

CSE 16 - Data Warehousing and Data Mining

UNIT 1

1. Introduction to Data Mining	1
2. Getting to know your data	2
3. Introduction to Data Warehousing.....	3
3.1 Basics of Data Warehouse	
3.2 Decision Support Systems	
3.3 Operational versus DSS	
4. Architecture of DWH.....	4
4.1 ETL, OLAP vs OLTP	
4.2. OLAP operations	
5. Dimensional Data modeling	5
5.1. Star Schema	
5.2. Snowflake Schema	
5.3. Fact constellation Schema	
6. OLAP Operations.....	6

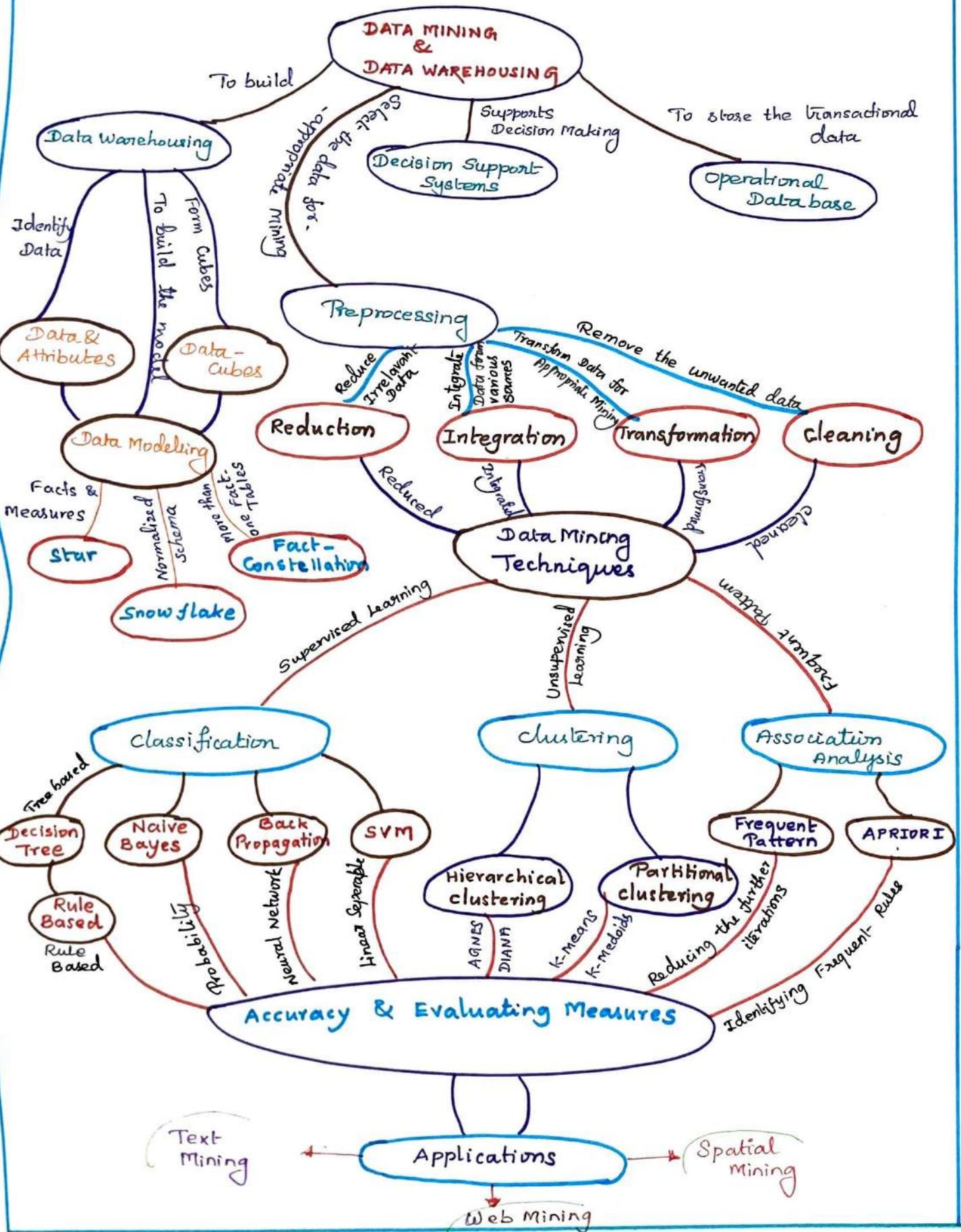
UNIT 2

10.Basic concepts in classification & prediction.....	10
11.Decision Tree	11
12.Bayesian Classification	12
13.Rule Based Classification	13
14.Back Propagation, SVM, Associate Classification	14
15.Accuracy and Error Measures	15
15.1. Cross validation (k-fold)	
15.2. Metrics for performance evaluation	
16.Req. for Cluster Analysis & Partitional clustering.....	16
16.1. K-Means clustering	
16.2. K-Medoids	
17.Hierarchical clustering	17
17.1. AGNES & DIANA	
18.Frequent Pattern Mining.....	18
18.1. FP Growth	
18.2. Apriori	

UNIT 4

19.Text Mining, Web Mining	19
20.Spatial Mining	20
21.Applications & Research aspects of Data Mining.....	21

UNIT 5



Data Mining :-

→ Tools to discover knowledge from data

→ Knowledge from data

→ Data platform analysis

→ Knowledge extraction

Kinds of data :-

Database data

- * Context of data

- * Primarily stored in database table

Transactional data

- * Organized by time stamps
will backs.

- * Repository of information
- * Analysis to make more decisions.

Data warehouse

Other kinds of data

- * Big data
- * Structured, semi structured, unstructured data
- * Time stamp data

- * Machine data

Kinds of pattern to be mined:

* Data to be associated

Class concept

Data objects

Anomaly mining

Outlier analysis

cluster analysis

* Used multi-dimensional data

Association analysis

pattern to be mined

* Minimizing intra class similarity

Frequent Pattern

* Pattern occurs frequently

Regression

- * Continuous valued function
- * Statistical method

* Finding a model

Ex: Decision tree neural networks

Data mining adopts techniques from many domains:

supervised

Unsupervised

Semi supervised

Searched data

* Queries are formed by keywords

↑

Machine learning

Statistics

* performs the analysis and help

* visualize data

Pattern Recognition

* automated recognition of patterns

* Representation, feature extraction, classification

Visualization

* Makes it easier to identify

* to identify patterns and outliers in large database

↓

Data mining

↓

Information Retrieval

↓

Algorithm

↓

SVM

APriori

CART

↓

KNN

Naïves bayes

↓

Data mining

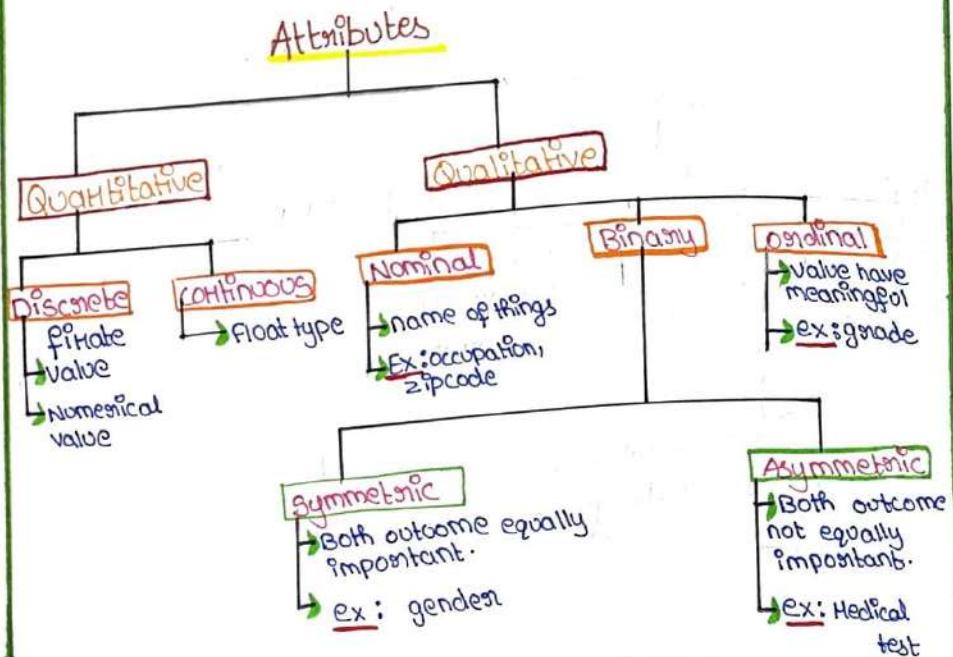
↓</

GETTING TO KNOW YOUR DATA

Data objects and Attribute Types

→ Data objects → represent entity

- Ex :- Sales Data base : customer, store items
- described attributes
- Data rows → Data objects
- columns → data attribute



Basic Statistical Descriptions of Data

Motivation - Better understand the data : central tendency, variation and spread.

Data dispersion characteristics - Median, max, min, outliers.

Numerical Dimensions - Data dispersion, Box plot.

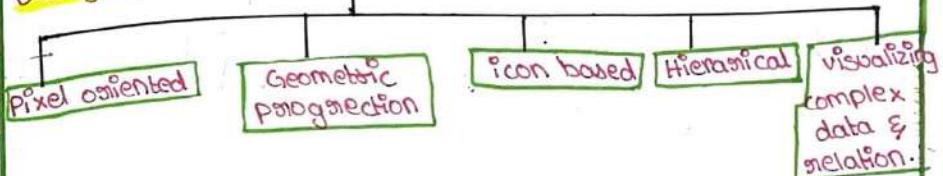
Dispersion analysis on computed measures -

- * Folding measures into numerical dimension.
- * Boxplot on the transformed cube.

Data Visualization

- Mapping data onto graphical primitives.
- provide qualitative overview of large datasets.
- Search for patterns.
- provide a visual proof.

