# CS6220: DATA MINING TECHNIQUES

## 1: Introduction

**Instructor: Yizhou Sun**

yzsun@ccs.neu.edu

September 9, 2013

# Course Information

- Course homepage:
  http://www.ccs.neu.edu/home/yzsun/classes/2013Fall_CS6220/index.htm
  - Class schedule
  - Slides
  - Announcement
  - Assignments
  - …
- Piazza:
  https://piazza.com/northeastern/fall2013/cs6220/home

- Prerequisites
  - CS 5800 or CS 7800, or consent of instructor
  - More generally
    - You are expected to have background knowledge in data structures, algorithms, basic linear algebra, and basic statistics.
    - You will also need to be familiar with at least one programming language, and have programming experiences.

# **Meeting Time and Location**

- When
  - Tuesdays, 6-9pm
- Where
  - Behrakis Health Sciences Center 310

# Instructor and TA Information

- Instructor: Yizhou Sun
  - Homepage: http://www.ccs.neu.edu/home/yzsun/
  - Email: yzsun@ccs.neu.edu
  - Office: 320 WVH
  - Office hour: Wednesdays 3-5pm
- TA:
  - Moonyoung (Moon) Kang
    - Email: yerihyo@gmail.com
    - Office hours: Tuesdays 2-4pm at 472 WVH
  - Qizhen Ruan
    - Email: ruan.qi@husky.neu.edu
    - Office hours: Mondays 4:30-6:30pm at 102 Main Lab WVH

# Grading

- Homework: 40%

- Midterm exam: 25%

- Course project: 30%

- Participation: 5%

# Grading: Homework

- Homework: 40%

  - Four assignments are expected
    - 2 paper-based assignments
    - 2 program-based assignments

  - Deadline: 11:59pm of the indicated due date via *Blackboard* or class system
    - within 1 hour late: 90% max; within 8 hours late: 60% max; otherwise: 0%

  - No copying or sharing of homework!
    - But you can discuss general challenges and ideas with others

# Grading: Midterm Exam

- Midterm exam: 25%
  - Closed book exam, but you can take a "cheating sheet" of A4 size

# Grading: Course Project

- Course project: 30%
  - Group project (3-4 people for one group)
  - Goal: Choose one interesting problem, formalize it as a data mining task, collect data, provide solutions, and evaluate and compare your solutions.
  - You are expected to submit one project proposal early this semester, and your datasets, code, and a project report at the end of the semester
  - You are expected to present your project at the end of the semester.

# Grading: Participation

- Participation (5%)
  - In-class participation
  - Online participation (piazza)

# Textbook

- Jiawei Han, Micheline Kamber, and Jian Pei. Data Mining: Concepts and Techniques, 3rd edition, Morgan Kaufmann, 2011
- References
  - "Data Mining" by Pang-Ning Tan, Michael Steinbach, and Vipin Kumar (http://www-users.cs.umn.edu/~kumar/dmbook/index.php)
  - "Machine Learning" by Tom Mitchell (http://www.cs.cmu.edu/~tom/mlbook.html)
  - "Introduction to Machine Learning" by Ethem ALPAYDIN (http://www.cmpe.boun.edu.tr/~ethem/i2ml/)
  - "Pattern Classification" by Richard O. Duda, Peter E. Hart, David G. Stork (http://www.wiley.com/WileyCDA/WileyTitle/productCd-0471056693.html)
  - "The Elements of Statistical Learning: Data Mining, Inference, and Prediction" by Trevor Hastie, Robert Tibshirani, and Jerome Friedman (http://www-stat.stanford.edu/~tibs/ElemStatLearn/)
  - "Pattern Recognition and Machine Learning" by Christopher M. Bishop (http://research.microsoft.com/en-us/um/people/cmbishop/prml/)

# Course Content

- By data types:
  - matrix data
  - set data
  - sequence data
  - time series
  - graph and network (Next Semester: Advanced Topics)
- By functions:
  - Classification
  - Clustering
  - Frequent pattern mining
  - Prediction
  - Similarity search
  - Ranking

# Goal of the Course

- Know what is data mining and the basic algorithms

- Know how to apply algorithms to real-world applications

- Provide a starting course for research in data mining

# 1. Introduction

- Why Data Mining?

- What Is Data Mining?

- A Multi-Dimensional View of Data Mining

  - What Kinds of Data Can Be Mined?

  - What Kinds of Patterns Can Be Mined?

  - What Kinds of Technologies Are Used?

  - What Kinds of Applications Are Targeted?

- Major Issues in Data Mining

# Why Data Mining?

- The Explosive Growth of Data: from terabytes to petabytes

  - Data collection and data availability

    - Automated data collection tools, database systems, Web, computerized society

  - Major sources of abundant data

    - Business: Web, e-commerce, transactions, stocks, …

    - Science: Remote sensing, bioinformatics, scientific simulation, …

    - Society and everyone: news, digital cameras, YouTube

- <u>We are drowning in data, but starving for knowledge!</u>

- "Necessity is the mother of invention"—Data mining—Automated analysis of massive data sets

# Big Data Challenges

- Video 1: Big Data Challenges (Ads by DataStax)
  - http://www.youtube.com/watch?v=or6Pse8fxD4


- Video 2: Explaining Big Data
  - http://www.youtube.com/watch?v=7D1CQ_LOizA

# 1. Introduction

- Why Data Mining?

- What Is Data Mining?

- A Multi-Dimensional View of Data Mining

  - What Kinds of Data Can Be Mined?

  - What Kinds of Patterns Can Be Mined?

  - What Kinds of Technologies Are Used?

