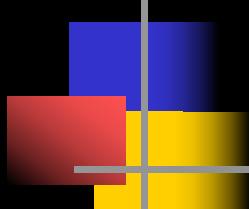


Data Mining: Concepts and Techniques

— Slides for Textbook —
— Chapter 1 —

©Jiawei Han and Micheline Kamber
Intelligent Database Systems Research Lab
School of Computing Science
Simon Fraser University, Canada

<http://www.cs.sfu.ca>

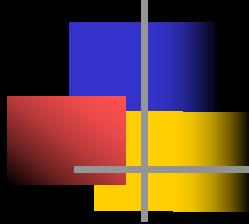


Acknowledgements

- n This work on this set of slides started with my (Han's) tutorial for UCLA Extension course in February 1998
- n Dr. **Hongjun Lu** from Hong Kong Univ. of Science and Technology taught jointly with me a Data Mining Summer Course in Shanghai, China in July 1998. He has contributed many excellent slides to it
- n Some graduate students have contributed many new slides in the following years. Notable contributors include **Eugene Belchev**, **Jian Pei**, and **Osmar R. Zaiane** (now teaching in Univ. of Alberta).

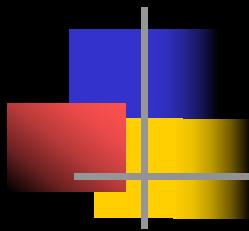
CMPT-459-00.3 Course Schedule

- n Chapter 1. Introduction {W1:L2, L3}
- n Chapter 2. Data warehousing and OLAP technology for data mining {W2:L1-3, W3:L1-2}
 - n Homework # 1 distribution (SQLServer7.0+ DBMiner2.0)
- n Chapter 3. Data preprocessing {W3:L3, W4: L1-L2}
- n Chapter 4. Data mining primitives, languages and system architectures {W4: L3, W5: L1}
 - n Homework #1 due, homework #2 distribution
- n Chapter 5. Concept description: Characterization and comparison {W5: L2, L3, W6: L2}
 - n W6:L1 Thanksgiving Day
- n Chapter 6. Mining association rules in large databases {W6: L3, W7: L1-3, W8: L2}
 - n Midterm {W8: L2}
- n Chapter 7. Classification and prediction {W8:L3, W9: L1-L3}
- n Chapter 8. Clustering analysis {W10: L1-L3}
 - n W10: L3 Homework #2 due
- n Chapter 9. Mining complex types of data {W11: L2-L3, W12:L1-L3}
 - n W11:L1 Remembrance Day, W12:L3 Course project due
- n Chapter 10. Data mining applications and trends in data mining {W13: L1-L3}



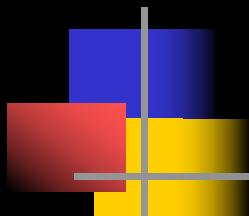
Where to Find the Set of Slides?

- n Tutorial sections (MS PowerPoint files):
 - n <http://www.cs.sfu.ca/~han/dmbook>
- n Other conference presentation slides (.ppt):
 - n <http://db.cs.sfu.ca/> or <http://www.cs.sfu.ca/~han>
- n Research papers, DBMiner system, and other related information:
 - n <http://db.cs.sfu.ca/> or <http://www.cs.sfu.ca/~han>



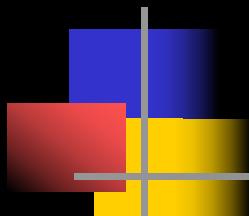
Chapter 1. Introduction

- n Motivation: Why data mining?
- n What is data mining?
- n Data Mining: On what kind of data?
- n Data mining functionality
- n Are all the patterns interesting?
- n Classification of data mining systems
- n Major issues in data mining



Motivation: “Necessity is the Mother of Invention”

- „ Data explosion problem
 - „ Automated data collection tools and mature database technology lead to tremendous amounts of data stored in databases, data warehouses and other information repositories
- „ We are drowning in data, but starving for knowledge!
- „ Solution: Data warehousing and data mining
 - „ Data warehousing and on-line analytical processing
 - „ Extraction of interesting knowledge (rules, regularities, patterns, constraints) from data in large databases



Evolution of Database Technology

(See Fig. 1.1)

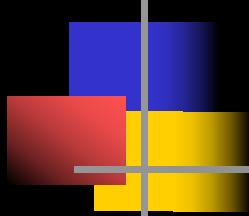
- n 1960s:
 - n Data collection, database creation, IMS and network DBMS
- n 1970s:
 - n Relational data model, relational DBMS implementation
- n 1980s:
 - n RDBMS, advanced data models (extended-relational, OO, deductive, etc.) and application-oriented DBMS (spatial, scientific, engineering, etc.)
- n 1990s–2000s:
 - n Data mining and data warehousing, multimedia databases, and Web databases



What Is Data Mining?