Category of Visualization Method



Similarity Measures

- real value function that quantifies the similarity between two objects.
- Measure how two data objects are alike
- often falls in the range $[0,1]$: 0 - no similarity
1 - completely similar

Dissimilarity Measures

- Numerical measure of how different two data objects.
- Minimum dissimilarity
- range $[0,1]$ or $[0,\infty]$

Data Warehousing:

- * provides architecture and tools for business executives to systematically organize, understand and use their data to make strategic decisions.
- * DW is subject oriented, integrated, time variant and non volatile.



Subject Oriented: Data warehouse is organised around major subjects such as customer, supplier, product and sales.

Integrated: A dw is usually constructed by integrating multiple heterogeneous sources such as flat files, Relational db and online transaction records.

Time variant: Historic data/information [5-10 yrs]

Non volatile: Application data found in Operational environment [permanent storage].

Difference b/w Datawarehouse & Operation database

Datawarehouse	Operation database
* Datawarehouse is repository for structured, filtered data.	* Database changes frequently.
* Denormalized schema.	* Normalized schema.
* Historical data	* Current transaction data.
* Online analytical processing.	* Online transaction processing.

Decision Support System (DSS):

- * It is computerized program used to support determinations, judgements and course of actions in organization.

Ex:

DSS: target (or) projected revenue, forecast the sales. It is completely computerized & powered by humans.

Characteristic	Operation data	DSS
Data currency/ granularity	Real time data/ atomic-data.	Historic data/ summarized data.
Data model	RDBMS	non-normalized model.
Transaction volumes	high	Summary
Query complexity	Simple to medium	very complex
Data volumes	Hundreds of GB	Tera (or) petabytes
Query activity	Low to medium	High

DATA WAREHOUSING ARCHITECTURE

query/repost

analyse

Data mining

option tool set

The data is transformed & consolidated from its intended analytical use case.

SQL or NoSQL scenarios → CRM & ERP systems → Flat files.

→ Email → webpages.

Extraction, Transformation & loading Extracted, ETL

pulling data extraction, raw data is copied, structured or unstructured. These sources of data include but are not limited to:

SQL or NoSQL scenarios → CRM & ERP systems → Flat files.

High performance, high availability, high flexibility, and user autonomy.

High performance, high availability, high flexibility, and user autonomy.

High performance, high availability, high flexibility, and user autonomy.

High performance, high availability, high flexibility, and user autonomy.

High performance, high availability, high flexibility, and user autonomy.

High performance, high availability, high flexibility, and user autonomy.

High performance, high availability, high flexibility, and user autonomy.

High performance, high availability, high flexibility, and user autonomy.

High performance, high availability, high flexibility, and user autonomy.

High performance, high availability, high flexibility, and user autonomy.

High performance, high availability, high flexibility, and user autonomy.

High performance, high availability, high flexibility, and user autonomy.

High performance, high availability, high flexibility, and user autonomy.

High performance, high availability, high flexibility, and user autonomy.

High performance, high availability, high flexibility, and user autonomy.

High performance, high availability, high flexibility, and user autonomy.

High performance, high availability, high flexibility, and user autonomy.

High performance, high availability, high flexibility, and user autonomy.

High performance, high availability, high flexibility, and user autonomy.

High performance, high availability, high flexibility, and user autonomy.

High performance, high availability, high flexibility, and user autonomy.

High performance, high availability, high flexibility, and user autonomy.

High performance, high availability, high flexibility, and user autonomy.

High performance, high availability, high flexibility, and user autonomy.

High performance, high availability, high flexibility, and user autonomy.

High performance, high availability, high flexibility, and user autonomy.

High performance, high availability, high flexibility, and user autonomy.

High performance, high availability, high flexibility, and user autonomy.

High performance, high availability, high flexibility, and user autonomy.

High performance, high availability, high flexibility, and user autonomy.

High performance, high availability, high flexibility, and user autonomy.

High performance, high availability, high flexibility, and user autonomy.

High performance, high availability, high flexibility, and user autonomy.

High performance, high availability, high flexibility, and user autonomy.

High performance, high availability, high flexibility, and user autonomy.

High performance, high availability, high flexibility, and user autonomy.

High performance, high availability, high flexibility, and user autonomy.

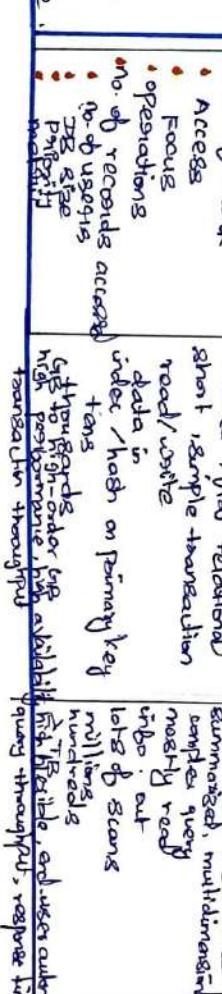
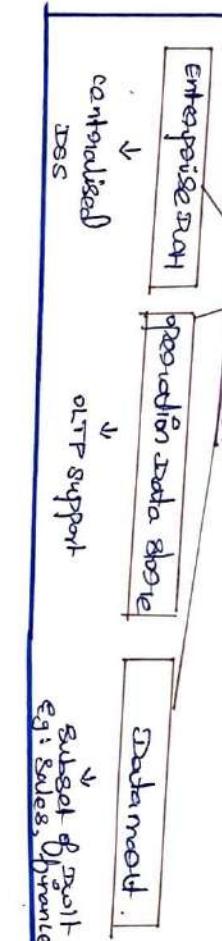
High performance, high availability, high flexibility, and user autonomy.

High performance, high availability, high flexibility, and user autonomy.

High performance, high availability, high flexibility, and user autonomy.

High performance, high availability, high flexibility, and user autonomy.

High performance, high availability, high flexibility, and user autonomy.



Dimensional Data Modelling

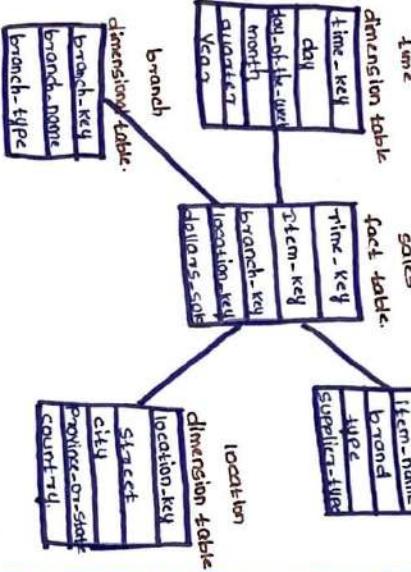
Cicago
New York
Toronto
Vancouver

	H	C	P	S
Q1	605	825	10	400
Q2	650	952	12	512
Q3	812	1023	31	573
Q4	927	700	21	545

Star, snowflake and fact constellation :-

* schemas for multi dimensional data model!

Star - schema :-



Fact constellation schema :-

* It consists of dim-table that are shared by several fact tables

→ bulk of data with no-redundancy.
→ a set of small & attendant table [1 - for each dimension]

* It contains both dimensional table and fact table.

→ also known as "Galaxy Schema"

Star schema definition :-

define cube sales_star [time, item, branch, location];
dollars_sold = sum (sales ~in-dollars), units_sold = count(*)

define dimension time as (time-key, day, day-of-week, month)

define dimension item as (item-key, it-name, item-type)

define dimension branch as (branch-key, bran-name, bran-type)

define dimension location as (loc-key, street, city, province ~or-state)

define dimension item as (item-key, it-name, item-type)

define dimension branch as (branch-key, bran-name, bran-type)

define dimension location as (loc-key, street, city, province ~or-state)

define dimension item as (item-key, it-name, item-type)

define dimension branch as (branch-key, bran-name, bran-type)

define dimension location as (loc-key, street, city, province ~or-state)

define dimension item as (item-key, it-name, item-type)

define dimension branch as (branch-key, bran-name, bran-type)

define dimension location as (loc-key, street, city, province ~or-state)

define dimension item as (item-key, it-name, item-type)

define dimension branch as (branch-key, bran-name, bran-type)

define dimension location as (loc-key, street, city, province ~or-state)

define dimension item as (item-key, it-name, item-type)

define dimension branch as (branch-key, bran-name, bran-type)

define dimension location as (loc-key, street, city, province ~or-state)

define dimension item as (item-key, it-name, item-type)

define dimension branch as (branch-key, bran-name, bran-type)

define dimension location as (loc-key, street, city, province ~or-state)

define dimension item as (item-key, it-name, item-type)

define dimension branch as (branch-key, bran-name, bran-type)

define dimension location as (loc-key, street, city, province ~or-state)

define dimension item as (item-key, it-name, item-type)

define dimension branch as (branch-key, bran-name, bran-type)

define dimension location as (loc-key, street, city, province ~or-state)

define dimension item as (item-key, it-name, item-type)

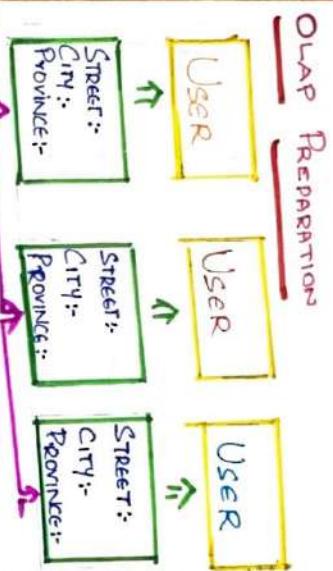
define dimension branch as (branch-key, bran-name, bran-type)

define dimension location as (loc-key, street, city, province ~or-state)

