

Data Mining

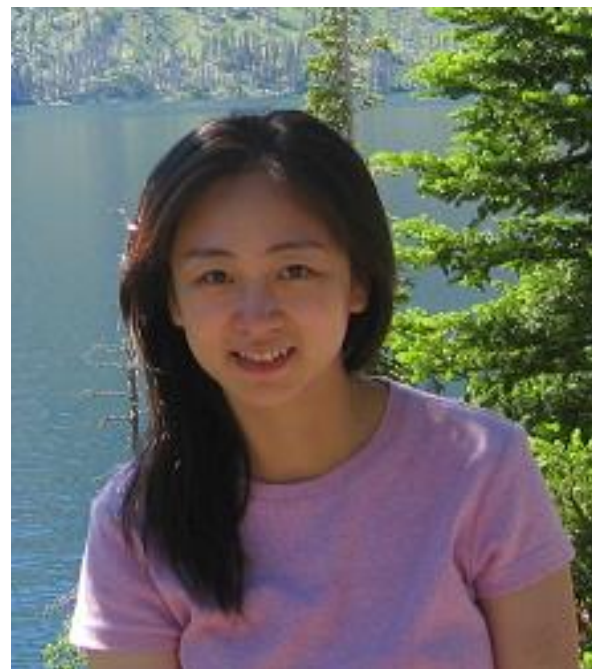
Ying Liu, Prof., Ph.D

*School of Computer Science and Technology
University of Chinese Academy of Sciences
The Key Lab of Big Data Mining and Knowledge Management*

Welcome

■ Ying Liu

- Computer Engineering, Ph.D, Northwestern University, USA
- Research interests
 - Data Mining
 - Artificial Intelligence
 - High Performance Computing
 - Big Data
- Email: yingliu@ucas.ac.cn



Useful Information

- Teaching Assistants

- ?

- ?

- ?

- Class: Monday & Wednesday 8:30 - 10:10, 教
1-101

- Website: <http://sepucas.ac.cn>

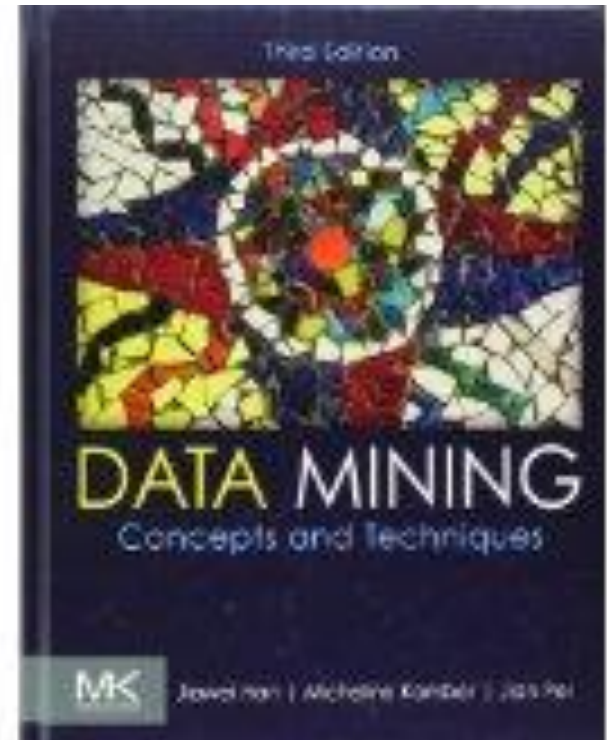
Textbook and References

■ Textbook

- Data Mining, Concepts and Techniques. Jiawei Han and Micheline Kamber, Morgan Kaufmann, 2011 (Third Edition)

■ References

- Research papers. To be announced in class.



Prerequisites

- Data Structure
- Algorithm
- Database
- Programming: C/C++ (preferred), Python, Java

A Mini Survey

- How many people were major in computer science?
- How many people took machine learning courses before?
- How many people took statistics courses before?
- How many people took database courses before?

Grading Scheme

- Assignments (30%)
 - 2 homework assignments
- Course Project (30%)
 - Group project (4 students/group)
 - Solve a real problem: propose an algorithm/approach and implement it
- Final Exam (40%)
 - In class, closed book

About the Project

■ Option 1:

- 2019 CCF大数据与计算智能大赛
(<https://www.datafountain.cn/competitions>)
- Choose a topic from the following topics
- Read through some related research papers and fully understand them
- Develop and Implement the method
- To be evaluated by the ranking or feedback from the contest

中国科学院邮件系统 x 赛题详情(Competition D x +

https://www.datafountain.cn/competitions/352

收藏 JD JD 我的订单

DataFountain 首页 竞赛 AI指北 数据集 | 我要办赛 人工智能教学实验室 (En) 登录 注册

京东

乘用车细分市场销量预测

主办方：中国计算机学会 & 深瞳云涂

预测回归 数据挖掘

队伍 / 人数 1065 / 1177 奖励 ¥100,000

【2019/09/01 20:38:15】DF1567341434998参加了乘用车...

开赛 初赛 08.17 ~ 10.25 A 榜 B 榜 复赛 10.30 ~ 11.11 A 榜 B 榜 决赛 11.23 ~ 11.24 A 榜

赛制规则 数据与评测 排行榜 参赛队伍 赛题讨论 常见问题 报名参赛

大赛介绍 赛题名称

大赛介绍

CCF大数据与计算智能大赛 (CCF Computing Intelligence Contest, 简称CCF BDCI) 是由中国计算机学会大数据专家委员会于2013年创办的

无痕迹浏览 下载

在这里输入你要搜索的内容

21:18 2019/9/1

中国科学院邮件系统 x 赛题详情(Competition D x +

https://www.datafountain.cn/competitions/350

收藏 JD JD 我的订单

DataFountain 首页 竞赛 AI指北 数据集 | 我要办赛 人工智能教学实验室 (En) 登录 注册

京东

互联网新闻情感分析

主办方：中国计算机学会 & 中移软件

自然语言处理 机器学习

队伍 / 人数 1053 / 1135 奖励 ¥ 20,000

【2019/09/01 20:43:51】DF1567341569283参加了互联网...