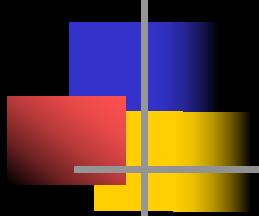
- n Data mining (knowledge discovery in databases):
 - n Extraction of interesting (non-trivial, implicit, previously unknown and potentially useful) information or patterns from data in large databases
- n Alternative names and their “inside stories”:
 - n Data mining: a misnomer?
 - n Knowledge discovery(mining) in databases (KDD), knowledge extraction, data/pattern analysis, data archeology, data dredging, information harvesting, business intelligence, etc.
- n What is not data mining?
 - n (Deductive) query processing.
 - n Expert systems or small ML/statistical programs





Why Data Mining? — Potential Applications

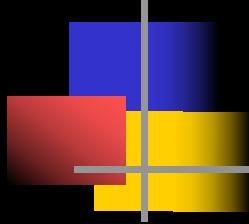
- n Database analysis and decision support
 - n Market analysis and management
 - n target marketing, customer relation management, market basket analysis, cross selling, market segmentation
 - n Risk analysis and management
 - n Forecasting, customer retention, improved underwriting, quality control, competitive analysis
 - n Fraud detection and management
- n Other Applications
 - n Text mining (news group, email, documents) and Web analysis.
 - n Intelligent query answering



Market Analysis and Management

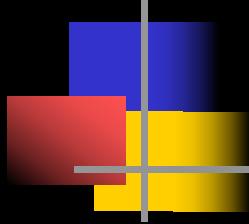
(1)

- n Where are the data sources for analysis?
 - n Credit card transactions, loyalty cards, discount coupons, customer complaint calls, plus (public) lifestyle studies
- n Target marketing
 - n Find clusters of “model” customers who share the same characteristics: interest, income level, spending habits, etc.
- n Determine customer purchasing patterns over time
 - n Conversion of single to a joint bank account: marriage, etc.
- n Cross-market analysis
 - n Associations/co-relations between product sales
 - n Prediction based on the association information



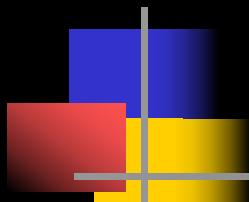
Market Analysis and Management (2)

- n Customer profiling
 - n data mining can tell you what types of customers buy what products (clustering or classification)
- n Identifying customer requirements
 - n identifying the best products for different customers
 - n use prediction to find what factors will attract new customers
- n Provides summary information
 - n various multidimensional summary reports
 - n statistical summary information (data central tendency and variation)



Corporate Analysis and Risk Management

- n Finance planning and asset evaluation
 - n cash flow analysis and prediction
 - n contingent claim analysis to evaluate assets
 - n cross-sectional and time series analysis (financial-ratio, trend analysis, etc.)
- n Resource planning:
 - n summarize and compare the resources and spending
- n Competition:
 - n monitor competitors and market directions
 - n group customers into classes and a class-based pricing procedure
 - n set pricing strategy in a highly competitive market



Fraud Detection and Management (1)

n Applications

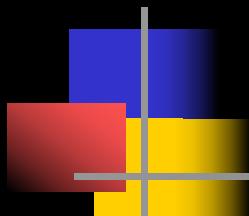
- n widely used in health care, retail, credit card services, telecommunications (phone card fraud), etc.

n Approach

- n use historical data to build models of fraudulent behavior and use data mining to help identify similar instances

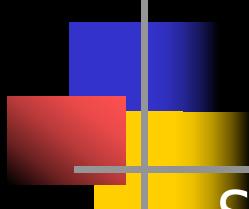
n Examples

- n auto insurance: detect a group of people who stage accidents to collect on insurance
- n money laundering: detect suspicious money transactions (US Treasury's Financial Crimes Enforcement Network)
- n medical insurance: detect professional patients and ring of doctors and ring of references



Fraud Detection and Management (2)

- n Detecting inappropriate medical treatment
 - n Australian Health Insurance Commission identifies that in many cases blanket screening tests were requested (save Australian \$1m/yr).
- n Detecting telephone fraud
 - n Telephone call model: destination of the call, duration, time of day or week. Analyze patterns that deviate from an expected norm.
 - n British Telecom identified discrete groups of callers with frequent intra-group calls, especially mobile phones, and broke a multimillion dollar fraud.
- n Retail
 - n Analysts estimate that 38% of retail shrink is due to dishonest employees.



