausreißer mit Ebenenänderung, der Effekt der Ausreißer verringert sich aber bei den nachfolgenden Beobachtungen exponentiell. Die Reihe kehrt schließlich auf ihre normale Ebene zurück.
- **Saisonal additive Ausreißer.** Ein saisonal additiver Ausreißer tritt als überraschend großer oder kleiner Wert auf, der wiederholt in regelmäßigen Intervallen erscheint.
- **Ausreißer mit lokalem Trend.** Ein Ausreißer mit lokalem Trend verursacht in den Reihen eine allgemeine Tendenz, die nach dem Entstehen des ersten Ausreißers durch ein Muster in den Ausreißern verursacht wird.

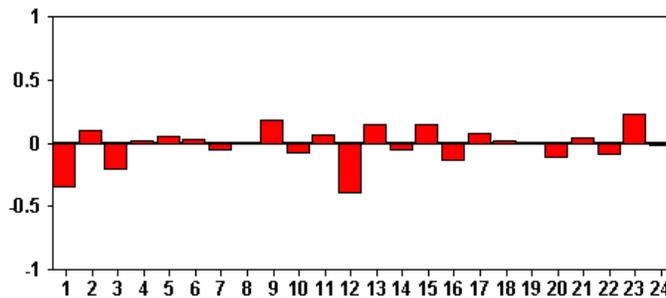
Beim Erkennen von Ausreißern in Zeitreihen müssen die Position, der Typ und der Umfang aller vorhandenen Ausreißer festgestellt werden. Tsay (1988) hat ein iteratives Verfahren zum Erkennen der durchschnittlichen Ebenenänderung vorgeschlagen, mit dem deterministische Ausreißer identifiziert werden. Bei diesem Prozess wird ein Zeitreihenmodell, in dem angenommen wird, dass keine Ausreißer vorhanden sind, mit einem anderen Modell verglichen, das Ausreißer berücksichtigt. Die Differenzen zwischen den Modellen ergeben Schätzungen dafür, welchen Effekt die Behandlung eines beliebigen Punkts als Ausreißer besitzt.

### ***Autokorrelation und partielle Autokorrelationsfunktionen***

Die Autokorrelation und die partielle Autokorrelation sind Maßstäbe für die Beziehung zwischen aktuellen und vergangenen Reihenwerten, die anzeigen, welche vergangenen Reihenwerte bei der Vorhersage zukünftiger Werte am nützlichsten sind. Mit diesem Wissen können Sie die Reihenfolge der Prozesse in einem ARIMA-Modell festlegen. Spezieller

- **Autokorrelationsfunktion (ACF).** Mit Abstand  $k$  ist dies die Korrelation zwischen Reihenwerten, die  $k$  Intervalle entfernt sind.
- **Partielle Autokorrelationsfunktion (PACF).** Mit Abstand  $k$  ist dies die Korrelation zwischen Reihenwerten, die  $k$  Intervalle entfernt sind, wobei die Werte der dazwischenliegenden Intervalle berücksichtigt werden.

Abbildung 13-6  
ACF-Plot für eine Reihe



Die x-Achse des ACF-Plots gibt den Abstand an, bei dem die Autokorrelation berechnet wurde. Die y-Achse gibt den Wert der Korrelation an (zwischen  $-1$  und  $1$ ). Beispiel: Eine Spitze bei Abstand 1 in einem ACF-Plot zeigt eine starke Korrelation zwischen allen Reihenwerten und dem vorherigen Wert an. Eine Spitze bei Abstand 2 zeigt eine starke Korrelation zwischen allen Werten und dem Wert an, der zwei Punkte vorher auftrat, usw.

- Eine positive Korrelation zeigt an, dass große aktuelle Werte großen Werten zum angegebenen Abstand entsprechen. Eine negative Korrelation zeigt an, dass große aktuelle Werte kleinen Werten zum angegebenen Abstand entsprechen.
- Der absolute Wert einer Korrelation ist ein Maßstab für die Stärke der Verbindung, wobei größere absolute Werte eine stärkere Beziehung anzeigen.

## Reihentransformationen

Transformationen sind häufig nützlich, um eine Reihe zu stabilisieren, bevor Modelle geschätzt werden. Dies ist insbesondere für ARIMA-Modelle wichtig, für die eine Reihe **feststehend** sein muss, bevor Modelle geschätzt werden. Eine Reihe ist feststehend, wenn die globale Ebene (Mittelwert) und die durchschnittliche Abweichung von der Ebene (Varianz) über die Reihen hinweg konstant sind.

Obwohl die meisten Reihen nicht feststehend sind, ist ARIMA effektiv, solange Reihen durch Transformationen, wie natürlicher Logarithmus, Differenz oder saisonale Differenz, feststehend gemacht werden können.

**Varianz stabilisierende Transformation.** Reihen, in denen sich die Varianz mit der Zeit ändert, können häufig mit einer natürlichen Logarithmus- oder einer Quadratwurzel-Transformation feststehend gemacht werden. Diese werden auch als funktionale Transformationen bezeichnet.

- **Natürlicher Logarithmus.** Der natürliche Logarithmus wird auf Reihenwerte angewendet.
- **Quadratwurzel.** Die Quadratwurzelfunktion wird auf die Reihenwerte angewendet.

Natürliche Logarithmus- und Quadratwurzel-Transformationen können für Reihen mit negativen Werten verwendet werden.

**Ebene stabilisierende Transformationen.** Ein langsamer Rückgang von Werten in der ACF zeigt an, dass alle Reihenwerte eine starke Korrelation mit dem vorherigen Wert besitzen. Durch die Analyse der Veränderung der Reihenwerte erhalten Sie eine stabile Ebene.

- **Einfache Differenz.** Die Differenzen zwischen den einzelnen Werten und dem vorherigen Wert in der Reihe werden berechnet, wobei der älteste Wert der Reihe ausgenommen wird. Dies bedeutet, dass die Differenzreihen einen Wert weniger als die Originalreihen besitzen.
- **Saisonale Differenz.** Identisch mit der einfachen Differenz, außer, dass die Differenzen zwischen den einzelnen Werten und den vorherigen saisonalen Werten berechnet werden.

Wenn entweder die einfache oder die saisonale Differenz gleichzeitig mit der Log- oder Quadratwurzel-Transformation eingesetzt wird, wird immer zuerst die Varianz stabilisierende Transformation angewendet. Wenn sowohl die einfache als auch die saisonale Differenz angewendet wird, erhalten Sie, unabhängig von der Reihenfolge, in der die Differenz gebildet wird, immer dieselben resultierenden Reihenwerte.

## Prädiktorreihen

Prädiktorreihen enthalten verwandte Daten, die das Ergebnis von Reihen erklären können, für die eine Vorhersage erstellt werden soll. Beispiel: Ein über das Internet oder mit Katalogen arbeitender Einzelhändler kann den Absatz anhand der versendeten Kataloge, der Anzahl verfügbarer Telefonverbindungen oder der Anzahl der Klicks auf die Website des Unternehmens vorhersagen.

Alle Reihen können als Prädiktor verwendet werden, wenn die Reihen so weit in die Zukunft reichen, wie sie eine Vorhersage erstellen möchten, und wenn vollständige Daten ohne fehlende Werte vorliegen.

Seien Sie vorsichtig, wenn Sie Prädiktoren zu einem Modell hinzufügen. Wenn eine große Anzahl von Prädiktoren hinzugefügt wird, steigt die für die Schätzung von Modellen erforderliche Zeit. Während das Hinzufügen von Prädiktoren die Anpassungsgüte des Modells für Vergangenheitsdaten verbessern kann, bedeutet dies nicht zwingend, dass das Modell bessere Vorhersagen liefert. Aus diesem Grund kann es sein, dass sich diese zusätzliche Komplexität nicht lohnt. Idealerweise sollte das Ziel die Identifizierung des einfachsten Modells darstellen, das eine gute Vorhersage ermöglicht.

Eine allgemeine Regel besagt, dass die Anzahl der Prädiktoren geringer sein sollte als der durch 15 geteilte Stichprobenumfang (höchstens ein Prädiktor für 15 Fälle).

**Prädiktoren mit fehlenden Daten.** Prädiktoren mit unvollständigen oder fehlenden Daten können nicht für Vorhersagen verwendet werden. Dies gilt sowohl für Vergangenheitsdaten als auch für zukünftige Werte. In einigen Fällen können Sie diese Einschränkung umgehen, indem Sie die Schätzspanne des Modells so einstellen, dass die ältesten Werte bei der Schätzung von Modellen ausgeschlossen werden.

## Zeitreihen – Modellierungsknoten

Der Zeitreihenknoten berechnet Schätzungen für exponentielle Glättung, univariate ARIMA-Modelle (ARIMA steht für Autoregressive Integrated Moving Average (autoregressiver integrierter gleitender Durchschnitt)) und multivariate ARIMA-Modelle (Transferfunktionsmodelle) für Zeitreihendaten und erstellt Vorhersagen auf der Grundlage der Zeitreihendaten.

**Exponentielles Glätten** ist ein Prognoseverfahren, bei dem gewichtete Werte aus früheren Beobachtungen der Zeitreihe verwendet werden, um zukünftige Werte vorherzusagen. An sich beruht das exponentielle Glätten nicht auf einem theoretischen Verständnis der Daten. Es wird

jeweils ein Punkt vorhergesagt und die Vorhersagen werden angepasst, wenn neue Daten verfügbar sind. Dieses Verfahren ist nützlich für die Vorhersage von Zeitreihen, die Trend und/oder Saisonalität aufweisen. Dabei können Sie zwischen verschiedenen Modellen mit exponentiellem Glätten wählen, die sich hinsichtlich der Behandlung von Trends und Saisonalität unterscheiden.

**ARIMA**-Modelle bieten feinere Methoden für die Modellierung von Trend- und saisonalen Komponenten als die Modelle der exponentiellen Glättung und weisen insbesondere den zusätzlichen Vorteil auf, dass unabhängige Variablen (Prädiktorvariablen) in das Modell integriert werden können. Hierfür müssen die die Ordnung der Autoregression, die Ordnung des gleitenden Durchschnitts und der Grad der Differenzbildung angegeben werden. Sie können Prädiktorvariablen einschließen und Transferfunktionen für bestimmte oder alle dieser Variablen definieren und die automatische Erkennung von Ausreißern oder einer bestimmten Gruppe von Ausreißern festlegen.

*Hinweis:* In der Praxis bedeutet dies, dass ARIMA-Modelle besonders nützlich sind, wenn Prädiktoren eingeschlossen werden sollen, die zur Erklärung des Verhaltens der prognostizierten Zeitreihe beitragen können, wie beispielsweise die Anzahl der versendeten Kataloge oder die Anzahl der Aufrufe einer Firmenwebseite. Modelle für das exponentielle Glätten beschreiben das Verhalten der Zeitreihen, ohne dass versucht wird zu verstehen, warum sich die Zeitreihe so verhält. Beispielsweise ist davon auszugehen, dass eine Zeitreihe, die bisher alle 12 Monate einen Höhepunkt erreicht hat, dies auch weiterhin tun wird, auch wenn Sie nicht wissen, warum.

Außerdem ist ein **Expert Modeler** verfügbar, der automatisch das am besten angepasste ARIMA-Modell bzw. das am besten angepasste Modell mit exponentiellem Glätten für eine oder mehrere Zielvariablen ermittelt, sodass das geeignete Modell nicht mehr nach dem Prinzip von Versuch und Irrtum ermittelt werden muss. In allen Fällen wählt der Expert Modeler jeweils das beste Modell für jede der angegebenen Zielvariablen. Im Zweifelsfall sollte der Expert Modeler verwendet werden.

Bei Angabe von Prädiktorvariablen wählt der Expert Modeler diejenigen Variablen zum Einschluss in ARIMA-Modelle aus, die eine statistisch signifikante Beziehung mit der abhängigen Zeitreihe aufweisen. Modellvariablen werden ggf. durch Differenzierung und/oder Quadratwurzeltransformation bzw. Transformation mit natürlichem Logarithmus transformiert. In der Standardeinstellung berücksichtigt der Expert Modeler alle Modelle für das exponentielle Glätten sowie alle ARIMA-Modelle und wählt für jedes Zielfeld das jeweils beste Modell aus. Sie können festlegen, dass der Expert Modeler nur das beste Modell mit exponentiellem Glätten oder nur das beste ARIMA-Modell auswählen soll. Sie können auch die automatische Erkennung von Ausreißern festlegen.

**Beispiel.** Ein Analyst eines Breitband-Providers soll Vorhersagen über die Vertragsabschlüsse mit Kunden erstellen, um die Auslastung der Bandbreite prognostizieren zu können. Es werden Vorhersagen für alle lokalen Märkte benötigt, die zusammen den landesweiten Kundenstamm ergeben. Mit der Zeitreihenmodellierung können Sie Vorhersagen für die nächsten drei Monate für eine Reihe von lokalen Märkten erstellen. [Für weitere Informationen siehe Thema Prognoseerstellung mit dem Zeitreihenknoten in Kapitel 15 in IBM SPSS Modeler 14.2-Anwendungshandbuch.](#)

## Voraussetzungen

Der Zeitreihenknoten weicht von anderen IBM® SPSS® Modeler-Knoten dahingehend ab, dass Sie ihn nicht einfach in einen Stream einfügen und den Stream ausführen können. Dem Zeitreihenknoten muss stets ein Zeitintervallknoten vorangehen, der Informationen angibt wie das zu verwendende Zeitintervall (Jahre, Quartale, Monate usw.), die für die Schätzung zu verwendenden Daten und wie weit in die Zukunft sich eine Vorhersage erstrecken soll, sofern verwendet.

Abbildung 13-7

Einem Zeitreihenknoten muss stets ein Zeitintervallknoten vorangehen



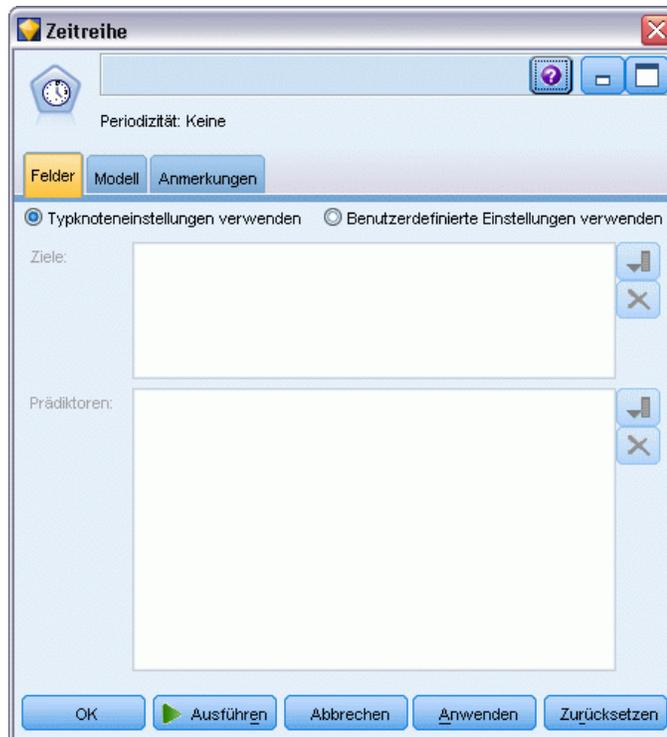
Die Zeitreihendaten müssen gleichmäßige Abstände aufweisen. Bei den Methoden zur Modellierung von Zeitreihendaten ist ein einheitliches Intervall zwischen den Messungen erforderlich; fehlende Werte werden durch leere Zeilen dargestellt. Falls Ihre Daten diese Bedingung nicht bereits erfüllen, können Sie sie mithilfe des Zeitintervallknotens entsprechend transformieren. [Für weitere Informationen siehe Thema Zeitintervallknoten in Kapitel 4 in IBM SPSS Modeler 14.2- Quellen-, Prozess- und Ausgabeknoten.](#)

Außerdem ist bei Zeitreihendaten zu beachten:

- Die Felder müssen numerisch sein.
- Datumsfelder können nicht als Eingaben verwendet werden.
- Partitionen werden ignoriert.

### Feldoptionen

Abbildung 13-8  
Dialogfeld des Zeitreihenknotens, Registerkarte "Felder"



Auf der Registerkarte "Felder" geben Sie an, welche Felder bei der Erstellung des Modells verwendet werden sollen. Bevor Sie ein Modell erstellen können, müssen Sie festlegen, welche Felder als Ziele und als Eingaben verwendet werden sollen. Normalerweise verwendet der Zeitreihenknoten Feldinformationen aus einem weiter oben liegenden Typknoten. Wenn Sie einen Typknoten benutzen, um Eingabe- und Zielfelder auszuwählen, brauchen Sie auf dieser Registerkarte keine Änderungen vorzunehmen.

**Typknoteneinstellungen verwenden.** Diese Option weist den Knoten an, die Feldinformationen von einem weiter oben liegenden Typknoten zu verwenden. Dies ist die Standardeinstellung.

**Benutzerdefinierte Einstellungen verwenden.** Diese Option weist den Knoten an, die hier angegebenen Feldinformationen anstelle der in einem weiter oben liegenden Typknoten angegebenen zu verwenden. Geben Sie nach Auswahl dieser Option die unten stehenden Felder an. Beachten Sie, dass als Datumswerte gespeicherte Felder nicht als Ziel- oder Eingabefelder zulässig sind.

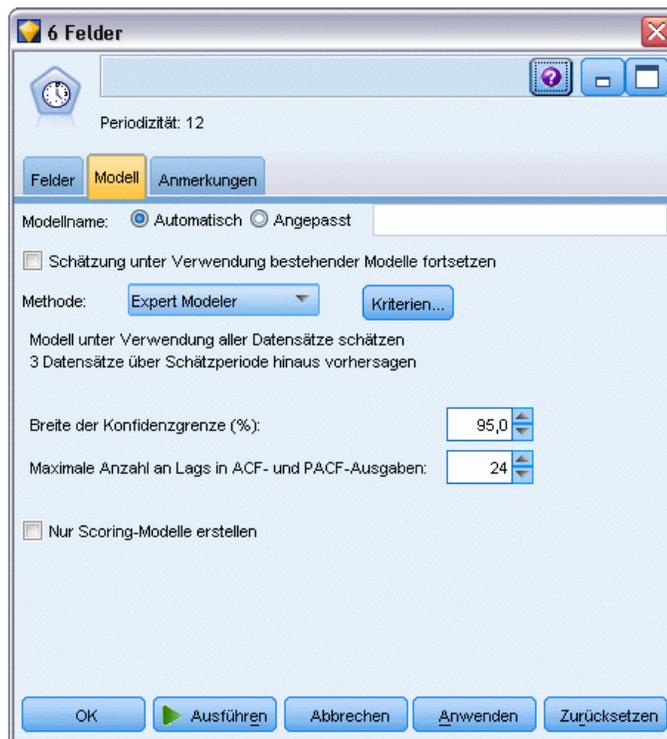
- Ziele.** Wählen Sie ein oder mehrere Zielfelder aus. Dies ist so, als würden Sie in einem Typknoten für die Rolle eines Felds den Wert *Ziel* festlegen. Zielfelder für Zeitreihenmodelle müssen ein Messniveau des Typs *Stetig* aufweisen. Für jedes Zielfeld wird ein separates Modell erstellt. Für Zielfelder kommen alle angegebenen *Eingabe*-Felder mit Ausnahme des jeweiligen Zielfelds selbst als mögliche Eingaben in Betracht. Daher kann dasselbe Feld in

beiden Listen vorkommen; ein solches Feld wird als mögliche Eingabe für alle Modelle verwendet, außer für das Modell, bei dem es ein Zielfeld ist.

- **Eingaben.** Wählen Sie die Eingabefelder aus. Dies ist so, als würden Sie in einem Typknoten für die Rolle eines Felds den Wert *Eingabe* festlegen. Eingabefelder für Zeitreihenmodelle müssen numerisch sein.

## Zeitreihenmodelle – Optionen

Abbildung 13-9  
Dialogfeld des Zeitreihenknotens, Registerkarte "Modell"



**Modellname.** Gibt den Namen des Modells an, das beim Ausführen des Knotens generiert wird.

- **Auto.** Generiert den Modellnamen automatisch basierend auf den Ziel- oder ID-Feldnamen oder dem Namen des Modelltyps in Fällen, in denen kein Ziel angegeben ist (z. B. Clustermodelle).
- **Benutzerdefiniert.** Hier können Sie einen benutzerdefinierten Namen für das Modell-Nugget angeben.

**Schätzung unter Verwendung bestehender Modelle fortsetzen.** Wenn Sie bereits ein Zeitreihenmodell generiert haben, wählen Sie diese Option aus, um die für dieses Modell angegebenen Kriterieneinstellungen wiederzuverwenden und einen neuen Modellknoten in der Modellpalette zu generieren, anstatt ein Modell völlig neu zu erstellen. Auf diese Weise können Sie Zeit sparen, indem Sie eine erneute Schätzung durchführen und eine neue Vorhersage auf der Grundlage derselben Modelleinstellungen erstellen wie zuvor, nur unter Verwendung aktuellerer Daten. Wenn also beispielsweise als ursprüngliches Modell für eine bestimmte Zeitreihe der lineare Trend

nach Holt verwendet wurde, wird derselbe Modelltyp für erneute Schätzungen und Vorhersagen für diese Daten verwendet; das System versucht nicht erneut, den besten Modelltyp für die neuen Daten zu ermitteln. Durch die Auswahl dieser Option werden die Steuerelemente Methode und Kriterien deaktiviert. [Für weitere Informationen siehe Thema Erneute Schätzung und Vorhersage auf S. 454.](#)

**Methode.** Sie haben die Wahl zwischen Expert Modeler, exponentiellem Glätten und ARIMA. [Für weitere Informationen siehe Thema Zeitreihen – Modellierungsknoten auf S. 438.](#) Wählen Sie Kriterien, um Optionen für die ausgewählte Methode anzugeben.

- **Expert Modeler.** Wählen Sie diese Option, um den Expert Modeler zu verwenden, der automatisch das jeweils am besten angepasste Modell für die einzelnen abhängigen Zeitreihen ermittelt.
- **Exponentielles Glätten.** Mit dieser Option können Sie ein benutzerdefiniertes Modell mit exponentiellem Glätten angeben.
- **ARIMA.** Mit dieser Option können Sie ein ARIMA-Modell angeben.

### ***Zeitintervallinformationen***

Dieser Bereich des Dialogfelds enthält Informationen zu Spezifikationen für Schätzer und Vorhersagen, die am Zeitintervallknoten vorgenommen werden. Beachten Sie, dass dieser Abschnitt nicht angezeigt wird, wenn Sie die Option Schätzung unter Verwendung bestehender Modelle fortsetzen auswählen.

Die erste Zeile der Informationen gibt an, ob Datensätze aus dem Modell ausgeschlossen oder als Holdouts verwendet werden. [Für weitere Informationen siehe Thema Schätzperiode in Kapitel 4 in IBM SPSS Modeler 14.2- Quellen-, Prozess- und Ausgabeknoten.](#)

Die zweite Zeile bietet Informationen zu den im Zeitintervallknoten angegebenen Vorhersageperioden. [Für weitere Informationen siehe Thema Vorhersagen in Kapitel 4 in IBM SPSS Modeler 14.2- Quellen-, Prozess- und Ausgabeknoten.](#)

Wenn in der ersten Zeile Kein Zeitintervall definiert steht, bedeutet dies, dass kein Zeitintervallknoten eingebunden ist. Dies führt zu einem Fehler beim Versuch, den Stream auszuführen; Sie müssen einen Zeitintervallknoten oberhalb des Zeitreihenknotens einfügen.

### ***Sonstige Informationen***

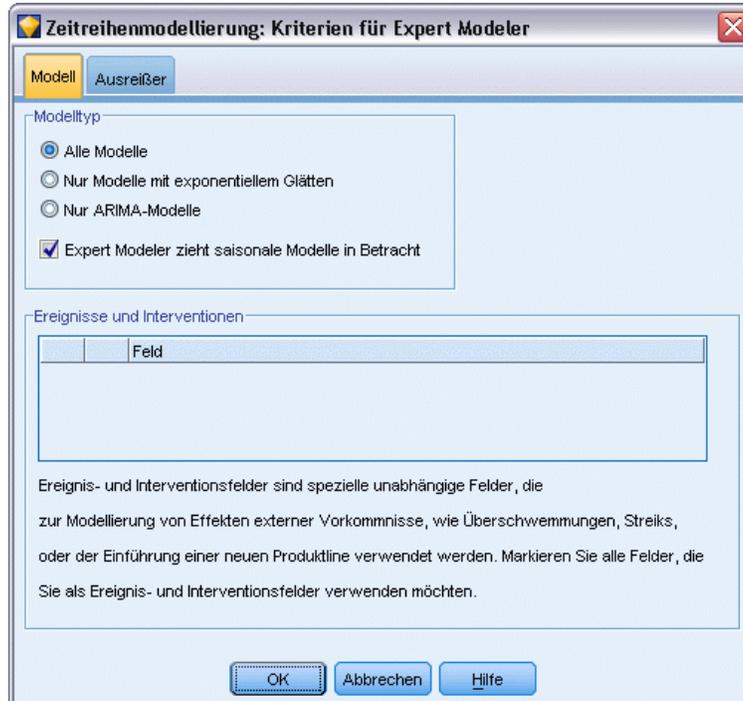
**Konfidenzniveau (%).** Konfidenzintervalle werden für die Modellvorhersagen und Residuen-Autokorrelationen berechnet. Es kann ein beliebiger positiver Wert unter 100 angegeben werden. In der Standardeinstellung wird ein Konfidenzintervall von 95 % verwendet.

**Maximale Anzahl an Lags in ACF- und PACF-Ausgaben.** Sie können die Höchstanzahl von Intervallen festlegen, die in Tabellen und Diagrammen für Autokorrelationen und partielle Autokorrelationen angezeigt werden.

**Nur Scoring-Modell aufbauen.** Markieren Sie dieses Kontrollkästchen, um die im Modell gespeicherte Datenmenge zu reduzieren. Dies kann die Leistung beim Erstellen von Modellen mit extrem vielen Zeitreihen (Zehntausende) verbessern. Wenn Sie diese Option wählen, werden die Registerkarten "Modell", "Parameter" und "Residuen" nicht im Zeitreihen-Modell-Nugget angezeigt, aber Sie können dennoch die Daten wie üblich scoren.

## Zeitreihen – Expert Modeler-Kriterien

Abbildung 13-10  
Dialogfeld "Kriterien für Expert Modeler," Registerkarte "Modell"



**Modelltyp.** Die folgenden Optionen sind verfügbar:

- **Alle Modelle.** Der Expert Modeler berücksichtigt sowohl ARIMA-Modelle als auch Modelle mit exponentiellem Glätten.
- **Nur Modelle mit exponentiellem Glätten.** Der Expert Modeler berücksichtigt nur Modelle mit exponentiellem Glätten.
- **Nur ARIMA-Modelle.** Der Expert Modeler berücksichtigt nur ARIMA-Modelle.

**Expert Modeler berücksichtigt saisonale Modelle.** Diese Option ist nur verfügbar, wenn für die Arbeitsdatei eine Periodizität definiert wurde. Wenn diese Option aktiviert ist, berücksichtigt der Expert Modeler sowohl saisonale als auch nichtsaisonale Modelle. Wenn diese Option deaktiviert ist, berücksichtigt der Expert Modeler nur nichtsaisonale Modelle.

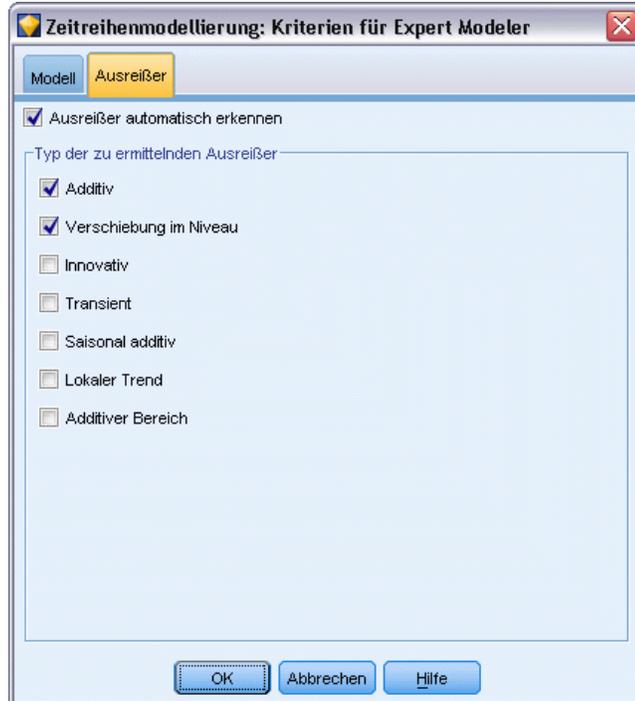
**Ereignisse und Interventionen.** Mit dieser Option können Sie bestimmte Eingabefelder als Ereignis- bzw. Interventionsfelder kennzeichnen. Dadurch wird angegeben, dass das betreffende Feld Zeitreihendaten enthält, die von Ereignissen (vorhersagbare wiederkehrende Situationen, z. B. Werbeaktionen) oder Interventionen (einmalige Vorfälle, z. B. Stromausfall, Streik) betroffen sind. Der Expert Modeler berücksichtigt nur einfache Regression und nicht frei wählbare Transferfunktionen für Eingaben, die als Ereignis- bzw. Interventionsfelder gekennzeichnet sind.

Eingabefelder müssen das Messniveau *Flag*, *Nominal* oder *Ordinal* aufweisen und müssen numerisch sein (z. B. "1"/"0" und nicht "Wahr"/"Falsch" für ein Flag-Feld), um in dieser Liste angezeigt zu werden. [Für weitere Informationen siehe Thema Impulse und Schritte auf S. 434.](#)

## Ausreißer

Abbildung 13-11

Dialogfeld "Kriterien für Expert Modeler," Registerkarte "Ausreißer"



**Ausreißer automatisch erkennen.** In der Standardeinstellung wird keine automatische Erkennung von Ausreißern durchgeführt. Aktivieren Sie diese Option, um die automatische Erkennung von Ausreißern durchzuführen, und wählen Sie anschließend die gewünschten Ausreißertypen aus.

Für weitere Informationen siehe [Thema Ausreißer auf S. 435](#).

## Zeitreihen – Kriterien für exponentielles Glätten

Abbildung 13-12  
Dialogfeld "Kriterien für exponentielles Glätten"



**Modelltyp.** Modelle für das exponentielle Glätten werden als saisonal oder nichtsaisonal klassifiziert. Saisonale Modelle stehen nur zur Verfügung, wenn die unter Verwendung des Zeitintervallknotens definierte Periodizität saisonal ist. Es gibt folgende saisonale Periodizitäten: zyklische Perioden, Jahre, Quartale, Monate, Tage pro Woche, Stunden pro Tag, Minuten pro Tag und Sekunden pro Tag. [Für weitere Informationen siehe Thema Zeitintervallknoten in Kapitel 4 in IBM SPSS Modeler 14.2- Quellen-, Prozess- und Ausgabeknoten.](#)

- **Einfach.** Dieses Modell eignet sich für Zeitreihen ohne Trend oder Saisonalität. Der einzige relevante Glättungsparameter für dieses Modell ist das Niveau. Einfaches exponentielles Glätten weist die größte Ähnlichkeit mit einem ARIMA-Modell mit Autoregression der Ordnung null, Differenzenbildung der Ordnung eins, gleitendem Durchschnitt der Ordnung eins und ohne Konstante auf.
- **Linearer Trend nach Holt.** Dieses Modell eignet sich für Zeitreihen mit linearem Trend und ohne Saisonalität. Die dafür relevanten Glättungsparameter sind Niveau und Trend, die bei diesem Modell nicht durch die Werte des jeweils anderen Parameters eingeschränkt sind. Das Holt-Modell ist allgemeiner als das Brown-Modell, die Berechnung von Schätzungen für große Zeitreihen kann allerdings mehr Zeit in Anspruch nehmen. Das exponentielle Glätten nach Holt weist die größte Ähnlichkeit mit einem ARIMA-Modell mit Autoregression der Ordnung null, Differenzenbildung der Ordnung zwei und gleitendem Durchschnitt der Ordnung zwei auf.
- **Linearer Trend nach Brown.** Dieses Modell eignet sich für Zeitreihen mit linearem Trend und ohne Saisonalität. Die dafür relevanten Glättungsparameter sind Niveau und Trend. Bei diesem Modell wird jedoch davon ausgegangen, dass diese gleich sind. Das Brown-Modell ist daher ein Spezialfall des Holt-Modells. Das exponentielle Glätten nach Brown weist die größte Ähnlichkeit mit einem ARIMA-Modell mit Autoregression der Ordnung null, Differenzenbildung der Ordnung zwei und gleitendem Durchschnitt der Ordnung zwei auf.

wobei der Koeffizient der zweiten Ordnung des gleitenden Durchschnitts die Hälfte des quadrierten Koeffizienten für die erste Ordnung beträgt.