  - What Kinds of Applications Are Targeted?

- Major Issues in Data Mining

- A Brief History of Data Mining and Data Mining Society

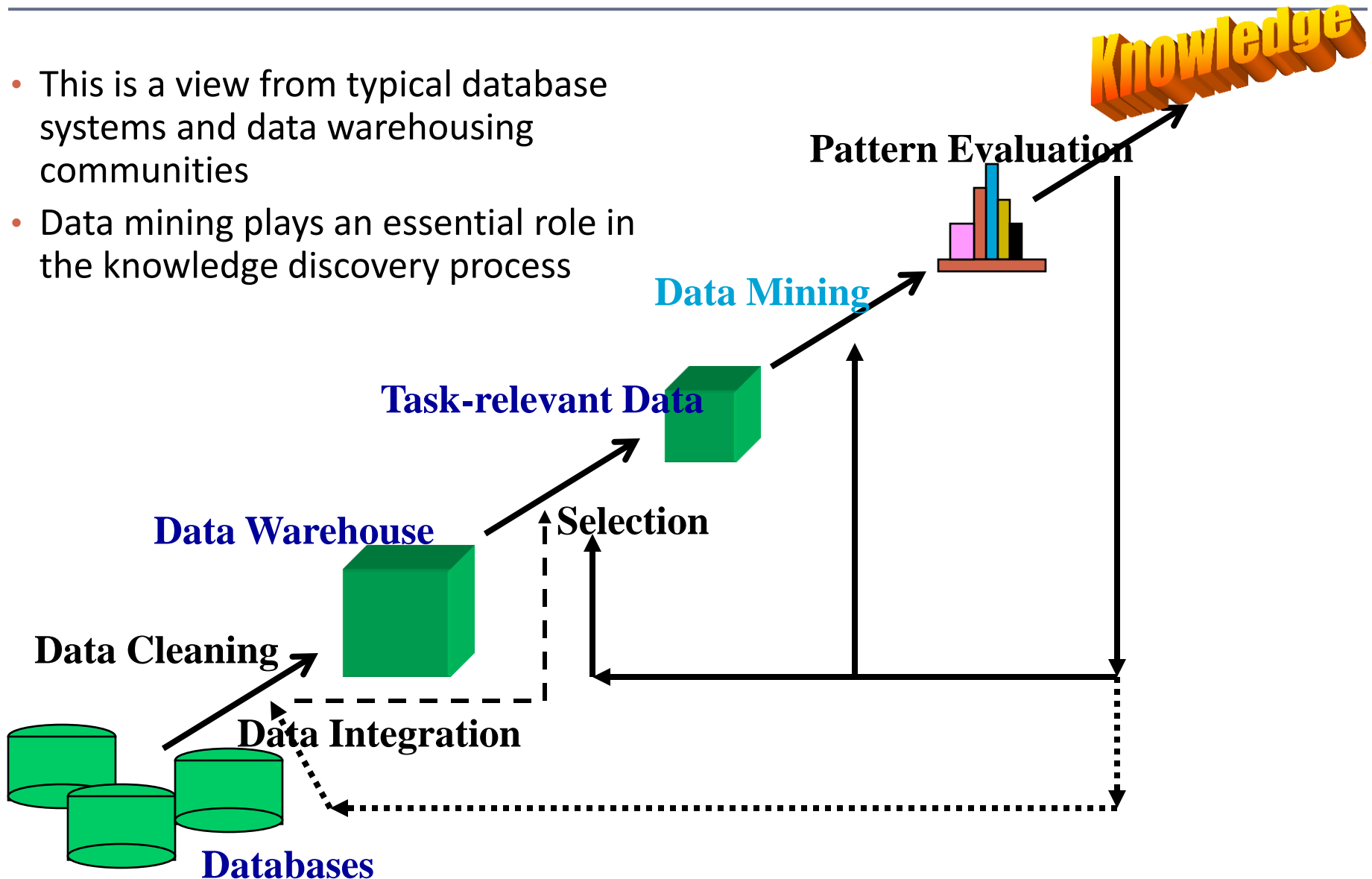- Summary

# What Is Data Mining?

- Data mining (knowledge discovery from data)
  - Extraction of interesting (non-trivial, implicit, previously unknown and potentially useful) patterns or knowledge from huge amount of data

