





Data Mining Metoder och tekniker

Mathias Lanner Sas Institute





What is data mining



Sas. HERE

"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 [is] the process of efficient discovery of nonobvious valuble information from large collection of data."

-Berson and Smith (1997), p.565

"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 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





Data Mining Is:

- Discovering patterns and relationships represented in data.
 - Developing models to understand and describe characteristics and activity based on these patterns.
- Using this understanding to help evaluate future options, gain insights and take decisions.
 - Deploy the results of the analysis to affect business change.





Two types of analysis



Pattern Discovery





Predictive Modeling



Pattern Discovery



The Essence of Data Mining?

"...the discovery of interesting, unexpected, or valuable structures in large data sets."

David Hand



Pattern Discovery Applications Data reduction



Novelty detection



Clustering



Market basket analysis



Sequence analysis



Predictive Modeling



The Essence of Data Mining

"Most of the big payoff [in data mining] has been in predictive modeling."

- Herb Edelstein



Data Mining Delivering Value Across Industries

Retail	Identify the most profitable customers and the underlying reasons for their loyalty
Life Sciences	Find promising new molecular drug compounds
Manufacturing	Save on downtime by applying predictive maintenance to manufacturing sites
Financial Services	Grow customer profitability and reduce risk exposure through more accurate credit scoring
Insurance	Rate making - set more profitable rates for insurance premiums
Telco	Help retain customers (prevent churn) and identify up-sell/cross- sell opportunities for individual customers
Government	Detect and deter fraudulent behavior





Process -methodology



Successful Data Mining – Iterative & Interactive



S.Sas. HERE



Explore data and find relationships



SSAS HERE

Modify variables/data





SSAS HERE



Modelling

SEMMA

Prediktive modelling

Decision tree



Regression



Neuralnetworks



Exploring techniques

Cluster analysis



Copyright © 2006, SAS Institute Inc. All rights reserved.

Basket analysis



Univariate





Assess –







Statistics & Data mining

Statistics

- Experimental
- Prior Hypothesis
 - Idea before data acquisition
 - Data acquisition planned
- Experimental Design
 - Sampling strategies
 - Factorial designs
 - Required confidence
 - Minimize model terms
- Inference
 - Hypothesis testing
 - Inference

Data mining

- Commercial
- Posterior Hypothesis
 - Idea after data acquisition
 - Data acquisition
 opportunistic
- No Experimental Design
 - Explore data
 - Create hypothesis
 - Generate query
 - Create models
- Prediction
 - Lift, Profit, Response
 - Prediction





Data Challenges





Data is collected differently

	Experimental	Opportunistic
Purpose	Research	Operational
Value	Scientific	Commercial
Generation	Actively controlled	Passively observed
Size	Small?	Massive
		(large N and p)
Hygiene	Clean	Dirty
State	Static	Dynamic

Where does mining data come from ?



Sas. HONOR

Data is becoming wider and wider

- Used to work with a couple of dozens of variables
- Nowadays at least a couple of hundreds
 - Data from different sources
 - Derived data (differences, rations, trends etc.)
 - Data from combined algorithms (market basket analysis, combined with clustering combined with predictive modeling)
- Can become thousands
 - Pharma: micro-array data
 - Interactions

as Here



New data sources

Extreme commercial data warehouses

- Many gigabytes of data
- Stores may have 100,000+ SKU items
- Sales histories for every item/basket saved
- Rollups can produce terms >> 10,000 terms
- Digital data acquisition
 - Biometrics: microarray, mass spectrometry
 - Chip fabs: 30,000 measurements per manufacturing run.
 - ISP: every page, server, router, switch, at timepoints
 - University: 50-60 GB / day
 - Regional telecom: 6 TB / day



Integration

- Integrate data access and management
 - Prepare data for analytics in enterprise warehouse
 - Join tables
 - Clean data
 - Create derived variables (aggregations, ratios, trends etc.)
 - Create samples
 - Create data mining metadata (targets, inputs, rejected)



Predictive Model Development Data



Marketing Data Warehouse



Data Quality

Data quality

- Data mining requires detail data
- New level of data quality is necessary
- Lot of time spent for data cleaning
- Use the warehouse to correct the errors

• META Group:

"10 to 20 percent of the raw data used is corrupt or incomplete in some way. It is not unusual to discover that as many as half the records in a database contain some type of information that needs to be corrected" META Group Program Director John Ladley



Data Quality

- Intelligent methods to deal with missing values
 - Use robust estimators for distribution
 - Predict missing values from remaining information with trees
 - Track the replaced values add degrees of freedom for missingness
 - Use clustering for replacing missing values
 - Use algorithms that can deal with missing values automatically





Techniques



Definition: What is a model?

• A function assigning an outcome

- Prediction, classification, cluster, rule,...
- input vars, output vars
- SCORECODE

Attributes of the function

- Type, formula, objective, target
- Attribute importance, ranking, p-value
- Model quality measures: stats, gains, distributions, centroids

History of the function

- Who, what, when, how
- Auditing, regulation, best practices
- Purpose of the function
 - Cross-sell, acquisition, approval, collection, churn, fraud,...
- Status of the model
 - Development, production, retired, champion, challenger



Algorithms



- A model is an abstraction of the data and belongs with the data
- There is nothing more in a model than what is already in the data



Algorithms

- There is no BEST algorithm per se
- Depends on
 - Nature of relationships in data
 - Data quality
 - Time available to build a model
 - Nature of model deployment
 - operational use
 - insights for business users
 - decision support etc.



Data Mining Algorithms

① predictive (supervised)

use data on past processes to *predict* future production



② descriptive (unsupervised) use data on past processes to <u>describe</u> current situation





Supervised Learning

Tries to find good rules for predicting the value of a target(s) from the values of the inputs variables.



Enterprise Miner

- Logistic and OLS
- Tree Classifiers
- Neural Networks
- Ensembles
- Memory Based Reasoning
- Two-stage modeling
- Fast Variable Selection
- Principal Components
- •PLS Regression
- •Support Vector Machine
- •Gradient Bosting
- SAS/STAT



Data Mining Model Training and Scoring





Unsupervised Learning

Tries to divide the data into groups such that the observations within a group have traits more similar than those assigned to different groups Enterprise Miner







- k-Means
- SOM/Kohonen Networks
- •Rule Builder
- SAS/STAT





Unsupervised Classification

Training Data

case 1: inputs, ?
case 2: inputs, ?
case 3: inputs, ?
case 4: inputs, ?
case 5: inputs, ?

Training Data

case 1: inputs, cluster 1
case 2: inputs, cluster 3
case 3: inputs, cluster 2
case 4: inputs, cluster 1
case 5: inputs, cluster 2



4



Market Basket Analysis



<u>Rule</u>	<u>Support</u>	<u>Confidence</u>
$A \Longrightarrow D$	2/5	2/3
$C \Rightarrow A$	2/5	2/4
$A \Rightarrow C$	2/5	2/3
$B \And C \Rightarrow D$	1/5	1/3





10,000

Support(SVG \Rightarrow CK) = 50% Confidence(SVG \Rightarrow CK) = 83% Expected Confidence(SVG \Rightarrow CK) = 85% Lift(SVG \Rightarrow CK) = 0.83/0.85 < 1

