

IBM SPSS Modeler 15 应用程序
指南



注意：使用本信息及其支持的产品之前，请阅读注意事项第 页码下的一般信息。

此版本适用于 IBM SPSS Modeler 15 及所有后续发布和修订，除非在新版本中另有说明。

Adobe 产品屏幕截图重印已获得 Adobe Systems Incorporated 的许可。

Microsoft 产品屏幕截图重印已获得 Microsoft Corporation 的许可。

受许可保护材料 - IBM 所有

Copyright IBM Corporation 1994, 2012.

美国政府用户受限权利 - 使用、复制或披露受与 IBM Corp. 签订的 GSA ADP Schedule Contract 的限制。

前言

IBM® SPSS® Modeler 是 IBM Corp. 企业级数据挖掘工作平台。SPSS Modeler通过深入的数据分析帮助组织改进与客户和市民的关系。组织通过借助源自 SPSS Modeler 的洞察力可以留住优质客户，识别交叉销售机遇，吸引新客户，检测欺诈，降低风险，促进政府服务交付。

SPSS Modeler’ 的可视化界面让用户可以应用他们自己的业务专长，这将生成更强有力的预测模型，缩减实现解决方案所需的时间。SPSS Modeler 提供了多种建模技术，例如预测、分类、细分和关联检测算法。模型创建成功后，通过 IBM® SPSS® Modeler Solution Publisher，在广泛的企业内交付给决策者，或通过数据库交付。

关于 IBM Business Analytics

IBM Business Analytics 软件为决策者提供可信赖的完整、一致和准确信息，以帮助其提升业务绩效。这一涵盖 [商务智能](#)、[预测分析](#)、[财务绩效与战略管理](#) 以及 [分析应用程序](#) 的全面组合可提供有关当前业务表现的清晰、立即和切实可行的深入见解，并能够有效预测未来结果。其中整合了丰富的行业解决方案、经过验证的做法与专业服务，以帮助各种规模的组织提升生产效率、自动化决策并取得卓越成果。

作为该软件组合的一部分，IBM SPSS Predictive Analytics 软件能够帮助各类组织有效地预测未来事件，并针对所得到的深入见解提前采取行动，以取得更优秀的业务成果。全球企业、政府和学院客户依赖 IBM SPSS 技术作为吸引、留住和增加客户数量的竞争优势，并降低欺诈和转移风险。通过将 IBM SPSS 软件融入其日常运营中，这些组织将成为“预测型”企业，即能够指引并自动化决策，以实现业务目标和取得可衡量的竞争优势。有关详细信息，或联系我们的代表，请访问 <http://www.ibm.com/spss>。

技术支持

我们提供有技术支持服务以维护客户。客户可就 IBM Corp. 产品使用或某一受支持硬件环境的安装帮助寻求技术支持。要获得技术支持，请访问 IBM Corp. 网站 <http://www.ibm.com/support>。在请求帮助时，请做好准备，以便识别您自己、您的组织以及您的支持协议。

内容

1 关于 IBM SPSS Modeler	1
IBM SPSS Modeler 产品	1
IBM SPSS Modeler	1
IBM SPSS Modeler Server	2
IBM SPSS Modeler Administration Console	2
IBM SPSS Modeler Batch	2
IBM SPSS Modeler Solution Publisher	2
IBM SPSS Modeler Server 适配器用于 IBM SPSS Collaboration and Deployment Services	2
IBM SPSS Modeler 版本	3
IBM SPSS Modeler 文档	3
SPSS Modeler Professional 文档	4
SPSS Modeler Premium 文档	4
应用程序示例	5
Demos 文件夹	6

部分 I: 简介和新手入门

2 IBM SPSS Modeler 概述	8
新手入门	8
启动 IBM SPSS Modeler	8
从命令行启动	9
连接到 IBM SPSS Modeler Server	9
更改 Temp 目录	12
启动多个 IBM SPSS Modeler 会话	13
IBM SPSS Modeler 界面一览	13
IBM SPSS Modeler 流工作区	14
节点选项板	14
IBM SPSS Modeler 管理器	15
IBM SPSS Modeler 工程	16
IBM SPSS Modeler 工具栏	17
自定义工具栏	18
自定义 IBM SPSS Modeler 窗口	19
更改流的图标尺寸	19
在 IBM SPSS Modeler 中使用鼠标	21
使用快捷键	21
打印	22

自动化 IBM SPSS Modeler	22
3 建模简介	24
构建流	25
浏览模型	31
评估模型	35
对记录评分	38
摘要	38
4 标志目标的自动建模	40
对客户响应建模（自动分类器）	40
历史数据	40
构建流	42
生成和比较模型	47
摘要	52
5 连续目标的自动建模	53
属性值（自动数值）	53
训练数据	53
构建流	54
比较模型	58
摘要	60
部分 II：数据准备示例	
6 自动数据准备（ADP）	62
构建流	62
比较模型准确性	68

7	准备分析数据（数据审核）	71
	构建流	71
	浏览统计量和图表	75
	处理离群值和缺失值	78
8	药物治疗（勘察表/C5.0）	83
	读取文本数据	83
	添加表	87
	创建分布图	88
	创建散点图	90
	创建网络图	92
	导出新字段	93
	构建模型	96
	浏览模型	99
	使用分析节点	101
9	筛选预测变量（特征选择）	103
	构建流	103
	构建模型	106
	比较结果	107
	摘要	108
10	减少输入数据字符串长度（重新分类节点）	109
	减少输入数据字符串长度（重新分类）	109
	重新分类数据	109

部分 III：建模示例

11 对客户响应建模（决策列表）	115
历史数据	116
构建流	117
创建模型	120
使用 Excel 计算自定义测量量	133
修改 Excel 模板	138
保存结果	140
12 电信业客户分类（多项 Logistic 回归）	142
构建流	143
浏览模型	147
13 电信客户流失（二项 Logistic 回归）	152
构建流	152
浏览模型	160
14 预测带宽利用率（时间序列）	165
使用时间序列节点进行预测	165
创建流	166
检查数据	167
定义日期	171
定义目标	173
设置时间区间	174
创建模型	176
检查模型	178
摘要	187
重新应用时间序列模型	187
检索流	188
检索保存的模型	190
生成建模节点	191

生成新模型	191
检查新模型	193
摘要	195
15 预测产品分类销售情况（时间序列）	196
创建流	196
检查数据	200
指数平稳	200
ARIMA	205
摘要	211
16 向客户报价（自学）	212
构建流	213
浏览模型	218
17 预测贷款拖欠者（贝叶斯网络）	223
构建流	223
浏览模型	227
18 每个月重新训练模型（贝叶斯网络）	232
构建流	232
评估模型	236
19 零售促销（神经网络/C&RT）	243
检查数据	243
学习和检验	246

20 状态监测（神经网络/C5.0）	248
检查数据	249
数据准备	251
学习	252
测试	252
21 电信客户分类（判别式分析）	254
创建流	254
检查模型	260
逐步判别式分析	262
有关逐步法的警告说明	263
检查模型拟合	263
结构矩阵	264
区域图	265
分类结果	266
摘要	266
22 分析区间型删失的生存数据（广义线性模型）	267
创建流	267
模型效应检验	272
拟合仅治疗模型	272
参数估计值	273
预测复发和生存的概率	274
按周期对复发概率进行建模	278
模型效应检验	284
拟合简化模型	284
参数估计值	285
预测复发和生存的概率	286
摘要	290
23 使用泊松回归来分析船只损坏率（广义线性模型）	292
拟合“高度离散的”泊松回归	292

拟合优度统计	297
Omnibus 检验	297
模型效应检验	298
参数估计值	298
拟合其他模型	299
拟合优度统计	302
摘要	303
24 将 Gamma 回归拟合至汽车保险理赔（广义线性模型）	304
创建流	304
参数估计值	308
摘要	308
25 细胞样本分类（SVM）	309
创建流	310
检查数据	315
尝试另一种函数	317
比较结果	319
摘要	320
26 将 Cox 回归用于客户流失时间模型	321
构建合适的模型	321
删失的观测值	325
分类变量编码	326
变量选择	327
协变量均值	329
生存曲线	330
风险曲线	331
评估	332
跟踪仍在的预期客户数	337
评分	351
摘要	356

27 市场购物篮分析（规则归纳/C5.0） 357

访问数据 357
发现购物篮内容的关系 358
描绘客户群的特征 362
摘要 363

28 评估新车辆产品（KNN） 364

创建流 365
检查输出 369
 预测变量空间 370
 对等图表 371
 邻元素和距离表 373
摘要 374

附录

A 注意事项 375

参考书目 377

索引 378

关于 IBM SPSS Modeler

IBM® SPSS® Modeler 是一组数据挖掘工具，通过这些工具可以采用商业技术快速建立预测性模型，并将其应用于商业活动，从而改进决策过程。SPSS Modeler 参照行业标准 CRISP-DM 模型设计而成，可支持从数据到更优商业成果整个数据挖掘过程。

SPSS Modeler 提供了各种借助机器学习、人工智能和统计学的建模方法。通过建模选项板中的方法，您可以根据数据生成新的信息以及开发预测模型。每种方法各有所长，同时适用于解决特定类型的问题。

SPSS Modeler 可以作为独立产品购买，或作为客户端与 SPSS Modeler Server 一起使用。同时提供了大量其他选项，以下各节将对这些选项进行概述。有关详细信息，请参阅 <http://www.ibm.com/software/analytics/spss/products/modeler/>。

IBM SPSS Modeler 产品

IBM® SPSS® Modeler 系列产品及其相关软件包括如下成员：

- IBM SPSS Modeler
- IBM SPSS Modeler Server
- IBM SPSS Modeler Administration Console
- IBM SPSS Modeler Batch
- IBM SPSS Modeler Solution Publisher
- IBM SPSS Modeler Server 适配器用于 IBM SPSS Collaboration and Deployment Services

IBM SPSS Modeler

SPSS Modeler 是在个人电脑上安装并运行的完整功能版本产品。您可以在本地模式下将 SPSS Modeler 作为独立产品来运行；也可以在分布式模式下与 IBM® SPSS® Modeler Server 协同使用，从而提高了对大数据集的处理速度。

使用 SPSS Modeler，您可以快速直观地建构精确的预测模型，无需编程。使用其独特的可视界面，您可以轻松地将数据挖掘过程可视化。通过内嵌在该产品中的高级分析支持，您可以发现之前在您的数据中隐藏的模式与趋势。您可以建构输出模型，并理解其影响因子，让您可以更好地利用业务机会、降低风险。

SPSS Modeler 提供两个版本：SPSS Modeler Professional 和 SPSS Modeler Premium。有关详细信息，请参阅 [IBM SPSS Modeler 版本中的 IBM SPSS Modeler 15 用户指南](#)。

IBM SPSS Modeler Server

SPSS Modeler 使用客户端/服务器体系结构将资源集约型操作的请求分发给功能强大的服务器软件，因而使大数据集的传输速度大大加快。

SPSS Modeler Server 是一项独立许可的产品，可以协同一份或多份 IBM® SPSS® Modeler 安装，以分布式分析模式在服务器主机上持续运行。因此，SPSS Modeler Server 在处理大型数据集时具有卓越的性能，因为内存密集型的操作可在服务器完成而无需将数据下载到客户端计算机。IBM® SPSS® Modeler Server 还提供了对 SQL 优化的支持与数据库内建模能力，有助于进一步提高性能与自动化程度。

IBM SPSS Modeler Administration Console

Modeler Administration Console 是一个图形化的应用程序，可管理 SPSS Modeler Server 的多项配置选项，配置也可通过编辑一个选项文件来进行。该应用程序提供了一个控制台用户界面，用以监控和配置所安装的 SPSS Modeler Server，SPSS Modeler Server 的现用户可以免费使用该应用程序。应用程序只能安装在 Windows 计算机上；但是它可以管理安装在任何受支持平台上的服务器。

IBM SPSS Modeler Batch

虽然数据挖掘通常是交互式的过程，但是也可以通过命令行来运行 SPSS Modeler 而无需图形用户界面。例如，您可能有些任务需长期或重复性地运行而无用户干预。SPSS Modeler Batch 是该产品一个特殊版本，无需通过常规的用户界面即可完整地实现 SPSS Modeler 的分析功能。使用 SPSS Modeler Batch 需要具备 SPSS Modeler Server 许可。

IBM SPSS Modeler Solution Publisher

SPSS Modeler Solution Publisher 是一个工具，让您可以创建 SPSS Modeler 流的打包版本，由外部运行时引擎运行，或嵌入一个外部应用程序。以此方式，您可以发布与部署完整的 SPSS Modeler 流，在未安装 SPSS Modeler 的环境下也能使用。SPSS Modeler Solution Publisher 是作为 IBM SPSS Collaboration and Deployment Services - Scoring 服务的一部分来发行的，需另行购买许可。获得许可后，您会收到 SPSS Modeler Solution Publisher Runtime，让您可以执行已发布的流。

IBM SPSS Modeler Server 适配器用于 IBM SPSS Collaboration and Deployment Services

有若干种 IBM® SPSS® Collaboration and Deployment Services 适配器可以让您通过 SPSS Modeler 和 SPSS Modeler Server 与 IBM SPSS Collaboration and Deployment Services 存储库交互。以此方式，部署在存储库中的 SPSS Modeler 流即可实现多用户共享，或通过瘦客户端应用程序 IBM SPSS Modeler Advantage 访问。适配器须安装在存储库的主机系统中。

IBM SPSS Modeler 版本

SPSS Modeler 提供下列版本。

SPSS Modeler Professional

SPSS Modeler Professional 提供您在处理大多数类型的结构化数据（例如 CRM 系统中跟踪的行为或交互活动、人口统计、购买行为与销售数据）时所需的所有工具。

SPSS Modeler Premium

SPSS Modeler Premium 是一项单独许可的产品，它扩展了 SPSS Modeler Professional 的功能，使其可以处理如实体分析或社会网络专门数据，以及非结构化的文本数据。SPSS Modeler Premium 包含下列组件。

IBM® SPSS® Modeler Entity Analytics 在 IBM® SPSS® Modeler 预测分析的基础上添加了全新的维度。预测分析会尝试根据过去数据预测未来行为，而实体分析侧重于通过解析记录自身的身份冲突，提高当前数据的连贯性和一致性。身份可以指个人、组织、对象或可能不确定的任何其他实体的身份。身份解析在许多领域中都非常重要，包括客户关系管理、检测、反洗钱以及国家与国际安全。

IBM SPSS Modeler Social Network Analysis 将关于关系的信息转换为字段，这些字段可描述个人和组社交行为的特征。使用介绍社交网络之下关系的数据，IBM® SPSS® Modeler Social Network Analysis 可识别影响网络中他人行为的社交领导。此外，可确定受其他网络参与者影响最大的人员。通过结合这些结果和其他测量，您可创建个人的综合配置文件，作为预测模型的基础。包括此社交信息的模型比不包括的模型执行效果更好。

Text Analytics for IBM® SPSS® Modeler 采用了先进语言技术和 Natural Language Processing (NLP)，以快速处理大量无结构文本数据，抽取和组织关键概念，以及将这些概念分为各种类别。抽取的概念和类别可以和现有结构化数据中进行组合（例如人口统计学），并且可用于借助 SPSS Modeler 的一整套数据挖掘工具来进行建模，以此实现更好更集中的决策。

IBM SPSS Modeler 文档

可以从 SPSS Modeler 的帮助菜单中获取在线帮助格式的文档。此文档包括 SPSS Modeler、SPSS Modeler Server 和 SPSS Modeler Solution Publisher 的文档以及《应用程序指南》和其他支持材料。

每个产品的完整文档（包括安装说明）也在每个产品 DVD 的 \Documentation 文件夹下以 PDF 格式提供。安装文档也可从以下网页中下载：
<http://www-01.ibm.com/support/docview.wss?uid=swg27023172>。

两种格式的文档均可从 SPSS Modeler 信息中心获取，其网址如下：
<http://publib.boulder.ibm.com/infocenter/spssmodl/v15r0m0/>。

SPSS Modeler Professional 文档

SPSS Modeler Professional 的文档套件（不含安装说明）如下：

- **IBM SPSS Modeler 用户指南。** 使用 SPSS Modeler 的一般使用介绍，包括如何构建数据流、处理缺失值、生成 CLEM 表达式、处理项目和报告以及将用于部署的流打包为 IBM SPSS Collaboration and Deployment Services、Predictive Application 或 IBM SPSS Modeler Advantage。
- **IBM SPSS Modeler 源、处理和输出节点。** 介绍用于以不同的格式读取、处理和输出数据的所有节点。实际上这表示所有节点而非建模节点。
- **IBM SPSS Modeler 建模节点。** 有关用于创建数据挖掘模型的所有节点的描述。IBM® SPSS® Modeler 可提供各种借助机器学习、人工智能和统计学的建模方法。 [有关详细信息，请参阅第 3 章中的建模节点概述中的 IBM SPSS Modeler 15 建模节点。](#)
- **IBM SPSS Modeler 算法指南。** 介绍 SPSS Modeler 中所用建模方法的数学基础。此指南仅提供 PDF 版。
- **IBM SPSS Modeler 应用程序指南。** 本指南中的示例旨在为具体的建模方法和技术提供具有针对性的简介。还可以在“帮助”菜单中查阅本指南的在线版本。 [有关详细信息，请参阅应用程序示例中的 IBM SPSS Modeler 15 用户指南。](#)
- **IBM SPSS Modeler 脚本编写与自动化。** 通过编写脚本实现系统自动化的相关信息，包括用于操作节点和流的属性信息。
- **IBM SPSS Modeler 部署指南。** 有关在 IBM® SPSS® Collaboration and Deployment Services Deployment Manager 中以处理作业的步骤形式运行 SPSS Modeler 流和方案的信息。
- **IBM SPSS Modeler CLEF 开发人员指南** CLEF 提供了将第三方程序（例如，数据处理例程或建模算法）作为节点集成到 SPSS Modeler 的功能。
- **IBM SPSS Modeler 数据库内数据挖掘指南。** 有关如何利用数据库的功能通过第三方法来改进性能并增强分析功能的信息。
- **IBM SPSS Modeler Server 管理和性能指南。** 有关如何配置和管理 IBM® SPSS® Modeler Server 的信息。
- **IBM SPSS Modeler Administration Console 用户指南。** 有关安装和使用控制台用户界面以监视和配置 SPSS Modeler Server 的信息。控制台实现为 Deployment Manager 应用程序的插件。
- **IBM SPSS Modeler Solution Publisher 指南。** SPSS Modeler Solution Publisher 是一个附加式组件，通过它组织可发布在标准 SPSS Modeler 环境之外使用的流。
- **IBM SPSS Modeler CRISP-DM 指南。** 借助 CRISP-DM 方法进行 SPSS Modeler 数据挖掘的分步指南。
- **IBM SPSS Modeler Batch 用户指南。** 在批处理模式下使用 IBM SPSS Modeler 的完整指南，包括批处理模式的执行与命令行参数的详细信息。此指南仅提供 PDF 版。

SPSS Modeler Premium 文档

SPSS Modeler Premium 的文档套件（不含安装说明）如下：

- **IBM SPSS Modeler Entity Analytics 用户指南。** 关于通过 SPSS Modeler 使用实体分析的信息，涵盖存储库的安装与配置、实体分析节点以及管理任务。

- **IBM SPSS Modeler Social Network Analysis 用户指南。** 通过 SPSS Modeler 进行社会网络分析的指南，包括群组分析与传播分析。
- **Text Analytics for SPSS Modeler 用户指南。** 关于通过 SPSS Modeler 使用文本分析的信息，涵盖文本发掘节点、交互式工作台、模板及其他资源。
- **Text Analytics for IBM SPSS Modeler Administration Console 用户指南。** 有关安装和使用控制台用户界面监视和配置 IBM® SPSS® Modeler Server 用于 Text Analytics for SPSS Modeler 的信息。控制台实现为 Deployment Manager 应用程序的插件。

应用程序示例

SPSS Modeler 中的数据挖掘工具可以帮助解决很多业务和组织问题，应用程序示例将提供有关特定建模方法和技术的简明的针对性说明。此处使用的数据集比某些数据挖掘器管理的大量数据存储要小得多，但涉及的概念和方法应可扩展到实际的应用程序。

可以通过在 SPSS Modeler 中的“帮助”菜单中单击[应用程序示例](#)来访问示例。数据文件和样本流安装在产品安装目录下的 Demos 文件夹中。[有关详细信息，请参阅Demos 文件夹中的IBM SPSS Modeler 15 用户指南。](#)

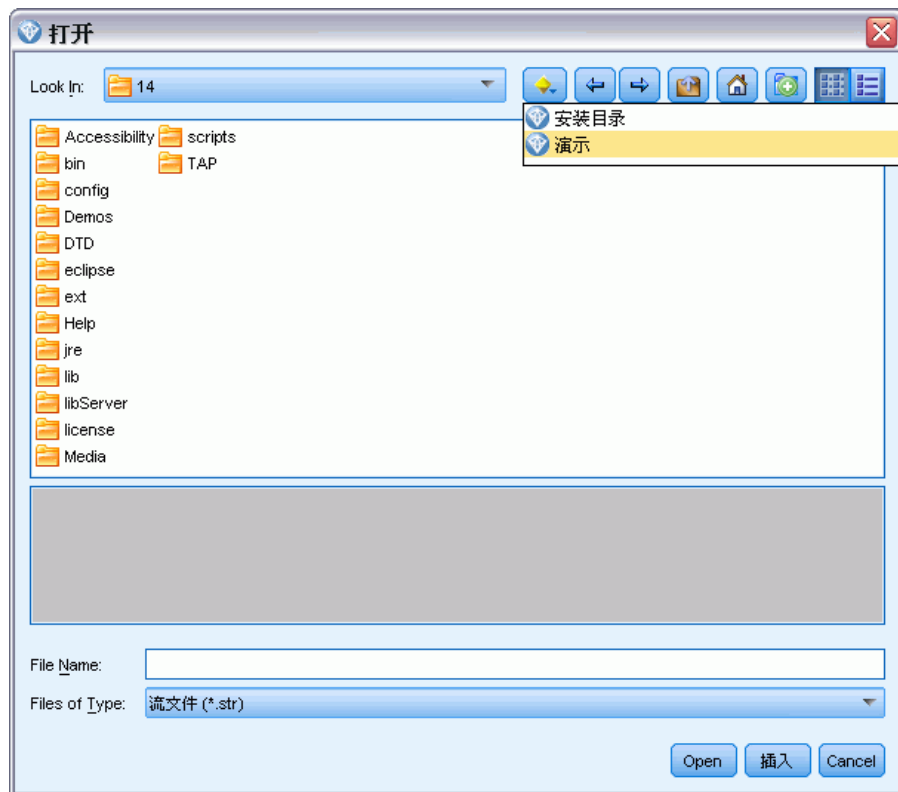
数据库建模示例。 请参阅 IBM SPSS Modeler 数据库内挖掘指南 中的示例。

编写示例脚本。 请参阅 IBM SPSS Modeler 脚本编写和自动化指南 中的示例。

Demos 文件夹

与应用程序示例一起使用的数据文件和样本流安装在产品安装目录下的 Demos 文件夹中。可从 Windows 的“开始”菜单中 IBM SPSS Modeler 15 程序组访问该文件夹，也可以在“文件打开”对话框中最近目录的列表中单击 Demos。

图片 1-1
在最近使用的目录列表中选择 Demos 文件夹



部分 I: 简介和新手入门

IBM SPSS Modeler 概述

新手入门

作为一种数据挖掘应用程序，IBM® SPSS® Modeler 提供了用以寻找大数据集中有用关系的策略性方法。与较为传统的统计方法有所不同，您在开始时不必知道您在寻找的是什么。您可以通过拟合不同的模型和研究不同的关系来探索您的数据，直至发现有用的信息。

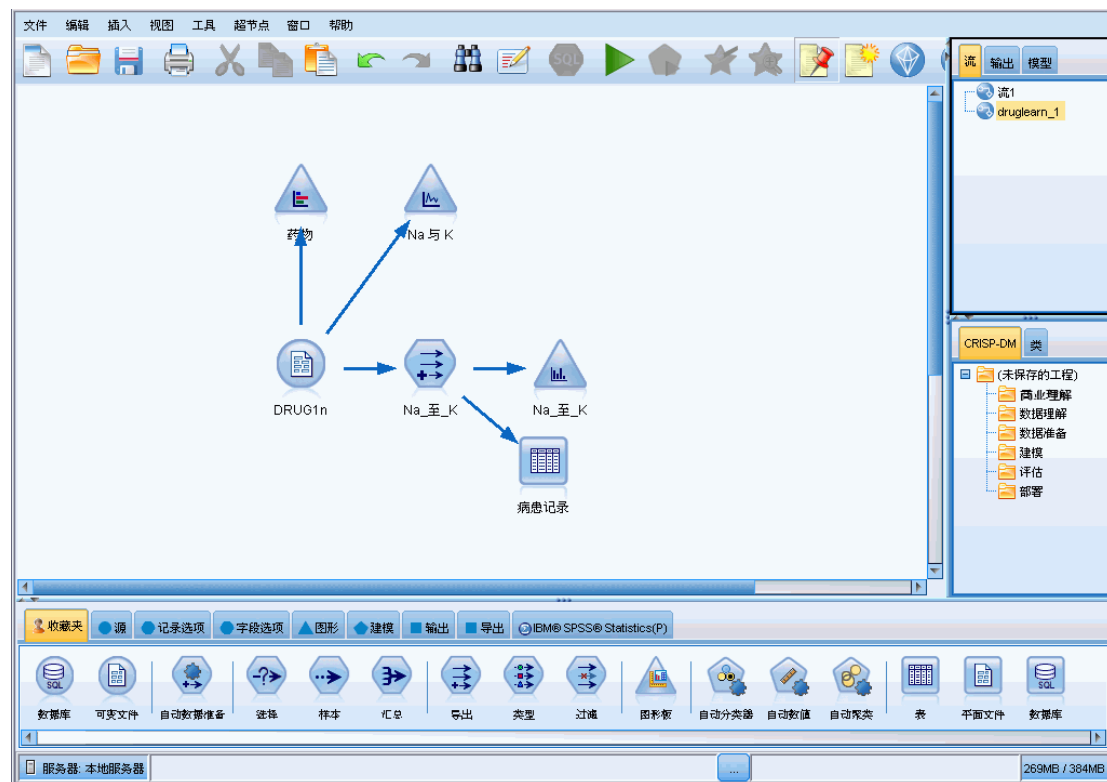
启动 IBM SPSS Modeler

要启动应用程序，单击：

开始 > [所有] 程序 > IBM SPSS Modeler15 > IBM SPSS Modeler15

几秒钟后，屏幕上将显示主窗口。

图片 2-1
IBM SPSS Modeler 主应用程序窗口



从命令行启动

您可以使用操作系统的命令行来如下启动 IBM® SPSS® Modeler:

- ▶ 在安装了 IBM® SPSS® Modeler 的计算机上，打开 DOS 或命令提示符窗口。
- ▶ 要采用互动模式启动 SPSS Modeler 界面，请输入后接所需参数的 `modelerclient` 命令；例如：

```
modelerclient -stream report.str -execute
```

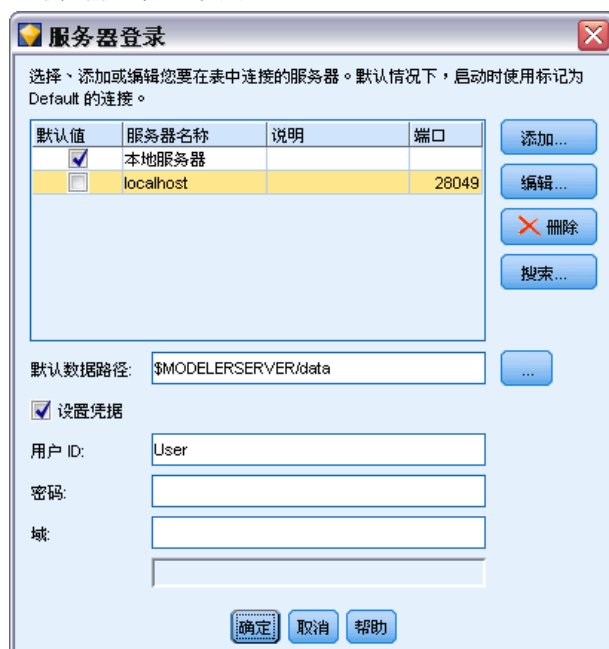
可用参数（标记）允许您连接到一个服务器、加载流、运行脚本或根据需要指定其他参数。

连接到 IBM SPSS Modeler Server

IBM® SPSS® Modeler 可作为独立的应用程序运行，或作为直接连接到 IBM® SPSS® Modeler Server 的客户端运行，或者作为通过进程协调器（COP）插件从 IBM® SPSS® Collaboration and Deployment Services 连接到 SPSS Modeler Server 或服务器群集的客户端运行。当前连接状态显示在 SPSS Modeler 窗口的左下角。

无论何时想连接到服务器，都请手动输入想要连接的服务器的名称或选择之前已定义的名称。但是，如果您拥有 IBM SPSS Collaboration and Deployment Services，则可以从“服务器登录”对话框搜索服务器列表或服务器群集列表。可以通过进程协调器执行浏览网络上运行的 Statistics 服务的功能。[有关详细信息，请参阅附录 D 中的使用服务器群集负载均衡中的 IBM SPSS Modeler Server 15 管理和性能指南。](#)

图片 2-2
“服务器登录”对话框



连接到服务器

- ▶ 在“工具”菜单上，单击**服务器登录**。将打开“服务器登录”对话框。或者，双击 SPSS Modeler 窗口的连接状态区域。
- ▶ 使用该对话框指定要连接到本地服务器计算机的选项或从表中选择连接。
 - 单击**添加或编辑**以添加或编辑连接。有关详细信息，请参阅[添加并编辑 IBM SPSS Modeler Server 连接中的 IBM SPSS Modeler 15 用户指南](#)。
 - 单击**搜索**以访问进程协调器中的服务器或服务器群集。有关详细信息，请参阅[搜索 IBM SPSS Collaboration and Deployment Services 中的服务器中的 IBM SPSS Modeler 15 用户指南](#)。

服务器表。该表包含已定义的服务器连接集。该表显示默认连接、服务器名称、说明和端口号。可以手动添加新的连接，以及选择或搜索现有连接。要将特定的服务器设置为默认连接，请在表中“默认”列中为此连接选择复选框。

默认数据路径。指定用于服务器计算机上的数据的路径。单击省略号按钮 (...)，以浏览至所需要的位置。

设置凭证。不选中此复选框可启用**单点登录**功能，该功能尝试使您使用本地计算机用户名和密码详细信息登录服务器。如果无法使用单点登录，或您选中此复选框以禁用单点登录（例如，登录管理员帐户），则启用以下字段让您输入您的凭证。

用户 ID。输入用于登录到服务器的用户名。

密码。输入与指定用户名关联的密码。

域。指定用于登录到服务器的域。只有服务器计算机与客户计算机处于不同的 Windows 域时，才需要域名。

- ▶ 单击**确定**以完成此连接。

断开与服务器的连接

- ▶ 在“工具”菜单上，单击**服务器登录**。将打开“服务器登录”对话框。或者，双击 SPSS Modeler 窗口的连接状态区域。
- ▶ 在此对话框中，选择“本地服务器”，然后单击**确定**。

添加并编辑 IBM SPSS Modeler Server 连接

可以在“服务器登录”对话框中手动编辑或添加服务器连接。单击“添加”可以访问空的“添加/编辑服务器”对话框，在此对话框中可以输入服务器连接的详细信息。在“服务器登录”对话框中选择现有连接并单击“编辑”，将打开“添加/编辑服务器”对话框，其中包含所选连接的详细信息，以便可以进行任何更改。

注意：不能编辑从 IBM® SPSS® Collaboration and Deployment Services 中添加的服务器连接，因为名称、端口及其他详细信息已在 IBM SPSS Collaboration and Deployment Services 中做过定义。

图片 2-3
“服务器登录：添加/编辑服务器”对话框



添加服务器连接

- ▶ 在“工具”菜单上，单击服务器登录。将打开“服务器登录”对话框。
- ▶ 在此对话框中，单击添加。将打开“服务器登录：添加/编辑服务器”对话框。
- ▶ 输入服务器连接的详细信息，然后单击确定保存此连接并返回“服务器登录”对话框。
 - **服务器。** 指定可用服务器或从列表选择一个服务器。服务器计算机的名称可以使用字母数字（例如 myserver）或指派给服务器计算机的 IP 地址（例如，202.123.456.78）。
 - **端口。** 指定服务器正在侦听的端口号。如果默认设置不可用，请向系统管理员索取正确的端口号。
 - **说明。** 输入此服务器连接的说明（可选）。
 - **确保安全连接（使用 SSL）。** 指定是否应使用 SSL（安全套接层）连接。SSL 是常用于确保网络发送数据的安全的协议。要使用此功能，必须在承载 IBM® SPSS® Modeler Server 的服务器中启用 SSL。必要时请联系本地管理员，以了解详细信息。

编辑服务器连接

- ▶ 在“工具”菜单上，单击服务器登录。将打开“服务器登录”对话框。
- ▶ 在此对话框中，选择希望编辑的连接，然后单击编辑。将打开“服务器登录：添加/编辑服务器”对话框。
- ▶ 更改服务器连接详细信息，然后单击确认保存更改内容并返回至“服务器登录”对话框。

搜索 IBM SPSS Collaboration and Deployment Services 中的服务器

在 IBM® SPSS® Collaboration and Deployment Services 中，可以使用进程协调器选择网络上可用的服务器或服务器群集，从而代替手动输入服务器连接。服务器群集是一组服务器，进程协调器从这组服务器中确定最适合对处理要求作出响应的服务器。有关详细信息，请参阅附录 D 中的使用服务器群集负载均衡中的 IBM SPSS Modeler Server 15 管理和性能指南。

尽管可在“服务器登录”对话框中手动添加服务器，但通过搜索可用的服务器，可在无需知道正确服务器名称和端口号的情况下连接到服务器。此信息是自动提供的。但仍需输入正确的登录信息，如用户名、域和密码。

注意：如果您没有访问进程协调器功能的权限，仍然可以手动输入希望连接的服务器名称或选择之前已定义的服务器名称。有关详细信息，请参阅[添加并编辑 IBM SPSS Modeler Server 连接中的 IBM SPSS Modeler 15 用户指南](#)。

图片 2-4
“搜索服务器”对话框



搜索服务器和服务器群集

- ▶ 在“工具”菜单上，单击服务器登录。将打开“服务器登录”对话框。
- ▶ 在此对话框中，单击搜索打开“搜索服务器”对话框。如果在尝试浏览进程协调器时未登录到 IBM SPSS Collaboration and Deployment Services，则系统会提示您执行此项操作。有关详细信息，请参阅第 9 章中的[连接至存储库中的 IBM SPSS Modeler 15 用户指南](#)。
- ▶ 从列表中选择服务器或服务器群集。
- ▶ 单击确定以关闭对话框，然后将此连接添加到“服务器登录”对话框的表中。

更改 Temp 目录

IBM® SPSS® Modeler Server 执行的某些操作可能需要创建临时文件。默认情况下，IBM® SPSS® Modeler 在系统临时目录下创建临时文件。可通过以下步骤更改临时目录的位置。

- ▶ 创建新目录 `spss` 及其子目录 `servertemp`。
- ▶ 编辑 `options.cfg`，该文件位于 SPSS Modeler 安装目录的 `/config` 目录下。将此文件中的参数 `temp_directory` 编辑为：`temp_directory, "C:/spss/servertemp"`。
- ▶ 完成此操作后，必须重新启动 SPSS Modeler Server 服务。可通过单击 Windows 控制面板中的服务选项卡进行此服务重启操作。只需停止该服务然后重启它即可激活所做的更改。重启计算机也可重新启动此服务。

所有临时文件此时将写入该新目录。

注意：上述操作中最常见的错误是斜杠的使用不正确。由于 SPSS Modeler 以 UNIX 为基础，因此采用正斜杠。

启动多个 IBM SPSS Modeler 会话

如果需要同时启动一个以上的 IBM® SPSS® Modeler 会话，则必须对 IBM® SPSS® Modeler 和 Windows 的设置做一些更改。例如，如果您有两个独立的服务器许可证，并且希望从同一台客户机针对两台不同的服务器运行两个流，则需要对上述设置做一些更改。

要启用多个 SPSS Modeler 会话：

- ▶ 单击：
开始 > [所有] 程序 > IBM SPSS Modeler15
- ▶ 在 IBM SPSS Modeler15 快捷键（带箭头的图标）上右键单击并选择属性。
- ▶ 在目标文本框中，将 `-noshare` 添加到该字符串的结尾。
- ▶ 在 Windows 资源管理器中选择：
工具 > 文件夹选项...
- ▶ 在“文件类型”选项卡上选择“SPSS Modeler 流”选项，然后单击高级。
- ▶ 在“编辑文件类型”对话框中，选择用 SPSS Modeler 打开，然后单击编辑。
- ▶ 在用于执行操作的应用程序文本框中，在 `-stream` 参数前添加 `-noshare`。

IBM SPSS Modeler 界面一览

在数据挖掘过程中的每一个阶段，均可通过 IBM® SPSS® Modeler 易于使用的界面来邀请特定业务的专家。建模算法（如预测、分类、细分和关联检测）可确保得到强大而准确的模型。模型结果可以方便地部署和读入到数据库、IBM® SPSS® Statistics 和各种其他应用程序中。

使用 SPSS Modeler 即处理数据的三个步骤。

- 首先，将数据读入 SPSS Modeler。
- 然后，通过一系列操纵运行数据。
- 最后，将数据发送到目标位置。

这一操作序列称为**数据流**，因为数据以一条条记录的形式，从数据源开始，依次经过各种操纵，最终到达目标（模型或某种数据输出）。

图片 2-5
简单流



IBM SPSS Modeler 流工作区

流工作区是 IBM® SPSS® Modeler 窗口的最大区域，也是您构建和操纵数据流的位置。

通过在界面的主工作区中绘制与业务相关的数据操作图表来创建流。每个操作都用一个图标或节点表示，这些节点通过流链接在一起，流表示数据在各个操作之间的流动。

在 SPSS Modeler 中，可以在同一流工作区或通过打开新的流工作区来一次处理多个流。会话期间，流存储在 SPSS Modeler 窗口右上角的“流”管理器中。

节点选项板

IBM® SPSS® Modeler 中的大部分数据和建模工具位于节点选项板中，该选项板位于流工作区下方窗口的底部。

例如，可以使用“记录选项”选项板选项卡中包含的节点对数据记录执行操作，如选择、合并和追加等。

要将节点添加到工作区，请在节点选项板中双击图标或将其拖放到工作区。随后可将各个图标连接以创建一个表示数据流动的流。

图片 2-6
节点选项板中的“记录选项”选项卡



每个选项板选项卡均包含一组不同的流操作阶段中使用的相关节点，如：

- **源**。此类节点可将数据引入 SPSS Modeler。
- **记录选项**。此类节点可对数据记录执行操作，如选择、合并和追加等。
- **字段选项**。此类节点可对数据字段执行操作，如过滤、导出新字段和确定给定字段的测量级别等。
- **图形**。此类节点可在建模前后以图表形式显示数据。图形包括散点图、直方图、网络节点和评估图表。
- **建模**。此类节点可使用 SPSS Modeler 中提供的建模算法，如神经网络、决策树、聚类算法和数据排序等。
- **数据库建模**。节点使用 Microsoft SQL Server、IBM DB2 和 Oracle 数据库中可用的建模算法。
- **输出**。节点生成可在 SPSS Modeler 中查看的数据、图表和模型等多种输出结果。
- **导出**。节点生成可在外部应用程序（如 IBM® SPSS® Data Collection 或 Excel）中查看的多种输出。
- **SPSS Statistics**。节点将数据导入 IBM® SPSS® Statistics 或从中导出数据，以及运行 SPSS Statistics 过程。

随着对 SPSS Modeler 的熟悉，您也可以自定义供自己使用的选项板内容。[有关详细信息，请参阅第 12 章中的自定义节点选项板中的 IBM SPSS Modeler 15 用户指南。](#)

在“节点选项板”的下方，报告窗格提供各种操作的进度反馈，如将数据读入数据流时的进度。同样在“节点选项板”的下方，状态窗格提供有关应用程序当前正在执行操作的信息以及需要用户反馈时的指示信息。

IBM SPSS Modeler 管理器

窗口右侧顶部为管理器窗格。其中包含三个选项卡，用于管理流、输出和模型。

可以使用“流”选项卡打开、重命名、保存和删除在会话中创建的流。

图片 2-7
“流”选项卡



“输出”选项卡中包含由 IBM® SPSS® Modeler 中的流操作生成的各类文件，如图形和表格。您可以显示、保存、重命名和关闭此选项上列出的表格、图形和报告。

图片 2-8
“输出”选项卡



“模型”选项卡是管理器选项卡中功能最强大的选项卡。该选项卡中包含所有模型块，这些模型块包含针对当前会话在 SPSS Modeler 中生成的模型。这些模型可以直接从“模型”选项卡上浏览或将其添加到工作区的流中。

图片 2-9
包含模型块的“模型”选项卡

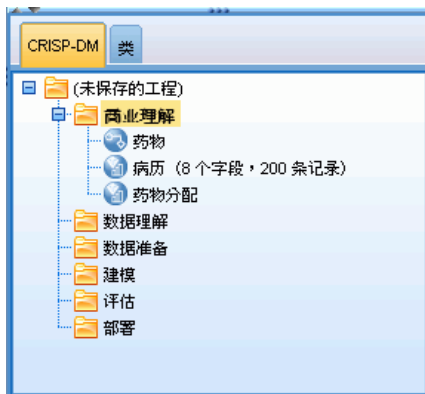


IBM SPSS Modeler 工程

窗口右侧底部是工程窗格，用于创建和管理数据挖掘工程（与数据挖掘任务相关的文件组）。有两种方式可查看您在 IBM® SPSS® Modeler 中创建的工程 - 类视图或 CRISP-DM 视图。

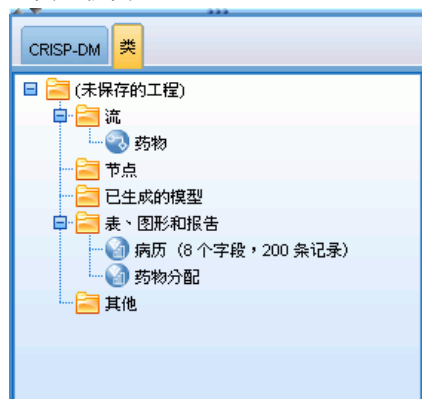
依据“跨行业数据挖掘过程标准”，这一业内认可且无所属的方法理论，“CRISP-DM”选项卡提供了一种组织工程的方式。不论是有经验的数据挖掘人员还是新手，使用 CRISP-DM 工具都会使您事半功倍。

图片 2-10
CRISP-DM 视图



“类”选项卡提供了一种在 SPSS Modeler 中按类别（即，按照所创建对象的类别）组织您工作的方式。此视图在获取数据、流、模型的详尽目录时十分有用。










图片 2-11
“类”视图



IBM SPSS Modeler 工具栏

IBM® SPSS® Modeler 窗口顶部有一个图标工具栏，其中包含许多有用功能。下表列出了这些工具栏按钮及其功能。

	创建新流		打开流
	保存流		打印当前流
	剪切并移到剪贴板		复制到剪贴板
	粘贴选择		撤消上一次操作
	重新		搜索节点
	编辑流属性		预览 SQL 生成
	运行当前流		运行流选择

	停止流（仅在流运行期间处于激活状态）		添加超节点
	放大（仅对超节点有效）		缩小（仅对超节点有效）
	流中无标记		插入注释
	隐藏流标记（如果有）		显示隐藏的流标记
	在 IBM® SPSS® Modeler Advantage 中打开流		

流标记由流注释、模型链接和评分分支指示组成。

有关流注释的详细信息，请参阅[添加注释和注解到节点和流](#)第 页码。

有关评分分支指示的详细信息，请参阅[评分分支](#)第 页码。

在《IBM SPSS 建模节点》指南中介绍了模型链接。

自定义工具栏

您可以更改工具栏的不同方面，如：

- 是否显示
- 图标是否有可用工具提示
- 使用大或小图标

要打开或关闭工具栏显示：

- ▶ 在主菜单中，单击：
视图 > 工具栏 > 输出

要更改工具提示或图标大小设置：

- ▶ 在主菜单中，单击：
视图 > 工具栏 > 自定义

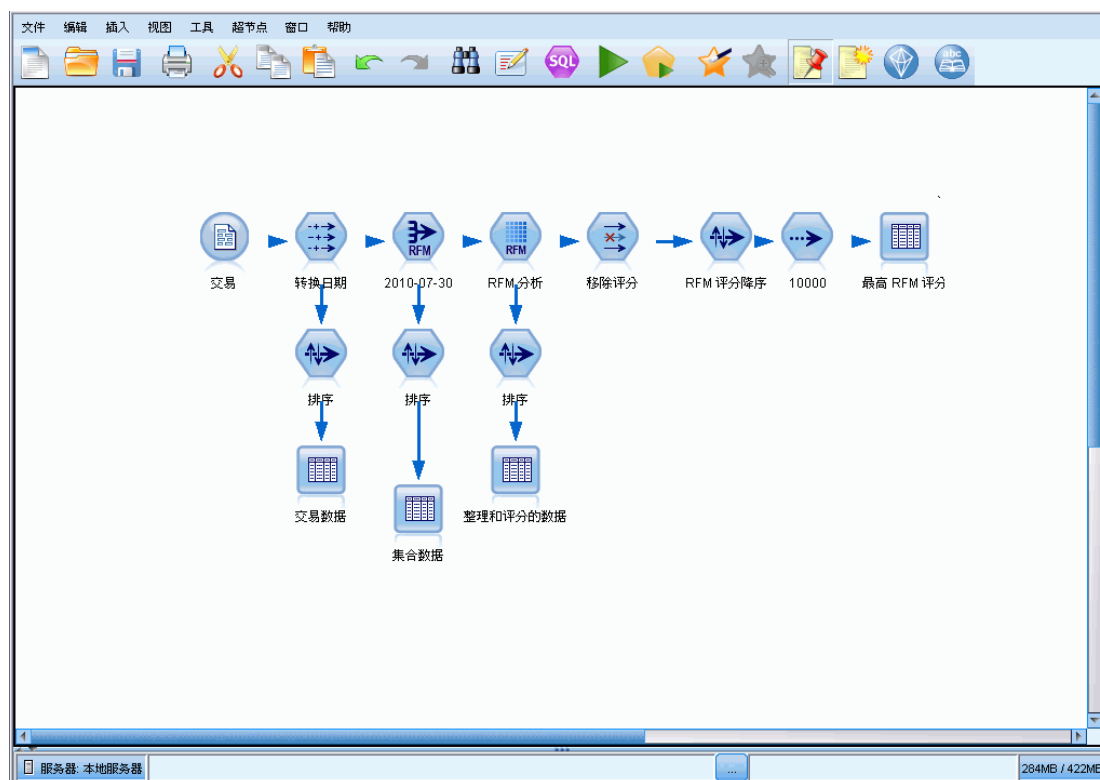
根据需要单击显示工具提示或大按钮。

自定义 IBM SPSS Modeler 窗口

使用 IBM® SPSS® Modeler 界面各部分之间的分界线，可以调整工具的大小或关闭某些工具以满足个人偏好。例如，如果处理的工作流很大，可以使用每条分界线上的小箭头关闭节点选项板、管理器窗格和工程窗格。这样可以最大化流工作区，为处理大型流或多个流提供足够工作空间。

此外，从“视图”菜单上，单击节点选项板、管理器或工程可打开或关闭这些项目的显示。

图片 2-12
最大化的流工作区



另外一种关闭节点选项板和管理器以及工程窗格的方法是：垂直或水平移动 SPSS Modeler 窗口侧面或底部的滚动条，将流工作区当作滚动页面使用。

您还可以控制屏幕标记的显示，该标记由流注释、模型链接和评分分支指示组成。要开启或关闭该显示，请单击：

视图 > 流标记

更改流的图标尺寸

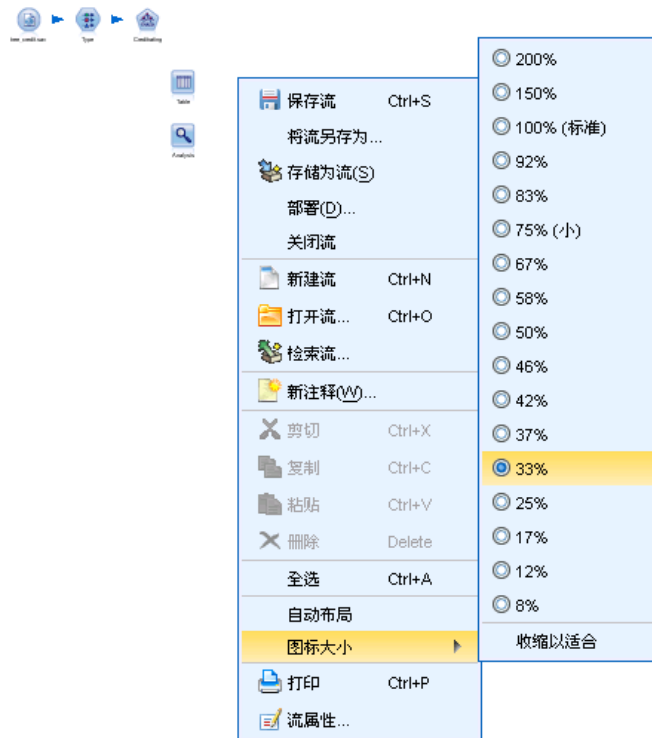
可通过以下方式更改流图标的尺寸。

- 流属性设置

- 流中的弹出菜单
- 使用键盘

您可将整个流视图放大到标准图标尺寸的 8% 至 200% 之间的某个尺寸。

图片 2-13
更改图标尺寸



缩放整个流（流属性方法）

- ▶ 在主菜单中，选择
工具 > 流属性 > 选项 > 布局。
- ▶ 从“图标尺寸”菜单中选择所需的大小。
- ▶ 单击应用以查看结果。
- ▶ 单击确定保存更改。

缩放整个流（菜单方法）

- ▶ 在工作区中右键单击流背景。
- ▶ 选择图标尺寸，并选择所需的大小。

缩放整个流（键盘方法）

- ▶ 在主键盘上按 Ctrl + [-] 以缩小至下一个较小的尺寸。

- ▶ 在主键盘上按 Ctrl + Shift + [+] 以放大至下一个较大的尺寸。

该功能尤其适合用于获得某个复杂流的总体视图。还可通过它来减少在打印某个流时所需的页数。

在 IBM SPSS Modeler 中使用鼠标

IBM® SPSS® Modeler 中最常见的鼠标用法如下所示：

- **单击。** 使用鼠标左键或右键选择菜单选项，打开弹出菜单以及访问其他各种标准控件和选项。单击并按住按键可移动和拖动节点。
- **双击。** 双击鼠标左键可将节点置于流工作区并编辑现有节点。
- **中键单击。** 单击鼠标中键并拖动光标可在流工作区中连接节点。双击鼠标中键可断开某个节点的连接。如果没有三键鼠标，可在单击并拖动鼠标时通过按 Alt 键来模拟此功能。

使用快捷键

IBM® SPSS® Modeler 中的许多可视化编程操作均有与之关联的快捷键。例如，可通过单击某个节点并按键盘上的 Delete 键将此节点删除。同样，可在按住 Ctrl 键的同时按 S 键快速保存某个流。类似控制命令用 Ctrl 键与另一个键的组合来表示，例如 Ctrl+S。

标准 Windows 操作中采用了许多快捷键，例如 Ctrl+X 表示剪切。SPSS Modeler 不仅支持这些快捷键，而且还支持下列应用程序特定的快捷键。

注意：某些时候，SPSS Modeler 中使用的旧快捷键会与标准 Windows 快捷键发生冲突。这些旧快捷键与 Alt 键组合使用仍然有效。例如，Ctrl+Alt+C 可用来切换高速缓存的开与关。

表 2-1
支持的快捷键

快捷键	函数
Ctrl+A	全选
Ctrl+X	剪切
Ctrl+N	新建流
Ctrl+O	打开流
Ctrl+P	打印
Ctrl+C	粘贴
Ctrl+V	粘贴
Ctrl+Z	撤消
Ctrl+Q	选择选定节点的所有下游节点
Ctrl+W	取消选择所有下游节点（使用 Ctrl+Q 切换）
Ctrl+E	从选定节点运行
Ctrl+S	保存当前流
Alt+箭头键	按所用箭头键的方向在流工作区中移动所选节点。
Shift+F10	打开选定节点的弹出菜单

表 2-2
支持的旧热键快捷键

快捷键	函数
Ctrl+Alt+D	复制节点
Ctrl+Alt+L	载入节点
Ctrl+Alt+R	重命名节点
Ctrl+Alt+U	创建用户输入节点
Ctrl+Alt+C	切换高速缓存开关
Ctrl+Alt+F	刷新高速缓存
Ctrl+Alt+X	扩展超节点
Ctrl+Alt+Z	放大/缩小
Delete	删除节点或连接

打印

可在 IBM® SPSS® Modeler 中打印下列对象：

- 流图表
- 图形
- 表
- 报告（来自报告节点和工程报告）
- 脚本（来自“流属性”、“独立脚本”或“超节点脚本”对话框）
- 模型（模型浏览器、包含当前内容的对话框选项卡、树查看器）
- 注解（使用输出的“注解”选项卡）

要打印对象：

- 要不预览就打印，请单击工具栏上的“打印”按钮。
- 要在打印前设置页面，请选择“文件”菜单中的页面设置。
- 要在打印前预览，请选择“文件”菜单中的打印预览。
- 要查看标准打印对话框中用于选择打印机以及指定外观的选项，请选择“文件”菜单中的打印。

自动化 IBM SPSS Modeler

由于高级数据挖掘往往是一个冗长的复杂过程，因此 IBM® SPSS® Modeler 包含对几种类型的编码和自动处理的支持。

- **表达式操作控制语言 (CLEM)** 是一种用于分析和操作在 SPSS Modeler 流中流动的数据的语言。数据挖掘人员可在流操作中广泛使用 CLEM 语言来执行根据成本和收入数据推导利润这样的简单任务，也可以执行将 Web 日志数据转换为具有有

用信息的一系列字段和记录这样的复杂任务。有关详细信息，请参阅第 7 章中的关于 CLEM 中的 IBM SPSS Modeler 15 用户指南。

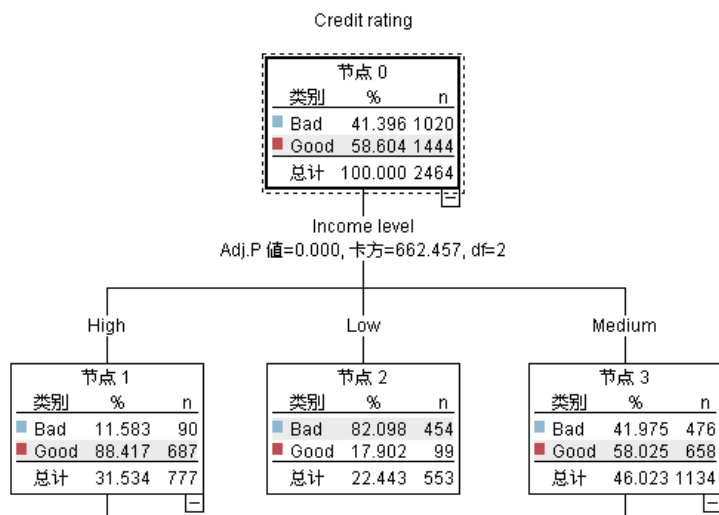
- **脚本编写**是用于在用户界面上实现过程自动化的强大工具。用户以鼠标或键盘实现的操作，也可以通过脚本来实现。可以为节点设置选项并使用 CLEM 子集实现派生。还可以指定输出并操纵生成的模型。有关详细信息，请参阅第 2 章中的脚本编写概述中的 IBM SPSS Modeler 15 脚本编写和自动化指南。

建模简介

模型是一组规则、公式或方程式，可以用它们根据一组输入或变量来预测输出。例如，一家财务机构可根据对过往申请人的已知信息，使用模型预测贷款申请人可能存在优良还是不良风险。

能够预测结果是预测性分析的中心目标，并且了解建模过程是使用 IBM® SPSS® Modeler 的关键。

图片 3-1
简单的决策树模型



本示例使用**决策树**模型，该模型使用一系列决策规则对记录进行分类（并预测响应），例如：

如果收入 = 中等
并且卡 < 5
则 -> “优良”

本示例使用 CHAID（卡方自动交互效应检测）模型时，旨在进行常规的介绍，大部分概念会广泛应用于 SPSS Modeler 中的其他建模类型。

无论要了解哪种模型，均需要首先了解进入该模型的数据。此示例中的数据包含有关银行客户的信息。其中使用了下列字段：

字段名	描述
Credit_rating	信用评价：0=不良，1=优良，9=丢失值
年龄	Age in years
收入	收入水平：1=低，2=中，3=高

字段名	描述
Credit_cards	持有的信用卡数量：1=少于五张， 2=五张或更多
教育	教育程度：1=高中，2=大学
Car_loans	贷款的汽车数量：1=没有或一辆， 2=超过两辆

银行可维护一个包含银行贷款客户历史信息，包括这些客户是正在还贷（信用评价 = 优良）还是在拖欠贷款（信用评价 = 不良）的数据库。银行希望使用现有的数据建立一个模型，允许他们预测未来贷款申请人拖欠贷款的可能性。

使用决策树模型，您可分析两组客户的特征，并预测拖欠贷款的可能性。

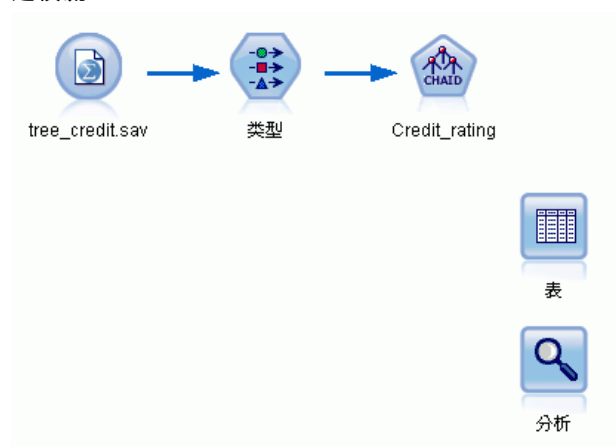
本示例使用了名为 modelingintro.str 的流，该流位于 streams 子文件夹下的 Demos 文件夹中。数据文件是 tree_credit.sav。有关详细信息，请参阅第 1 章中的 Demos 文件夹中的 IBM SPSS Modeler 15 用户指南。

我们来看一下流。

- ▶ 从主菜单中选择下列选项：
文件 > 打开流
- ▶ 单击“打开”对话框的工具栏上的金色模型块图标，然后选择 Demos 文件夹。
- ▶ 双击 streams 文件夹。
- ▶ 双击名为 modelingintro.str 的文件。

构建流

图片 3-2
建模流



要构建流以创建模型，至少需要三个元素：

- 一个从某些外部源读取数据的源节点，在本示例中为 IBM® SPSS® Statistics 数据文件。
- 一个指定字段属性的源节点或“类型”节点，字段属性包括测量级别（字段包含的数据类型）以及每个字段在建模过程中的角色是目标还是输入等。
- 一个在运行流时生成模型块的建模节点。

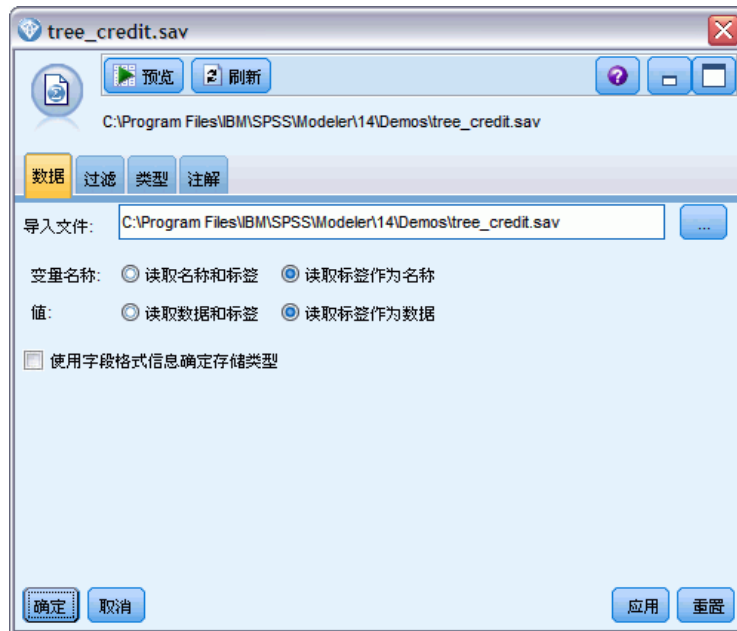
在本例中，我们使用 CHAID 建模节点。CHAID，或卡方自动交互效应检测，是一种通过使用称作卡方统计量的特定统计类型识别决策树中的最优分割来构建决策树的分类方法。

如果在源节点中指定了测量级别，则可以去掉单独的“类型”节点。从功能上来说，结果是一样的。

该流中还包含“表”节点和“分析”节点，创建模型块并将其添加到流中后，将使用这两个节点查看评分结果。

Statistics 文件源节点从 tree_credit.sav 数据文件读取 SPSS Statistics 格式数据，该文件安装在 Demos 文件夹中。（名为 \$CLEO_DEMOS 的特殊变量用于引用位于当前 IBM® SPSS® Modeler 安装下的该文件。这样，无论当前的安装文件夹或版本是什么，均可以确保路径有效。）

图片 3-3
使用 Statistics 文件源节点读取数据



类型节点指定每个字段的**测量级别**。测量级别是一种指示字段中数据类型的类别。我们的源数据文件使用三种不同的测量级别。

连续字段（例如年龄字段）包含连续的数字值，而**名义**字段（例如信用评价字段）有两个或多个不同值，例如不良、优良或无信用历史。**有序**字段（例如收入水平字段）用于描述具有顺序固定的不同值的数据，在本例中为低、中和高。

图片 3-4
用类型节点设置目标和输入字段



对于每个字段，类型节点还指定**角色**，以指示每个字段在建模中扮演的部分。将字段信用评价的角色设置为目标，此字段指示指定的客户是否拖欠贷款。这是**目标**，或者是要预测其值的字段。

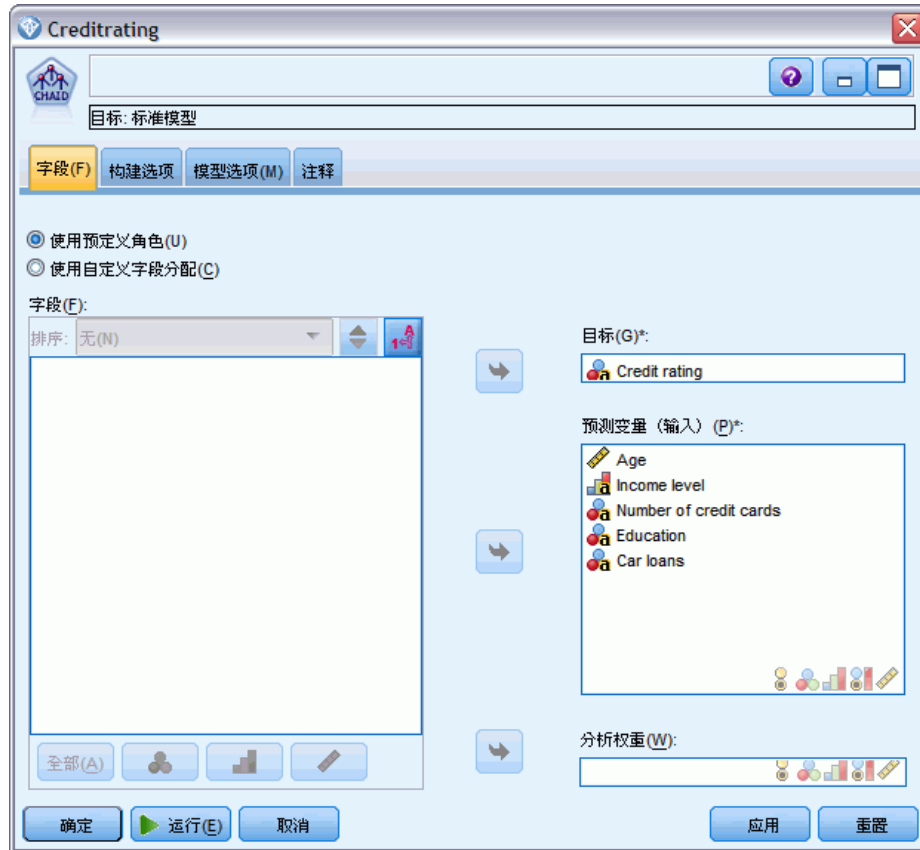
对于其他字段，将角色设置为输入。输入字段有时也称为**预测变量**，或建模算法用其值来预测目标字段值的字段。

CHAID 建模节点生成模型。

在建模节点的“字段”选项卡中，已选中使用预定义角色，这意味着将按在类型节点中的指定使用目标和输入。我们可以在此处更改字段角色，但本例中我们不做任何更改使用这些角色。

- ▶ 单击“构建选项”选项卡。

图片 3-5
CHAID 建模节点、“字段”选项卡



此处包含的选项可以用于指定要构建的模型类型。

由于我们想要一个全新的模型，因此使用默认选项构建新模型。

我们还要求它为单个标准决策树模型，并且不包含任何增强，因此保留默认目标选项构建单个树。

我们可以选择启动允许对模型进行微调的交互建模会话，本示例只使用默认设置生成模型来生成模型。

图片 3-6
CHAID 建模节点、“构建选项”选项卡



对于此示例，我们希望保持树的结构简单，因此通过增加用于父节点和子节点的最小个案数限制树的生长。

- ▶ 在“构建选项”选项卡上，从左侧的导航器窗格选择停止规则。
- ▶ 选择使用绝对值选项。
- ▶ 将父分支中的最小记录数设置为 400。

- ▶ 将子分支中的最小记录数设置为 200。

图片 3-7
为构建决策树设置停止标准

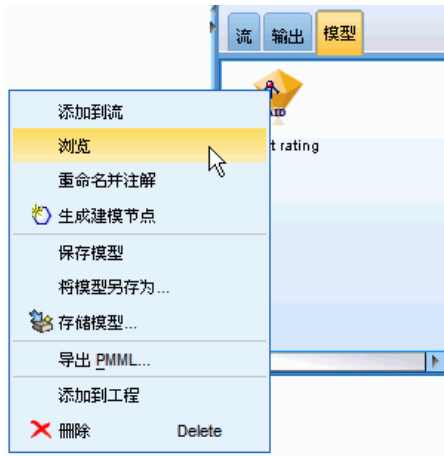


在本例中，我们可以使用所有其他默认选项，因此单击运行以创建模型。（另外，也可以右键单击该节点，然后从上下文菜单中选择运行，或选择节点，并从“工具”菜单中选择运行。）

浏览模型

执行完成后，模型块将添加到应用程序窗口右上角的“模型”选项板中，它还会置于流工作区中，并带有指向创建它的建模节点的链接。要查看模型的详细信息，右键单击模型块并选择浏览（在模型选项板上）或编辑（在工作区上）。

图片 3-8
“模型”选项板



对于 CHAID 模型块，“模型”选项卡以规则集的形式显示详细信息，规则集实际上是根据不同输入字段的值将各个记录分配给子节点的一组规则。

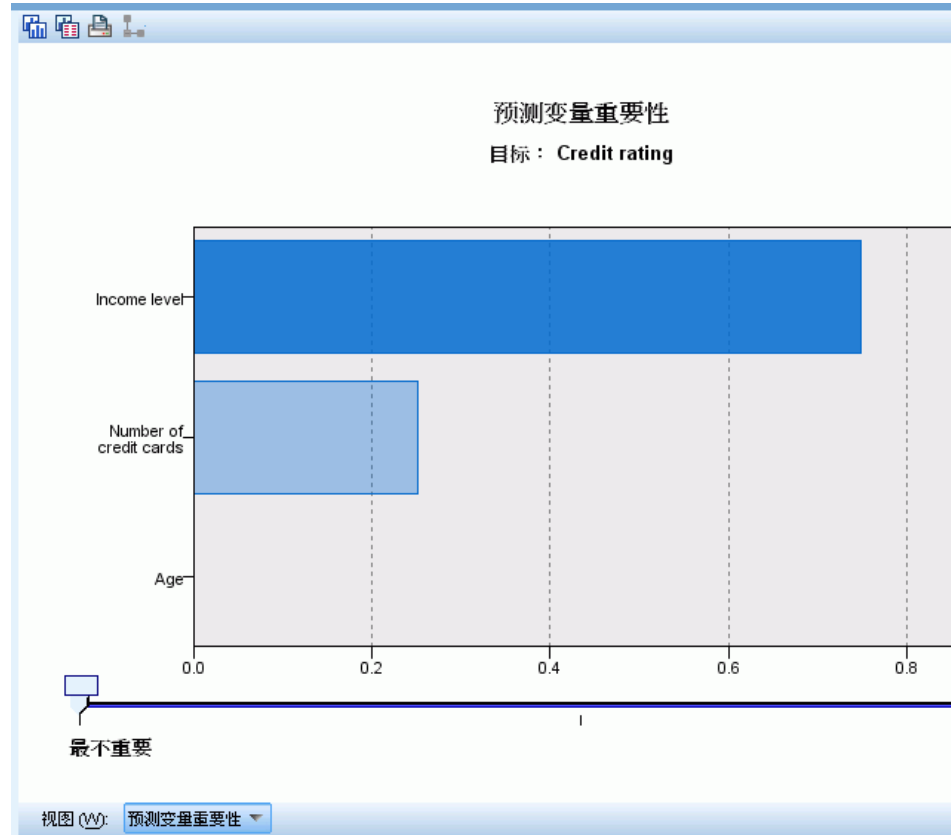
图片 3-9
CHAID 模型块、规则集



对于每个决策树终端节点一意味着那些树节点没有进一步拆分—返回优良或不良的预测值。对于落在该节点内的记录，所有个案中的预测均由**模式**或最常见的响应决定。

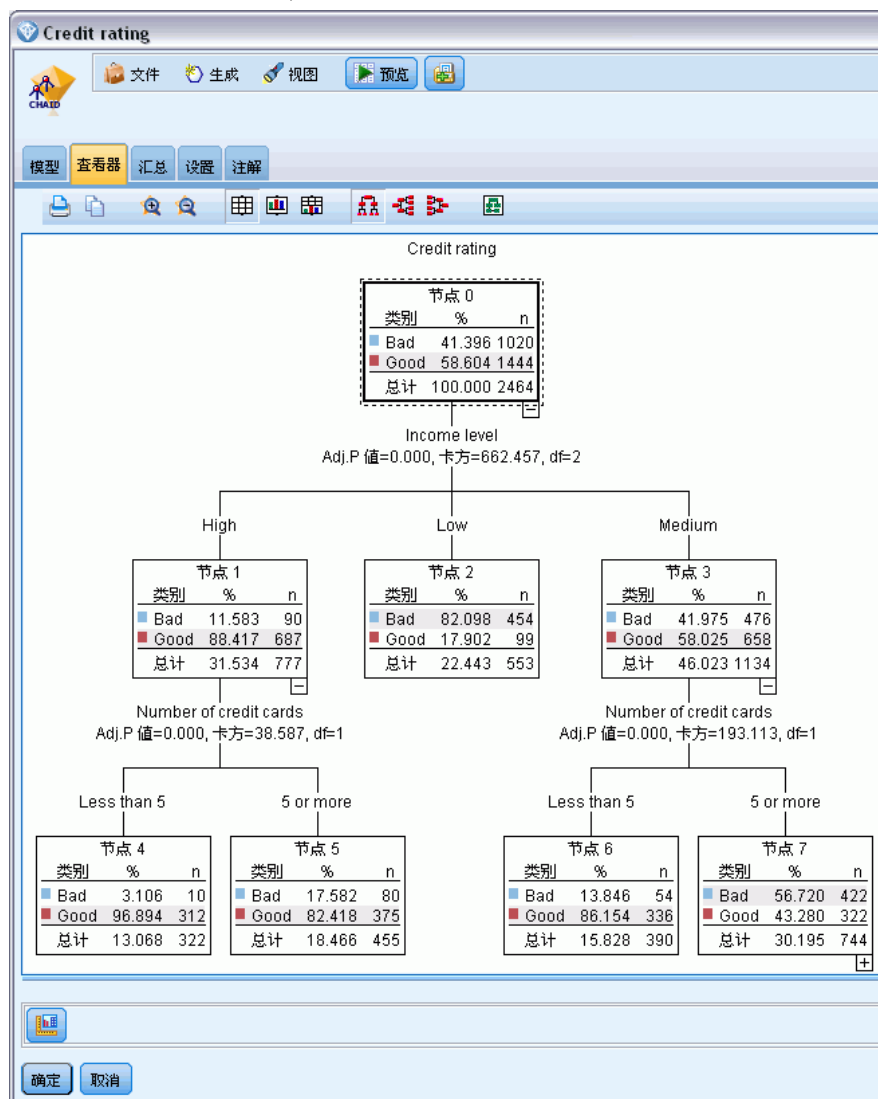
在规则集的右侧，“模型”选项卡显示预测变量重要性图表，该图表显示评估模型时每个预测变量的相对重要性。通过这一点，我们看到收入水平在此个案中最显著，而其他唯一显著的因子是信用卡数量。

图片 3-10
预测变量重要性图表



模型块中的“查看器”选项卡以树的形式显示相同的模型，每个决策点上都有一个节点。可使用工具栏上的缩放控件放大特定节点，或缩小节点以查看更完整的树。

图片 3-11
模型块中的查看器选项卡，已选择缩小



查看树的上部分，第一个节点（节点 0）为我们提供数据集中所有记录的摘要。数据集中超过 40% 的个案分类为不良风险。这是相当高的比例，因此让我们看看树是否能为我们提供哪些因子负责的任何线索。

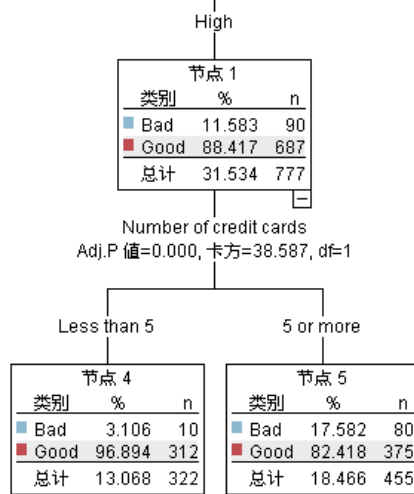
我们可以看到第一个分割是根据收入水平。收入水平位于低类别的记录被指定到节点 2，并且看到此类别包含贷款拖欠人的最高百分比不足为奇。我们可以很明显地了解，此类别中的客户具有高风险。

但是，此类别中的 16% 客户实际上没有拖欠，因此预测并非始终准确。没有模型能够预测每一个响应，但好的模型能够根据可用数据预测对每一个记录作出的最常见的响应。

同样，如果我们查看高收入客户（节点 1），我们看到绝大部分（89%）是优良风险。但是在这些客户中 10 位中有超过 1 位也会拖欠。我们能精炼自己的贷款标准以便将此处风险最小化吗？

注意模型如何根据持有的信用卡数量，将这些客户分成两个子类别（节点 4 和节点 5）。对于高收入客户，如果我们只向那些信用卡少于 5 张的客户贷款，则可以将我们的成功率从 89% 提高到 97%—甚至更满意的结果。

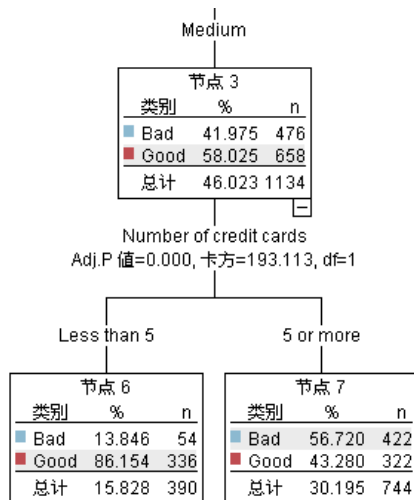
图片 3-12
高收入客户的树状视图



但中等收入类别（节点 3）中的那些客户是什么情况？他们更加均匀地划分为优良和不良评价。

子类别（此情况中是节点 6 和 7）仍然能帮助我们。这次，只向那些信用卡少于 5 张的中等收入客户贷款，可将优良评价的百分比从 58% 提高到 85%，这是显著的改进。

图片 3-13
中等收入客户的树状视图



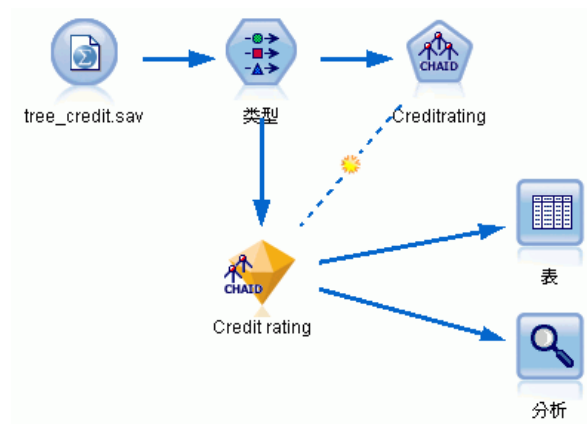
因此，我们了解到输入此模型的每项记录都将被分配到一个特定节点，并且根据该节点最常见的响应分配优良或不良的预测值。

为各个记录分配预测值的这一过程称为**评分**。通过对用于评估该模型的相同记录进行评分，可以评估该模型执行训练数据（我们知道结果的数据）的准确度。让我们看看如何做到这一点。

评估模型

我们浏览了模型以了解评分方式。但是，如果要评估模型的准确度，则需要对一些记录进行评分，并将模型预测的响应与实际结果进行比较。接下来对用于评估该模型的相同记录进行评分，以将观察到的响应与预测响应进行比较。

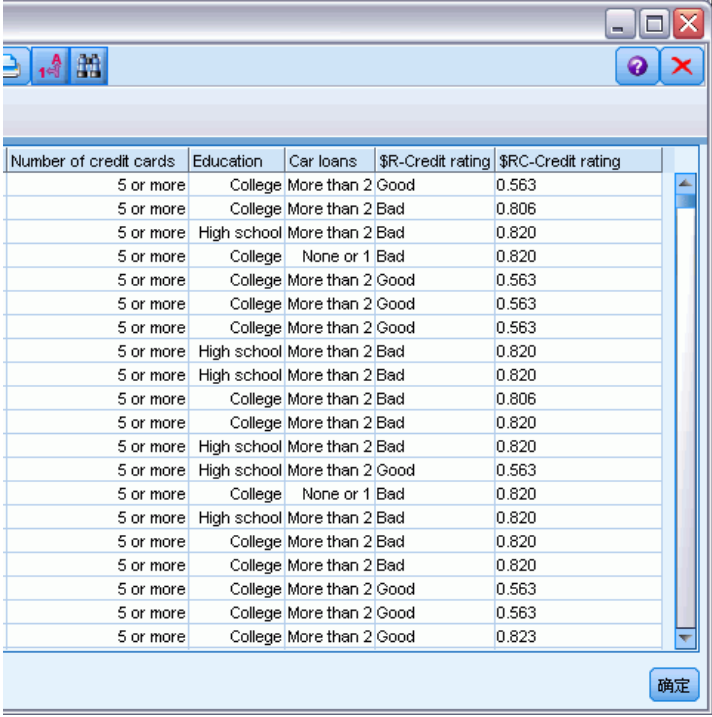
图片 3-14
将模型块附加到输出节点用于模型评估



- ▶ 要查看分数或预测值，请将表节点添加到模型块，然后双击“表”节点，并单击运行。
表在名为 \$R-Credit rating 的字段中显示预测分数，该字段由模型创建。我们可以将这些值与包含实际响应的原始信用评价字段进行比较。

按照惯例，在评分过程中生成的字段的名称基于目标字段，但是要加上标准前缀，例如 \$R- 表示预测值，\$RC- 表示置信度值。不同的模型类型使用不同的前缀集。**置信度值**是模型自己的评估，尺度从 0.0 到 1.0，表示每个预测值的精确程度。

图片 3-15
表格显示生成的分数和置信度值



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	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

与预期的一样，预测值与大多数（并非全部）记录的实际响应相匹配。原因是每个 CHAID 终端节点均有混合响应。预期值与最常见 的响应相匹配，但对于该节点中的其他响应，该预期值是错误的。（记住，16% 的少部分低收入客户没有拖欠。）

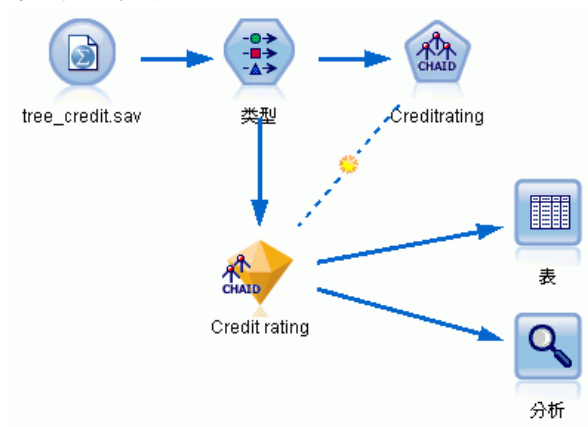
若要避免出现此情况，应继续将树分割为更小的分支，直到每个节点都是不含混合响应的 100% 纯节点为止——即全部为优良或不良。但是，这样的模型可能会非常复杂，并且不易推广到其他数据集。

要查看具体有多少预测值正确，我们可通读表格，并计算预测字段 \$R-Credit rating 的值匹配信用评价的值的记录数量。幸运的是，这里有更简单的方式——我们可使用分析节点，它自动进行此项操作。