Other Applications

n Sports

- IBM Advanced Scout analyzed NBA game statistics (shots blocked, assists, and fouls) to gain competitive advantage for New York Knicks and Miami Heat

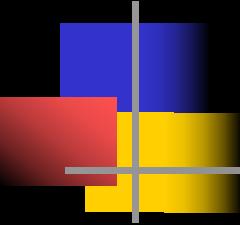
n Astronomy

- JPL and the Palomar Observatory discovered 22 quasars with the help of data mining

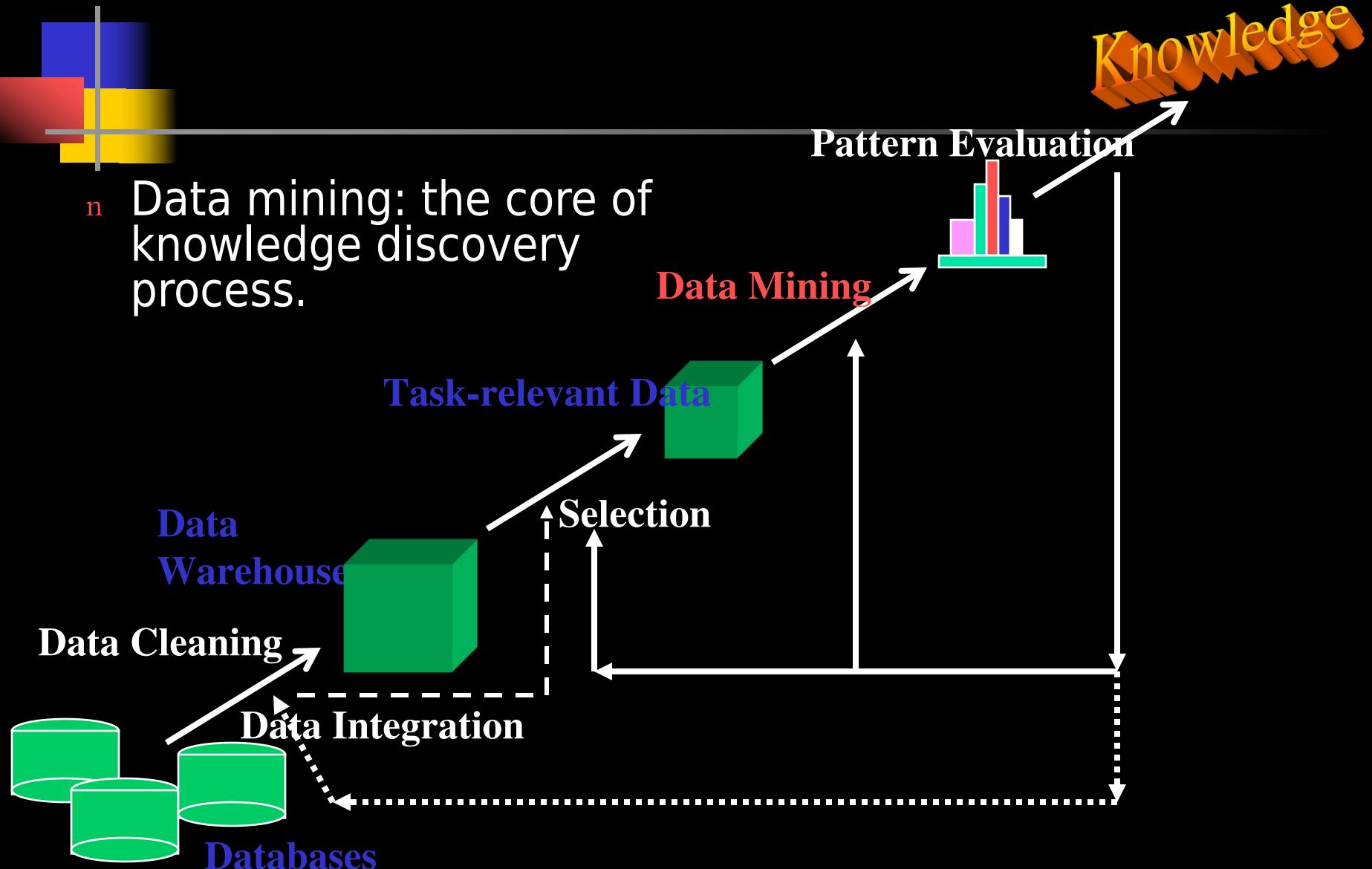
n Internet Web Surf-Aid

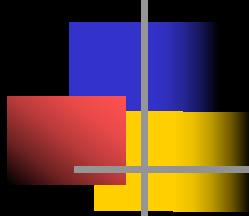
- IBM Surf-Aid applies data mining algorithms to Web access logs for market-related pages to discover customer preference and behavior pages, analyzing effectiveness of Web marketing, improving Web site organization, etc.

Data Mining: A KDD Process



- n Data mining: the core of knowledge discovery process.





