# CS6220: DATA MINING TECHNIQUES

## 1: Introduction

**Instructor: Yizhou Sun**

yzsun@ccs.neu.edu

September 28, 2015

# Course Information

- Course homepage:
  [http://www.ccs.neu.edu/home/yzsun/classes/2015Fall_CS6220/index.htm](http://www.ccs.neu.edu/home/yzsun/classes/2015Fall_CS6220/index.htm)

  - Class schedule
  - Slides
  - Announcement
  - Assignments
  - …

- 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
  - Monday, 6-9pm
- Where
  - Forsyth Building 236

# Instructor and TA Information

- Instructor: Yizhou Sun
  - Homepage: http://www.ccs.neu.edu/home/yzsun/
  - Email: yzsun@ccs.neu.edu
  - Office: 358 WVH
  - Office hour: Tuesdays 10-12pm
- TA: Monisha Singh
  - Email: msingh28@ccs.neu.edu
  - Office hours: Thursdays 10:00-12:00pm at 462 WVH

# Grading

- Homework: 40%

- Midterm exam: 25%

- Course project: 30%

- Participation: 5%

# Grading: Homework

- Homework: 40%

  - Six assignments are expected
  - Deadline: 11:59pm of the indicated due date via *Blackboard* or class system
    - *No Late Submission!*
  - No copying or sharing of homework!
    - But you can discuss general challenges and ideas with others
    - *Suspicious cases will be reported to OSCCR (*Office of Student Conduct and Conflict Resolution*)*

# 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: Solve an open data mining problem
  - You are expected to submit a project report and your code at the end of the semester

# Grading: Participation

- Participation (5%)
  - In-class participation
  - quizzes
  - Online participation (piazza)
    - piazza.com/northeastern/fall2014/cs6220

# 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/)

# 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?

- Content covered by this course

# 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

- We are drowning in data, but starving for knowledge!

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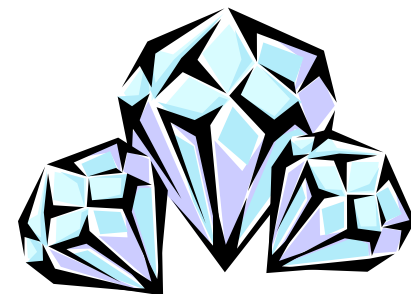
# 1. Introduction

- 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?

