# DB2 Intelligent Miner for Data and
# DB2 Intelligent Miner Scoring, Modeling, & Visualization

Gregor Meyer
gregorm@us.ibm.com

IBM Software Group

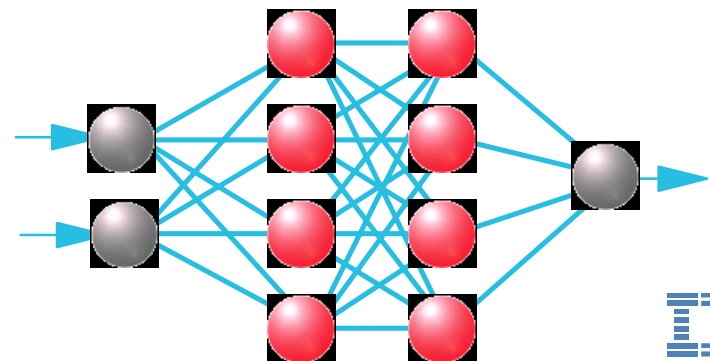# Agenda

- **Data Mining, quick overview**

  ► **Concepts, Business uses**

- **DB2 Intelligent Miner Technology**

  ► **Intelligent Miner for Data, the workbench**

  ► **IM Scoring, Modeling, Visualization**

- **Trends, directions**

  ► **Are we on the right track?**

IBM

# We've got plenty of data.  What we need are answers.

personalization

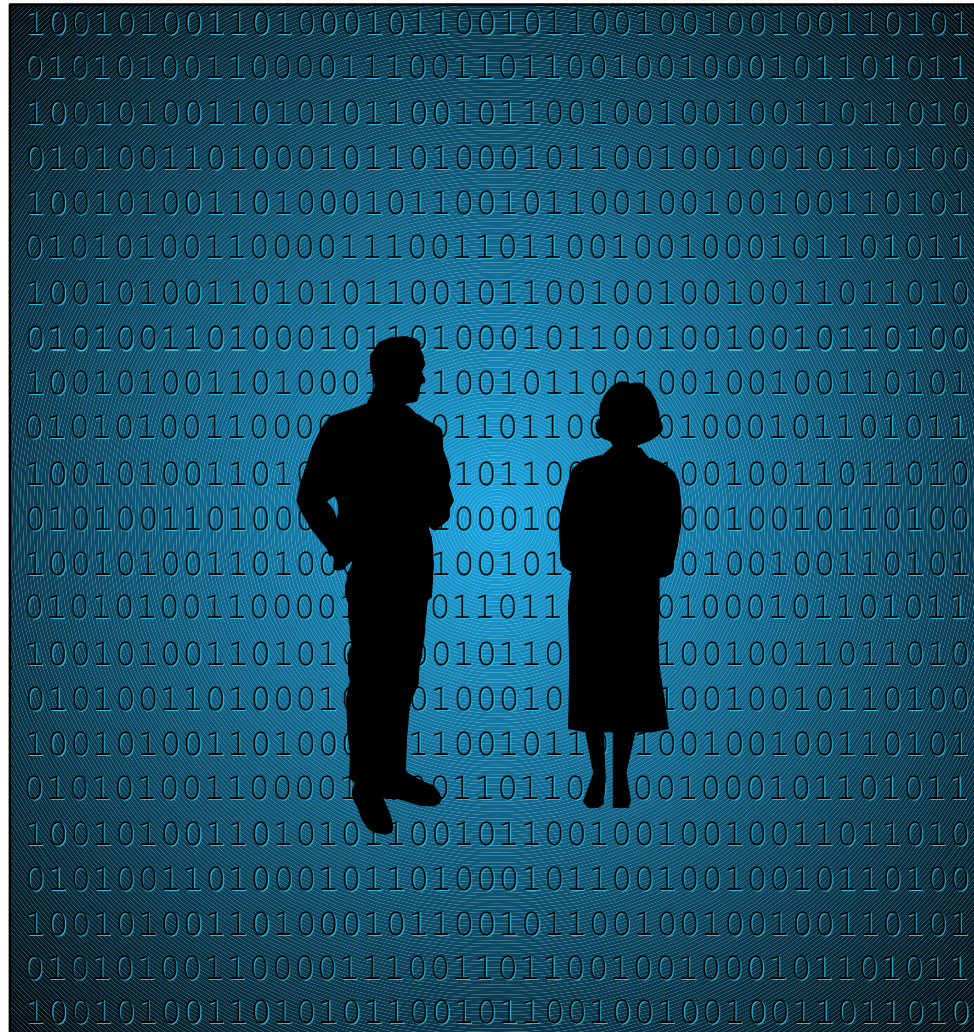find important attributes

Marketing automation

what is our risk?



