

CS 412 Intro. to Data Mining

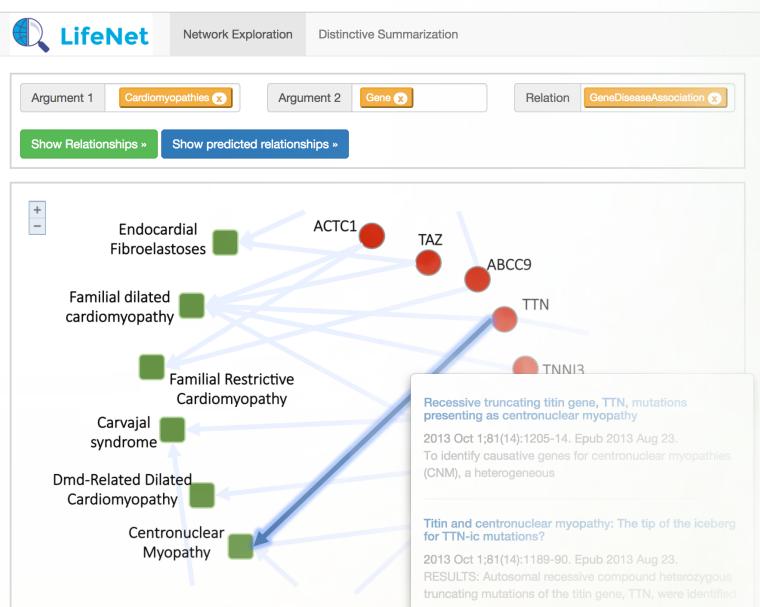
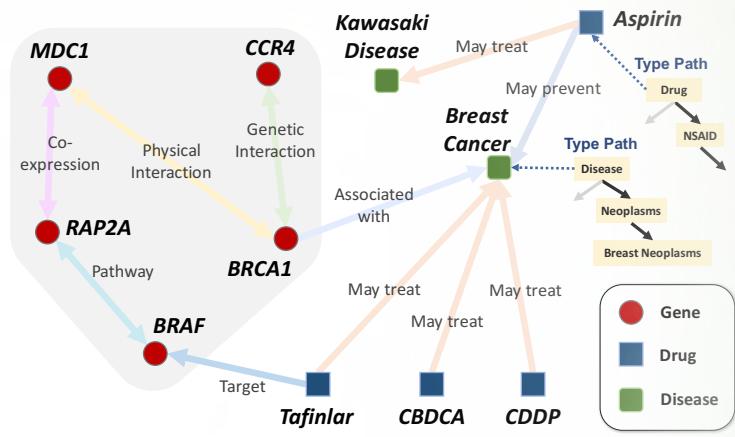
Chapter 1. Introduction

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



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

- What we are doing? *ផ្តល់ព័ត៌មានទៅសារប៊ិនការងារនៃបច្ចេកវិទ្យា function*
 - A scalable system that transforms biomedical papers into a knowledge graph & supports various search/analytics functions *ការអាជីវការណ៍ Analysis នៃព័ត៌មាន*
- What we already have? *មិនបានព័ត៌មាន*
 - 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** *សាមគរកំណើនហាទីតាមលក្ខណៈ*
- 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?

កំណត់ការអវិជ្ជមាន

ការពេញ

- 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 **លេខធម្មន៍ភ័ត៌មានរបស់គេ**
 - 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