- **Gedämpfter Trend.** Dieses Modell eignet sich für Zeitreihen mit auslaufendem linearem Trend und ohne Saisonalität. Die dafür relevanten Glättungsparameter sind Niveau, Trend und gedämpfter Trend. Gedämpftes exponentielles Glätten weist die größte Ähnlichkeit mit einem ARIMA-Modell mit Autoregression der Ordnung eins, Differenzenbildung der Ordnung eins und gleitendem Durchschnitt der Ordnung zwei auf.
- **Einfach saisonal.** Dieses Modell eignet sich für Zeitreihen ohne Trend und mit einem saisonalen Effekt, der im Zeitverlauf konstant bleibt. Die dafür relevanten Glättungsparameter sind Niveau und Saison. Saisonales exponentielles Glätten weist die größte Ähnlichkeit mit einem ARIMA-Modell mit Autoregression der Ordnung null, Differenzenbildung der Ordnung eins, saisonaler Differenzenbildung der Ordnung eins und den Ordnungen 1,  $p$  und  $p+1$  für den gleitenden Durchschnitt auf, wobei  $p$  die Anzahl der Perioden in einem saisonalen Intervall ist. Für Monatsdaten gilt:  $p = 12$ .
- **Additives Winters-Modell.** Dieses Modell eignet sich für Zeitreihen mit linearem Trend und mit einem saisonalen Effekt, der im Zeitverlauf konstant bleibt. Die dafür relevanten Glättungsparameter sind Niveau, Trend und Saison. Das exponentielle Glätten nach dem additiven Winters-Modell weist die größte Ähnlichkeit mit einem ARIMA-Modell mit Autoregression der Ordnung null, Differenzenbildung der Ordnung eins, saisonaler Differenzenbildung der Ordnung eins und  $p+1$  Ordnungen für den gleitenden Durchschnitt auf, wobei  $p$  die Anzahl der Perioden in einem saisonalen Intervall ist. Für Monatsdaten gilt:  $p = 12$ .
- **Multiplikatives Winters-Modell.** Dieses Modell eignet sich für Zeitreihen mit linearem Trend und mit einem saisonalen Effekt, der sich mit der Größenordnung der Zeitreihe ändert. Die dafür relevanten Glättungsparameter sind Niveau, Trend und Saison. Exponentielles Glätten mit dem multiplikativen Winters-Modell weist keine Ähnlichkeit zu irgendeinem ARIMA-Modell auf.

**Zieltransformation.** Sie können für jede abhängige Variable eine Transformation angeben, die vor deren Modellierung durchgeführt werden soll. [Für weitere Informationen siehe Thema Reihentransformationen auf S. 437.](#)

- **Keine.** Es wird keine Transformation durchgeführt.
- **Quadratwurzel.** Es wird eine Quadratwurzeltransformation durchgeführt.
- **Natürlicher Logarithmus.** Es wird eine Transformation mit natürlichem Logarithmus durchgeführt.

## **Zeitreihen – ARIMA-Kriterien**

Mit dem Zeitreihenknoten können Sie benutzerdefinierte nichtsaisonale oder saisonale ARIMA-Modelle – auch als Box-Jenkins-Modelle bekannt – mit oder ohne festes Set von Eingabevariablen (Prädiktorvariablen) erstellen. Sie können Transferfunktionen für bestimmte oder alle Eingabevariablen definieren und die automatische Erkennung von Ausreißern oder einer bestimmten Gruppe von Ausreißern festlegen.

Alle angegebenen Eingabevariablen werden explizit in das Modell aufgenommen. Im Gegensatz dazu werden beim Expert Modeler Eingabevariablen nur aufgenommen, wenn sie eine statistisch signifikante Beziehung zu der Zielvariablen aufweisen.

### Modell

Auf der Registerkarte "Modelle" können Sie die Struktur eines benutzerdefinierten ARIMA-Modells festlegen.

Abbildung 13-13  
Dialogfeld "ARIMA-Kriterien," Registerkarte "Modell"

The screenshot shows a dialog box titled "Zeitreihenmodellierung: ARIMA-Kriterien" with a close button (X) in the top right corner. It has three tabs: "Modell" (selected), "Übertragungsfunktionen", and "Ausreißer".

Under the "Arima-Ordnungen" section, there is a "Struktur:" label above a table:

	Nichtsaisonal	Saisonal
Autoregressiv (p)	0	0
Differenz(d)	0	0
Gleitender Durchschnitt (q)	0	0

Below the table is a "Zieltransformation" section with three radio buttons: "Keine" (selected), "Quadratwurzel", and "Natürlicher Logarithmus".

At the bottom left, there is a checked checkbox labeled "Konstante in Modell einbeziehen".

At the bottom center, there are three buttons: "OK", "Abbrechen", and "Hilfe".

**ARIMA-Ordnungen.** Geben Sie Werte für die verschiedenen ARIMA-Komponenten des Modells in die entsprechenden Zellen des Strukturgitters ein. Alle Werte müssen nichtnegative Ganzzahlen sein. Bei autoregressiven Komponenten und Komponenten des gleitenden Durchschnitts stellt der Wert die höchste Ordnung dar. Alle positiven niedrigeren Ordnungen werden in das Modell eingeschlossen. Wenn Sie beispielsweise 2 angeben, enthält das Modell die Ordnungen 2 und 1. Die Zellen in der Spalte "Saisonal" sind nur verfügbar, wenn für die Arbeitsdatei eine Periodizität definiert wurde.

- **Autoregressiv (p).** Die Anzahl autoregressiver Ordnungen im Modell. Autoregressive Ordnungen geben die zurückliegenden Werte der Zeitreihe an, die für die Vorhersage der aktuellen Werte verwendet werden. Eine autoregressive Ordnung von 2 gibt beispielsweise an, dass die Werte der Zeitreihe, die zwei Zeitperioden zurückliegt, für die Vorhersage der aktuellen Werte verwendet wird.
- **Differenz (d).** Gibt die Ordnung der Differenzierung an, die vor dem Schätzen der Modelle auf die Zeitreihe angewendet wurde. Differenzierung ist erforderlich, wenn Trends vorhanden sind. (Zeitreihen mit Trends sind normalerweise nichtstationär, und bei der ARIMA-Modellierung wird Stationarität angenommen.) Mithilfe der Differenzierung werden

die Effekte der Trends entfernt. Die Ordnung der Differenzierung entspricht dem Grad des Trends der Zeitreihe: Differenzierung erster Ordnung erklärt lineare Trends, Differenzierung zweiter Ordnung erklärt quadratische Trends usw.

- **Gleitender Durchschnitt (q).** Die Anzahl von Ordnungen des gleitenden Durchschnitts im Modell. Ordnungen des gleitenden Durchschnitts geben an, wie Abweichungen vom Mittelwert der Zeitreihe für zurückliegende Werte zum Vorhersagen der aktuellen Werte verwendet werden. Ordnungen des gleitenden Durchschnitts von 1 und 2 geben beispielsweise an, dass beim Vorhersagen der aktuellen Werte der Zeitreihe Abweichungen vom Mittelwert der Zeitreihe von den beiden letzten Zeitperioden berücksichtigt werden sollen.

**Saisonale Ordnungen.** Saisonale autoregressive Komponenten, Komponenten des gleitenden Durchschnitts und Differenzierungskomponenten entsprechen im Prinzip ihren nichtsaisonalen Gegenstücken. Bei saisonalen Ordnungen werden die Werte der aktuellen Zeitreihe jedoch von Werten zurückliegender Zeitreihen beeinflusst, die um eine oder mehrere saisonalen Perioden getrennt sind. Bei monatlichen Daten (saisonale Periode von 12) beispielsweise bedeutet eine saisonale Ordnung von 1, dass der Wert der aktuellen Zeitreihe durch den Zeitreihenwert beeinflusst wird, der 12 Perioden vor dem aktuellen liegt. Eine saisonale Ordnung von 1 entspricht bei monatlichen Daten einer nichtsaisonalen Ordnung von 12.

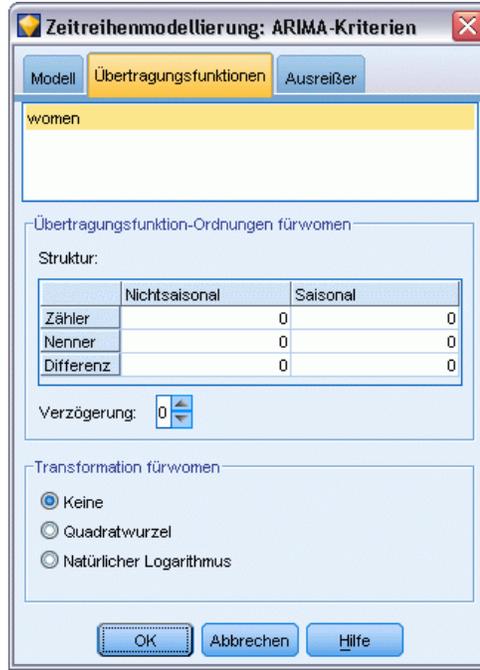
**Zieltransformation.** Sie können für jede Zielvariable eine Transformation angeben, die vor deren Modellierung durchgeführt werden soll. [Für weitere Informationen siehe Thema Reihentransformationen auf S. 437.](#)

- **Keine.** Es wird keine Transformation durchgeführt.
- **Quadratwurzel.** Es wird eine Quadratwurzeltransformation durchgeführt.
- **Natürlicher Logarithmus.** Es wird eine Transformation mit natürlichem Logarithmus durchgeführt.

**Konstante in Modell einschließen.** Der Einschluss einer Konstanten ist das Standardverfahren, sofern Sie nicht sicher wissen, dass der Gesamtmittelwert der Zeitreihe 0 ist. Bei der Anwendung von Differenzierung empfiehlt es sich, die Konstante auszuschließen.

## Transferfunktionen

Abbildung 13-14  
Dialogfeld "ARIMA-Kriterien", Registerkarte "Übertragungsfunktionen"



Auf der Registerkarte "Übertragungsfunktionen" können Sie Transferfunktionen für einige oder alle Eingabefelder definieren. Mithilfe von Transferfunktionen können Sie angeben, auf welche Weise frühere Werte der betreffenden Felder für die Vorhersage zukünftiger Werte der Ziel-Zeitreihe verwendet werden sollen.

Die Registerkarte wird nur dann angezeigt, wenn Eingabefelder (mit der Rolle *Eingabe*) angegeben wurden, entweder auf dem Typknoten oder auf der Registerkarte "Felder" des Zeitreihenknotens (wählen Sie Benutzerdefinierte Einstellungen verwenden – Eingaben). [Für weitere Informationen siehe Thema Festlegen der Feldrolle in Kapitel 4 in IBM SPSS Modeler 14.2-Quellen-, Prozess- und Ausgabeknoten.](#)

In der Liste oben werden alle Eingabefelder angezeigt. Die übrigen Informationen in diesem Dialogfeld hängen davon ab, welches Eingabefeld in der Liste ausgewählt wurde.

**Übertragungsfunktion-Ordnungen.** Geben Sie Werte für die verschiedenen Komponenten der Übertragungsfunktion in die entsprechenden Zellen des Strukturgitters ein. Alle Werte müssen nichtnegative Ganzzahlen sein. Bei Zähler- und Nennerkomponenten stellt der Wert die höchste Ordnung dar. Alle positiven niedrigeren Ordnungen werden in das Modell eingeschlossen. Darüber hinaus wird die Ordnung 0 bei Zählerkomponenten immer eingeschlossen. Wenn Sie beispielsweise 2 für den Zähler angeben, enthält das Modell die Ordnungen 2, 1 und 0. Wenn Sie 3 für den Nenner angeben, enthält das Modell die Ordnungen 3, 2 und 1. Die Zellen in der Spalte "Saisonal" sind nur verfügbar, wenn für die Arbeitsdatei eine Periodizität definiert wurde.

**Zähler.** Die Zählerordnung der Transferfunktion gibt an, welche zurückliegenden Werte aus der ausgewählten unabhängigen Zeitreihe (Prädiktor-Zeitreihe) zum Vorhersagen der aktuellen Werte der abhängigen Zeitreihe verwendet werden. Ein Zähler-Term von 1 gibt beispielsweise an, dass der Wert einer unabhängigen Zeitreihe, die eine Periode zurückliegt, und der aktuelle Wert der unabhängigen Zeitreihe zum Vorhersagen des aktuellen Werts der einzelnen abhängigen Zeitreihen verwendet werden.

**Nenner.** Die Nennerordnung der Transferfunktion gibt an, wie Abweichungen vom Mittelwert der Zeitreihe für zurückliegende Werte der ausgewählten unabhängigen Zeitreihe (Prädiktor-Zeitreihe) zum Vorhersagen der aktuellen Werte der abhängigen Zeitreihe verwendet werden. Ein Nenner-Term von 1 gibt beispielsweise an, dass beim Vorhersagen der aktuellen Werte für die einzelnen abhängigen Zeitreihen Abweichungen vom Mittelwert einer unabhängigen Zeitreihe berücksichtigt werden sollen, die eine Zeitperiode zurückliegt.

**Differenz.** Gibt die Ordnung der Differenzierung an, die vor dem Schätzen der Modelle auf die ausgewählte unabhängige Zeitreihe (Prädiktoren) angewendet wurde. Wenn Trends vorhanden sind, ist die Differenzierung erforderlich, um die Effekte der Trends zu entfernen.

**Saisonale Ordnungen.** Saisonale Zähler-, Nenner- und Differenzierungskomponenten entsprechen im Prinzip ihren nichtsaisonalen Gegenstücken. Bei saisonalen Ordnungen werden die Werte der aktuellen Zeitreihe jedoch von Werten zurückliegender Zeitreihen beeinflusst, die um eine oder mehrere saisonalen Perioden getrennt sind. Bei monatlichen Daten (saisonale Periode von 12) beispielsweise bedeutet eine saisonale Ordnung von 1, dass der Wert der aktuellen Zeitreihe durch den Zeitreihenwert beeinflusst wird, der 12 Perioden vor dem aktuellen liegt. Eine saisonale Ordnung von 1 entspricht bei monatlichen Daten einer nichtsaisonalen Ordnung von 12.

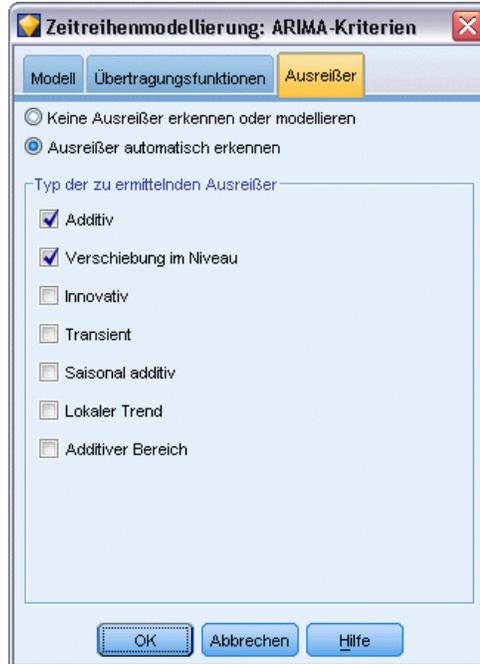
**Verzögerung** Wenn eine Verzögerung festgelegt wird, verzögert sich der Einfluss des Eingabefelds um die Anzahl der angegebenen Intervalle. Bei einer Verzögerung mit dem Wert 5 beeinflusst der Wert des Eingabefelds zum Zeitpunkt  $t$  die Vorhersagen erst nach dem Ablauf von fünf Perioden ( $t + 5$ ).

**Transformation.** Die Angabe einer Transferfunktion für ein Set von unabhängigen Variablen enthält auch eine optionale Transformation, die für diese Variablen ausgeführt werden soll.

- **Keine.** Es wird keine Transformation durchgeführt.
- **Quadratwurzel.** Es wird eine Quadratwurzeltransformation durchgeführt.
- **Natürlicher Logarithmus.** Es wird eine Transformation mit natürlichem Logarithmus durchgeführt.

## Umgang mit Ausreißern

Abbildung 13-15  
Dialogfeld "ARIMA-Kriterien," Registerkarte "Ausreißer"



Auf der Registerkarte "Ausreißer" ist eine Reihe von Möglichkeiten für die Behandlung von Ausreißern in den Daten verfügbar .

**Ausreißer nicht erkennen oder modellieren.** In der Standardeinstellung werden Ausreißer weder erkannt noch modelliert. Aktivieren Sie diese Option, um die Erkennung und Modellierung von Ausreißern zu deaktivieren.

**Ausreißer automatisch erkennen.** Aktivieren Sie diese Option, um eine automatische Erkennung von Ausreißern durchzuführen, und wählen Sie mindestens einen der gezeigten Ausreißertypen aus

**Typ der zu ermittelnden Ausreißer.** Wählen Sie die Ausreißertypen aus, die erkannt werden sollen. Folgende Typen werden unterstützt:

- Additiv (Standard)
- Verschiebung im Niveau (Standard)
- Neuerung
- Kurzlebig
- Saisonal additiv
- Lokaler Trend
- Additive Gruppe

Für weitere Informationen siehe Thema Ausreißer auf S. 435.

## **Generieren von Zeitreihenmodellen**

Dieser Abschnitt bietet einige allgemeine Informationen über bestimmte Aspekte beim Generieren von Zeitreihenmodellen:

- Generieren mehrerer Modelle
- Verwenden von Zeitreihenmodellen bei der Prognoseerstellung
- Erneute Schätzung und Vorhersage

Das generierte Modell-Nugget wird in einem separaten Kapitel beschrieben. [Für weitere Informationen siehe Thema Zeitreihen-Modell-Nugget auf S. 454.](#)

### **Generieren mehrerer Modelle**

Bei der Zeitreihenmodellierung in IBM® SPSS® Modeler wird für jedes Zielfeld ein einzelnes Modell (entweder ARIMA oder exponentielle Glättung) erstellt. Wenn mehrere Zielfelder vorhanden sind, generiert SPSS Modeler also mehrere Knoten in einem einzigen Vorgang. Dadurch wird Zeit gespart und Sie erhalten gleichzeitig die Möglichkeit, die Einstellungen für die einzelnen Modelle zu vergleichen.

Wenn Sie ein ARIMA-Modell und ein Modell mit exponentiellem Glätten für dasselbe Zielfeld vergleichen möchten, können Sie den Zeitreihenknoten separat ausführen und jeweils ein anderes Modell angeben.

### **Verwenden von Zeitreihenmodellen bei der Prognoseerstellung**

Bei der Erstellung von Zeitreihen wird eine bestimmte Reihe von geordneten Fällen, die so genannte Schätzungsspanne, verwendet, um ein Modell zu erstellen, das zur Vorhersage der zukünftigen Werte der Zeitreihe verwendet werden kann. Dieses Modell enthält Informationen zur verwendeten Zeitspanne, einschließlich des Intervalls. Um mithilfe dieses Modells Prognosen zu erstellen, müssen für Zielvariablen und Prädiktorvariablen jeweils dieselbe Zeitspanne und dieselben Intervallinformationen mit derselben Zeitreihe verwendet werden.

Beispiel: Angenommen, Anfang Januar möchten Sie den monatlichen Absatz von Produkt 1 für die ersten drei Monate des Jahres vorhersagen. Sie erstellen ein Modell unter Verwendung der tatsächlichen monatlichen Absatzdaten für Produkt 1 von Januar bis Dezember des Vorjahrs ("Jahr 1"), wobei Sie als Zeitintervall "Monate" verwenden. Anschließend können Sie mit diesem Modell den Absatz von Produkt 1 für die ersten drei Monate von Jahr 2 prognostizieren.

Im Prinzip könnten Sie eine Vorhersage für beliebig viele Monate erstellen, allerdings lässt die Effektivität des Modells immer mehr nach, je weiter sich die Prognose in die Zukunft erstreckt. Es wäre jedoch nicht möglich, eine Prognose für die ersten drei Wochen von Jahr 2 zu erstellen, da für die Erstellung des Modells das Intervall "Monate" verwendet wurde. Ebenso wäre es sinnlos, dieses Modell für die Vorhersage des Absatzes von Produkt 2 heranzuziehen. Zeitreihenmodelle sind immer nur für die Daten relevant, mit denen sie definiert wurden. [Für weitere Informationen siehe Thema Prognoseerstellung mit dem Zeitreihenknoten in Kapitel 15 in IBM SPSS Modeler 14.2- Anwendungshandbuch.](#)

## Erneute Schätzung und Vorhersage

Die Schätzperiode ist untrennbar mit dem generierten Modell verbunden. Daher werden bei Anwendung des aktuellen Modells auf neue Daten alle Werte außerhalb der Schätzperiode ignoriert. Zeitreihenmodelle müssen also jedes Mal, wenn neue Daten verfügbar sind, neu geschätzt werden; andere IBM® SPSS® Modeler-Modelle können dagegen zu Scoring-Zwecken ohne Veränderungen erneut angewendet werden.

Setzen wir das vorherige Beispiel fort: Angenommen, Anfang April von Jahr 2 liegen die tatsächlichen monatlichen Absatzdaten für Januar bis März vor. Wenn Sie jedoch das Modell, das Sie Anfang Januar erstellt haben, erneut anwenden, wird erneut eine Vorhersage für Januar bis März erstellt und die bekannten Absatzdaten für diesen Zeitraum werden ignoriert.

Die Lösung besteht darin, ein neues Modell auf der Grundlage der aktualisierten Ist-Daten zu erstellen. Sofern Sie die Prognoseparameter nicht ändern, kann das neue Modell zur Vorhersage der nächsten drei Monate, April bis Juni, verwendet werden. Wenn Sie noch Zugriff auf den Stream haben, der zur Erzeugung des ursprünglichen Modells verwendet wurde, können Sie einfach den Verweis auf die Quelldatei in diesem Stream durch einen Verweis auf die Datei ersetzen, die die aktualisierten Daten enthält, und den Stream dann erneut ausführen, um das neue Modell zu generieren. Wenn Sie nur noch das ursprüngliche Modell als Datei besitzen, können Sie damit einen Zeitreihenknoten erstellen, den Sie anschließend einem neuen Stream hinzufügen können, der einen Verweis zur aktualisierten Quelldatei enthält. Vorausgesetzt, dieser neue Stream stellt dem Zeitreihenknoten einen Zeitintervallknoten voran, bei dem das Intervall auf "Monate" gesetzt ist, wird durch die Ausführung dieses neuen Streams das erforderliche neue Modell erstellt. [Für weitere Informationen siehe Thema Erneutes Anwenden eines Zeitreihenmodells in Kapitel 15 in IBM SPSS Modeler 14.2- Anwendungshandbuch.](#)

## Zeitreihen-Modell-Nugget

Die Zeitreihenmodellierung erstellt eine Reihe neuer Felder mit dem Präfix \$TS-, nämlich:

\$TS-Spaltenname	Der vom Modell für die einzelnen Ziel-Zeitreihen vorhergesagte Wert.
\$TSLCI-Spaltenname	Die unteren Konfidenzintervalle für die einzelnen vorhergesagten Zeitreihen.*
\$TSUCI-Spaltenname	Die oberen Konfidenzintervalle für die einzelnen vorhergesagten Zeitreihen.*
\$TSNR-Spaltenname	Der Wert des Restrauschens für die einzelnen Spalten der Daten des generierten Modells.*
\$TS-Total	Der Gesamtwert der \$TS-Spaltenname-Werte für die betreffende Zeile.
\$TSLCI-Total	Der Gesamtwert der \$TSLCI-Spaltenname-Werte für die betreffende Zeile.*
\$TSUCI-Total	Der Gesamtwert der \$TSUCI-Spaltenname-Werte für die betreffende Zeile.*
\$TSNR-Total	Der Gesamtwert der \$TSNR-Spaltenname-Werte für die betreffende Zeile.*

\* Die Sichtbarkeit dieser Felder (z. B. in der Ausgabe von einem angegliederten Tabellenknoten) hängt von Optionen in der Registerkarte "Einstellungen" des Zeitreihenmodell-Nuggets ab. [Für weitere Informationen siehe Thema Zeitreihen – Modelleinstellungen auf S. 461.](#)

Abbildung 13-16  
Zeitreihen-Modell-Nugget, Registerkarte "Modell"

Anzahl der bei der Schätzung verwendeten Datensätze: 60

	Ziel	Modell	Prädiktoren	StationaryR**2	Q	df	Sig.
<input checked="" type="checkbox"/>	Market_1	Linearer Tre...	0	0,264	8,53	16,0	0,931
<input checked="" type="checkbox"/>	Market_2	Linearer Tre...	0	0,121	35,9	16,0	0,003
<input checked="" type="checkbox"/>	Market_3	Linearer Tre...	0	0,258	15,76	16,0	0,47
<input checked="" type="checkbox"/>	Market_4	Linearer Tre...	0	0,25	27,714	16,0	0,034
<input checked="" type="checkbox"/>	Market_5	Additives Wl...	0	0,544	11,888	15,0	0,688
<input checked="" type="checkbox"/>	Total	Linearer Tre...	0	0,049	27,616	16,0	0,035

Übersichtsstatistiken

	Statistik	StationaryR**2	Q	df	Sig.
ÜBERSICHT	MITTELVERT	0,247	21,235	15,833	0,36
ÜBERSICHT	SE	0,169	10,738	0,408	0,396
ÜBERSICHT	MINIMUM	0,049	8,53	15	0,003
ÜBERSICHT	MAXIMUM	0,544	35,9	16	0,931
ÜBERSICHT	PERZENTIL 5	0,049	8,53	15	0,003
ÜBERSICHT	PERZENTIL 10	0,049	8,53	15	0,003
ÜBERSICHT	PERZENTIL 25	0,103	11,048	15,75	0,026
ÜBERSICHT	PERZENTIL 50	0,254	21,688	16	0,252
ÜBERSICHT	PERZENTIL 75	0,334	29,761	16	0,749
ÜBERSICHT	PERZENTIL 90	0,544	35,9	16	0,931
ÜBERSICHT	PERZENTIL 95	0,544	35,9	16	0,931

Das Modell-Nugget vom Typ "Zeitreihe" zeigt Details der verschiedenen Modelle an, die für die einzelnen Zeitreiheneingaben im Erstellungsknoten der Zeitreihe ausgewählt wurden. Die Eingabe mehrerer Zeitreihen (z. B. Daten zu Produktlinien, Regionen oder Filialen) ist möglich und für jede Ziel-Zeitreihe wird ein separates Modell erstellt. Wenn der Umsatz in der Region Ost beispielsweise für ein ARIMA-Modell geeignet ist, die Region West sich jedoch nur für einen einfachen gleitenden Durchschnitt eignet, wird jede Region mit dem entsprechenden Modell gesort.

Die Standardausgabe zeigt für die einzelnen erstellten Modelle jeweils den Modelltyp, die Anzahl der angegebenen Prädiktoren und das Maß für die Anpassungsgüte (Standard: stationäres  $R$ -Quadrat) an. Wenn Sie Ausreißermethoden angegeben haben, ist eine Spalte vorhanden, in der die Anzahl der ermittelten Ausreißer angezeigt wird. Die Standardausgabe beinhaltet außerdem Spalten für Ljung-Box  $Q$ , Freiheitsgrade und Signifikanzwerte.

Außerdem können Sie die erweiterte Ausgabe auswählen; hierbei werden zusätzlich folgende Spalten angezeigt:

- $R$ -Quadrat
- RMSE (Root Mean Square Error, Wurzel der mittleren quadratischen Abweichung)

- MAPE (Mean Absolute Percentage Error, Mittlerer absoluter Fehler in Prozent)
- MAE (Mean Absolute Error, Mittlerer absoluter Fehler)
- MaxAPE (Maximum Absolute Percentage Error, Maximaler absoluter Fehler in Prozent)
- MaxAE (Maximum Absolute Error, Maximaler absoluter Fehler)
- Norm. BIC (Normalized Bayesian Information Criterion, Normalisiertes Bayes'sches Informationskriterium)

**Generieren.** Ermöglicht die Erzeugung eines Zeitreihen-Modellierungsknotens im Stream oder eines Modell-Nuggets in der Palette.

- **Erzeugen eines Modellierungsknotens.** Platziert einen Zeitreihen-Modellierungsknoten in einen Stream – mit den Einstellungen, die zum Erstellen dieses Modells verwendet wurden. Dies ist beispielsweise dann sinnvoll, wenn Sie einen Stream haben, in dem Sie diese Modelleinstellungen verwenden möchten, aber nicht mehr über den Modellierungsknoten verfügen, mit dem Sie sie generiert haben.
- **Modell zur Palette hinzufügen.** Platziert ein Modell-Nugget mit allen Zielen im Modell-Manager.

### Modell

Abbildung 13-17  
Schaltflächen "Alles markieren" und "Alle Markierungen aufheben"



**Kontrollkästchen aktivieren.** Wählen Sie aus, welche Modelle Sie beim Scoring verwenden möchten. Standardmäßig sind alle Kontrollkästchen aktiviert. Mit den Schaltflächen Alles markieren und Alle Markierungen aufheben werden alle Kontrollkästchen in einem einzigen Vorgang bearbeitet.

**Sortieren nach.** Mit dieser Option können Sie die Ausgabezeilen in aufsteigender oder absteigender Reihenfolge einer bestimmten Spalte der Anzeige sortieren. Mit der Option "Ausgewählt" wird die Ausgabe anhand einer oder mehrerer Zeilen sortiert, die über Kontrollkästchen ausgewählt wurden. Dies ist beispielsweise sinnvoll, um Zielfelder wie "Markt\_1" bis "Markt\_9" vor "Markt\_10" anzeigen zu lassen; bei der standardmäßigen Sortierreihenfolge wird nämlich "Markt\_10" unmittelbar nach "Markt\_1" angezeigt.

**Ansicht.** In der Standardansicht ("Einfach") wird die Grundmenge der Ausgabespalten angezeigt. Bei der Option "Erweitert" werden zusätzliche Spalten für die Maße der Anpassungsgüte angezeigt.

**Anzahl der bei der Schätzung verwendeten Datensätze.** Die Anzahl der Zeilen in der ursprünglichen Quelldatendatei.

**Ziel.** Die im Typknoten als Zielfelder (mit der Rolle *Ziel*) gekennzeichneten Felder.

**Modell.** Der für dieses Zielfeld verwendete Modelltyp.

**Prädiktoren.** Die Anzahl der für dieses Zielfeld verwendeten Prädiktoren (mit der Rolle *Eingabe*).

**Ausreißer.** Diese Spalte wird nur angezeigt, wenn Sie die automatische Erkennung von Ausreißern angefordert haben (bei den Expert Modeler- oder ARIMA-Kriterien). Der angezeigte Wert gibt die Anzahl der ermittelten Ausreißer an.