- Content covered by this course
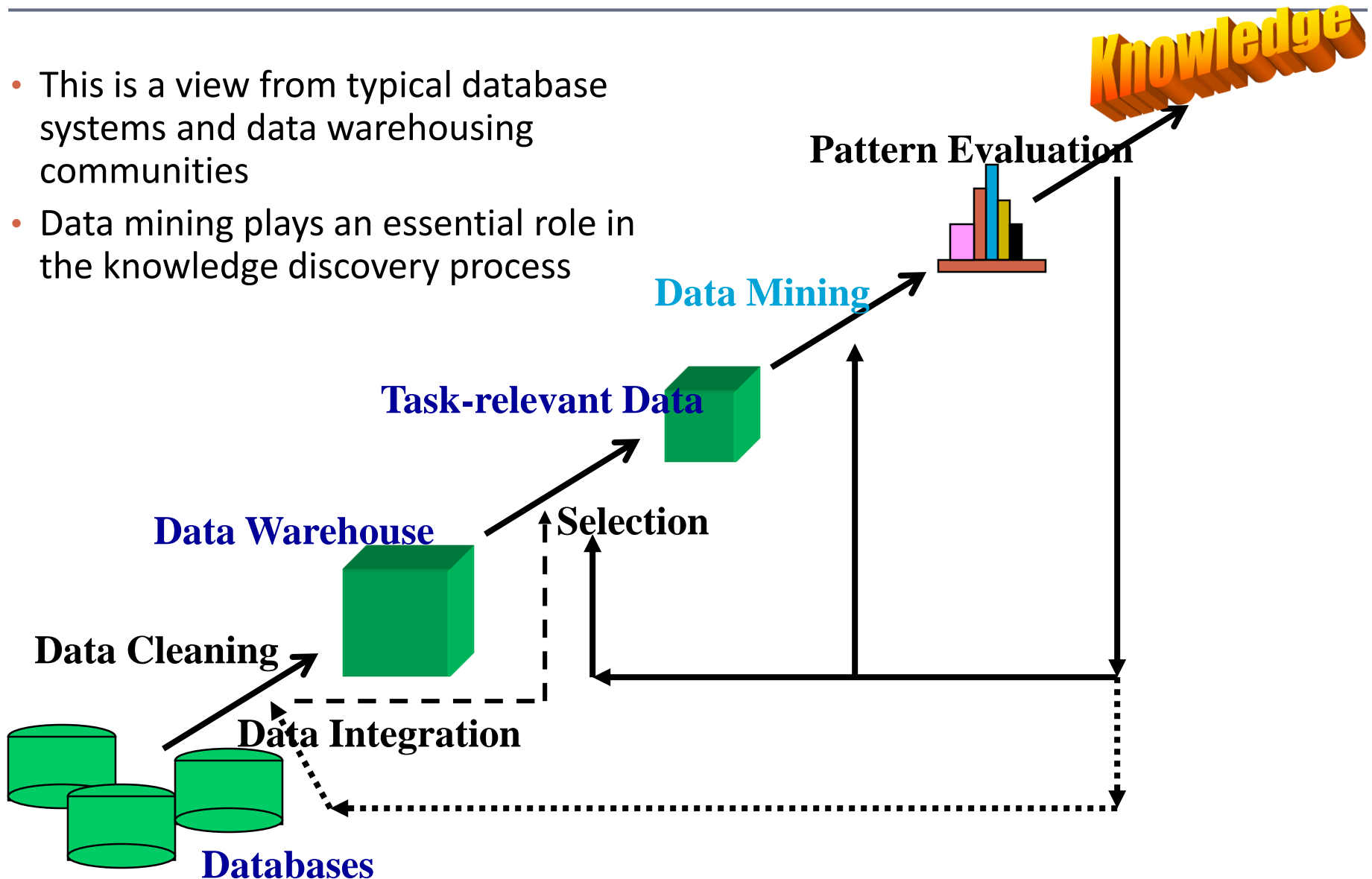
# What Is Data Mining?