- 将模型块连接到分析节点。

- ▶ 双击“分析”节点，然后单击运行。

图片 3-16
添加分析节点



分析表明，2464 个记录中有 1899 个记录（超过 77%）的模型预测值与实际响应相匹配。

图片 3-17
观察到的响应与预测的响应的比较分析结果

输出字段 Credit rating 的结果		
比较 \$R-Credit rating 与 Credit rating		
正确	1,960	79.55%
错误	504	20.45%
总计	2,464	

此结果受到评分的记录和用于评估模型的记录相同的事实限制。在真实情况中，可使用分区节点将数据分割为培训和评估的单独示例。

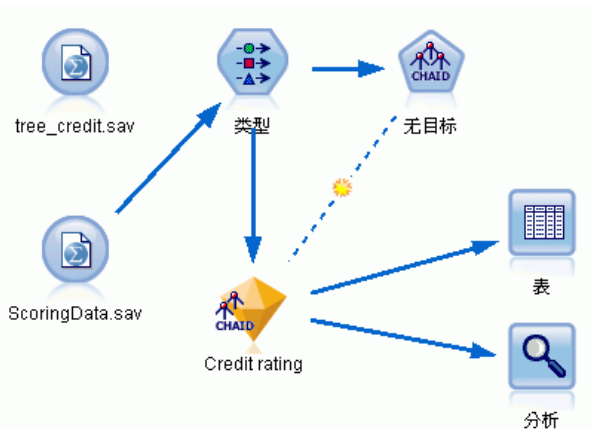
通过使用一个样本分区生成模型并使用另一个样本对模型进行检验，您会得到该模型推广到其他数据集的情况。

通常分析节点，可以针对已知道其实际结果的记录来检验模型。下一阶段介绍如何使用模型对我们不知道结果的记录进行评分。例如，这可能包括当前不是银行客户的人员，但他们是促销邮寄的潜在目标。

对记录评分

之前，我们对用于评估模型的相同记录进行了评分，以评价模型准确程度。现在，我们要查看如何对和用于创建模型不同的记录集进行评分。这是使用目标字段进行建模的目标：研究已知道其结果的记录，以标识您可以从中预测未知结果的模式。

图片 3-18
附加用于评分的新数据



可以更新 Statistics 文件源节点，使它指向其他数据文件，也可以添加一个新的源节点，从它读取要评分的数据。无论采用哪种方式，新数据集包含的输入字段必须与模型（年龄、收入水平、教育等）所使用的相同，但不包含目标字段信用评价。

另外，也可以将模型块添加到包含预期的输入字段的任何流中。无论是读取文件还是数据库，只要字段名和类型与模型使用的相匹配，源类型都无关紧要。

也可以将模型块保存为单独的文件、将模型导出为 PMML 格式以用于其他支持此格式的应用程序，或将模型存储到 IBM® SPSS® Collaboration and Deployment Services 存储库中，这样可以在企业范围对模型进行部署、评分和管理。

无论使用何种基础结构，模型自身都按相同的方式工作。

摘要

本示例演示创建、评估模型以及对模型评分的基本步骤。

-
- 建模节点通过研究已知道其结果的记录来评估模型，并创建模型块。这有时称为训练模型。
 - 可将模型块添加到包含预期字段的任何流中，以对记录进行评分。通过对已知道其结果的记录（如现有客户）进行评分，可以评估模型的运行情况。
 - 如果您对模型的运行情况感到满意，则可以对新数据（如潜在客户）进行评分，以预测他们的响应。
 - 用于训练或评估模型的数据可以称为分析数据或历史数据；评分数据也可以称为操作数据。

标志目标的自动建模

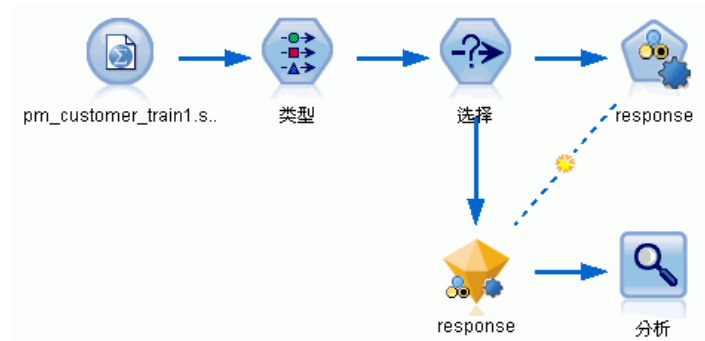
对客户响应建模（自动分类器）

通过“自动分类器”节点，您可以为标志（例如某个客户是否很可能拖欠贷款或者是否会对特定的报价做出响应）或名义（集合）目标自动创建和比较大量的不同模型。在本例中，我们将查找标志（是或否）结果。在一个相对简单的流中，节点生成一组候选模型并对它们进行排序，选择最有效的模型，然后将它们合并为一个汇总（整体）模型。此方法将自动化操作的方便性与组合多个模型的优势融为一体，从而产生任何单一模型所不能带来的更为准确的预测。

本示例以某虚构的公司为例，该公司希望通过为每个客户提供最适用的报价以获取更丰厚的收益。

此方法突出了自动操作的优势。有关使用连续（数值范围）目标的类似示例，请参阅第 5 章第 53 页码。

图片 4-1
自动分类器样本流



本示例使用安装在 streams 目录下 Demo 文件夹中的流 pm_binaryclassifier.str。所使用的数据文件为 pm_customer_train1.sav。有关详细信息，请参阅第 6 页码第 1 章中的 Demos 文件夹。

历史数据

文件 pm_customer_train1.sav 的历史数据可跟踪过去的营销活动中为特定客户提供的报价，由 campaign 字段的值表示。Premium account 活动中的记录数最大。

campaign 字段的值在数据中实际编码为整数（例如 2 = Premium account）。稍后，您可为这些值定义标签以用于给出更有意义的输出。

图片 4-2
以前促销活动的相关数据

表 (31 个字段, 21,927 条记录)

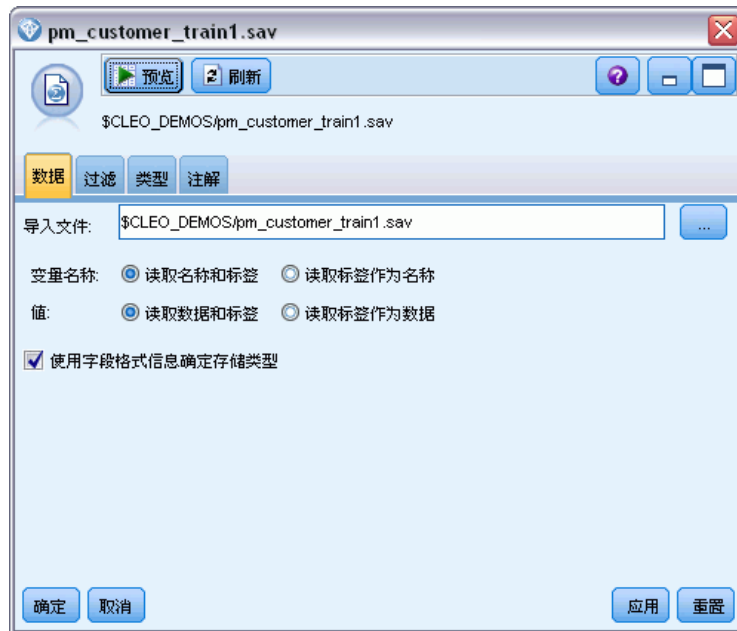
	customer_id	campaign	response	response_date	purchase	purchase_date	product_id	Rowid
1	7	2	0	\$null\$	0	\$null\$	\$null\$	1
2	13	2	0	\$null\$	0	\$null\$	\$null\$	2
3	15	2	0	\$null\$	0	\$null\$	\$null\$	3
4	16	2	1	2006-07-05 00:00:00	0	\$null\$	183	761
5	23	2	0	\$null\$	0	\$null\$	\$null\$	4
6	24	2	0	\$null\$	0	\$null\$	\$null\$	5
7	30	2	0	\$null\$	0	\$null\$	\$null\$	6
8	30	3	0	\$null\$	0	\$null\$	\$null\$	7
9	33	2	0	\$null\$	0	\$null\$	\$null\$	8
10	42	3	0	\$null\$	0	\$null\$	\$null\$	9
11	42	2	0	\$null\$	0	\$null\$	\$null\$	10
12	52	2	0	\$null\$	0	\$null\$	\$null\$	11
13	57	2	0	\$null\$	0	\$null\$	\$null\$	12
14	63	2	1	2006-07-14 00:00:00	0	\$null\$	183	1501
15	74	2	0	\$null\$	0	\$null\$	\$null\$	13
16	74	3	0	\$null\$	0	\$null\$	\$null\$	14
17	75	2	0	\$null\$	0	\$null\$	\$null\$	15
18	82	2	0	\$null\$	0	\$null\$	\$null\$	16
19	89	3	0	\$null\$	0	\$null\$	\$null\$	17
20	89	2	0	\$null\$	0	\$null\$	\$null\$	18

此文件还包含一个响应字段，该字段表明所提供的报价是否被接受（0 = 否，1 = 是）。这将是您希望预测的**目标字段**或值。此外，其中还包括若干包含每位客户的相关人口统计和财务信息的字段。这些字段可用于构建或“训练”一个可基于类似收入、年龄或每月交易次数等特征预测单个用户或用户群响应率的模型。

构建流

- ▶ 添加指向 pm_customer_train1.sav 的 Statistics 文件源节点，该文件位于 IBM® SPSS® Modeler 安装程序的 Demos 文件夹中。（您可以在文件路径中指定 \$CLEO_DEMOS/ 作为引用此文件夹的快捷方式。请注意，路径中必须使用正斜线而非反斜线，如上文所示。）

图片 4-3
读取数据



- ▶ 添加类型节点，然后选择响应作为目标字段（“角色”为目标）。将此字段的“测量”设置为标志。

图片 4-4
设置测量级别和角色



- ▶ 对于以下字段，应将角色设置为无：customer_id、campaign、response_date、purchase、purchase_date、product_id、Rowid 和 X_random。当您构建模型时，将忽略这些字段。
- ▶ 单击类型节点的读取值按钮以确保值获得实例化。

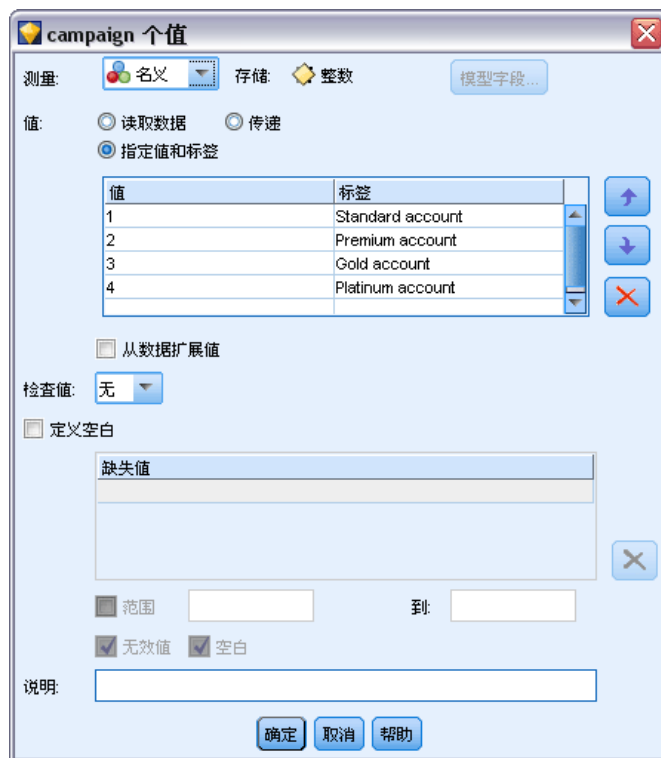
从前文看出，我们的源数据包含有关四项不同活动的信息，每个活动针对不同类型的客户帐户。这些活动在数据中编码为整数，以方便记住每个整数所代表的帐户类型，让我们为每一个都定义标签。

图片 4-5
选择以指定字段值



- ▶ 在活动字段的行上，单击值列中的条目。
- ▶ 从下拉列表选择指定。

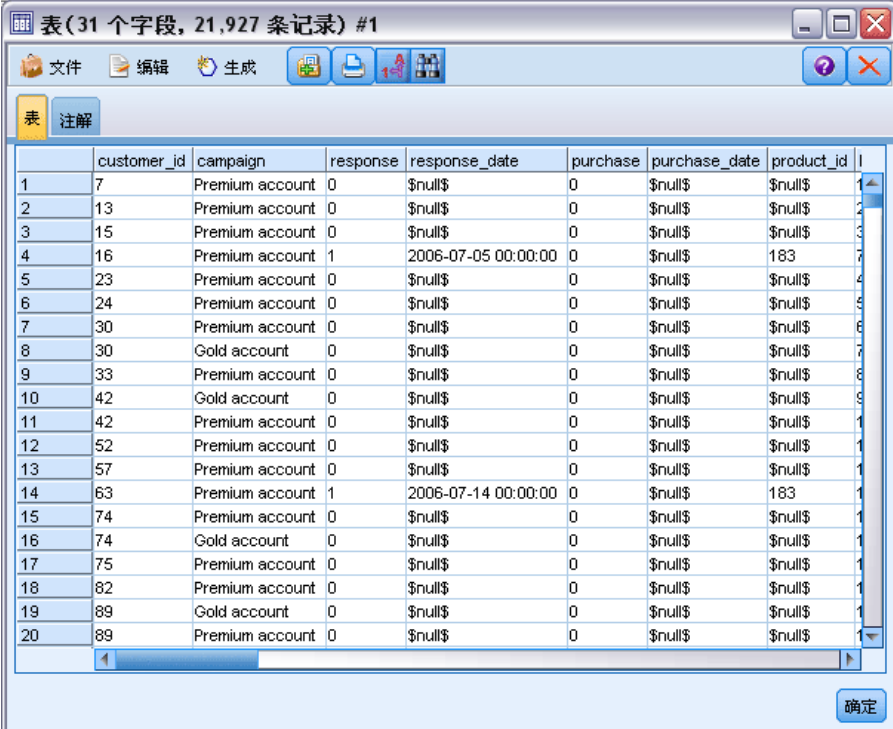
图片 4-6
定义字段值的标签



- ▶ 在标签列中，键入活动字段四个值中每个值所显示的标签。
- ▶ 单击确定。

现在您可在输出窗口中显示标签而非整数了。

图片 4-7
显示字段值标签



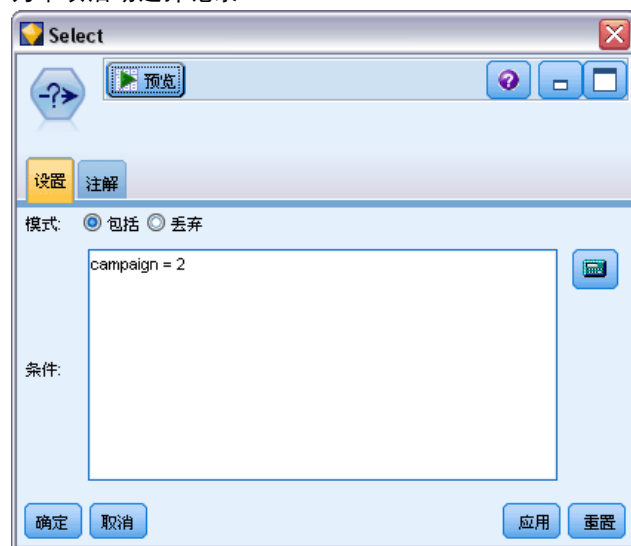
表(31 个字段, 21,927 条记录) #1

	customer_id	campaign	response	response_date	purchase	purchase_date	product_id
1	7	Premium account	0	\$null\$	0	\$null\$	\$null\$
2	13	Premium account	0	\$null\$	0	\$null\$	\$null\$
3	15	Premium account	0	\$null\$	0	\$null\$	\$null\$
4	16	Premium account	1	2006-07-05 00:00:00	0	\$null\$	183
5	23	Premium account	0	\$null\$	0	\$null\$	\$null\$
6	24	Premium account	0	\$null\$	0	\$null\$	\$null\$
7	30	Premium account	0	\$null\$	0	\$null\$	\$null\$
8	30	Gold account	0	\$null\$	0	\$null\$	\$null\$
9	33	Premium account	0	\$null\$	0	\$null\$	\$null\$
10	42	Gold account	0	\$null\$	0	\$null\$	\$null\$
11	42	Premium account	0	\$null\$	0	\$null\$	\$null\$
12	52	Premium account	0	\$null\$	0	\$null\$	\$null\$
13	57	Premium account	0	\$null\$	0	\$null\$	\$null\$
14	63	Premium account	1	2006-07-14 00:00:00	0	\$null\$	183
15	74	Premium account	0	\$null\$	0	\$null\$	\$null\$
16	74	Gold account	0	\$null\$	0	\$null\$	\$null\$
17	75	Premium account	0	\$null\$	0	\$null\$	\$null\$
18	82	Premium account	0	\$null\$	0	\$null\$	\$null\$
19	89	Gold account	0	\$null\$	0	\$null\$	\$null\$
20	89	Premium account	0	\$null\$	0	\$null\$	\$null\$

- ▶ 将表节点附加到类型节点。
- ▶ 打开“表”节点，然后单击运行。
- ▶ 在输出窗口上，单击显示字段和值标签工具栏按钮以显示标签。
- ▶ 单击确定关闭输出窗口。

尽管数据包含有关四项不同活动的信息，但每一次的分析应集中关注其中一项活动。由于 Premium account 活动（在数据中编码为 campaign=2）中的记录数最大，因此可以使用选择节点实现仅在流中包含这些记录。

图片 4-8
为单项活动选择记录

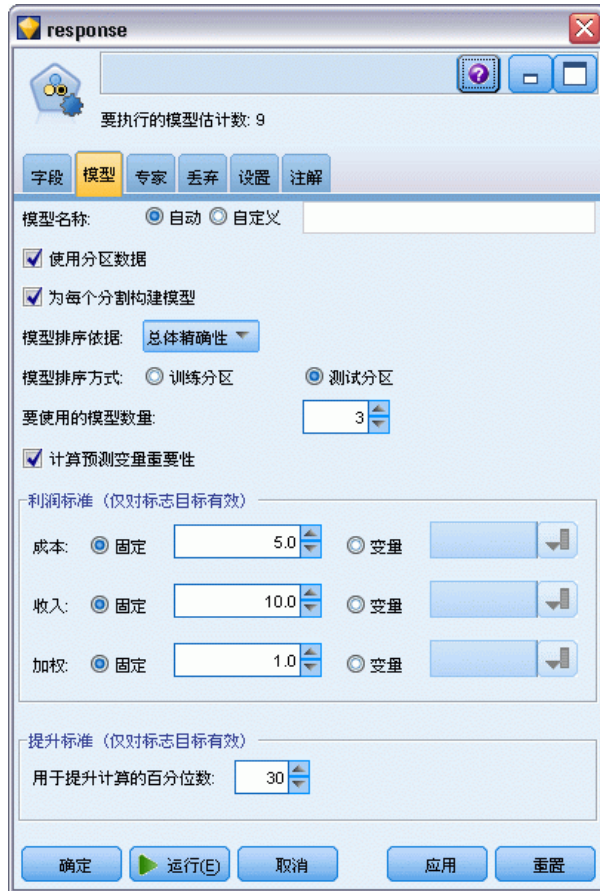


生成和比较模型

- ▶ 附加一个自动分类器节点，然后选择总体精确性作为对模型进行排序的度量。

- ▶ 将要使用的模型数设置为 3。这意味着在执行节点时将构建三个最佳模型。

图片 4-9
自动分类器节点“模型”选项卡



在“专家”选项卡上，可从最多 11 种不同模型算法中进行选择。

- ▶ 取消选择判别式和 SVM 模型类型。（这些模型需要花费更多时间培训这些数据，因此取消选中它们将可以加快示例的执行速度。如果您不介意稍等一下，也可以保留它们的选中状态。）

由于在“模型”选项卡上将要使用的模型数设置为 3，因此节点将计算余下九个算法的准确性，并构建包含三个最准确算法的单个模型块。

图片 4-10
自动分类器节点“专家”选项卡



- ▶ 在“设置”选项卡上，对于整体方法，请选择置信度加权投票。此选项确定如何为每条记录生成一个汇总得分。

使用简单投票方式时，若三个模型中有两个模型均预测是，则是将以 2 比 1 的投票结果取胜。在使用置信度加权投票方式的情况下，将基于各预测的置信度值进行加权投票。因此，如果一个预测否的模型的置信度比两个预测是的模型合在一起的置信度还高，则否取胜。

图片 4-11
自动分类器节点：“设置”选项卡



- ▶ 单击运行。

几分钟后，将构建生成的模型块，并放到工作区和窗口右上角的“模型”选项板中。您可浏览模型块，或以多种其他方式将其保存或部署。

打开模型块，它将列出在运行期间所创建的每个模型的详细信息。（实际情况中，由于大型数据集往往需要创建数百个模型，这可能会花费数小时的时间。）请参阅第 40 页码中的[图片 4-1](#)。

如果需要进一步探索任何单独的模型，可在**模型**列中双击此模型块图标，以向下浏览至单独模型结果，您可以从中生成建模节点、模型块或评估图表。在**图形**列中，可以双击缩略图生成标准大小的图形。

图片 4-12
自动分类器结果



是否...	图形	模型	构建时间 (分钟)	最大 利润	最大利润 发生比率 (%)	提升(Top 30%)	总体 精确性 (%)	使用的 字段编号	曲线下方 区域
<input checked="" type="checkbox"/>		C51	< 1	4,906.667	8	2.203	92.861	10	0.777
<input checked="" type="checkbox"/>		C&R...	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

默认情况下，模型会基于总体精确性排序，因为这是您在自动分类器节点“模型”选项卡中选择的度量。根据这一度量，C51 模型的精确性最高，但 C&R 树和 CHAID 模型的精确性与之相差不大。

您可以通过单击其他列的标题对该列进行排序，或者也可以从工具栏的排序方式下拉列表中选择所需的度量。

基于这些结果，您可决定使用所有三个最准确的模型。通过结合多个模型的预测，可以避免单个模型的局限性，从而使整体准确性更高。

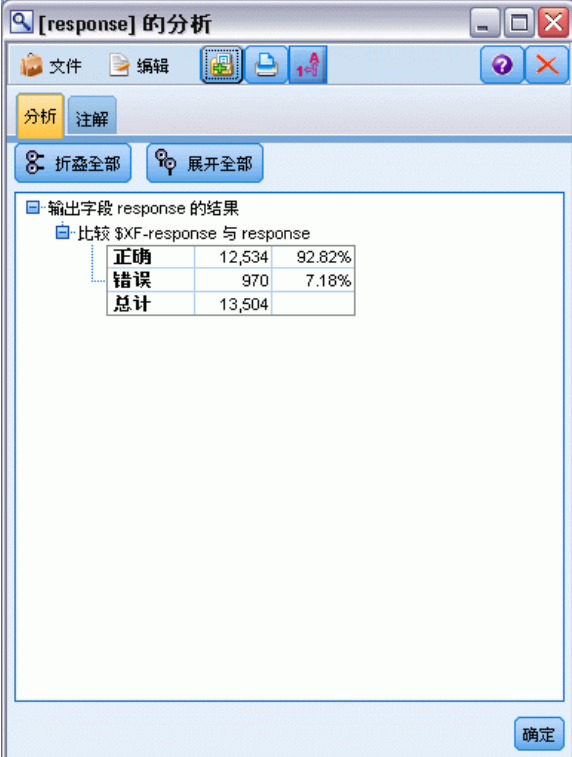
在使用? 列中，选择 C51，C&R 树和 CHAID 模型。

在模型块后附加一个分析节点（“输出”选项板）。右键单击分析节点，然后选择运行以运行流。

由整体模型生成的汇总得分将显示在名为 \$XF-response 的字段中。当根据训练数据度量时，预测值与实际响应（如原始响应字段中的记录所示）匹配的总精确性为 92.82%。

尽管这不如本例中三个模型的最高精确性高（C51 为 92.86%），但它们之间的差距小得可以忽略不计。一般来说，整体模型应用到除训练数据之外的数据集中时，通常更有可能效果较好。

图片 4-13
对三个整体模型的分析



The screenshot shows a software window titled "[response] 的分析" (Analysis of [response]). The window contains a menu bar with "文件" (File) and "编辑" (Edit), and a toolbar with icons for file operations. Below the toolbar are buttons for "分析" (Analyze) and "注解" (Annotations). There are also buttons for "折叠全部" (Collapse All) and "展开全部" (Expand All). The main content area displays a tree view under "输出字段 response 的结果" (Results for output field response), with a sub-entry "比较 XF-response 与 response" (Compare XF-response with response). A table is shown with the following data:

正确	12,534	92.82%
错误	970	7.18%
总计	13,504	

A "确定" (OK) button is located at the bottom right of the window.

摘要

综上所述，您使用自动分类器节点比较了多种不同的模型，然后使用三个最准确的模型并将它们添加到位于一个整体自动分类器模型块内的流中。

- 基于总体精确性，“C51”、“C&R 树”和 CHAID 模型对于训练数据效果最佳。
- 整体模型与最好的单个模型相比效果相差不大，而且当应用到其他数据集时可以起到更好的效果。如果您的目标是尽可能多地自动执行这一过程，您可以通过此方法获得在大多数情况下都很稳健的模型，而无需深入挖掘任意一个模型的细节。

连续目标的自动建模

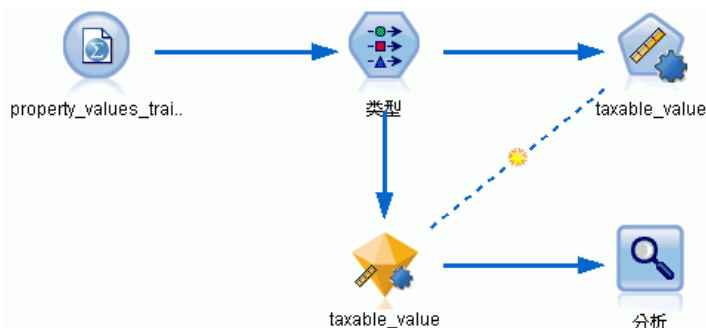
属性值（自动数值）

通过自动数值节点，您可以自动创建和比较连续（数值范围）结果的不同模型，例如，预测属性的应征税值。借助于单独节点，可以估计和比较一组候选模型，并生成一个模型子集以进一步分析。这类节点与自动分类器节点工作方式相同，但连续不仅限于标志或名义目标。

该节点将候选模型的最佳结果合并到单个汇总（整体）模型块中。此方法将自动化操作的方便性与组合多个模型的优势融为一体，从而产生任何单一模型所不能带来的更为准确的预测。

本示例重点关注一个虚拟的负责调整和评估房地产税的市政机构。为使预测更为准确，他们将构建一个根据建筑类型、周边状况、占地面积以及其他已知因素预测属性值的模型。

图片 5-1
自动数值样本流



此示例使用安装在 streams 下 Demos 文件夹中的流 property_values_numericpredictor.str。所使用的数据文件为 property_values_train.sav。有关详细信息，请参阅第 6 页码第 1 章中的 Demos 文件夹。

训练数据

数据文件包含一个名为 taxable_value 的字段，该字段就是要预测的目标字段或值。其他字段所包含的信息有周边情况、建筑类型以及内部体积，它们均可以用作预测变量。

字段名	Label
property_id	属性 ID
周边状况	城市内的区域
building_type	建筑物的类型

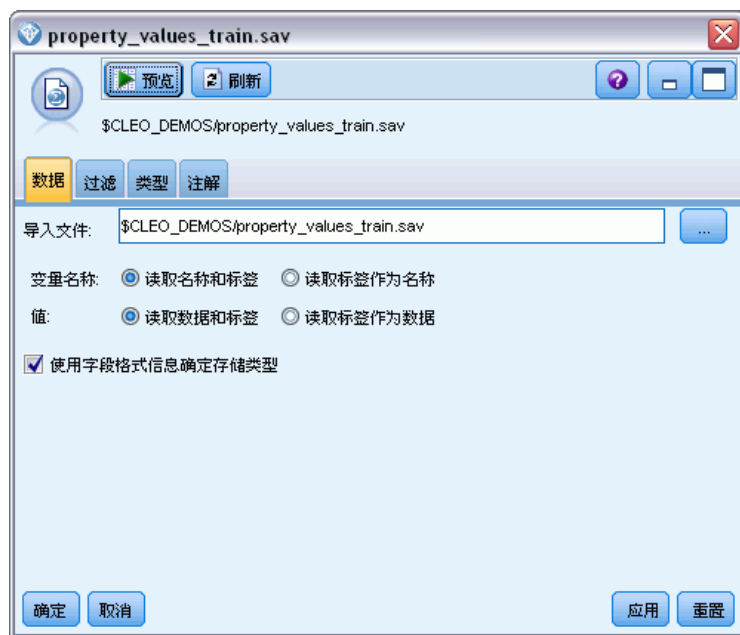
字段名	Label
year_built	建造年代
volume_interior	内部体积
volume_other	车库和其他建筑所占的体积
lot_size	占地面积
taxable_value	应征税值

在 Demos 文件夹中还包括一个名为 property_values_score.sav 的评分数据文件。该文件中除了没有 taxable_value 字段之外，剩下的字段与数据文件相同。在训练模型使用已知应征税值的数据集之后，您就可以对仍不知晓应征税值的记录进行评分。

构建流

- ▶ 添加指向 property_values_train.sav 的 Statistics 文件源节点，该文件位于 IBM® SPSS® Modeler 安装程序的 Demos 文件夹中。（您可以在文件路径中指定 \$CLEO_DEMOS/ 作为引用此文件夹的快捷方式。请注意，路径中必须使用正斜线而非反斜线，如上文所示。）

图片 5-2
读取数据



- ▶ 添加类型节点，然后选择 `taxable_value` 作为目标字段（角色为目标）。所有其他字段的角色均应设置为输入，从而指示这些字段将用作预测变量。

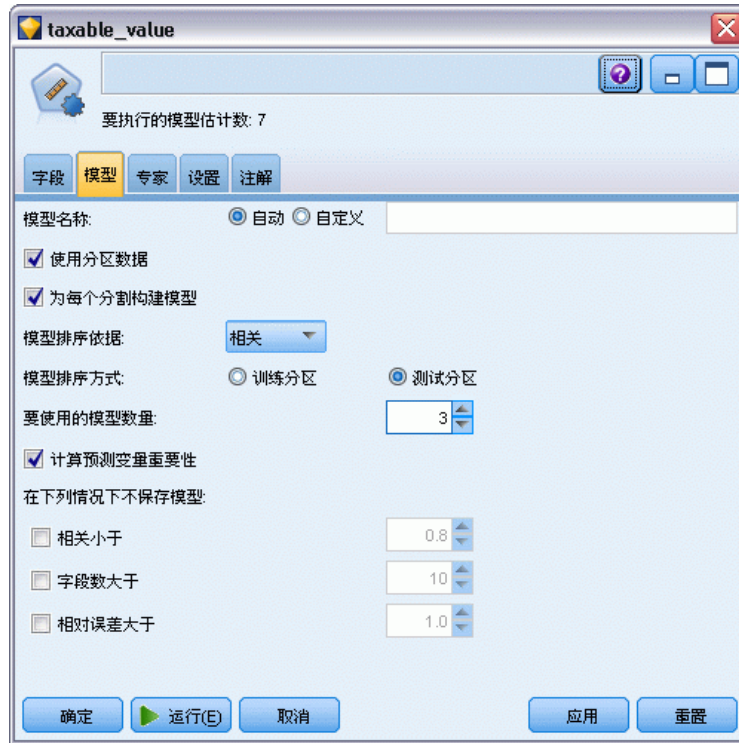
图片 5-3
设置目标字段



- ▶ 附加自动数值节点，并选择相关性作为对模型排序的方法。

- ▶ 将要使用的模型数设置为 3。这意味着在执行节点时将构建三个最佳模型。

图片 5-4
自动数值节点“模型”选项卡



- ▶ 在“专家”选项卡中，保留默认设置。节点将为每个算法评估单个模型（共七个模型）。（或者，您可以修改这些设置，以对每个模型类型的多个变量进行比较。）

由于在“模型”选项卡上将要使用的模型数设置为 3，因此节点将计算七个算法的准确性，并构建包含三个最准确算法的单个模型块。

图片 5-5
自动数值节点“专家”选项卡



- ▶ 在“设置”选项卡中，保留默认设置。由于这是一个连续目标，因此会由各个模型的平均得分生成整体得分。

图片 5-6
自动数值节点“设置”选项卡



比较模型

- ▶ 单击“运行”按钮。

构建模型块，并放到工作区和窗口右上角的“模型”选项板中。您可浏览模型块，或以多种其他方式将其保存或部署。

打开模型块，它将列出在运行期间所创建的每个模型的详细信息。（实际情况中，由于大型数据集往往需要评估数百个模型，这可能会花费数小时的时间。）请参阅第 53 页码中的 [图片 5-1](#)。

如果需要进一步探索任何单独的模型，可在模型列中双击此模型块图标，以向下浏览至单独模型结果，您可以从中生成建模节点、模型块或评估图表。

图片 5-7
自动数值结果



是否使用?	图形	模型	构建时间 (分钟)	相关	使用的 字段编号	相对错误
<input checked="" type="checkbox"/>		Generalize...	< 1	0.915	7	0.162
<input checked="" type="checkbox"/>		Regressio...	< 1	0.9	5	0.19
<input checked="" type="checkbox"/>		CHAID Tre...	< 1	0.892	5	0.204

默认情况下，模型将按相关性进行排序，这是因为您在自动数值节点中将相关性选作了测量量。出于排序的目的，使用了相关性的绝对值，值越接近于 1 则说明关系越强。在本测量中，广义线性模型的排序位置最高，但是还有几个模型也近乎准确。除此之外，广义线性模型还具有最低的相对错误。

通过单击列标题或从工具栏上的排序方式列表中选择所需的测量，您可以对不同的列进行排序。

每个图形都显示了相对于模型预测值的观测值散点图，从而可以快速直观地表示模型之间的相关性。对一个好的模型来说，所有的点都应聚集在对角线附近，在本例中所有模型都是如此。

在图形列中，可以双击缩略图生成标准大小的图形。

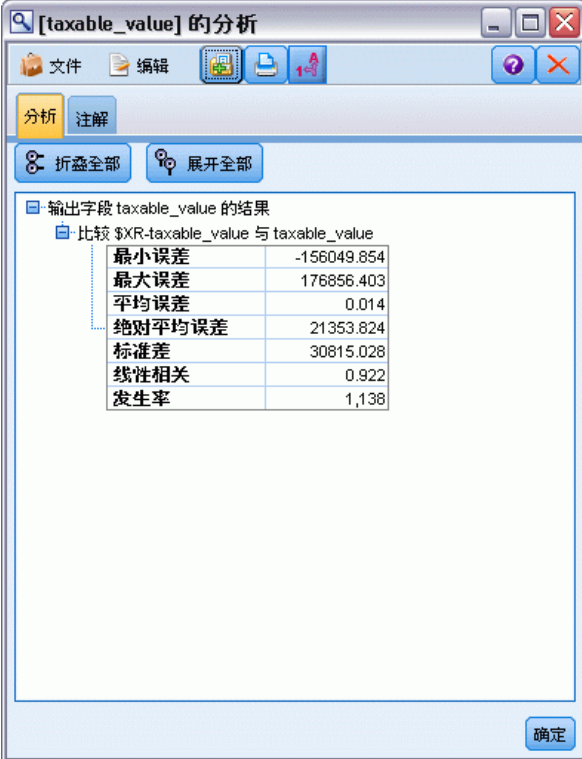
基于这些结果，您可决定使用所有三个最准确的模型。通过结合多个模型的预测，可以避免单个模型的局限性，从而使整体准确性更高。

在使用?列中，确保选中了所有三个模型。

在模型块后附加一个分析节点（“输出”选项板）。右键单击分析节点，然后选择运行以运行流。

由整体模型生成的平均得分会添加到名为 `$XR-taxable_value` 且相关性为 0.922 的字段中，该相关性值高于三个单独模型中的这些相关性值。该整体节点还显示了较低的绝对平均误差，因此与任何单独模型相比，在应用到其他数据集时，执行效果可能会更好。

图片 5-8
自动数值样本流



The screenshot shows a software window titled "[taxable_value] 的分析". It contains a table with the following data:

输出字段 taxable_value 的结果	
比较 \$XR-taxable_value 与 taxable_value	
最小误差	-156049.854
最大误差	176856.403
平均误差	0.014
绝对平均误差	21353.824
标准差	30815.028
线性相关	0.922
发生率	1,138

摘要

综上所述，您使用自动数值节点比较了多种不同的模型，然后选定三个最准确的模型并将它们添加到位于一个整体自动数值模型块内的流中。

- 基于总体准确性的考虑，广义线性模型、回归模型以及 CHAID 模型对训练数据的执行效果最好。
- 整体模型显示出优于两个单独模型的效果，在应用到其他数据库时，执行效果可能要更好一些。如果您的目标是尽可能多地自动执行这一过程，您可以通过此方法获得在大多数情况下都很稳健的模型，而无需深入挖掘任意一个模型的细节。

部分 II: 数据准备示例

自动数据准备 (ADP)

准备分析数据是任何数据挖掘项目中最重要的一步，而从传统来说也是最耗时的步骤之一。自动数据准备 (ADP) 您可以完全自动化地使用节点，允许节点选择并应用修正，或者也可在修正前预览更改，按照需要接受或拒绝。

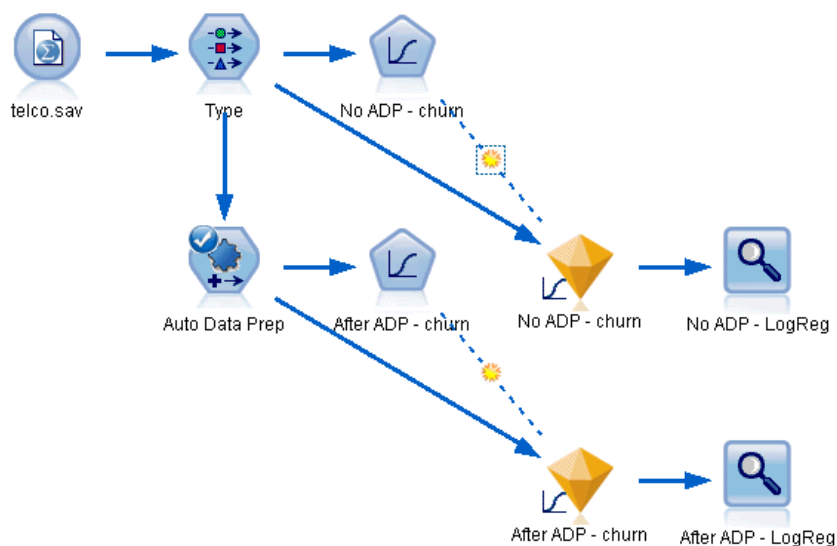
通过 ADP 节点，可以快速、方便地准备数据以供数据挖掘，无需具备相关统计概念的预备知识。如果使用默认设置运行节点，您将可以更快地构建模型并进行评分。

本示例使用名为 ADP_basic_demo.str 的流，该流引用名为 telco.sav 的数据文件来演示构建模型时通过使用默认的 ADP 节点设置所增加的准确性。这些文件可在任何 IBM® SPSS® Modeler 安装程序的 Demos 目录中找到。此目录可通过 WindowsIBM® SPSS® Modeler 程序组进行访问。文件 ADP_basic_demo.str 位于 streams 目录下。

构建流

- ▶ 要构建流，请添加指向 telco.sav 的 Statistics 文件源节点，telco.sav 位于 IBM® SPSS® Modeler 安装程序的 Demos 目录中。

图片 6-1
构建流



- ▶ 将一个类型节点附加到源节点，将 churn 字段的测量级别设置为标志，并将角色设置为目标。将所有其他字段的角色设置为 Input。

图片 6-2
选择目标



- ▶ 将 Logistic 节点附加到“类型”节点。

- ▶ 在 Logistic 节点上，单击“模型”选项卡并选择二项过程。在模型名称字段中，选择自定义并输入 No ADP - churn。

图片 6-3
选择模型选项



- ▶ 将 ADP 节点附加到“类型”节点。在“目标”选项卡上保留默认设置，以均衡速度与准确性的方式分析和准备数据。
- ▶ 在“目标”选项卡顶部，单击分析数据以分析和处理数据。

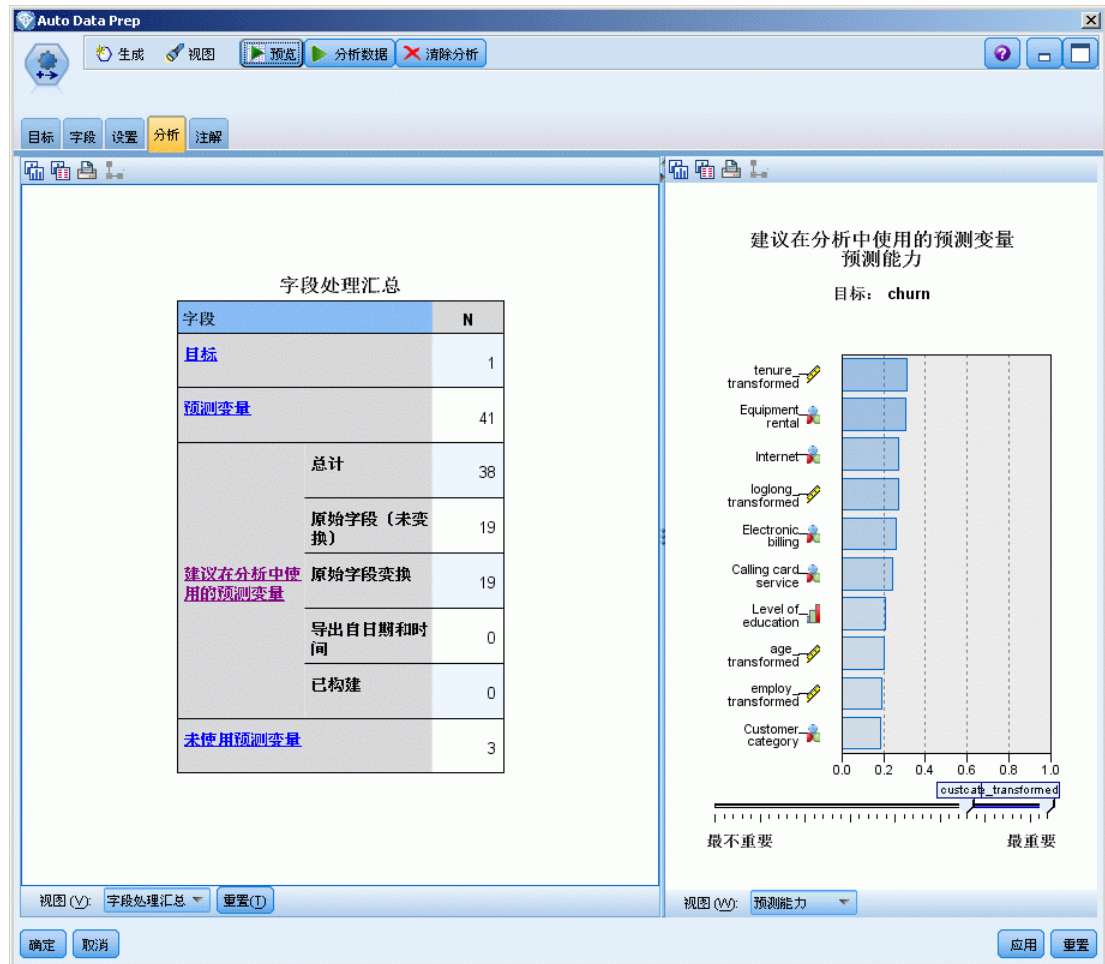
ADP 节点上的其他选项允许您指定优先考虑准确性还是处理速度，或者对许多数据准备处理步骤进行微调。

图片 6-4
ADP 默认目标



数据处理的结果显示在“分析”选项卡上。字段处理摘要显示 41 个被导入到 ADP 节点的数据特征中，19 个被转换为辅助处理，3 个因未使用被丢弃。

图片 6-5
数据处理摘要



- 将 Logistic 节点附加到 ADP 节点。

- ▶ 在 Logistic 节点上，单击“模型”选项卡并选择二项过程。在建模名称字段中，选择自定义并输入 After ADP - churn。

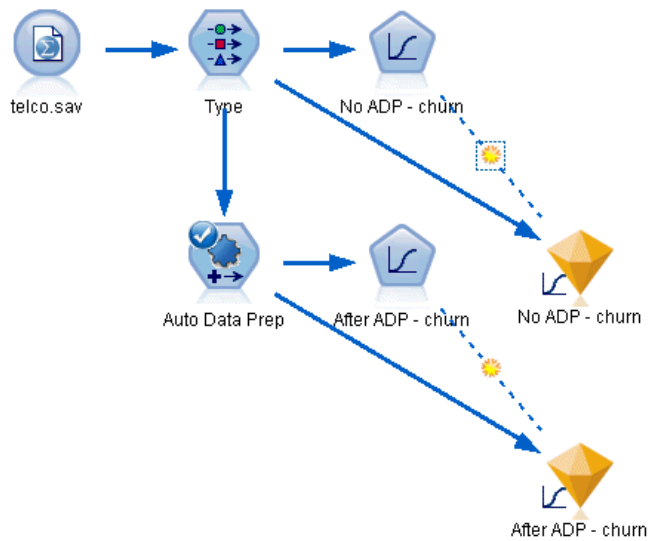
图片 6-6
选择模型选项



比较模型准确性

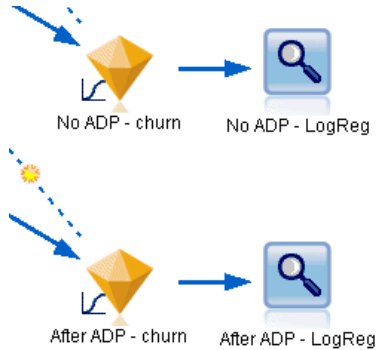
- ▶ 运行两个 Logistic 节点以创建模型块，该模型块将被添加到流和位于右上角的“模型”选项板中。

图片 6-7
附加模型节点



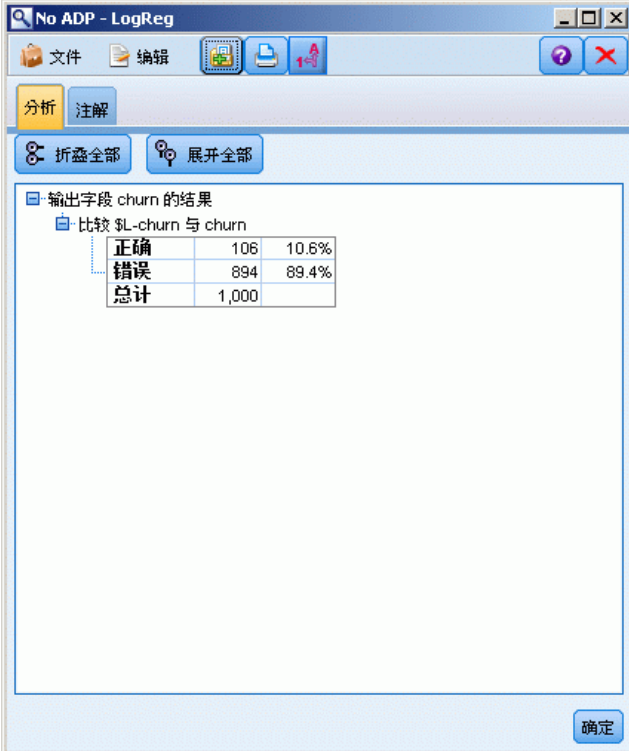
- ▶ 将“分析”节点附加到模型块，并使用其默认设置运行“分析”节点。

图片 6-8
附加分析节点



对非 ADP 派生模型的“分析”显示，在“Logistic 回归”节点中使用默认设置运行数据会使模型的准确性较低 - 仅为 10.6%。

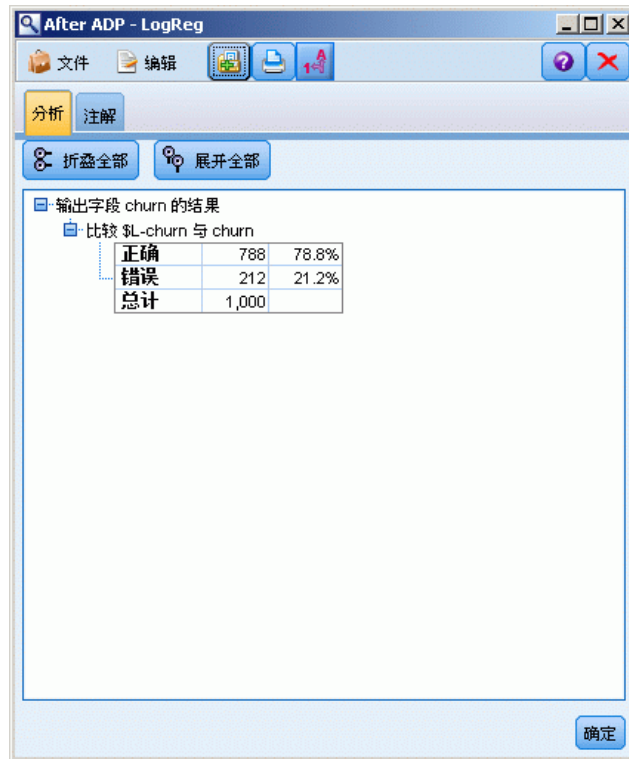
图片 6-9
非 ADP 派生模型结果



类别	数量	百分比
正确	106	10.6%
错误	894	89.4%
总计	1,000	

对 ADP 派生模型的“分析”显示，使用默认的 ADP 设置运行数据，您可构建正确率为 78.8% 的非常准确的模型。

图片 6-10
ADP 派生模型结果



总之，通过运行 ADP 节点对数据处理过程进行微调，您只需很少的直接数据操作便可构建更准确的模型。

当然，如果您有兴趣证明或推翻这一论断，或想构建特定模型，您会发现直接使用模型设置很有用；然而，对于那些构建时间短或有大量数据要处理的模型，ADP 节点可为您带来优势。

有关 IBM® SPSS® Modeler 中所用建模方法的数学原理的说明，请参阅《SPSS Modeler 算法指南》，该指南位于安装光盘的 \Documentation 目录中。

请注意，本示例中的结果仅基于训练数据产生。要评估模型对实际应用中的其他数据的拟合程度，可使用“分区”节点保留部分记录，以便于测试和验证。[有关详细信息，请参阅第 4 章中的分区节点中的 IBM SPSS Modeler 15 源、过程和输出节点。](#)

准备分析数据（数据审核）

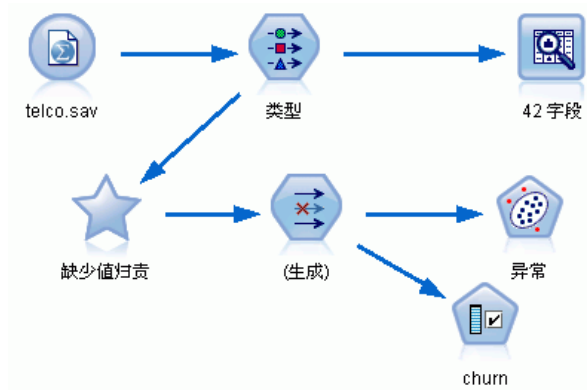
数据审核节点将首先全面检查用户导入到 IBM® SPSS® Modeler 的数据。我们在初始数据研究过程中经常会使用数据审核报告来显示汇总统计量以及每个数据字段的直方图和分布图，还可以确定如何处理缺失值、离群值和极值。

本例使用名为 telco_dataaudit.str 的流，该流引用名为 telco.sav 的数据文件。这些文件可在任何 IBM® SPSS® Modeler 安装程序的 Demos 目录中找到。此目录可通过 Windows SPSS Modeler 程序组进行访问。文件 telco_dataaudit.str 位于 streams 目录下。

构建流

- ▶ 要构建流，请添加指向 telco.sav 的 Statistics 文件源节点，telco.sav 位于 IBM® SPSS® Modeler 安装程序的 Demos 目录中。

图片 7-1
构建流



- ▶ 添加“类型”节点以定义字段，并将 churn 指定为目标字段（角色为目标）。为了使此字段成为唯一目标字段，应将所有其他字段的角色设置为输入。

图片 7-2
设置目标



- ▶ 确认已正确定义字段的测量级别。例如，大多数值为 0 和 1 的字段都可以用作标志字段，但某些字段，比如性别，作为包含两个值的名义字段会更加准确。

图片 7-3
设置测量级别



提示：要更改值（比如 0/1）相似的多个字段的属性，请单击值列标题按该列对字段进行排序，然后使用 Shift 键选择要更改的所有字段。然后可以右键单击选定项，更改所有选定字段的测量级别或其他属性。

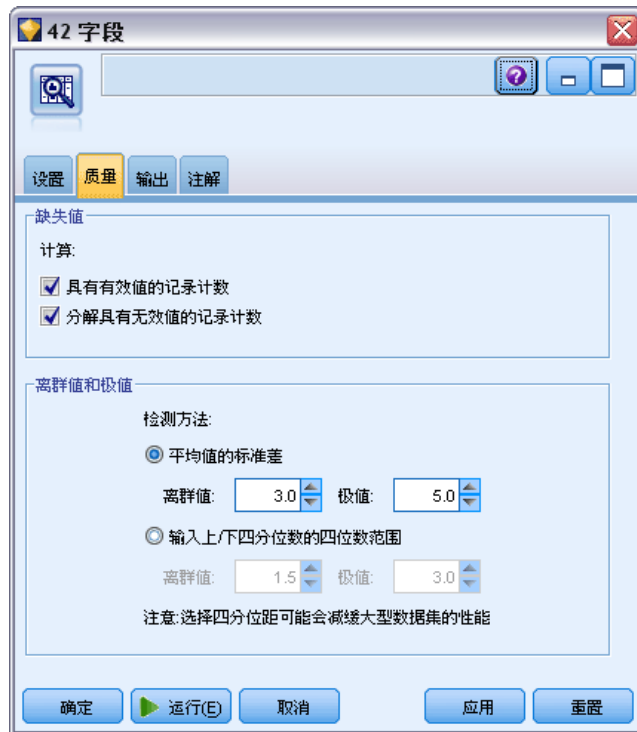
- ▶ 将数据审核节点附加到流。保留“设置”选项卡上的默认设置，以将所有字段包含在报告中。由于 churn 是类型节点中定义的唯一目标字段，系统会自动将其用作交叠字段。

图片 7-4
数据审核节点，“设置”选项卡



在“质量”选项卡上，保留检测缺失值、离群值和极值的所有默认设置，然后单击运行。

图片 7-5
数据审核节点，“质量”选项卡



浏览统计量和图表

将打开“数据审核”浏览器，其中显示各个字段的缩略图和描述性统计量。

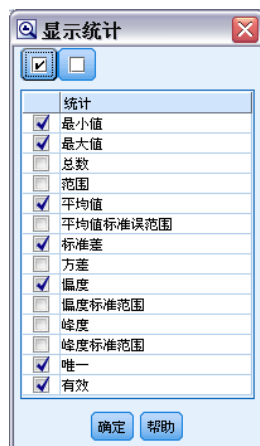
图片 7-6
数据审核浏览器



使用工具栏显示字段和值标签，并将图表从水平对齐切换为垂直对齐（仅适用于分类字段）。

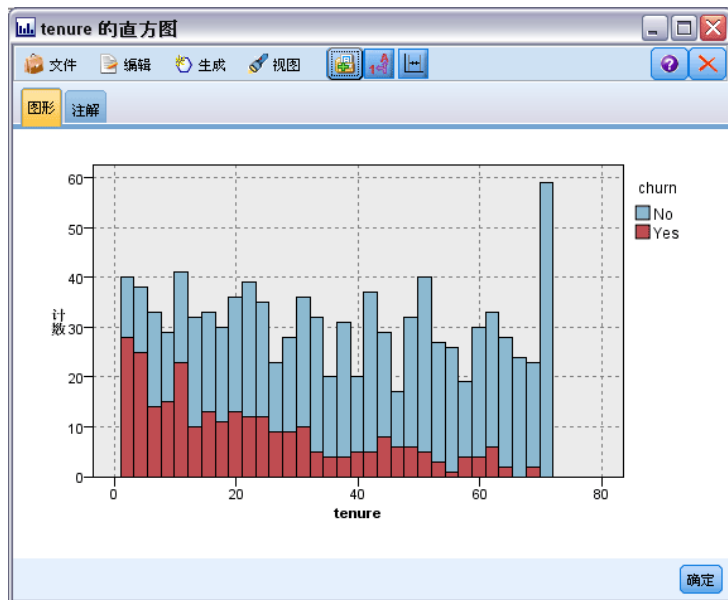
- 也可以使用工具栏或“编辑”菜单选择要显示的统计量。

图片 7-7
显示统计



双击审核报告中的任意缩略图，即可查看全尺寸图表。由于 churn 是流中的唯一目标字段，系统会将其自动用作交叠字段。可以使用图形窗口工具栏来切换字段和值标签的显示方式，也可以单击“编辑模式”按钮进一步自定义图表。

图片 7-8
工龄直方图



您也可以选择一个或多个缩略图并生成每个缩略图的图形节点。生成的节点放置在流工作区上，您可以将其添加到流以重新创建该特定图表。

图片 7-9
生成图形节点

The screenshot shows the 'Data Audit' window for 42 fields. The window is divided into several sections:

- Field List:** A list of fields including 'region', 'tenure', 'age', 'marital', 'address', and 'income'. Each field has a small thumbnail chart next to it.
- Statistics Table:** A table with columns for '最小值' (Minimum), '最大值' (Maximum), '平均值' (Average), '标准差' (Standard Deviation), '偏度' (Skewness), '唯一' (Unique), and '有效' (Valid). The data for the visible fields is as follows:

字段	最小值	最大值	平均值	标准差	偏度	唯一	有效
region		1	3	--	--	--	3
tenure		1	72	35.526	21.360	0.112	--
age		18	77	41.684	12.559	0.357	--
marital		0	1	--	--	--	2
address		0	55	11.551	10.087	1.106	--
income		9.000	1668.000	77.535	107.044	6.643	--
- Menu:** A menu is open over the 'age' field, listing various analysis options. The '图形节点' (Graphical Node) option is highlighted.
- Footer:** A status bar at the bottom indicates '* 指示多方式结果' and '* 指示取样结果'. A '确定' (OK) button is located in the bottom right corner.

处理离群值和缺失值

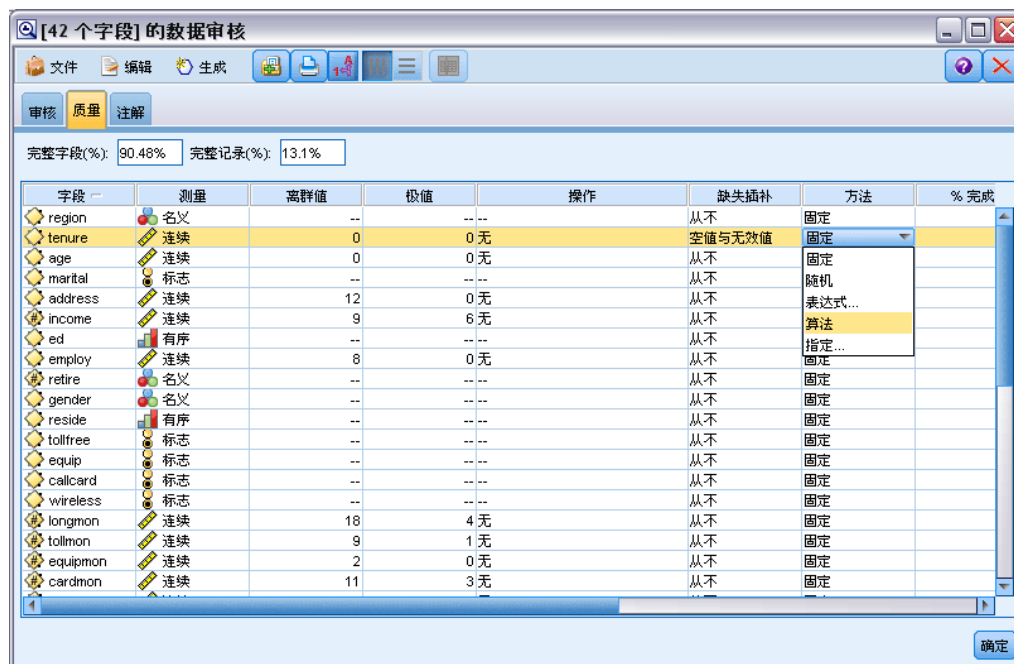
审核报告中的“质量”选项卡显示离群值、极值和缺失值的相关信息。

图片 7-10
“数据审核”浏览器，“质量”选项卡

字段	测量	离群值	极值	操作	缺失插补	方法	% 完成
region	名义	--	--	--	从不	固定	
tenure	连续	0	0 无		从不	固定	
age	连续	0	0 无		从不	固定	
marital	标志	--	--	--	从不	固定	
address	连续	12	0 无		从不	固定	
income	连续	9	6 无		从不	固定	
ed	有序	--	--	--	从不	固定	
employ	连续	8	0 无		从不	固定	
retire	名义	--	--	--	从不	固定	
gender	名义	--	--	--	从不	固定	
reside	有序	--	--	--	从不	固定	
tollfree	标志	--	--	--	从不	固定	
equip	标志	--	--	--	从不	固定	
callcard	标志	--	--	--	从不	固定	
wireless	标志	--	--	--	从不	固定	
longmon	连续	18	4 无		从不	固定	
tollmon	连续	9	1 无		从不	固定	
equipmon	连续	2	0 无		从不	固定	
cardmon	连续	11	3 无		从不	固定	

也可以指定处理这些值的方法并生成超节点，以自动应用各种变换。例如，您可以使用各种方法（包括 C&RT 算法）选择一个或多个字段并选择填补或替换这些字段的缺失值。

图片 7-11
选择归因方法



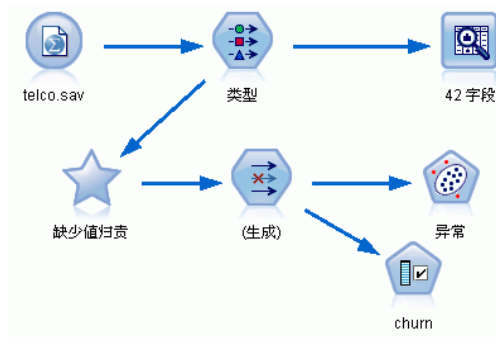
指定用于一个或多个字段的归因方法后，要生成缺失值超节点，请从菜单中选择：
生成 > 缺失值超节点

图片 7-12
生成超节点



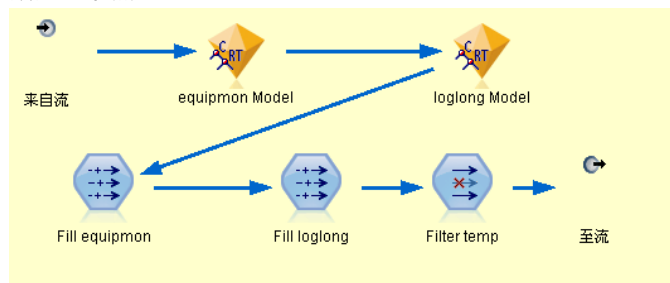
生成的超节点将添加到流工作区中，随后您可在流工作区中将超节点附加到流中，以应用各种变换。

图片 7-13
具有缺失值超节点的流



实际上，超节点包含执行所需变换的一系列节点。要了解超节点的工作方式，可编辑超节点并单击 放大。

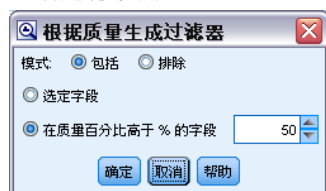
图片 7-14
放大超节点



对于使用算法归因的每个字段，将对应一个独立的 C&RT 模型，以及一个填充节点，该节点将用按模型预测的值替换空值和无效值。用户可在超节点中添加、编辑或删除特定节点，以进一步自定义节点行为。

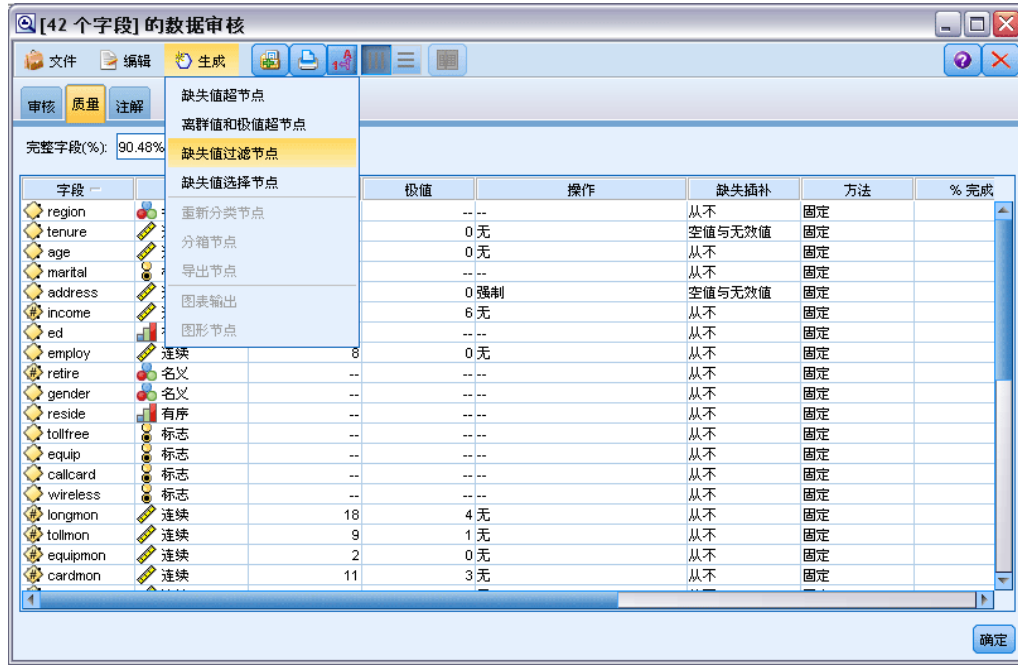
您也可以生成选择节点或过滤节点，以删除包含缺失值的字段或记录。例如，您可以过滤质量百分比低于指定阈值的任何字段。

图片 7-15
生成过滤节点



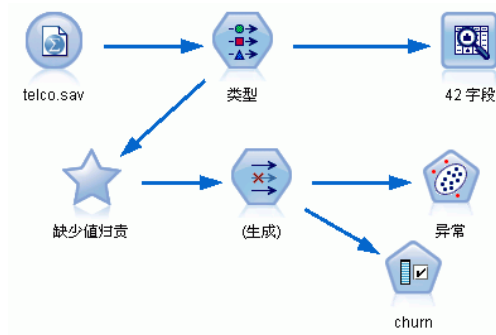
也可以用类似的方法来处理离群值和极值。指定要对各个字段执行的操作—强制、丢弃或取消—并生成超节点，以应用各种变换。

图片 7-16
生成过滤节点



审核完成并将生成的节点添加到流后，即可继续进行分析。您可能会选择使用“异常检测”、“特征选择”或其他多种方法来进一步筛选数据。

图片 7-17
具有缺失值超节点的流



药物治疗（勘察表/C5.0）

在本章中，假设您是一位正在汇总研究数据的医学研究员。您已收集了一组患有同一疾病的患者的数据。在治疗过程中，每位患者均对五种药物中的一种有明显反应。您的任务就是通过数据挖掘找出适合治疗此疾病的药物。

此示例使用名为 `druglearn.str` 的流，此流引用名为 `DRUG1n` 的数据文件。这些文件可在任何 IBM® SPSS® Modeler 安装程序的 Demos 目录中找到。此目录可通过 Windows IBM® SPSS® Modeler 程序组进行访问。文件 `druglearn.str` 位于 `streams` 目录中。

此 demo 中使用的数据字段包括：

数据字段	描述
年龄	（数值）
性别	男或女
BP	血压：高、正常或低
Cholesterol	血液中的胆固醇含量：正常或高
Na	血液中钠的浓度
K	血液中钾的浓度
Drug	对患者有效的处方药

读取文本数据

您可以使用**变量文件节点**读取定界文本数据。可以从选项板中添加变量文件节点，方法是单击**源**选项卡找到此节点，或者使用**收藏夹**选项卡（默认情况下，其中包含此节点）。然后，双击新添加的节点以打开相应的对话框。

图片 8-1
添加变量文件节点



可变文件



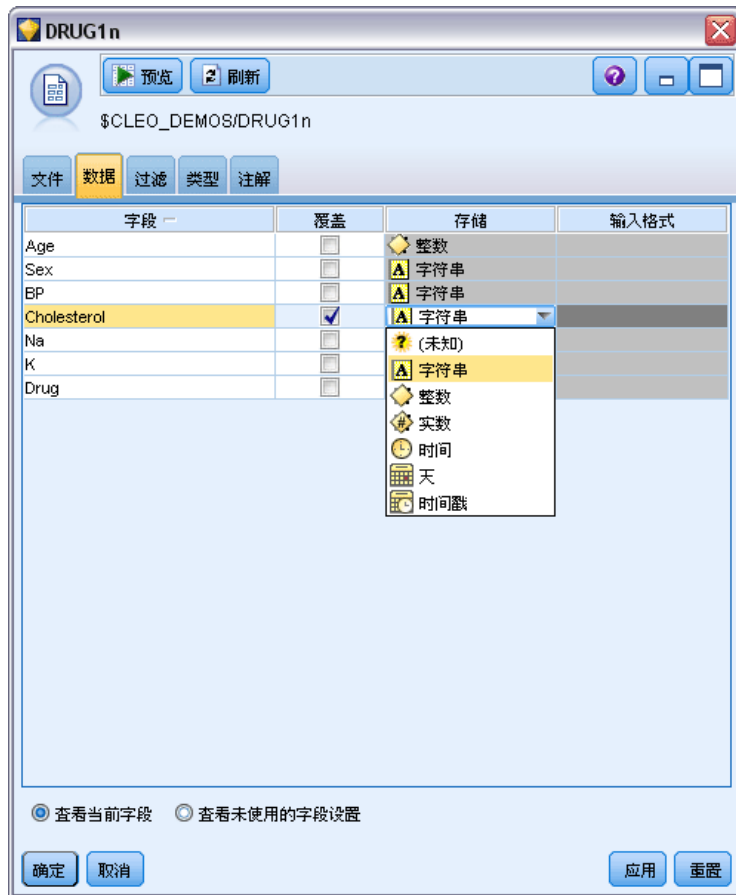
单击紧挨“文件”框右边以省略号“...”标记的按钮，浏览到您系统中的 IBM® SPSS® Modeler 安装目录。打开 Demos 目录，然后选择名为 DRUG1n 的文件。

确保选中了从文件读取字段名称，注意已载入此对话框中的字段和值。

图片 8-2
“变量文件”对话框



图片 8-3
更改字段的存储类型



图片 8-4
选择“类型”选项卡中的“值”选项



单击数据选项卡，覆盖和更改某个字段的**存储**。注意，存储不同于**测量**，即，数据字段的测量级别（或用途类型）。**类型**选项卡可帮助您了解数据中的更多字段类型。还可以选择读取值来查看各个字段的实际值，具体取决于您在值列中的选择。此过程称为**实例化**。

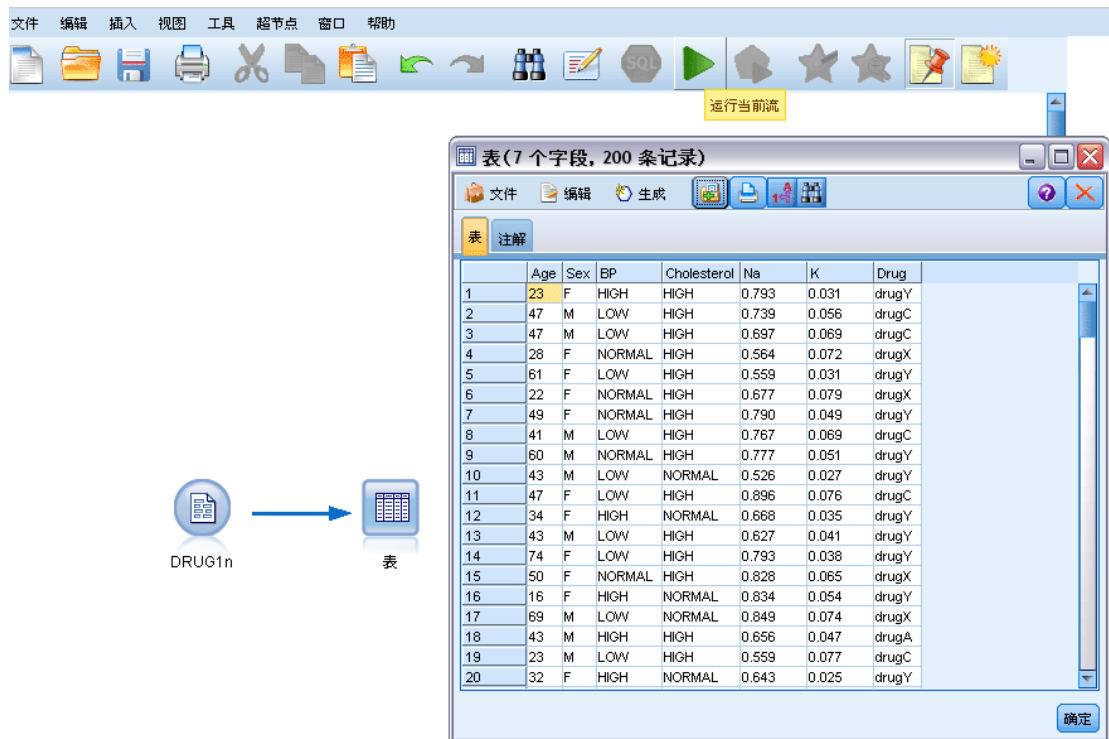
添加表

现在您已载入数据文件，可以浏览一下某些记录的值。其中一个方法就是构建一个包含“表”节点的流。要将表节点添加到流中，可双击选项板中的表节点图标或将其拖放到工作区。

图片 8-5
表节点已连接至数据源



图片 8-6
从工具栏运行流



双击选项板中的某个节点后，该节点将自动与流工作区中的选定节点相连接。此外，如果尚未连接节点，则可以使用鼠标中键将“源”节点与“表”节点相连接。要模拟鼠标中键操作，请在使用鼠标时按下 Alt 键。要查看表，请单击工具栏上的绿色箭头按钮运行流，或者右键单击“表”节点，然后选择运行。

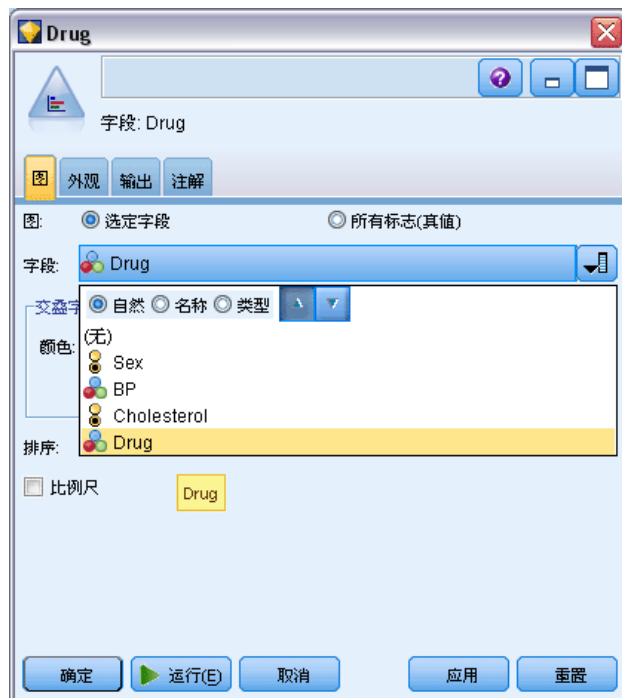
创建分布图

数据挖掘过程中，创建汇总视图通常有助于研究数据。IBM® SPSS® Modeler 提供了若干不同类型的图表供您选择，具体取决于您要汇总分析的数据类型。例如，要找出每种药物的对症患者的比例，请使用“分布”节点。

将“分布”节点添加到流，并将其与“源”节点相连接，然后双击该节点以编辑要显示的选项。

选择药品作为要显示其分布的目标字段。然后，在对话框中单击运行。

图片 8-7
选择药品作为目标字段



最终图表将有助于您查看数据的“结构”。结果表明，药品 Y 的对症患者最多，而药品 B 和药品 C 的对症患者最少。

图片 8-8
对症药品类型分布



图片 8-9
数据审核结果



此外，您还可以添加并执行“数据审核”节点，同时快速浏览所有字段的分布图和直方图。可以在“输出”选项卡中找到“数据审核”节点。

创建散点图

现在我们来了解一下有哪些因素会对药品（目标变量）产生影响。作为研究员，您一定知道钠和钾的浓度在血液中有着重大的影响。由于两者都是数值，您可以用颜色区分药品，创建一个关于钠和钾的散点图。

将“散点图”节点放在工作区中，并将其与“源”节点相连接，然后双击该节点对其进行编辑。

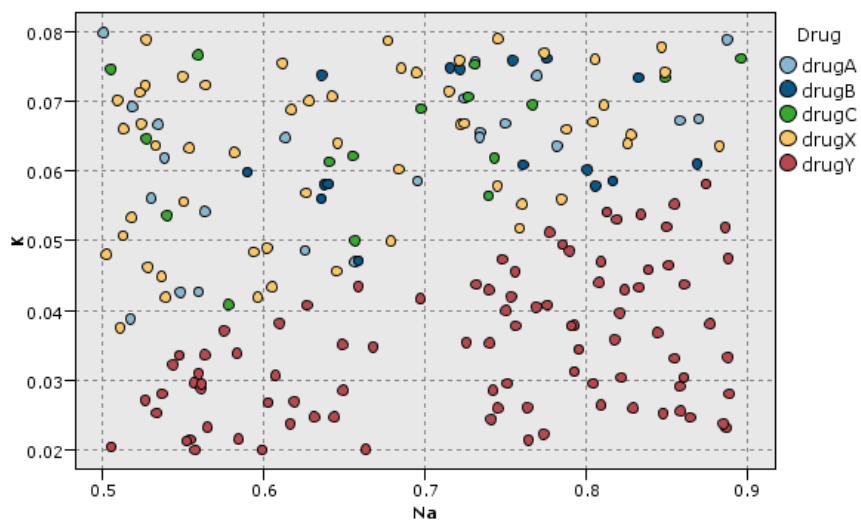
在“散点图”选项卡中，选择 Na 作为 X 字段，选择 K 作为 Y 字段，并选择 Drug 作为交叠字段。然后，单击运行。

图片 8-10
创建散点图



此散点图清楚地显示了一个阈值，在此阈值上方，对症药品始终是 Y，在此阈值下方，对症药品均不是 Y。此阈值等于钠 (Na) 和钾 (K) 的比。

图片 8-11
药品分布散点图

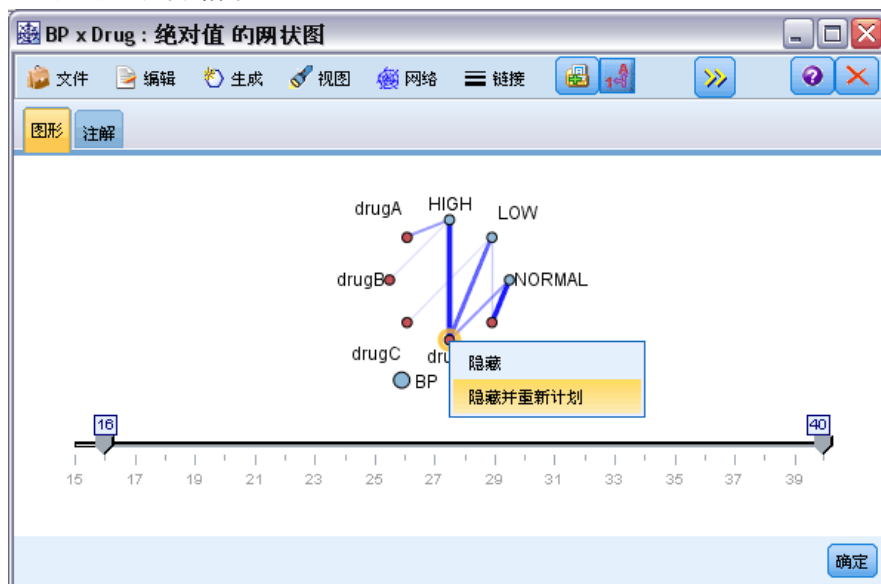


创建网络图

因为很多数据字段均可分类，您也可尝试绘制网络图，此图表将反映不同类别之间的联系。首先，将网络节点与您工作区中的源节点相连接。在“网络节点”对话框中，选择 BP（血压）和药品。然后，单击运行。

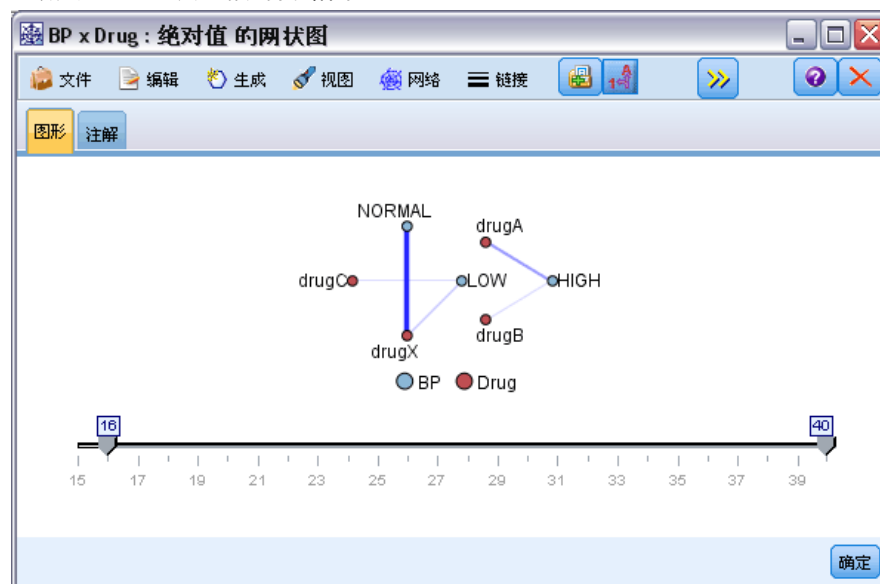
此图显示，药品 Y 与三种级别的血压均相关。这并不奇怪，因为您早已看出 Y 是最佳药品。要关注其他药品，您可隐藏药品 Y。在视图菜单上，选择编辑模式，然后右键单击药品 Y 点并选择隐藏并重新计划。

