



Introduction
to Data
Mining

Jun Huang

Introduction to Data Mining

Lecture1 Introduction

Jun Huang

Anhui University of Technology

Spring 2018

huangjun_cs@163.com



Welcome

Introduction
to Data
Mining

Jun Huang

Welcome

Syllabus

Introduction

What
Motivated
Data Mining?

What is Data
Mining?

Patterns of
Data Mining

What Kinds of
Data?

KDD Process

Research
Issues

- Jun Huang

- Technology of Computer Application, Ph.D
- University of Chinese Academy of Sciences
- Research Interests: Machine Learning, Data Mining, Multimedia Content Analysis, etc.

- Office: Room 227, Yifu Building

- Phone: 17353766628

- Wechat





Syllabus(Tentative)

Introduction
to Data
Mining

Jun Huang

Welcome

Syllabus

Introduction

What
Motivated
Data Mining?

What is Data
Mining?

Patterns of
Data Mining

What Kinds of
Data?

KDD Process

Research
Issues

- Introduction
- Data pre-processing
- Classification
- Clustering
- Association rules
- Advanced Topics in Data Mining and Machine Learning
 - Multi-Label Classification
 - Multi-View/Modal Learning
 - Multiple Clustering
 - Paralleled and Distributed Computing



Introduction

Introduction
to Data
Mining

Jun Huang

Welcome

Syllabus

Introduction

What
Motivated
Data Mining?

What is Data
Mining?

Patterns of
Data Mining

What Kinds of
Data?

KDD Process

Research
Issues

- Textbook and references
- Prerequisites
- Grading scheme
- What is data mining?
- Tasks of data mining
- Procedure of data mining



Textbook and References

Introduction
to Data
Mining

Jun Huang

Welcome

Syllabus

Introduction

What
Motivated
Data Mining?

What is Data
Mining?

Patterns of
Data Mining

What Kinds of
Data?

KDD Process

Research
Issues

● Textbook

- The Top 10 Algorithms in Data Mining. Xindong Wu, Kumar Vipin. Chinese Edition (translated by Wenbo Li and Suyan Wu), Tsinghua University Press.
- Data Mining, Concepts and Techniques. Jiawei Han and Micheline Kamber, Morgan Kaufmann, 2011(Third Edition).

● References

- Research paper. To be announced in class.



Prerequisites

Introduction
to Data
Mining

Jun Huang

Welcome

Syllabus

Introduction

What
Motivated
Data Mining?

What is Data
Mining?

Patterns of
Data Mining

What Kinds of
Data?

KDD Process

Research
Issues

- Data structure
- Algorithm
- Database
- Programming skills: Python, Java, Matlab,.etc.
- A mini survey
 - How many people were major in computer science?
 - How many people took machine learning courses before?
 - How many people took database courses before?



Grading Scheme

Introduction
to Data
Mining

Jun Huang

Welcome

Syllabus

Introduction

What
Motivated
Data Mining?

What is Data
Mining?

Patterns of
Data Mining

What Kinds of
Data?

KDD Process

Research
Issues

- Assignments (20%)
- Course Project (30%)
 - One Project
 - Develop an algorithm and hand in a project report
 - To be evaluated in technical innovation, performance, thoroughness of the work, clarity of presentation
- Final Exam (50%)



About the Project

Introduction
to Data
Mining

Jun Huang

Welcome

Syllabus

Introduction

What
Motivated
Data Mining?

What is Data
Mining?

Patterns of
Data Mining

What Kinds of
Data?

KDD Process

Research
Issues

- Choose a topic from a list of selected topics
- Read through some related research papers and fully understand them
- Implement and experimentally evaluate the major method
- Identify advantages and disadvantages
- Improve the method in effectiveness or efficiency, implement and experimentally evaluate your improvement(plus)
- Write a technical report and give presentation



How to do a good project

Introduction
to Data
Mining

Jun Huang

Welcome

Syllabus

Introduction

What
Motivated
Data Mining?

What is Data
Mining?

Patterns of
Data Mining

What Kinds of
Data?

KDD Process

Research
Issues

- Start early
 - It takes time to understand and think
- Discuss with me
 - Maybe I can give some suggestions or ideas
- Implement concretely
 - Understand the advantages and disadvantages
- Think creatively



Why take this course?

Introduction
to Data
Mining

Jun Huang

Welcome

Syllabus

Introduction

What
Motivated
Data Mining?

What is Data
Mining?

Patterns of
Data Mining

What Kinds of
Data?

KDD Process

Research
Issues

- Data Mining is so hot
- Solve many interesting problems in real applications, e.g. business management, WWW, science exploration
- Turn raw data into knowledge
- Promising in research of many disciplines
- Data miners' job market: many well-paid positions



Objectives of this course

Introduction
to Data
Mining

Jun Huang

Welcome

Syllabus

Introduction

What
Motivated
Data Mining?

What is Data
Mining?

Patterns of
Data Mining

What Kinds of
Data?

KDD Process

Research
Issues

- Introduce the motivation of data mining
- Outline principles, major algorithms
- Introduce applications
- Introduce advanced topics
- Enhance independent research capability



Policies

Introduction
to Data
Mining

Jun Huang

Welcome

Syllabus

Introduction

What
Motivated
Data Mining?

What is Data
Mining?

Patterns of
Data Mining

What Kinds of
Data?

KDD Process

Research
Issues

- Students are expected to attend all classes
- No late homework will be accepted
- All work must be efforts of your own
- **No Plagiarism!**



What motivated data mining?

Introduction
to Data
Mining

Jun Huang

Welcome

Syllabus

Introduction

What
Motivated
Data Mining?

What is Data
Mining?

Patterns of
Data Mining

What Kinds of
Data?