**R-Quadrat für stationären Teil.** Ein Maß, das den stationären Teil des Modells mit einem einfachen Mittelwert-Modell vergleicht. Dieses Maß ist dem gewöhnlichen R-Quadrat vorzuziehen, wenn ein Trend oder ein saisonales Muster vorliegt. R-Quadrat für den stationären Teil kann auch negativ sein, es nimmt Werte zwischen minus unendlich und 1 an. Negative Werte bedeuten, dass das betrachtete Modell schlechter ist als das Basismodell. Positive Werte bedeuten, dass das betrachtete Modell besser ist als das Basismodell.

**R-Quadrat.** Ein Maß für die Güte der Anpassung eines linearen Modells. Wird auch als Bestimmtheitsmaß bezeichnet. Es gibt den Anteil der Variation der abhängigen Variablen an, der durch das Regressionsmodell erklärt wird. Er liegt zwischen 0 und 1. Kleine Werte zeigen an, dass das Modell nicht gut zu den Daten passt.

**RMSE.** Steht für Root Mean Square Error, die Wurzel des mittleren quadratischen Fehlers. Die Quadratwurzel des mittleren Fehlerquadrats. Ein Maß dafür, wie stark eine abhängige Zeitreihe von ihrem durch das Modell vorhergesagten Niveau abweicht, und zwar ausgedrückt in derselben Maßeinheit wie die abhängige Zeitreihe.

**MAPE.** Mittlerer absoluter Fehler in Prozent. Ein Maß dafür, wie stark eine abhängige Zeitreihe von ihrem durch das Modell vorhergesagten Niveau abweicht. Es ist unabhängig von den verwendeten Maßeinheiten und kann daher verwendet werden, um Zeitreihen mit unterschiedlichen Einheiten zu vergleichen.

**MAE.** Mean Absolute Error, also mittlerer absoluter Fehler bzw. mittlerer Betrag des Fehlers. Er misst, wie stark die Zeitreihe von ihrem durch das Modell vorhergesagten Niveau abweicht. MAE wird in derselben Maßeinheit angegeben wie die ursprüngliche Zeitreihe.

**MaxAPE.** Maximaler absoluter Fehler in Prozent (Maximum Absolute Percentage Error, also maximaler Betrag des relativen Fehlers). Dies ist der größte vorhergesagte Fehler, ausgedrückt in Prozent. Dieses Maß hilft dabei, sich ein Worst-Case-Szenario für die Vorhersagen vorzustellen.

**MaxAE.** Maximaler absoluter Fehler (Maximum Absolute Error, also maximaler Betrag des Fehlers). Dies ist der größte vorhergesagte Fehler, ausgedrückt in derselben Maßeinheit wie die abhängige Zeitreihe. Genau wie MaxAPE hilft er dabei, sich ein Worst-Case-Szenario für die Vorhersagen vorzustellen. Der maximale absolute Fehler und der maximale absolute Fehler in Prozent können an verschiedenen Punkten in der Zeitreihe auftreten, beispielsweise wenn der absolute Fehler für einen großen Zeitreihenwert geringfügig größer ist als der absolute Fehler für einen kleinen Zeitreihenwert. In diesem Fall tritt der maximale absolute Fehler beim größeren Zeitreihenwert und der maximale absolute Fehler in Prozent beim kleineren Zeitreihenwert auf.

**Normalisiertes BIC.** Normalisiertes Bayes-Informationskriterium (BIC). Ein allgemeines Maß der insgesamt erreichten Güte der Anpassung, das auch die Komplexität des Modells zu berücksichtigen versucht. Es ist ein Wert, der auf dem mittleren quadratischen Fehler beruht und eine Penalisierung für die Anzahl der Modellparameter und die Länge der Zeitreihe enthält. Die Penalisierung neutralisiert die Überlegenheit von Modellen mit einer größeren Anzahl von Parametern und macht die Statistik damit gut vergleichbar für verschiedene Modelle derselben Zeitreihe.

**Q.** Die Ljung-Box-Q-Statistik. Ein Test der Zufälligkeit der Restfehler in diesem Modell.

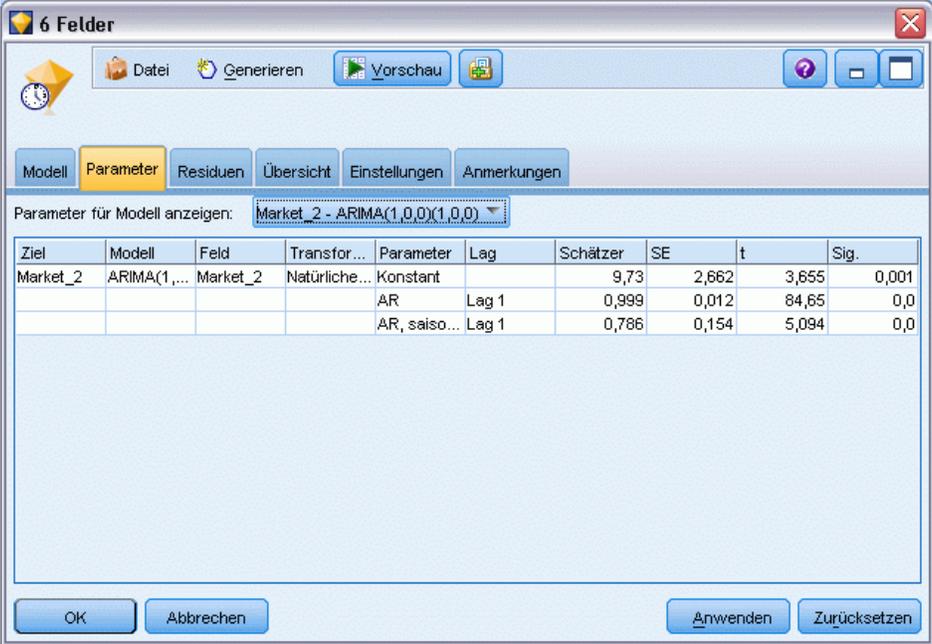
**df.** Freiheitsgrade. Die Anzahl der Modellparameter, die bei der Schätzung eines bestimmten Ziels frei variieren können.

**Sig.** Signifikanzwert der Ljung-Box-Statistik. Ein Signifikanzwert von weniger als 0,05 deutet darauf hin, dass die Restfehler nicht zufällig sind.

**Auswertungsstatistik.** Dieser Abschnitt enthält verschiedene Auswertungsstatistiken für die verschiedenen Spalten, darunter Mittelwert, Minimum, Maximum und Perzentilwerte.

## Zeitreihen – Modellparameter

Abbildung 13-18  
Zeitreihen – Registerkarte “Parameter”



Parameter für Modell anzeigen: Market\_2 - ARIMA(1,0,0)(1,0,0)

Ziel	Modell	Feld	Transfor...	Parameter	Lag	Schätzer	SE	t	Sig.
Market_2	ARIMA(1,...	Market_2	Natürliche...	Konstant		9,73	2,662	3,655	0,001
				AR	Lag 1	0,999	0,012	84,65	0,0
				AR, saiso...	Lag 1	0,786	0,154	5,094	0,0

Auf der Registerkarte “Parameter” sind Details verschiedener Parameter aufgeführt, die zum Erstellen eines ausgewählten Modells verwendet wurden.

**Parameter für Modell anzeigen.** Wählen Sie das Modell aus, für das die Parameterdetails angezeigt werden sollen.

**Ziel.** Der Name des von diesem Modell vorhergesagten Zielfelds (mit der Rolle *Ziel*).

**Modell.** Der für dieses Zielfeld verwendete Modelltyp.

**Feld (Nur ARIMA-Modelle).** Enthält genau einen Eintrag für jede der im Modell verwendeten Variablen. Dabei wird zunächst das Ziel angegeben und dann die Prädiktoren (sofern vorhanden).

**Transformation.** Zeigt ggf. an, welche Art von Transformation für dieses Feld festgelegt wurde, bevor das Modell erstellt wurde.

**Parameter.** Der Modellparameter, für den die folgenden Details angezeigt werden:

- **Feld (Nur ARIMA-Modelle).** Gibt ggf. die für diesen Parameter im Modell berücksichtigten Lags an.
- **Schätzer.** Der Parameterschätzer. Dieser Wert wird bei der Berechnung des Vorhersagewerts und der Konfidenzintervalle für das Zielfeld verwendet.
- **SE.** Der Standardfehler des Parameterschätzers.
- **t.** Der Wert des Parameterschätzers dividiert durch den Standardfehler.
- **Sig.** Das Signifikanzniveau des Parameterschätzers. Werte über 0,05 werden als statistisch nicht signifikant betrachtet.

## Zeitreihen – Modellresiduen

Abbildung 13-19

Zeitreihenmodell, Registerkarte "Residuen," Anzeige "ACF" und "PACF"

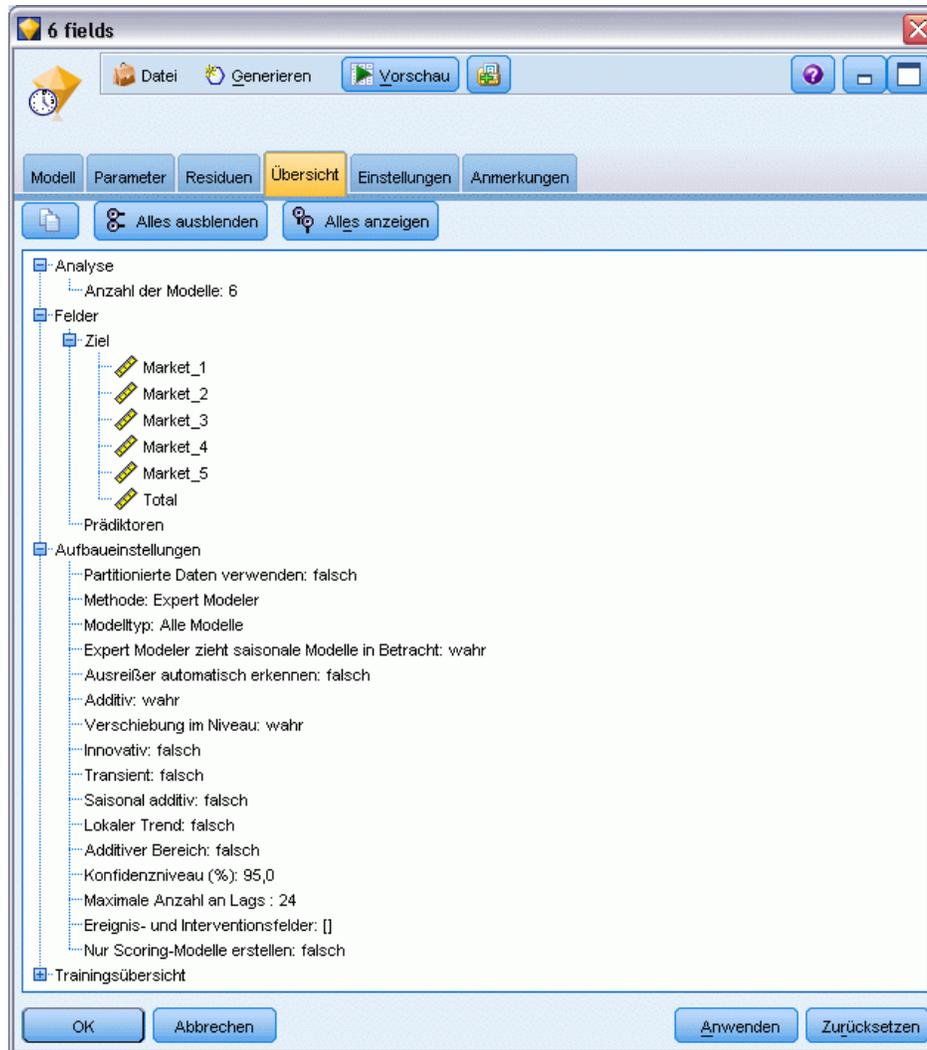


Auf der Registerkarte "Residuen" werden die Autokorrelationsfunktion (ACF) und die partielle Autokorrelationsfunktion (PACF) der Residuen (der Differenzen zwischen den erwarteten Werten und den Ist-Werten) für die einzelnen erstellten Modelle angezeigt. [Für weitere Informationen siehe Thema Autokorrelation und partielle Autokorrelationsfunktionen auf S. 436.](#)

**Diagramm für Modell anzeigen.** Wählen Sie das Modell, für das Sie die die Residuen-ACF und -PACF anzeigen möchten.

## Zeitreihen – Modellübersicht

Abbildung 13-20  
Zeitreihenmodell – Registerkarte "Übersicht"



Auf der Registerkarte "Übersicht" eines Modell-Nuggets finden Sie Informationen zum Modell selbst (*Analyse*), den im Modell verwendeten Feldern (*Felder*), den beim Erstellen des Modells verwendeten Einstellungen (*Aufbaueinstellungen*) und zum Modelltraining (*Trainingsübersicht*).

Beim ersten Durchsuchen des Knotens sind die Ergebnisse der Registerkarte "Übersicht" reduziert. Um die für Sie relevanten Ergebnisse anzuzeigen, vergrößern Sie sie mithilfe des Erweiterungssteuerelements links neben dem betreffenden Element oder klicken Sie auf die Schaltfläche Alles anzeigen, um alle Ergebnisse anzuzeigen. Um die Ergebnisse nach der Betrachtung wieder auszublenden, können Sie mit dem Erweiterungssteuerelement die

gewünschten Ergebnisse reduzieren. Alternativ können Sie mit der Schaltfläche Alles ausblenden alle Ergebnisse ausblenden.

**Analyse.** Zeigt Informationen zum jeweiligen Modell an.

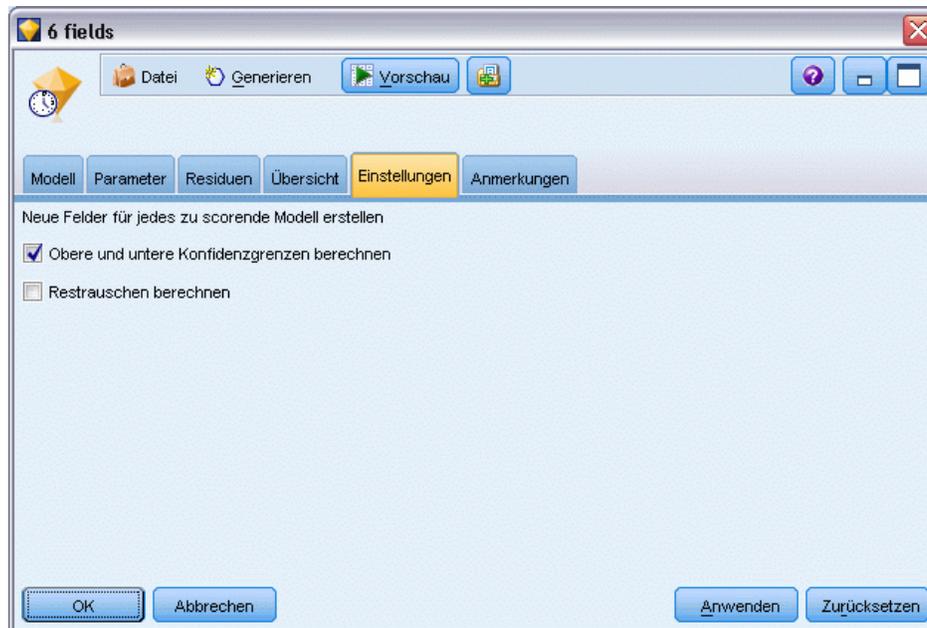
**Felder.** Listet die Felder auf, die bei der Erstellung des Modells als Ziel bzw. als Eingaben verwendet werden.

**Aufbaueinstellungen.** Enthält Informationen zu den beim Aufbau des Modells verwendeten Einstellungen.

**Trainingsübersicht.** In diesem Abschnitt werden folgende Informationen angezeigt: der Typ des Modells, der für seine Erstellung verwendete Stream, der Benutzer, der ihn erstellt hat, der Erstellungszeitpunkt und die für den Aufbau des Modells benötigte Zeit.

## Zeitreihen – Modelleinstellungen

Abbildung 13-21  
Zeitreihenmodell – Registerkarte "Einstellungen"



Auf der Registerkarte "Einstellungen" können Sie angeben, welche zusätzlichen Felder durch den Modellierungsvorgang erstellt werden sollen.

**Neue Felder für jedes zu scorende Modell erstellen.** Mit dieser Option können Sie die neuen Felder angeben, die für jedes zu scorende Modell erstellt werden sollen.

- **Obere und untere Konfidenzgrenzen berechnen.** Wenn diese Option aktiviert ist, werden neue Felder (mit den Standardpräfixen \$TSLCI- und \$TSUCI-) für die obere bzw. untere Grenze des Konfidenzintervalls für die einzelnen Zielfelder erstellt, zusammen mit den Gesamtsummen dieser Werte.
- **Restrauschen berechnen.** Wenn diese Option aktiviert ist, wird ein neues Feld (mit dem Standardpräfix \$TSNR-) für die Modellresiduen der einzelnen Zielfelder erstellt, zusammen mit einer Gesamtsumme dieser Werte.

# ***Selbstlern-Antwortknotenmodelle***

## ***SLRM-Knoten***

Mit dem Knoten für das Selbstlern-Antwortmodell (**Self-Learning Response Model**, SLRM) können Sie ein Modell erstellen, das Sie während der Erweiterung des Daten-Sets ständig aktualisieren bzw. neu schätzen können, ohne dass Sie das Modell jedes Mal anhand des vollständigen Daten-Sets neu erstellen müssen. Dies ist beispielsweise dann nützlich, wenn Sie mehrere Produkte führen und ermitteln möchten, welches Produkt ein Kunde mit der größten Wahrscheinlichkeit kauft, wenn Sie es ihm anbieten. Mit diesem Modell können Sie prognostizieren, welche Angebote für die Kunden am geeignetsten sind und mit welcher Wahrscheinlichkeit die Angebote angenommen werden.

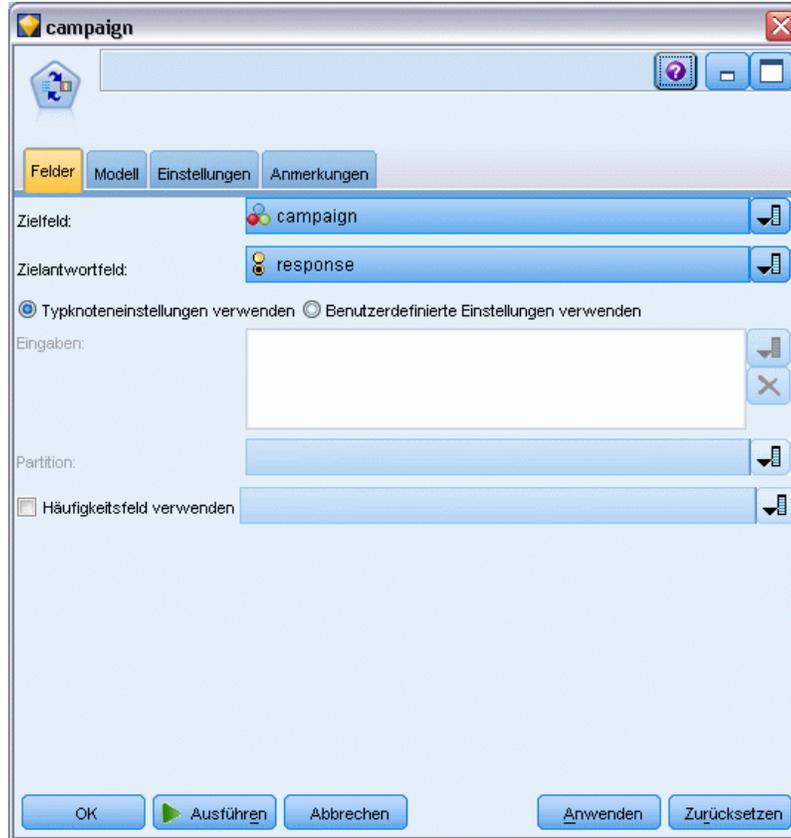
Das Modell kann zunächst mit einem kleinen Daten-Set mit zufällig ausgewählten Angeboten und den Reaktionen auf diese Angebote erstellt werden. Wenn das Daten-Set größer wird, kann das Modell aktualisiert werden, wodurch es besser in der Lage ist, die geeignetsten Angebote für Kunden und die Wahrscheinlichkeit, dass diese angenommen werden, auf der Grundlage anderer Eingabefelder, wie Alter, Geschlecht, Beruf und Einkommen, zu prognostizieren. Die verfügbaren Angebote können geändert werden, indem sie im Knotendialogfeld hinzugefügt bzw. entfernt werden. Es ist also nicht erforderlich, das Zielfeld des Daten-Sets zu ändern.

Bei einer Kombination mit IBM® SPSS® Collaboration and Deployment Services können Sie automatische regelmäßige Aktualisierungen des Modells einrichten. Dieser Vorgang, der keine Überwachung oder Eingriffe von menschlicher Seite erfordert, stellt eine flexible und kostengünstige Lösung für Organisationen und Anwendungen dar, bei denen individuelle Eingriffe durch einen Data-Mining-Experten nicht möglich oder erforderlich sind.

**Beispiel.** Ein Finanzinstitut möchte profitablere Ergebnisse erzielen, indem jedem Kunden das Angebot zugeordnet wird, das mit der größten Wahrscheinlichkeit angenommen wird. Mit dem Selbstlernmodell können Sie auf der Grundlage früherer Werbeaktionen die Eigenschaften der Kunden ermitteln, die mit hoher Wahrscheinlichkeit positiv reagieren werden, und das Modell in Echtzeit auf der Grundlage der jeweils aktuellsten Kundenreaktionen aktualisieren. [Für weitere Informationen siehe Thema Erstellen von Angeboten für Kunden \(Selbstlernfunktion\) in Kapitel 17 in IBM SPSS Modeler 14.2- Anwendungshandbuch.](#)

## Feldoptionen für den SLRM-Knoten

Abbildung 14-1  
Dialogfeld des SLRM-Knotens, Registerkarte "Felder"



Vor der Ausführung eines SLRM-Knotens müssen Sie die Ziel- und Zielantwortfelder auf der Registerkarte "Felder" des Knotens angeben.

**Zielfeld.** Wählen Sie das Zielfeld aus der Liste, z. B. ein nominales (Set-) Feld, das die verschiedenen Produkte enthält, die Sie den Kunden anbieten möchten.

*Anmerkung:* Das Zielfeld muss den Speichertyp "Zeichenkette" aufweisen (nicht "Numerisch").

**Zielantwortfeld.** Wählen Sie das Zielantwortfeld aus der Liste aus. Beispiele: "Angenommen" oder "Abgelehnt".

*Anmerkung:* Bei diesem Feld muss es sich um ein Flag-Feld handeln. Der Wahr-Wert des Flag zeigt die Annahme, der Falsch-Wert die Ablehnung des Angebots an.

Die restlichen Felder in diesem Dialogfeld sind Standardfelder, die überall in IBM® SPSS® Modeler verwendet werden. [Für weitere Informationen siehe Thema Feldoptionen der Modellierungsknoten in Kapitel 3 auf S. 36.](#)

*Anmerkung:* Wenn die Quelldaten Bereiche beinhalten, die als stetige Eingabefelder (numerischer Bereich) verwendet werden sollen, müssen Sie sicherstellen, dass die Metadaten die Details für Maximum und Minimum für jeden Bereich beinhalten.

## Modelloptionen für den SLRM-Knoten

Abbildung 14-2  
Dialogfeld des SLRM-Knotens, Registerkarte "Modell"



**Modellname.** Sie können den Modellnamen automatisch basierend auf den Ziel- oder ID-Feldnamen (oder dem Modelltyp in Fällen, in denen kein solches Feld angegeben ist) generieren oder einen benutzerdefinierten Namen eingeben.

**Partitionierte Daten verwenden.** Wenn ein Partitionsfeld definiert ist, gewährleistet diese Option, dass nur Daten aus der Trainingspartition für die Modellerstellung verwendet werden. [Für weitere Informationen siehe Thema Partitionsknoten in Kapitel 4 in IBM SPSS Modeler 14.2- Quellen-, Prozess- und Ausgabeknoten.](#)

**Training des bestehenden Modells fortsetzen.** In der Standardeinstellung wird bei jeder Ausführung eines Modellierungsknotens ein völlig neues Modell erstellt. Bei Auswahl dieser Option wird das Training mit dem letzten, vom Knoten erfolgreich aufgebauten Modell fortgesetzt. Dies ermöglicht die Aktualisierung eines bestehenden Modells, ohne dass ein Zugriff auf die ursprünglichen Daten erforderlich ist, und kann zu einer wesentlich schnelleren Leistung führen, da *ausschließlich* die neuen bzw. aktualisierten Datensätze in den Stream eingespeist werden. Details zum vorherigen Modell werden zusammen mit dem Modellierungsknoten gespeichert, wodurch diese Option auch dann verwendet werden kann, wenn das vorherige Modell-Nugget nicht mehr im Stream oder in der Modelpalette verfügbar ist.

**Werte des Zielfelds** Standardmäßig ist dieses Feld auf Alle verwenden gesetzt, was bedeutet, dass ein Modell erstellt wird, das alle Angebote enthält, die dem ausgewählten Zielfeldwert zugeordnet sind. Wenn Sie ein Modell erstellen möchten, das nur einige der Angebote des Zielfelds enthält, klicken Sie auf Angeben und verwenden Sie die Schaltflächen Hinzufügen, Bearbeiten und Löschen, um die Namen der Angebote, für die ein Modell erstellt werden soll, einzugeben bzw. abzuändern. Wenn Sie beispielsweise ein Ziel ausgewählt haben, das alle von Ihnen gelieferten Produkte auflistet, können Sie die angebotenen Produkte mit diesem Feld auf einige wenige einschränken, die Sie hier eingeben.

**Modellauswertung.** Die Felder in diesem Fenster sind dahingehend unabhängig vom Modell, dass sie das Scoring nicht beeinflussen. Stattdessen ermöglichen Sie Ihnen, eine visuelle Darstellung davon zu erstellen, wie gut das Modell Ergebnisse vorhersagt.

*Anmerkung:* Um die Ergebnisse der Modellauswertung im Modell-Nugget anzuzeigen, müssen Sie auch das Kontrollkästchen Modellevaluation anzeigen aktivieren.

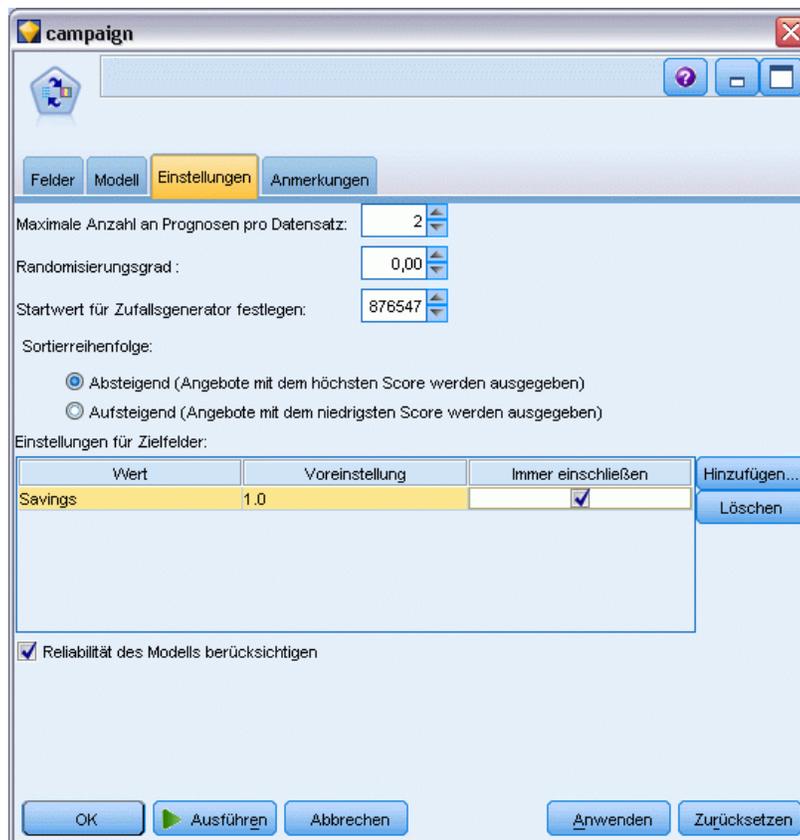
- **Modellauswertung einschließen.** Wählen Sie dieses Feld aus, um Diagramme zu erstellen, die die prognostizierte Genauigkeit des Modells für jedes ausgewählte Angebot zeigen.
- **Startwert für Zufallsgenerator festlegen.** Bei der Schätzung der Genauigkeit eines Modells auf der Grundlage eines Zufallsprozentsatzes können Sie mit dieser Option dieselben Ergebnisse in einer anderen Sitzung replizieren. Wenn Sie den vom Zufallszahlengenerator verwendeten Startwert angeben, stellen Sie sicher, dass bei jeder Ausführung des Knotens dieselben Datensätze zugewiesen werden. Gewünschten Startwert eingeben. Wenn diese Option nicht ausgewählt ist, wird bei jeder Ausführung des Knotens eine andere Stichprobe generiert.
- **Simulierte Stichprobengröße.** Geben Sie die Anzahl der Datensätze an, die bei der Bewertung des Modells in der Stichprobe verwendet werden sollen. Der Standardwert lautet 100.
- **Anzahl der Iterationen.** Mit dieser Option können Sie die Erstellung der Modellauswertung nach der angegebenen Anzahl an Iterationen beenden. Dient zur Angabe der maximalen Anzahl der Iterationen; der Standardwert lautet "20".

*Anmerkung:* Beachten Sie, dass große Stichproben und eine hohe Anzahl von Iterationen den Zeitbedarf für die Erstellung des Modells erhöhen.

**Modellevaluation anzeigen.** Wählen Sie diese Option, um eine grafische Darstellung der Ergebnisse im Modell-Nugget anzuzeigen.

## Einstellungsoptionen für den SLRM-Knoten

Abbildung 14-3  
Dialogfeld des SLRM-Knotens, Registerkarte "Einstellungen"



Mit den Einstellungsoptionen für den Knoten können Sie die Feinabstimmung der Modellerstellung vornehmen.

**Maximale Anzahl an Prognosen pro Datensatz.** Mit dieser Option können Sie die Anzahl der Vorhersagen für die einzelnen Datensätze im Daten-Set einschränken. Der Standardwert ist 3.

Sie könnten beispielsweise sechs Angebote haben (z. B. Sparbuch, Hypothek, Autokredit, Rentensparplan, Kreditkarte und Versicherung), möchten jedoch nur die beiden empfehlenswertesten Angebote ermitteln. In diesem Fall würden Sie das Feld auf 2 setzen. Wenn Sie das Modell erstellen und mit einer Tabelle verknüpfen, sehen Sie zwei Prognosespalten (und die zugehörige Konfidenz für die Wahrscheinlichkeit, dass das Angebot angenommen wird) pro Datensatz. Die Vorhersagen können jedes der sechs möglichen Angebote enthalten.

**Randomisierungsgrad.** Um Verzerrungen zu vermeiden – beispielsweise in einem kleinen oder unvollständigen Daten-Set – und alle potenziellen Angebote gleich zu behandeln, können Sie einen Randomisierungsgrad für die Angebotsauswahl und die Wahrscheinlichkeit, dass sie als empfohlene Angebote aufgenommen werden, angeben. Die Randomisierung wird als Prozentsatz ausgedrückt, der als Dezimalwert zwischen 0,0 (keine Randomisierung) und 1,0 (vollständig zufällig) angezeigt wird. Der Standardwert ist 0,0.

**Startwert für Zufallsgenerator festlegen.** Wenn Sie einen Randomisierungsgrad für die Auswahl eines Angebots angeben, können Sie mit dieser Option dieselben Ergebnisse in einer anderen Sitzung replizieren. Wenn Sie den vom Zufallszahlengenerator verwendeten Startwert angeben, stellen Sie sicher, dass bei jeder Ausführung des Knotens dieselben Datensätze zugewiesen werden. Gewünschten Startwert eingeben. Wenn diese Option nicht ausgewählt ist, wird bei jeder Ausführung des Knotens eine andere Stichprobe generiert.

