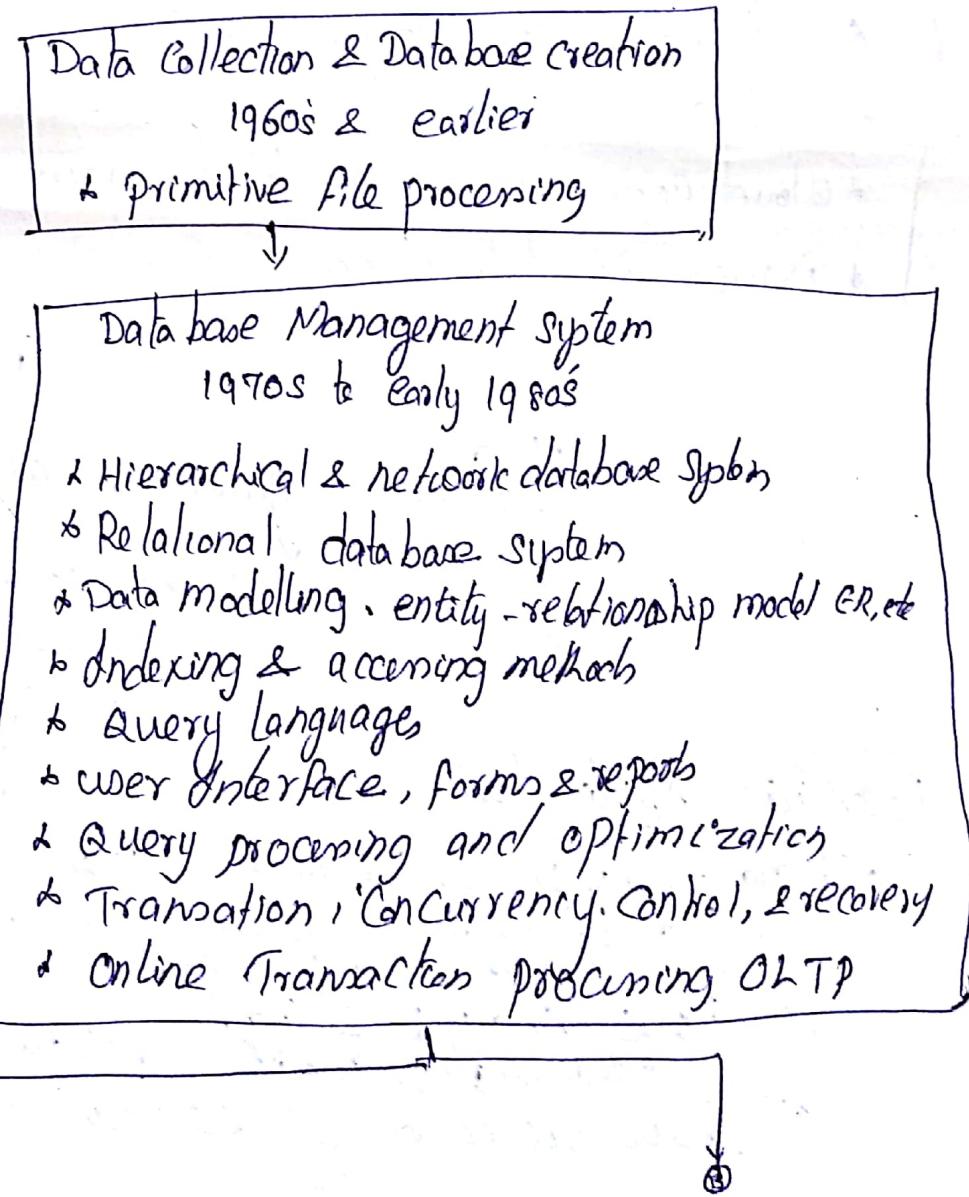
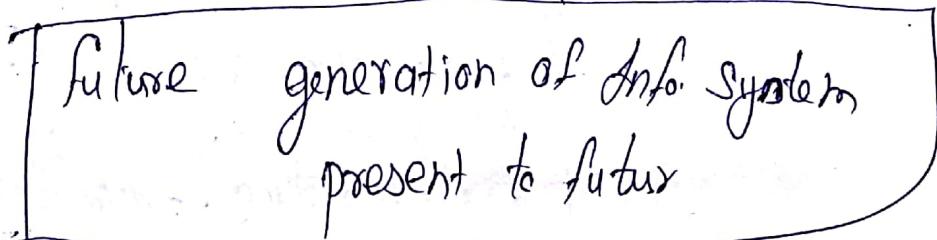
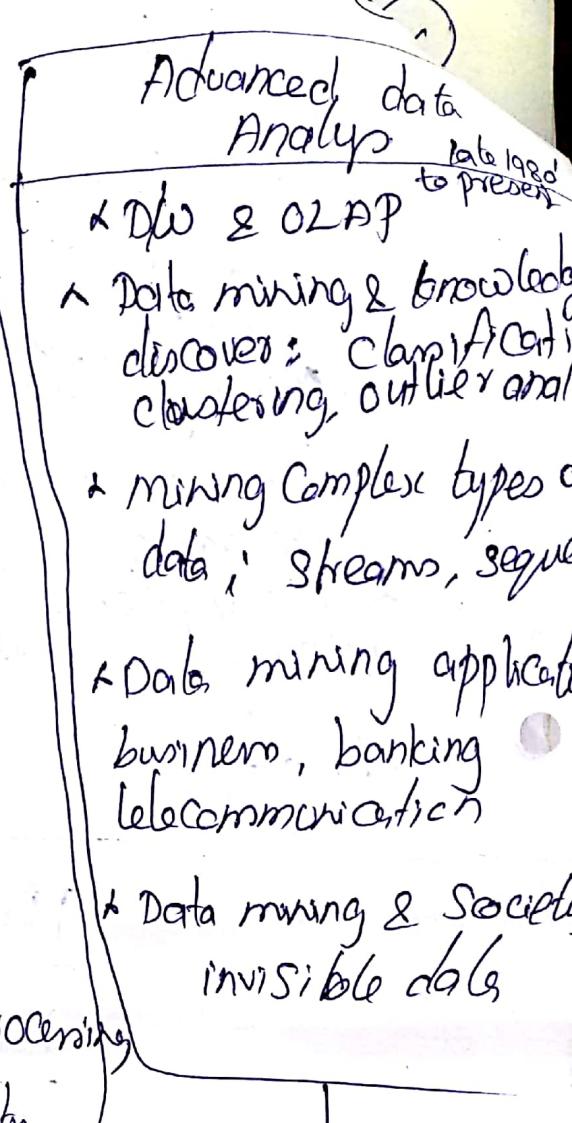
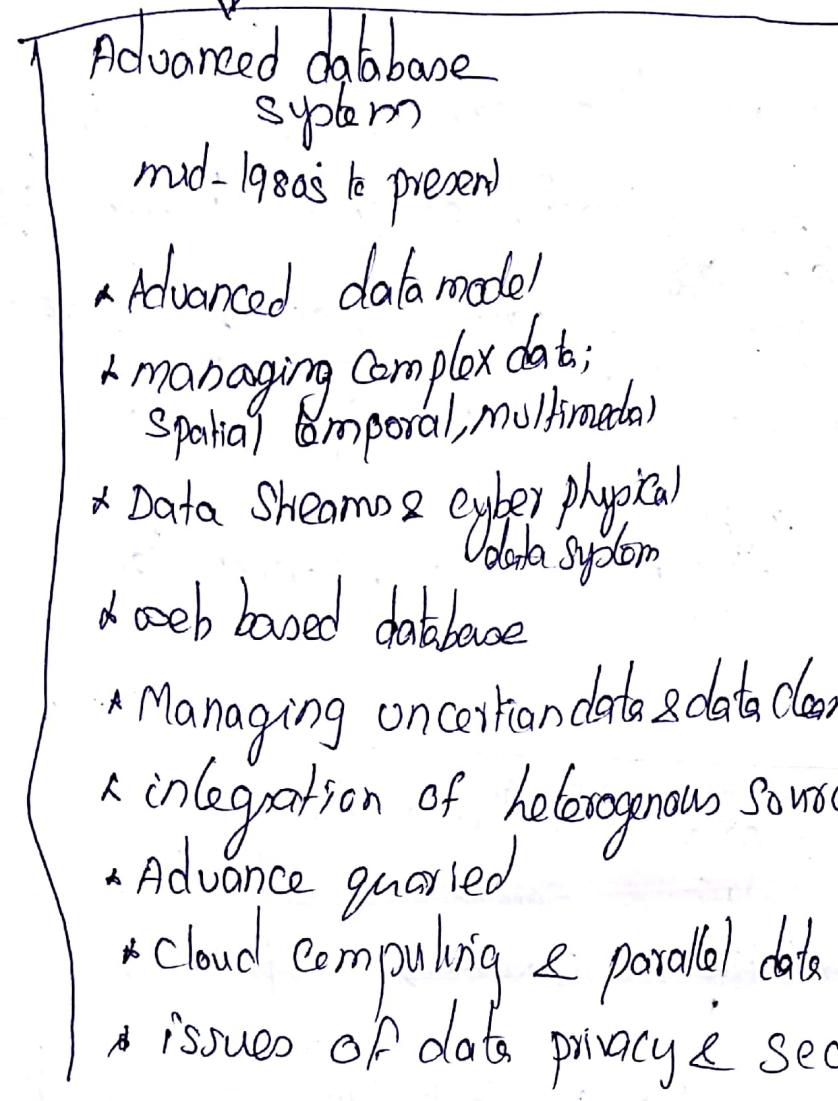