KDD Process

Research
Issues

- The explosive growth of data
- Data collection and data availability
 - Computer hardware & software develop dramatically
 - The amount of data collected and stored doubles/triples per year vs. CPU speed increases 15% per year(till 2003)
- Many types of databases
 - Object-oriented, spatial, temporal, time-series, text, multimedia, Web



What motivated data mining?

Business World

Introduction to Data Mining

Jun Huang

Welcome

Syllabus

Introduction

What Motivated Data Mining?

What is Data Mining?

Patterns of Data Mining

What Kinds of Data?

KDD Process

Research Issues

- Tremendous of data being collected and stored
 - E-commerce
 - Transactions
 - Stocks
 - Credit card transactions
 - Strong competitive pressure to extract and use knowledge hidden in the data to provide customized CRM





What motivated data mining?

Scientific World

Introduction
to Data
Mining

Jun Huang

Welcome

Syllabus

Introduction

What
Motivated
Data Mining?

What is Data
Mining?

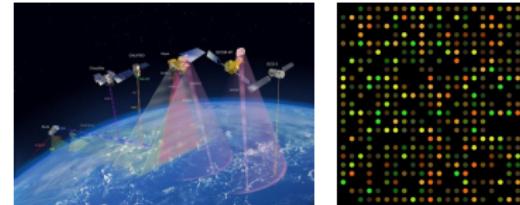
Patterns of
Data Mining

What Kinds of
Data?

KDD Process

Research
Issues

- Tremendous of data being collected and stored
 - Remote sensing
 - Bioinformatics(Microarrays)
 - Scientific simulation
- Scientists need strong data analysis to assist research, such as classification, segmetation, etc.





What motivated data mining?

Introduction
to Data
Mining

Jun Huang

Welcome

Syllabus

Introduction

What
Motivated
Data Mining?

What is Data
Mining?

Patterns of
Data Mining

What Kinds of
Data?

KDD Process

Research
Issues

- We are drowning in data, but starving for knowledge
 - Data rich, knowledge poor
 - Decision makers, domain experts have biases or errors
- Automated analysis of massive data sets



What is data mining?

Introduction
to Data
Mining

Jun Huang

Welcome

Syllabus

Introduction

What
Motivated
Data Mining?

What is Data
Mining?

Patterns of
Data Mining

What Kinds of
Data?

KDD Process

Research
Issues

- **Data Mining:** Discover valid, novel, useful, and understandable patterns in massive database





What is data mining?

Introduction
to Data
Mining

Jun Huang

Welcome

Syllabus

Introduction

What
Motivated
Data Mining?

What is Data
Mining?

Patterns of
Data Mining

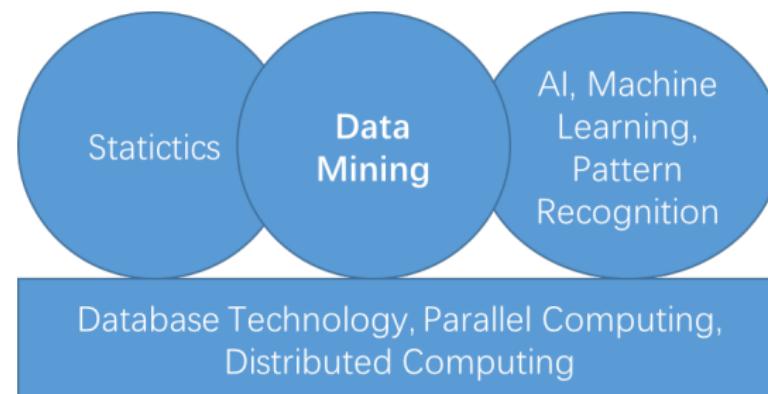
What Kinds of
Data?

KDD Process

Research
Issues

- Cross disciplines

- Databases
- Machine Learning: decision tree, Bayesian Classifier, etc.
- Statistics: regression, etc.
- Neural networks





Why not Traditional Data Analysis?

Introduction
to Data
Mining

Jun Huang

Welcome

Syllabus

Introduction

What
Motivated
Data Mining?

What is Data
Mining?

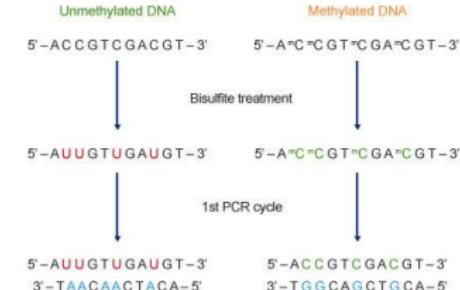
Patterns of
Data Mining

What Kinds of
Data?

KDD Process

Research
Issues

- Tremendous amount of data
 - Algorithms must be highly scalable to handle such as tera-bytes of data
- High-dimensionality of data
 - DNA sequences may have tens of thousands of dimensions





Why not traditional data analysis?

Introduction
to Data
Mining

Jun Huang

Welcome

Syllabus

Introduction

What
Motivated
Data Mining?

What is Data
Mining?

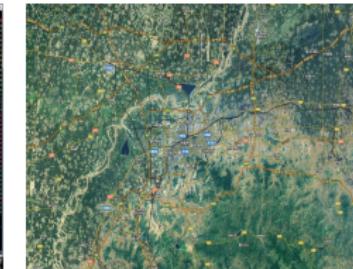
Patterns of
Data Mining

What Kinds of
Data?

KDD Process

Research
Issues

- High complexity of data
 - Data streams and sensor data
 - Time-series data, temporal data, sequence data
 - Graphs, social networks
 - Spatial, temporal, multimedia, text, and web data
- New and sophisticated applications





Why not traditional data analysis?

Introduction
to Data
Mining

Jun Huang

Welcome

Syllabus

Introduction

What
Motivated
Data Mining?

What is Data
Mining?