Steps of a KDD Process

- n Learning the application domain:
 - n relevant prior knowledge and goals of application
- n Creating a target data set: data selection
- n **Data cleaning** and preprocessing: (may take 60% of effort!)
- n **Data reduction and transformation:**
 - n Find useful features, dimensionality/variable reduction, invariant representation.
- n Choosing functions of data mining
 - n summarization, classification, regression, association, clustering.
- n Choosing the mining algorithm(s)
- n **Data mining:** search for patterns of interest
- n **Pattern evaluation and knowledge presentation**
 - n visualization, transformation, removing redundant patterns, etc.
- n Use of discovered knowledge

Data Mining and Business Intelligence

Increasing potential
to support
business decisions

Making Decisions

Data Presentation

Visualization Techniques

Data Mining
Information Discovery

Data Exploration

Statistical Analysis, Querying and Reporting

Data Warehouses / Data Marts

OLAP, MDA

Data Sources

Paper, Files, Information Providers, Database Systems, OLTP

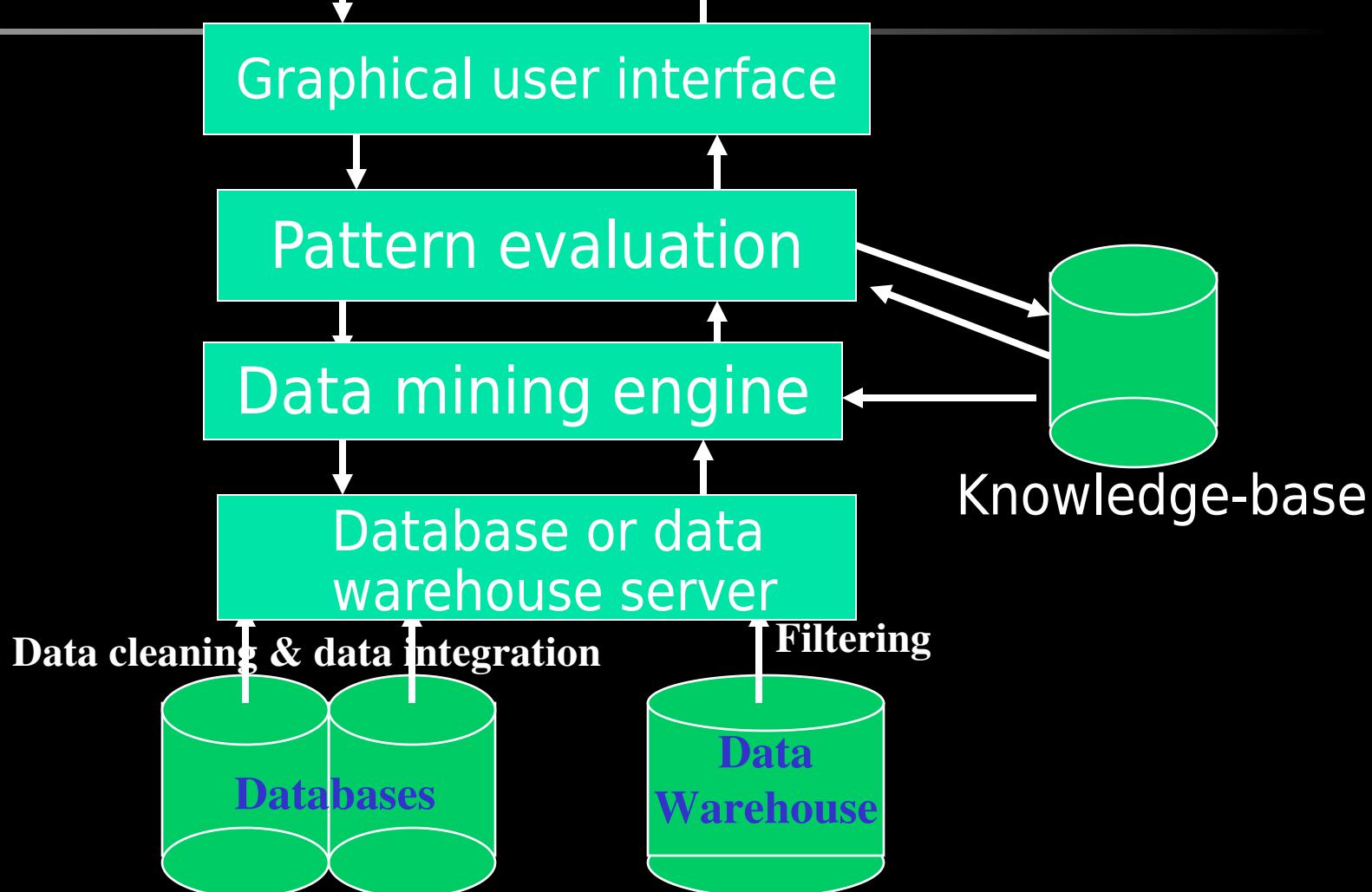
End User

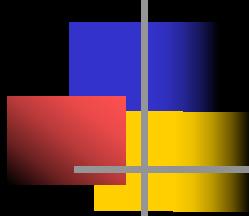
Business Analyst

Data Analyst

DBA

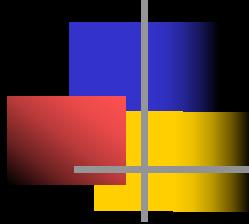
Architecture of a Typical Data Mining System





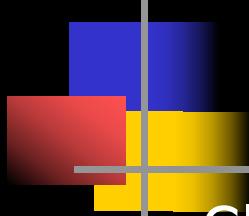
Data Mining: On What Kind of Data?