图片 8-12
药品和血压网络图



简图中隐藏了药品 Y 及其所有链接。现在您可以清楚地看到，只有药品 A 和 B 与高血压有关。只有药品 C 和 X 与低血压有关。而药品 X 与正常血压有关。此时，您仍然无法在药品 A 与 B 或药品 C 与 X 之间为指定患者作出选择。此时 建模可以助您一臂之力。

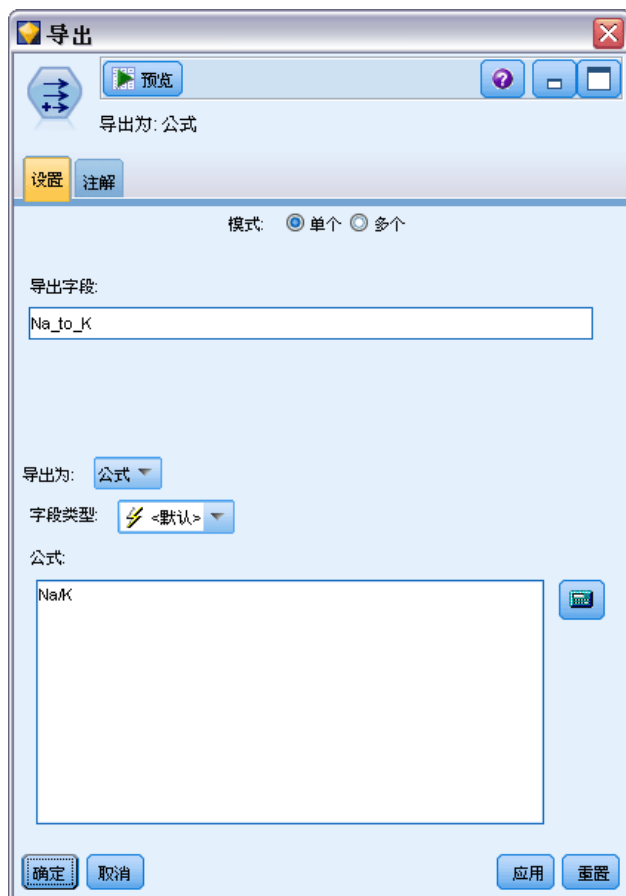
图片 8-13
隐藏药品 Y 及其链接的网络图



导出新字段

由于钠与钾的比似乎可以用来预测何时可以使用药品 Y，因此您可以为每条记录导出一个包含此比值的字段。该字段稍后可用于构建模型以预测何时可使用五种药品中的每一种药品。为了简化流布局，首先删除除 DRUG1n 源节点之外的所有节点。将派生节点（字段选项选项卡）附加到 DRUG1n，然后双击派生节点进行编辑。

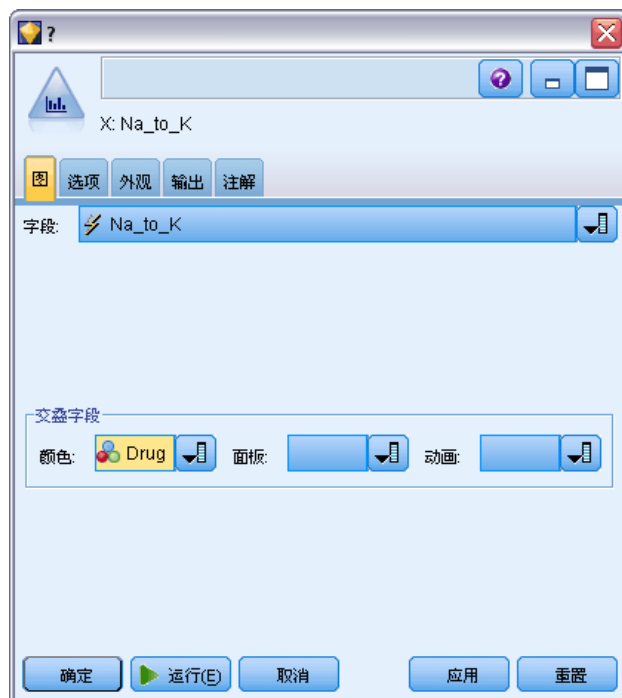
图片 8-14
编辑导出节点



将新字段命名为 Na_to_K。由于是通过将钠值除以钾值获取新字段，所以请在公式中输入 Na/K。您还可通过单击紧挨该字段右侧的图标来创建公式。此操作将打开“表达式构建器”，这是一种使用函数、操作数、字段及其字段值的内置列表交互式创建表达式的方式。

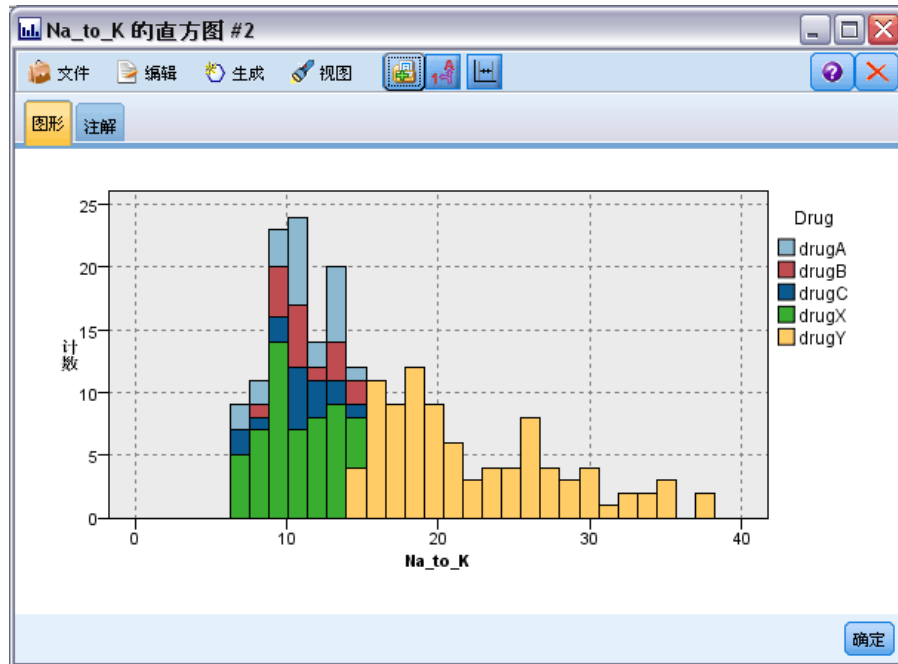
您可以通过将直方图节点添加到导出节点来检查新字段的分布情况。在直方图节点对话框中，将 Na_to_K 指定为要绘制的字段，并将药品指定为交叠字段。

图片 8-15
编辑直方图节点



运行流时，将在此显示图表。您可以根据显示结果得出以下结论：当 Na_to_K 字段的值等于或大于 15 时，应选择药品 Y。

图片 8-16
直方图显示



构建模型

通过研究和操作数据，您能够得出某些假设结论。血液中钠与钾的比例以及血压似乎都会影响药品的选择。但您还不能完全解释清楚所有关系。此时似乎可以通过建模找出某些答案。此种情况下，您可以尝试使用规则构建模型（C5.0）来拟合数据。

由于使用的是导出字段 Na_to_K，您可以过滤掉原始字段 Na 和 K，以避免在建模算法中重复操作。上述操作可通过过滤节点完成。

图片 8-17
编辑过滤节点



在“过滤”选项卡上，单击 Na 和 K 旁边的箭头。如果箭头上显示红色的 X，则表示该字段已被过滤。

然后，附加一个连接到“过滤”节点的“类型”节点。“类型”节点允许您指出要使用的字段类型以及如何使用这些字段预测结果。

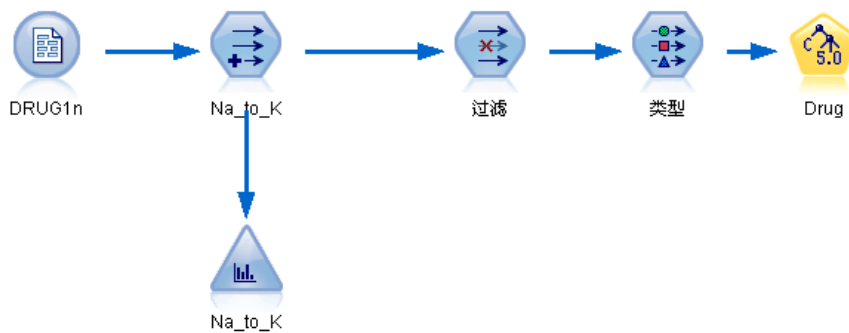
在“类型”选项卡上，将药品字段的角色设置为目标，表明您要预测该药品字段。将其他字段的角色设置为输入，表示这些字段将用作预测变量。

图片 8-18
编辑类型节点



要评估此模型，请将节点 C5.0 置于工作空间，然后将此节点附加到流的末端（如图所示）。单击绿色运行工具栏按钮运行流。

图片 8-19
添加 C5.0 节点



浏览模型

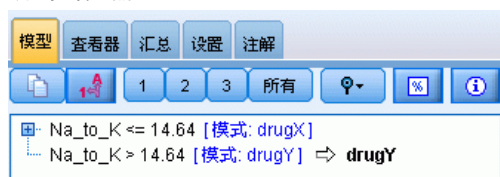
执行 C5.0 节点时，模型块被添加到流，同时添加到位于窗口右上角的“模型”选项板。要浏览模型，右键单击任一图标并从上下文菜单选择编辑或浏览。

图片 8-20
浏览模型



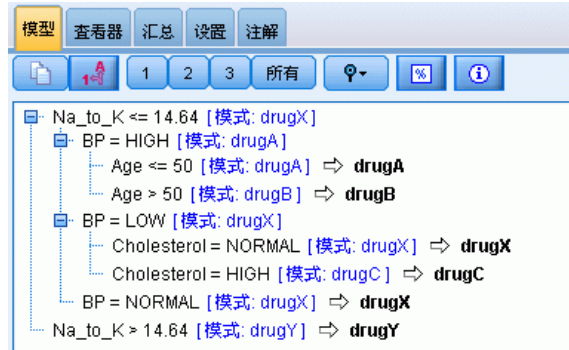
“规则”浏览器以决策树形式显示 C5.0 节点所生成的规则集。最初，决策树处于折叠状态。要展开决策树，请单击所有按钮显示所有层。

图片 8-21
规则浏览器



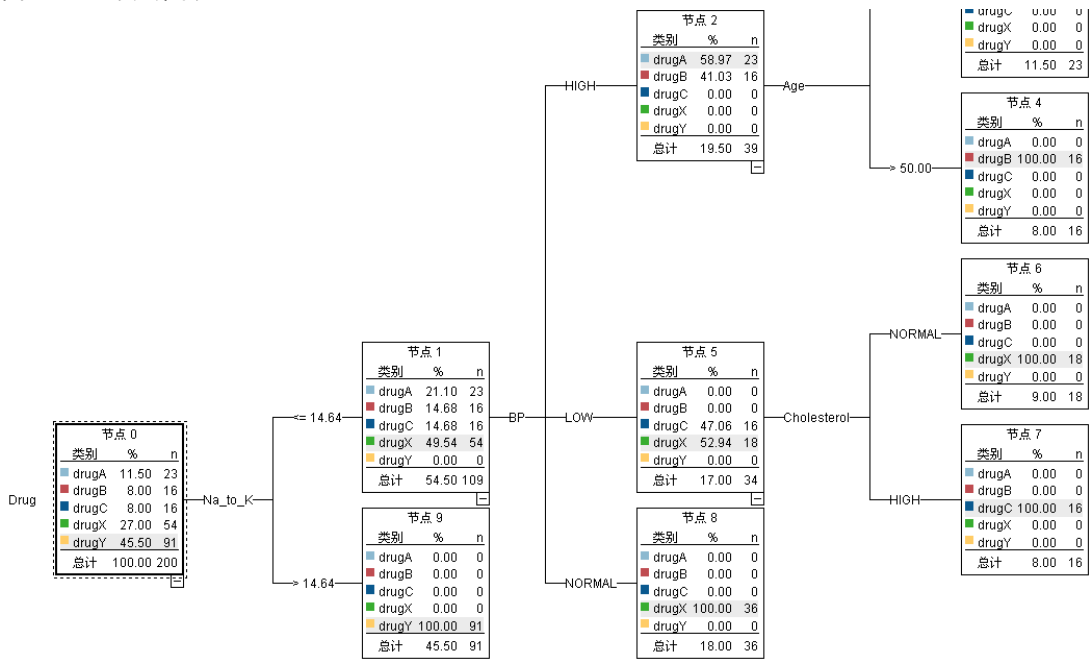
谜团将因此而解开。对于 Na 与 K 的比小于 14.64 的高血压患者，年龄将决定如何选择药品。对于低血压患者，胆固醇含量似乎是最有力的预测变量。

图片 8-22
完全展开的规则浏览器



通过单击查看器选项卡，还可以更复杂的图表形式查看同一决策树。通过此图表形式，您可以更轻松地查看各个血压类别的观测值数量以及各个观测值的百分比。

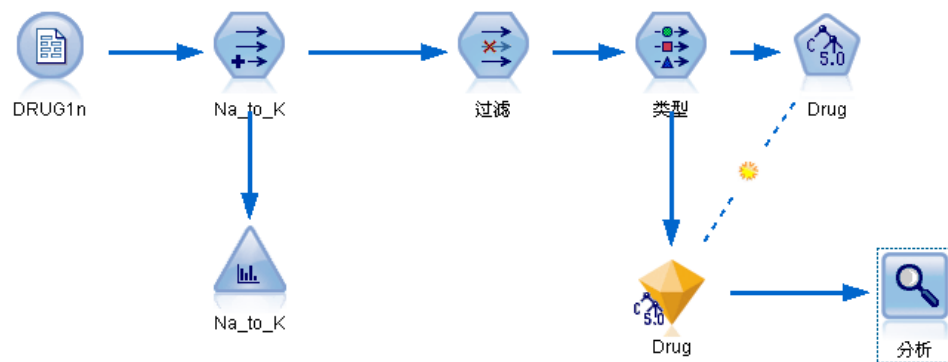
图片 8-23
图形形式的决策树



使用分析节点

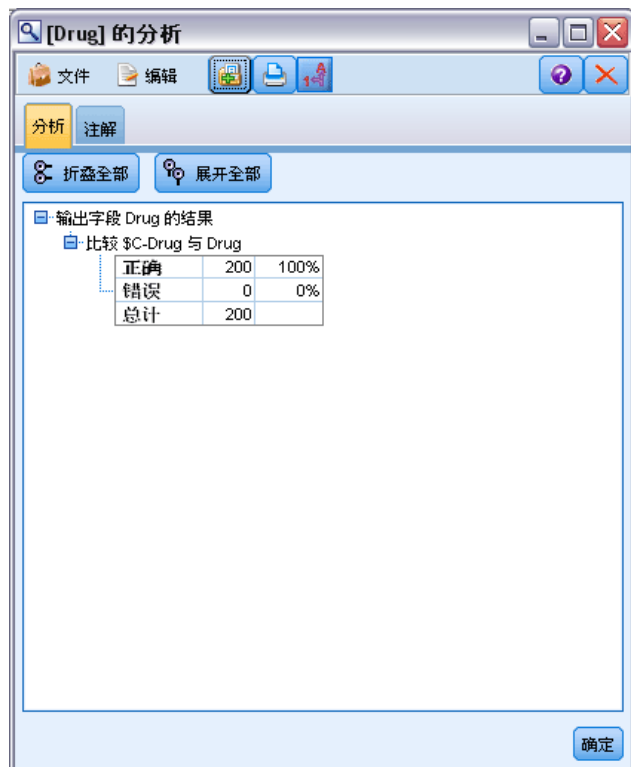
可以使用分析节点评估模型的精确度。将分析节点（从输出节点选项板）附加到模型块，打开分析节点并单击运行。

图片 8-24
添加分析节点



分析节点输出显示：使用该假设数据集，此模型已正确预测该数据集中每个记录的药品选择。在真正的数据集中，未必能做到完全准确，但分析节点可帮您确定模型的精确度能否满足特殊使用要求。

图片 8-25
分析节点输出



筛选预测变量（特征选择）

特征选择节点有助于识别用于预测特定结果的最重要的字段。特征选择节点可对成百乃至上千个预测变量进行筛选、排序，并选择出可能是最重要的预测变量。最后，会生成一个执行地更快且更加有效的模型—此模型使用较少的预测变量，执行地更快且更易于理解。

本示例中使用的数据由某虚构电话公司的数据仓库提供，且包含有关该公司的 5000 名客户对特定促销活动的响应的信息。该数据包含大量的字段，其中有客户年龄、职业、收入、电话使用情况等统计量。其中有三个“目标”字段，显示客户是否响应这三种促销。该公司想利用这些数据来预测哪些客户最可能在将来对类似报价做出响应。

此示例使用名为 featureselection.str 的流，此流引用名为 customer_dbase.sav 的数据文件。这些文件可在任何 IBM® SPSS® Modeler 安装程序的 Demos 目录中找到。此目录可通过 Windows IBM® SPSS® Modeler 程序组进行访问。文件 featureselection.str 位于 streams 目录下。

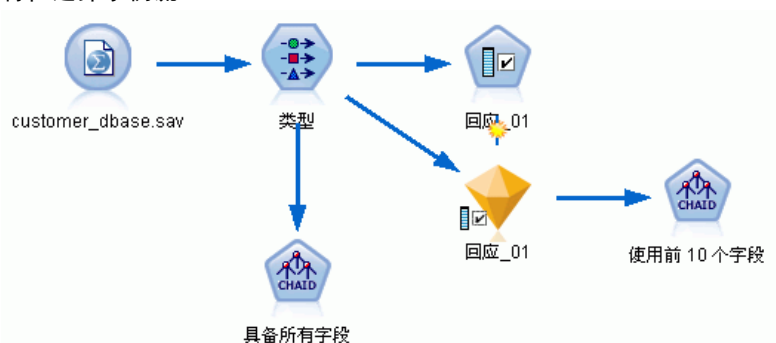
本示例仅关注其中一种促销活动，并将其作为目标。本示例使用 CHAID 树构建节点来开发模型，用以说明最有可能响应促销活动的客户。其中对以下两种方法作了对比：

- 不使用特征选择。数据集中的所有预测变量字段均可用作 CHAID 树的输入。
- 使用特征选择。使用特征选择节点选择最佳的 10 个预测变量。然后将其输入到 CHAID 树中。

通过比较两个生成的树模型，可以看到特征选择如何产生有效的结果。

构建流

图片 9-1
特征选择示例流



- ▶ 在一个空流工作区中，放置一个 Statistics 文件源节点。将此节点指向示例数据文件 customer_dbase.sav，该文件位于 IBM® SPSS® Modeler 安装程序的 Demos 目录下。（或者，可打开位于 streams 目录下的示例流文件 featureselection.str。）
- ▶ 添加类型节点。在“类型”选项卡上，向下滚动到底部并将 response_01 的角色更改为目标。将其他响应字段（response_02）和（response_03）以及客户 ID（列表

顶部的 custid) 的角色更改为无。将所有其他字段的角色设置为输入，并单击读取值按钮，然后单击确定。

- ▶ 为流添加“特征选择”建模节点。在此节点上，您可以指定要筛选的规则和标准，或要筛选的字段。
- ▶ 运行流以创建特征选择模型块。
- ▶ 右键单击流上或“模型”选项板中的模型块并选择编辑或浏览以查看结果。

图片 9-2
特征选择模型块中的“模型”选项卡



顶部面板显示了所找到的对预测非常有用的字段。这些字段基于重要性排序。底部面板显示了从分析中筛选出来的字段及筛选的原因。通过检查顶部面板中的字段，可以确定在随后的建模会话中要使用哪些字段。

- ▶ 现在，可以选择要在下游使用的字段。虽然最初已将 34 个字段识别为重要字段，但我们希望进一步减少预测变量集合的数目。
- ▶ 通过使用第一列上的复选标记仅选中前 10 个预测变量，可取消选择不需要的预测变量。（单击行 11 中的选中标记，按住 Shift 键并单击行 34 中的选中标记。）关闭模型块。
- ▶ 要在未选中特征的情况下比较结果，则必须向流中添加以下两个 CHAID 建模节点：一个模型使用特征选择，另一个模型不使用特征选择。
- ▶ 将一个 CHAID 节点连接到类型节点，将另一个连接到特征选择模型块。
- ▶ 打开每个 CHAID 节点，选择“构建选项”选项卡，确保在“目标”窗格中选中了选项构建新模型、构建单个树和启动交互会话。

在“基本”窗格上，确保将最大树深度设置为 5。

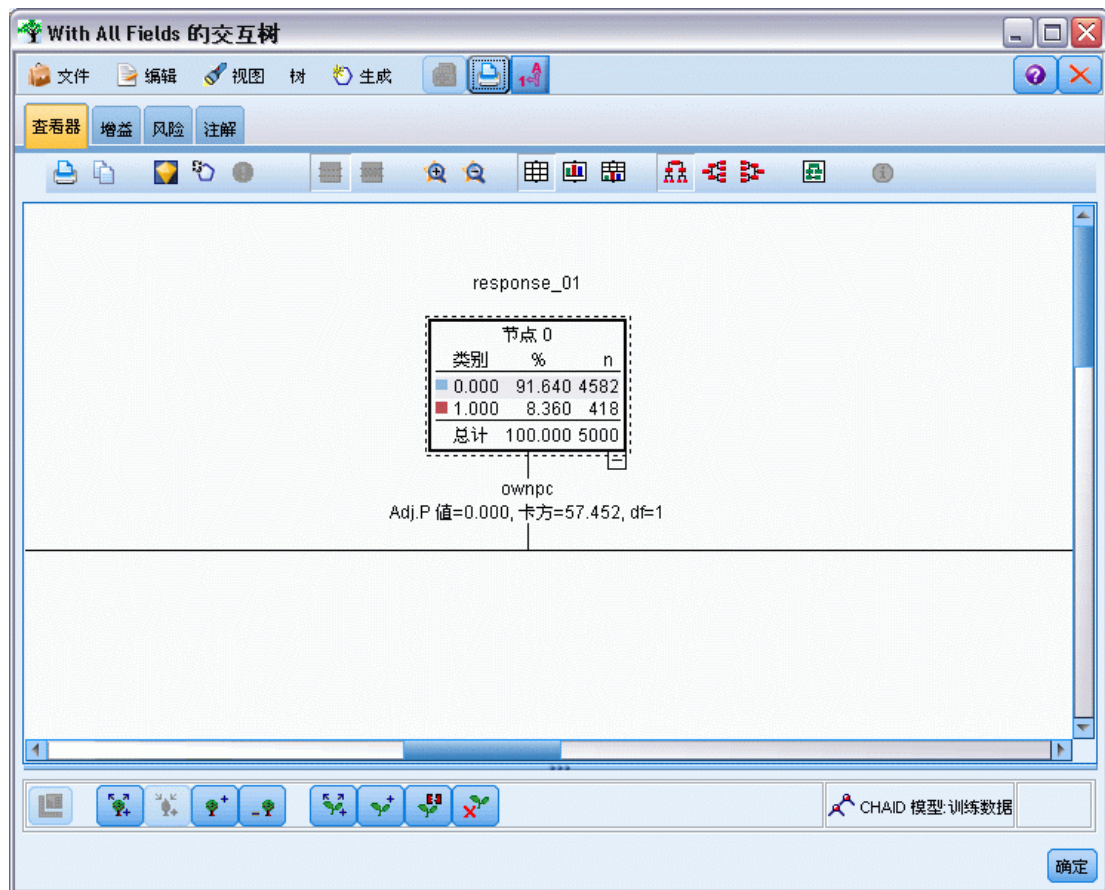
图片 9-3
所有预测变量字段的 CHAID 建模节点的“目标”设置



构建模型

- ▶ 执行使用数据集中所有预测变量的 CHAID 节点（即连接到“类型”节点的节点）。当节点执行时，注意观察执行节点所用的时间。表会显示在结果窗口中。
- ▶ 从菜单中，选择树 > 生长树，可生成并显示展开的树。

图片 9-4
在树构建器中生成树



- ▶ 现在，对另一个 CHAID 节点（此节点仅使用 10 个预测变量）执行相同的操作。打开树构建器后再次生成树。

第二个模型的执行速度应比第一个模型快。因为此数据集相当小，所以在执行时间上差别可能只有几秒钟；但对于更大的实际应用中的数据集，此差别可能非常明显（几分钟甚至几小时）。使用特征选择可以显著加快处理速度。

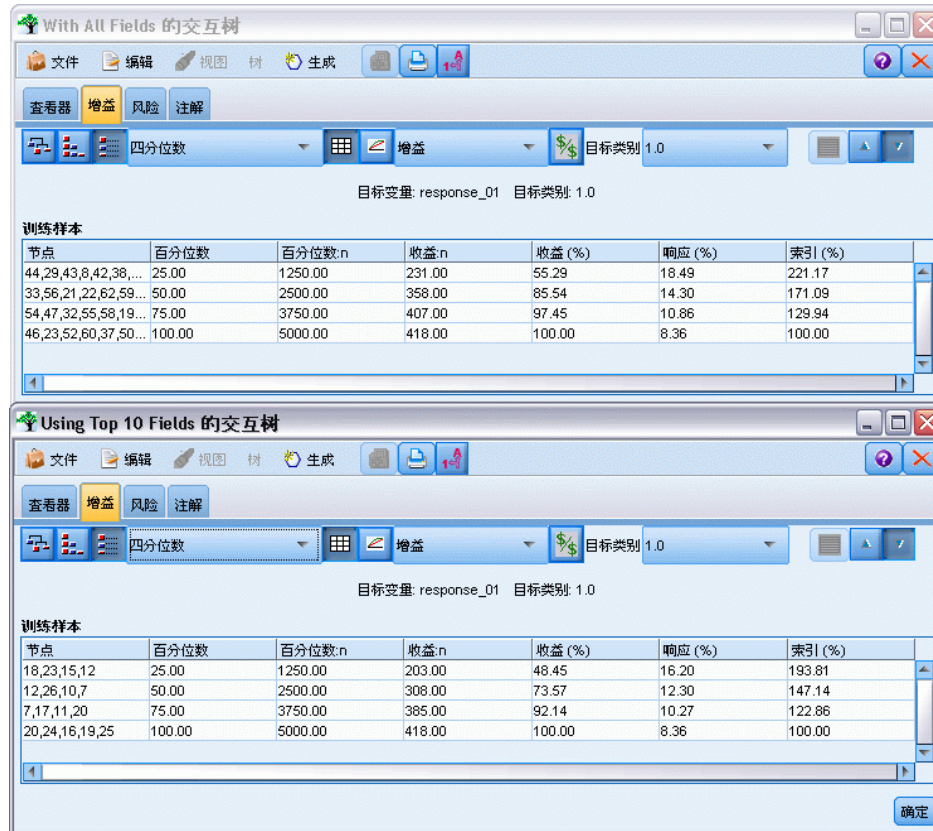
第二个树比第一个树包含的树节点也要少。因此更易于理解。但在决定使用此模型之前，需要查明此模型是否有效，并查明其与使用所有预测变量的模型相比较的结果。

比较结果

要比较两个结果，需要进行有效性测量。为此，可在树构建器中使用“收益”选项卡。我们将查看**提升**，该图可测量节点中的记录与数据集中的所有记录相比时，其落入目标类别的可能性究竟提升多少。例如，提升值 148% 表示与数据集中的所有记录相比，节点中的记录落在目标类别的可能性是其 1.48 倍。提升值显示在“收益”选项卡的指数列中。

- ▶ 在全部预测变量集合的树构建器中，单击“收益”选项卡。将“目标类别”更改为 1.0。首先单击“分位数”工具栏按钮将显示更改为四分位数。然后从此按钮右侧的下拉列表中选择四分位数。
- ▶ 在具有 10 个预测变量集合的树构建器中重复此步骤，就可以对两个类似的收益表进行比较，如下图中所示。

图片 9-5
两个 CHAID 模型的收益图表



在每个收益表中，都将其树的终端节点分组为四分位数。要比较两个模型的有效性，可查看每个表中 25% 分位数的提升（指数值）。

包括所有预测变量时，模型显示提升值 221%。即，具有这些节点中的特征的客户，其响应目标促销活动的可能性是其他客户的 2.2 倍。要查看这些具体特征，可单击以选定顶部的行。然后切换到“查看器”选项卡，其中相应的节点正以黑色框突出显示。沿树往下寻找每个突出显示的终端节点以查明这些预测变量是如何分割的。25%

分位数自身包括 10 个节点。如果转换为实际应用中的评分模型，则客户的 10 个不同特征是很难管理的！

如果仅包括前 10 个预测变量（由特征选择识别），则提升值接近为 194%。虽然此模型不如使用所有预测变量的模型那样有效，但它无疑也是有用的。并且此模型 25% 分位数仅包括四个节点，因此它更简单。因此，我们可以确定特征选择模型比使用所有预测变量的模型更优越。

摘要

让我们来检查特征选择的优势。使用较少的预测变量会降低成本。这意味着要收集、处理和输入模型的数据减少。并且节省了计算时间。在本示例中，即使增加了额外的特征选择步骤，但因具有较小的预测变量集合，模型构建的速度也明显提高。如果使用较大的实际数据集，则节省的时间应大大增加。

使用数目较少的预测变量会使评分更加简单。如示例所示，可能仅需识别有可能响应促销活动的客户的 4 个特征。请注意，如果预测变量数越多，则过度拟合模型的风险越大。生成的模型越简单，则对其他数据集会越有利（尽管可能需要通过测试确定该模型）。

您可能已经使用树构建算法来进行特征选择，使树能够识别最重要的预测变量。实际上，会经常使用 CHAID 算法来完成此操作，而且使用此算法还可以逐层生成树以控制树的深度和复杂性。但是，使用特征选择节点会更快更容易。此节点可通过一次性步骤迅速对所有预测变量进行排序，使您能够迅速识别最重要的字段。使用此节点还可以更改要包括的预测变量数。可以使用前 15 个或 20 个而不是 10 个预测变量再次轻松运行此示例，并比较其结果以确定最佳模型。

减少输入数据字符串长度（重新分类节点）

减少输入数据字符串长度（重新分类）

对于二项 logistic 回归模型和包含二项 logistic 回归模型的自动分类器模型，字符串字段被限制为最多不得超过八个字符。如果字符串超过八个字符，则可以使用重新分类节点对其重新编码。

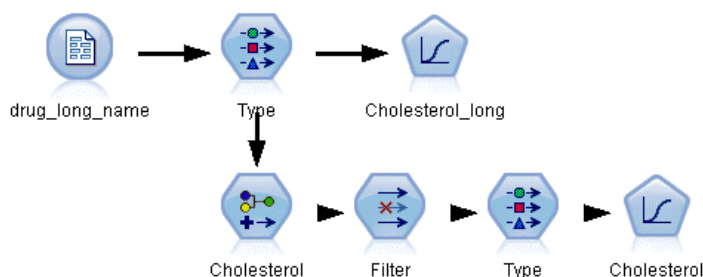
本示例使用名为 `reclassify_strings.str` 的流，该流所引用的数据文件名为 `drug_long_name`。这些文件可在任何 IBM® SPSS® Modeler 安装程序的 Demos 目录中找到。此目录可通过 Windows IBM® SPSS® Modeler 程序组进行访问。文件 `reclassify_strings.str` 位于 `streams` 目录下。

本示例主要集中在流中的一小部分，以显示可能因字符串过长而生成的某种错误；本示例还解释了如何使用重新分类节点更改字符串的详细信息，使其符合长度要求。尽管示例中使用了二项 Logistic 回归节点，但这相当于使用自动分类器节点来生成二项 Logistic 回归模型。

重新分类数据

- ▶ 使用一个变量文件源节点，连接到 Demos 文件夹下的数据集 `drug_long_name`。

图片 10-1
显示了对二项 logistic 回归的字符串重新分类的样本流



- ▶ 将类型节点添加至源节点，然后选择 `Cholesterol_long` 作为目标。
- ▶ 将 Logistic 回归节点添加到类型节点中。

- ▶ 在 Logistic 回归节点上，单击“模型”选项卡并选择二项过程。

图片 10-2

“Cholesterol_long” 字段中的长字符串详细信息

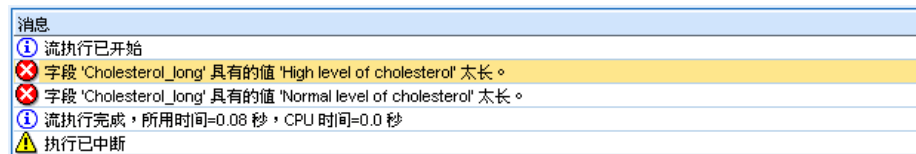


- ▶ 在 reclassify_strings.str 中执行 Logistic 回归节点时，会显示一个错误消息，警告您 Cholesterol_long 字符串值过长。

如果遇到此类型的错误消息，请按照本示例其他部分所说明的步骤来修改数据。

图片 10-3

执行二项 logistic 回归节点时所显示的错误消息



- ▶ 为类型节点添加一个重新分类节点。
- ▶ 在“重新分类”字段中，选择 Cholesterol_long。
- ▶ 键入 Cholesterol 作为新的字段名称。
- ▶ 单击获取按钮，将 Cholesterol_long 值添加至原始值列。

- ▶ 在新的值列中，在高胆固醇水平原始值旁边，键入高，在正常胆固醇水平原始值的旁边，键入正常。

图片 10-4
重新分类长字符串



- ▶ 为重新分类节点添加一个过滤节点。

- ▶ 在“过滤”列中，单击以删除 Cholesterol_long。

图片 10-5
过滤数据中的“Cholesterol_long”字段



- ▶ 将类型节点添加至过滤节点并选择 Cholesterol 作为目标。

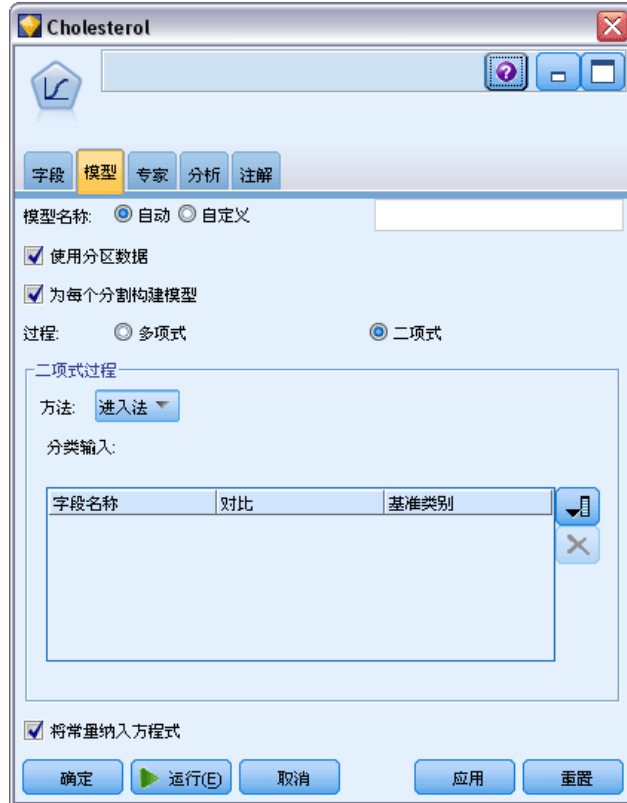
图片 10-6
“Cholesterol” 字段中的短字符串详细信息



- ▶ 将一个 Logistic 节点添加到类型节点中。
- ▶ 在 Logistic 节点上，单击“模型”选项卡并选择二项过程。

- 您现在可以执行二项 Logistic 节点，并生成一个不会显示错误消息的模型。

图片 10-7
选择“二项”作为程序



该示例仅显示了部分流。如果您需要更多有关该类型流（可能需要重新分类长字符串）的信息，请参阅以下示例：

- 自动分类器节点。有关详细信息，请参阅第 40 页码第 4 章中的对客户响应建模（自动分类器）。
- 二项 logistic 回归节点。有关详细信息，请参阅第 152 页码第 13 章中的电信客户流失（二项 Logistic 回归）。

有关如何使用 IBM® SPSS® Modeler（如《用户指南》、《节点参考》和《算法指南》）的详细信息，可从安装光盘的 \Documentation 目录下找到。

部分 III: 建模示例

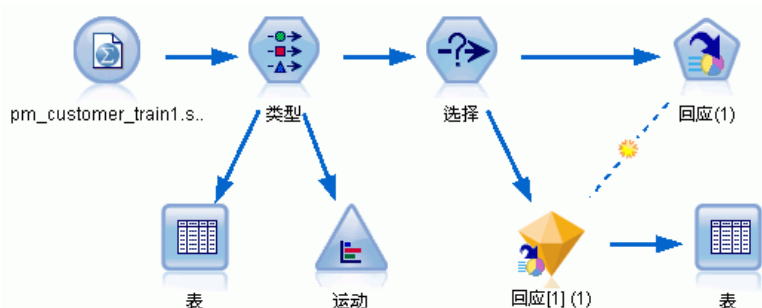
对客户响应建模（决策列表）

决策列表算法可生成表示给定的二元结果（是/否）的可能性上限和下限规则。决策列表模型在客户关系管理（例如呼叫中心或市场营销）中得到了广泛的应用。

本示例以某金融公司为例，该公司希望通过为每个客户提供最适合的报价以在未来的商业竞争中取得更大收益。特别地，该示例可使用决策列表模型，根据以前的促销活动，识别积极响应当前活动的客户的特征，并根据识别结果生成邮件发送清单。

决策列表模型尤其适用于交互建模，可允许您调整模型中的参数并立即看到结果。对于允许您自动创建多个不同模型并对结果进行排序的其他方法，可以改用自动分类器节点。

图片 11-1
决策列表示例流

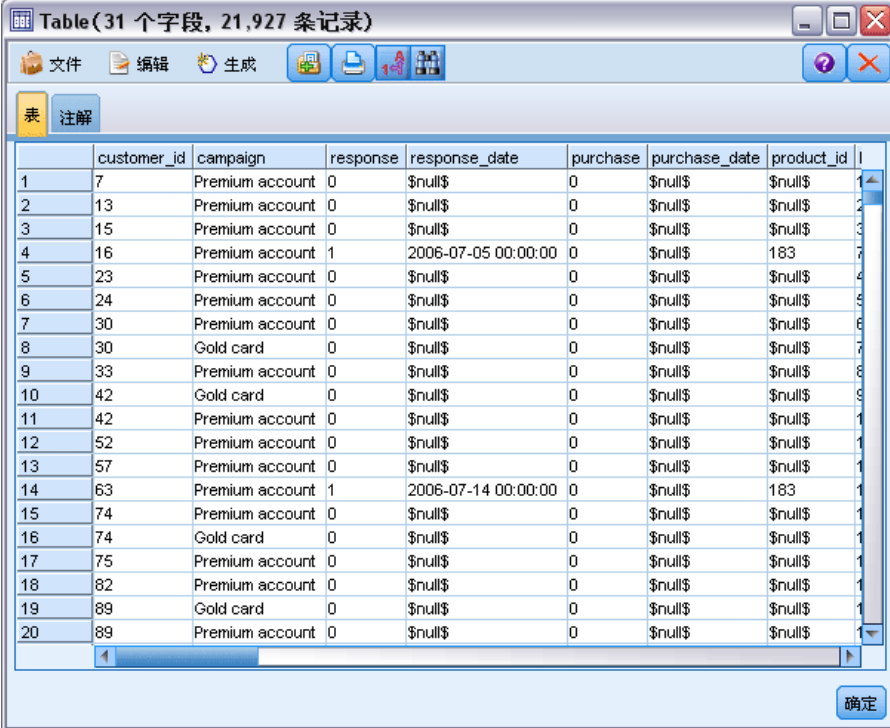


本示例使用流 `pm_decisionlist.str`，该流引用数据文件 `pm_customer_train1.sav`。这些文件可在任何 IBM® SPSS® Modeler 安装程序的 Demos 目录中找到。此目录可通过 Windows IBM® SPSS® Modeler 程序组进行访问。文件 `pm_decisionlist.str` 位于 `streams` 目录下。

历史数据

文件 pm_customer_train1.sav 的历史数据可跟踪过去的营销活动中为特定客户提供的报价，由 campaign 字段的值表示。Premium account 活动中的记录数最大。

图片 11-2
以前促销活动的相关数据



The screenshot shows a window titled "Table (31 个字段, 21,927 条记录)". The table contains the following data:

	customer_id	campaign	response	response_date	purchase	purchase_date	product_id
1	7	Premium account	0	\$null\$	0	\$null\$	\$null\$
2	13	Premium account	0	\$null\$	0	\$null\$	\$null\$
3	15	Premium account	0	\$null\$	0	\$null\$	\$null\$
4	16	Premium account	1	2006-07-05 00:00:00	0	\$null\$	183
5	23	Premium account	0	\$null\$	0	\$null\$	\$null\$
6	24	Premium account	0	\$null\$	0	\$null\$	\$null\$
7	30	Premium account	0	\$null\$	0	\$null\$	\$null\$
8	30	Gold card	0	\$null\$	0	\$null\$	\$null\$
9	33	Premium account	0	\$null\$	0	\$null\$	\$null\$
10	42	Gold card	0	\$null\$	0	\$null\$	\$null\$
11	42	Premium account	0	\$null\$	0	\$null\$	\$null\$
12	52	Premium account	0	\$null\$	0	\$null\$	\$null\$
13	57	Premium account	0	\$null\$	0	\$null\$	\$null\$
14	63	Premium account	1	2006-07-14 00:00:00	0	\$null\$	183
15	74	Premium account	0	\$null\$	0	\$null\$	\$null\$
16	74	Gold card	0	\$null\$	0	\$null\$	\$null\$
17	75	Premium account	0	\$null\$	0	\$null\$	\$null\$
18	82	Premium account	0	\$null\$	0	\$null\$	\$null\$
19	89	Gold card	0	\$null\$	0	\$null\$	\$null\$
20	89	Premium account	0	\$null\$	0	\$null\$	\$null\$

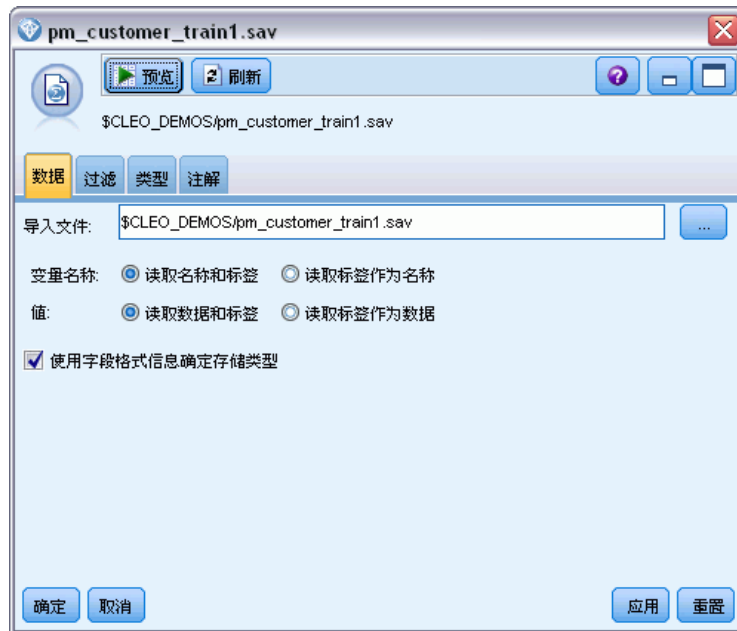
campaign 字段的值在数据中实际编码为整数，并带有类型节点中定义的标签（例如，2 = Premium account）。可以使用工具栏切换表中值标签的显示。

该文件还包括若干包含每位客户的相关人口统计和金融信息的字段，这些字段可用于构建或“训练”依据特定特征针对不同组预测响应率的模型。

构建流

- ▶ 添加指向 pm_customer_train1.sav 的 Statistics 文件节点，该文件位于 IBM® SPSS® Modeler 安装程序的 Demos 文件夹中。（您可以在文件路径中指定 \$CLEO_DEMOS/ 作为引用此文件夹的快捷方式。）

图片 11-3
读取数据



- ▶ 添加类型节点，然后选择响应作为目标字段（“角色”为目标）。将此字段的测量级别设置为标志。

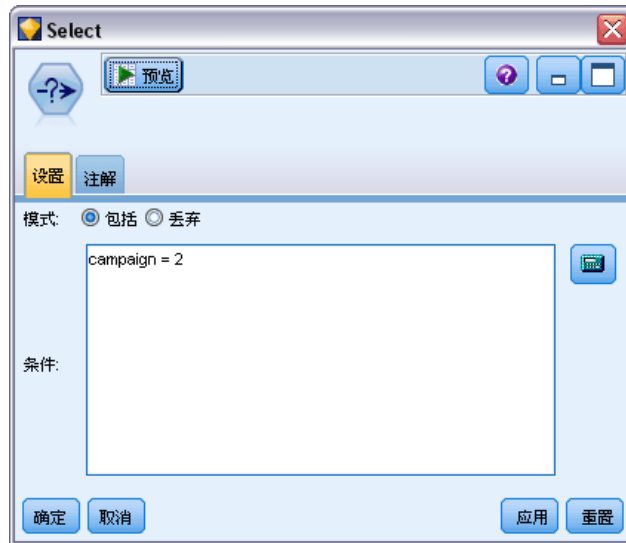
图片 11-4
设置测量级别和角色



- ▶ 对于以下字段，应将角色设置为无：customer_id、campaign、response_date、purchase、purchase_date、product_id、Rowid 和 X_random。这些字段在数据中均有用途，但不会在实际模型的构建中使用。
- ▶ 单击类型节点的读取值按钮以确保值获得实例化。

尽管数据包含有关四项不同活动的信息，但每一次的分析应集中关注其中一项活动。由于 Premium 活动（在数据中编码为 `campaign=2`）中的记录数最大，因此可以使用选择节点实现仅在流中包含这些记录。

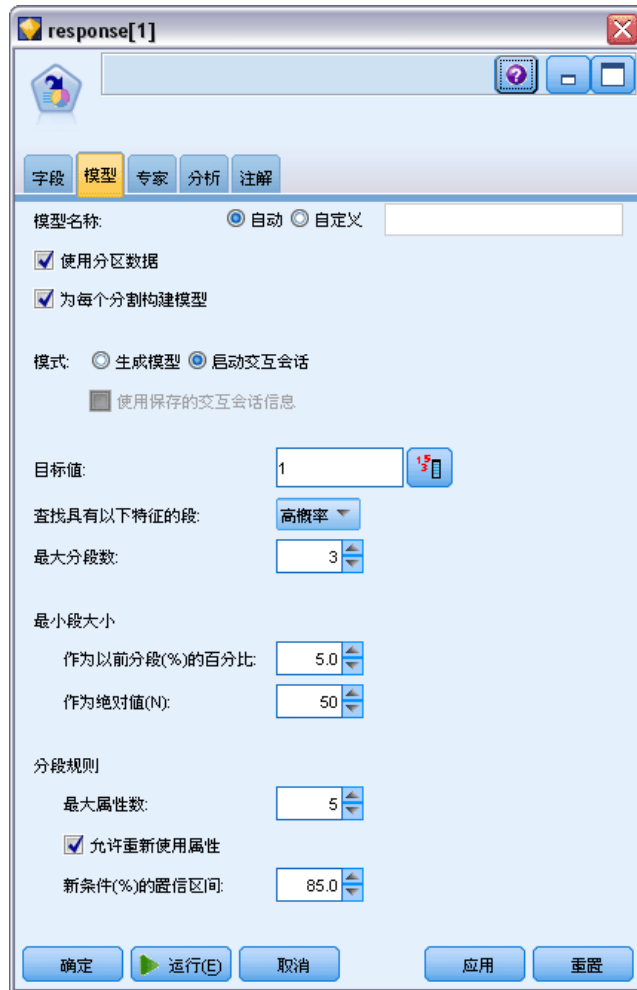
图片 11-5
为单项活动选择记录



创建模型

- ▶ 在流中添加“决策列表”节点。在“模型”选项卡中，将目标值设为 1 以表示要搜索的结果。在这种情况下，您正在搜索对以前的报价发出是响应的客户。

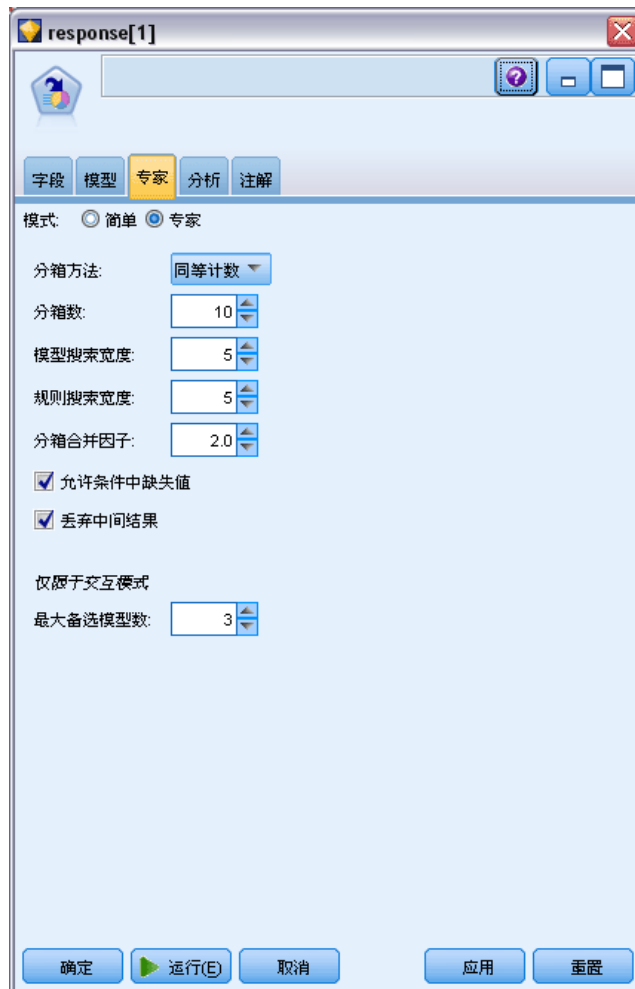
图片 11-6
决策列表节点，“模型”选项卡



- ▶ 选择启动交互会话。
- ▶ 为简化本例中的模型，请将最大段数设为 3。
- ▶ 将新条件的置信区间更改为 85%。

- ▶ 在“专家”选项卡上，将模式设置为专家。

图片 11-7
决策列表节点，“专家”选项卡



- ▶ 将最大替代值数设为 3。此选项与在“模型”选项卡上所选的启动交互会话设置一起使用。
- ▶ 单击运行显示“交互列表”查看器。

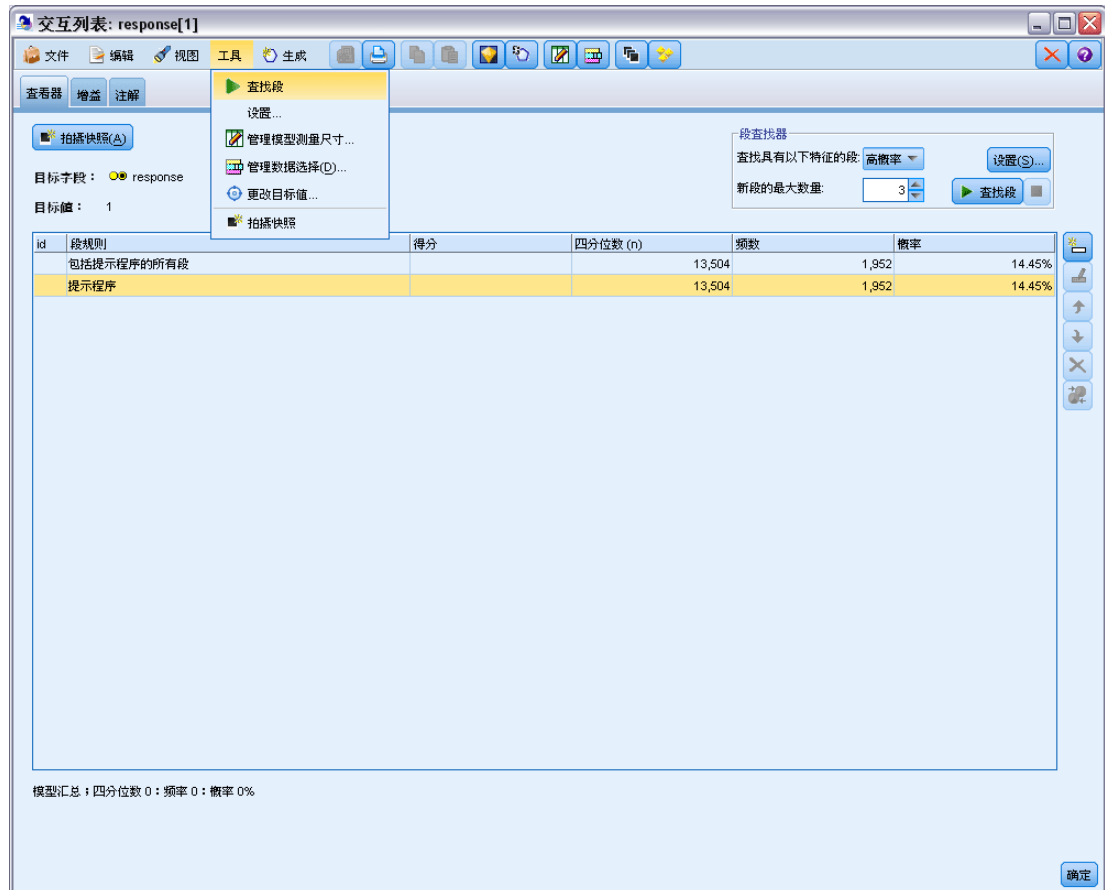
图片 11-8
交互列表查看器

id	段规则	得分	四分位数 (n)	频数	概率
	包括提示程序的所有段		13,504	1,952	14.45%
	提示程序		13,504	1,952	14.45%

尚未定义任何段，因此所有记录都位于其余段中。在示例的 13,504 个记录中，有 1,952 个记录的响应是是，总匹配率为 14.45%。可通过识别更多可能（或较少可能）作积极响应的客户段来提高此匹配率。

- ▶ 在交互列表查看器中，从菜单中选择以下选项：
工具 > 查找段

图片 11-9
交互列表查看器



此操作将根据在决策列表节点中指定的设置运行默认的挖掘任务。任务完成后将获得三个替代模型，这三个模型列在“模型作品集”对话框的“替代”选项卡中。

图片 11-10
可用替代模型

模型作品集

名称	目标	段数	四分位数	频数	概率
替换方法 1	1		3	2,375	1,348 56.76%
替换方法 2	1		3	2,368	1,326 56.00%
替换方法 3	1		3	2,380	1,329 55.84%

替换方法预览

id	段规则	得分	四分位数 (n)	频数	概率
	包括提示程序的所有段		13,504	1,952	14.45%
1	income, number_products income > 55267.000 和 number_products > 1.000	1	912	795	87.17%
2	rfm_score, number_transactions rfm_score > 12.333 和 number_transactions > 2.000	1	737	360	48.85%
3	number_transactions, income number_transactions > 0.000 和 number_transactions <= 1.000 和 income > 46072.000	1	731	174	23.80%

↑ 加载

替换方法 快照

确定 取消 帮助

- 从列表选择第一个替代模型；其详细信息显示在“替代预览”窗格中。

图片 11-11
所选替代模型

名称	目标	段数	四分位数	频数	概率	
替换方法 1	1		3	2,375	1,348	56.76%
替换方法 2	1		3	2,368	1,326	56.00%
替换方法 3	1		3	2,380	1,329	55.84%

id	段规则	得分	四分位数 (n)	频数	概率
	包括提示程序的所有段		13,504	1,952	14.45%
	income, number_products				
1	income > 55267.000 和 number_products > 1.000	1	912	795	87.17%
	rfm_score, number_transactions				
2	rfm_score > 10.535 和 number_transactions > 3.000	1	725	357	49.24%
	average#balance#feed#index, numbe				
3	average#balance#feed#index > 0.000 并 average#balance#feed#index <= 349.0(1 number_products <= 2.000 和 rfm_score > 9.239		738	196	26.56%
	提示程序		11,129	604	5.43%

↑ 加载

替换方法 快照

确定 取消 帮助

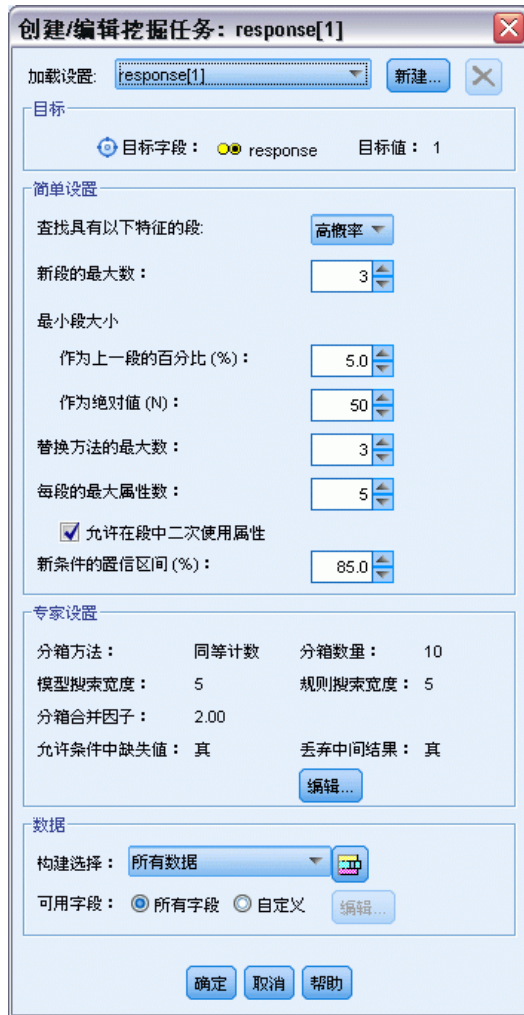
使用“替代预览”窗格，可以在不更改工作模型的情况下快速浏览任意数量的替代模型，从而简化尝试不同方法的过程。

注意：为了更好地查看模型，可能需要最大化对话框中的“替代预览”窗格，方法如下。您可以通过拖动窗格边框来进行此操作。

使用基于预测变量（如收入、每月事务数和 RFM 得分）的规则，模型可识别响应率高于总体样本响应率的段。组合段后，该模型会提示您可以将匹配率提高至 56.76%。但该模型只占总体样本的一小部分，还有将近 11,000 条记录—其中有几百条匹配记录一位于一其余段中。希望模型在仍排除低响应率段的情况下捕获更多的匹配记录。

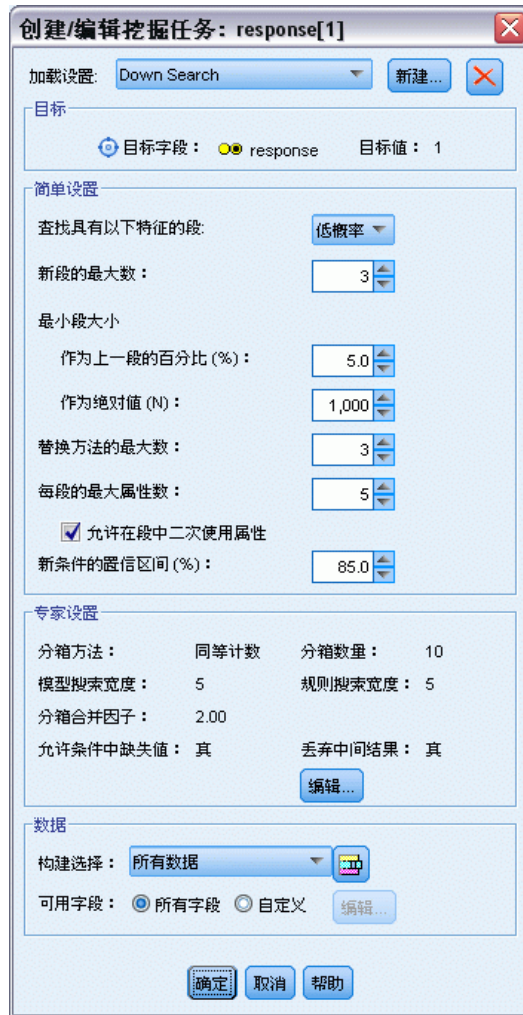
- ▶ 要尝试使用不同的建模方法，可从菜单中选择以下项：
工具 > 设置

图片 11-12
“创建/编辑挖掘任务”对话框



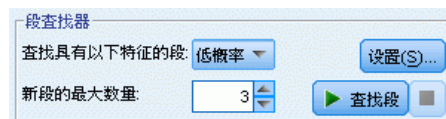
- ▶ 单击新建按钮（右上角）以创建第二个挖掘任务，并指定向下搜索作为“新建设置”对话框中的任务名称。

图片 11-13
“创建/编辑挖掘任务”对话框



- ▶ 将任务的搜索方向更改为低概率。此操作将使算法搜索具有最低（而不是最高）响应率的段。
- ▶ 增加最小段大小到 1,000。单击确定返回到“交互列表”查看器。
- ▶ 在“交互列表”查看器中，确保段查找器窗格正在显示新任务详细信息并单击查找段。

图片 11-14
在新挖掘任务中查找段



该任务返回一组新的替代模型，显示在“模型作品集”对话框的“替代”选项卡中，且可用与以前结果相同的方法进行预览。

图片 11-15
向下搜索模型结果



这一次，每个模型都可识别具有低响应率而不是高响应率的段。浏览第一个替代模型，只需排除这些低响应率段就可以将其余段中的匹配率提高到 39.81%。此值低于前面浏览的模型的值，但其覆盖率更高（总匹配率更高）。

将这两种方法组合使用—使用“低概率”搜索剔除不需要的记录，然后使用“高概率”搜索—可以改进此结果。

- ▶ 单击加载以使其（第一个向下搜索替代模型）成为工作模型并单击确定以关闭“模型作品集”对话框。

图片 11-16
排除段

交互列表: response[1]

目标字段: response
目标值: 1

段查找器
查找具有以下特征的段: 低概率
新段的最大数量: 3
查找段

id	段规则	得分	四分位数 (n)	频数	概率
	包括提示程序的所有段		13,504	1,952	14.45%
1	months_customer months_customer = "0"	1		1,747	0.00%
2	rfm_score rfm_score <= 0.000	1		6,003	0.00%
3	income, rfm_score income > 40297.000 和 income <= 55267.000 和 rfm_score > 0.000 和 rfm_score <= 10.535	1		1,433	232 16.19%
	提示程序			4,321	1,720 39.81%

模型汇总: 四分位数 9,183: 频数 232: 概率 2.53%

确定

- ▶ 分别右键单击前两个段，然后选择排除段。这些段共同捕获将近 8,000 条记录，但之间的匹配率为零，因此可以将其从未来的报价中排除。（排除的段的得分将为空，并以此来表示这些段。）
- ▶ 右键单击第三个段并选择删除段。此段的匹配率 16.19% 与基准匹配率 14.45% 的差别不是很大，因此不足以证明应将其保留。

注意：删除段与排除段是两种不同的操作。排除段只是更改其得分方式，而删除段则将段完全从模型中删除。

排除最低响应率段后，便可以在剩余段中搜索高响应率段。

- ▶ 单击表中的其余行以将其选中，使下一项挖掘任务仅应用于该行。

图片 11-17
选择段

The screenshot shows a software window titled "交互列表: response[1]". It features a menu bar with options like "文件", "编辑", "视图", "工具", and "生成". Below the menu is a toolbar with icons for various actions. The main area contains a "查看器" (Viewer) section with buttons for "增益" (Gain) and "注解" (Annotation). A "快照快捷(A)" (Snapshot Shortcut) button is also present. The "目标字段" (Target Field) is set to "response" and the "目标值" (Target Value) is "1". A "段查找器" (Segment Finder) panel is visible, with a dropdown for "低概率" (Low Probability) and a "设置(S)..." (Settings...) button. The "新段的最大数量" (Maximum number of new segments) is set to 3, and a "查找段" (Find Segments) button is available. The main table displays the following data:

id	段规则	得分	四分位数 (n)	频数	概率
1	包括提示程序的所有段 months_customer months_customer = "0"	1	13,504	1,952	14.45%
2	rfm_score rfm_score <= 0.000	1	1,747	0	0.00%
	提示程序		6,003	0	0.00%
			5,754	1,952	33.92%

At the bottom of the window, there is a summary: "模型汇总: 四分位数 7,750; 频率 0; 概率 0%". A "确定" (OK) button is located in the bottom right corner.

- ▶ 使用所选剩余模型，单击设置以重新打开“创建/编辑挖掘任务”对话框。
- ▶ 在顶部的加载设置中，选择默认挖掘任务：响应[1]。
- ▶ 编辑简单设置以将新段数增加到 5 并将最小段大小增加到 500。

- ▶ 单击确定返回到“交互列表”查看器。

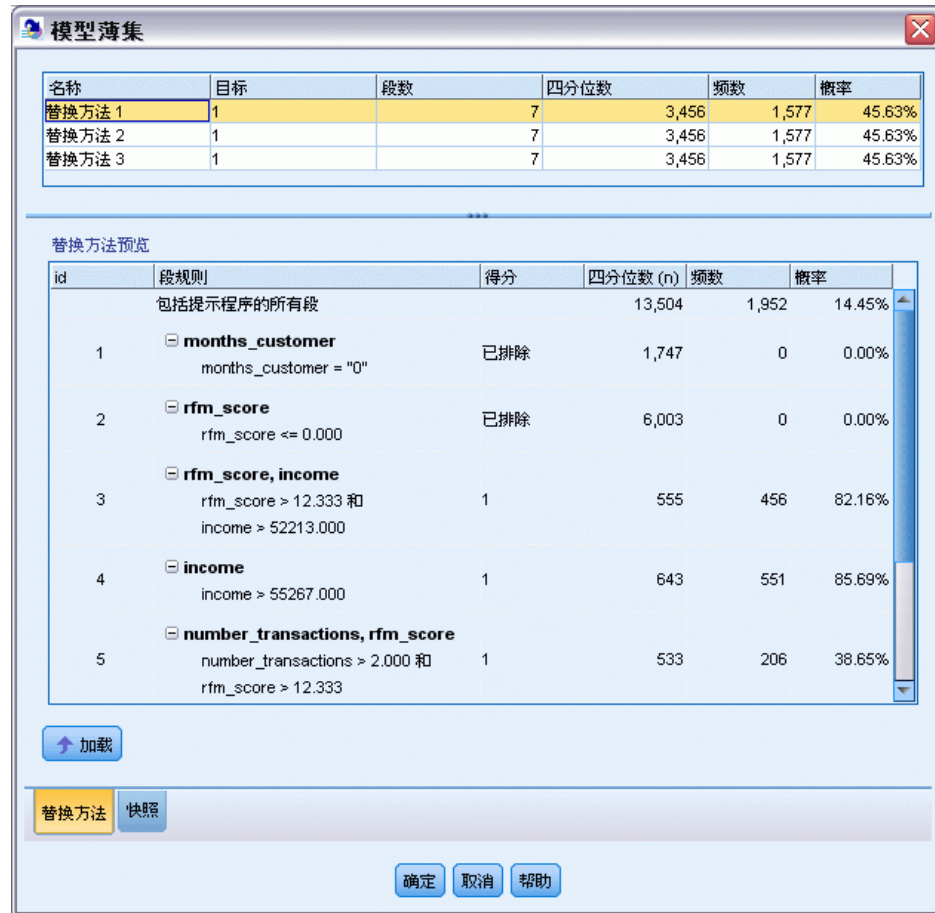
图片 11-18
选择默认挖掘任务。

- ▶ 单击查找段。

这会显示另一组替代模型。通过将一项挖掘任务的结果反馈给另一项挖掘任务，这些最新模型将同时包含高响应率段和低响应率段。具有低响应率的段将被排除，这意味着其得分为 Null，而未排除的段得分将为 1。总体统计量反映了上述排除结果，其中第

一个替代模型具有 45.63% 的匹配率，覆盖率（3,456 条记录中有 1,577 条匹配记录）高于以前任何一个模型。

图片 11-19
组合模型的替代模型



- 预览首个替代模型然后单击加载以使其成为工作模型。

使用 Excel 计算自定义测量量

- 为更清楚地了解模型的实际性能，可以在“工具”菜单中选择组织模型测量。

图片 11-20
组织模型测量

The screenshot shows the 'Interactive List: response[1]' window. The 'Tools' menu is open, highlighting 'Manage Model Measurement Dimensions...'. The 'Segment Finder' panel is visible, showing search criteria: 'High Probability' and 'Maximum Number of Segments: 5'. The main table displays the following data:

id	段规则	得分	四分位数 (n)	频数	概率
	包括提示程序的所有段		13,504	1,952	14.45%
1	months_customer months_customer = "0"	已排除		1,747	0.00%
2	rfm_score rfm_score <= 0.000	已排除		6,003	0.00%
3	rfm_score, income rfm_score > 12.333 和 income > 52213.000	1		555	456 82.16%
4	income income > 55267.000	1		643	551 85.69%
5	number_transactions, rfm_score number_transactions > 2.000 和 rfm_score > 12.333	1		533	206 38.65%

模型汇总：四分位数 3,456；频数 1,577；概率 45.63%

使用“组织模型测量”对话框，可以选择要在交互列表查看器中显示的测量量（或列）。还可以指定是根据所有记录还是根据选定子集计算测量量，如有可能，您可选择显示饼图而不是数值。

图片 11-21
“组织模型测量”对话框



此外，如果已安装 Microsoft Excel，可以链接至 Excel 模板，该模板将用于计算自定义测量量并将其添加到交互显示中。

- ▶ 在“组织模型测量”对话框中，将计算 Excel 中的自定义测量量设置成是。
- ▶ 单击连接到 Excel (TM)
- ▶ 选择 template_profit.xlt 工作簿（位于 IBM® SPSS® Modeler 安装位置的 Demos 文件夹中的 streams 下），然后单击打开启动电子表格。

图片 11-22
Excel 模型测量工作表

	A	B	C	D	E	F	G	
1								
2								Profitability A
3	#	Use	Metric: Frequency	Imported Metric: Cover	Calculated Metric: Profit Margin	Calculated Metric: Cumulative Profit	Target	Segment
4	1					-2,500.00		
5	2							

Excel 模板包含以下三个工作表：

- 模型测量用于显示从模型中导出的模型测量，并计算要重新导出至模型的自定义测量量。
- 设置包含要用于计算自定义测量量的参数。
- 配置用于定义要导入模型和从其中导出的测量量。

重新导出至模型的矩阵如下：

- **毛利**。来自段的净收入
- **累积利润**。来自活动的总利润

通过以下公式定义：

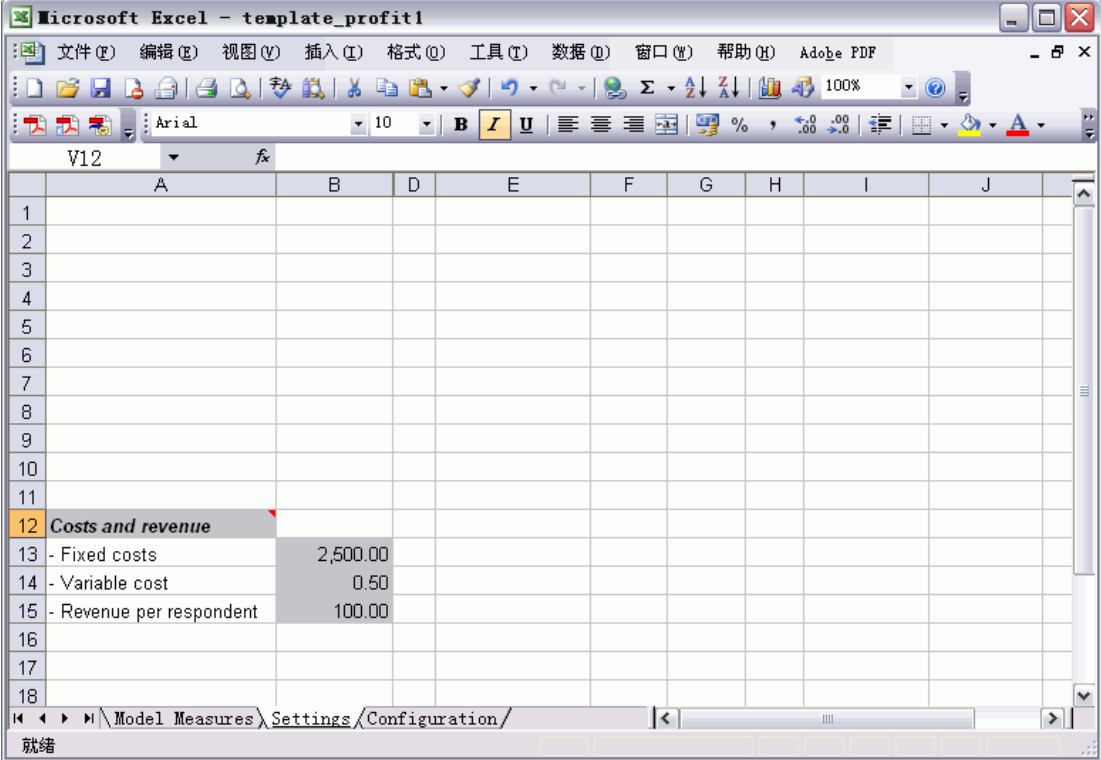
毛利 = 频率 * 每个响应者的收入 - 涉及范围 * 可变成本

累积利润 = 总毛利 - 固定成本

请注意，“频率和涉及范围”从模型中导入。

成本和收入参数由用户在“设置”工作表中指定。

图片 11-23
Excel 设置工作表



The screenshot shows a Microsoft Excel window titled "Microsoft Excel - template_profit1". The spreadsheet has columns A through J and rows 1 through 18. A table is visible starting at row 12, column A. The table has the following data:

	A	B	D	E	F	G	H	I	J
12	Costs and revenue								
13	- Fixed costs	2,500.00							
14	- Variable cost	0.50							
15	- Revenue per respondent	100.00							
16									
17									
18									

The status bar at the bottom shows the path: \Model Measures\Settings\Configuration/ and the text "就绪".

