# Data Mining with SAS

Mathias Lanner
mathias.lanner@swe.sas.com

§sas | THE POWER TO KNOW.

# Agenda

- Data mining Introduction

- Data mining applications

- Data mining techniques

- SEMMA mythology

- Survival data mining

- Time series data mining

# What is data mining?

*"Data mining uses sophisticated statistical analysis and modeling techniques to uncover patterns and relationships hidden in organizational databases-patterns that ordinary methods might miss."*
-Two Crows Corporation (1998),p.1

*"Data Mining, as we use the term, is the exploration and analysis by automatic or semiautomatic means, of large quantities of data in order to discover meaningsful patterns and rules."*
-Berry and Linoff(1997), p.5

*"Data Mining [is] the process of efficient discovery of nonobvious valuble information from large collection of data."*
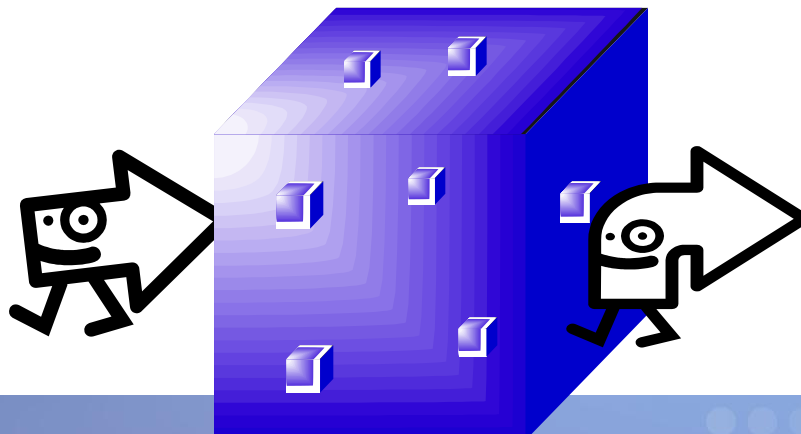-Berson and Smith (1997), p.565

*"Data Mining is the process of discovering meaningful new correlations, patterns and trends by sifting through large amounts of data stored in repositories, using pattern recognation technologies as well as statistical and mathematical techniques."*
-Erick Brethnoux, Gartner Group

# Data Mining Definition :

The *process of selecting, exploring, and modeling* large amounts of data to uncover previously unknown information for a *business advantage*
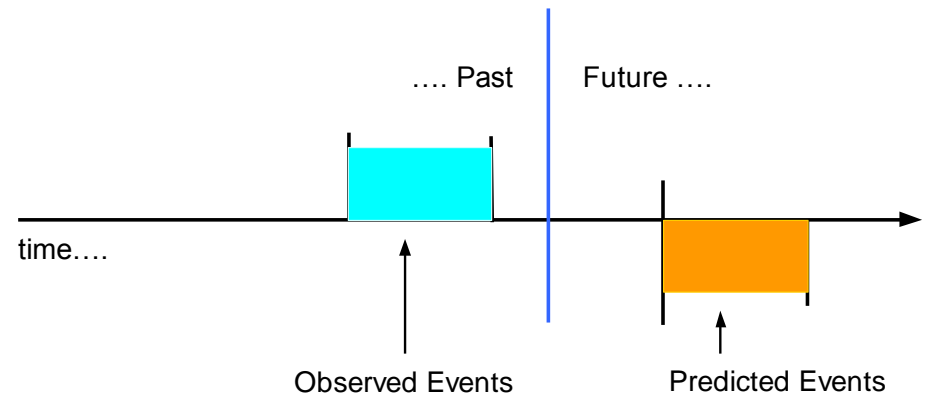
**Operational Systems & Data Warehouse**

