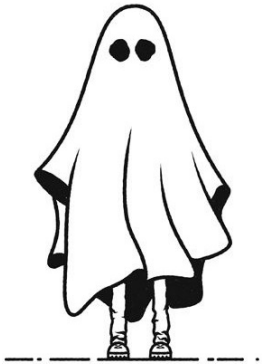


CS 412 Intro. to Data Mining

Chapter 1. Introduction

Jiawei Han, Computer Science, Univ. Illinois at Urbana-Champaign, 2017



ระบบสำรวจและวิเคราะห์ข้อมูลงานเคสช่วยที่มีโครงสร้าง

Help Needed: LifeNet—A Structured Network-Based Knowledge Exploration and Analytics System for Life Sciences

- What we are doing?

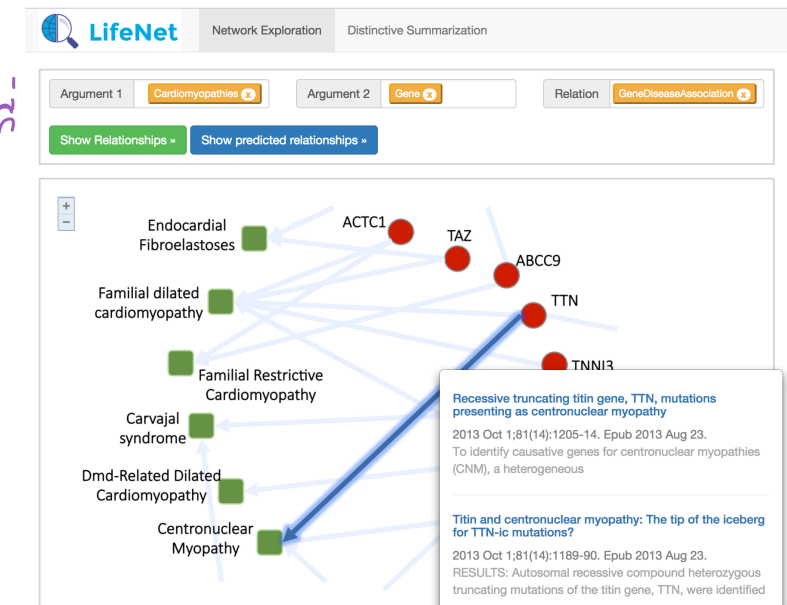
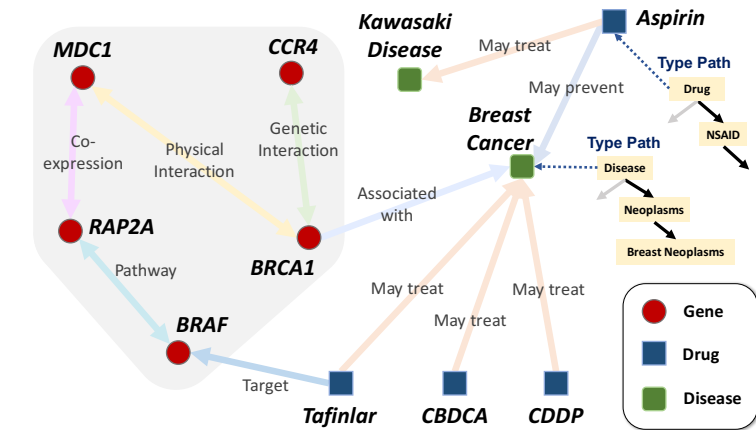
เปลี่ยนเอกสารเป็นกราฟและรองรับ function
การค้นหาคือ Analysis ข้อมูล

 - A scalable system that transforms biomedical papers into a knowledge graph & supports various search/analytics functions
- What we already have?

ระบบต้นแบบ

 - A working prototype system & an ACL demo paper
- What we are looking for?
 - Students with expertise on HTML/CSS & JavaScript
 - Experiences on web frameworks and databases
 - System design experience will be a big plus

สามารถศึกษาได้จากหลายที่
- What you will gain?
 - Hourly pay (\$12-\$15 per hour, 6-20 hours per week)
 - Possible research publications & a good thesis topic



Send us your resume if interested: Jiaming Shen (mickeysjm@gmail.com)

Why Data Mining? ^{ทำไมต้องทำเหมืองข้อมูล}

^{การเติบโต}

- The Explosive ^{การเติบโต} Growth of Data: from terabytes to petabytes

- Data ^{รวบรวม} collection and data ^{พร้อมใช้งาน} availability

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

- Major sources of abundant data ^{แหล่งข้อมูลที่สำคัญมากมาย}

- 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!

- ^{ความจำเป็น} “Necessity is the mother of invention” — ^{เอนกต้นกำเนิดของการประดิษฐ์} Data mining — ^{ชุดข้อมูล} Automated analysis of ^{การวิเคราะห์ชุดข้อมูลขนาดใหญ่} massive data sets



อะไรคือเหมืองข้อมูล What Is Data Mining?