- Alternative names
  - Knowledge discovery (mining) in databases (KDD), knowledge extraction, data/pattern analysis, data archeology, data dredging, information harvesting, business intelligence, etc.
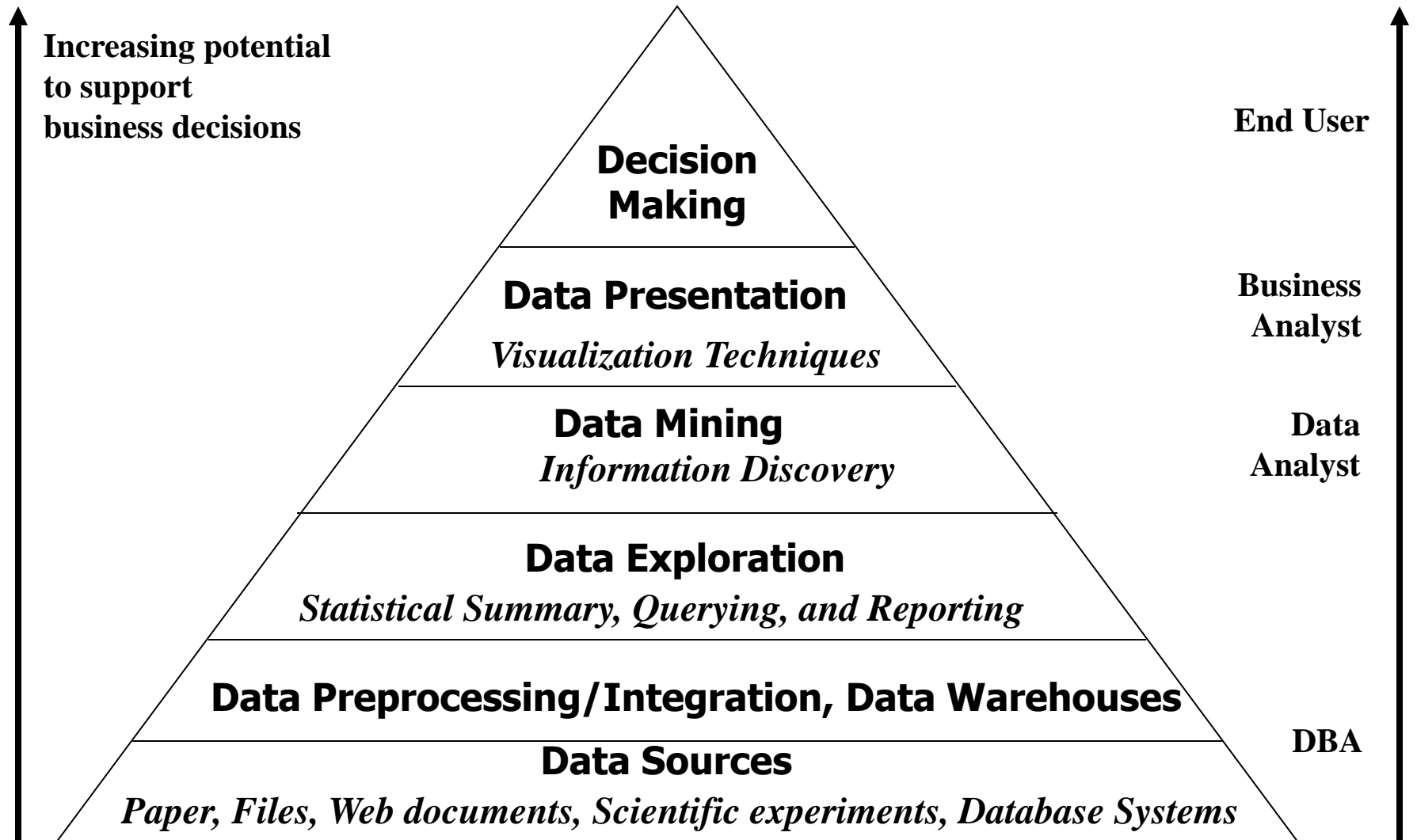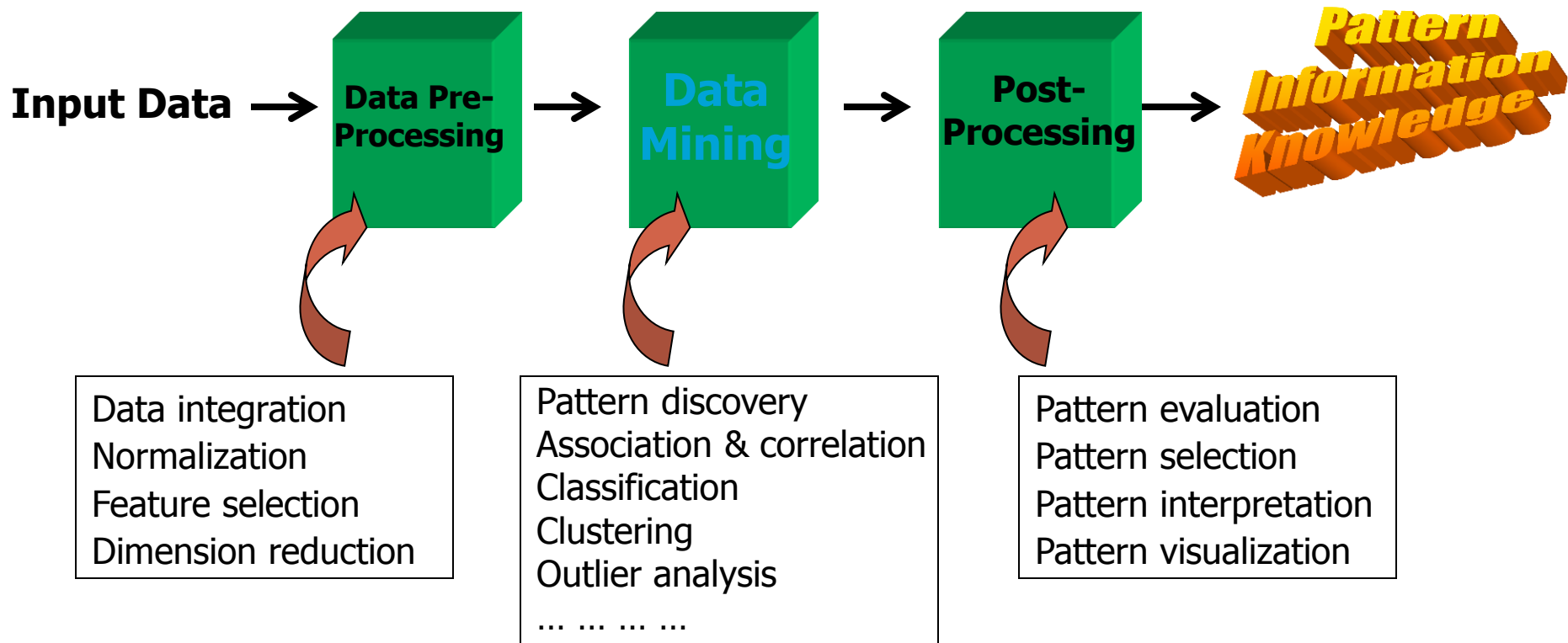
# Knowledge Discovery (KDD) Process