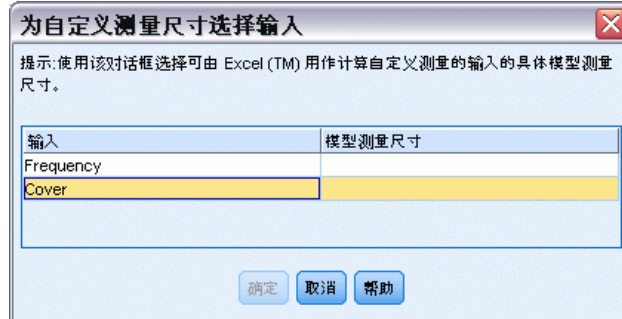
固定成本是活动的启动成本，如设计和计划的成本。

可变成本是将报价送达每位客户的成本，如信封和邮票的成本。

每个响应者的收入是响应报价的客户的净收入。

- ▶ 要完成返回模型的链接，可使用 Windows 任务栏（或按 Alt+Tab）返回交互列表查看器。

图片 11-24
选择自定义测量的输入



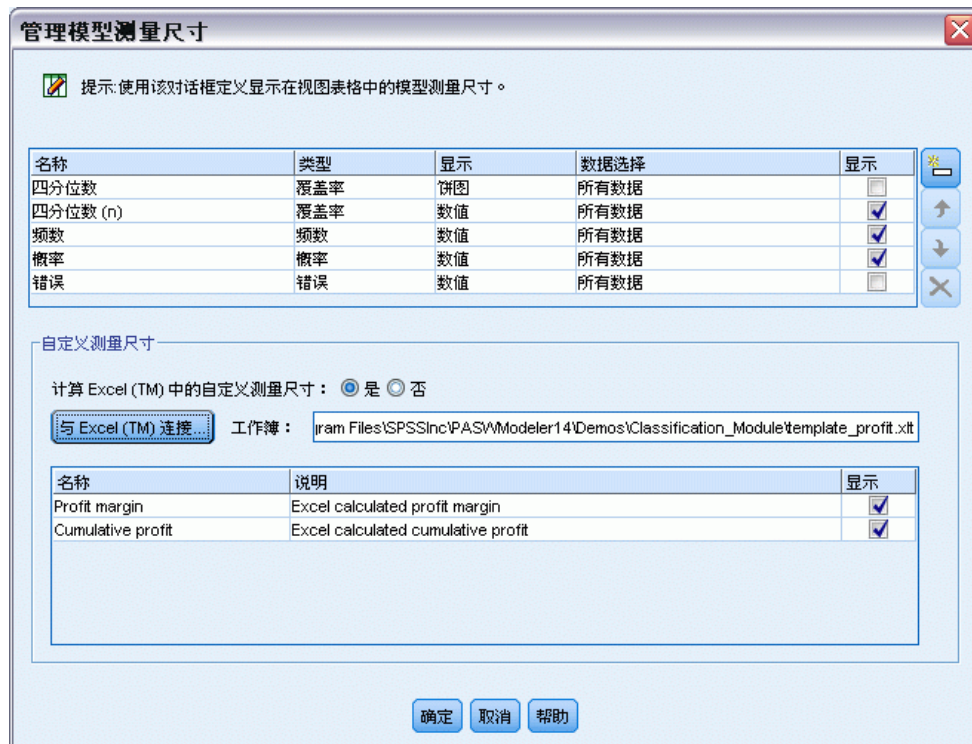
此时将显示“选择自定义测量的输入”对话框，您可在其中将模型中的输入映射为模板中定义的特定参数。左列列出可用测量量，右列将这些测量量映射为“配置”工作表中定义的电子表格参数。

- ▶ 在模型测量量列中，选择相对于各自输入的频率和涉及范围（n）并单击确定。

在本例中，模板—频率和涉及范围（n）—中的参数名恰好与输入匹配，但也可以使用其他名称。

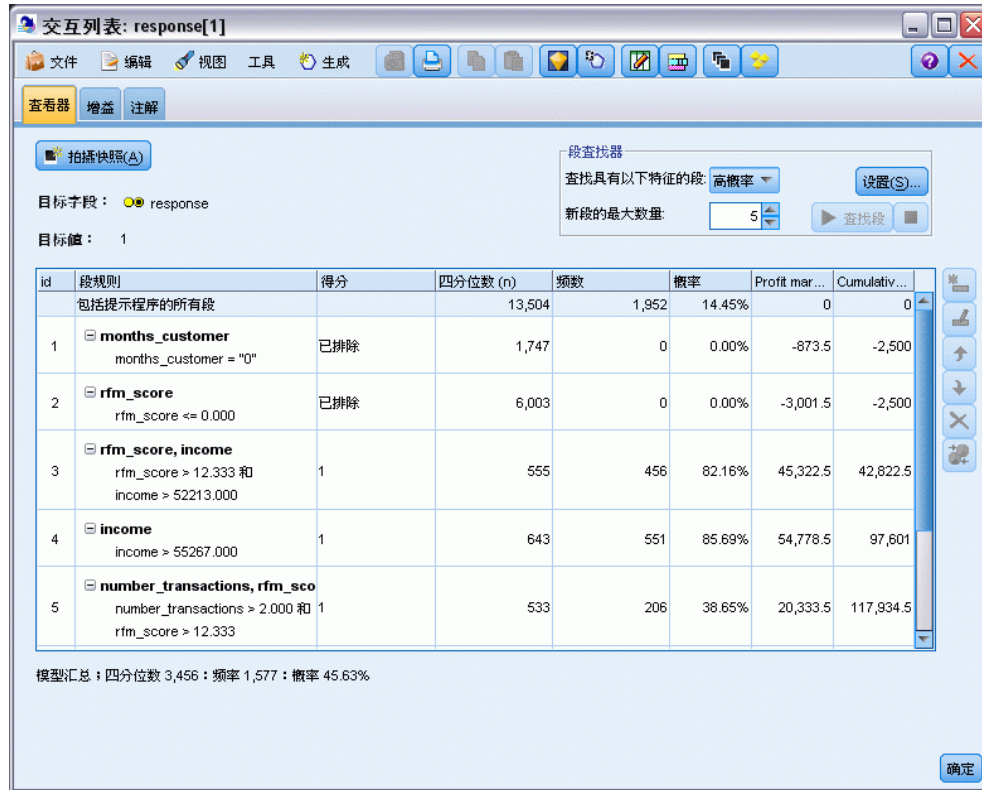
- ▶ 在“组织模型测量”对话框中单击确定以更新交互列表查看器。

图片 11-25
显示 Excel 中的自定义测量量的“组织模型测量”对话框



这时新的测量量将作为新列添加到窗口中，并在模型每次更新时进行重新计算。

图片 11-26
交互列表查看器中显示的 Excel 自定义测量量



通过编辑 Excel 模板，可以创建任意数量的自定义测量量。

修改 Excel 模板

尽管 IBM® SPSS® Modeler 提供了可以与交互列表查看器一起使用的默认 Excel 模板，但您可能希望更改设置或添加自己的设置。例如，对您的组织而言，模板中的成本可能并不正确并且需要修改。

注意：如果确实要修改现有模板或创建自己的模板，请记住应以 Excel 2003 .xlt 为后缀保存文件。

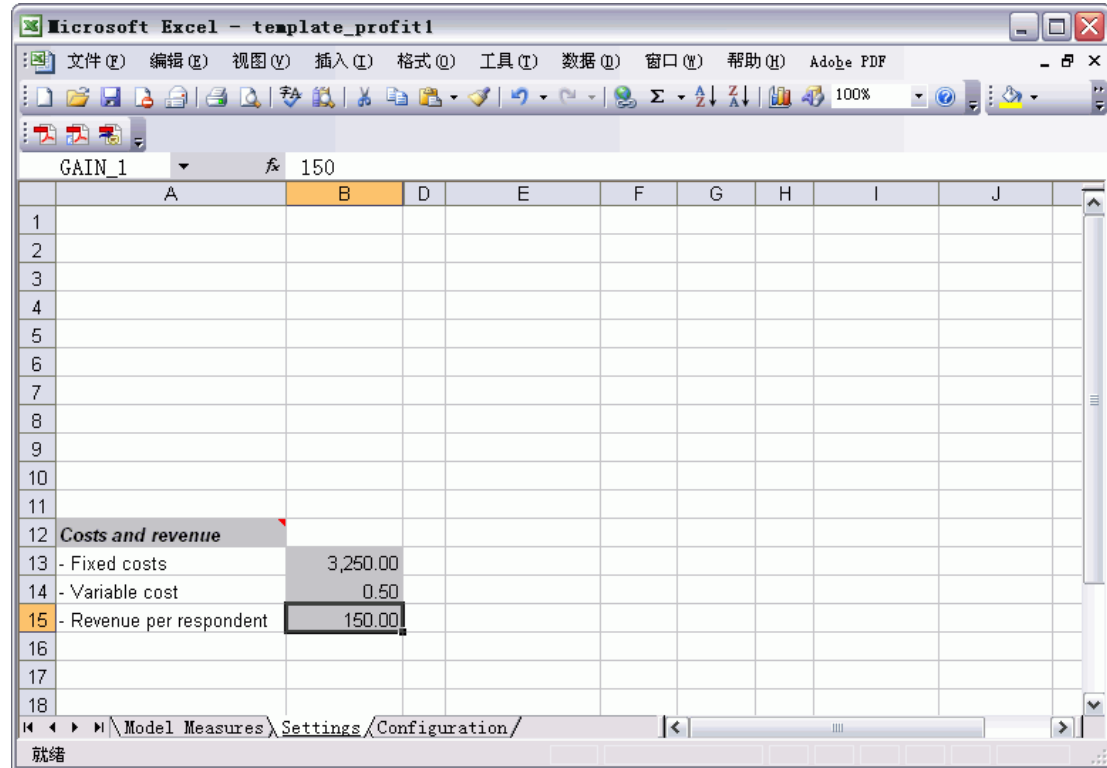
要使用新成本和收入细节修改默认模板，并使用新数据更新交互列表查看器：

- ▶ 请在交互列表查看器中，选择“工具”菜单中的组织模型测量。
- ▶ 在“组织模型测量”对话框中，单击连接到 Excel。
- ▶ 选择 template_profit.xlt 工作簿并单击打开以启动电子表格。
- ▶ 选择“设置”工作表。

- ▶ 将固定成本编辑为 3,250.00，并将每个响应者的收入编辑为 150.00。

图片 11-27

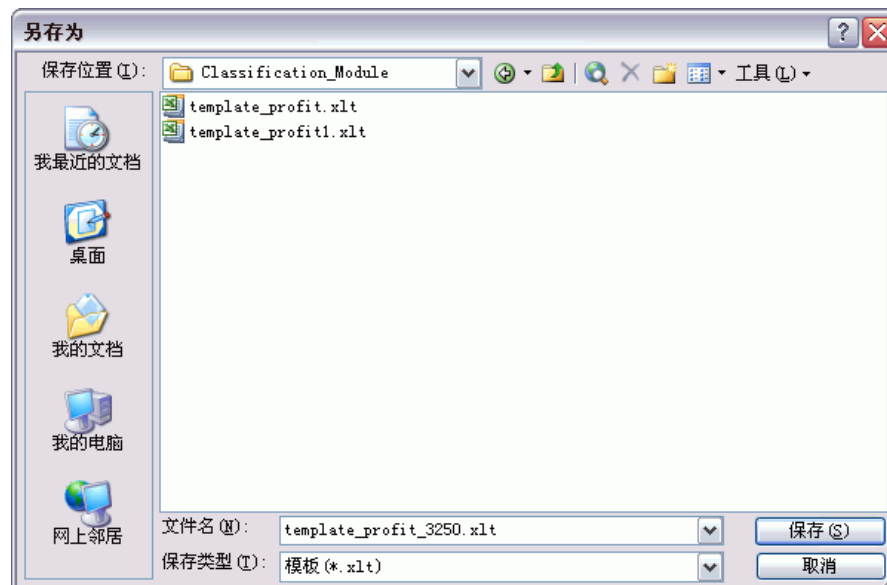
Excel 设置工作表中的已修改值



- ▶ 使用唯一且相关的文件名保存已修改的模板。请确保文件的扩展名为 Excel 2003 .xlt。

图片 11-28

保存已修改的 Excel 模板



- ▶ 使用 Windows 任务栏（或按 Alt+Tab）返回交互列表查看器。
在“选择自定义测量的输入”对话框中，选择要显示的测量量并单击确定。
- ▶ 在“组织模型测量”对话框中，单击确定以更新交互列表查看器。

显然，本模型仅显示了一种修改 Excel 模板的简单方式；您还可以进一步更改，例如，从交互列表查看器提取数据或向其中传递数据；也可以在 Excel 中生成其他输出（如图形）。

图片 11-29
交互列表查看器中显示的已修改的 Excel 自定义测量量



保存结果

为了保存模型以供日后在交互会话过程中使用，可以对模型进行快照，此快照将列在“快照”选项卡上。在交互会话过程中可以随时返回任何已保存的快照。

按此方式继续，可以尝试使用其他挖掘任务搜索其他段。还可以编辑现有段、根据自己的业务规则插入自定义段、创建数据选择以优化特定组的模型以及以多种其他方式自定义模型。最后，可以明确地包括或排除每个段，同时正确指定每个段的得分方式。

如果对结果感到满意，则可以使用“生成”菜单生成模型，此模型可添加到流或部署此模型以用于评分。

此外，要保存交互会话的当前状态以备他日使用，可以在“文件”菜单中选择**更新建模节点**。此操作将“决策列表”建模节点更新为当前设置，包括挖掘任务、模型快照、数据选择和自定义测量量。下次运行流时，只需确保在决策列表建模节点中选中了**使用保存的会话信息**，就可将会话恢复到其当前状态。[有关详细信息，请参阅第 9 章中的决策表中的 IBM SPSS Modeler 15 建模节点。](#)

电信业客户分类（多项 Logistic 回归）

Logistic 回归是一种统计方法，它可根据输入字段的值对记录进行分类。这种技术与线性回归类似，但用分类目标字段代替了数值字段。

例如，假设某个电信提供商根据服务使用模式对它的客户群进行了分段，将这些客户分为了四个组。如果人口统计学数据可用于预测组成员资格，则可以为各个潜在客户自定义服务。

此示例使用名为 telco_custcat.str 的流，此流引用名为 telco.sav 的数据文件。这些文件可在任何 IBM® SPSS® Modeler 安装程序的 Demos 目录中找到。此目录可通过 WindowsIBM® SPSS® Modeler 程序组进行访问。文件 telco_custcat.str 位于 streams 目录下。

该示例主要关注于使用人口统计数据预测使用模式。目标字段 custcat 有四个可能的值对应于四个客户组，如下所示：

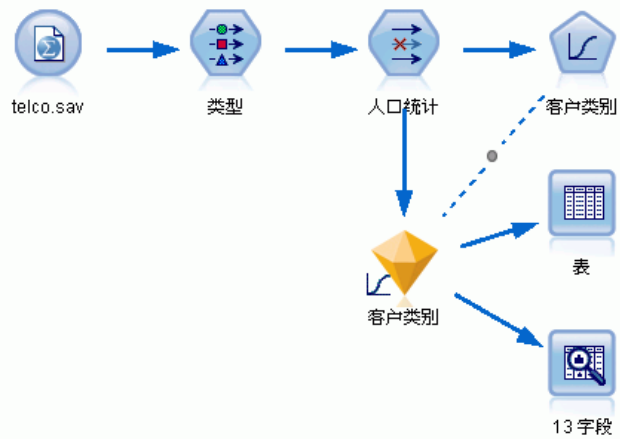
值	Label
1	基本服务
2	电子服务
3	增值服务
4	全套服务

因为目标含有多个类别，因此将使用多项模型。如果目标含有两个截然不同的分类，例如是/否，真/假，或流失/保持，则会转而创建二项模型。[有关详细信息，请参阅第 152 页码第 13 章中的电信客户流失（二项 Logistic 回归）。](#)

构建流

- ▶ 在 Demos 文件夹中添加指向 telco.sav 的 Statistics 文件源节点。

图片 12-1
使用多项 logistic 回归分类客户的示例流



- ▶ 添加类型节点并单击读取值，确保所有测量水平设置正确。例如，具有值 0 和 1 的多数字段可视为标志。

图片 12-2
设置多个字段的测量级别



提示：要更改具有相似值（如 0/1）的多个字段，请单击值列标题，以便按值对字段进行排序，然后按住 Shift 键的同时使用鼠标或箭头键选择所有要更改的字段。然后可以右键单击选定的内容以更改选定字段的测量级别或其他属性。

注意，性别更准确而言应视为具有两个值的集合的字段，而不是标志，所以将其测量值保留为名义。

- ▶ 将客户类别字段的角色设置为目标。将所有其他字段的角色设置为 Input。

图片 12-3
设置字段角色



因为此示例主要关注人口统计，所以请使用过滤节点以选取相关字段（地区、年龄、婚姻状况、地址、收入、教育程度、行业、退休、性别、居住地和客户类别）。其他字段可以排除在此分析之外。

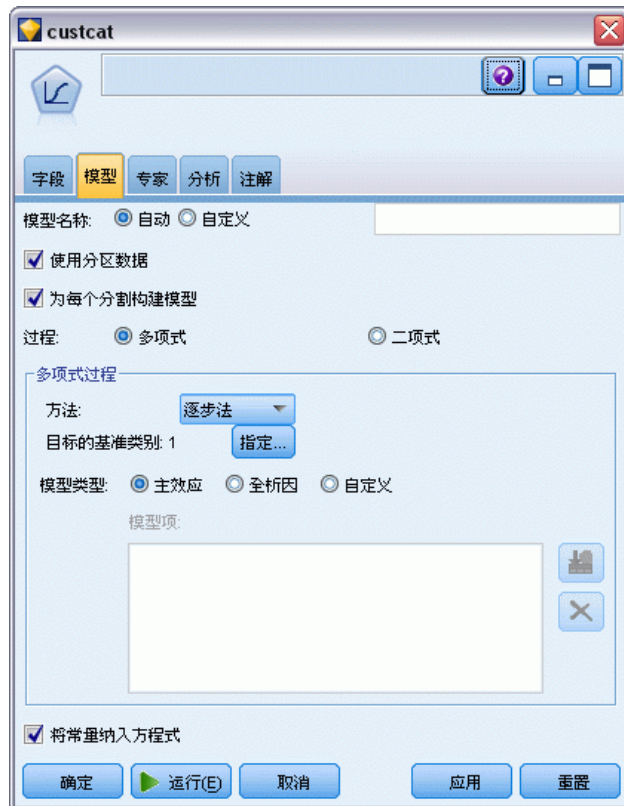
图片 12-4
过滤人口统计字段



（另外，您可以将这些字段的角色更改为无，而不要排除这些字段，或者选择要在建模节点中使用的字段。）

- ▶ 在 Logistic 节点上，单击模型选项卡并选择逐步法。选中多项、主效应和将常量纳入方程式。

图片 12-5
选择模型选项



将目标的底数类别保留为 1。模型将对其他客户与预订基本服务的客户进行比较。

- ▶ 在“专家”选项卡上，选中专家模式，选中输出，然后在“高级输出”对话框中选中分类表。

图片 12-6
选择输出选项

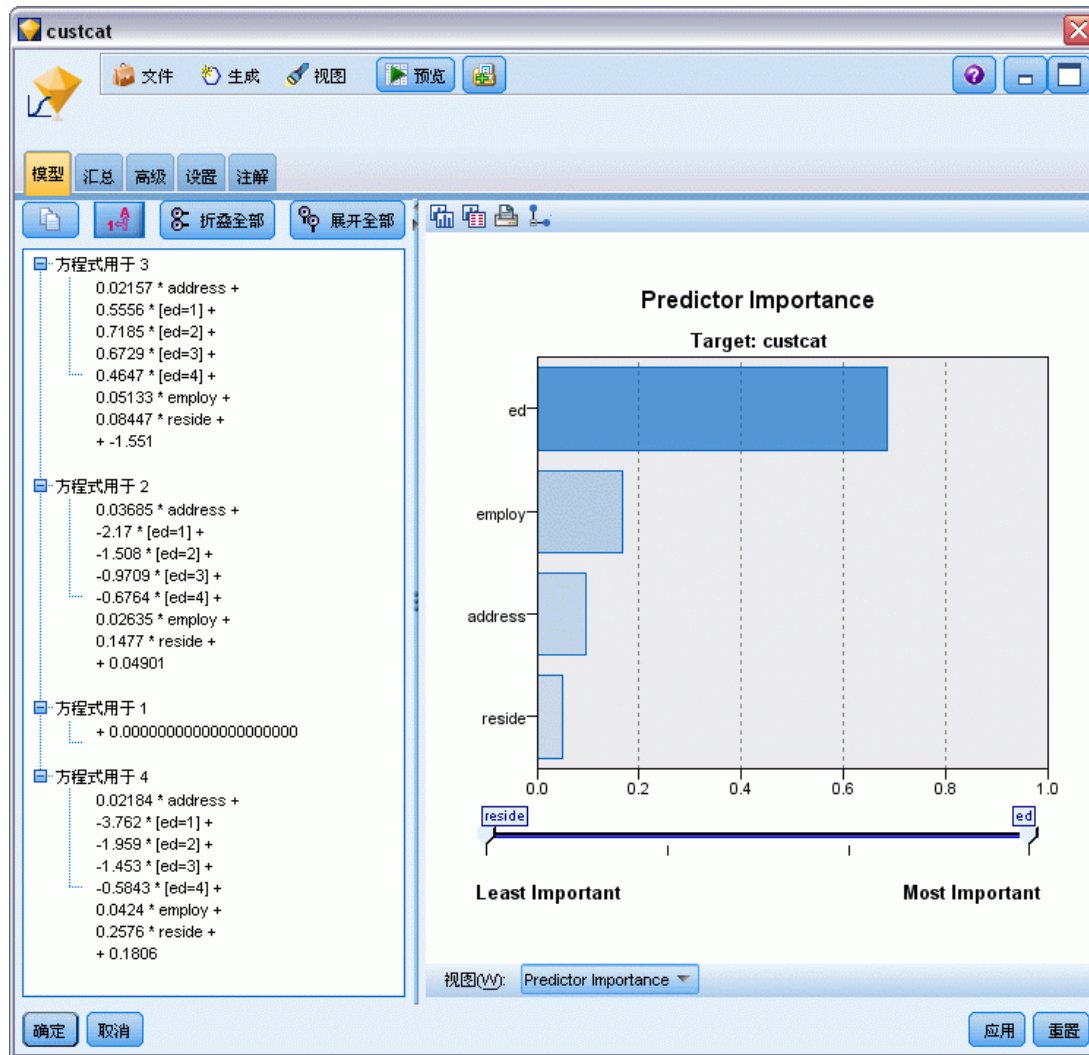


浏览模型

- ▶ 执行这个节点以生成模型，它被添加至右上角的模型调色板中。要查看其详细信息，请在生成的模型节点上用右键单击并选择浏览。

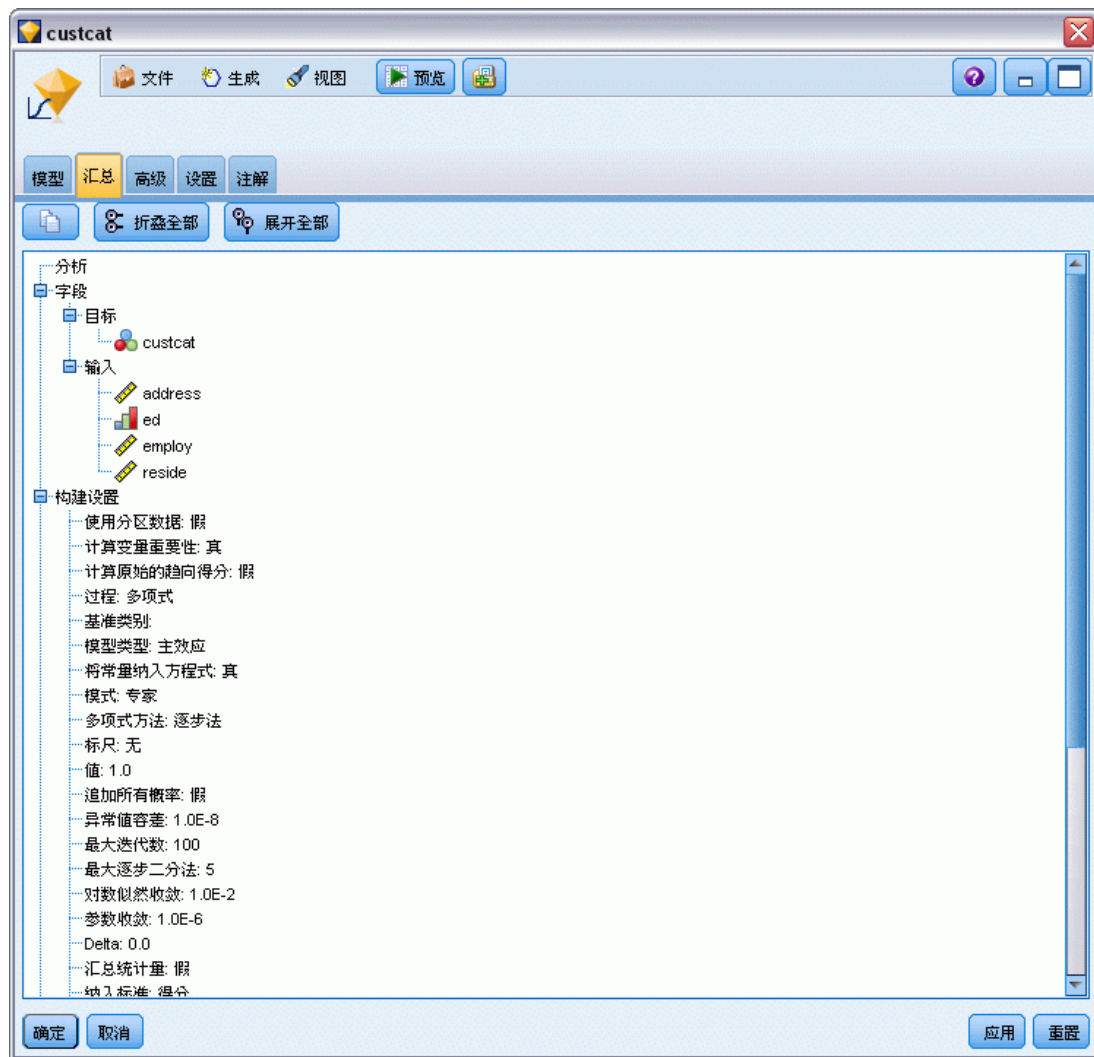
“模型”选项卡中显示了用于将记录分配到目标字段的每个类别的方程式。有四个可能的类别，其中之一就是基准类别，在该类别中没有显示方程式的详细信息。在余下的三个方程式中显示了详细信息，其中类别 3 表示附加服务，依此类推。

图片 12-7
浏览模型结果



“汇总”选项卡显示了（包括其他内容）模型中使用的目标字段和输入字段（预测变量字段）。注意，这些字段是根据逐步法实际选择出来的字段，不是为进行分析而提交的完整列表。

图片 12-8
显示目标和输入字段的模型汇总

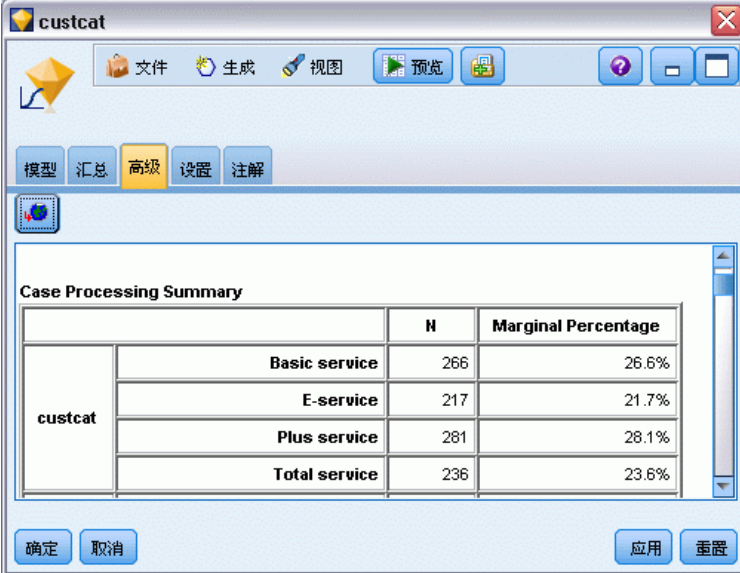


“高级”选项卡上显示的项目取决于在建模节点的“高级输出”对话框中选中的选项。

其中通常显示的一个项目是观测值处理概要，它显示了落在目标字段每个类别中的记录的百分比。这将生成一个空模型用作比较的基础。

在不构建使用预测变量的模型的情况下，最好的预计结果可能是将所有的客户分配到最普通的组（附加服务组）中。

图片 12-9
个案处理摘要



The screenshot shows a software window titled 'custcat' with a menu bar containing '文件', '生成', '视图', '预览', and a help icon. Below the menu bar are tabs for '模型', '汇总', '高级', '设置', and '注解'. The main area displays a 'Case Processing Summary' table with the following data:

		N	Marginal Percentage
custcat	Basic service	266	26.6%
	E-service	217	21.7%
	Plus service	281	28.1%
	Total service	236	23.6%

At the bottom of the window are buttons for '确定', '取消', '应用', and '重置'.

如果基于训练数据将所有客户分配到空模型，则得到的正确率将是 $281/1000 = 28.1\%$ 。“高级”选项卡还包括其他信息，使您能够检查模型的预测。然后，可将这些预测与空模型的结果相比，以查看使用此数据的模型的执行效果。

在“高级”选项卡底部，分类表显示了此模型的结果，其正确率为 39.9% 。

特别是，此模型在识别全套服务客户（类别 4）时表现优异，而在识别电子服务客户（类别 2）时表现很差。如果想提高预测类别 2 中客户的准确性，可能需要再找到一个预测变量来识别此类客户。

图片 12-10
分类表

Observed	Predicted				Percent Correct
	Basic service	E-service	Plus service	Total service	
Basic service	122	8	75	61	45.9%
E-service	58	10	68	81	4.6%
Plus service	89	8	133	51	47.3%
Total service	47	12	43	134	56.8%
Overall Percentage	31.6%	3.8%	31.9%	32.7%	39.9%

依赖于您所期望的预测，模型可充分满足您的需求。例如，如果您对识别类别 2 中的客户并不关心，那么该模型的准确性足以满足需求。这种情况可能是，电子服务仅是一种为吸引顾客而出售且获利微薄的产品。

例如，如果投资的最高回报来自于落在类别 3 或类别 4 中的客户，则该模型能够提供所需的信息。

当构建模型时，可使用“高级输出”对话框中的大量诊断信息来评估模型实际拟合数据的程度。有关详细信息，请参阅第 10 章中的 [Logistic 模型块高级输出中的 IBM SPSS Modeler 15 建模节点](#)。有关 IBM® SPSS® Modeler 中所用建模方法的数学原理的说明，请参阅《SPSS Modeler 算法指南》，该指南位于安装光盘的 \Documentation 目录中。

还请注意，这些结果仅基于训练数据产生。要评估模型对实际应用中的其他数据的拟合程度，可使用分区节点保留部分记录，以便于测试和验证。有关详细信息，请参阅第 4 章中的分区节点中的 [IBM SPSS Modeler 15 源、过程和输出节点](#)。

电信客户流失（二项 Logistic 回归）

Logistic 回归是一种统计方法，它可根据输入字段的值对记录进行分类。这种技术与线性回归类似，但用分类目标字段代替了数值字段。

此示例使用名为 telco_churn.str 的流，此流引用名为 telco.sav 的数据文件。这些文件可在任何 IBM® SPSS® Modeler 安装程序的 Demos 目录中找到。此目录可通过 WindowsIBM® SPSS® Modeler 程序组进行访问。文件 telco_churn.str 位于 streams 目录下。

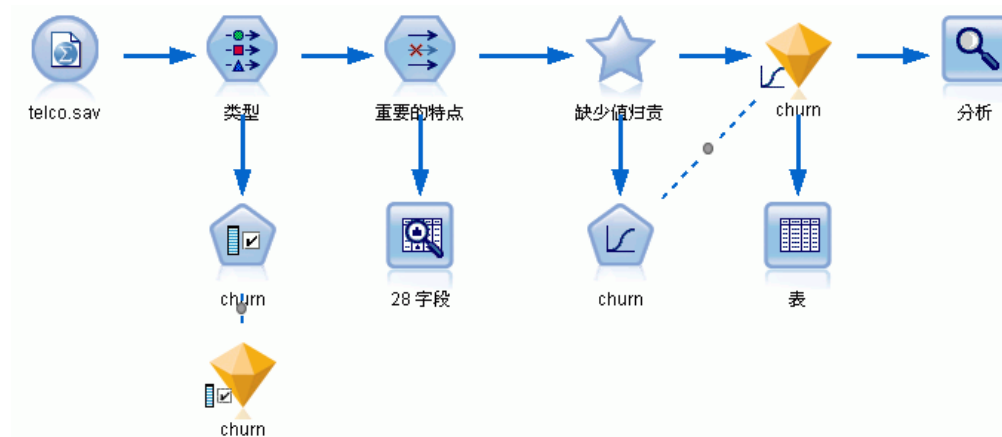
例如，假设某个电信服务提供商非常关心流失到竞争对手那里的客户数。如果可以和服务使用数据预测有可能转移到其他提供商的客户，则可通过自定义服务使用数据来尽可能多地保留这些客户。

本示例将焦点集中于利用使用数据预测客户的丢失（流失）。因为目标含有两个截然不同的类别，因此将使用二项模型。如果目标中含有多个类别，则会转而创建多项模型。有关详细信息，请参阅第 142 页码第 12 章中的电信业客户分类（多项 Logistic 回归）。

构建流

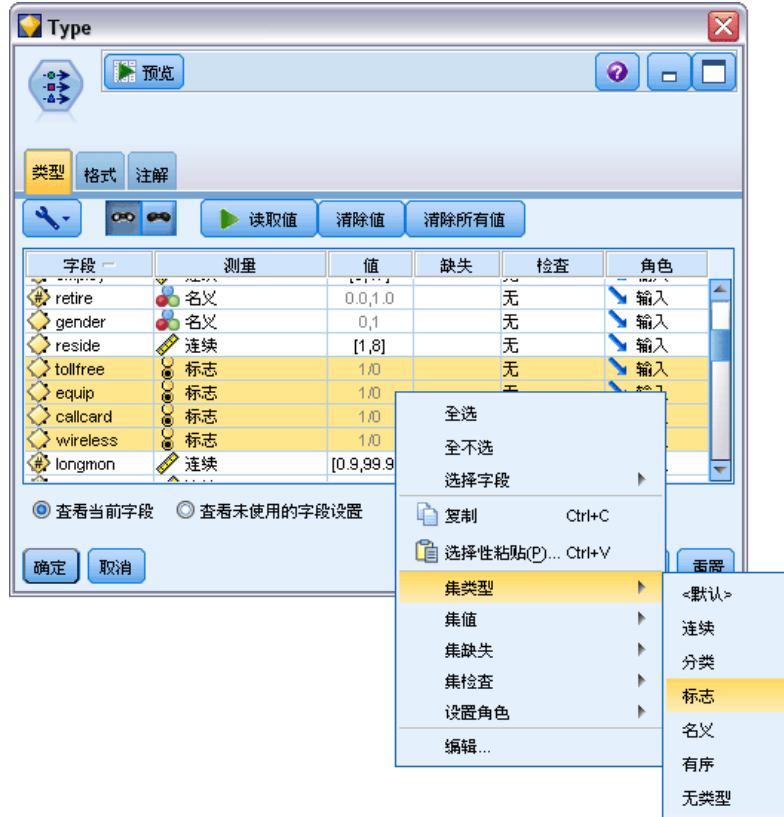
- 在 Demos 文件夹中添加指向 telco.sav 的 Statistics 文件源节点。

图片 13-1
使用二项 logistic 回归分类客户的示例流



- ▶ 添加“类型”节点以定义字段，确保所有测量级别都已正确设置。例如，大多数值为 0 和 1 的字段都可以用作标志字段，但某些字段，比如性别，作为包含两个值的名义字段会更加准确。

图片 13-2
设置多个字段的测量级别



提示：要更改具有类似值（例如 0/1）的多个字段的属性，可单击值列标题以便按值对字段进行排序，然后在使用鼠标或箭头键时按住 Shift 键以选定所有要更改其属性的字段。然后可以右键单击选定的内容以更改选定字段的测量级别或其他属性。

- ▶ 将流失字段的测量级别设置为标志，并将其角色设置为目标。将所有其他字段的角色设置为 Input。

图片 13-3
设置流失字段的测量水平和角色



- ▶ 为类型节点添加“特征选择”建模节点。

通过使用“特征选择”节点，对于不能为预测变量/目标之间的关系添加任何有用信息的预测变量或数据，可以将其删除。

- ▶ 运行流。

- ▶ 打开结果模型块，并从生成菜单中，选择过滤以创建过滤节点。

图片 13-4
从“特征选择”节点中生成“过滤”节点



不是 telco.sav 文件中的所有数据都对预测客户流失有用。可以使用过滤器仅选择被认为很重要的数据来用作预测变量。

- ▶ 在“生成过滤”对话框中，选择所有标记的字段：重要并单击确定。

- ▶ 将生成过滤器节点附加到类型节点。

图片 13-5
选择重要字段



- ▶ 将数据审核节点附加到生成的“过滤”节点。
打开数据审核节点，然后单击运行。
- ▶ 在“数据审核”浏览器的“质量”选项卡上，单击 % 完成列以便按数值升序顺序对此列进行排序。这样就可以识别所有含有大量缺失数据的字段；在本示例中，唯一需要修改的字段是 logtoll，其完成值比例小于 50%。

- ▶ 在 logtoll 的归因于缺失列中，单击指定。

图片 13-6
logtoll 的归因于缺失值



- ▶ 对于归因条件，选择空白值和Null 值。对于固定为，选择平均值，然后单击确定。

选择平均值可确保归因值不会反过来影响总数据中所有值的平均值。

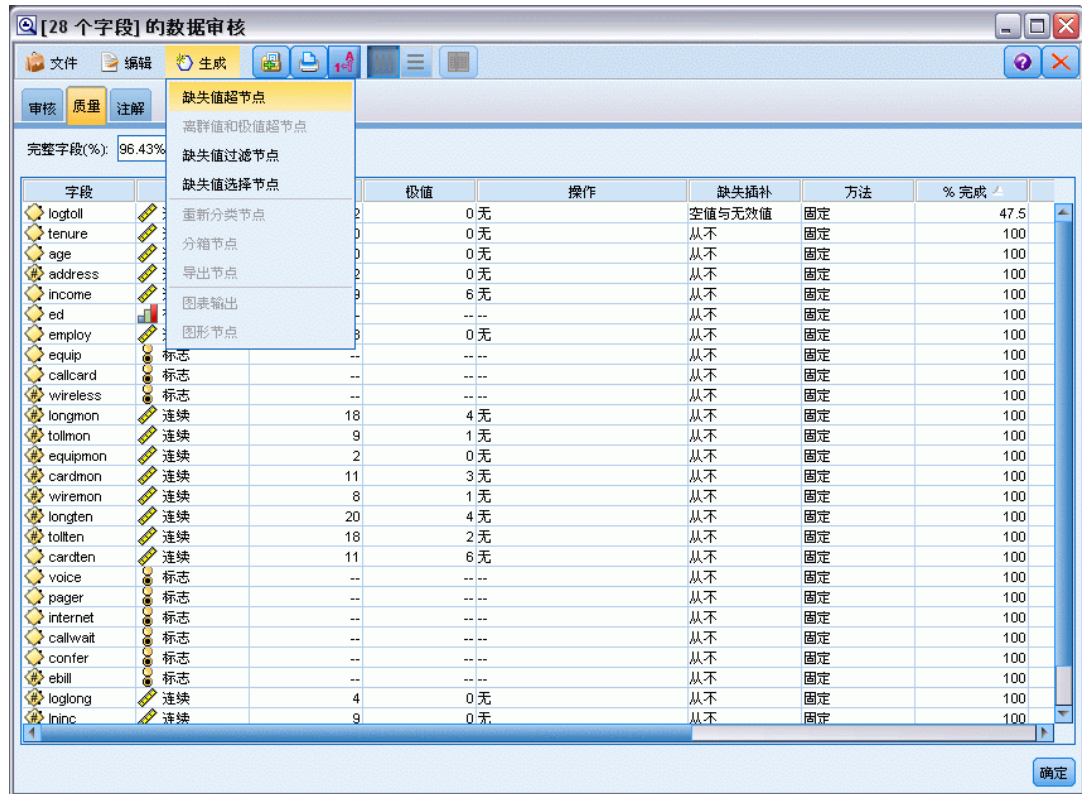
图片 13-7
选择归因设置



- ▶ 在“数据审核”浏览器的“质量”选项卡上，生成缺失值超节点。为此，请从菜单中选择：

生成 > 缺失值超节点

图片 13-8
生成缺失值超节点

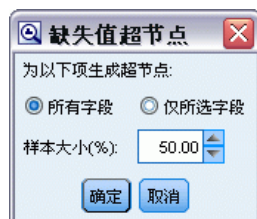


在“缺失值超节点”对话框中，将样本大小增加到 50%，然后单击确定。

此时超节点将显示在流工作区中，其标题为：缺失值归因。

- ▶ 将超节点附加到过滤节点。

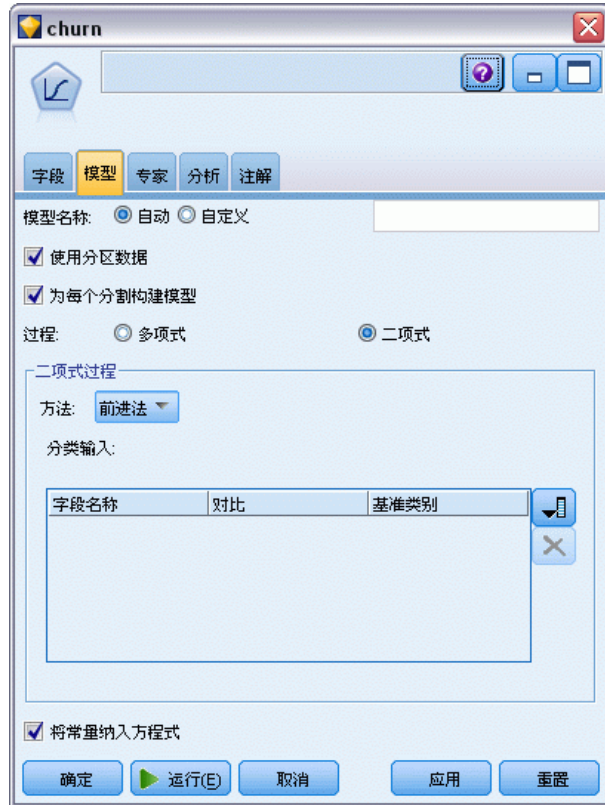
图片 13-9
指定样本大小



- ▶ 将 Logistic 节点添加到超节点。

- ▶ 在 Logistic 节点上，单击“模型”选项卡并选择二项过程。在二项过程区域，选择前进法。

图片 13-10
选择模型选项



- ▶ 在“专家”选项卡上，选择专家模式，然后单击输出。此时显示“高级输出”对话框。
- ▶ 在“高级输出”对话框中，选择在每个步骤作为显示类型。选择迭代历史和参数估计，然后单击确定。

图片 13-11
选择输出选项



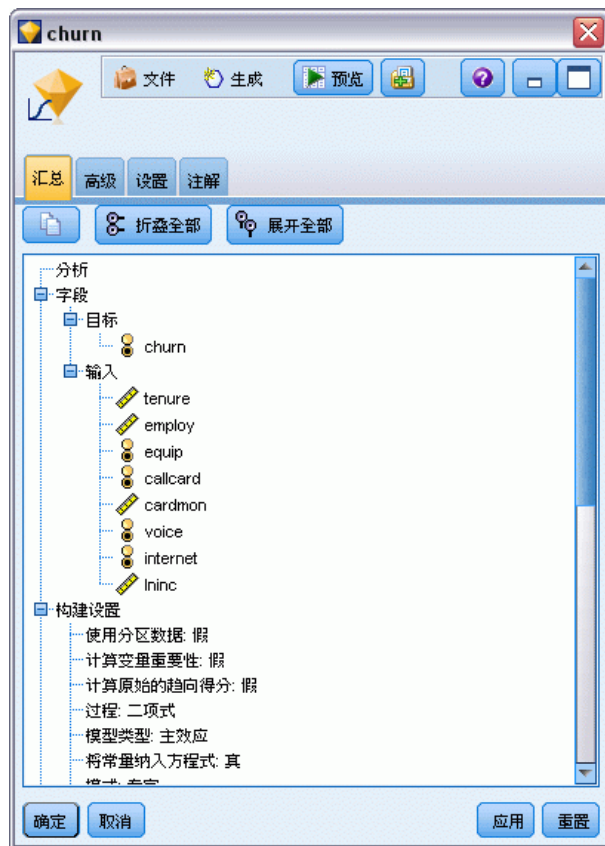
浏览模型

- ▶ 在 Logistic 节点上，单击运行创建模型。

模型块被添加到流工作区，同时添加到右上角的“模型”选项板。要查看其详细信息，右键单击模型块并选择编辑或浏览。

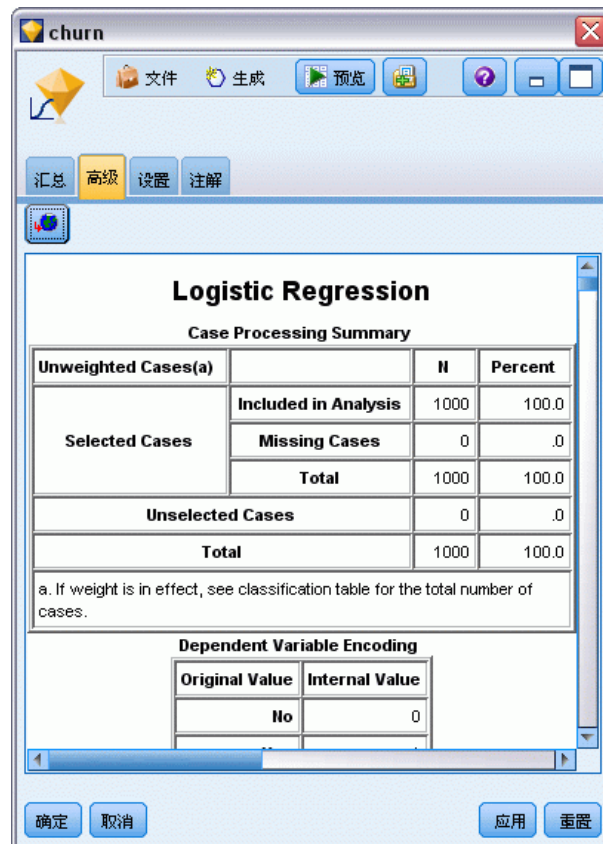
“汇总”选项卡显示了（包括其他内容）模型中使用的目标字段和输入字段（预测变量字段）。注意，这些字段是根据前进法实际选择出来的字段，不是为进行分析而提交的完整列表。

图片 13-12
显示目标和输入字段的模型汇总



“高级”选项卡上显示的项目取决于在 Logistic 节点的“高级输出”对话框中选中的选项。其中通常显示的一个项目是观测值处理概要，它显示了包括在分析中的记录数及百分比。此外，在此汇总中还列出了其中有一个或多个输入字段不可用的缺失观测值的数目（如果有的话），及所有未选定的观测值数。

图片 13-13
个案处理摘要



- 向下滚动观测值处理概要以显示块 0 下的分类表：起始块。

开始使用前进逐步法时会有一个空模型（即，没有预测变量的模型），可将此空模型用作与最终构建的模型进行比较的基础。空模型按常规将所有值预测为 0，因此空模型的准确性为 72.6%，这仅仅是因为已正确预测到有 726 个没有流失的客户。但是，根本没有正确预测到已流失的客户。

图片 13-14
开始分类表 - 块 0

b. Initial -2 Log Likelihood: 1174.394

c. Estimation terminated at iteration number 4 because parameter estimates changed by less than .000.

Classification Table(a,b)

	Observed	churn	Predicted		Percentage Correct
			churn		
			No	Yes	
Step 0	churn	No	726	0	100.0
		Yes	274	0	.0
Overall Percentage					72.6

a. Constant is included in the model.

b. The cut value is .500

Variables in the Equation

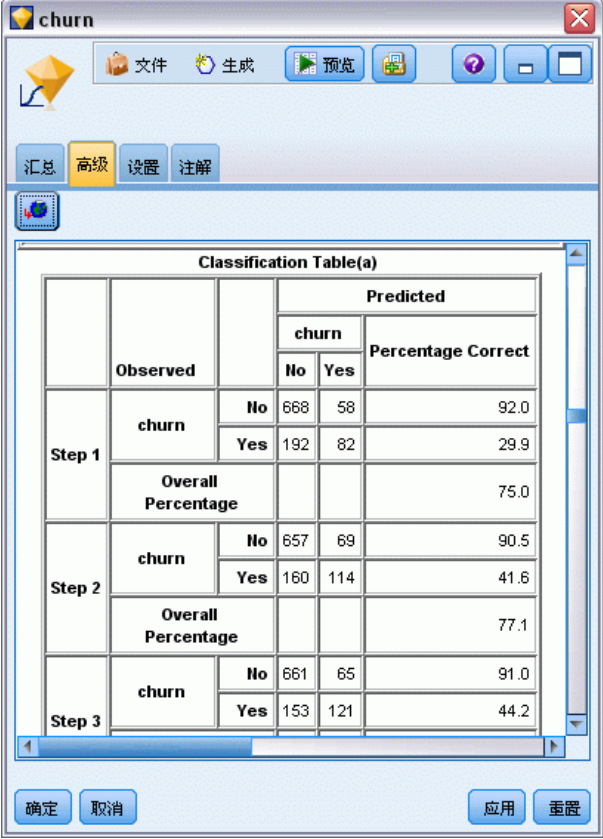
	B	S.E.	Wald	df	Sig.	Exp(B)
--	---	------	------	----	------	--------

确定 取消 应用 重置

- ▶ 现在向下滚动以显示块 1（方法 = 前进逐步）下的分类表。

此分类表显示了模型在每个步骤中添加的预测变量。在第一个步骤中（在仅使用了一个预测变量之后），模型预测流失的准确性就已从 0.0% 增加到 29.9%。

图片 13-15
分类表 - 块 1



		Observed	Predicted		Percentage Correct
			churn		
			No	Yes	
Step 1	churn	No	668	58	92.0
		Yes	192	82	29.9
	Overall Percentage				75.0
Step 2	churn	No	657	69	90.5
		Yes	160	114	41.6
	Overall Percentage				77.1
Step 3	churn	No	661	65	91.0
		Yes	153	121	44.2

- 向下滚动到此分类表的底部。

分类表在步骤 8 之后结束。在此步骤中，算法已确定不用再向模型添加任何其他预测变量。虽然预测非流失客户的准确性有所下降，达到了 91.2%，但预测已流失客户的准确性却从原来的 0% 上升到了 47.1%。这相比原来不使用任何预测变量的空模型其有效性显著提高。

图片 13-16
分类表 - 块 1

		Overall Percentage				
						78.7
Step 7	churn	No	657	69		90.5
		Yes	144	130		47.4
	Overall Percentage					
Step 8	churn	No	662	64		91.2
		Yes	145	129		47.1
	Overall Percentage					

a. The cut value is .500

Variables in the Equation							
		B	S.E.	Wald	df	Sig.	Exp(B)
Step 1(a)	tenure	-.046	.004	123.346	1	.000	.955
	Constant	.462	.136	11.574	1	.001	1.587

对于希望减少流失的客户，能够将流失率减少接近一半将会成为保护其收入流的主要步骤。

注意：此示例还显示出将总体百分比看作判断模型准确性的依据在某些情况下易引起错误。原来空模型的总准确性为 72.6%，而最终预测模型的总准确性为 79.1%；但是，正如我们所看到的，其实际单个类别的预测准确性的差别极大。

当构建模型时，可使用“高级输出”对话框中的大量诊断信息来评估模型实际拟合数据的程度。有关详细信息，请参阅第 10 章中的 Logistic 模型块高级输出中的 IBM SPSS Modeler 15 建模节点。有关 IBM® SPSS® Modeler 中所用建模方法的数学原理的说明，请参阅《SPSS Modeler 算法指南》，该指南位于安装光盘的 \Documentation 目录中。

还请注意，这些结果仅基于训练数据产生。要评估模型对实际应用中的其他数据的拟合程度，可使用分区节点保留部分记录，以便于测试和验证。有关详细信息，请参阅第 4 章中的分区节点中的 IBM SPSS Modeler 15 源、过程和输出节点。

预测带宽利用率（时间序列）

使用时间序列节点进行预测

全国宽带提供商的分析员需要对用户订阅进行预测，以便预测带宽的利用率。分析师需要对各地市场进行预测，才能得出全国注册用户数量。分析师将使用时间序列建模来得到未来三个月若干地区市场的预测数字。第二个示例则说明源数据的格式不适合作为时间序列节点的输入时应如何转换源数据。

这两个示例均使用名为 broadband_create_models.str 的流，该流引用名为 broadband_1.sav 的数据文件。这些文件可在任何 IBM® SPSS® Modeler 安装程序的 Demos 文件夹中找到。此目录可通过 WindowsIBM® SPSS® Modeler 程序组进行访问。文件 broadband_create_models.str 位于 streams 文件夹中。

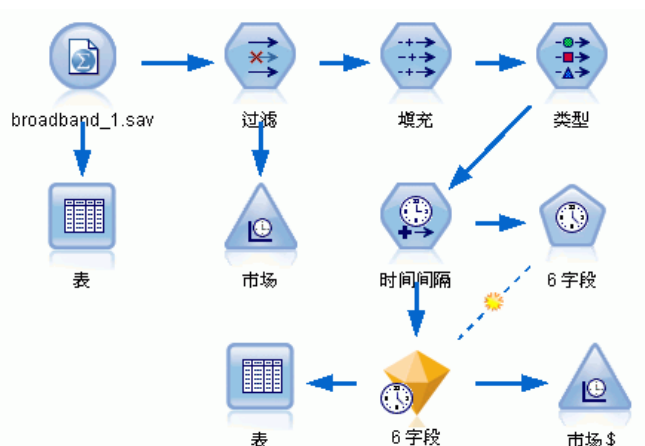
最后一个示例演示如何将保存的模型应用于更新过的数据集，以将预测时间延长三个月。

在 SPSS Modeler 中，可以在单一操作中生成多个时间序列模型。将要使用的源文件具有 85 个不同市场的时间序列数据，但为简便起见，只为其中五个市场以及总体市场的数据建模。

broadband_1.sav 数据文件具有全部 85 个地区市场的月度带宽使用率数据。在本示例中，将只使用前五个序列；并且将为这五个序列各创建一个模型以及为总序列创建一个模型。

该文件还包含指明每个记录的年度和月度的日期字段。此字段将在时间区间节点中用来标记记录。日期字段会以字符串格式读入到 SPSS Modeler 中，但为了在 SPSS Modeler 中使用该字段，必须使用填充节点将存储类型转换为数字日期格式。

图片 14-1
显示时间序列建模的样本流



时间序列节点要求每个序列各占一列，每个区间各占一行。SPSS Modeler 提供用于变换数据的方法，以在必要时使数据符合此格式。

图片 14-2
宽带地区市场的月度预订数据

The screenshot shows a data table window titled "Table (89 个字段, 60 条记录) #1". The table contains 20 columns labeled Market_1 through Market_20 and 20 rows numbered 1 through 20. The data represents monthly booking values for 20 different broadband markets.

	Market_1	Market_2	Market_3	Market_4	Market_5	Market_6	Market_7	Market_8	Market_9	Market_10	Market_11	Market_12	Market_13	Market_14	Market_15	Market_16	Market_17	Market_18	Market_19	Market_20			
1	3750	11489	11659	4571	2205	5488	6144	2363	5041	2352	5802	6670	2469	5231	4010	12801	13716	5211	2490	5899	6929	2574	5401
2	3846	11984	12228	4825	2301	5672	6390	2404	5160	2352	5802	6670	2469	5231	4010	12801	13716	5211	2490	5899	6929	2574	5401
3	3894	12266	12897	5041	2352	5802	6670	2469	5231	4010	12801	13716	5211	2490	5899	6929	2574	5401	2352	5802	6670	2469	5231
4	4010	12801	13716	5211	2490	5899	6929	2574	5401	2352	5802	6670	2469	5231	4010	12801	13716	5211	2490	5899	6929	2574	5401
5	4147	13291	14647	5383	2534	6017	7312	2654	5541	2352	5802	6670	2469	5231	4010	12801	13716	5211	2490	5899	6929	2574	5401
6	4335	13828	15419	5496	2664	6137	7493	2699	5771	2352	5802	6670	2469	5231	4010	12801	13716	5211	2490	5899	6929	2574	5401
7	4554	14273	16108	5747	2738	6250	7702	2786	5901	2352	5802	6670	2469	5231	4010	12801	13716	5211	2490	5899	6929	2574	5401
8	4744	14664	16958	5885	2754	6439	7965	2847	6031	2352	5802	6670	2469	5231	4010	12801	13716	5211	2490	5899	6929	2574	5401
9	4885	15130	17642	6053	2874	6701	8107	2967	6151	2352	5802	6670	2469	5231	4010	12801	13716	5211	2490	5899	6929	2574	5401
10	5020	15851	18453	6229	2975	6957	8366	3099	6341	2352	5802	6670	2469	5231	4010	12801	13716	5211	2490	5899	6929	2574	5401
11	5208	16509	19181	6320	3042	7111	8684	3195	6631	2352	5802	6670	2469	5231	4010	12801	13716	5211	2490	5899	6929	2574	5401
12	5379	17225	19885	6499	3095	7275	8997	3341	6761	2352	5802	6670	2469	5231	4010	12801	13716	5211	2490	5899	6929	2574	5401
13	5574	18173	20565	6593	3199	7380	9326	3376	7021	2352	5802	6670	2469	5231	4010	12801	13716	5211	2490	5899	6929	2574	5401
14	5828	19287	21155	6680	3207	7633	9543	3443	7331	2352	5802	6670	2469	5231	4010	12801	13716	5211	2490	5899	6929	2574	5401
15	5942	20171	21655	6757	3298	7985	9673	3617	7491	2352	5802	6670	2469	5231	4010	12801	13716	5211	2490	5899	6929	2574	5401
16	6139	21379	21964	6804	3387	8236	9934	3732	7711	2352	5802	6670	2469	5231	4010	12801	13716	5211	2490	5899	6929	2574	5401
17	6244	22067	22756	6915	3450	8464	10211	3831	7941	2352	5802	6670	2469	5231	4010	12801	13716	5211	2490	5899	6929	2574	5401
18	6274	23074	23464	7035	3528	8575	10440	3886	8291	2352	5802	6670	2469	5231	4010	12801	13716	5211	2490	5899	6929	2574	5401
19	6347	23729	24324	7151	3546	8817	10763	3938	8581	2352	5802	6670	2469	5231	4010	12801	13716	5211	2490	5899	6929	2574	5401
20	6399	24803	25351	7304	3604	9041	11012	3953	8711	2352	5802	6670	2469	5231	4010	12801	13716	5211	2490	5899	6929	2574	5401

创建流

- ▶ 新建流并添加指向 broadband_1.sav 的 Statistics 文件源节点。
- ▶ 使用过滤节点过滤掉 Market_6 至 Market_85 字段以及 MONTH_ 和 YEAR_ 字段，以简化模型。

提示：要一次选定多个相邻字段，请单击 Market_6 字段，然后按住鼠标左键并向下拖至 Market_85 字段。选定字段将以蓝色突出显示。要添加其他字段，请按住 Ctrl 键，然后单击 MONTH_ 和 YEAR_ 字段。

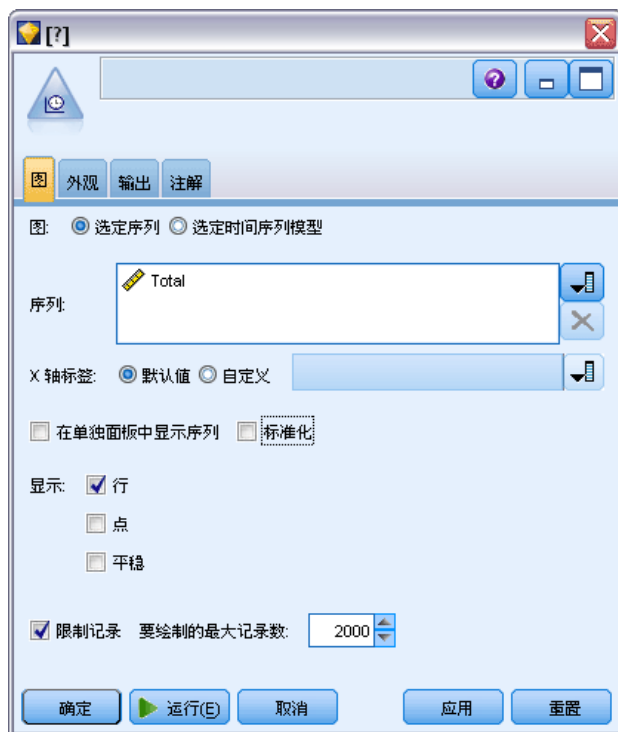
图片 14-3
简化模型



检查数据

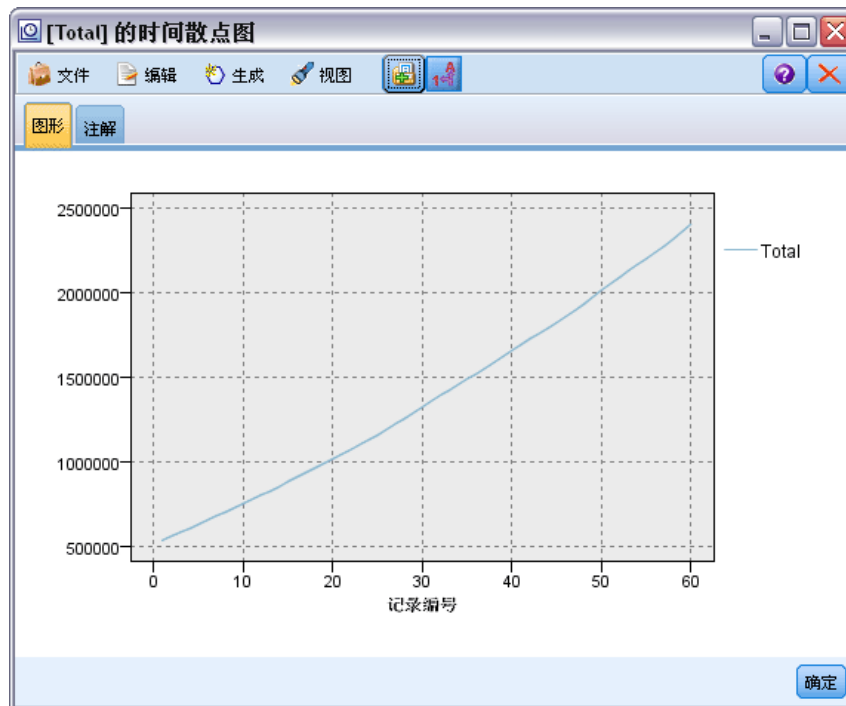
在构建模型前了解数据的性质始终是一种好的做法。数据是否呈现季节性变化？虽然 Expert Modeler 可以自动找出每个序列的最佳季节性或非季节性模型，但是当数据中不存在季节性时，通常可以通过将搜索对象限制为非季节性模型，从而更快速地获得结果。虽然未检查各地区市场的数据，但我们通过标绘这五个市场的总订户数，可大体了解是否存在季节性的因素。

图片 14-4
标绘总订户数



- ▶ 通过“图形”选项板，可将时间散点图节点附加到过滤节点中。
- ▶ 将总计字段添加到“序列”列表。
- ▶ 取消选择在单独面板中显示序列和标准化复选框。
- ▶ 单击运行。

图片 14-5
“总计”字段的时间散点图

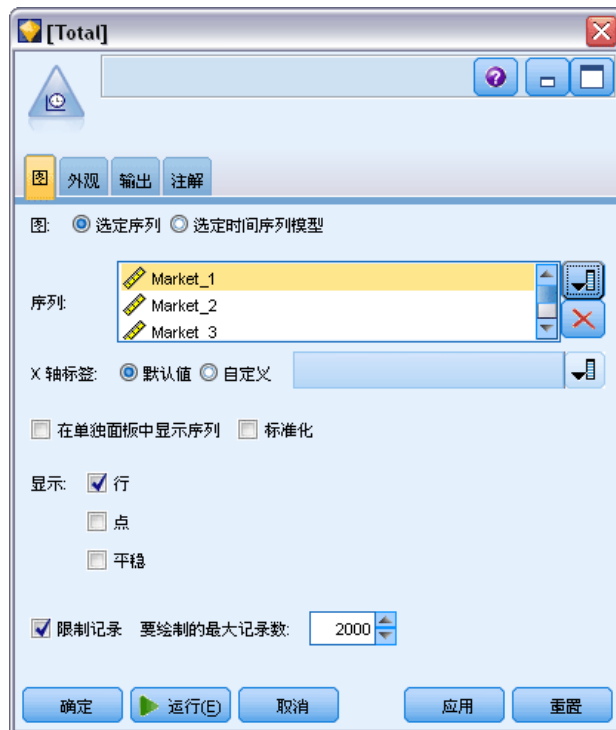


此序列表现出了一个很平滑的上升趋势，没有出现季节性变化的迹象。可能在个别序列中存在季节性，但总的看来，季节性不是数据的显著特征。

当然，应该在排除季节性模型之前检查每个序列。然后分离出表现季节性的序列，并分别对这些序列建模。

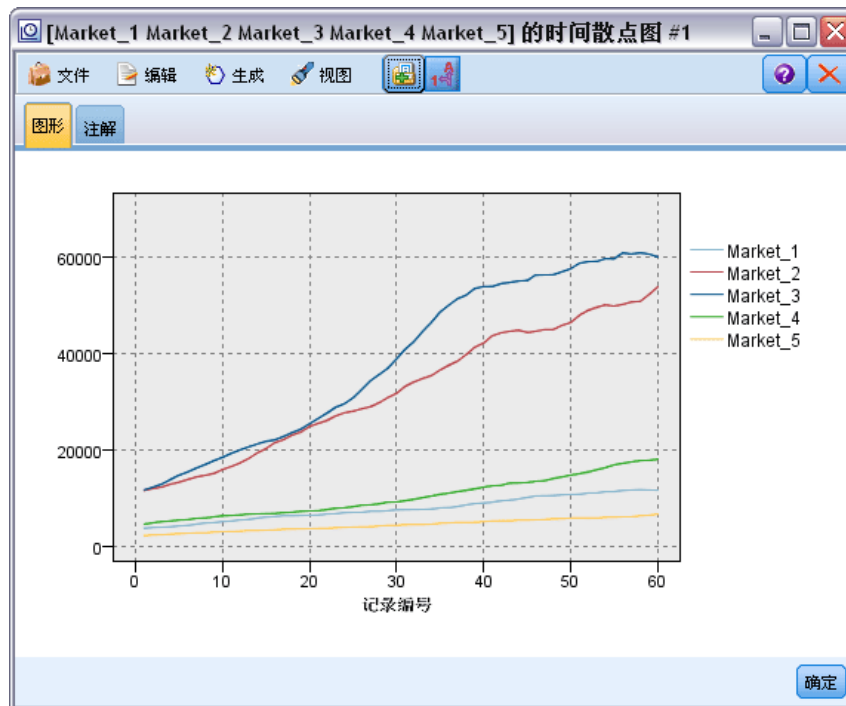
通过 IBM® SPSS® Modeler，可轻松地同时标绘多个序列。

图片 14-6
标绘多个时间序列



- ▶ 重新打开时间散点图节点。
- ▶ 从“序列”列表中删除总计字段（将其选中，然后单击红色的 X 按钮）。
- ▶ 将 Market_1 至 Market_5 字段添加到列表中。
- ▶ 单击运行。

图片 14-7
多个字段的时间散点图



审视各个市场后发现每个市场的曲线均呈稳定上升趋势。虽然一些市场的曲线上升不如其他市场那么稳定，但也未表现出任何季节性趋势。

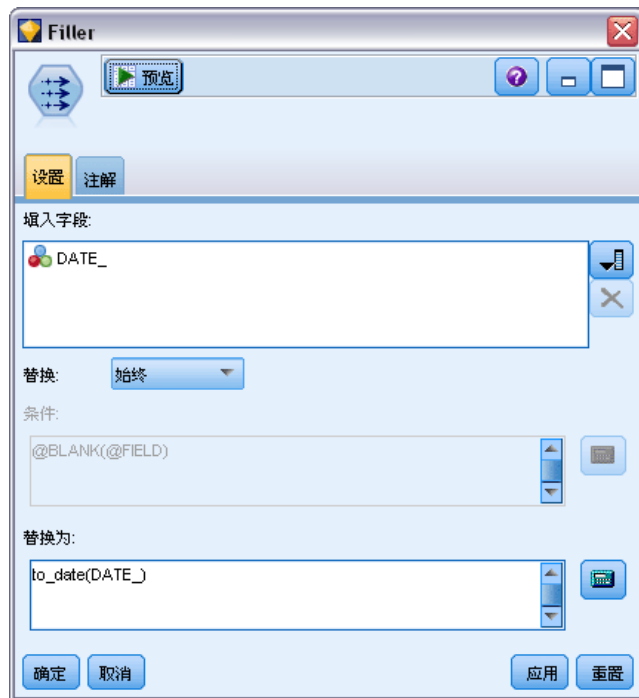
定义日期

现在需要将 DATE_ 字段的存储类型更改为日期格式。

- ▶ 将填充节点附加到过滤节点。
- ▶ 打开填充节点并单击字段选择器按钮。
- ▶ 选择 DATE_ 并将它添加到填入字段。
- ▶ 将替换条件设置为始终。

- ▶ 将替换为的值设置为 `to_date(DATE_)`。

图片 14-8
设置日期存储类型

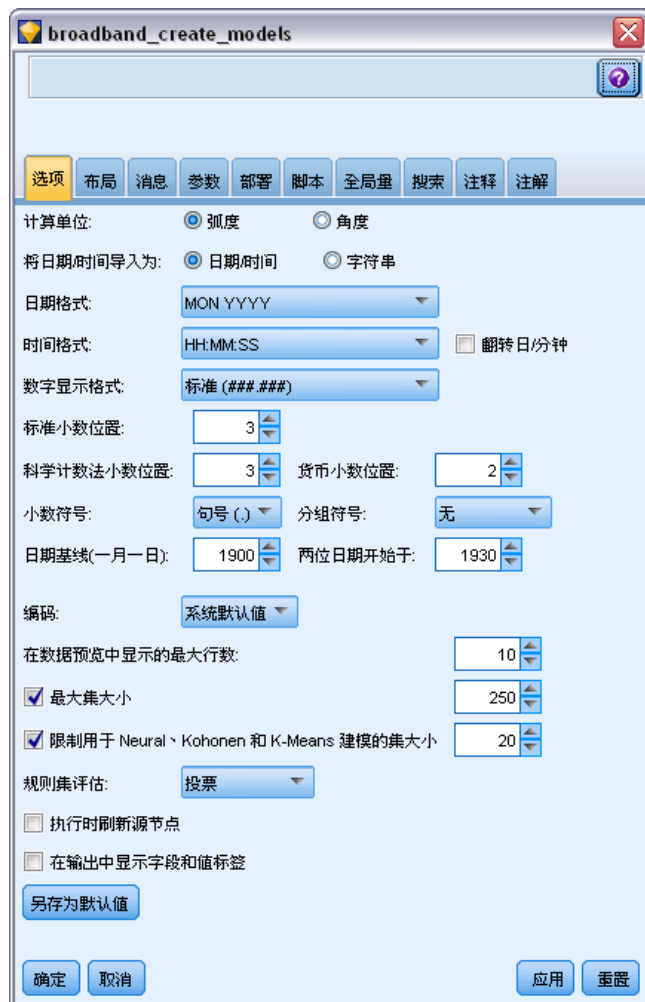


更改默认日期格式以匹配“日期”字段的格式。这对于要按预期转换“日期”字段是必需的。

- ▶ 在菜单上，选择工具 > 流属性 > 选项，以显示“流选项”对话框。

- ▶ 将默认日期格式设置为 MON YYYY。

图片 14-9
设置日期格式



定义目标

- ▶ 添加“类型”节点并将 DATE_ 字段的角色设置为无。将所有其他字段（Market_n 字段以及合计字段）的角色设置为目标。

- ▶ 单击读取值按钮以填充“值”列。

图片 14-10
为多个字段设置角色

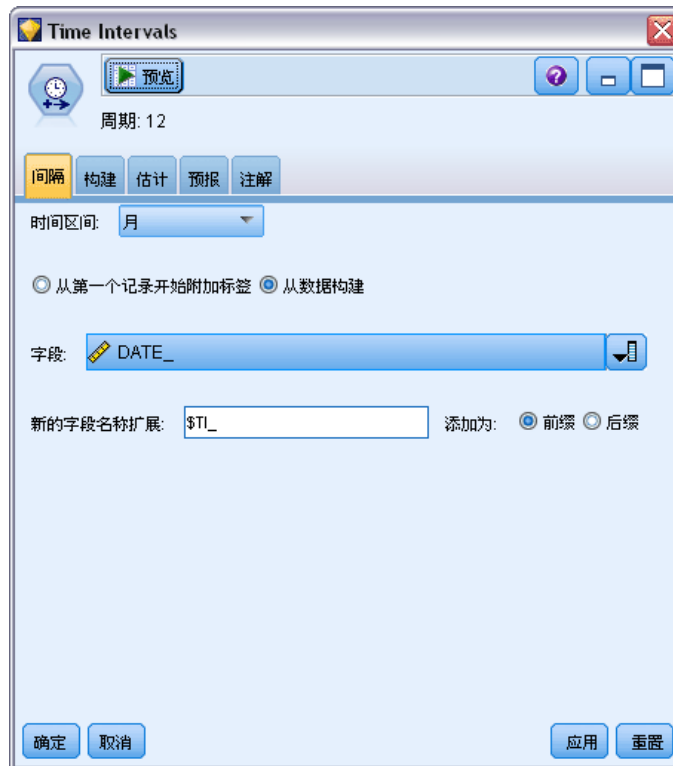


设置时间区间

- ▶ 添加“时间区间”节点（通过“字段操作”选项板）。
- ▶ 在“区间”选项卡上，选择月作为时间区间。
- ▶ 选中根据数据构建选项。

- ▶ 选择 DATE_ 作为构建字段。

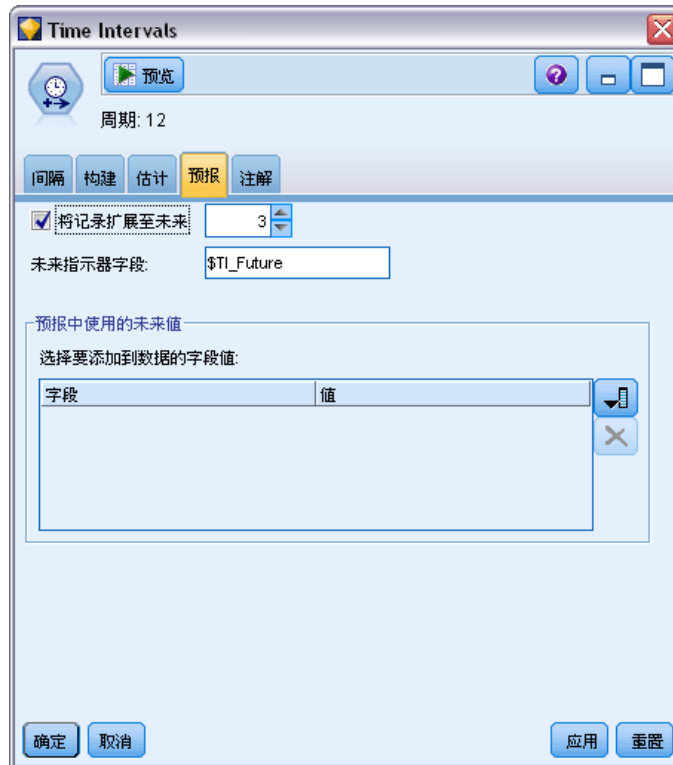
图片 14-11
设置时间区间



- ▶ 在“预测”选项卡上，选中将记录扩展到未来复选框。
- ▶ 将值设置为 3。

- ▶ 单击确定。

图片 14-12
设置预测时限

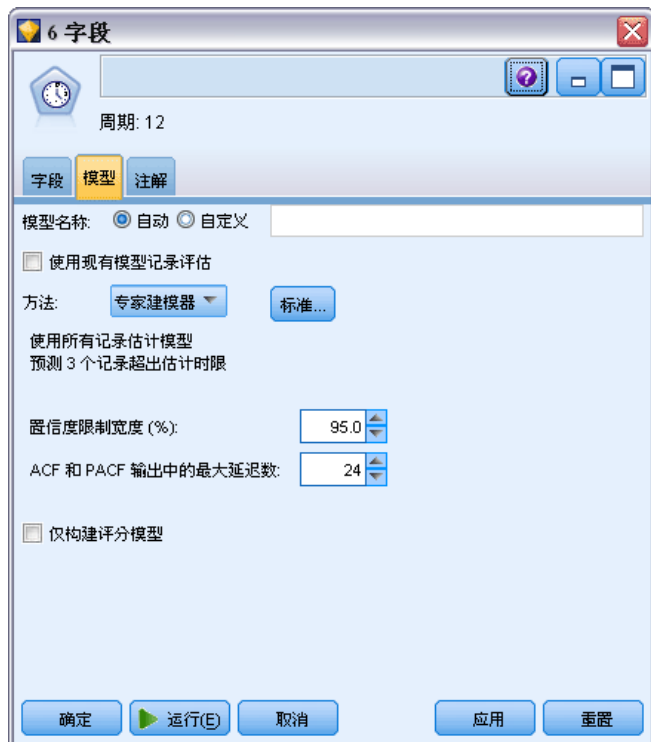


创建模型

- ▶ 从“建模”选项板中，将一个时间序列节点添加到流，并将它附加到时间区间节点。

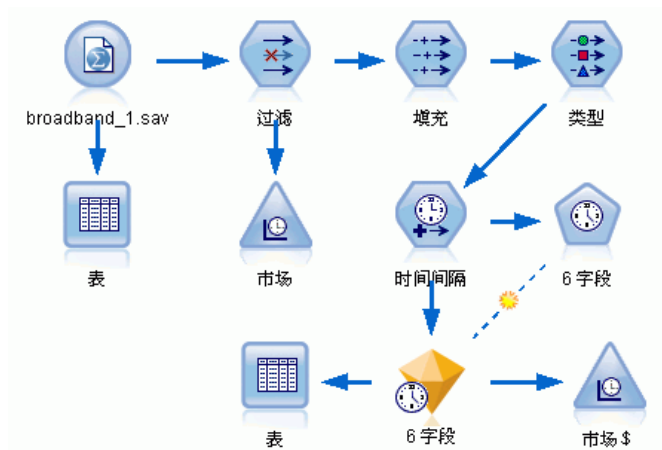
- ▶ 单击使用全部默认设置的时间序列节点上的运行。此操作使 Expert Modeler 能够将最合适的模型用于每个时间序列。

图片 14-13
为时间序列选择 Expert Modeler



- ▶ 将时间序列模型块附加到时间区间节点。
- ▶ 将“表”节点附加到时间序列模型并单击运行。

图片 14-14
显示时间序列建模的样本流



现在有三个新行（第 61 至 63 行）附加到原始数据中。这三行用于预测时限，在本例中为 2004 年 1 至 3 月。

现在还有几个新列，即时间区间节点添加的若干 \$TI_ 列和时间序列节点添加的若干 \$TS_ 列。这些列表示每行（也就是时间序列数据中的每个区间）的以下内容：

Column	描述
\$TI_TimeIndex	此行的时间区间索引值。
\$TI_TimeLabel	此行的时间区间标签。
\$TI_Year	此行中生成数据的年份和月份指示符。
\$TI_Month	
\$TI_Count	确定此行的新数据时所涉及记录的数量。
\$TI_Future	指明此行是否包含预测数据。
\$TS-colname	由每列原始数据生成的模型数据。
\$TSLCI-colname	每列生成的模型数据中的置信区间下限值。
\$TSUCI-colname	每列生成的模型数据中的置信区间上限值。
\$TS-Total	此行中 \$TS-colname 的总值。
\$TSLCI-Total	此行中 \$TSLCI-colname 的总值。
\$TSUCI-Total	此行中 \$TSUCI-colname 的总值。

对预测操作最重要的列是 \$TS-Market_n、\$TSLCI-Market_n 和 \$TSUCI-Market_n。特别是这些列的第 61 至 63 行，它们包含各个地区市场的用户预订预测数据和置信区间。

检查模型

- ▶ 双击时间序列模型块，以显示有关为每个市场生成的模型的数据。

请注意 Expert Modeler 如何选择通过为其他市场生成的类型来为市场 5 生成不同类型的模型。

图片 14-15
为市场生成的时间序列模型

6 字段

文件 生成 预览

模型 参数 残差 汇总 设置 注解

排序方式 已选择 视图: 简单

评估中使用的记录数:60

	目标	模型	预测变量	固定 R**2	Q	df	Sig.
<input checked="" type="checkbox"/>	Market_1	Holts 线性趋势	0	0.264	8.53	16.0	0.931
<input checked="" type="checkbox"/>	Market_2	Holts 线性趋势	0	0.121	35.9	16.0	0.003
<input checked="" type="checkbox"/>	Market_3	Holts 线性趋势	0	0.258	15.76	16.0	0.47
<input checked="" type="checkbox"/>	Market_4	Holts 线性趋势	0	0.25	27.714	16.0	0.034
<input checked="" type="checkbox"/>	Market_5	Winters 加法	0	0.544	11.888	15.0	0.688
<input checked="" type="checkbox"/>	Total	Holts 线性趋势	0	0.049	27.616	16.0	0.035

汇总统计

	统计	固定 R**2	Q	df	Sig.
SUMMARY	MEAN	0.247	21.235	15.833	0.36
SUMMARY	SE	0.169	10.738	0.408	0.396
SUMMARY	MINIMUM	0.049	8.53	15	0.003
SUMMARY	MAXIMUM	0.544	35.9	16	0.931
SUMMARY	PERCENTILE 5	0.049	8.53	15	0.003
SUMMARY	PERCENTILE ...	0.049	8.53	15	0.003
SUMMARY	PERCENTILE ...	0.103	11.048	15.75	0.026
SUMMARY	PERCENTILE ...	0.254	21.688	16	0.252
SUMMARY	PERCENTILE ...	0.334	29.761	16	0.749
SUMMARY	PERCENTILE ...	0.544	35.9	16	0.931
SUMMARY	PERCENTILE ...	0.544	35.9	16	0.931

确定 取消 应用 重置

“预测变量”列显示有多少个字段作为每个目标的预测变量—在本例中为 0。

此视图中余下的列显示的是每个模型的拟合优度测量值。StationaryR**2 列显示的是固定的 R 平方值。此统计量是序列中由模型解释的总变异所占比例的估计值。该值越高（最大值为 1.0），则模型拟合会越好。

Q、df 和 Sig. 列与 Ljung-Box 统计量相关联，该检验是对模型中残差错误的随机检验；错误的随机性越大，则模型会变得越好。Q 是 Ljung-Box 统计量本身，而 df（自由度）则表示评估特定目标时可任意更改的模型参数数量。

Sig. 列给出了 Ljung-Box 统计量的显著性值，从而以另一种方式来表示指定的模型是否正确。显著性值小于 0.05 表示残差误差不是随机的，则意味着所观测的序列中存在模型无法解释的结构。

如果将固定 R 平方值和显著性值考虑在内，则 Expert Modeler 为 Market_1、Market_3 和 Market_5 选择的模型完全可以接受。Market_2 和 Market_4 的 Sig. 值均小于 0.05，表明可能还必须进行一些实验，以便为这些市场找到拟合度更好的模型。

屏幕下方的汇总值提供了有关这些统计量在所有模型中的分布情况的信息。例如，所有模型的固定 R 平方均值为 0.247，而此值的最小值为 0.049（总计 模型的值），最大值为 0.544（Market_5 的值）。

SE 表示每个统计量在所有模型中标准误。例如，固定 R 平方值在所有模型中的标准误差为 0.169。

汇总部分还包括百分位数值，它们提供有关统计量在模型中的分布情况的信息。对于每个百分位，模型百分比具有一个小于规定值的拟合统计量值。

例如，仅 25% 的模型的固定 R 平方值低于 0.121。

- ▶ 单击“视图”下拉列表并选择高级。

屏幕上即会显示若干其他拟合优度测量值。 R^2 为 R 平方值，即，时间序列中可由模型解释的总变异估计值。由于此统计量的最大值为 1.0，因此在这一点上，我们的模型表现不错。

图片 14-16
时间序列模型的高级显示

评估中使用的记录数:60

	MAPE	MAE	MaxAPE	MaxAE	标准化BIC	Q	df	Sig.
17	0.94	73.869	2.147	224.517	9.15	8.53	16.0	0.931
76	0.94	314.721	1.867	927.949	12.059	35.9	16.0	0.003
33	0.776	306.877	1.918	1,030.105	12.1	15.76	16.0	0.47
38	0.78	79.49	1.942	233.544	9.329	27.714	16.0	0.034
32	0.936	39.963	2.481	137.633	8.114	11.888	15.0	0.688
74	0.094	1,326.071	0.299	7,062.662	15.243	27.616	16.0	0.035

汇总统计

MAPE	MAE	MaxAPE	MaxAE	标准化BIC	Q	df	Sig.
0.744	356.832	1.776	1,602.735	10.999	21.235	15.833	0.36
0.328	490.119	0.758	2,702.397	2.641	10.738	0.408	0.396
0.094	39.963	0.299	137.633	8.114	8.53	15	0.003
0.94	1,326.071	2.481	7,062.662	15.243	35.9	16	0.931
0.094	39.963	0.299	137.633	8.114	8.53	15	0.003
0.094	39.963	0.299	137.633	8.114	8.53	15	0.003
0.605	65.393	1.475	202.796	8.891	11.048	15.75	0.026
0.858	193.183	1.93	580.747	10.694	21.688	16	0.252
0.94	567.559	2.231	2,538.245	12.886	29.761	16	0.749
0.94	1,326.071	2.481	7,062.662	15.243	35.9	16	0.931
0.94	1,326.071	2.481	7,062.662	15.243	35.9	16	0.931

RMSE 是指均方根误差，它是一种测量序列实际值与模型预测值之间差异的度量方法，采用与序列本身所用的相同单位表示。由于这是误差测量值，因此我们希望该值应尽可能地低。乍一看，Market_2 和 Market_3 的模型的成功率低于其他三个市场，但根据目前所观察的统计量仍可以接受。

其他的拟合优度测量值包括均值绝对百分比误差 (MAPE) 和其最大值 (MaxAPE)。绝对百分比误差用于度量目标序列与其模型预测水平的差异程度，用百分比值表示。通过检查所有模型的均值和最大值，可获得预测中不确定性的指示。

MAPE 值显示所有模型的平均不确定性都低于 1%，这是一个很低的数字。MaxAPE 值显示最大绝对百分比误差，对设想预测的最坏情形很有帮助。它显示，每个模型的最大百分比误差大约在 1.8% 至 2.5% 之间，仍然是一组很低的数字。

MAE（绝对平均误差）值用于显示预测误差绝对值的均值。如 RMSE 值，使用与序列本身所用的相同单位表示。MaxAE 用于以相同的单位显示最大预测误差，并指示预测的最坏情形。