*Hinweis:* Bei Verwendung der Option Startwert für Zufallsgenerator festlegen mit Datensätzen, die aus einer Datenbank eingelesen wurden, ist möglicherweise vor der Stichprobenziehung ein Sortierknoten erforderlich, um zu gewährleisten, dass bei jeder Ausführung des Knotens dasselbe Ergebnis erzielt wird. Dies liegt daran, dass der Startwert für den Zufallsgenerator von der Reihenfolge der Datensätze abhängt, die in relationalen Datenbanken nicht unbedingt gleich bleibt. [Für weitere Informationen siehe Thema Sortierknoten in Kapitel 3 in IBM SPSS Modeler 14.2- Quellen-, Prozess- und Ausgabeknoten.](#)

**Sortierreihenfolge.** Wählen Sie die Reihenfolge aus, in der die Angebote im erstellten Modell angezeigt werden sollen:

- **Absteigend.** Das Modell zeigt die Angebote mit den höchsten Scores zuerst an. Hierbei handelt es sich um die Angebote, die mit der größten Wahrscheinlichkeit angenommen werden.
- **Aufsteigend.** Das Modell zeigt die Angebote mit den niedrigsten Scores zuerst an. Hierbei handelt es sich um die Angebote, die mit der größten Wahrscheinlichkeit abgelehnt werden. Dies kann beispielsweise nützlich sein, wenn Sie ermitteln möchten, welche Kunden aus einer Marketingkampagne für ein bestimmtes Angebot gestrichen werden sollen.

**Einstellungen für Zielfelder.** Beim Erstellen eines Modells kann es bestimmte Datenaspekte geben, die Sie aktiv begünstigen bzw. entfernen möchten. Wenn Sie beispielsweise ein Modell erstellen, das das beste Finanzangebot auswählt, für das beim Kunden geworben werden soll, möchten Sie möglicherweise sicherstellen, dass ein bestimmtes Angebot immer aufgenommen wird, unabhängig davon, wie gut sein Score bei den einzelnen Kunden ist.

Um ein Angebot in diesem Fenster einzuschließen und seinen Präferenzgrad zu bearbeiten, klicken Sie auf Hinzufügen, geben Sie den Namen des Angebots ein (z. B. "Sparbuch" oder "Hypothek" und klicken Sie auf OK.

- **Wert.** Hier wird der Name des hinzugefügten Angebots angezeigt.
- **Voreinstellung.** Gibt den Präferenzgrad an, der auf das Angebot angewendet werden soll. Die Präferenz wird als Prozentsatz ausgedrückt, der als Dezimalwert zwischen 0,0 (keine Präferenz) und 1,0 (höchste Präferenz) angezeigt wird. Der Standardwert ist 0,0.
- **Immer einschließen.** Aktivieren Sie dieses Kontrollkästchen, um sicherzustellen, dass ein bestimmtes Angebot immer in die Vorhersagen eingeschlossen wird.

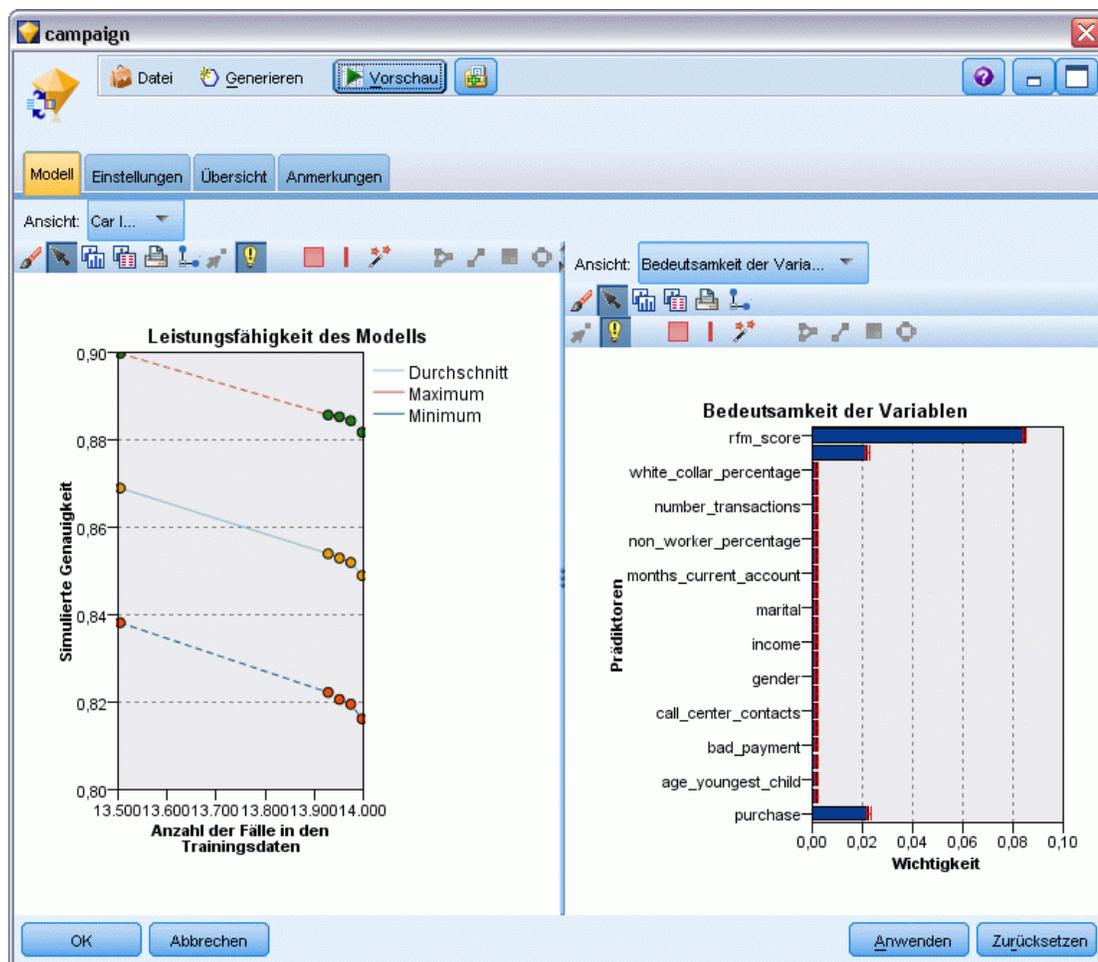
*Hinweis:* Wenn das Feld Voreinstellung auf 0.0 gesetzt ist, wird die Einstellung Immer einschließen ignoriert.

**Reliabilität des Modells berücksichtigen.** Ein gut strukturiertes Modell, das reich an Daten ist und durch mehrere erneute Generierungen eine Feinabstimmung erfahren hat, sollte stets genauere Ergebnisse liefern als ein ganz neues Modell mit wenigen Daten. Wenn Sie die höhere Reliabilität des ausgereifteren Modells nutzen möchten, aktivieren Sie dieses Kontrollkästchen.

## SLRM-Modell-Nuggets

*Anmerkung:* Auf dieser Registerkarte werden nur dann Ergebnisse angezeigt, wenn auf der Registerkarte “Modelloptionen” sowohl Modellauswertung einschließen als auch Modellevaluation anzeigen ausgewählt wurde.

Abbildung 14-4  
Anzeige des SLRM-Modell-Nuggets



Wenn Sie einen Stream ausführen, der ein SLRM-Modell enthält, berechnet der Knoten Schätzungen für die Genauigkeit der Vorhersagen für jeden Wert des Zielfelds (Angebot) und die Wichtigkeit der einzelnen verwendeten Prädiktoren.

*Anmerkung:* Wenn Sie auf der Registerkarte “Modell” des Modellknotens die Option Training des bestehenden Modells fortsetzen ausgewählt haben, werden die im Modell-Nugget angezeigten Informationen bei jeder erneuten Generierung des Modells aktualisiert.

Bei mit IBM® SPSS® Modeler 12.0 oder höher erstellten Modellen ist die Registerkarte “Modell” des Modell-Nuggets in zwei Spalten unterteilt:

**Linke Spalte.**

- **Ansicht.** Wenn Sie mehrere Angebote haben, wählen Sie das Angebot aus, für das Ergebnisse angezeigt werden sollen.
- **Leistungsfähigkeit des Modells.** Zeigt die geschätzte Modellgenauigkeit für jedes Angebot an. Das Test-Set wird durch Simulation generiert.

**Rechte Spalte.**

- **Ansicht.** Wählen Sie aus, ob Details zu Assoziation mit Antwort oder zu Bedeutsamkeit der Variablen angezeigt werden sollen.
- **Assoziation mit Antwort.** Zeigt die Assoziation (Korrelation) der einzelnen Prädiktoren mit der Zielvariablen an.
- **Bedeutsamkeit des Prädiktors.** Gibt die relative Wichtigkeit der einzelnen Prädiktoren für die Schätzung des Modells an. Normalerweise ist es sinnvoll, die Modellierungsbemühungen auf die wichtigsten Prädiktoren zu konzentrieren und diejenigen Prädiktoren zu verwerfen bzw. zu ignorieren, die die geringste Bedeutung haben. Dieses Diagramm kann auf dieselbe Weise interpretiert werden wie für andere Modelle, die die Bedeutsamkeit der Prädiktoren anzeigen. Allerdings wird bei SLRM das Diagramm vom SLRM-Algorithmus durch Simulation erzeugt. Dies geschieht, indem jeder Prädiktor nacheinander aus dem Modell entfernt und angezeigt wird, wie dies die Modellgenauigkeit beeinflusst. [Für weitere Informationen siehe Thema Bedeutsamkeit des Prädiktors in Kapitel 3 auf S. 54.](#)

## ***SLRM-Modell – Einstellungen***

Auf der Registerkarte “Einstellungen” für ein SLRM-Modell-Nugget werden Optionen zum Ändern des erstellten Modells angegeben. Beispielsweise können Sie mit dem SLRM-Knoten unter Verwendung derselben Daten und Einstellungen mehrere verschiedene Modelle erstellen und anschließend diese Registerkarte bei jedem Modell verwenden, um die Einstellungen leicht abzuändern und zu ermitteln, welche Auswirkungen dies auf die Ergebnisse hat.

*Anmerkung:* Diese Registerkarte ist erst verfügbar, nachdem das Modell-Nugget einem Stream hinzugefügt wurde.

Abbildung 14-5  
Dialogfeld des SLRM-Modell-Nuggets, Registerkarte "Einstellungen"

campaign

Datei Generieren Vorschau

Modell **Einstellungen** Übersicht Anmerkungen

Maximale Anzahl an Prognosen pro Datensatz: 2

Randomisierungsgrad: 0,00

Startwert für Zufallsgenerator festlegen: 876547

Sortierreihenfolge:

Absteigend (Angebote mit dem höchsten Score werden ausgegeben)

Aufsteigend (Angebote mit dem niedrigsten Score werden ausgegeben)

Einstellungen für Zielfelder:

Wert	Voreinstellung	Immer einschließen
Savings	1.0	<input checked="" type="checkbox"/>

Hinzufügen...  
Löschen

Reliabilität des Modells berücksichtigen

OK Abbrechen Anwenden Zurücksetzen

**Maximale Anzahl an Prognosen pro Datensatz.** Mit dieser Option können Sie die Anzahl der Vorhersagen für die einzelnen Datensätze im Daten-Set einschränken. Der Standardwert ist 3.

Sie könnten beispielsweise sechs Angebote haben (z. B. Sparbuch, Hypothek, Autokredit, Rentensparplan, Kreditkarte und Versicherung), möchten jedoch nur die beiden empfehlenswertesten Angebote ermitteln. In diesem Fall würden Sie das Feld auf 2 setzen. Wenn Sie das Modell erstellen und mit einer Tabelle verknüpfen, sehen Sie zwei Prognosespalten (und die zugehörige Konfidenz für die Wahrscheinlichkeit, dass das Angebot angenommen wird) pro Datensatz. Die Vorhersagen können jedes der sechs möglichen Angebote enthalten.

**Randomisierungsgrad.** Um Verzerrungen zu vermeiden – beispielsweise in einem kleinen oder unvollständigen Daten-Set – und alle potenziellen Angebote gleich zu behandeln, können Sie einen Randomisierungsgrad für die Angebotsauswahl und die Wahrscheinlichkeit, dass sie als empfohlene Angebote aufgenommen werden, angeben. Die Randomisierung wird als Prozentsatz ausgedrückt, der als Dezimalwert zwischen 0,0 (keine Randomisierung) und 1,0 (vollständig zufällig) angezeigt wird. Der Standardwert ist 0,0.

**Startwert für Zufallsgenerator festlegen.** Wenn Sie einen Randomisierungsgrad für die Auswahl eines Angebots angeben, können Sie mit dieser Option dieselben Ergebnisse in einer anderen Sitzung replizieren. Wenn Sie den vom Zufallszahlengenerator verwendeten Startwert angeben, stellen Sie sicher, dass bei jeder Ausführung des Knotens dieselben Datensätze zugewiesen werden. Gewünschten Startwert eingeben. Wenn diese Option nicht ausgewählt ist, wird bei jeder Ausführung des Knotens eine andere Stichprobe generiert.

*Hinweis:* Bei Verwendung der Option Startwert für Zufallsgenerator festlegen mit Datensätzen, die aus einer Datenbank eingelesen wurden, ist möglicherweise vor der Stichprobenziehung ein Sortierknoten erforderlich, um zu gewährleisten, dass bei jeder Ausführung des Knotens dasselbe Ergebnis erzielt wird. Dies liegt daran, dass der Startwert für den Zufallsgenerator von der Reihenfolge der Datensätze abhängt, die in relationalen Datenbanken nicht unbedingt gleich bleibt. [Für weitere Informationen siehe Thema Sortierknoten in Kapitel 3 in IBM SPSS Modeler 14.2- Quellen-, Prozess- und Ausgabeknoten.](#)

**Sortierreihenfolge.** Wählen Sie die Reihenfolge aus, in der die Angebote im erstellten Modell angezeigt werden sollen:

- **Absteigend.** Das Modell zeigt die Angebote mit den höchsten Scores zuerst an. Hierbei handelt es sich um die Angebote, die mit der größten Wahrscheinlichkeit angenommen werden.
- **Aufsteigend.** Das Modell zeigt die Angebote mit den niedrigsten Scores zuerst an. Hierbei handelt es sich um die Angebote, die mit der größten Wahrscheinlichkeit abgelehnt werden. Dies kann beispielsweise nützlich sein, wenn Sie ermitteln möchten, welche Kunden aus einer Marketingkampagne für ein bestimmtes Angebot gestrichen werden sollen.

**Einstellungen für Zielfelder.** Beim Erstellen eines Modells kann es bestimmte Datenaspekte geben, die Sie aktiv begünstigen bzw. entfernen möchten. Wenn Sie beispielsweise ein Modell erstellen, das das beste Finanzangebot auswählt, für das beim Kunden geworben werden soll, möchten Sie möglicherweise sicherstellen, dass ein bestimmtes Angebot immer aufgenommen wird, unabhängig davon, wie gut sein Score bei den einzelnen Kunden ist.

Um ein Angebot in diesem Fenster einzuschließen und seinen Präferenzgrad zu bearbeiten, klicken Sie auf Hinzufügen, geben Sie den Namen des Angebots ein (z. B. "Sparbuch" oder "Hypothek" und klicken Sie auf OK.

- **Wert.** Hier wird der Name des hinzugefügten Angebots angezeigt.
- **Voreinstellung.** Gibt den Präferenzgrad an, der auf das Angebot angewendet werden soll. Die Präferenz wird als Prozentsatz ausgedrückt, der als Dezimalwert zwischen 0,0 (keine Präferenz) und 1,0 (höchste Präferenz) angezeigt wird. Der Standardwert ist 0,0.
- **Immer einschließen.** Aktivieren Sie dieses Kontrollkästchen, um sicherzustellen, dass ein bestimmtes Angebot immer in die Vorhersagen eingeschlossen wird.

*Hinweis:* Wenn das Feld Voreinstellung auf 0.0 gesetzt ist, wird die Einstellung Immer einschließen ignoriert.

**Reliabilität des Modells berücksichtigen.** Ein gut strukturiertes Modell, das reich an Daten ist und durch mehrere erneute Generierungen eine Feinabstimmung erfahren hat, sollte stets genauere Ergebnisse liefern als ein ganz neues Modell mit wenigen Daten. Wenn Sie die höhere Reliabilität des ausgereifteren Modells nutzen möchten, aktivieren Sie dieses Kontrollkästchen.

# Modelle vom Typ "Support Vector Machine"

## Informationen zu SVM

Support Vector Machine (SVM) ist ein robustes Klassifikations- und Regressionsverfahren, das die Vorhersagegenauigkeit von Modellen erhöht, ohne dass es zu einer Überanpassung an die Trainingsdaten kommt. SVM ist insbesondere für die Analyse von Daten mit einer sehr großen Anzahl (z. B. mehrere Tausend) an Prädiktorfeldern geeignet.

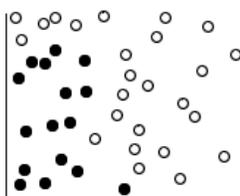
Für SVM gibt es Anwendungsbereiche in zahlreichen Fachgebieten, wie Customer Relationship Management (CRM), Gesichts- und sonstige Bilderkennung, Bioinformatik, Konzeptextrahierung beim Text-Mining, Intrusion Detection, Proteinstrukturvorhersage und Stimm- und Spracherkennung.

## So funktioniert SVM:

SVM funktioniert durch Zuordnung von Daten zu einem hochdimensionalen Merkmalsraum, sodass Datenpunkte kategorisiert werden können, selbst wenn die Daten anderweitig nicht linear getrennt werden können. Eine Trennlinie zwischen den Kategorien wird ermittelt. Anschließend werden die Daten derart transformiert, dass die Trennlinie als Hyperebene gezeichnet werden könnte. Danach kann anhand der Eigenschaften neuer Daten die Gruppe vorhergesagt werden, zu der ein neuer Datensatz gehören sollte.

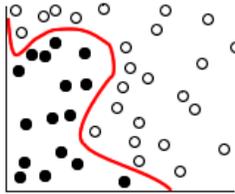
Betrachten Sie beispielsweise folgende Abbildung, bei der die Punkte in zwei verschiedene Kategorien entfallen:

Abbildung 15-1  
*Ursprüngliches Daten-Set*



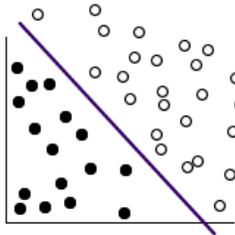
Die beiden Kategorien können dann durch eine Kurve getrennt werden:

Abbildung 15-2  
Daten nach Hinzufügen der Trennlinie



Nach der Transformation lässt sich die Grenze zwischen den Kategorien durch eine Hyperebene definieren:

Abbildung 15-3  
Transformierte Daten



Die für die Transformation verwendete mathematische Funktion wird als **Kernel**-Funktion bezeichnet. SVM in IBM® SPSS® Modeler unterstützt folgende Kernel-Typen:

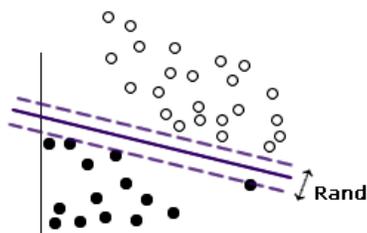
- Linear
- Polynomial
- Radiale Basisfunktion (RBF)
- Sigmoid

Eine lineare Kernel-Funktion wird empfohlen, wenn die lineare Trennung der Daten unproblematisch ist. In anderen Fällen sollte eine der folgenden Funktionen verwendet werden. Sie sollten die verschiedenen Funktionen ausprobieren, um in jedem Fall das beste Modell zu erzielen, da jede davon andere Algorithmen und Parameter verwendet.

## Feinabstimmung von SVM-Modellen

Neben der Trennlinie zwischen den Kategorien ermittelt ein SVM-Modell zur Klassifikation auch Randlinien, die den Raum zwischen den beiden Kategorien definieren:

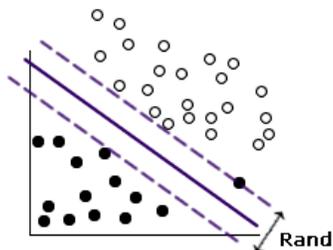
Abbildung 15-4  
Daten mit einem vorläufigen Modell



Die in den Rändern liegenden Datenpunkte werden als **Support-Vektoren** bezeichnet.

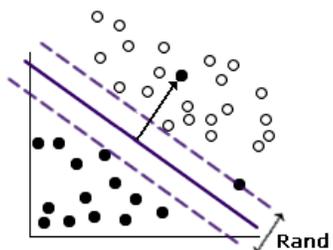
Je breiter der Rand zwischen den beiden Kategorien ist, desto besser ist das Modell bei der Vorhersage der Kategorie für neue Datensätze. Im vorherigen Beispiel ist der Rand nicht besonders breit und das Modell gilt als **überangepasst**. Ein geringes Maß an Fehlklassifizierung kann hingegenommen werden, um den Rand zu verbreitern. Beispiel:

Abbildung 15-5  
Daten mit einem verbesserten Modell



In einigen Fällen ist die lineare Trennung schwieriger. Beispiel:

Abbildung 15-6  
Problem bei linearer Trennung



In solchen Fällen besteht das Ziel darin, die optimale Balance zwischen einem möglichst breiten Rand und einer möglichst kleinen Zahl fehlklassifizierter Datenpunkte zu finden. Die Kernel-Funktion weist einen **Regularisierungsparameter** (als "C" bekannt) auf, der

den Ausgleich zwischen diesen beiden Werten steuert. Es wird vermutlich erforderlich sein, verschiedene Werte für diesen und andere Kernel-Parameter auszuprobieren, um das beste Modell zu finden.

## **SVM-Knoten**

Mit dem SCM-Knoten können Sie eine Support Vector Machine zum Klassifizieren von Daten verwenden. SCM eignet sich insbesondere für umfangreiche Daten-Sets, also solche mit einer großen Anzahl an Prädiktorfeldern. Mit den Standardeinstellungen im Knoten können Sie in relativ kurzer Zeit ein Grundmodell erstellen. Alternativ können Sie mithilfe der Experteneinstellungen verschiedene Typen von SVM-Modellen ausprobieren.

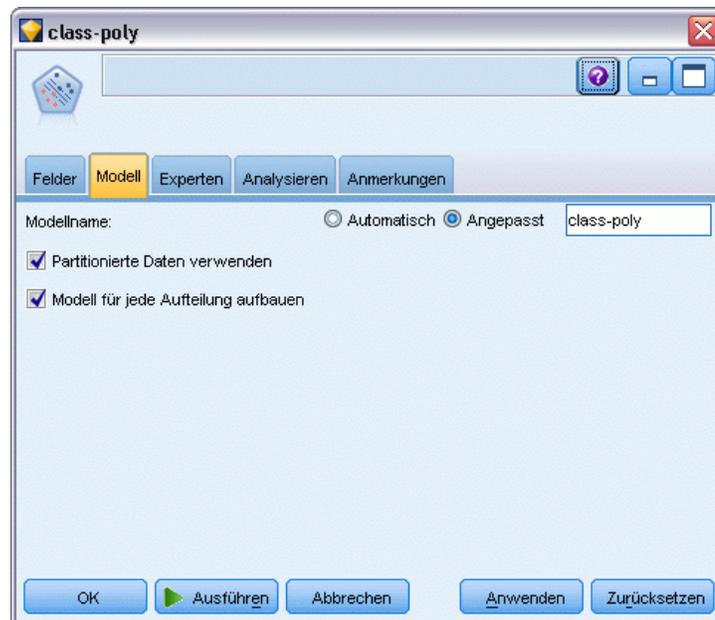
Nach der Erstellung des Modells haben Sie folgende Möglichkeiten:

- Durchsuchen des Modell-Nuggets zur Anzeige der relativen Bedeutsamkeit der Eingabefelder bei der Erstellung des Modells.
- Anfügen eines Tabellenknotens zum Modell-Nugget zur Anzeige der Modellausgabe.

**Beispiel.** Ein medizinischer Forscher hat ein Daten-Set mit den Eigenschaften einer Reihe von Stichproben menschlicher Zellen erstellt, die von Patienten stammen, bei denen ein Krebsrisiko angenommen wurde. Die Analyse der ursprünglichen Daten ergab, dass bei vielen der Eigenschaften deutliche Unterschiede zwischen den gutartigen und den bösartigen Proben bestehen. Der Forscher möchte ein SVM-Modell entwickeln, das die Werte ähnlicher Zelleneigenschaften in Proben von anderen Patienten verwenden kann, um eine Frühindikation dafür abzugeben, ob die Proben vermutlich gutartig oder bösartig sind. [Für weitere Informationen siehe Thema Klassifikation von Zellproben \(SVM\) in Kapitel 26 in IBM SPSS Modeler 14.2-Anwendungshandbuch.](#)

## Modelloptionen für SVM-Knoten

Abbildung 15-7  
Modelloptionen für SVM-Knoten



**Modellname.** Sie können den Modellnamen automatisch basierend auf den Ziel- oder ID-Feldnamen (oder dem Modelltyp in Fällen, in denen kein solches Feld angegeben ist) generieren oder einen benutzerdefinierten Namen eingeben.

**Partitionierte Daten verwenden.** Wenn ein Partitionsfeld definiert ist, gewährleistet diese Option, dass nur Daten aus der Trainingspartition für die Modellerstellung verwendet werden. [Für weitere Informationen siehe Thema Partitionsknoten in Kapitel 4 in IBM SPSS Modeler 14.2- Quellen-, Prozess- und Ausgabeknoten.](#)

**Aufteilungsmodelle erstellen.** Erstellt ein separates Modell für jeden möglichen Wert von Eingabefeldern, die als Aufteilungsfelder angegeben werden. [Für weitere Informationen siehe Thema Erstellung von aufgeteilten Modellen in Kapitel 3 auf S. 31.](#)

## Expertenoptionen für SVM-Knoten

Wenn Sie über umfassende Kenntnisse im Bereich Support Vector Machines verfügen, können Sie mithilfe der Expertenoptionen eine Feinabstimmung des Trainingsvorgangs vornehmen. Für den Zugriff auf die Expertenoptionen setzen Sie den Modus auf der Registerkarte "Experten" auf Experten.

Abbildung 15-8  
Expertenoptionen für SVM-Knoten



**Alle Wahrscheinlichkeiten anhängen (nur gültig für kategoriale Ziele).** Wenn diese Option ausgewählt ist, werden für jeden Datensatz, der vom Knoten verarbeitet wird, die Wahrscheinlichkeiten für jeden möglichen Wert eines Zielfelds vom Typ “Flag” oder “Nominal” angezeigt. Wenn diese Option nicht ausgewählt ist, wird für Zielfelder vom Typ “Flag” oder “Nominal” nur die Wahrscheinlichkeit des vorhergesagten Werts angezeigt. Die für dieses Kontrollkästchen festgelegte Einstellung bestimmt die Standardeinstellung des entsprechenden Kontrollkästchens in der Anzeige für das Modell-Nugget.

**Grenzkriterien.** Bestimmt, wann der Optimierungsalgorithmus gestoppt werden soll. Die Werte liegen im Bereich von  $1,0E-1$  bis  $1,0E-6$ ; der Standardwert ist  $1,0E-3$ . Eine Verringerung des Werts führt zu einem genaueren Modell bei gleichzeitiger Verlängerung des Zeitaufwands für das Trainieren des Modells.

**Regularisierungsparameter (C).** Steuert den Ausgleich zwischen der Maximierung des Rands und der Minimierung des Fehlerterms für das Training. Der Wert sollte normalerweise im Bereich von 1 bis 10 liegen; der Standardwert ist 10. Bei einer Erhöhung des Werts wird die Klassifizierungsgenauigkeit für die Trainingsdaten verbessert (bzw. der Regressionsfehler verringert), dies kann jedoch auch zu einer Überanpassung führen.

**Regressionsgenauigkeit (Epsilon).** Nur verwendet, wenn es sich beim Messniveau des Zielfelds um *Continuous* (Stetig) handelt. Führt dazu, dass Fehler hingenommen werden, vorausgesetzt, dass sie unter dem hier angegebenen Wert liegen. Eine Erhöhung des Werts kann zu schnellerer Modellierung führen, jedoch geht dies auf Kosten der Genauigkeit.

**Kernel-Typ.** Bestimmt den Typ der für die Transformation verwendeten Kernel-Funktion. Unterschiedliche Kernel-Typen führen dazu, dass die Trennlinie auf unterschiedliche Weise berechnet wird, weshalb es ratsam ist, mit den verschiedenen Optionen zu experimentieren. Der Standardwert lautet RBF (Radiale Basisfunktion).

**RBF-Gamma.** Nur aktiviert, wenn der Kernel-Typ auf RBF gesetzt ist. Der Wert sollte normalerweise zwischen  $3/k$  und  $6/k$  liegen, wobei  $k$  die Anzahl der Eingabefelder ist. Bei 12 Eingabefeldern beispielsweise wären somit Werte im Bereich von 0,25 bis 0,5 einen Versuch wert. Bei einer Erhöhung des Werts wird die Klassifizierungsgenauigkeit für die Trainingsdaten verbessert (bzw. der Regressionsfehler verringert), dies kann jedoch auch zu einer Überanpassung führen.

**Gamma.** Nur aktiviert, wenn der Kernel-Typ auf Polynomial oder Sigmoid. gesetzt ist. Bei einer Erhöhung des Werts wird die Klassifizierungsgenauigkeit für die Trainingsdaten verbessert (bzw. der Regressionsfehler verringert), dies kann jedoch auch zu einer Überanpassung führen.

**Verzerrung.** Nur aktiviert, wenn der Kernel-Typ auf Polynomial oder Sigmoid. gesetzt ist. Legt den Wert `coef0` in der Kernel-Funktion fest. Der Standardwert 0 ist in den meisten Fällen geeignet.

**Grad.** Nur aktiviert, wenn der Kernel-Typ auf Polynomial gesetzt ist. Steuert die Komplexität (Dimension) des Zuordnungsraums. Normalerweise werden nur Werte bis maximal 10 verwendet.

## SVM-Modell-Nugget

Das SVM-Modell erstellt eine Reihe neuer Felder. Das wichtigste dieser Felder ist das Feld `$$-feldname`, das den vom Modell vorhergesagten Wert für das Zielfeld anzeigt.

Anzahl und Namen der vom Modell erstellten neuen Felder hängen vom Messniveau des Zielfelds ab (in den folgenden Tabellen durch *Feldname* angegeben).

Um diese Felder und die zugehörigen Werte anzuzeigen, müssen Sie einen Tabellenknoten zum SVM-Modell-Nugget hinzufügen und den Tabellenknoten ausführen.

Tabelle 15-1  
Messniveau des Zielfelds ist 'Nominal' oder 'Flag'

Neuer Feldname	Beschreibung
<code>\$\$-Feldname</code>	Vorhergesagter Wert des Zielfelds.
<code>\$\$P-Feldname</code>	Wahrscheinlichkeit des vorhergesagten Werts.
<code>\$\$P-Wert</code>	Wahrscheinlichkeit jedes möglichen Werts von "Nominal" oder "Flag" (nur angezeigt, wenn auf der Registerkarte "Einstellungen" des Modell-Nuggets die Option Alle Wahrscheinlichkeiten anhängen aktiviert ist).
<code>\$\$SRP-Wert</code>	(Nur bei Flag-Zielen) Scores für Rohneigung (SRP) und korrigierte Neigung (SAP), die die Likelihood eines Ergebnisses vom Typ "wahr" für das Zielfeld angeben. Diese Scores werden nur angezeigt, wenn vor der Generierung des Modells die entsprechenden Kontrollkästchen auf der Registerkarte "Analysieren" des SVM-Modellierungsknotens ausgewählt wurden. <a href="#">Für weitere Informationen siehe Thema Analyseoptionen bei Modellierungsknoten in Kapitel 3 auf S. 41.</a>
<code>\$\$SAP-Wert.</code>	

Tabelle 15-2  
Messniveau des Zielfelds ist 'Continuous' (Stetig)

Neuer Feldname	Beschreibung
<code>\$\$-Feldname</code>	Vorhergesagter Wert des Zielfelds.