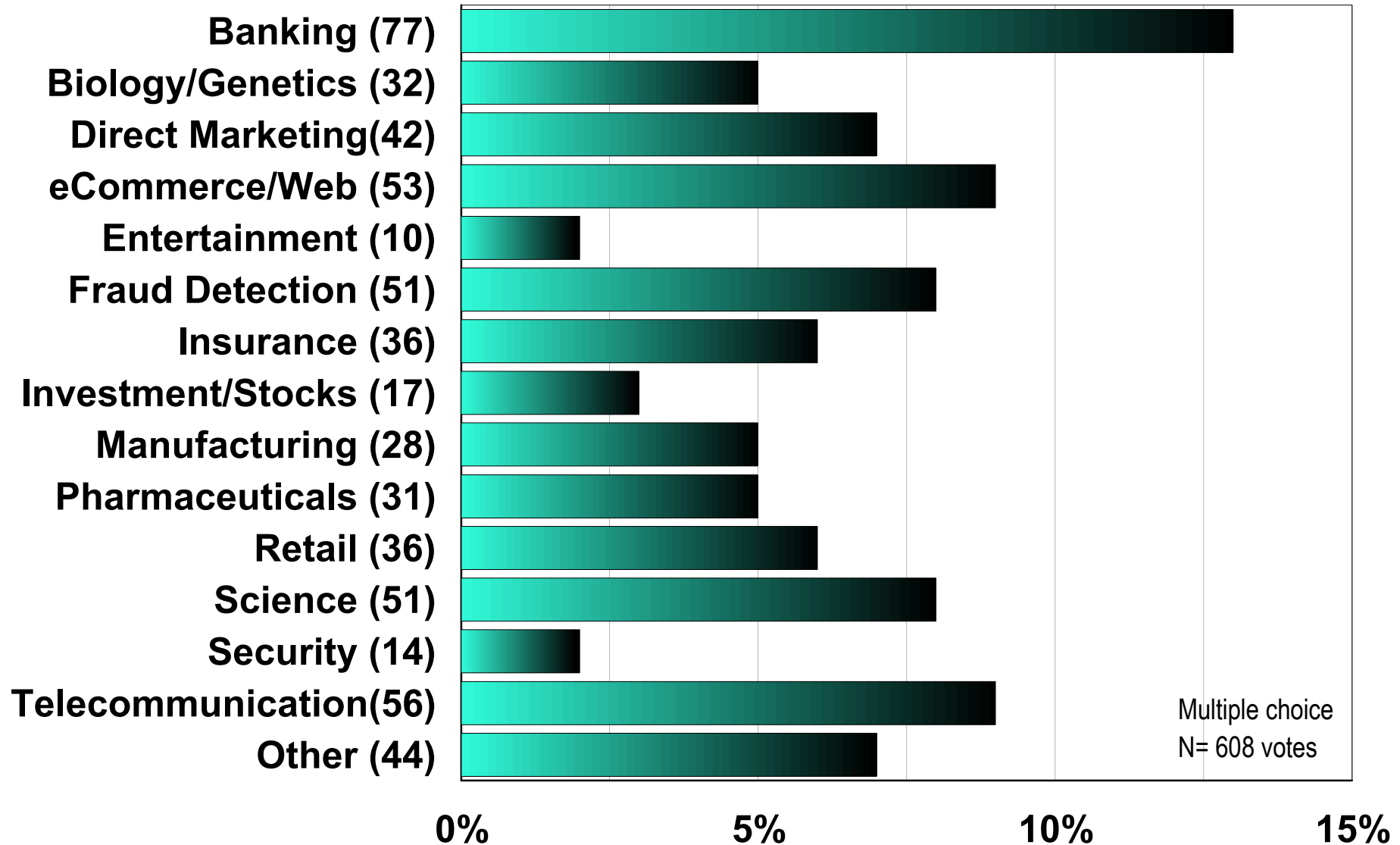give the right message to the right person

predict customer behavior

detect anomalies

quality analysis

1-to-1 Marketing

IBM

# Where do you plan to use data mining in 2002?

| Category | Percentage |
|---|---|
| Banking (77) | ~13% |
| Biology/Genetics (32) | ~5% |
| Direct Marketing(42) | ~7% |
| eCommerce/Web (53) | ~9% |
| Entertainment (10) | ~2% |
| Fraud Detection (51) | ~8% |
| Insurance (36) | ~6% |
| Investment/Stocks (17) | ~3% |
| Manufacturing (28) | ~5% |
| Pharmaceuticals (31) | ~5% |
| Retail (36) | ~6% |
| Science (51) | ~8% |
| Security (14) | ~2% |
| Telecommunication(56) | ~9% |
| Other (44) | ~7% |

Multiple choice
N= 608 votes

0%    5%    10%    15%

source: KDnuggets : Polls : Data Mining Applications in June 2002
http://www.kdnuggets.com/polls/current_application_fields_2002.htm

# Primary Data Mining Techniques

## Clustering

- What are the ages, ethnicity, location, family size, and affluence of our clients?
- What are the ages, ethnicity, location, family size, and affluence of our clients by profit, number of products, and product groups?

## Prediction

- Which prospects are most likely to buy a specific product or service?
- Which clients are at risk of defecting to competitors? What are the attributes associated with them leaving?

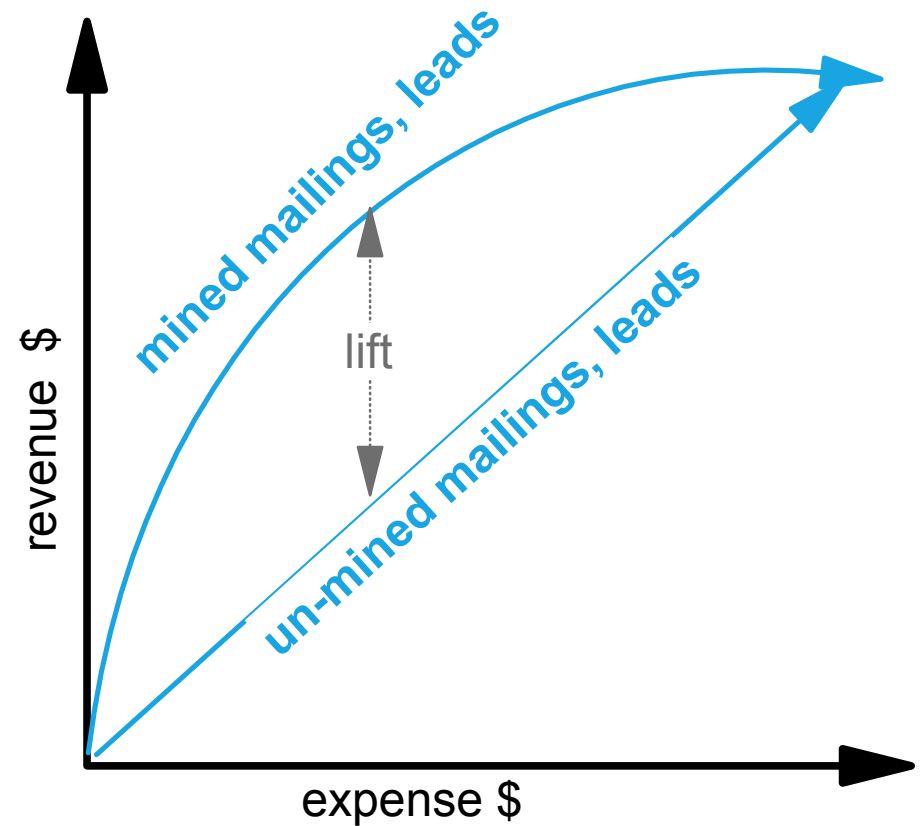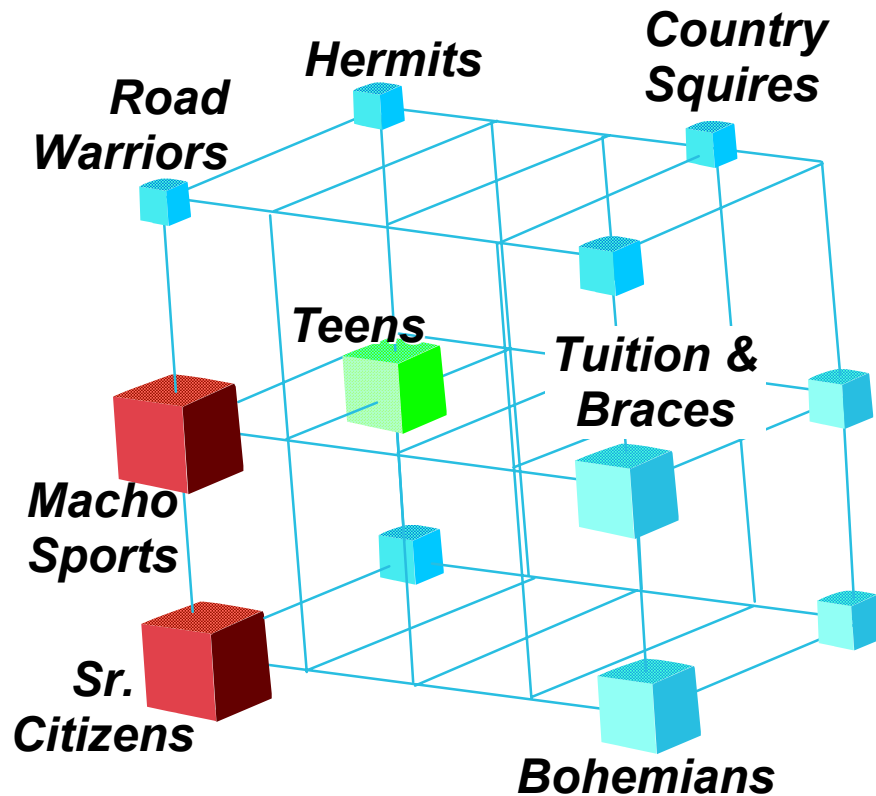## Affinity