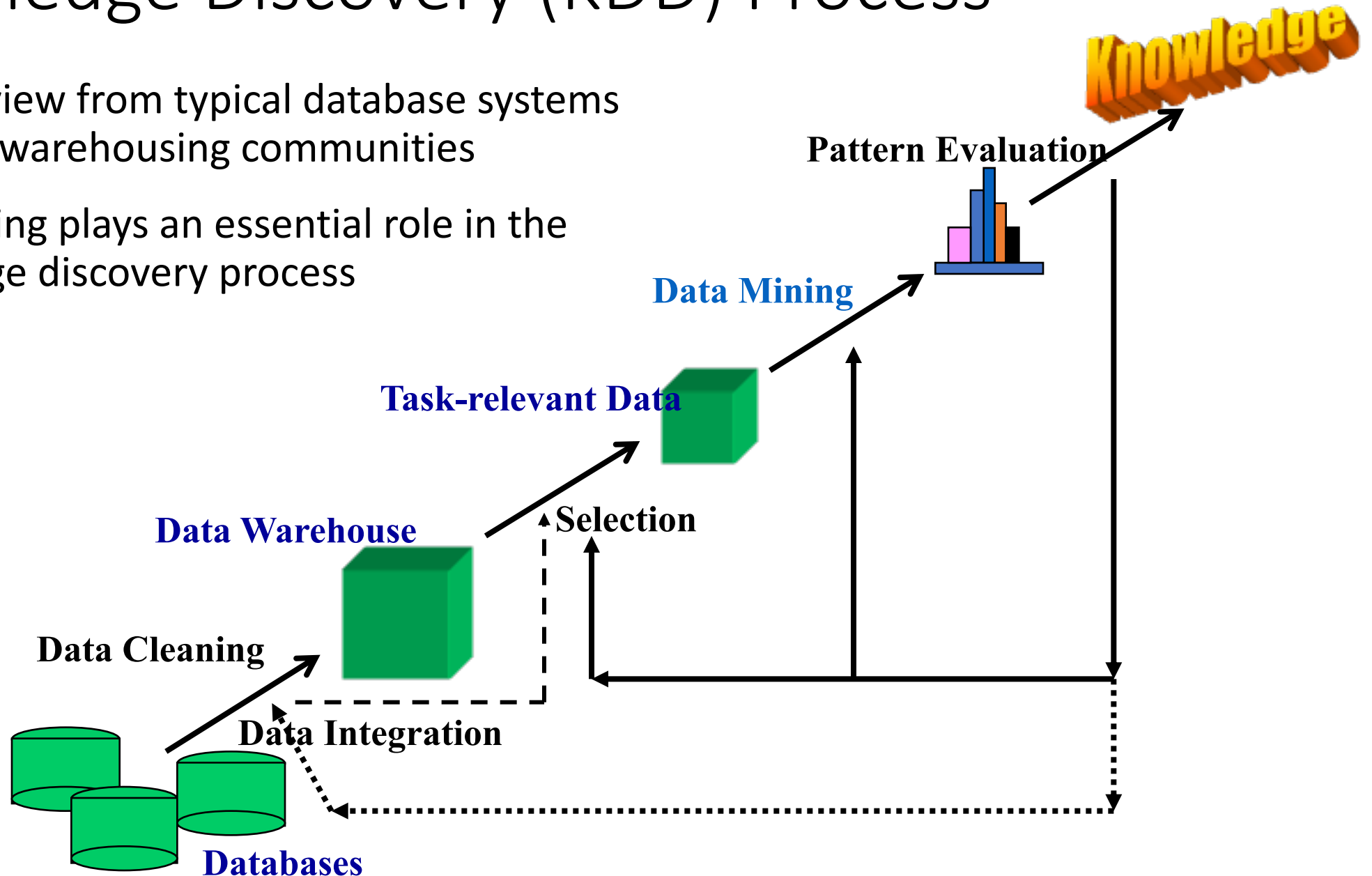


- Data mining (knowledge discovery from data)
 - Extraction of interesting (non-trivial, implicit, previously unknown and potentially useful) patterns or knowledge from huge amount of data
 - Data mining: a misnomer?
- Alternative names:
 - Knowledge discovery (mining) in databases (KDD), knowledge extraction, data/pattern analysis, data archeology, data dredging, information harvesting, business intelligence, etc.
- Watch out: Is everything “data mining”?
 - Simple search and query processing
 - (Deductive) expert systems