- n Relational databases
- n Data warehouses
- n Transactional databases
- n Advanced DB and information repositories
 - n Object-oriented and object-relational databases
 - n Spatial databases
 - n Time-series data and temporal data
 - n Text databases and multimedia databases
 - n Heterogeneous and legacy databases
 - n WWW



Data Mining Functionalities (1)

- n Concept description: Characterization and discrimination
 - n Generalize, summarize, and contrast data characteristics, e.g., dry vs. wet regions
- n Association (correlation and causality)
 - n Multi-dimensional vs. single-dimensional association
 - n $\text{age}(X, "20..29") \wedge \text{income}(X, "20..29K") \rightarrow \text{buys}(X, "PC")$ [support = 2%, confidence = 60%]
 - n $\text{contains}(T, "computer") \rightarrow \text{contains}(x, "software")$ [1%, 75%]



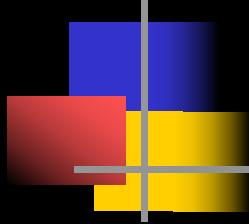
Data Mining Functionalities (2)

n Classification and Prediction

- n Finding models (functions) that describe and distinguish classes or concepts for future prediction
- n E.g., classify countries based on climate, or classify cars based on gas mileage
- n Presentation: decision-tree, classification rule, neural network
- n Prediction: Predict some unknown or missing numerical values

n Cluster analysis

- n Class label is unknown: Group data to form new classes, e.g., cluster houses to find distribution patterns
- n Clustering based on the principle: maximizing the intra-class similarity and minimizing the interclass similarity



Data Mining Functionalities (3)

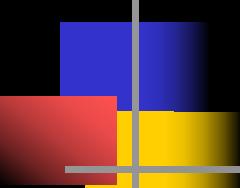
n Outlier analysis

- n Outlier: a data object that does not comply with the general behavior of the data
- n It can be considered as noise or exception but is quite useful in fraud detection, rare events analysis

n Trend and evolution analysis

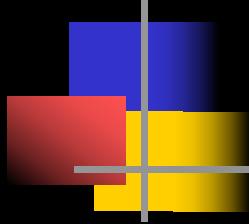
- n Trend and deviation: regression analysis
- n Sequential pattern mining, periodicity analysis
- n Similarity-based analysis

n Other pattern-directed or statistical analyses



Are All the “Discovered” Patterns Interesting?

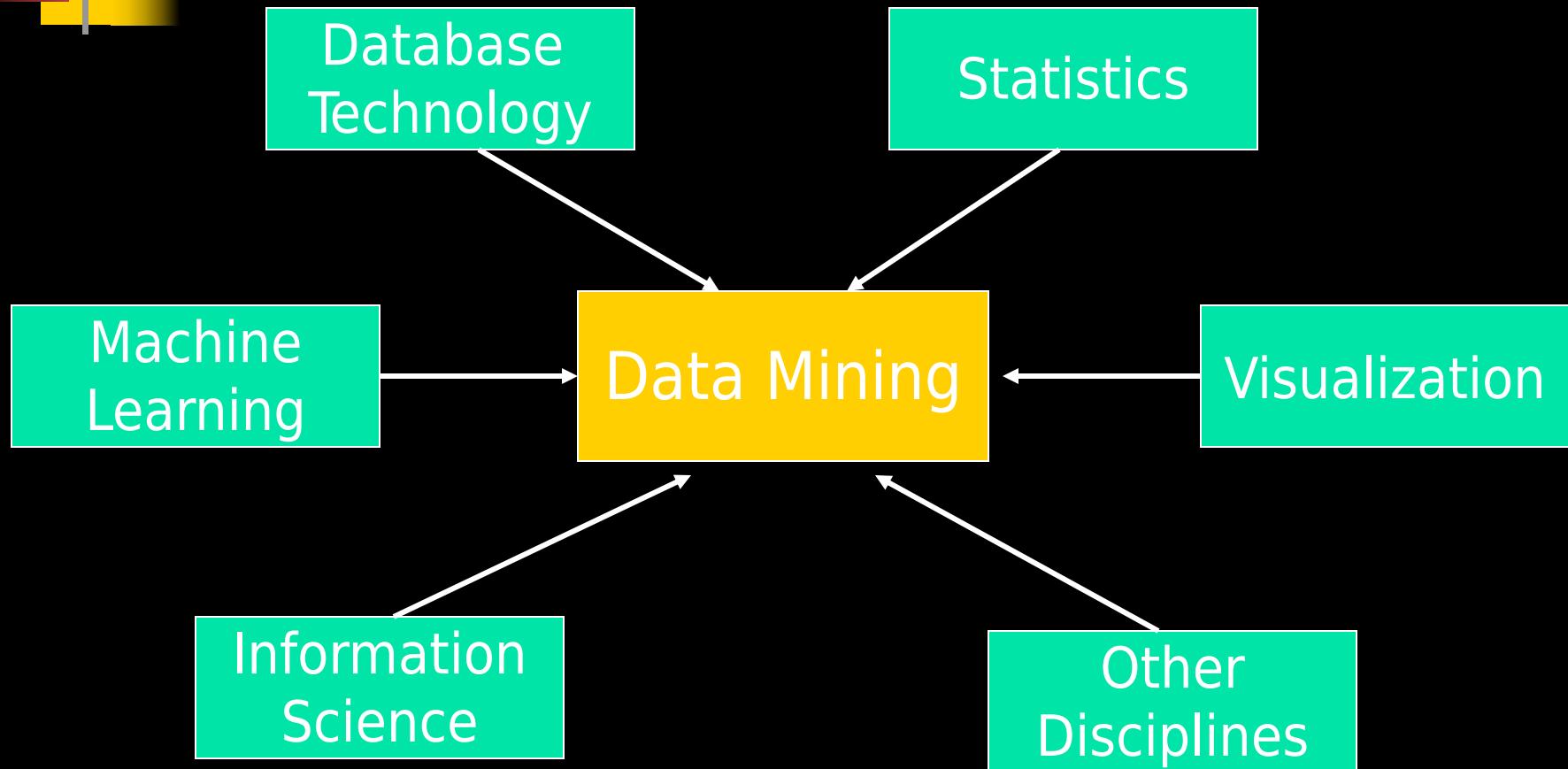
- n A data mining system/query may generate thousands of patterns, not all of them are interesting.
 - n Suggested approach: Human-centered, query-based, focused mining
 - n **Interestingness measures**: A pattern is **interesting** if it is easily understood by humans, valid on new or test data with some degree of certainty, potentially useful, novel, or validates some hypothesis that a user seeks to confirm
 - n **Objective vs. subjective interestingness measures:**
 - n Objective: based on statistics and structures of patterns, e.g., support, confidence, etc.
 - n Subjective: based on user’s belief in the data, e.g., unexpectedness, novelty, actionability, etc.