- Which products are most often purchased by the most profitable clients?
- Are there any product purchases that often trigger additional purchases later?

## Sequences & Pattern Detection

- Are there any repeating patterns for customers who stop buying or cancel service?
- Where are there patterns of purchases and returns?

# Typical Business Uses



**Clustering & Segmentation**

**Prediction**
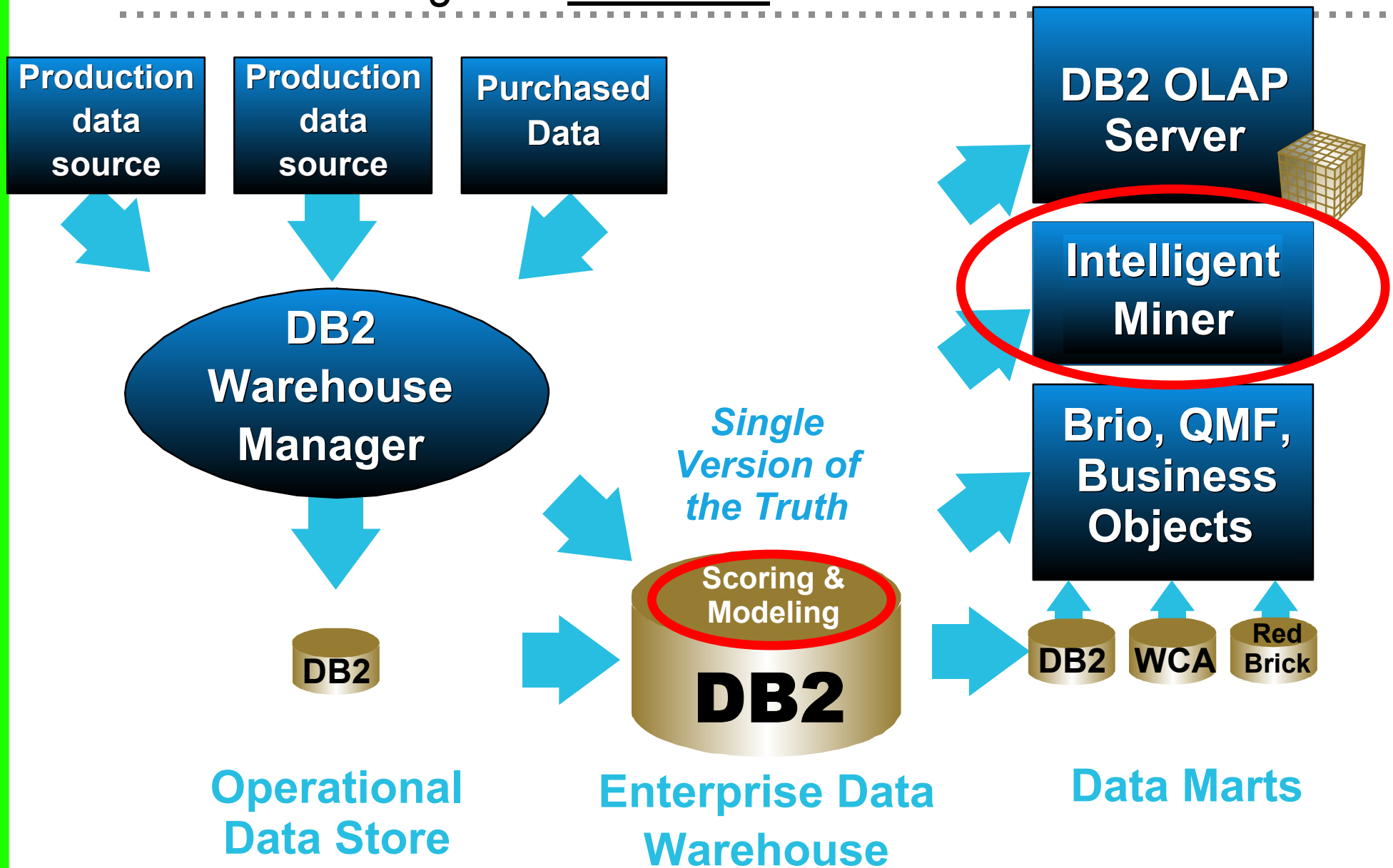
IBM

# Clustering (banking example)

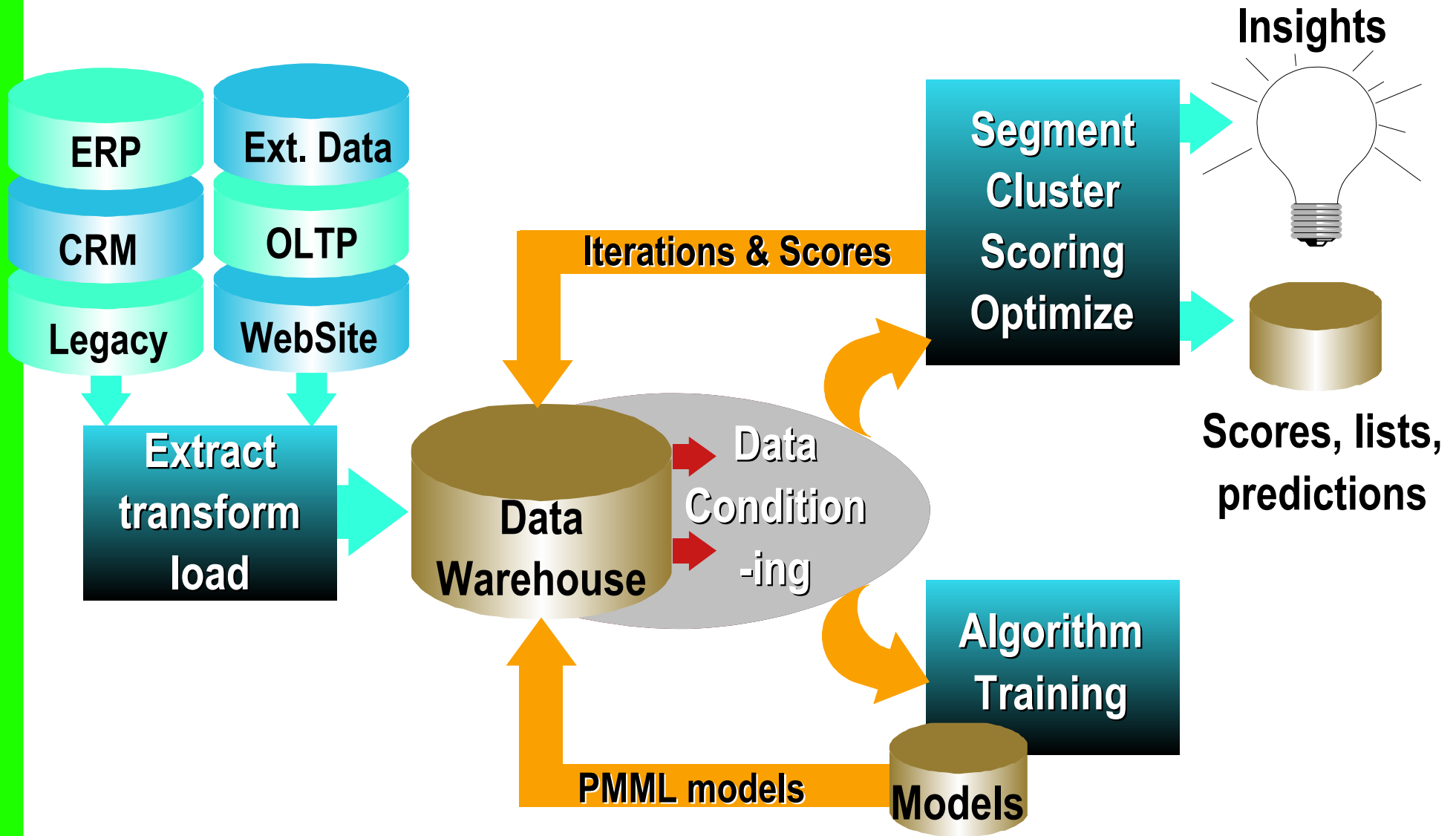| Segments | Segment Services used | Primary differentiator | Secondary |
|---|---|---|---|
| Bohemians | checking, credit card | high restaurant cc use | 3X credit card usage |
| Road Warrior | Business credit line, checking, college funds, mortgage | business credit card & $5.5K monthly balance | out of country ATM use |
| Senior Citizens | checking, credit card | 3+ savings accounts > $20K | mortgage principle < 18% loan |
| Country Squires | mortgage, credit line, checking, porfolio, platinum credit card | 7X more likely to have boat loan | rental property |
| Tuition and Braces | mortgage, checking, credit card, college savings | Mortgage < 8 years old | 2X credit card usage |
| • • • | • • • | • • • | • • • |
| Hermits | checking, credit card, mortgage | cc usage < 8 per month | electronics purchases |

# Prediction (banking example)

| Segments | Will Buy | Will Attrite |
|---|---|---|
| **Bohemians** | ▪ **Merchant co-sponsored credit card 28% of time** | ▪ **balance of 33 cc's reaches 85% of limit** |
| **Road Warrior** | ▪ **Airline co-sponsored cc 43% of time**<br>▪ **2nd home mortgage 11% of time** | ▪ **business cc balance drops 70% for 3 consecutive months** |
| **Senior Citizens** | ▪ **Bond fund ABC 21% of time**<br>▪ **Bond fund XYZ 17% of time** | ▪ **3 call center complaints < 4 months** |
| **Country Squires** | ▪ **foreign country home mortgage 21% of time**<br>▪ **Bond fund XYZ 31% of time** | ▪ **relocation & nearest branch office > 4 miles** |
| **Tuition and Braces** | ▪ **Refi mortgage when 1.8% delta** | ▪ **College fund < $10K and cc balance < $3K and call center complaints in last two months** |
| • • • | • • • | • • • |
| **Hermits** | | |