Online Analytical Processing

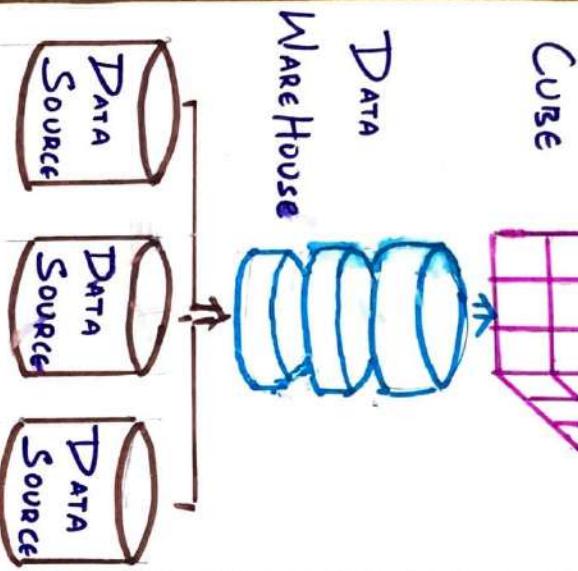
OLAP - Multidimensional

Information - Manage, Analysis, Interactive



App

OLAP Cube



location	Toronto	305
Vancouver	605	

location	Vancouver	605
Canada	1000	1000

OLAP - Operations

Roll Up : Aggregation on data

ways :- (1) climbing up on a concept hierarchy for a dimension

"street < city < province < country"

Aggregated hierarchy location from "city" to "country"

Roll-Down : stepping down a hierarchy

* New dimension

Slice : Create a new sub-cube "day < Month < Quarter < Year"

work : hierarchy for dimension time from one - particular cube

way : Dimension "time" -> time = "Q,"

Dice : Three dimension

Location = "Toronto" or "Vancouver"

Time = "Q₁" or "Q₂"

Item = "Mobile" or "TV"

Pivot : PIVOT - Rotation Operation

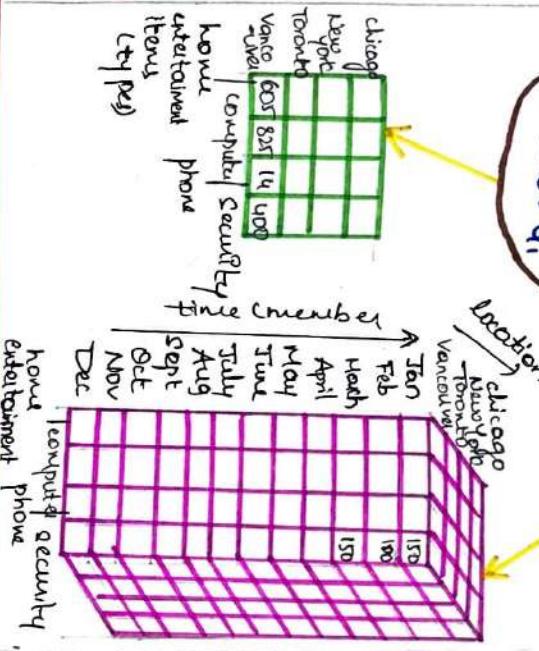
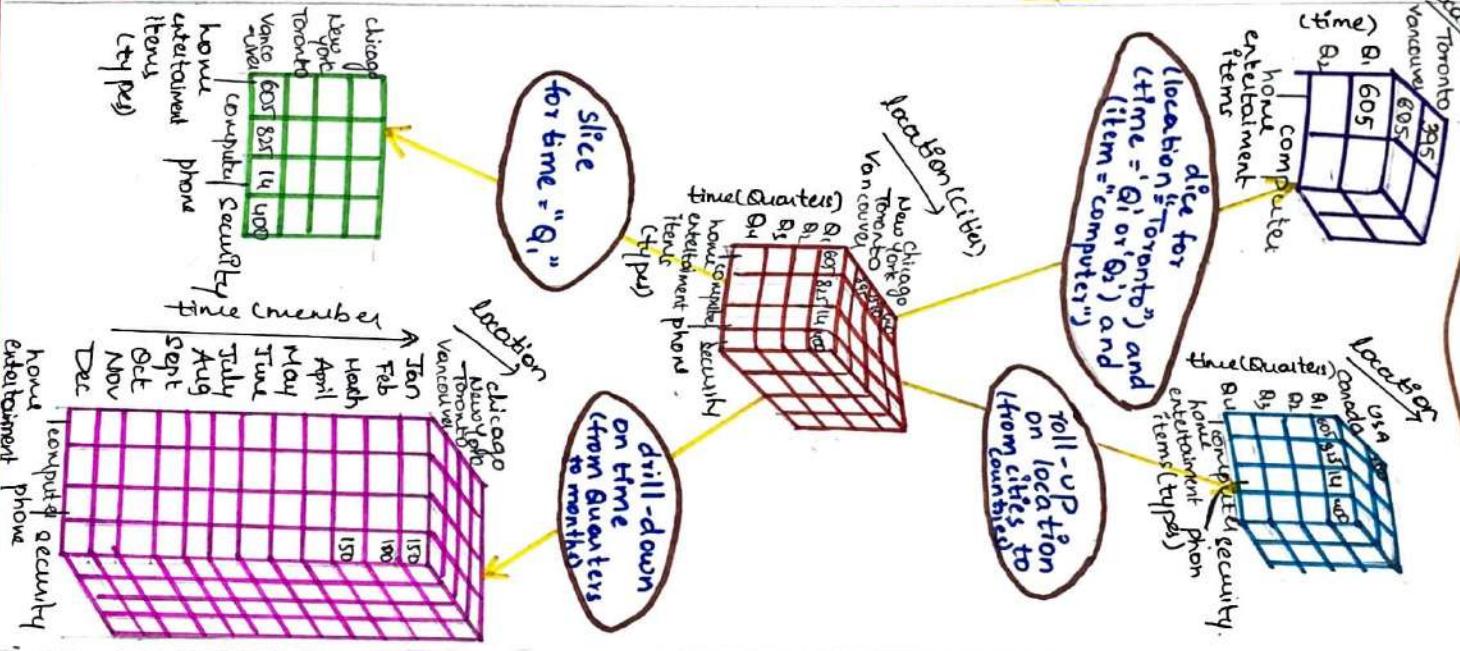
TOTAL NO OF CUBOIDS :-

