

Unit II – Data Mining for Business Intelligence

Prepared by
Dr. Lakshmi Priya GG

Overview

- Motivation for Data Mining
- Data Mining-
 - Functionalities
 - Classification of Data Mining systems,
 - Data Mining task primitives,
 - Data Streaming,
 - Integration of a Data Mining system with a Database or a Data Warehouse,
 - Major issues in Data Mining.
- Data pre- processing
- BI lifecycle- implementation of BI

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

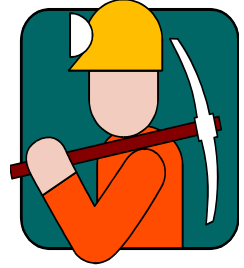
Evolution of Sciences

- 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

- 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?

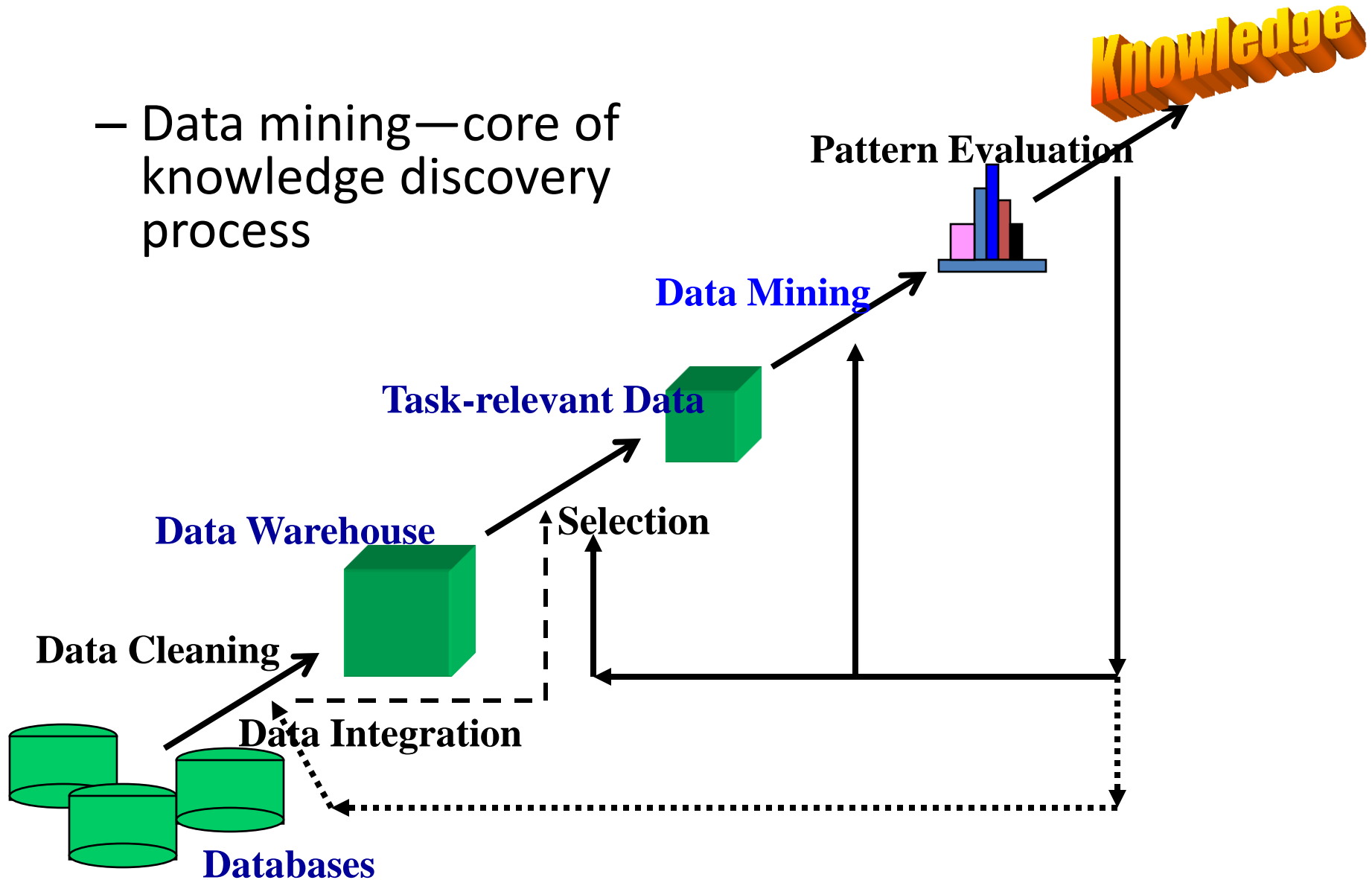


- Data mining (knowledge discovery from data)
 - Extraction of interesting (non-trivial, implicit, previously unknown and potentially useful) patterns or knowledge from huge amount of data
 - Data mining: a misnomer?
- Alternative names
 - Knowledge discovery (mining) in databases (KDD), knowledge extraction,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