Data mining

- * It is also called as knowledge discovery from data or KDD
 - * Extraction of interesting patterns or knowledge from huge amount of data
 - * Mining knowledge from data, knowledge extraction
 - * data/pattern analysis, data archaeology & data dredging
 - * Mining of knowledge or data.
- (ii) Data mining as the Evolution of Information Technology.

Technology.





Applications

Potential Application

- * Data Analytics & Decision Support
- (i) Market Analysis & management
Target marketing, customer relationship management (CRM)
- (ii) Risk Analysis & Mgt
* forecasting, quality control, etc
- (iii) Fraud detection & detection of unusual patterns (outliers)

Other Applications

- * Text mining
- * co-emerging
- * Stream data mining
- * Bioinformatics

Market Analysis & mgt

* Data come from

- * Credit Card transactions, discount Coupons, Customer Complaint calls etc

* Target Marketing

- * find cluster of 'Mode', customer who share the characteristics such as interest, income level etc

- * Determine customer purchasing patterns over time

* Gross Market analysis

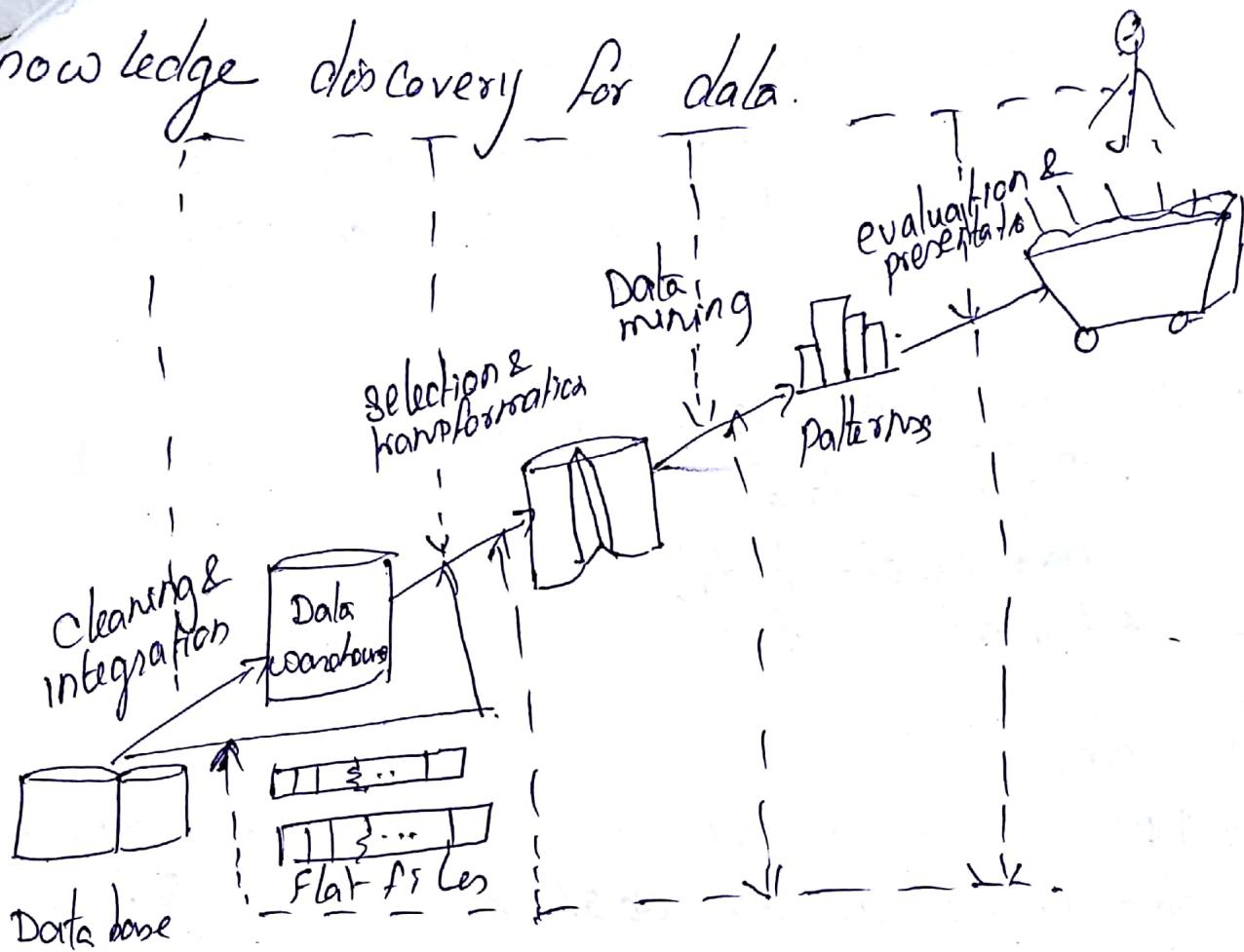
- * Association / Correlations b/w product sales & prediction based on such associations