ការងារ
សេវាទីនៃការបង្កើតទូទៅ

ទូទៅ

ការងារ: សេវាទីនៃការបង្កើតទូទៅ

ឧបកដីលប់ពេទ្យជាបន្ទូល

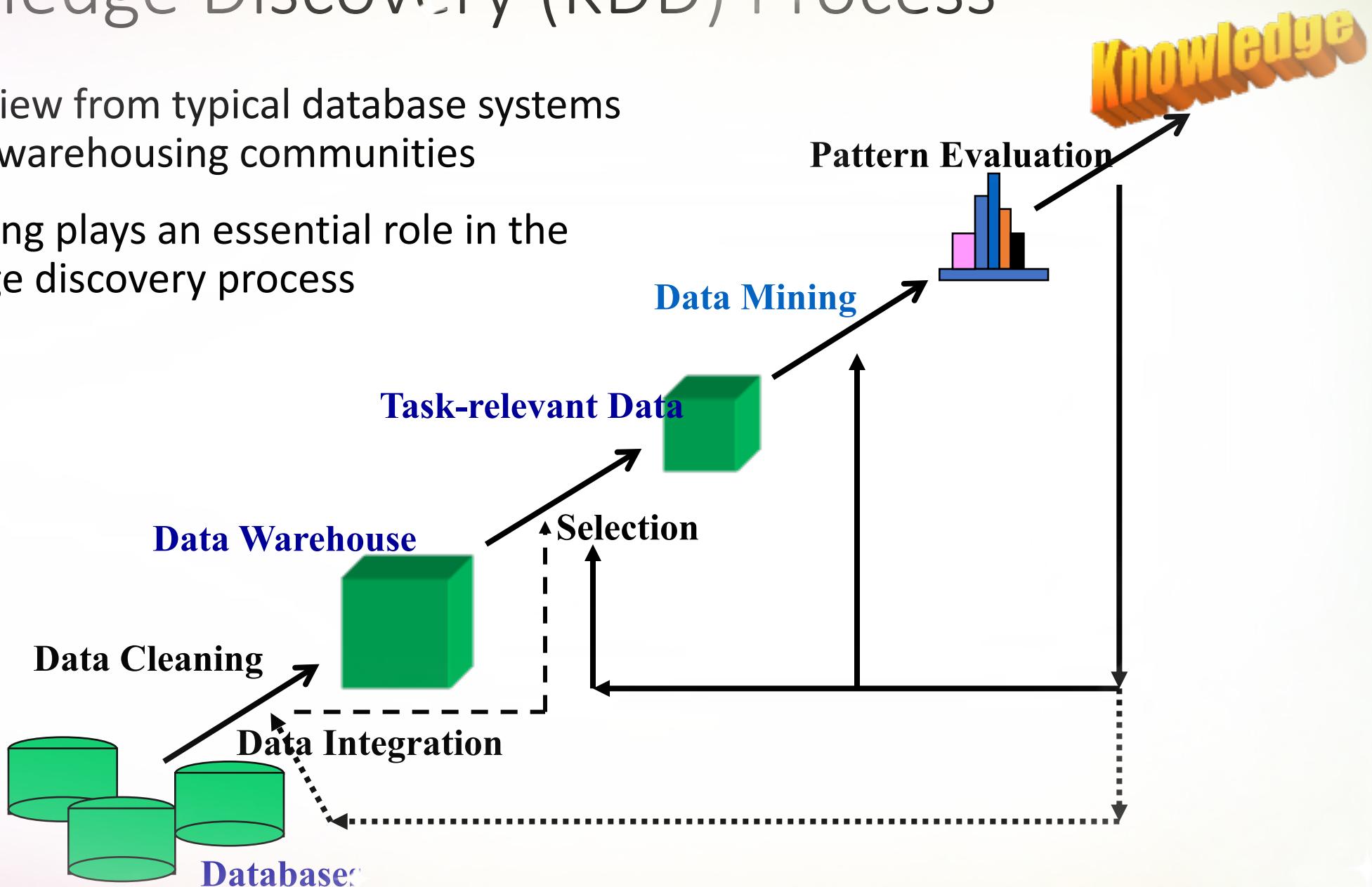
What Is Data Mining?

- Data mining (knowledge discovery from data)
ការបានឃុំពិភាក្សាដែលមានតម្លៃខ្លួន
ការរក្សា
• Extraction of interesting (non-trivial, implicit, previously unknown and potentially useful) patterns or knowledge from huge amount of data
ស្មើកម្រិតណ៍
• Data mining: a misnomer?
សំគាល់ឡើង
- Alternative names:
ការកំណត់ពាណិជ្ជកម្ម
• 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”?
ការច្បាប់អាណាពលការបានឃុំ នៅពេលមិនបានឃុំ ការរក្សាដែលមានតម្លៃខ្លួន
• Simple search and query processing
ទំនួរក្នុងបញ្ជាផល
• (Deductive) expert systems