Patterns of
Data Mining

What Kinds of
Data?

KDD Process

Research
Issues

Database	Data Mining
Storage-oriented	Discover knowledge from data in databases
Provide simple queries	
Data warehouse	Data Mining
Subject-oriented	Advanced data analysis tools
A multidimensional view of data	
Operations to access summarized data	
Statistical algorithms	Data Mining
Based on many hypothesis	Less hypothesis; Find patterns in large number of samples; Abnormal patterns
Find patterns in small number of samples	



Characteristics of Data Mining

Introduction
to Data
Mining

Jun Huang

Welcome

Syllabus

Introduction

What
Motivated
Data Mining?

What is Data
Mining?

Patterns of
Data Mining

What Kinds of
Data?

KDD Process

Research
Issues

- Massive dataset
- Automatically searching for interesting patterns from history data
- Fast
- Scalable
- Update easily
- Practical
- Decision support



Exercises

Introduction
to Data
Mining

Jun Huang

Welcome

Syllabus

Introduction

What
Motivated
Data Mining?

What is Data
Mining?

Patterns of
Data Mining

What Kinds of
Data?

KDD Process

Research
Issues

- ① Could you present an application of data mining in business domain?
- ② Could you present an application of data mining in scientific domain?



What kinds of Patterns?

Association rules

Introduction
to Data
Mining

Jun Huang

Welcome

Syllabus

Introduction

What
Motivated
Data Mining?

What is Data
Mining?

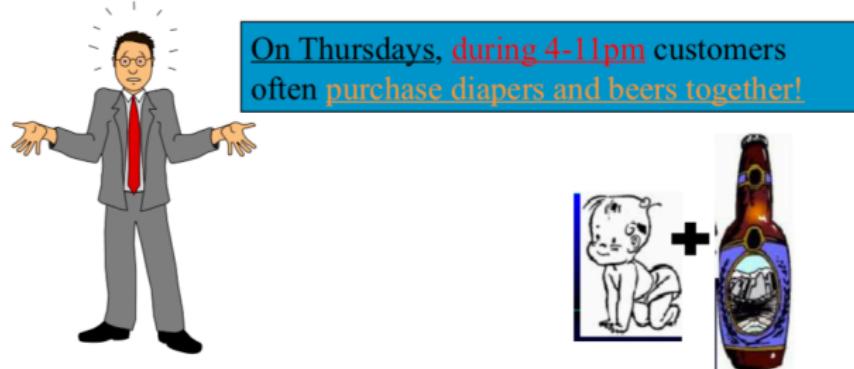
Patterns of
Data Mining

What Kinds of
Data?

KDD Process

Research
Issues

- **Association rules:** Detect sets of attributes or items that frequently co-occur in many database records and rules among them





Ex.1 Market basket analysis and management

Introduction
to Data
Mining

Jun Huang

Welcome

Syllabus

Introduction

What
Motivated
Data Mining?

What is Data
Mining?

Patterns of
Data Mining

What Kinds of
Data?

KDD Process

Research
Issues

- Where does the data come from?
 - supermarket transactions, membership cards, discount coupons, customer complaint calls
- Cross-marketing analysis
 - What products were often purchased together? Purchase recommendation, cross selling
 - What are the subsequent purchases after buying a given product?
- Target-marketing
 - What types of customers buy what products
- Catalog design





What kinds of Patterns?

Classification

Introduction
to Data
Mining

Jun Huang

Welcome

Syllabus

Introduction

What
Motivated
Data Mining?

What is Data
Mining?

Patterns of
Data Mining

What Kinds of
Data?

KDD Process

Research
Issues

- **Classification:** Build a model of classes on training dataset, and then, assign a new record to one of several predefined classes



- Decision Tree
- Rule1: if **Income<=40K** and **Debt=0** then **good**
- Rule2: if **Income>40K** and **Debt < 10% of Income** then **good**



Ex.2 Credit Scoring

Introduction
to Data
Mining

Jun Huang

Welcome

Syllabus

Introduction

What
Motivated
Data Mining?

What is Data
Mining?

Patterns of
Data Mining

What Kinds of
Data?

KDD Process

Research
Issues

- Where does the data come from?
 - credit card transactions, credit card payments, loan payments, demographic data
- Predict the probability to bankrupt or charge-off
- Reduce the credit risk to the banks
- Increase the profitability of the banks



What kinds of Patterns?

Clustering

Introduction
to Data
Mining

Jun Huang

Welcome

Syllabus

Introduction

What
Motivated
Data Mining?

What is Data
Mining?

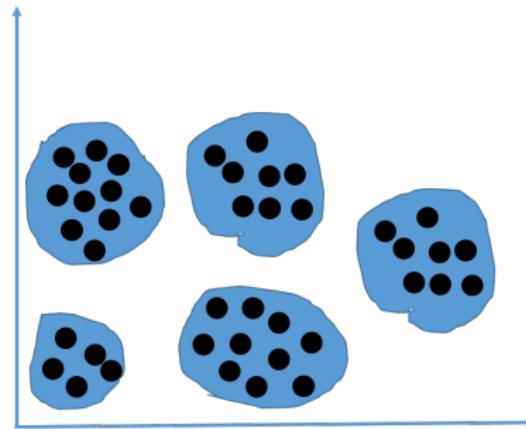
Patterns of
Data Mining

What Kinds of
Data?

KDD Process

Research
Issues

- **Clustering:** Partition the dataset into groups such that elements in a group have lower inter-group similarity and higher intra-group similarity





Ex.3 Scientific simulation

Introduction
to Data
Mining

Jun Huang

Welcome

Syllabus

Introduction

What
Motivated
Data Mining?

What is Data
Mining?

