

---

# **Data Mining**

# About Myself

---

- Contact Details
  - Room: 28
  - Email: [mghouri@numl.edu.pk](mailto:mghouri@numl.edu.pk)
  - Contact Hours: Mon 8:00 AM – 1:00 PM

# Textbook

---

- Han, J., Kamber , M., & Pei, J. “Data Mining: Concepts and Techniques”, Latest Edition, Morgan Kaufmann

# Evaluation Criteria


---

- Subject to change as per policy

Activity	Evaluation Percentage
Assignments	15
Quizzes	10
Project	10
Mid Term	25
End Term	40

# 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


# 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

# 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

# What Is Data Mining?



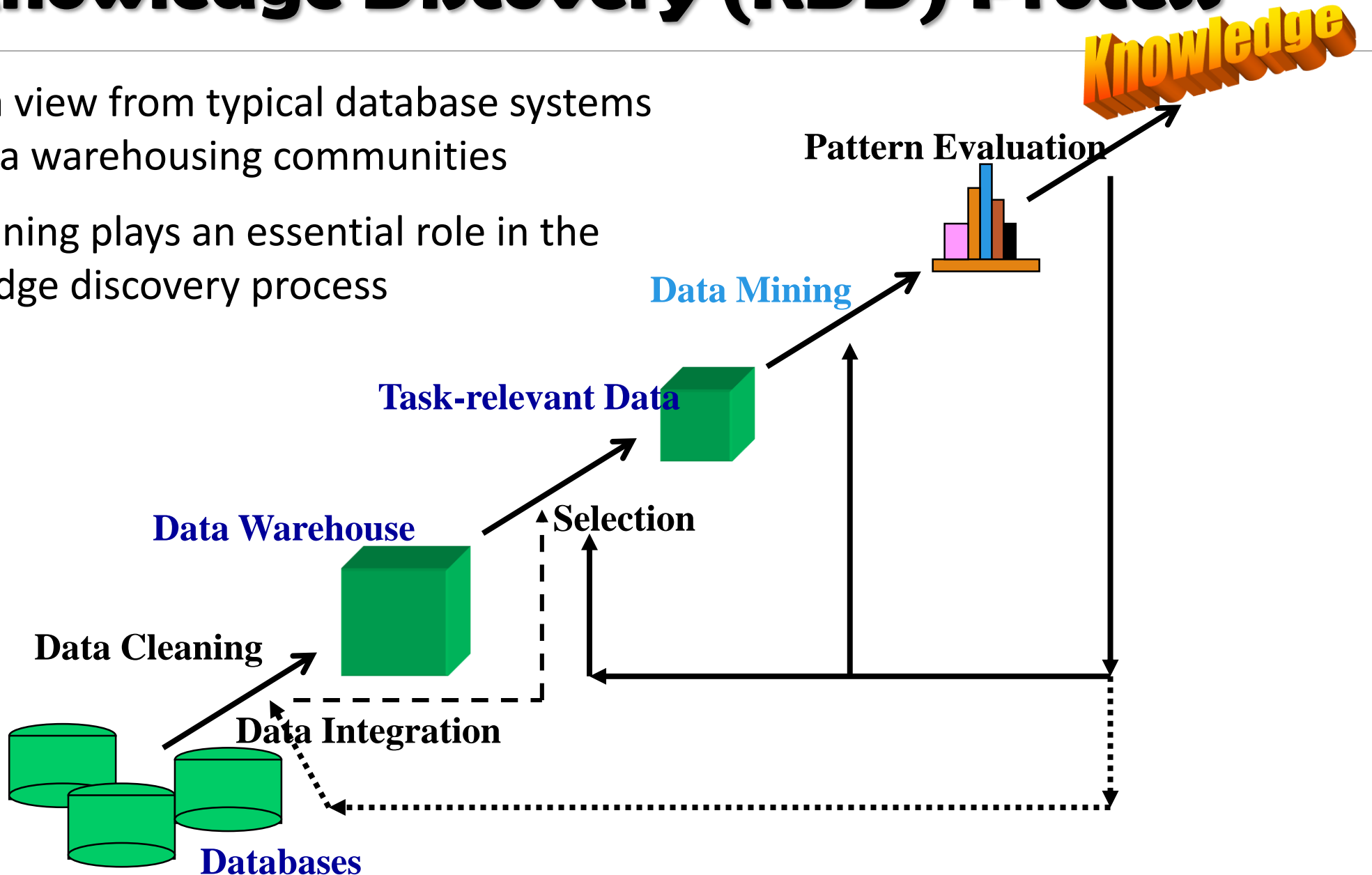
- ❑ Data mining (knowledge discovery from data)
  - ❑ Extraction of interesting (non-trivial, hidden, previously unknown and potentially useful) patterns or knowledge from huge amount of data
- ❑ Alternative names
  - ❑ Knowledge discovery (mining) in databases (KDD), knowledge extraction, data/pattern analysis, data archeology, data dredging, information harvesting, business intelligence, etc.
- ❑ 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