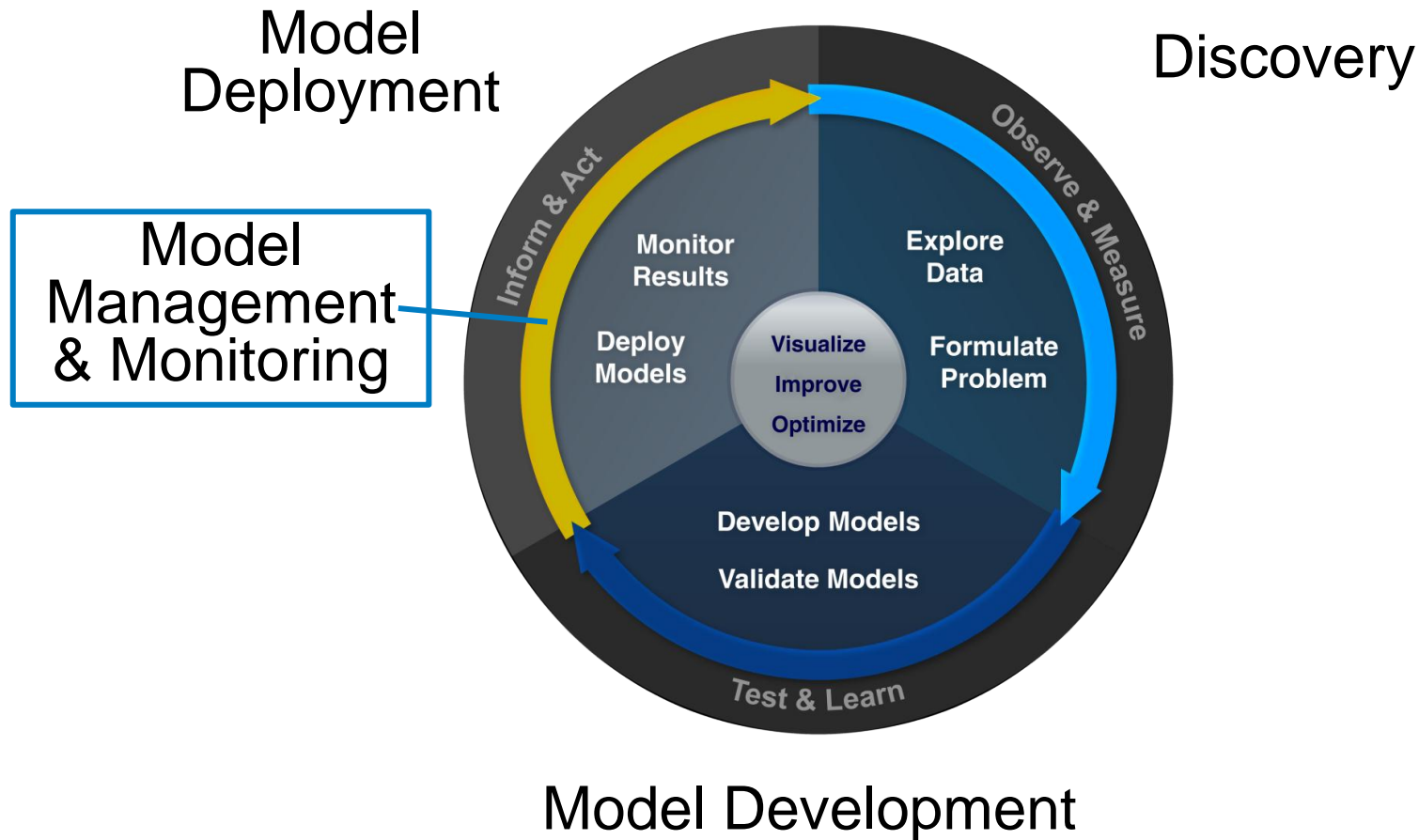**Significant & Actionable Information**

4

§sas | THE POWER TO KNOW.

# Data Mining Is:

- Discovering patterns, trends and relationships represented in data

- Developing models to understand and describe characteristics and activity based on these patterns

- Use insights to help evaluate future options and take fact-based decisions

- Deploy scores and results for timely, appropriate action

.... Past    Future ....

time....

Observed Events    Predicted Events

# Predictive Analytics and Data Mining
*Key Components*

Model Deployment

Discovery

Model Management & Monitoring



Model Development

# Cross-Industry Data Mining Applications
## *Customer Analytics*

| Application | What is Predicted? | Driven Business Decision |
|---|---|---|
| Profiling and Segmentation | Customer's behaviors and needs by segment | How to create better-targeted product/service offers? |
| Cross-sell and Up-Sell | Identify what will customer's buy? | Which product/service to recommend? |
| Acquisition and Retention | Customer's preferences and purchase patterns | How to grow and maintain valuable customers? |
| Campaign Management | Evaluate the success of customer communications | How to direct right offer to right person at the right time? |
| Profitability and Life-time Value | Understand the drivers of future value (margin and retention) | Identify economically valuable channels/demographics and incremental benefits? |

§sas | THE POWER TO KNOW.