开赛 初赛 08.17 ~ 10.25 A 榜 B 榜 复赛 10.30 ~ 11.11 A 榜 B 榜 结束

赛制规则 数据与评测 排行榜 参赛队伍 赛题讨论 常见问题 报名参赛

大赛介绍 赛题名称

大赛介绍

CCF大数据与计算智能大赛 (CCF Computing Intelligence Contest, 简称CCF BDCI) 是由中国计算机学会大数据专家委员会于2013年创办的

在这里输入你要搜索的内容

无痕浏览 下载

21:18 2019/9/1

中国科学院邮件系统

赛题详情(Competition D x

+

https://www.datafountain.cn/competitions/362

收藏 JD JD 我的订单

2019

智慧气象服务创新大赛

"神气"大数据算法与应用赛

科技创新 气象惠民

2019

智慧气象服务创新大赛

"神气"大数据算法与应用赛

科技创新 气象惠民

算法组：卫星云图+地面观测云图预测辐照量

主办方：中国气象服务协会 & 中国气象网

图像识别

【2019/09/01 21:07:32】胡搞瞎搞参加了算法组：卫星云图...

队伍 / 人数

29 / 31

奖金

¥ 70,000

开赛

A 榜

初赛 08.30 ~ 10.30

B 榜

A 榜

结束赛

11.12

~

11.12

赛制规则

数据与评测

排行榜

参赛队伍

神气大数据

常见问题

报名参赛

大赛背景

大赛主题

赛程赛制

赛题任务

大赛奖项

大赛背景

第二届智慧气象服务创新大赛——2019"神气"大数据算法与应用赛是由中国气象局指导、中国气象服务协会主办、中国天气网承办的气象领域高规格赛事。通过气象大数据与行业数据的融合，激发高校、科研院所、社会各行业等对气象数据应用的想象力和创新力，促进产学研一体化发展，推动气象数据共享、开放和科研成果的转化与应用，打造气象服务创新平台，推动气象大数据价值链的形成和产业生态发展，提升气象部门的社会影响力，助力气象服务现代化发展。

在这里输入你要搜索的内容

无痕浏览

下载

21:19

2019/9/1

About the Project

■ Option 2:

- Contest in class
- Assigned a topic (to be announced)
- Read through some related research papers and fully understand them
- Develop and Implement the method
- To be evaluated by the ranking in class

How to Do a Good Project?

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

Why Take This Course ?

- Data mining is 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

➤ *Data Mining is very useful!*

Syllabus (Tentative)

- Introduction
- Data warehouse
- Data pre-processing
- Classification
- Association rules
- Clustering
- Applications
- Big data mining

Objectives of This Course

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

Policies

- Students are expected to attend all classes
- No late homework will be accepted
- All work must be efforts of your own (individual assignment) or of your approved team (group assignment)

No Plagiarism!

What Motivated Data Mining?

- 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

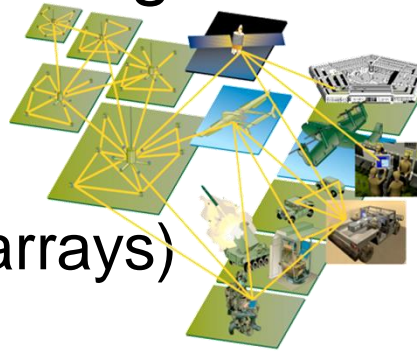
- Tremendous of data being collected and stored
 - E-commerce
 - Transactions
 - Stocks
 - Credit card transactions
- Strong competitive pressure to extract and use the knowledge hidden in the data to provide customized CRM



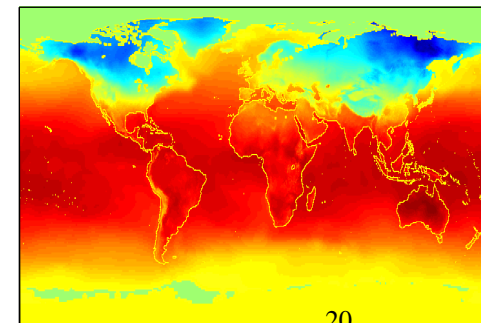
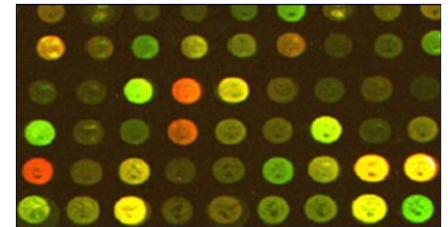
What Motivated Data Mining – Scientific World

- Tremendous of data being collected and stored

- Remote sensing
- Bioinformatics (Microarrays)
- Scientific simulation



- Scientists need strong data analysis to assist research, such as classification, segmentation, etc.



What Motivated Data Mining?

- 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?

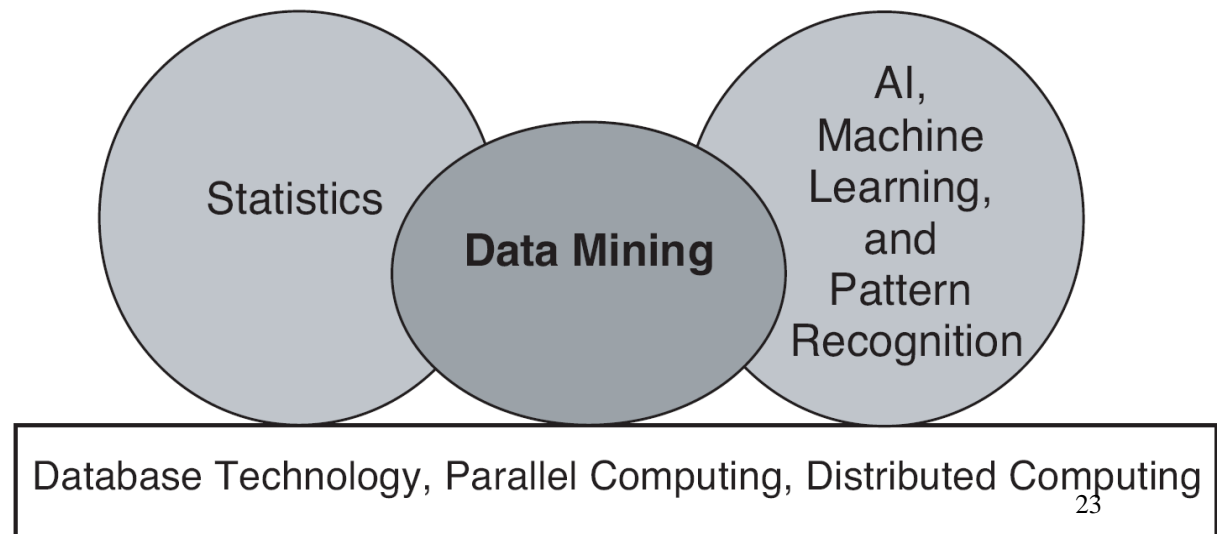
- Data mining — Discover valid, novel, useful, and understandable patterns in massive datasets



What is Data Mining?

■ Cross Disciplines

- Databases
- Machine learning: decision tree, Bayesian classifier, etc.
- Statistics: regression, etc.
- Neural networks
- Parallel/Distributed computing



Why Not Traditional Data Analysis?