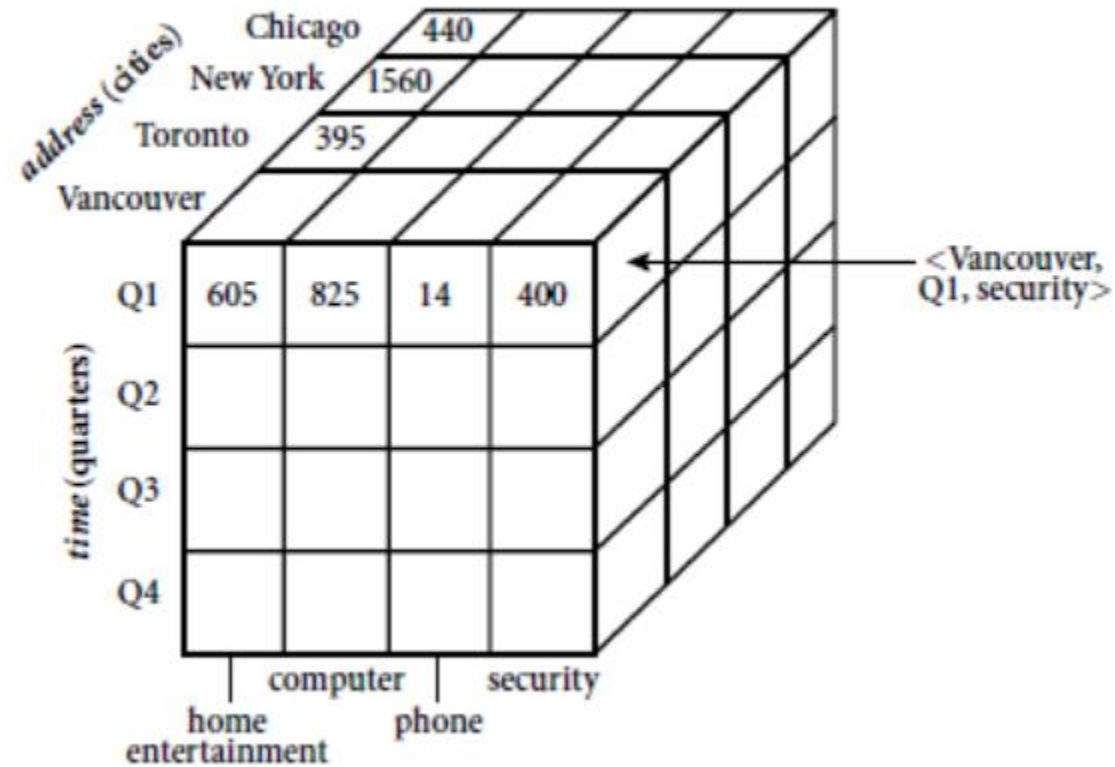


# Example: A Web Mining Framework

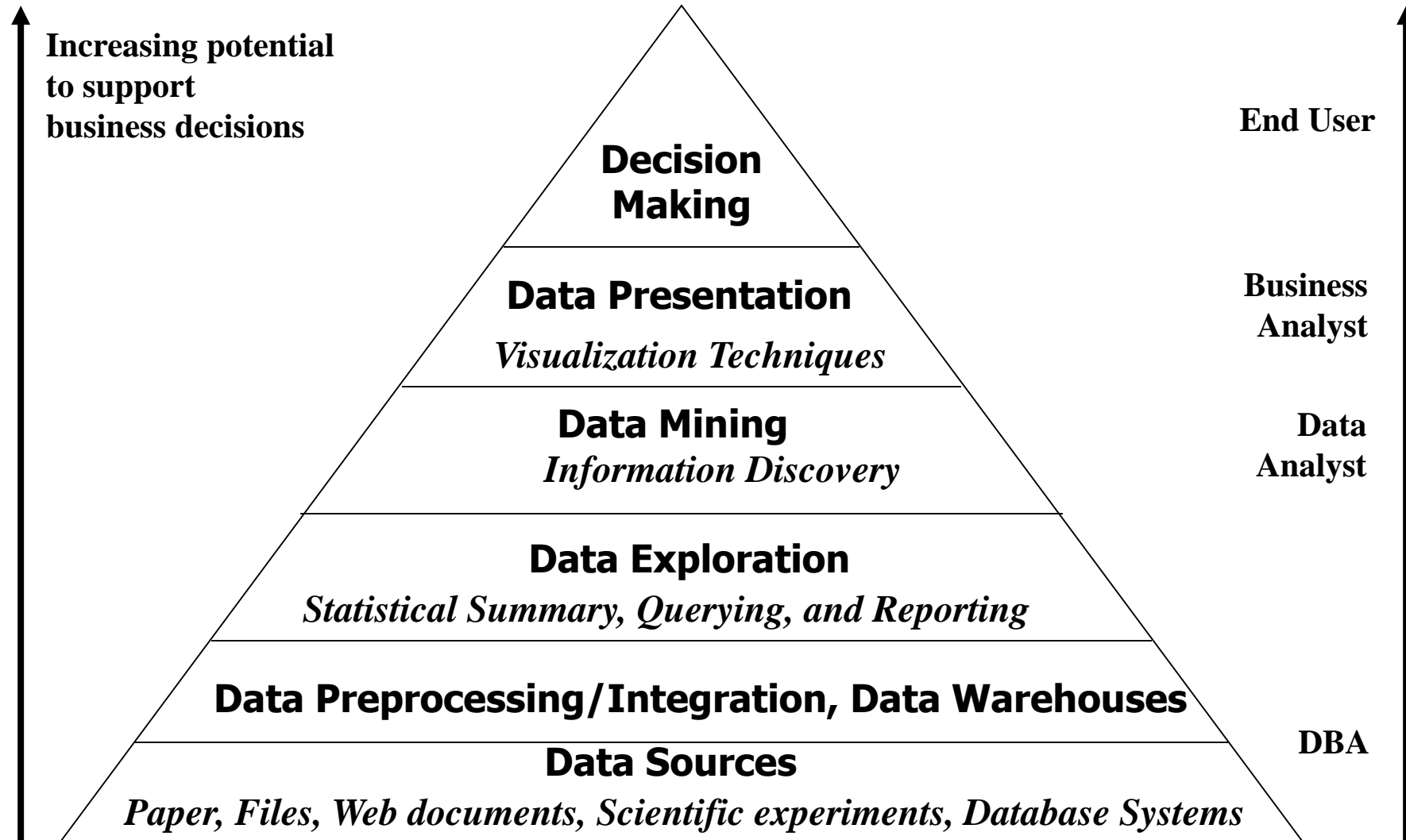
---

- ❑ 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

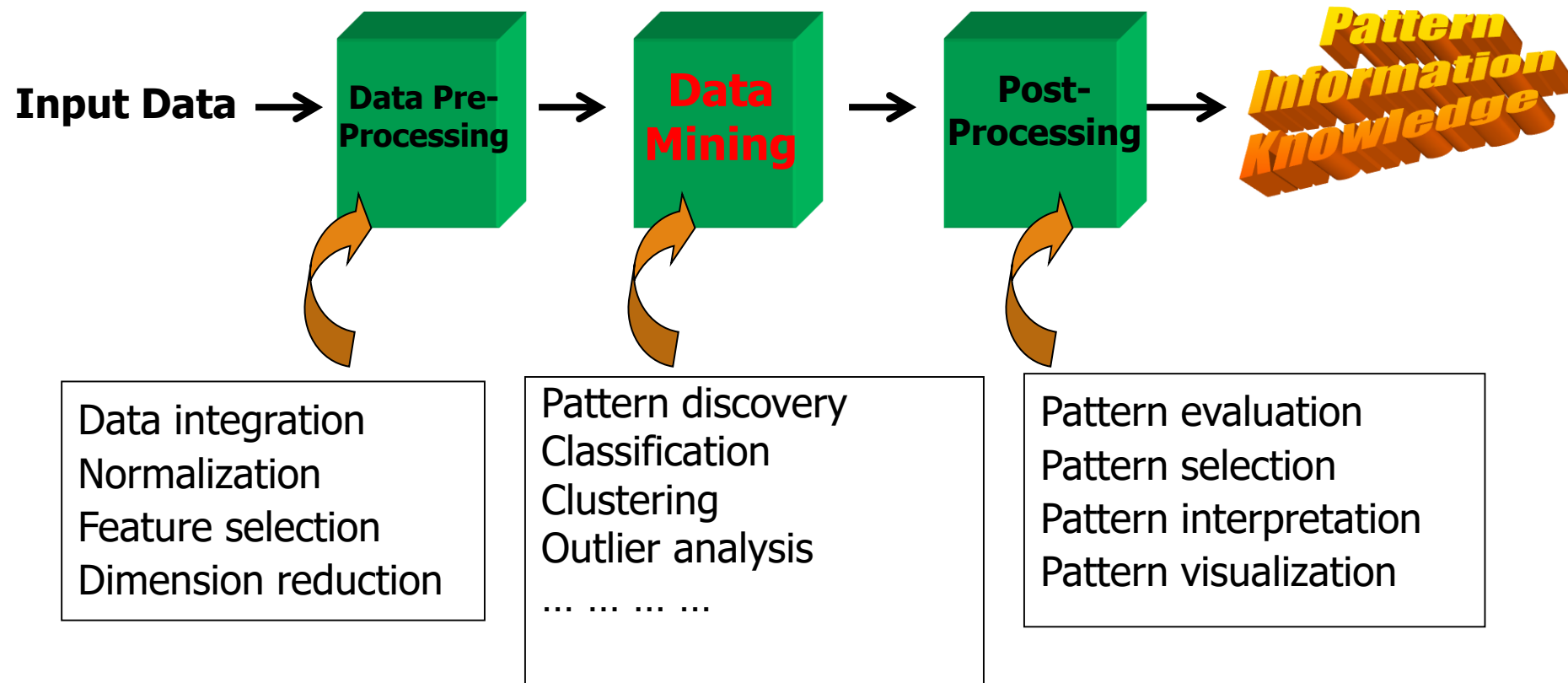
# Data Cube



# Data Mining in Business Intelligence




# KDD Process: A View from ML and Statistics



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

# 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

# Multi-Dimensional View of Data Mining

---

## ❑ Data to be mined

- ❑ Database data (extended-relational, object-oriented, heterogeneous), data warehouse, transactional data, stream (e.g. video), spatiotemporal (e.g. maps), time-series and sequence (e.g. stock), text (customer sentiment) and web, multi-media, graphs & social and information networks

## ❑ Knowledge to be mined (or: Data mining functions)

- ❑ Characterization (summarization of general features of a target class. E.g. characteristics of software products whose sales were increased by 10% last year)
- ❑ Discrimination (comparison of features of contrasting classes. Eg. Comparing features of dataset where software were increased by 10% Vs. which were decreased by 30%)
- ❑ Association, classification, clustering, trend/deviation, outlier analysis
- ❑ Descriptive (characterize properties of data) vs. predictive data mining (predictions)

# Multi-Dimensional View of Data Mining

---

## ❑ Techniques utilized

- ❑ Data-intensive, data warehouse (OLAP), machine learning, statistics, pattern recognition, visualization, high-performance, etc.

## ❑ Applications adapted

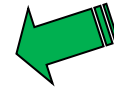
- ❑ Retail, telecommunication, banking, fraud analysis, bio-data mining, stock market analysis, text mining, Web mining, etc.