IBM

# Business Intelligence **Products**

**Production data source**

**Production data source**

**Purchased Data**

**DB2 OLAP Server**

**DB2 Warehouse Manager**

**Intelligent Miner**

*Single Version of the Truth*

**Brio, QMF, Business Objects**

DB2

**Scoring & Modeling**

**DB2**

DB2   WCA   Red Brick

**Operational Data Store**

**Enterprise Data Warehouse**

**Data Marts**

**DB2** **Data Management Software**
Copyright IBM 2002

IBM

# The Data Mining process

ERP

Ext. Data

CRM

OLTP

Legacy

WebSite

**Extract transform load**

**Data Warehouse**

**Data Condition -ing**

**Iterations & Scores**

**Segment Cluster Scoring Optimize**

**Insights**

**Scores, lists, predictions**

**Algorithm Training**

**PMML models**

**Models**

# Data Mining products

## Workbench

**Statistician**

**DB2**

extract

**data warehouse**

## RDBMS Extenders

**SQL invokes extender**

**Programmer**

**DB2 Instance**

**DB2 Instance**

**DB2 Instance**

**DB2 Instance**

## Application Embedded

**Consumer**

**Packaged Applications**

**algorithm**

SAS   Unica   hnc   Retek

IBM

# Visualization of database attributes

**Categorical Variable**
**inner circle - this segment**
**outer circle - total population**

**Unknown**

**Female**

**Male**

**This segment contains more males and fewer females than the overall population sample**

GENDER

DB2 **Data Management Software**
Copyright IBM 2002

IBM

# Continuous Variable Visualization



**Total sample Population**

**This Segment**

bin 1     bin 2     bin 3     bin 4     bin 5     bin 6

Low  ←——————→ **Annual Income** ——————→ High

IBM

# Intelligent Miner Segments Visualization



PM_usage     multi_lines     wireless_Rm     T_rev98     call_card     .47

MM_usage     multi_lines     ISP_usa     Call_waiting     wireless_3T     .44

MM_usage     T_rev98     Call_waiting     call_card     PM_usage     .43

res_bldg     call_card     wireless_3T     Sys_rx_ft     multi_lines     .36

ISP_amt     ISP_amt     MM_usage     T_rev98     Rtx     .33

rural_zip     wireless_Rm     ISP_amt     PM_usage     call_card     .29

Call_waiting     multi_lines     PM_usage     MM_usage     ISP_amt     .27

MM_usage     T_rev98     rural_zip     multi_lines     PM_usage     .19

# Associations Visualizer

List    Edit    View      Help

| Support % | Confidence% | Lift | Item Set |
|:---:|:---:|:---:|:---|
| 3.00 | 100.0 | 26.70 | [cheddar cheese]+[rn crackers] |
| 3.37 | 100.0 | 26.70 | [cheddar cheese]+[wtr crackers] |
| 3.00 | 100.0 | 26.70 | [mineral water]+[limes]+[napkins] |
| 3.00 | 80.0 | 21.36 | [rn crackers]+[Wines] |
| 2.55 | 82.0 | 21.36 | [mineral water]+[lemons]+[wtr crackers] |
| 3.07 | 100.0 | 19.42 | [soft drink]+[salty snacks] |
| 3.37 | 100.0 | 19.42 | [mineral water]+[fruit juice] |
| 3.00 | 72.7 | 19.07 | [Wines]+[brie cheese] |
| 2.52 | 69.3 | 19.07 | [Film photo]+[salty snacks]+[soft drink] |
| 2.44 | 100.0 | 19.07 | [Fruit juice]+[pop tarts] |
| 3.12 | 58.4 | 19.07 | [Dish soap]+[beers]+[salty snacks] |
| 3.00 | 63.3 | 19.07 | [baby carrots]+[blue cheese] |
| 3.09 | 66.7 | 19.07 | [Motor oils]+[Tide detergent] |
| 3.21 | 90.2 | 19.07 | [Beer]+[salty snacks] |
| 3.28 | 74.2 | 18.48 | [mineral water]+[household] |
| 3.00 | 80.0 | 17.80 | [eggs]+[confetti] |
| 3.00 | 82.0 | 17.16 | [soft drink]+[cheese platter]+[film] |