Fraud detection & mining unusual patterns

* Application:

- Healthcare, retail, Credit Card service, telecomm.

knowledge discovery for data



Data Cleaning : Remove noise & inconsistent data

Data integration : Multiple data source may be combined

Data Selection : data relevant to the analysis task are retrieved from the database

Data transformation : data are transformed & consolidated into forms by performing summary or aggregation operations

Data mining : Intelligent methods are applied to extract data patterns

Pattern evaluation : to identify the truly interesting patterns representing knowledge based on interestingness measures

knowledge presentation; where visualization & knowledge representation techniques are used to present mined knowledge to users

what kind of Data Can be mined

- * Database data
- * data warehouse data
- * transactional data

Database data

- * DBMS consists of a collection of interrelated data known as database, and a set of s/w pgms to manage and access the data.
- * s/w pgm provide mechanisms for database structures and data storage for specifying & managing concurrent, shared and distributed data access.
- * A Relational database is a collection of tables each of which is assigned a unique name.
- * Each table consists of a set of attributes (columns or fields) and stores a large set of tuples (records).
- * Each tuple in a relational table represents an object identified by a unique key and described by a set of attribute values.

A Semantic data model such as Entity relationship (ER) data model is constructed for relational data base

* ER model represents the database as a set of entities & their relationships.

eg. A relational database for All Electronics.

It consists of following relational tables

customer, item, employee & branch.

The relational table Customer consists of set of attributes describing the customer information including customer cust-id, name, address, age, occupation, annual-income, etc.

item - item-id, brand, category, type, price, place-made, supplier, cost.

employee - emp-id, name, category, group, salary, ...

branch - branch-id, name, address, ..

The tables can be used to represent the relationships between or among multiple entities: eg,

purchases - customer purchase items, creating a sales transaction handled by an employee.

item sold - list of item sold in a given transaction.

work-at - employee work at a branch of AllElectronics.

(6)

Relational schema for a relational database.

customer (cust-ID, name, address, age, ...)

item (item-ID, brand, Category, type, price...)

employee (empl-ID, name, Category, group, salary...)

branch (branch-ID, name, address, ...)

purchases (trans-ID, cust-ID, empl-ID, date, time, method.paid, amount)

item.sold (trans-ID, item-ID, qty)

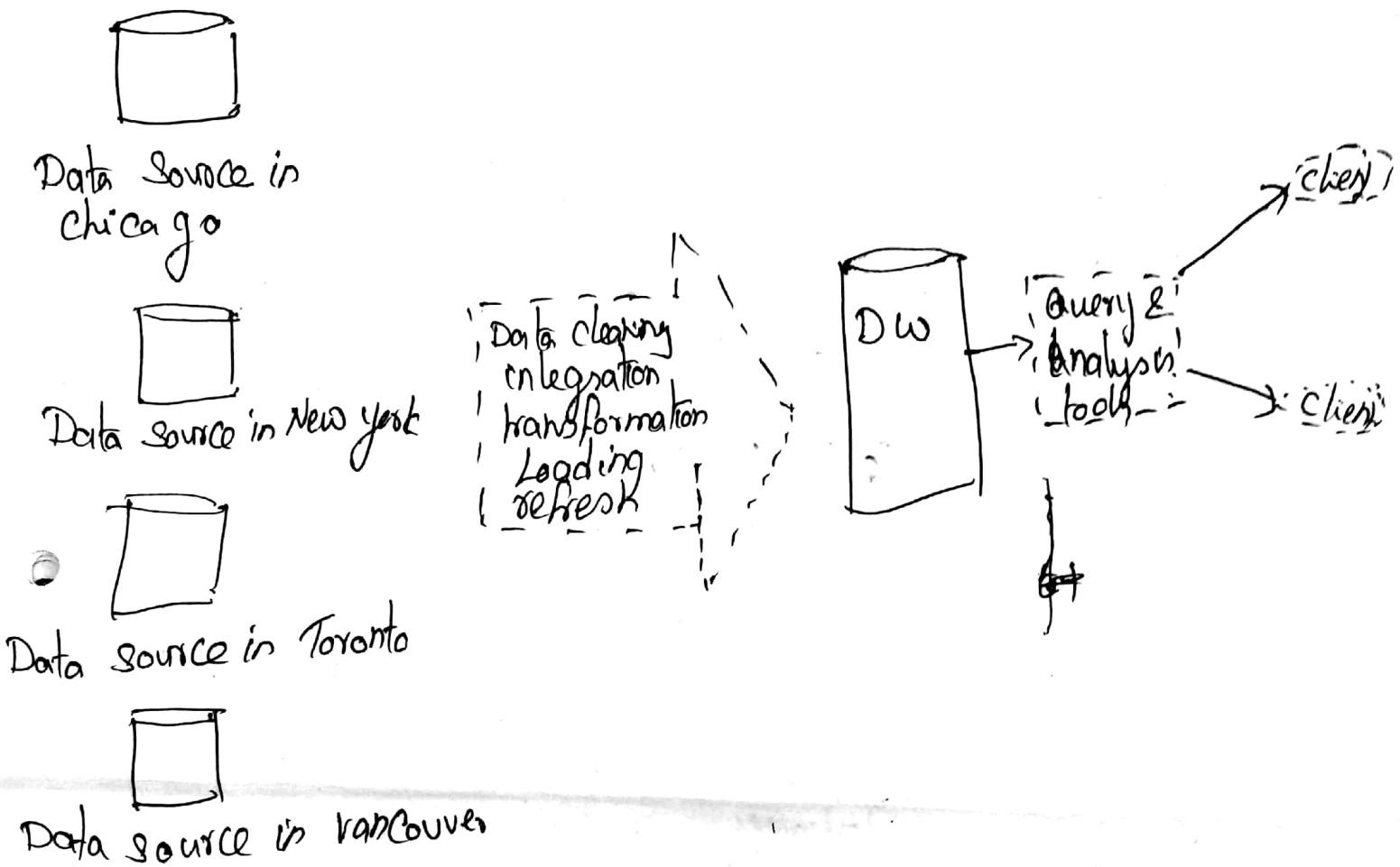
works-at (empl-ID, branch-ID)