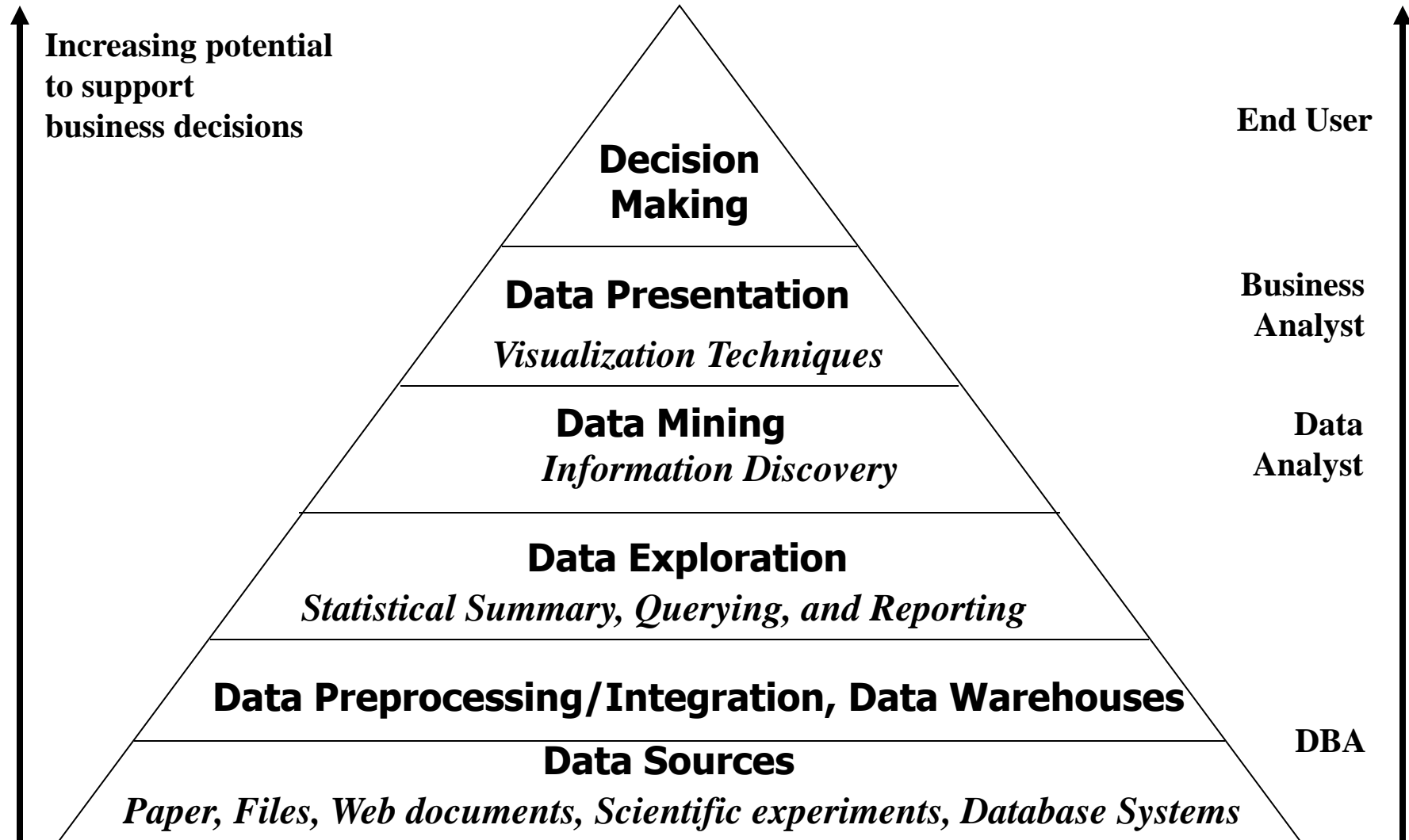


Knowledge Discovery (KDD) Process

- Data mining—core of knowledge discovery process



Data Mining and Business Intelligence



What is (not) Data Mining?

● What is not Data Mining?