Internet

# The PMML Interoperability Standard

consume
pmml

emit
pmml

**Unica**

**SPSS**

**Epiphany**

**Startup
de Jure**

**Walker**

**WebSphere
apps**

**Intelligent Miner**

**PMML
model(s)**

**SAS**

**SAP**

**Kana**

**Fair Isaac**

**Chordiant**

**Query Tools**

**WebTrends**

**Hyperion**

**JDA**

emit
pmml

consume
pmml

open standards enable innovation, partnerships, and choice

**DB2** **Data Management Software**
Copyright IBM 2002

IBM

# What is PMML?

- **Predictive Model Markup Language**
- **PMML 2.0 is a standard for XML documents which express trained instances of analytic models. The following classes of model are addressed**
  - ► Association Rules, Decision Trees , Center-Based & Distribution-Based Clustering , Polynomial Regression , General Regression, Neural Networks, Naive Bayes
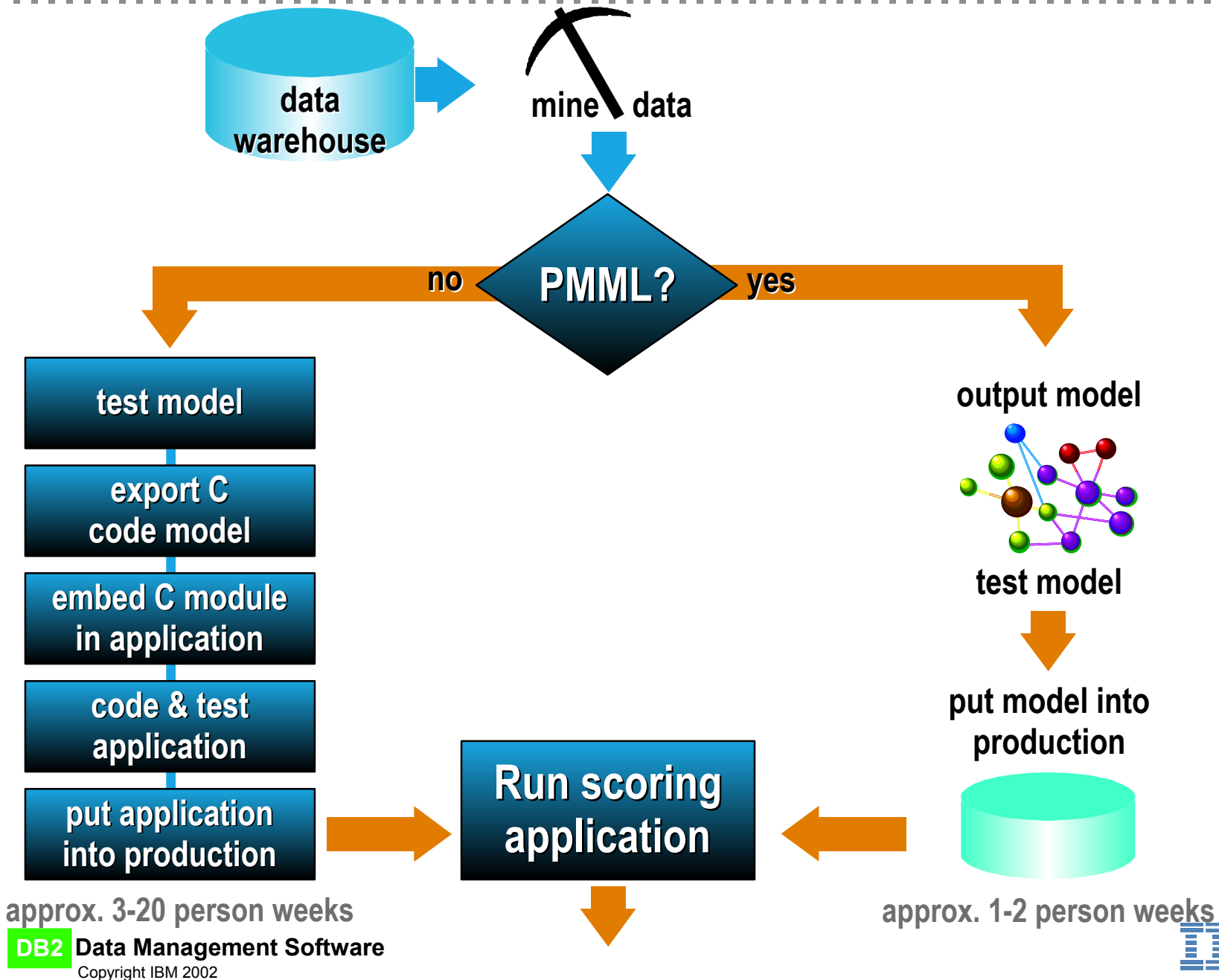- **A standard for developing a data mining model -- the training set -- to be used by one or more consuming software tools**
- **PMML Standards supporters**
  - ► Angoss Software, IBM,  Magnify, Oracle, National Center for Data Mining  University of Illinois, SAS, SPSS Inc., Xchange, MINEit Software, NCR, KXEN, ...
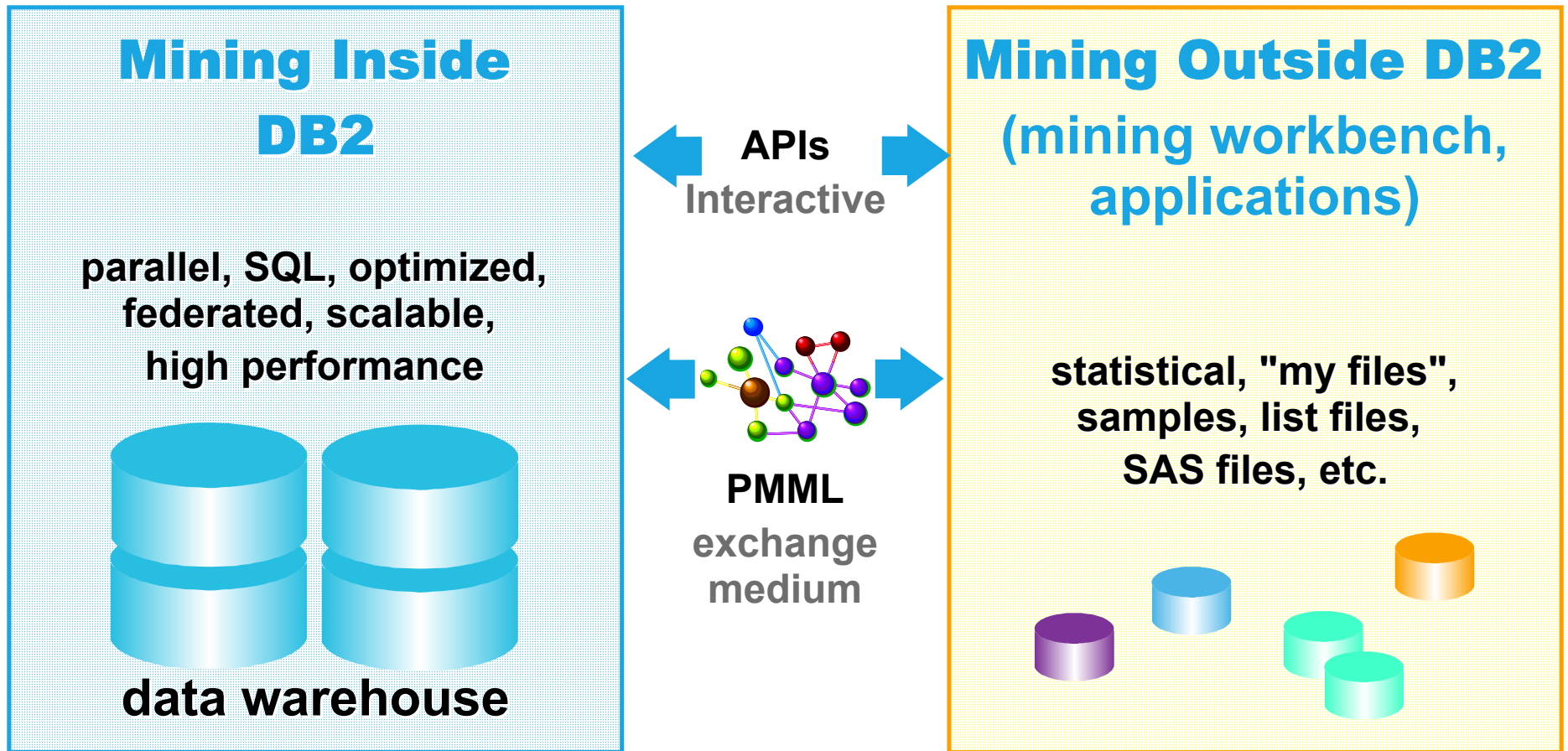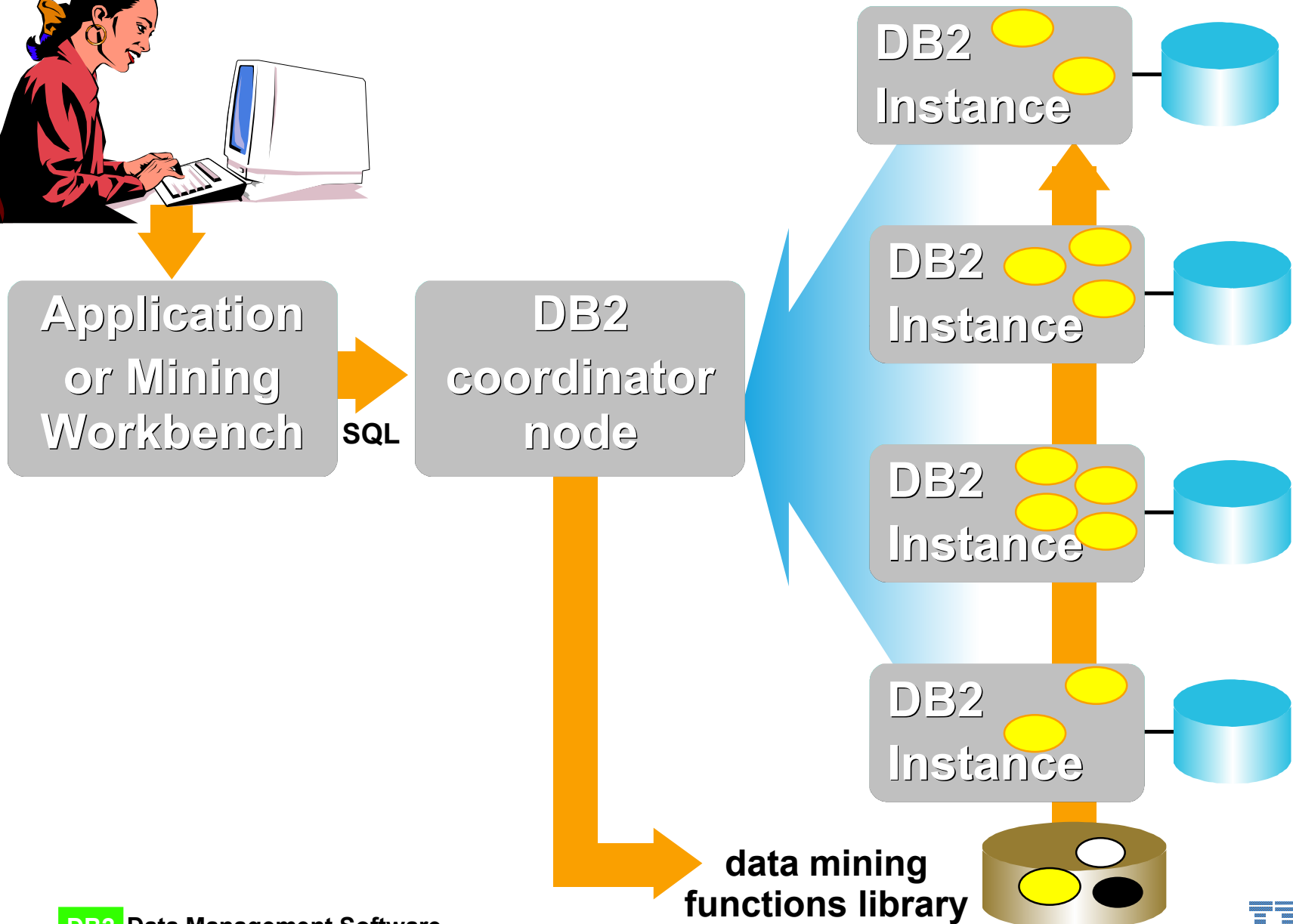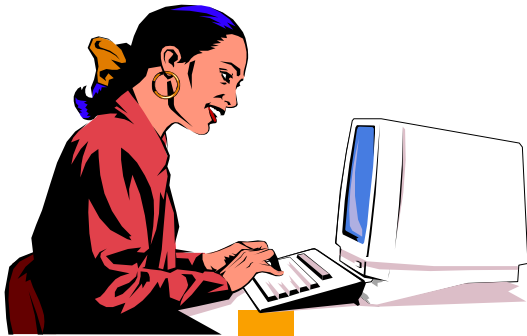- **http://www.dmg.org**