Relational data can be accessed by using database query written in ~~SQL~~ relational query language, eg selection, join

Data warehouses;

- A repository of information collected from multiple sources, stored under a unified schema and residing at a single site.
- Data warehouses are constructed via a process of data cleaning, data integration, data transformation, data loading & data refreshing.

typical framework for all Electronics



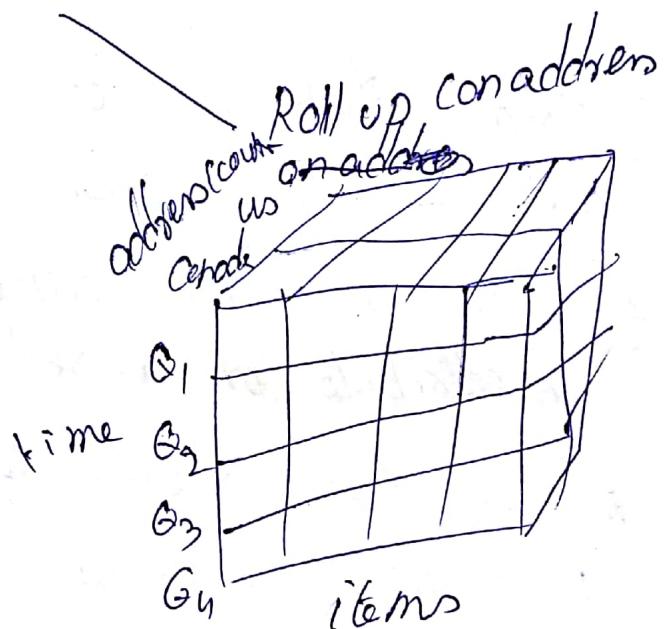
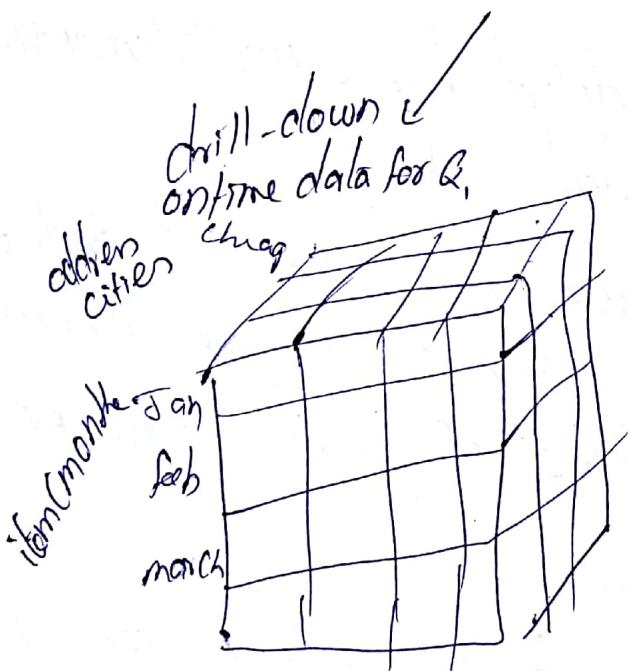
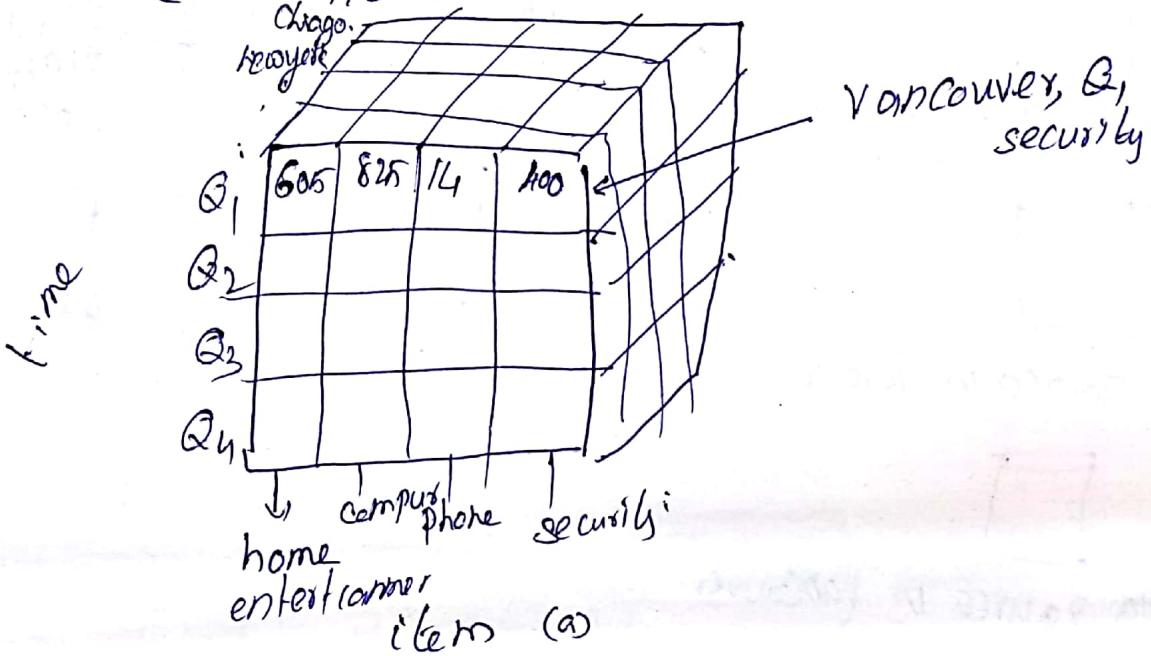
- A data ~~warehouse~~ is modeled by a multidimensional data structure called a data cube.
- in which each dimension correspond to set of an attribute or a set of attributes in the schema!
- each cell stores the value of some aggregation measures count or sum.
- A data cube provides a multidimensional view of data and allows the pre-computation and fast access of summarized data