$$\prod_{i=1}^n (L_i + 1)$$

L_i = Levels associated with dimension

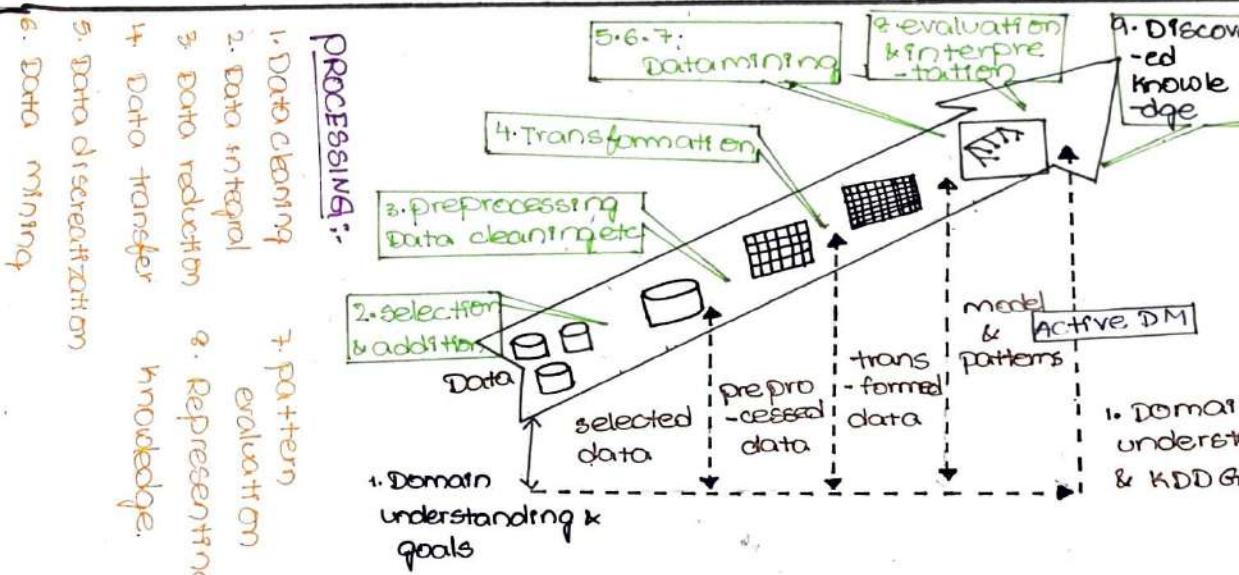
1 = Virtual top level

DATA Cube Computation



item	Jan	Feb	Mar
computer	100	150	100
security	200	250	200
phone	300	350	300

KNOWLEDGE DISCOVERY IN DATABASE (KDD):-

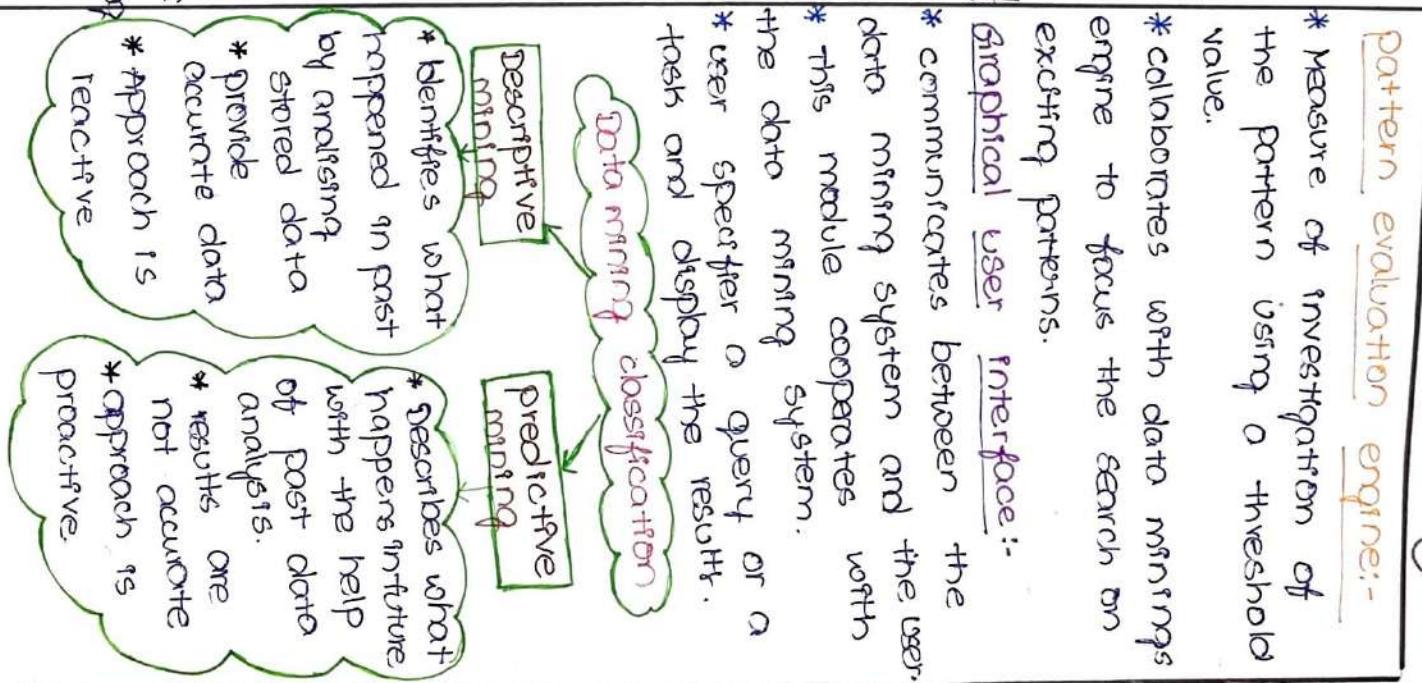
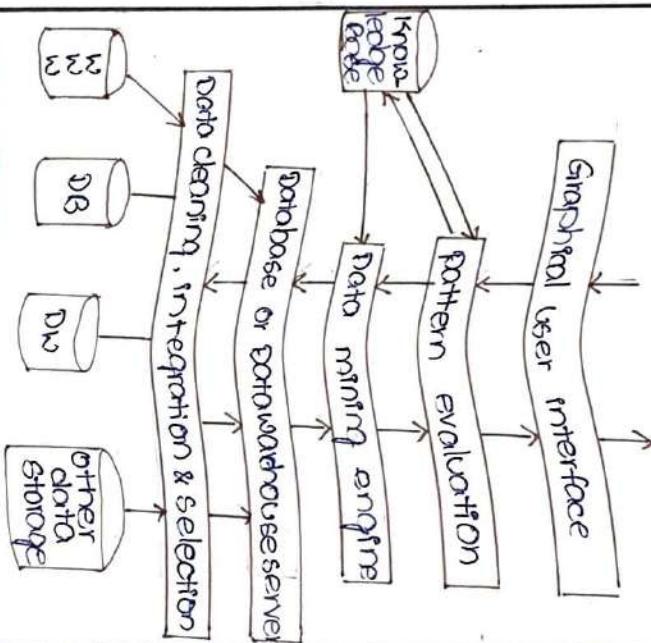


DATA SOURCE:-
Database, www, text files and other documents.

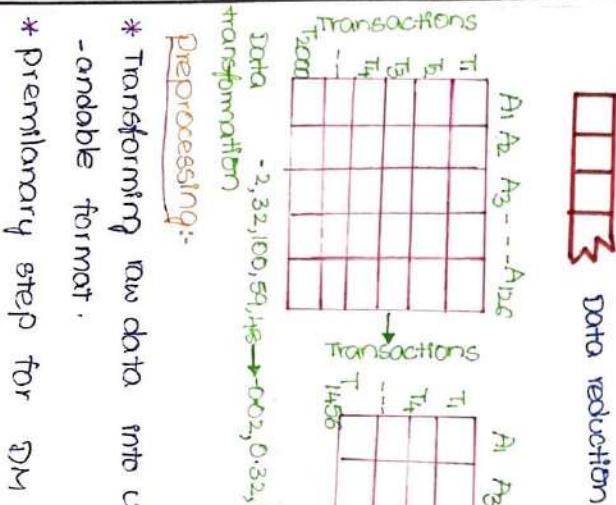
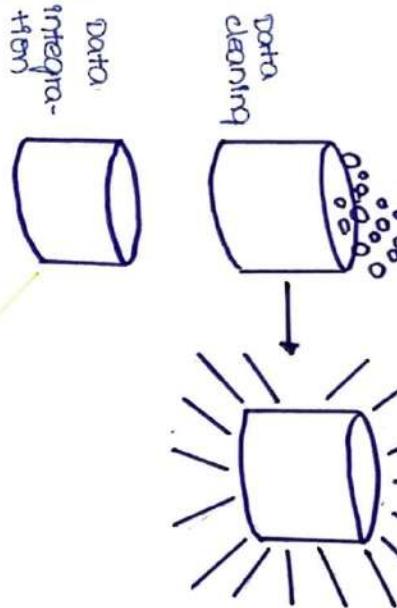
Organization store data in datawarehouse different process:-

- * Data must be cleaned
- * Data must be integrated & selected.
- * Data of interest to be selected & passed to the server.

Data mining engine: contains modules like association, characterization, classification, clustering prediction, time-series analysis



STEPS IN DATA PREPROCESSING



- * check data quality
- 1. Data cleaning:-
* remove incorrect data.
- 2. Data integration:-
* combining multiple sources into a single data set.
- 3. Data reduction:
* Reduction of data volume.
- 4. Data transformation- Mining task
* convert data into required format.

DATA CLEANING:-

- * Data - real world is dirty.
- * incomplete-locking attribute values eq: age = "NOISY"

* contain noise or error salary = "10"
INCONSISTENT: contain discrepancies, age = 10
How to handle missing data?

Ignore the tuple:

- * fill the missing value manually.
- * fill automatically: global constant (mean value).

How to handle noisy data?

Binning:-

- 1. sort data & partition into (equal frequency) bin.
- 2. smooth bin means, median by boundaries

Regression:-

- * Transforming raw data into understandable format.
- * preliminary step for DM

Preprocessing:-

- * Transforming raw data into understandable format.
- * smooth by fitting data into regression clustering: detect & remove outliers

DATA CLEANING AS A PROCESS:-

DATA DISCREPENCY DETECTION:-

- * use metadata
- * check uniqueness rule, consecutive rule

DATA MIGRATION & INTEGRATION:-

- * Allows transformation to be specified
- ETL:- allows users to specify transformations through query.

DATA INTEGRATION:-

- * combines data from multiple sources into a coherent store.
- Handling redundancy in data integration:-

- * redundant data occur often when integration of multiple database.

X² (chi-square) test.

- * Object identification

* Derivable data.

- * Redundant attributes detect using correlation analysis and covariance analysis.
- Correlation analysis:-

$$\chi^2 = \frac{\sum (\text{observed} - \text{expected})^2}{\text{expected}}$$

large χ^2 value, more likely the variables are related.

$$\text{Covariance: } \text{cov}(A, B) = E((A - \bar{A})(B - \bar{B})) = \frac{1}{n} \sum_{i=1}^n (A_i - \bar{A})(B_i - \bar{B})$$

$$VA, VB = \text{cov}(A, B) / \sigma_A \sigma_B$$

positive covariance:- if $\text{cov}_A, B > 0$.
negative covariance:- if $\text{cov}_A, B < 0$.

larger than its expected value, B is small than expected value.

Independence : $\text{cov } A, B = 0$

Data Processing : Data Reduction & Transformation

(9)

Data Reduction :- Reduced representation
Data set that smaller in volume.

Produces the same analytical result.

Data Reduction Strategies

Dimensionality reduction

Wavelet Transform :- To preserve relative distance between objects at different levels of resolution.

Principal Component analysis :- (PCA)

Original data are projected onto a much smaller space, resulting in dimensionality reduction.

Find eigen vectors of the covariance matrix, these eigen vector defines space.

Numerosity reduction :- Regression and log-linear

Linear :- Data modelled to fit a straight line.

Histogram analysis :- Divide data into buckets and store average (sum) of each bucket population size. Equal bucket range and equal frequency.

Data Transformation :- Which maps the entire set of values of altitude to new set of values.

Smoothing :- Remove noise from data altitude. New attribute - constructed from the given ones.

Aggregation :- Summarization, data cube construction.

Normalization :- Fall within a small min - max normalization:-

$$v_i = \frac{v - v_{\min}}{v_{\max} - v_{\min}} \quad (\text{new_max} - \text{new_min})$$

Z score - Normalization :-

$$v_i = \frac{v - \bar{v}_A}{\sigma_A}$$

\bar{v} = Mean σ = standard deviation

Normalization by decimal scaling :- $v'_i = v / 10^j$; j = smallest integer such that $\max(v'_i) \leq 1$

Discretization :- Divide the range of continuous attribute into intervals.

Interval labels can be used to replace data values.

Reduce data size by discretization method.

