# DATA MINING LECTURE 1

Introduction

Heba AL.Marwai

# What is data mining?

- After years of data mining there is still no unique answer to this question.

- A tentative definition:

Data mining is the use of efficient techniques for the analysis of very large collections of data and the extraction of useful and possibly unexpected patterns in data.
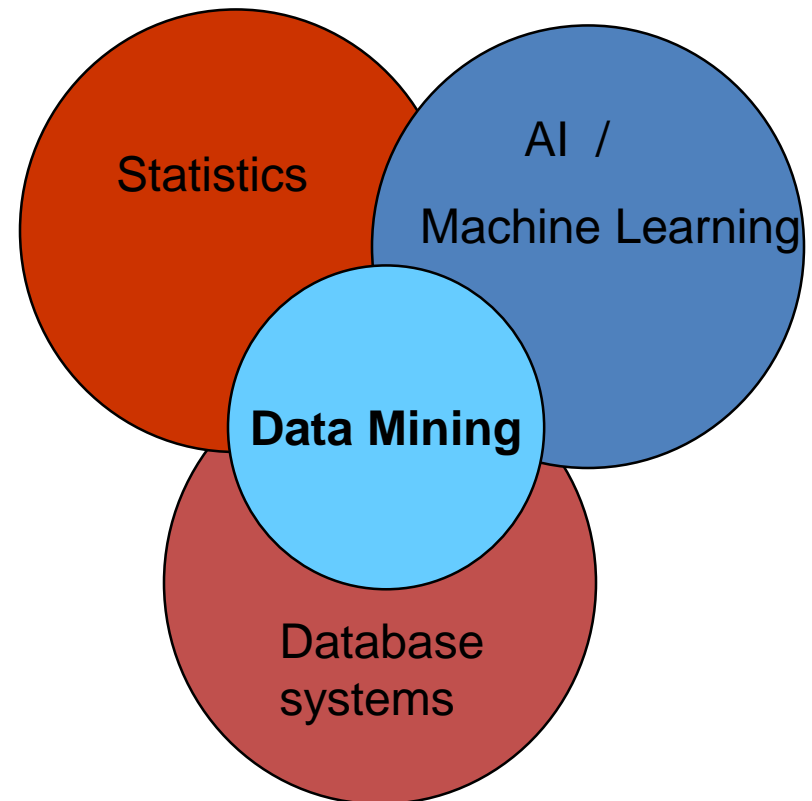
# What is Data Mining?

- Data Mining is:

  (1) The efficient discovery of previously unknown, valid, potentially useful, understandable patterns in large datasets

  (2) The analysis of (often large) observational data sets to find unsuspected relationships and to summarize the data in novel ways that are both understandable and useful to the data owner

# Overview of terms

- <u>Data:</u> a set of facts (items) D, usually stored in a database

- <u>Pattern:</u> an expression E in a language L, that describes a subset of facts

- <u>Attribute:</u> a field in an item i in D.

- <u>Interestingness:</u> a function ID,L that maps an expression E in L into a measure space M