(8)

(2)

e.g. A data cube for All electronics

- * Cube has three dimensions. (i) address (with city value Chicago, New York, Toronto, Vancouver)
- (ii) time (with Quarter values Q₁, Q₂, Q₃, Q₄)
- (iii) item (values home entertainment, computer, phone security)



item types

Transactional data:

- each record in a transactional database captures a transaction such as customer's purchase etc.
- transaction includes a unique transaction identity number and list of the items making up the transaction.

Eg

Trans id	List-of-item-ids
T 100	I ₁ , I ₃ , I ₈ , I ₁₆
T 200	I ₂ , I ₈
⋮	⋮ ⋮ ⋮ ⋮

Other kind of Data.

- time related or sequence data
(eg) historical data records, stock exchange data, etc.
- data streams:
(eg) video surveillance & sensor data, etc
- spatial data
eg, maps.
- engineering design data
design of buildings, system components etc
- hyper text & multimedia data.
(eg) text, image, video & audio data
- graph & N/c data
social & N/c data
- web
huge & widely distributed data

Data mining functionalities: or what kinds of patterns can be mined

* It is used to specify the kind of patterns to be found in data mining task.

* Data mining task are classified into two categories
 (i) descriptive
 (ii) predictive.

• Descriptive data mining task

* it ~~can~~ characterize the general properties of data in the database

Predictive data mining task

* it perform inference on the current data in order to make predictions

Functionalities

(i) Concept/Class description: characterization and discrimination.

* Data can be associated with classes or concepts.

e.g. In All electronic stores

classes of item for sale include
 Computers & printers

Concepts include big spenders or budget spenders.

(11)

data characterization,

- * by summarizing the data of the class under study (target class)

data discrimination

- * by comparison of the target class with one or a set of comparative classes (contrasting class)

Data characterization:

- * is a summarization of the general characteristics or features of a target class of data.
- * the data cube based OLAP rollup operation can be used to perform user-controlled data summarization
- * An attribute-oriented induction technique can be used to perform
- * the o/p of data characterization can be of various forms
eg. pie charts, bar charts, curves, Multidimensional data cubes & multidimensional tables

Data discrimination.

- * It is a comparison of the general features of the target class data objects against the general features of object from one or multiple contrasting class
- * It is expressed in the form of rules is called discriminant rule

Mining Frequent patterns, Associations, and correlation

Frequent patterns:

- * patterns that occur frequently in data.
- * Many kinds of frequent patterns like, frequent itemsets, frequent Subsequence, frequent Substructure
- * Frequent itemset :
it refers to a set of items that often appear together in a transactional data set (eg) Milk etc.
- * Frequent subsequence
* the patterns that the customers tend to purchase
eg laptop → digital camera — memory card

13

for association rule $X \rightarrow Y$, the support of the rule is denoted as $Sup(X \rightarrow Y)$ and is no. of transactions where $X \rightarrow Y$ appears divided by the total number of transaction.

- Substructure refers different structural forms eg. graphs, trees or lattices. They may combined with itemset or sequence
- if the substructure occur frequently it is called (frequent) structured patterns.
- Mining frequent patterns leads to the discovery of interesting associations and correlation with data

eg. Association Analysis

eg If the manager wants to know which items are frequently purchased together or co-appear in same transaction

$\text{buys}(X, \text{"computer"}) \Rightarrow \text{buys}(X, \text{"software"})$ [Support = 1%, Confidence = 50 %]

where X is a variable representing customer.

A Confidence or Certainty of 50% means that if a customer buys a computer there is a 50% chance that customer will buy s/w also.

A 1% Support means that 1% transaction under analysis show that computer & s/w purchase together

• Here association rule involves a single attribute

or predicate \in buys that repeat..

- Association rule that contains a single attribute or predicate are referred as Single dimensional association rule

e.g. $\text{age}(x, "20..29") \wedge \text{income}(x, "40k..49k") \Rightarrow \text{buys}(x, "laptop")$ [support = 2%, confidence 60%]

2% of 20 to 29 year old customer with income 40k to 49k have purchased laptop.

60% probability that a customer in this age and income group will purchase a laptop.

- More than one attribute or predicate is used in age, income, buy. So

- More than one attribute is referred as Multidimensional association rule

3. Classification & Regression for predictive Analysis

Classification:

- is the process of finding a model or function that describes and distinguishes data classes or concepts.

- The models are derived based on the analysis of a set of training data

(1) a data object for which the class labels are known.

- The derived model may be represented in various forms such as classification rules.
 - IF-THEN, decision trees, mathematical formulae or neural nets

A decision tree

- is a flow chart like tree structure where each node denotes a test on an attribute value, each branch represents an outcome of the test, and tree leaves represent classes or classification
- it can be converted to classification rules

Neural Net:

used for classification, is a collection of neuron like processing units with weighted connections between the units

Regression

- are continuous valued functions
- it is used to predict missing or unavailable numerical data values rather

- than class label prediction.

* CE is used for numeric prediction

then.

