

# CSE 626: Data mining

Instructor: Sargur N. Srihari

E-mail: [srihari@cedar.buffalo.edu](mailto:srihari@cedar.buffalo.edu)

Phone: 645-6164, ext. 113

# What is Data Mining?

- Different perspectives: CSE, Business, IT
- As a field of research in CSE:
  - Science of extracting useful information from large data sets or databases
- Also known as
  - Knowledge Discovery and Data Mining (KDD)
  - Knowledge Discovery in Databases (KDD)

# Data Mining Definitions

1. Analysis of datasets to find unsuspected relationships
2. Summarize data in novel ways that are  
understandable  
useful to data owner
3. Extraction of knowledge from data  
non-trivial extraction of implicit, previously unknown & potentially useful  
knowledge from data
4. Process of discovering patterns:  
automatically or semi-automatically, in large quantities of data  
  
Patterns discovered must be useful:  
meaningful in that they lead to some advantage, usually economic

# Why Data Mining?

1. Large datasets are common: due to advances in digital data acquisition and storage technology

## Business

Supermarket transactions  
Credit card usage records  
Telephone call details  
Government statistics

## Scientific

Images of astronomical bodies  
Molecular databases  
Medical records

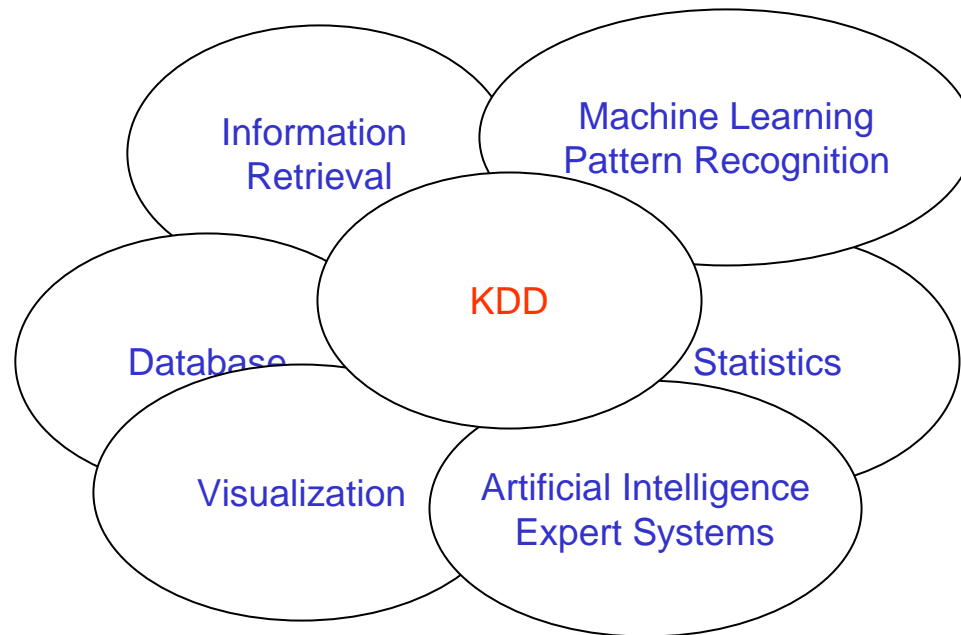
International organizations produce more information in a week than many people could read in a lifetime

2. *Automatic data production* leads to need for *automatic data consumption*

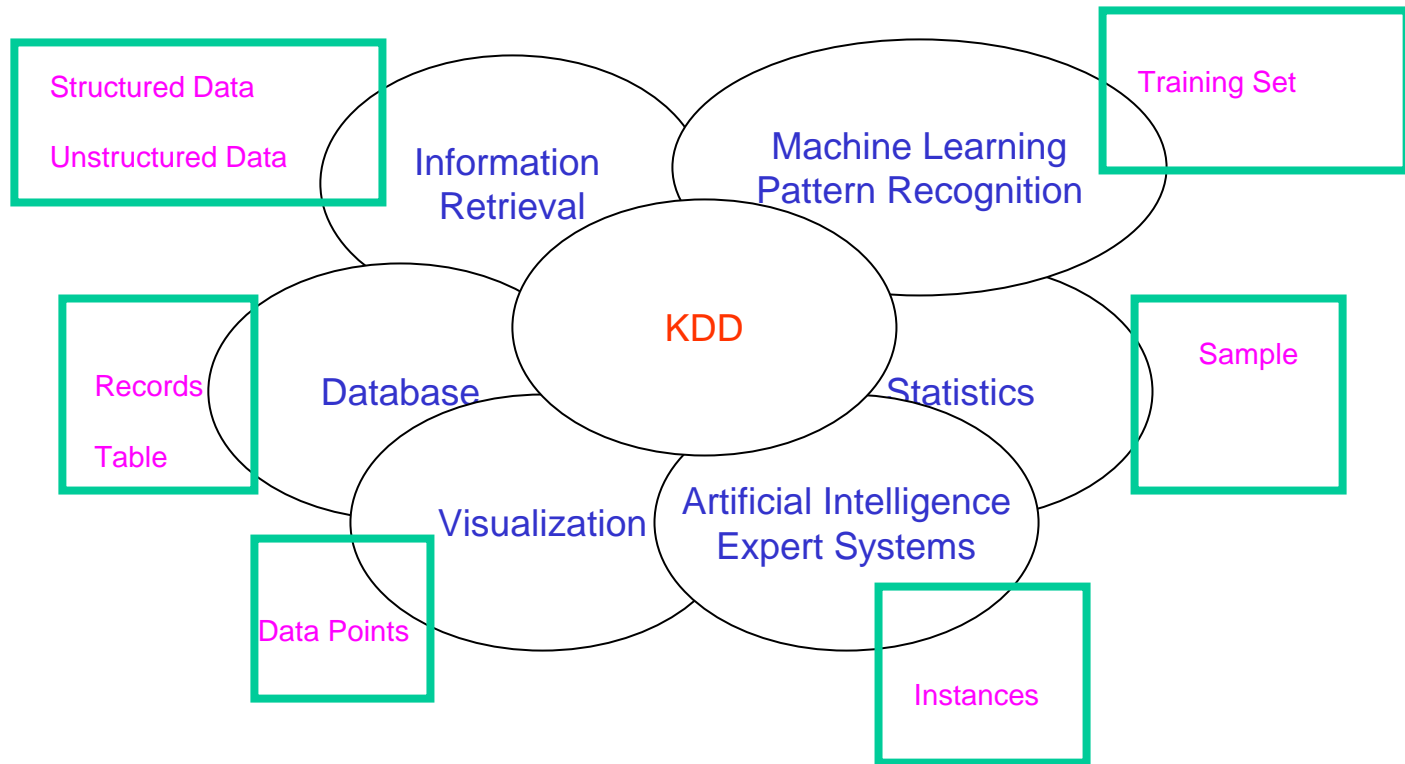
3. Large databases mean vast amounts of information

