

CSC5312 DATA MINING AND DATA WAREHOUSING

Introduction

2

Data Mining: Concepts and Techniques

Lecture 1. Introduction

3

- Motivation: Why data mining?
- What is data mining?
- Data Mining: On what kind of data?
- Data mining functionality
- Classification of data mining systems
- Major issues in data mining
- Overview of the course

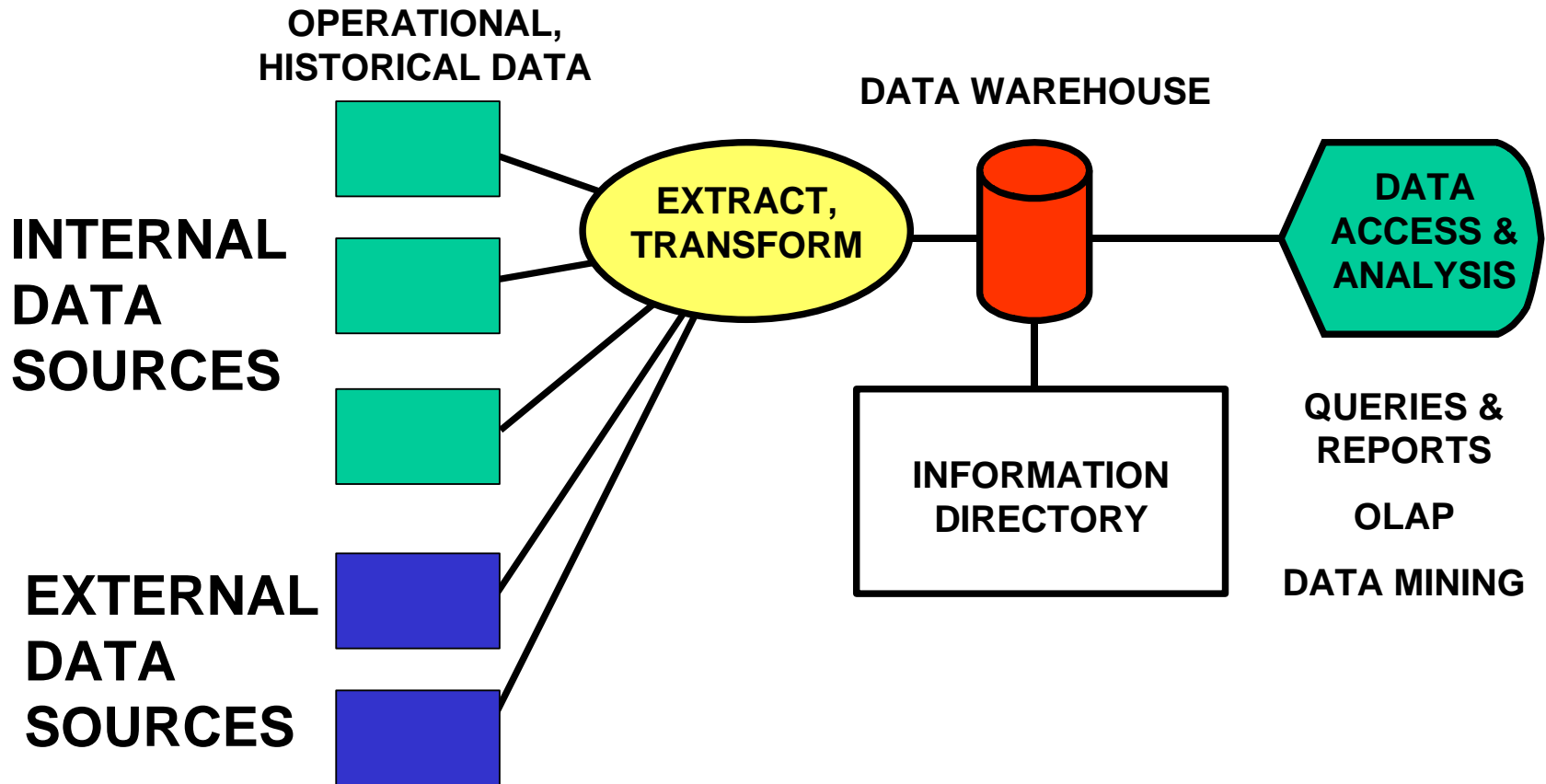
Why Data Mining?