尽管关注的是这些绝对值，但由于目标序列表示的是不同规模市场的订户量，因此在此情况下百分比误差值（MAPE 和 MaxAPE）更有用。

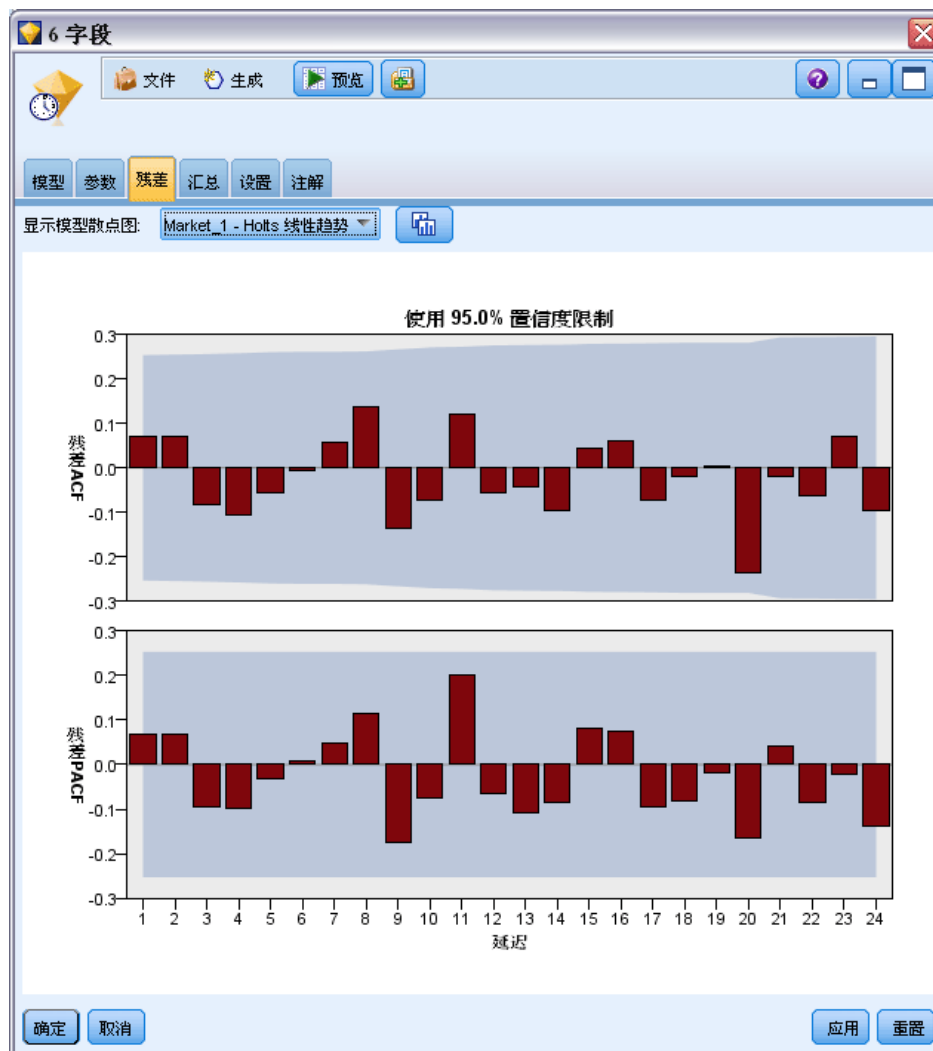
MAPE 和 MaxAPE 值可以使用模型显示可接受的不确定性吗？这些值确实很低。由于可接受风险因问题的不同而异，因此商业意识可以在此派上用场。假设拟合优度统计量在可接受的范围之内，然后继续查看残差错误。

检查模型残差的自相关函数（ACF）和偏自相关函数（PACF）的值比只查看拟合优度统计量能更多地从量化角度来了解模型。

合理指定的时间模型将捕获所有非随机的变异，其中包括季节性、趋势、循环周期以及其他重要的因素。如果是这种情况，则任何误差都不会随着时间的推移与其自身相关联（自关联）。这两个自相关函数中的显著结构都可以表明基础模型不完整。

- ▶ 单击“残差”选项卡以显示第一个地区市场模型中残差错误的自相关函数（ACF）值和偏自相关函数（PACF）值。

图片 14-17
市场的 ACF 值和 PACF 值



在这些散点图中，为了解误差变量的原始值是否与时间变化相关，这些值已滞后多达 24 个时限并与原始值进行了比较。对于可接受的模型，上（ACF）散点图中的条形块在正（上）方向或负（下）方向均不应扩展到阴影区之外。

如果出现此情况，您需要检查下（PACF）散点图，以了解是否已确认此处的结构。PACF 散点图主要关注在控制插入时间点的序列值之后的相关性。

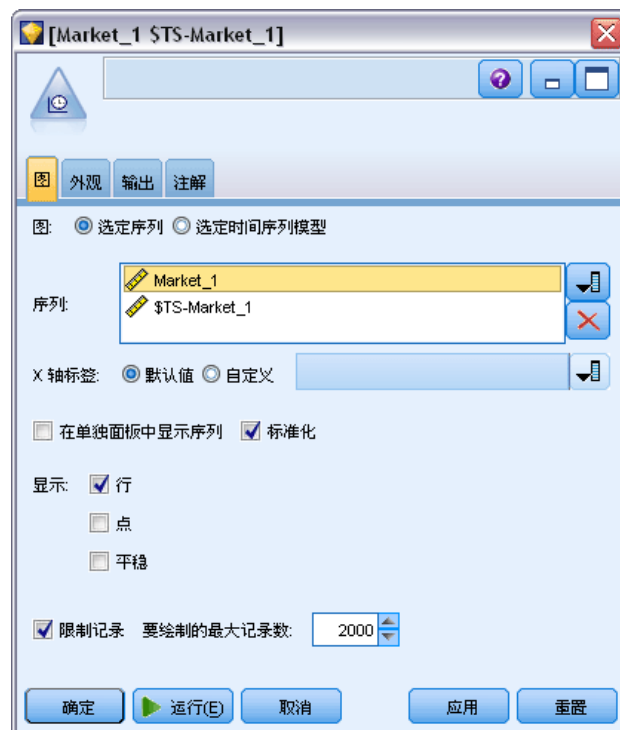
Market_1 的值都位于阴影区之内，因此我们可以继续检查其他市场的值。

- ▶ 单击显示模型散点图下拉列表可显示其他地区市场和总体市场的 ACF 值和 PACF 值。

由于 Market_2 和 Market_4 的值不太重要，因此可以确认原先对 Sig. 值所产生的可疑因素。我们需要在某些时间试验这些市场的其他一些不同的模型，以了解是否可以获取更佳的拟合，但对于该示例的剩余部分，应更多地考虑可以从 Market_1 模型中获取的其他内容。

- ▶ 通过“图形”选项板，可将时间散点图节点附加到时间序列模型块中。
- ▶ 在“散点图”选项卡上，取消选中在单独面板中显示系列复选框。
- ▶ 在系列列表上，单击字段选择器按钮，选定 Market_1 和 \$TS-Market_1 字段，然后单击确定将它们添加到列表中。
- ▶ 单击运行，以显示第一个地区市场的实际数据和预测数据的线图。

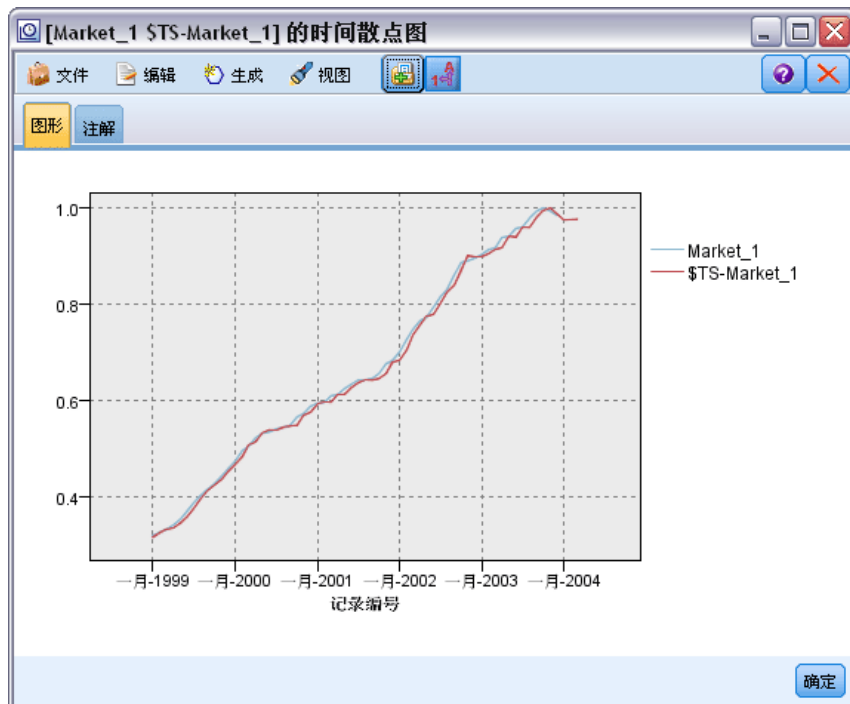
图片 14-18
选择要标绘的字段



请注意预测（\$TS-Market_1）线如何通过实际数据的末端向外延伸。现已得出对此市场未来三个月预期需求的预测。

整个时间序列上的实际数据线和预测数据线在图上非常接近，表明对此特定时间序列这是一个可靠的模型。

图片 14-19
Market_1 的实际数据和预测数据的时间散点图



将模型保存在文件中，以便在将来的示例中使用：

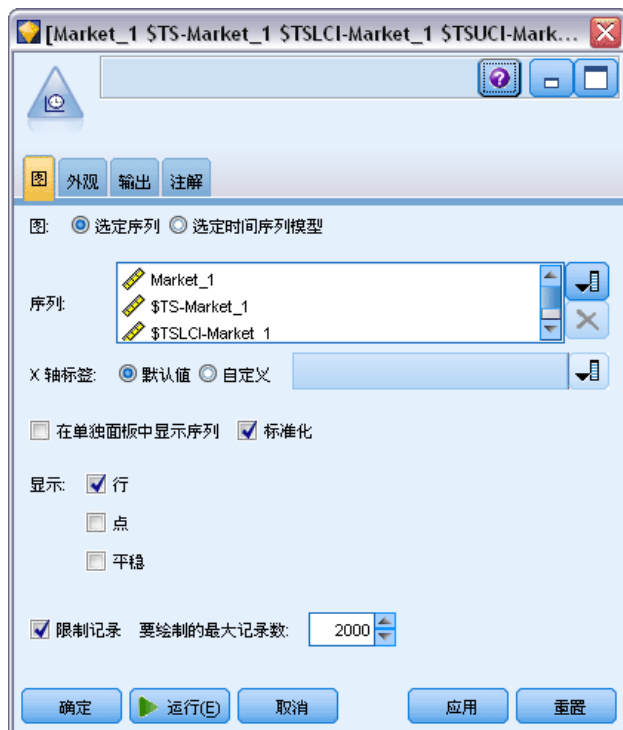
- ▶ 单击确定关闭当前图形。
- ▶ 打开时间序列模型块。
- ▶ 选择文件 > 保存节点并指定文件位置。
- ▶ 单击保存。

现在虽然有了此特定市场的可靠模型，但该预测的误差到底有多大呢？可通过检查置信区间得到预测的误差大小。

- ▶ 双击流中最后的时间散点图（标注为 `Market_1 $TS-Market_1`），以重新打开该节点的对话框。
- ▶ 单击字段选择器按钮并将 `$TSLCI-Market_1` 和 `$TSUCI-Market_1` 字段添加到系列列表中。

- ▶ 单击运行。

图片 14-20
添加更多要标绘的字段

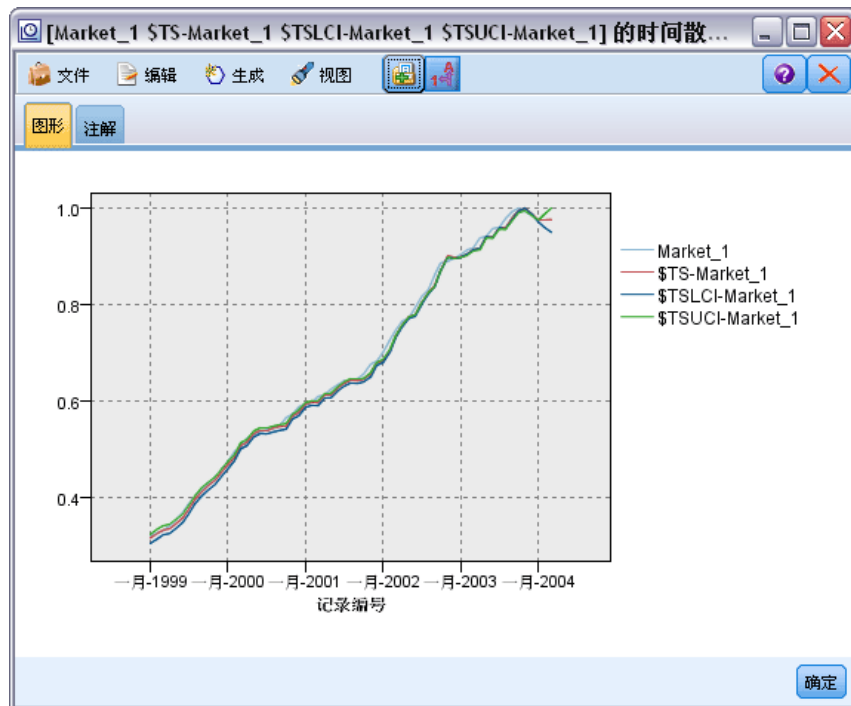


现在有了与以前一样的图形，但添加了置信区间上限（\$TSUCI）和下限（\$TSLCI）。

请注意置信区间的边界如何随预测时限而分叉，这表示预测越指向更远的将来，不确定性就變得越来越大。

但是，随着每个时段的流逝，您就会多一个时段（在本例中为月）的实际使用率数据作为预测的依据。您可以将这些新数据读入到流中，并再次应用您的模型，因为您知道它是可靠的。[有关详细信息，请参阅第 187 页码重新应用时间序列模型。](#)

图片 14-21
添加了置信区间的时间散点图



摘要

您已学习了如何使用 Expert Modeler 为多个时间序列生成预测，并已将得到的模型保存到外部文件中。

在下一个示例中，您将看到如何将非标准时间序列数据转换为适合输入到时间序列节点的格式。

重新应用时间序列模型

本示例将应用来自第一个时间序列示例但也可以独立使用的时间序列模型。[有关详细信息，请参阅第 165 页码使用时间序列节点进行预测。](#)

与原来的情形一样，全国宽带提供商的一位分析师为了预测带宽需求，需要为若干地区市场中的每一个市场得出用户预订的月度预测数据。您已经使用专家建模器创建了模型，要对未来的三个月进行预测。

您的数据仓库现已使用原来预测时限的实际数据进行了更新，因而您想使用这些数据对另外三个月做出预测。

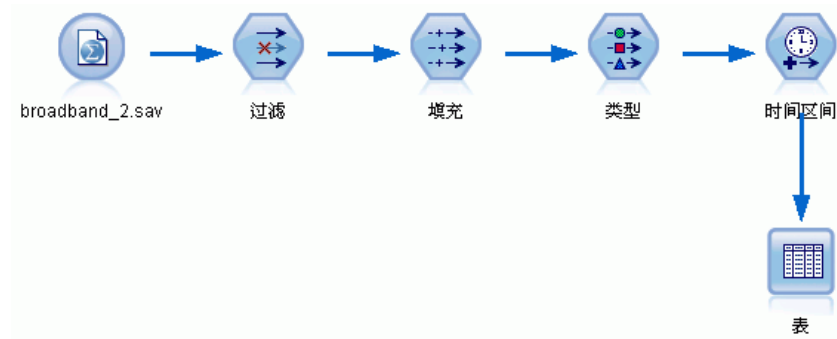
此示例使用名为 broadband_apply_models.str 的流，该流引用名为 broadband_2.sav 的数据文件。这些文件可在任何 IBM® SPSS® Modeler 安装程序的 Demos 文件夹中找到。此目录可通过 WindowsIBM® SPSS® Modeler 程序组进行访问。文件 broadband_apply_models.str 位于 streams 文件夹中。

检索流

在此示例中，将从保存在第一个示例中的时间序列模型中重新创建一个时间序列节点。如果未保存模型也不用担心，因为 Demos 文件夹中提供了一个模型。

- ▶ 从 Demos 下的 streams 文件夹中打开流 broadband_apply_models.str。

图片 14-22
打开流



图片 14-23
更新的销售数据

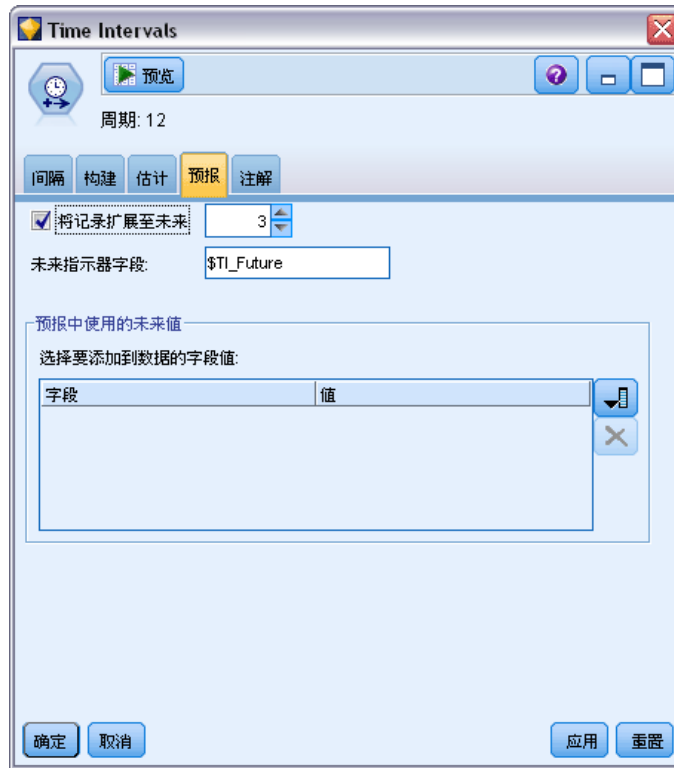
	Market_82	Market_83	Market_84	Market_85	Total	YEAR_	MONTH_	DATE_
44	58820	20482	14326	16935	17917...	2002	8	AUG 2002
45	60119	21211	14349	17179	18249...	2002	9	SEP 2002
46	61320	21893	14333	17601	18601...	2002	10	OCT 2002
47	63099	22471	14229	17816	18945...	2002	11	NOV 2002
48	64687	23112	14514	17937	19343...	2002	12	DEC 2002
49	65518	23686	14856	18003	19752...	2003	1	JAN 2003
50	65570	24669	15182	17875	20148...	2003	2	FEB 2003
51	66567	25469	15709	18214	20540...	2003	3	MAR 2003
52	67527	25868	16155	18557	20922...	2003	4	APR 2003
53	67724	26284	16521	19190	21300...	2003	5	MAY 2003
54	68644	26468	16567	19938	21669...	2003	6	JUN 2003
55	69878	26781	16618	20876	22004...	2003	7	JUL 2003
56	71538	27566	16553	21514	22398...	2003	8	AUG 2003
57	73162	28164	16597	21779	22773...	2003	9	SEP 2003
58	74167	28693	16669	22266	23160...	2003	10	OCT 2003
59	76036	28922	16748	22559	23616...	2003	11	NOV 2003
60	76630	29811	16798	23018	24067...	2003	12	DEC 2003
61	79002	30034	17122	23160	24509...	2004	1	JAN 2004
62	81123	30091	17581	23698	24968...	2004	2	FEB 2004
63	83909	30162	17894	24355	25383...	2004	3	MAR 2004

更新过的月度数据收集在 broadband_2. sav 中。

- ▶ 将表节点附加到 IBM® SPSS® Statistics 文件源节点，打开表节点并单击运行。
注意：数据文件已经用 2004 年 1 至 3 月份（第 61 至 63 行）的实际销售数据进行了更新。
- ▶ 打开流上的时间区间节点。
- ▶ 单击预测选项卡。

- 确保将记录扩展到未来设置为 3。

图片 14-24
检查预测时限的设置

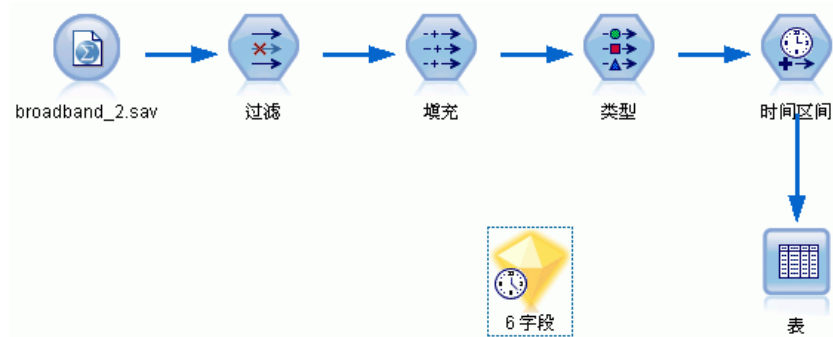


检索保存的模型

- 在 IBM® SPSS® Modeler 菜单上，选择插入 > 来自文件的节点，并从 Demos 文件夹中选择 TSmodel.nod 文件（或使用在第一个时间序列示例中保存的时间序列模型）。

此文件包含来自上一个示例的时间序列模型。插入操作将把对应的时间序列模型块放置在工作区上。

图片 14-25
添加模型块

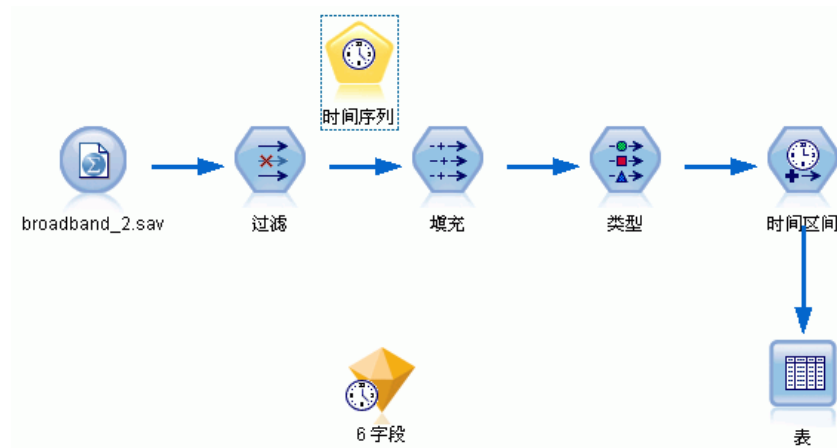


生成建模节点

- ▶ 打开时间序列模型块并选择生成 > 生成建模节点。

此操作将把时间序列建模节点放置在工作区上。

图片 14-26
从模型块生成建模节点



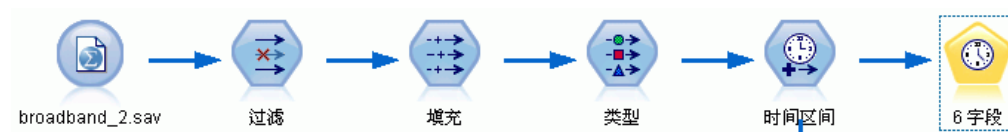
生成新模型

- ▶ 关闭时间序列模型块并将它从工作区中删除。

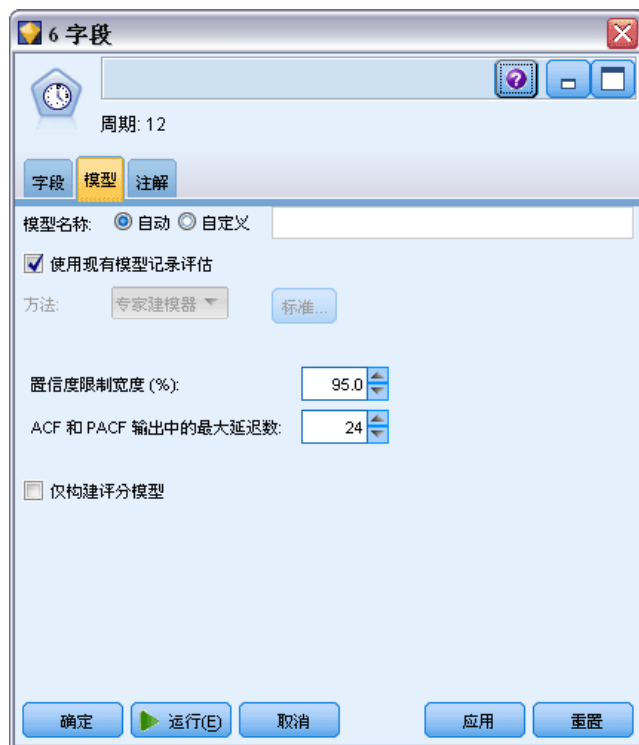
旧模型以 60 行数据为基础构建。现在需要基于更新过的销售数据（63 行）来生成新模型。

- ▶ 将新生成的时间序列构建节点附加到流。

图片 14-27
将建模节点附加到流中



图片 14-28
重新使用已为时间序列模型存储的设置



- ▶ 打开时间序列节点。
- ▶ 在模型选项卡上，请确保已选中使用现有模型继续评估。
- ▶ 单击运行以将新的模型块放在工作区及“模型”选项板中。

检查新模型

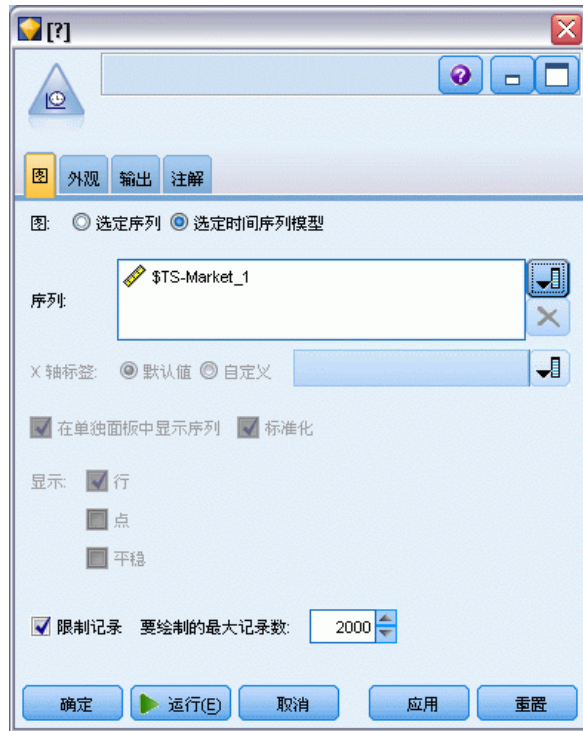
图片 14-29
显示新预测的表

	\$TI_TimeLabel	\$TI_Year	\$TI_Month	\$TI_Count	\$TI_Future	\$TS-Market_1	\$TSLCI-Market_1
47	十一月 2002	2002	11	1	0	10552	10365
48	十二月 2002	2002	12	1	0	10593	10406
49	一月 2003	2003	1	1	0	10653	10466
50	二月 2003	2003	2	1	0	10740	10553
51	三月 2003	2003	3	1	0	10851	10664
52	四月 2003	2003	4	1	0	10909	10722
53	五月 2003	2003	5	1	0	11153	10966
54	六月 2003	2003	6	1	0	11178	10991
55	七月 2003	2003	7	1	0	11382	11195
56	八月 2003	2003	8	1	0	11408	11221
57	九月 2003	2003	9	1	0	11627	11440
58	十月 2003	2003	10	1	0	11795	11608
59	十一月 2003	2003	11	1	0	11869	11682
60	十二月 2003	2003	12	1	0	11793	11607
61	一月 2004	2004	1	1	0	11686	11500
62	二月 2004	2004	2	1	0	11896	11710
63	三月 2004	2004	3	1	0	11996	11810
64	四月 2004	2004	4	0	1	12278	12056
65	五月 2004	2004	5	0	1	12416	12100
66	六月 2004	2004	6	0	1	12553	12167

- ▶ 将表节点附加到工作区中新的时间序列模型块。
- ▶ 打开“表”节点，然后单击运行。

由于使用的是已存储的设置，因此新模型仍然向前预测三个月的需求。不过，此次的预测时限为 4 至 6 月，因为估计时限（在时间区间节点上指定）现在是在 3 月而不是 1 月结束。

图片 14-30
指定要标绘的字段

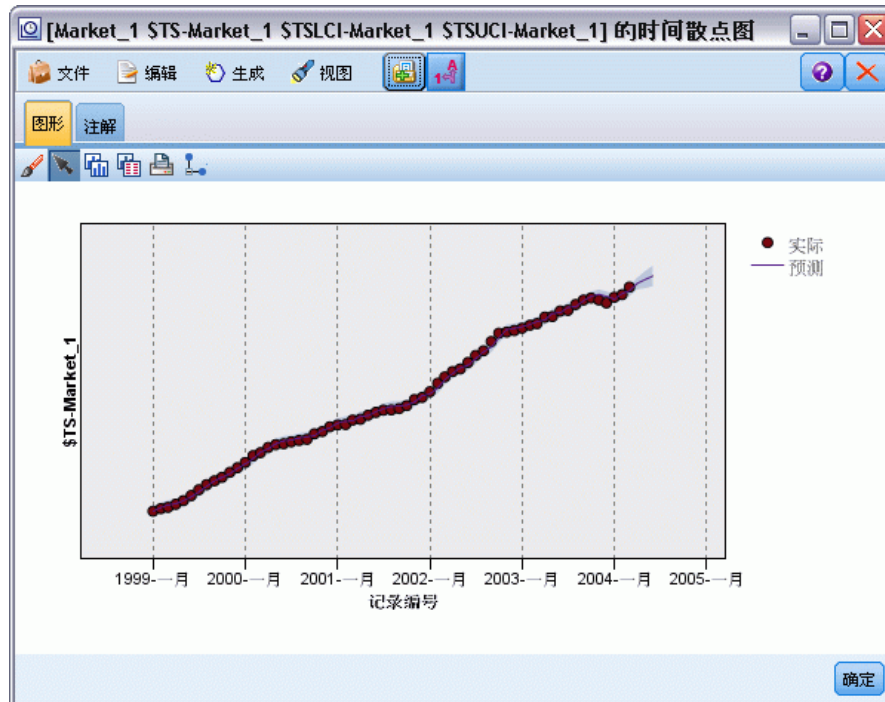


- ▶ 将时间散点图节点附加到时间序列模型块。
此次将使用专为时间序列模型设计的时间散点图显示方式。
- ▶ 在“散点图”选项卡上，选择选定的时间序列模型选项。
- ▶ 在系列列表上，单击字段选择器按钮，选定 \$TS-Market_1 字段，然后单击确定将它添加到列表中。
- ▶ 单击运行。

现在的图形显示的是截止 2004 年 3 月 Market_1 的实际销售量，以及截止 2004 年 6 月的预测销售量及置信区间（蓝色阴影区）。

与第一个示例中一样，在整个预测时限内，预测值与实际数据贴得很紧，再次表明构建的是一个比较理想的模型。

图片 14-31
扩展到 6 月的预测



摘要

上面学习了有更多当前数据可用时如何应用保存的模型以扩展以前的预测，并可在无需重新构建模型的情况下实现这一点。当然，如果有理由认为模型已改变，则应重新构建模型。

预测产品分类销售情况（时间序列）

一个产品分类销售公司会根据过去 10 年的销售数据来预测其男装生产线的月销售情况。

此示例使用名为 catalog_forecast.str 的流，此流引用名为 catalog_seasfac.sav 的数据文件。这些文件可在任何 IBM® SPSS® Modeler 安装程序的 Demos 目录中找到。此目录可通过 WindowsIBM® SPSS® Modeler 程序组进行访问。文件 catalog_forecast.str 位于 streams 目录中。

在前一个示例中我们了解了如何通过 Expert Modeler 找到最适合您的时间序列的模型。现在我们来深入了解两种可用于选择模型的方法 — 指数平滑与 ARIMA。

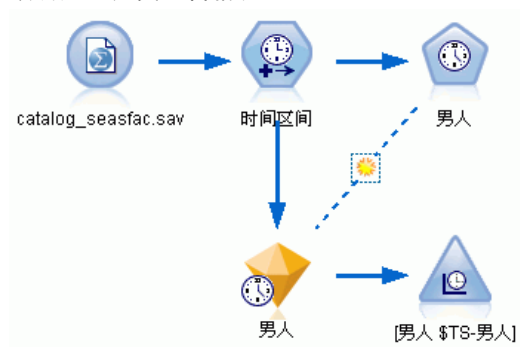
为了帮您找到适当的模型，最好先绘制时间序列。时间序列的可视化检查通常可以很好地指导并帮助您进行选择。另外，您需要弄清以下几点：

- 此序列是否存在整体趋势？如果是，趋势是显示持续存在还是显示将随时间而消逝？
- 此序列是否显示季节变化？如果是，那么这种季节的波动是随时间而加剧还是持续稳定存在？

创建流

- ▶ 新建流并添加指向 catalog_seasfac.sav 的 Statistics 文件源节点。

图片 15-1
预测产品分类销售情况

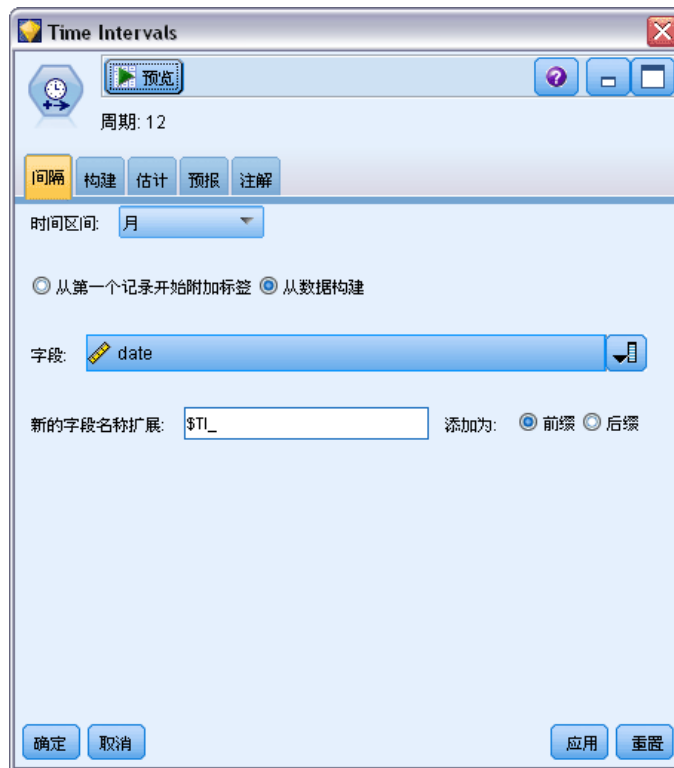


图片 15-2
指定目标字段



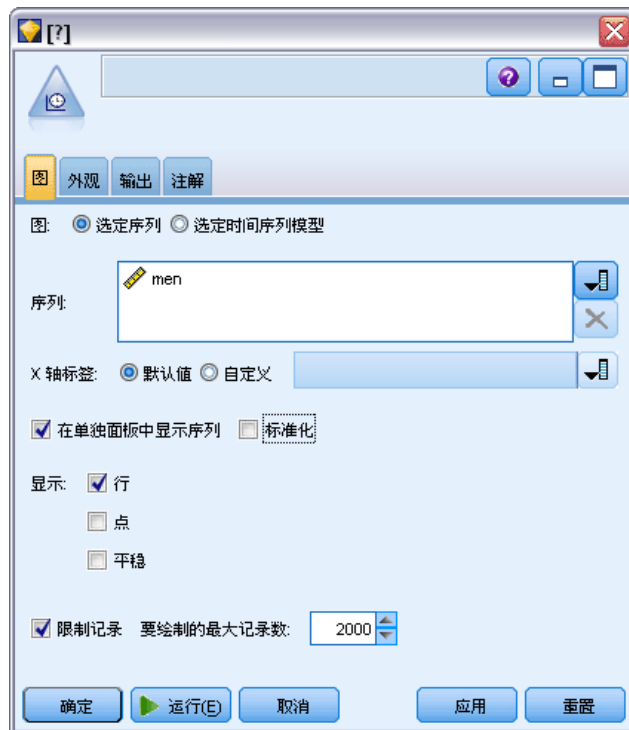
- ▶ 打开 IBM® SPSS® Statistics 文件源节点并选择“类型”选项卡。
- ▶ 单击读取值，然后单击确定。
- ▶ 单击角色列（在男字段中），将角色设置为目标。
- ▶ 将所有其他字段的角色设置为无，然后单击确定。

图片 15-3
设置时间区间



- ▶ 将时间区间节点添加到 SPSS Statistics 文件源节点。
- ▶ 打开时间区间节点，然后将时间区间设置为月。
- ▶ 选择从数据构建。
- ▶ 将字段设置为日期，然后单击确定。

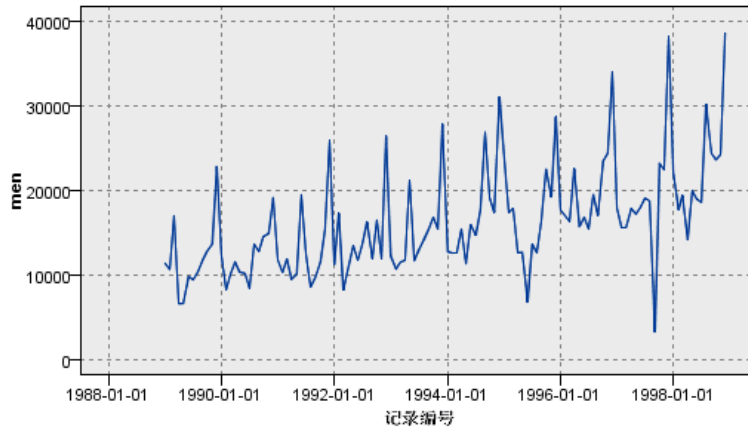
图片 15-4
绘制时间序列



- ▶ 将“时间散点图”节点添加到“时间区间”节点。
- ▶ 在“散点图”选项卡上，将男添加到序列列表。
- ▶ 取消选择标准化复选框。
- ▶ 单击运行。

检查数据

图片 15-5
男装的实际销售情况



此序列显示整体上升趋势，即序列值随时间而增加。上升趋势似乎将持续，即为线性趋势。

此序列还有一个明显的季节特征，即年度高点在十二月（根据图表中的垂直线可以看出）。季节变化显示随上升序列而增长的趋势，表明是乘法季节模型而不是加法季节模型。

- ▶ 单击确定以关闭此散点图。

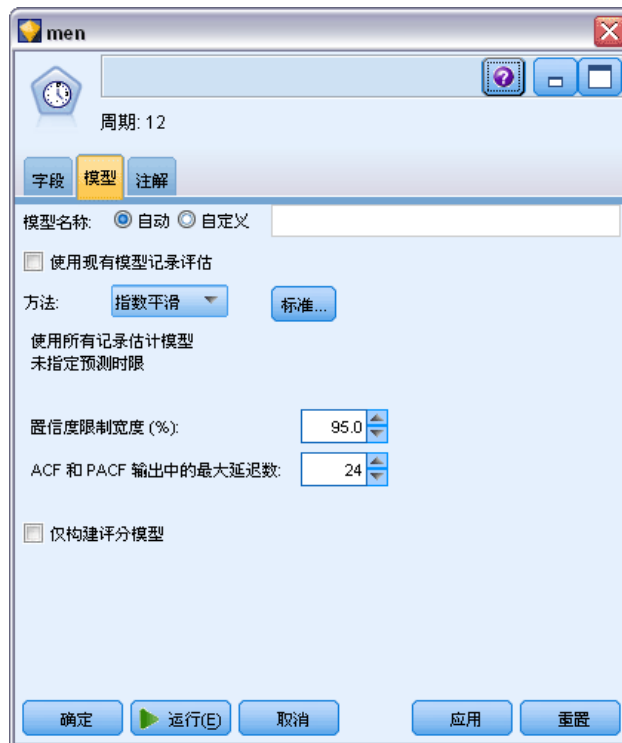
您已了解此序列的特征，便可以开始尝试构建其模型。指数平滑法有助于预测存在趋势和/或季节的序列。如您所见，此处数据同时体现上述两种特征。

指数平稳

创建最适当的指数平滑模型包括确定模型类型—此模型是否需要包含趋势和/或季节—，然后获取最适合选定模型的参数。

随着时间的推移，男装销售散点图建议您使用同时包含线性趋势和乘法季节的模型。这里暗指 Winters 模型。首先我们将开发一个简单模型（即无趋势也无季节），然后开发一个 Holt 模型（存在线性趋势但无季节）。此操作将让您了解模型在什么时候不适合数据，这是成功构建模型的基本技巧。

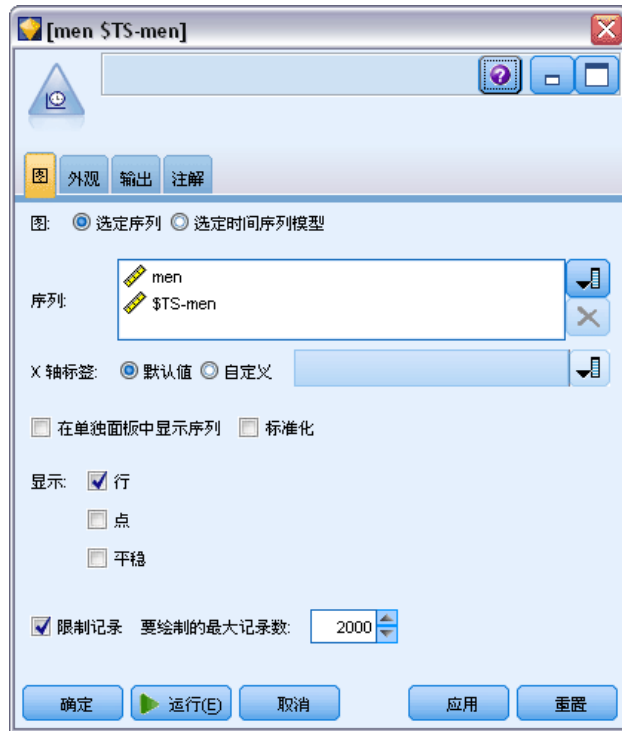
图片 15-6
指定指数平滑



下面我们将开始构建一个简单的指数平滑模型。

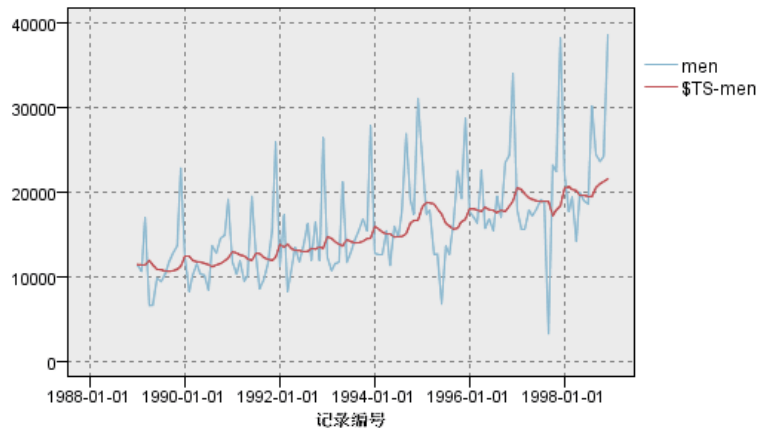
- ▶ 将时间序列节点添加到时间区间节点。
- ▶ 在模型选项卡中，将方法设置为指数平滑。
- ▶ 单击运行创建模型块。

图片 15-7
绘制时间序列模型



- ▶ 将时间散点图节点附加到模型块。
- ▶ 在散点图选项卡中，将男和 \$TS-men 添加到序列列表。
- ▶ 取消选择在单独面板中显示序列和标准化复选框。
- ▶ 单击运行。

图片 15-8
简单指数平滑模型

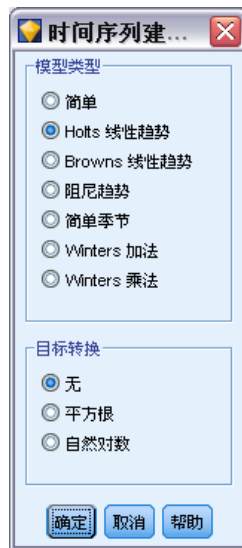


在散点图中，男 表示实际数据，\$TS-men 则表示时间序列模型。

虽然简单模型确实显示了渐进（十分冗长）上升趋势，但它并未考虑季节。您完全可以拒绝此模型。

- ▶ 单击确定关闭时间散点图窗口。

图片 15-9
选择 Holt 模型

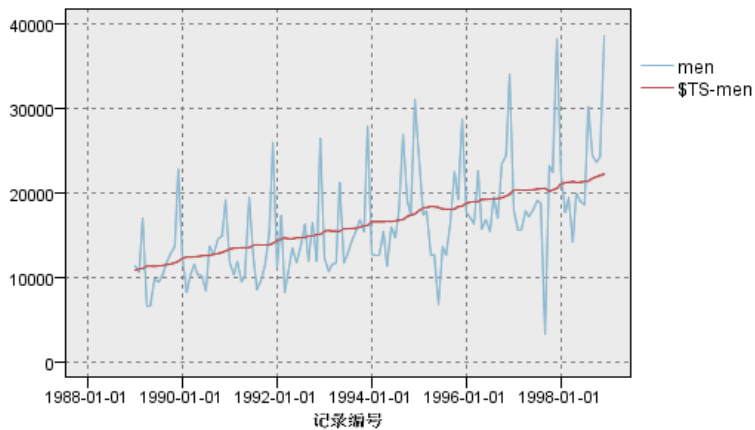


下面试着做一个 Holt 线性模型。虽然，此模型的趋势性会比简单模型稍强，但它同样无法捕捉季节。

- ▶ 重新打开时间序列节点。
- ▶ 在模型选项卡中（依然选择指数平滑方法），单击标准。
- ▶ 在“指数平滑标准”对话框中，选择 Holt 线性趋势。

- ▶ 单击**确定**关闭此对话框。
- ▶ 单击**运行**以再次创建模型块。
- ▶ 再次打开时间散点图节点并单击**运行**。

图片 15-10
Holt 线性趋势模型

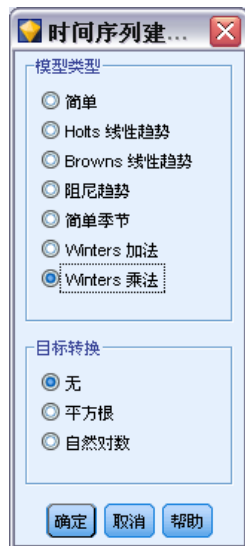


虽然，Holt 模型显示比简单模型更强的平滑趋势，但它仍未考虑季节，所以还应放弃此模型。

- ▶ 关闭时间散点图窗口。

随着时间的推移，男装销售散点图建议您使用同时包含线性趋势和乘法季节的模型，您可以恢复最初的男装销售散点图。因此 Winters 模型才是更适合的备选方案。

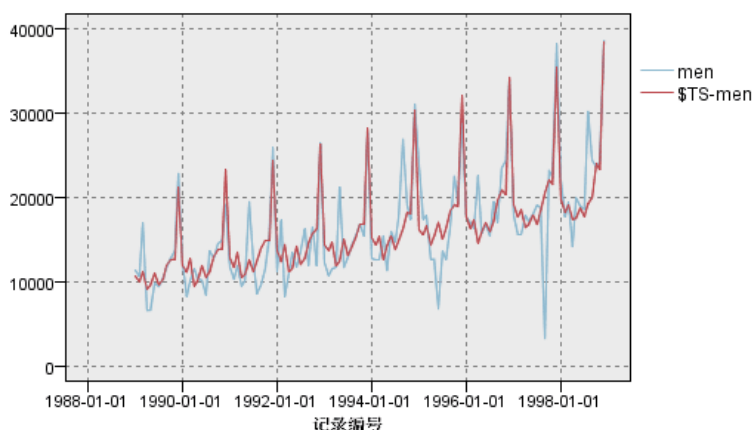
图片 15-11
选择 Winters 模型



- ▶ 重新打开时间序列节点。

- ▶ 在模型选项卡中（依然选择指数平滑方法），单击标准。
- ▶ 在“指数平滑标准”对话框中，选择 Winters 乘法。
- ▶ 单击确定关闭此对话框。
- ▶ 单击运行以再次创建模型块。
- ▶ 打开时间散点图节点并单击运行。

图片 15-12
Winters 乘法模型



此模型较好，因为它同时反映了数据的趋势和季节。

此数据集的时间跨度为 10 年，并且包含 10 个季节峰值（出现在每年十二月份）。这 10 个峰值表示与实际数据中的 10 个年度峰值完全匹配的预测结果。

但此结果同时指出了“指数平滑”步骤的局限性。看看上升和下降的峰值，还有一些重要结构没有得到解释。

如果您旨在构建一个长期包含季节变化趋势的模型，那么可以选择指数平滑模型。要构建此类结构较复杂的模型，则需要考虑使用 ARIMA 步骤。

ARIMA

ARIMA 步骤允许您创建一个适用于时间序列微调建模的自回归整合移动平均 (ARIMA) 模型。ARIMA 模型构建趋势和季节模型的方法比构建指数平滑模型更复杂，并新增了包含预测变量的功能。

继续以想要开发预测模型的产品分类销售公司为例，我们了解了公司如何通过若干可用于解释某些销售变化情况的序列来收集男装月销售数据。预测变量可包括：邮递的产品目录数和产品目录的页数、开通的订购热线数目、印刷广告投入额以及客户服务代表人数。

这些预测变量是否对预测都有用？包含预测变量的模型是否真的比不包含预测变量的模型要好？通过 ARIMA 步骤，我们可以创建一个包含预测变量的预测模型，然后看一下此模型是否与不包含预测变量的指数平滑模型在预测能力上存在巨大差别。

ARIMA 方法使您能够通过指定自动回归、差分和移动平均的顺序以及相应的季节产物来对模型进行微调。由于手动确定上述各部分的最佳值时需要有大量的试错，从而可能变得十分耗时，因此在本例中，我们将为 Expert Modeler 选择 ARIMA 模型。

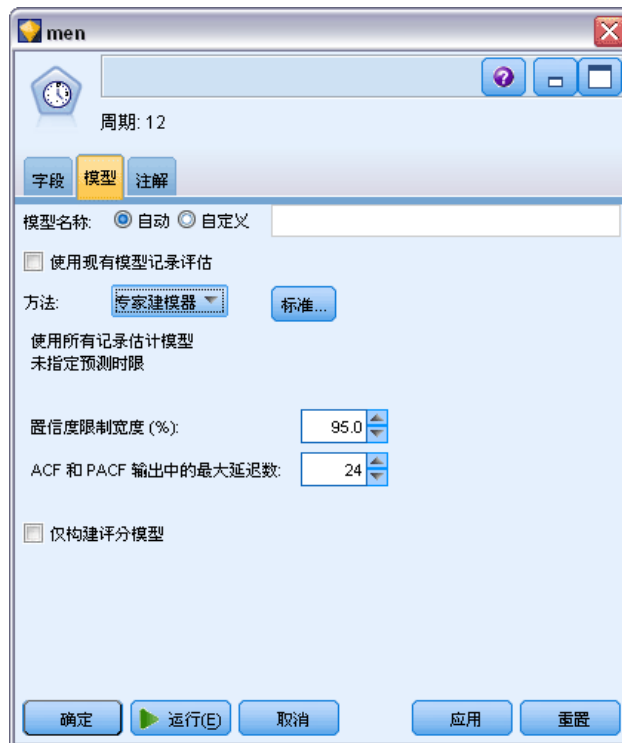
我们会将数据集中的某些其他变量视为预测变量，从而尝试构建更好的模型。最适合作为预测值包含在内的变量有：邮递的产品目录数（邮件）、产品目录的页数（页）、开通的订购热线数目（电话）、印刷广告投入额（印刷）和客户服务代表人数（服务）。

图片 15-13
设置预测值字段



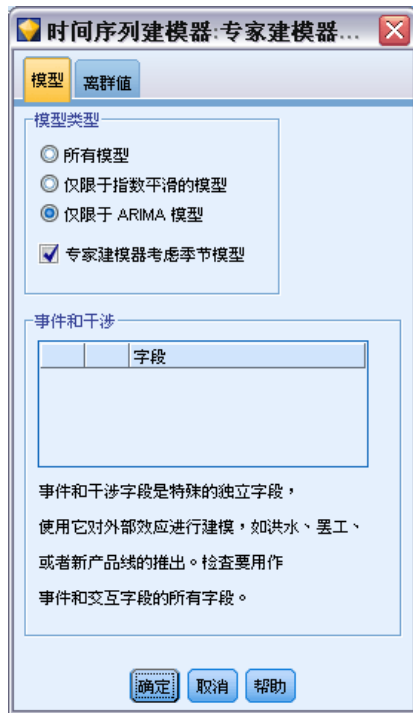
- ▶ 打开 IBM® SPSS® Statistics 文件源节点。
- ▶ 在“类型”选项卡中，将邮件、页、电话、印刷和服务的角色设置为输入。
- ▶ 确保将男的角色设置为目标，并将所有剩余字段的角色设置为无。
- ▶ 单击确定。

图片 15-14
选择 Expert Modeler



- ▶ 打开时间序列节点。
- ▶ 在“模型”选项卡中，将方法设置为 Expert Modeler，然后单击标准。

图片 15-15
仅选择 ARIMA 模型



- ▶ 在“Expert Modeler 标准”对话框中，选择仅 ARIMA 模型选项，并确保选中 Expert Modeler 考虑季节模型。
- ▶ 单击确定关闭此对话框。
- ▶ 单击“模型”选项卡上的运行以再次创建模型块。

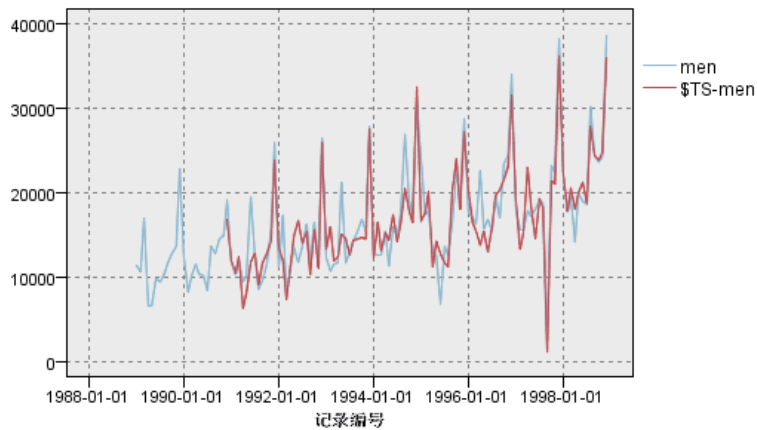
图片 15-16
Expert Modeler 可选择两个预测值



- ▶ 打开该模型块。
注意 Expert Modeler 如何仅选择了 5 个指定预测值中的 2 个作为模型的重大预测值。
- ▶ 单击确定关闭模型块。

- ▶ 打开时间散点图节点并单击运行。

图片 15-17
包含指定预测值的 ARIMA 模型



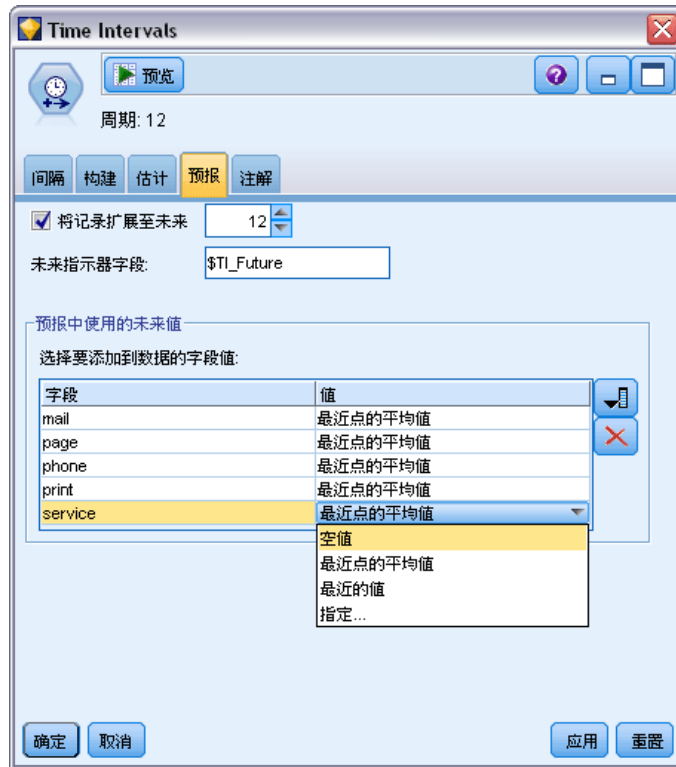
比前一模型有所改进，此模型可以捕捉大型的下降峰值，并将其保持为当前最适合的值。

我们可以进一步优调此模型，但从此以后所进行的任何调整都将微乎其微。我们已创建了最适合的包含预测值的 ARIMA 模型，只需使用刚刚构建的这一模型即可。此示例的目的是预测明年的销售情况。

- ▶ 单击确定关闭时间散点图窗口。
- ▶ 打开时间区间节点，然后选择预测选项卡。
- ▶ 选择将记录扩展到未来复选框，然后将值设置为 12。

预测时使用预测值时，您需要为预测时使用的字段指定估计值，这样建模器可以较为准确地预测目标字段。

图片 15-18
为预测值字段指定未来值



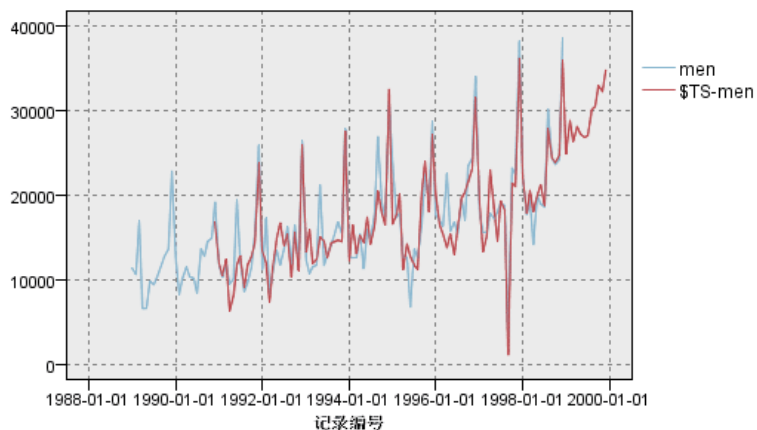
- ▶ 在要在预测中使用的未来值组中，单击“值”列右侧的“字段选择器”按钮。
- ▶ 在“选择字段”对话框中，通过服务选择邮件，然后单击确定。

在实际操作环境中，您可以在此手动指定未来值，因为这五个预测值与您所控制的项都有关。此示例旨在使用某个预定义函数来保存“必须为每个预测值指定 12 个值”选项。（如果对此示例比较熟悉，您可能会试着测试不同的未来值来了解其对模型有何影响。）

- ▶ 对于每个字段，依次单击值字段以显示可能值列表，然后选择最近点的平均值。此选项将计算该字段后三个数据点的平均值，并将其用作每个案例的估计值。
- ▶ 单击确定。
- ▶ 打开时间序列节点并单击运行以再次创建模型块。
- ▶ 打开时间散点图节点并单击运行。

1999 年的预测形势良好：如预期的那样，销售水平继十二月高峰期后再次回复正常，下半年一直保持平稳的上升趋势，整体销售情况比去年有明显的好转。

图片 15-19
以指定预测值进行销售预测



摘要

您已成功构建了一个复杂的时间序列，不仅包含上升趋势，还包含季节变化以及其他变化。通过试错，您还知道如何一步步得到精确模型，然后借助此模型预测未来的销售情况。

事实上，在更新实际销售数据时（例如每月或每季度），您需要重新应用此模型并生成最新的预测。有关详细信息，请参阅第 187 页码第 14 章中的重新应用时间序列模型。

向客户报价（自学）

自学响应模型（SLRM）节点用于生成和启用模型更新，该模型可用于预测哪些报价对客户最合适及报价被接受的概率。这种模型在客户关系管理中非常有用，比如用于市场营销或客户中心。

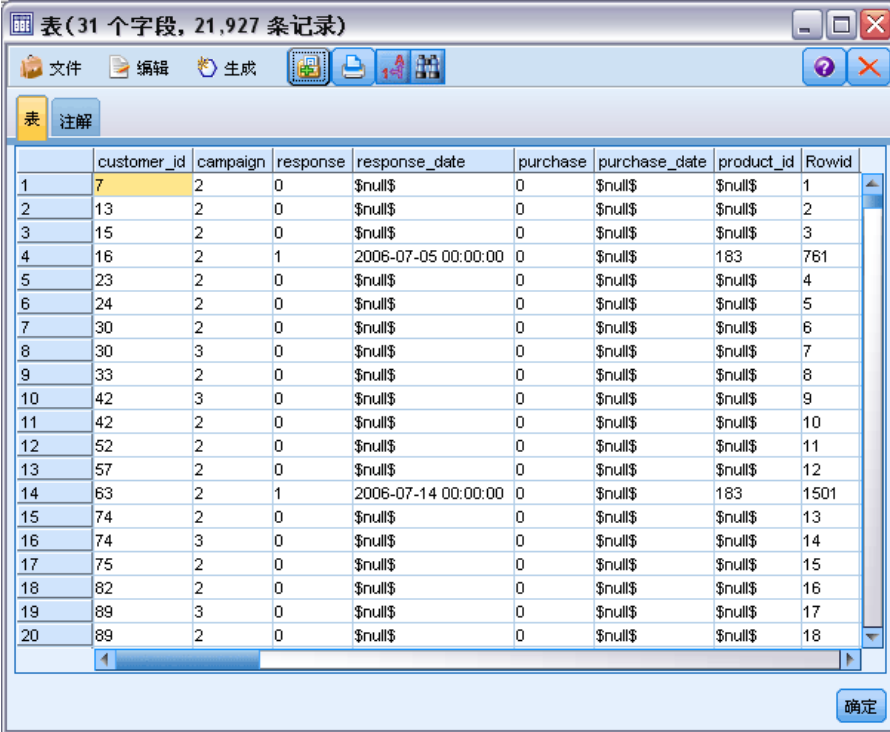
此示例中提到的金融公司纯属虚构。市场营销部希望通过为各个客户匹配合适的报价，在未来的营销活动中创造更好的业绩。示例具体使用自学响应模型，根据以前的报价和响应确定最可能做出正面响应的客户的特征，并根据结果提高当前报价被接受的概率。

本示例使用流 `pm_selflearn.str`，该流引用数据文件 `pm_customer_train1.sav`、`pm_customer_train2.sav` 和 `pm_customer_train3.sav`。这些文件可在任何 IBM® SPSS® Modeler 安装程序的 Demos 文件夹中找到。此目录可通过 Windows IBM® SPSS® Modeler 程序组进行访问。文件 `pm_selflearn.str` 位于 `streams` 文件夹中。

现有数据

公司拥有追踪以前营销活动中向客户做出的报价及客户对报价的响应的历史数据。这些数据还含有可用于预测不同客户响应率的人口统计和金融信息。

图片 16-1
以前报价的响应

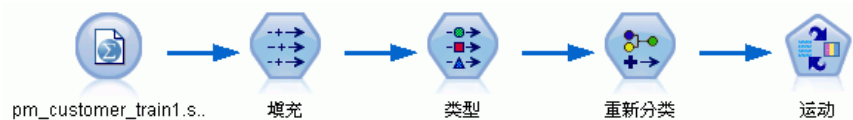


	customer_id	campaign	response	response_date	purchase	purchase_date	product_id	Rowid
1	7	2	0	\$null\$	0	\$null\$	\$null\$	1
2	13	2	0	\$null\$	0	\$null\$	\$null\$	2
3	15	2	0	\$null\$	0	\$null\$	\$null\$	3
4	16	2	1	2006-07-05 00:00:00	0	\$null\$	183	761
5	23	2	0	\$null\$	0	\$null\$	\$null\$	4
6	24	2	0	\$null\$	0	\$null\$	\$null\$	5
7	30	2	0	\$null\$	0	\$null\$	\$null\$	6
8	30	3	0	\$null\$	0	\$null\$	\$null\$	7
9	33	2	0	\$null\$	0	\$null\$	\$null\$	8
10	42	3	0	\$null\$	0	\$null\$	\$null\$	9
11	42	2	0	\$null\$	0	\$null\$	\$null\$	10
12	52	2	0	\$null\$	0	\$null\$	\$null\$	11
13	57	2	0	\$null\$	0	\$null\$	\$null\$	12
14	63	2	1	2006-07-14 00:00:00	0	\$null\$	183	1501
15	74	2	0	\$null\$	0	\$null\$	\$null\$	13
16	74	3	0	\$null\$	0	\$null\$	\$null\$	14
17	75	2	0	\$null\$	0	\$null\$	\$null\$	15
18	82	2	0	\$null\$	0	\$null\$	\$null\$	16
19	89	3	0	\$null\$	0	\$null\$	\$null\$	17
20	89	2	0	\$null\$	0	\$null\$	\$null\$	18

构建流

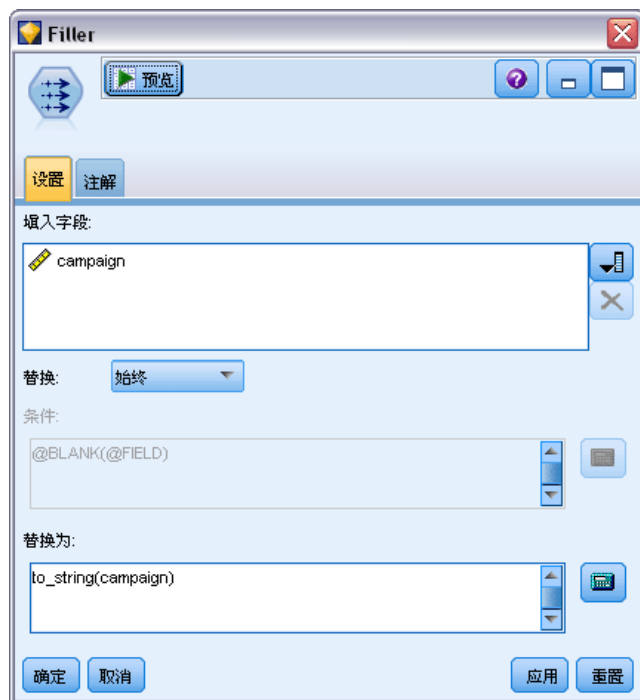
- ▶ 添加指向 pm_customer_train1.sav 的 Statistics 文件源节点，该文件位于 IBM® SPSS® Modeler 安装程序的 Demos 文件夹中。

图片 16-2
SLRM 样本流



- ▶ 添加一个填充节点，并选择 `campaign` 作为字段的填充内容。
- ▶ 选择始终作为替换类型。
- ▶ 在“替换为”文本框中输入 `to_string(campaign)` 并单击确定。

图片 16-3
导出 `campaign` 字段



- ▶ 添加一个类型节点，然后将 customer_id、response_date、purchase_date、product_id、Rowid 和 X_random 字段的角色设置为无。

图片 16-4
更改类型节点设置



- ▶ 将 campaign 和 response 字段的角色设置为目标。这些字段将作为预测的基准。将响应字段的测量设置为标志。

- ▶ 单击读取值，然后单击确定。

由于 campaign 字段数据显示为一列数字（1、2、3 和 4），因此可以对字段进行重新分类以显示更有意义的标题。

- ▶ 为类型节点添加一个重新分类节点。
- ▶ 在重新分类为字段中，选择现有字段。
- ▶ 在重新分类字段列表中，选择 campaign。
- ▶ 单击获取按钮；campaign 的值将添加到原始值列。
- ▶ 在新值列中，在前四行中输入以下活动名称：
 - 抵押
 - 汽车贷款
 - 储蓄
 - 退休金

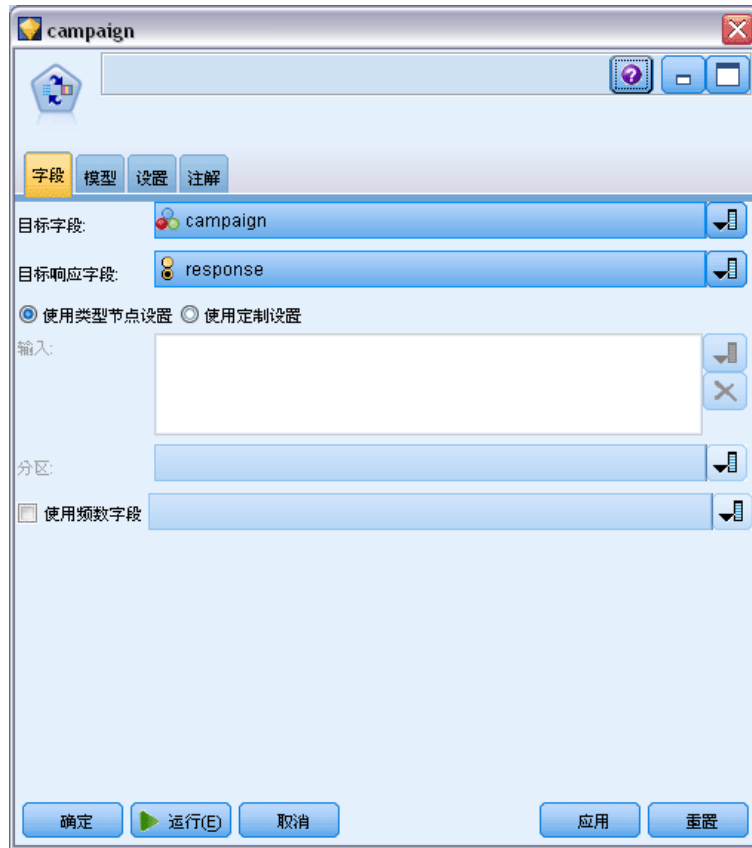
- ▶ 单击确定。

图片 16-5
对活动名称进行重新分类



- ▶ 将 SLRM 建模节点附加到重新分类节点。在“字段”选项卡上，将 `campaign` 和 `response` 分别选择为目标字段和目标响应字段。

图片 16-6
选择目标和目标响应



- ▶ 在“设置”选项卡的“每条记录的最大预测数”字段中，将数字减为 2。这表示对于每位客户，将确定两项具有最高接受概率的报价。

- ▶ 确保选中了考虑模型可靠性，并单击运行。

图片 16-7
SLRM 节点设置



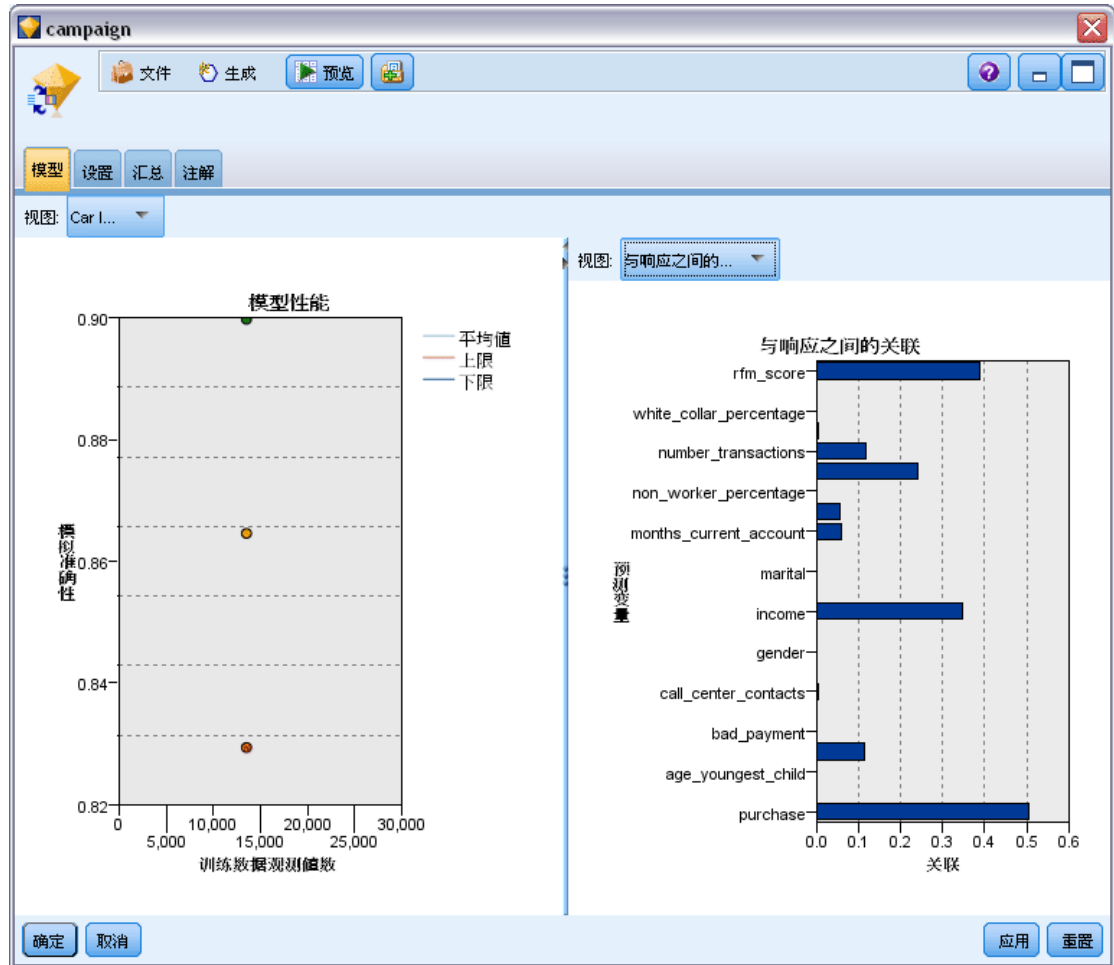
浏览模型

- ▶ 打开该模型块。“模型”选项卡最初显示每项报价的预测准确性估计值，以及每个预测变量在估计模型时的相对重要性。

要显示每个预测变量与目标变量的相关性，从右侧窗格的视图列表中选择与响应关联。

- ▶ 要在具有预测值的四个报价之间进行切换，从左侧窗格的视图列表中选择所需报价。

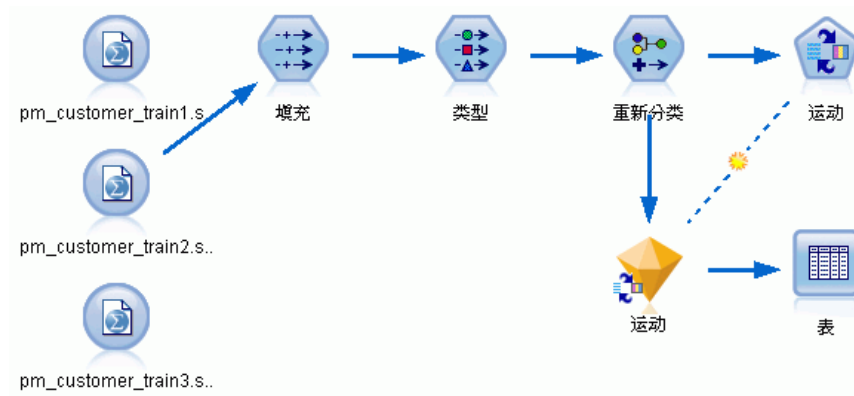
图片 16-8
SLRM 模型块



- ▶ 关闭模型块窗口。
- ▶ 在流工作区上，将指向 pm_customer_train1.sav 的 IBM® SPSS® Statistics 文件源节点断开连接。

- ▶ 添加指向 pm_customer_train2.sav 的 Statistics 文件源节点（该文件位于 IBM® SPSS® Modeler 安装目录的 Demos 文件夹中），并将其连接到过滤节点。

图片 16-9
将第二个数据源附加到 SLRM 流



- ▶ 在 SLRM 节点的“模型”选项卡中，选择继续训练现有模型。

图片 16-10
继续训练模型

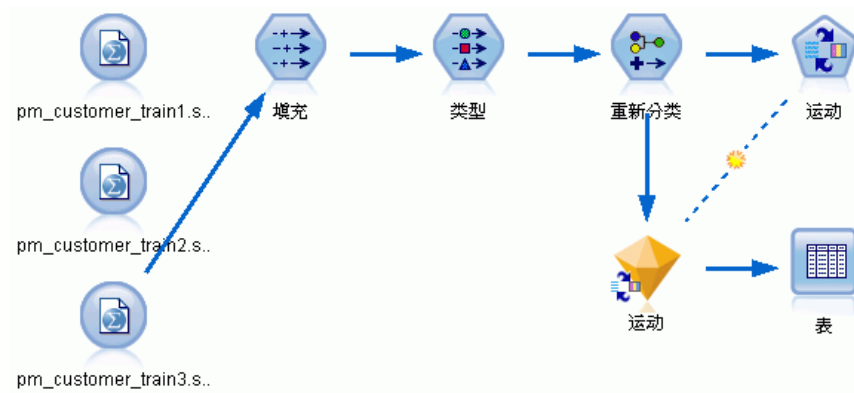


- ▶ 单击运行以再次创建模型块。要查看其详细信息，双击工作区上的模型块。

此时“模型”选项卡将显示每项报价的预测准确性修正估计值。

- ▶ 添加指向 pm_customer_train3.sav 的 Statistics 文件源节点（该文件位于 SPSS Modeler 安装目录的 Demos 文件夹中），并将其连接到过滤节点。

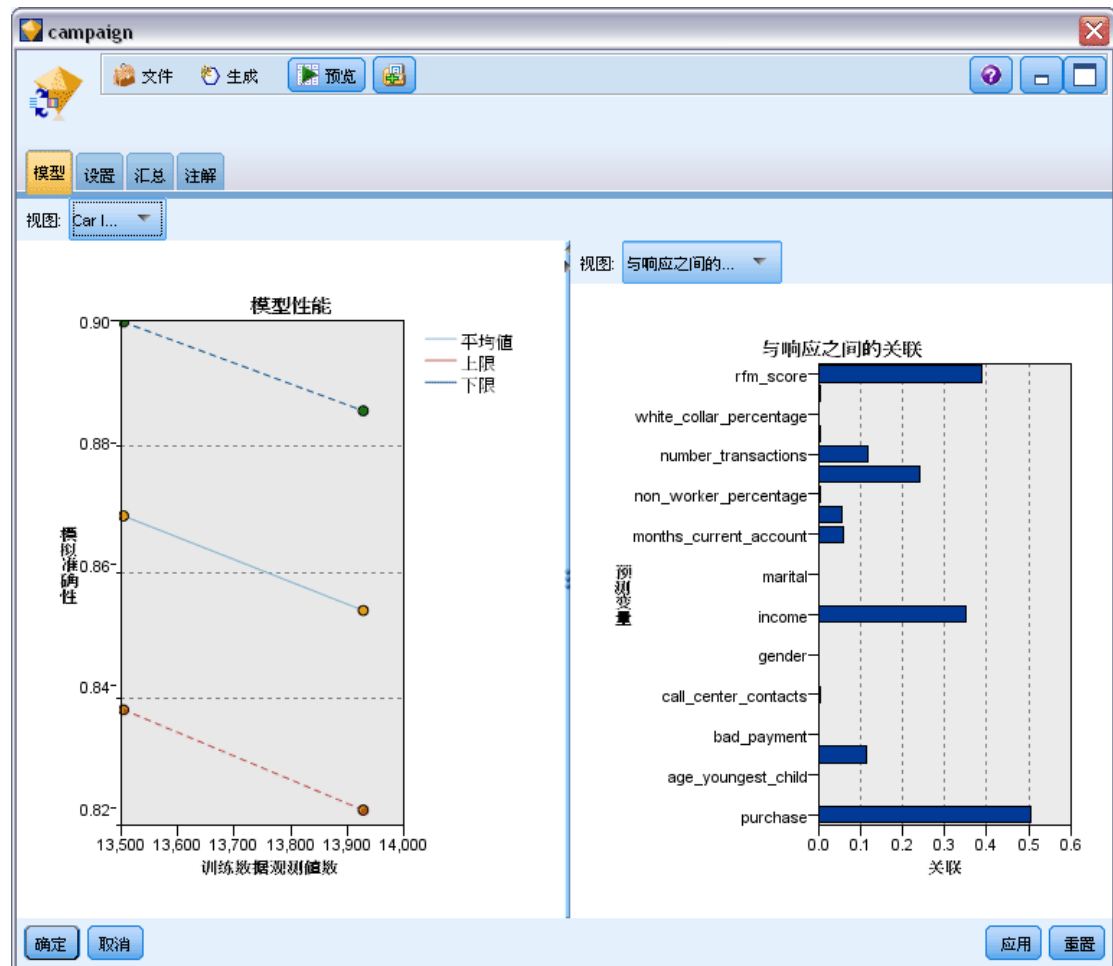
图片 16-11
将第三个数据源附加到 SLRM 流



- ▶ 单击运行以再次创建模型块。要查看其详细信息，双击工作区上的模型块。
- ▶ 此时“模型”选项卡将显示每项报价的预测准确性最终估计值。

可以看到，当您添加其他数据源时，平均准确性稍有下降（从 86.9% 降至 85.4%）；但这种波动程度很小，可能归因于可用数据中的细微异常。

图片 16-12
更新 SLRM 模型块



- ▶ 将表节点附加到生成的最后一个（第三个）模型，然后执行该表节点。
- ▶ 滚动至表的右边。预测将显示客户最有可能接受哪些报价，以及他们将会接受的置信度，具体取决于每位客户的详细信息。

例如，在显示表格的第一行中，以前取得汽车贷款的某位客户将根据提供的报价接受退休金的置信度比率只有 13.2%（表示为 \$SC-campaign-1 列中的 0.132）。但是，第二行和第三行显示了另外两位也曾取得汽车贷款的客户；在他们的案例中，他们以及具有

类似历史记录的其他客户将根据提供的储蓄报价开立储蓄帐户的置信度为 95.7%，并且接受退休金的置信度高于 80%。

图片 16-13
模型输出 - 预测报价和置信度

	X_random	\$S-campaign-1	\$SC-campaign-1	\$S-campaign-2	\$SC-campaign-2
1	1	Pension	0.076	Mortgage	0.061
2	1	Savings	0.977	Mortgage	0.875
3	1	Savings	0.977	Pension	0.686
4	3	Pension	0.076	Mortgage	0.061
5	1	Pension	0.694	Savings	0.101
6	3	Pension	0.076	Mortgage	0.061
7	2	Pension	0.076	Mortgage	0.061
8	3	Pension	0.076	Mortgage	0.061
9	1	Pension	0.076	Mortgage	0.061
10	1	Pension	0.076	Mortgage	0.061
11	2	Pension	0.076	Mortgage	0.061
12	2	Pension	0.076	Mortgage	0.061
13	2	Savings	0.977	Mortgage	0.901
14	2	Pension	0.076	Mortgage	0.061
15	2	Savings	0.977	Pension	0.924
16	2	Pension	0.076	Mortgage	0.061
17	3	Pension	0.076	Mortgage	0.061
18	3	Pension	0.076	Mortgage	0.061
19	3	Savings	0.103	Pension	0.076
20	2	Pension	0.076	Mortgage	0.061

有关 SPSS Modeler 中所用建模方法的数学原理的说明，请参阅《SPSS Modeler 算法指南》，该指南位于产品 DVD 的 \Documentation 目录中。