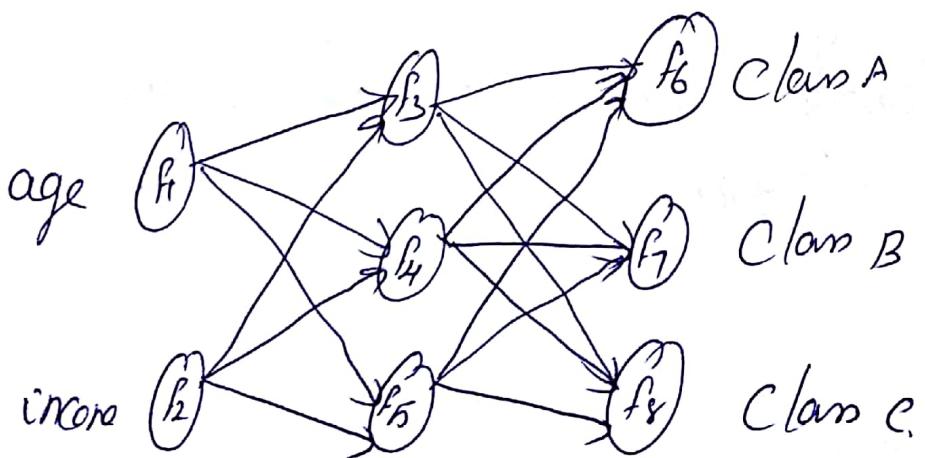
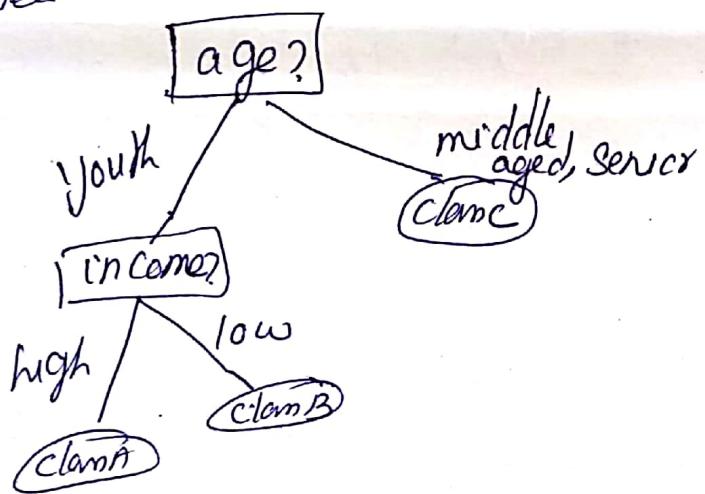
age (x , "youth") And income (x , "high") \rightarrow class (x , "A")

age (x , "youth") And income (x , "low") \rightarrow class (x , "B")

age (x , "middle-aged") \rightarrow class (x , "C")

age (x , "senior") \rightarrow class (x , "C")

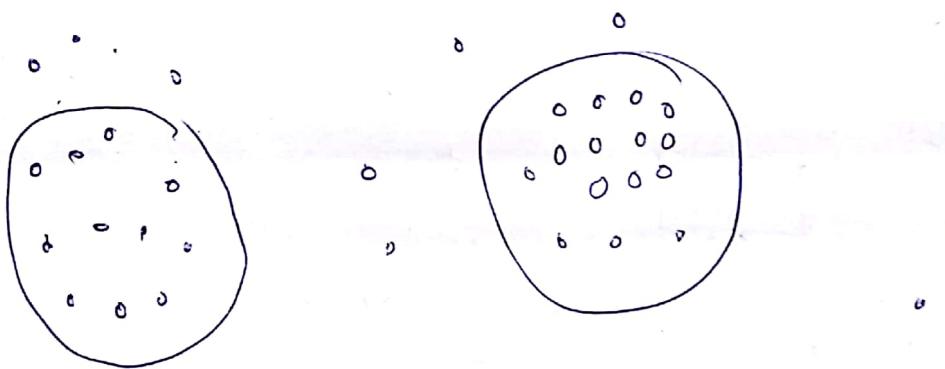
decision tree



(17)

4. Cluster Analysis

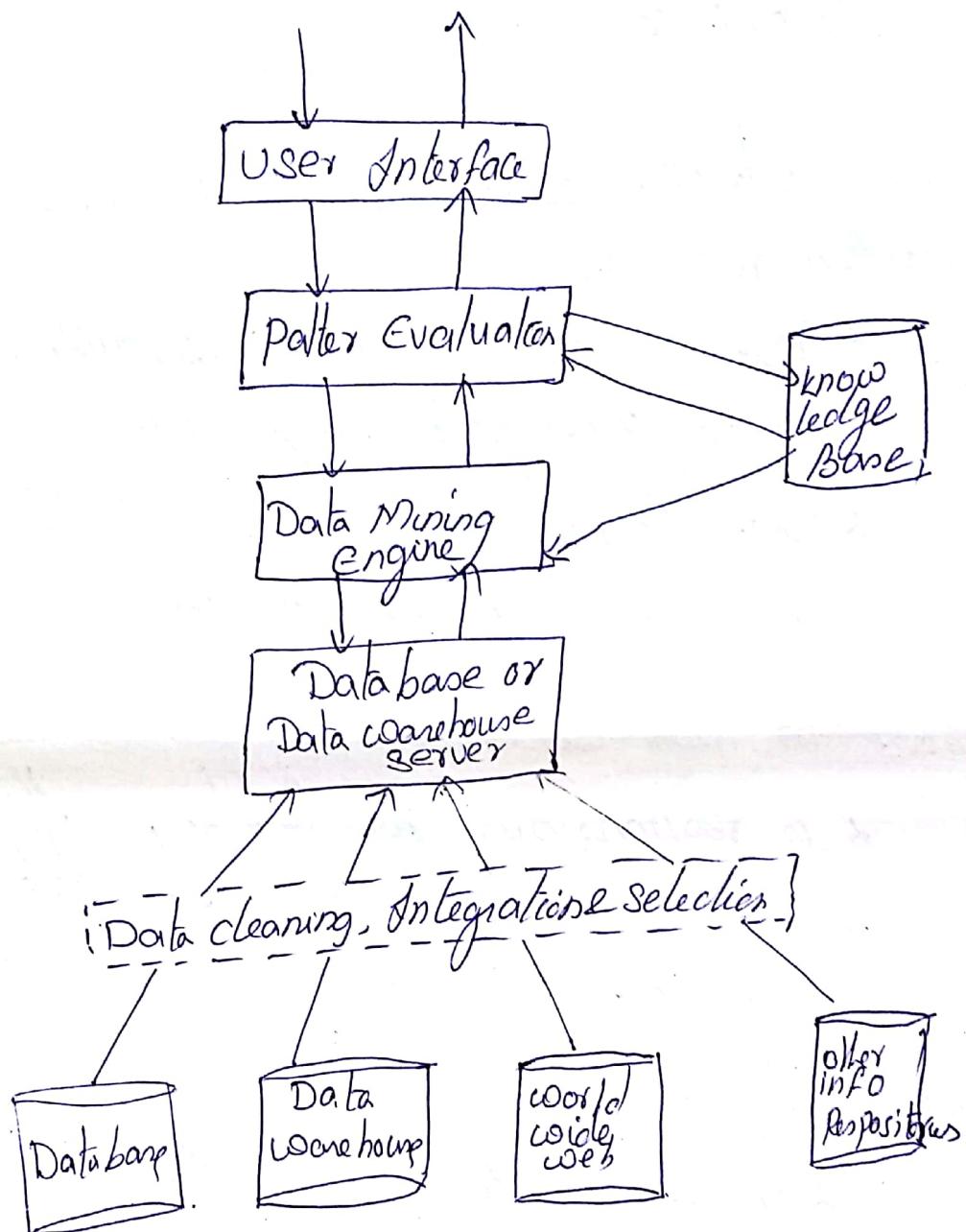
- clustering is used to generate class labels for a group of data.
- Objects are clustered or grouped based on the principle of maximizing the intra-class similarity & minimizing the inter-class similarity.