# 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
- ❑ Summary



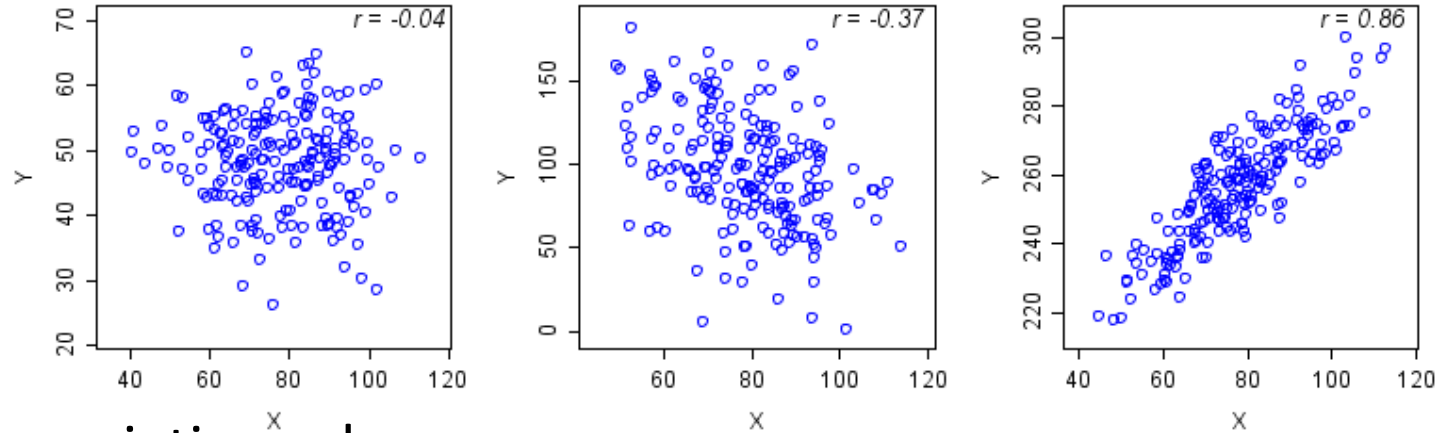
# Data Mining Functions: (1) Generalization

- ❑ Information integration and data warehouse construction
  - ❑ Data cleaning, transformation, integration, and multidimensional data model
- ❑ Data cube technology
  - ❑ Multidimensional aggregates
  - ❑ OLAP (online analytical processing)
- ❑ Multidimensional concept description: Characterization and discrimination
  - ❑ Generalize, summarize, and contrast data characteristics, e.g., dry vs. wet region



# Data Mining Functions: (2) Pattern Discovery

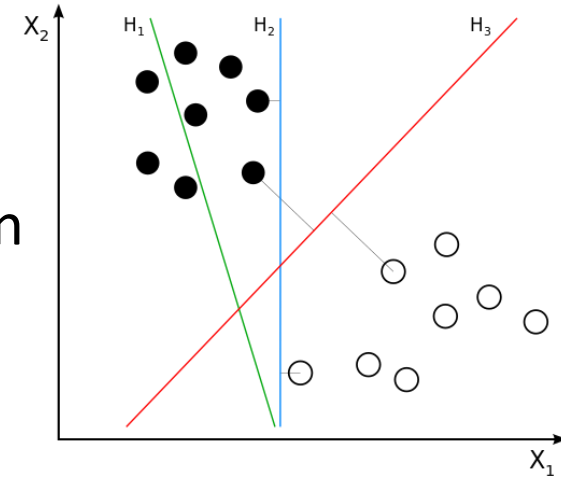
- Frequent patterns (or frequent itemsets)
  - What items are frequently purchased together in your Walmart?
- Association and Correlation Analysis



- A typical association rule
  - Diaper  $\rightarrow$  Juice [0.5%, 75%] (support, confidence) ( $P(XUY), P(X|Y)$ )
- How to mine such patterns and rules efficiently in large datasets?
- How to use such patterns for classification, clustering, and other applications?