- This is a view from typical database systems and data warehousing communities

- Data mining plays an essential role in the knowledge discovery process

**Knowledge**

**Pattern Evaluation**

**Data Mining**

**Task-relevant Data**

**Data Warehouse**

**Selection**

**Data Cleaning**

**Data Integration**

**Databases**

# Data Mining in Business Intelligence

**Increasing potential
to support
business decisions**

**End User**

### Decision Making

**Business Analyst**

### Data Presentation
*Visualization Techniques*

**Data Analyst**

### Data Mining
*Information Discovery*

### Data Exploration
*Statistical Summary, Querying, and Reporting*

### Data Preprocessing/Integration, Data Warehouses

**DBA**

### Data Sources
*Paper, Files, Web documents, Scientific experiments, Database Systems*

# KDD Process: A Typical View from ML and Statistics

**Input Data** → Data Pre-Processing → **Data Mining** → Post-Processing → *Pattern Information Knowledge*

| Data integration<br>Normalization<br>Feature selection<br>Dimension reduction | Pattern discovery<br>Association & correlation<br>Classification<br>Clustering<br>Outlier analysis<br>… … … … | Pattern evaluation<br>Pattern selection<br>Pattern interpretation<br>Pattern visualization |

- This is a view from typical machine learning and statistics communities

# 1. Introduction

- Why Data Mining?

- What Is Data Mining?

- A Multi-Dimensional View of Data Mining

  - What Kinds of Data Can Be Mined?

  - What Kinds of Patterns Can Be Mined?

  - What Kinds of Technologies Are Used?

  - What Kinds of Applications Are Targeted?

- Major Issues in Data Mining

# Multi-Dimensional View of Data Mining

- **Data to be mined**
  - Database data (extended-relational, object-oriented, heterogeneous, legacy), data warehouse, transactional data, stream, spatiotemporal, time-series, sequence, text and web, multi-media, graphs & social and information networks
- **Knowledge to be mined (or: Data mining functions)**
  - Characterization, discrimination, association, classification, clustering, trend/deviation, outlier analysis, etc.
  - Descriptive vs. predictive data mining
  - Multiple/integrated functions and mining at multiple levels
- **Techniques utilized**
  - Data-intensive, data warehouse (OLAP), machine learning, statistics, pattern recognition, visualization, high-performance, etc.
- **Applications adapted**
  - Retail, telecommunication, banking, fraud analysis, bio-data mining, stock market analysis, text mining, Web mining, etc.

# 1. Introduction

- Why Data Mining?

- What Is Data Mining?

- A Multi-Dimensional View of Data Mining

  - What Kinds of Data Can Be Mined?

  - What Kinds of Patterns Can Be Mined?

  - What Kinds of Technologies Are Used?

  - What Kinds of Applications Are Targeted?

- Major Issues in Data Mining

# Data Mining: On What Kinds of Data?

- Database-oriented data sets and applications

  - Relational database, data warehouse, transactional database

- Advanced data sets and advanced applications

  - Data streams and sensor data

  - Time-series data, temporal data, sequence data (incl. bio-sequences)

  - Structure data, graphs, social networks and multi-linked data

  - Object-relational databases

  - Heterogeneous databases and legacy databases

  - Spatial data and spatiotemporal data

  - Multimedia database

  - Text databases

  - The World-Wide Web

# Matrix Data

|       | Sex | Race | Height | Income | Marital Status | Years of Educ. | Liberal-ness |
|-------|-----|------|--------|--------|----------------|----------------|--------------|
| R1001 | M   | 1    | 70     | 50     | 1              | 12             | 1.73         |
| R1002 | M   | 2    | 72     | 100    | 2              | 20             | 4.53         |
| R1003 | F   | 1    | 55     | 250    | 1              | 16             | 2.99         |
| R1004 | M   | 2    | 65     | 20     | 2              | 16             | 1.13         |
| R1005 | F   | 1    | 60     | 10     | 3              | 12             | 3.81         |
| R1006 | M   | 1    | 68     | 30     | 1              | 9              | 4.76         |
| R1007 | F   | 5    | 66     | 25     | 2              | 21             | 2.01         |
| R1008 | F   | 4    | 61     | 43     | 1              | 18             | 1.27         |
| R1009 | M   | 1    | 69     | 67     | 1              | 12             | 3.25         |