还请注意，这些结果仅基于训练数据产生。要评估模型对实际应用中的其他数据的拟合程度，可使用分区节点保留部分记录，以便于测试和验证。有关详细信息，请参阅第 4 章中的分区节点中的 IBM SPSS Modeler 15 源、过程和输出节点。有关 SLRM 节点的详细信息，请参阅节点参考中的第 14 章。

预测贷款拖欠者（贝叶斯网络）

使用贝叶斯网络，可以通过将观察到并记录下的证据与实际常识结合起来构建概率模型，以通过使用表面看上去不相关的属性确定发生的可能性。

此示例使用名为 bayes_bankloan.str 的流，它引用名为 bankloan.sav 的数据文件。这些文件位于任意 IBM® SPSS® Modeler 安装程序中的 Demos 目录下，并可从 Windows “开始” 菜单上的 IBM® SPSS® Modeler 程序组进行访问。文件 bayes_bankloan.str 位于 streams 目录下。

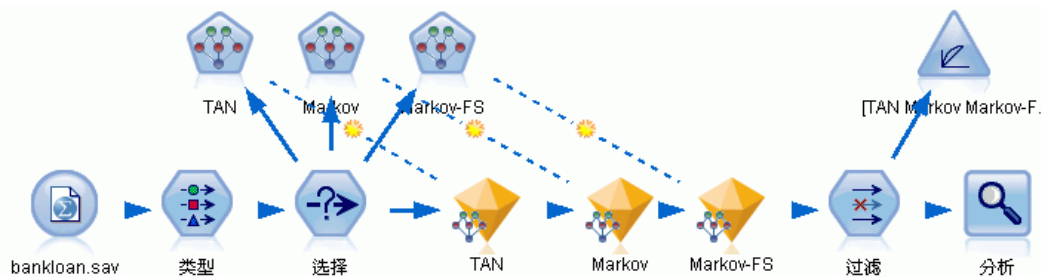
例如，假设某个银行希望了解不偿还贷款的潜在情况。如果先前的贷款拖欠数据可用于预测哪些潜在客户可能在偿还贷款时有问题，则可以对这些“不良风险”的客户减少贷款或者为他们提供其他产品。

此示例主要使用了现有贷款拖欠数据来预测今后出现的潜在贷款拖欠者，并观察了三个不同的贝叶斯网络模型类型，从而确定在这种情况下哪个类型的预测效果更好。

构建流

- ▶ 在 Demos 文件夹中添加指向 bankloan.sav 的 Statistics 文件源节点。

图片 17-1
贝叶斯网络样本流



- ▶ 将类型节点添加到源节点，并将默认字段的角色设为目标。将所有其他字段的角色设置为 Input。

- ▶ 单击读取值按钮以填充值列。

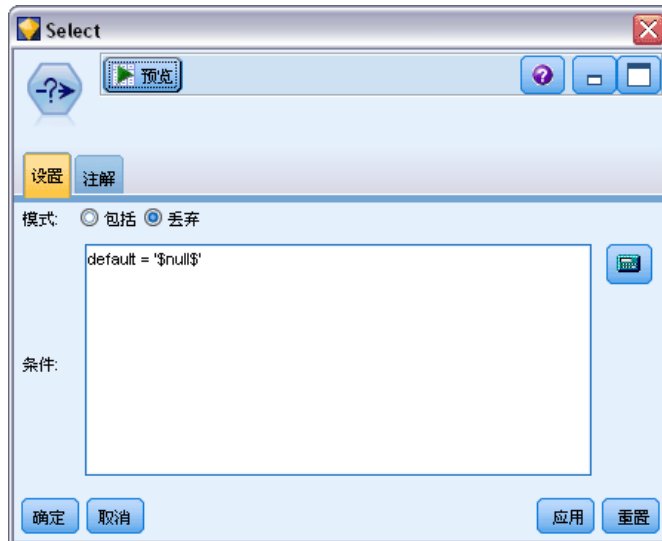
图片 17-2
选择目标字段



构建模型时其目标字段有空值的观测值没有意义。您可以排除这些观测值以防止在模型评估中使用它们。

- ▶ 为类型节点添加一个选择节点。
- ▶ 对于模型，请选择丢弃。
- ▶ 在“条件”框中，输入 `default = '$null$'`。

图片 17-3
丢弃空的目标字段



因为您构建了多个不同类型的贝叶斯网络，所以最好对它们进行比较，从而明确哪些模型可提供最好的预测。第一个要创建的模型是树扩展朴素贝叶斯（TAN）模型。

- ▶ 将贝叶斯网络节点附加到选择节点上。
- ▶ 对于“模型”选项卡上的模型名称，请选择自定义，并在文本框中输入 TAN。
- ▶ 对于结构类型，请选择 TAN 并单击确定。

图片 17-4
创建树扩展朴素贝叶斯模型



第二个要创建的模型具有马尔可夫覆盖结构。

- ▶ 将第二个贝叶斯网络节点添加到选择节点上。
- ▶ 对于“模型”选项卡上的模型名称，请选择自定义，并在文本框中输入 Markov。

- ▶ 对于结构类型，请选择 **Markov Blanket** 并单击确定。

图片 17-5
创建马尔可夫覆盖模型



要创建的第三个模型具有马尔可夫覆盖结构，同时也使用了特征选择预处理来选择与目标变量有重大关联的输入。

- ▶ 将第三个贝叶斯网络节点添加到选择节点上。
- ▶ 对于“模型”选项卡上的模型名称，请选择自定义，并在文本框中输入 **Markov-FS**。
- ▶ 对于结构类型，请选择 **Markov Blanket**。

- ▶ 选择包括特征选择预处理步骤并单击确定。

图片 17-6
使用特征选择预处理来创建马尔可夫覆盖模型



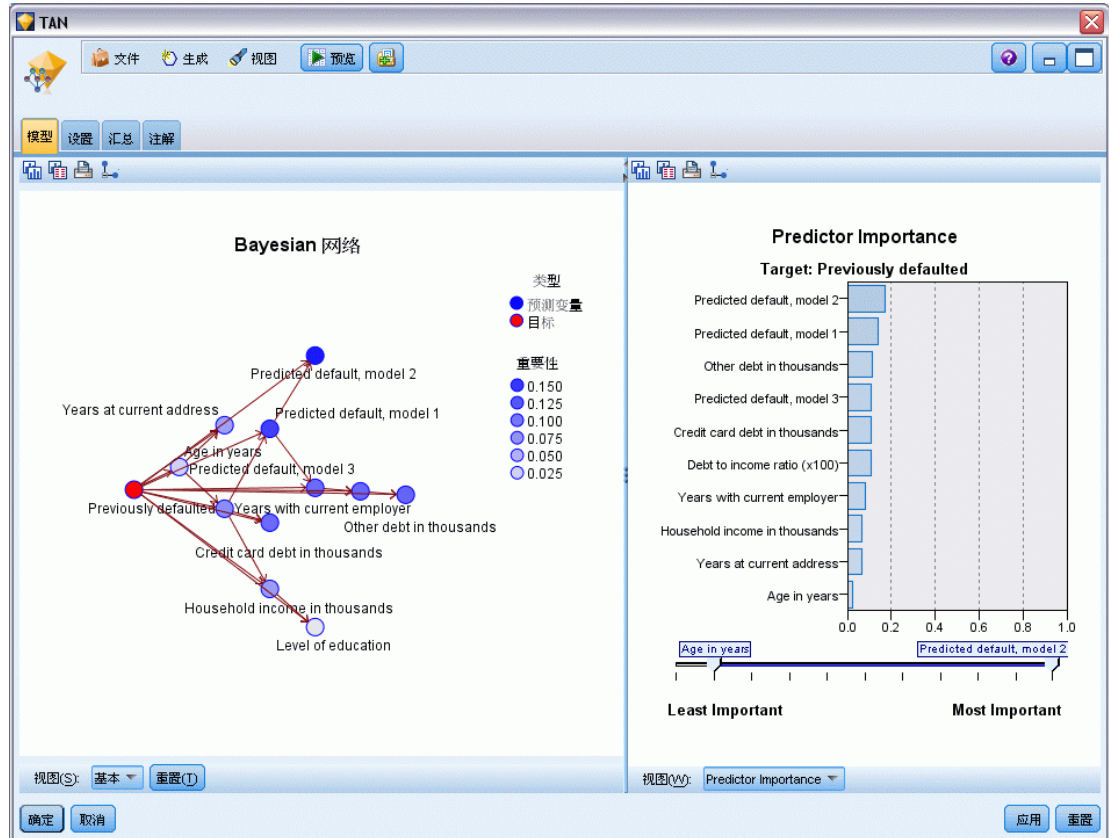
浏览模型

- ▶ 运行流以创建模型块，该模型块将被添加到流和位于右上角的“模型”选项板中。要查看其详细信息，请双击流中的任一模型块。

模型块“模型”选项卡分为两个窗格。左窗格包含节点网络图，可显示目标与其最重要预测变量之间的关系，以及各预测变量之间的关系。

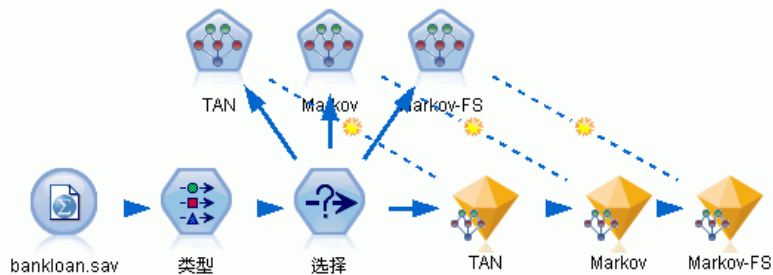
右窗格可能显示预测变量重要性，它表示评估模型时每个预测变量的相对重要性，右窗格也可能显示条件概率，它包含各个节点值的条件概率值，以及各节点的父节点中的所有值组合。

图片 17-7
查看树扩展朴素贝叶斯模型



- ▶ 将 TAN 模型块连接到 Markov 模型块（选择警告对话框上的替换）。
- ▶ 将 Markov 模型块连接到 Markov-FS 模型块（选择警告对话框上的替换）。
- ▶ 通过选择节点对齐三个模型块以方便查看。

图片 17-8
对齐流中的模型块



- ▶ 要重新命名您正在创建的评估图形上的模型输出以避免混淆，请将过滤节点附加到 Markov-FS 模型块。
- ▶ 在右侧的字段栏中，将 \$B-default、\$B1-default 和 \$B2-default 分别重新命名为 TAN、马尔可夫和 Markov-FS。

图片 17-9
重新命名模型字段名

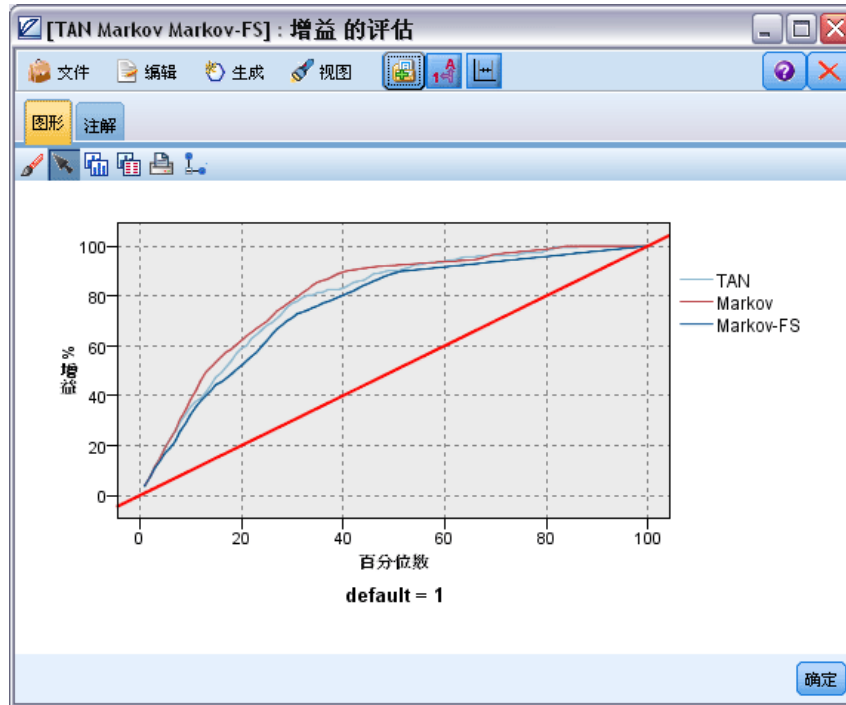


要比较模型的预测准确性，您可以构建一个收益图表。

- ▶ 将评估图形节点附加到过滤节点上，然后使用图形节点的默认设置来执行它。

该图形显示，每个模型类型都生成了相似的结果，但是马尔可夫模型要稍微好一些。

图片 17-10
评估模型准确性

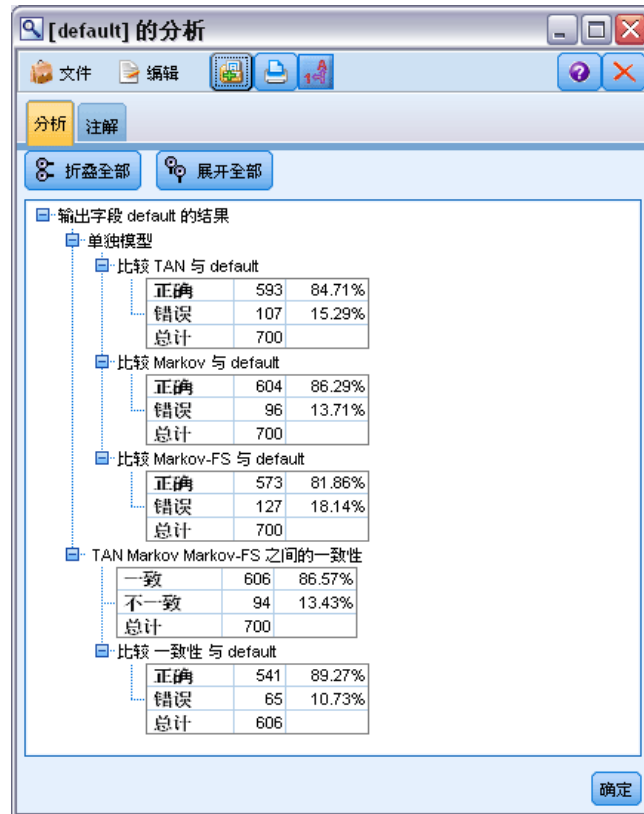


要检查每个模型的预测效果，您可能会使用分析节点而不是评估图形。下图显示了依据正确和不正确的预测百分比得出的准确性。

- ▶ 将分析节点附加到过滤节点上，然后使用分析节点的默认设置来执行它。

与评估图形一样，本图显示 Markov 模型在正确预测方面稍微好一些；但 Markov-FS 模型仅落后 Markov 模型几个百分点。这可能就意味着使用 Markov-FS 模型要更为方便一些，因为它计算结果所需的输入更少，因此节省了数据收集和输入的时间以及处理时间。

图片 17-11
分析模型准确性



有关 IBM® SPSS® Modeler 中所用建模方法的数学原理的说明，请参阅《SPSS Modeler 算法指南》，该指南位于安装光盘的 \Documentation 目录中。

还请注意，这些结果仅基于训练数据产生。要评估模型对实际应用中的其他数据的拟合程度，可使用分区节点保留部分记录，以便于测试和验证。有关详细信息，请参阅第 4 章中的分区节点中的 IBM SPSS Modeler 15 源、过程和输出节点。

每个月重新训练模型（贝叶斯网络）

使用贝叶斯网络，可以通过将观察到并记录下的证据与实际常识结合起来构建概率模型，以通过使用表面看上去不相关的属性确定发生的可能性。

此示例使用名为 bayes_churn_retrain.str 的流，此流引用名为 telco_Jan.sav 和 telco_Feb.sav 的数据文件。这些文件位于任意 IBM® SPSS® Modeler 安装程序中的 Demos 目录下，并可从 Windows “开始” 菜单上的 IBM® SPSS® Modeler 程序组进行访问。文件 bayes_churn_retrain.str 位于 streams 目录下。

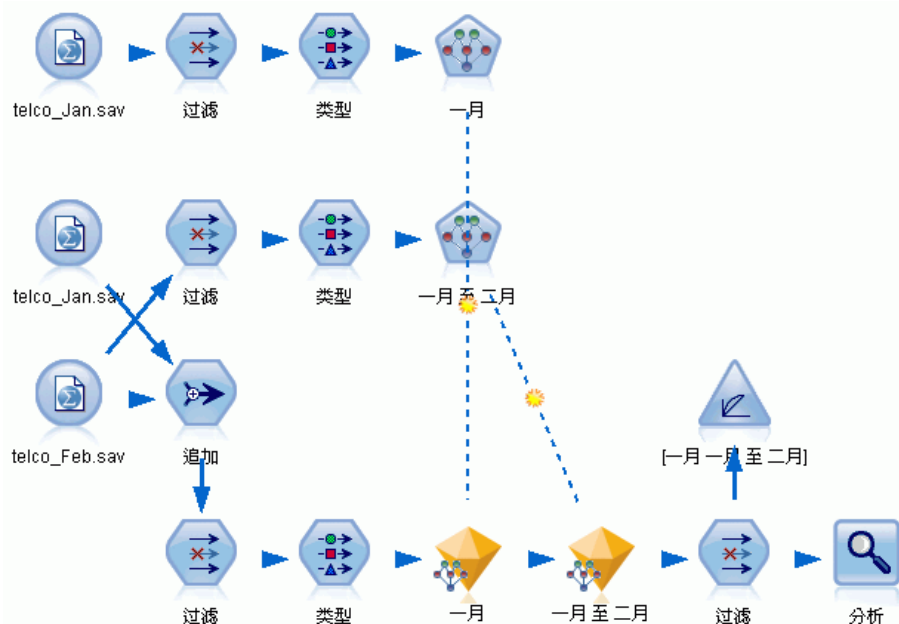
例如，假设某个电信服务提供商非常关心流失到竞争对手那里的客户数(流失)。如果历史的客户数据用作预测哪些客户更可能在今后流失，那么这些客户可能划定为出于刺激性的动机或其他一些使它们失去信心的意图而倒向了另一个服务提供商。

本示例主要使用现有的每月流失数据来预测哪些客户在今后流失，然后将这些数据添加到模型中，从而精炼和重新训练模型。

构建流

- ▶ 在 Demos 文件夹中添加指向 telco_Jan.sav 的 Statistics 文件源节点。

图片 18-1
贝叶斯网络样本流



前面的分析已经显示，预测流失时，有几个数据字段并不太重要。这些字段可从数据集中滤出，从而在构建模型以及对模型评分时提高处理速度。

- ▶ 为源节点添加一个过滤节点。
- ▶ 排除除地址、年龄、流失、客户类别、教育程度、行业、性别、婚姻状况、居住地、退休和工龄外的所有字段。
- ▶ 单击确定。

图片 18-2
过滤不必要的字段



- ▶ 为过滤节点添加一个类型节点。
- ▶ 打开类型节点并单击读取值按钮以填充值列。

- 为了使“评估”节点可以评估值的真假，需要将流失字段的测量级别设置为标志，并将其角色设置为目标。单击确定。

图片 18-3
选择目标字段



您可以构建多个不同类型的贝叶斯网络；但在本示例中，您将构建的是树扩展朴素贝叶斯（TAN）模型。该模型创建了一个大型网络，可以确保其中已经囊括了数据变量间所有可能存在的连接关系，从而构建一个稳健的初始模型。

- 将贝叶斯网络节点附加到类型节点上。
- 对于“模型”选项卡上的模型名称，请选择自定义，并在文本框中输入 Jan。
- 对于参数学习方法，请选择对小单元格计数的贝叶斯调整。

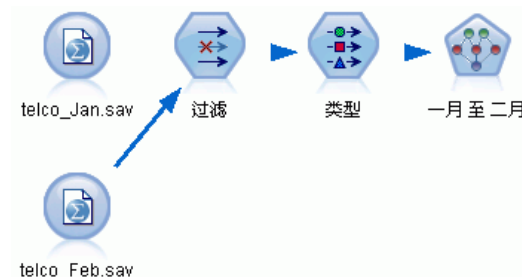
- ▶ 单击运行。模型块被添加到流，同时添加到右上角的?模型?选项板。

图片 18-4
创建树扩展朴素贝叶斯模型



- ▶ 在 Demos 文件夹中添加指向 telco_Feb.sav 的 Statistics 文件源节点。
- ▶ 将此新的源节点附加到过滤节点（在警告对话框上，选择替换以替换到前一源节点的连接）。

图片 18-5
添加第二个月的数据



- ▶ 对于贝叶斯网络节点的“模型”选项卡上的模型名称，请选择自定义，并在文本框中输入 Jan-Feb。
- ▶ 选择继续训练现有模型。

- ▶ 单击运行。模型块覆盖流中的现有模型块，但同时也将添加到右上角的“模型”选项板。

图片 18-6
重新训练模型



评估模型

要对模型进行比较，必须将两个数据集合并到一起。

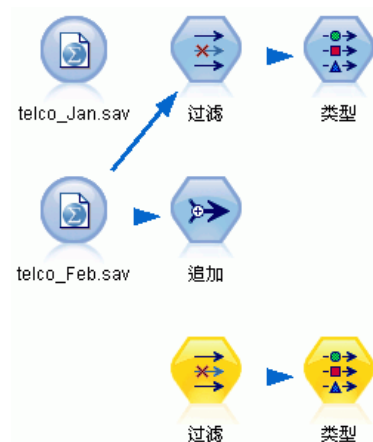
- ▶ 添加一个追加节点并将telco_Jan.sav 和 telco_Feb.sav 源节点都附加到该节点上。

图片 18-7
追加两个数据源



- ▶ 从早期的流中复制过滤节点和类型节点，然后将它们粘贴到流工作区上。
- ▶ 将追加节点附加到最新复制的过滤节点上。

图片 18-8
将复制的节点粘贴到流中

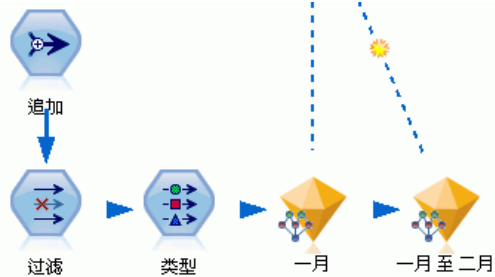


两个贝叶斯网络模型的模型块位于右上角的“模型”选项板中。

- ▶ 双击 Jan 模型块以将其导入流，并将其附加到新复制的类型节点。
- ▶ 将流中现有的 Jan-Feb 模型块附加到 Jan 模型块。

- ▶ 打开 Jan 模型块。

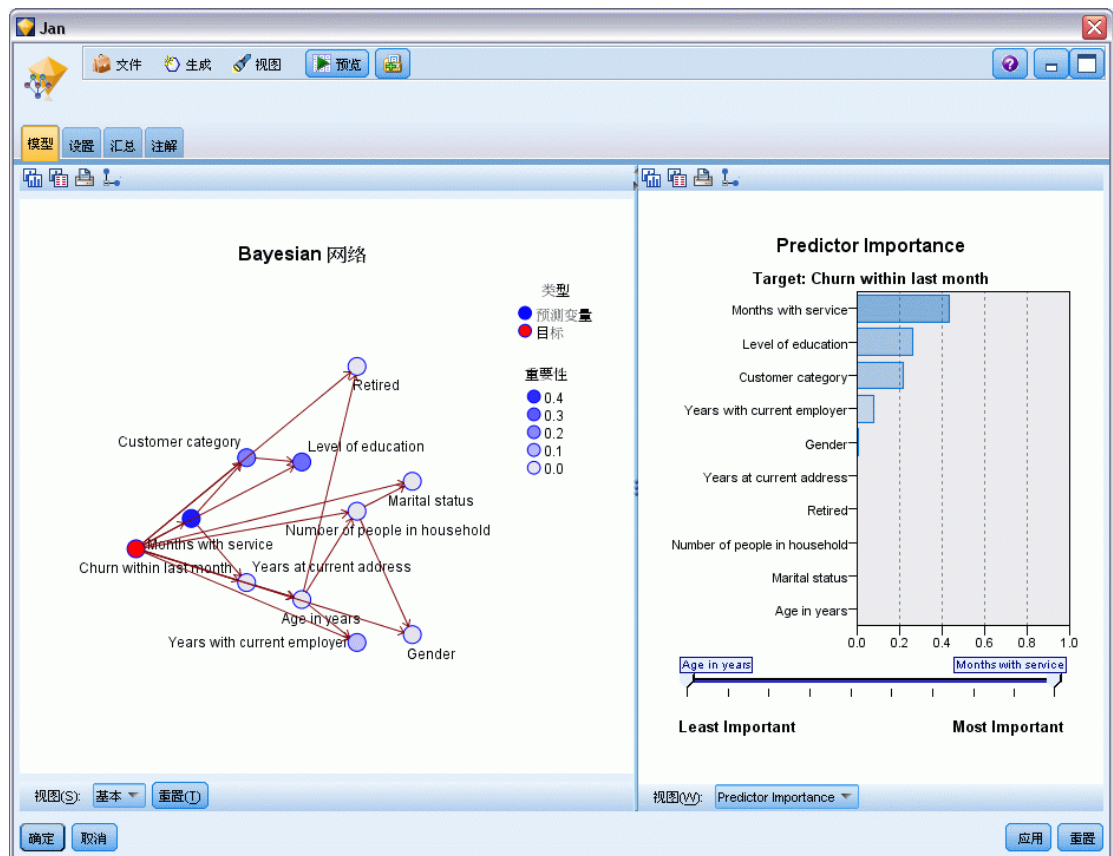
图片 18-9
将模型块添加到流中



贝叶斯网络模型块“模型”选项卡分为两列。左列包含节点网络图，可显示目标与其最重要预测变量之间的关系，以及各预测变量之间的关系。

右列可能显示预测变量重要性，它表示评估模型时每个预测变量的相对重要性，右窗格也可能显示条件概率，它包含各个节点值的条件概率值，以及各节点的父节点中的所有值组合。

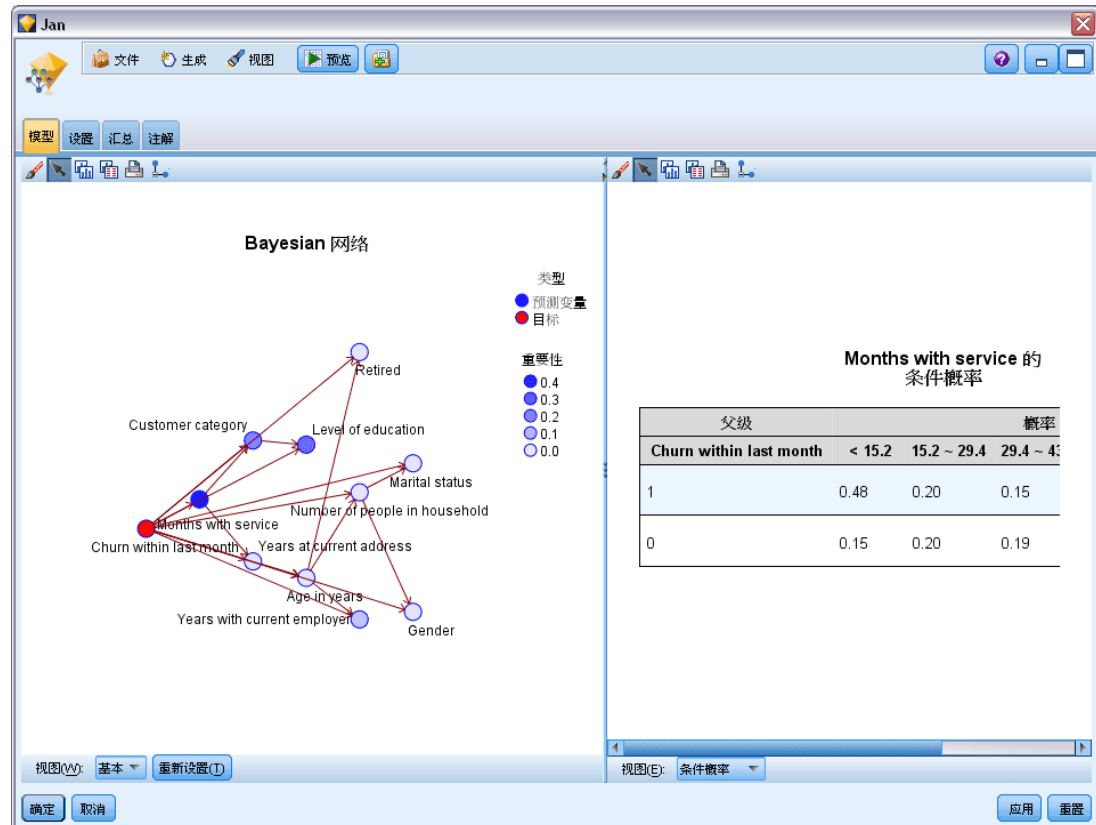
图片 18-10
贝叶斯网络模型显示了预测变量重要性



要显示所有节点的条件概率，请单击左列中的节点。相应的右列会更新以显示所需的详细信息。

显示每个分级的条件概率，这些分级中的数据值已划分为与该节点的父节点和同胞节点相关。

图片 18-11
贝叶斯网络模型显示了条件概率



- 要重新命名模型输出以避免混淆，请将过滤节点附加到 Jan-Feb 模型块。

- ▶ 在右侧的字段列，将 \$B-churn 和 \$B1-churn 分别重新命名为 Jan 和 Jan-Feb。

图片 18-12
重新命名模型字段名

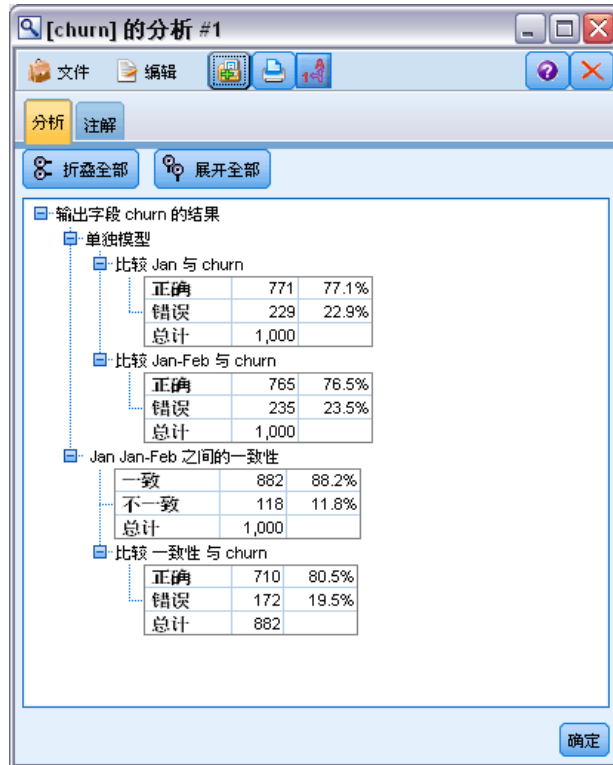


要检查每个模型预测流失的好坏，请使用分析节点；这样将显示依据正确和不正确的预测百分比得出的准确性。

- ▶ 将分析节点附加到过滤节点。
- ▶ 打开分析节点并单击运行。

这表明两个模型在预测流失时具有类似精确度。

图片 18-13
分析模型准确性

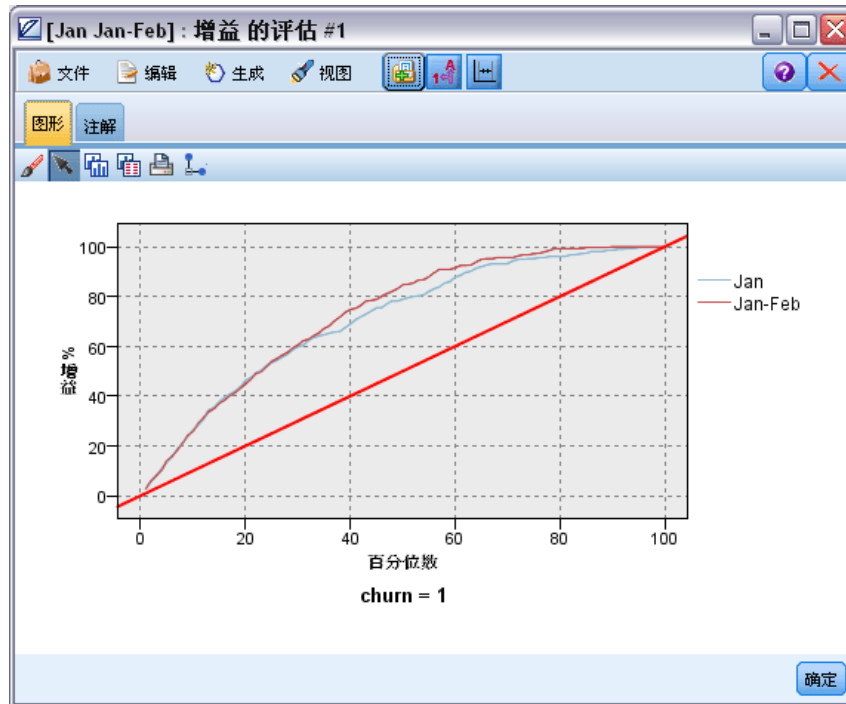


您可以使用评估图形代替分析节点，通过构建一个收益图表比较模型的预测准确性。

- 将评估图形节点附加到过滤节点。
并使用其默认设置执行图形节点。

与分析节点相同，该图形显示两个模型类型都生成了相似的结果；但是，因为使用两个月数据的重新训练模型在预测中具有更高水平的置信度，所以要稍微好一些。

图片 18-14
评估模型准确性



有关 IBM® SPSS® Modeler 中所用建模方法的数学原理的说明，请参阅《SPSS Modeler 算法指南》，该指南位于安装光盘的 \Documentation 目录中。

还请注意，这些结果仅基于训练数据产生。要评估模型对实际应用中的其他数据的拟合程度，可使用分区节点保留部分记录，以便于测试和验证。有关详细信息，请参阅第 4 章中的分区节点中的 IBM SPSS Modeler 15 源、过程和输出节点。

零售促销（神经网络/C&RT）

此示例使用数据来说明零售产品线 and 促销对销售的影响。（此数据纯为虚构。）此示例的目的在于预测未来促销活动的影响。与条件监视示例类似，数据挖掘过程包括探索、数据准备、训练和检验阶段。

此示例使用名为 `goodsplot.str` 和 `goodslearn.str` 的流，这些流引用名为 `GOODS1n` 和 `GOODS2n` 的数据文件。这些文件可在任何 IBM® SPSS® Modeler 安装程序的 Demos 目录中找到。此目录可通过 WindowsIBM® SPSS® Modeler 程序组进行访问。流 `goodsplot.str` 在 `streams` 文件夹中，而 `goodslearn.str` 文件在 `streams` 目录中。

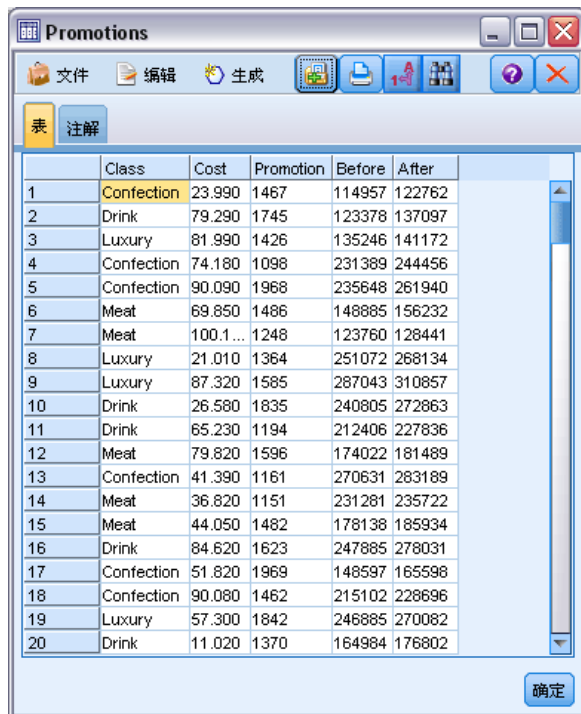
检查数据

每条记录含有：

- `Class`. 模型类型。
- `Cost`. 单价。
- `Promotion`. 特定促销上所花费金额的指数。
- `Before`. 促销之前的收入。
- `After`. 促销之后的收入。

流 `goodsplot.str` 含有一个用于在表格中显示数据的简单流。两个收入字段（即 `Before` 和 `After`）用绝对值来表示；但是，可能促销后收入的增长量（并假定收入增长源于促销）是更有用的数据。

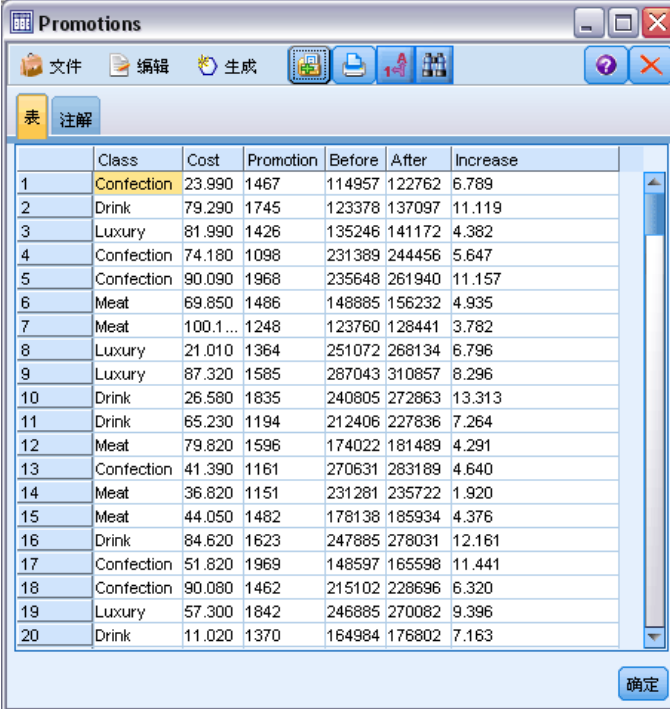
图片 19-1
促销对产品销售的影响



	Class	Cost	Promotion	Before	After
1	Confection	23.990	1467	114957	122762
2	Drink	79.290	1745	123378	137097
3	Luxury	81.990	1426	135246	141172
4	Confection	74.180	1098	231389	244456
5	Confection	90.090	1968	235648	261940
6	Meat	69.850	1486	148885	156232
7	Meat	100.1...	1248	123760	128441
8	Luxury	21.010	1364	251072	268134
9	Luxury	87.320	1585	287043	310857
10	Drink	26.580	1835	240805	272863
11	Drink	65.230	1194	212406	227836
12	Meat	79.820	1596	174022	181489
13	Confection	41.390	1161	270631	283189
14	Meat	36.820	1151	231281	235722
15	Meat	44.050	1482	178138	185934
16	Drink	84.620	1623	247885	278031
17	Confection	51.820	1969	148597	165598
18	Confection	90.080	1462	215102	228696
19	Luxury	57.300	1842	246885	270082
20	Drink	11.020	1370	164984	176802

goodsplot.str 也包含引导出该值的节点，然后在名称为增长量的字段中用促销前的收入百分比来表达该值，并显示一个带有该字段的表格。

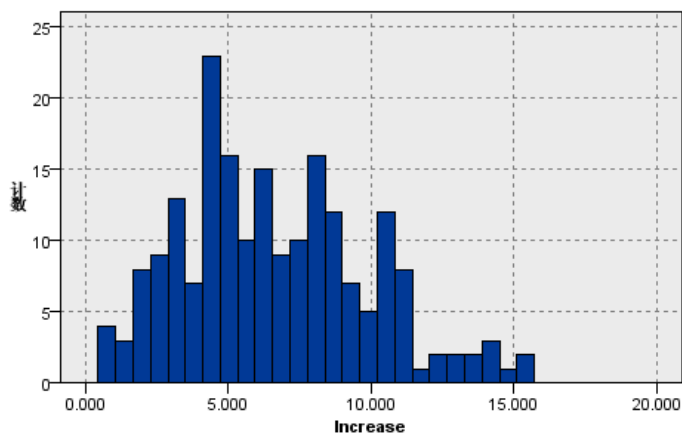
图片 19-2
促销之后的收入增长量



	Class	Cost	Promotion	Before	After	Increase
1	Confection	23.990	1467	114957	122762	6.789
2	Drink	79.290	1745	123378	137097	11.119
3	Luxury	81.990	1426	135246	141172	4.382
4	Confection	74.180	1098	231389	244456	5.647
5	Confection	90.090	1968	235648	261940	11.157
6	Meat	69.850	1486	148885	156232	4.935
7	Meat	100.1...	1248	123760	128441	3.782
8	Luxury	21.010	1364	251072	268134	6.796
9	Luxury	87.320	1585	287043	310857	8.296
10	Drink	26.580	1835	240805	272863	13.313
11	Drink	65.230	1194	212406	227836	7.264
12	Meat	79.820	1596	174022	181489	4.291
13	Confection	41.390	1161	270631	283189	4.640
14	Meat	36.820	1151	231281	235722	1.920
15	Meat	44.050	1482	178138	185934	4.376
16	Drink	84.620	1623	247885	278031	12.161
17	Confection	51.820	1969	148597	165598	11.441
18	Confection	90.080	1462	215102	228696	6.320
19	Luxury	57.300	1842	246885	270082	9.396
20	Drink	11.020	1370	164984	176802	7.163

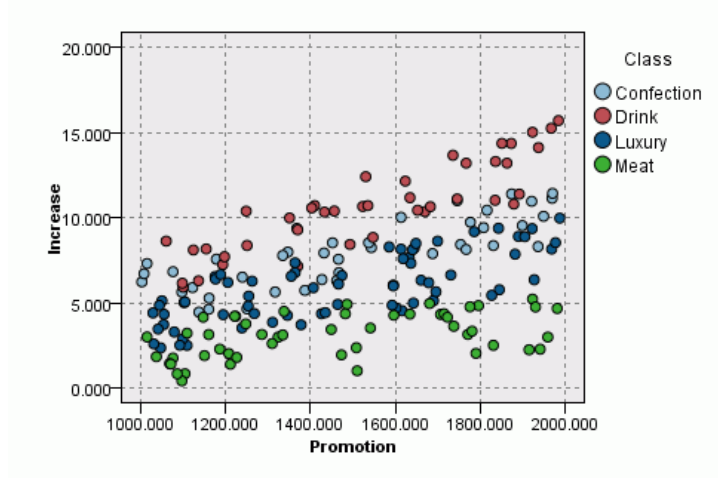
另外，流将显示一个增长量的直方图和一个以促销费用为参照的增长量的散点图，产品的各个类别的散点图将叠放在一起。

图片 19-3
收入增长量直方图



散点图显示对于每类产品，收入增长量和促销费用之间存在几乎线性的关系。因此，决策树或神经网络似乎可以合理和准确地预测其他可用字段上的收入增长量。

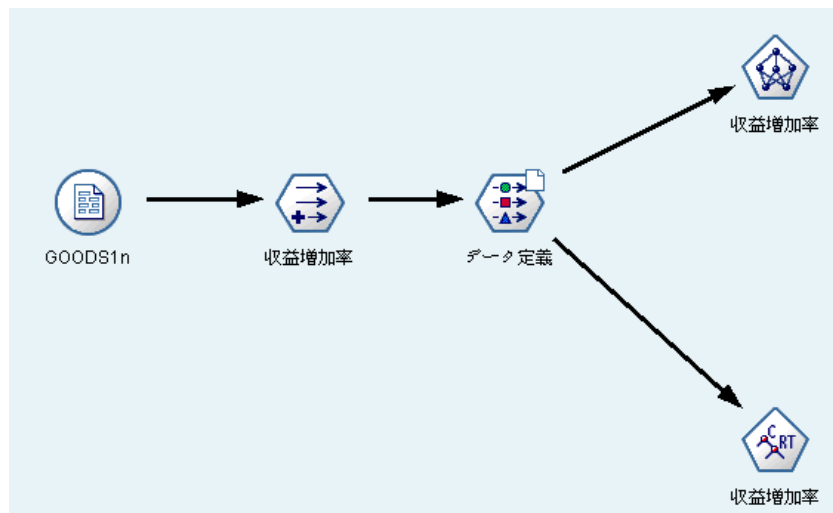
图片 19-4
收入增长量与促销费用



学习和检验

流 goodslearn.str 将训练神经网络和决策树，以对收入增长量做出预测。

图片 19-5
构建流 goodslearn.str



执行模型节点和生成实际模型之后，即可检验学习过程的效果。检验方法如下：将“类型”节点和新“分析”节点之间的决策树和网络串联起来，接着将输入（数据）文件更改为 GOODS2n，然后执行“分析”节点。按照此节点的输出数据，特别是按照

预测的增长量与正确答案之间的线性相关进行判断，可以发现已训练系统对收入增长量的预测成功率颇高。

进一步的探索应该集中在那些与已训练系统的预测有较大差别的案例上；通过收入的预测增长量与真实增长量的对比图，可标识出这些案例。可使用 IBM® SPSS® Modeler 的迭代图来选择图上的离群值，而依据离群值的属性，通过调整数据说明和学习过程，提高预测的准确性将成为可能。

状态监测（神经网络/C5.0）

本示例涉及如何监测计算机的状态信息及识别和预测故障状态相关问题。其中的数据通过虚构模拟创建得到并包括大量按时间测量的连续序列。每个记录都是与计算机的以下方面相关的快照报告：

- 时间。整数。
- 功率。整数。
- 温度。整数。
- 压力。0 表示正常，1 表示瞬时压力报警。
- 正常工作时间。上次运行时间。
- 状态。正常情况下是 0，发生错误时更改为错误代码（101、202 或 303）。
- 结果。在此时间序列中显示错误代码，如果没有发生错误，则显示为 0。（提供这些代码唯一的好处是可在事后了解出现的错误。）

此示例使用名为 condplot.str 和 condlearn.str 的流，这些流引用名为 COND1n 和 COND2n 的数据文件。这些文件可在任何 IBM® SPSS® Modeler 安装程序的 Demos 目录中找到。此目录可通过 WindowsIBM® SPSS® Modeler 程序组进行访问。文件 condplot.str 和 condlearn.str 都位于目录 streams 下。

对于每个时间序列，都会对应一系列正常运行期间产生的记录，后跟一系列非正常运行期间产生的故障记录，如下表所示：

Time	幂	温度	压力	正常工作 时间	状态	结果
0	1059	259	0	404	0	0
1	1059	259	0	404	0	0
...						
51	1059	259	0	404	0	0
52	1059	259	0	404	0	0
53	1007	259	0	404	0	303
54	998	259	0	404	0	303
...						
89	839	259	0	404	0	303
90	834	259	0	404	303	303
0	965	251	0	209	0	0
1	965	251	0	209	0	0
...						
51	965	251	0	209	0	0
52	965	251	0	209	0	0
53	938	251	0	209	0	101
54	936	251	0	209	0	101
...						

Time	幂	温度	压力	正常工作 时间	状态	结果
208	644	251	0	209	0	101
209	640	251	0	209	101	101

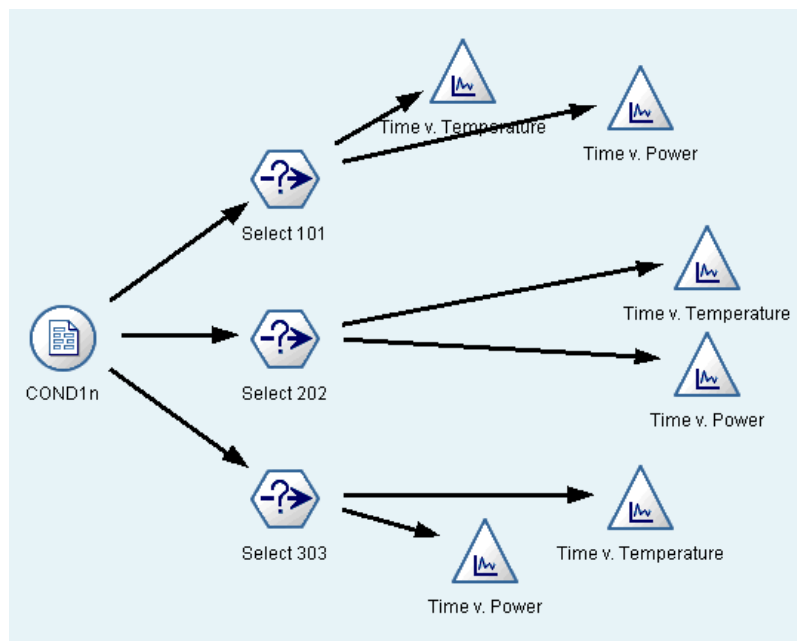
通常，大多数数据挖掘工程都会经历以下过程：

- 检查数据以确定哪些属性可能与相关状态的预测或识别有关。
- 保留这些属性（如果已存在），或者在必要时导出这些属性并将其添加到数据中。
- 使用结果数据训练规则和神经网络。
- 使用独立测试数据测试经过训练的系统。

检查数据

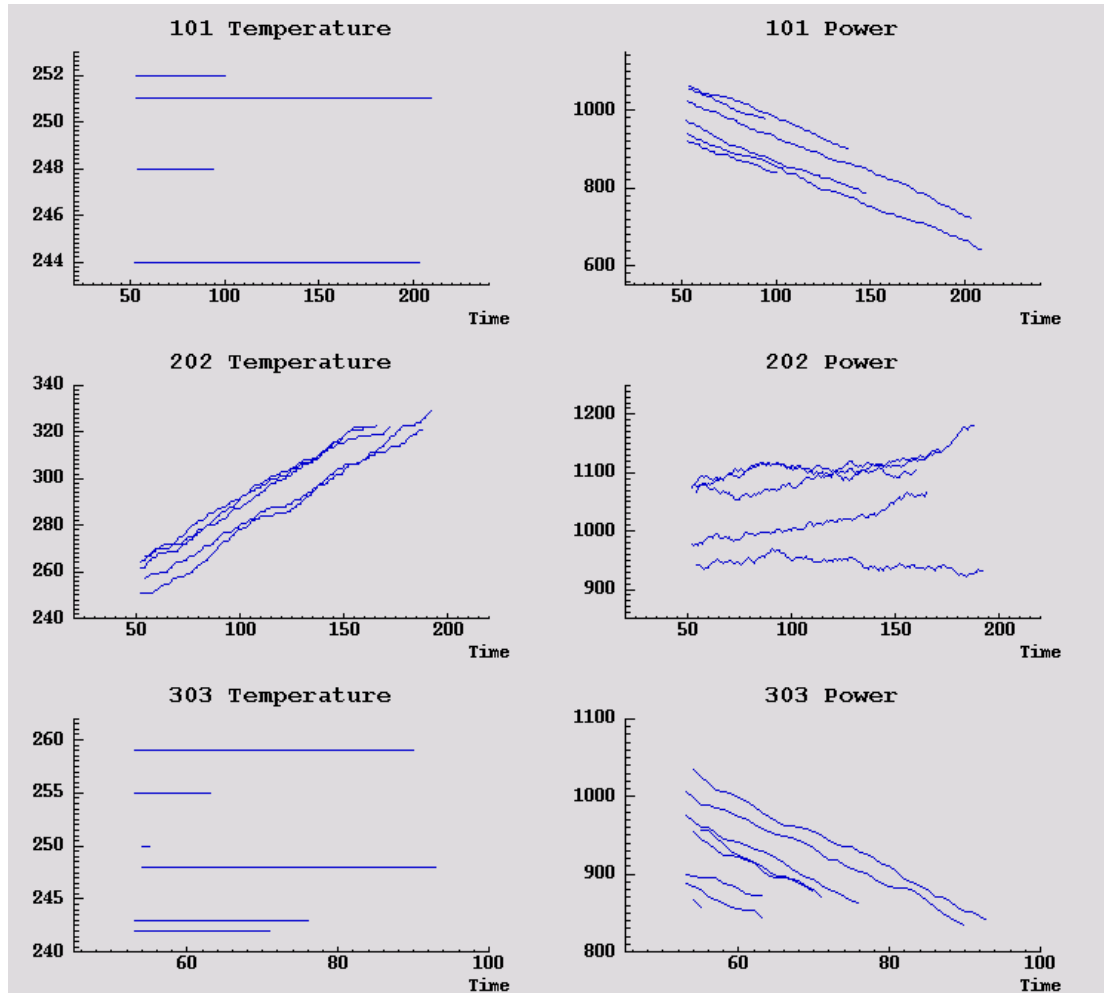
文件 condplot.str 说明上述过程的第一部分。该文件包含可绘制大量图形的流。如果温度或功率的时间序列中包含可见特性曲线，则可以识别即将发生的错误状态，还可以预测要发生的错误状态。对于温度和功率，下面的流可绘制与单独图形中的三个不同错误代码相关联的时间序列，并生成六个图形。选择节点可分隔与不同错误代码关联的数据。

图片 20-1
Condplot 流



该流的结果显示在下图中。

图片 20-2
按时间测量的温度和功率



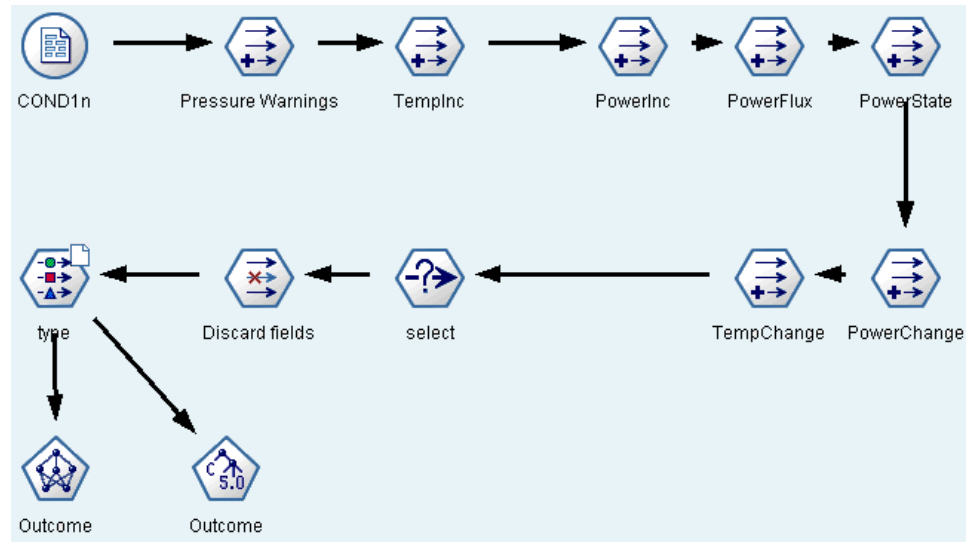
这些图形可清楚地显示将 202 错误与 101 和 303 错误区分开来的特性曲线。202 错误显示了随着时间的推移温度不断上升且功率发生波动；而其他两个错误则未显示。但是，用于区分 101 和 303 错误的特性曲线却不是很清晰。这两个错误图都显示了平滑的温度曲线和功率的下降，但 303 错误图中的功率下降显得更加急剧一些。

从上述图形可以看出，温度和功率的变化和变化率以及波动的存在和波动程度，都与预测和区别故障相关。因此应先将这些属性添加到数据，然后再应用学习系统。

数据准备

根据数据研究结果，流 condlearn.str 可导出相关数据并学习如何预测故障。

图片 20-3
Condlearn 流



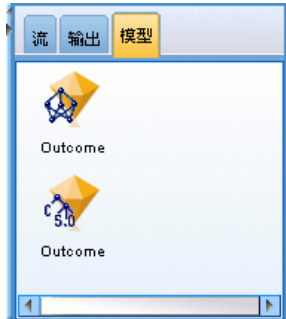
该流使用大量导出节点准备建模数据。

- **变量文件节点。** 读取数据文件 COND1n。
- **导出压力报警。** 计算瞬时压力报警数。当时间返回到 0 时重置。
- **导出 TempInc。** 使用 @DIFF1 计算温度瞬时变化率。
- **导出 Power Inc。** 使用 @DIFF1 计算功率瞬时变化率。
- **导出 PowerFlux。** 这是一个标志，如果在上一个记录和本记录中功率按相反的方向变化（即功率的峰值或波谷值），则此值为真。
- **导出 PowerState。** 起始为稳定，当检测到两个连续的功率波动时，切换为波动的状态。仅当五个时间区间内都没有出现功率波动或重置时间时，才切换回稳定状态。
- **PowerChange。** 最近五个时间区间内 PowerInc 的平均值。
- **TempChange。** 最近五个时间区间内 TempInc 的平均值。
- **丢弃初始（选择）。** 丢弃每个时间序列的第一个记录，以避免在功率和温度的边界处出现大的（不正确的）跳跃。
- **丢弃字段。** 削减记录字段，只保留正常工作时间、状态、结果、压力报警、PowerState、PowerChange 和 TempChange 字段。
- **类型。** 将结果的角色定义为**目标**（要预测的字段）。此外，将结果的测量级别定义为**名义**、将压力报警的测量级别定义为**连续**，将 PowerState 的测量级别定义为**标志**。

学习

运行 `condlearn.str` 中的流可训练 C5.0 规则和神经网络。训练网络需要一段时间，但可以早些中断训练以保存生成合理结果的网络。学习完成后，位于管理器窗口右上角的“模型”选项卡将闪烁，通知您已创建两个新的模型块：一个节点表示神经网络，一个节点表示规则。

图片 20-4
带有模型块的模型管理器



还可以将模型块添加到现有的流中，这允许我们测试系统，或导出模型结果。在此示例中，将测试模型结果。

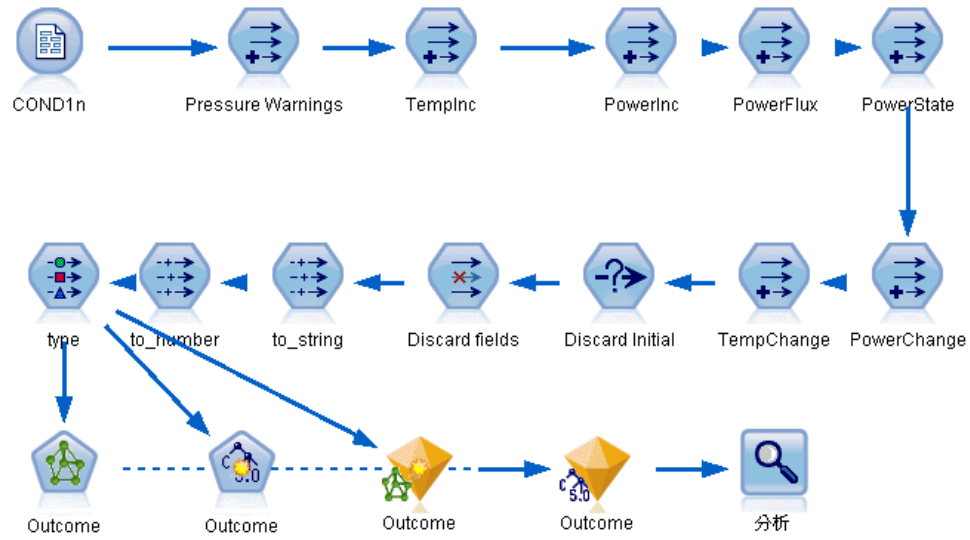
测试

模型块可以添加到流，两者均连接到类型节点。

- ▶ 如图所示，重新定位模型块，以使类型节点连接到神经网络模型块，后者连接到 C5.0 模型块。
- ▶ 将分析节点附加到 C5.0 模型块。

- ▶ 编辑原始源节点以读取文件 COND2n（而不是 COND1n），因为 COND2n 包含隐藏的测试数据。

图片 20-5
测试经过训练的网络



- ▶ 打开分析节点并单击“运行”。

这将生成可反映经过训练的网络和规则的准确性的图表。

电信客户分类（判别式分析）

判别式分析是一项根据输入字段值对记录进行分类的统计技术。这种技术与线性回归类似，但用分类目标字段代替了数值字段。

例如，假设某个电信提供商根据服务使用模式对它的客户群进行了分段，将这些客户分为了四个组。如果人口统计学数据可用于预测组成员资格，则可以为各个潜在客户自定义服务。

本示例使用的流名为 `telco_custcat_discriminant.str`，该流引用名为 `telco.sav` 的数据文件。这些文件可在任何 IBM® SPSS® Modeler 安装程序的 Demos 目录中找到。此目录可通过 WindowsIBM® SPSS® Modeler 程序组进行访问。文件 `telco_custcat_discriminant.str` 位于 `streams` 目录下。

该示例主要关注于使用人口统计数据预测使用模式。目标字段客户类别具有四个可能的值，分别对应四个客户组，如下所示：

值	Label
1	基本服务
2	电子服务
3	增值服务
4	全套服务

创建流

- ▶ 首先，设置流属性，以便在输出中显示变量和值标签。从菜单中选择：
文件 > 流属性... > 选项 > 通常

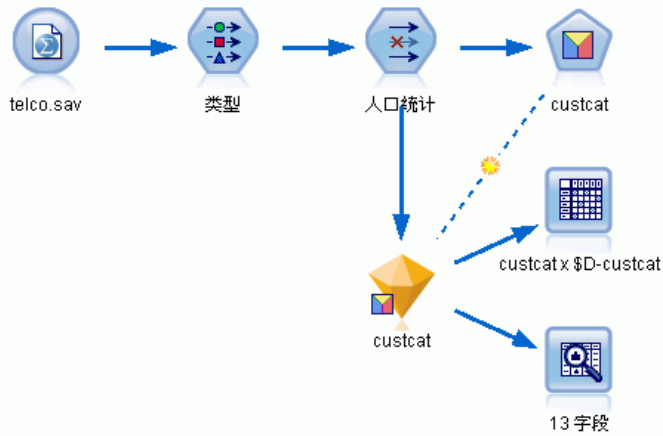
- ▶ 确保选择在输出中显示字段和值标签，然后单击确定。

图片 21-1
流属性



- ▶ 在 Demos 文件夹中添加指向 telco.sav 的 Statistics 文件源节点。

图片 21-2
使用判别式分析对客户进行分类的流示例



- ▶ 添加类型节点并单击读取值，确保所有测量水平设置正确。例如，具有值 0 和 1 的多数字段可视为标志。

图片 21-3
设置多个字段的测量级别



提示：要更改具有相似值（如 0/1）的多个字段，请单击值列标题，以便按值对字段进行排序，然后按住 Shift 键的同时使用鼠标或箭头键选择所有要更改的字段。然后可以右键单击选定的内容以更改选定字段的测量级别或其他属性。

注意，性别更准确而言应视为具有两个值的集合的字段，而不是标志，所以将其测量值保留为名义。

- ▶ 将客户类别字段的角色设置为目标。将所有其他字段的角色设置为 Input。

图片 21-4
设置字段角色



因为此示例主要关注人口统计，所以请使用过滤节点以选取相关字段（地区、年龄、婚姻状况、地址、收入、教育程度、行业、退休、性别、居住地和客户类别）。其他字段可以排除在此分析之外。

图片 21-5
过滤人口统计字段



（另外，您可以将这些字段的角色更改为无，而不要排除这些字段，或者选择要在建模节点中使用的字段。）

- ▶ 在判别节点中，单击“模型”选项卡，然后选择逐步法。

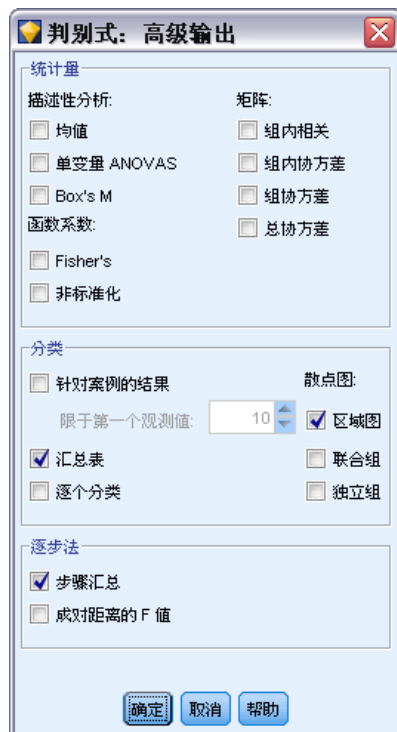
图片 21-6
选择模型选项



- ▶ 在“专家”选项卡上，将模式设置为专家，然后单击输出。

- ▶ 在“高级输出”对话框中，选择汇总表、区域图和步骤汇总，然后单击确定。

图片 21-7
选择输出选项



检查模型

- ▶ 单击运行以创建模型，该模型将添加到流和位于右上角的“模型”选项板中。要查看其详细信息，双击流中的模型块。

“汇总”选项卡显示目标（还有其他内容），以及针对考虑事项提交的完整输入（预测变量字段）列表。

图片 21-8
显示目标和输入字段的模型汇总



有关判别式分析结果的详细信息：

- ▶ 单击高级选项卡。
- ▶ 单击“在外部浏览器中启动”按钮（就在“模型”选项卡下）以在您的 Web 浏览器中查看结果。

逐步判别式分析

图片 21-9
未包含在分析步骤 0 中的变量

Step		Tolerance	Min. Tolerance	F to Enter	Wilks' Lambda
0	Age in years	1.000	1.000	7.521	.978
	Marital status	1.000	1.000	3.500	.990
	Years at current address	1.000	1.000	8.433	.975
	Household income in thousands	1.000	1.000	6.689	.980
	Level of education	1.000	1.000	61.454	.844
	Retired	1.000	1.000	3.005	.991
	Years with current employer	1.000	1.000	16.976	.951
	Gender	1.000	1.000	.373	.999
	Number of people in household	1.000	1.000	3.976	.988

拥有大量预测变量时，逐步法有助于自动选择“最适合的”用于模型的变量。逐步法的最初模型不包括任何预测变量。在每个步骤中，会将具有超出输入标准值（默认为 3.84）的最大 F to Enter 值的预测变量添加到模型中。

图片 21-10
未包含在分析步骤 3 中的变量

Step		Tolerance	Min. Tolerance	F to Enter	Wilks' Lambda
3	Age in years	.535	.535	.252	.795
	Marital status	.605	.593	1.507	.792
	Years at current address	.776	.771	3.514	.787
	Household income in thousands	.688	.657	.687	.794
	Retired	.917	.880	.353	.795
	Gender	.997	.931	.395	.795

在最后一个步骤中保留在分析之外的变量具有的 F to Enter 值都小于 3.84，因此不再向分析中添加其他变量。

图片 21-11
分析中包括的变量

Step		Tolerance	F to Remove	Wilks' Lambda
1	Level of education	1.000	61.454	
2	Level of education	.953	59.108	.951
	Years with current employer	.953	14.933	.844
3	Level of education	.951	60.046	.940
	Years with current employer	.934	15.824	.834
	Number of people in household	.979	4.841	.807

此表显示了每个步骤中包括在分析中的变量的统计信息。容差指该变量的方差中不能由方程式的其他自变量解释的部分所占比例。容差很小的变量可以向模型提供的信息很少，并且可导致计算问题。

F to Remove 值有助于说明从当前模型中删除某个变量（假设其他变量仍保留）时可能发生的情况。输入变量的 F to Remove 与上述步骤中的 F to Enter 相同（显示于“不包括在分析中的变量”表）。

有关逐步法的警告说明

逐步法很方便，但也有其局限。请注意，因为逐步法仅根据统计意义选择模型，所以它有可能选择不具有**实际意义**的预测变量。如果您比较熟悉数据并对有重要意义的预测变量有所预期，那么应该利用您的经验而不使用逐步法。但是，如果存在多个预测变量而您不知道从何处着手，则运行逐步分析法并调整选定的模型比完全没有模型要好。

检查模型拟合

图片 21-12
特征值

Function	Eigenvalue	% of Variance	Cumulative %	Canonical Correlation
1	.198	80.2	80.2	.407
2	.048	19.4	99.6	.214
3	.001	.4	100.0	.031

几乎所有由模型解释的方差都源于前两个判别函数。三个函数可自动拟合，但由于第三个函数特征值极小，可以完全忽视此函数而不用担心安全性。

图片 21-13
Wilk 的 lambda

Test of Function(s)	Wilks' Lambda	Chi-square	df	Sig.
1 through 3	.796	227.345	9	.000
2 through 3	.953	47.486	4	.000
3	.999	.929	1	.335

Wilks' lambda 认同仅有前两个函数是有用的。对于每一个函数集合，该判别式检验各组所列函数的均值相等的假设。对函数 3 的检验具有的显著性值大于 0.10，因此该函数对模型而言意义甚微。

结构矩阵

图片 21-14
结构矩阵

	Function		
	1	2	3
Level of education	.966*	-.090	-.244
Years with current employer	-.182	.964*	-.193
Age in years ^a	-.162	.598*	-.285
Household income in thousands ^a	.109	.514*	-.190
Years at current address ^a	-.151	.394*	-.214
Retired ^a	-.108	.230*	-.137
Gender ^a	.008	.054*	.009
Number of people in household	.232	.097	.968*
Marital status ^a	.132	.134	.600*

Pooled within-groups correlations between discriminating variables and standardized canonical discriminant functions
Variables ordered by absolute size of correlation within function.

*. Largest absolute correlation between each variable and any discriminant function

a. This variable not used in the analysis.

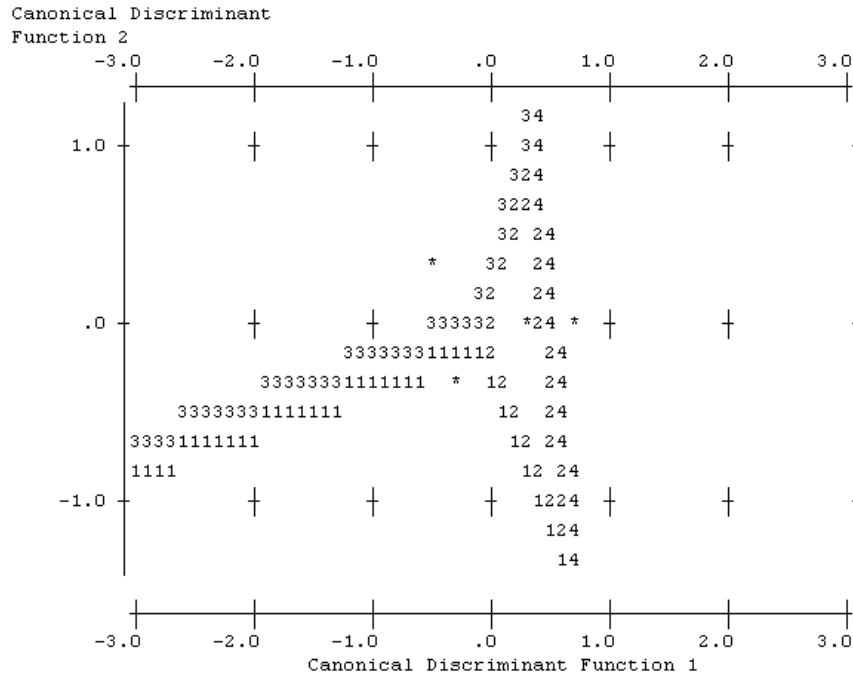
当存在多个判别式函数时，用星号来标记每个变量与某典范函数的最大绝对相关度。在每个函数内部，这些标记星号 (*) 的变量将按相关度大小排序。

- 教育程度与第一个函数具有最强相关度，并且它是与该函数具有最强相关度的唯一变量。
- 工作时间、年龄、家庭收入（以千计）、现住址居住时间、是否退休、以及性别与第二个函数具有最强相关度，而性别和是否退休与该函数的相关度比其他变量要弱许多。其他变量将该函数标记为“稳定”函数。
- 家庭成员数和婚姻状况与第三个函数具有最强相关度，但该函数是无用函数，因此这些变量是几乎无用的预测变量。

区域图

图片 21-15

区域图



区域图有助于研究组与判别式函数之间的关系。结合结构矩阵的结果，区域图能够对预测变量和组之间的关系提供图形化的解释。第一个函数，显示在水平轴上，将组 4（全套服务用户）从其他组中区分开来。因为教育程度与第一个函数具有很强的明确的关联度，这表明全套服务用户通常具有最高的教育程度。第二个函数将组 1 和 3（基本服务和附加服务用户）区分开来。附加服务用户倾向于比基本服务用户具有更长工作时间和更大的年龄。尽管区域图表明电子服务用户受过良好教育并且具有中等工作经验，但无法很好地将它与其他组区分开来。

总体而言，标记有星号（*）的组的矩心靠近区域边界时，则表明所有组间的分隔不是非常强。

区域图仅绘制了前两个判别式函数，但由于第三个函数无关紧要，因此区域图提供了判别式模型的全面视图。

分类结果

图片 21-16
分类结果

Original	Count	Customer category	Predicted Group Membership				Total
			Basic service	E-service	Plus service	Total service	
		Basic service	125	11	61	69	266
		E-service	49	15	58	95	217
		Plus service	102	14	112	53	281
		Total service	40	16	37	143	236
%		Basic service	47.0	4.1	22.9	25.9	100.0
		E-service	22.6	6.9	26.7	43.8	100.0
		Plus service	36.3	5.0	39.9	18.9	100.0
		Total service	16.9	6.8	15.7	60.6	100.0

a. 39.5% of original grouped cases correctly classified.

根据 Wilks' lambda 检验，可以得知模型的预测能力比猜测要强大，但需要借助于分类结果才能确定其强大的程度。对于给定的观测数据，“空”模型（即不包括任何预测变量的模型）将把所有用户分类到附加服务模型组。因此，空模型的正确率将是 $281/1000 = 28.1\%$ 。模型可获得较之空模型多 11.4% 即 39.5% 的用户。特别是模型在鉴别全套服务用户时表现优异。然而，它在对电子服务用户进行分类时表现得格外糟糕。可能需要寻找新的预测变量来区分这些用户。

摘要

已创建了一个判别式模型，用以基于每个用户的人口统计学信息将用户分类到四个预定义的“服务使用”组之一。利用结构矩阵和区域图，能够鉴别出那些最有助于分割客户群的变量。最后，分类结果显示模型对电子服务用户进行分类时表现欠佳。需要进一步研究来确定另一个预测变量，以便更好地对这些用户进行分类，但该模型可能完全能够满足您的需求，这取决于您希望预测的内容。例如，如果您对电子服务用户的鉴别并不关心，那么该模型可足以满足需求。这种情况可能是，将电子服务作为一种仅为吸引顾客而出售并产生微薄利润的产品。例如，如果投资的最高回报来自于附加服务或全套服务用户，则该模型能够提供所需的信息。

还请注意，这些结果仅基于训练数据。要评估该模型适用于其他数据的程度，可以使用分区节点保留部分记录，用于测试和验证。[有关详细信息，请参阅第 4 章中的分区节点中的 IBM SPSS Modeler 15 源、过程和输出节点。](#)

《IBM® SPSS® Modeler 算法指南》中列出了对 SPSS Modeler 中用到的建模方法的数学原理的说明。此文件可在安装光盘的 \Documentation 目录中找到。

分析区间型删失的生存数据（广义线性模型）

当分析区间型删失的生存数据时—即不知道所关注事件的准确时间，而只知道事件发生在给定的时间间隔内—则可将 Cox 模型应用到时间间隔内的事件危险性并生成互补重对数回归模型。

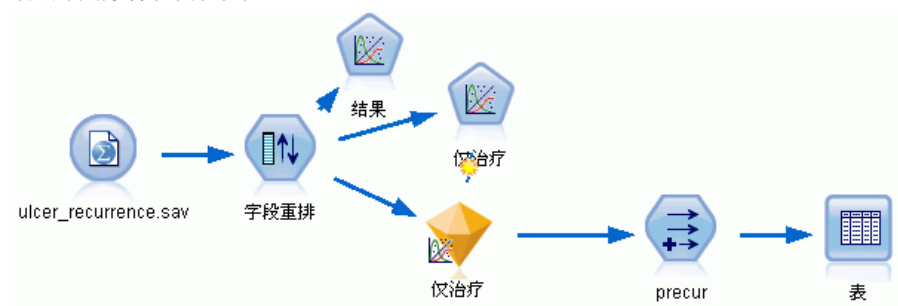
从研究（此研究的设计目的是比较两种防止溃疡复发的疗法的功效）中获取的部分信息位于 ulcer_recurrence.sav 中。此数据集已在其他地方给出并分析。使用广义线性模型，可以复制互补重对数回归模型的结果。

此示例使用名称为 ulcer_genlin.str 的流，此流参考的是数据文件 ulcer_recurrence.sav。数据文件和流文件分别位于 Demos 文件夹和 streams 子文件夹中。有关详细信息，请参阅第 1 章中的 Demos 文件夹中的 IBM SPSS Modeler 15 用户指南。

创建流

- ▶ 在 Demos 文件夹中添加一个指向 ulcer_recurrence.sav 的 Statistics 文件源节点。

图片 22-1
预测溃疡复发的样本流



- ▶ 在源节点的“过滤”选项卡上，过滤掉 id 和时间。

图片 22-2
过滤不需要的字段



- ▶ 在源节点的“类型”选项卡上，将 result 字段的角色设置为 Target，将其测量级别设置为 Flag。结果为 1 表示溃疡已复发。将所有其他字段的角色设置为 Input。
- ▶ 单击读取值以实例化数据。

图片 22-3
设置字段角色



- ▶ 添加字段重排节点并指定持续时间、治疗和年龄作为输入的顺序。此操作将确定在模型中输入字段的顺序并会帮助您尝试复制 Collett 的结果。

图片 22-4

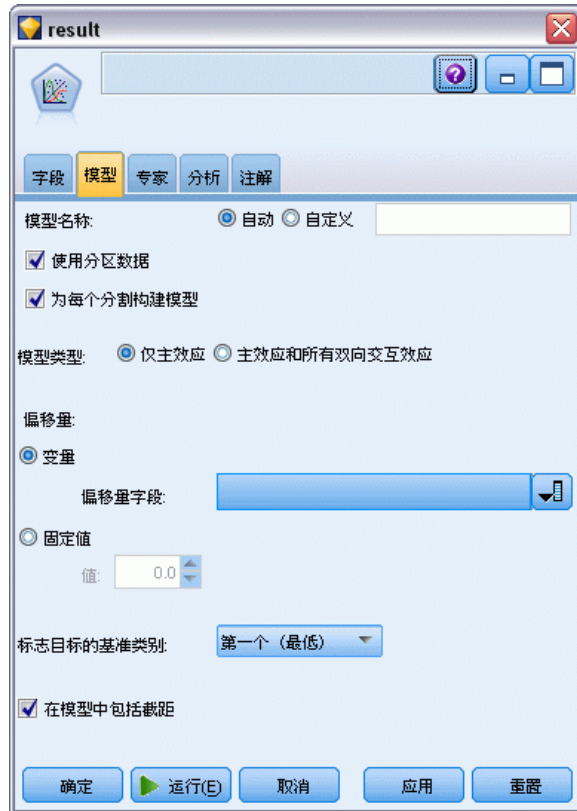
重排字段使得能够按照需要将字段输入到模型中



- ▶ 将 GenLin 节点附加到源节点；在 GenLin 节点中，单击模型选项卡。
- ▶ 选择第一个（最低值）作为目标的参考类别。此操作表示第二个类别是所关注的事件，它对模型的影响在参数估计中进行解释。系数为正的连续预测变量表示复发概率随着

预测变量值的增加而增加；相对于其他设置的类别，系数越大的名义预测变量的类别表示复发概率越大。

图片 22-5
选择模型选项



- ▶ 单击专家选项卡并选择专家以激活专家建模选项。
- ▶ 选择二项作为分布，互补重对数作为连接函数。
- ▶ 选择固定值作为估计尺度参数的方法，并选择默认值 1.0。

- ▶ 选择降序作为因子的类别顺序。这指示每个因子的第一个类别将是其参考类别；模型中此项选择的效应由参数估计值解释。

图片 22-6
选择专家选项



- ▶ 运行流以创建模型块，此模型块将被添加到流工作区和位于右上角的“模型”选项板中。要查看模型详细信息，请右键单击此模型块并选择编辑或浏览。

模型效应检验

图片 22-7
对主效应模型的模型效果测试

Source	Type III		
	Wald Chi-Square	df	Sig.
(Intercept)	.536	1	.464
duration	.003	1	.958
treatment	.382	1	.537
age	.358	1	.550

Dependent Variable: Result

Model: (Intercept), duration, treatment, age

没有任何的模型效果是在统计意义下显著的；但是，任何可观察到的治疗效果上的差异都具有临床意义，因此我们将仅以治疗作为模型项拟合一个简化模型。

拟合仅治疗模型

- ▶ 在 GenLin 节点的“字段”选项卡上，单击使用自定义设置。
- ▶ 选择结果作为目标。

- ▶ 选择治疗作为唯一的输入。

图片 22-8
选择字段选项



- ▶ 运行流，并打开生成的模型块。

在模型块上，选择高级选项卡，并滚动到底部。

参数估计值

图片 22-9
仅治疗模型的参数估计

Parameter	B	Std. Error	95% Wald Confidence Interval		Hypothesis Test		
			Lower	Upper	Wald Chi-Square	df	Sig.
(Intercept)	-1.442	.5012	-2.425	-.460	8.282	1	.004
[treatment=1]	.378	.6288	-.855	1.610	.361	1	.548
[treatment=0]	0 ^a						
(Scale)	1 ^b						

Dependent Variable: Result
Model: (Intercept), treatment

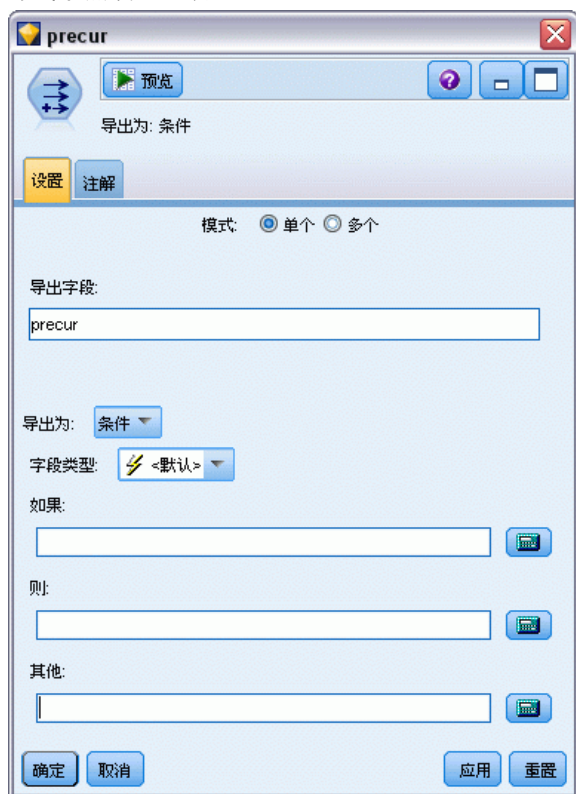
a. Set to zero because this parameter is redundant.

b. Fixed at the displayed value.

