

①

1. Introduction

2. Data

3. Exploring Data

4. Classification

5. Association analysis

6. Cluster analysis

7. Anomaly Detection

8. Applications before research topics

text mining

graph mining

web mining

1. Introduction to Data Mining

Pang Ning Tan / Vipin Kumar

Michael Steinbach

2. Data mining Concepts and Techniques

Jianwei Han / Micheline Kamber / Jian Pei

Introduction

Rapid advances in data collection and storage technology have enabled organizations to accumulate vast amount of data.

However, extracting useful information has proven extremely challenging.

Often, traditional data analysis tools and techniques cannot be used because of the massive size of a data set.

Sometimes, the non-traditional nature of the data, traditional approaches cannot be applied even if the data set is relatively small.

Data mining is a technology that blends traditional data analysis methods with sophisticated algorithms for processing large volumes of data.

It has also opened up exciting opportunities for exploring and analyzing new types of data and for analyzing old types of data in new ways.

Business applications

Retailers are allowed to collect up to the minute date about customer purchases at the checkout counters of their stores.

web logs from e-commerce.

to help ~~customers~~ better understand the needs of their customers and make more informed business decisions.

Data mining techniques can be used to support a wide range of business intelligence applications

such as

- customer profiling
- targeted marketing
- work-flow management
- store layout,
- fraud detection

Help retailers:

- who are the most profitable customers?
- what products can be cross sold?
- what is the revenue outlook of the company for next year?

Medicine, Science and Engineering

- An important step towards improving our understanding of the Earth's Climate System.
- Data mining can and
 - what is the relationship between the frequency and intensity of ecosystems disturbances such as droughts and hurricanes to global warming?
 - How is land surface precipitation and temperature affected by ocean surface temperature?
 - How well can we predict the beginning and end of the rainy season for a region?
- Researchers in molecular biology hope to use the large amounts of genomic data currently being gathered to better understand the structure and function of genes.
- In addition, data mining can also be used to address other important biological challenges such as protein structure prediction, multiple sequence alignment, the modeling of biochemical pathways, and phylogenetics.

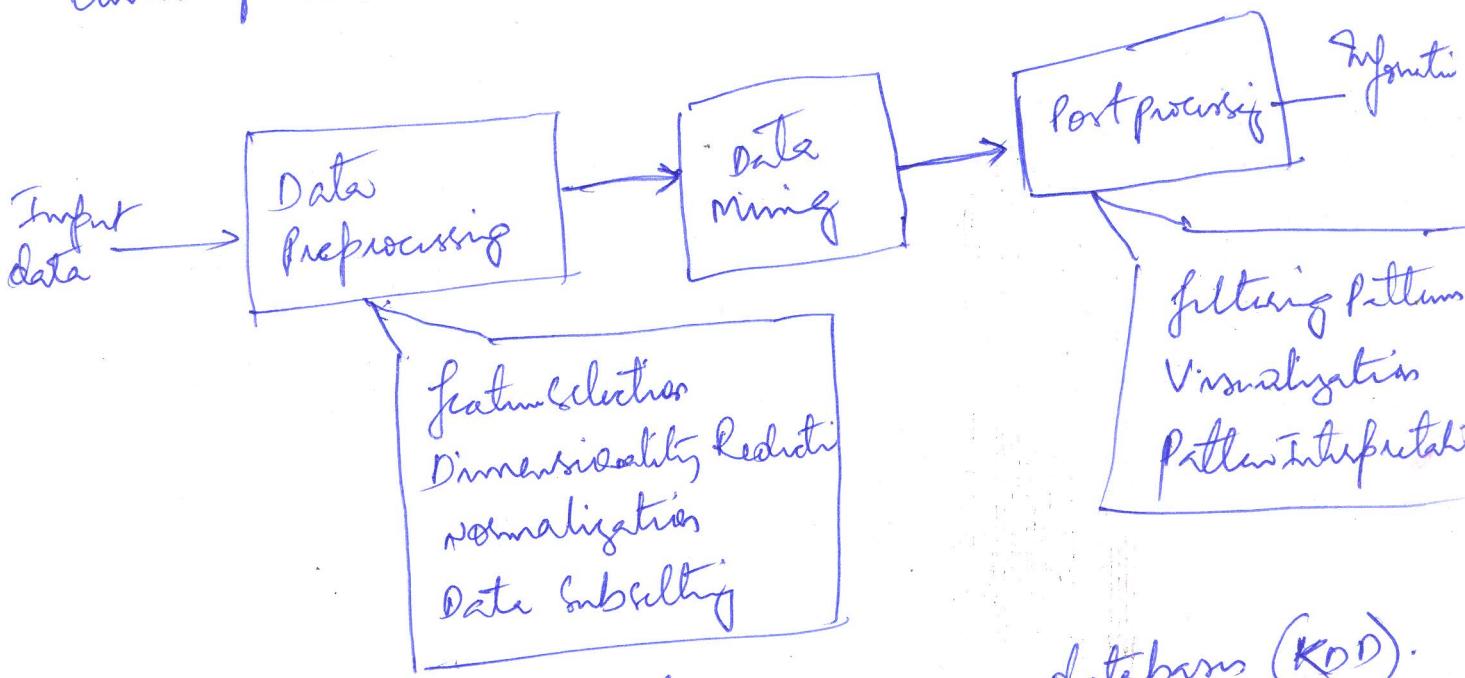
What is Data Mining?