Binning :- Top down split

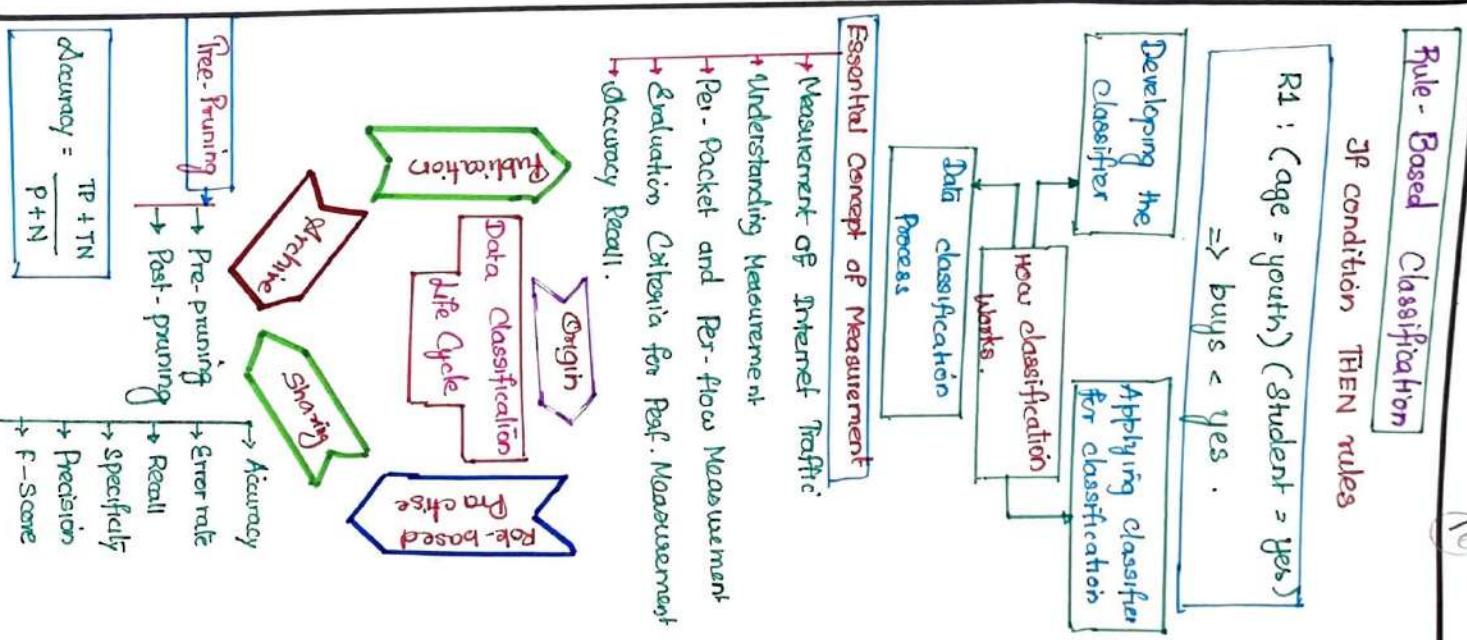
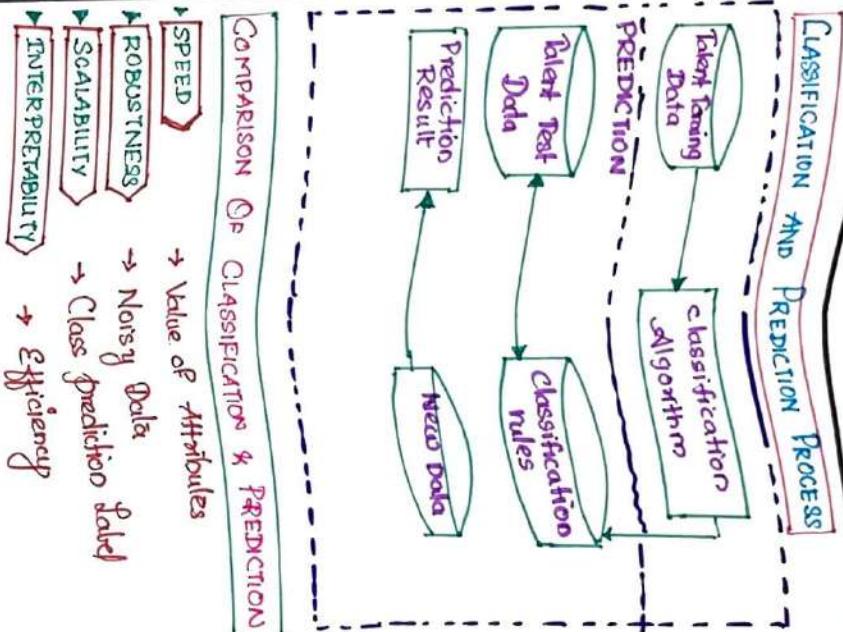
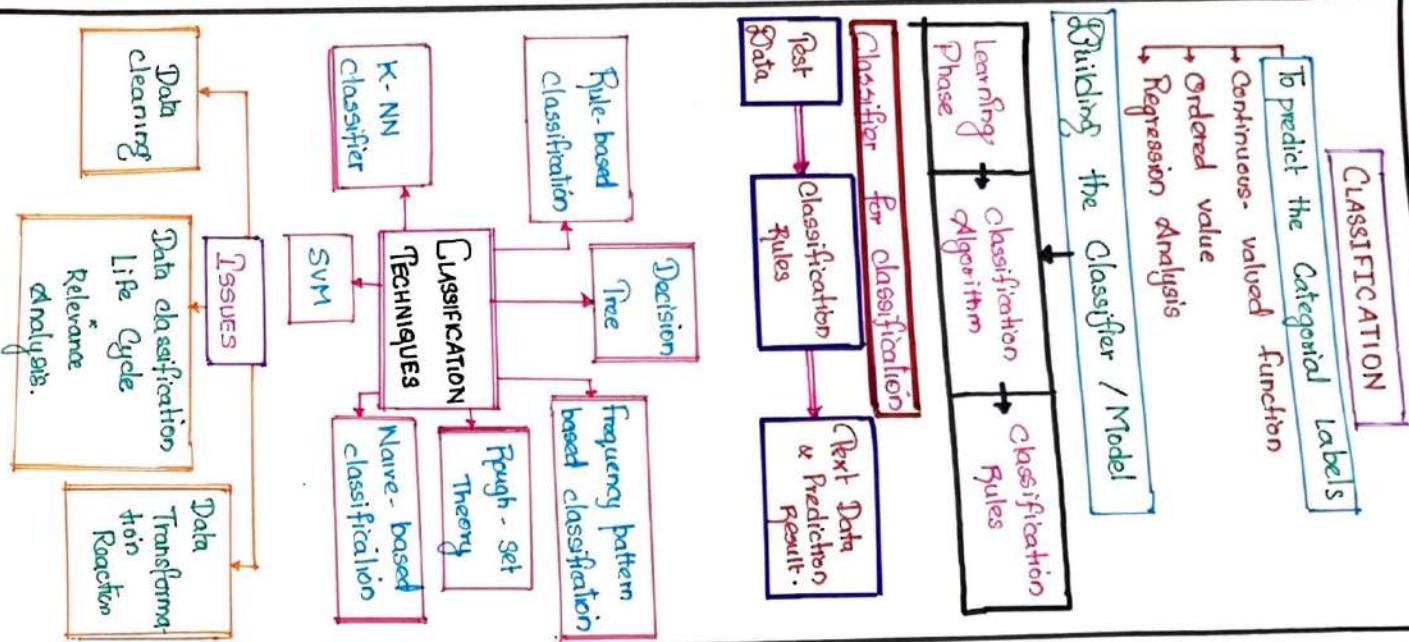
Histogram analysis :- Top down split.

Equal-width (distance) Partitioning :- Divide the range into N intervals of equal size.

$\Rightarrow A \& B$ are the lowest & highest values of attribute wide $w = (B - A) / N$.

Equal-depth (Frequency) Partitioning :- Divide the range into N intervals, each containing approximately same number of samples.

\Rightarrow Good data scaling :- Concept hierarchy organized, concept hierarchically and is usually associated with each dimension in a data warehouse.



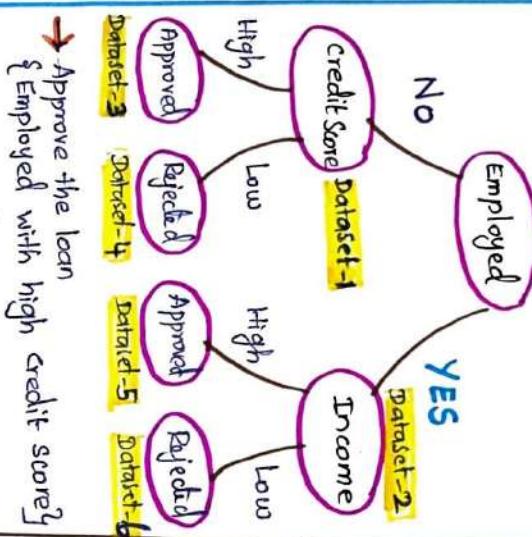
DECISION TREE

* Tree based classification & regression can be done using 'Decision Tree'.

* Process :
Dataset \rightarrow D.T Algorithm \rightarrow classifies the data

Example :

Decides if the loan should be Approved / Rejected



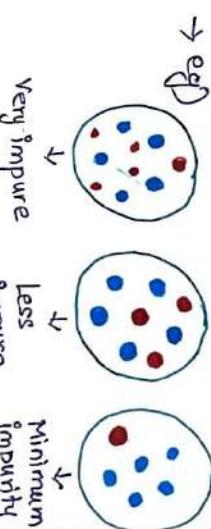
- Approve the loan & { Employed with high credit score }
- Approve the Loan & { Employed with high income }
- Reject the loan - { Unemployed with low credit score }
- Reject the loan - { Employed with low income }
- What is the source to repay the loan

challenges :

* Attribute Selection
* Popular Attribute Selection

1. Information Gain :- (IG):
IG is calculated for a split by subtracting the weighted "entropies" of each branch from the original entropy.

2. GINI Index
Entropy : \rightarrow Measures the change in entropy i.e measure the uncertainty in a group of observations.
 \rightarrow e.g)



$$\rightarrow \text{Formula :- } E = - \sum_{i=1}^N p_i \cdot \log_2 p_i$$

Example:
 $X = \{a, a, a, b, b, b, b, b\}$

Total instances = 8

Instance a = 3

$$\text{Entropy} = - \left[\frac{3}{8} \log \frac{3}{8} + \frac{5}{8} \log \frac{5}{8} \right] = 0.954$$

(ii) GINI Index :-

\rightarrow Gini Index & entropy are the measures for calculating information gain (IG).

\rightarrow D.T. algorithm uses IG to split a node. Both Gini Index & Entropy are the measures of impurity of a node.

$$\text{Gini} = 1 - \sum_{i=1}^n p_i^2 (c_i)$$

$$P(c_i) \rightarrow \text{The probability of class } c_i \text{ in a node.}$$

$$\text{IG} = (\text{Entropy of Parent node}) - \left(\sum_{i=1}^n \text{Weighted entropies of child node} \right)$$

Decision Tree Advantages & Disadvantages :-

- Advantages :-
They are very fast & efficient
Compare to KNN
- Easy to understand, interpret & visualize
- All type of data such as numerical & categorical are possible.
- Disadvantages :-
→ Training the model take higher time
→ Inadequate for applying regression & predicting continuous values.

