# Information Processing Technology of Internet of Things

## Lecture 1: Introduction

Wu Liu

Beijing University of Posts and Telecommunications

# *Course Information*

- Prerequisites:
  - Data structures, Database, Linear Algebra

- Required Textbooks:
  - Data Mining-Concepts and Techniques (Third Edition). Jiawei Han, Micheline Kamber, Jian Pei
  - Modern Information Retrieval-The Concepts and Technology behind Search (second edition). Ricardo Baeza-Yates, Berthier Ribeiro-Neto
  - Image Processing, Analysis, and Machine Vision (Third Edition). Milan Sonka, Vaclav Hlavac, Roger Boyle

- References:
  - …

# *Course outlines*

- Chapter 0: Introduction 1
- Chapter 1: Data Preprocessing 2
- Chapter 2: Data Mining 4
- Chapter 3: Information Retrieval 4
- Chapter 4: Visual Information Processing 3
- Mid-term Test 1
- Review 1

- Homework 15%，Mid-Term 15%，Final 70%

# *Chapter 0 Introduction*

# *What Is Data Mining?*

- Data mining (knowledge discovery from data)
  - Extraction of interesting (<u>non-trivial,</u> <u>implicit</u>, <u>previously unknown</u> and <u>potentially useful)</u> 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
  - Expert systems

