

Knowledge Discovery with Data and Text Mining



Dursun Delen, Ph.D.

Associate Professor of MSIS
William S. Spears School of Business
Oklahoma State University - Tulsa

Introduction

- Why data mining?
- What is data mining?
- Data Mining Process
- Data Mining Techniques
- Text Mining
- Data & Text Mining Examples
- A Demonstration of Data Mining in Action

2 © D. Delen KPM Symposium – Tulsa, OK – October 3-4, 2007



Data Deluge

Why Now?

hospital patient registries
electronic point-of-sale data
stock trades OLTP telephone calls
catalog orders bank transactions tax returns
remote sensing images credit card charges
airline reservations

Barcodes
RFIDs
Magnetic Strips
Sensors...



Motivation:
"Necessity is the Mother of Invention"

- Competitive pressure to make better decisions
- Data explosion problem (or opportunity)
 - Why do we have more data now, then we had before?
 - Technology driven reasons: Automated data collection tools and techniques...
 - Software driven reasons: Mature database technology...
 - Cost driven reasons (both hardware and software)

We are drowning in data, but starving for knowledge!

Solution: Data and data mining

4 © D. Delen KPM Symposium – Tulsa, OK – October 3-4, 2007 OSU

What Is Data Mining?

- Data mining :
 - Extraction of interesting (non-trivial, implicit, previously unknown and potentially useful) information (or patterns) from data
- Alternative names...
 - Data mining: a misnomer?
 - Knowledge discovery, knowledge extraction, data/pattern analysis, data archeology, data dredging, information harvesting, business intelligence, etc.

5 © D. Delen KPM Symposium – Tulsa, OK – October 3-4, 2007 OSU

Standardized DM Processes

- CRISP-DM
 - Cross-Industry Standard Process for Data Mining
 - www.crisp-dm.org
- SEMMA
 - ✓ Sample the data
 - ✓ Explore the data
 - ✓ Modify the data
 - ✓ Model
 - ✓ Assess
 - <http://www.sas.com/technologies/analytics/datamining/miner/semma.html>

6 © D. Delen KPM Symposium – Tulsa, OK – October 3-4, 2007 OSU

Standardized Data Mining Processes

Step 1: Business Understanding

- Determine the business objectives
- Assess the situation
- Determine the data mining goals
- Produce a project plan

Cross-Industry Standard Process for Data Mining CRISP-DM

7 © D. Delen KPM Symposium – Tulsa, OK – October 3-4, 2007 **OSU**

Standardized Data Mining Processes

Step 2: Data Understanding

- Collect the initial data
- Describe the data
- Explore the data
- Verify the data

Cross-Industry Standard Process for Data Mining CRISP-DM

8 © D. Delen KPM Symposium – Tulsa, OK – October 3-4, 2007 **OSU**

Standardized Data Mining Processes

Step 3: Data Preparation

- Select data
- Clean data
- Construct data
- Integrate data
- Format data

Cross-Industry Standard Process for Data Mining CRISP-DM

9 © D. Delen KPM Symposium – Tulsa, OK – October 3-4, 2007 **OSU**

Standardized Data Mining Processes

Step 4: Modeling

- Select the modeling technique
- Generate test design
- Build the model
- Assess the model

Cross-Industry Standard Process for Data Mining **CRISP-DM**

10 © D. Delen KPM Symposium – Tulsa, OK – October 3-4, 2007 **OSU**

Standardized Data Mining Processes

Step 5: Evaluation

- Evaluate results
- Review process
- Determine next step

Cross-Industry Standard Process for Data Mining **CRISP-DM**

11 © D. Delen KPM Symposium – Tulsa, OK – October 3-4, 2007 **OSU**

Standardized Data Mining Processes

Step 6: Deployment

- Plan deployment
- Plan monitoring and maintenance
- Produce final report
- Review the project

Cross-Industry Standard Process for Data Mining **CRISP-DM**

12 © D. Delen KPM Symposium – Tulsa, OK – October 3-4, 2007 **OSU**



Data Mining Techniques (1)