# Origins of Data Mining

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

- Must address:
  - Enormity of data
  - High dimensionality of data
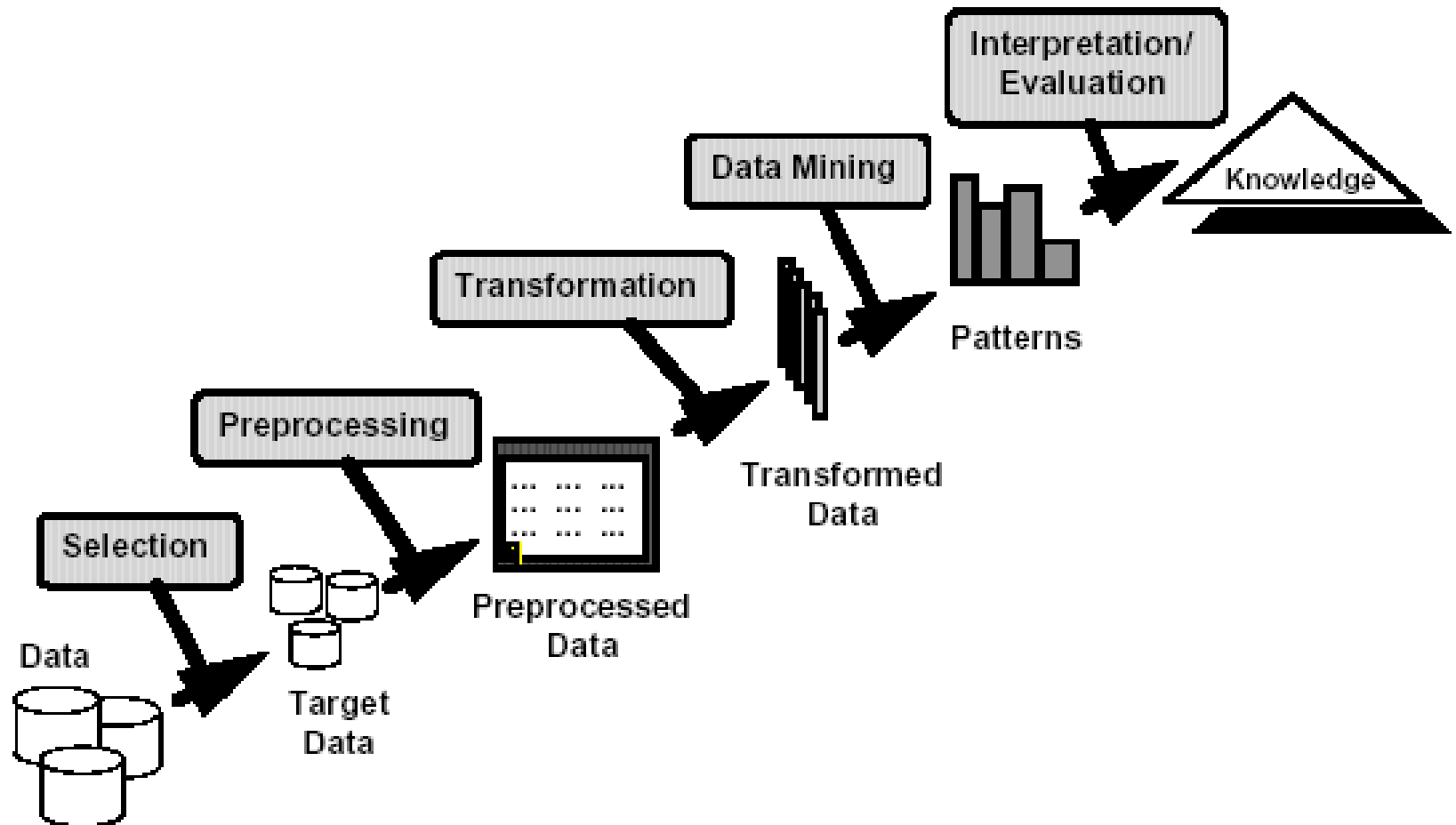  - Heterogeneous, distributed nature of data

# Why do we need data mining?

- Really, really huge amounts of raw data!!
- In the digital age, TB of data is generated by the second
- Mobile devices, digital photographs, web documents.
- Facebook updates, Tweets, Blogs, User-generated content
- Transactions, sensor data.
- Queries, clicks, browsing
- Cheap storage has made possible to maintain this data
- Need to analyze the raw data to extract knowledge

# Why do we need data mining?

- "The data is the computer"
  - Large amounts of data can be more powerful than complex algorithms and models
    - Google has solved many Natural Language Processing problems, simply by looking at the data
    - Example: misspellings, synonyms
  - Data is power!
    - Today, the collected data is one of the biggest assets of an online company
    - Query logs of Google
    - The friendship and updates of Facebook
    - Tweets and follows of Twitter
  - We need a way to harness the collective intelligence

# Knowledge Discovery

# The Data Mining Process

1. Understand the domain

2. Create a dataset:

   - Select the interesting attributes
   - Data cleaning and preprocessing

3. Choose the data mining task and the specific algorithm

4. Interpret the results, and possibly return to 2

# Data Mining Tasks

1. Classification: learning a function that maps an item into one of a set of predefined classes

2. Regression: learning a function that maps an item to a real value

3. Clustering: identify a set of groups of similar items

# Data Mining Tasks

4. Dependencies and associations:

   identify significant dependencies between data attributes

5. Summarization: find a compact description of the dataset or a subset of the dataset

# The data is also very complex

- Multiple types of data: tables, time series, images, graphs, etc

- Interconnected data of different types:
  - From the mobile phone we can collect, location of the user, friendship information, opinions through twitter, images though cameras, queries to search engines

# So, what is Data?

- Collection of data objects and their attributes

- An attribute is a property or characteristic of an object
  - Examples: eye color of a person, temperature, etc.
  - Attribute is also known as variable, field, characteristic, or feature
- A collection of attributes describe an object
  - Object is also known as record, point, case, sample, entity, or instance