Patterns of
Data Mining

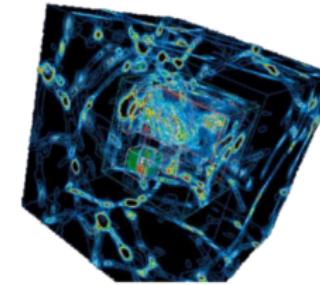
What Kinds of
Data?

KDD Process

Research
Issues

● Cosmological simulation

- Simulate the formation of the galaxy
- Enormous particles at each evolution stage, beyond the capability of human being to analyze





What kinds of Patterns?

Sequence mining

Introduction
to Data
Mining

Jun Huang

Welcome

Syllabus

Introduction

What
Motivated
Data Mining?

What is Data
Mining?

Patterns of
Data Mining

What Kinds of
Data?

KDD Process

Research
Issues

- **Sequence mining:** Give a set of sequences, find the complete set of frequent subsequences



Marking strategy: recommend a new CPU for the customer 9 months after his first purchase



What kinds of Patterns?

Anomaly detection

Introduction
to Data
Mining

Jun Huang

Welcome

Syllabus

Introduction

What
Motivated
Data Mining?

What is Data
Mining?

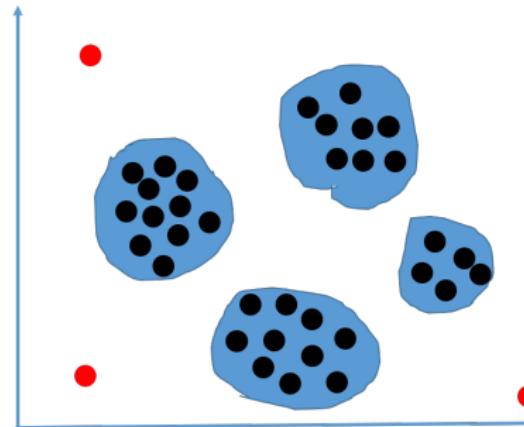
Patterns of
Data Mining

What Kinds of
Data?

KDD Process

Research
Issues

- **Anomaly detection:** Give a set of n objects, and k , the number of expected anomalies, find the top k objects that are considerably dissimilar or inconsistent with the remaining data





What kinds of Patterns?

Community Analysis

Introduction
to Data
Mining

Jun Huang

Welcome

Syllabus

Introduction

What
Motivated
Data Mining?

What is Data
Mining?

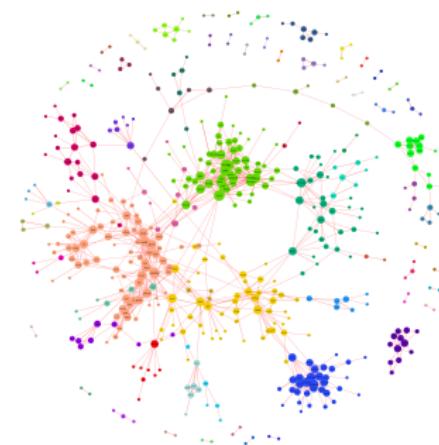
Patterns of
Data Mining

What Kinds of
Data?

KDD Process

Research
Issues

- **Community Analysis:** In social media mining, analyzing communities is essential.
 - How to detect communities?
 - How do communities evolve and how to study evolving communities?
 - How to evaluate detected communities?





What kinds of Patterns?

Recommender systems

Introduction
to Data
Mining

Jun Huang

Welcome

Syllabus

Introduction

What
Motivated
Data Mining?

What is Data
Mining?

Patterns of
Data Mining

What Kinds of
Data?

KDD Process

Research
Issues

● Recommender systems:

- Recommend products that would be interesting to individuals
- Build a function, $f: U \times I \rightarrow \mathbb{R}$, for user set U and item set I

Product



amazon
JD.com
天猫 Tmall.com

Customers Who Viewed This Item Also Viewed



Nivea UV Whitening Extra
Moisture Antiperspirant
Deodorant Roll-On 50ml
★★★★★ (142)
\$8.33



Nivea UV Whitening Extra
Cell Repair and Renewal
Body Lotion 400ml
★★★★★ (14)
\$20.00



Nivea UV Whitening Extra
Cell Repair Milk Repair
200ml
★★★★★ (8)
\$5.95



热映推荐
一路·英·冲·梦
中国式立功歌合集
向南的心的爱歌

Movie

Music



Exercises

Introduction
to Data
Mining

Jun Huang

Welcome

Syllabus

Introduction

What
Motivated
Data Mining?

What is Data
Mining?

Patterns of
Data Mining

What Kinds of
Data?

KDD Process

Research
Issues

- ① Please present an example where data mining is crucial to the success of the business. What data mining techniques are used? (What kinds of patterns are mined?)
- ② Can you describe other possible kind of knowledge that needs to be discovered by data mining methods but not been mentioned in class yet?



On What Kinds of Data?

Introduction
to Data
Mining

Jun Huang

Welcome

Syllabus

Introduction

What
Motivated
Data Mining?

What is Data
Mining?

Patterns of
Data Mining

What Kinds of
Data?

KDD Process

Research
Issues

① Database-oriented data sets and applications

- Relational database, data warehouse, transactional database

② Advanced database applications

- Data streams
- Spatial data
- Text database
- Multimedia data
- Time-series
- Bio-medical data
- Network traffic data



Relational Databases

Introduction
to Data
Mining

Jun Huang

Welcome

Syllabus

Introduction

What
Motivated
Data Mining?

What is Data
Mining?

Patterns of
Data Mining

What Kinds of
Data?