# Set Data

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

# Sequence Data

## SYNTENIC ASSEMBLIES FOR CG15386

```
MD106   ATGCTTAGTAATCCCTACTTTAAGTCCGTTTTGTGGCTGATTGGCTTCGGAGGAATGGG
NEWC    ATGCTTAGTAATCCTTACTTTAAATCCGTTTTGTGGCTGATTGGCTTCGGAGGAATGGG
W501    ATGCTTAGTAATCCCTACTTTAAGTCCGTTTTGTGGCTGATTGGCTTCGGAGGAATGGG
MD199   ATGCTTAGTAATCCCTACTTTAAGTCCGTTTTGTGGCTGATTGGCTTCGGAGGAATGGG
C1674   ATGCTTAGTAATCCCTACTTTAAGTCCGTTTTGTGGCTGATTGGCTTCGGAGGAATGGG
SIM4    ATGCTTAGTAATCCCTACTTTAAGTCCGTTTTGTGGCTGATTGGCTTCGGAGGAATGGG

MD106   CTACGGCCTAATGGTGCTAACAGAGCCGAACGTCGACAAAATAGAGCGCATCAAAGCCT
NEWC    CTACGGCCTAATGGTGCTAACCGAGCCGAACGTCGACAAAATAGAGCGCATCAAAGCCT
W501    CTACGGCCTAATGGTGCTAACCGAGCCGAACGTCGACAAAATAGAGCGCATCAAAGCCT
MD199   CTACGGCCTAATGGTGCTAACCGAGCCGAACGTCGACAAAATAGAGCGCATCAAAGCCT
C1674   CTACGGCCTAATGGTGCTAACCGAGCCGAACGTCGACAAAATAGAGCGCATCAAAGCCT
SIM4    CTACGGCCTAATGGTGCTAACCGAGCCGAACGTCGACAAAATAGAGCGCATCAAAGCCT