# *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
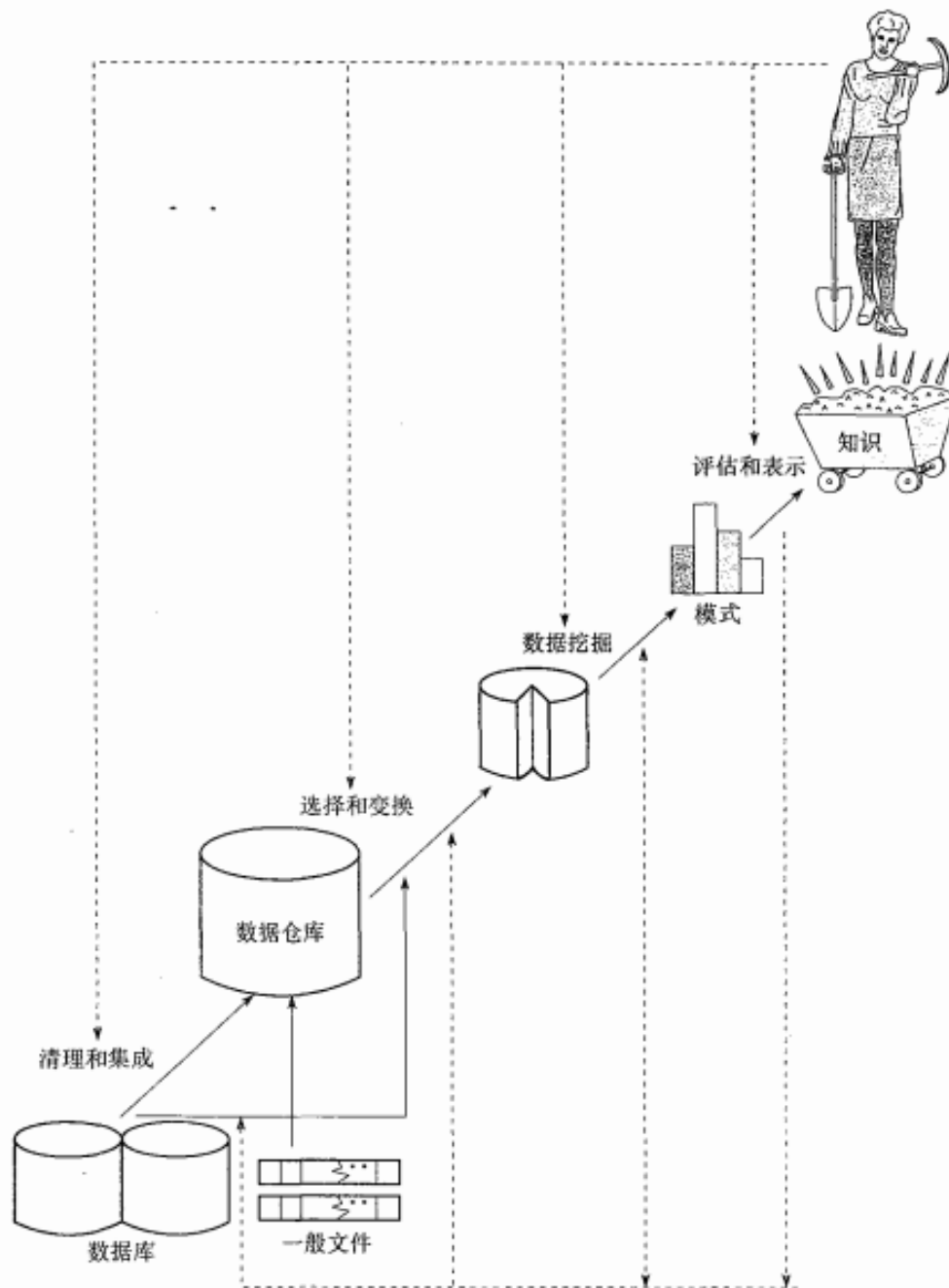  - Patterns and knowledge to be used or stored into knowledge-base
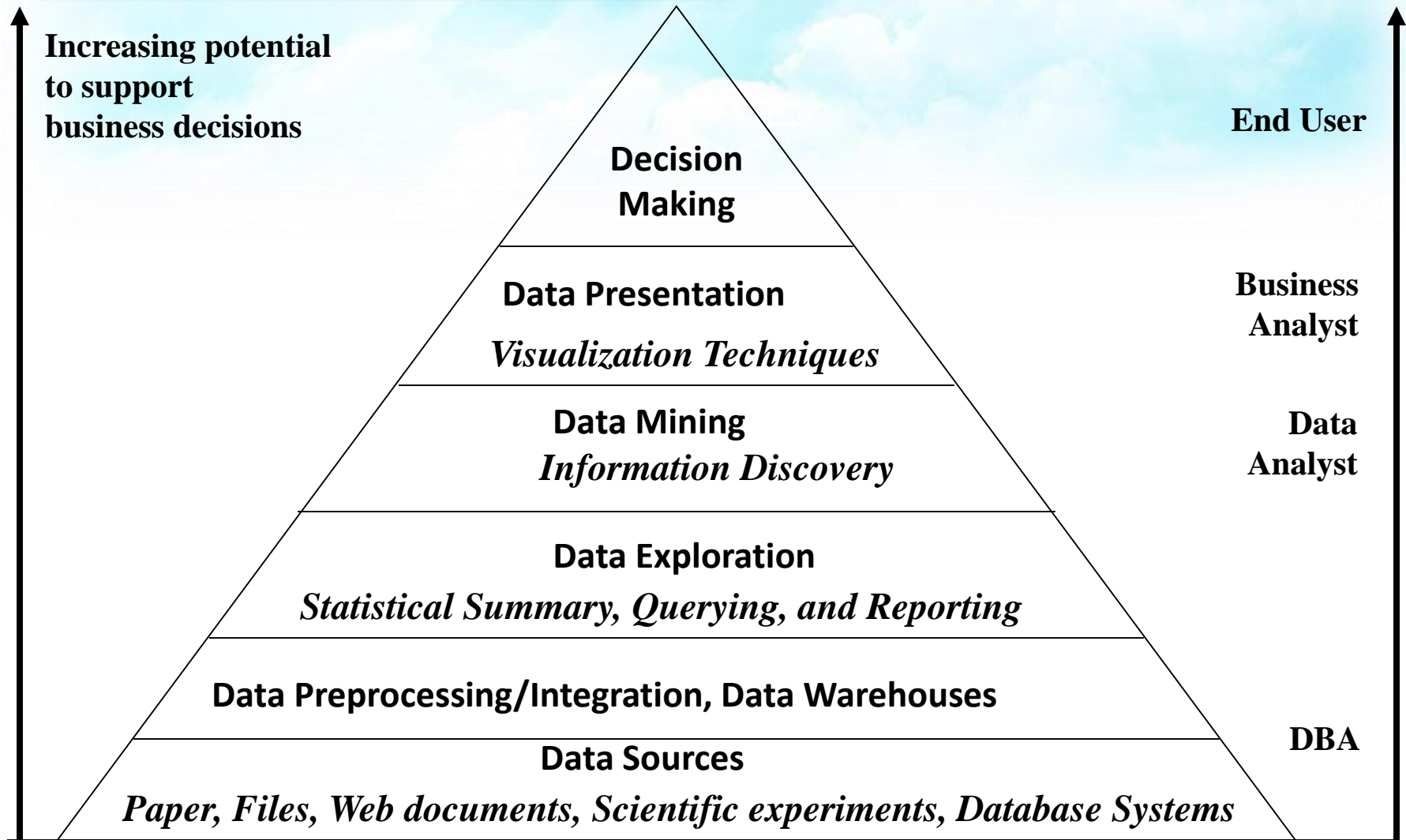
# *Data Mining*