BAYESIAN CLASSIFICATION

* Classification technique based on "Bayes' theorem" with an assumption of independent among features.

* The presence of one feature does not affect the other feature. All the features are independent of each other.

* Parametric model [Assumptions about a form of a function to ease the learning process]

Bayes' Formula :-

$$P(A|B) = \frac{P(B|A) \cdot P(A)}{P(B)}$$

Bayesian classification uses Bayes' theorem to predict the occurrence of any event.

* A & B → Events $P(A) \neq 0$

$P(A|B) \rightarrow$ Conditional Probability

[Occurrence of event A is given that B is true]

$P(B|A) \rightarrow$ Conditional Probability → Likelihood

[Occurrence of event B is given that A is true]

* $P(A) \& P(B) \rightarrow$ Probabilities of observing A & B independently of each other.

Example : $\begin{cases} \text{Marginal Probability} \\ \text{Fruit} = \{\text{Yellow, Sweet, Long}\} \end{cases}$

FRUIT	Yellow	Sweet	Long	TOTAL
Orange	350	0	650	
Banana	400	300	350	400
Others	50	100	50	150
TOTAL	800	850	400	1200

→ To predict & classify which one has the most probability of yellow, sweet & long.

$$\therefore \text{Probability (Yellow | Orange)} = \frac{P(\text{Orange}|\text{Yellow})}{P(\text{Yellow})}$$

$$= \frac{350}{800} \times \frac{800}{1200}$$

NOTE

{ • 800 → Yellow Fruit }

{ • 1200 → Total Fruits }

$$= \frac{0.4375 \times 0.667}{0.541667}$$

• 350 → Orange

• 800 → Total

• 1200 → Total

$$\therefore \text{Probability of Yellow | Orange} = 0.5387 //$$

$$P(\text{Fruit | Others}) = \frac{P(\text{Yellow})}{\text{others}} \times P(\text{Sweet}) \times P(\text{Long})$$

$$= 0.33 \times 0.667 \times 0.33$$

$$= 0.072$$

From the above,

→ Fruit of Orange is zero,
∴ eliminated

→ Fruit of Banana & } Which fruit of other fruits } is the most probability of yellow, sweet & long.

$$\text{Conclusion :-}$$

$$\begin{aligned} \text{Prob. (Sweet | Orange)} &= \frac{P(\text{Orange}) \cdot P(\text{Sweet})}{P(\text{Orange})} \\ &= \frac{450}{850} \times \frac{850}{1200} \\ &= 0.529 \times 0.7083 \\ &= 0.541667 \\ &= 0.6917 // \end{aligned}$$

Banana is the fruit which is having yellow, sweet & long compare to other all given fruits.

$$P(\text{Fruit | Orange}) = P(\text{Yellow | Orange}) \times P(\text{Sweet | Orange}) \times P(\text{Long | Orange})$$

$$\text{Note : } \{ P(\text{Long | Orange}) = 0 \}$$

$$= 0.53 \times 0.69 \times 0 = 0$$

Rule Based classification

- * It uses set of If-THEN rules for classifications
 - * Example:
If age = "youth" and student = "yes" then buys computer = "yes"
 - * Rule extraction from decision tree
- Example:
- ```

graph TD
 L30[≤ 30] --> Age[age?]
 Age -- No --> Student[student?]
 Age -- Yes --> 340[3 <= 40]
 Student -- No --> Rating[credit rating?]
 Student -- Yes --> 240[2 <= 40]
 Rating -- No --> Fair[fair]
 Rating -- Yes --> 340[3 <= 40]

```

| <u>Assessment of a rule:</u> |                                                         |
|------------------------------|---------------------------------------------------------|
| * Coverage &                 | bit the target function correctly                       |
| * Accuracy &                 | confidence? $\rightarrow$ how much R $\rightarrow$ Rule |
| N-covers =                   | No. of tuples correctly classified by R                 |

\* If more than one rule are triggered need conflict resolution.

size ordering: Assign the highest priority to the triggered value.

class based ordering: Decreasing order of prevalence

(or) misclassification cost per class

Rule based ordering:

Rules are organized into one long priority list, according to some measure of quality (by experts).

## Rule induction

- \* Sequential covering alg/method used
- \* Seq covering  $\rightarrow$  Extracts rules from dataset directly
- \* Seq covr  $\rightarrow$  Mainly based on one of the Evaluation measure
  - i) Accuracy
  - ii) coverage

Rules are learned sequentially  
Rules are learned one at a time.  
i.e sequentially covered every rule  
(one by one)

## Algorithm:

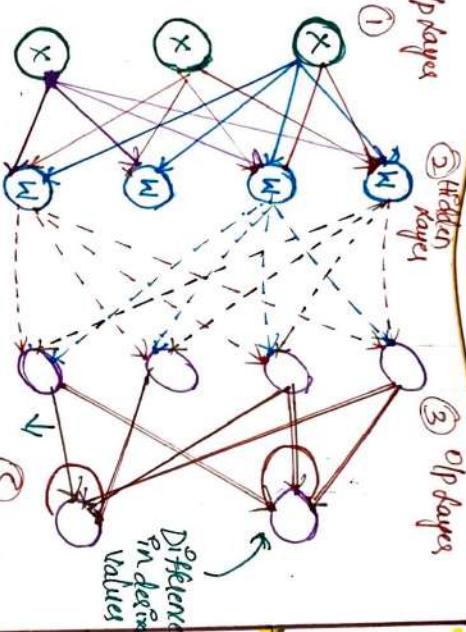
- 1) creates an empty set of decision (rules) list.
  - 2) A function called "learned one-rule" is used
  - 3) If all training  $\in$  class  $y \Rightarrow$  +ve records (Accepted)
  - 4) If all training  $\notin$  class  $y \Rightarrow$  -ve records (Rejected)
  - 5) Get only desirable values (only +ve)
  - 6) Eliminate records (-ve)
  - 7) New rule is added to the bottom of R
- Example:
- 
- Records  $R_1, R_2$  &  $R_3$  located in the list sequentially.

## CLASSIFICATION BY BACK-PROPAGATION

### Back propagation :-

- \* Neurobiological - to develop and list.
- Computation analogous of neuron
- \* Artificial Neural Net (ANN) uses back propagation as a learning algorithm
- “to compute a gradient descent with respect to weight.”
- \* Defined O/Ps are compared to achieved system output & then the system are turned by adjusting connection weights This process to narrow the difference b/w two as much as possible
- \* perception consists of 2 types of nodes:-  
 ↘ I/P-node : Represent I/P attributes  
 ↘ O/P-node : Represent model O/P
- \* each node connected with weight to O/P node
- \* back propagation contains the following layers:-  
 ↘ I/P layer: → Receives I/Ps → X  
 ↘ O/P layer: → Difference in defined values  
 ↘ hidden layers: → calculate the O/P + data to ready @ the O/P layer.  
 The gradient-loss function calculates the difference b/w the O/P + its probable O/P

## I/P layer      ② Hidden layer      ③ O/P layer



### perception model?

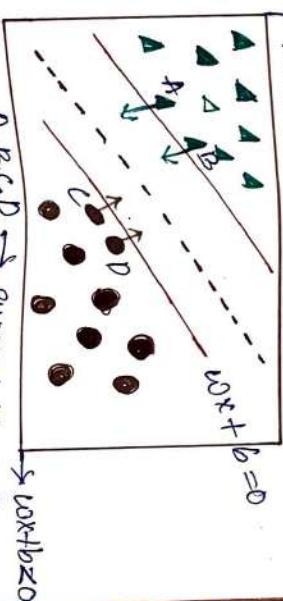
$$Y = \text{sign}(\sum w_i x_i - t)$$

### SUPPORT VECTOR MACHINE

(SVM)

- \* SVM - supervised machine learning algorithm  
→ Used for both classification & prediction

- \* hyperplane - used to separate the data
- \* MMC → Maximum Margin classifier helps to pick the best hyperplane



$$A, B, C, D \rightarrow \text{support vectors}$$

- \* Increase the max-Margin width, if need.
- \* Sometimes, deleting the support-vectors easily change / pick the position of the optimal hyperplane.
- \* Association Classification
- \* Mine data → find strong association b/w frequent pattern
- \* Association Rules:-
 
$$P_1 \wedge P_2 \wedge \dots \wedge P_n \rightarrow "A \text{ class } = c"$$

(confidence support)
- \* Association classification methods more accurate than other method
- \* Mine possible association-rules
- \* Classification based on multiple association rules.
- \* → statistical analysis on multiple ratios
- \* Classification based on predictive association rule
- \* Generation of predictive rules but allow covered rules to retain with reduced weight.



## CLUSTER ANALYSIS / CLUSTERING

\* Clustering is the process of partitioning a set of data into subcluster / subsets

### Application Areas:-

- Image Pattern Recognition, web search, Market - Basket Analysis etc.,

### Basic clustering Methods:-

**Partitioning Methods**

→ A divisive data-objects into non-overlapping "subsets" s.t. each data-object is in exactly one subset.



**Hierarchical Methods**

→ A set of nested clusters organized as a hierarchical tree.



k-Means clustering is an example for partitioning method

chameleon is an example for hierarchical clustering.

## K-Means clustering:-

\* Partitioning clustering approach

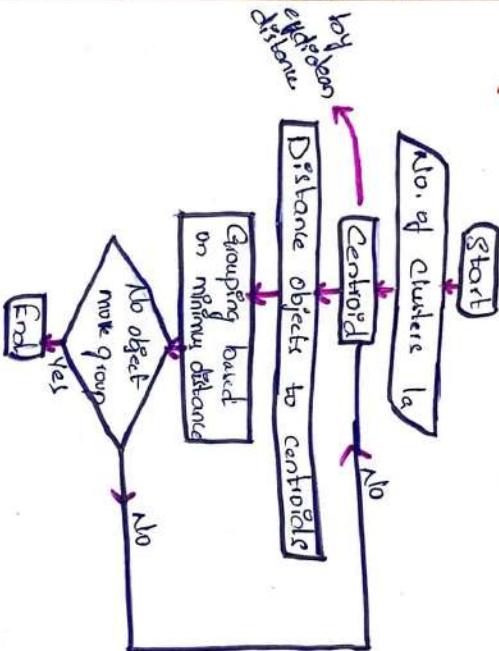
\* Each cluster is associated with a "centroid" i.e. centre point

\* Each point is assigned to the cluster with the "closest centroid".