4

- The Explosive Growth of Data: from terabytes to petabytes
- Data are any facts, numbers, images or text that can be processed by a computer.
 - 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!
 - An urgent need for transforming data into useful **information** and **knowledge**.

Components of Data Warehouse

5



Evolution of Sciences

6 Before 1600, **empirical science**

- 1600-1950s, **theoretical science**
 - ▣ Each discipline has grown a *theoretical* component. Theoretical models often motivate experiments and generalize our understanding.
- 1950s-1990s, **computational science**
 - ▣ Over the last 50 years, most disciplines have grown a third, *computational* branch (e.g. empirical, theoretical, and computational ecology, or physics, or linguistics.)
 - ▣ Computational Science traditionally meant simulation. It grew out of our inability to find closed-form solutions for complex mathematical models.
- 1990-now, **data science**
 - ▣ The flood of data from new scientific instruments and simulations
 - ▣ The ability to economically store and manage petabytes of data online
 - ▣ The Internet and computing Grid that makes all these archives universally accessible
 - ▣ Scientific info. management, acquisition, organization, query, and visualization tasks scale almost linearly with data volumes. **Data mining** is a major new challenge!
- Jim Gray and Alex Szalay, *The World Wide Telescope: An Archetype for Online Science*, Comm. ACM, 45(11): 50-54, Nov. 2002

Evolution of Database Technology

7

1960s:

- ▣ Data collection, database creation, IMS and network DBMS

1970s:

- ▣ Relational data model, relational DBMS implementation

1980s:

- ▣ RDBMS, advanced data models (extended-relational, OO, deductive, etc.)
- ▣ Application-oriented DBMS (spatial, scientific, engineering, etc.)

1990s:

- ▣ Data mining, data warehousing, multimedia databases, and Web databases

2000s

- ▣ Stream data management and mining
- ▣ Data mining and its applications
- ▣ Web technology (XML, data integration) and global information systems

What Is Data Mining?



8

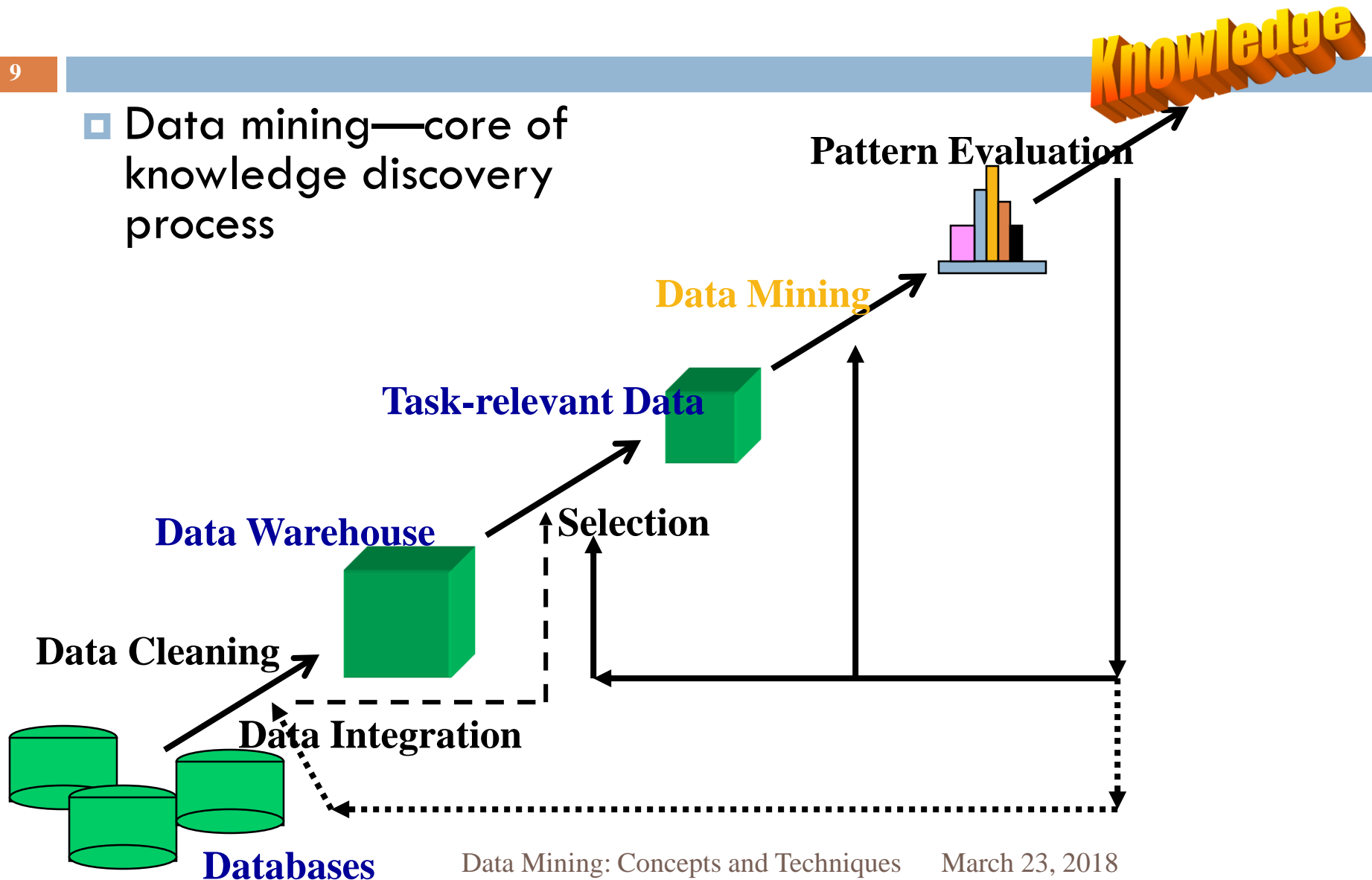
- Data mining (knowledge discovery from data)
 - ▣ Extraction of interesting (non-trivial, implicit, previously unknown and potentially useful) patterns or knowledge from huge amount of data
 - ▣ Data mining: a misnomer?
- Alternative names
 - ▣ Knowledge discovery (mining) in databases (KDD), knowledge 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