DM (4)

Data mining is the process of automatically discovering useful information in large data repositories.

Data mining and Knowledge Discovery

Data mining is an integral part of knowledge discovery in databases (KDD), which is the overall process of converting raw data into useful information.



The process of knowledge discovery in databases (KDD).

The input data - Variety of formats

- flat files

- spreadsheets

- relational tables

Central repository
or
distributed across
multiple sites.

DM-5

The purpose ~~for~~ of Preprocessing is to transform the raw input data into an appropriate format for subsequent analysis.

Motivating Challenges:

1. Scalability:

Because of advances in data generation and collection, data sets with sizes of gigabytes, terabytes, or even petabytes are becoming common.

If data mining algorithms are to handle these massive data sets, then they must be scalable.

2. High Dimensionality:

~~Thousands~~ thousands ~~handful~~ instead of the handful

- thousands of attributes common a few decades ago.

- in bioinformatics, progress in microarray technology has produced gene expression data involving thousand of features.

- temporal data set

- spatial data set

Heterogeneous and Complex data :

Heterogeneous attributes :

- webpages containing semi-structured text and hyperlinks;
- DNA data with sequences and their structural structure ;

Data ownership and Distribution :

The data is geographically distributed among resources belonging to multiple entities .

The origins of Data mining :

by the goal of meeting the challenges ,
 researchers from different disciplines began to focus
 on developing more efficient and scalable tools that
 could handle diverse types of data .
 This work , which culminated in the field of data mining
 methodology ~~at~~
 In particular , data mining draws upon ideas , such
 as ① sampling , estimation , and hypothesis testing
 from statistics . and .
 ② search algorithm , modelling techniques .

learning theories from artificial intelligence, pattern recognition, and machine learning.

Data mining has also been quick to adopt ideas from other areas, including optimization, evolutionary computing, information theory, signal processing, visualization and information retrieval.

Data mining as a confluence of many disciplines



- Database systems are needed to provide support for efficient storage, indexing and query processing.
- The techniques from high performance (Parallel) computing are often important in addressing the massive size of some of date sets.
- Distributed techniques can also help address the issue of size and are essential when the data cannot be gathered in one location.