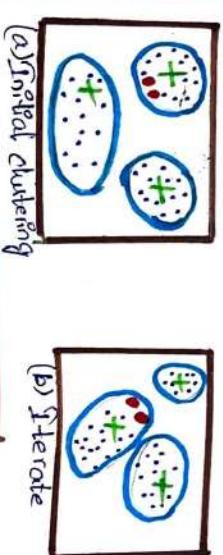
\* closeness is measured by "Euclidean Distance"

\* Initial centroid are often chosen "randomly"

### Algorithm | Flow chart Steps:



### Example:-



(a) Initial clustering  
(b) Iterate  
(c) Final clusters

## K-Medoids

K-Medoids also called "Partitioning Around Medoid"

\* The cost in K-Medoids algorithm is given as

$$C = \sum_{i=1}^k \sum_{j \in P_i} \| p_j - c_i \|$$

### Algorithm:-

1. Initialize: Select k random points out of the n data points as the medoids.

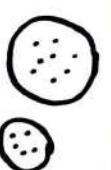
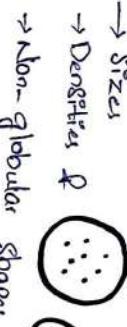
2. Associate each data point to the closest medoid by using any common distance metric methods.

3. While the cost decreases: For each medoid m, for each data O point which is not a medoid:  
→ Swap m and O, associate each data point to the closest medoid and recompute the cost.

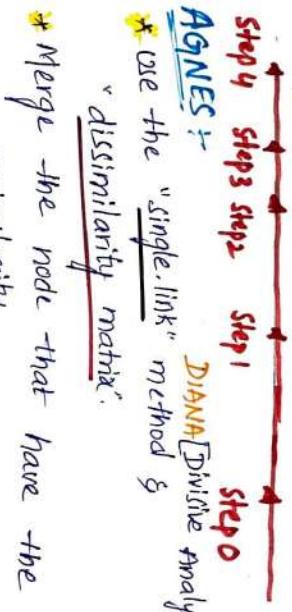
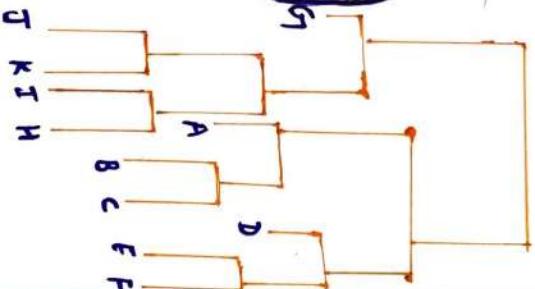
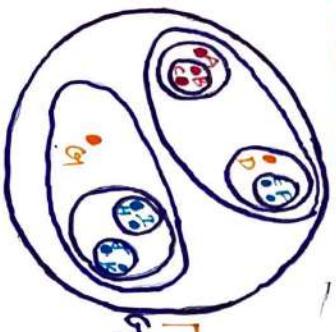
→ If the total cost is more than that in the previous step, undo the swap.

## LIMITATIONS OF K-MEANS:-

\* K-Means has problems when clusters are of differing sizes



## Hierarchical clustering



Ex: AGNES

- Agglomerative nesting  
DIANA

Dissolve Analysis.

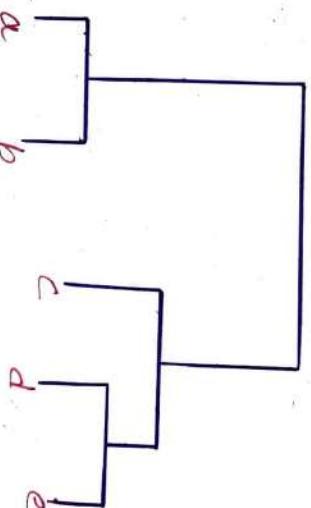
- \* Grouping data into a tree of clusters
- \* It begins by treating every data point as a separate cluster.
- Steps:**
- 1) Identify the clusters which can be closest together.
- 2) Merge the 2 maximum comparable clusters. we need to continue these steps until all the clusters are merged together.

- \* Bottom up Approach
- \* It successively merges the objects (or group close to one another until all the groups.
- \* One another until all the objects are in one cluster, or a termination condition hold

- Agglomerative nesting  
DIANA
- AGNES [Agglomerative nesting]

Step 0 Step 1 Step 2 Step 3 Step 4

Dendogram:



AGNES:

DIANA [Divisive analysis]

- \* use the "single.link" method & dissimilarity matrix.

- \* Merge the node that have the least dissimilarity

- \* Go on in a non-descending fashion

- \* Eventually all nodes belongs to same cluster.

## AGGLOMERATIVE vs DIVISIVE

### DIANA:



- \* Top down approach
- \* In each successive iteration a cluster is split into smaller clusters.
- Steps:**
- 1) Initially all points to the dataset
- 2) partition the cluster into two least similar clusters
- 3) proceed recursively to form new cluster until the desired number of cluster is obtain.

- \* Inverse order of AGNES
- \* Eventually each nodes forms a cluster on its own.

- \* A tree structure called a "dendrogram" is commonly used to represent the process of hierarchical clustering
- \* Either agglomerative/divisive can be used.

- \* A tree structure called a "dendrogram" is commonly used to represent the process of hierarchical clustering
- \* Either agglomerative/divisive can be used.

## FREQUENT PATTERN (FP) GROWTH

### ALGORITHM

FP-Growth Algorithm in Data Mining

→ FP-Growth is an improved version of the Apriori algorithm which is used for frequent-pattern mining (also known as Association Rule Mining).

→ Association Rule Mining can be viewed as a 2-step process:

1) Find all frequent item sets

→ Apriori Algorithm

2) Generate Strong association rules from the frequent itemsets.

→ These rules must satisfy the following:

• Minimum support &  
• Maximum confidence

FP-Tree :

→ FP-growth alg. represents the database in the form of a tree called "FP-Tree".

→ Purpose: To mine the frequent pattern

→ The root node represents null, while the lower-node repres. itemsets.

→ It is a compact data structure that stores quantitative info. about freq-patterns in a database.

Algorithm / Pseudocode :-

Procedure FP-Growth \* ( $T$ )

Input: A conditional FP-tree  $T$

Output: The complete set of all FPs corresponding to  $T$ .

Method:

1. If  $T$  only contains a single branch  $B$
2. For each subset  $y$  of the set of items in  $B$
3. Output itemset  $y$   $\cup$   $T$ .base with count = smallest count of nodes in  $y$ ,
4. else
  - for each  $i$  in  $T$ .header do begin
  5. Output  $y = T$ .base  $\cup$   $\{i\}$  with  $i$ .count;
  6. if  $T$ .FP-array is defined
    - Construct a new header table for  $y$ 's - FP-tree from  $T$ .FP-array
  7. Construct  $y$ 's - FP-tree from  $T$ .FP-array
  8. else
    - construct a new header-table from  $T$ .
  9. Construct  $y$ 's conditional FP-tree  $T_y$  & Possibly if FP-array  $A_y$ ;
  10. If  $T_y \neq \emptyset$
  11. call FP-growth \* ( $T_y$ ) ;
  12. End.

## APRIORI ALGORITHM

(18)

Steps in Apriori :-

Step-1 : Determine the support of itemsets in the transactional database & select the "minimum support" as "confidence".

Step-2 : Take all supports in the transaction with higher support value than the minimum / selected support value.

Step-3 : Find all the rules of these subsets that have higher confidence value than the threshold [ $\sigma$ ] min-confidence.

Step-4 : Sort the rules as the decreasing order to fit.

Flow chart :-

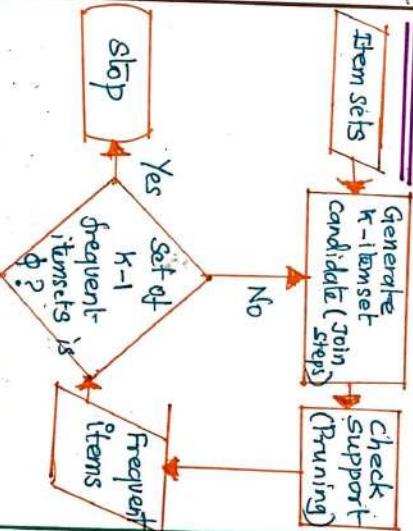
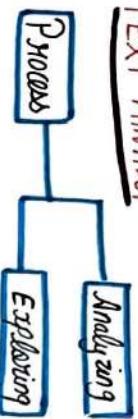


Fig. Flowchart of Apriori

## TEXT MINING



Finding / Outputs  $\Rightarrow$  concepts, pattern, topic, etc....