Attributes

| Tid | Refund | Marital Status | Taxable Income | Cheat |
|-----|--------|----------------|----------------|-------|
| 1 | Yes | Single | 125K | No |
| 2 | No | Married | 100K | No |
| 3 | No | Single | 70K | No |
| 4 | Yes | Married | 120K | No |
| 5 | No | Divorced | 95K | Yes |
| 6 | No | Married | 60K | No |
| 7 | Yes | Divorced | 220K | No |
| 8 | No | Single | 85K | Yes |
| 9 | No | Married | 75K | No |
| 10 | No | Single | 90K | Yes |

Objects

Size: Number of objects
Dimensionality: Number of attributes
Sparsity: Number of populated object-attribute pairs

# Types of Attributes

- There are different types of attributes
    - Categorical
        - Examples: eye color, zip codes, words, rankings (e.g, good, fair, bad), height in {tall, medium, short}
        - Nominal (no order or comparison) vs Ordinal (order )
    - Numeric
        - Examples: dates, temperature, time, length, value, count.
        - Discrete (counts) vs Continuous (temperature)
        - Special case: Binary attributes (yes/no, exists/not exists)

# Numeric Record Data

- If data objects have the same fixed set of numeric attributes, then the data objects can be thought of as points in a multi-dimensional space, where each dimension represents a distinct attribute

- Such data set can be represented by an n-by-d data matrix, where there are n rows, one for each object, and d columns, one for each attribute

| Projection of x Load | Projection of y load | Distance | Load | Thickness |
|---|---|---|---|---|
| 10.23 | 5.27 | 15.22 | 2.7 | 1.2 |
| 12.65 | 6.25 | 16.22 | 2.2 | 1.1 |

# Categorical Data

- Data that consists of a collection of records, each of which consists of a fixed set of categorical attributes

| Tid | Refund | Marital Status | Taxable Income | Cheat |
|-----|--------|----------------|----------------|-------|
| 1 | Yes | Single | High | No |
| 2 | No | Married | Medium | No |
| 3 | No | Single | Low | No |
| 4 | Yes | Married | High | No |
| 5 | No | Divorced | Medium | Yes |
| 6 | No | Married | Low | No |
| 7 | Yes | Divorced | High | No |
| 8 | No | Single | Medium | Yes |
| 9 | No | Married | Medium | No |
| 10 | No | Single | Medium | Yes |

# Document Data

- Each document becomes a `term` vector,
  - each term is a component (attribute) of the vector,
  - the value of each component is the number of times the corresponding term occurs in the document.
  - Bag-of-words representation – no ordering

| | team | coach | play | ball | score | game | win | lost | timeout | season |
|---|---|---|---|---|---|---|---|---|---|---|
| Document 1 | 3 | 0 | 5 | 0 | 2 | 6 | 0 | 2 | 0 | 2 |
| Document 2 | 0 | 7 | 0 | 2 | 1 | 0 | 0 | 3 | 0 | 0 |
| Document 3 | 0 | 1 | 0 | 0 | 1 | 2 | 2 | 0 | 3 | 0 |

# Transaction Data

- Each record (transaction) is a set of items.

| TID | Items |
|-----|-------|
| 1 | Bread, Coke, Milk |
| 2 | Beer, Bread |
| 3 | Beer, Coke, Diaper, Milk |
| 4 | Beer, Bread, Diaper, Milk |
| 5 | Coke, Diaper, Milk |

- A set of items can also be represented as a binary vector, where each attribute is an item.

- A document can also be represented as a set of words (no counts)

