Data mining Theory and examples

Mathias Lanner, SAS Institute Sverige





THE POWER TO KNOW.

Agenda

- Introduction to data mining
- Process methodology
- Data Challenges
- Modelling Techniques



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 [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 Applications



Data reduction



Novelty detection



Clustering



Market basket analysis



Sequence analysis



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?

8

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 <i>(Retail)</i>	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. <i>(Health</i> <i>Insurance)</i>	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?



9

POWER TO KNOW

Successful Analytics – Iterative and Interactive



Sas He POWER TO KNOW

10

Copyright © 2010, SAS Institute Inc. All rights reserved.

Successful Data Mining – Iterative & Interactive



Where does mining data come from ?



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



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
- 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
 - Social Media



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



Sas. THE POWER TO KNOW.

16

Copyright © 2010, SAS Institute Inc. All rights reserved.

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



Predictive modeling and scoring



Copyright © 2010, SAS Institute Inc. All rights reserved.

POWER TO KNOW

Algorithms

- There is no BEST algorithm
- 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.



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

Predictive Modeling Essentials







Optimize complexity



21

Copyright © 2010, SAS Institute Inc. All rights reserved.

Curse of Dimensionality



22

THE POWER TO KNOW。

SSas







POWER TO KNOW.





THE POWER TO KNOW。

Data Splitting







The Model Comparison tool provides



Summary statistics



Statistical graphics



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



Market Basket Analysis



<u>Rule</u>	<u>Support</u>	<u>Confidence</u>
$A \Rightarrow D$	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

Copyright © 2010, SAS Institute Inc. All rights reserved.

THE POWER TO KNOW。

S.Sa

Implication?



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



Copyright © 2010 SAS Institute Inc. All rights reserved