

IT in Finance

MINING OF FINANCIAL DATABASES

INTRODUCTION

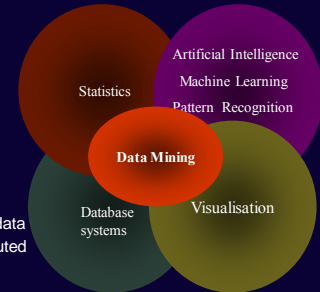
Jerzy KORCZAK

email: jerzy.korczak@ue.wroc.pl
<http://www.korczak-leliwa.pl>

1

Origins of Data Mining

- Draws ideas from AI/machine learning, pattern recognition, statistics, and database systems



- Traditional Techniques may be unsuitable due to
 - Enormity of data
 - High dimensionality of data
 - Heterogeneous, distributed nature of data

2

Data Mining

- Data mining (knowledge discovery in databases, KDD)
 - Extraction of interesting, non-trivial, implicit, previously unknown and potentially useful information (knowledge) or patterns from data in large databases or other information repositories
- Scientific point of view: data abstraction and KDD
- Commercial point of view: competitive pressure
- Necessity is the mother of invention
 - *Data is everywhere — data mining should be everywhere, too!*
 - Understand and use data — an imminent task!

3

Statistics vs Data Mining

- Statistics: a discipline dedicated to data analysis
- What are the differences?
 - Huge amount of data—in Giga to Tera bytes
 - Fast computer—quick response, interactive analysis
 - Multi-dimensional, powerful, thorough analysis
 - High-level, “declarative”—user’s ease and control
 - Automated or semi-automated—mining functions hidden or built-in in many systems

4

Types of Decision-Support Systems (DSS)

Model-driven DSS:

- Primarily stand-alone systems
- Use a strong theory or model to perform “what-if” analyses

Data-driven DSS:

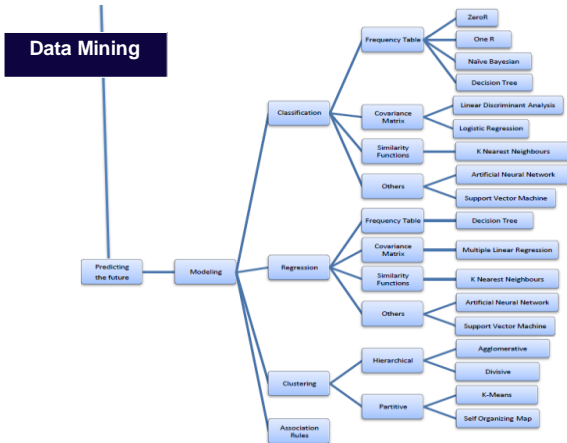
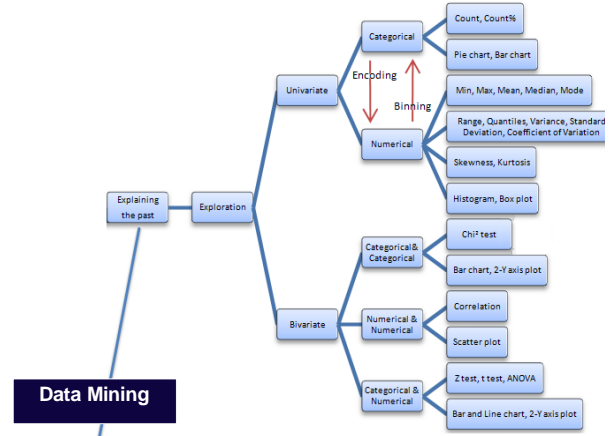
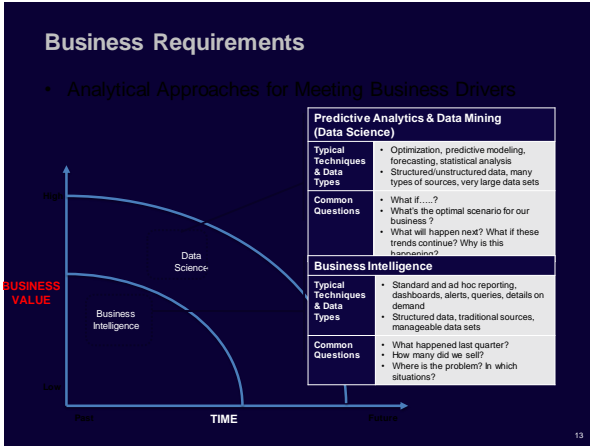
- Integrated with large pools of data in major enterprise systems and Web sites
- Support decision making by enabling user to extract useful information
- Data mining: can obtain types of information such as associations, sequences, classifications, clusters, and forecasts

5

Data Sets, Database, Images

- Relational database — A commodity of every enterprise
- Huge data warehouses are under construction
- POS (Point of Sales): Transactional DBs in terabytes
- Object-relational databases, distributed, heterogeneous, and legacy databases
- Spatial databases (GIS), remote sensing database (EOS), and scientific/engineering databases
- Time-series data (e.g., stock trading) and temporal data
- Text (documents, emails) and multimedia databases
- WWW: A huge, hyper-linked, dynamic, global information system

6



- ## Research Progress
- Multi-dimensional data analysis: Data Warehouse and OLAP
 - Association, correlation, and causality analysis
 - Classification: scalability and new approaches
 - Clustering and outlier analysis
 - Sequential patterns and time-series analysis
 - Similarity analysis: curves, trends, images, texts, etc.
 - Text mining, Web mining and Weblog analysis
 - Social networks, link analysis
 - Spatial, multimedia, scientific data analysis
 - Smart sensors: IoT
 - Image classification and interpretation
 - Data preprocessing and database compression
 - Data visualization and visual data mining
 - Many others, e.g., collaborative filtering
- 16