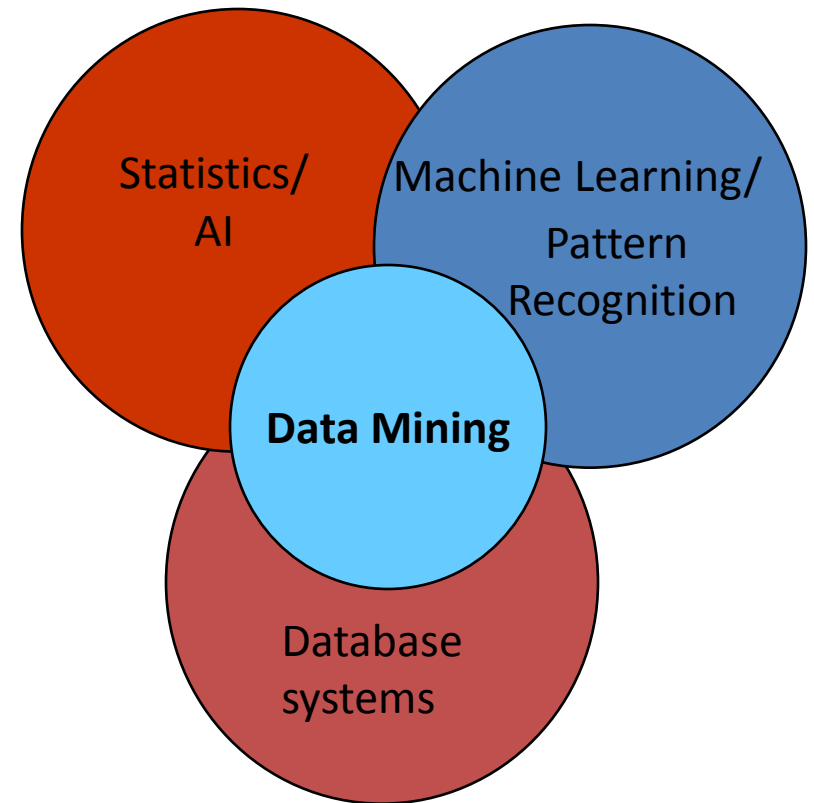
- Look up phone number in phone directory
- Query a Web search engine for information about “Amazon”

● What is Data Mining?

- Certain names are more prevalent in certain US locations (O’Brien, O’Rourke, O’Reilly... in Boston area)
- Group together similar documents returned by search engine according to their context (e.g. Amazon rainforest, Amazon.com,)

Origins of Data Mining

- **Draws ideas from machine learning/AI, pattern recognition, statistics, and database systems**
- **Traditional Techniques may be unsuitable due to**
 - Enormity of data
 - High dimensionality of data
 - Heterogeneous, distributed nature of data



Are All the “Discovered” Patterns Interesting?

- A data mining system/query 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.

Can We Find All and Only Interesting Patterns?

- Find all the interesting patterns: Completeness
 - Can a data mining system find all the interesting patterns?
 - Association vs. classification vs. clustering
- Search for only interesting patterns: Optimization
 - Can a data mining system find only the interesting patterns?
 - Highly desirable, progress has been made, but still a challenge
 - Approaches
 - First generate all the patterns and then filter out the uninteresting ones.
 - Generate only the interesting patterns—mining query optimization

Data Mining: Classification Schemes

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

Data Mining Tasks

- Prediction Methods
 - Use some variables to predict unknown or future values of other variables.
- Description Methods
 - Find human-interpretable patterns that describe the data.

A Multi-Dimensional View of Data Mining Classification

- **Databases to be mined**
 - Relational, transactional, object-oriented, object-relational, active, spatial, time-series, text, multi-media, heterogeneous, legacy, WWW, etc.
- **Knowledge to be mined**
 - Characterization, discrimination, association, classification, clustering, trend, deviation and outlier analysis, etc.
 - Multiple/integrated functions and mining at multiple levels
- **Techniques utilized**
 - Database-oriented, data warehouse (OLAP), machine learning, statistics, visualization, neural network, etc.
- **Applications adapted**
 - Retail, telecommunication, banking, fraud analysis, DNA mining, stock market analysis, Web mining, Weblog analysis, etc.

Data Mining Tasks...

- Classification [Predictive]
- Clustering [Descriptive]
- Association Rule Discovery [Descriptive]
- Sequential Pattern Discovery [Descriptive]
- Regression [Predictive]
- Deviation Detection [Predictive]