### **Bedeutsamkeit des Prädiktors**

Optional kann auf der Registerkarte “Modell” auch ein Diagramm, das die relative Bedeutsamkeit der einzelnen Prädiktoren für die Schätzung des Modells angibt, angezeigt werden. Normalerweise ist es sinnvoll, die Modellierungsbemühungen auf die wichtigsten Prädiktoren zu konzentrieren und diejenigen Prädiktoren zu verwerfen bzw. zu ignorieren, die die geringste Bedeutung haben. Beachten Sie, dass dieses Diagramm nur verfügbar ist, wenn vor dem Generieren des Modells Bedeutsamkeit der Prädiktoren berechnet auf der Registerkarte “Analysieren” ausgewählt wurde. [Für weitere Informationen siehe Thema Bedeutsamkeit des Prädiktors in Kapitel 3 auf S. 54.](#)

*Anmerkung:* Die Berechnung der Bedeutsamkeit der Prädiktoren kann bei SVM länger dauern als bei anderen Modelltypen und ist nicht standardmäßig auf der Registerkarte “Analysieren” aktiviert. Die Auswahl dieser Option kann die Leistung verlangsamen, insbesondere bei großen Daten-Sets.

### **Einstellungen beim SVM-Modell**

Abbildung 15-9  
SVM-Modell, Registerkarte “Einstellungen”



In der Registerkarte “Einstellungen” können Sie zusätzliche Felder angeben, die bei der Anzeige der Ergebnisse verwendet werden sollen (z. B. durch Ausführen eines Tabellenknotens, der mit dem Nugget verknüpft ist). Sie können den Effekt jeder dieser Optionen sehen, indem Sie sie auswählen und auf die Schaltfläche “Vorschau” klicken. Führen Sie in der Vorschau-Ausgabe einen Bildlauf nach rechts durch, um die zusätzlichen Felder zu sehen.

**Alle Wahrscheinlichkeiten anhängen (nur gültig für kategoriale Ziele).** Wenn diese Option ausgewählt ist, werden für jeden Datensatz, der vom Knoten verarbeitet wird, die Wahrscheinlichkeiten für jeden möglichen Wert eines Zielfelds vom Typ “Nominal” oder “Flag” angezeigt. Wenn diese Option nicht ausgewählt ist, werden für Zielfelder vom Typ “Flag” oder “Nominal” nur der vorhergesagte Wert und seine Wahrscheinlichkeit angezeigt.

Die Standardeinstellung dieses Kontrollkästchens richtet sich nach dem entsprechenden Kontrollkästchen des Modellierungsknotens.

**Scores für Rohneigung berechnen.** Bei Modellen mit einem Flag-Ziel (das als Vorhersage "Ja" bzw. "Nein" ausgibt) können Sie Neigungs-Scores anfordern, die die Wahrscheinlichkeit des für das Zielfeld angegebenen wahren Ergebnisses angeben. Diese Werte werden zusätzlich zu den anderen Vorhersage- und Konfidenzwerten angegeben, die ggf. während des Scorings erstellt werden.

**Scores für korrigierte Neigung berechnen.** Rohneigungs-Scores beruhen ausschließlich auf den Trainingsdaten und sind aufgrund der Neigung vieler Modelle zur übermäßigen Anpassung an die Trainingsdaten möglicherweise zu optimistisch. Bei korrigierten Neigungen wird versucht, dies durch Evaluation der Modellleistung anhand einer Test- bzw. Validierungspartition zu kompensieren. Bei dieser Option muss im Stream ein Partitionsfeld definiert sein und die Scores für die korrigierte Neigung müssen im Modellierungsknoten aktiviert werden, bevor das Modell generiert wird.

# Nächste-Nachbarn-Modelle

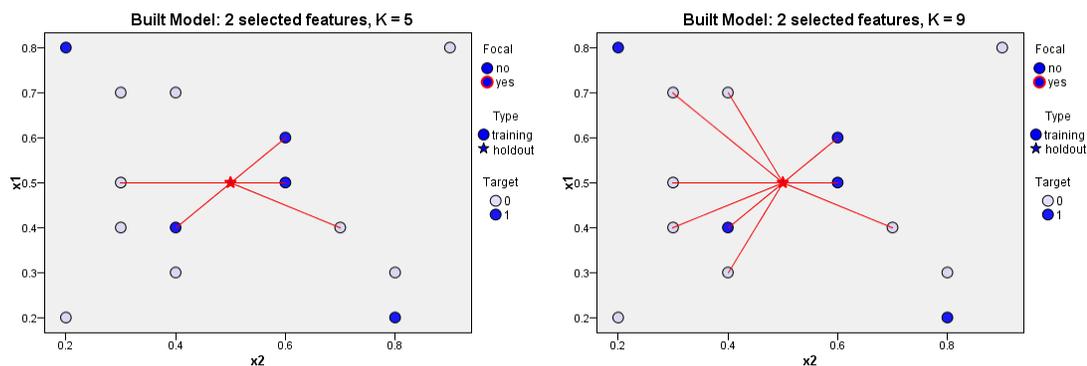
## KNN-Knoten

Die Nächste-Nachbarn-Analyse ist eine Methode zur Klassifizierung von Fällen anhand ihrer Ähnlichkeit zu anderen Fällen. Im Maschinenlernen wurde es entwickelt, um Datenmuster zu erkennen, ohne dass eine exakte Übereinstimmung mit gespeicherten Mustern oder Fällen benötigt wird. Ähnliche Fälle befinden sich nahe beieinander und unterschiedliche Fälle sind voneinander entfernt. Somit gilt die Distanz zwischen zwei Fällen als Maß für ihre Unähnlichkeit.

Befinden sich Fälle nahe beieinander, werden sie als "Nachbarn" bezeichnet. Wenn ein neuer Fall (Holdout) angegeben wird, wird seine Distanz zu jedem der Fälle im Modell berechnet. Die Klassifizierungen der ähnlichsten Fälle – die nächsten Nachbarn – werden gezählt und der neue Fall wird einer Kategorie zugeordnet, die die größte Anzahl der nächsten Nachbarn enthält.

Sie können die Zahl der zu untersuchenden nächsten Nachbarn festlegen; dieser Wert wird  $k$  genannt. Die Abbildungen zeigen, wie ein neuer Fall mit Hilfe von zwei verschiedenen Werten von  $k$  klassifiziert würde. Ist  $k = 5$ , wird der neue Fall der Kategorie 1 zugeordnet, da der Großteil der nächsten Nachbarn der Kategorie 1 angehört. Ist jedoch  $k = 9$ , wird der neue Fall der Kategorie 0 zugeordnet, da der Großteil der nächsten Nachbarn der Kategorie 0 angehört.

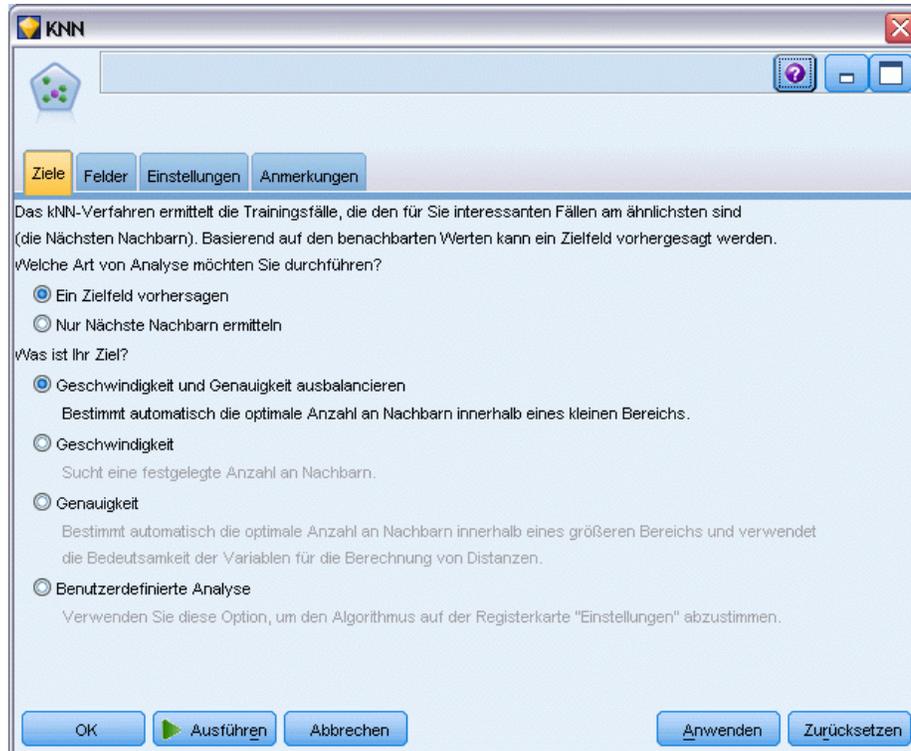
Abbildung 16-1  
Auswirkungen der Änderung von "k" bei der Klassifizierung



Die Nächste-Nachbarn-Analyse kann auch zur Berechnung von Werten für ein stetiges Ziel verwendet werden. Dabei wird der durchschnittliche oder Median-Zielwert der nächsten Nachbarn verwendet, um den vorhergesagten Wert für den neuen Fall zu beziehen.

## Zieloptionen für KNN-Knoten

Abbildung 16-2  
Zieloptionen für KNN-Knoten



In der Registerkarte “Ziele” können Sie wählen, ob Sie ein Modell erstellen möchten, das den Wert eines Zielfelds in Ihren Eingabedaten auf der Basis der Werte seiner nächsten Nachbarn vorhersagt, oder ob Sie einfach die nächsten Nachbarn für einen bestimmten Fall von Interesse herausfinden möchten.

### **Welche Art der Analyse möchten Sie ausführen?**

**Zielfeld vorhersagen.** Wählen Sie diese Option, wenn Sie den Wert eines Zielfelds auf der Grundlage der Werte seiner nächsten Nachbarn vorhersagen möchten.

**Nur nächste Nachbarn identifizieren.** Wählen Sie diese Option, wenn Sie nur die nächsten Nachbarn für ein bestimmtes Eingabefeld sehen möchten.

Wenn Sie sich entscheiden, nur einen der nächsten Nachbarn zu identifizieren, werden die übrigen Optionen in dieser Registerkarte (für Genauigkeit und Geschwindigkeit) deaktiviert, da sie nur für die Prognose von Zielen relevant sind.

**Wie lautet Ihr Ziel?**

Mit dieser Optionsgruppe können Sie entscheiden, ob Geschwindigkeit und/oder Genauigkeit die wichtigsten Faktoren beim Vorhersagen eines Zielfelds sind. Alternativ können Sie die Einstellungen selbst anpassen.

Wenn Sie die Option “Balancieren”, “Geschwindigkeit” oder “Genauigkeit” wählen, trifft der Algorithmus die am besten geeignete Einstellungskombination für diese Option als Voreinstellung. Erfahrene Benutzer möchten diese Voreinstellungen eventuell überschreiben; dies ist in den verschiedenen Bereichen der Registerkarte “Einstellungen” möglich.

**Geschwindigkeit und Genauigkeit ausgleichen.** Wählt die beste Anzahl an Nachbarn in einem kleinen Bereich aus.

**Geschwindigkeit.** Findet eine feste Anzahl an Nachbarn.

**Genauigkeit.** Wählt die beste Anzahl an Nachbarn in einem größeren Bereich aus und berechnet Entfernungen anhand der Prädiktorenbedeutung.

**Analyse anpassen** Wählen Sie diese Option, um den Algorithmus in der Registerkarte “Einstellungen” genauer einzustellen.

*Anmerkung:* Anders als bei anderen Modellen steigt die Größe des Ergebnis-KNN-Modells linear mit der Menge an Trainingsdaten an. Wird beim Versuch, ein KNN-Modell zu erstellen, eine Fehlermeldung angezeigt, dass nicht genügend Speicher vorhanden ist, versuchen Sie, den maximalen von IBM® SPSS® Modeler verwendeten Systempeicher zu erhöhen. Wählen Sie dazu

Werkzeuge > Optionen > Systemoptionen

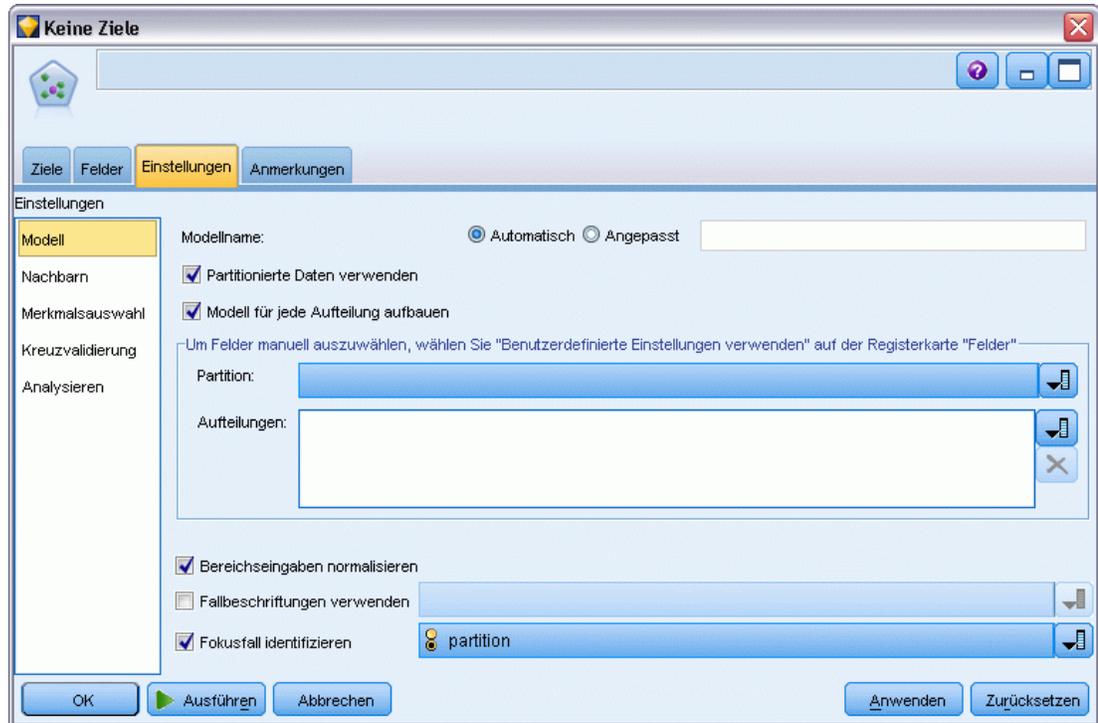
und geben Sie die neue Größe im Feld Maximale Speichergröße ein. Im Dialogfeld “Systemoptionen” vorgenommene Änderungen werden erst nach einem Neustart von SPSS Modeler wirksam.

**KNN-Knoten - Einstellungen**

In der Registerkarte “Einstellungen” geben Sie die spezifischen Optionen für die Nächste-Nachbarn-Analyse an. Die linke Randleiste am Bildschirm listet die Bereiche auf, die Sie zur Festlegung der Optionen verwenden.

## Modell

Abbildung 16-3  
Modelloptionen für KNN-Knoten



Das Modellfenster bietet Optionen, die steuern, wie das Modell erstellt werden soll, z. B. ob Partitionierung oder Aufteilungsmodelle benutzt werden, ob numerische Eingabefelder so transformiert werden, dass sie alle im selben Bereich liegen, und wie bestimmte Fälle verwaltet werden. Sie können auch einen benutzerdefinierten Namen für das Modell angeben.

**Modellname.** Sie können den Modellnamen automatisch basierend auf den Ziel- oder ID-Feldnamen (oder dem Modelltyp in Fällen, in denen kein solches Feld angegeben ist) generieren oder einen benutzerdefinierten Namen eingeben.

**Partitionierte Daten verwenden.** Wenn ein Partitionsfeld definiert ist, gewährleistet diese Option, dass nur Daten aus der Trainingspartition für die Modellerstellung verwendet werden. [Für weitere Informationen siehe Thema Partitionsknoten in Kapitel 4 in IBM SPSS Modeler 14.2- Quellen-, Prozess- und Ausgabeknoten.](#)

**Aufteilungsmodelle erstellen.** Erstellt ein separates Modell für jeden möglichen Wert von Eingabefeldern, die als Aufteilungsfelder angegeben werden. [Für weitere Informationen siehe Thema Erstellung von aufgeteilten Modellen in Kapitel 3 auf S. 31.](#)

**Felder manuell auswählen...** Standardmäßig verwendet der Knoten die Partitions- und Aufteilungsfeldeinstellungen (falls vorhanden) des Typknotens. Sie können diese Einstellungen jedoch hier überschreiben. Um die Partitions- und Aufteilungsfelder zu aktivieren, wählen Sie die Registerkarte Felder und dann die Option Benutzerdefinierte Einstellungen verwenden. Kehren Sie anschließend hierher zurück.

- **Partition.** In diesem Feld können Sie ein Feld angeben, das verwendet wird, um die Daten in getrennte Stichproben für die Trainings-, Test- und Validierungsphase der Modellbildung aufzuteilen. Indem Sie mit einer Stichprobe das Modell erstellen und es mit einer anderen Stichprobe testen, erhalten Sie einen guten Hinweis dafür, wie gut das Modell sich für größere Datenmengen generalisieren lässt, die den aktuellen Daten ähneln. Wenn mehrere Partitionsfelder mithilfe von Typ- oder Partitionsknoten erstellt wurden, muss in jedem Modellierungsknoten, der die Partitionierung verwendet, auf der Registerkarte “Felder” ein einzelnes Partitionsfeld ausgewählt werden. (Wenn nur eine einzige Partition vorhanden ist, wird diese immer automatisch verwendet, wenn die Partitionierung aktiviert ist.) [Für weitere Informationen siehe Thema Partitionsknoten in Kapitel 4 in IBM SPSS Modeler 14.2-Quellen-, Prozess- und Ausgabeknoten.](#) Beachten Sie außerdem, dass die Partitionierung auch auf der Registerkarte “Modelloptionen” des Knotens aktiviert sein muss, wenn die ausgewählte Partitionierung in Ihrer Analyse angewendet werden soll. (Wenn die Auswahl der Option aufgehoben ist, kann die Partitionierung ohne Änderung der Feldeinstellungen deaktiviert werden.)
- **Aufteilen.** Wählen Sie für Aufteilungsmodelle das Aufteilungsfeld bzw. die Aufteilungsfelder. Dies ist so, als würden Sie in einem Typknoten für die Rolle eines Felds den Wert *Aufteilung* festlegen. Sie können nur Felder vom Typ Flag, Nominal oder Ordinal als Aufteilungsfelder festlegen. Als Aufteilungsfelder gewählte Felder können nicht als Ziel-, Prädiktor-, Partitions-, Häufigkeits- oder Gewichtungsfelder verwendet werden. [Für weitere Informationen siehe Thema Erstellung von aufgeteilten Modellen in Kapitel 3 auf S. 31.](#)

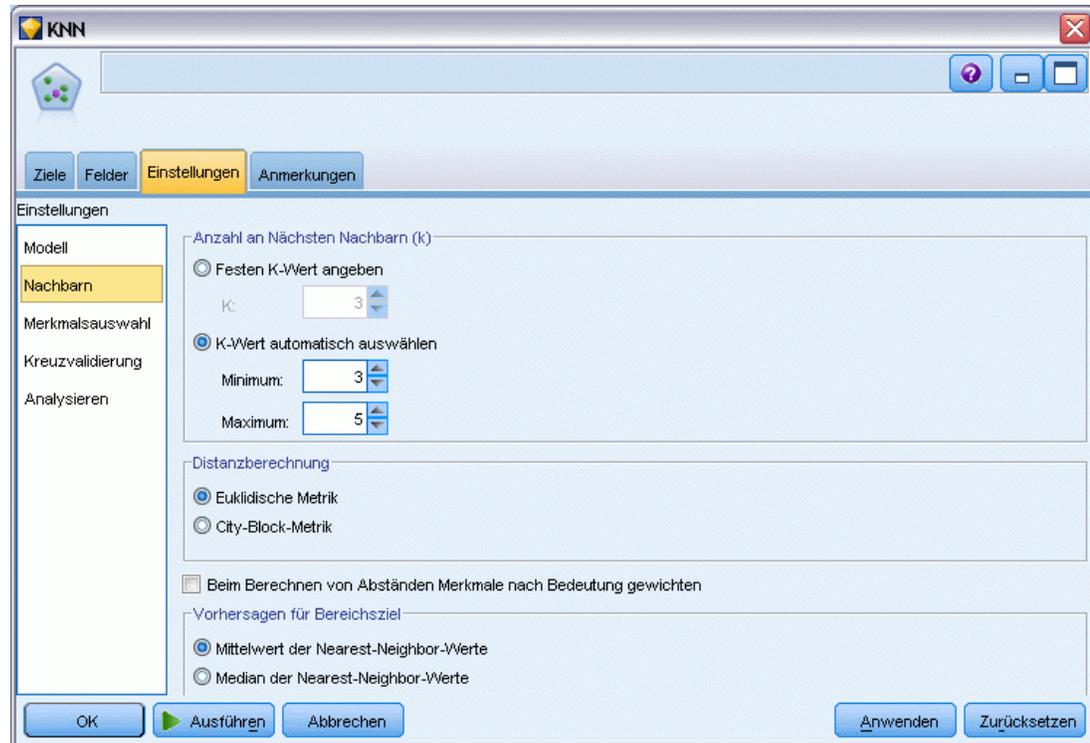
**Bereichseingaben normalisieren.** Markieren Sie dieses Kontrollkästchen, um die Werte für stetige Eingabefelder zu normalisieren. Normalisierte Merkmale umfassen denselben Wertebereich, was die Leistung des Schätzalgorithmus verbessern kann. Es wird eine korrigierte Normalisierung,  $[2*(x-\min)/(\max-\min)]-1$ , angewendet. Korrigierte, normalisierte Werte liegen im Bereich zwischen -1 und 1.

**Falllabels verwenden.** Markieren Sie dieses Kontrollkästchen, um die Dropdown-Liste zu aktivieren, in der Sie ein Feld wählen können, dessen Werte als Beschriftungen verwendet werden, um die Fälle von Interesse im Prädiktorraumdiagramm, Vergleichsdiagramm und in Quadrantenkarten im Modell-Viewer zu identifizieren. Sie können ein beliebiges Feld mit einem Messniveau von *Nominal*, *Ordinal* oder *Flag* als Beschriftungsfeld wählen. Wenn Sie hier kein Feld wählen, werden Datensätze in den Modell-Viewer-Diagrammen mit den nächstgelegenen Nachbarn angezeigt, wobei diese anhand der Zeilennummer in den Quelldaten identifiziert werden. Wenn Sie die Daten überhaupt nach der Erstellung des Modells ändern, verwenden Sie Labels, damit Sie nicht jedesmal auf die Quelldaten zurückverweisen müssen, um die Fälle in der Anzeige zu identifizieren.

**Fokusdatensatz identifizieren.** Markieren Sie dieses Kontrollkästchen, um die Dropdown-Liste zu aktivieren, in der Sie ein Eingabefeld von besonderem Interesse kennzeichnen können (nur für Flag-Felder). Wenn Sie hier ein Feld festlegen, sind die Punkte, die dieses Feld repräsentieren, anfangs im Modell-Viewer ausgewählt, während das Modell erstellt wird. Die Auswahl eines Fokusdatensatzes an dieser Stelle ist optional. Jeder Punkt kann temporär zu einem Fokusdatensatz werden, wenn er manuell im Modell-Viewer ausgewählt wird.

## Nachbarn

Abbildung 16-4  
Nachbarnoptionen für KNN-Knoten



Das Fenster “Nachbarn” enthält ein Set an Optionen, die steuern, wie die Anzahl der nächsten Nachbarn berechnet wird.

**Anzahl der nächstgelegenen Nachbarn (k).** Geben Sie die Anzahl der nächsten Nachbarn für einen bestimmten Fall ein. Beachten Sie dabei, dass eine höhere Anzahl an Nachbarn nicht unbedingt ein präziseres Modell hervorbringt.

Wenn geplant ist, dass ein Ziel vorhergesagt werden soll, stehen zwei Optionen zur Auswahl:

- **Festen Wert für k angeben.** Verwenden Sie diese Option, wenn Sie eine feste Anzahl von nächsten Nachbarn angeben möchten, die gefunden werden sollen.
- **k automatisch auswählen.** Sie können alternativ mithilfe der Felder Minimum und Maximum einen Wertebereich festlegen und es der Prozedur überlassen, die “beste” Anzahl an Nachbarn innerhalb dieses Bereichs zu wählen. Die Methode zur Bestimmung der Anzahl der nächsten Nachbarn hängt davon ab, ob die Merkmalsauswahl im Merkmalsauswahlfenster verlangt wird.

Wenn die Funktionsauswahl aktiviert wurde, wird für jeden Wert von  $k$  im angegebenen Bereich eine Funktionsauswahl durchgeführt und  $k$  und die zugehörige Funktionsgruppe mit der niedrigsten Fehlerrate (oder dem geringsten Quadratsummen-Fehler, falls das Ziel stetig ist) werden ausgewählt.

Wenn die Funktionsauswahl nicht aktiviert ist, wird eine  $V$ -fache Kreuzvalidierung angewendet, um die "beste" Anzahl an Nachbarn zu ermitteln. Siehe den Bereich "Kreuzvalidierung" für die Steuerung der Zuweisung von Aufteilungen.

**Distanzberechnung.** Mit diesem Wert wird das Längenmaßsystem für die Messung der Ähnlichkeit von Fällen festgelegt.

- **Euklidisch.** Der Abstand zwischen zwei Fällen,  $x$  und  $y$ , ergibt sich aus der Quadratwurzel der Summe, über alle Dimensionen, der quadrierten Differenzen zwischen den Werten für die Fälle.
- **Stadtblock.** Die Distanz zwischen zwei Fällen ergibt sich aus der Summe, über alle Dimensionen, der absoluten Differenzen zwischen den Werten der Fälle. Dies wird auch als Manhattan-Distanz bezeichnet.

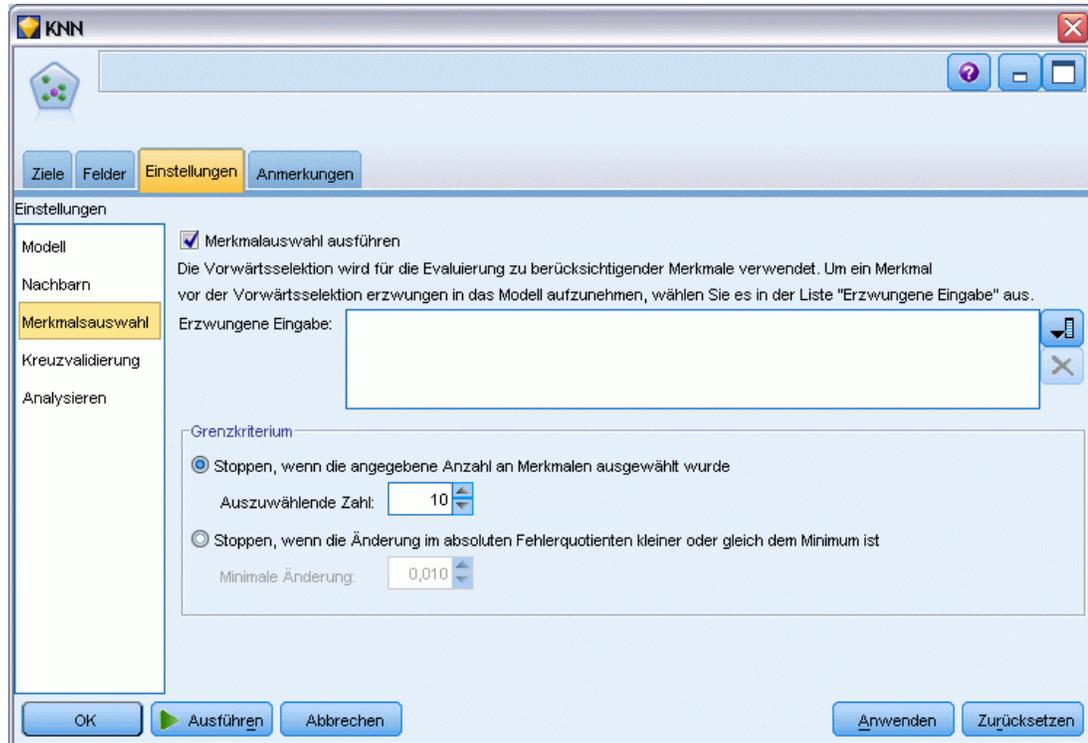
Optional: Wenn geplant ist, ein Ziel vorherzusagen, können Sie beim Berechnen von Distanzen die Merkmale nach ihrer normalisierten Bedeutung gewichten. Die Merkmalswichtigkeit für einen Prädiktor wird durch das Verhältnis der Fehlerquote oder des Quadratsummenfehlers des Modells (wobei der Prädiktor vom Modell entfernt wird) zur Fehlerquote bzw. zum Quadratsummenfehler für das vollständige Modell berechnet. Die normalisierte Wichtigkeit wird durch die Neugewichtung der Werte der Funktionswichtigkeit berechnet, so dass deren Summe 1 ergibt.

**Gewichtungsmerkmale nach Wichtigkeit beim Berechnen von Abständen.** (Wird nur angezeigt, wenn ein Ziel vorhergesagt werden soll.) Markieren Sie dieses Kontrollkästchen, damit die Bedeutsamkeit des Prädiktors bei der Distanzberechnung zwischen Nachbarn verwendet wird. Die Bedeutsamkeit des Prädiktors wird dann im Modell-Nugget angezeigt und in Vorhersagen verwendet (und beeinflusst damit das Scoring). [Für weitere Informationen siehe Thema Bedeutsamkeit des Prädiktors in Kapitel 3 auf S. 54.](#)

**Vorhersagen für das Bereichsziel.** (Wird nur angezeigt, wenn ein Ziel vorhergesagt werden soll.) Wenn ein stetiges Ziel (numerischer Bereich) angegeben ist, definiert dies, ob der vorhergesagte Wert auf der Basis des Mittel- oder Medianwerts der nächsten Nachbarn berechnet wird.

## Funktionsauswahl

Abbildung 16-5  
Merkmalsauswahloptionen für KNN-Knoten



Dieser Bereich wird nur aktiviert, wenn ein Ziel vorhergesagt werden soll. Hier können Sie Optionen zur Merkmalsauswahl anfordern und angeben. Standardmäßig werden bei der Funktionsauswahl alle Funktionen berücksichtigt, Sie können optional aber auch eine Untergruppe von Funktionen auswählen, die in das Modell aufgenommen werden sollen.

**Funktionsauswahl durchführen.** Markieren Sie dieses Kontrollkästchen, um die Optionen zur Merkmalsauswahl zu aktivieren.

- **Erwungene Aufnahme.** Klicken Sie auf die Feldauswahlschaltfläche neben diesem Feld und wählen Sie ein oder mehrere Merkmale, deren Aufnahme in das Modell erzwungen werden soll.

**Stoppkriterien.** Bei jedem Schritt wird die Funktion, deren Integration in das Modell den geringsten Fehler hervorruft (für kategoriale Ziele als Fehlerrate und für stetige Ziele als Quadratsummenfehler berechnet), für die Integration in das Modell in Betracht gezogen. Die Vorwärtsselektion wird fortgesetzt, bis die angegebene Bedingung erfüllt wird.