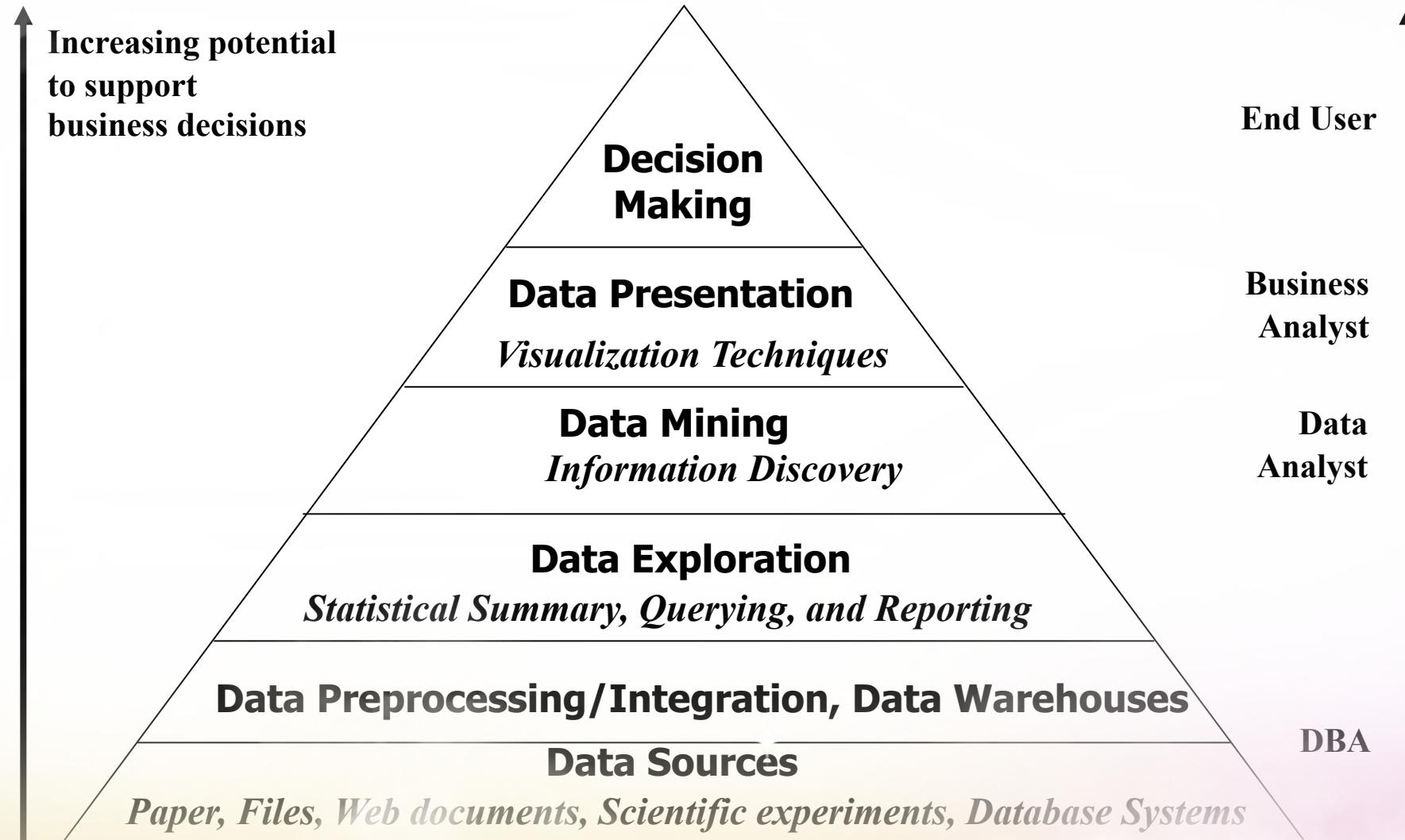


Knowledge Discovery (KDD) Process

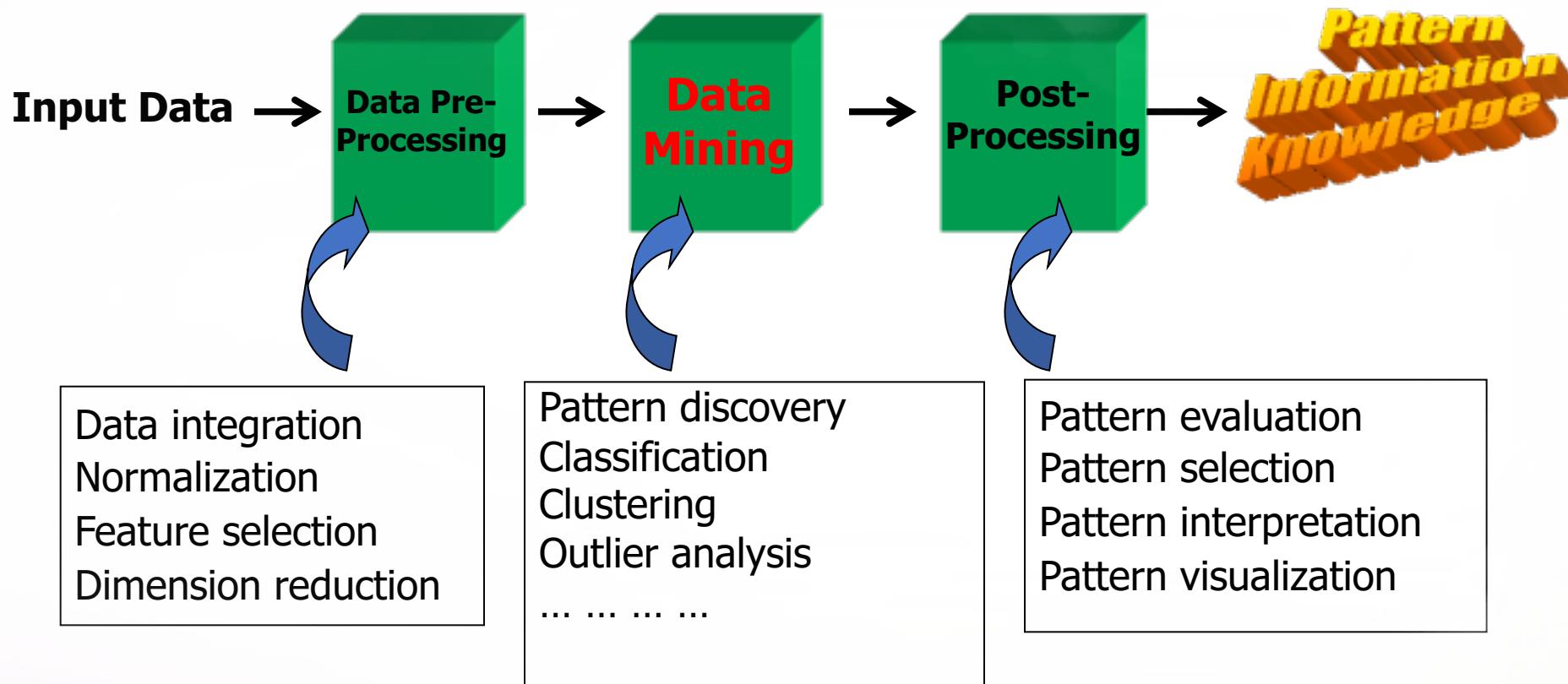
- This is a view from typical database systems and data warehousing communities
- Data mining plays an essential role in the knowledge discovery process



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

ធម្មលសាបនិតិ

Multi-Dimensional View of Data Mining

- Data to be mined ធម្មលដែលត្រួតពិនិត្យ

- 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) ការងារដែលត្រួតពិនិត្យ

- 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 ការគាំទ្ររបស់ខ្លួន

- Retail, telecommunication, banking, fraud analysis, bio-data mining, stock market analysis, text mining, Web mining, etc.

Data Mining: On What Kinds of Data? ជាមួយទៅ?

ចំណាំ

- Database-oriented data sets and applications
 - Relational database, data warehouse, transactional database
 - Object-relational databases, Heterogeneous databases and legacy databases
- Advanced data sets and advanced applications
 - Data streams and sensor data
 - Time-series data, temporal data, sequence data (incl. bio-sequences)
 - Structure data, graphs, social networks and information networks
 - Spatial data and spatiotemporal data
 - Multimedia database
 - Text databases
 - The World-Wide Web

Data Mining Functions: (1) Generalization

ກາງຈົນຈາກແກ້ໄຂ:ສ່ວນຄລົງຂອງ

- Information integration and data warehouse construction

- Data cleaning, transformation, integration, and multidimensional data model ຈຳລັງເບຄນຊື່ຕີ