KDD Process

Research
Issues

- Structured data

- Table-records-attributes
- Accessed by queries, SQL

- Online transactional processing (OLTP)

- Insert a student “Jun Huang” into class “Introduction to Data Mining”, Spring 2018

Name	Time	Course	Score	Room
Jun Huang	Spring 2018	Introduction to Data Mining	90	002
Lucy	Spring 2017	Introduction to Data Mining	70	002
Merlisa	Spring 2017	Math	80	001
Jack	Fall 2016	Machine Learning	80	001



Data Warehouse

Introduction
to Data
Mining

Jun Huang

Welcome

Syllabus

Introduction

What
Motivated
Data Mining?

What is Data
Mining?

Patterns of
Data Mining

What Kinds of
Data?

KDD Process

Research
Issues

- A subject-oriented, integrated, cleaned collection of data in support of management's decision making process
- Data from multiple databases
- Consistency checking in data warehouses
- Data warehouses can answer OLAP queries efficiently
 - Online analytical processing (OLAP)
 - Find the average class score of "Jun Huang" in the last 3 years, grouped by semesters
- Many patterns are summarization of data
 - Roll-up, drill-down



Data Warehouse

Introduction
to Data
Mining

Jun Huang

Welcome

Syllabus

Introduction

What
Motivated
Data Mining?

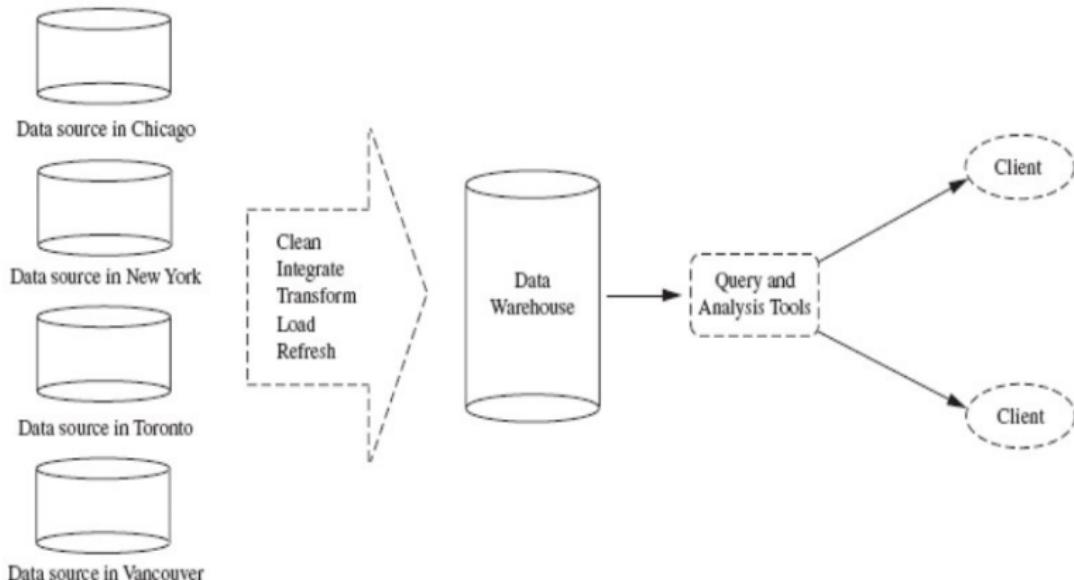
What is Data
Mining?

Patterns of
Data Mining

What Kinds of
Data?

KDD Process

Research
Issues





Transactional Databases

Introduction
to Data
Mining

Jun Huang

Welcome

Syllabus

Introduction

What
Motivated
Data Mining?

What is Data
Mining?

Patterns of
Data Mining

What Kinds of
Data?

KDD Process

Research
Issues

- $I = \{x_1, x_2, \dots, x_n\}$ is the set of items
- An **itemset** is a subset of I
- A **transaction** is a tuple (tid, X)
 - tid : Transaction Id
 - X : Itemset
- A **transactional database** is a set of transactions

Tid	Itemset
T100	milk, bread, beer, diaper
T200	beer, cook, fish, potato, orange, apple
...	...



Spatial Data

Introduction
to Data
Mining

Jun Huang

Welcome

Syllabus

Introduction

What
Motivated
Data Mining?

What is Data
Mining?

Patterns of
Data Mining

What Kinds of
Data?

KDD Process

Research
Issues

● Spatial Information

- Geographical database
- VLSI chip design databases
- Satellite/remote sensing image database
- Medical image database

● Spatial Patterns

- Find characteristics of homes near a give location
- Change in trend of metropolitan poverty rates based on distances from major highways

编号	中心	正右方	右上方	面积
1	居民地	绿地	水体	120
2	绿地	水体	水体	40
3	水体	居民地	居民地	610
...



Time Series

Introduction
to Data
Mining

Jun Huang

Welcome

Syllabus

Introduction

What
Motivated
Data Mining?

What is Data
Mining?

Patterns of
Data Mining

What Kinds of
Data?

KDD Process

Research
Issues

- A sequence of values that change over time
 - Sequences of stock price at every 5 minutes
 - Daily temperature
 - Power supply
 - Electrocardiogram
- Typical Operations
 - Similarity search
 - Trend analysis
 - Periodic pattern discovery





Text Databases & Multimedia Databases

Introduction
to Data
Mining

Jun Huang

Welcome

Syllabus

Introduction

What
Motivated
Data Mining?

What is Data
Mining?

Patterns of
Data Mining

What Kinds of
Data?