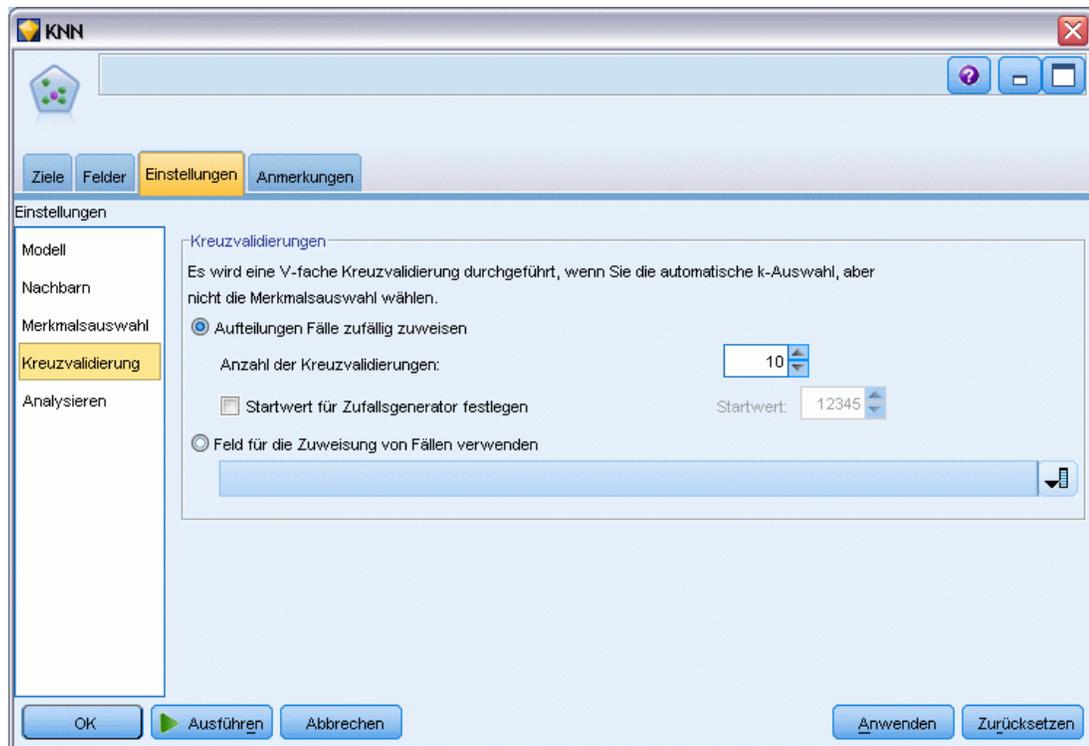
- **Stoppen, wenn die angegebene Anzahl an Merkmalen ausgewählt wurde.** Der Algorithmus fügt neben den erzwungenen Funktionen eine feste Anzahl an Funktionen in das Modell ein. Geben Sie eine positive Ganzzahl ein. Eine geringere Anzahl an Werten führt zu einem sparsameren Modell. Dabei läuft man allerdings Gefahr, wichtige Funktionen zu

vernachlässigen. Bei einer höheren Anzahl an Werten werden alle wichtigen Funktionen erfasst, dafür läuft man aber Gefahr, Funktionen einzufügen, die den Modellfehler erhöhen.

- **Stoppen, wenn die Änderung im absoluten Fehlerverhältnis kleiner oder gleich dem Minimum ist.** Der Algorithmus wird beendet, wenn die Änderung im absoluten Fehlerquotienten vermuten lässt, dass das Modell durch Hinzufügen weiterer Funktionen nicht mehr weiter optimiert werden kann. Geben Sie eine positive Zahl an. Geringere Werte für die Mindeständerung berücksichtigen mehr Merkmale und bergen das Risiko, dass Merkmale aufgenommen werden, die dem Modell keinen zusätzlichen Wert bescheren. Bei einem höheren Wert für die minimale Änderungen werden mehr Funktionen ausgeschlossen, was dazu führen kann, dass Funktionen ausgeschlossen werden, die wichtig für das Modell wären. Der "optimale" Wert für die minimale Änderung hängt von den jeweiligen Daten und dem Anwendungsbereich ab. Informationen dazu, wie Sie beurteilen, welche Funktionen am wichtigsten sind, finden Sie im Protokoll über die Funktionsauswahlfehler in der Ausgabe. [Für weitere Informationen siehe Thema Prädiktor-Auswahlfehler-Protokoll auf S. 501.](#)

## Kreuzvalidierung

Abbildung 16-6  
Kreuzvalidierungsoptionen für KNN-Knoten



Dieser Bereich wird nur aktiviert, wenn ein Ziel vorhergesagt werden soll. Die Optionen in diesem Bereich steuern, ob beim Berechnen der nächsten Nachbarn Kreuzvalidierung verwendet werden soll.

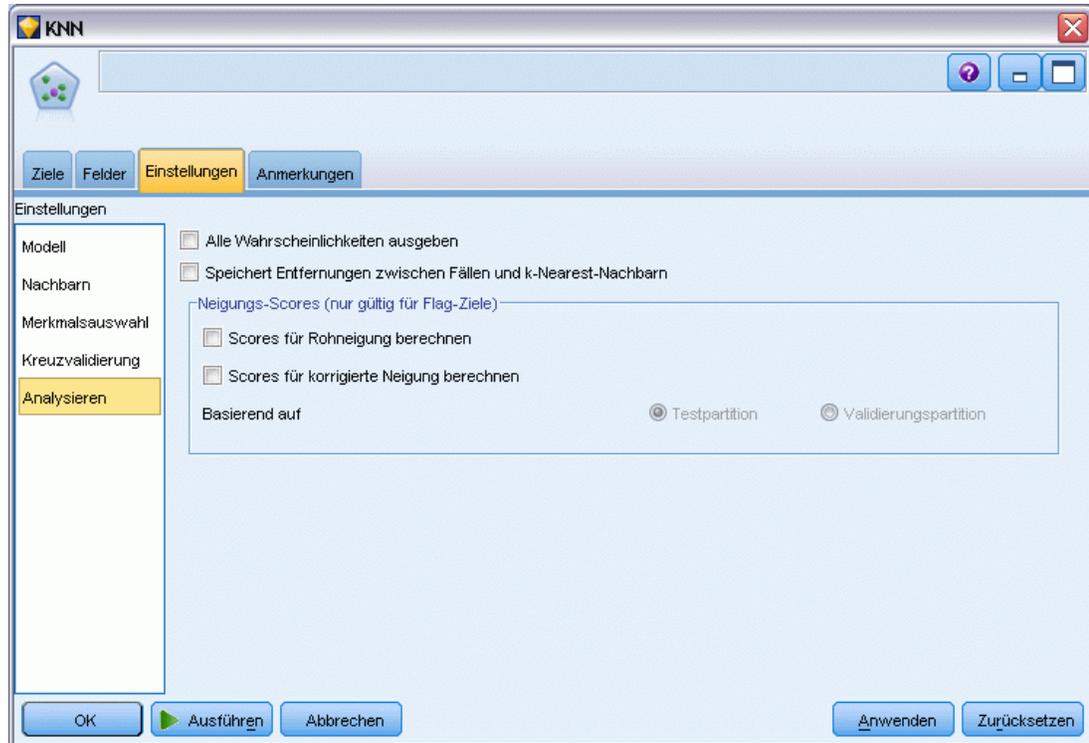
Bei der Kreuzvalidierung wird die Stichprobe in mehrere Teilstichproben oder **Aufteilungen** gegliedert. Anschließend werden Nächste-Nachbarn-Modelle erzeugt; dabei werden nacheinander die Daten der einzelnen Stichproben ausgeschlossen. Das erste Modell beruht auf allen Fällen mit Ausnahme der Fälle in der ersten Stichprobenaufteilung, das zweite Modell auf allen Fällen mit Ausnahme der Fälle in der zweiten Stichprobenaufteilung usw. Bei jedem Modell wird jeweils der Fehler geschätzt. Hierzu wird das Modell auf die Teilstichprobe angewendet, die beim Erstellen des Modells ausgeschlossen war. Die “beste” Anzahl an nächstgelegenen Nachbarn ist die Anzahl, die die wenigsten Fehler für alle Aufteilungen erzeugt.

**Kreuzvalidierungsaufteilungen.** Um die “beste” Anzahl an Nachbarn zu ermitteln wird eine  $V$ -fache Vergleichsprüfung durchgeführt. Bei Funktionsauswahl ist sie aus Leistungsgründen nicht verfügbar.

- **Aufteilungen willkürlich Fälle zuweisen.** Geben Sie die Anzahl an Aufteilungen an, die für die Vergleichsprüfung herangezogen werden sollen. Die Prozedur weist Fälle willkürlich Aufteilungen zu und nummeriert sie von 1 bis  $V$ , die Anzahl an Aufteilungen.
- **Startwert für Zufallsgenerator festlegen.** Bei der Schätzung der Genauigkeit eines Modells auf der Grundlage eines Zufallsprozentsatzes können Sie mit dieser Option dieselben Ergebnisse in einer anderen Sitzung replizieren. Wenn Sie den vom Zufallszahlengenerator verwendeten Startwert angeben, stellen Sie sicher, dass bei jeder Ausführung des Knotens dieselben Datensätze zugewiesen werden. Gewünschten Startwert eingeben. Wenn diese Option nicht ausgewählt ist, wird bei jeder Ausführung des Knotens eine andere Stichprobe generiert.
- **Feld zum Zuweisen von Fällen verwenden.** Geben Sie ein numerisches Feld an, das jeden Fall im aktiven Daten-Set einer Unterteilung zuweist. Das Feld muss numerisch sein und Werte von 1 bis  $V$  annehmen. Wenn Werte in diesem Bereich fehlen und wenn für Aufteilungsfelder etwaige Aufteilungsmodelle wirksam sind, wird ein Fehler ausgelöst.

## Analysieren

Abbildung 16-7  
Analyseoptionen für KNN-Knoten



Der Bereich “Analyse” wird nur aktiviert, wenn ein Ziel vorhergesagt werden soll. Sie können diesen Bereich verwenden, um anzugeben, ob das Modell zusätzliche Variablen enthalten soll, die Folgendes enthalten:

- Wahrscheinlichkeiten für jeden möglichen Wert eines Zielfelds
- Distanzen zwischen einem Fall und seinen nächsten Nachbarn
- Scores für Rohneigung und angepasste Neigung (nur für Flag-Ziele)

**Alle Wahrscheinlichkeiten ausgeben.** Wenn diese Option ausgewählt ist, werden für jeden Datensatz, der vom Knoten verarbeitet wird, die Wahrscheinlichkeiten für jeden möglichen Wert eines Zielfelds vom Typ “Nominal” oder “Flag” angezeigt. Wenn diese Option nicht ausgewählt ist, werden für Zielfelder vom Typ “Flag” oder “Nominal” nur der vorhergesagte Wert und seine Wahrscheinlichkeit angezeigt.

**Distanzen zwischen Fällen und k-Nächsten-Nachbarn speichern.** Für jeden Fokusdatensatz wird eine separate Variable für jeden  $k$ -Nächsten-Nachbarn (aus dem Schulungsbeispiel) und die entsprechenden  $k$ -Nächsten-Distanzen erstellt.

### **Neigungs-Scores**

Neigungs-Scores können im Modellierungsknoten oder auf der Registerkarte “Einstellungen” im Modell-Nugget aktiviert werden. Diese Funktionalität ist nur verfügbar, wenn das ausgewählte Ziel ein Flag-Feld ist. [Für weitere Informationen siehe Thema Neigungsbewertungen in Kapitel 3 auf S. 44.](#)

**Scores für Rohneigung berechnen.** Rohneigungs-Scores werden ausschließlich auf der Grundlage der Trainingsdaten aus dem Modell abgeleitet. Wenn das Modell den Wert *wahr* (wird antworten) vorhersagt, ist die Neigung mit P identisch. Dabei ist P die Wahrscheinlichkeit der Vorhersage. Wenn das Modell den Wert “falsch” vorhersagt, wird die Neigung als  $(1 - P)$  berechnet.

- Wenn Sie bei der Modellerstellung diese Option auswählen, werden standardmäßig Neigungs-Scores im Modell-Nugget aktiviert. Sie können jedoch immer festlegen, dass Rohneigungs-Scores im Modell-Nugget aktiviert werden sollen, unabhängig davon, ob Sie sie im Modellierungsknoten auswählen.
- Beim Scoring des Modells werden Rohneigungs-Scores in einem Feld hinzugefügt, bei dem die Buchstaben *RP* an das Standard-Präfix angehängt sind. Wenn sich die Vorhersagen beispielsweise im Feld *\$R-Abwanderung* befinden, lautet der Name des Felds für den Neigungs-Score *\$RRP-Abwanderung*.

**Scores für korrigierte Neigung berechnen.** Rohneigungen basieren ausschließlich auf vom Modell angegebenen Schätzern. Beim Modell kann jedoch eine Überanpassung vorliegen, was zu übermäßig optimistischen Schätzern für die Neigung führt. Korrigierte Neigungen versuchen, dies zu kompensieren, indem untersucht wird, wie leistungsfähig das Modell bei den Test- bzw. Validierungspartitionen ist, und die Neigungen entsprechend angepasst werden, um einen besseren Schätzer zu erzeugen.

- Diese Einstellung ist nur möglich, wenn ein gültiges Partitionsfeld im Stream vorhanden ist. [Für weitere Informationen siehe Thema Partitionsknoten in Kapitel 4 in IBM SPSS Modeler 14.2- Quellen-, Prozess- und Ausgabeknoten.](#)
- Anders als rohe Konfidenz-Scores müssen Scores für die korrigierte Neigung bei der Erstellung des Modells berechnet werden; anderenfalls stehen Sie beim Scoring des Modell-Nuggets nicht zur Verfügung.
- Beim Scoring des Modells werden Scores für die korrigierte Neigung in einem Feld hinzugefügt, bei dem die Buchstaben *AP* an das Standard-Präfix angehängt sind. Wenn sich die Vorhersagen beispielsweise im Feld *\$R-Abwanderung* befinden, lautet der Name des Felds für den Neigungs-Score *\$RAP-Abwanderung*. Scores für die korrigierte Neigung stehen bei logistischen Regressionsmodellen nicht zur Verfügung.
- Bei der Berechnung der Scores für die korrigierte Neigung darf die für die Berechnung verwendete Test- bzw. Validierungspartition nicht ausbalanciert worden sein. Um dies zu vermeiden, müssen Sie darauf achten, dass in etwaigen weiter oben im Stream befindlichen Balancierungsknoten die Option Balancierung nur für Trainingsdaten durchführen ausgewählt wurde. [Für weitere Informationen siehe Thema Festlegen der Optionen für den Balancierungsknoten in Kapitel 3 in IBM SPSS Modeler 14.2- Quellen-, Prozess- und](#)

*Ausgabeknoten.* Außerdem gilt: Wenn weiter oben im Stream eine komplexe Stichprobe gezogen wurde, werden dadurch die Scores für die korrigierte Neigung ungültig.

- Scores für die korrigierte Neigung stehen bei verstärkten Baum- und Regelmengenmodellen nicht zur Verfügung. [Für weitere Informationen siehe Thema Verbesserte C5.0-Modelle in Kapitel 6 auf S. 190.](#)

**Basierend auf.** Um Scores für die angepasste Neigung berechnen zu können, muss im Stream ein Partitionsfeld vorhanden sein. Sie können angeben, ob die Test- bzw. Validierungspartition für diese Berechnung verwendet werden soll. Um bestmögliche Ergebnisse zu erzielen, sollte die Test- bzw. Validierungspartition mindestens so viele Datensätze enthalten wie die Partition, die zum Trainieren des ursprünglichen Modells verwendet wurde. [Für weitere Informationen siehe Thema Partitionsknoten in Kapitel 4 in IBM SPSS Modeler 14.2- Quellen-, Prozess- und Ausgabeknoten.](#)

## KNN-Modell-Nugget

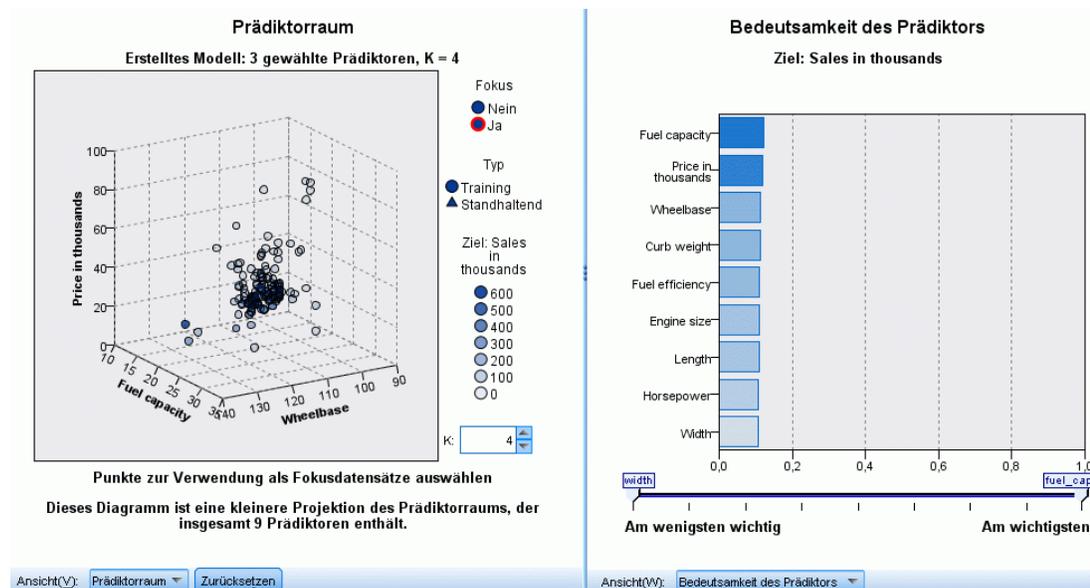
Das KNN-Modell erstellt eine Reihe neuer Felder wie in der folgenden Tabelle gezeigt. Um diese Felder und die zugehörigen Werte anzuzeigen, müssen Sie einen Tabellenknoten zum KNN-Modell-Nugget hinzufügen und den Tabellenknoten ausführen oder auf die Schaltfläche “Vorschau” am Nugget klicken.

Tabelle 16-1  
KNN-Modellfelder

Neuer Feldname	Beschreibung
\$KNN-Feldname	Vorhergesagter Wert des Zielfelds.
\$KNNP-Feldname	Wahrscheinlichkeit des vorhergesagten Werts.
\$KNNP-Wert	Wahrscheinlichkeit jedes möglichen Werts von nominalen oder Flag-Felds. Wird nur angezeigt, wenn auf der Registerkarte “Einstellungen” des Modell-Nuggets die Option Alle Wahrscheinlichkeiten anhängen aktiviert ist.
\$KNN-Nachbar- <i>n</i>	Der Name des <i>n.</i> nächsten Nachbarn zum Fokusdatensatz. Wird nur eingeschlossen, wenn für Anzeigen von Nearest in der Registerkarte “Einstellungen” des Modell-Nuggets auf ein Wert ungleich null angegeben ist.
\$KNN-Distanz- <i>n</i>	Die relative Distanz des Fokusdatensatzes vom <i>n.</i> nächsten Nachbarn. Wird nur eingeschlossen, wenn für Anzeigen von Nearest in der Registerkarte “Einstellungen” des Modell-Nuggets auf ein Wert ungleich null angegeben ist.

## Modellansicht

Abbildung 16-8  
Modellansicht für die Analyse nächstgelegener Nachbar



Das Fenster der Modellansicht setzt sich aus zwei Bereichen zusammen:

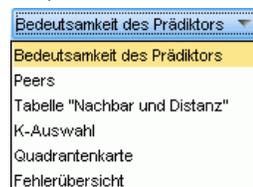
- Im ersten Bereich wird eine Übersicht des Modells, die sogenannte Hauptansicht, angezeigt.
- Im zweiten Bereich wird eine der beiden folgenden Ansichten angezeigt:

Die Hilfsmodellansicht enthält mehr Informationen zum Modell, ist dafür aber weniger stark auf das Modell an sich konzentriert.

Die verknüpfte Ansicht zeigt Details zu einer bestimmten Funktion des Modells an, wenn der Benutzer einen Teil der Hauptansicht ansteuert.

Standardmäßig wird im ersten Bereich der Prädiktorbereich und im zweiten Bereich das Diagramm für die Bedeutsamkeit der Prädiktoren angezeigt. Wenn das Diagramm für die Bedeutsamkeit der Prädiktoren nicht verfügbar ist (d. h. wenn Merkmale nach Bedeutung gewichten nicht im Bereich "Nachbarn" der Registerkarte "Einstellungen" ausgewählt wurde), wird die erste verfügbare Ansicht aus dem Dropdown-Menü "Ansicht" angezeigt.

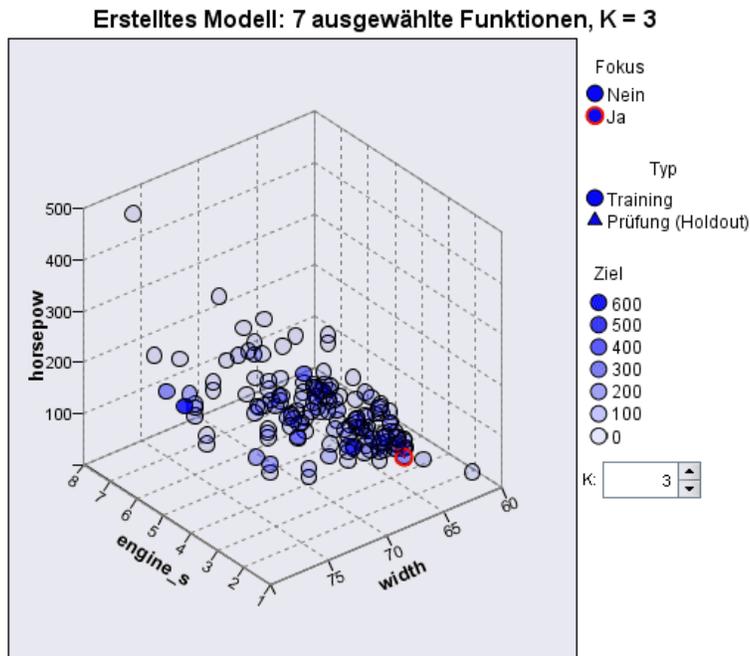
Abbildung 16-9  
Dropdown-Liste "Modellansicht" im Dialogfeld "Analyse nächstgelegener Nachbar"



Wenn für eine Ansicht keine Informationen verfügbar sind, wird sie im Dropdown-Menü "Ansicht" nicht angezeigt.

## Prädiktorbereich

Abbildung 16-10  
Prädiktorbereich



Das Prädiktorbereichsdiagramm ist ein interaktives Diagramm für den Prädiktorbereich (bzw. -unterbereich, bei mehr als drei Prädiktoren). Jede Achse stellt einen Prädiktor im Modell dar und die Position der Punkte in der Tabelle gibt die Werte dieser Prädiktoren für Fälle in den Trainings- und Holdout-Partitionen an.

**Erläuterungen** Neben den Prädiktorwerten liefern die Punkte im Diagramm weitere Informationen.

- Die Form gibt die Partition an, zu der ein Punkt gehört (Training oder Holdout).
- Die Farbe/Schattierung eines Punkts gibt den Wert des Ziels für diesen Fall an. Dabei entsprechen eindeutige Farbwerte den Kategorien eines kategorialen Ziels und Schattierungen dem Wertebereich eines stetigen Ziels. Für Trainings-Partitionen ist der angegebene Wert der festgestellte Wert. Für Holdout-Partitionen handelt es sich um den vorhergesagten Wert. Wenn kein Ziel angegeben ist, wird diese Erläuterung nicht angezeigt.
- Kräftigere Umrisse weisen auf Fokusfälle hin. Fokusdatensätze werden im Zusammenhang mit ihren  $k$  nächstgelegenen Nachbarn angezeigt.

**Steuerelemente und Interaktivität.** Sie können den Prädiktorbereich mit einer Reihe an Steuerelementen im Diagramm untersuchen.

- Sie können festlegen, welche Untermenge an Prädiktoren im Diagramm angezeigt werden soll, und ändern, welche Prädiktoren in den Dimensionen dargestellt werden.

- “Fokusdatensätze” sind Punkte, die im Prädiktorbereichsdiagramm ausgewählt wurden. Wenn Sie eine Fokusdatensatzvariable angegeben haben, werden zuerst die Punkte ausgewählt, die die Fokusdatensätze darstellen. Es kann jedoch jeder Punkt vorübergehend ein Fokusdatensatz werden, wenn Sie ihn auswählen. Die gängigen Steuerelemente für Punkte sind verfügbar: Wenn Sie auf einen Punkt klicken, wird dieser Punkt ausgewählt und die Auswahl aller anderen Punkte aufgehoben. Wenn Sie die Strg-Taste drücken und auf einen Punkt klicken, wird er der Menge an gewählten Punkten hinzugefügt. Verknüpfte Ansichten wie das Peers-Diagramm werden automatisch mit den Fällen aktualisiert, die im Prädiktorbereich ausgewählt werden.
- Sie können die Anzahl an für Fokusdatensätze anzuzeigenden nächstgelegenen Nachbarn ( $k$ ) ändern.
- Wenn Sie die Maus über einen Punkt im Diagramm bewegen, wird eine QuickInfo mit dem Wert der Fallbeschriftung oder, wenn keine Fallbeschriftungen definiert sind, der Fallnummer und dem festgestellten und vorhergesagten Zielwert angezeigt.
- Sie können den Prädiktorbereich über die Schaltfläche “Zurücksetzen” wieder in seinen Originalzustand versetzen.

### **Ändern der Achsen im Prädiktorbereichsdiagramm**

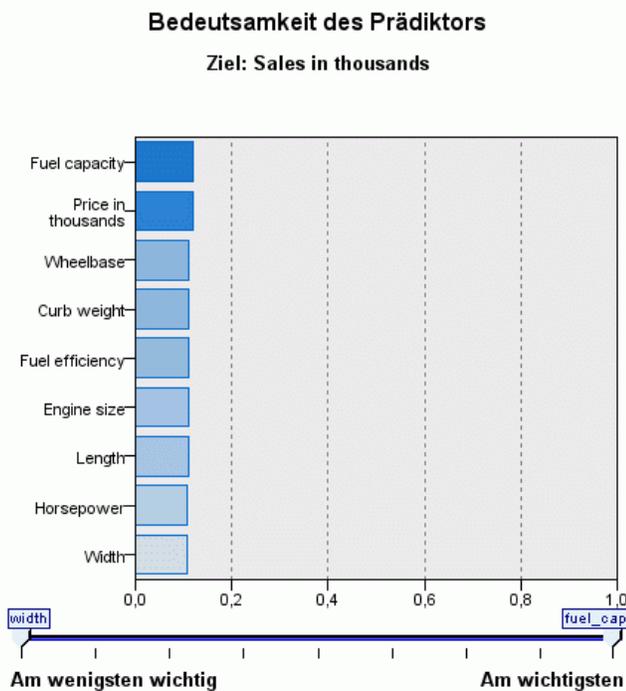
Sie können steuern, welche Funktionen an den Achsen im Prädiktorbereichsdiagramm angezeigt werden.

#### **So ändern Sie die Achseneinstellungen:**

- ▶ Klicken Sie auf die Schaltfläche für den Bearbeitungsmodus (Pinselsymbol) im linken Bereich, um den Modus “Bearbeiten” für den Prädiktorbereich zu wählen.
- ▶ Ändern Sie die Ansicht (beliebig) im rechten Bereich. Der Bereich Zonen anzeigen wird zwischen den beiden Hauptbereichen angezeigt.
- ▶ Klicken Sie auf das Kontrollkästchen Zonen anzeigen.
- ▶ Klicken Sie auf einen beliebigen Datenpunkt im Prädiktorbereich.
- ▶ So ersetzen Sie einen Prädiktor durch eine Funktion desselben Datentyps:
  - ▶ Ziehen Sie den neuen Prädiktor über das Zonenlabel (mit der kleinen X-Schaltfläche), das Sie ersetzen möchten.
- ▶ So ersetzen Sie eine Achse durch einen Prädiktor eines anderen Datentyps:
  - Klicken Sie am Zonenlabel des zu ersetzenden Prädiktors auf die kleine X-Schaltfläche. Der Prädiktorbereich ändert sich in eine zweidimensionale Ansicht.
  - Ziehen Sie den neuen Prädiktor auf das Zonenlabel Dimension hinzufügen.
- ▶ Klicken Sie auf die Schaltfläche für den Interaktionsmodus (Pfeilspitzensymbol) im linken Bereich, um den Modus “Bearbeiten” zu verlassen.

### Bedeutsamkeit des Prädiktors

Abbildung 16-11  
Bedeutsamkeit des Prädiktors



Normalerweise ist es sinnvoll, die Modellierungsbemühungen auf die wichtigsten Prädiktorfelder zu konzentrieren und diejenigen Prädiktoren zu verwerfen bzw. zu ignorieren, die die geringste Bedeutung haben. Dabei unterstützt Sie das Diagramm für die Bedeutsamkeit der Prädiktoren, da es die relative Bedeutsamkeit der einzelnen Prädiktoren für das Modell angibt. Da die Werte relativ sind, beträgt die Summe der Werte aller Prädiktoren im Diagramm 1,0. Die Bedeutsamkeit der Prädiktoren bezieht sich nicht auf die Genauigkeit des Modells. Sie bezieht sich lediglich auf die Bedeutsamkeit der einzelnen Prädiktoren für eine Vorhersage und nicht auf die Genauigkeit der Vorhersage.

### Abstände zwischen nächstgelegenen Nachbarn

Abbildung 16-12  
Abstände zwischen nächstgelegenen Nachbarn

**k-Nächste-Nachbarn und Distanzen**  
Angezeigt für Anfangsfokusdatensätze

Fokusdatensatz	Nächste Nachbarn				Kürzeste Distanzen			
	1	2	3	4	1	2	3	4
101	106	76	70	61	0,157	0,161	0,171	0,173

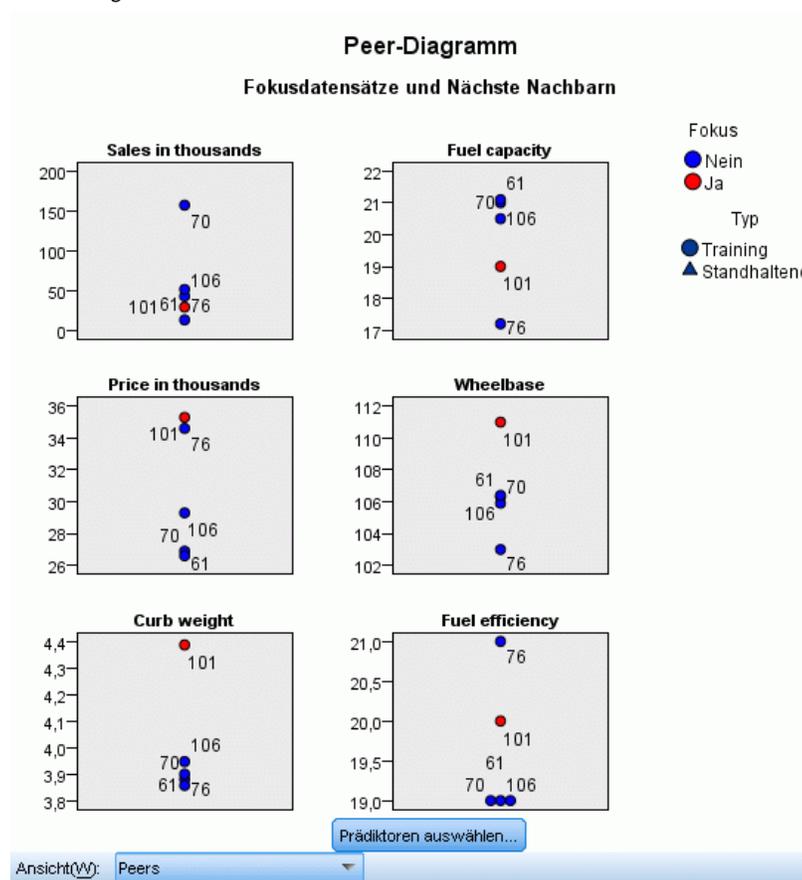
Diese Tabelle zeigt nur die  $k$  nächstgelegenen Nachbarn und Abstände für Fokusdatensätze an. Sie ist verfügbar, wenn eine Fokusdatensatz-ID im Modellierungsknoten angegeben wurde, und zeigt nur Fokusdatensätze, die durch diese Variable identifiziert werden.