Classification: Definition

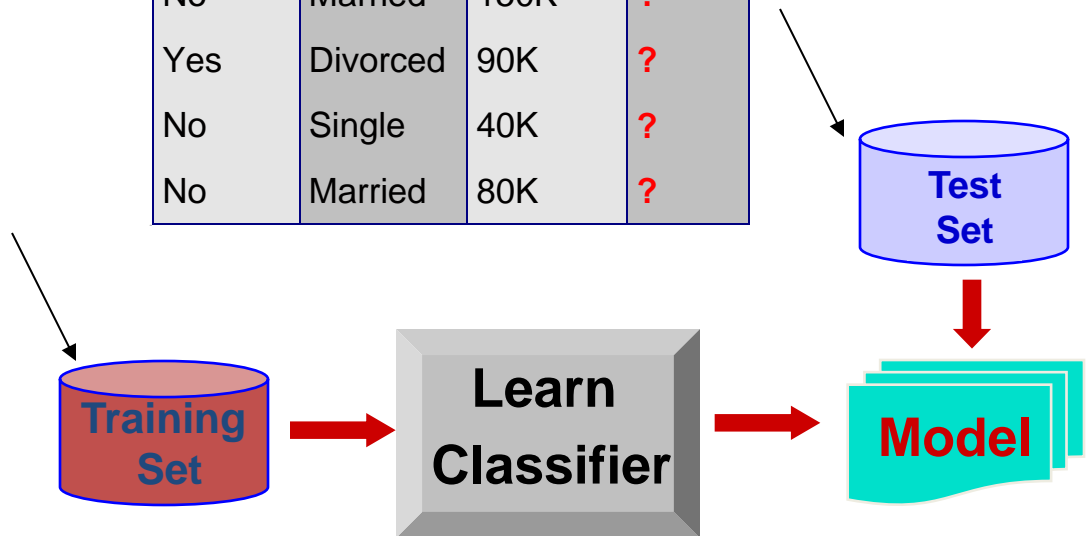
- Given a collection of records (*training set*)
 - Each record contains a set of *attributes*, one of the attributes is the *class*.
- Find a *model* for class attribute as a function of the values of other attributes.
- Goal: previously unseen records should be assigned a class as accurately as possible.
 - A *test set* is used to determine the accuracy of the model. Usually, the given data set is divided into training and test sets, with training set used to build the model and test set used to validate it.

Classification Example

categorical
categorical
continuous
class

<i>Tid</i>	Refund	Marital Status	Taxable Income	Cheat
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes

Refund	Marital Status	Taxable Income	Cheat
No	Single	75K	?
Yes	Married	50K	?
No	Married	150K	?
Yes	Divorced	90K	?
No	Single	40K	?
No	Married	80K	?



Classification: Application 1

- Direct Marketing
 - Goal: Reduce cost of mailing by *targeting* a set of consumers likely to buy a new cell-phone product.
 - Approach:
 - Use the data for a similar product introduced before.
 - We know which customers decided to buy and which decided otherwise. This *{buy, don't buy}* decision forms the *class attribute*.
 - Collect various demographic, lifestyle, and company-interaction related information about all such customers.
 - Type of business, where they stay, how much they earn, etc.
 - Use this information as input attributes to learn a classifier model.

Classification: Application 2

- Fraud Detection
 - Goal: Predict fraudulent cases in credit card transactions.
 - Approach:
 - Use credit card transactions and the information on its account-holder as attributes.
 - When does a customer buy, what does he buy, how often he pays on time, etc
 - Label past transactions as fraud or fair transactions. This forms the class attribute.
 - Learn a model for the class of the transactions.
 - Use this model to detect fraud by observing credit card transactions on an account.

Classification: Application 3

- Customer Attrition/Churn:
 - Goal: To predict whether a customer is likely to be lost to a competitor.
 - Approach:
 - Use detailed record of transactions with each of the past and present customers, to find attributes.
 - How often the customer calls, where he calls, what time-of-the day he calls most, his financial status, marital status, etc.
 - Label the customers as loyal or disloyal.
 - Find a model for loyalty.

Classification: Application 4

- Sky Survey Cataloging
 - Goal: To predict class (star or galaxy) of sky objects, especially visually faint ones, based on the telescopic survey images (from Palomar Observatory).
 - 3000 images with 23,040 x 23,040 pixels per image.
 - Approach:
 - Segment the image.
 - Measure image attributes (features) - 40 of them per object.
 - Model the class based on these features.
 - Success Story: Could find 16 new high red-shift quasars, some of the farthest objects that are difficult to find!

Clustering Definition

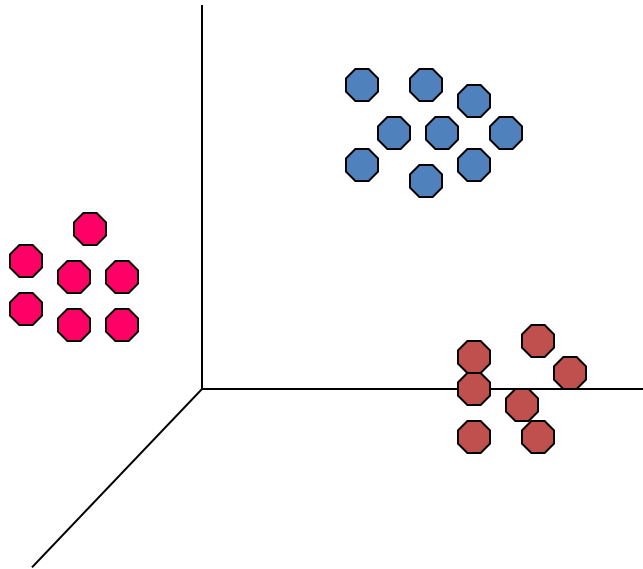
- Given a set of data points, each having a set of attributes, and a similarity measure among them, find clusters such that
 - Data points in one cluster are more similar to one another.
 - Data points in separate clusters are less similar to one another.
- Similarity Measures:
 - Euclidean Distance if attributes are continuous.
 - Other Problem-specific Measures.