# Industry Specific Data Mining Applications

| Application | What is Predicted? | Driven Business Decision |
|---|---|---|
| Credit Scoring (Banking) | Measure credit worthiness of new and existing set of customers | How to assess and control risk within existing (or new) consumer portfolios? |
| Market Basket Analysis (Retail) | Which products are likely to purchased together? | How to increase sales with cross-sell/up-sell, loyalty programs, promotions? |
| Asset Maintenance (Utilities, Mfg., Oil & Gas) | Identify real drivers of asset or equipment failure | How to minimize operational disruptions and maintenance costs? |
| Health & Condition Mgmt. (Health Insurance) | Identify patients at risk of a chronic illness & offer treatment program | How can we reduce healthcare costs and satisfy patients? |
| Fraud Mgmt. (Govt., Insurance, Banks) | Detect unknown fraud cases and future risks | How to decrease fraud losses and lower false positives? |
| Drug Discovery (Life Science) | Find compounds that have desirable effects & detect drug behavior during trials | How to bring drugs quickly and effectively to the marketplace? |

§sas | THE POWER TO KNOW.

# Data mining techniques - Connect business problems with the right analytical technique
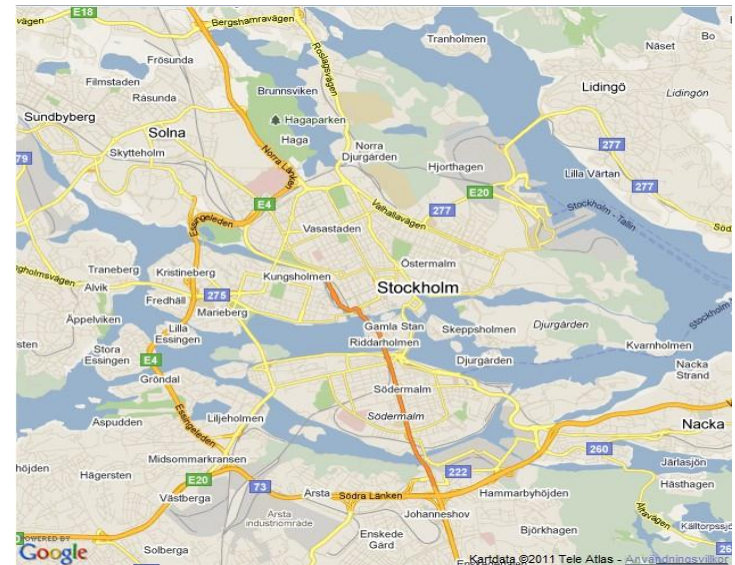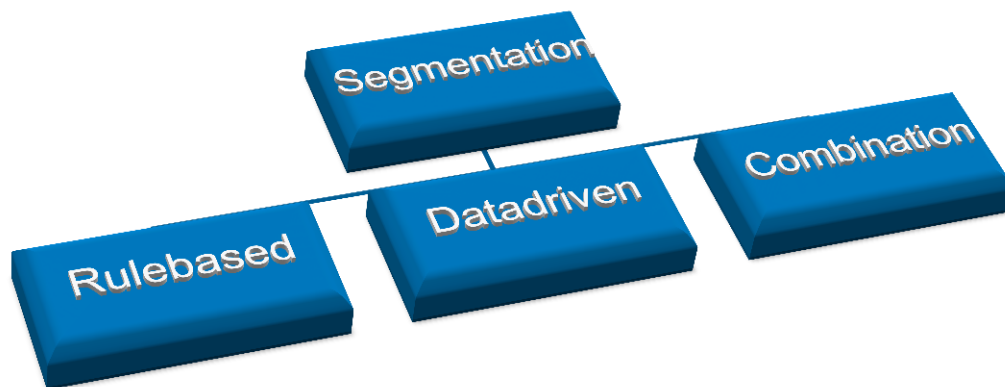
- Basket-/Sequence analysis
  - Find the association between products purchased together or sequentially

- Clustering/Segmentation
  - Divide the customers into different groups for different campaigns

- Predictive modeling
  - Classification – Find potential buyers
  - Prediction – Predict the future value

SAS | THE POWER TO KNOW.

# Predictiv Modeling- how does it works

- Predict how someone / something will behave (prediction / scoring)
  - Will this customer respond to the offer?
  - Will this borrower manage their interest payments?