- **Association** (correlation and causality)
 - $\text{age}(X, "20..29") \wedge \text{income}(X, "20..29K") \rightarrow \text{buys}(X, "PC")$
[support = 2%, confidence = 60%]
 - $\text{contains}(T, "computer") \rightarrow \text{contains}(x, "software")$
[support = 1%, confidence = 75%]
- **Prediction** (Classification & Regression)
 - **Classification**: developing models (functions) that describe and distinguish classes or concepts for future prediction
 - Predicting whether it will rain tomorrow
 - Classifying loan applicant as "good" or "bad"



Data Mining Techniques (2)

- **Prediction** (Classification & Regression)
 - **Regression**: developing models (functions) that forecast the value of a continuous numerical variable
 - Predicting what the temperature will be tomorrow
 - Forecasting the future value of a stock a year from now
- **Cluster analysis**
 - Class label is unknown: Group the samples of objects based on their quantifiable characteristics
 - Cluster the customers who share the same characteristics: interest, income level, spending habits, etc.
 - Clustering is done by maximizing the intra-class similarity and minimizing the interclass similarity

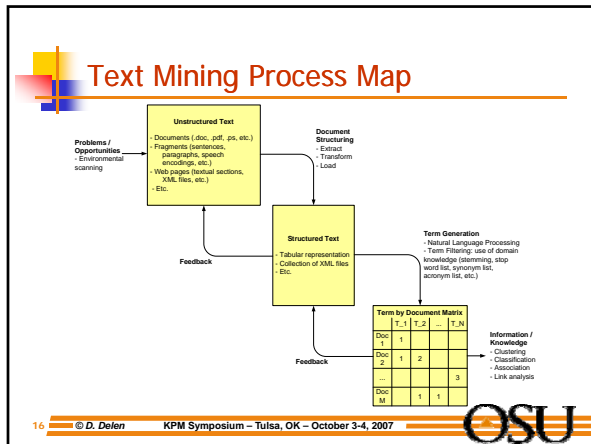


Text Mining

Statistical
Natural
Language
Processing

Data
Mining

- Text mining is a process that employs
 - a set of algorithms for converting unstructured text into structured data objects, and
 - the quantitative methods that analyze these data objects to discover knowledge
- Text Mining = Statistical NLP + DM



Data Mining Applications... Military Health System (1/3)

- KBSI's Phase II SBIR research project
 - Funded through SBIR program by the Offices of the Secretary of Defense
 - SBIR: Phase I ⇒ Phase II ⇒ Phase III

- DM in Healthcare
 - Managerial
 - Clinical

Reference

- Delen, D. and S. Ramachandran (2003). A Hybrid Approach to Knowledge Discovery from Military Health Systems. *Journal of Neural, Parallel and Scientific Computing*, Volume 11, Number 1&2, pp. 161-183.

17 © D. Delen KPM Symposium – Tulsa, OK – October 3-4, 2007 OSU

Data Mining Applications... Military Health System (2/3)

- One of the largest health systems in the US
 - 90+ hospitals
 - 100s of outpatient clinics, treatment facilities
 - 180,000 employees (doctors, nurses, other staff)
 - 8 million beneficiaries
 - \$20 Billion/year budget
- Purpose/mission is to
 - Provide healthcare to eligible veterans
 - Provide education and training opportunities to health profession (residency, practical training, etc.)
 - Conduct medical research, create innovation
 - Provide public health service at the time of natural disasters
- Problem: ↑ demand, ↓ budget

18 © D. Delen KPM Symposium – Tulsa, OK – October 3-4, 2007 OSU



Forecasting Box Office Success of Hollywood Movies

Ramesh Sharda and Dursun Delen

Institute for Research in Information Systems
Oklahoma State University



Forecasting Box-Office Receipts: A Tough Problem!

“... No one can tell you how a movie is going to do in the marketplace... not until the film opens in darkened theatre and sparks fly up between the screen and the audience”

Mr. Jack Valenti
*President and CEO
of the Motion Picture Association of America*



Current work on prediction...

- A lot of research have been done
 - Behavioral models
 - Analytical models
 - 🕒 Predict after the initial release
- Our approach
 - Use a data mining approach
 - Use as much historical data as possible
 - Make it web-enabled
 - 🕒 Predict before the initial release