- 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.
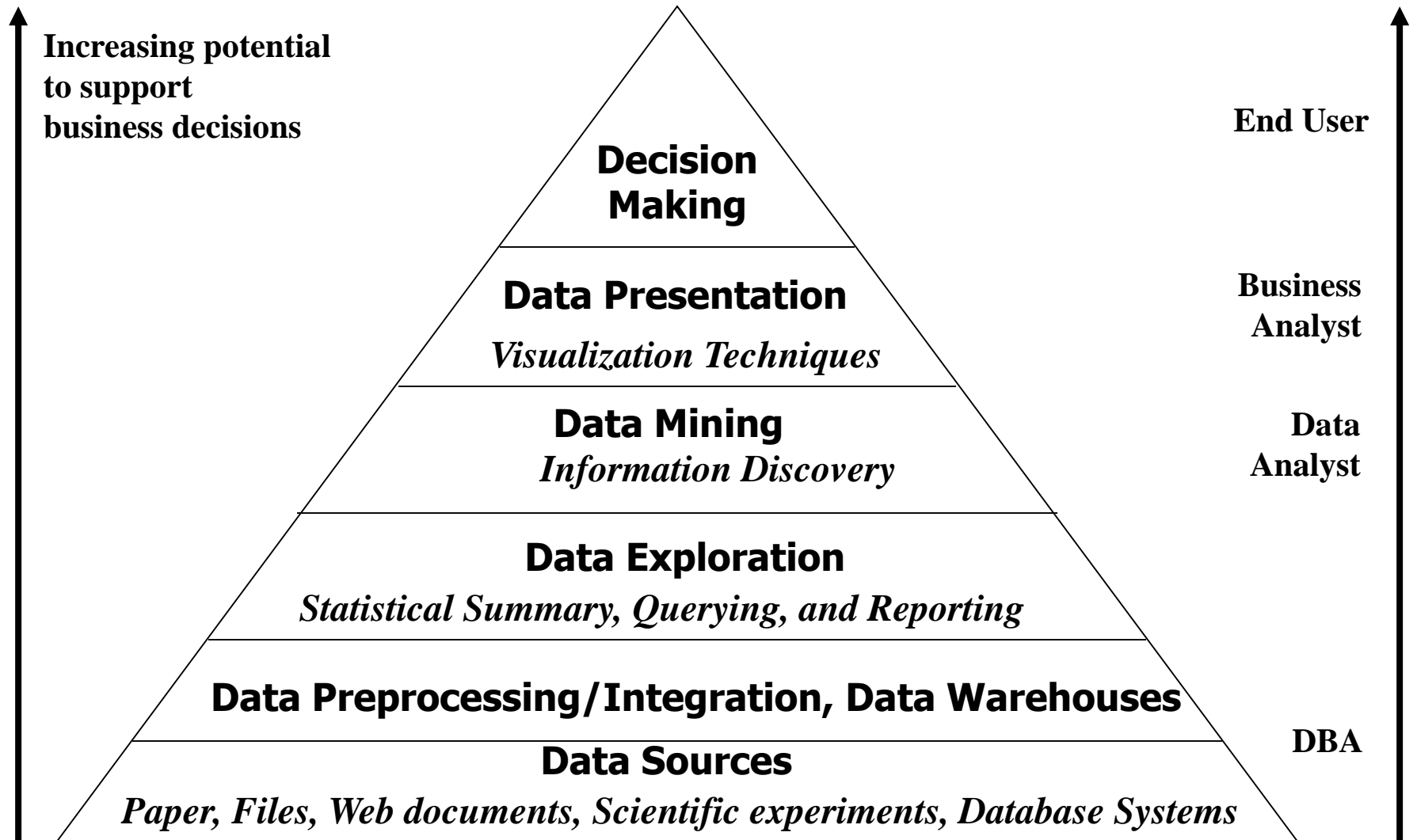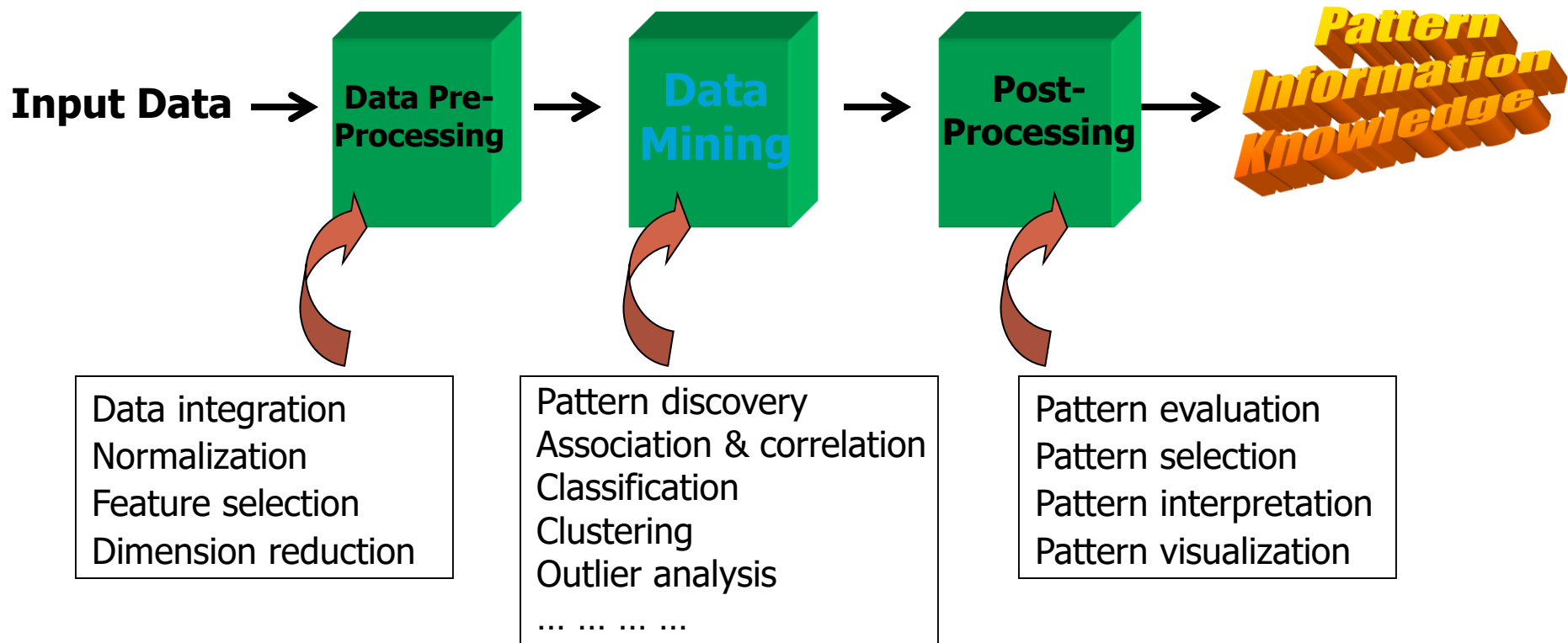
# Knowledge Discovery (KDD) Process