- Tremendous amount of data
 - Algorithms must be highly scalable to handle such as tera-bytes of data

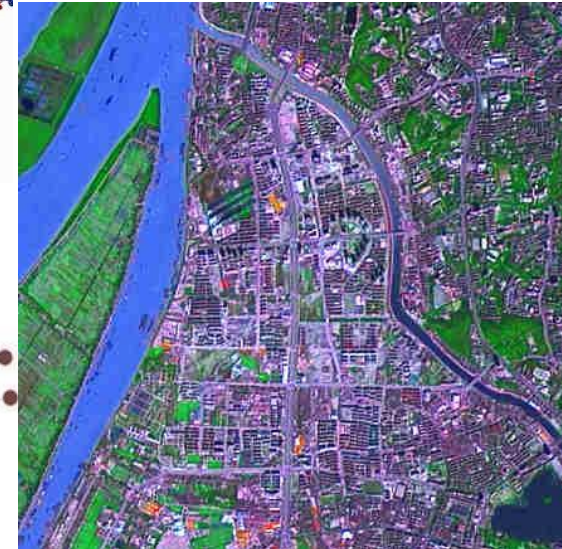
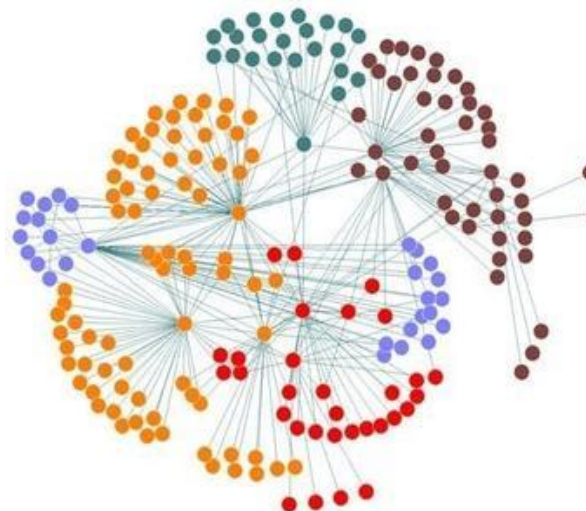
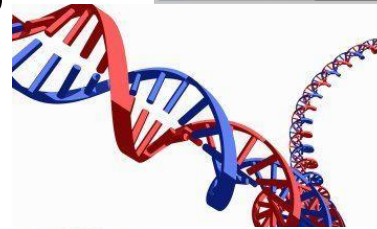


- High-dimensionality of data
 - DNA sequences may have tens of thousands of dimensions

TRFE_CHICK	WHLICLTNLSLBIAVCFAP	PKSVIRICTISSPEEXCHNLDTODERIS	LTQVQKATLDCIKAIANNEADATSLGGQVFEADLAPINLKPVAEYEH	
TRFE_HUMAN	MRLAYGALLYGAYLQLCLAYP	OKTVRICAVSDEATKODSFHMKSVIPDQGSVACVKKASTLDCIRAIANNEADAVTLQGLYIDA LAPHLNKPVAEYFG		
TRFE_XENLA	WFLSLRYALQLHMLALCLATG	YKQVRRKCKVSNELKXCKLYOTCKNE	IKLSCEYKSNTECSTATGDAICYQGGYKQSLQFYNLKPVAEYFG	
TRFE_RABIT	MRLAQLLACAAQLCLAYT	EXTVRICAVNHKSKCANFRDSMKVLPEDQPI	ICVKKASTLDCIKAIANNEADAVTLQGLYHEDLTPHLNKPVAEYFG	
TRFE_BOVIN	MSPAYRALLACAYLQLCLADP	ERTVRICITISTHANNICASFRENILRI	LESG-PFYSCNKTSTHWDIKAIANNEADAVTLQGLYHEDLTPHLNKPVAEYFG	
TRFE_PIG		YA	OKTVRICTISNDEANICSSPFRENKAYKING-PLYSCKVKSSTLDCIKAIANNEADAVTLQGLYHEDLTPHLNKPVAEYFG	
TRFE_HORSE	MRLAIRALLACAYLQLCLA	EDTVRICTVSNHNSKASFQDQWSTIVAP-PLVACVKTSTLDCIKAIANNEADAVTLQGLYHEDLTPHLNKPVAEYFG		
TRFE_ANPL		AP	PKTTVRICTISSAEDKXCHLKHQDERVT	LSQVQKATLDCIKAIANNEADATSLGGQVFEADLAPINLKPVAEYEH
TRF1_SALSA	WLLLLSALLQDLATAYAP	AEGIVKVKYKSEDELKCHDLANKVAEFS	CYKQGSFQCTQAKGGEADATLGGQVITAGLTNYLQLOPIIAEDYQ	
TRF2_SALSA	WLLLLSALLQDLATAYAP	AEGIVKVKYKSEDELKCHDLANKVAEFS	CYKQGSFQCTQAKGGEADATLGGQVITAGLTNYLQLOPIIAEDYQ	
NRL_ILFG		QRRSVQVCAVSNPEATKCFQWQRMKVRG	PPYSCIKRQSPICQCIQAIANNEADAVTLQGLYHEDLTPHLNKPVAEYFG	
TRF_BLAD1	WLLQLTLISABAVLHMTPEQSPH	IKVQVPEALES-CHNGGE	SQHMTCYAAERITLQDKIKHREADAPYQEDHMYAAKIPQDPIIIEVIRTK	
TRF_HANSE	WMLKLLTILALTDAAANAKSS	YNLCVPAATNKD-CEHLEVPK	SKYALECYPARDVDELSPYQGRADAPYQEDHMYAAKIPQDPIIIEVIRTK	
TRF1_HUMAN	WMLVLLVLLGALQLCLAGR	RRSVQVCAVSOPEATKCFQWQRMKVRG	PPYSCIKRQSPICQCIQAIANNEADAVTLQGLYHEDLTPHLNKPVAEYFG	
TRF1_BOVIN	WMLVYALLSLGALQLCLAP	RINVRICITISQPEVFCRQWQRMKVLDA	PSITCYRRAFALEDICRAIANNEADAVTLQGLYHEDLTPHLNKPVAEYFG	
TRF1_HUMAN	WROPSGALWLLALRTVLDG	VEYRVKATSPQEHKCNSEAFREAD	IGPOLLCHRTSADHCVOLIAADADATLQGLYHEDLTPHLNKPVAEYFG	
TRF1_HOUSE	WMLLIPSLIFLEALQLCLA	KATTYQVCAVSNSEEDCLQWQRMKVRG	PPLSCVKSSTROCIQAIANNEADATLQGLYHEDLTPHLNKPVAEYFG	
SAX_RANCA	NAPTFTALFFTIISLBFAAP	NAKTVRICAISLEBKXCHLYSSCNFD	ITLVCYLRSSTEDCMTAKDQADHFLSGEYKQSLNKPVAEYEH	

Why Not Traditional Data Analysis?