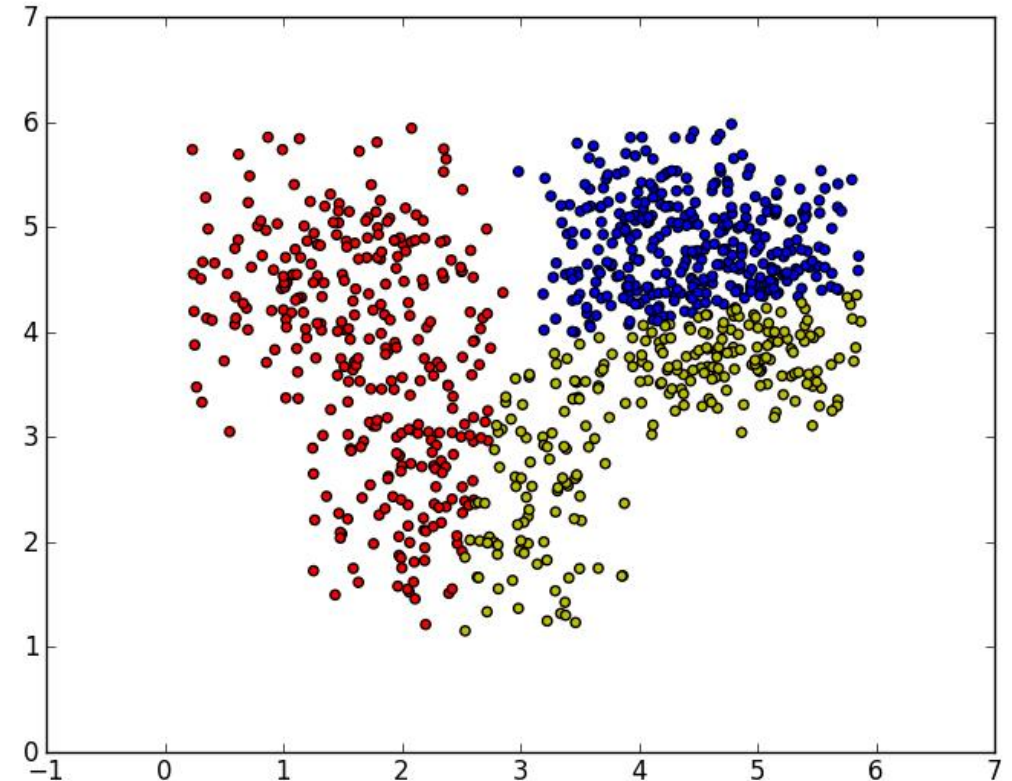
# Data Mining Functions: (3) Classification

- ❑ 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 diseases, web-pages, ...



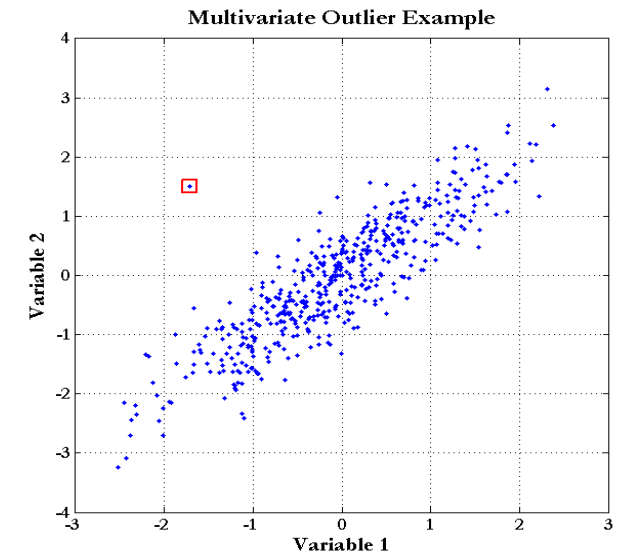
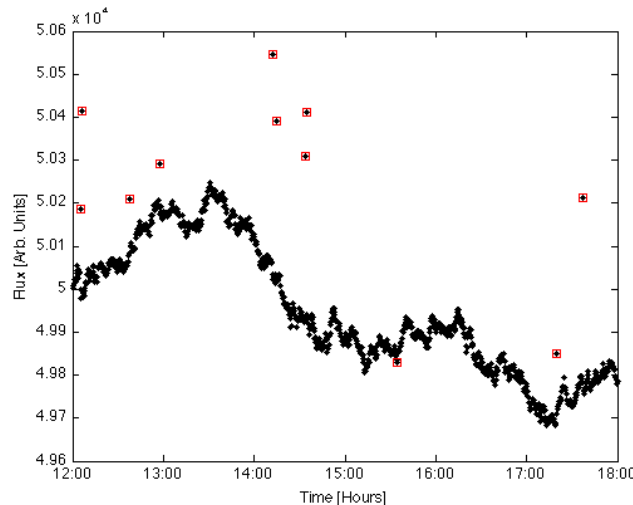
# Data Mining Functions: (4) Cluster Analysis

- ❑ 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
  - ❑ Like market segmentation, community detection, improving search results, ...