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

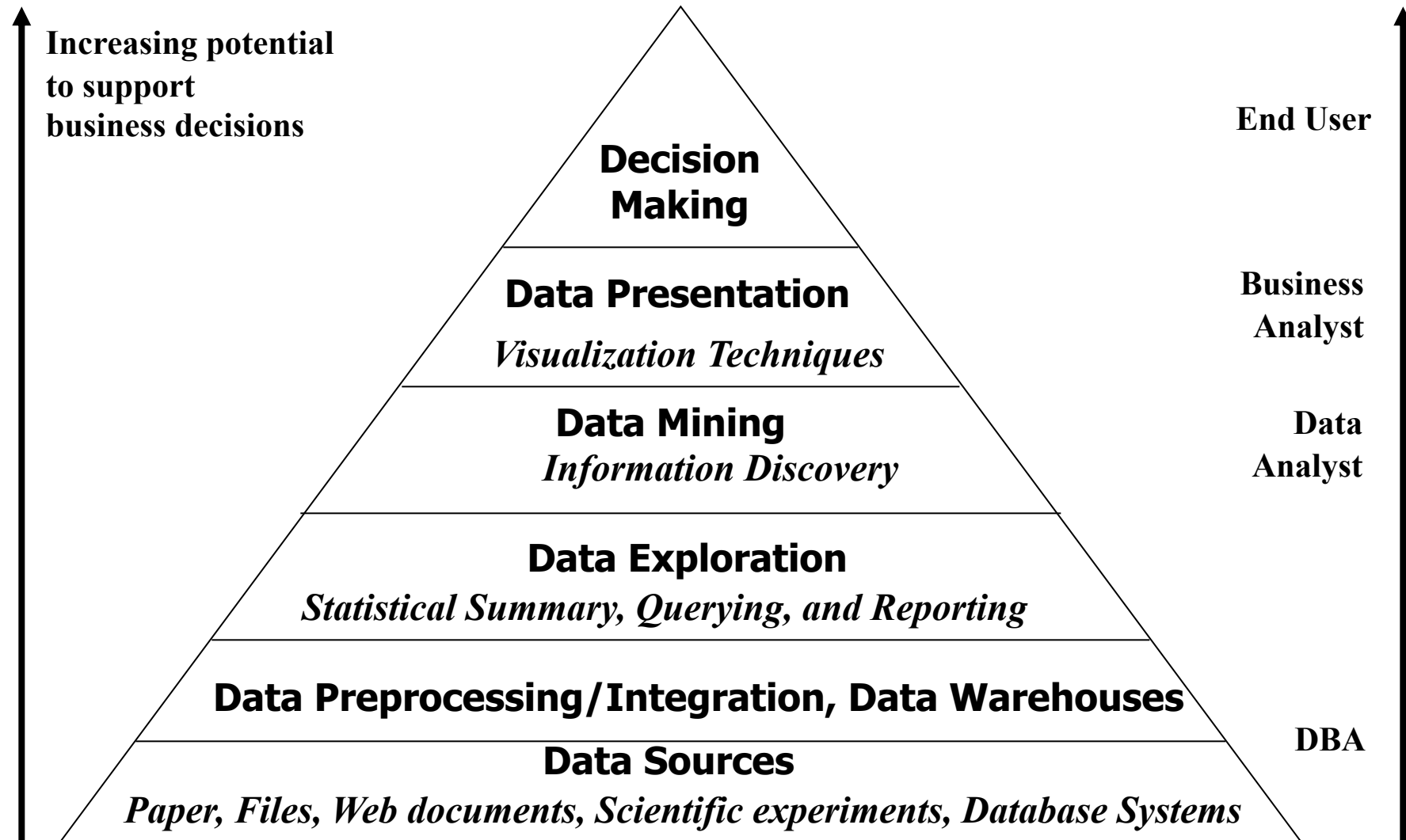


Example: A Web Mining Framework

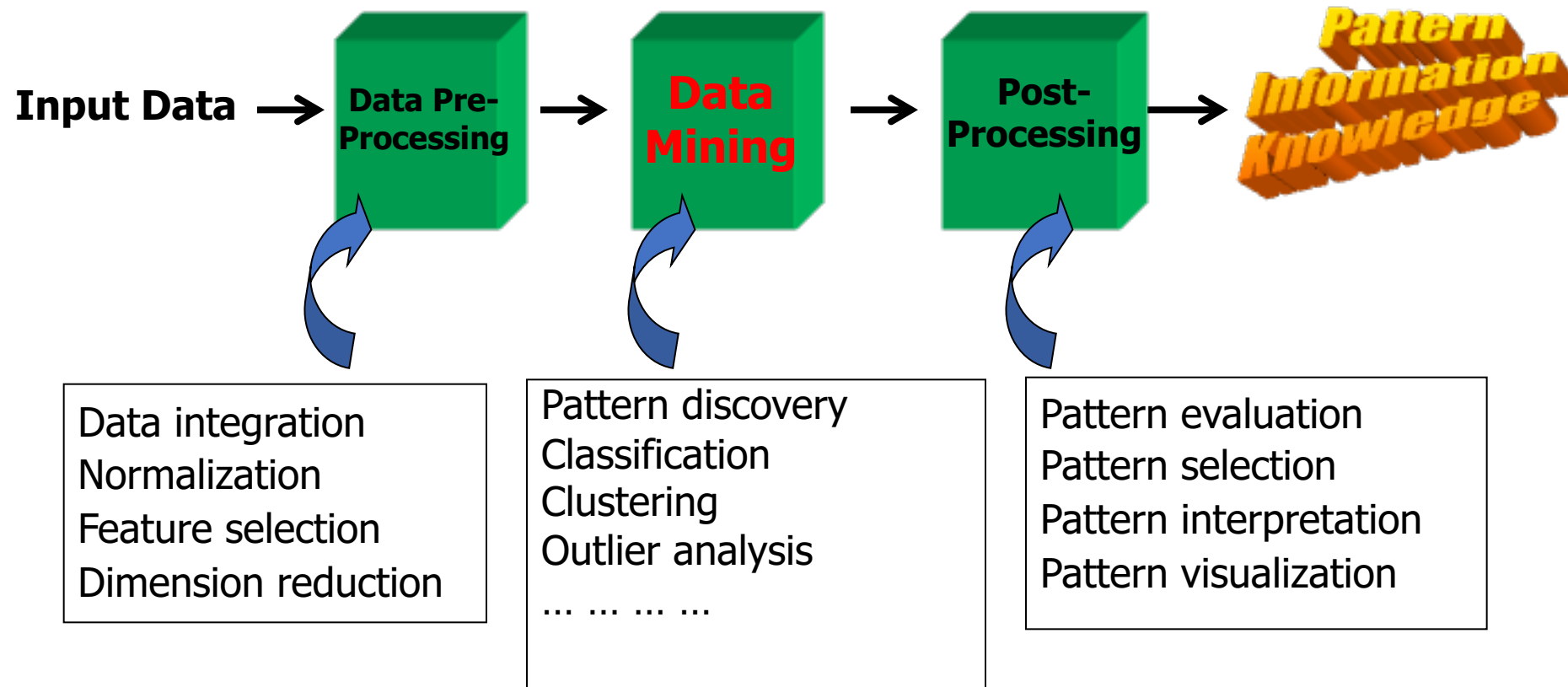
- Web mining usually involves
 - Data cleaning
 - Data integration from multiple sources
 - Warehousing the data
 - Data cube construction
 - Data selection for data mining
 - Data mining
 - Presentation of the mining results
 - Patterns and knowledge to be used or stored into knowledge-base



Data Mining in Business Intelligence



KDD Process: A View from ML and Statistics



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

Data Mining vs. Data Exploration

- Which view do you prefer?
 - KDD vs. ML/Stat. vs. Business Intelligence
 - Depending on the data, applications, and your focus
- Data Mining vs. Data Exploration
 - Business intelligence view
 - Warehouse, data cube, reporting but not much mining
 - Business objects vs. data mining tools
 - Supply chain example: mining vs. OLAP vs. presentation tools
 - Data presentation vs. data exploration



Multi-Dimensional View of Data Mining

- **Data to be mined** ข้อมูลที่ขุด
 - Database data (extended-relational, object-oriented, heterogeneous), 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)** ฟังก์ชันการทำเหมืองข้อมูล
 - Characterization, discrimination, association, classification, clustering, trend/deviation, outlier analysis, ...
 - 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** การดัดแปลง
 - Retail, telecommunication, banking, fraud analysis, bio-data mining, stock market analysis, text mining. Web mining. etc.