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



Why Not Traditional Data Analysis?

■ Database

- Storage-oriented
- Provide simple queries

Data mining

Discover knowledge from data in databases

■ Data warehouse

- Subject-oriented
- A multidimensional view of data
- Operations to access summarized data

Advanced data analysis tools

■ Statistical algorithms

- Based on many hypothesis
- Find patterns in small number of samples

Less hypothesis

Find patterns in large number of samples

Abnormal patterns

Characteristics of Data Mining

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

Exercises

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



On Thursdays, during 4-11pm customers often purchase diapers and beers together!



Ex. 1: Market Basket Analysis and Management

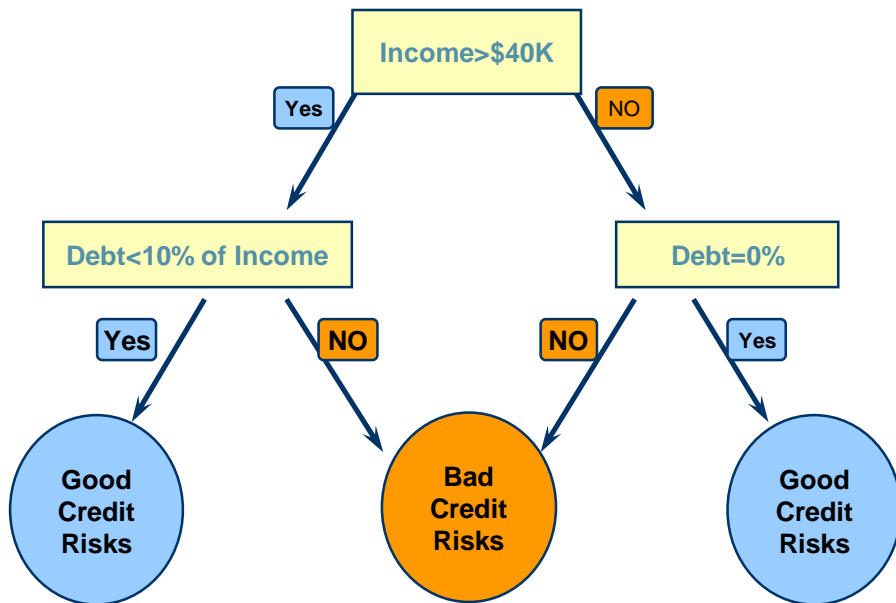
- Where does the data come from?
 - Supermarket transactions, membership cards, discount coupons, customer complaint calls
- Cross-marketing analysis
 - What products were often purchased together?
Purchase recommendation, cross selling
 - What are the subsequent purchases after buying a given product?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



What Kinds of Patterns?

■ Classification

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



• Decision Tree

rule 1: if (Income ≤ \$40k) and (Debt = 0) then “good”

rule 2: if (Income > \$40K) and (Debt < 10% of Income) then “good”

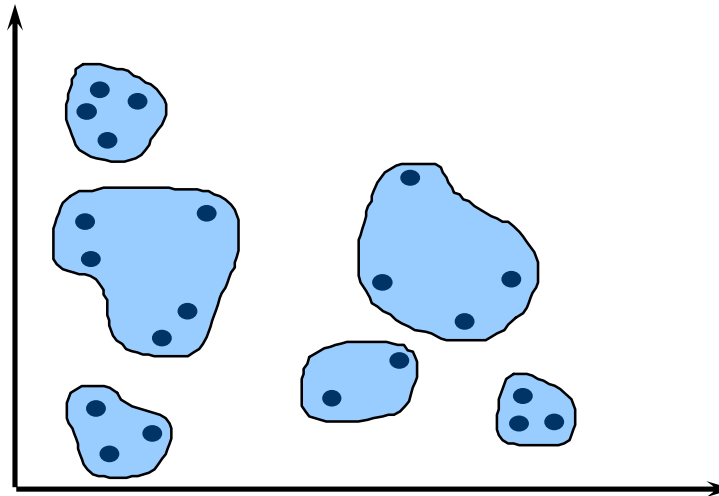
Ex.2 Credit Scoring

- 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?

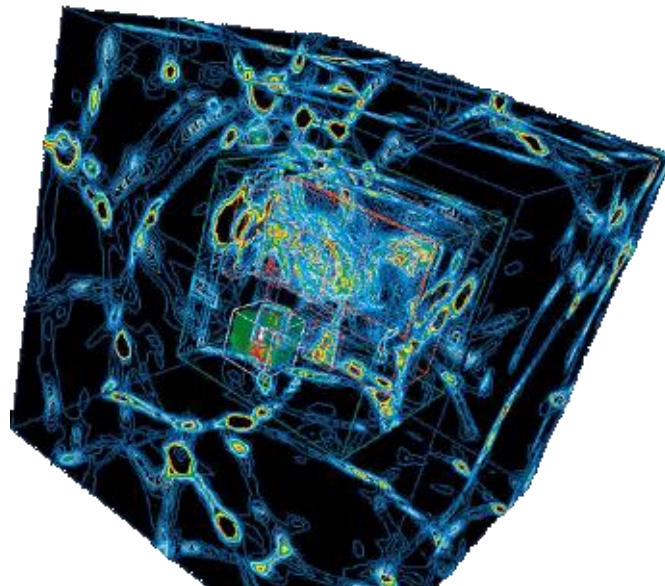
■ Clustering

- 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

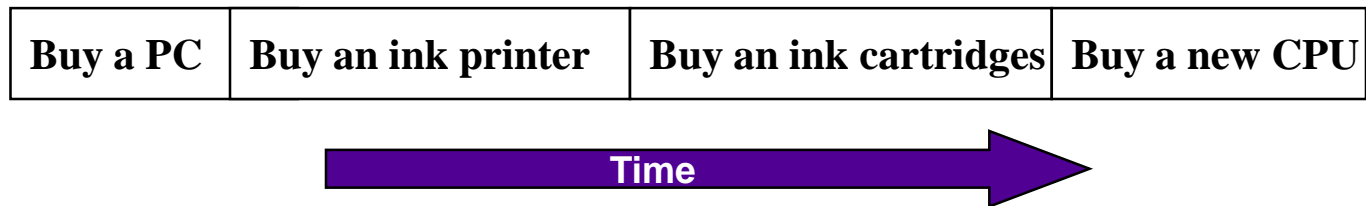
- Cosmological simulation
 - Simulate the formation of the galaxy
 - Enormous particles at each evolution stage, beyond the capability of human being to analyze



What Kinds of Patterns?