IBM

# PMML Benefits - Reduced Labor costs



data warehouse → mine data

PMML?

no → test model → export C code model → embed C module in application → code & test application → put application into production

yes → output model → test model → put model into production

Run scoring application

approx. 3-20 person weeks

approx. 1-2 person weeks

IBM

# Mining Inside & Outside of the RDBMS



**Mining Inside DB2**

parallel, SQL, optimized, federated, scalable, high performance

**data warehouse**

**APIs**
Interactive

**PMML**
exchange medium

**Mining Outside DB2 (mining workbench, applications)**
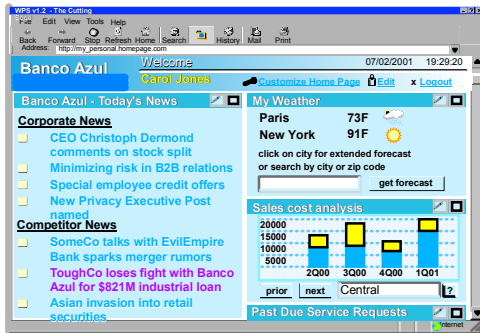
statistical, "my files", samples, list files, SAS files, etc.

IBM

# End User activates Mining Extenders

**Application or Mining Workbench**

**SQL**

**DB2 coordinator node**

**DB2 Instance**

**DB2 Instance**

**DB2 Instance**

**DB2 Instance**

**data mining functions library**

IBM

# Applications Connect to Real Time Data Mining



**BI Reports**
**Portal KPIs**

**Call Center**
**Web Site**
Triggers or Agents

**PMML**

**Data Mining**

**Data Warehouse**

**Production Database**

IBM

# Benefits of DB2 Extenders Data Mining

- **Data mining can be performed using parallel processing**
  - ► Workload scalability --big tasks
  - ► SMP and/or cluster parallelism = response time reduction
    - – More "runs" per day for deeper accuracy and more discoveries
- **Avoiding extracts & loads -- less data movement & redundancy**
  - ► less labor costs, less processing = lower costs overall
- **EDW Server has the bigger faster CPUs, memory, disks, etc.**
  - ► CPU intensive workloads finish quicker
- **More data can be mined resulting in more detailed analysis**
  - ► Sometimes "sample" subsets aren't enough for accuracy
    - – e.g. deviation detection
- **Scoring Extenders can be leveraged in online applications**

IBM

# SQL Modeling, example

```
call IDMMX.BuildClusModel(
        'CustomerSegments',
        IDMMX.DM_MiningData()
        ..DM_defMiningData('BANKING_CUST'),
        IDMMX.genClusSettings(
                IDMMX.DM_LogicalDataSpec()
                ..DM_genLogicalDataSpec('BANKING_CUST')
                ..DM_remDataSpecFld( 'GENDER' ) )
        ..DM_setMaxNumClus(9)
        ..DM_setFldUsageType( 'PRODUCT', 2 )
        ..expClusSettings()    );
```