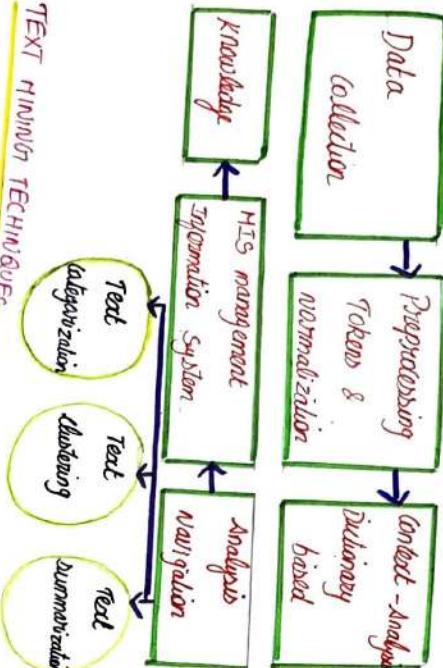
### METHODS:-

#### STEP 1: Gathering unstructured information

#### STEP 2: Preprocessing - Cleaning

#### STEP 3: Analysis - Analyze the patterns

Clustering & summarization



### KEY WORD BASED ASSOCIATION

collected set of keywords / items that one

occurred frequently

removing the stop words  $\Rightarrow$  by passing, stemming,

$\Rightarrow$  implement association mining algorithm

$\Rightarrow$  view a set of keywords in document as a set of item

in the transaction

### AUTOMATIC DOCUMENT CLASSIFICATION

$\Rightarrow$  A.K.A. categorization

$\Rightarrow$  process of managing text & unstructured information

by categorization (or) clustering text

$\Rightarrow$  Enables users to organize content quickly

### APPLICATIONS

$\Rightarrow$  Find useful information from the

user interactions

$\Rightarrow$  understanding consumers behaviour, need & buying patterns

### TEXT MINING TECHNIQUES

\* Information extraction

\* Information retrieval

\* Natural language processing

\* Text clustering analysis

\* Text summarization

## APPLICATION AREAS

\* Digital library

\* Academic & research field

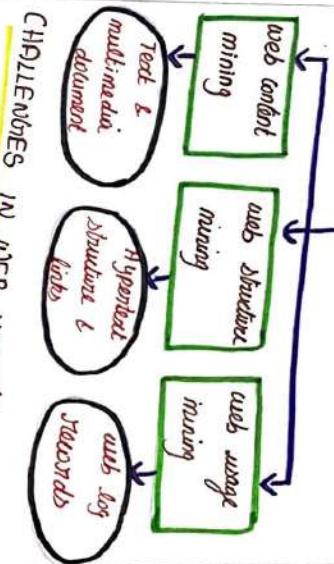
\* Life science

\* Social media

\* Business intelligence

## WEB MINING

\* Data acquired through "web crawler"  
[or] "web analysis".



### CHALLENGES IN WEB MINING

\* The complexity of web pages

\* Diversity of web data source

\* Relevancy of data

\* The web is too broad

### ADVANTAGES

$\Rightarrow$  A.K.A. categorization

$\Rightarrow$  process of managing text & unstructured information

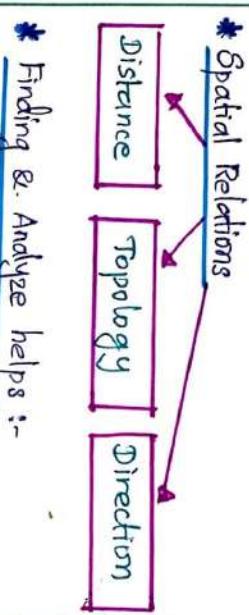
by categorization (or) clustering text

$\Rightarrow$  Enables users to organize content quickly

## MINING SPATIAL DATA

- Geographical [or] spatial information to produce business intelligence [or] other results.

- Data related to spatial description of the objects such as



- \* Finding & Analyze helps :-

→ Earthquake points  
→ Climate / Weather predictions

→ Google Map  
GPS (Global Positioning System)

→ Trend Analysis

- \* Clustering Analysis Methods :-

→ PAM → Partitioning Around Medoids

[*M*arla's K-Means clustering]  
→ PAM divides data into groups based on medoids.

Methods [Others] :-

→ Spatial Auto Correlation - Measure of dependency among points in a spatial neighborhood

## Partitioning Around Medoids (PAM)

20

(b) Spatial Heterogeneity - variation in events, features & relationships across a region.

- \* Geographic Warehouse :-

- Built to collect data from various resources.

- \* Spatial Prediction :-

- Identify the relationship between variables in different datasets.

### Types of Spatial Data :

#### 1) Feature Data :

- Follow the vector data model

- Represents the entity of the real world. i.e., roads, trees,

buildings etc.,

This information can be visually represented in the form of a point, line (or) polygon.

#### 2) Coverage Data :

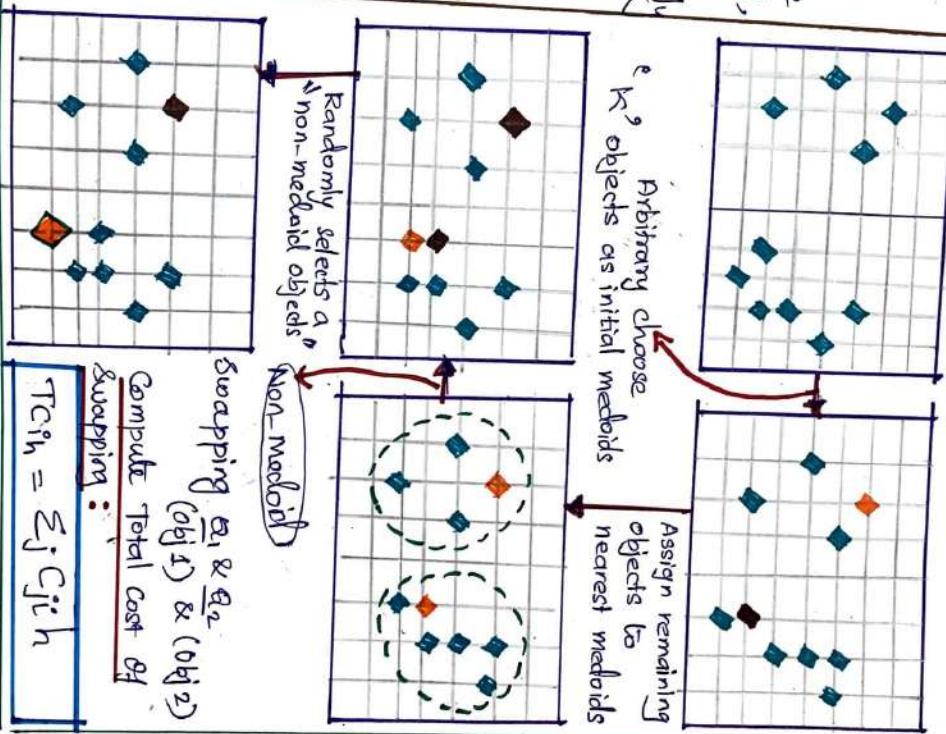
- Follows the raster data model

- Data contains the mapping of continuous data in space.

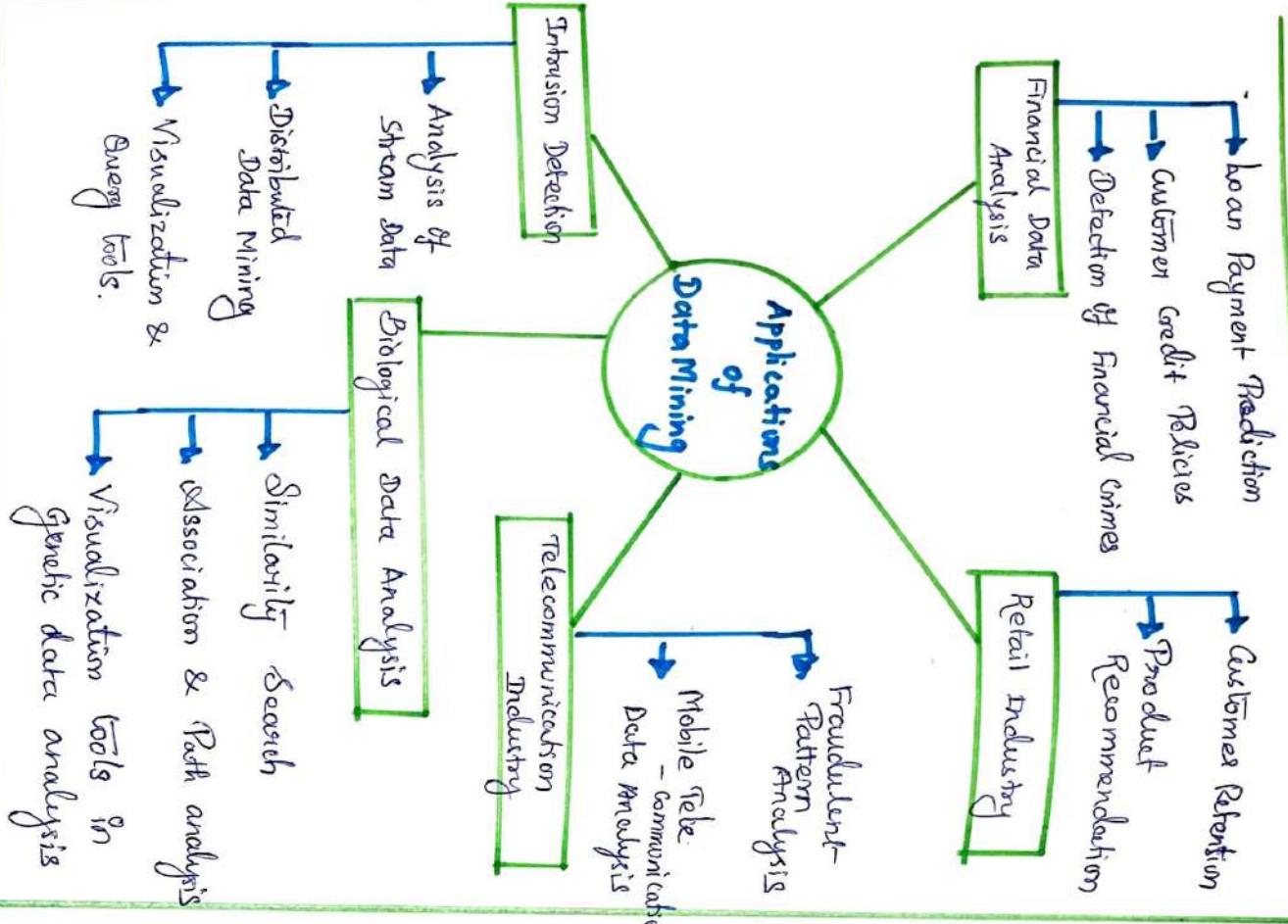
→ Represented as a satellite image, digital surface model etc.,

The visual representation of coverage data is in the form of a "Grid" [or]

triangulated irregular network



## APPLICATIONS IN DATA MINING



## TRENDS IN DATA MINING