Data Mining: On What Kinds of Data? ^{ชุดข้อมูลประเภทใด?}

- ^{ชุดข้อมูล} Database-oriented data sets and applications
 - ^{ฐานข้อมูลเชิงสัมพันธ์} Relational database, ^{คลังข้อมูล} data warehouse, ^{ระบบธุรกรรม} transactional database
 - Object-relational databases, ^{ต่างกัน} Heterogeneous databases and ^{สืบทอด} legacy databases
- ^{ชุดข้อมูลขั้นสูง} Advanced data sets and advanced applications
 - Data streams and sensor data
 - ^{ข้อมูลตามเวลา} Time-series data, ^{ข้อมูลทาง} temporal data, ^{ลำดับ} sequence data (incl. bio-sequences)
 - Structure data, graphs, social networks and information networks
 - ^{เชิงพื้นที่} Spatial data and spatiotemporal data
 - Multimedia database
 - Text databases
 - The World-Wide Web

ลักษณะทั่วไป

Data Mining Functions: (1) Generalization

การรวบรวมและสร้างคลังข้อมูล

- Information integration and data warehouse construction

- Data ^{ล้าง}cleaning, ^{แปลง}transformation, ^{รวม}integration, and multidimensional data model ^{จำลองลักษณะที่}

- Data cube technology

การนำข้อมูลมาวิเคราะห์ตามความต้องการ

- Scalable methods for computing (i.e., materializing) multidimensional aggregates ^{การรวบรวมข้อมูลเป็น}

- OLAP (online analytical processing) ^{เชิงวิเคราะห์}

- Multidimensional concept description: Characterization and discrimination ^{การนิยามลักษณะเฉพาะและการเลือกประเภท}

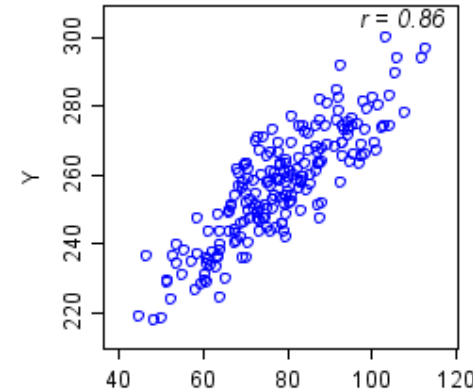
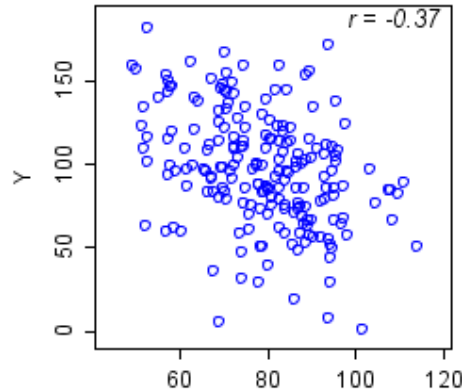
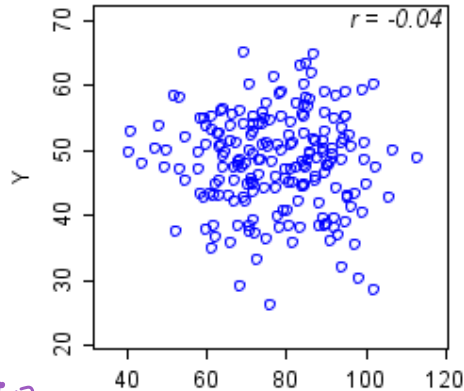
- Generalize, summarize, and contrast data characteristics, e.g., dry vs. wet region

^{สรุป} ^{เปรียบเทียบคุณลักษณะข้อมูล}



Data Mining Functions: (2) Pattern Discovery

- Frequent patterns (or frequent itemsets)
 - What items are frequently purchased together in your Walmart?
- Association and Correlation Analysis



□ A typical association rule

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

□ Are strongly associated items also strongly correlated?

□ How to mine such patterns and rules efficiently in large datasets?

□ How to use such patterns for classification, clustering, and other applications?

การค้นพบรูปแบบ → ที่สามารถช่วยอะไร

บ่อย ๆ รูปแบบ ชุดสินค้าที่บ่อย

ช่วยตัดสินใจ

วิเคราะห์ความสัมพันธ์

การเชื่อมโยง

Ex. ถ้าคุณซื้อสินค้าไปทางไปทางใกล้ ๆ มักจะพบได้เหมือนกัน

Data Mining Functions: (3) Classification

การจำแนกประเภท: ทศ

การทำนายผล

- Classification and label prediction

สร้างแบบจำลอง

- Construct models (functions) based on some training examples

อธิบาย: อธิบาย

- Describe and distinguish classes or concepts for future prediction

แนวคิด

ทำนายอนาคต

- Ex. 1. Classify countries based on (climate)

จำแนกตามภูมิอากาศ

- Ex. 2. Classify cars based on (gas mileage)

รถตามระยะ: ย. กม

- Predict some unknown class labels

รู้ล่วงหน้า

- Typical methods

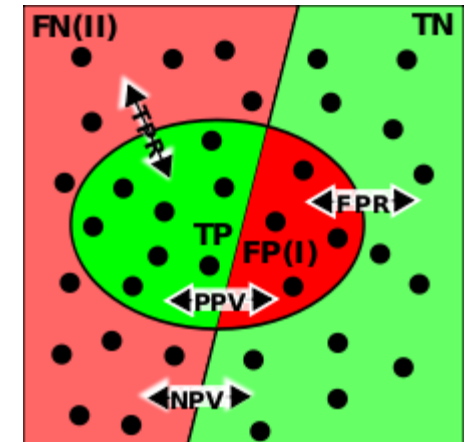
- Decision trees, naïve Bayesian classification, support vector machines, neural networks, rule-based classification, pattern-based classification, logistic regression, ...