4. Difficulty lies in accessing it

# KDD is a multidisciplinary field



# Terminology for Data



# Course Textbook

Hand, David, Heikki Mannila, and Padhraic Smyth,  
*Principles of Data Mining*, MIT Press 2001.

## Approach:

Fundamental principles

Emphasis on Theory and Algorithms

## Many other textbooks:

Emphasize business applications, case studies

# Many Other Textbooks

1. Han and Kamber, *Data Mining Concepts and Techniques*, Morgan Kaufmann, 2000 (Data Base Perspective)
2. Witten, I. H., and E. Frank, *Data Mining: Practical Machine Learning Tools and Techniques with Java Implementations*, Morgan Kaufmann, 2000. (Machine Learning Perspective)
3. Adriaans, P., and D. Zantinge, *Data Mining*, Addison- Wesley, 1998. (Layman Perspective)
4. Groth, R., *Data Mining: A Hands-on Approach for Business Professionals*, Prentice-Hall PTR, 1997. (Business Perspective)
5. Kennedy, R., Y. Lee, et al., *Solving Data Mining Problems through Pattern Recognition*, Prentice-Hall PTR, 1998. (Pattern Recognition Perspective)
6. Weiss, S., and N. Indurkha, *Predictive Data Mining: A Practical Guide*, Morgan Kaufmann, 1998. (Statistical Perspective)

# More Data Mining Textbooks

7. S.Chakrabarti, *Mining the web*, Morgan Kaufman, 2003 (Emphasis on webpages and hyperlinks)
8. T. Dasu and T. Johnson, *Exploratory Data Mining and Data Cleaning*, Wiley, 2003 (Focus on data quality)
9. K. Cios, W. Pedrycz and R. Swiniarski, *Data Mining Methods for Knowledge Discovery*, Kluwer, 1998, (Focus on Mathematical issues, e.g., rough sets)
10. M. Kantardzic, *Data Mining: Concepts, Models and Algorithms*, IEEE-Wiley, 2003 (Focus on Machine Learning)
11. A. K. Pujari, *Data Mining Techniques*, Universities Press, 2001, (Data Base Perspective)
12. R. Groth, *Data Mining: A hands-on approach for business professionals*, Prentice Hall, 1998 (Business user perspective including software CD)

# Data Mining *vs* Statistics

- Objective of data mining exercise plays **no role** in data collection strategy
- In this way it **differs** from much of statistics
- For this reason, data mining is referred to as **secondary** data analysis
- KDD more complicated than initially thought
  - 80% preparing data
  - 20% mining data

# Query: Data Base vs Data Mining

Data Base: When you know **exactly** what you are looking for

- Query Tool: SQL (Structured Query Language) example

Table called **Persons**

<b>LastName</b>	<b>FirstName</b>	<b>Address</b>	<b>City</b>
Hansen	Ola	Timoteivn 10	Sandnes
Svendson	Tove	Borgvn 23	Sandnes
Pettersen	Kari	Storgt 20	Stavanger

- Query: **SELECT LastName FROM Persons** results in

**LastName**

Hansen
Svendson
Pettersen

Data Mining: When you only **vaguely** know what you are looking for

# Data Mining Tasks and Techniques

- Not so much a single technique
- Idea that there is more knowledge hidden in the data than shows itself on the surface
- Any technique that helps to extract more out of data is useful
- Five major task types:
  1. Exploratory Data Analysis (Visualization)
  2. Descriptive Modeling (Density estimation, Clustering)
  3. Predictive Modeling (Classification and Regression)
  4. Discovering Patterns and Rules (Association rules)
  5. Retrieval by Content (Retrieve items similar to pattern of interest)

Model  
building

# Topics in Data Mining

## 1. Fundamentals

- Nature of Data
- Measurement
- Summarizing and Visualization (includes PCA)
- Uncertainty and Inference

## 2. Data Mining Components

- Models
- Score Functions
- Optimization and Search

## 3. Data Mining Tasks and Algorithms

- Density Estimation and Clustering
- Classification (decision trees, neural networks, genetic algorithms)
- Regression
- Pattern Discovery (association rules)
- Retrieval by Content (includes Image Retrieval and Text Analytics)