Sparsity: average number of products bought by a customer

# Ordered Data

- Genomic sequence data

    GGTTCCGCCTTCAGCCCCGCGCC
    CGCAGGGCCCGCCCCGCGCCGTC
    GAGAAGGGCCCGCCTGGCGGGCG
    GGGGGAGGCGGGGCCGCCCGAGC
    CCAACCGAGTCCGACCAGGTGCC
    CCCTCTGCTCGGCCTAGACCTGA
    GCTCATTAGGCGGCAGCGGACAG
    GCCAAGTAGAACACGCGAAGCGC
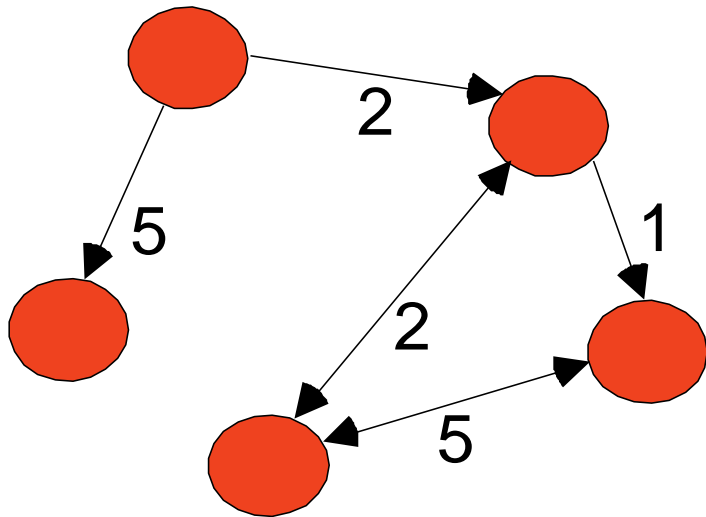    TGGGCTGCCTGCTGCGACCAGGG

- Data is a long ordered string

# Ordered Data

- Time series
  - Sequence of ordered (over "time") numeric values.

# Graph Data

- Examples: Web graph and HTML Links



<a href="papers/papers.html#bbbb">
Data Mining </a>
<li>
<a href="papers/papers.html#aaaa">
Graph Partitioning </a>
<li>
<a href="papers/papers.html#aaaa">
Parallel Solution of Sparse Linear System of Equations </a>
<li>
<a href="papers/papers.html#ffff">
N-Body Computation and Dense Linear System Solvers

# Types of data

- **Numeric data**: Each object is a point in a multidimensional space

- **Categorical data**: Each object is a vector of categorical values

- **Set data**: Each object is a set of values (with or without counts)
  - Sets can also be represented as binary vectors, or vectors of counts

- **Ordered sequences**: Each object is an ordered sequence of values.

- **Graph data**

# What can you do with the data?

- Suppose that you are the owner of a supermarket and you have collected billions of market basket data. What information would you extract from it and how would you use it?

| TID | Items |
|-----|-------|
| 1 | Bread, Coke, Milk |
| 2 | Beer, Bread |
| 3 | Beer, Coke, Diaper, Milk |
| 4 | Beer, Bread, Diaper, Milk |
| 5 | Coke, Diaper, Milk |

Product placement

Catalog creation

Recommendations

- What if this was an online store?

# What can you do with the data?

- Suppose you are a search engine and you have a toolbar log consisting of
  - pages browsed,
  - queries,
  - pages clicked,
  - ads clicked

Ad click prediction

Query reformulations

- each with a user id and a timestamp. What information would you like to get our of the data?

# What can you do with the data?

- You are the owner of a social network, and you have full access to the social graph, what kind of information do you want to get out of your graph?

- Who is the most important node in the graph?
- What is the shortest path between two nodes?
- How many friends two nodes have in common?
- How does information spread on the network?

# Why data mining?

- Commercial point of view
- Data has become the key competitive advantage of companies
  - Examples: Facebook, Google, Amazon
  - Being able to extract useful information out of the data is key for exploiting them commercially.
- Scientific point of view
  - Scientists are at an unprecedented position where they can collect TB of information
    - Examples: Sensor data, astronomy data, social network data, gene data
  - We need the tools to analyze such data to get a better understanding of the world and advance science
- Scale (in data size and feature dimension)
  - Why not use traditional analytic methods?
  - Enormity of data, curse of dimensionality
  - The amount and the complexity of data does not allow for manual processing of the data. We need automated techniques.

# Why Data Preprocessing?

- Data in the real world is dirty
  - incomplete: lacking attribute values, lacking certain attributes of interest .
  - noisy: containing errors or outliers
  - inconsistent: containing discrepancies in codes or names
- No quality data, no quality mining results!
  - Quality decisions must be based on quality data
  - Data warehouse needs consistent integration of quality data
  - Required for both OLAP and Data Mining!

# Why can Data be Incomplete?

- Attributes of interest are not available (e.g., customer information for sales transaction data)
- Data were not considered important at the time of transactions, so they were not recorded!
- Data not recorder because of misunderstanding or malfunctions
- Data may have been recorded and later deleted!
- Missing/unknown values for some data

# Data Cleaning

- Data cleaning tasks
    - Fill in missing values
    - Identify outliers and smooth out noisy data
    - Correct inconsistent data

# What can we do with data mining?

- Some examples:

  - Frequent itemsets and Association Rules extraction

  - Clustering

  - Classification

  - Ranking

  - Exploratory analysis