- Typical applications:

- Credit card fraud detection, direct marketing, classifying stars, diseases, web-pages, ...

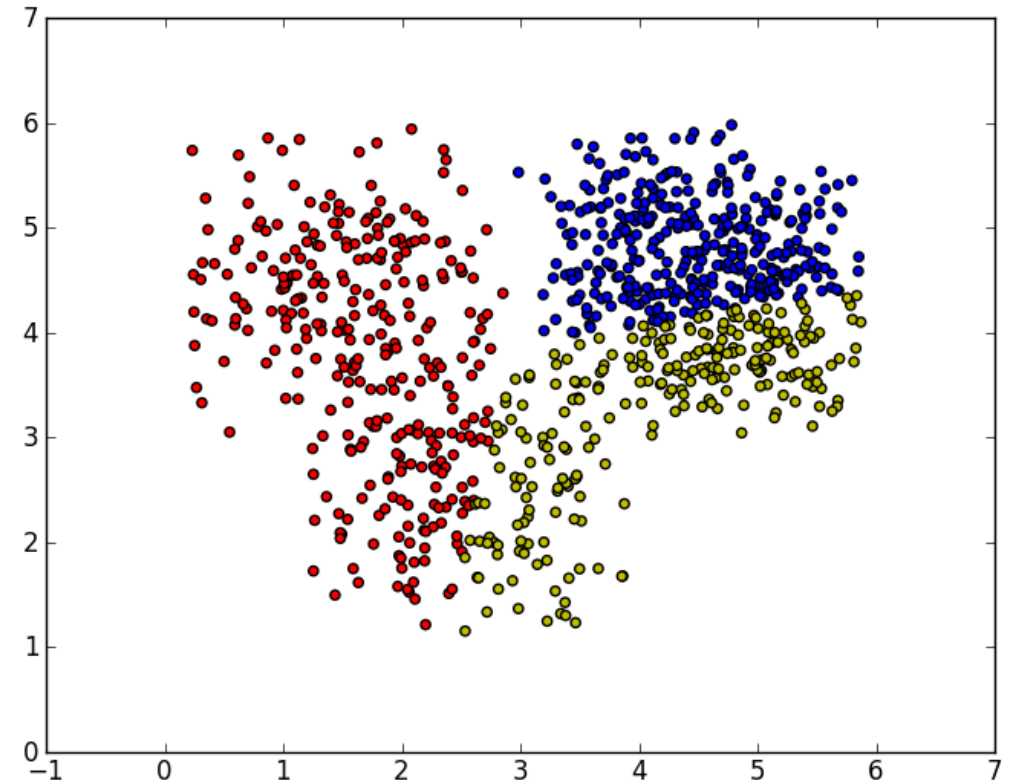
การตรวจจับการฉ้อโกงบัตรเครดิต

การตลาดทางตรง



Data Mining Functions: (4) Cluster Analysis

- การรู้แบบไม่มีผู้ดูแล เช่น ไม่รู้จัดเข้ากับ
• Unsupervised learning (i.e., Class label is unknown)
- เพื่อสร้างหมวดหมู่ใหม่
• Group data to form new categories (i.e., clusters), e.g., cluster houses to find distribution patterns
การจับ
- เพื่อหาความคล้ายคลึง และลด
• Principle: Maximizing intra-class similarity & minimizing interclass similarity
- Many methods and applications



Data Mining Functions: (5) Outlier Analysis

- Outlier analysis

- Outlier: A data object that does not comply with the general behavior of the data

ไม่สอดคล้องกับพฤติกรรมทั่วไปของข้อมูล

- Noise or exception?—One person's garbage could be another person's treasure

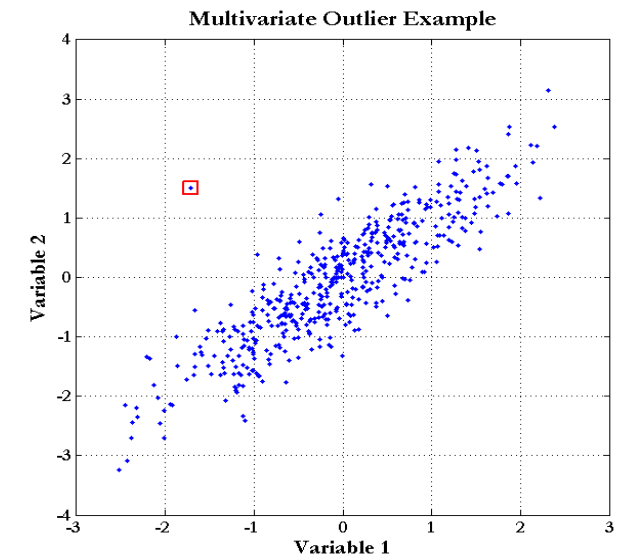
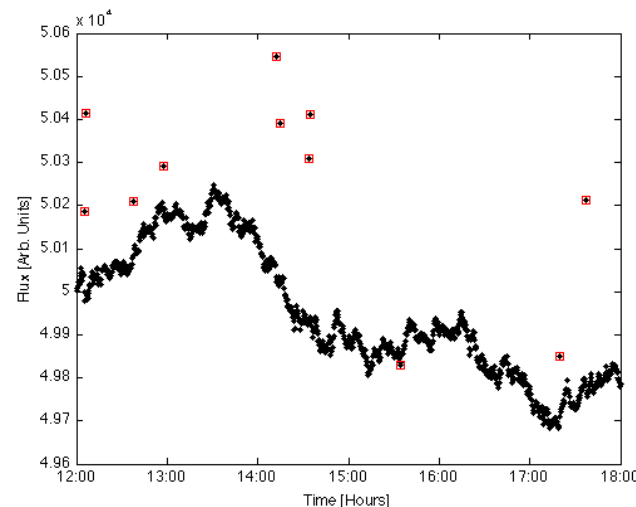
ของใครบางคนเสีย