- 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**

**Decision
Making**

**End User**

**Data Presentation**

*Visualization Techniques*

**Business
Analyst**

**Data Mining**

*Information Discovery*

**Data
Analyst**

**Data Exploration**

*Statistical Summary, Querying, and Reporting*

**Data Preprocessing/Integration, Data Warehouses**

**Data Sources**

*Paper, Files, Web documents, Scientific experiments, Database Systems*

**DBA**

# KDD Process: A Typical View from ML and Statistics

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

Data integration
Normalization
Feature selection
Dimension reduction

Pattern discovery
Association & correlation
Classification
Clustering
Outlier analysis
... ... ... ...

Pattern evaluation
Pattern selection
Pattern interpretation
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?

- Content covered by this course

# 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?

- Content covered by this course

# 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         |

# Text Data

- "Text mining, also referred to as text data mining, roughly equivalent to text analytics, refers to the process of deriving high-quality information from text. High-quality information is typically derived through the devising of patterns and trends through means such as statistical pattern learning. Text mining usually involves the process of structuring the input text (usually parsing, along with the addition of some derived linguistic features and the removal of others, and subsequent insertion into a database), deriving patterns within the structured data, and finally evaluation and interpretation of the output. 'High quality' in text mining usually refers to some combination of relevance, novelty, and interestingness. Typical text mining tasks include text categorization, text clustering, concept/entity extraction, production of granular taxonomies, sentiment analysis, document summarization, and entity relation modeling (i.e., learning relations between named entities)." –from wiki