Jede Zeile der:

- Die Spalte Fokusdatensatz enthält den Wert der Fallbeschriftungsvariablen für den Fokusdatensatz. Wenn keine Fallbeschriftungen angegeben wurden, enthält diese Spalte die Fallnummer des Fokusdatensatzes.
- Die  $i$ -te Spalte unter der Gruppe Nächste Nachbarn enthält den Wert der Fallbeschriftungsvariablen für den  $i$ -ten nächsten Nachbarn des Fokusdatensatzes. Wenn keine Fall-Labels definiert sind, enthält diese Spalte die Fallnummer des  $i$ -ten nächsten Nachbarn des Fokusdatensatzes.
- Die  $i$ -te Spalte unter der Gruppe Nächste Distanzen enthält die Distanz des  $i$ -ten nächsten Nachbarn zum Fokusdatensatz.

## Gruppen

Abbildung 16-13  
Peers-Diagramm



Dieses Diagramm enthält die Fokusfälle und ihre  $k$  nächstgelegenen Nachbarn für jeden Prädiktor im Ziel. Es ist verfügbar, wenn ein Fokusfall im Prädiktorbereich ausgewählt ist.

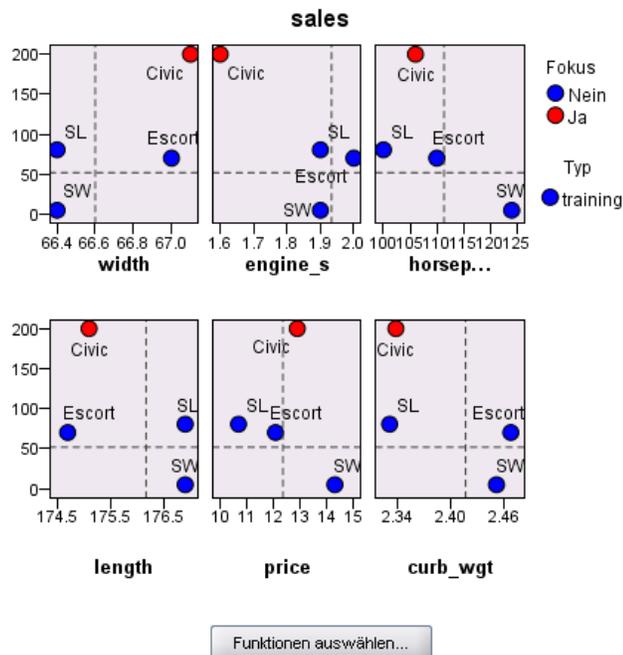
Das Vergleichsdiagramm ist auf zwei Arten mit dem Prädiktorbereich verknüpft.

- Im Peers-Diagramm werden die im Prädiktorbereich gewählten Fokusfälle sowie ihre  $k$  nächstgelegenen Nachbarn angezeigt.
- Der Wert  $k$  wird im Prädiktorbereich gewählt und im Peers-Diagramm herangezogen.

**Prädiktoren auswählen.** Ermöglicht Ihnen, die Prädiktoren für die Anzeige im Vergleichsdiagramm auszuwählen.

### Quadrantenkarte

Abbildung 16-14  
Quadrantenkarte



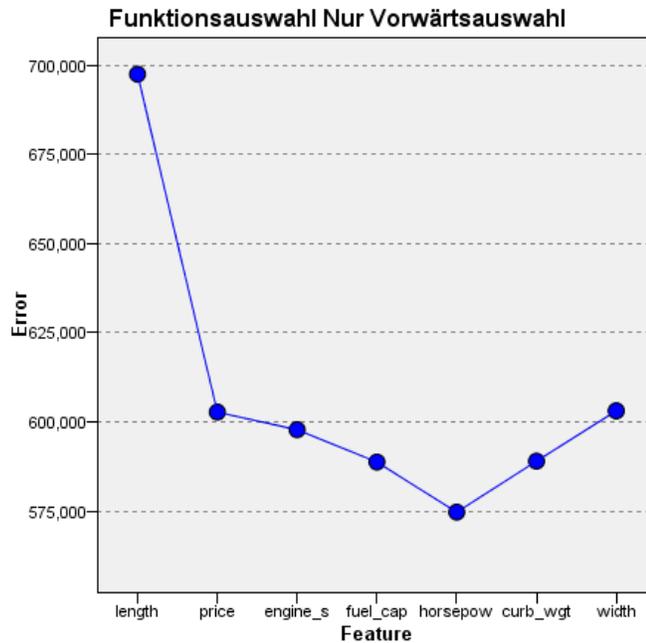
Dieses Diagramm zeigt die Fokusfälle und ihre  $k$  nächstgelegenen Nachbarn als Streudiagramm (oder Punktdiagramm, je nach Messniveau des Ziels) mit dem Ziel auf der  $y$ -Achse und einem metrischen Prädiktor auf der  $x$ -Achse nach Prädiktoren in einzelne Felder unterteilt an. Es ist verfügbar, wenn ein Ziel vorhanden und ein Fokusfall im Prädiktorbereich ausgewählt ist.

- Für stetige Variablen werden bei den Mittelwerten der Variablen in der Trainingspartition Referenzlinien gezogen.

**Prädiktoren auswählen.** Ermöglicht Ihnen, die Prädiktoren für die Anzeige in der Quadrantenkarte auszuwählen.

### Prädiktor-Auswahlfehler-Protokoll

Abbildung 16-15  
Prädiktorauswahl



Punkte im Diagramm zeigen den Fehler (je nach Messniveau des Ziels entweder die Fehlerrate oder den Quadratsummenfehler) auf der  $y$ -Achse für das Modell mit dem Prädiktor auf der  $x$ -Achse an (plus allen Prädiktoren weiter links auf der  $x$ -Achse). Dieses Diagramm ist verfügbar, wenn ein Ziel und eine Funktionsauswahl aktiviert sind.

### Klassifikationsmatrix

Abbildung 16-18  
Klassifikationsmatrix

Partition		Vorhergesagt		
		0	1	Prozent korrekt
Training	0	111	1	99.11%
	1	7	33	82.50%
	Prozent (insgesamt)	77.64%	22.37%	94.74%

Diese Tabelle enthält die Kreuzklassifikation der festgestellten Werte im Vergleich zu den vorhergesagten Werten des Ziels nach Partitionen. Verfügbar, wenn ein Ziel vorhanden ist und es kategorial ist (Flag, nominal oder ordinal).

- Die Zeile (Fehlend) in der Holdout-Partition enthält Holdout-Fälle mit fehlenden Werten im Ziel. Diese Fälle tragen zur Prüfstichprobe bei: Gesamtprozentwerte, aber nicht die Werte für "Prozent korrekt".

### Fehlerzusammenfassung

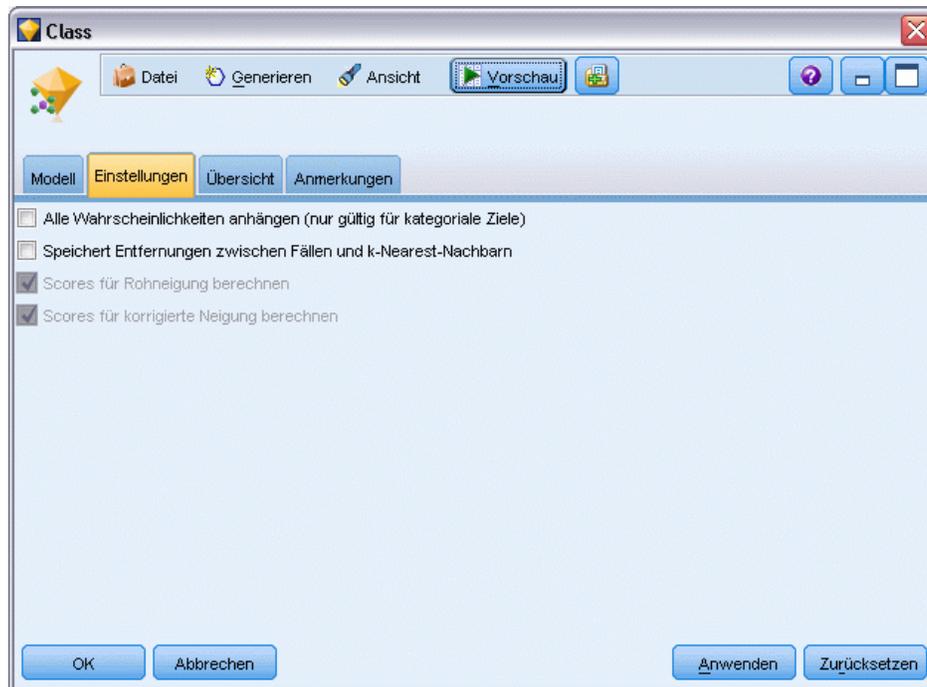
Abbildung 16-19  
Fehlerzusammenfassung

Partition	Sum-of-Squares Error
Training	622043

Diese Tabelle ist verfügbar, wenn eine Zielvariable vorhanden ist. Sie enthält die Fehler für das Modell, Quadratsummenfehler für stetige Ziele und die Fehlerrate (100 % – Gesamtprozent korrekt) für kategoriale Ziele.

### KNN-Modell-Einstellungen

Abbildung 16-20  
KNN-Modell-Nugget - Einstellungen



In der Registerkarte “Einstellungen” können Sie zusätzliche Felder angeben, die bei der Anzeige der Ergebnisse verwendet werden sollen (z. B. durch Ausführen eines Tabellenknotens, der mit dem Nugget verknüpft ist). Sie können den Effekt jeder dieser Optionen sehen, indem Sie sie auswählen und auf die Schaltfläche “Vorschau” klicken. Führen Sie in der Vorschau-Ausgabe einen Bildlauf nach rechts durch, um die zusätzlichen Felder zu sehen.

**Alle Wahrscheinlichkeiten anhängen (nur gültig für kategoriale Ziele).** Wenn diese Option ausgewählt ist, werden für jeden Datensatz, der vom Knoten verarbeitet wird, die Wahrscheinlichkeiten für jeden möglichen Wert eines Zielfelds vom Typ “Nominal” oder “Flag” angezeigt. Wenn diese

Option nicht ausgewählt ist, werden für Zielfelder vom Typ “Flag” oder “Nominal” nur der vorhergesagte Wert und seine Wahrscheinlichkeit angezeigt.

Die Standardeinstellung dieses Kontrollkästchens richtet sich nach dem entsprechenden Kontrollkästchen des Modellierungsknotens.

**Scores für Rohneigung berechnen.** Bei Modellen mit einem Flag-Ziel (das als Vorhersage “Ja” bzw. “Nein” ausgibt) können Sie Neigungs-Scores anfordern, die die Wahrscheinlichkeit des für das Zielfeld angegebenen wahren Ergebnisses angeben. Diese Werte werden zusätzlich zu den anderen Vorhersage- und Konfidenzwerten angegeben, die ggf. während des Scorings erstellt werden.

**Scores für korrigierte Neigung berechnen.** Rohneigungs-Scores beruhen ausschließlich auf den Trainingsdaten und sind aufgrund der Neigung vieler Modelle zur übermäßigen Anpassung an die Trainingsdaten möglicherweise zu optimistisch. Bei korrigierten Neigungen wird versucht, dies durch Evaluation der Modellleistung anhand einer Test- bzw. Validierungspartition zu kompensieren. Bei dieser Option muss im Stream ein Partitionsfeld definiert sein und die Scores für die korrigierte Neigung müssen im Modellierungsknoten aktiviert werden, bevor das Modell generiert wird.

**Anzeigen von Nearest.** Wenn Sie diesen Wert auf  $n$  einstellen (mit  $n$  als positive Ganzzahl ungleich null), werden die  $n$  nächsten Nachbarn des Fokusdatensatzes zusammen mit ihren Distanzen vom Fokusdatensatz in das Modell aufgenommen.

## **Hinweise**

This information was developed for products and services offered worldwide.

IBM may not offer the products, services, or features discussed in this document in other countries. Consult your local IBM representative for information on the products and services currently available in your area. Any reference to an IBM product, program, or service is not intended to state or imply that only that IBM product, program, or service may be used. Any functionally equivalent product, program, or service that does not infringe any IBM intellectual property right may be used instead. However, it is the user's responsibility to evaluate and verify the operation of any non-IBM product, program, or service.

IBM may have patents or pending patent applications covering subject matter described in this document. The furnishing of this document does not grant you any license to these patents. You can send license inquiries, in writing, to:

*IBM Director of Licensing, IBM Corporation, North Castle Drive, Armonk, NY 10504-1785, U.S.A.*

For license inquiries regarding double-byte character set (DBCS) information, contact the IBM Intellectual Property Department in your country or send inquiries, in writing, to:

*Intellectual Property Licensing, Legal and Intellectual Property Law, IBM Japan Ltd., 1623-14, Shimotsuruma, Yamato-shi, Kanagawa 242-8502 Japan.*

**Der folgende Absatz gilt nicht für Großbritannien oder andere Länder, in denen derartige Bestimmungen nicht mit dem dort geltenden Recht vereinbar sind.** SPSS INC., AN IBM COMPANY, ÜBERNIMMT FÜR DIE VORLIEGENDE DOKUMENTATION KEINERLEI GEWÄHRLEISTUNG IRGENDWELCHER ART, WEDER AUSDRÜCKLICH NOCH STILLSCHWEIGEND, EINSCHLIESSLICH (JEDOCH NICHT DARAUF BEGRENZT) DER STILLSCHWEIGENDEN GEWÄHRLEISTUNGEN IN BEZUG AUF DIE NICHTVERLETZUNG VON RECHTEN DRITTER, AUF HANDELSÜBLICHKEIT ODER DIE EIGNUNG FÜR EINEN BESTIMMTEN ZWECK. Einige Staaten lassen bei bestimmten Transaktionen keine Ausschlussklauseln ausdrücklicher oder stillschweigender Gewährleistungen zu, sodass diese Erklärung möglicherweise nicht auf Sie zutrifft.

Diese Informationen können technische Ungenauigkeiten oder typografische Fehler enthalten. An den hierin enthaltenen Informationen werden in regelmäßigen Abständen Änderungen vorgenommen, die in spätere Ausgaben der Publikation eingearbeitet werden. SPSS Inc. kann jederzeit ohne Vorankündigung Verbesserungen und/oder Veränderungen an den in dieser Publikation beschriebenen Produkten und/oder Programmen vornehmen.

Alle in diesen Ausführungen enthaltenen Verweise auf Websites, die nicht zu SPSS bzw. IBM gehören, dienen lediglich der Information. Die Nennung bedeutet nicht, dass SPSS bzw. IBM den Inhalt dieser Websites unterstützen. Das Material auf diesen Websites ist kein Bestandteil des Materials für dieses SPSS Inc.-Produkt. Sie verwenden diese Websites auf eigenes Risiko.

Wenn Sie Informationen an IBM bzw. SPSS senden, räumen Sie IBM und SPSS das nicht ausschließliche Recht ein, die Informationen in jeglicher Form zu verwenden bzw. weiterzugeben, die dem Unternehmen geeignet erscheint, ohne dass ihm daraus Verbindlichkeiten Ihnen gegenüber entstehen.

Informationen zu Nicht-SPSS-Produkten stammen von den Herstellern dieser Produkte, ihren veröffentlichten Verlautbarungen oder aus anderen öffentlich verfügbaren Quellen. SPSS hat diese Produkte nicht getestet und kann daher die Richtigkeit der Angaben zu Leistung und Kompatibilität oder anderer Behauptungen in Bezug auf Nicht-SPSS-Produkte nicht bestätigen. Fragen zu den Fähigkeiten von Nicht-SPSS-Produkten sind an die Hersteller dieser Produkte zu richten.

Licensees of this program who wish to have information about it for the purpose of enabling: (i) the exchange of information between independently created programs and other programs (including this one) and (ii) the mutual use of the information which has been exchanged, should contact:

*IBM Software Group, Attention: Licensing, 233 S. Wacker Dr., Chicago, IL 60606, USA.*

Such information may be available, subject to appropriate terms and conditions, including in some cases, payment of a fee.

The licensed program described in this document and all licensed material available for it are provided by IBM under terms of the IBM Customer Agreement, IBM International Program License Agreement or any equivalent agreement between us.

Any performance data contained herein was determined in a controlled environment. Therefore, the results obtained in other operating environments may vary significantly. Some measurements may have been made on development-level systems and there is no guarantee that these measurements will be the same on generally available systems. Furthermore, some measurements may have been estimated through extrapolation. Actual results may vary. Users of this document should verify the applicable data for their specific environment.

Information concerning non-IBM products was obtained from the suppliers of those products, their published announcements or other publicly available sources. IBM has not tested those products and cannot confirm the accuracy of performance, compatibility or any other claims related to non-IBM products. Questions on the capabilities of non-IBM products should be addressed to the suppliers of those products.

All statements regarding IBM's future direction or intent are subject to change or withdrawal without notice, and represent goals and objectives only.

Diese Informationen enthalten Beispiele für Daten und Berichte, die in alltäglichen Betriebsabläufen verwendet werden. Um sie möglichst umfassend darzulegen, enthalten die Beispiele Namen von Einzelpersonen, Unternehmen, Marken und Produkten. Alle diese Namen sind frei erfunden und jegliche Ähnlichkeit mit Namen und Adressen, die von einem tatsächlichen Handelsunternehmen verwendet werden, ist rein zufällig.

If you are viewing this information softcopy, the photographs and color illustrations may not appear.

**Marken**

IBM, das IBM-Logo und ibm.com sind Marken von IBM Corporation, die in vielen Ländern weltweit eingetragen sind. Eine aktuelle Liste der IBM-Marken finden Sie im Internet unter <http://www.ibm.com/legal/copytrade.shtml>.

SPSS ist eine Marke von SPSS Inc., an IBM Company, die in vielen Ländern weltweit eingetragen sind.

Adobe, das Adobe-Logo, PostScript und das PostScript-Logo sind entweder registrierte Marken oder Marken von Adobe Systems Incorporated in den USA und/oder anderen Ländern.

IT Infrastructure Library ist eine eingetragene Marke der Central Computer and Telecommunications Agency, die nun zum Office of Government Commerce gehört.

Intel, das Intel-Logo, Intel Inside, das Intel Inside-Logo, Intel Centrino, das Intel Centrino-Logo, Celeron, Intel Xeon, Intel SpeedStep, Itanium und Pentium sind Marken oder eingetragene Marken der Intel Corporation oder ihrer Tochtergesellschaften in den USA und anderen Ländern.

Linux ist eine eingetragene Marke von Linus Torvalds in den USA und/oder anderen Ländern.

Microsoft, Windows, Windows NT und das Windows-Logo sind Marken von Microsoft Corporation in den USA und/oder anderen Ländern.

ITIL ist eine eingetragene Marke und eine eingetragene Gemeinschaftsmarke des Office of Government Commerce und ist beim U.S. Patent and Trademark Office eingetragen.

UNIX ist eine eingetragene Marke von The Open Group in den USA und anderen Ländern.

Cell Broadband Engine ist eine Marke von Sony Computer Entertainment, Inc. in den USA und/oder anderen Ländern und wird im Rahmen einer Lizenz dieses Unternehmens verwendet.

Java und alle Java-basierten Marken und Logos sind Marken von Sun Microsystems, Inc. in den USA und/oder anderen Ländern.

Linear Tape-Open, LTO, the LTO Logo, Ultrium, and the Ultrium logo are trademarks of HP, IBM Corp. and Quantum in the U.S. and other countries.

Weitere Produkt- oder Servicenamen können Marken von IBM, SPSS oder anderen Unternehmen sein.

- A Priori-Modelle
  - Evaluierungsmaße, 392
  - Expertenoptionen, 392
  - Modellierungsknoten, 390
  - Modellierungsknotenoptionen, 391
  - Tabellen- im Vergleich zu Transaktionsdaten, 37
- A-priori-Wahrscheinlichkeit
  - Entscheidungsbäume, 167
- A-priori-Wahrscheinlichkeiten, 167
- absolute Konfidenz-Differenz zum Vorgänger
  - A Priori-Evaluierungsmaß, 393
- Abstand
  - ACF und PACF, 436
- Abstände zwischen nächstgelegenen Nachbarn
  - in der Nächste-Nachbarn-Analyse, 498
- Additive Ausreißer, 436
  - Patches, 436
  - Zeitreihenmodellierung, 452
- Akaike-Informationskriterium
  - in linearen Modellen, 271
- Aktualisieren von Maßen, 257
- Aktualisieren von Modellen
  - Selbstlern-Antwortmodelle, 465
- Algorithmen, 46, 127
- Allgemeine schätzbare Funktion
  - Verallgemeinerte lineare Modelle, 338
- Alternative Modelle, 253
- Alternativen, Registerkarte, 240
- Analyse Nächstgelegener Nachbar
  - Modellansicht, 495
- Anomalieerkennungmodelle, 89
  - Anomalie-Index, 85
  - Anomaliefelder, 85, 91
  - Anpassungskoeffizient, 86
  - Bewertung, 88, 91
  - Cutoff-Wert, 85, 90
  - Ebene des Rauschens, 86
  - Fehlende Werte, 86
  - Gruppen, 86, 90
  - Modellierungsknoten, 83
- ANOVA
  - in linearen Modellen, 281
- Anpassen eines Modells, 253
- Anpassungsgüte des Modells
  - Logistische Regressionsmodelle, 309
- Antezedens
  - Regeln ohne, 399
- Anwendungsbeispiele, 3
- Arbeitsmodellbereich, 238
- ARIMA-Modelle, 439
  - Ausreißer, 452
  - autoregressive Ordnungen, 448
  - Differenzierungsordnungen, 448
  - Konstante, 448
  - Kriterien für Zeitreihenmodelle, 447
  - Ordnungen des gleitenden Durchschnitts, 448
  - saisonale Ordnungen, 448
  - Übertragungsfunktionen, 450
- Assoziationsregelmodelle, 188, 194, 196, 388, 422, 424, 427
  - Angeben von Filtern, 404
  - apriori, 390
  - bereitstellen, 413
  - CARMA, 394
  - Diagrammerstellung, 405
  - Einstellungen, 407
  - Erstellen eines gefilterten Modells, 411
  - Für Sequenzen, 416
  - Generieren einer Regelmenge, 410
  - IBM InfoSphere Warehouse, 37
  - Modell-Nugget, 400
  - Modell-Nugget-Übersicht, 409
  - Nähere Informationen zum Modell-Nugget, 401
  - Scoring-Regeln, 412
  - Transponieren von Scores, 413
- Asymptotische Korrelation
  - Logistische Regressionsmodelle, 298, 309
- Asymptotische Kovarianz
  - Logistische Regressionsmodelle, 298
- Aufbaueregeln, 161
  - C&R-Baum, Knoten, 148
  - CHAID-Knoten, 148–149
  - Entscheidungsbäume, 150
  - QUEST-Knoten, 148
- Aufgeteilte Modelle, 485
  - Betroffene Merkmale, 35
  - Erstellen, 31
  - im Vergleich zu partitionierten Modellen, 33
  - Modellierungsknoten, 34
- Aufteilungen
  - Entscheidungsbäume, 132, 134
- Aufteilungen, Kreuzvalidierung, 490
- Ausführen einer Mining-Aufgabe, 244
- Ausreißer, 435
  - Additive Patches, 436
  - ARIMA-Modelle, 452
  - deterministisch, 435
  - Ebenenänderung, 436
  - erkennen, 83
  - Expertenmodellierung, 445
  - in Reihen, 434
  - in Zeitreihenmodellen, 452
  - innovatorisch, 436
  - lokaler Trend , 436
  - saisonal additiv, 436
  - vorübergehende Änderung, 436
- Ausreißer mit lokalem Trend , 436
  - Zeitreihenmodellierung, 452
- Ausreißer mit vorübergehender Änderung, 436
- Auswahlknoten
  - Generieren von Entscheidungsbäumen, 151
- Auswerten eines Modells, 256

- Auswertung in Excel, 258
- Auto-Numerisch, Modelle, 93
  - Abbruchregeln, 95, 111
  - Algorithmuseinstellungen, 94
  - Einstellungen, 113
  - Ergebnis-Browser, Fenster, 120
  - Evaluationsdiagramme, 123–124
  - Generieren von Modellierungsknoten und Nuggets, 123
  - Modell-Nugget, 120
  - Modellierungsknoten, 107–108
  - Modellierungsoptionen, 108
  - Modelltypen, 111
- Autokorrelationsfunktion
  - Reihen, 436
- Autom. Cluster, Modelle, 93
  - Abbruchregeln, 95
  - Algorithmuseinstellungen, 94
  - Ergebnis-Browser, Fenster, 120
  - Evaluationsdiagramme, 123
  - Generieren von Modellierungsknoten und Nuggets, 123
  - Modell-Nugget, 120
  - Modellierungsknoten, 114–115
  - Modelltypen, 117
  - Partitionen, 117
  - Rangenteilung von Modellen, 115
  - Verwerfen von Modellen, 119
- Automatische Datenaufbereitung
  - in linearen Modellen, 276
- Automatischer Klassifizierer, Modelle, 93
  - Abbruchregeln, 95
  - Algorithmuseinstellungen, 94
  - Einführung, 96
  - Einstellungen, 105
  - Ergebnis-Browser, Fenster, 120
  - Evaluationsdiagramme, 123–124
  - Generieren von Modellierungsknoten und Nuggets, 123
  - Modell-Nugget, 120
  - Modellierungsknoten, 96, 98
  - Modelltypen, 100
  - Partitionen, 100
  - Rangenteilung von Modellen, 98
  - Verwerfen von Modellen, 104
- Automatisierte Modellierung, Knoten
  - Auto-Numerisch, Modelle, 93
  - Autom. Cluster, Modelle, 93
  - Automatischer Klassifizierer, Modelle, 93
- Autoregression
  - ARIMA-Modelle, 448
- Basiskategorie
  - Logistiksknoten, 293
- Baumbasierte Analyse
  - Allgemeine Verwendung, 128
- Baumstruktur
  - Diagrammerstellung, 191
  - Entscheidungsbaum-Modelle, 187
- Baumtiefe, 162
- Bayes-Netzwerk-Modelle
  - Expertenoptionen, 203
  - Modell-Nugget, 205
  - Modell-Nugget-Einstellungen, 206
  - Modell-Nugget-Übersicht, 208
  - Modellierungsknoten, 198
  - Modelloptionen, 200
- Bedeutsamkeit des Prädiktors
  - Diskriminanzmodelle, 325
  - Entscheidungsbaum-Modelle, 183
  - Filtern von Feldern, 57
  - in der Nächste-Nachbarn-Analyse, 498
  - Lineare Modelle, 277
  - Logistische Regressionsmodelle, 303
  - Modellergebnisse, 41, 54, 57
  - Neuronale Netzwerke, 220
  - Verallgemeinerte lineare Modelle, 339
- Beispiele
  - Anwendungshandbuch, 3
  - Übersicht, 4
- Benutzerdefinierte Aufteilungen
  - Entscheidungsbaume, 132, 134
- Beschriftungen
  - value, 71
  - Variablen, 71
- Beste Untergruppen
  - in linearen Modellen, 271
- Betrugserkennung
  - Anomalieerkennung, 83
- Binomiale logistische Regression, Modelle, 287–288
- Bonferroni-Korrektur
  - CHAID-Knoten, 173
- Box' M-Test
  - Diskriminanzknoten, 322
- C&R-Baum-Modelle
  - A-priori-Wahrscheinlichkeit, 167
  - Baumtiefe, 162
  - beschneiden, 162
  - Diagrammerstellung aus dem Modellnugget, 191
  - Ensembles, 165
  - Fallgewichtungen, 39
  - Fehlklassifizierungskosten, 166
  - Feldoptionen, 158
  - Grenzooptionen, 164
  - Häufigkeitsgewichtung, 39
  - Modell-Nugget, 181
  - Modellierungsknoten, 129, 154–155, 187–188
  - Surrogate, 163
  - Unreinheitsmaße, 169
  - Ziele, 159
- C5.0-Modelle, 127
  - beschneiden, 179
  - Diagrammerstellung aus dem Modellnugget, 191
  - Fehlklassifizierungskosten, 179
  - Leistung, 178, 181
  - Modell-Nugget, 181, 194, 196

- Modellierungsknoten, 177, 179, 187–188, 190
- Optionen, 179
- Parallele Verarbeitung, 178, 181
- Verbesserung, 179, 190
- CARMA-Modelle
  - Datenformate, 395
  - Expertenoptionen, 399
  - Feldoptionen, 395
  - ID-Feld, 395
  - Inhaltsfeld(er), 395
  - Mehrere Sukzedenzien, 413
  - Modellierungsknoten, 394
  - Modellierungsknotenoptionen, 398
  - Tabellen- im Vergleich zu Transaktionsdaten, 399
  - Zeitfeld, 395
- CHAID-Modelle, 127
  - Baumtiefe, 162
  - Diagrammerstellung aus dem Modellnugget, 191
  - Ensembles, 165
  - Exhaustive CHAID, 162
  - Fehlklassifizierungskosten, 168
  - Feldoptionen, 158
  - Grenzooptionen, 164
  - Modell-Nugget, 181
  - Modellierungsknoten, 129, 154, 156, 187–188
  - Ziele, 159
- Chi-Quadrat
  - CHAID-Knoten, 174
  - Funktionsauswahl, 79
- Chi-Quadrat nach Pearson
  - CHAID-Knoten, 174
  - Funktionsauswahl, 79
- Cluster-Viewer
  - Ansicht Bedeutsamkeit des Prädiktors im Cluster, 378
  - Ansicht Clustervergleich, 381
  - Ansicht Clusterzentrum, 375
  - Ansicht Zellverteilung, 380
  - Anzeige Zelleninhalt, 377
  - Basisansicht, 378
  - Bedeutsamkeit des Prädiktors, 378
  - Cluster sortieren, 377
  - Cluster und Merkmale transponieren, 376
  - Cluster und Merkmale vertauschen, 376
  - Clusteransicht, 375
  - Clusteranzeige sortieren., 377
  - Clustergrößen, 379
  - Clustergrößenansicht, 379
  - Diagrammerstellung, 385
  - Merkmalanzeige sortieren, 377
  - Merkmale sortieren, 377
  - Modellübersicht, 374
  - über Clustermodelle, 371
  - Übersicht, 372
  - Übersichtsansicht, 374
  - Vergleich von Clustern, 381
  - Verteilung der Zellen, 380
  - verwenden, 382
  - Zelleninhalt sortieren, 377
- Clusteranalyse
  - Anomalieerkennung, 86
  - Anzahl der Cluster, 368
- clustering, 354
- Clustering, 355, 362, 365–367, 370–371
  - Cluster anzeigen, 372
  - Gesamtanzeige, 372
- confidence
  - A Priori-Knoten, 391
  - Assoziationsregeln, 402, 404
  - CARMA-Knoten, 398
  - Entscheidungsbaum-Modelle, 186
  - Für Sequenzen, 424
  - Sequenzknoten, 419
- Cox-Regressionsmodelle, 353
  - Einstellungsoptionen, 351
  - Erweiterte Ausgabe, 349, 353
  - Expertenoptionen, 347
  - Feldoptionen, 343
  - Konvergenzkriterien, 348
  - Modell-Nugget, 352
  - Modellierungsknoten, 342
  - Modelloptionen, 344
  - Schrittkriterien, 350
- Cramér-*V*
  - Funktionsauswahl, 79
- Datenauswahl organisieren, 249
- Datenreduktion, 128
  - PCA-/Faktor-Modelle, 311
- Deskriptive Statistiken
  - Verallgemeinerte lineare Modelle, 338
- Diagrammerstellung
  - Assoziationsregeln, 405
- Diagrammoptionen, 263
- Differenz des Konfidenz-Quotienten zur 1
  - A Priori-Evaluierungsmaß, 393
- Dimensionsreduzierung, 355
- Direkte Oblimin-Rotation
  - PCA-/Faktor-Modelle, 314
- Diskriminanzmodelle
  - Bewertung, 325
  - Erweiterte Ausgabe, 322, 326
  - Expertenoptionen, 321
  - Konvergenzkriterien, 321
  - Modell-Nugget, 325–327
  - Modellformat, 320
  - Modellierungsknoten, 319
  - Neigungs-Scores, 326
  - Schrittkriterien (Feldauswahl), 324
- Dokumentation, 3
- Doppelköpfige Regeln, 399
- DTD, 71