MD106   CCGTTTCAAGTACCAAACTGAGTGCGGATGAGCAGCGAAAGGCTCTGTTTATGAAGAAG
NEWC    CCGTTTCAAGTACCAAACTGAGTGCGGATGAGCAGCGAAAGGCTCTGTTTATGAAGAAG
W501    CCGTTTCAAGTACCAAACTGAGTGCGGATGAGCAGCGAAAGGCTCTGTTTATGAAGAAG
MD199   CCGTTTCAAGTACCAAACTGAGTGCGGATGAGCAGCGAAAGGCTCTGTTTATGAAGAAG
C1674   CCGTTTCAAGTACCAAACTGAGTGCGGATGAGCAGCGAAAGGCTCTGTTTATGAAGAAG
SIM4    CCGTTTCAAGTACCAAACTGAGTGCGGATGAGCAGCGAAAGGCTCTGTTTATGAAGAAG

MD106   CTGCAGGAGGCGTCCACCACCAGTGCCCCAATCTACAGGTCAGCGGCCGAGAAATAG
NEWC    CTGCAGGAGGCGTCCACCACCAGTGCCCCAATCTACAGGTCATCGGCCGAGAAATAG
W501    CTGCAGGAGGCGTCCACCACCACTGCCCCAATCTACAGGTCATCGGCCGAGAAATAG
MD199   CTGCAGGAGGCGTCCACCACCAGTGCCCCAATCTACAGGTCAGCGGCCGAGAAATAG
C1674   CTGCAGGAGGCGTCCACCACCAGTGCCCCAATCTACAGGTCAGCGGCCGAGAAATAG
SIM4    CTGCAGGAGGCGTCCACCACCAGTGCCCCAATCTACAGGTCAGCGGCCGAGAAATAG
```
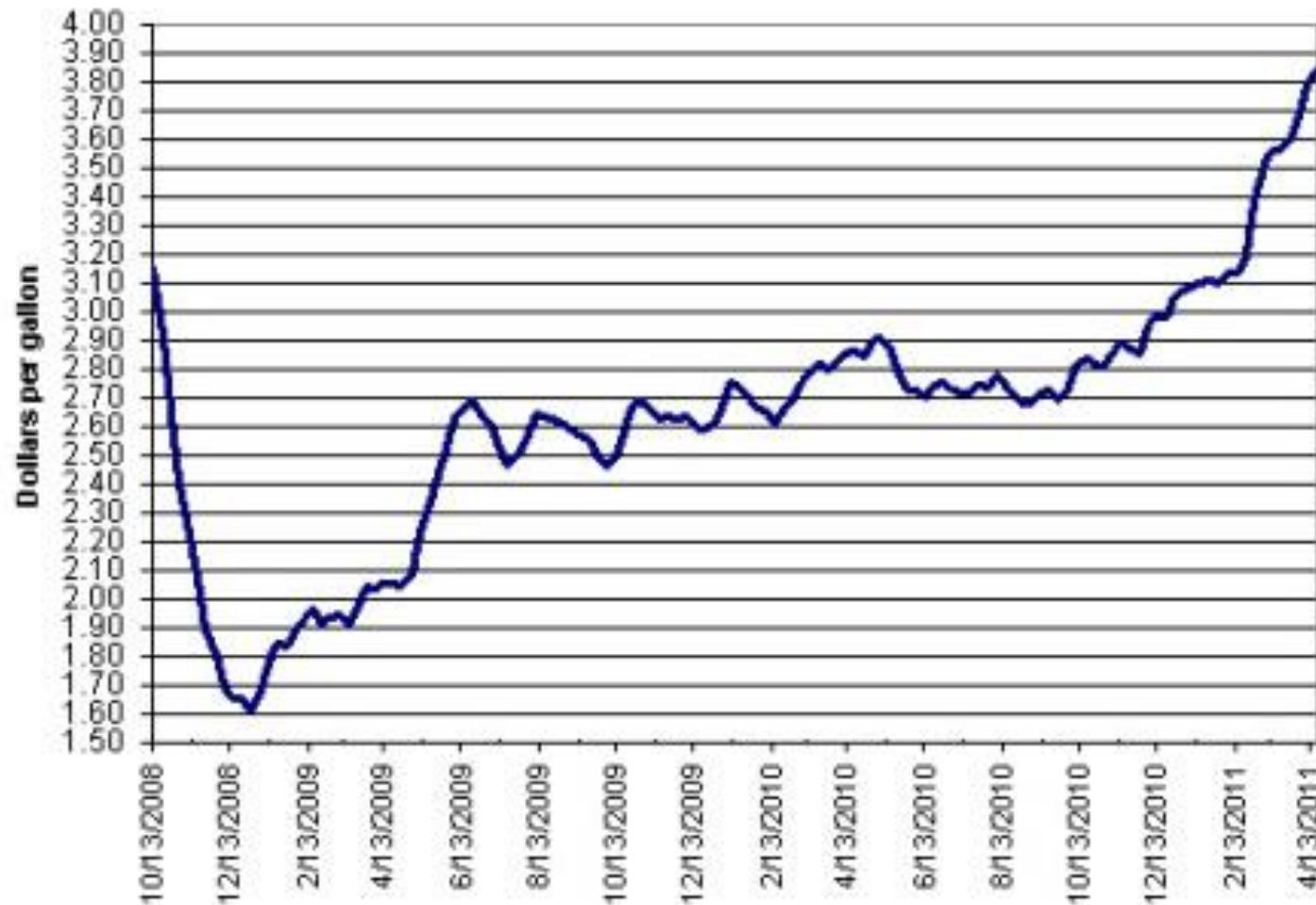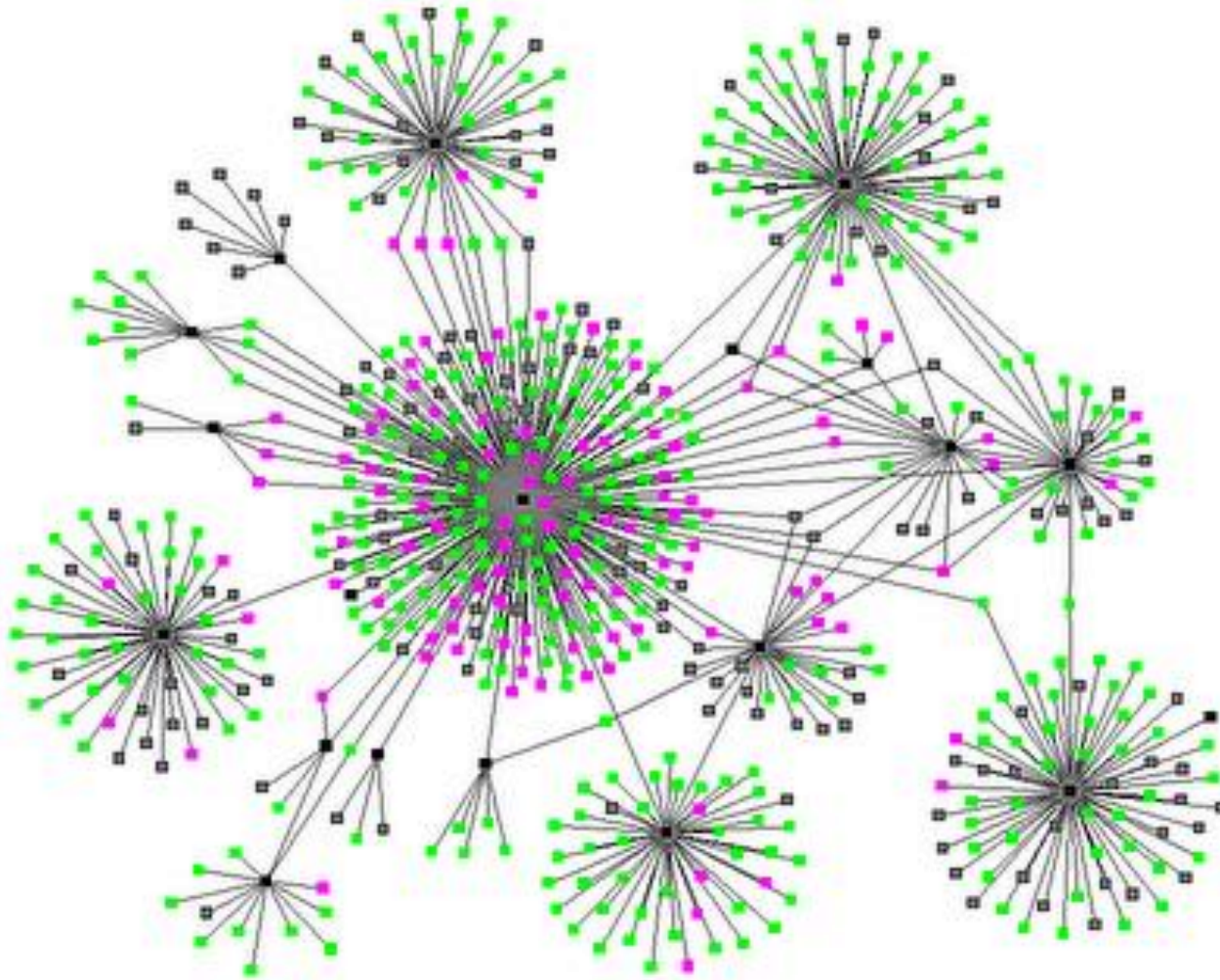
# Time Series