Illustrating Clustering

✗ Euclidean Distance Based Clustering in 3-D space.

Intracuster distances
are minimized

Intercluster distances
are maximized



Clustering: Application 1

- Market Segmentation:
 - Goal: subdivide a market into distinct subsets of customers where any subset may conceivably be selected as a market target to be reached with a distinct marketing mix.
 - Approach:
 - Collect different attributes of customers based on their geographical and lifestyle related information.
 - Find clusters of similar customers.
 - Measure the clustering quality by observing buying patterns of customers in same cluster vs. those from different clusters.

Clustering: Application 2

- Document Clustering:
 - Goal: To find groups of documents that are similar to each other based on the important terms appearing in them.
 - Approach: To identify frequently occurring terms in each document. Form a similarity measure based on the frequencies of different terms. Use it to cluster.
 - Gain: Information Retrieval can utilize the clusters to relate a new document or search term to clustered documents.

Illustrating Document Clustering

- Clustering Points: 3204 Articles of Los Angeles Times.
- Similarity Measure: How many words are common in these documents (after some word filtering).

<i>Category</i>	<i>Total Articles</i>	<i>Correctly Placed</i>
<i>Financial</i>	555	364
<i>Foreign</i>	341	260
<i>National</i>	273	36
<i>Metro</i>	943	746
<i>Sports</i>	738	573
<i>Entertainment</i>	354	278

Clustering of S&P 500 Stock Data

- ⌘ Observe Stock Movements every day.
- ⌘ Clustering points: Stock-{UP/DOWN}
- ⌘ Similarity Measure: Two points are more similar if the events described by them frequently happen together on the same day.
 - ⌘ We used association rules to quantify a similarity measure.

	<i>Discovered Clusters</i>	<i>Industry Group</i>
1	Applied-Matl-DOWN,Bay-Net work-Down,3-COM-DOWN, Cabletron-Sys-DOWN,CISCO-DOWN,HP-DOWN, DSC-Comm-DOWN,INTEL-DOWN,LSI-Logic-DOWN, Micron-Tech-DOWN,Texas-Inst-Down,Tellabs-Inc-Down, Natl-Semiconduct-DOWN,Orac1-DOWN,SGI-DOWN, Sun-DOWN	Technology1-DOWN
2	Apple-Comp-DOWN,Autodesk-DOWN,DEC-DOWN, ADV-Micro-Device-DOWN,Andrew-Corp-DOWN, Computer-Assoc-DOWN,Circuit-City-DOWN, Compaq-DOWN, EMC-Corp-DOWN, Gen-Inst-DOWN, Motorola-DOWN,Microsoft-DOWN,Scientific-Atl-DOWN	Technology2-DOWN
3	Fannie-Mac-DOWN,Fed-Home-Loan-DOWN, MBNA-Corp-DOWN,Morgan-Stanley-DOWN	Financial-DOWN
4	Baker-Hughes-UP,Dresser-Inds-UP,Halliburton-HLD-UP, Louisiana-Land-UP,Phillips-Petro-UP,Unocal-UP, Schlumberger-UP	Oil-UP

Association Rule Discovery: Definition

- Given a set of records each of which contain some number of items from a given collection;
 - Produce dependency rules which will predict occurrence of an item based on occurrences of other items.

<i>TID</i>	<i>Items</i>
1	Bread, Coke, Milk
2	Cola, Bread
3	Cola, Coke, Diaper, Milk
4	Cola, Bread, Diaper, Milk
5	Coke, Diaper, Milk

Rules Discovered:

{Milk} --> {Coke}

{Diaper, Milk} --> {Cola}

Association Rule Discovery: Application 1

- Marketing and Sales Promotion:
 - Let the rule discovered be
{Bagels, ... } --> {Potato Chips}
 - Potato Chips as consequent => Can be used to determine what should be done to boost its sales.
 - Bagels in the antecedent => Can be used to see which products would be affected if the store discontinues selling bagels.
 - Bagels in antecedent *and* Potato chips in consequent => Can be used to see what products should be sold with Bagels to promote sale of Potato chips!

Association Rule Discovery: Application 2

- Supermarket shelf management.
 - Goal: To identify items that are bought together by sufficiently many customers.
 - Approach: Process the point-of-sale data collected with barcode scanners to find dependencies among items.
 - A classic rule --
 - If a customer buys diaper and milk, then he is very likely to buy Cola.
 - So, don't be surprised if you find six-packs stacked next to diapers!