■ Sequence mining

- Given a set of sequences, find the complete set of frequent subsequences

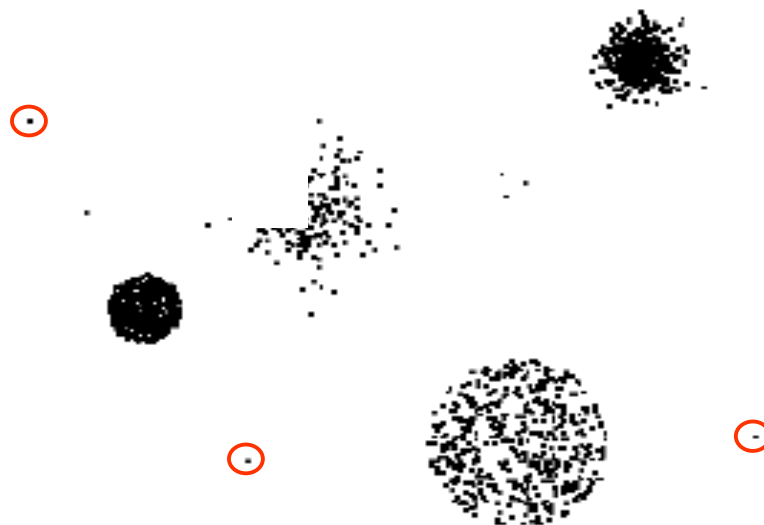


Marketing strategy: recommend a new CPU for the customer 9 months after his first purchase

What Kinds of Patterns?

■ Anomaly detection

- Given 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



Anomalies may be valuable!

What Kinds of Patterns?

■ 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 京东

天猫 Tmall.com

iqiyi 爱奇艺

youku 优酷

腾讯视频 V.qq.com

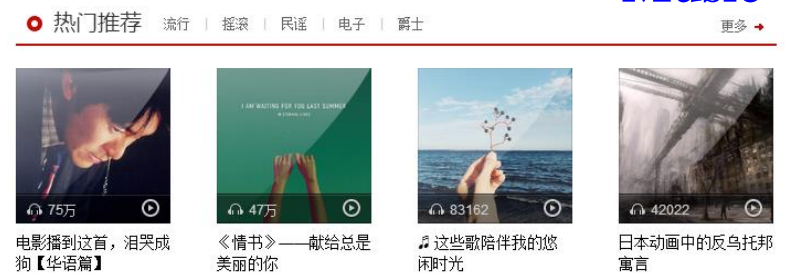
Customers Who Viewed This Item Also Viewed



Movie



Music



Exercises

1. 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?

- Database-oriented data sets and applications
 - Relational database, data warehouse, transactional database
- Advanced database applications
 - Data streams
 - Spatial data
 - Text database
 - Multimedia data
 - Time-series
 - Bio-medical data
 - Network traffic data

Relational Databases

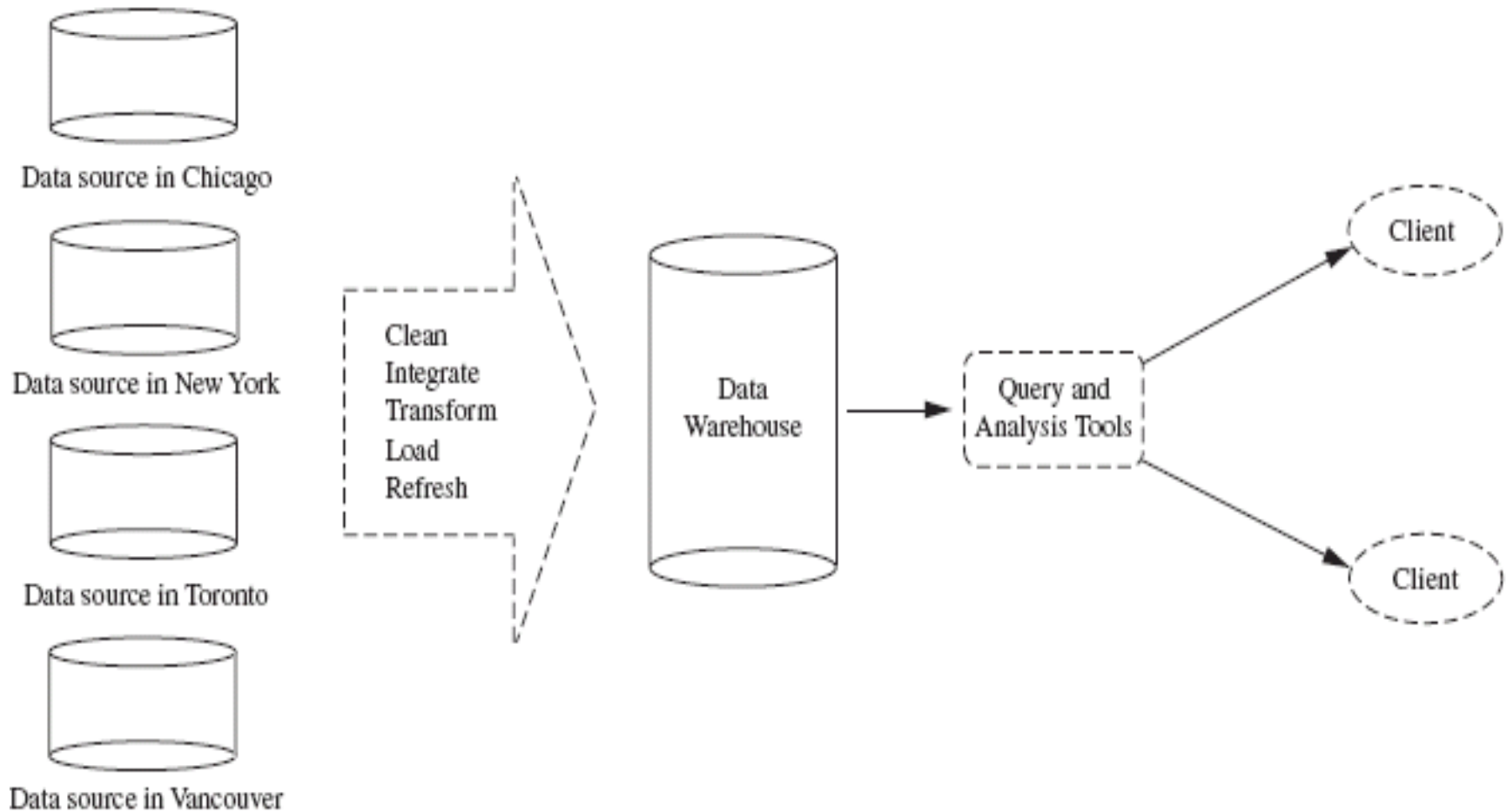
- Structured data
 - Table – records – attributes
 - Accessed by queries, SQL
- Online transactional processing (OLTP)
 - Insert a student “Ying Liu” into class “Introduction to Data Mining”, fall 2014