- By starting from how other individuals / units behaved
  - Database – loyalty card
    - » Information about customers (age, gender, buying amount…)
    - » Information about their previous responds
  - Database –  Credit application
    - » Information about customers (age, gender, income, occupation, etc.)
    - » Information about how they have performed theirs interest payments

# Clustring/Segmentation

- A strategic map of the customer base

- Communicate it into the business

- The division can be done in different ways

- What information will be used to generate the segments

# Clustring/Segmentation

*"Create groups that have similar characteristics - also provides a measure of how big the difference is between the different groups*

**Cluster B**
Average income, new customers, tenant

**Cluster A**
High value, high income, home owner

**Cluster C**
Low value, low income, inactive customer

# Association / Sequence Analysis

*"Identify events that occur in association with each other, possibly in a particular order "*

What products customers purchase together or in a certain sequence?

Can we predict what you will have in your shopping cart?

§sas | THE POWER TO KNOW.

# SAS® Enterprise Miner™ 7.1
## *Model Development Process*

# S ample    E xplore    M odify    M odel    A ssess

| Sample | Explore | | Modify | Model | | Assess |
|---|---|---|---|---|---|---|
| Input Data | Association | DMDB | Transform Variables | Decision Tree | Neural Network | Model Comparison |
| File Import | Cluster | SOM/Kohonen | Impute | AutoNeural | SVM | Score |
| Sample | Variable Selection | Graph Explore | Replacement | Dmine Regression | Partial Least Squares | Segment Profile |
| Data Partition | Market Basket | | Interactive Binning | DMNeural | Regression | Decisions |
| Merge | StatExplore | | Rules Builder | Ensemble | Rule Induction | Cutoff |
| Filter | Variable Clustering | | Drop | Gradient Boosting | TwoStage | |
| Append | MultiPlot | | Principal Components | LARS | Model Import | |
| Time Series | Path Analysis | | | MBR | | |

§sas | THE POWER TO KNOW.

# SAS® Enterprise Miner™ 7.1
## *Model Development Process*

### Utility

- Metadata
- SAS Code
- Start Groups
- End Groups
- Control Point
- Reporter
- Score Code Export
- Ext Demo