Weekly U.S. Retail Gasoline Prices, Regular Grade

Source: Energy Information Administration

29

# Graph / Network

# 1. Introduction

- Why Data Mining?

- What Is Data Mining?

- A Multi-Dimensional View of Data Mining

  - What Kinds of Data Can Be Mined?

  - What Kinds of Patterns Can Be Mined?

  - What Kinds of Technologies Are Used?

  - What Kinds of Applications Are Targeted?

- Major Issues in Data Mining

# Data Mining Function: Association and Correlation Analysis

- Frequent patterns (or frequent itemsets)

  - What items are frequently purchased together in your Walmart?

- Association, correlation vs. causality

  - A typical association rule

    - Diaper → Beer [0.5%, 75%]  (support, confidence)

  - Are strongly associated items also strongly correlated?

# Data Mining Function: Classification

- Classification and label prediction
  - Construct models (functions) based on some training examples
  - Describe and distinguish classes or concepts for future prediction
    - E.g., classify countries based on (climate), or classify cars based on (gas mileage)
  - Predict some unknown class labels
- Typical methods
  - Decision trees, naïve Bayesian classification, support vector machines, neural networks, rule-based classification, pattern-based classification, logistic regression, …
- Typical applications:
  - Credit card fraud detection, direct marketing, classifying stars, diseases, web-pages, …

# Data Mining Function: Cluster Analysis

- Unsupervised learning (i.e., Class label is unknown)

- Group data to form new categories (i.e., clusters), e.g., cluster houses to find distribution patterns

- Principle: Maximizing intra-class similarity & minimizing interclass similarity

- Many methods and applications

# Data Mining Function: Others

- Prediction

- Similarity search
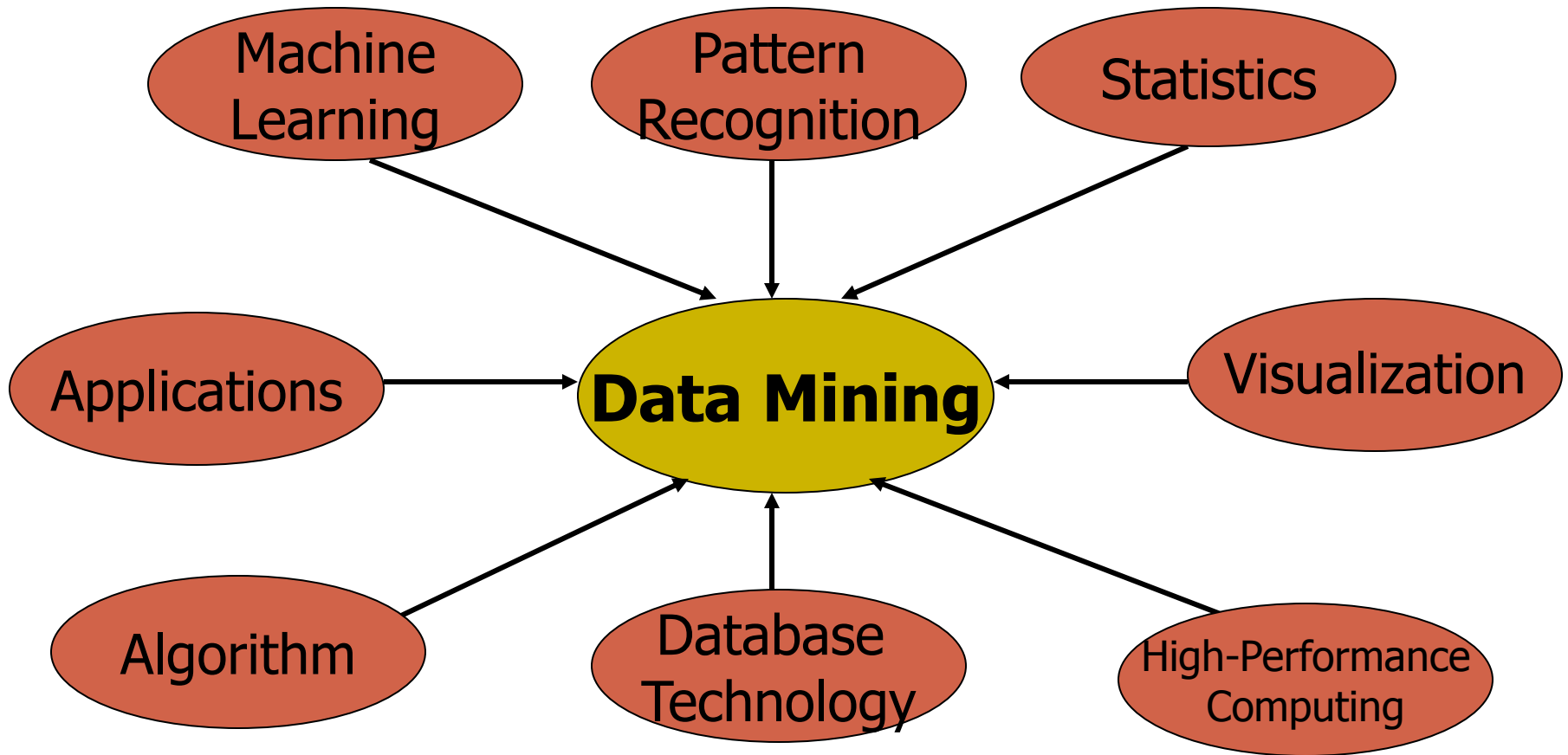
- Ranking

- Outlier detection

- ...

# Evaluation of Knowledge

- Are all mined knowledge interesting?
  - One can mine tremendous amount of "patterns" and knowledge
  - Some may fit only certain dimension space (time, location, …)
  - Some may not be representative, may be transient, …

- Evaluation of mined knowledge → directly mine only interesting knowledge?
  - Descriptive vs. predictive
  - Coverage
  - Typicality vs. novelty
  - Accuracy
  - Timeliness
  - …

# 1. Introduction

- Why Data Mining?

- What Is Data Mining?

- A Multi-Dimensional View of Data Mining

  - What Kinds of Data Can Be Mined?

  - What Kinds of Patterns Can Be Mined?

  - What Kinds of Technologies Are Used?