# Data Mining Functions: (5) Outlier Analysis

- ❑ Outlier analysis
  - ❑ Outlier: A data object that does not comply with the general behavior of the data
  - ❑ Noise or exception?—One person's garbage could be another person's treasure
  - ❑ Methods: by product of clustering or regression analysis, ...
  - ❑ Useful in fraud detection, rare events analysis



# Evaluation of Knowledge

---

- ❑ Are all mined knowledge interesting?
  - ❑ One can mine tremendous amount of “patterns”
  - ❑ Some may fit only certain dimension space (time, location, ...)
  - ❑ Some may not be representative, may be temporary, ...
- ❑ Evaluation of mined knowledge
  - ❑ Descriptive vs. predictive
  - ❑ Coverage
  - ❑ Typicality vs. novelty
  - ❑ Accuracy
  - ❑ Relevance
  - ❑ ...





# Chapter 1. Introduction

---

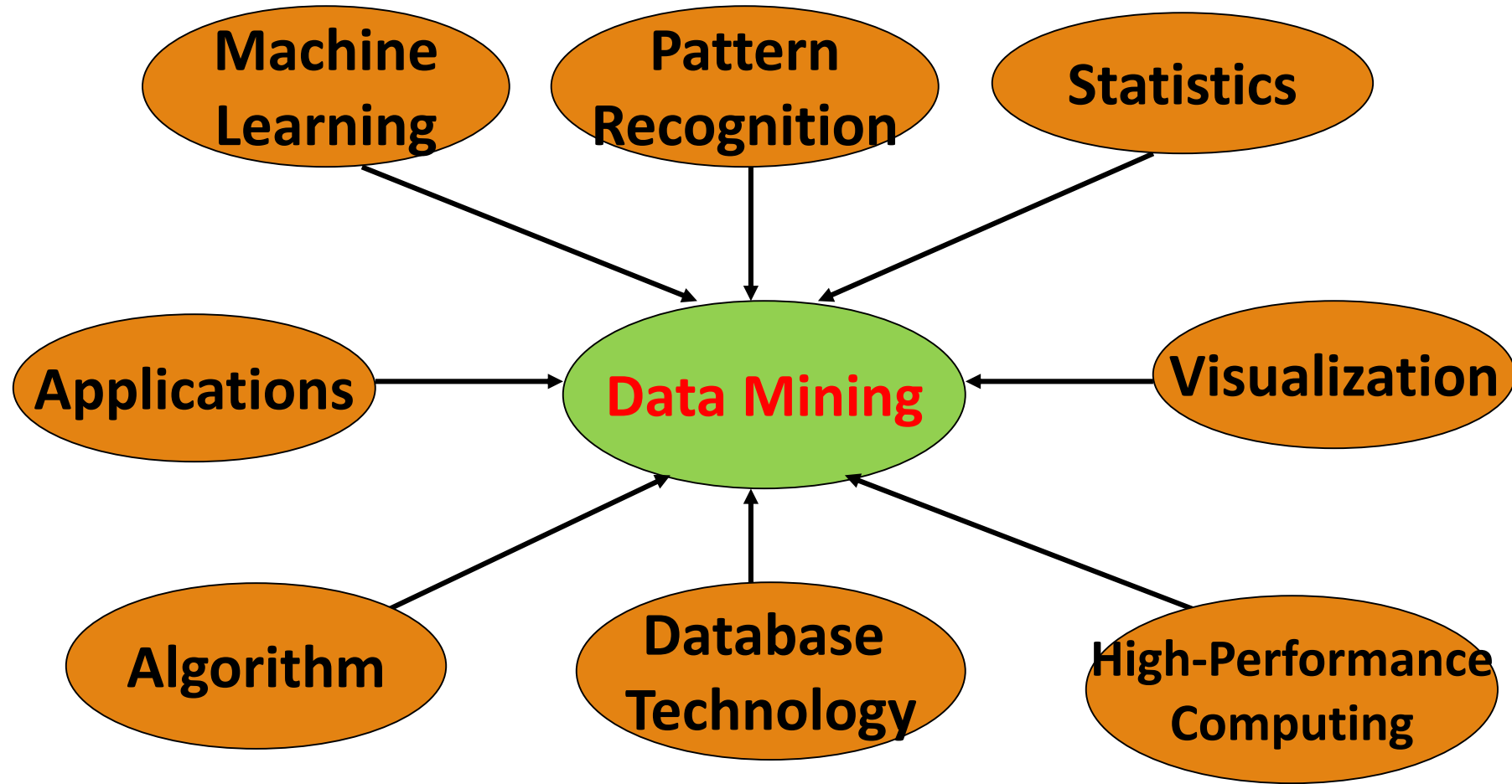
- ❑ Why Data Mining?
- ❑ What Is Data Mining?
- ❑ A Multi-Dimensional View of Data Mining
- ❑ What Kinds of Data Can Be Mined?
- ❑ What Kinds of Patterns Can Be Mined?
- ❑ What Kinds of Technologies Are Used?
- ❑ What Kinds of Applications Are Targeted?
- ❑ Major Issues in Data Mining
- ❑ Summary