- Methods: by product of clustering or regression analysis, ...

โดยผลคูณของการจัดกลุ่มและการวิเคราะห์ถดถอย

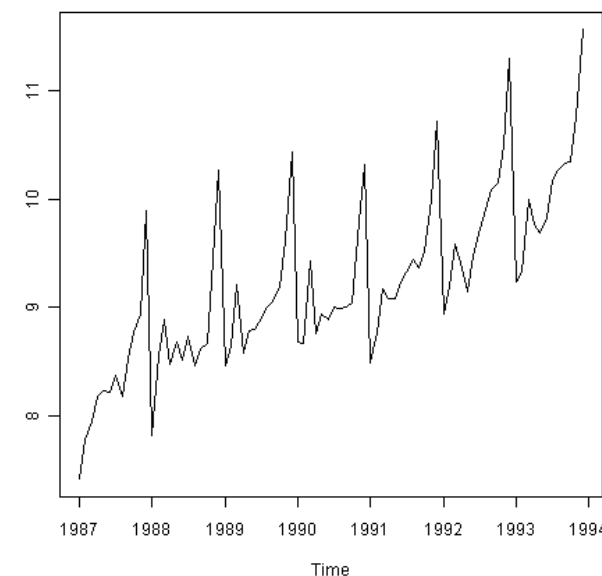
สเปซโอบเจกต์ในการตรวจสอบการผิดปกติก่อน, การวิเคราะห์เหตุการณ์

- Useful in fraud detection, rare events analysis



Data Mining Functions: (6) Time and Ordering: Sequential Pattern, Trend and Evolution Analysis

- Sequence, trend and evolution analysis
 - Trend, time-series, and deviation analysis
 - e.g., regression and value prediction
 - Sequential pattern mining
 - e.g., buy digital camera, then buy large memory cards
 - Periodicity analysis
 - Motifs and biological sequence analysis
 - Approximate and consecutive motifs
 - Similarity-based analysis
- Mining data streams
 - Ordered, time-varying, potentially infinite, data streams



Data Mining Functions: (7) Structure and Network Analysis

การวิเคราะห์โครงสร้างและเครือข่าย

- Graph mining

การขุดค้นกราฟ

 - Finding frequent subgraphs (e.g., chemical compounds), trees (XML), substructures (web fragments)

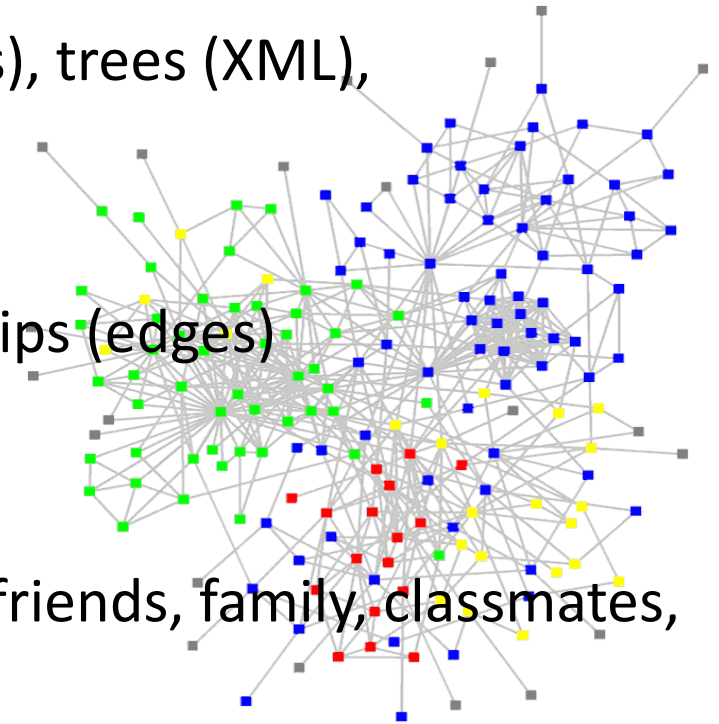
ค้นหากราฟบ่อยๆ
- Information network analysis

การวิเคราะห์เครือข่ายข้อมูล

 - Social networks: actors (objects, nodes) and relationships (edges)

เครือข่ายสังคม

 - e.g., author networks in CS, terrorist networks
 - Multiple heterogeneous networks
 - A person could be multiple information networks: friends, family, classmates, ...
 - Links carry a lot of semantic information: Link mining
- Web mining
 - Web is a big information network: from PageRank to Google
 - Analysis of Web information networks