Name	Time	Course	score	Room
Ying Liu	Fall 2014	Introduction to Data Mining	90	002
Tom	Fall 2014	Math	85	001
Merlisa	Spring 2014	Compiler	70	001
George	Fall 2014	Graphics	92	001

Data Warehouses

- 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 “Ying Liu” in the last 3 years, grouped by semesters
- Many patterns are summarization of data
 - Roll-up, drill-down

Data Warehouses



Transactional Databases

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

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

Spatial Data

■ Spatial information

- Geographic databases (map)
- VLSI chip design databases
- Satellite/remote sensing image databases
- Medical image database

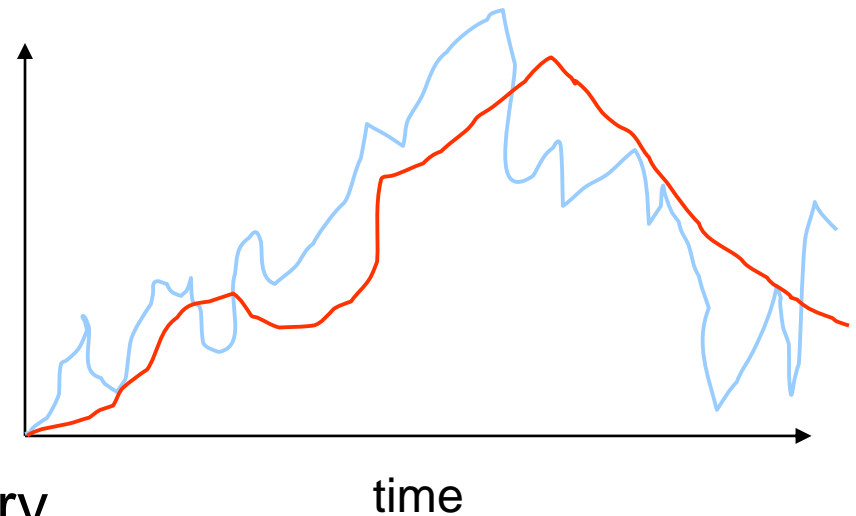
编号	中心	正右方	右上方	面积
1	居民地	绿地	水体	100
2	绿地	水体	水体	50
3	水体	居民地	居民地	600
4	水体	绿地	绿地	54
...

■ Spatial patterns

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

Time Series

- 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

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



Data Streams

- 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 networks, etc.

Biomedical Data

■ Bio-sequences

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



World-Wide Web

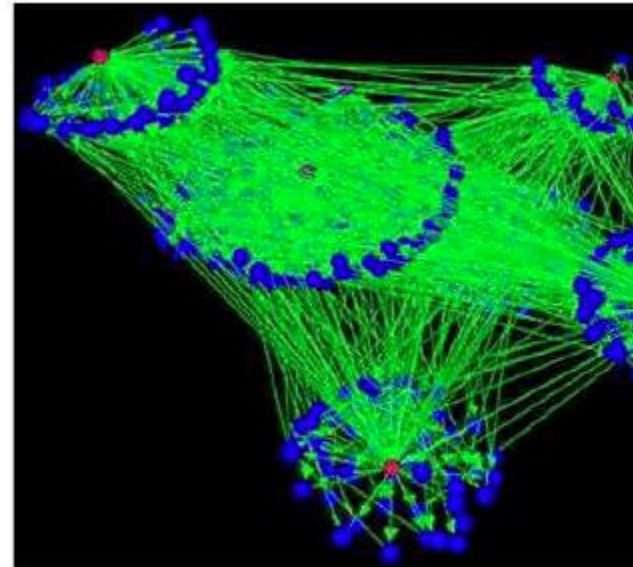
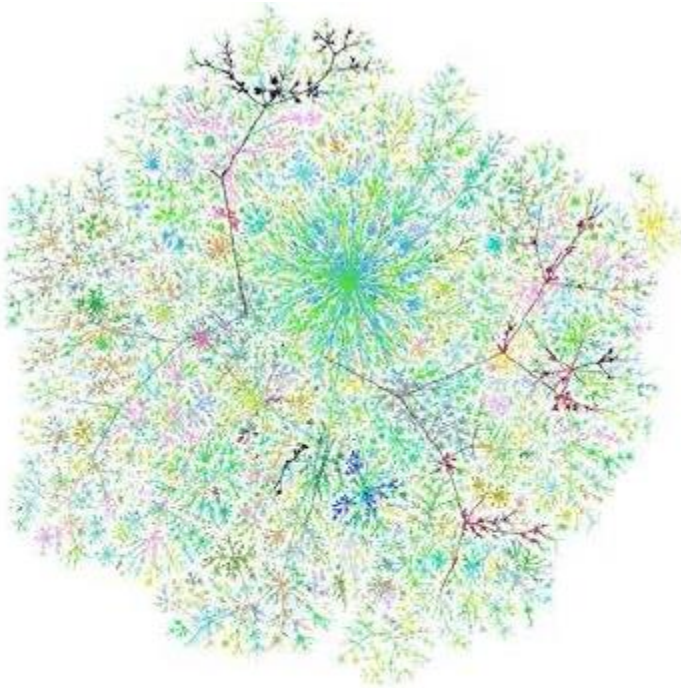
- The WWW is huge, widely distributed, global information service center for
 - Information services: news, advertisements, consumer information, financial management, education, government, e-commerce, etc.
 - Hyper-link information
 - Access and usage information
- WWW provides rich sources for data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

World-Wide Web

- Web Usage: Logs and IP package header streams
 - Mine Weblog records to discover user accessing patterns of Web pages
- Web Content
 - Extract knowledge from a Web documents, automatic categorization
- Web Structure
 - Identifying interesting graph patterns among different Web pages