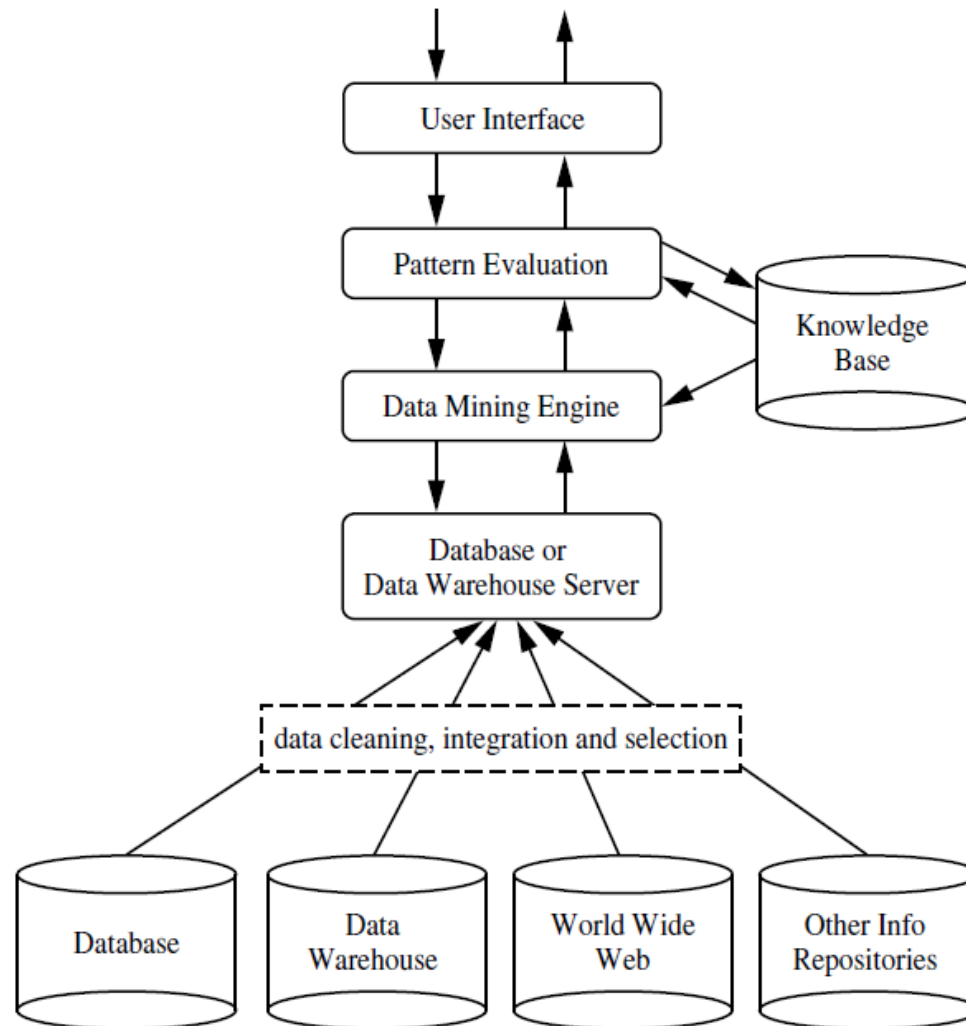
Knowledge Discovery (KDD) Process

9

- Data mining—core of knowledge discovery process

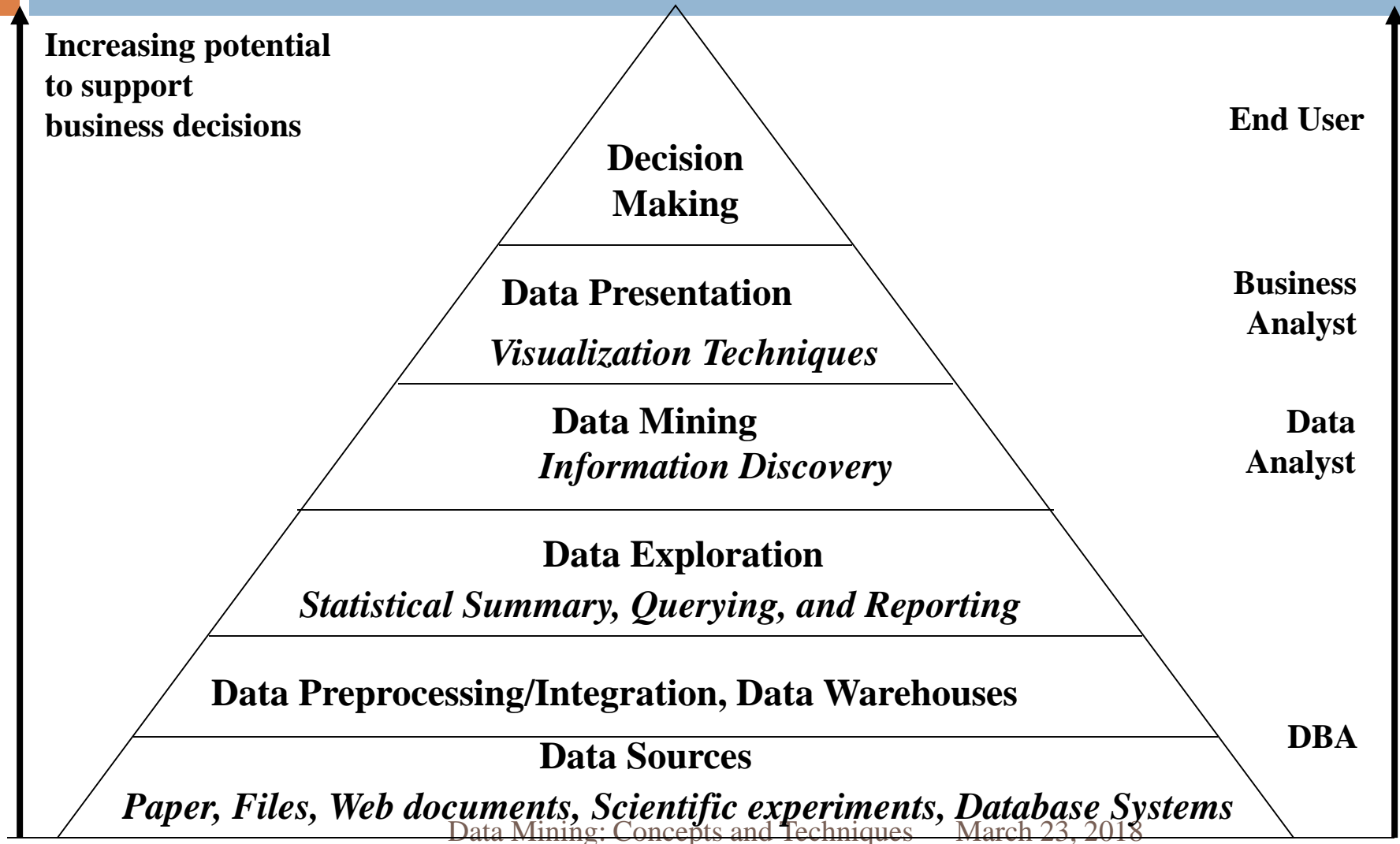


10



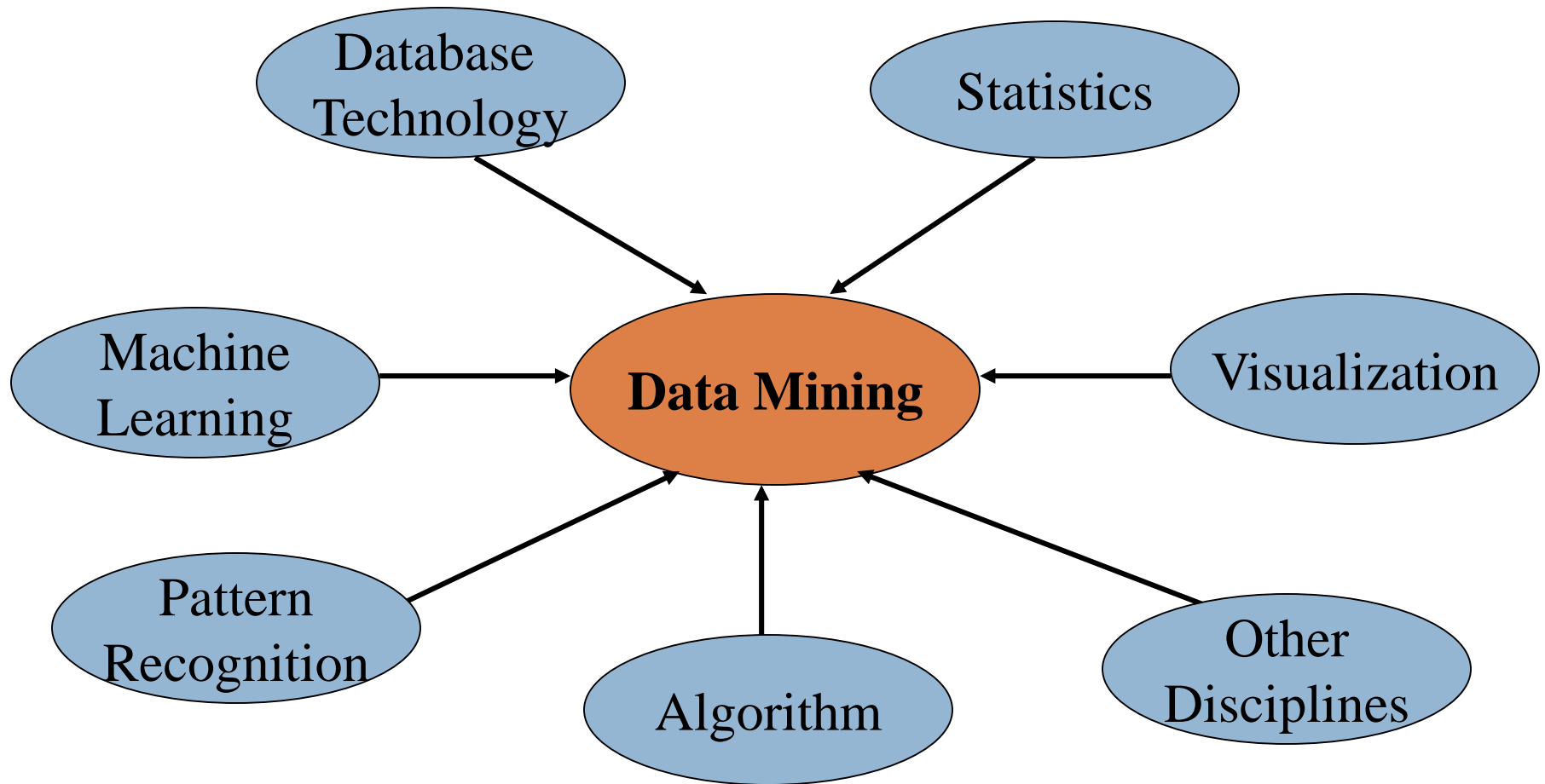
Data Mining and Business Intelligence

11



Data Mining: Confluence of Multiple Disciplines

12



Why Not Traditional Data Analysis?

13

- Tremendous amount of data
 - Algorithms must be highly scalable to handle such as tera-bytes of data
- High-dimensionality of data
 - Micro-array may have tens of thousands of dimensions
- High complexity of data
 - Data streams and sensor data
 - Time-series data, temporal data, sequence data
 - Structure data, graphs, social networks and multi-linked data
 - Heterogeneous databases and legacy databases
 - Spatial, spatiotemporal, multimedia, text and Web data
 - Software programs, scientific simulations
- New and sophisticated applications

Multi-Dimensional View of Data Mining

14

Data to be mined