- Data cube technology

ກາປັບພະດີຕໍ່ຕາມຕາມຄົ່ນຫາ

- Scalable methods for computing (i.e., materializing) multidimensional aggregates
 - OLAP (online analytical processing)

- Multidimensional concept description: Characterization and discrimination ດິນທີ ລັ້ງຄະ: ເດວາ: ແກ້:ກາງ / ລື້ອກປັບປຸງ

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

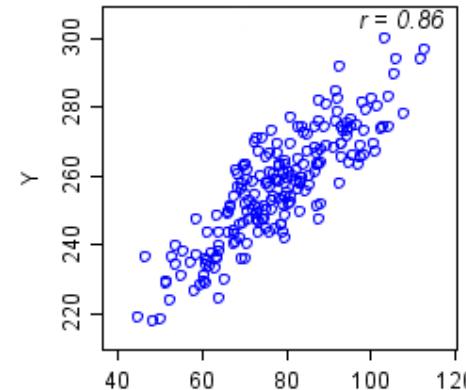
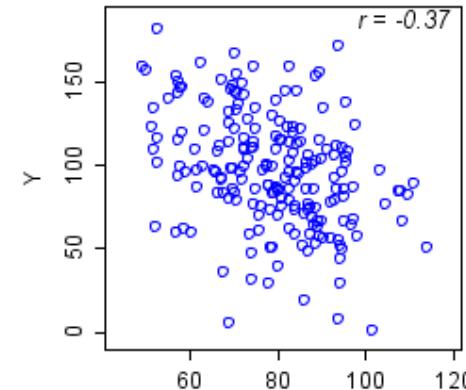
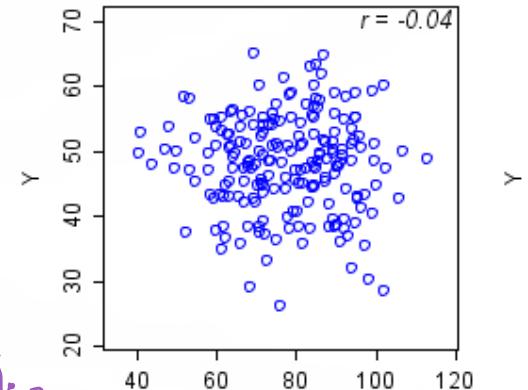
ສຽງ ເນັບທີ່ຍຸດຖານກພ່ອນູກ



Data Mining Functions: (2) Pattern Discovery

ទីផ្សារ ស្តីពី ទុកដាក់ និង ការងារ

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



- A typical association rule  Ex. តាមឱ្យលានឯកសារអេឡិចត្រូនុកម្ពស់បានក្នុងការងារ
- Diaper \rightarrow 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?

Data Mining Functions: (3) Classification

ការរំអាយកតាករ

- Classification and label prediction

- Construct models (functions) based on some training examples
- Describe and distinguish classes or concepts for future prediction
 - Ex. 1. Classify countries based on (climate)
 - Ex. 2. Classify cars based on (gas mileage)
- Predict some unknown class labels

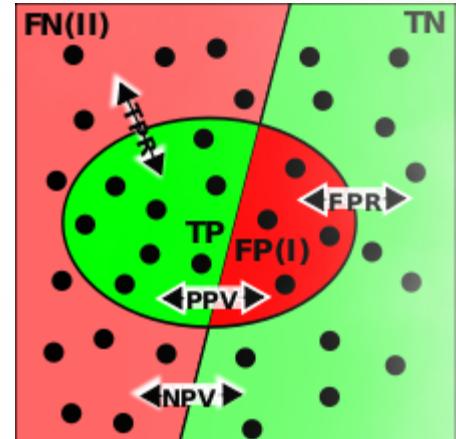
- Typical methods

- 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, ...

ការគ្រប់គ្រងការស្នើសុំជាមួយកម្រិតទី៣
ការគ្រប់គ្រងការស្នើសុំជាមួយកម្រិតទី២



Data Mining Functions: (4) Cluster Analysis

ការបង្កើរបែងចែកពួកគេ នៃ មុខរៀងប្រើប្រាស់

- 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