治疗效果（两种治疗水平之间的线性预测变量的差异；即[治疗 =1] 的系数）仍然不是统计意义下显著的，而是仅能估计出治疗 A [治疗 =0]可能比 B [治疗 =1]的效果好，因为治疗 B 的参数估计大于 A 的参数估计，因而与前 12 个月内复发概率增加相关联。线性预测变量（截距 + 治疗效果）是 $\log(-\log(1-P(\text{recur}_{12,t})))$ 的估计值，其中 $P(\text{recur}_{12,t})$ 是治疗 (t=A 或 B)12 个月后复发概率。可为数据集中的每个观测数据生成这些预测的可能性。

预测复发和生存的概率

图片 22-10
导出节点设置选项



- ▶ 对于每位患者，模型都可对预测结果和该预测结果的概率进行评分。为查看预测的复发概率，可将生成的模型复制到选项板并附加导出节点。
- ▶ 在“设置”选项卡中，键入 `precur` 作为导出字段。
- ▶ 选择将它作为条件导出。
- ▶ 单击计算器按钮可打开 `if` 条件的表达式构建器。

图片 22-11
导出节点：If 条件的表达式构建器



- ▶ 将 $\$G\text{-result}$ 字段插入表达式中。
- ▶ 单击确定。

导出字段 precur 在 $\$G\text{-result}$ 等于 1 时和 0 时分别取 Then 表达式和 Else 表达式的值。

图片 22-12
导出节点：Then 表达式的表达式构建器



- ▶ 单击计算器按钮可打开 Then 表达式的表达式构建器。
- ▶ 将 \$GP-result 字段插入表达式中。
- ▶ 单击确定。

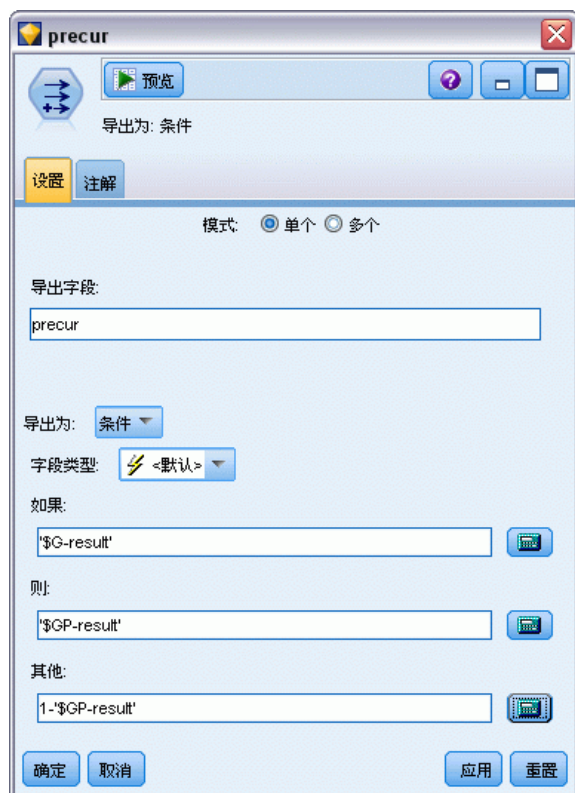
图片 22-13
导出节点：Else 表达式的表达式构建器



- ▶ 单击计算器按钮可打开 Else 表达式的表达式构建器。

- ▶ 在表达式中键入 1-, 然后将 \$GP-result 字段插入表达式中。
- ▶ 单击确定。

图片 22-14
导出节点设置选项



- ▶ 将表节点附加到导出节点并执行它。

图片 22-15
预测概率(R)

	result	duration	treatment	age	\$G-result	\$GP-result	\$GP-0	\$GP-1
1	1	2	1	48	0	0.708	0.708	0.292
2	0	1	1	73	0	0.708	0.708	0.292
3	0	1	1	54	0	0.708	0.708	0.292
4	0	2	1	58	0	0.708	0.708	0.292
5	0	1	0	56	0	0.789	0.789	0.211
6	0	2	0	49	0	0.789	0.789	0.211
7	0	1	1	71	0	0.708	0.708	0.292
8	0	1	0	41	0	0.789	0.789	0.211
9	0	1	1	23	0	0.708	0.708	0.292
10	1	1	1	37	0	0.708	0.708	0.292
11	0	1	1	38	0	0.708	0.708	0.292
12	0	2	1	76	0	0.708	0.708	0.292
13	0	2	0	38	0	0.789	0.789	0.211
14	1	1	0	27	0	0.789	0.789	0.211
15	1	1	1	47	0	0.708	0.708	0.292
16	0	1	0	54	0	0.789	0.789	0.211
17	1	1	1	38	0	0.708	0.708	0.292
18	1	2	1	27	0	0.708	0.708	0.292
19	0	2	0	58	0	0.789	0.789	0.211
20	0	1	1	75	0	0.708	0.708	0.292

对于分配到治疗 A 的患者，其在前 12 个月内病情复发概率的估计值为 0.211；对于分配到治疗 B 的患者，其概率的估计值为 0.292。注意， $1-P(\text{recur}_{12}, t)$ 是 12 个月后的生存概率，对于生存分析的专家来说这部分具有更重要的意义。

按周期对复发概率进行建模

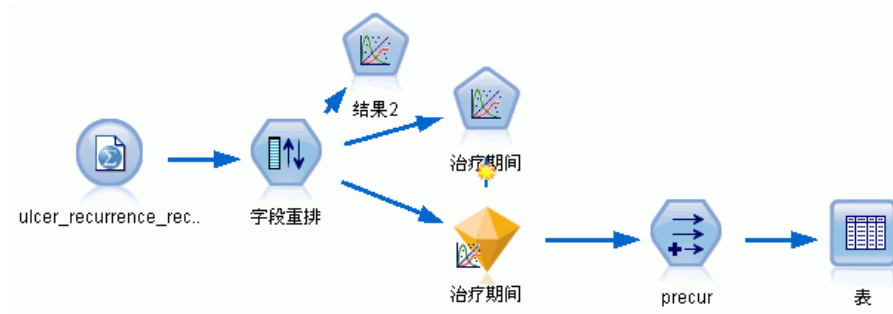
模型建立时出现的一个问题是它忽略了在第一次检查时所收集的信息；即，许多患者在前六个月内没有病情的复发。“更理想”的模型会模拟二元响应，该响应可记录事件是否会在每个时间间隔内发生。使用此模型拟合需要对原始数据集进行重建，此数据集位于 ulcer_recurrence_recoded.sav 中。[有关详细信息，请参阅第 1 章中的 Demos 文件夹中的 IBM SPSS Modeler 15 用户指南。](#) 此文件包含两个附加变量：

- 周期，它可记录病历符合第一个检查周期还是第二个检查周期。
- 按周期得到的结果，它可记录给定患者在给定周期内是否出现病情复发。

在风险集中，每个原始病历（病人）都为其所存在的每个时间间隔提供一个实例。因而，例如，患者 1 提供两个实例；一个在第一次检查周期内，此时病情没有复发，另一个在第二次检查周期内，此时记录到一次病情的复发。另一方面，患者 10 仅提供了一个实例，因为已在第一个周期内记录到病情的复发。患者 16、28 和 34 在六个月后放弃参加研究，因此仅向新数据集提供一个实例。

- ▶ 在 Demos 文件夹中添加一个指向 ulcer_recurrence_recoded.sav 的 Statistics 文件源节点。

图片 22-16
预测溃疡复发的样本流



- ▶ 在源节点的“过滤”选项卡上，过滤掉 id、时间和结果。

图片 22-17
过滤不需要的字段



- ▶ 在源节点的“类型”选项卡上，将 result2 字段的角色设置为 Target，将其测量级别设置为 Flag。将所有其他字段的角色设置为 Input。

图片 22-18
设置字段角色



- ▶ 添加字段重排节点并指定周期、持续时间、治疗和年龄作为输入的顺序。将周期作为第一个输入（不包括模型中的截距项）使您能够拟合完整的虚设变量集以捕获周期效果。

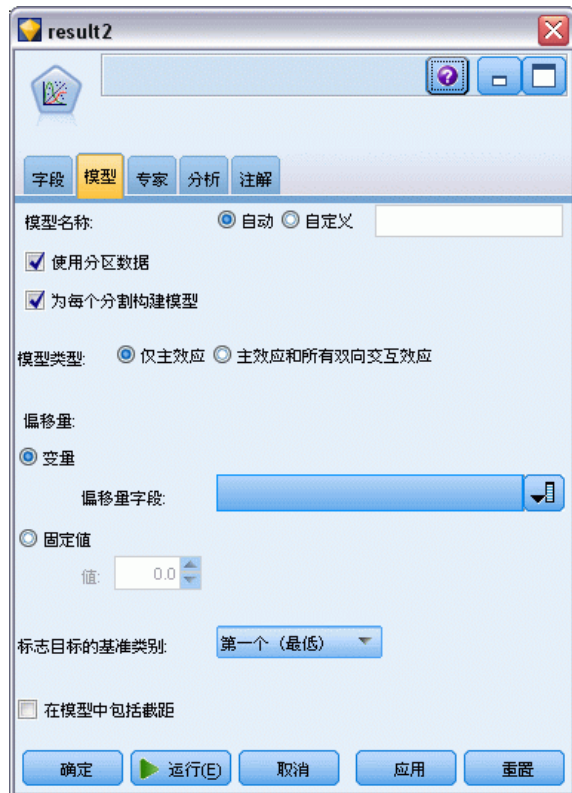
图片 22-19

重排字段使得能够按照需要将字段输入到模型中



- ▶ 在 GenLin 节点中，单击模型选项卡。

图片 22-20
选择模型选项



- ▶ 选择第一个（最低值）作为目标的参考类别。此操作表示第二个类别是所关注的事件，它对模型的影响在参数估计中进行解释。
- ▶ 取消选择 Include intercept in model。

- ▶ 单击专家选项卡并选择专家以激活专家建模选项。

图片 22-21
选择专家选项



- ▶ 选择二项作为分布，互补重对数作为连接函数。
- ▶ 选择固定值作为估计尺度参数的方法，并选择默认值 1.0。
- ▶ 选择降序作为因子的类别顺序。这指示每个因子的第一个类别将是其参考类别；模型中此项选择的效应由参数估计值解释。
- ▶ 运行流以创建模型块，此模型块将被添加到流工作区和位于右上角的“模型”选项板中。要查看模型详细信息，请右键单击此模型块并选择编辑或浏览。

模型效应检验

图片 22-22
对主效应模型的模型效果测试

Source	Type III		
	Wald Chi-Square	df	Sig.
period	.464	1	.496
duration	.000	1	.988
treatment	.117	1	.732
age	.314	1	.575

Dependent Variable: Result by period
Model: period, duration, treatment, age

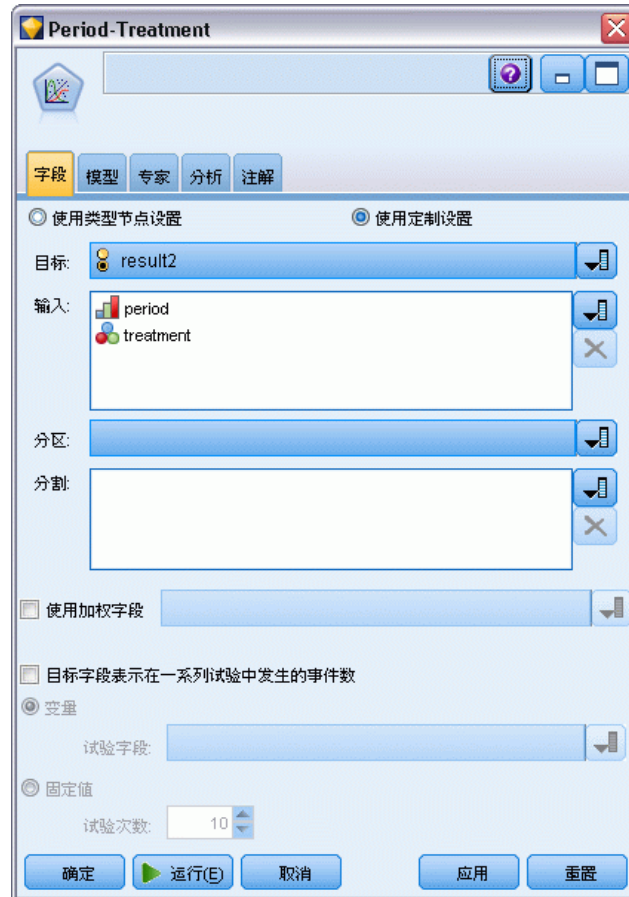
没有任何的模型效果是在统计意义下显著的；但是，任何可观察到的周期和治疗效果上的差异都具有临床意义，因此我们将仅为这些模型项拟合一个简化模型。

拟合简化模型

- ▶ 在 GenLin 节点的“字段”选项卡上，单击使用自定义设置。
- ▶ 选择结果 2 作为目标。

- ▶ 选择周期和治疗作为输入。

图片 22-23
选择字段选项



- ▶ 运行节点并浏览生成的模型，然后将生成的模型复制到选项板，附加表节点后再次运行。

参数估计值

图片 22-24
仅治疗模型的参数估计

Parameter	B	Std. Error	95% Wald Confidence Interval		Hypothesis Test		
			Lower	Upper	Wald Chi-Square	df	Sig.
[period=2]	-1.794	.5792	-2.929	-.659	9.597	1	.002
[period=1]	-2.206	.5912	-3.365	-1.047	13.926	1	.000
[treatment=1]	.195	.6279	-1.035	1.426	.097	1	.756
[treatment=0]	0 ^a
(Scale)	1 ^b

Dependent Variable: Result by period
Model: period, treatment

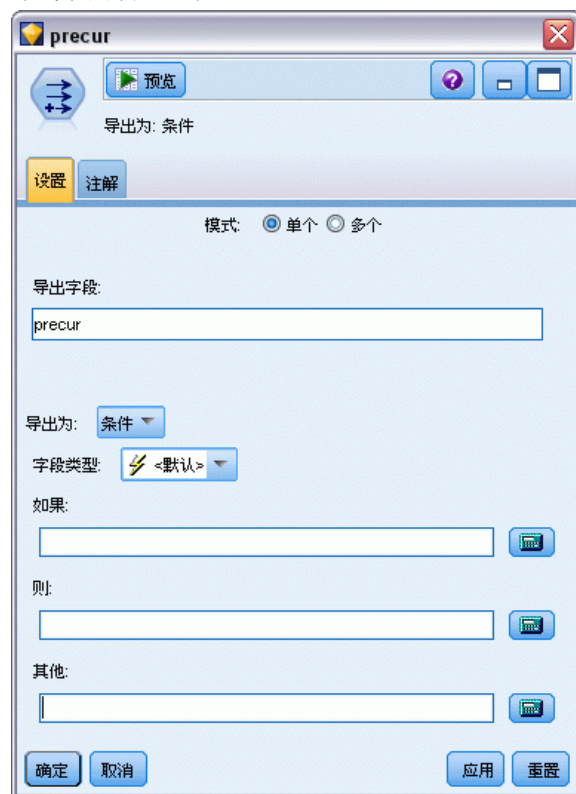
a. Set to zero because this parameter is redundant.

b. Fixed at the displayed value.

治疗效果仍然不是统计意义下显著的，而是仅能估计出治疗 A 可能比 B 的效果好，因为治疗 B 的参数估计与前 12 个月内复发的概率增加相关联。周期值在统计意义下显著地不为 0，但这是因为截距项没有拟合的缘故。周期效果（[周期 =1]和 [周期 =2] 的线性预测变量的值之间的差异）不是统计意义下显著的，这一点可以在模型效果测试中看到。线性预测变量（周期效果 + 治疗效果）是 $\log(-\log(1-P(\text{recur}_{p,t})))$ 的估计值，其中 $P(\text{recur}_{p,t})$ 是在给定治疗（ $t=A$ 或 B ）的周期（ $p=1$ 或 2 ，表示 6 个月或 12 个月）内复发的概率。可为数据集中的每个观测数据生成这些预测的可能性。

预测复发和生存的概率

图片 22-25
导出节点设置选项



- ▶ 对于每位患者，模型都可对预测结果和该预测结果的概率进行评分。为查看预测的复发概率，可将生成的模型复制到选项板并附加导出节点。
- ▶ 在“设置”选项卡中，键入 `precur` 作为导出字段。
- ▶ 选择将它作为条件导出。
- ▶ 单击计算器按钮可打开 `if` 条件的表达式构建器。

图片 22-26
导出节点：If 条件的表达式构建器



- ▶ 将 \$G-result2 字段插入表达式中。
- ▶ 单击确定。

导出字段 precur 在 \$G-result2 等于 1 时和 0 时分别取 Then 表达式和 Else 表达式的值。

图片 22-27
导出节点：Then 表达式的表达式构建器



- ▶ 单击计算器按钮可打开 Then 表达式的表达式构建器。
- ▶ 将 \$GP-result2 字段插入表达式中。
- ▶ 单击确定。

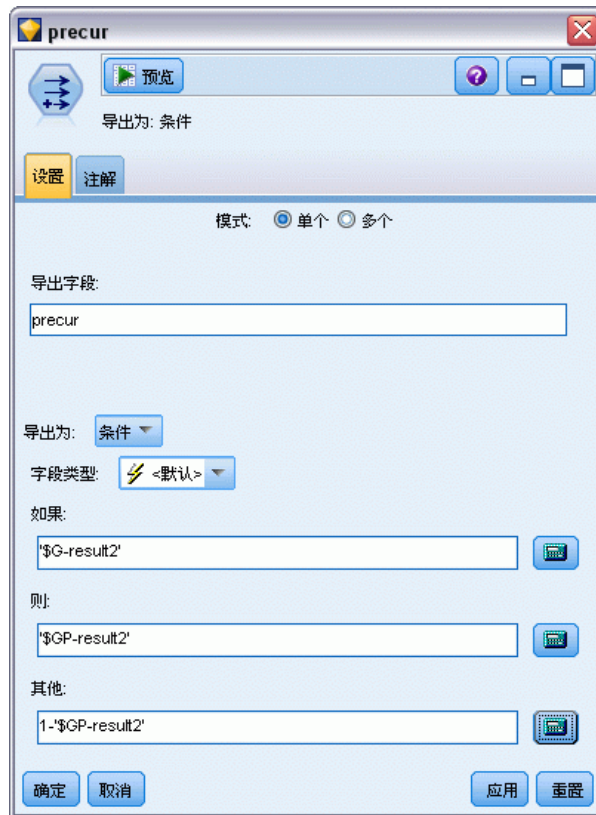
图片 22-28
导出节点：Else 表达式的表达式构建器



- ▶ 单击计算器按钮可打开 Else 表达式的表达式构建器。

- ▶ 在表达式中键入 1-, 然后将 \$GP-result2 字段插入表达式中。
- ▶ 单击确定。

图片 22-29
导出节点设置选项



- ▶ 将表节点附加到导出节点并执行它。

图片 22-30
预测概率(R)

	result2	period	duration	treatment	age	\$G-result2	\$GP-result2	\$GP-0	\$GP-1
1	0	1	2	1	48	0	0.875	0.875	0.125
2	1	2	2	1	48	0	0.817	0.817	0.183
3	0	1	1	1	73	0	0.875	0.875	0.125
4	0	2	1	1	73	0	0.817	0.817	0.183
5	0	1	1	1	54	0	0.875	0.875	0.125
6	0	2	1	1	54	0	0.817	0.817	0.183
7	0	1	2	1	58	0	0.875	0.875	0.125
8	0	2	2	1	58	0	0.817	0.817	0.183
9	0	1	1	0	56	0	0.896	0.896	0.104
10	0	2	1	0	56	0	0.847	0.847	0.153
11	0	1	2	0	49	0	0.896	0.896	0.104
12	0	2	2	0	49	0	0.847	0.847	0.153
13	0	1	1	1	71	0	0.875	0.875	0.125
14	0	2	1	1	71	0	0.817	0.817	0.183
15	0	1	1	0	41	0	0.896	0.896	0.104
16	0	2	1	0	41	0	0.847	0.847	0.153
17	0	1	1	1	23	0	0.875	0.875	0.125
18	0	2	1	1	23	0	0.817	0.817	0.183
19	1	1	1	1	37	0	0.875	0.875	0.125
20	0	1	1	1	38	0	0.875	0.875	0.125

估计的复发概率可总结如下：

治疗	6 个月	12 个月
A	0.104	0.153
B	0.125	0.183

根据这些数据，可将 12 个月后的生存概率估计为 $1 - (P(\text{recur}_{1, t}) + P(\text{recur}_{2, t}) \times (1 - P(\text{recur}_{1, t})))$ ；因此，对于每种治疗有如下结果：

$$A: 1 - (0.104 + 0.153 \times 0.896) = 0.759$$

$$B: 1 - (0.125 + 0.183 \times 0.875) = 0.715$$

再一次显示出 A 作为更理想的治疗不是统计意义下显著的。

摘要

使用广义线性模型，已通过一系列互补重对数回归模型拟合了区间型删失的生存数据。虽然对于选择治疗 A 显示出一定的支持，但要取得统计意义下显著的结果还需要更大量的研究。不过，研究现有的数据还有一些其他方法。

- 值得一试的是使用模型重新拟合交互效应，尤其是周期和治疗组之间的交互效应。

《SPSS Modeler 算法指南》中列出了对 IBM® SPSS® Modeler 中用到的建模方法的数学原理的说明。

使用泊松回归来分析船只损坏率 (广义线性模型)

广义线性模型能够用来为计数数据的分析拟合泊松回归。例如，在别处被提出和分析的关于波浪对货船造成的损坏的数据集。如果有预测变量的值，便可将事件计数的模型建为以泊松比率发生，而且结果模型可以帮助您确定哪种类型的船最容易损坏。

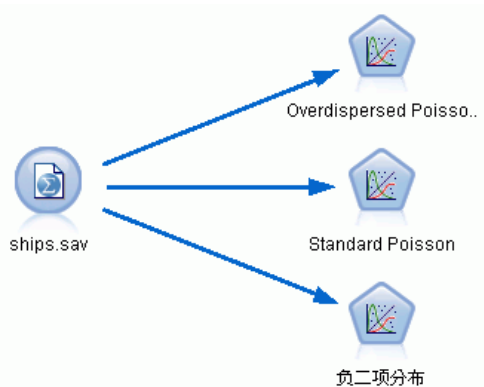
本示例使用流 `ships_genlin.str`，该流引用了数据文件 `ships.sav`。数据文件和流文件分别位于 `Demos` 文件夹和 `streams` 子文件夹中。有关详细信息，请参阅第 1 章中的 `Demos` 文件夹中的 `IBM SPSS Modeler 15 用户指南`。

由于分类汇总服务月数会随船只类型而变化，因此，在这种情况下为原始单元格计数建模会使人产生误解。像这种测量“遭遇”风险的数量的变量在广义线性模型中被按照偏移变量来处理。此外，泊松回归假设因变量的对数在预测变量中为线性。因此，要使用广义线性模型来将泊松回归拟合到事故率，您需要使用 `Logarithm of aggregate months of service`。

拟合“高度离散的”泊松回归

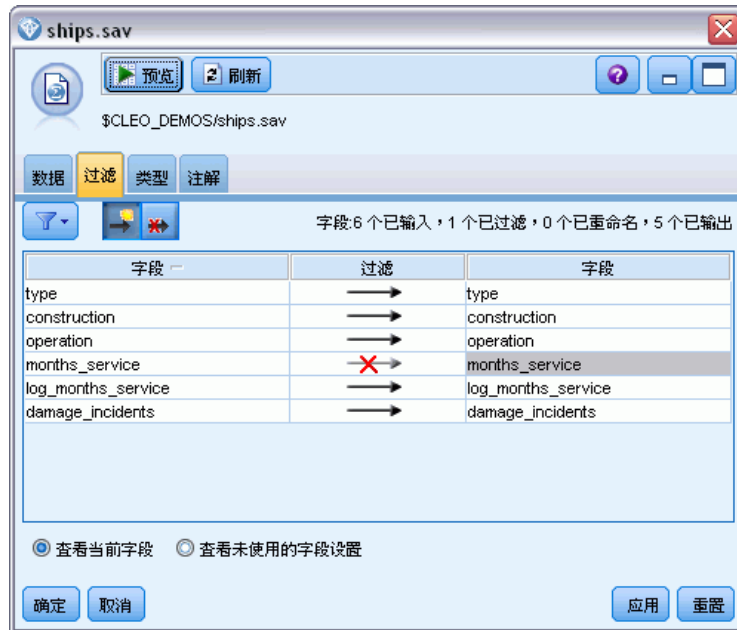
- ▶ 在 `Demos` 文件夹中添加指向 `ships.sav` 的 `Statistics` 文件源节点。

图片 23-1
取样流来分析损坏率



- ▶ 在源节点的“过滤”选项卡上，排除字段 `months_service`。该变量的经对数转换的值包含在 `log_months_service` 中，这些值将在分析中使用。

图片 23-2
过滤不需要的字段



（或者，也可以将“类型”选项卡上该字段的角色改为无而不是排除该字段，或选择您要在模型节点中使用的字段。）

- ▶ 在源节点的类型选项卡中，将 `damage_incidents` 字段的角色设置为 `Target`。将所有其他字段的角色设置为 `Input`。

- ▶ 单击读取值以实例化数据。

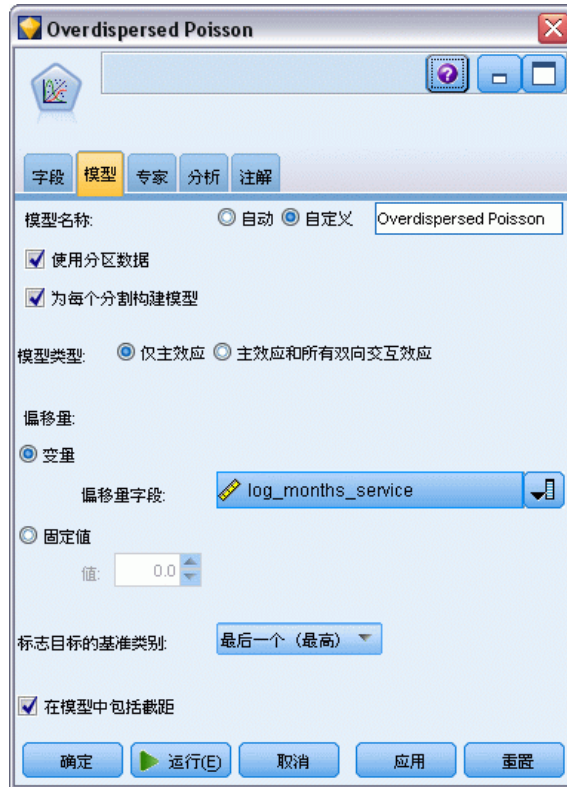
图片 23-3
设置字段角色



- ▶ 将 GenLin 节点附加到源节点；在 GenLin 节点中，单击模型选项卡。

- ▶ 选择 `log_months_service` 作为偏移变量。

图片 23-4
选择模型选项



- ▶ 单击专家选项卡并选择专家以激活专家建模选项。

图片 23-5
选择专家选项



- ▶ 选择泊松作为响应的分布，并选择对数作为关联函数。
- ▶ 选择 Pearson 卡方作为估计尺度参数的方法。尺度参数在泊松回归中通常假设为 1，但 McCullagh 和 Nelder 却用 Pearson 卡方估计来获得更保守的方差估计值和显著性水平。
- ▶ 选择降序作为因子的类别顺序。这指示每个因子的第一个类别将是其参考类别；模型中此项选择的效应由参数估计值解释。
- ▶ 单击运行以创建模型块，此模型块将被添加到流工作区和位于右上角的“模型”选项板中。要查看模型详细信息，请右键单击此模型块并选择编辑或浏览，然后单击高级选项卡。

拟合优度统计

图片 23-6
拟合优度统计量

	Value	df	Value/df
Deviance	38.695	25	1.548
Scaled Deviance	22.883	25	
Pearson Chi-Square	42.275	25	1.691
Scaled Pearson Chi-Square	25.000	25	
Log Likelihood ^a	-68.281		
Akaike's Information Criterion (AIC)	154.562		
Finite Sample Corrected AIC (AICC)	162.062		
Bayesian Information Criterion (BIC)	168.299		
Consistent AIC (CAIC)	177.299		

Dependent Variable: Number of damage incidents
Model: (Intercept), type, construction, operation, offset = log_months_
service

- a. The full log likelihood function is displayed and used in computing information criteria.
b. Information criteria are in small-is-better form.

拟合优度统计表提供了对于比较竞争模型很有用的度量。此外，偏差和 Pearson 卡方统计量的 Value/df 值给出了对尺度参数的相应估计。泊松回归中的这些值应该接近 1.0；这些值大于 1.0 的事实表明拟合高度离散的模型也许是合理的。

Omnibus 检验

图片 23-7
Omnibus 检验

Likelihood Ratio Chi-Square	df	Sig.
107.633	8	.000

Dependent Variable: Number of damage incidents
Model: (Intercept), type, construction, operation, offset = log_
months_service

- a. Compares the fitted model against the intercept-only model.

Omnibus 检验是当前模型与零（此个案中为截距）模型的似然比卡方检验。小于 0.05 的显著性值表明当前模型的性能要高于零模型的性能。

模型效应检验

图片 23-8
模型效应检验

Source	Type III		
	Wald Chi-Square	df	Sig.
(Intercept)	2138.657	1	.000
type	15.415	4	.004
construction	17.242	3	.001
operation	6.249	1	.012

Dependent Variable: Number of damage incidents
Model: (Intercept), type, construction, operation, offset = log_months_service

检验模型中的每一项是否具有效应。显著性值小于 0.05 的项具有一定可辨别效应。每个主效应项都对模型有贡献。

参数估计值

图片 23-9
参数估计

Parameter	B	Std. Error	95% Wald Confidence Interval		Hypothesis Test		
			Lower	Upper	Wald Chi-Square	df	Sig.
(Intercept)	-6.406	.2828	-6.960	-5.852	513.238	1	.000
[type=5]	.326	.3067	-.276	.927	1.127	1	.288
[type=4]	-.076	.3779	-.817	.665	.040	1	.841
[type=3]	-.687	.4279	-1.526	.151	2.581	1	.108
[type=2]	-.543	.2309	-.996	-.091	5.536	1	.019
[type=1]	0 ^a
[construction=75]	.453	.3032	-.141	1.048	2.236	1	.135
[construction=70]	.818	.2208	.386	1.251	13.743	1	.000
[construction=65]	.697	.1946	.316	1.079	12.835	1	.000
[construction=60]	0 ^a
[operation=75]	.384	.1538	.083	.686	6.249	1	.012
[operation=60]	0 ^a
(Scale)	1.691 ^b

Dependent Variable: Number of damage incidents
Model: (Intercept), type, construction, operation, offset = log_months_service

a. Set to zero because this parameter is redundant.

b. Computed based on the Pearson chi-square.

参数估计值表汇总了每个预测变量的作用。关联函数的性质决定了此模型的系数难于解释，但是，通过协变量系数的符号和因子水平系数的相对值，还是可以获得模型预测变量效应的重要信息。

- 对于协变量，正（负）系数表示预测变量和结果的正（反）关系。系数为正的协变量值的增加对应于损坏事件比率的增加。
- 对于因子，系数越大的因子水平表示损坏的发生率越高。因子水平的系数符号取决于因子水平相对于参考类别的效应。

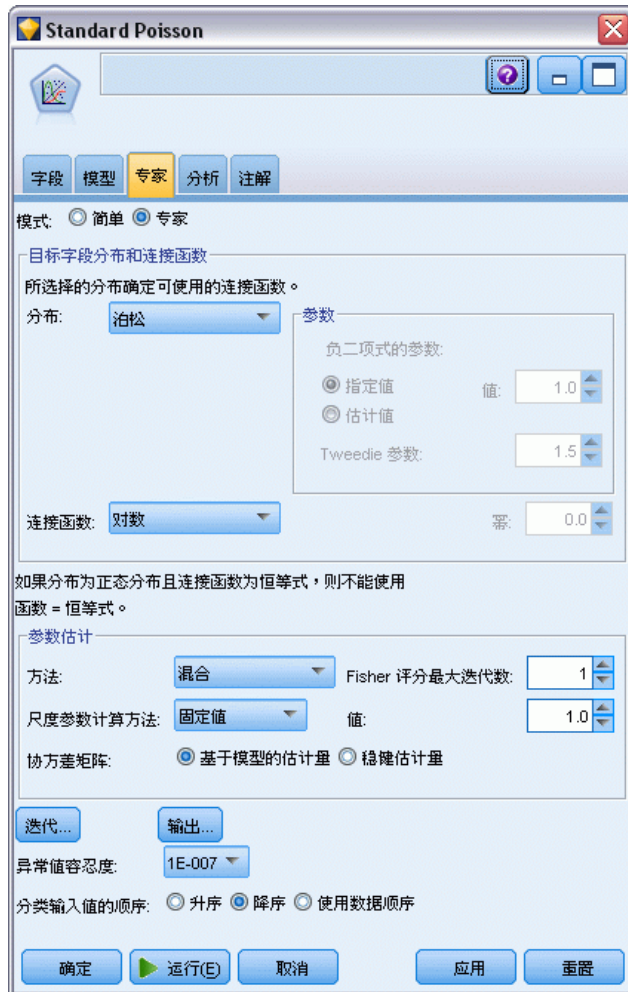
根据参数估计值，可以进行以下解释：

- 船只类型 B [type=2] 比类型 A [type=1]（参考类别）具有统计显著性上（p 值为 0.019）较低的损坏率（估计系数为 -0.543）。类型 C [type=3] 的估计参数实际要低于 B，但是 C 的估计值中的变异性遮盖了效应。有关因子水平间所有关系的解释，请参阅估计边际均值。
- 1965 - 69 [construction=65] 和 1970 - 74 [construction=70] 年间建造的船只比 1960 - 64 [construction=60]（参考类别）年间建造的船只具有统计显著性上（p 值为 <0.001）较高的损坏率（估计系数分别为 0.697 和 0.818）。有关因子水平间所有关系的解释，请参阅估计边际均值。
- 航运时间为 1975 - 79 [航运时间 =75] 的船只的损伤率（估计系数为 0.384）在统计意义下显著（p 为 0.012）高于航运时间为 1960 - 1974 [航运时间 =60] 的船只。

拟合其他模型

“高度离散的”泊松回归存在一个问题，即没有正式的方式来检验它与“标准”泊松回归。但是，建议的用来确定是否具有高度离散的正式检验是在所有其他设置相同的情况下执行“标准”泊松回归和负二项式回归之间的似然比检验。如果泊松回归中不具有高度离散，那么统计量 $-2 \times (\text{泊松模型的对数似然} - \text{负二项式模型的对数似然})$ 应该具有混合分布（其一半的概率集中在 0，其余的则在卡方分布中且自由度为 1）。

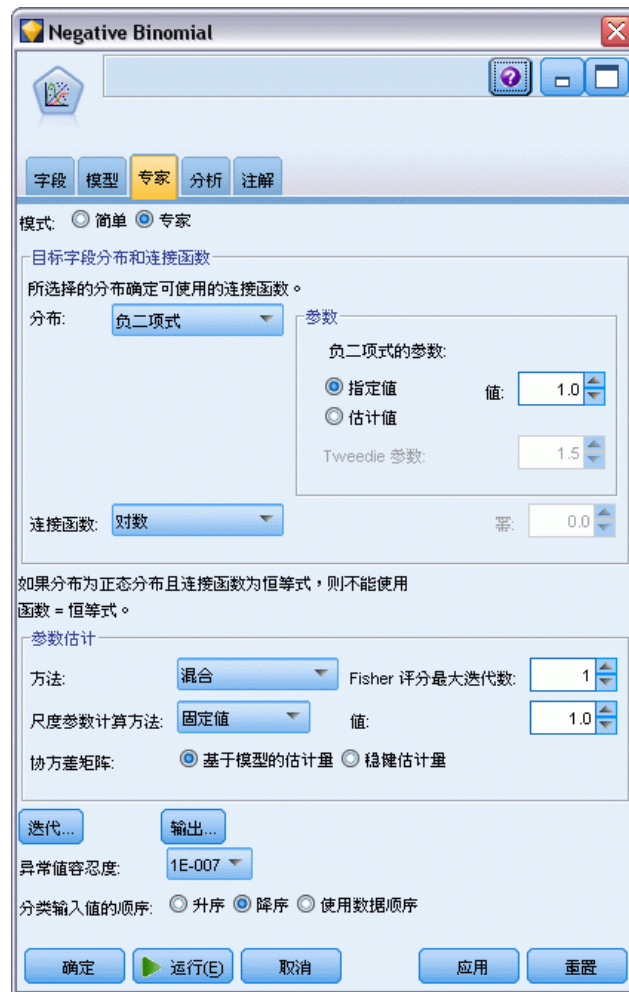
图片 23-10
“专家”选项卡



要拟合“标准”泊松回归，可复制并粘贴 GenLin 节点，将其附加到源节点，然后打开新节点，并单击专家选项卡。

- ▶ 选择固定值作为估计尺度参数的方法。默认情况下，该值为 1。

图片 23-11
“专家”选项卡



- ▶ 要拟合负二项回归，可复制并粘贴 GenLin 节点，将其附加到源节点，然后打开新节点，并单击专家选项卡。
- ▶ 选择负二项式作为分布。保留辅助参数的默认值为 1。
- ▶ 运行流，并在新创建的模型块上浏览“高级”选项卡。

拟合优度统计

图片 23-12
标准泊松回归的拟合优度统计量

	Value	df	Value/df
Deviance	38.695	25	1.548
Scaled Deviance	38.695	25	
Pearson Chi-Square	42.275	25	1.691
Scaled Pearson Chi-Square	42.275	25	
Log Likelihood ^a	-68.281		
Akaike's Information Criterion (AIC)	154.562		
Finite Sample Corrected AIC (AICC)	162.062		
Bayesian Information Criterion (BIC)	168.299		
Consistent AIC (CAIC)	177.299		

Dependent Variable: Number of damage incidents
Model: (Intercept), type, construction, operation, offset = log_months_
service

- a. The full log likelihood function is displayed and used in computing information criteria.
- b. Information criteria are in small-is-better form.

报告的标准泊松回归的对数似然为 -68.281。将该值与负二项式模型进行比较。

图片 23-13
负二项式回归的拟合优度统计量

	Value	df	Value/df
Deviance	11.145	25	.446
Scaled Deviance	11.145	25	
Pearson Chi-Square	8.815	25	.353
Scaled Pearson Chi-Square	8.815	25	
Log Likelihood ^a	-83.725		
Akaike's Information Criterion (AIC)	185.450		
Finite Sample Corrected AIC (AICC)	192.950		
Bayesian Information Criterion (BIC)	199.187		
Consistent AIC (CAIC)	208.187		

Dependent Variable: Number of damage incidents
Model: (Intercept), type, construction, operation, offset = log_months_
service

- a. The full log likelihood function is displayed and used in computing information criteria.
- b. Information criteria are in small-is-better form.

报告的负二项式回归的对数似然为 -83.725。而实际上负二项式回归的对数似然要小于泊松回归的对数似然，这就表示该负二项式回归没有提供优于泊松回归的改进。

但是，选择负二项式分布的辅助参数为 1 也许不是非常适合该数据集。检验高度离散的另一种方法是以辅助参数为 0 来拟合负二项式模型，并请求专家选项卡的输出对话框上的 Lagrange 乘数检验。如果该检验不显著，则过散布对于该数据集不是问题。

摘要

通过使用“广义线性模型”，您为计数数据拟合了三个不同的模型。负二项回归表明它没有对泊松回归进行任何改进。高度离散的泊松回归似乎可以合理地替代标准泊松模型，但对于两者间的选择还没有正式的检验。

《SPSS Modeler 算法指南》中列出了对 IBM® SPSS® Modeler 中用到的建模方法的数学原理的说明。

将 Gamma 回归拟合至汽车保险理赔 (广义线性模型)

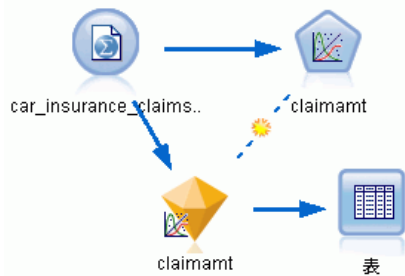
广义线性模型可以被应用于拟合正范围数据分析的 Gamma 回归。例如，在其他地方出现和分析的数据集与汽车的损伤理赔有关。平均理赔金额可以当作其具有 gamma 分布来建模，通过使用逆联接函数将因变量的均值与预测值的线性组合关联。出于考虑到用于计算平均理赔金额的理赔数目不同，指定理赔数作为尺度权重。

本示例使用命名为 `car-insurance_genlin.str` 的流，它引用命名为 `car_insurance_claims.sav` 的数据文件。数据文件和流文件分别位于 Demos 文件夹和 streams 子文件夹中。有关详细信息，请参阅第 1 章中的 Demos 文件夹中的 IBM SPSS Modeler 15 用户指南。

创建流

- ▶ 在 Demos 文件夹中添加指向 `car_insurance_claims.sav` 的 Statistics 文件源节点。

图片 24-1
预测汽车保险理赔的示例流



- ▶ 在源节点的“类型”选项卡中，将 `claimamt` 字段的角色设置为 **Target**。将所有其他字段的角色设置为 **Input**。

- ▶ 单击读取值以实例化数据。

图片 24-2
设置字段角色



- ▶ 将 GenLin 节点附加到源节点；在 GenLin 节点中，单击“字段”选项卡。

- ▶ 选择 nclaims 作为尺度权重字段。

图片 24-3
选择字段选项



- ▶ 单击专家选项卡并选择专家以激活专家建模器选项。

图片 24-4
选择专家选项



- ▶ 选择 Gamma 作为响应分布。
- ▶ 选择幂作为关联函数，并键入 -1.0 作为幂函数的指数。这是一个逆联接。
- ▶ 选择 Pearson 卡方作为估计尺度参数的方法。这是 McCullagh 和 Nelder 应用的方法，因此我们在此沿用它来精确重现其结果。
- ▶ 选择降序作为因子的类别顺序。这指示每个因子的第一个类别将是其参考类别；模型中此项选择的效应由参数估计值解释。
- ▶ 单击运行以创建模型块，此模型块将被添加到流工作区和位于右上角的“模型”选项板中。要查看模型详细信息，请右键单击此模型块并选择编辑或浏览，然后选择“高级”选项卡。

参数估计值

图片 24-5
参数估计

Parameter	B	Std. Error	95% Wald Confidence Interval		Hypothesis Test		
			Lower	Upper	Wald Chi-Square	df	Sig.
(Intercept)	.003411	.000418	.002591	.004230	66.593	1	.000
[holderage=8]	.000920	.000416	.000105	.001735	4.898	1	.027
[holderage=7]	.000916	.000408	.000117	.001716	5.046	1	.025
[holderage=6]	.000969	.000405	.000176	.001763	5.740	1	.017
[holderage=5]	.001370	.000419	.000548	.002192	10.682	1	.001
[holderage=4]	.000462	.000411	-.000342	.001267	1.268	1	.260
[holderage=3]	.000350	.000412	-.000458	.001158	.720	1	.396
[holderage=2]	.000101	.000436	-.000754	.000956	.054	1	.816
[holderage=1]	.000000 ^a
[vehiclegroup=4]	-.001421	.000181	-.001775	-.001067	61.883	1	.000
[vehiclegroup=3]	-.000614	.000170	-.000947	-.000281	13.039	1	.000
[vehiclegroup=2]	.000038	.000169	-.000293	.000368	.050	1	.823
[vehiclegroup=1]	.000000 ^a
[vehicleage=4]	.004154	.000442	.003287	.005021	88.175	1	.000
[vehicleage=3]	.001651	.000227	.001207	.002096	53.013	1	.000
[vehicleage=2]	.000366	.000101	.000169	.000564	13.191	1	.000
[vehicleage=1]	.000000 ^a
(Scale)	1.209 ^b	.	.	.001	.0004	.000	.002

Dependent Variable: Average cost of claims

Model: (Intercept), holderage, vehiclegroup, vehicleage

a. Set to zero because this parameter is redundant.

b. Computed based on the Pearson chi-square.

Omnibus 检验和模型效应检验（未显示）表明这个模型的性能要高于“零模型”的性能而且每个主效应项都对此模型起作用。参数估计值表显示与 McCullagh 和 Nelder 获得的因子水平相同的值和刻度参数。

摘要

通过使用广义线性模型，可以将 gamma 回归拟合至理赔数据。请注意：gamma 分布的典型关联函数被应用于此模型中的同时，对数链接也将给出合理结果。通常，很难甚至根本就不可能直接将模型与不同的关联函数进行比较；不过，对数链接是一个特例，在对数链接中，指数为零，因此，您可以将一个模型中的偏差与一个对数链接以及一个具有幂关联模型进行比较，以确定哪一个能更好的拟合（例如参阅 McCullagh 和 Nelder 的第 11.3 部分）。

《SPSS Modeler 算法指南》中列出了对 IBM® SPSS® Modeler 中用到的建模方法的数学原理的说明。

细胞样本分类（SVM）

Support Vector Machine (SVM) 是一种特别适用于大型数据集的分类和回归技术。大型数据集具有大量预测变量，例如可能会在生物信息学领域遇到（对生物化学和生物学数据应用信息技术）的预测变量。

一位医学研究人员获得了一个包含大量人体细胞样本特征的数据集，这些样本是从极有可能患上癌症的患者身上提取的。通过对原始数据进行分析，发现良性样本与恶性样本之间的许多特征显著不同。该研究人员希望开发一个 SVM 模型，使该模型可以使用其他患者样本中的这些细胞特征值尽早发现他们的样本是良性还是恶性。

本示例使用了名为 `svm_cancer.str` 的流，该流位于 Demos 文件夹下的 `streams` 子文件夹中。数据文件为 `cell_samples.data`。有关详细信息，请参阅第 1 章中的 Demos 文件夹中的 IBM SPSS Modeler 15 用户指南。

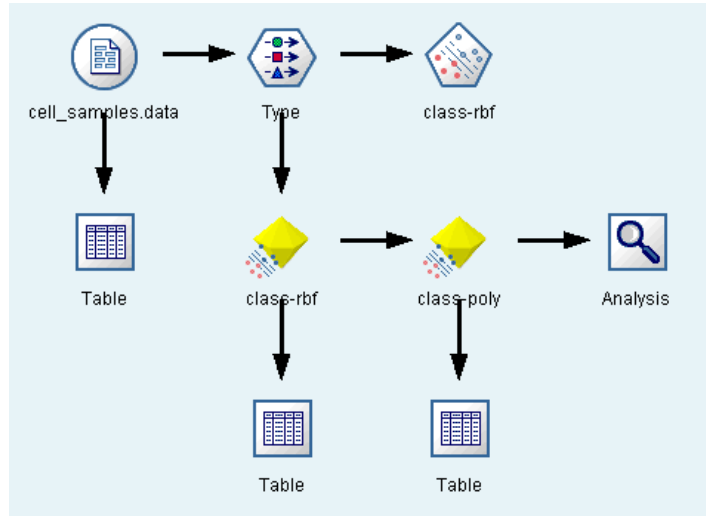
本示例基于可从 UCI Machine Learning Repository (Asuncion 和 Newman, 2007) 上公开获取的数据集。数据集由数百条人体细胞样本记录组成，每条记录均包含一组细胞特征值。每条记录中包含的字段包括：

字段名	描述
ID	患者标识符
肿块	肿块的厚度
UnifSize	细胞大小的均匀度
UnifShape	细胞大小的均匀度
MargAdh	边缘的粘连
SingEpiSize	单层上皮细胞的大小
BareNuc	裸核
BlandChrom	温和的染色质
NormNucl	正常的核仁
Mit	有丝分裂
Class	良性或恶性

为达到本示例的目的，我们使用的是每条记录包含相对较少预测变量的数据集。

创建流

图片 25-1
显示 SVM 建模的样本流



- ▶ 创建一个新流，然后在 IBM® SPSS® Modeler 安装程序的 Demos 文件夹中添加一个指向 cell_samples.data 的变量文件源节点。

让我们看一下源文件中的数据。

- ▶ 为流添加表节点。
- ▶ 将表节点附加到变量文件节点并运行流。

图片 25-2
SVM 的源数据

	Clump	Uniformity of Cell Size	Uniformity of Cell Shape	Marginal Adhesion	Single Epithelial Cell Size	Bare Nuclei	Chromatin Clumping	Normal Nuclei	Mitoses	Class
1	1	1	1	2	1	3	1	1	2	
2	4	5	7	10	3	2	1	2		
3	1	1	2	2	3	1	1	2		
4	8	1	3	4	3	7	1	2		
5	1	3	2	1	3	1	1	2		
6	10	8	7	10	9	7	1	4		
7	1	1	2	10	3	1	1	2		
8	2	1	2	1	3	1	1	2		
9	1	1	2	1	1	1	5	2		
10	1	1	2	1	2	1	1	2		
11	1	1	1	1	3	1	1	2		
12	1	1	2	1	2	1	1	2		
13	3	3	2	3	4	4	1	4		
14	1	1	2	3	3	1	1	2		
15	5	10	7	9	5	5	4	4		
16	6	4	6	1	4	3	1	4		
17	1	1	2	1	2	1	1	2		
18	1	1	2	1	3	1	1	2		
19	7	6	4	10	4	1	2	4		
20	1	1	2	1	3	1	1	2		

ID 字段包含患者的标识符。来自每位患者的细胞样本特征包含在从 Clump 到 Mit 的字段中。这些字段的值按照 1 到 10 进行分级，值为 1 表示最接近于良性。

Class 字段包含诊断，由多步独立的医疗程序确认，用于表明样本是良性（值 = 2）还是恶性（值 = 4）。

图片 25-3
类型节点设置



▶ 添加一个类型节点并将它附加到变量文件节点。

▶ 打开该类型节点。

我们希望模型预测 Class 的值（即，良性 (=2) 还是恶性 (=4)）。由于此字段可以为仅有的这两个可能值之一，我们需要更改其测量级别以反映这一情况。

▶ 在类字段的测量列（列表中最后一个），单击值连续并将其更改为标志。

▶ 单击读取值。

▶ 在角色列中，将 ID（患者的标识符）的角色设置为无，因为此字段将不会用作预测变量或模型的目标。

▶ 将目标 Class 的角色设置为目标，并将所有其他字段（预测变量）的角色保留为输入。

▶ 单击确定。

SVM 节点提供多个可选的核函数用于执行处理过程。由于无法轻易知道哪个函数对于任意给定的数据集性能最佳，我们将依次选择不同的函数并对结果进行比较。我们从默认函数开始，即 RBF（径向基函数）。

图片 25-4
模型选项卡设置



- ▶ 在“建模”选项板中，将 SVM 节点附加到类型节点。
- ▶ 打开该 SVM 节点。在模型选项卡中，单击模型名称的定制选项，然后在相邻的文本字段中键入 class-rbf。

图片 25-5
默认专家选项卡设置



- ▶ 在专家选项卡上，将模式设为专家以获得可靠性，但保持所有默认选项不变。注意，Kernel 类型默认设为 RBF。在简单模式下所有选项均为灰显。

图片 25-6
“分析”选项卡设置

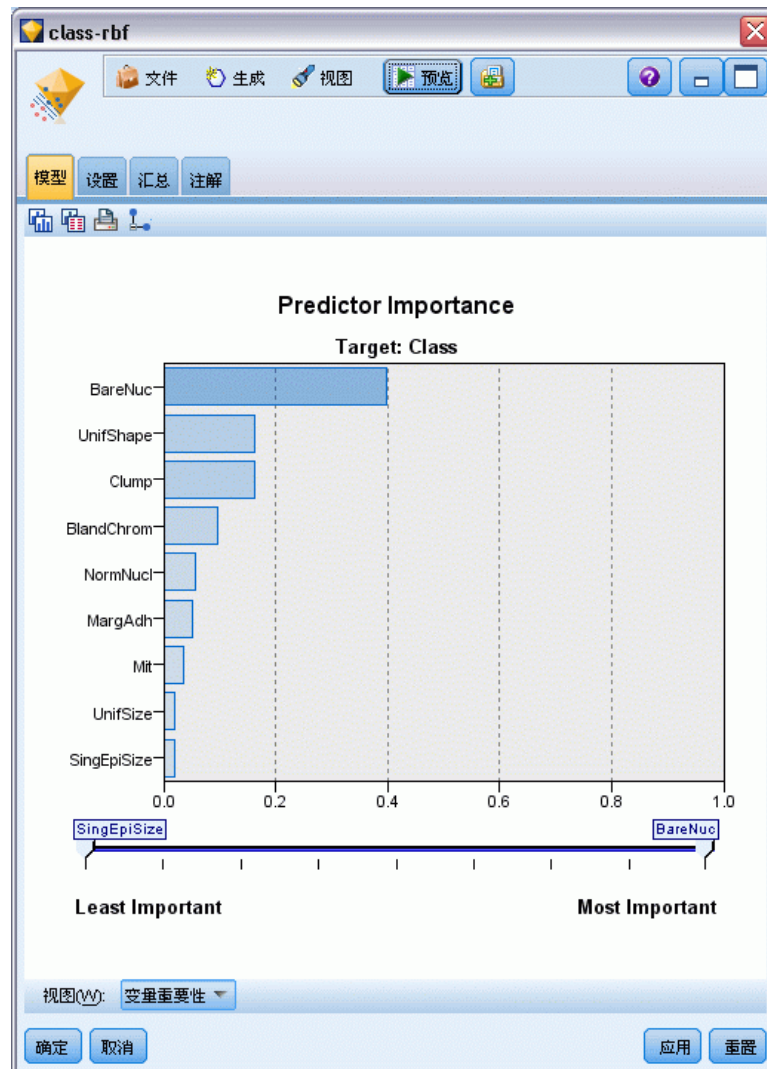


- ▶ 在分析选项卡上，选中计算变量重要性复选框。
- ▶ 单击运行。模型块放在流中，同时还放在屏幕右上角上的“模型”选项板中。

- ▶ 双击流中的模型块。

检查数据

图片 25-7
预测变量重要性图形



在“模型”选项卡上，预测变量重要性图显示了不同字段对预测的相应影响。此图向我们显示了 BareNuc 无疑具有最大的影响，而 UnifShape 和 Clump 的影响也很大。

- ▶ 单击确定。
- ▶ 将表节点附加到 class-rbf 模型块。
- ▶ 打开“表”节点，然后单击运行。

图片 25-8
为预测值和置信度值添加的字段

	gEpiSize	BareNuc	BlandChrom	NormNucl	Mit	Class	\$S-Class	\$SP-Class
1		1	3	1	1	2	2	0.992
2		10	3	2	1	2	4	0.899
3		2	3	1	1	2	2	0.994
4		4	3	7	1	2	4	0.915
5		1	3	1	1	2	2	0.992
6		10	9	7	1	4	4	0.999
7		10	3	1	1	2	2	0.907
8		1	3	1	1	2	2	0.997
9		1	1	1	5	2	2	0.997
10		1	2	1	1	2	2	0.996
11		1	3	1	1	2	2	0.999
12		1	2	1	1	2	2	0.999
13		3	4	4	1	4	2	0.514
14		3	3	1	1	2	2	0.989
15		9	5	5	4	4	4	0.991
16		1	4	3	1	4	4	0.691
17		1	2	1	1	2	2	0.997
18		1	3	1	1	2	2	0.995
19		10	4	1	2	4	4	0.996
20		1	3	1	1	2	2	0.986

- 此模型创建了两个附加字段。向右滚动表输出可看到这两个字段：

新字段名	描述
\$S-Class	由模型预测的 Class 值。
\$SP-Class	此预测值的倾向得分（即此预测值为真的似然，其值介于 0.0 到 1.0 之间）。

只需查看上表，我们就可以看到大多数记录的倾向得分（在 \$SP-Class 列）都相当高。

但是，也存在一些明显的例外情况；例如，位于第 13 行的患者 1041801 的记录，其倾向得分为不可接受的低分 0.514。同时，通过比较 Class 和 \$S-Class，可以清楚地看到此模型作出了许多不正确的预测，即使是倾向得分相对高的地方也是如此（例如，第 2 行和第 4 行）。

让我们看一下是否可以通过选择另一种函数类型获得较好的效果。

尝试另一种函数

图片 25-9
为模型设置一个新的名称



- ▶ 关闭“表格”输出窗口。
- ▶ 将第二个 SVM 建模节点附加到类型节点。
- ▶ 打开新的 SVM 节点。
- ▶ 在模型选项卡上，选择自定义和类型 class-poly 作为模型名称。

图片 25-10
“专家”选项卡的“多项式”设置



- ▶ 在专家选项卡上，将模式设置为专家。
- ▶ 将核类型设置为多项式并单击运行。class-poly 模型块被添加到流，同时还添加到屏幕右上角的“模型”选项板。
- ▶ 将 class-rbf 模型块连接到 class-poly 模型块（在警告对话框上选择替换）。
- ▶ 将表节点附加到 class-poly 模型块。
- ▶ 打开“表”节点，然后单击运行。

比较结果

图片 25-11
为多项式函数添加的字段

	ormNucl	Mit	Class	\$S-Class	\$SP-Class	\$S1-Class	\$SP1-Class
78		1	2	2	0.992	2	0.998
79		1	2	2	0.968	2	0.967
80		1	2	2	0.998	2	0.994
81		1	2	2	0.986	2	0.991
82		1	2	2	0.996	2	0.997
83		1	2	2	0.991	2	0.998
84		1	2	2	0.970	2	0.998
85	0	7	4	4	0.992	4	1.000
86	0	10	4	4	0.974	4	1.000
87		1	4	4	0.786	4	0.958
88		3	4	4	0.988	4	0.935
89		1	2	2	0.995	2	0.997
90		1	2	2	0.998	2	0.991
91		1	2	2	0.999	2	0.993
92		1	2	2	0.998	2	0.996
93		1	2	2	0.995	2	0.997
94		1	2	2	0.999	2	0.994
95		1	2	2	0.998	2	0.995
96		1	2	2	0.999	2	0.993
97		1	2	2	0.999	2	0.995

- 向右滚动表输出可看到新添加的字段。

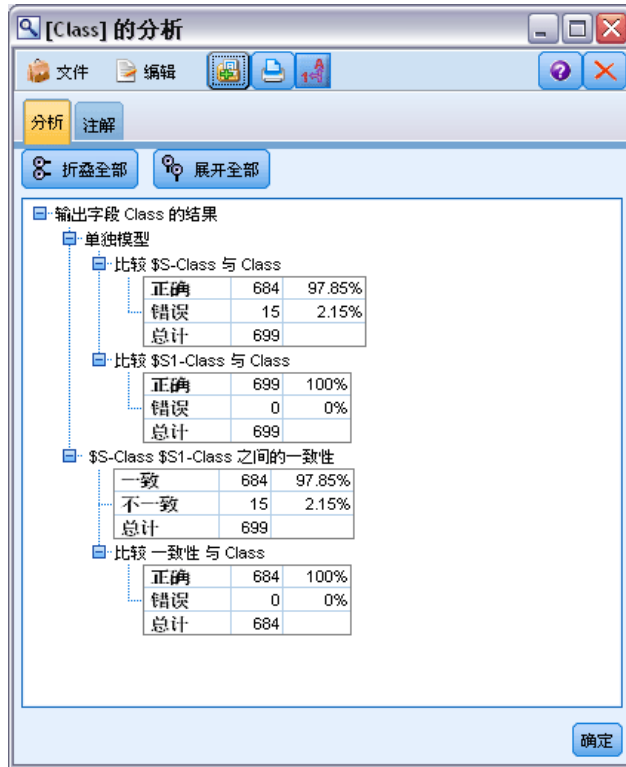
为多项式函数类型生成的字段分别名为 \$S1-Class 和 \$SP1-Class。

多项式的结果看起来好多了。许多倾向得分均为 0.995 或更高，这些结果非常令人鼓舞。

- 要确认此模型的性能更为优异，请将分析节点附加到 class-poly 模型块。

打开分析节点并单击运行。

图片 25-12
分析节点



在此方法中使用分析节点，您可以比较同类型的两个或更多模型块。来自分析节点的输出显示 RBF 函数可以正确地预测 97.85% 的观测值，这仍是一个不错的结果。但是，输出显示多项式函数可以正确预测每个观测值中的诊断。实际使用中，未必能做到完全准确，但分析节点可帮您确定模型的精确度能否满足特殊使用要求。

事实上，对于这个特定的数据集，其他任意两种函数类型（Sigmoid 和线性）的效果均不如多项式函数。但用于其他数据集时，其结果可能会明显不同，因此始终应该尝试所有选项。

摘要

您使用了多种不同类型的 SVM 核函数针对若干参数预测分类情况。此外您还了解了不同类型的核对于相同的数据集给出的结果存在多大差异，以及如何相对其他模型来度量某个模型是否性能更佳。

将 Cox 回归用于客户流失时间模型

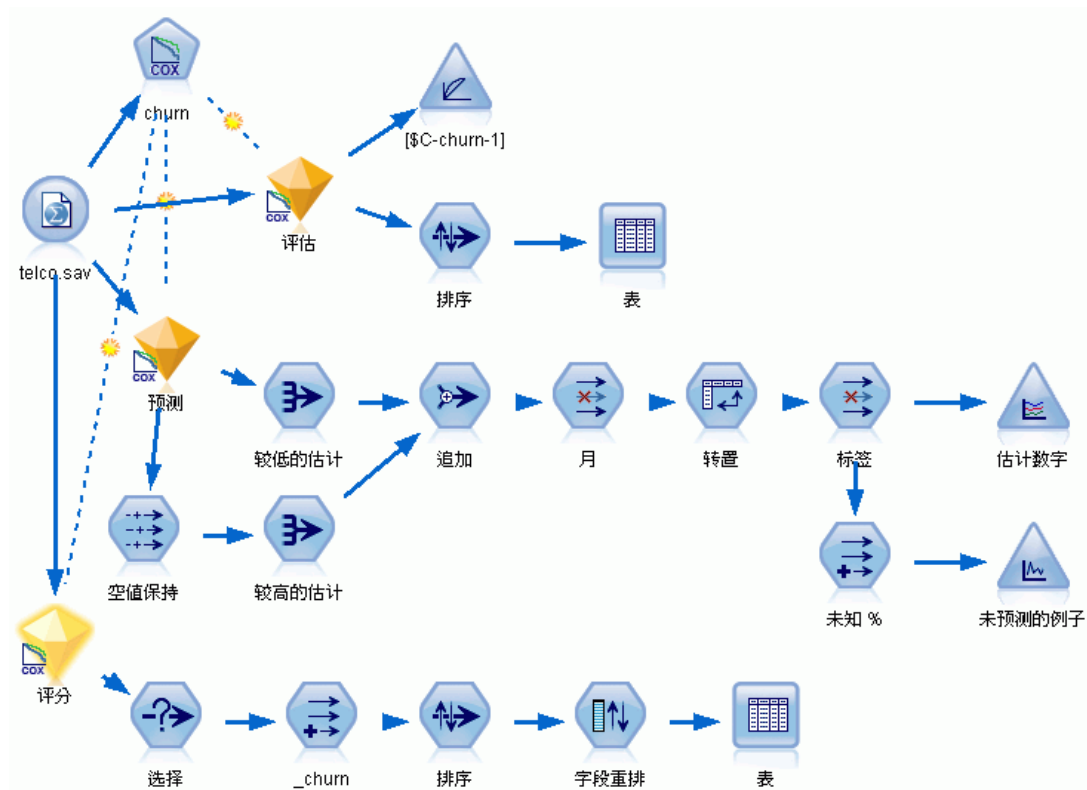
作为其减少客户流失所做工作的一部分，电信公司对“流失时间”很感兴趣，借此他们可以确定哪些因素导致客户在很短的时间内更换使用其他电信服务。为此，随机选取了一些客户样本，和他们作为客户所花费的时间（无论他们是否仍为活动客户）以及从数据库中抽取的其他各种字段。

此示例使用流 telco_coxreg.str，此流参考的是数据文件 telco.sav。数据文件和流文件分别位于 Demos 文件夹和 streams 子文件夹中。[有关详细信息，请参阅第 1 章中的 Demos 文件夹中的 IBM SPSS Modeler 15 用户指南。](#)

构建合适的模型

- ▶ 在 Demos 文件夹中添加指向 telco.sav 的 Statistics 文件源节点。

图片 26-1
分析流失时间的样本流



- ▶ 在源节点的“过滤”选项卡上，排除区域、收入字段、从 longten 到 wireten 的字段，以及从 loglong 到 logwire 的字段。

图片 26-2
过滤不需要的字段



（或者，可以在“类型”选项卡上将这些字段的角色更改为无，而不用排除它，或者在建模节点中选择要使用的字段。）

- ▶ 在源节点的“类型”选项卡上，将流失字段的角色设置为目标，将其测量级别设置为标志。将所有其他字段的角色设置为 Input。

- ▶ 单击读取值以实例化数据。

图片 26-3
设置字段角色



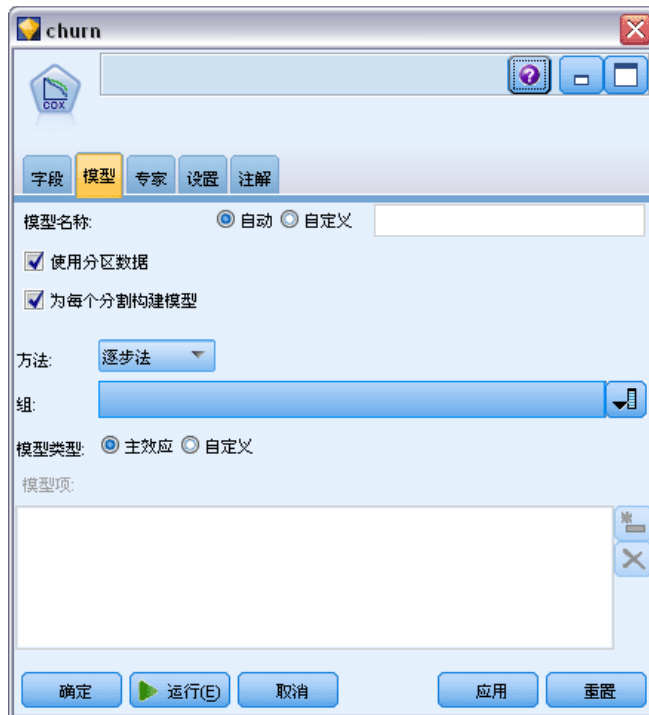
- ▶ 将 Cox 节点添加到源节点中；在字段选项卡中，选择工龄作为生存时间变量。

图片 26-4
选择字段选项



- ▶ 单击模型选项卡。
- ▶ 选择逐步法作为变量选择方法。

图片 26-5
选择模型选项



- ▶ 单击专家选项卡并选择专家以激活专家建模选项。

- ▶ 单击输出。

图片 26-6
选择高级输出选项



- ▶ 选择生存和风险作为要生成的散点图，然后单击确定。
- ▶ 单击运行以创建模型块，该模型块将被添加到流和位于右上角的“模型”选项板中。要查看其详细信息，双击流上的模型块。首先，请查看“高级输出”选项卡。

删失的观测值

图片 26-7
个案处理摘要

		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

状态变量确定是否已发生给定观测值的事件。如果事件尚未发生，则观测值被称为已删失。已删失的观测值不能用于计算回归系数，但可用于计算基线风险。观测值处理概要显示 726 个观测值已删失。该数字表示尚未流失的客户量。

分类变量编码

图片 26-8
分类变量编码

		Frequency	(1) ^b	(2)	(3)	(4)
marital ^a	0=Unmarried	505	1			
	1=Married	495	0			
ed ^a	1=Did not complete high school	204	1	0	0	0
	2=High school degree	287	0	1	0	0
	3=Some college	209	0	0	1	0
	4=College degree	234	0	0	0	1
	5=Post-undergraduate degree	66	0	0	0	0
retire ^a	.00=No	953	1			
	1.00=Yes	47	0			
gender ^a	0=Male	483	1			
	1=Female	517	0			
tollfree ^a	0=No	526	1			
	1=Yes	474	0			
equip ^a	0=No	614	1			
	1=Yes	386	0			
callcard ^a	0=No	322	1			
	1=Yes	678	0			
wireless ^a	0=No	704	1			
	1=Yes	296	0			
multiline ^a	0=No	525	1			
	1=Yes	475	0			
voice ^a	0=No	696	1			
	1=Yes	304	0			
pager ^a	0=No	739	1			
	1=Yes	261	0			
internet ^a	0=No	632	1			
	1=Yes	368	0			
callid ^a	0=No	519	1			
	1=Yes	481	0			
callwait ^a	0=No	515	1			
	1=Yes	485	0			
forward ^a	0=No	507	1			
	1=Yes	493	0			
confer ^a	0=No	498	1			
	1=Yes	502	0			
ebill ^a	0=No	629	1			
	1=Yes	371	0			
custcat ^a	1=Basic service	266	1	0	0	
	2=E-service	217	0	1	0	
	3=Plus service	281	0	0	1	
	4=Total service	236	0	0	0	

分类变量编码是解释分类协变量（特别是二元型变量）回归系数的有用参考。默认情况下，参考类别是“最后一个”类别。例如，即使已婚客户在数据文件中的变量值为 1，但为了回归的目的，这些变量值都被编码为 0。

变量选择

图片 26-9
公用检验

Step	-2 Log Likelihood	Overall (score)			Change From Previous Step			Change From Previous Block		
		Chi-square	df	Sig.	Chi-square	df	Sig.	Chi-square	df	Sig.
1 ^a	3392.536	162.303	1	.000	133.828	1	.000	133.828	1	.000
2 ^b	3087.314	249.392	2	.000	305.222	1	.000	439.050	2	.000
3 ^c	3027.085	328.426	3	.000	60.229	1	.000	499.279	3	.000
4 ^d	2990.790	347.197	4	.000	36.294	1	.000	535.574	4	.000
5 ^e	2973.790	362.673	5	.000	17.000	1	.000	552.574	5	.000
6 ^f	2958.796	376.140	6	.000	14.994	1	.000	567.568	6	.000
7 ^g	2945.503	384.717	7	.000	13.293	1	.000	580.861	7	.000
8 ^h	2936.993	417.341	8	.000	8.510	1	.004	589.371	8	.000
9 ⁱ	2926.000	423.911	9	.000	10.994	1	.001	600.364	9	.000
10 ^j	2917.551	428.078	10	.000	8.449	1	.004	608.813	10	.000
11 ^k	2913.308	436.837	11	.000	4.243	1	.039	613.056	11	.000
12 ^l	2908.078	440.158	12	.000	5.230	1	.022	618.286	12	.000

a. Variable(s) Entered at Step Number 1: callcard
 b. Variable(s) Entered at Step Number 2: longmon
 c. Variable(s) Entered at Step Number 3: equip
 d. Variable(s) Entered at Step Number 4: employ
 e. Variable(s) Entered at Step Number 5: multline
 f. Variable(s) Entered at Step Number 6: voice
 g. Variable(s) Entered at Step Number 7: address
 h. Variable(s) Entered at Step Number 8: equipmon
 i. Variable(s) Entered at Step Number 9: ebill
 j. Variable(s) Entered at Step Number 10: callid
 k. Variable(s) Entered at Step Number 11: internet
 l. Variable(s) Entered at Step Number 12: reside
 m. Beginning Block Number 0, initial Log Likelihood function: -2 Log likelihood: 3526.364
 n. Beginning Block Number 1. Method = Forward Stepwise (Likelihood Ratio)

模型构建过程采用前进逐步算法。公用检验是指对模型执行情况测量的。上一步骤的卡方更改是上一步骤和当前步骤模型的 -2 对数似然之间的差值。如果在某一步中要添加变量，则在更改的显著性小于 0.05 时才能进行此包含操作。如果某一步中要移除变量，则在更改的显著性大于 0.10 时才能进行此排除操作。在 12 个步骤中，12 个变量将添加到模型中。

图片 26-10
公式中的变量（仅适用于步骤 12）

		B	SE	Wald	df	Sig.	Exp(B)
Step 12	address	-.035	.009	14.543	1	.000	.966
	employ	-.051	.010	25.767	1	.000	.950
	reside	-.103	.046	5.037	1	.025	.902
	equip	-1.948	.381	26.180	1	.000	.143
	callcard	.777	.151	26.451	1	.000	2.175
	longmon	-.233	.022	115.619	1	.000	.792
	equipmon	-.042	.011	15.377	1	.000	.959
	multline	.612	.145	17.854	1	.000	1.844
	voice	-.501	.157	10.197	1	.001	.606
	internet	-.362	.160	5.114	1	.024	.697
	callid	-.464	.148	9.790	1	.002	.629
	ebill	-.399	.156	6.557	1	.010	.671

最终的模型包含地址、雇用状况、居住地址、equip、电话卡、longmon、equipmon、multline、声音、因特网、callid，以及电子帐单。要了解单个预测变量的效果，请查看 Exp(B)，可将 Exp(B) 解释为预测变量中单元增量风险中的预测更改。

- 地址的 Exp(B) 的值表示, 对于住在同一地址一年的用户, 流失风险会减少 $100\% - (100\% \times 0.966) = 3.4\%$ 。在同一地址居住五年的客户的流失风险百分比减少 $100\% - (100\% \times 0.966^5) = 15.88\%$ 。
- 电话卡的 Exp(B) 值表示没有订购电话卡服务的客户流失的风险比率是订购此服务的客户的 2.175 倍。重新调用分类变量编码, 其回归的 No = 1。
- 因特网的 Exp(B) 值表示未订购因特网服务的客户流失的风险比率是订购此服务的客户的 0.697 倍。这一点让人有些苦恼, 因为这表明使用该种服务的客户比不使用的客户取消公司服务的速度更快。

图片 26-11
模型中没有的变量 (仅适用于 12 个步骤)

	Score	df	Sig.
Step 12	age	.122	1 .726
	marital	.648	1 .421
	income	1.476	1 .224
	ed	6.328	4 .176
	ed(1)	.007	1 .934
	ed(2)	.203	1 .652
	ed(3)	.835	1 .361
	ed(4)	5.773	1 .016
	retire	.013	1 .908
	gender	.214	1 .644
	tollfree	3.243	1 .072
	wireless	.668	1 .414
	tollmon	.000	1 .987
	cardmon	3.163	1 .075
	wiremon	1.084	1 .298
	pager	1.808	1 .179
	callwait	.266	1 .606
	forward	2.201	1 .138
	confer	2.568	1 .109
	custcat	.864	3 .834
custcat(1)	.466	1 .495	
custcat(2)	.450	1 .502	
custcat(3)	.019	1 .889	

模型左侧变量的得分统计量的显著性值均大于 0.05。但是, tollfree 和 cardmon 的显著性值不小于 0.05, 且与该值很接近。二者在今后的研究中有待进一步考证。

协变量均值

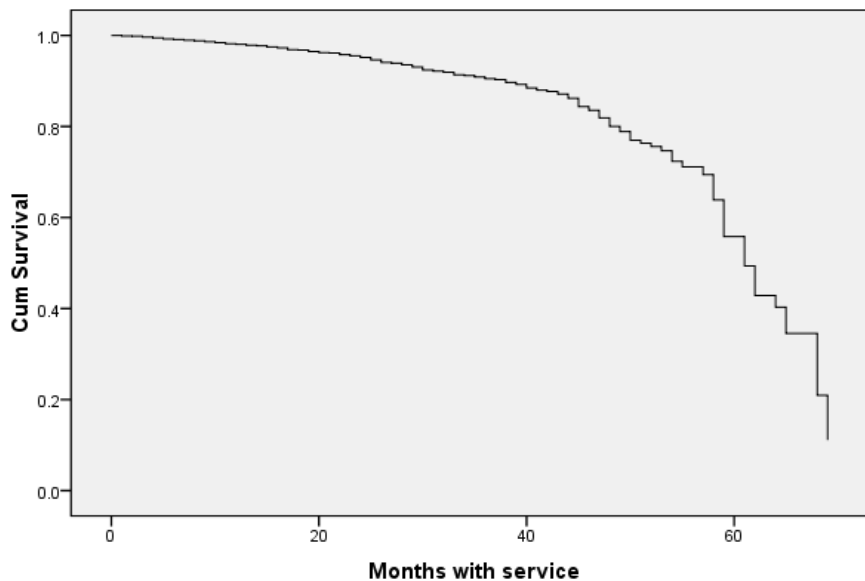
图片 26-12
协变量均值

	Mean
age	41.684
marital	.505
address	11.551
income	77.535
ed(1)	.204
ed(2)	.287
ed(3)	.209
ed(4)	.234
employ	10.987
retire	.953
gender	.483
reside	2.331
tollfree	.526
equip	.614
callcard	.322
wireless	.704
longmon	11.723
tollmon	13.274
equipmon	14.220
cardmon	13.781
wiremon	11.584
multiline	.525
voice	.696
pager	.739
internet	.632
callid	.519
callwait	.515
forward	.507
confer	.498
ebill	.629
custcat(1)	.266
custcat(2)	.217
custcat(3)	.281

此表格显示了每个预测变量的平均值。在查看为均值构建的生存散点图时，此表格是很有用的参考。但请注意，在您查看分类预测变量的指示符变量均值时，“平均”客户实际上并不存在。即使使用所有尺度预测变量，您也无法找到一位其所有协变量值都接近均值的客户。如果要查看特殊观测值的生存曲线，可以更改协变量值，这些协变量用于在“散点图”对话框中绘制生存曲线。如果要查看特殊观测值的生存曲线，可以更改协变量值，这些值用于在“高级输出”对话框的散点图组中绘制生存曲线。

生存曲线

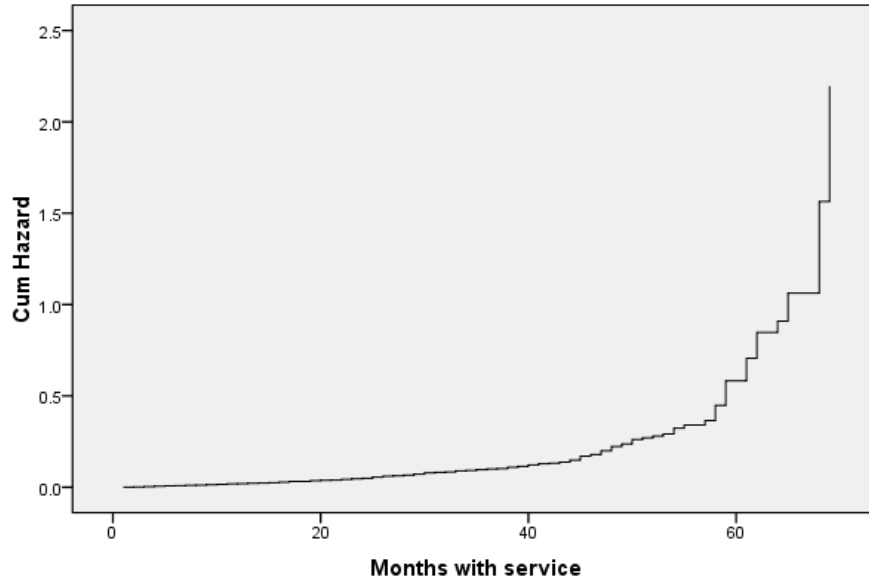
图片 26-13
“平均”客户的生存曲线



基本生存曲线是“平均”客户的模型预测流失时间的可视化显示。水平轴显示事件发生的时间。垂直轴显示生存概率。所以，生存曲线上的任何一点表示“平均”客户经过某段时间仍未流失的概率。55 个月过后，生存曲线变得不平滑。很少有客户那么长的时间使用公司服务，这样可获取的信息变少，导致曲线变成块状。

风险曲线

图片 26-14
“平均”客户的风险曲线



基本风险曲线是“平均”客户的累积模型预测流失可能性的可视化显示。水平轴显示事件发生的时间。垂直轴显示累积风险，等于生存概率的负对数。55 个月过后，同生存曲线一样，风险曲线也变得不平滑，变化原因相同。

评估

逐步选择法保证模型仅包含“统计意义上显著”的预测变量，但不保证模型实际上非常适用于预测目标。为此，需要分析已评分的记录。

图片 26-15
Cox 块：“设置”选项卡



- ▶ 将模型块放置在工作区上并将其附加到源节点，然后打开该模型块并单击“设置”选项卡。
- ▶ 选择时间字段并指定工龄。将根据每个记录的工龄长度对其进行评分。
- ▶ 选择追加所有概率。

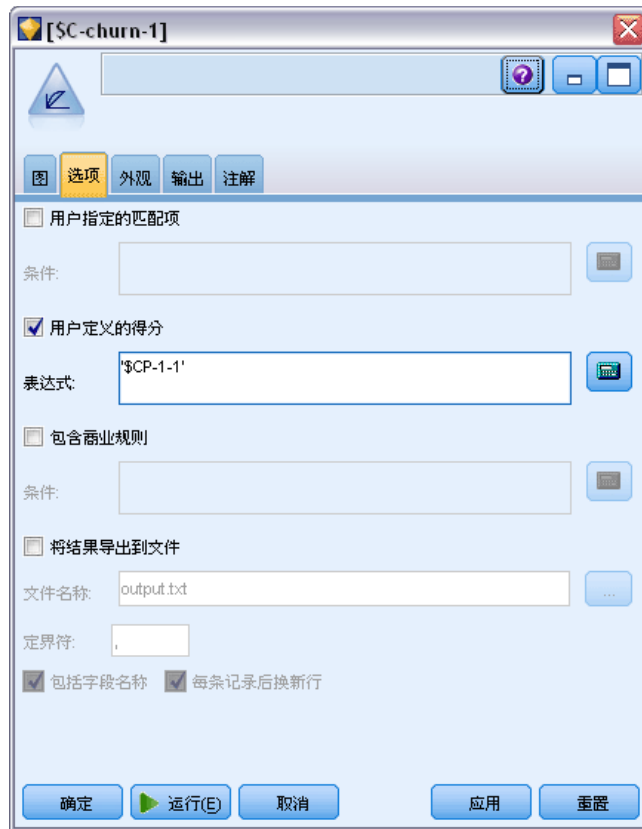
这样可使用 0.5 作为截断值区分客户是否流失来创建得分；如果客户流失的倾向值大于 0.5，则这些客户将被标记为“流失者”。该数值并不特殊，但不同的截断值可能产生不同的期望结果。在考虑选择截断值时，有一种方法是使用评估节点。