- ▣ Relational, data warehouse, transactional, stream, object-oriented/relational, active, spatial, time-series, text, multi-media, heterogeneous, legacy, WWW

Knowledge to be mined

- ▣ Characterization, discrimination, association, classification, clustering, trend/deviation, outlier analysis, etc.
- ▣ Multiple/integrated functions and mining at multiple levels

Techniques utilized

- ▣ Database-oriented, data warehouse (OLAP), machine learning, statistics, visualization, etc.

Applications adapted

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

Data Mining: Classification Schemes

15

- General functionality
 - ▣ Descriptive data mining
 - ▣ Predictive data mining
- Descriptive data mining: describes concepts or task relevant data sets in concise, summarative, informative, discriminative forms
- Predictive mining: Based on data and analysis, constructs models for the database, and predicts the trend and properties of unknown data

Data Mining: Classification Schemes

16

- Different views lead to different classifications
 - ▣ **Data** view: Kinds of data to be mined
 - ▣ **Knowledge** view: Kinds of knowledge to be discovered
 - ▣ **Method** view: Kinds of techniques utilized
 - ▣ **Application** view: Kinds of applications adapted

Data Mining: On What Kinds of Data?

17

- 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 multi-linked data
 - Object-relational databases
 - Heterogeneous databases and legacy databases
 - Spatial data and spatiotemporal data
 - Multimedia database
 - Text databases
 - The World-Wide Web

Data Mining Functionalities

18

- Multidimensional concept description: Characterization and discrimination
 - ▣ Generalize, summarize, and contrast data characteristics, e.g., dry vs. wet regions
- Frequent patterns, association, correlation vs. causality
 - ▣ Diaper → Beer [0.5%, 75%] (Correlation or causality?)
- Classification and prediction
 - ▣ Construct models (functions) that 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 or missing numerical values

Data Mining Functionalities (2)

19

Cluster analysis

- ▣ Class label is unknown: Group data to form new classes, e.g., cluster houses to find distribution patterns
- ▣ Maximizing intra-class similarity & minimizing interclass similarity

Outlier analysis

- ▣ Outlier: Data object that does not comply with the general behavior of the data
- ▣ Noise or exception? Useful in fraud detection, rare events analysis

Trend and evolution analysis

- ▣ Trend and deviation: e.g., regression analysis
- ▣ Sequential pattern mining: e.g., digital camera → large SD memory
- ▣ Periodicity analysis
- ▣ Similarity-based analysis

Other pattern-directed or statistical analyses

Summary

20

- Data mining: Discovering interesting patterns from large amounts of data
- A natural evolution of database 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 information repositories
- Data mining functionalities: characterization, discrimination, association, classification, clustering, outlier and trend analysis, etc.
- Data mining systems and architectures
- Major issues in data mining