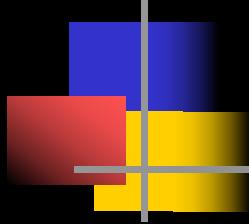


Can We Find All and Only Interesting Patterns?

- n Find all the interesting patterns: Completeness
 - n Can a data mining system find all the interesting patterns?
 - n Association vs. classification vs. clustering
- n Search for only interesting patterns: Optimization
 - n Can a data mining system find only the interesting patterns?
 - n Approaches
 - n First generate all the patterns and then filter out the uninteresting ones.
 - n Generate only the interesting patterns—mining query optimization

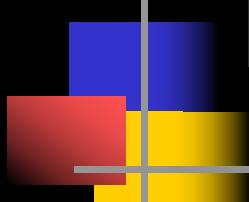
Data Mining: Confluence of Multiple Disciplines





Data Mining: Classification Schemes

- n General functionality
 - n Descriptive data mining
 - n Predictive data mining
- n Different views, different classifications
 - n Kinds of databases to be mined
 - n Kinds of knowledge to be discovered
 - n Kinds of techniques utilized
 - n Kinds of applications adapted



A Multi-Dimensional View of Data Mining Classification

n **Databases to be mined**

- n Relational, transactional, object-oriented, object-relational, active, spatial, time-series, text, multi-media, heterogeneous, legacy, WWW, etc.

n **Knowledge to be mined**

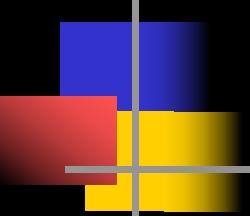
- n Characterization, discrimination, association, classification, clustering, trend, deviation and outlier analysis, etc.
- n Multiple/integrated functions and mining at multiple levels

n **Techniques utilized**

- n Database-oriented, data warehouse (OLAP), machine learning, statistics, visualization, neural network, etc.