Association Rule Discovery: Application 3

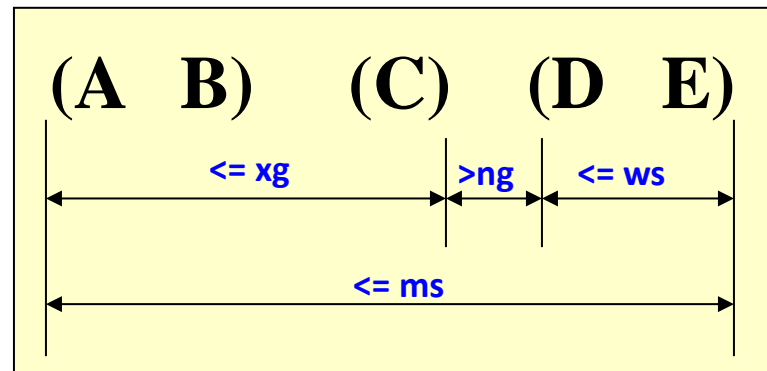
- Inventory Management:
 - Goal: A consumer appliance repair company wants to anticipate the nature of repairs on its consumer products and keep the service vehicles equipped with right parts to reduce on number of visits to consumer households.
 - Approach: Process the data on tools and parts required in previous repairs at different consumer locations and discover the co-occurrence patterns.

Sequential Pattern Discovery: Definition

- Given is a set of *objects*, with each object associated with its own *timeline of events*, find rules that predict strong **sequential dependencies** among different events.

(A B) (C) \longrightarrow (D E)

- Rules are formed by first discovering patterns. Event occurrences in the patterns are governed by timing constraints.



Sequential Pattern Discovery: Examples

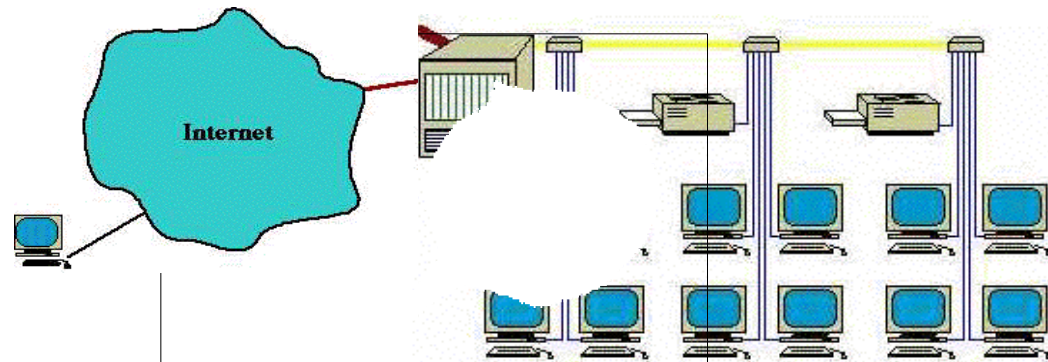
- In telecommunications alarm logs,
 - (Inverter_Problem Excessive_Line_Current)
(Rectifier_Alarm) --> (Fire_Alarm)
- In point-of-sale transaction sequences,
 - Computer Bookstore:
(Intro_To_Visual_C) (C++_Primer) -->
(Perl_for_dummies,Tcl_Tk)
 - Athletic Apparel Store:
(Shoes) (Racket, Racketball) --> (Sports_Jacket)

Regression

- Predict a value of a given continuous valued variable based on the values of other variables, assuming a linear or nonlinear model of dependency.
- Greatly studied in statistics, neural network fields.
- Examples:
 - Predicting sales amounts of new product based on advertising expenditure.
 - Predicting wind velocities as a function of temperature, humidity, air pressure, etc.
 - Time series prediction of stock market indices.

Deviation/Anomaly Detection

- Detect significant deviations from normal behavior
- Applications:
 - Credit Card Fraud Detection
 - Network Intrusion Detection



Typical network traffic at University level may reach over 100 million connections per day

DBMS, OLAP, and Data Mining

	DBMS	OLAP	Data Mining
Task	Extraction of detailed and summary data	Summaries, trends and forecasts	Knowledge discovery of hidden patterns and insights
Type of result	Information	Analysis	Insight and Prediction
Method	Deduction (Ask the question, verify with data)	Multidimensional data modeling, Aggregation, Statistics	Induction (Build the model, apply it to new data, get the result)
Example question	Who purchased mutual funds in the last 3 years?	What is the average income of mutual fund buyers by region by year?	Who will buy a mutual fund in the next 6 months and why?