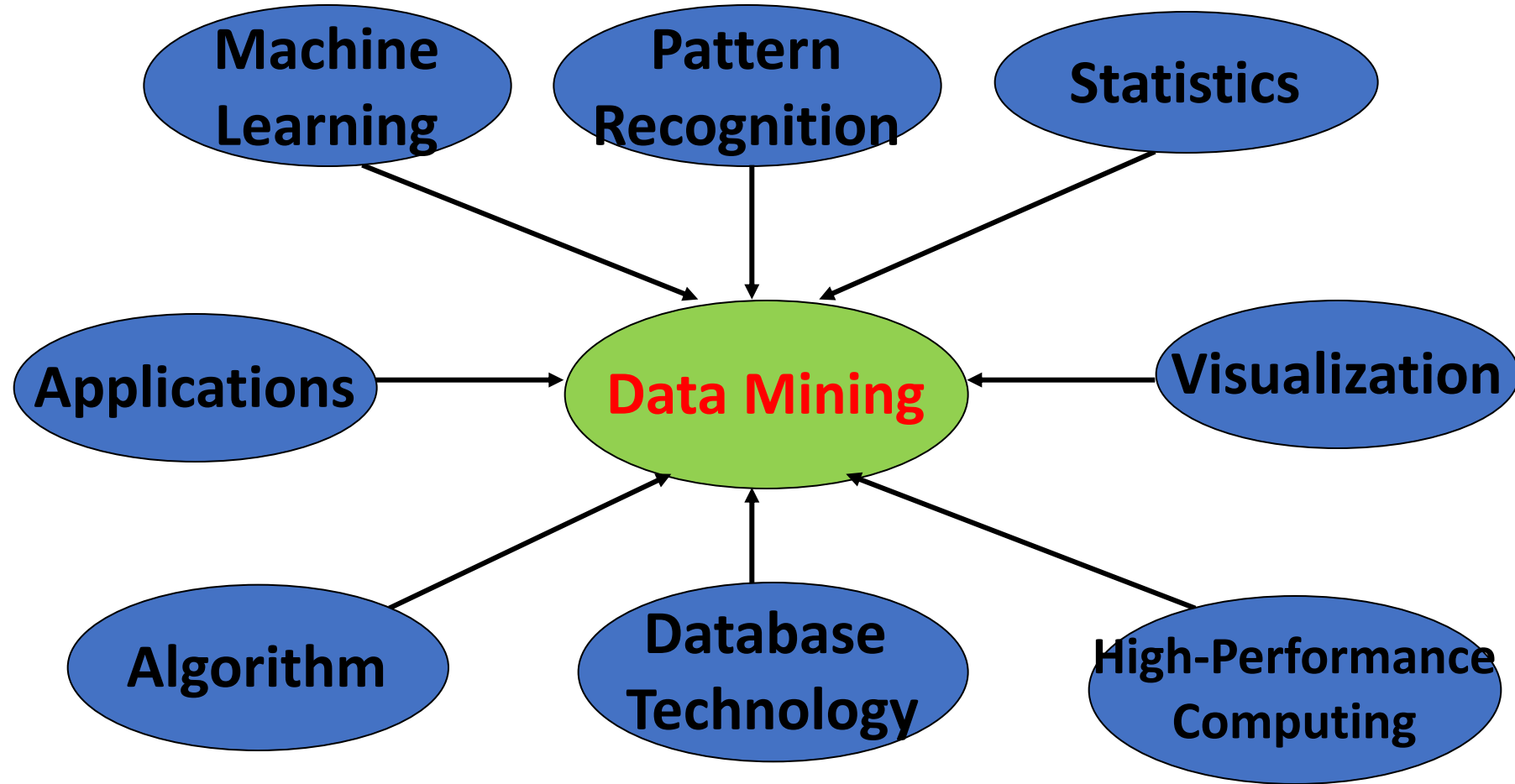
การประเมินความรู้

Evaluation of Knowledge

- Are all mined knowledge interesting? ข้อมูลที่ขุดได้มีความน่าสนใจหรือไม่
 - One can mine tremendous amount of “patterns” สามารถขุดมาได้จำนวนมาก
 - Some may fit only certain dimension space (time, location, ...) บางส่วนอาจสอดคล้องกับพื้นที่บางมิติเท่านั้น
 - Some may not be representative, may be transient, ... ชั่วคราว
- Evaluation of mined knowledge → directly mine only interesting knowledge? การประเมินความรู้ที่ถูกต้อง
 - Descriptive vs. predictive ความครอบคลุมเชิงพรรณนา
 - Coverage การคาดการณ์
 - Typicality vs. novelty แปลกใหม่
 - Accuracy ค.แม่นยำ
 - Timeliness ทันเวลา, ทันสมัย
 - ...



Data Mining: Confluence of Multiple Disciplines



ပါးပါး

အသံပူဆူသော

Why Confluence of Multiple Disciplines?

အလွန်အမင်းများပြားသော

- Tremendous amount of data

အလွန်အမင်းများပြားသော

- Algorithms must be scalable to handle big data

အလွန်အမင်း

- High-dimensionality of data

- Micro-array may have tens of thousands of dimensions

အလွန်အမင်း

- High complexity of data

- Data streams and sensor data
 - Time-series data, temporal data, sequence data
 - Structure data, graphs, social and information networks
 - Spatial, spatiotemporal, multimedia, text and Web data
 - Software programs, scientific simulations

- New and sophisticated applications



การประยุกต์ใช้

Applications of Data Mining

- Web page analysis: classification, clustering, ranking
การวิเคราะห์หน้าเว็บ จำแนก จัดกลุ่ม จัดอันดับ
- Collaborative analysis & recommender systems
วิเคราะห์ร่วมกัน แนะนำ
- Basket data analysis to targeted marketing
เพื่อการตลาดเป้าหมาย
- Biological and medical data analysis
- Data mining and software engineering
- Data mining and text analysis
- Data mining and social and information network analysis
- Built-in (invisible data mining) functions in Google, MS, Yahoo!, Linked, Facebook, ...
ฝังไว้ในตัว การขุดข้อมูลที่ไม่มองเห็น
- Major dedicated data mining systems/tools
เครื่องมือที่เน้นเฉพาะขุด
 - SAS, MS SQL-Server Analysis Manager, Oracle Data Mining Tools)