评估和表示

知识

模式

数据挖掘

选择和变换

数据仓库

清理和集成

数据库     一般文件

图 1.4 数据挖掘视为知识发现过程的一个步骤

# Data Mining in Business Intelligence

Increasing potential
to support
business decisions

**Decision Making**

**Data Presentation**
*Visualization Techniques*

**Data Mining**
*Information Discovery*

**Data Exploration**
*Statistical Summary, Querying, and Reporting*

**Data Preprocessing/Integration, Data Warehouses**

**Data Sources**
*Paper, Files, Web documents, Scientific experiments, Database Systems*

**End User**

**Business Analyst**

**Data Analyst**

**DBA**

# Multi-Dimensional View of Data Mining

- **Data to be mined**
  - Database data (extended-relational, object-oriented, heterogeneous, legacy), 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)**
  - Characterization, discrimination, association, classification, clustering, trend/deviation, outlier analysis, etc.
  - Descriptive vs. predictive data mining
- **Techniques utilized**
  - 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.

# *Data Mining: On What Kinds of Data?*

- Database-oriented data sets and applications

  - Relational database, data warehouse, transactional database

- 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 Function: (1) Generalization*

- Information integration and data warehouse construction
  - Data cleaning, transformation, integration, and multidimensional data model
- Data cube technology
  - Scalable methods for computing (i.e., materializing) 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 Function: (2) Association and Correlation Analysis*

- Frequent patterns (or frequent itemsets)

  - What items are frequently purchased together in your Walmart?

- Association, correlation vs. causality

  - A typical association rule

    - 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 Function: (3) Classification*

- Classification and label prediction
  - Construct models (functions) based on some training examples
  - Describe and distinguish classes or concepts for future prediction
    - E.g., classify countries based on (climate), or 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, …

# *Data Mining Function: (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

# *Data Mining Function: (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

# *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., first buy digital camera, then buy large SD 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

# *Structure and Network Analysis*

- Graph mining
  - Finding frequent subgraphs (e.g., chemical compounds), trees (XML), substructures (web fragments)
- Information network analysis
  - Social networks: actors (objects, nodes) and relationships (edges)
    - 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
    - Web community discovery, opinion mining, usage mining, …

# 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 transient, …
- Evaluation of mined knowledge → directly mine only interesting knowledge?
  - Descriptive vs. predictive
  - Coverage
  - Typicality vs. novelty
  - Accuracy
  - Timeliness
  - …

# *Applications of Data Mining*

- Web page analysis: from web page classification, clustering to PageRank & HITS algorithms

- Collaborative analysis & recommender systems

- Basket data analysis to targeted marketing

- Biological and medical data analysis: classification, cluster analysis (microarray data analysis),  biological sequence analysis, biological network analysis

- …