

Introduction: Necessity Is the Mother of Invention

- Data explosion problem
 - Automated data collection tools, widely used database systems, computerized society, and the Internet lead to tremendous amounts of data accumulated and/or to be analyzed in databases, data warehouses, WWW, 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
 - Mining interesting knowledge (rules, regularities, patterns, constraints) from data in large databases

1

What Is Data Mining?



- 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"?
 - (Deductive) query processing.
 - Expert systems or small ML/statistical programs



2

Why Data Mining?—Potential Applications

- Data analysis and decision support
 - Market analysis and management
 - Target marketing, customer relationship management (CRM), market basket analysis, cross selling, market segmentation
 - Risk analysis and management
 - Forecasting, customer retention, improved underwriting, quality control, competitive analysis
 - Fraud detection and detection of unusual patterns (outliers)
- Other Applications
 - Text mining (news group, email, documents) and Web mining
 - Stream data mining
 - Bioinformatics and bio-data analysis

3

Example 1: Market Analysis and Management

- Where does the data come from?—Credit card transactions, loyalty cards, discount coupons, customer complaint calls, plus (public) lifestyle studies
- Target marketing
 - Find clusters of "model" customers who share the same characteristics: interest, income level, spending habits, etc.,
 - Determine customer purchasing patterns over time
- Cross-market analysis—Find associations/co-relations between product sales, & predict based on such association
- Customer profiling—What types of customers buy what products (clustering or classification)
- Customer requirement analysis
 - Identify the best products for different customers
 - Predict what factors will attract new customers
- Provision of summary information
 - Multidimensional summary reports
 - Statistical summary information (data central tendency and variation)

4

Example 2: Corporate Analysis & Risk Management

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

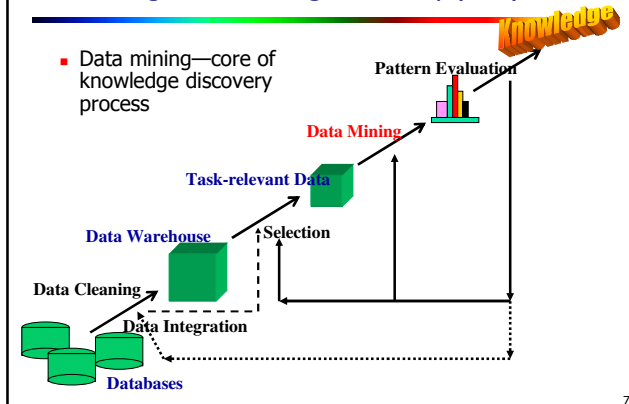
5

Example 3: Fraud Detection & Mining Unusual Patterns

- Approaches: Clustering & model construction for frauds, outlier analysis
- Applications: Health care, retail, credit card service, telecomm.
 - Auto insurance: ring of collisions
 - Money laundering: suspicious monetary transactions
 - Medical insurance
 - Professional patients, ring of doctors, and ring of references
 - Unnecessary or correlated screening tests
 - Telecommunications: phone-call fraud
 - Phone call model: destination of the call, duration, time of day or week. Analyze patterns that deviate from an expected norm
 - Retail industry
 - Analysts estimate that 38% of retail shrink is due to dishonest employees
 - Anti-terrorism

6

Data Mining and Knowledge Discovery (KDD) Process



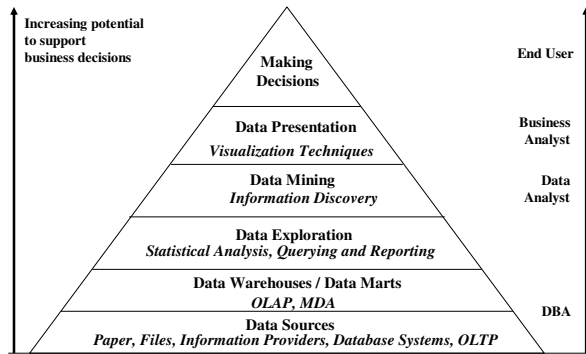
7

Steps of a KDD Process

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

8

Data Mining and Business Intelligence



9

Data Mining: On What Kinds of Data?

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

10

Data Mining Functionalities

- Multidimensional concept description: Characterization and discrimination
 - Generalize, summarize, and contrast data characteristics, e.g., dry vs. wet regions
- Frequent patterns, association, correlation and 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

11

Data Mining Functionalities (2)

- 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? No! useful in fraud detection, rare events analysis
- Trend and evolution analysis
 - Trend and deviation: e.g., regression analysis
 - Sequential pattern mining, periodicity analysis
 - Similarity-based analysis
- Other pattern-directed or statistical analyses

12

Are All the "Discovered" Patterns Interesting?

- Data mining may generate thousands of patterns: Not all of them are interesting
 - Suggested approach: Human-centered, query-based, focused mining
- 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
- Objective vs. subjective interestingness measures**
 - Objective**: based on **statistics and structures of patterns**, e.g., support, confidence, etc.
 - Subjective**: based on **user's belief** in the data, e.g., unexpectedness, novelty, actionability, etc.