KDD Process

Research
Issues

- HTML web documents
- XML documents
- Digital libraries
- Annotated multimedia databases
 - Image, audio and video data
- Typical Operations
 - Similarity-based pattern matching
 - Image classification



电影排行榜

1. 《夏洛特烦恼》
2. 《捉妖记》

解忧杂货店

记忆大师

INFORMATION RESEARCH
an international electronic journal
ISSN 1360-1613

Volume 29 No 4 December, 2013

Editorial

Editor
The impact of information technology capability, information sharing and knowledge management on the operational performance of emergency incident management systems

Mihirshi Saito-Malina, Nenca Karapicikci, Perizzi Yalcin and Katerina Zafeiri
Linking empathy effects information users on a questionnaire study in city administration

David Jofre-Sastre and Manuel Sánchez-Pérez
Developing an e-government capacity in measurement scale

Jung Sun Hahn, De Yoon Kim, Heon Choi Kim and Min Sung Son
Simplification into the extension of the lexicon effect in key phrase extraction

Orla Kiely
A comparative theoretical exploration of information sharing and trust in a dispersed community of design scholars

David P. Nettleton, Ricardo Basurto-Tobas and Mart-Carrión Marzo
Information retrieval and document classification in terms of the Library of Congress classification and key publishers

Hilke Mev
Identifying information-seeking behaviour of highly dedicated online research visitors



Data Streams

Introduction
to Data
Mining

Jun Huang

Welcome

Syllabus

Introduction

What
Motivated
Data Mining?

What is Data
Mining?

Patterns of
Data Mining

What Kinds of
Data?

KDD Process

Research
Issues

- Data in the form of continuous arrival in multiple, rapid, time-varying, possibly unpredictable and unbounded streams
 - Dynamically changing patterns, high volume, infinite, quick response, no re-scan
- Many applications
 - Stock exchange, network monitoring, telecommunications data management, web application, sensor network, etc.



Biomedical Data

Introduction
to Data
Mining

Jun Huang

Welcome

Syllabus

Introduction

What
Motivated
Data Mining?

What is Data
Mining?

Patterns of
Data Mining

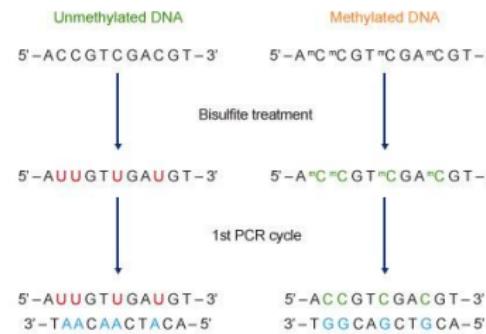
What Kinds of
Data?

KDD Process

Research
Issues

● Bio-sequences

- DNA; very long sequences of nucleotides
- Similarity search
- Identify sequential patterns that play roles in various diseases
- Association analysis: co-occurring gene





World-Wide Web

Introduction
to Data
Mining

Jun Huang

Welcome

Syllabus

Introduction

What
Motivated
Data Mining?

What is Data
Mining?

Patterns of
Data Mining

What Kinds of
Data?

KDD Process

Research
Issues

- The WWW is huge, widely distributed, global information service center for
 - Information service: news, advertisements, consumer information, financial management, education, government, e-commerce, etc.
 - Hyper-link information
 - Access and usage information
- WWW provides rich sources for data mining
- Challenges
 - Too huge for effective data warehousing and data mining
 - Too complex and heterogeneous: no standards and structure



World-Wide Web

Introduction
to Data
Mining

Jun Huang

Welcome

Syllabus

Introduction

What
Motivated
Data Mining?

What is Data
Mining?

Patterns of
Data Mining

What Kinds of
Data?

KDD Process

Research
Issues

- Web-Usage: Logs and IP package head streams
 - Mining weblog records to discover users accessing patterns of Web pages
- Web Content
 - Extract knowledge from Web documents, automatic categorization
- Web Structure
 - Identifying interesting graph patterns among different Web pages



Graph

Introduction
to Data
Mining

Jun Huang

Welcome

Syllabus

Introduction

What
Motivated
Data Mining?

What is Data
Mining?

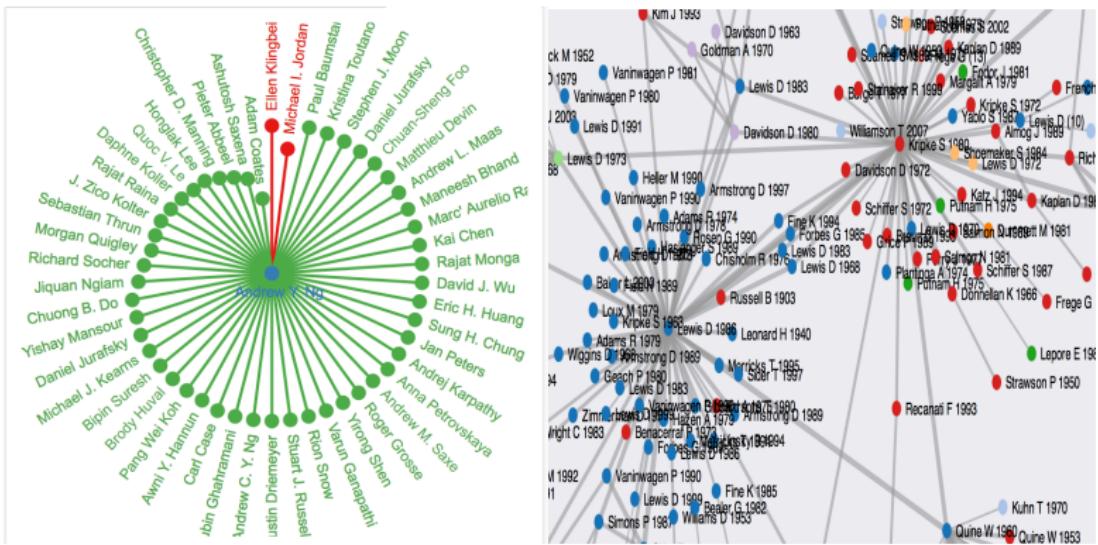
Patterns of
Data Mining

What Kinds of
Data?