n **Applications adapted**

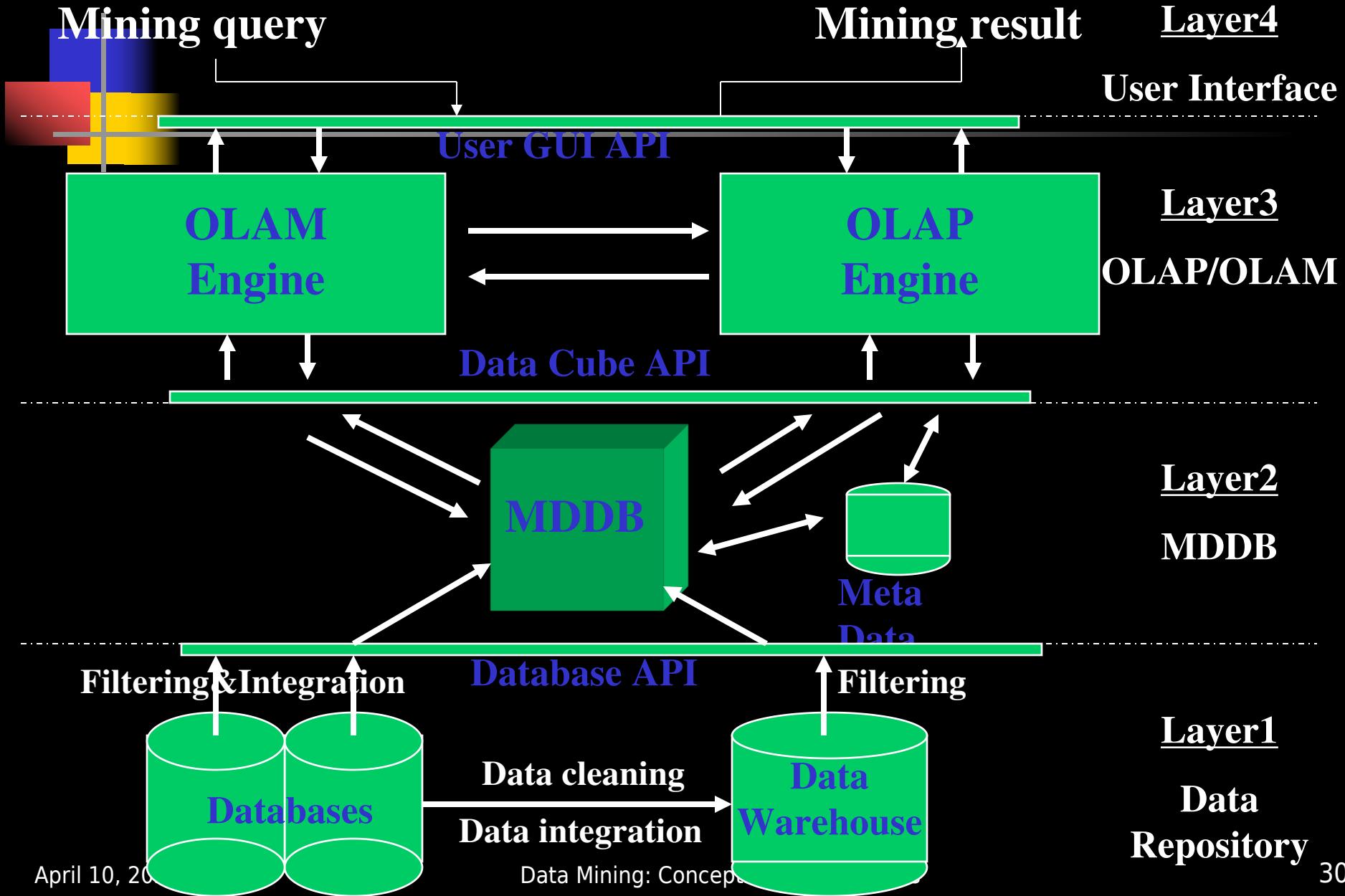
- n Retail, telecommunication, banking, fraud analysis, DNA mining, stock market analysis, Web mining, Weblog analysis, etc.

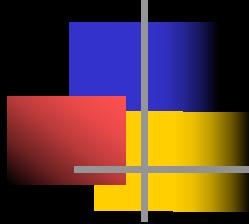


OLAP Mining: An Integration of Data Mining and Data Warehousing

- n **Data mining systems, DBMS, Data warehouse systems coupling**
 - n No coupling, loose-coupling, semi-tight-coupling, tight-coupling
- n **On-line analytical mining data**
 - n integration of mining and OLAP technologies
- n **Interactive mining multi-level knowledge**
 - n Necessity of mining knowledge and patterns at different levels of abstraction by drilling/rolling, pivoting, slicing/dicing, etc.
- n **Integration of multiple mining functions**
 - n Characterized classification, first clustering and then

An OLAM Architecture





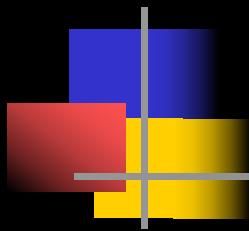
Major Issues in Data Mining (1)

n Mining methodology and user interaction

- n Mining different kinds of knowledge in databases
- n Interactive mining of knowledge at multiple levels of abstraction
- n Incorporation of background knowledge
- n Data mining query languages and ad-hoc data mining
- n Expression and visualization of data mining results
- n Handling noise and incomplete data
- n Pattern evaluation: the interestingness problem

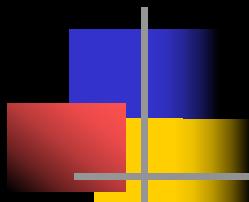
n Performance and scalability

- n Efficiency and scalability of data mining algorithms
- n Parallel, distributed and incremental mining methods