Major Issues in Data Mining (1) ประเด็นหลักในการขุดข้อมูล

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

Major Issues in Data Mining (2)

- Efficiency and Scalability ပြု: ပုံစံကောင်းမှု နှင့်: ချောမွေ့မှု
 - Efficiency and scalability of data mining algorithms ပြုမည်ကောင်း
 - Parallel, distributed, stream, and incremental mining methods
- Diversity of data types ဂ. ပုံစံကောင်းမှု
 - Handling complex types of data
 - Mining dynamic, networked, and global data repositories
- Data mining and society အကျိုးကျေးဇူးတင်မှု
 - Social impacts of data mining
 - Privacy-preserving data mining
 - Invisible data mining အမြင်မရသော အချက်အလက်များကို ရှာဖွေခြင်း



A Brief History of Data Mining Society

- 1989 IJCAI Workshop on Knowledge Discovery in Databases
 - Knowledge Discovery in Databases (G. Piatetsky-Shapiro and W. Frawley, 1991)
- 1991-1994 Workshops on Knowledge Discovery in Databases
 - Advances in Knowledge Discovery and Data Mining (U. Fayyad, G. Piatetsky-Shapiro, P. Smyth, and R. Uthurusamy, 1996)
- 1995-1998 International Conferences on Knowledge Discovery in Databases and Data Mining (KDD'95-98)
 - Journal of Data Mining and Knowledge Discovery (1997)
- ACM SIGKDD conferences since 1998 and SIGKDD Explorations
- More conferences on data mining
 - PAKDD (1997), PKDD (1997), SIAM-Data Mining (2001), (IEEE) ICDM (2001), WSDM (2008), etc.
- ACM Transactions on KDD (2007)

Conferences and Journals on Data Mining

- KDD Conferences

- ACM SIGKDD Int. Conf. on Knowledge Discovery in Databases and Data Mining (**KDD**)
- SIAM Data Mining Conf. (**SDM**)
- (IEEE) Int. Conf. on Data Mining (**ICDM**)
- European Conf. on Machine Learning and Principles and practices of Knowledge Discovery and Data Mining (**ECML-PKDD**)
- Pacific-Asia Conf. on Knowledge Discovery and Data Mining (**PAKDD**)
- Int. Conf. on Web Search and Data Mining (**WSDM**)

- Other related conferences

- DB conferences: ACM SIGMOD, VLDB, ICDE, EDBT, ICDT, ...
- Web and IR conferences: WWW, SIGIR, WSDM
- ML conferences: ICML, NIPS
- PR conferences: CVPR,

- Journals

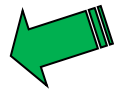
- Data Mining and Knowledge Discovery (DAMI or DMKD)
- IEEE Trans. On Knowledge and Data Eng. (TKDE)
- KDD Explorations
- ACM Trans. on KDD

Where to Find References? DBLP, CiteSeer, Google

- Data mining and KDD (SIGKDD)
 - Conferences: ACM-SIGKDD, IEEE-ICDM, SIAM-DM, PKDD, PAKDD, etc.
 - Journal: Data Mining and Knowledge Discovery, KDD Explorations, ACM TKDD
- Database systems (SIGMOD)
 - 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.
- AI & Machine Learning
 - Conferences: Machine learning (ML), AAAI, IJCAI, COLT (Learning Theory), CVPR, NIPS, etc.
 - Journals: Machine Learning, Artificial Intelligence, Knowledge and Information Systems, IEEE-PAMI, etc.
- Web and IR
 - Conferences: SIGIR, WWW, CIKM, etc.
 - Journals: WWW: Internet and Web Information Systems,
- Statistics
 - Conferences: Joint Stat. Meeting, etc.
 - Journals: Annals of statistics, etc.
- Visualization
 - Conference proceedings: CHI, ACM-SIGGraph, etc.
 - Journals: IEEE Trans. visualization and computer graphics, etc.

Chapter 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



Summary

การขุดข้อมูล: การค้นพบรูปแบบและความรู้ที่น่าสนใจจากข้อมูลจำนวนมาก

- Data mining: Discovering interesting patterns and knowledge from massive amount of data
- A natural evolution of science and information technology, in great demand, with wide applications
วิวัฒนาการทางวิทยาศาสตร์และเทคโนโลยี
ซึ่งมีความต้องการอย่างมาก
- A KDD process includes data cleaning, data integration, data selection, transformation, data mining, pattern evaluation, and knowledge presentation
- Mining can be performed in a variety of data
- Data mining functionalities: characterization, discrimination, association, classification, clustering, trend and outlier analysis, etc.
- Data mining technologies and applications
- Major issues in data mining



Recommended Reference Books

- Charu C. Aggarwal, *Data Mining: The Textbook*, Springer, 2015
- E. Alpaydin. *Introduction to Machine Learning*, 2nd ed., MIT Press, 2011
- R. O. Duda, P. E. Hart, and D. G. Stork, *Pattern Classification*, 2ed., Wiley-Interscience, 2000
- 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
- T. M. Mitchell, *Machine Learning*, McGraw Hill, 1997
- P.-N. Tan, M. Steinbach and V. Kumar, *Introduction to Data Mining*, Wiley, 2005 (2nd ed. 2016)
- I. H. Witten and E. Frank, *Data Mining: Practical Machine Learning Tools and Techniques with Java Implementations*, Morgan Kaufmann, 2nd ed. 2005
- Mohammed J. Zaki and Wagner Meira Jr., *Data Mining and Analysis: Fundamental Concepts and Algorithms* 2014