KDD Process

Research
Issues

● Citation Graph





Graph

Introduction to Data Mining

Jun Huang

Welcome

Syllabus

Introduction

What Motivated Data Mining?

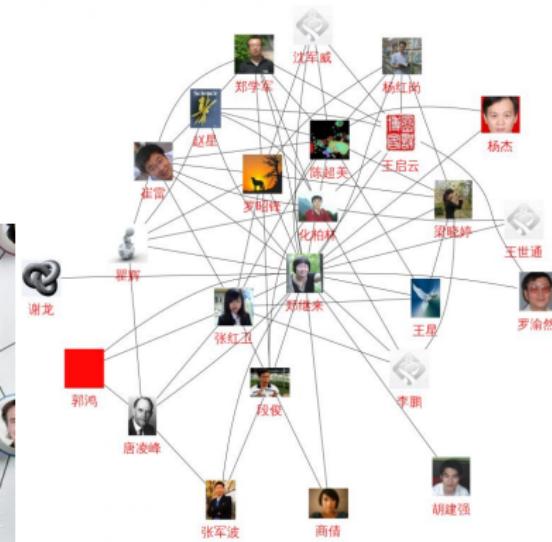
What is Data Mining?

Patterns of Data Mining

What Kinds of Data?

KDD Process

Research Issues





Graph

Introduction to Data Mining

Jun Huang

Welcome

Syllabus

Introduction

What Motivated Data Mining?

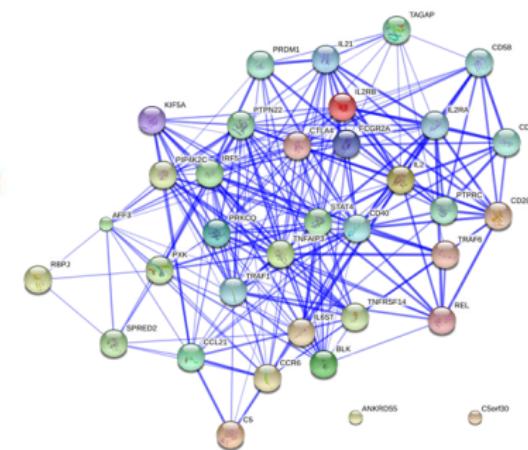
What is Data Mining?

Patterns of Data Mining

What Kinds of Data?

KDD Process

Research Issues





Graph

Graph Mining

Introduction
to Data
Mining

Jun Huang

Welcome

Syllabus

Introduction

What
Motivated
Data Mining?

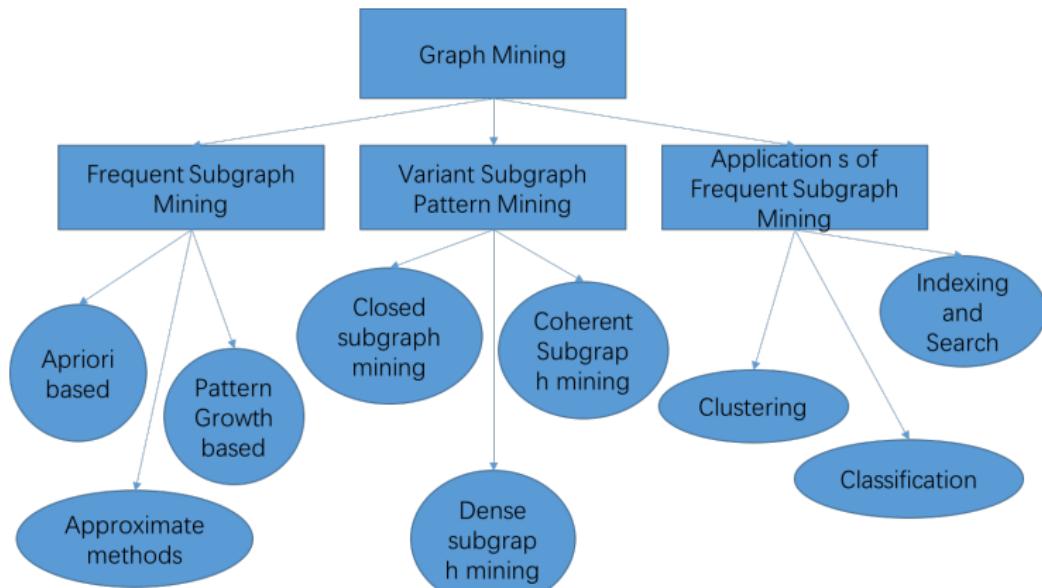
What is Data
Mining?

Patterns of
Data Mining

What Kinds of
Data?

KDD Process

Research
Issues





Knowledge Discovery(KDD) Process

Introduction
to Data
Mining

Jun Huang

Welcome

Syllabus

Introduction

What
Motivated
Data Mining?

What is Data
Mining?