5) Outlier Analysis

- Data set may contain objects that do not comply with the general behavior or model of data.
- The analysis of outlier data is referred to as outlier Analysis or anomaly mining.

Architecture of Typical Data Mining Systems



Components One

- * Database, data warehouse, cww or other information Repository,
- * Database or Data warehouse server
- * knowledge base
- * Data mining Engine

(19)

* Pattern Evaluation Module

* User Interface

(i) Database, data warehouse, world wide web, or other information repository:

* This is one or set of db, dwh, www or other information repository

* Data cleaning, integration selection techniques may be performed on the data

(ii) Database or data warehouse Server:

* is responsible for fetching the relevant data, based on the user's data mining request.

(iii) knowledge base:

* is used to guide the search or evaluate the interestingness of resulting patterns. ~~such knowledge can include~~

* Concept hierarchies

* knowledge can include concept hierarchies used to organize attributes or attribute values into different levels

- * Data mining Engine;

- * Set of functional Modules for task such as Characterization, association and Correlation analysis, Classification, Prediction, cluster analysis, outliers analysis & evolution analysis

- * Pattern evolution Module;

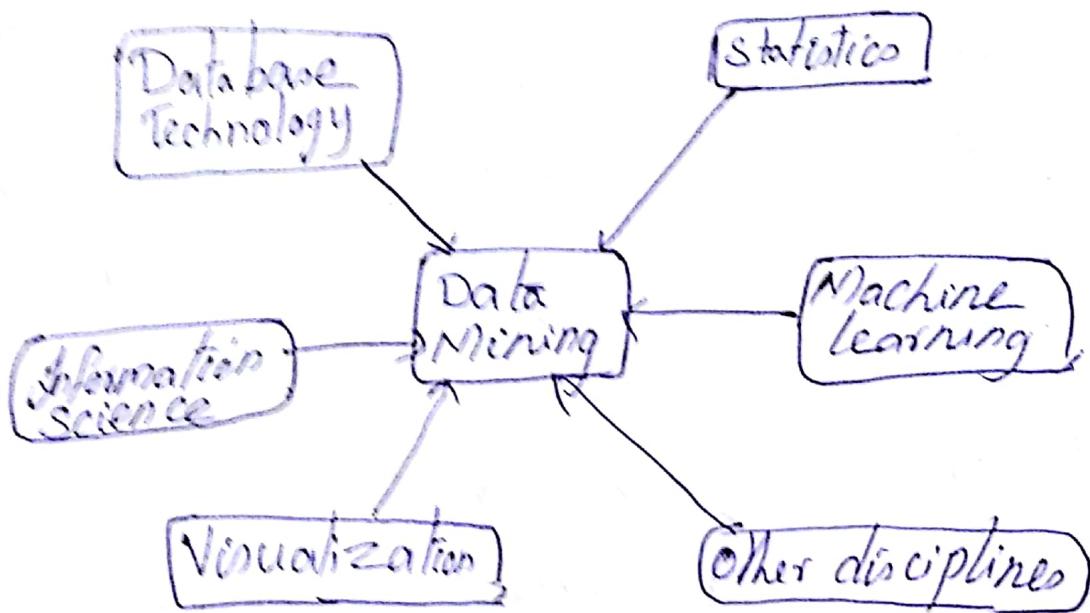
- * employs interestingness measures and interacts with the data mining modules so as to focus the search towards interesting patterns.

- * Pattern Evolution module is integrated with the mining module, depending on the implementation of the data mining method used

- * User Interface;

- * It communicates b/w users and the data mining system along the user to interact with the system by specifying a data mining query or task providing,

Classification of Data Mining System



- * DM consist of set of disciplines including
 - * database systems,
 - * Statistics,
 - * N/l/e learning
 - * visualization
 - * information science
- * Other disciplines include
Neural nets, fuzzy and/or rough set theory, knowledge representation, inductive logic programming &
high-performance computing

Classification according to the kind of database mined.

- * Data mining system can be classified according to the kinds of database mined.
- * Database systems can be classified according to criteria such as data models or types of data or application involved.
- * If classifying according to data model then we have relational, transactional, object-relational, or data warehouse mining system.
- * If according to types of data, then spatial, time series, text, stream data, multimedia or video mining system

Classification according to the kinds of knowledge Mined.

- * Based on the data mining functionalities such as Characterization, discrimination, association and Correlation analysis, Classification, Prediction, Clustering, Outlier and evolution analysis.

Classification according to the kinds of techniques utilized:

- The techniques can be described according to the degree of user interaction involved (autonomous systems, interactive, exploratory systems etc) or methods of data analysis employed (database oriented or dw-oriented techniques, M/c learning)

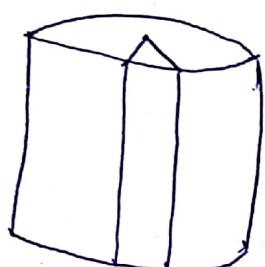
Classification according to the application adopted.

- DM can be categorized according to the applications they adapt.
e.g. for finance, telecommunication, DNA, stock markets, e-mail..

Data mining task Primitives

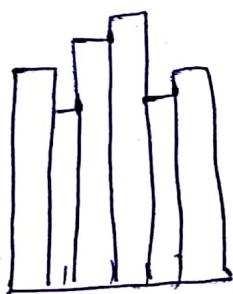
- Data mining task can be specified as the form of a data mining Query, which is the input to the data mining system.
- Data mining Query is defined in terms of data mining primitives.
- primitives allow the user to actively communicate with the data mining system during discovery of data

The data mining primitives specify the following.