# Data Mining: Confluence of Multiple Disciplines

---



# Why Confluence of Multiple Disciplines?

---

- ❑ Tremendous amount of data
  - ❑ Algorithms must be scalable to handle big data
- ❑ High-dimensionality of data
  - ❑ May have tens of thousands of dimensions
- ❑ High complexity of data
  - ❑ Data streams and sensor data (temp, humidity, air pressure, gps, heart rate etc.)
  - ❑ Time-series data, temporal data, sequence data
  - ❑ Structure data, graphs, social and information networks
  - ❑ Spatial (3D), spatiotemporal (maps), multimedia, text and Web data
  - ❑ Software programs, scientific simulations
- ❑ New and sophisticated applications


# 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
- ❑ Summary



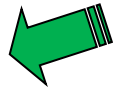
# Applications of Data Mining

- ❑ Web page analysis: classification, clustering, ranking
  - ❑ Recommender systems
  - ❑ Basket data analysis to targeted marketing
  - ❑ Biological and medical data analysis
  - ❑ Data mining and text analysis
  - ❑ Data mining and social and information network analysis
  - ❑ Data mining and software engineering (e.g. bug mining, i.e. mining of software bugs in large programs)
  - ❑ Built-in (invisible data mining) functions in Google, MS, Yahoo!, Linked, Facebook, ...
- 



# 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 
- ❑ Summary

# Major Issues in Data Mining (1)

---

- ❑ Mining Methodology
  - ❑ Mining various and new kinds of knowledge
  - ❑ Mining knowledge in multi-dimensional space (e.g. cube)
  - ❑ Data mining: An interdisciplinary effort
  - ❑ Boosting the power of discovery in a networked environment
  - ❑ Handling noise, uncertainty, and incompleteness of data
  - ❑ Pattern evaluation and pattern- or constraint-guided mining
- ❑ User Interaction
  - ❑ Interactive mining
  - ❑ Incorporation of background knowledge
  - ❑ Presentation and visualization of data mining results

# Major Issues in Data Mining (2)

---

- ❑ Efficiency and Scalability
  - ❑ Efficiency and scalability of data mining algorithms
  - ❑ Parallel and incremental mining methods (dealing with new input data)
- ❑ Diversity of data types
  - ❑ Handling complex types of data
  - ❑ Mining dynamic, networked, and global data repositories
- ❑ Data mining and society
  - ❑ Social impacts of data mining (good use in society vs misuse)
  - ❑ Privacy-preserving data mining
  - ❑ Invisible data mining (web search engines, internet-based stores)

# 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
- ❑ Summary 



# Summary

---

- ❑ 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
- ❑ 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, 3<sup>rd</sup> ed. , 2011
- ❑ T. Hastie, R. Tibshirani, and J. Friedman, The Elements of Statistical Learning: Data Mining, Inference, and Prediction, 2<sup>nd</sup> 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 (2<sup>nd</sup> ed. 2016)
- ❑ I. H. Witten and E. Frank, Data Mining: Practical Machine Learning Tools and Techniques with Java Implementations, Morgan Kaufmann, 2<sup>nd</sup> ed. 2005
- ❑ Mohammed J. Zaki and Wagner Meira Jr., Data Mining and Analysis: Fundamental Concepts and Algorithms 2014