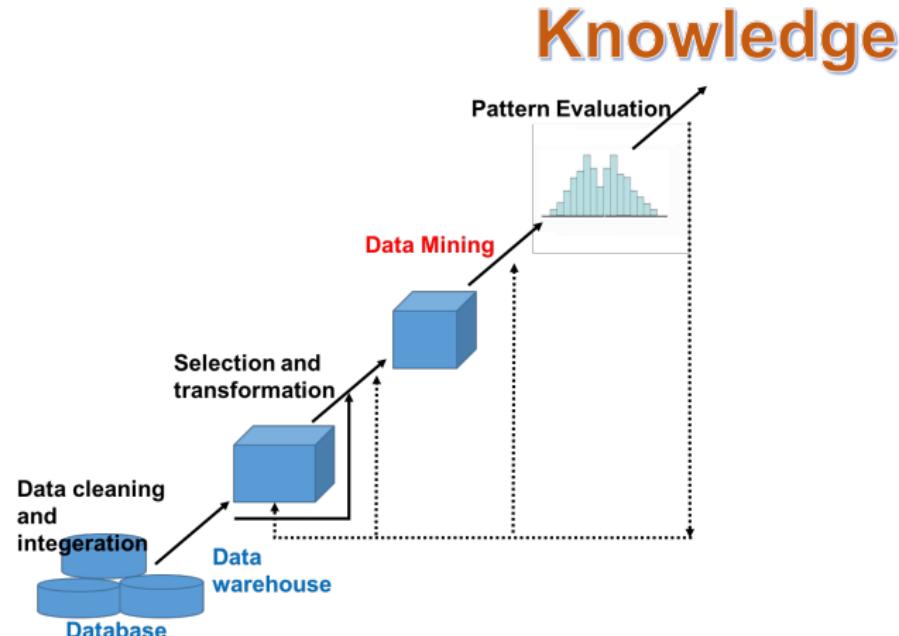
Patterns of
Data Mining

What Kinds of
Data?

KDD Process

Research
Issues

- Data Mining-Core of knowledge discovery process





Key Steps in KDD Process

Introduction
to Data
Mining

Jun Huang

Welcome

Syllabus

Introduction

What
Motivated
Data Mining?

What is Data
Mining?

Patterns of
Data Mining

What Kinds of
Data?

KDD Process

Research
Issues

- Learning the application domain
 - Relevant prior knowledge and goals of application
- Creating a target data resource
- Data cleaning and preprocessing: (may take 60% of effort!)
- Data reduction and transformation
 - Find useful features, dimensionality/variable reduction, invariant representation
- Choosing the mining algorithm(s) to search for patterns of interest
- Pattern evaluation and knowledge presentation
 - visualization, transformation, removing redundant patterns, etc.
- Use of discovered knowledge



Are All the “Discovered” Patterns Interesting?

Introduction
to Data
Mining

Jun Huang

Welcome

Syllabus

Introduction

What
Motivated
Data Mining?

What is Data
Mining?

Patterns of
Data Mining

What Kinds of
Data?

KDD Process

Research
Issues

- Data mining may generate thousands of patterns: Not all of them are interesting
- Interestingness measures
 - A pattern is interesting if it is easily understood by humans, valid on new or test data with some degree of certainty, potentially useful, novel, or validates some hypothesis that a user seeks to confirm
- Objective vs. subjective interestingness measures
 - Objective: based on statistics and structures of patterns, e.g., support, confidence, etc.
 - Subjective: based on user's belief in the data, e.g., unexpectedness, novelty, actionability, etc.



Find All and Only Interesting Patterns?

Introduction
to Data
Mining

Jun Huang

Welcome

Syllabus

Introduction

What
Motivated
Data Mining?

What is Data
Mining?

Patterns of
Data Mining

What Kinds of
Data?

KDD Process

Research
Issues

- Find all the interesting patterns: Completeness
 - Can a data mining system find all the interesting patterns?
Do we need to find all of the interesting patterns?
 - Heuristic vs. exhaustive searching
- Search for only interesting patterns: An optimization problem - Challenging
 - Can a data mining system find only the interesting patterns?
 - Approaches
 - First general all the patterns and then filter out the uninteresting ones
 - Guide and constrain the discovery process



Research Issues in Data Mining

Introduction
to Data
Mining

Jun Huang

Welcome

Syllabus

Introduction

What
Motivated
Data Mining?

What is Data
Mining?

Patterns of
Data Mining

What Kinds of
Data?

KDD Process

Research
Issues

● Mining methodology

- Mining different kinds of knowledge from diverse data types, e.g., Web, graph, bioinformatics, stream
- Performance: efficiency, effectiveness, and scalability
- Parallel, distributed and incremental mining methods
- Pattern evaluation: the interesting problem
- Handling noise and incomplete data
- Incorporation of background knowledge

● User Interaction

- Data mining query languages
- Expression and visualization of data mining results

● Application and social impacts

- Domain-specific data mining
- Protection of data security, integrity, and privacy



Top 10 Challenging Problems in Data Mining Research

Introduction
to Data
Mining

Jun Huang

Welcome

Syllabus

Introduction

What
Motivated
Data Mining?

What is Data
Mining?

Patterns of
Data Mining

What Kinds of
Data?

KDD Process

Research
Issues

- ① Developing a Unifying Theory of Data Mining
- ② Scaling Up for High Dimensional Data and High Speed Data Streams
- ③ Mining Sequence Data and Time Series Data
- ④ Mining Complex Knowledge from Complex Data
- ⑤ Data Mining in a Network Setting
- ⑥ Distributed Data Mining and Mining Multi-agent Data
- ⑦ Data Mining for Biological and Environmental Problems
- ⑧ Data-Mining-Process Related Problems
- ⑨ Security, Privacy and Data Integrity
- ⑩ Dealing with Non-static, Unbalanced and Cost-sensitive Data



Important Resources

Introduction
to Data
Mining

Jun Huang

Welcome

Syllabus

Introduction

What
Motivated
Data Mining?

What is Data
Mining?

Patterns of
Data Mining

What Kinds of
Data?

KDD Process

Research
Issues

- Data mining Conferences
 - **KDD**, ICDM, WWW, WSDM, SDM, ECML-PKDD, PAKDD
- Database Conferences
 - **SIGMOD**, VLDB, PODS, ICDE, EDBT, ICDT
- Important Journals
 - IEEE Transactions on Knowledge and Data Engineering (TKDE)
 - Knowledge and Information Systems (KAIS)
 - Data Mining and Knowledge Discovery: An International Journal (DMKD)
 - ACM Data Mining and Knowledge Discovery (TKDD)
- Useful Website
 - <https://www.kdnuggets.com/>
 - <https://www.kaggle.com/>