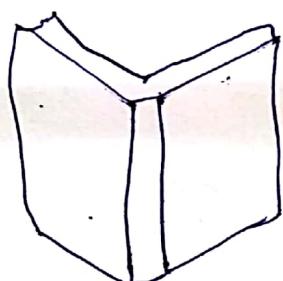
Task-relevant data

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



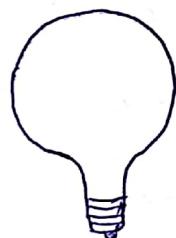
knowledge type to be mined

- * Characterization
- * Discrimination
- * Association/Correlation
- * Classification/Prediction
- * Clustering



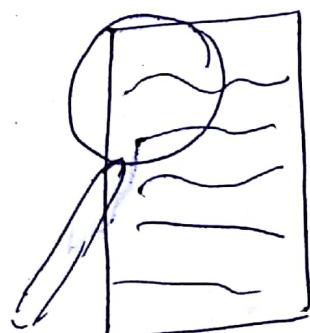
background knowledge

- * Concept hierarchies
- * user beliefs about relationships in the data



Pattern interestingness measures

- * Simplicity
- * Certainty (e.g. confidence)
- * Utility (e.g. support)
- * Novelty



Visualization of discovered patterns

- * Rules, tables, reports, charts, graphs, decision trees, & cubes
- * Drill-down & roll-up

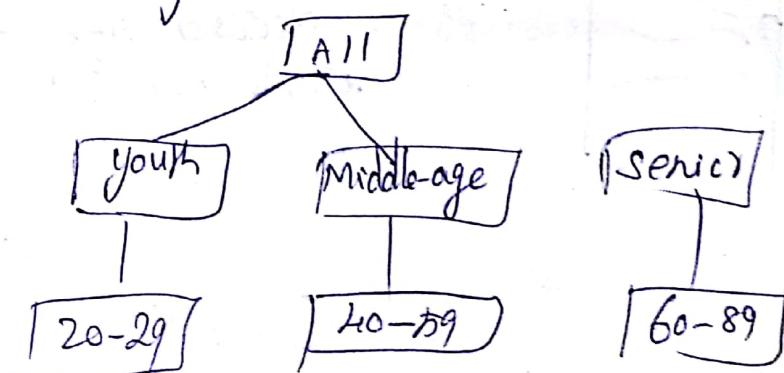
* Th Data mining primitives

- * The set of task relevant data to be mined
- * The kind of knowledge to be mined
- * The background knowledge to be used in the discovery process
- * The interestingness measures and thresholds for pattern evaluation.
- * The expected representation for visualizing the discovered patterns.

Data Mining Language [DMQL]

The design of an effective data mining Query language requires a deep understanding of the power, limitation and underlying mechanisms of the various kinds of data mining task.

(eg) Concept hierarchy for attribute (or dimension) age



Major issues of Data Mining

- (i) Mining Methodology & User Interaction issues
- (ii), Performance issues (iii), Efficiency & Scalability
- (iv) diversity of database type issues
- (v) Data Mining & Society
 - (i) Mining Methodology & User Interaction issues;
 - Mining Methodology issues
 - * Mining different kinds of knowledge in database;
 - * users can be interested in different kinds of knowledge, data mining cover world wide spectrum of data analysis and knowledge discovery tasks including data characterization, discriminative association and Correlation analysis, classification, prediction, clustering, outlier analysis
 - * ~~Interactive~~
 - * Mining knowledge in Multidimensional space
 - * Can search for interesting patterns among Combinations of dimensions (attributes) at various levels of abstraction
 - * Such mining is known as (exploratory) Multidimensional data mining.
 - * Mining knowledge in cube space
Can enhance the power & flexibility of data minin

Data Mining - a interdisciplinary effort:

- the power of Data Mining can be substantially enhanced by integrating new methods from multiple disciplines.
 - eg mining data with natural language text
 - (ii) the mining of S/w bugs in large programs

Boosting the power of discovery in Networked environment

- data objects reside in a linked or interconnected environment, it may be web, database relations, files or documents.
- knowledge derived in one set of objects can be used to boost the discovery of knowledge in a related or semantically linked set of objects
- Handling of Data
 - * Data contain noise, error, exceptions etc.
 - * Errors and noise may confuse the data mining process leading to the derivation of erroneous pattern.

- * Data Cleaning, data preprocessing, outlier are examples of techniques need to be integrated with data mining process
- * Pattern Evaluation and pattern or constraint-guided Mining :

By using interestingness measures or user specified ~~measures~~^{constraint} to guide the discovery process.

User Interaction issues

* Interactive Mining

- * DM is highly interactive
- Interactive mining should allow users to dynamically change the focus of a search, to refine mining request based on return results and to drill, dice and pivot thro the data and knowledge space interactively

* Incorporation of background knowledge:

- * Background knowledge, constraints, rules and other information regarding the domain should be incorporated into the knowledge discovery process

Adhoc data mining and DMQL:

- * Query languages have played an important role in flexible searching because they allow users to pose adhoc queries
- * High level DMQL or high level flexible user interfaces will give users to define adhoc data mining tasks

Presentation & visualization of data mining results:

- * How to a data mining system present data mining results, vividly & flexibly, so that the discovered knowledge can be easily understood and usable by humans.

3. Efficiency & scalability:

- * Efficiency & Scalability of data mining algorithm
 - * Data Mining algorithms must be efficient & scale in order to effectively extract information from huge amounts of data on many data repositories.
- * Efficiency, Scalability, performance, optimization and the ability to execute in real time are key criteria for the development of new data mining algorithms

- parallel, distributed and incremental mining algorithms
- Size of data set, distribution of data & Computational complexity are factors that motivate the development of parallel & distributed data-intensive mining algorithms
- Algorithms first partition the data into pieces. Each piece is processed by searching for patterns
- cloud computing & cluster computing use computers on a distributed and collaborative way to solve very large-scale computational task

④ Diversity of Database type.

- Handling complex types of data
- it generate wide spectrum of new data types, from structured data such as relational & data warehouse data to semi structured and unstructured data; from stable data repositories to dynamic data from simple data to temporal data object etc.

- * Mining dynamic, networked and global data repositories:
 - * Multiple source of data are connected by the internet & various kinds of N/w forming distributed & heterogeneous global information systems & N/w.

b. Data Mining & Society

- * Social impact of data mining:
 - * Based on the everyday lives it is important to study the impact of data mining on society.
- * Privacy - preserving data mining:
 - * DM will help scientific discovery, business management, economy and security protection.
- * Invisible data mining:
 - * Intelligent search engines & Internet based stores perform such invisible data mining by incorporating data mining into their components to improve their functionality & performance