Example of DBMS, OLAP and Data Mining: Weather Data

DBMS:

Day	outlook	temperature	humidity	windy	play
1	sunny	85	85	false	no
2	sunny	80	90	true	no
3	overcast	83	86	false	yes
4	rainy	70	96	false	yes
5	rainy	68	80	false	yes
6	rainy	65	70	true	no
7	overcast	64	65	true	yes
8	sunny	72	95	false	no
9	sunny	69	70	false	yes
10	rainy	75	80	false	yes
11	sunny	75	70	true	yes
12	overcast	72	90	true	yes
13	overcast	81	75	false	yes
14	rainy	71	91	true	no

Example of DBMS, OLAP and Data Mining: Weather Data

- By querying a DBMS containing the above table we may answer questions like:
- What was the temperature in the sunny days? {85, 80, 72, 69, 75}
- Which days the humidity was less than 75? {6, 7, 9, 11}
- Which days the temperature was greater than 70? {1, 2, 3, 8, 10, 11, 12, 13, 14}
- Which days the temperature was greater than 70 and the humidity was less than 75? The intersection of the above two: {11}

Example of DBMS, OLAP and Data Mining: Weather Data

OLAP:

- Using OLAP we can create a **Multidimensional Model** of our data (**Data Cube**).
- For example using the dimensions: **time**, **outlook** and **play** we can create the following model.

9 / 5	sunny	rainy	overcast
Week 1	0 / 2	2 / 1	2 / 0
Week 2	2 / 1	1 / 1	2 / 0

Example of DBMS, OLAP and Data Mining: Weather Data

Data Mining:

- Using the ID3 algorithm we can produce the following decision tree:
 - **outlook = sunny**
 - humidity = high: no
 - humidity = normal: yes
 - **outlook = overcast: yes**
 - **outlook = rainy**
 - windy = true: no
 - windy = false: yes

Major Issues in Data Warehousing and Mining

- Mining methodology and user interaction
 - Mining different kinds of knowledge in databases
 - Interactive mining of knowledge at multiple levels of abstraction
 - Incorporation of background knowledge
 - Data mining query languages and ad-hoc data mining
 - Expression and visualization of data mining results
 - Handling noise and incomplete data
 - Pattern evaluation: the interestingness problem
- Performance and scalability
 - Efficiency and scalability of data mining algorithms
 - Parallel, distributed and incremental mining methods

Major Issues in Data Warehousing and Mining

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

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.
- Classification of data mining systems
- Major issues in data mining

Mining Data Streams

- What is stream data?
- Why Stream Data Systems?
- Stream data management systems: Issues and solutions
- Research issues

Characteristics of Data Streams

- Data Streams
 - Data streams—continuous, ordered, changing, fast, huge amount
 - Traditional DBMS—data stored in finite, persistent data sets
- Characteristics
 - Huge volumes of continuous data, possibly infinite
 - Fast changing and requires fast, real-time response
 - Data stream captures nicely our data processing needs of today
 - Random access is expensive—single scan algorithm (*can only have one look*)
 - Store only the summary of the data seen thus far
 - Most stream data are at pretty low-level or multi-dimensional in nature, needs multi-level and multi-dimensional processing

Stream Data Applications

- Telecommunication calling records
- Business: credit card transaction flows
- Network monitoring and traffic engineering
- Financial market: stock exchange
- Engineering & industrial processes: power supply & manufacturing
- Sensor, monitoring & surveillance: video streams, RFIDs
- Security monitoring
- Web logs and Web page click streams
- Massive data sets (even saved but random access is too expensive)

DBMS versus DSMS

- Persistent relations
- One-time queries
- Random access
- “Unbounded” disk store
- Only current state matters
- No real-time services
- Relatively low update rate
- Data at any granularity
- Assume precise data
- Access plan determined by query processor, physical DB design
- Transient streams
- Continuous queries
- Sequential access
- Bounded main memory
- Historical data is important
- Real-time requirements
- Possibly multi-GB arrival rate
- Data at fine granularity
- Data stale/imprecise
- Unpredictable/variable data arrival and characteristics

Ack. From Motwani's PODS tutorial slides

Architecture: Stream Query Processing

SDMS (Stream Data Management System)

Continuous Query

Multiple streams

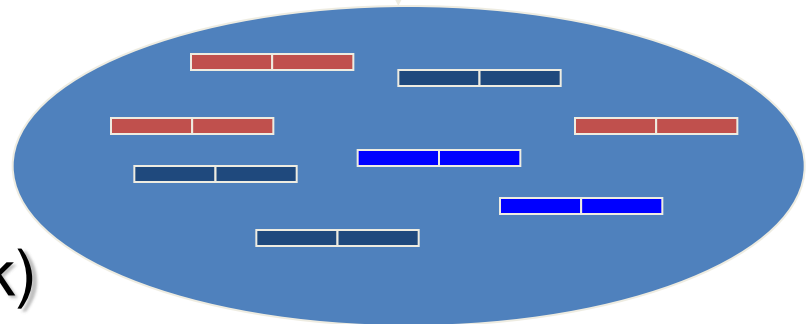


User/Application

Stream Query Processor

Results

Scratch Space
(Main memory and/or Disk)