Date mining tasks

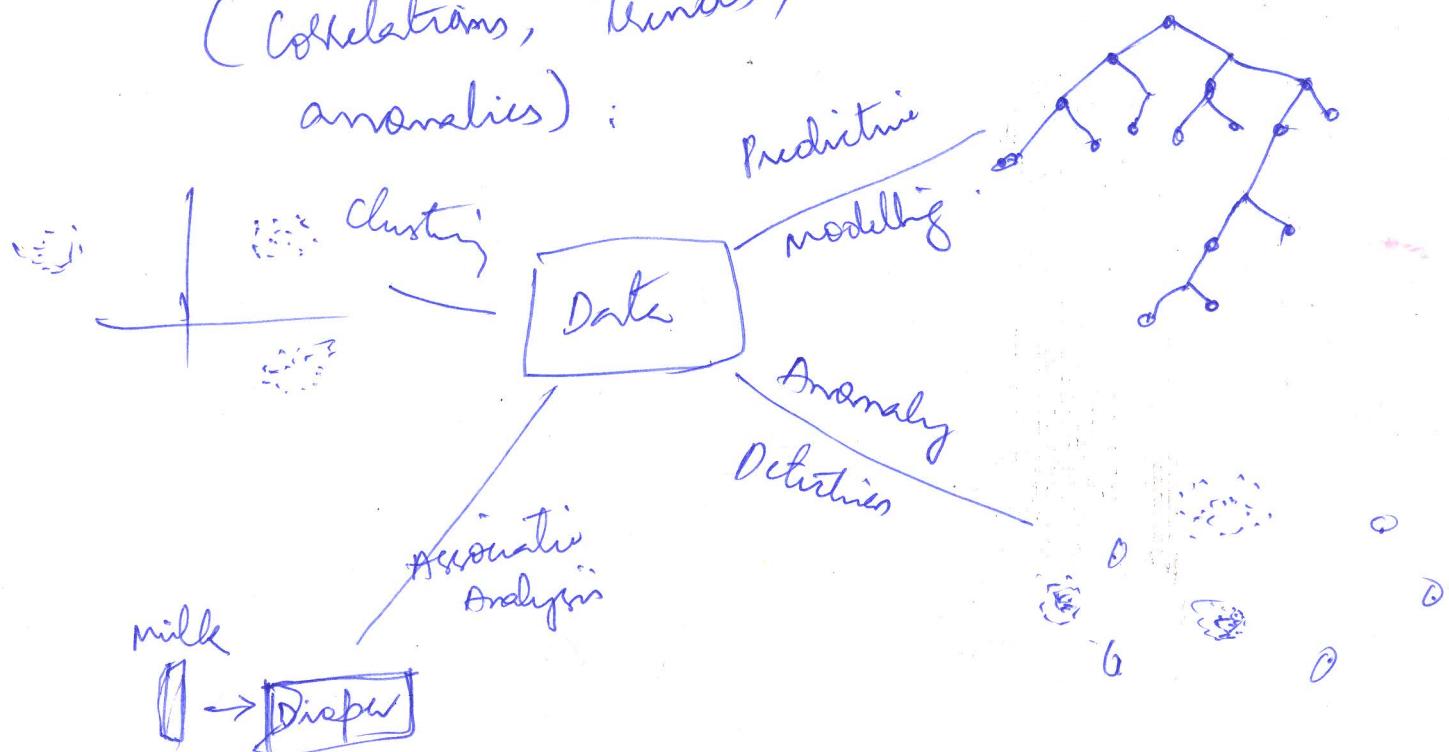
two major categories

Predictive Tasks

Predict the value of a particular ~~task~~ attribute based on the values of other attributes.

Descriptive Tasks

The objective is to derive patterns that summarize the underlying relationship in data (correlations, trends, clusters, trajectories and anomalies).



Predictive modeling

Refers to the task of building a model for the target variable as a function of the explanatory variables.

- Classification : which is used for discrete target variables.
- Regression : which is used for continuous target variables.

For example:

- Predicting whether a web user will make a purchase at an online bookstore. (Classification)
- predicting the future price of a stock in a regression task. (because price is a continuous-valued attribute).

Association analysis

Find the discernible patterns that describe strongly associated features in the data.

Eg: → identifying web pages that are accessed together.
 → ~~sold~~ to find items that are frequently bought together.

Cluster analysis

DM - 10

Seeks to find groups of closely related observations to those

Eg:- find areas of the ocean that have a significant impact on the earth's climate.

- Document clustering (collection of news articles)

Anomaly detection:

is the task of identifying observations whose characteristics are significantly different from the rest of the data.

A good anomaly detector must have a high detection rate and a low false alarm rate.

Applications

- fraud detection
- network intrusions
- unusual patterns of disease
- ecosystem disturbances.

Exercises

1. Discuss whether or not each of the following activities is a data mining task.
- Dividing the customers of a company according to their gender: **NO** database query
 - Dividing the customer of a company according to their profitability: **NO**
 - Computing the total sales of a company: **NO**
 - Sorting a student database based on student identification numbers: **NO**
 - Predicting the outcome of tossing a (fair) pair of dice
 - Predicting the future stock price of a company using historical records: **YES**
 - Monitoring the heart rate of a patient for abnormalities: **YES**
 - Monitoring seismic waves for earthquake activities: **YES**
 - Monitoring seismic waves for early warning: **YES**
 - Extracting the frequency of a sound wave: **NO** signal processing