Graph

■ Internet graph



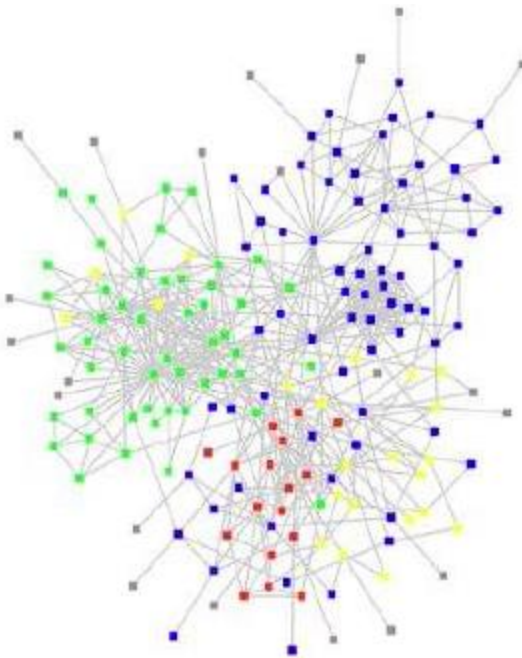
The images are downloaded from
<http://www.maths.bris.ac.uk/~maarw/graphs/graph.html>
and <http://www.netdimes.org/new/?q=node/17>

- Citation graph



Graph

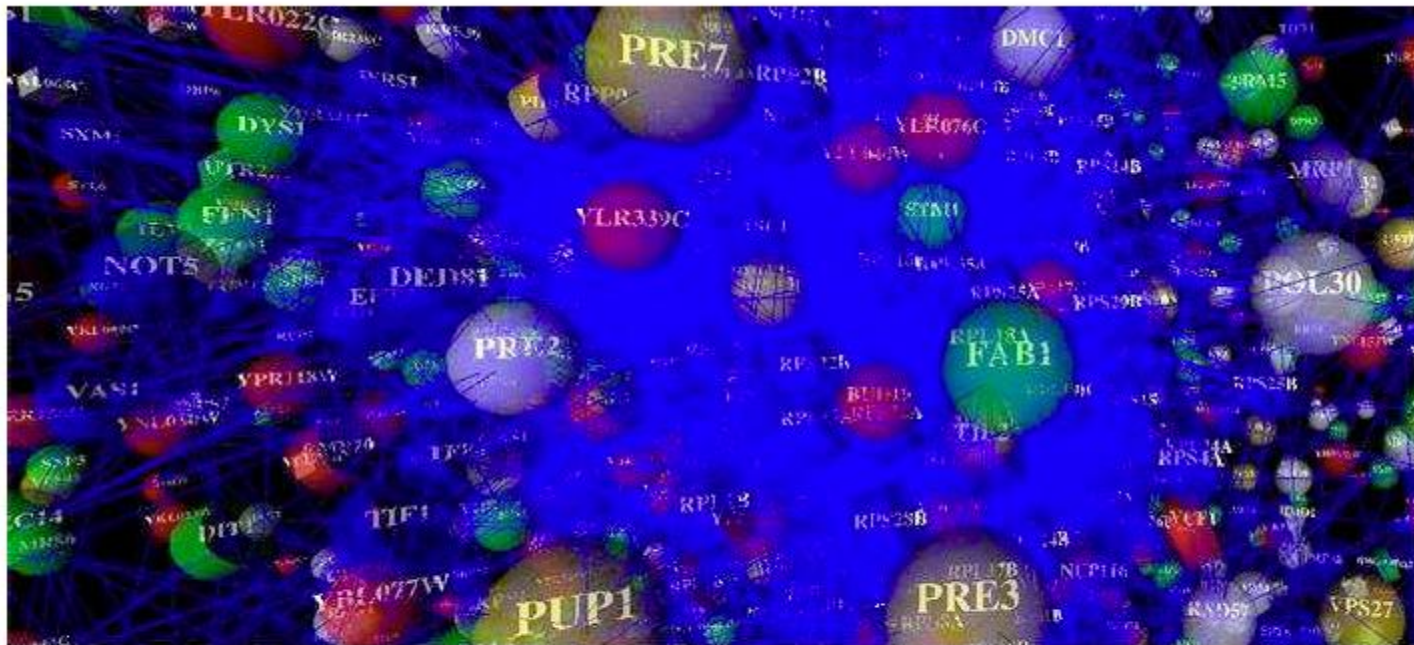
■ Friendship graph



The images are downloaded from
<http://www.thenetworkthinker.com/>
and [http://myweb20list.com/blog/2008/03/23/
new-amazing-facebook-photo-mapper/my-facebook-friend-graph/](http://myweb20list.com/blog/2008/03/23/new-amazing-facebook-photo-mapper/my-facebook-friend-graph/)