# 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
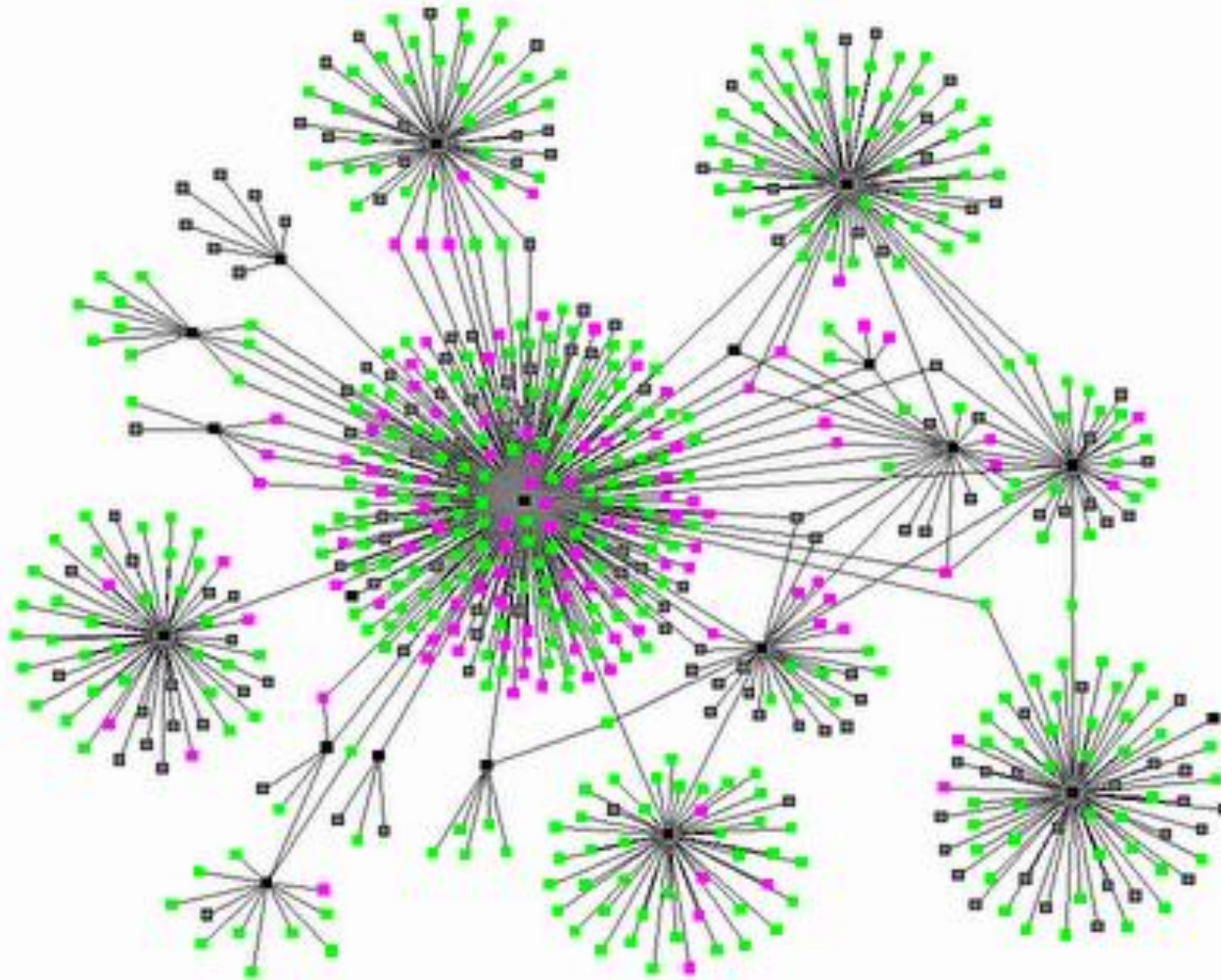
# Graph / Network

# Image Data

# 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?

- Content covered by this course

# 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 Functions: Others

- Prediction

- Similarity search

- Ranking

- Outlier detection

- …

# 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?

- Content covered by this course

# 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?

- Content covered by this course

# 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

# Google Flu Trends

- https://www.youtube.com/watch?v=6111nS66Dpk

Annual U.S. Flu Activity - Mid-Atlantic Region

ILI percentage

● Google Flu Trends   ● CDC Data

# NetFlix Prize

- https://www.youtube.com/watch?v=4_e2sNYYfxA

# Facebook MyPersonality App

- https://www.youtube.com/watch?v=GOZArvMMHKs

## Private traits and attributes are predictable from digital records of human behavior

Michal Kosinski[a,1], David Stillwell[a], and Thore Graepel[b]