图片 26-16
评估节点：“散点图”选项卡



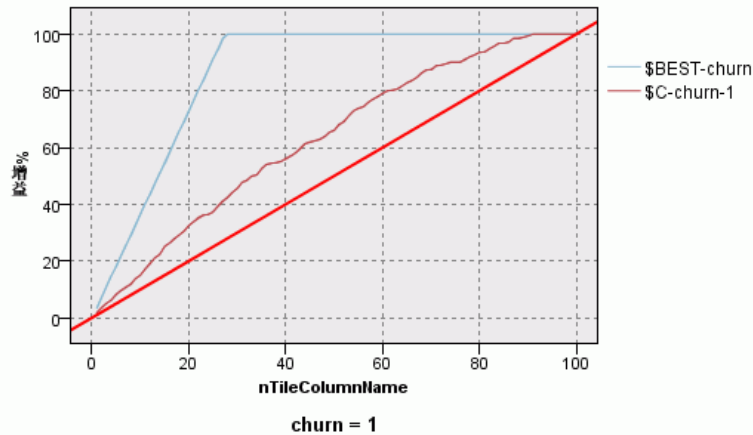
- ▶ 将评估节点附加到模型块中；在“散点图”选项卡上，选择包含最佳线。
- ▶ 单击选项选项卡。

图片 26-17
评估节点：“选项”选项卡



- ▶ 选择用户定义的得分，然后键入 '\$CP-1-1' 作为表达式。这是与流失倾向相对应的模型生成的字段。
- ▶ 单击运行。

图片 26-18
收益图



累积增益图会在给定的类别中显示通过把个案总数的百分比作为目标而“增益”的个案总数的百分比。例如，曲线上的一点（10%，15%），表示当您使用模型对数据集进行评分，并使用预测的流失倾向对所有观测值进行排序时，您可能希望前 10% 中包含实际类别为 1（流失者）的所有观测值中约 15% 的观测值。同样，前 60% 包含大约 79.2% 的流失者。如果选择所有已评分的数据集，则将获得数据集中所有流失者。

对角线是“基线”曲线；如果从已评分数据集中随机选取 20% 的记录，则可能从实际类别为 1 的所有记录中“获得”约 20% 的记录。曲线离基线的上方越远，增益越大。“最佳”线显示将较高的倾向得分分配到每个流失者（而非每个非流失者）的“完美”模型的曲线。您可以使用累积增益图帮助通过选择对应于大量收益的百分比选择分类标准值，然后将百分比与适当分界值映射。

生成“大量”收益取决于类型 I 与类型 II 错误的成本。也就是说，将流失者划分为非流失者（第一类）的成本是多少？将非流失者分类为流失者（第二类）的成本是多少？当客户保持成为首要考虑的问题时，您会希望降低第一类错误；在累积收益表上，这对应于预测倾向值为 1 的前 60% 中的客户已增加的客户维护，这将捕获 79.2% 的可能流失者，但将花费时间和资源用于获得新客户。如果降低维护当前客户群成本是首先要考虑的事，则您会希望降低第二类错误。在图表上，这对应于前 20% 增加的客户维护，这将捕获 32.5% 的流失者。通常，这两方面都是很重要的事，这样您必须选择一个决策规则用于对具有最大敏感度和特征性的客户进行分类。

图片 26-19
排序节点：“设置”选项卡



- ▶ 假定您所定的期望收益为 45.6%，则相应会显示前面 30% 的记录。要找到合适的分类截断值，请将排序节点附加到模型块。
- ▶ 在“设置”选项卡上，选择以降序按 \$CP-1-1 排序，然后单击确定。

图片 26-20
Table

rn	\$C-churn-1	\$CP-churn-1	\$CP-0-1	\$CP-1-1
292	0	0.744	0.744	0.256
293	0	0.745	0.745	0.255
294	0	0.745	0.745	0.255
295	0	0.746	0.746	0.254
296	0	0.748	0.748	0.252
297	0	0.749	0.749	0.251
298	0	0.749	0.749	0.251
299	0	0.750	0.750	0.250
300	0	0.752	0.752	0.248
301	0	0.752	0.752	0.248
302	0	0.754	0.754	0.246
303	0	0.754	0.754	0.246
304	0	0.755	0.755	0.245
305	0	0.756	0.756	0.244
306	0	0.757	0.757	0.243
307	0	0.757	0.757	0.243
308	0	0.758	0.758	0.242
309	0	0.759	0.759	0.241
310	0	0.761	0.761	0.239
311	0	0.762	0.762	0.238

- ▶ 将表节点附加到排序节点。
- ▶ 打开“表”节点，然后单击运行。

将输出向下滚动，就会看到第 300 个记录的 $\$CP-1-1$ 值为 0.248。使用 0.248 作为分类截断值，就会导致约 30% 的客户被评定为流失者，由此捕获实际总流失者的个数约占 45%。

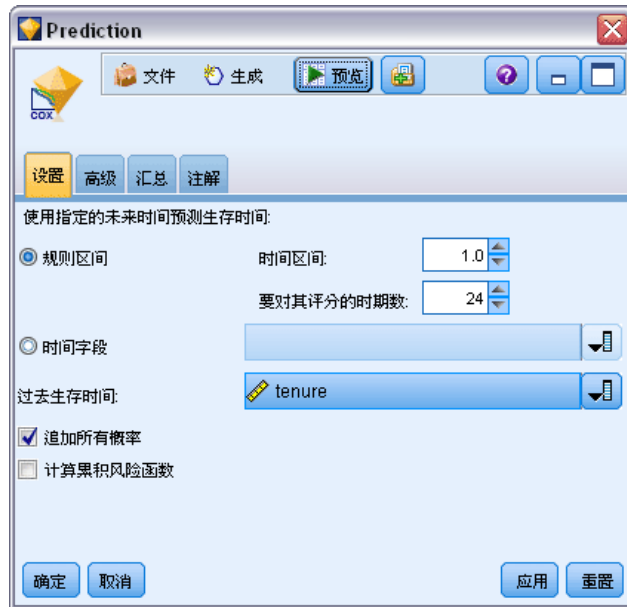
跟踪仍在的预期客户数

对模型满意后，您会希望跟踪数据集中未来两年内会保留的预期客户数。空值代表客户的总工龄（未来时间 + 工龄）不在用于训练模型的数据中的生存时间范围内，呈现出一种很有趣的挑战。处理空值的一种方法是创建两个预测集合，其中一个集合中

的空值被假定为已流失，另一个集合中的空值被假定为已保留。使用这种方法，可以创建保留的预测客户数的上限和下限。

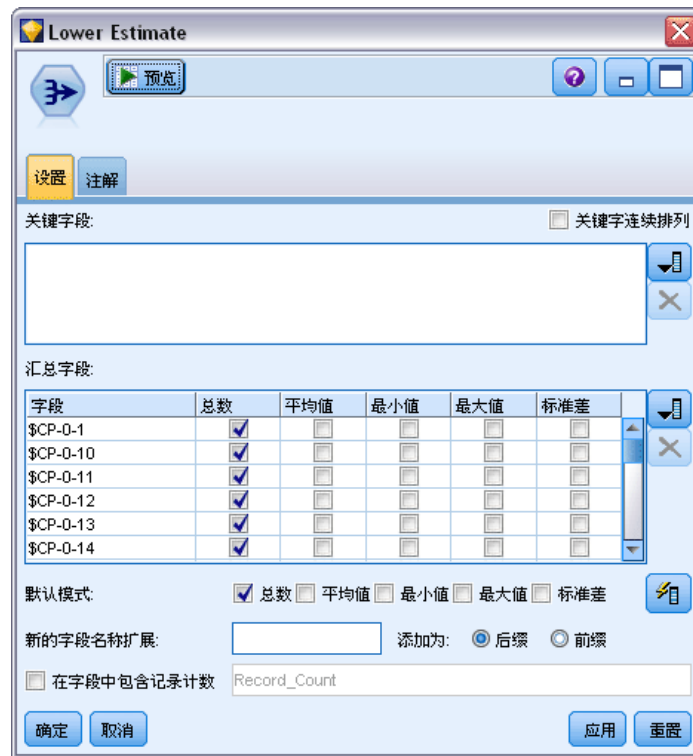
图片 26-21

Cox 块：“设置”选项卡



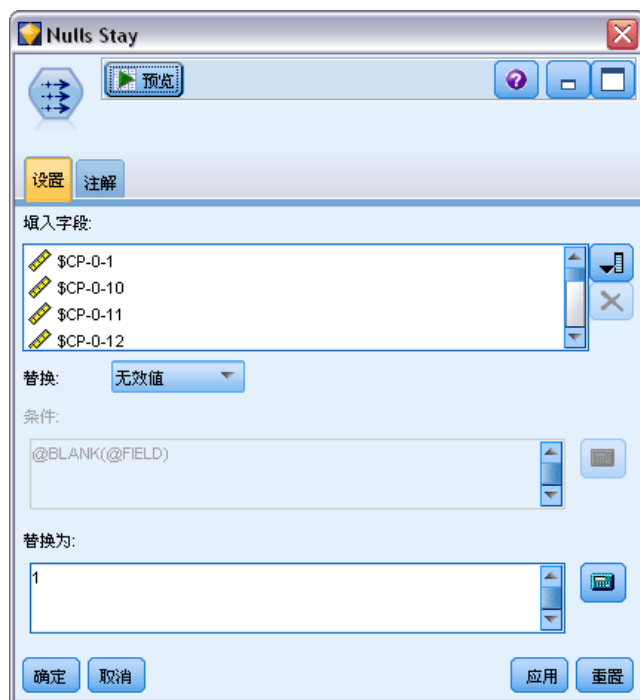
- ▶ 双击“模型”选项板中的模型块（或复制和粘贴流工作区上的模型块），并将新的模型块附加到源节点。
- ▶ 打开模型块的“设置”选项卡。
- ▶ 确保选中规则区间，并指定 1.0 为时间区间，24 为要对其评分的时段数。此项操作指定在未来 24 个月中每个月都会对每个记录进行评分。
- ▶ 选择 tenure 作为指定过去生存时间的字段。得分算法将考虑每个客户持续作为公司客户的时间长度。
- ▶ 选择追加所有概率。

图片 26-22
 汇总节点：“设置”选项卡



- ▶ 将汇总节点添加到模型块；在“设置”选项卡上，取消选中均值作为默认模式。
- ▶ 选择从 \$CP-0-1 到 \$CP-0-24 的字段（即格式为 \$CP-0-n 的字段）作为要汇总的字段。如果在“选择字段”对话框上，按名字（即字母顺序）对字段进行排序，则这种方法最容易。
- ▶ 取消选中在字段中包含记录计数。
- ▶ 单击确定。此节点可创建“下限”预测。

图片 26-23
填充节点：“设置”选项卡



- ▶ 将填充节点附加到 Coxreg 块（我们刚刚附加汇总节点的位置）；在“设置”选项卡上，选择从 \$CP-0-1 到 \$CP-0-24 的字段（即格式为 \$CP-0-n 的字段）作为要填充的字段。如果在“选择字段”对话框上，按名字（即字母顺序）对字段进行排序，则这种方法最容易。
- ▶ 选择使用值 1 替换空值。
- ▶ 单击确定。

图片 26-24
 汇总节点：“设置”选项卡



- ▶ 将汇总节点附加到填充节点；在“设置”选项卡上，取消选中均值作为默认模式。
- ▶ 选择从 \$CP-0-1 到 \$CP-0-24 的字段（即格式为 \$CP-0-n 的字段）作为要汇总的字段。如果在“选择字段”对话框上，按名字（即字母顺序）对字段进行排序，则这种方法最容易。
- ▶ 取消选中在字段中包含记录计数。
- ▶ 单击确定。此节点可创建“上限”预测。

图片 26-25
过滤节点：“设置”选项卡



- ▶ 将追加节点附加到两个汇总节点；然后将过滤节点附加到追加节点。
- ▶ 在过滤节点的“设置”选项卡上，将字段重新命名为 1 到 24。使用转置节点，这些字段名称将变为图表下游中 x 轴上的值。

图片 26-26
转置节点：“设置”选项卡



- ▶ 将转置节点附加到过滤节点。
- ▶ 第 2 类作为新字段数。

图片 26-27
过滤节点：“过滤”选项卡



- ▶ 将过滤节点附加到转置节点。
- ▶ 在过滤节点的“设置”选项卡上，将 ID 重新命名为月，将字段1重新命名为较低估计，将字段2命名为较高估计。

图片 26-28
多重散点图节点：“散点图”选项卡



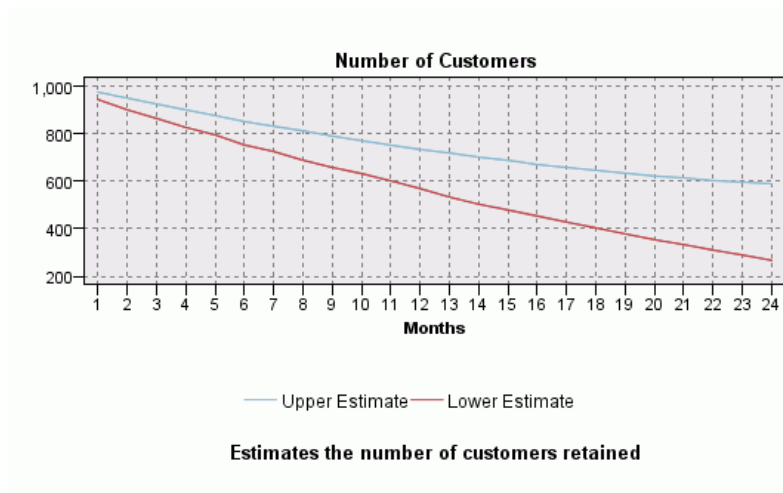
- ▶ 将多重散点图节点附加到过滤节点。
- ▶ 在“散点图”选项卡上，月是 X 字段，较低估计和较高估计是 Y 字段。

图片 26-29
多重散点图节点：“外观”选项卡



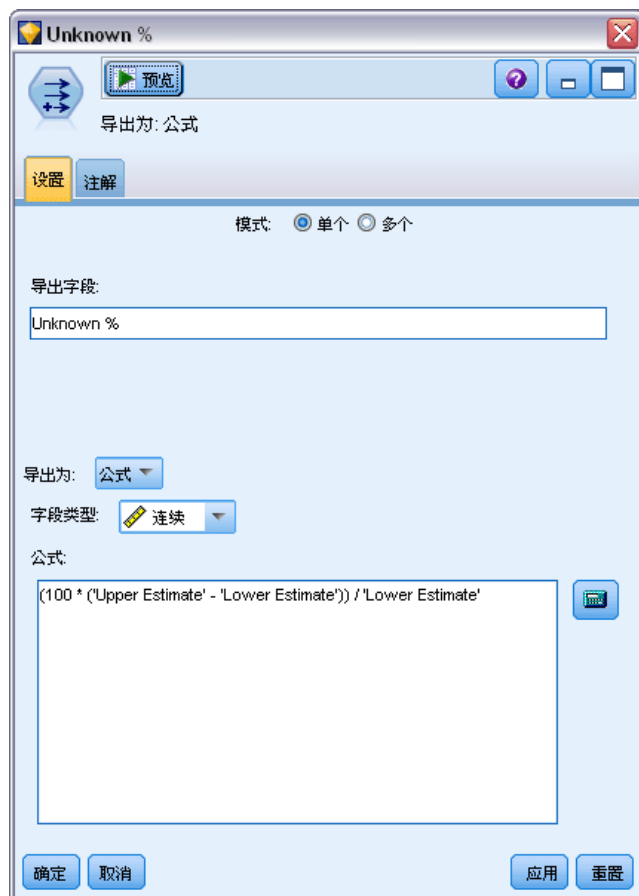
- ▶ 单击“外观”选项卡。
- ▶ 键入客户数作为标题。
- ▶ 键入估计保留的客户数作为标注。
- ▶ 单击运行。

图片 26-30
多重散点图可预测保留的客户数。



已绘制出估计的保留客户数的上限和下限。这两条线的差值是得分为空值的客户数，因此其状态很难确定。这些客户数量随时间增加。12 个月过后，可以预计数据集中保留的原始客户数介于 601 到 735 之间；24 个月过后，该值介于 288 到 597 之间。

图片 26-31
导出节点：“设置”选项卡



- ▶ 要查看保留的客户数估计的不确定程度，请将导出节点附加到过滤节点。
- ▶ 在导出节点的“设置”选项卡上，键入未知百分比作为导出字段。
- ▶ 选择连续作为字段类型。
- ▶ 键入 $(100 * ('Upper Estimate' - 'Lower Estimate')) / 'Lower Estimate'$ 作为公式。未知百分比是作为较低估计百分比的“质疑”客户数。
- ▶ 单击确定。

图片 26-32
散点图节点：“散点图”选项卡



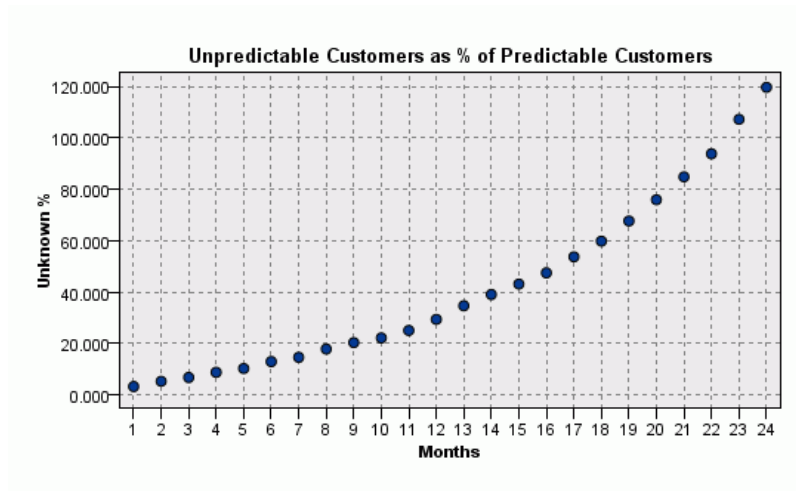
- ▶ 将散点图节点附加到导出节点。
- ▶ 在散点图节点的“散点图”选项卡上，选择月作为 X 字段，未知百分比作为 Y 字段。
- ▶ 单击“外观”选项卡。

图片 26-33
散点图节点：“外观”选项卡



- ▶ 键入 Unpredictable Customers as % of Predictable Customers 作为标题。
- ▶ 执行节点。

图片 26-34
无法预测客户的散点图

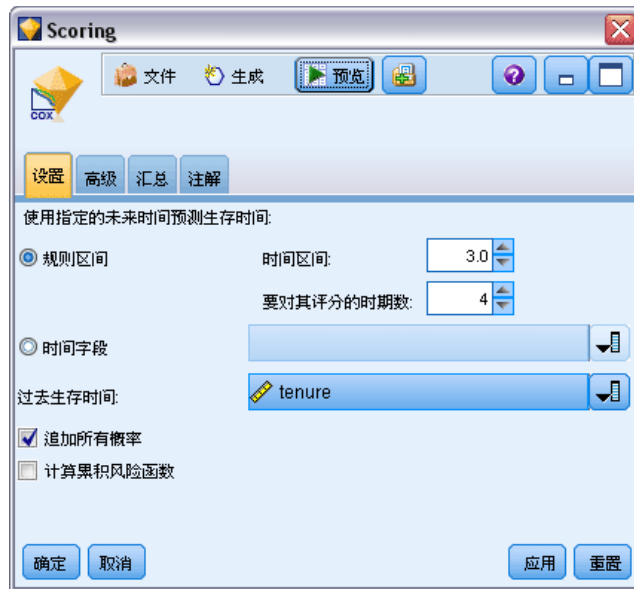


第一年里，无法预测客户的百分比以明显的线性速率增长。但第二年增长速率猛增，一直到第 23 个月，具有空值的客户数超过了保留的预测客户数。

评分

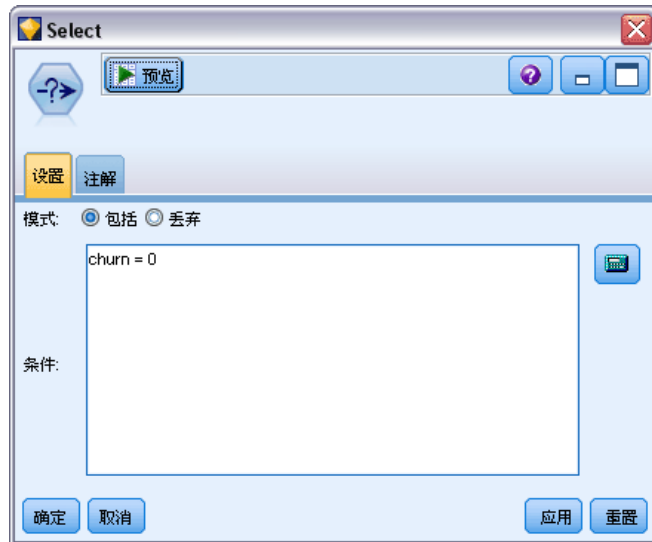
如果对模型满意，则您会希望对客户进行评分以确认下一年一个季度内最可能流失的客户。

图片 26-35
Coxreg 块：“设置”选项卡



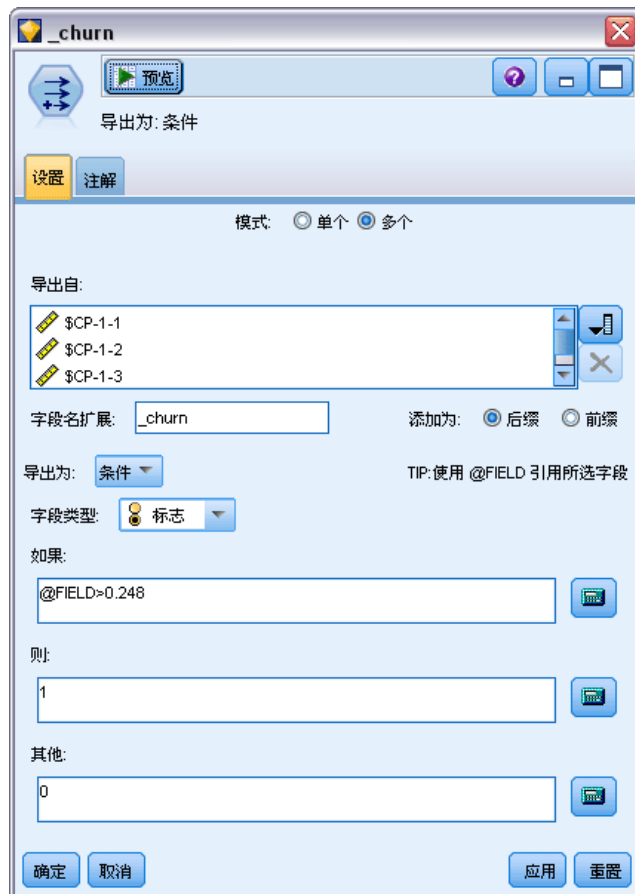
- ▶ 将第三个模型块附加到源节点并打开模型块。
- ▶ 确保选中规则区间，并指定 3.0 为时间区间，4 为要对其评分的时段数。此项操作指定在未来四个季度中将对每个记录进行评分。
- ▶ 选择 tenure 作为指定过去生存时间的字段。得分算法将考虑每个客户持续作为公司客户的时间长度。
- ▶ 选择追加所有概率。使用这些附加字段更容易对表中要查看的记录进行排序。

图片 26-36
选择节点：“设置”选项卡



- ▶ 将选择节点附加到模型块；在“设置”选项卡上，键入 `churn=0` 作为条件。此项操作将移除已从结果表中流失的客户。

图片 26-37
导出节点：“设置”选项卡



- ▶ 将导出节点附加到选择节点；在“设置”选项卡上，选择多个作为模式。
- ▶ 选择导出从 \$CP-1-1 到 \$CP-1-4 的字段（即格式为 \$CP-1-n 的字段），然后键入 _churn 作为要添加的后缀。如果在“选择字段”对话框上，按名字（即字母顺序）对字段进行排序，则这种方法最容易。
- ▶ 选择将字段导出为 条件。
- ▶ 选择标志为测量级别。
- ▶ 键入 @FIELD>0.248 作为 If 条件。请记住，这是评估期间确定的分类截断值。
- ▶ 键入 1 作为 Then 表达式。
- ▶ 键入 0 作为 Else 表达式。
- ▶ 单击确定。

图片 26-38
排序节点：“设置”选项卡



- ▶ 将排序节点附加到导出节点；在“设置”选项卡上，选择先按 \$CP-1-1_churn 到 \$CP-1-4_churn、再按 \$CP-1-1 到 \$CP-1-4 进行排序，所有顺序都按降序排列。预测要流失的客户会出现在顶端。

图片 26-39
字段重排节点：“重排”选项卡



- ▶ 将字段重排节点附加到排序节点；在“重排”选项卡上，选择将 \$CP-1-1_churn 到 \$CP-1-4 放在其他字段之前。此项操作会使结果表更易于读取，且是可选的。您将需要使用按钮将字段移动到图中显示的位置。

图片 26-40
显示客户得分的表

	\$CP-1-1_churn	\$CP-1-1	\$CP-1-2_churn	\$CP-1-2	\$CP-1-3_churn	\$CP-1-3	\$CP-1-4_churn	\$CP-1-4
255	0	0.032	0	0.075	0	0.147	1	0.298
256	0	0.027	0	0.064	0	0.127	1	0.260
257	0	0.023	0	0.130	0	0.233	1	0.308
258	0	0.021	0	0.127	0	0.239	1	0.320
259	0	0.021	0	0.125	0	0.237	1	0.318
260	0	0.021	0	0.053	0	0.198	1	0.331
261	0	0.021	0	0.053	0	0.196	1	0.329
262	0	0.020	0	0.050	0	0.189	1	0.317
263	0	0.017	0	0.043	0	0.163	1	0.278
264	0	0.015	0	0.039	0	0.148	1	0.253
265	0	0.197	0	0.197	0	\$null\$	0	\$null\$
266	0	0.109	0	0.109	0	\$null\$	0	\$null\$
267	0	0.101	0	0.214	0	\$null\$	0	\$null\$
268	0	0.081	0	0.137	0	0.194	0	0.245
269	0	0.074	0	0.159	0	\$null\$	0	\$null\$
270	0	0.070	0	0.116	0	0.158	0	0.237
271	0	0.070	0	0.128	0	0.189	0	0.234
272	0	0.062	0	0.105	0	0.151	0	0.191
273	0	0.062	0	0.130	0	0.163	0	0.212
274	0	0.061	0	0.123	0	0.182	0	0.241

- 将表节点附加到字段重排节点并运行。

预计年末将流失 264 位客户，第三个季度末流失 184 位，第二个季度末流失 103 位，第一个季度末流失 31 位。请注意，如果指定两类客户，一类在第一季度中具有较高的流失倾向，但在以后的季度中并非就具有较高的流失倾向；例如，请查看记录 256 和 260。这可能是由于客户当前工龄之后的几个月中的风险函数类型所致。例如，因升职而加入的客户与因个人推荐而加入的客户相比，前者可能会更早地离开。但如果不是这样，则实际上这些客户在其剩余的工龄内会成为更加忠实的用户。您可能希望对这些客户重新排序以获得最可能流失的客户的不同视图。

图片 26-41
显示具有空值的客户的表

The screenshot shows a window titled "Table (50 个字段, 726 条记录) #1". The table contains 726 rows and 9 columns. The columns are labeled as follows: an unlabeled ID column, \$CP-1-1_churn, \$CP-1-1, \$CP-1-2_churn, \$CP-1-2, \$CP-1-3_churn, \$CP-1-3, \$CP-1-4_churn, and \$CP-1-4. The data shows a pattern where the 'churn' columns contain '0' and the corresponding variable columns contain '\$null\$'. This pattern repeats for every row from 707 to 726.

	\$CP-1-1_churn	\$CP-1-1	\$CP-1-2_churn	\$CP-1-2	\$CP-1-3_churn	\$CP-1-3	\$CP-1-4_churn	\$CP-1-4
707	0	\$null\$	0	\$null\$	0	\$null\$	0	\$null\$
708	0	\$null\$	0	\$null\$	0	\$null\$	0	\$null\$
709	0	\$null\$	0	\$null\$	0	\$null\$	0	\$null\$
710	0	\$null\$	0	\$null\$	0	\$null\$	0	\$null\$
711	0	\$null\$	0	\$null\$	0	\$null\$	0	\$null\$
712	0	\$null\$	0	\$null\$	0	\$null\$	0	\$null\$
713	0	\$null\$	0	\$null\$	0	\$null\$	0	\$null\$
714	0	\$null\$	0	\$null\$	0	\$null\$	0	\$null\$
715	0	\$null\$	0	\$null\$	0	\$null\$	0	\$null\$
716	0	\$null\$	0	\$null\$	0	\$null\$	0	\$null\$
717	0	\$null\$	0	\$null\$	0	\$null\$	0	\$null\$
718	0	\$null\$	0	\$null\$	0	\$null\$	0	\$null\$
719	0	\$null\$	0	\$null\$	0	\$null\$	0	\$null\$
720	0	\$null\$	0	\$null\$	0	\$null\$	0	\$null\$
721	0	\$null\$	0	\$null\$	0	\$null\$	0	\$null\$
722	0	\$null\$	0	\$null\$	0	\$null\$	0	\$null\$
723	0	\$null\$	0	\$null\$	0	\$null\$	0	\$null\$
724	0	\$null\$	0	\$null\$	0	\$null\$	0	\$null\$
725	0	\$null\$	0	\$null\$	0	\$null\$	0	\$null\$
726	0	\$null\$	0	\$null\$	0	\$null\$	0	\$null\$

该表的底端是具有预测空值的客户。这些客户的总工龄（未来时间 + 工龄）在用于训练模型的数据中的生存时间范围内。

摘要

使用 Cox 回归，您已找到流失时间的适用模型，该模型绘制下两年中保留的预测客户数并且识别下一年中最可能流失的客户。请注意，这只是适用模型，不一定是最佳模型。理想的作法是，您至少应当将使用前进逐步法获取的该模型与使用后退逐步法创建的模型进行比较。

《SPSS Modeler 算法指南》中列出了对 IBM® SPSS® Modeler 中用到的建模方法的数学原理的说明。

市场购物篮分析（规则归纳/C5.0）

本示例处理描述超级市场购物篮内容（即，所购买的全部商品的集合）的虚构数据，以及购买者的相关个人数据（可通过忠诚卡方案获得）。目的是寻找购买相似产品并且可按人口统计学方式（如按年龄、收入等）刻画其特征的客户群。

本示例说明了数据挖掘的两个阶段：

- 关联规则建模和一个揭示所购买商品之间联系的 Web 显示
- C5.0 规则归纳（描绘已标识产品组的购买者的特征）

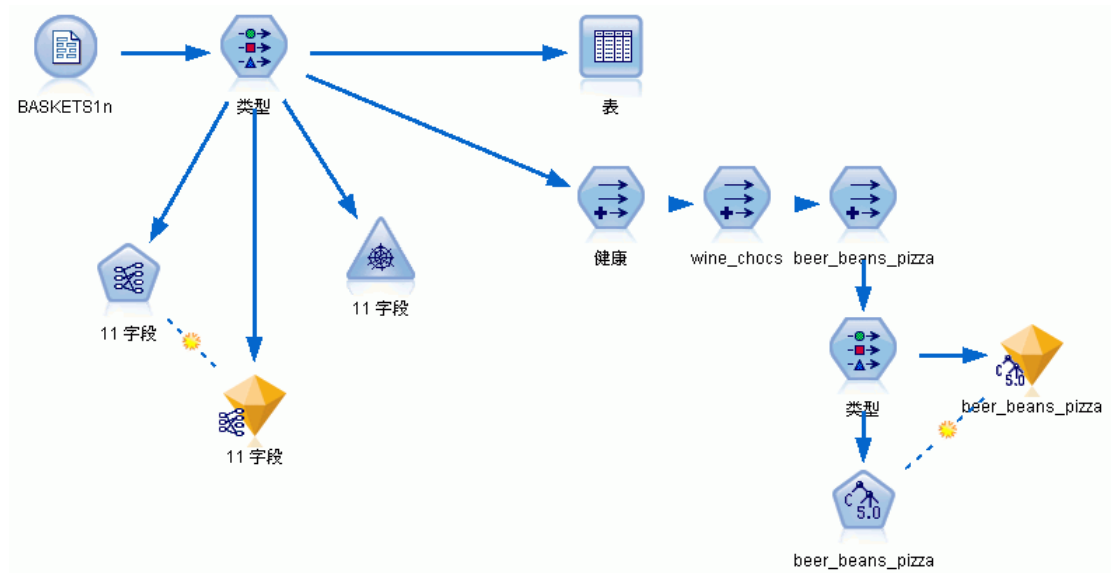
注意：此应用不直接使用预测建模，因此，不对最终模型进行准确性度量，在数据挖掘过程中也不存在与之相关的训练/检验两个步骤的区分。

本例使用名为 baskrule 的流，该流引用名为 BASKETS1n 的数据文件。这些文件可在任何 IBM® SPSS® Modeler 安装程序的 Demos 目录中找到。此目录可通过 WindowsIBM® SPSS® Modeler 程序组进行访问。文件 baskrule 位于 streams 目录下。

访问数据

使用“变量文件”节点连接到数据集 BASKETS1n，选择要从该文件读取的字段名称。将“类型”节点连接到数据源，然后将该节点连接到“表”节点。将字段卡 ID 的测量级别设置为无类型（因为每个忠诚卡 ID 在数据集中只出现一次，因此对于建模没有用处）。选择名义作为字段性别的测量级别（这是为了确保 Apriori 建模算法不会将性别视为标志）。

图片 27-1
购物规则流



现在，运行该流以将“类型”节点实例化并显示表。数据集包含 18 个字段，其中每条记录表示一个购物篮。

下列标题中会显示 18 个字段。

购物篮摘要：

- cardid. 购买此篮商品的客户的忠诚卡标识符。
- value. 购物篮的总购买价格。
- pmethod. 购物篮的支付方法。

卡持有者的个人详细信息：

- sex
- homeown. 卡持有者是否拥有住房。
- income
- age

购物篮内容—产品类别的出现标志：

- fruitveg
- freshmeat
- dairy
- cannedveg
- cannedmeat
- frozenmeal
- beer
- wine
- softdrink
- fish
- confectionery

发现购物篮内容的关系

首先，需要使用 Apriori 大致了解购物篮内容的关系（关联）以生成关联规则。选择要在此建模过程中使用的字段，方法是：编辑“类型”节点，将所有产品类别的角色设置为两者，并将所有其他角色设置为无。（双向表示该字段可以是结果模型的输入或输出。）

注意：通过按住 Shift 键并单击以选择多个字段，然后指定列中的选项，可为多个字段设置选项。

图片 27-2
选择用于建模的字段



指定了用于建模的字段后，请将 Apriori 节点附加到“类型”节点，编辑它，选择选项只显示值为真的标志变量，然后在 Apriori 节点上单击“运行”。结果（管理器窗口右上角“模型”选项卡上的模型）包含您可以查看（使用上下文菜单，然后选择浏览）的关联规则。

图片 27-3
关联规则

后项	前项	支持度 %	置信度 %
frozenmeal	beer cannedveg	16.7	87.425
cannedveg	beer frozenmeal	17.0	85.882
beer	frozenmeal cannedveg	17.3	84.393

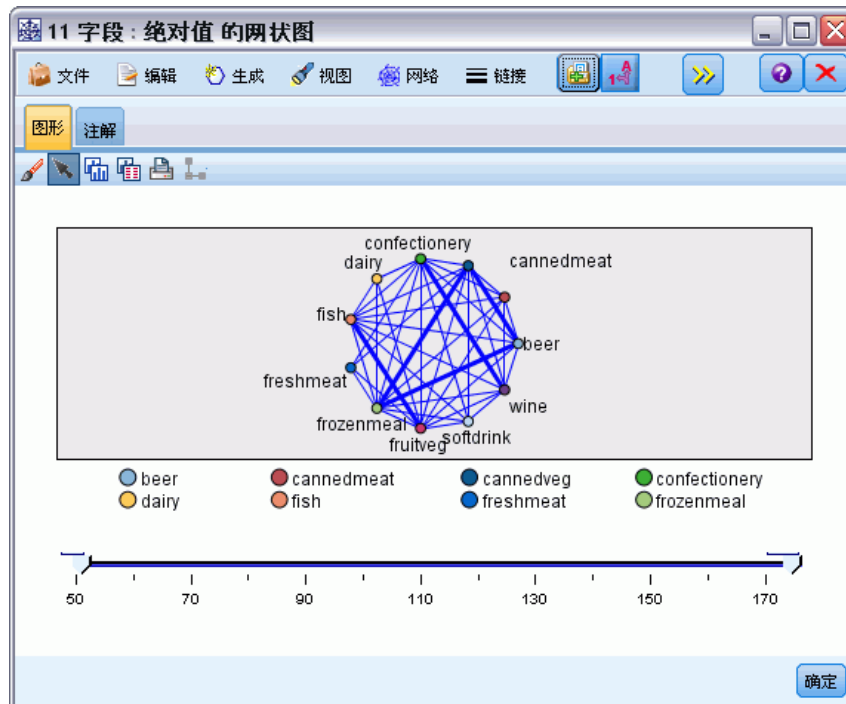
这些规则显示冻肉、罐装蔬菜和啤酒之间存在多种关联。出现双向关联规则（如：

```
frozenmeal -> beer  
beer -> frozenmeal
```

提示：Web 显示（只显示双向关联）可能会突出显示此数据中的一些模式。

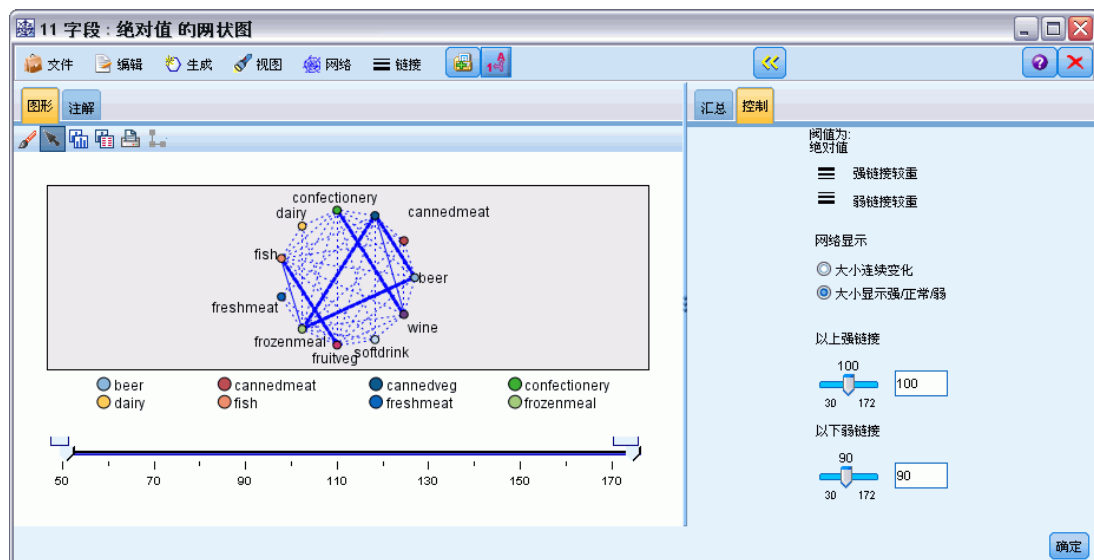
将 Web 节点附加到“类型”节点，编辑 Web 节点，选择所有购物篮内容字段，选择仅显示 true 标志，然后在 Web 节点上单击“运行”。

图片 27-4
产品关联的 Web 显示



因为大多数产品类别组合都会出现在多个购物篮中，所以此 Web 上的强链接太多，无法显示模型表示的客户群。

图片 27-5
限制性 Web 显示



- ▶ 要指定弱连接和强连接，请单击工具栏上的黄色双箭头按钮。这会展开显示 Web 输出摘要和控件的对话框。
- ▶ 选择大小表示强/正常/弱。
- ▶ 将弱链接设置为低于 90。
- ▶ 将强链接设置为高于 100。

在最终显示中，会有三个客户群突出显示：

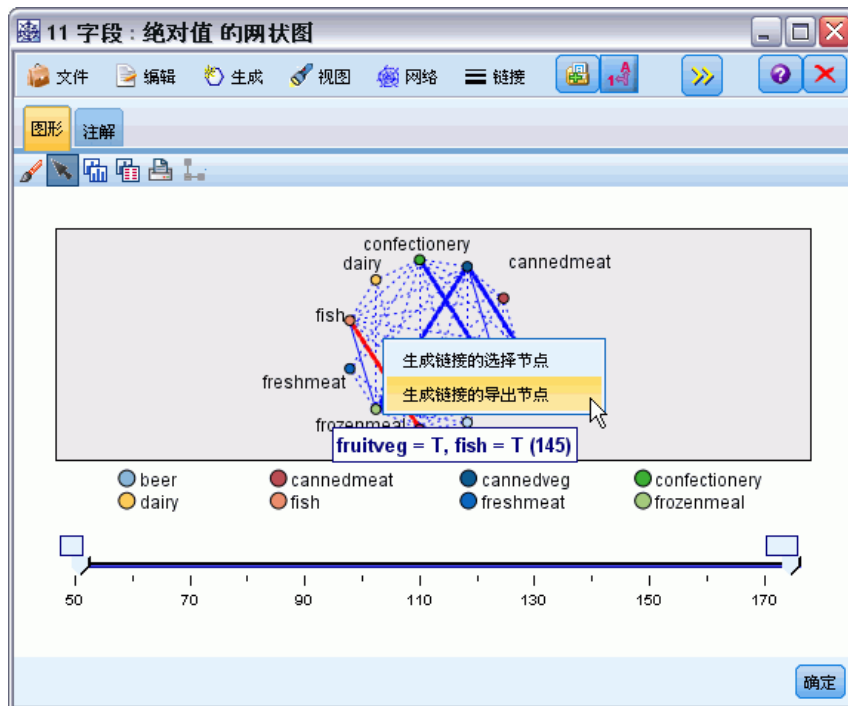
- 购买鱼和果蔬的客户，可将这类客户称为“健康食客”
- 购买酒和粮果的客户
- 购买啤酒、冻肉和罐装蔬菜（“啤酒、豆类和比萨”）的客户

描绘客户群的特征

现在，您已经根据客户购买的产品类型标识了三个客户群，但是还要知道这些客户是谁，即，他们的人口统计学特征。通过为每个群中的每名客户添加标志，并使用规则归纳 (C5.0) 来基于规则描绘这些标志的特征，可以实现这一点。

首先，必须获取每个群的标志。使用刚刚创建的 Web 显示，可以自动生成每个群的标志。使用鼠标右键，单击 fruitveg 和 fish 之间的链接以突出显示该链接，然后右键单击并选择为链接生成“派生”节点。

图片 27-6
获取每个客户群的标志



编辑最终的“派生”节点以将“派生”字段名称更改为健康。使用从wine到confectionery的链接重复该练习，并将最终的“派生”字段命名为 wine_chocs。

对于第三个群（涉及三个链接），首先要确保未选择任何链接。然后，在按住 shift 键的同时单击鼠标左键，从而选择cannedveg、beer和frozenmeal中的全部三个链接。（您一定要处于“交互”模式而不是“编辑”模式。）然后，从 Web 显示菜单中选择：

生成 > 导出节点（“和”）

将最终“派生”字段的名称更改为 beer_beans_pizza。

要描绘这些客户群的特征，请连续将现有的类型节点连接到这三个导出节点，然后附加另一个类型节点。在新“类型”节点中，请将除以下字段外的所有字段的角色都设置为无：value、pmethod、sex、homeown、income 和 age（这些字段的角色应该设置为输入），以及相关的客户群（例如，beer_beans_pizza，它们的角色应该设置为目标）。附加 C5.0 节点，将输出类型设置为规则集，然后在节点上单击“运行”。最终模型（用于 beer_beans_pizza）包含此客户群的明确人口统计学特征：

```
Rule 1 for T:  
if sex = M  
and income <= 16,900  
then T
```

通过在第二个类型节点中选择其他客户群标志作为输出，可将同一方法应用到这些标志。通过在此上下文中使用 Apriori 代替 C5.0，可生成更多替代特征描绘；Apriori 也可用于同时描绘所有客户群标志的特征，原因是，Apriori 并非被限制到一个输出字段。

摘要

此示例说明如何使用 IBM® SPSS® Modeler 通过建模（使用 Apriori）和直观化（使用 Web 显示）发现数据库中的关系（即链接）。这些链接与数据中的案例组相对应，并且，通过建模（使用 C5.0 规则集）可详细研究这些组并描绘其特征。

例如，在零售领域，可能会使用这种客户组确定特殊优惠目标，以提高直接邮寄的响应率，或自定义某分部的存货产品范围以与其人口统计学基础的需求匹配。

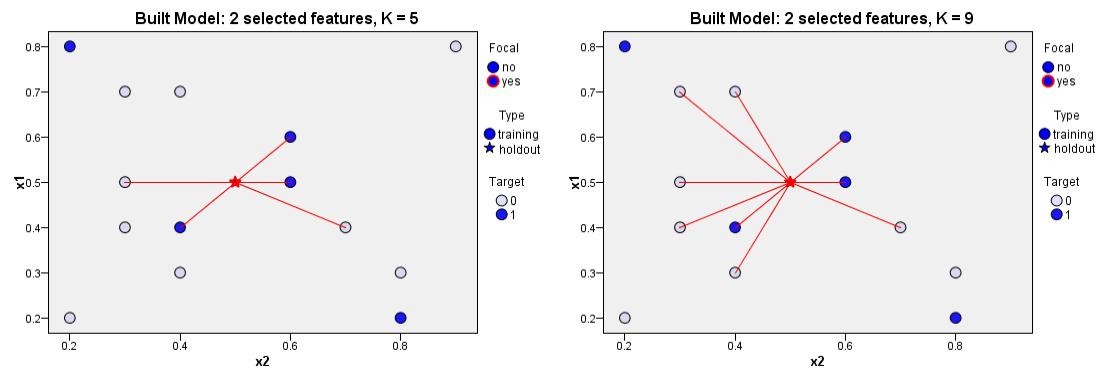
评估新车辆产品 (KNN)

“最近相邻元素分析”是根据观测值与其他观测值的类似程度分类观测值的方法。在机器学习中，将其开发为识别数据模式的一种方法，而不需要与任何存储模式或观测值完全匹配。相似个案相互邻近，非相似个案则相互远离。因此，两个观测值之间的距离是其不相似性的测量。

将靠近彼此的个案视为“相邻元素。”当提出新的观测值（保留观测值）时，计算其到模型中每个观测值的距离。计算最相似观测值 - 最近相邻元素 - 的分类并将新观测值放在包含最多最近相邻元素的类别中。

您可以规定需要检验的最近相邻元素的数量；此值叫做k。图片显示如何使用两个不同的 k 值分类新观测值。当 k = 5 时，新观测值将被置于类别 1 中，因为大多数最近相邻元素属于类别 1。但当 k = 9 时，新观测值将被置于类别 0 中，因为大多数最近相邻元素属于类别 0。

图片 28-1
更改 k 对分类的影响



最近相邻元素分析也可用于计算连续目标的值。在此情况下，最近相邻元素的平均值或中间目标值用于获得新观测值的预测值。

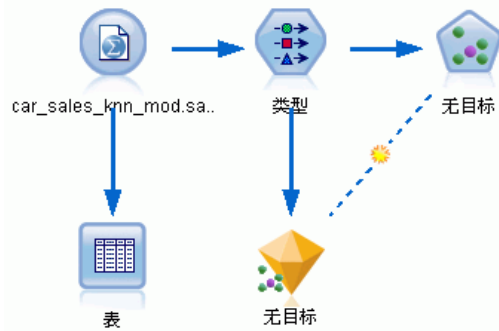
一家汽车制造商开发了两种新车辆（分别为轿车和货车）的原型。在将新车型引入其产品系列之前，该制造商想确定市场上的现有车辆与产品原型有多相像，即哪些车辆是它们的最近相邻元素，并以此确定它们将与哪些车型展开竞争。

制造商收集了有关现有车型的不同类别的数据，并添加了其原型产品的详细信息。需要在不同车型间进行比较的类别包括以千为单位的价 (price)、发动机尺寸 (engine_s)、马力 (horsepow)、轴距 (wheelbas)、车宽 (width)、车长 (length)、整车重量 (curb_wgt)、油箱容量 (fuel_cap) 和燃油效率 (mpg)。

本示例使用了名为 car_sales_knn.str 的流，该流位于 Demos 文件夹下的 streams 子文件夹中。数据文件为 car_sales_knn_mod.sav。有关详细信息，请参阅第 1 章中的 Demos 文件夹中的 IBM SPSS Modeler 15 用户指南。

创建流

图片 28-2
KNN 建模样本流



创建一个新流，然后在 IBM® SPSS® Modeler 安装程序的 Demos 文件夹中添加一个指向 car_sales_knn_mod.sav 的 Statistics 文件源节点。

现在，我们来看制造商收集的数据。

- ▶ 将“表”节点附加到 Statistics 文件源节点。
- ▶ 打开“表”节点，然后单击运行。

图片 28-3
轿车与货车的源数据

	manufact	model	sales	resale	type	price	engine_s	horsepow	wheelbas	width
140	Toyota	Celica	33.269	15.445	0.0...	16...	1.800	140.000	102.400	68.3...
141	Toyota	Tacoma	84.087	9.575	1.0...	11...	2.400	142.000	103.300	66.5...
142	Toyota	Sienna	65.119	\$null\$	1.0...	22...	3.000	194.000	114.200	73.4...
143	Toyota	RAV4	25.106	13.325	1.0...	16...	2.000	127.000	94.900	66.7...
144	Toyota	4Run...	68.411	19.425	1.0...	22...	2.700	150.000	105.300	66.5...
145	Toyota	Land ...	9.835	34.080	1.0...	51...	4.700	230.000	112.200	76.4...
146	Volksw...	Golf	9.761	11.425	0.0...	14...	2.000	115.000	98.900	68.3...
147	Volksw...	Jetta	83.721	13.240	0.0...	16...	2.000	115.000	98.900	68.3...
148	Volksw...	Passat	51.102	16.725	0.0...	21...	1.800	150.000	106.400	68.5...
149	Volksw...	Cabrio	9.569	16.575	0.0...	19...	2.000	115.000	97.400	66.7...
150	Volksw...	GTI	5.596	13.760	0.0...	17...	2.000	115.000	98.900	68.3...
151	Volksw...	Beetle	49.463	\$null\$	0.0...	15...	2.000	115.000	98.900	67.9...
152	Volvo	S40	16.957	\$null\$	0.0...	23...	1.900	160.000	100.500	67.6...
153	Volvo	V40	3.545	\$null\$	0.0...	24...	1.900	160.000	100.500	67.6...
154	Volvo	S70	15.245	\$null\$	0.0...	27...	2.400	168.000	104.900	69.3...
155	Volvo	V70	17.531	\$null\$	0.0...	28...	2.400	168.000	104.900	69.3...
156	Volvo	C70	3.493	\$null\$	0.0...	45...	2.300	236.000	104.900	71.5...
157	Volvo	S80	18.969	\$null\$	0.0...	36...	2.900	201.000	109.900	72.1...
158		newC...	\$null\$	\$null\$	\$n...	21...	1.500	76.000	106.300	67.9...
159		newT...	\$null\$	\$null\$	\$n...	34...	3.500	167.000	109.800	75.2...

两个原型（分别名为 newCar 和 newTruck）的详细信息，已被添加到文件末尾。

从源数据中我们可以看到，制造商使用了“货车”分类（在类型列中值为 1）而不是粗略地指代任何非客车车辆。

最后一列分区是必需的，以便我们在确定两个原型的最近相邻元素时可以将其指定为保留。在这种情况下，其数据不会影响计算，因为它们是我们要考虑的其余市场部分。将两个保留记录的分区值设置为 1，同时所有其他记录在此字段上为 0 值，这允许我们稍后在设置焦点记录时使用该字段，焦点记录即我们要为其计算最近相邻元素的记录。

现在保持表输出窗口处于打开，稍后会引用它。

图片 28-4
类型节点设置

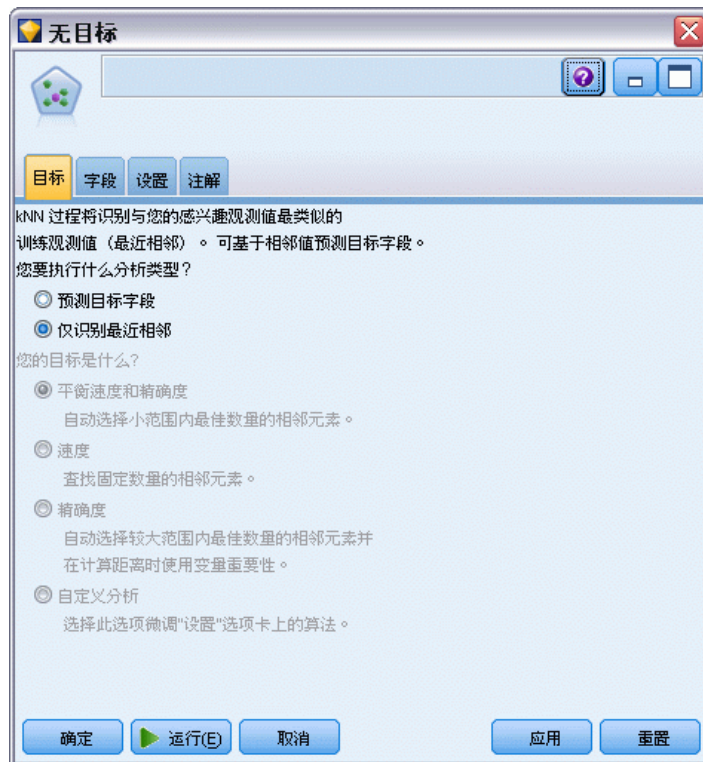


- ▶ 为流添加类型节点。
- ▶ 将“类型”节点附加到 Statistics 文件源节点。
- ▶ 打开该类型节点。

我们只想在字段 price 至 mpg 上进行比较，因此保持所有这些字段的角色设置为输入。

- ▶ 将所有其他字段（manufact 至 type，以及 lnsales）的角色设置为无。
- ▶ 将最后一个字段分区的测量级别设置为标志。确保将其角色设置为输入。
- ▶ 单击读取值以读取数据值到流中。
- ▶ 单击确定。

图片 28-5
选择识别最近相邻元素



- ▶ 将 KNN 节点附加到“类型”节点。
- ▶ 打开 KNN 节点。

由于我们只想为两个原型寻找最近相邻元素，因此这次不会预测目标字段。

- ▶ 在目标选项卡上，选择只识别最近相邻元素。
- ▶ 单击设置选项卡。

图片 28-6
使用分区字段识别焦点记录



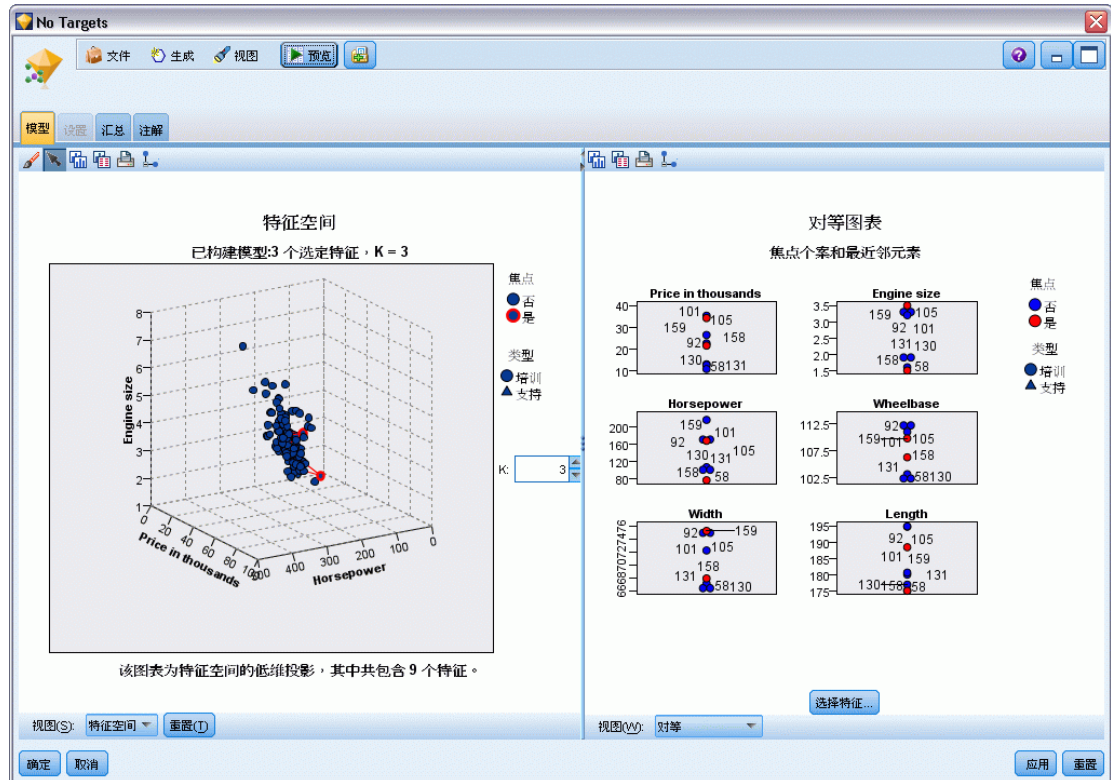
现在，可以使用分区字段来识别焦点记录，即我们要为其确定最近相邻元素的记录。通过使用标志字段，可以确保此字段值设置为 1 的记录成为焦点记录。

如您所见，该字段值为 1 的记录只有 newCar 和 newTruck，因此它们将作为我们的焦点记录。

- ▶ 在设置选项卡的模型面板上，选中识别焦点记录复选框。
- ▶ 从该字段的下拉列表中，选择分区。
- ▶ 单击运行按钮。

检查输出

图片 28-7
“模型查看器”窗口

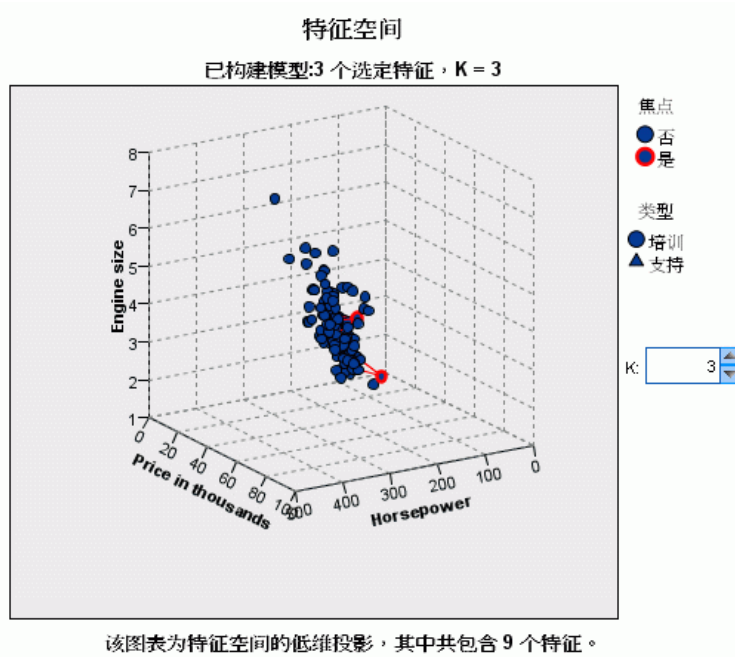


在流工作区和“模型”选项板中已创建了模型块。打开任一模型块即可显示“模型查看器”，它包含一个双面板窗口：

- 第一个面板显示模型概览，称为主视图。“最近相邻元素”模型的主视图称为**预测变量空间**。
- 第二个面板显示两种视图类型之一：
辅助模型视图显示有关模型的更多信息，但并不专注于模型本身。
当您深入查看主视图某个部分时，链接视图显示有关某个模型特征的详细信息。

预测变量空间

图片 28-8
预测变量空间图表



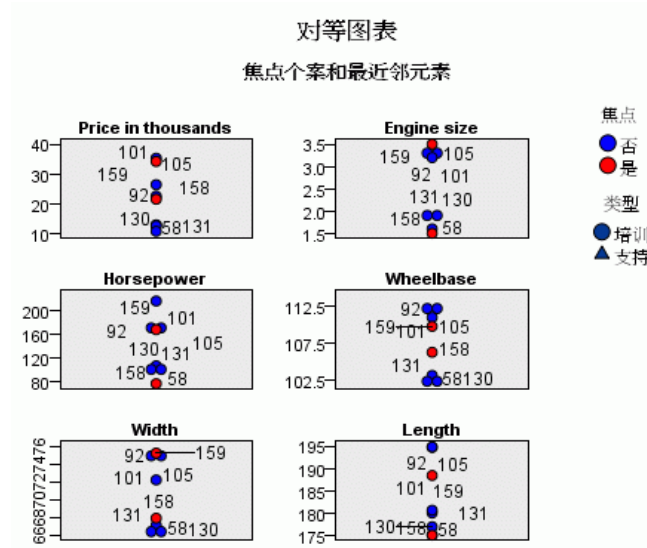
预测变量空间图表为交互式三维图形, 它绘制三个特征 (实际上是源数据的前三个输入字段) 的数据点, 分别代表价格、发动机尺寸和马力。

两个焦点记录突出显示为红色, 通过线条连接到其 k 个最近相邻元素。

通过单击并拖动图表, 可将其旋转以更方便地查看预测变量空间中的数据点分布。单击重置按钮, 将其恢复到默认视图。

对等图表

图片 28-9
对等图表



默认辅助视图为对等图表，其中突出显示在预测变量空间中选中的两个焦点记录，及其在六个特征（源数据的前六个输入字段）上的 k 个最近相邻元素。

车辆通过它们在源数据中的记录编号来表示。我们需要从“表”节点获得输出以帮助识别它们。

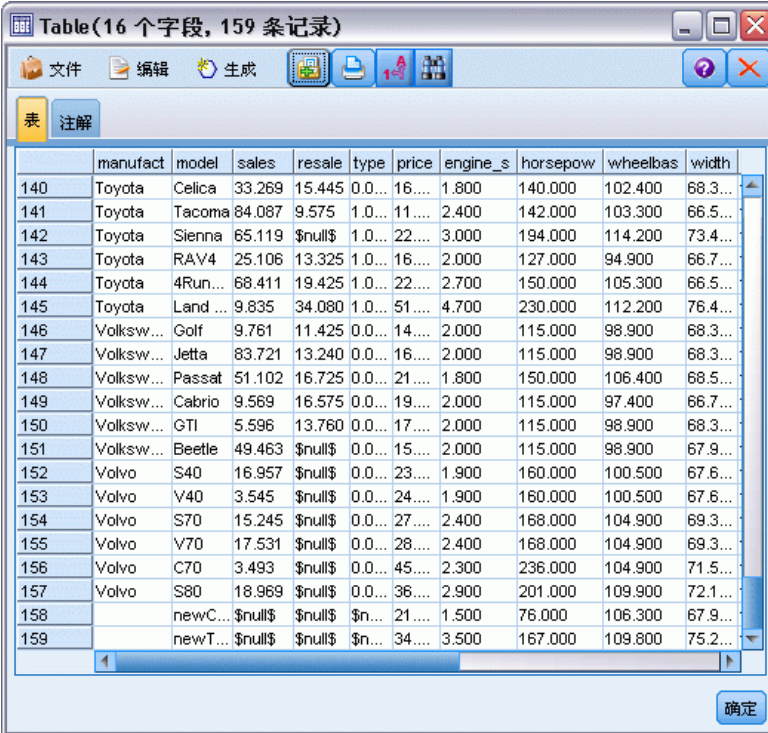
如果“表”节点输出仍然可用：

- ▶ 单击位于主 IBM® SPSS® Modeler 窗口右上角的管理器窗格的输出选项卡。
- ▶ 双击条目表（16 个字段，159 个记录）。

如果表输出不再可用：

- ▶ 在主 SPSS Modeler 窗口上，打开“表”节点。
- ▶ 单击运行。

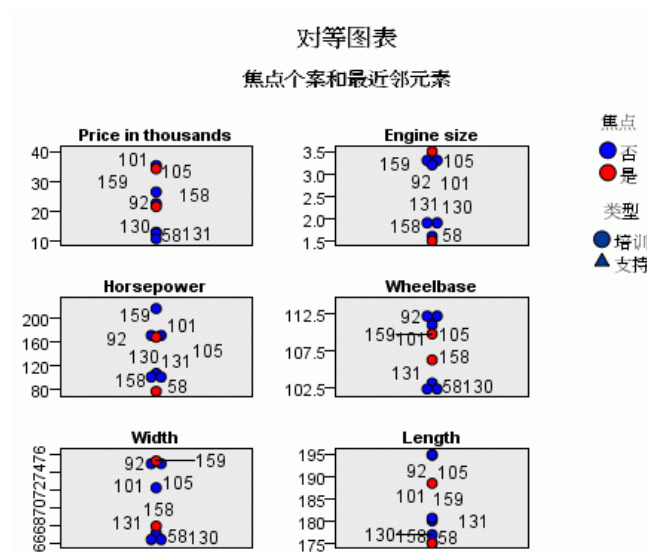
图片 28-10
按记录编号标识记录



	manufact	model	sales	resale	type	price	engine_s	horsepow	wheelbas	width
140	Toyota	Celica	33.269	15.445	0.0...	16....	1.800	140.000	102.400	68.3...
141	Toyota	Tacoma	84.087	9.575	1.0...	11....	2.400	142.000	103.300	66.5...
142	Toyota	Sienna	65.119	\$null\$	1.0...	22....	3.000	194.000	114.200	73.4...
143	Toyota	RAV4	25.106	13.325	1.0...	16....	2.000	127.000	94.900	66.7...
144	Toyota	4Run...	68.411	19.425	1.0...	22....	2.700	150.000	105.300	66.5...
145	Toyota	Land ...	9.835	34.080	1.0...	51....	4.700	230.000	112.200	76.4...
146	Volksw...	Golf	9.761	11.425	0.0...	14....	2.000	115.000	98.900	68.3...
147	Volksw...	Jetta	83.721	13.240	0.0...	16....	2.000	115.000	98.900	68.3...
148	Volksw...	Passat	51.102	16.725	0.0...	21....	1.800	150.000	106.400	68.5...
149	Volksw...	Cabrio	9.569	16.575	0.0...	19....	2.000	115.000	97.400	66.7...
150	Volksw...	GTI	5.596	13.760	0.0...	17....	2.000	115.000	98.900	68.3...
151	Volksw...	Beetle	49.463	\$null\$	0.0...	15....	2.000	115.000	98.900	67.9...
152	Volvo	S40	16.957	\$null\$	0.0...	23....	1.900	160.000	100.500	67.6...
153	Volvo	V40	3.545	\$null\$	0.0...	24....	1.900	160.000	100.500	67.6...
154	Volvo	S70	15.245	\$null\$	0.0...	27....	2.400	168.000	104.900	69.3...
155	Volvo	V70	17.531	\$null\$	0.0...	28....	2.400	168.000	104.900	69.3...
156	Volvo	C70	3.493	\$null\$	0.0...	45....	2.300	236.000	104.900	71.5...
157	Volvo	S80	18.969	\$null\$	0.0...	36....	2.900	201.000	109.900	72.1...
158		newC...	\$null\$	\$null\$	\$n...	21....	1.500	76.000	106.300	67.9...
159		newT...	\$null\$	\$null\$	\$n...	34....	3.500	167.000	109.800	75.2...

向下滚动到表的底部，可以看到 newCar 和 newTruck 为数据的最后两条记录，编号分别为 158 和 159。

图片 28-11
在对等图表上比较特征



例如，在这里，可以从对等图表上看到 newTruck (159) 的发动机尺寸大于它的任何最近相邻元素，而 newCar (158) 的发动机尺寸则小于它的任何最近相邻元素。

对于六个特征中的每个特征，可以将鼠标移至单独点上，查看特定个案的每个特征的实际值。

但哪些车辆是 newCar 和 newTruck 的最近相邻元素呢？

对等图表稍微有点拥挤，现在我们转换到较简单的视图。

- ▶ 单击位于对等图表底部的视图下拉列表（当前名为对等的条目）。
- ▶ 选择邻元素和距离表。

邻元素和距离表

图片 28-12
邻元素和距离表

k 最近邻元素和距离
针对初始焦点个案显示

焦点个案	最近邻元素			最近距离		
	1	2	3	1	2	3
158	131	130	58	0.979	0.990	1.011
159	105	92	101	0.580	0.634	0.644

很好。现在可以看到在市场上与两种原型最接近的三种车型。

对于 newCar（焦点记录 158），它们是 Saturn SC（131）、Saturn SL（130）和 Honda Civic（58）。

很正常，它们均为中型轿车，因此 newCar 的市场定位不错，特别是它具有优秀的燃油效率。

对于 newTruck（焦点记录 159），最近相邻元素为 Nissan Quest（105）、Mercury Villager（92）和 Mercedes M-Class（101）。

如前面所提及的那样，它们并不一定是传统意义上的货车，而只是分类为非客车的车辆。来看最近相邻元素的“表”节点输出，可以看到 newTruck 相对较为昂贵，并且是同类中最重的车辆。不过，其燃油效率优于最接近的对手，因此值得青睐。

摘要

我们介绍了如何使用最近相邻元素分析来比较来自特定数据集的个案的多个特征。我们还为两个明显不同的保留记录计算了个案，它们最接近地呈现了这些保留。

注意事项

这些信息开发用于在全球提供的产品和服务。

IBM 可能在其他国家/地区中不提供在本文档中讨论的产品、服务或功能。请咨询您当地的 IBM 代表以了解有关您所在地区当前可用产品和服务的信息。任何对 IBM 产品、程序或服务的引用，并不意味着仅可使用这些 IBM 产品、程序或服务。作为替代，可以使用任何功能相当的产品、程序或服务，前提是不侵犯任何 IBM 知识产权。不过，用户应自行负责评估和验证任何非 IBM 产品、程序或服务的工作情况。

在本文档中介绍的主题可能涉及 IBM 的专利或申请中的专利。您获得该文档并不意味着授予您任何这些专利许可。您可以将书面许可查询函件发送至：

IBM Director of Licensing, IBM Corporation, North Castle Drive, Armonk, NY 10504-1785, U. S. A.

要查询双字节字符集 (DBCS) 相关许可证信息，请联系所在国家/地区中的 IBM 知识产权部门，或者以书面形式将查询函件发送至：

Intellectual Property Licensing, Legal and Intellectual Property Law, IBM Japan Ltd., 1623-14, Shimotsuruma, Yamato-shi, Kanagawa 242-8502 Japan.

以下段落不适用于英国或此类条款与当地法律不符的其他国家/地区： INTERNATIONAL BUSINESS MACHINES 公司“按原样”提供本出版物，不保证任何明示或暗示，包括但不限于对非侵权性、适销性或对特定用途适用性的暗示担保。有些州/省不允许针对某些交易的明示或暗示免责条款，因此本声明可能不适用于您。

本信息可能含有技术误差或排版错误。此处的信息会定期进行更改；这些更改会体现在本出版物的新版本中。IBM 可能随时对本出版物中所述的产品和/或程序进行改进和/或更改，恕不另行通知。

在本信息材料中对任何非 IBM 网站的引用仅为了方便用户，并不以任何方式表明对这些网站的认可。这些网站上的材料并非本 IBM 产品材料的一部分，您对这些网站的使用需自担风险。

IBM 可以自认为适当并且不会对您构成任何约束的任何方式使用或分发您提供的任何信息。

如果本程序的受许可方试图了解有关程序的信息以启用：(i) 在独立创建的程序和其他程序（包括本程序）之间交换信息；(ii) 相互使用交换的信息，则应联系：

IBM Software Group, Attention:Licensing, 233 S. Wacker Dr., Chicago, IL 60606, USA.

此类信息的提供应遵照相关条款和条件，其中包括在某些情况下支付适当费用。

在本文档中介绍的受许可保护程序，及其所有受许可保护材料由 IBM 在双方签署的“IBM 客户协议”、“IBM 国际程序许可证协议”或任何其他等同协议下提供。

此处所含的性能数据均在受控环境下决定。因此，在其他操作环境中获得的结果可能差异较大。有些测量可能在开发级的系统中进行，不保证这些测量结果与常用系统上的测量结果相同。此外，有些测量结果可能通过推断来估计得出。实际结果可能有所差异。此文档的用户应针对其具体环境验证适用的数据。

非 IBM 产品的相关信息来自这些产品的供应商，及其发布的公告或其他公开来源。IBM 并未测试过这些产品，无法确认有关非 IBM 产品的性能准确性、兼容性或任何其他声明。有关非 IBM 产品功能的问题应由这些产品的供应商负责。

有关 IBM 未来方向或意向的所有声明均可能未经通知即变更或撤销，并且仅代表目标和目的。

本信息包含日常业务运营中使用的数据和报告的示例。为了尽可能详尽地对其进行说明，示例中包含了人员的姓名、公司、品牌和产品的名称。所有这些名称均为虚构，与真实商业企业使用的名称和地址的任何雷同纯属巧合。

如果您正在查阅此信息的软拷贝，照片和彩色插图可能不会显示。

商标

IBM、IBM 徽标、ibm.com 和 SPSS 是 IBM Corporation 的商标，在全球许多司法辖区注册。有关最新的 IBM 商标列表，请访问网页 <http://www.ibm.com/legal/copytrade.shtml>。

Intel、Intel 徽标、Intel Inside、Intel Inside 徽标、Intel Centrino、Intel Centrino 徽标、Celeron、Intel Xeon、Intel SpeedStep、Itanium 和 Pentium 是 Intel Corporation 或其子公司在美国和其他国家或地区的商标或注册商标。

Linux 是 Linus Torvalds 在美国、其他国家或地区或两者的注册商标。

Microsoft、Windows、Windows NT 和 Windows 徽标是 Microsoft Corporation 在美国、其他国家或地区或两者的商标。

UNIX 是 The Open Group 在美国和其他国家或地区的注册商标。

Java 和所有基于 Java 的商标和徽标是 Sun Microsystems, Inc. 在美国、其他国家或地区或两者的商标。

其他产品和服务名称可能是 IBM 或其他公司的商标。



参考书目

Asuncion, A., 和 D. Newman. 2007. "UCI 机器学习存储库." Available at <http://mllearn.ics.uci.edu/MLRepository.html>.

- CLEM
 - 简介, 22
 - COP, 11
 - Cox 回归
 - 分类变量编码, 326
 - 变量选择, 327
 - 已审查的个案数, 325
 - 生存曲线, 330
 - 风险曲线, 331
 - CRISP-DM, 16
 - Excel
 - 与决策列表模型连接, 133
 - 修改决策列表模板, 138
 - fields
 - 筛选, 103
 - 选择分析, 103
 - 重要性排序, 103
 - gamma 回归
 - 在“广义线性模型”中, 304
 - IBM SPSS Modeler, 1, 13
 - 从命令行运行, 9
 - 文档, 3
 - 新手入门, 8
 - 概述, 8
 - IBM SPSS Modeler Server
 - password, 9
 - 主机名, 9-10
 - 域名 (Windows), 9
 - 用户 ID, 9
 - 端口号, 9-10
 - Microsoft Excel
 - 与决策列表模型连接, 133
 - 修改决策列表模板, 138
 - Omnibus 检验
 - 在“广义线性模型”中, 297
 - password
 - IBM SPSS Modeler Server, 9
 - SLRM 节点
 - 应用示例, 212
 - 构建流, 213
 - 流构建示例, 213
 - 浏览模型, 218
 - SPSS Modeler Server, 2
 - temp 目录, 12
 - Wilk 的 lambda
 - 判别式分析, 263
-
- 主机名
 - IBM SPSS Modeler Server, 9-10
 - 主窗口, 14
-
- 交互列表查看器
 - 使用, 120
 - 应用示例, 120
 - 预览窗格, 120
-
- 低概率搜索
 - 决策列表模型, 126
 - 余数
 - 决策列表模型, 120
 - 停止执行, 17
-
- 公用检验
 - 在“Cox 回归”中, 327
-
- 决策列表查看器, 120
 - 决策列表模型
 - 与 Excel 连接, 133
 - 使用 Excel 自定义测量量, 133
 - 保存会话信息, 140
 - 修改 Excel 模板, 138
 - 应用示例, 115
 - 生成, 140
 - 决策列表节点
 - 应用示例, 115
 - 准备, 93
-
- 分析节点, 101
 - 分类变量编码
 - 在“Cox 回归”中, 326
 - 分类表
 - 判别式分析, 266
 - 分组生存数据
 - 在“广义线性模型”中, 267
 - 判别式分析
 - Wilk 的 lambda, 263
 - 分类表, 266
 - 特征值, 263
 - 结构矩阵, 264
 - 逐步方法, 262
 - 面积图, 265
 - 剪切, 17
-
- 区间型删失的生存数据
 - 在“广义线性模型”中, 267
-
- 协变量均值
 - 在“Cox 回归”中, 329
 - 单点登录, 10
-
- 参数估计值
 - 在“广义线性模型”中, 273, 285, 298, 308
-
- 变量文件节点, 83
-
- 可视化编程, 13

- 向下搜索
 - 决策列表模型, 126
- 命令行
 - 启动 IBM SPSS Modeler, 9
- 商标, 376

- 图形节点, 92
- 图标
 - 设置选项, 19

- 块
 - 定义, 15
- 域名 (Windows)
 - IBM SPSS Modeler Server, 9

- 复制, 17

- 多个 IBM SPSS Modeler 会话, 13

- 导出节点, 93

- 工作区, 14
- 工具栏, 17
- 工程, 16

- 已审查的个案数
 - 在“Cox 回归”中, 325

- 市场购物篮分析, 357

- 广义线性模型
 - Omnibus 检验, 297
 - 参数估计值, 273, 285, 298, 308
 - 拟合优度, 297, 302
 - 模型效应检验, 272, 284, 298
 - 泊松回归, 292
- 应用程序示例, 3

- 建模, 96, 99, 101

- 快捷键
 - 键盘, 21

- 打印, 22
 - 流, 19
- 拟合优度
 - 在“广义线性模型”中, 297, 302
- 挖掘任务
 - 决策列表模型, 120
- 排序预测变量, 103
- 撤销, 17

- 数据
 - 建模, 96, 99, 101
 - 操作, 93
 - 查看, 87
 - 读取, 83

- 文档, 3

- 最小化, 19

- 服务器
 - 添加连接, 10
 - 登录, 9
 - 通过 COP 搜索服务器, 11

- 模型效应检验
 - 在“广义线性模型”中, 272, 284, 298

- 泊松回归
 - 在“广义线性模型”中, 292
- 法律注意事项, 375
- 流, 8, 14
 - 构建, 83
 - 缩放以查看, 19
- 添加 IBM SPSS Modeler Server 连接, 10 - 11
- 源节点, 83

- 热键, 21

- 片段
 - 从评分中排除, 129
 - 决策列表模型, 120

- 特征值
 - 判别式分析, 263
- 特征选择模型, 103
- 特征选择节点
 - 排序预测变量, 103
 - 筛选预测变量, 103
 - 重要性, 103

- 状态监测, 248

- 生存曲线
 - 在“Cox 回归”中, 330
- 生成的模型选项板, 15

- 用户 ID
 - IBM SPSS Modeler Server, 9

- 登录到 IBM SPSS Modeler Server, 9

索引

示例

- KNN, 364
- SVM, 309
- 产品分类销售, 196
- 判别式分析, 254
- 多项 logistic 回归, 142, 152
- 字符串长度减少, 109
- 市场购物篮分析, 357
- 应用程序指南, 3
- 新车辆产品评估, 364
- 概述, 5
- 状态监测, 248
- 电信, 142, 152, 165, 187, 254
- 细胞样本分类, 309
- 贝叶斯网络, 223, 232
- 输入字符串长度减少, 109
- 重新分类节点, 109
- 零售分析, 243

端口号

- IBM SPSS Modeler Server, 9 - 10

筛选预测变量, 103

简介

- IBM SPSS Modeler, 8

管理器, 15

类, 16

粘贴, 17

结构矩阵

- 判别式分析, 264

缩放, 17

缩放流以查看, 19

网络节点, 92

脚本编写, 22

自学响应模型节点

- 应用示例, 212
- 构建流, 213
- 流构建示例, 213
- 浏览模型, 218

节点, 8

表节点, 87

表达式构建器, 93

调整大小, 19

调色板, 14

负二项式回归

- 在“广义线性模型”中, 299

输出, 15

过滤, 96

进程协调器, 11

连接

- 服务器群集, 11

- 至 IBM SPSS Modeler Server, 9 - 11

逐步方法

- 判别式分析, 262

- 在“Cox 回归”中, 327

通过 COP 搜索连接, 11

重要性

- 排序预测变量, 103

零售分析, 243

面积图

- 判别式分析, 265

预测变量

- 筛选, 103

- 选择分析, 103

- 重要性排序, 103

风险曲线

- 在“Cox 回归”中, 331

鼠标

- 在 IBM SPSS Modeler 中使用, 21

鼠标中键

- 模拟, 21