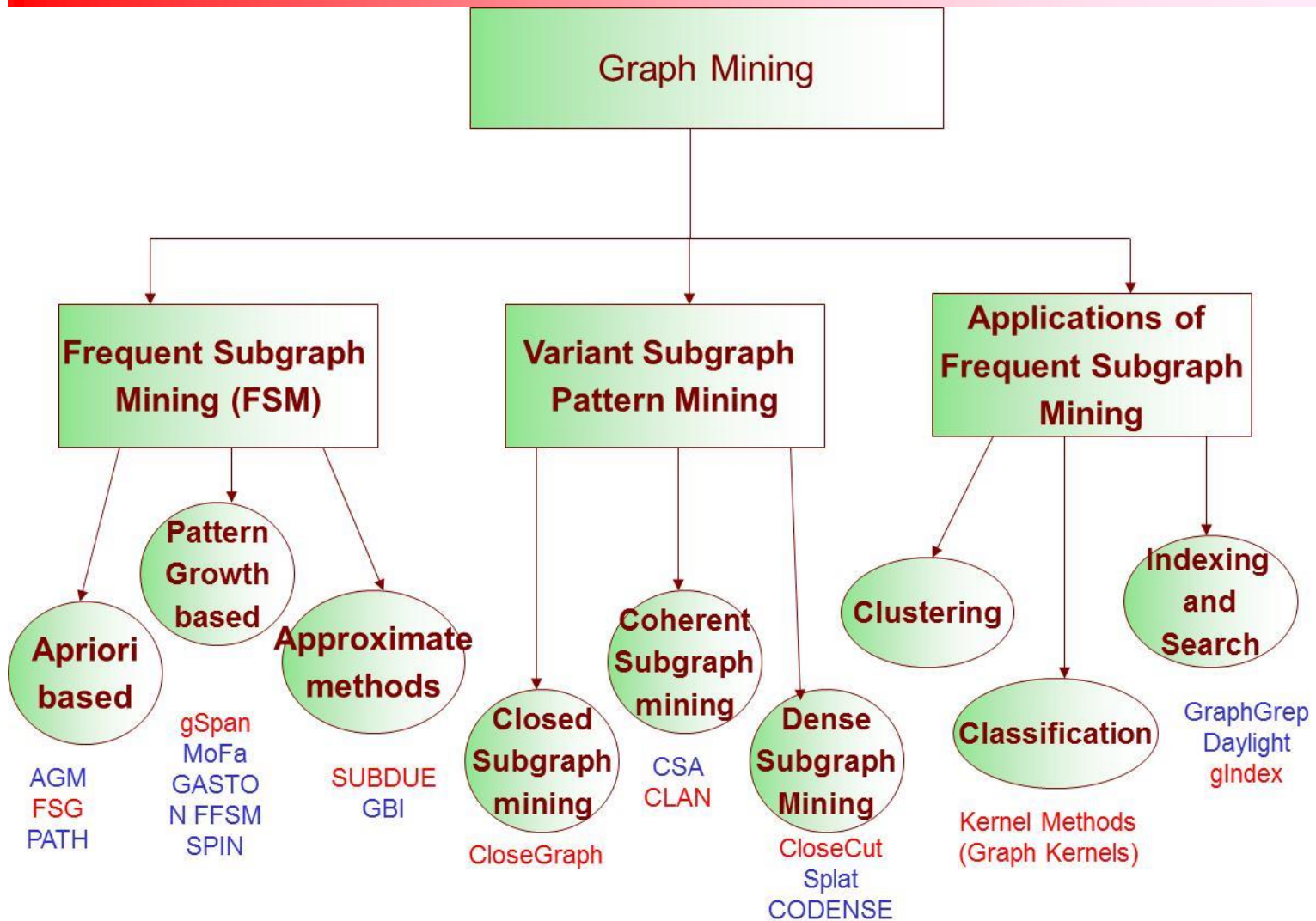
Graph

■ Protein interaction graph



The images are downloaded from
<http://bioinformatics.icmb.utexas.edu/lgl/Images/rsomZoom.jpg>

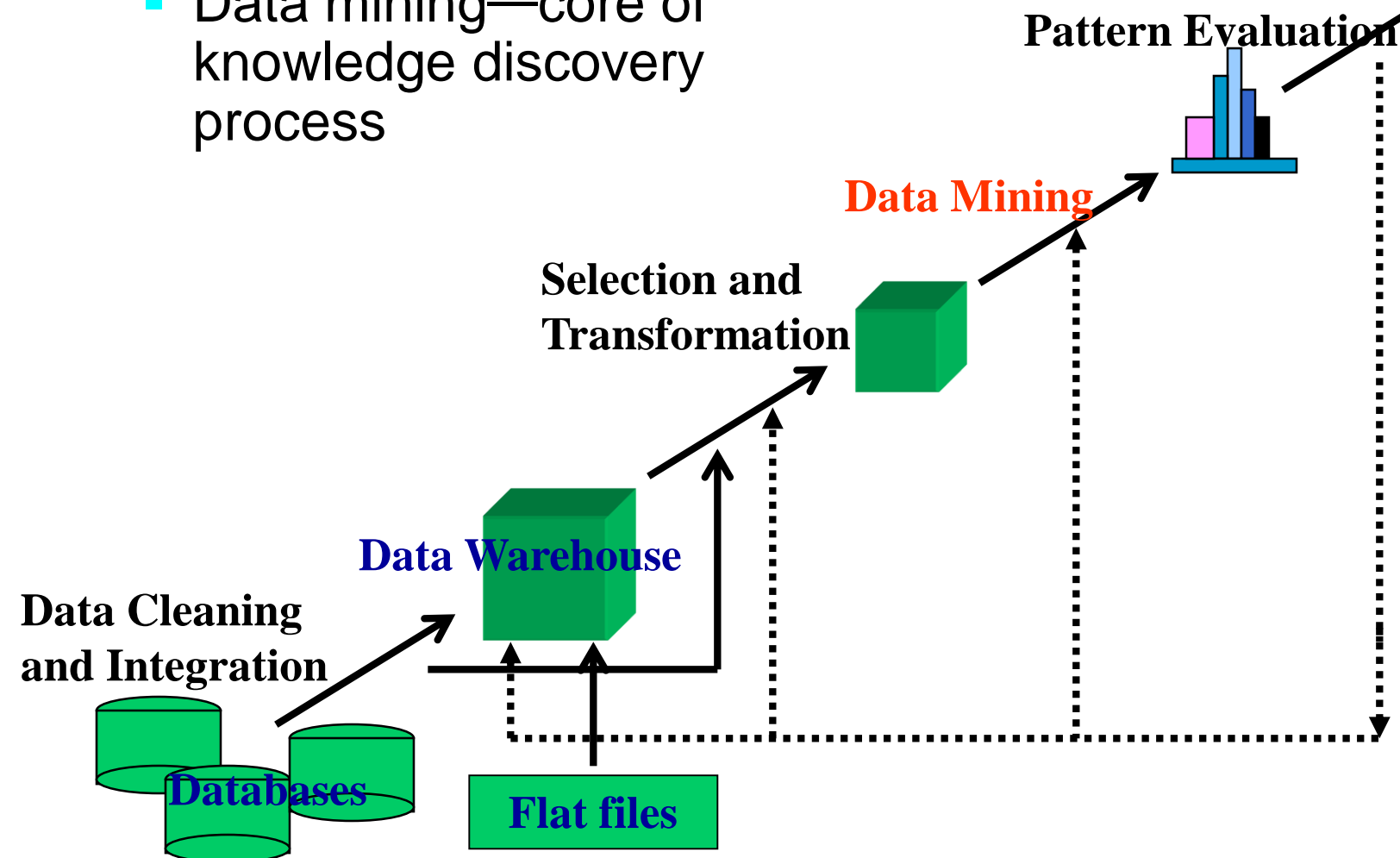
Graph



Knowledge Discovery (KDD) Process

Knowledge

- Data mining—core of knowledge discovery process



Key Steps in KDD Process

- 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?

- 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?

- 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 search
- Search for only interesting patterns: An optimization problem — Challenging
 - Can a data mining system find only the interesting patterns?
 - Approaches
 - First generate all the patterns and then filter out the uninteresting ones
 - Guide and constrain the discovery process

Research Issues in Data Mining

■ Mining methodology

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

Research Issues in Data Mining

- User interaction
 - Data mining query languages
 - Expression and visualization of data mining results
- Applications and social impacts
 - Domain-specific data mining
 - Protection of data security, integrity, and privacy

Important Resources

- Data mining conferences
 - ACM SIGKDD, IEEE ICDM, SIAM DM, PKDD, PAKDD
- Database conferences
 - ACM SIGMOD, VLDB, ACM PODS, IEEE ICDE, EDBT, ICDT
- Important journals
 - ACM Data Mining and Knowledge Discovery
 - IEEE Transactions on Knowledge and Data Engineering
 - Knowledge and Information Systems