### Applications

- Survival
- Ratemaking

### Time Series

- TS Similarity
- TS Exponential Smoothing
- TS Data Preparation

# SAS® Enterprise Miner™ 7.1
## *SEMMA in Action – Repeatable Process*

# Sample and Explore

- Data selection
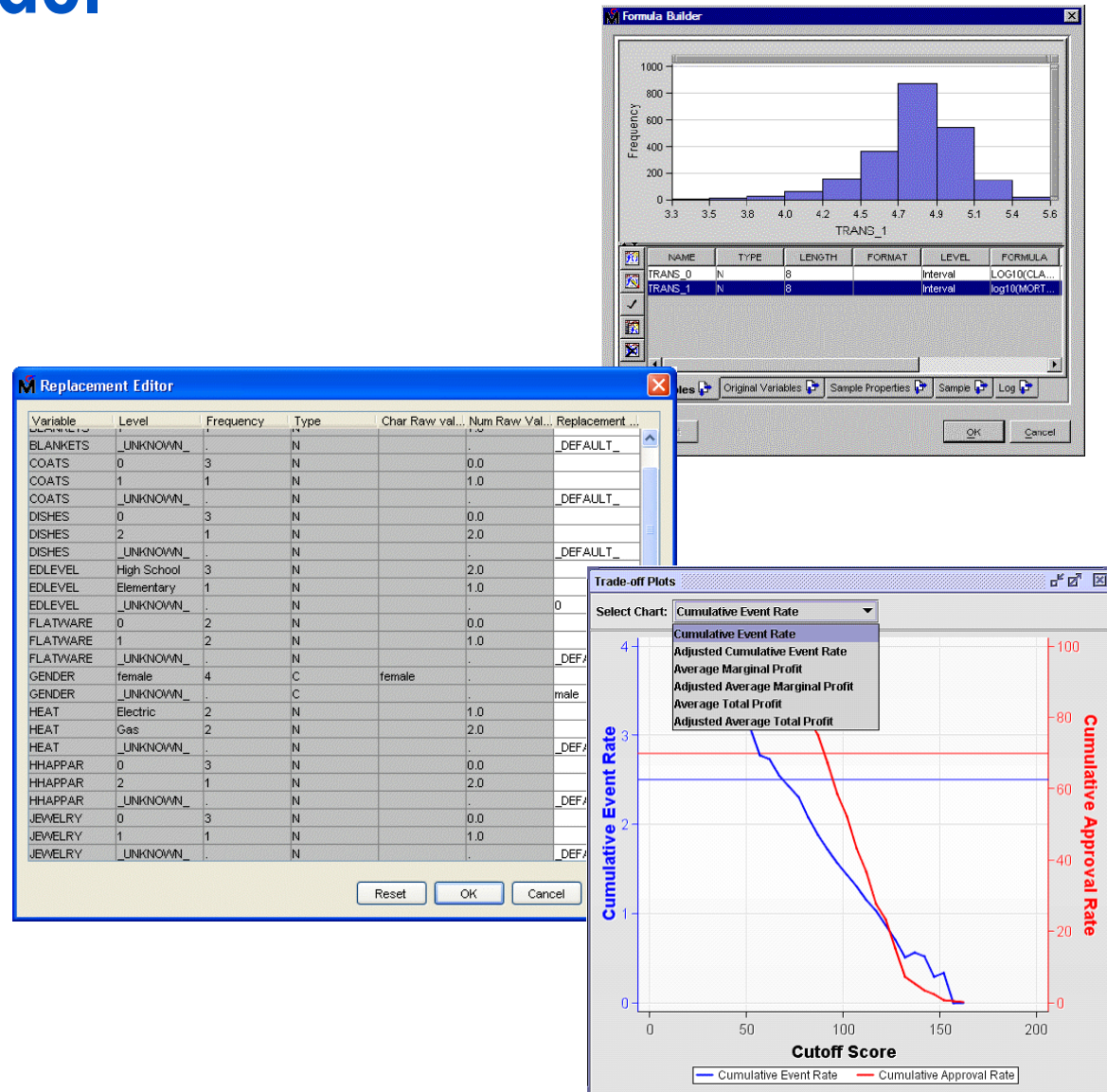  - Required & excluded fields
  - Sample balancing
  - Data partitioning

- Data evaluation
  - Statistical measures
  - Visualization
  - Identifying outliers
  - Analytical segmentation
  - Variable creation & selection

# Modify and Model

- Binning
- Scaling
- Imputation
- Replacement/recoding
- Modeling Policies
  - Prediction functions
  - Classification functions
- Modeling Methods

# Assess

- Compare training performance

- Champion / Challenger
  - Training and monitoring

- Ensure Generalization
  - Prevent over fitting

- Estimate deployment performance
  - Acquire target measures

- Select final model

# Survival Data Mining Analysis

- Prediction when event will happen, not just if it will happen

- Very popular for customer behavior modeling, such as
  - When will customer churn
  - When will customer upgrade

- Predicts event probability for time intervals for each customer
  - i.e. Customer has 50% chance to cancel next month but 75% chance to cancel the month after

- Can take external factors into account
  - customers with more than 2 products tend to stay longer

§sas | THE POWER TO KNOW.

# Survival Data Mining Analysis
## *Approach*

- Look at probability of hazard (Event) at discrete time points

- Time to event is most important feature

- Covariates can be integrated – need to be categorized

  - Customer demographics – age group

  - Product details – bronze, silver, gold tariff

  - Usage history – high, normal, low usage etc.

# Survival Data Mining Analysis Node
## *Reports*

- Discrete time to event regression with additive logistic regression.

- Data preparation to define time interval and time training range (tenure view)

- Time effect is modeled with cubic splines to allow for flexible shapes of the hazard functions.

- Proportional hazard function is fitted with constant covariates.

# Survival Data Mining Node
## *Scoring*

- Mean Residual Lifetime: expected time till event occurs based on projecting hazard function into the future
  - Projection based on constant hazard function
  - Projection based on continuing trend of hazard function

# Time Series Data Mining Nodes (experimental )

- **Integrate time dimension into analysis**

- **Data is often stored as transactional data with time stamp or in form of time series**

- **Nodes in SAS Enterprise Miner 7.1**

  - Data Preparation

    » Provides a tool of aggregation, differencing, summarization, etc.

  - Exponential Smoothing

    » Fits ESM to interval variables

  - Similarity

    » Computes several similarity measures among time series