- Ebene stabilisierende Transformation, 437
- Ebene verändernde Ausreißer, 436
  - Zeitreihenmodellierung, 452
- edit
  - Erweiterte Parameter, 247
- Eigenwerte
  - PCA-/Faktor-Modelle, 313
- Einflussvariablen
  - Auswahl für die Analyse, 75, 78, 80, 82
  - Entscheidungsbäume, 134
  - Rangordnung der Wichtigkeit, 75, 78–80, 82
  - Screening, 75, 80, 82
  - Surrogate, 134
- Eingabefelder
  - Auswahl für die Analyse, 77
  - Screening, 77
- Einstellungsoptionen
  - Cox-Regressionsmodelle, 351
  - SLRM-Knoten, 467
- Einteilen von kontinuierlichen Variablen, 128
- Ensemble-Viewer, 57
  - Automatische Datenaufbereitung, 65
  - Bedeutsamkeit des Prädiktors, 61
  - Komponentenmodelldetails, 64
  - Komponentenmodellgenauigkeit, 63
  - Modellzusammenfassung, 60
  - Prädiktorhäufigkeit, 62
- Ensembles
  - in linearen Modellen, 273
  - in neuronalen Netzwerken, 216
- Entfernen von Modellverknüpfungen, 46
- Entscheidungsbäume beschneiden, 155, 162
- Entscheidungsbaum-Modelle, 126, 129, 131, 136, 154–158, 177, 181, 187, 191
  - Baumregeln, 183
  - Bedeutsamkeit des Prädiktors, 183
  - Benutzerdefinierte Aufteilungen, 132
  - Diagrammerstellung, 191
  - Einflussvariablen, 134
  - erzeugen, 146–147
  - Exportieren von Ergebnissen, 150
  - Fehlklassifizierungskosten, 166, 168
  - Fenster mit zusätzlichen Informationen, 186
  - Gewinne, 137–138, 141, 144
  - Grenzoptionen, 164
  - Modell-Nugget, 183
  - Modellierungsknoten, 151
  - Profite, 140
  - Regelhäufigkeiten, 186
  - ROI, 140
  - Surrogate, 134, 186
  - Viewer, 187
- Entscheidungslistenmodelle
  - Alternativen, Registerkarte, 240
  - Arbeiten mit Viewer, 244
  - Arbeitsmodellbereich, 238
  - Ausschließen von Segmenten, 227
  - Bewertung, 227, 235
  - Einstellungen, 237
  - Expertenoptionen, 234
  - Klassiermethode, 234
  - Mailinglisten, 226
  - Modellierungsknoten, 226
  - Modelloptionen, 232
  - PMML, 235
  - requirements, 231
  - Segmente, 235
  - Snapshots, Registerkarte, 242
  - SQL-Erzeugung, 237
  - Suchbreite, 234
  - Suchrichtung, 232
  - Viewer-Arbeitsbereich, 238
  - Zielwert, 232
- Epsilon für Konvergenz
  - CHAID-Knoten, 174
- Equamax-Rotation
  - PCA-/Faktor-Modelle, 314
- Ersetzen von Modellen, 49
- Erste Schritte, 238
- Erstellungsauswahl
  - definieren, 245
- Erster Treffer, Regelmenge, 194
- Erweiterte Ausgabe
  - Cox-Regressionsmodelle, 349
  - Faktor/PCA-Knoten, 315
- Erweiterte Parameter, 247
- Evaluationsdiagramme
  - aus Modellen vom Typ “Auto-Numerisch”, 123–124
  - aus Modellen vom Typ “Autom. Cluster”, 123
  - Automatischer Klassifizierer, Modelle, 123–124
- Evaluierungsmaße
  - A Priori-Knoten, 392
- events
  - erkennen, 434
- Exhaustive CHAID, 127, 129, 162
- Expertenausgabe
  - Cox-Regressionsmodelle, 349
- Expertenmodellierung
  - Ausreißer, 445
  - Kriterien für Zeitreihenmodelle, 444
- Expertenoptionen
  - A Priori-Knoten, 392
  - Bayes-Netzwerk-Knoten, 203
  - CARMA-Knoten, 399
  - Cox-Regressionsmodelle, 347
  - K-Means-Modelle, 364
  - Kohonen-Modelle, 359
  - Sequenzknoten, 420
- exponentielles Glätten, 438
  - Kriterien für Zeitreihenmodelle, 446
- exportieren
  - Modell-Nuggets, 50
  - SQL, 53

- Exportieren
  - PMML, 70, 73
- F*-Statistik
  - Funktionsauswahl, 79
  - in linearen Modellen, 271
- Faktor-Modelle
  - Anzahl der Faktoren, 313
  - Eigenwerte, 313
  - Erweiterte Ausgabe, 318
  - Expertenoptionen, 313
  - Faktorwerte, 313
  - Fehlende Werte, Behandlung, 313
  - Gleichungen, 316
  - Iteration, 313
  - Modell-Nugget, 315–316, 318
  - Modellierungsknoten, 311
  - Modelloptionen, 312
  - Rotation, 314
- Fehlende Daten
  - Prädiktorreihen, 438
- Fehlende Werte
  - Ausschließen aus SQL, 189
  - CHAID-Knoten, 133
  - Screening von Feldern, 77
- Fehlerzusammenfassung
  - in der Nächste-Nachbarn-Analyse, 502
- Fehlklassifizierungskosten
  - C5.0-Knoten, 179
  - Entscheidungsbaume, 103, 167, 169
- Feldoptionen
  - Cox-Knoten, 343
  - Modellierungsknoten, 36
  - SLRM-Knoten, 464
- Feldwichtigkeit
  - Filtern von Feldern, 57
  - Modellergebnisse, 41, 54, 57
  - Zuweisen von Rängen zu Feldern, 75, 78–80, 82
- Fenster “Alternative Regeln”, 250
- Fenster mit zusätzlichen Informationen
  - Entscheidungsbaum-Modelle, 186
- Filterknoten
  - Generieren von Entscheidungsbaumen, 151
- Filterregeln, 403, 426
  - Assoziationsregeln, 404
- Fokusdatensätze, 486
- Funktionale Transformation, 437
- Generierte Sequenzregelmenge, 411
- Gewichtete kleinste Quadrate, 39
- Gewichtungsfelder, 39
- Gewinnbasierte Auswahl, 144
- Gewinne
  - Diagramm, 262
  - Entscheidungsbaume, 137–138, 141
  - exportieren, 150
- Gini-Unreinheitsmaß, 170
- Gleitender Durchschnitt
  - ARIMA-Modelle, 448
- Grenzooptionen
  - Entscheidungsbaume, 164
- Gruppen
  - Anomalieerkennung, 86
  - in der Nächste-Nachbarn-Analyse, 499
- Häufigkeiten
  - Entscheidungsbaum-Modelle, 186
- Häufigkeitsfelder, 39
- Haupteffekte
  - Logistische Regressionsmodelle, 294
- Hauptkomponentenanalyse. *Siehe* PCA-Modelle, 311, 315
- Hits
  - Entscheidungsbaum-Gewinne, 137
- Hosmer-Lemeshow-Anpassungsgüte
  - Logistische Regressionsmodelle, 309
- IBM InfoSphere Warehouse (ISW)
  - PMML-Export, 73
- IBM ISW-Assoziationsmodelle
  - Transaktionsdaten, 37
- IBM SPSS Modeler, 1
  - Dokumentation, 3
- IBM SPSS Statistics-Modelle, 30
- IBM SPSS Text Analytics, 2
- ID-Feld
  - CARMA-Knoten, 395
  - Sequenzknoten, 417
- Importieren
  - PMML, 50, 71, 73
- Impulse
  - in Reihen, 434
- Index
  - Entscheidungsbaum-Gewinne, 137
- Informationsdifferenz
  - A Priori-Evaluierungsmaß, 394
- Informationskriterien
  - in linearen Modellen, 271
- Inhaltsfeld(er)
  - CARMA-Knoten, 395
  - Sequenzknoten, 417
- innovatorische Ausreißer, 436
  - Zeitreihenmodellierung, 452
- Instanzen, 402, 426
  - Entscheidungsbaum-Modelle, 186
- Integration
  - ARIMA-Modelle, 448
- Interaktionen
  - Logistische Regressionsmodelle, 294
- Interaktionsidentifizierung, 128
- Interaktive Bäume, 126, 129, 131, 134, 136
  - Benutzerdefinierte Aufteilungen, 132
  - Diagrammerstellung, 191
  - Exportieren von Ergebnissen, 150
  - Generieren von Modellen, 146–147

- Gewinne, 137–138, 141, 144
- Profite, 140
- ROI, 140
- Surrogate, 134
- Interventionen
  - erkennen, 434
- Iterationsprotokoll
  - Logistische Regressionsmodelle, 298
  - Verallgemeinerte lineare Modelle, 338
- K*-Means-Modelle, 354, 362–364
  - Clustering, 362, 366
  - Diagrammerstellung aus dem Modellnugget, 385
  - Distanzfeld, 363
  - Expertenoptionen, 364
  - Grenzkriterien, 364
  - Modell-Nugget, 365–366
  - Verschlüsselungswert für Sets, 364
- Kassenrollendaten, 389, 412–413
- Kategoriezusammenführung, 128
- Kernel-Funktionen
  - Support Vector Machine, Modelle, 473
- Klassifikationstabelle
  - in der Nächste-Nachbarn-Analyse, 501
  - Logistische Regressionsmodelle, 298
- Klassifizierungsbäume, 155–157, 177
- Klassifizierungsgewinne
  - Entscheidungsbaume, 138, 141
- KNN. *Siehe* Nächste-Nachbarn-Modelle, 482
- Knoten “neuralnetwork”, 210
- Kohonen, Modelle, 354
- Kohonen-Modelle, 355, 357, 359
  - Binäre Setverschlüsselungsoption (entfernt), 357
  - Diagrammerstellung aus dem Modellnugget, 385
  - Expertenoptionen, 359
  - Feedback-Diagramm, 357
  - Grenzkriterien, 357
  - Lernrate, 359
  - Modell-Nugget, 361
  - Modellierungsknoten, 355
  - neighborhood, 355, 359
  - Neuronale Netzwerke:, 355, 361
- Kombinieren der Regeln
  - in linearen Modellen, 273
  - in neuronalen Netzwerken, 216
- Konfidenz
  - Assoziationsregeln, 426
- Konfidenz-Differenz
  - A Priori-Evaluierungsmaß, 393
- Konfidenz-Quotient
  - A Priori-Evaluierungsmaß, 393
- Konfidenz-Scores, 44
- Konfidenzen
  - Entscheidungsbaum-Modelle, 188
  - Logistische Regressionsmodelle, 306
  - Regelmengen, 188
- Konfidenzintervalle
  - Logistische Regressionsmodelle, 298
- Kontinuierliche Variablen
  - Segmentierung, 128
- Kontrastkoeffizienten-Matrix
  - Verallgemeinerte lineare Modelle, 338
- Konvergenzoptionen
  - CHAID-Knoten, 174
  - Cox-Regressionsmodelle, 348
  - Logistische Regressionsmodelle, 297
  - Verallgemeinerte lineare Modelle, 336
- Kopieren von Modellverknüpfungen, 48
- Korrelationsmatrix
  - Verallgemeinerte lineare Modelle, 338
- Korrigierte Neigung, Scores
  - Balancieren von Daten, 44
  - Diskriminanzmodelle, 326
  - Entscheidungslistenmodelle, 237
  - Verallgemeinerte lineare Modelle, 340
- Korrigiertes R-Quadrat
  - in linearen Modellen, 271
- Kosten
  - Entscheidungsbaume, 166, 168
- Kovarianzmatrix
  - Verallgemeinerte lineare Modelle, 338
- Kriterium zur Verhinderung übermäßiger Anpassung
  - in linearen Modellen, 271
- L*-Matrix
  - Verallgemeinerte lineare Modelle, 338
- Laden
  - Modell-Nuggets, 50
- Lagrange-Multiplikator-Test
  - Verallgemeinerte lineare Modelle, 338
- Lambda
  - Funktionsauswahl, 79
- Leistung
  - C5.0-Modelle, 178, 181
- Leistungsverbesserungen, 301, 359, 364, 392
- Lift, 403
  - Assoziationsregeln, 404
  - Entscheidungsbaum-Gewinne, 137
- Lift Chart
  - Entscheidungsbaum-Gewinne, 143
- Likelihood-Quotienten-Chi-Quadrat
  - CHAID-Knoten, 174
  - Funktionsauswahl, 79
- Likelihood-Quotiententest
  - Logistische Regressionsmodelle, 298, 309
- Lineare Modelle, 267
  - ANOVA-Tabelle, 281
  - Ausreißer, 280
  - Automatische Datenaufbereitung, 270, 276
  - Bedeutsamkeit des Prädiktors, 277
  - Ensembles, 273
  - Ergebnisse reproduzieren, 274
  - geschätzte Mittel, 285

- Informationskriterium, 275
- Koeffizienten, 283
- Kombinieren der Regeln, 273
- Konfidenzniveau, 270
- Modellauswahl, 271
- Modellerstellungsübersicht, 286
- Modelloptionen, 274
- Modellübersicht, 275
- Nugget-Einstellungen, 287
- R-Quadrat-Statistik, 275
- Residuen, 279
- Vorhersage nach Beobachtung, 278
- Ziele, 268
- Lineare Regression, Modelle, 265
  - Gewichtete kleinste Quadrate, 39
  - Modellierungsknoten, 266
- lineare Trends
  - erkennen, 432
- Linearer Kernel
  - Support Vector Machine, Modelle, 473
- linearnode-Knoten, 267
- Log-Odds
  - Logistische Regressionsmodelle, 303
- Log-Transformation, 437
  - Zeitreihenmodellierung, 450
- Logistische Regressionsmodelle, 265
  - Bedeutsamkeit des Prädiktors, 303
  - Binomiale Modelle, Optionen, 288
  - Erweiterte Ausgabe, 298, 309
  - Expertenoptionen, 296
  - Haupteffekte, 294
  - Hinzufügen von Fachausdrücken, 294
  - Interaktionen, 294
  - Konvergenzoptionen, 297
  - Modell-Nugget, 302–303, 305–306
  - Modellgleichungen, 303
  - Modellierungsknoten, 287
  - Multinomiale Modelle, Optionen, 288
  - Schrittoptionen, 301
- löschen
  - Modellverknüpfungen, 46
- Mailinglisten
  - Entscheidungslistenmodelle, 226
- Manager
  - Registerkarte “Modelle”, 50
- Marken, 506
- Mehrschicht-Perzeptron (MLP)
  - in neuronalen Netzwerken, 214
- Merkmalsauswahlmodelle, 80, 82
  - Generieren von Filterknoten, 82
  - Rangordnung von Prädiktoren, 75–76, 80
  - Screening von Prädiktoren, 75–76, 80
  - Wichtigkeit, 75–76, 80
- Mining-Aufgabe
  - starten, 245
- Mining-Aufgaben, 244
  - bearbeiten, 246
  - Entscheidungslistenmodelle, 226
  - erstellen, 245
- MLP (Mehrschicht-Perzeptron)
  - in neuronalen Netzwerken, 214
- Modell-Nuggets, 45, 74, 181, 188, 190, 194, 196, 341
  - Aufgeteilte Modelle, 66
  - Drucken, 52
  - Ensemble-Modelle, 57
  - Erstellung von Verarbeitungsknoten, 68
  - exportieren, 50, 52
  - Menüs, 52
  - Registerkarte “Übersicht”, 54
  - Scoring von Daten, 68
  - speichern, 52
  - Speichern und Laden, 50
  - Verwendung in Streams, 68
- Modellaktualisierung
  - Selbstlern-Antwortmodelle, 465
- Modellansicht
  - in der Nächste-Nachbarn-Analyse, 495
- Modelle
  - ARIMA (X11 ARIMA), 448
  - aufteilen, 31, 33–35
  - Ersetzen, 49
  - Importieren, 50
  - Registerkarte “Übersicht”, 54
- Modellierungsknoten, 24, 83, 177, 198, 355, 362, 367, 390, 416, 463
- Modellinformationen
  - Verallgemeinerte lineare Modelle, 338
- Modellmaße
  - Aktualisieren, 257
  - definieren, 256
- Modelloptionen
  - Bayes-Netzwerk-Knoten, 200
  - Cox-Regressionsmodelle, 344
  - SLRM-Knoten, 465
- Modellregeln hinzufügen, 250
- Modellverknüpfungen, 46
  - Definieren und Entfernen, 46
  - Kopieren und Einfügen, 48
  - und Superknoten, 48
- MS Excel Setup, Integrationsformat, 259
- Multinomiale logistische Regression, Modelle, 288
- Multinomiale logistische Regression, Modelle, 287
- Nächste-Nachbarn-Modelle
  - Analyseoptionen, 492
  - Einstellungsoptionen, 484
  - Info zu, 482
  - Kreuzvalidierungsoptionen, 490
  - Merkmalsauswahl, Optionen, 489
  - Modellierungsknoten, 482
  - Modelloptionen, 485
  - Nachbarnoptionen, 487

- Zieloptionen, 483
- Neigungs-Scores
  - Balancieren von Daten, 44
  - Diskriminanzmodelle, 326
  - Entscheidungslistenmodelle, 237
  - Verallgemeinerte lineare Modelle, 340
- Neues Modell erzeugen, 255
- Neuronale Netzwerk-Modelle
  - Feldoptionen, 36
- Neuronale Netzwerke:, 210
  - Abbruchregeln, 215
  - Bedeutsamkeit des Prädiktors, 220
  - Ensembles, 216
  - Ergebnisse reproduzieren, 217
  - Fehlende Werte, 217
  - Klassifikation, 222
  - Kombinieren der Regeln, 216
  - Mehrschicht-Perzeptron (MLP), 214
  - Modelloptionen, 218
  - Modellübersicht, 219
  - Netz, 223
  - Nugget-Einstellungen, 224
  - Radiale Basisfunktion (RBF), 214
  - Verdeckte Schichten, 214
  - Verhinderung übermäßiger Anpassung, 217
  - Vorhersage nach Beobachtung, 221
  - Ziele, 212
- Nicht überwachtetes Lernen, 354–355
- Nicht verfeinerte Modelle, 74, 80, 82, 388
- Nicht verfeinerte Regelmodelle, 400–401, 409–410
- nichtlineare Trends
  - erkennen, 432
- nichtseasonale Zyklen, 433
- Nominale Regression, 287
- Normalisiertes Chi-Quadrat
  - A Priori-Evaluierungsmaß, 394
- Nuggets für aufgeteilte Modelle, 66
  - Registerkarte “Übersicht”, 54
  - Viewer, 66
- Optimieren der Leistung, 359, 364, 392
- p*-Wert, 79
- Palette der Modelle, 45, 50
- Parallele Verarbeitung
  - C5.0-Modelle, 178, 181
- Parameter
  - in Zeitreihenmodellen, 458
- Parameterschätzer
  - Logistische Regressionsmodelle, 309
  - Verallgemeinerte lineare Modelle, 338
- Partielle Autokorrelationsfunktion
  - Reihen, 436
- Partitionen, 38, 396, 418, 485
  - auswählen, 38, 396, 418, 485
  - Modellerstellung, 99, 109, 116, 179, 200, 232, 289, 312, 320, 331, 345, 357, 363, 369, 419, 465, 477, 485
- PCA-Modelle
  - Anzahl der Faktoren, 313
  - Eigenwerte, 313
  - Erweiterte Ausgabe, 318
  - Expertenoptionen, 313
  - Faktorwerte, 313
  - Fehlende Werte, Behandlung, 313
  - Gleichungen, 316
  - Iteration, 313
  - Modell-Nugget, 315–316, 318
  - Modellierungsknoten, 311
  - Modelloptionen, 312
  - Rotation, 314
- Periodizität
  - Zeitreihenmodellierung, 450
- PMML
  - Exportieren von Modellen, 50, 70, 73
  - Importieren von Modellen, 50, 71, 73
- Prädiktorauswahl
  - in der Nächste-Nachbarn-Analyse, 501
- Prädiktorbereichsdiagramm
  - in der Nächste-Nachbarn-Analyse, 496
- Prädiktorreihen, 438
  - Fehlende Daten, 438
- Profite
  - Entscheidungsbaum-Gewinne, 140
- Promax-Rotation
  - PCA-/Faktor-Modelle, 314
- Pseudo-*R*-Quadrat
  - Logistische Regressionsmodelle, 309
- Punkt-Interventionen
  - erkennen, 434
- Quadrantenkarte
  - in der Nächste-Nachbarn-Analyse, 500
- Quadratwurzeltransformation, 437
  - Zeitreihenmodellierung, 450
- Quartimax-Rotation
  - PCA-/Faktor-Modelle, 314
- QUEST-Modelle, 127
  - A-priori-Wahrscheinlichkeit, 167
  - Baumtiefe, 162
  - beschneiden, 162
  - Diagrammerstellung aus dem Modellnugget, 191
  - Ensembles, 165
  - Fehlklassifizierungskosten, 166
  - Feldoptionen, 158
  - Grenzooptionen, 164
  - Modell-Nugget, 181
  - Modellierungsknoten, 129, 154, 157, 187–188
  - Surrogate, 163
  - Ziele, 159
- R-Quadrat
  - in linearen Modellen, 275
- Radiale Basisfunktion (RBF)
  - in neuronalen Netzwerken, 214

- Rangordnung von Prädiktoren, 75, 78–80, 82
- RBF (Radiale Basisfunktion)  
in neuronalen Netzwerken, 214
- Rechtliche Hinweise, 504
- Referenzkategorie  
Logistikknoten, 293
- Regel-ID, 403
- Regel-Superknoten  
Generieren aus Sequenzregeln, 428
- Regelbedingungen  
Entscheidungslistenmodelle, 226
- Regelerstellungsknoten, 181
- Regelinduktion, 126, 155–157, 177, 390
- Regelmenge, 151, 188, 194, 196, 407, 410–411  
Generieren von Entscheidungsbäumen, 151
- Regeln  
Assoziationsregeln, 390, 394  
Regelunterstützung, 403, 426
- Registerkarte “Snapshots”, 242
- Regressionsbäume, 155–157
- Regressionsgewinne  
Entscheidungsbaum, 141, 144
- Regressionsmodelle  
Modellierungsknoten, 266
- Reihen  
transformieren, 437
- Residuen  
in Zeitreihenmodellen, 459
- Risiken  
exportieren, 150
- Risikoschätzung  
Entscheidungsbaum-Gewinne, 146
- Rohneigung, Scores, 44
- ROI  
Entscheidungsbaum-Gewinne, 140
- Rotation  
PCA-/Faktor-Modelle, 314
- Saisonabhängigkeit, 433  
erkennen, 433
- saisonal additive Ausreißer, 436  
Zeitreihenmodellierung, 452
- saisonale Ordnungen  
ARIMA-Modelle, 448
- Schichtung, 128
- Schritt-Interventionen  
erkennen, 434
- Schrittoptionen  
Cox-Regressionsmodelle, 350  
Logistische Regressionsmodelle, 301
- Schrittweise Feldauswahl  
Diskriminanzknoten, 324
- Schrittweise vorwärts  
in linearen Modellen, 271
- Score-Statistik, 299, 301
- Scoring von Daten, 68
- Screening von Eingabefeldern, 77
- Screening von Prädiktoren, 75, 80, 82
- Segmente  
Ausschließen, 254  
bearbeiten, 251  
Einfügen, 250  
Entscheidungslistenmodelle, 226  
Kopieren, 252  
löschen, 254  
Löschen von Regelbedingungen, 252  
Prioritäten zuweisen, 254
- Segmentierung, 128
- Segmentregelerstellung, 244
- Selbstlern-Antwortmodelle  
Einstellungen, 467, 470–471  
Einstellungen für Zielfelder, 468, 472  
Feldoptionen, 464  
Modell-Nugget, 469  
Modellaktualisierung, 465  
Modellierungsknoten, 463  
Randomisierung der Ergebnisse, 467, 471  
Variablenwichtigkeit, 469
- Selbstorganisierende Karten, 355
- Sequenz-Browser, 427
- Sequenzerkennung, 388, 416
- Sequenzmodelle  
Datenformate, 417  
Expertenoptionen, 420  
Feldoptionen, 417  
Generieren eines Regel-Superknotens, 428  
ID-Feld, 417  
Inhaltsfeld(er), 417  
Modell-Nugget, 422, 424, 427  
Modell-Nugget-Einstellungen, 427  
Modell-Nugget-Übersicht, 427  
Modellierungsknoten, 416  
Nähere Informationen zum Modell-Nugget, 424  
Optionen, 419  
Sequenz-Browser, 427  
sortieren, 427  
Tabellen- im Vergleich zu Transaktionsdaten, 420  
Vorhersagen, 422  
Zeitfeld, 417
- Signifikanzniveau  
für Aufteilung, 172–173  
für das Zusammenführen, 173
- SLRM. *Siehe* Selbstlern-Antwortmodelle, 463
- Snapshot  
erstellen, 242
- Sortiertes Twoing-Unreinheitsmaß, 170
- SPSS Modeler Server, 1
- SQL  
export, 53  
Logistische Regressionsmodelle, 307  
Regelmengen, 189
- Statistics-Modelle, 30
- Statistik für Anpassungsgüte  
Logistische Regressionsmodelle, 309

- Verallgemeinerte lineare Modelle, 338
- Statistische Modelle, 265
- Sukzedens
  - Mehrere Sukzedenzen, 399
- Superknoten
  - und Modellverknüpfungen, 48
- Support Vector Machine, Modelle
  - Einstellungen, 480
  - Expertenoptionen, 477
  - Feinabstimmung, 475
  - Info zu, 473
  - Kernel-Funktionen, 473
  - Modell-Nugget, 479, 494
  - Modellierungsknoten, 476
  - Modelloptionen, 477
  - Überanpassung, 475
- Surrogate
  - Entscheidungsbaum-Modelle, 186
  - Entscheidungsbäume, 134, 163
- SVM. *Siehe* Support Vector Machine-Modelle, 473
  
- t*-Statistik
  - Funktionsauswahl, 79
- Tabellendaten, 389, 412
  - A Priori-Knoten, 37
  - CARMA-Knoten, 395
  - Sequenzknoten, 417
  - transponieren, 413
- Territorien
  - Diskriminanzknoten, 322
- Transaktionsdaten, 389, 412–413
  - A Priori-Knoten, 37
  - CARMA-Knoten, 395
  - IBM ISW-Assoziationsmodelle, 37
  - MS-Assoziationsregel-Knoten, 37
  - Sequenzknoten, 417
- Transformation der Differenz, 437
  - ARIMA-Modelle, 448
- Transformation der saisonalen Differenz, 437
  - ARIMA-Modelle, 448
- Transformation mit natürlichem Logarithmus, 437
  - Zeitreihenmodellierung, 450
- Transformieren von Reihen, 437
- Transiente Ausreißer
  - Zeitreihenmodellierung, 452
- Transponieren einer Tabellenausgabe, 413
- Tree Builder, 129, 131, 136
  - Benutzerdefinierte Aufteilungen, 132
  - Diagrammerstellung, 191
  - Einflussvariablen, 134
  - Exportieren von Ergebnissen, 150
  - Generieren von Modellen, 146–147
  - Gewinne, 137–138, 141, 144
  - Profite, 140
  - ROI, 140
  - Surrogate, 134
- Trefferdiagramme
  - Entscheidungsbaum-Gewinne, 137, 143
- Trends
  - erkennen, 432
- Twoing-Unreinheitsmaß, 170
- TwoStep-Cluster-Modelle, 354
- TwoStep-Clustermodelle, 368, 370
  - Anzahl der Cluster, 368
  - Clustering, 370
  - Diagrammerstellung aus dem Modellnugget, 385
  - Modell-Nugget, 370
  - Modellierungsknoten, 367
  - Optionen, 368
  - Standardisierung der Felder, 368
  - Umgang mit Ausreißern, 368
- Überanpassung von SVM-Modellen, 475
- Übertragungsfunktionen, 450
  - Nenner-Terme, 450
  - Ordnung der Differenzen, 450
  - saisonale Ordnungen, 450
  - Verzögerung, 450
  - Zähler-Terme, 450
- Unreinheitsmaße
  - C&R-Baum, Knoten, 170
  - Entscheidungsbäume, 169
- Unterstützung
  - A Priori-Knoten, 391
  - Antezedens-Unterstützung, 402, 426
  - Assoziationsregeln, 404
  - CARMA-Knoten, 398–399
  - Für Sequenzen, 424
  - Regelunterstützung, 403, 426
  - Sequenzknoten, 419
- Variablen
  - Screening, 128
- Variablenwichtigkeit
  - Selbstlern-Antwortmodelle, 469
- Varianz stabilisierende Transformation, 437
- Varianzkoeffizient
  - Screening von Feldern, 77
- Varimax-Rotation
  - PCA-/Faktor-Modelle, 314
- Verallgemeinerte lineare Modelle
  - Erweiterte Ausgabe, 337, 340
  - Expertenoptionen, 332
  - Felder, 330
  - Konvergenzoptionen, 336
  - Modell-Nugget, 339, 341
  - Modellformat, 331
  - Modellierungsknoten, 328
  - Neigungs-Scores, 340
- Verbesserung, 161, 179, 190
  - in linearen Modellen, 268
  - in neuronalen Netzwerken, 212
- verfügbare Felder, 248

- 
- Verhinderung übermäßiger Anpassung
    - in neuronalen Netzwerken, 217
  - Verlauf
    - Entscheidungsbaum-Modelle, 186
  - Verstärkung, 161
    - in linearen Modellen, 268
    - in neuronalen Netzwerken, 212
  - Verwendbarkeitsmaß, 403
  - Viewer, Registerkarte
    - Diagrammerstellung, 191
    - Entscheidungsbaum-Modelle, 187
  - Visualisieren eines Modells, 262
  - Visualisierung
    - Clustermodelle, 372
    - Diagrammerstellung, 191, 385, 405
    - Entscheidungsbäume, 187
  - Vorhersagen
    - Prädiktorreihen, 438
    - Übersicht, 431
  - Vorschau
    - Modellinhalt, 54
  - Voting-Regelmenge, 194
  
  - Wahrheitstabellendaten, 389, 412–413
  - Wahrscheinlichkeiten
    - Logistische Regressionsmodelle, 303
  - Wald-Statistik, 299, 301
  - Warenkorbdaten, 389, 412–413
  - Wichtigkeit
    - Filtern von Feldern, 57
    - Prädiktoren in Modellen, 41, 54, 57
    - Rangordnung von Prädiktoren, 75, 78–80, 82
  
  - Zeitfeld
    - CARMA-Knoten, 395
    - Sequenzknoten, 417
  - Zeitreihenmodelle
    - ARIMA-Kriterien, 447
    - ARIMA-Modelle, 439
    - Ausreißer, 445, 452
    - Expert Modeler-Kriterien, 444
    - exponentielles Glätten, 438
    - Kriterien für exponentielles Glätten, 446
    - Modell-Nugget, 454
    - Modellierungsknoten, 438
    - Modellparameter, 458
    - Periodizität, 450
    - requirements, 440
    - Residuen, 459
    - Übertragungsfunktionen, 450
    - Zeitreihentransformation, 450
  - Zielwert ändern, 255
  - Zusammenhänge
    - Modell, 46