  - What Kinds of Applications Are Targeted?

- Major Issues in Data Mining

# Data Mining: Confluence of Multiple Disciplines

# 1. Introduction

- Why Data Mining?

- What Is Data Mining?

- A Multi-Dimensional View of Data Mining

  - What Kinds of Data Can Be Mined?

  - What Kinds of Patterns Can Be Mined?

  - What Kinds of Technologies Are Used?

  - What Kinds of Applications Are Targeted?

- Major Issues in Data Mining

# Applications of Data Mining

- Web page analysis: from web page classification, clustering to PageRank & HITS algorithms

- Collaborative analysis & recommender systems

- Basket data analysis to targeted marketing

- Biological and medical data analysis: classification, cluster analysis (microarray data analysis),  biological sequence analysis, biological network analysis

- Data mining and software engineering (e.g., IEEE Computer, Aug. 2009 issue)

- Social media

- Game

# Example

- Street Bump Boston Project
  - http://www.cityofboston.gov/doit/apps/streetbump.asp

# 1. Introduction

- Why Data Mining?

- What Is Data Mining?

- A Multi-Dimensional View of Data Mining

  - What Kinds of Data Can Be Mined?

  - What Kinds of Patterns Can Be Mined?

  - What Kinds of Technologies Are Used?

  - What Kinds of Applications Are Targeted?

- Major Issues in Data Mining

# Major Issues in Data Mining (1)

- Mining Methodology
  - Mining various and new kinds of knowledge
  - Mining knowledge in multi-dimensional space
  - Data mining: An interdisciplinary effort
  - Boosting the power of discovery in a networked environment
  - Handling noise, uncertainty, and incompleteness of data
  - Pattern evaluation and pattern- or constraint-guided mining
- User Interaction
  - Interactive mining
  - Incorporation of background knowledge
  - Presentation and visualization of data mining results

# Major Issues in Data Mining (2)

- Diversity of data types

  - Handling complex types of data

  - Mining dynamic, networked, and global data repositories

- Efficiency and Scalability

  - Efficiency and scalability of data mining algorithms

  - Parallel, distributed, stream, and incremental mining methods

- Data mining and society

  - Social impacts of data mining

  - Privacy-preserving data mining

# Where to Find References? DBLP, CiteSeer, Google

- <u>Data mining and KDD (SIGKDD: CDROM)</u>
  - Conferences: ACM-SIGKDD, IEEE-ICDM, SIAM-DM, PKDD, PAKDD, etc.
  - Journal: Data Mining and Knowledge Discovery, KDD Explorations, ACM TKDD
- <u>Database systems (SIGMOD: ACM SIGMOD Anthology—CD ROM)</u>
  - Conferences: ACM-SIGMOD, ACM-PODS, VLDB, IEEE-ICDE, EDBT, ICDT, DASFAA
  - Journals: IEEE-TKDE, ACM-TODS/TOIS, JIIS, J. ACM, VLDB J., Info. Sys., etc.
- <u>AI & Machine Learning</u>
  - Conferences: Machine learning (ML), AAAI, IJCAI, COLT (Learning Theory), CVPR, NIPS, etc.
  - Journals: Machine Learning, Artificial Intelligence, Knowledge and Information Systems, IEEE-PAMI, etc.
- <u>Web and IR</u>
  - Conferences: SIGIR, WWW, CIKM, etc.
  - Journals: WWW: Internet and Web Information Systems,
- <u>Statistics</u>
  - Conferences: Joint Stat. Meeting, etc.
  - Journals: Annals of statistics, etc.
- <u>Visualization</u>
  - Conference proceedings: CHI, ACM-SIGGraph, etc.
  - Journals: IEEE Trans. visualization and computer graphics, etc.

# Recommended Reference Books

- **E. Alpaydin. Introduction to Machine Learning, 2nd ed., MIT Press, 2011**

- S. Chakrabarti. Mining the Web: Statistical Analysis of Hypertex and Semi-Structured Data. Morgan Kaufmann, 2002

- R. O. Duda, P. E. Hart, and D. G. Stork, Pattern Classification, 2ed., Wiley-Interscience, 2000

- T. Dasu and T. Johnson.  Exploratory Data Mining and Data Cleaning. John Wiley & Sons, 2003

- U. M. Fayyad, G. Piatetsky-Shapiro, P. Smyth, and R. Uthurusamy. Advances in Knowledge Discovery and Data Mining. AAAI/MIT Press, 1996

- U. Fayyad, G. Grinstein, and A. Wierse, Information Visualization in Data Mining and Knowledge Discovery, Morgan Kaufmann, 2001

- J. Han, M. Kamber, and J. Pei, Data Mining: Concepts and Techniques. Morgan Kaufmann, 3$^{rd}$ ed. , 2011

- T. Hastie, R. Tibshirani, and J. Friedman, The Elements of Statistical Learning: Data Mining, Inference, and Prediction, 2$^{nd}$ ed., Springer, 2009

- B. Liu, Web Data Mining, Springer 2006

- T. M. Mitchell, Machine Learning, McGraw Hill, 1997

- Y. Sun and J. Han, Mining Heterogeneous Information Networks, Morgan & Claypool, 2012

- P.-N. Tan, M. Steinbach and V. Kumar, Introduction to Data Mining, Wiley, 2005

- S. M. Weiss and N. Indurkhya, Predictive Data Mining, Morgan Kaufmann, 1998

- I. H. Witten and E. Frank,  Data Mining: Practical Machine Learning Tools and Techniques with Java Implementations, Morgan Kaufmann, 2$^{nd}$ ed. 2005