Challenges of Stream Data Processing

- Multiple, continuous, rapid, time-varying, ordered streams
- Main memory computations
- Queries are often continuous
 - Evaluated continuously as stream data arrives
 - Answer updated over time
- Queries are often complex
 - Beyond element-at-a-time processing
 - Beyond stream-at-a-time processing
 - Beyond relational queries (scientific, data mining, OLAP)
- Multi-level/multi-dimensional processing and data mining
 - Most stream data are at low-level or multi-dimensional in nature

Processing Stream Queries

- Query types
 - One-time query vs. **continuous query** (being evaluated continuously as stream continues to arrive)
 - **Predefined query** vs. ad-hoc query (issued on-line)
- Unbounded memory requirements
 - For real-time response, **main memory algorithm** should be used
 - Memory requirement is unbounded if one will join future tuples
- Approximate query answering
 - With bounded memory, it is not always possible to produce exact answers
 - **High-quality approximate answers** are desired
 - Data reduction and synopsis construction methods
 - Sketches, random sampling, histograms, wavelets, etc.

Methodologies for Stream Data Processing

- Major challenges
 - Keep track of a large universe, e.g., pairs of IP address, not ages
- Methodology
 - Synopses (trade-off between accuracy and storage)
 - Use *synopsis data structure*, much smaller ($O(\log^k N)$ space) than their base data set ($O(N)$ space)
 - Compute an *approximate answer* within a *small error range* (factor ϵ of the actual answer)
- Major methods
 - Random sampling
 - Histograms
 - Sliding windows
 - Multi-resolution model
 - Sketches
 - Radomized algorithms

Stream Data Processing Methods (1)

- Random sampling (but without knowing the total length in advance)
 - *Reservoir sampling*: maintain a set of s candidates in the reservoir, which form a true random sample of the element seen so far in the stream. As the data stream flows, every new element has a certain probability (s/N) of replacing an old element in the reservoir.
- Sliding windows
 - Make decisions based only on *recent data* of sliding window size w
 - An element arriving at time t expires at time $t + w$
- Histograms
 - Approximate the frequency distribution of element values in a stream
 - Partition data into a set of contiguous buckets
 - Equal-width (equal value range for buckets) vs. V-optimal (minimizing frequency variance within each bucket)
- Multi-resolution models
 - Popular models: balanced binary trees, micro-clusters, and wavelets

Stream Data Processing Methods (2)

- Sketches

- Histograms and wavelets require multi-passes over the data but sketches can operate in a single pass

$$F_k = \sum_{i=1}^v m_i^k$$

- Frequency moments of a stream $A = \{a_1, \dots, a_N\}$, F_k :

where v : the universe or domain size, m_i : the frequency of i in the sequence

- Given N elts and v values, sketches can approximate F_0, F_1, F_2 in $O(\log v + \log N)$ space

- Randomized algorithms

- Monte Carlo algorithm: bound on running time but may not return correct result

- Chebyshev's inequality: $P(|X - \mu| > k) \leq \frac{\sigma^2}{k^2}$

- Let X be a random variable with mean μ and standard deviation σ

- Chernoff bound: $P[X < (1 - \delta)\mu] < e^{-\mu\delta^2/4}$

- Let X be the sum of independent Poisson trials X_1, \dots, X_n , δ in $(0, 1]$

- The probability decreases exponentially as we move from the mean

Stream Data Mining: Research Issues

- Mining sequential patterns in data streams
- Mining partial periodicity in data streams
- Mining notable gradients in data streams
- Mining outliers and unusual patterns in data streams
- Stream clustering
 - Multi-dimensional clustering analysis?
 - Cluster not confined to 2-D metric space, how to incorporate other features, especially non-numerical properties
 - Stream clustering with other clustering approaches?
 - Constraint-based cluster analysis with data streams?

Summary: Stream Data Mining

- Stream data mining: A rich and on-going research field
 - Current research focus in database community:
 - DSMS system architecture, continuous query processing, supporting mechanisms
 - Stream data mining and stream OLAP analysis
 - Powerful tools for finding general and unusual patterns
 - Effectiveness, efficiency and scalability: lots of open problems
- Our philosophy on stream data analysis and mining
 - A multi-dimensional stream analysis framework
 - Time is a special dimension: Tilted time frame
 - What to compute and what to save?—Critical layers
 - partial materialization and precomputation
 - Mining dynamics of stream data