[a]Free School Lane, The Psychometrics Centre, University of Cambridge, Cambridge CB2 3RQ United Kingdom; and [b]Microsoft Research, Cambridge CB1 2FB, United Kingdom

We show that easily accessible digital records of behavior, Facebook Likes, can be used to automatically and accurately predict a range of highly sensitive personal attributes including: sexual orientation, ethnicity, religious and political views, personality traits, intelligence, happiness, use of addictive substances, parental separation, age, and gender. The analysis presented is based on a dataset of over 58,000 volunteers who provided their Facebook Likes, detailed demographic profiles, and the results of several psychometric tests. The proposed model uses dimensionality reduction for preprocessing the Likes data, which are then entered into logistic/linear regression to predict individual psychodemographic profiles from Likes. The model correctly discriminates between homosexual and heterosexual men in 88% of cases, African Americans and Caucasian Americans in 95% of cases, and between Democrat and Republican in 85% of cases. For the personality trait "Openness," prediction accuracy is close to the test–retest accuracy of a standard personality test. We give examples of associations between attributes and Likes and discuss implications for online personalization

browsing logs (11–15). Similarly, it has been shown that personality can be predicted based on the contents of personal Web sites (16), music collections (17), properties of Facebook or Twitter profiles such as the number of friends or the density of friendship networks (18–21), or language used by their users (22). Furthermore, location within a friendship network at Facebook was shown to be predictive of sexual orientation (23).

This study demonstrates the degree to which relatively basic digital records of human behavior can be used to automatically and accurately estimate a wide range of personal attributes that people would typically assume to be private. The study is based on Facebook Likes, a mechanism used by Facebook users to express their positive association with (or "Like") online content, such as photos, friends' status updates, Facebook pages of products, sports, musicians, books, restaurants, or popular Web sites. Likes represent a very generic class of digital records, similar to Web search queries, Web browsing histories, and credit card purchases. For example, observing users' Likes related to music

41

# 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?

- Content covered by this course

# Course Content

- By data types:
  - matrix data
  - text data
  - set data
  - sequence data
  - time series
  - graph and network
  - Image data
- By functions:
  - Classification
  - Clustering
  - Frequent pattern mining
  - Prediction
  - Similarity search
  - Ranking

# Methods to Learn

| | Matrix Data | Text Data | Set Data | Sequence Data | Time Series | Graph & Network | Images |
|---|---|---|---|---|---|---|---|
| **Classification** | Decision Tree; Naïve Bayes; Logistic Regression SVM; kNN | | | HMM | | Label Propagation* | Neural Network |
| **Clustering** | K-means; hierarchical clustering; DBSCAN; Mixture Models; kernel k-means* | PLSA | | | | SCAN*; Spectral Clustering* | |
| **Frequent Pattern Mining** | | | Apriori; FP-growth | GSP; PrefixSpan | | | |
| **Prediction** | Linear Regression | | | | Autoregression | | |
| **Similarity Search** | | | | | DTW | P-PageRank | |
| **Ranking** | | | | | | PageRank | |

# Where to Find References? DBLP, CiteSeer, Google

- <u>Data mining and KDD</u> (SIGKDD: CDROM)
  - Conferences: ACM-SIGKDD, IEEE-ICDM, SIAM-DM, PKDD, PAKDD, etc.
  - Journal: Data Mining and Knowledge Discovery, KDD Explorations, ACM TKDD
- <u>Database systems</u> (SIGMOD: ACM SIGMOD Anthology—CD ROM)
  - 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, 3rd ed. , 2011

- T. Hastie, R. Tibshirani, and J. Friedman, The Elements of Statistical Learning: Data Mining, Inference, and Prediction, 2nd 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, 2nd ed. 2005