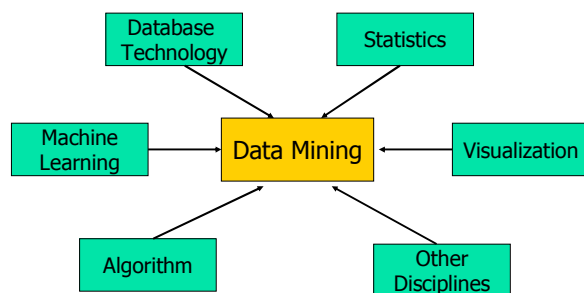
13

Can We Find All and Only Interesting Patterns?

- Find all the interesting patterns: Completeness**
 - Can a data mining system find **all** the interesting patterns?
 - Heuristic vs. exhaustive search
 - Association vs. classification vs. clustering
- Search for only interesting patterns: An optimization problem**
 - Can a data mining system find **only** the interesting patterns?
 - Approaches
 - First general all the patterns and then filter out the uninteresting ones.
 - Generate only the interesting patterns—mining query optimization

14

Data Mining: Confluence of Multiple Disciplines



15

Data Mining: Classification Schemes

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

16

Multi-Dimensional View of Data Mining

- **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.

17

Primitives that Define a Data Mining Task

- Task-relevant data
- Type of knowledge to be mined
- Background knowledge
- Pattern interestingness measurements
- Visualization/presentation of discovered patterns

18

Primitive 1: Task-Relevant Data

- Database or data warehouse name
- Database tables or data warehouse cubes
- Condition for data selection
- Relevant attributes or dimensions
- Data grouping criteria

19

Primitive 2: Types of Knowledge to Be Mined

- Characterization
- Discrimination
- Association
- Classification/prediction
- Clustering
- Outlier analysis
- Other data mining tasks

20

Primitive 3: Background Knowledge (Concept Hierarchies)

- Schema hierarchy
 - E.g., street < city < province_or_state < country
- Set-grouping hierarchy
 - E.g., {20-39} = young, {40-59} = middle_aged
- Operation-derived hierarchy
 - email address: hagonzal@cs.uiuc.edu
login-name < department < university < country
- Rule-based hierarchy
 - $\text{low_profit_margin}(X) \leq \text{price}(X, P_1) \text{ and } \text{cost}(X, P_2) \text{ and } (P_1 - P_2) < \50

21

Primitive 4: Measurements of Pattern Interestingness

- Simplicity
 - e.g., (association) rule length, (decision) tree size
- Certainty
 - e.g., confidence, $P(A|B) = \#(A \text{ and } B) / \#(B)$, classification reliability or accuracy, certainty factor, rule strength, rule quality, discriminating weight, etc.
- Utility
 - potential usefulness, e.g., support (association), noise threshold (description)
- Novelty
 - not previously known, surprising (used to remove redundant rules, e.g., Illinois vs. Champaign rule implication support ratio)

22

Primitive 5: Presentation of Discovered Patterns

- Different backgrounds/usages may require **different forms of representation**
 - E.g., rules, tables, crosstabs, pie/bar chart, etc.
- **Concept hierarchy** is also important
 - Discovered knowledge might be more understandable when represented at **high level of abstraction**
 - Interactive **drill up/down, pivoting, slicing and dicing** provide different perspectives to data
- Different kinds of **knowledge** require different representation: association, classification, clustering, etc.

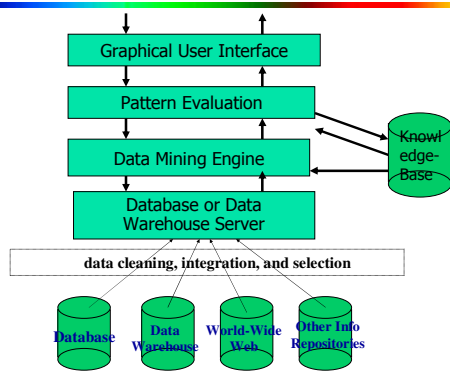
23

Integration of Data Mining and Data Warehousing

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

24

Architecture: Typical Data Mining System



25

Coupling Data Mining with DB/DW Systems

- No coupling—flat file processing, not recommended
- Loose coupling
 - Fetching data from DB/DW
- Semi-tight coupling—enhanced DM performance
 - Provide efficient implement a few data mining primitives in a DB/DW system, e.g., sorting, indexing, aggregation, histogram analysis, multiway join, precomputation of some stat functions
- Tight coupling—A uniform information processing environment
 - DM is smoothly integrated into a DB/DW system, mining query is optimized based on mining query, indexing, query processing methods, etc.

26

Major Issues in Data Mining

- Mining methodology
 - Mining different kinds of knowledge from diverse data types, e.g., bio, stream, Web
 - Performance: efficiency, effectiveness, and scalability
 - Pattern evaluation: the interestingness problem
 - Incorporation of background knowledge
 - Handling noise and incomplete data
 - Parallel, distributed and incremental mining methods
 - Integration of the discovered knowledge with existing one: knowledge fusion
- User interaction
 - Data mining query languages and ad-hoc mining
 - Expression and visualization of data mining results
 - Interactive mining of knowledge at multiple levels of abstraction
- Applications and social impacts
 - Domain-specific data mining & invisible data mining
 - Protection of data security, integrity, and privacy

27

Summary

- 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

28