# SQL Scoring, example

```
SELECT
    d.name, d.age,
    IDMMX.DM_getClusterID(
        IDMMX.DM_applyClusModel(
            cm.model,
            IDMMX.DM_impApplData(
                REC2XML(1,'COLATTVAL','',
                    d.age, d.salary,
                    d.region,d.product,
                    ..., d.goldcard)))
    )
FROM ClusterModels cm, MyData d
WHERE cm.modelname='CustomerSegments';
```

IBM

# SQL mining 'macros', ease of use

```
Call CreateClassifier(
      'Campaign',        -- name of new classification model
      'Customer',        -- input data
      'Response' );      -- target field


Call CreateScoringView(
    'CustomerWithScore',   -- name of new SQL VIEW
    'Customer'             -- source table

    'Campaign');           -- prediction model
```

IBM

# SQL mining 'macros', ease of use

```
Select ColumnName, Rank
From Table ( IDMMX.InfluenceFactors(
                'CustomerView', 'SalesPerMonth') ) F
order by rank descending;
```

-- returns table of colunms in 'CustomerView'

-- ranked by how much they are related to 'SalesPerMonth'

```
  ColumName,    Rank
----------------------------------
  Location      0.64
  Age           0.43
  Gender        0.142

  ...           ...
```
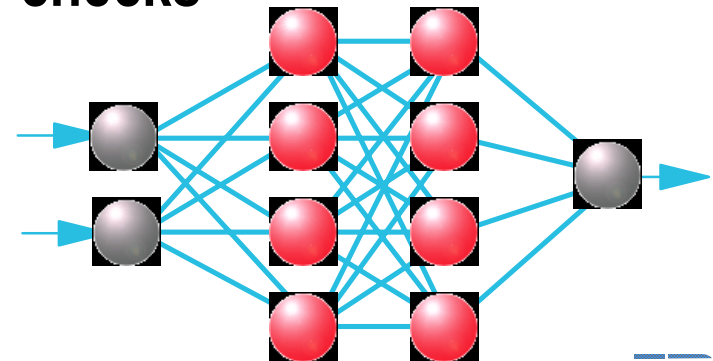
IBM

# Trends, directions for integrated mining

- **Ease of use, vertical integration**
  - ► Mining integrated into domain specific applications
    - − E.g. prediction in campaign management
  - ► End-user does not (have to) know mining
  - ► Mining expertise encapsulated in app.design
  - ► Robust, smart mining algorithms, visualization components
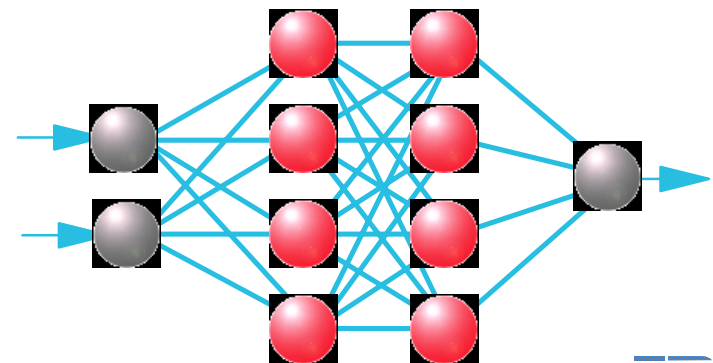- **Ease of use, intelligent helpers**
  - ► Find important influence factors
  - ► Outliers, anomalies, data quality checks
  - ► Guided OLAP, 'drill mine'

# Trends, directions for mining in the database

- **High performance, scalable**

- **Simple maintenance**
  - ► **common API, no separate mining tool admin**
- **Ease of use, intelligent helpers**
  - ► **'just another' simple database function**

# Trends, directions for standards

- **PMML**
  - **Model exchange format**
    - **e.g., SAS Modeling -> DB2 Scoring**
  - **Challenge: 'my algorithm is better than yours'**
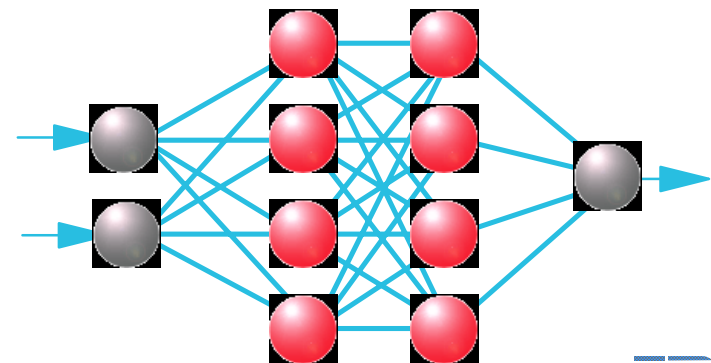- **SQL:  ISO SQL/MM, Microsoft OLE DB**
  - **"2 + 1" major players**
- **Java**
  - **New standard, extensive API**
- **Web Services, XMLA**
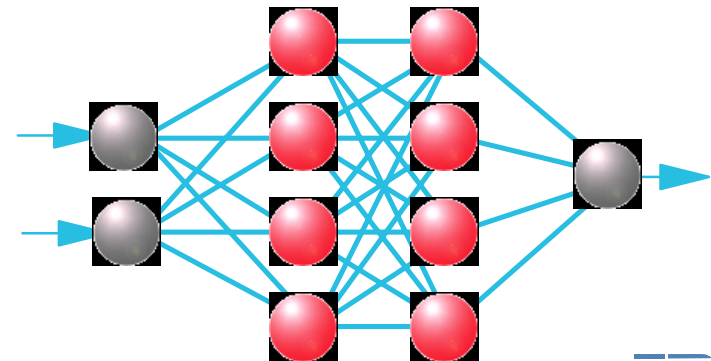  - **Business uses? Scoring.**

# Trends, directions for mining

- **End of mining workbenches for experts?**
  - ► No, still needed for flexible data analysis
  - ► But, also some support for 'mining' in other BI tools.
- **Database API with mining**
  - ► For heavy-lifting, back-end to mining tools & apps
  - ► Integrated with other transformations & information flow
  - ► Easy to use, intelligent helper functions that are packaged with the data warehouse server at no additional cost.

# Summary

- **Data Mining is vital to most industries, especially in CRM-Analytics**

- **Businesses use Data Mining for**
  - ► Lead generation & mailing lists
  - ► Cross-sell and Up-sell predictions
  - ► Anomaly detection

- **Intelligent Miner For Data & Extenders:**
  - ► Run Inside DB2 --fast and in parallel
  - ► Have the most algorithms to fit the most business needs
  - ► Setting the pace with Standards

- **Easy mining**