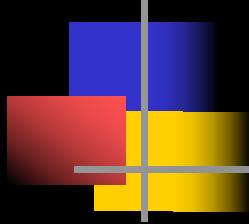
Major Issues in Data Mining (2)

- n Issues relating to the diversity of data types
 - n Handling relational and complex types of data
 - n Mining information from heterogeneous databases and global information systems (WWW)
- n Issues related to applications and social impacts
 - n Application of discovered knowledge
 - n Domain-specific data mining tools
 - n Intelligent query answering
 - n Process control and decision making
 - n Integration of the discovered knowledge with existing knowledge: A knowledge fusion problem
 - n Protection of data security, integrity, and privacy



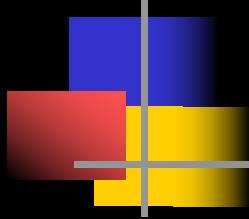
Summary

- n Data mining: discovering interesting patterns from large amounts of data
- n A natural evolution of database technology, in great demand, with wide applications
- n A KDD process includes data cleaning, data integration, data selection, transformation, data mining, pattern evaluation, and knowledge presentation
- n Mining can be performed in a variety of information repositories
- n Data mining functionalities: characterization, discrimination, association, classification, clustering, outlier and trend analysis, etc.
- n Classification of data mining systems
- n Major issues in data mining



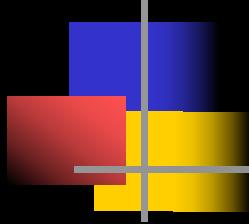
A Brief History of Data Mining Society

- 1989 IJCAI Workshop on Knowledge Discovery in Databases (Piatetsky-Shapiro)
 - 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)
- 1998 ACM SIGKDD, SIGKDD'1999-2001 conferences, and SIGKDD Explorations
- More conferences on data mining
 - PAKDD, PKDD, SIAM-Data Mining, (IEEE) ICDM, etc.



Where to Find References?

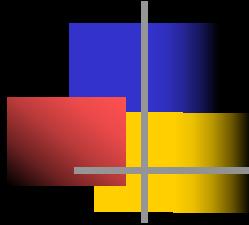
- n Data mining and KDD (SIGKDD member CDROM):
 - n Conference proceedings: KDD, and others, such as PKDD, PAKDD, etc.
 - n Journal: Data Mining and Knowledge Discovery
- n Database field (SIGMOD member CD ROM):
 - n Conference proceedings: ACM-SIGMOD, ACM-PODS, VLDB, ICDE, EDBT, DASFAA
 - n Journals: ACM-TODS, J. ACM, IEEE-TKDE, JIIS, etc.
- n AI and Machine Learning:
 - n Conference proceedings: Machine learning, AAAI, IJCAI, etc.
 - n Journals: Machine Learning, Artificial Intelligence, etc.
- n Statistics:
 - n Conference proceedings: Joint Stat. Meeting, etc.
 - n Journals: Annals of statistics, etc.
- n Visualization:
 - n Conference proceedings: CHI, etc.
 - n Journals: IEEE Trans. visualization and computer graphics, etc.



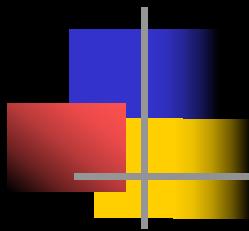
References

- U. M. Fayyad, G. Piatetsky-Shapiro, P. Smyth, and R. Uthurusamy. Advances in Knowledge Discovery and Data Mining. AAAI/MIT Press, 1996.
- J. Han and M. Kamber. Data Mining: Concepts and Techniques. Morgan Kaufmann, 2000.
- T. Imielinski and H. Mannila. A database perspective on knowledge discovery. Communications of ACM, 39:58-64, 1996.
- G. Piatetsky-Shapiro, U. Fayyad, and P. Smith. From data mining to knowledge discovery: An overview. In U.M. Fayyad, et al. (eds.), Advances in Knowledge Discovery and Data Mining, 1-35. AAAI/MIT Press, 1996.
- G. Piatetsky-Shapiro and W. J. Frawley. Knowledge Discovery in Databases. AAAI/MIT Press, 1991.

<http://www.cs.sfu.ca/~han>



Thank you !!!



CMPT-843 Course Arrangement

- n 1st week: full instructor teaching
- n 2nd to 11th week: 1/2 graduate student + 1/2 instructor teaching
- n 12-13th week: full student graduate project presentation
- n Course evaluation:
 - n presentation (quality of presentation slides 7% + presentation 8%) 15%
 - n midterm exam 35%
 - n project (presentation 5% + report 25%) total 30%
 - n homework (2): 20%
- n Deadline for the selection of your work in the semester:
 - n selection of course presentation: at the end of the 1st week
 - n selection of the course project: at the end of the 3rd week
 - n project proposal due date: at the end of the 4th week
 - n homework due dates:
 - n project due date: end of the semester
 - n Your presentation slides due date: one day before the presentation
 - n midterm date: end of the 8th week