



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

7. Architecture of Data Mining systems.....	7
7.1. KDD	
7.2. Classification of Data Mining Systems	
8. Steps in Preprocessing.....	8
8.1. Cleaning	
8.2. Integration	
9. Data Mining Processing	9
9.1. Reduction	
9.2. Transformation	
9.3. Discretization	

UNIT 3

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 and Bootstrap	
15.2. Metrics for performance evaluation	

UNIT 4

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 5

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

Data Mining :-

→ Tools to discover knowledge from data

→ Knowledge from data

→ Data platform analysis

→ Knowledge extraction

Kinds of data :-

Database data

* Context of data

* Positively stored in database table

Transactional data

* Organized by time stamps
will back.

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

Data warehouse

kinds of data

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

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

* Used in query language
query processing

* Data mining tasks and data warehousing techniques

Machine learning

Information Retrieval

Database system

Statistics

Data Mining

Algorithm

* performs the analysis and help

* visualize data

pattern
Recognition

Visualization

Data housing

* Makes it easier to identify

* automated recognition patterns and outliers in large database

* Central repositories

* Core components

* Business intelligence

Major Issues:

* Enhance the power and flexibility of datamining

* Handling uncertainty, noise
(or) incompleteness of data

* User interaction

* Efficiency and Scalability

* Social Impact of Datamining

* Privacy of datamining

Business Intelligence:-

* Historical

* OLAP predictive analysis.

Web Search Engines:-

* Search of information on web

* Web directories maintained by human editors.

(2)

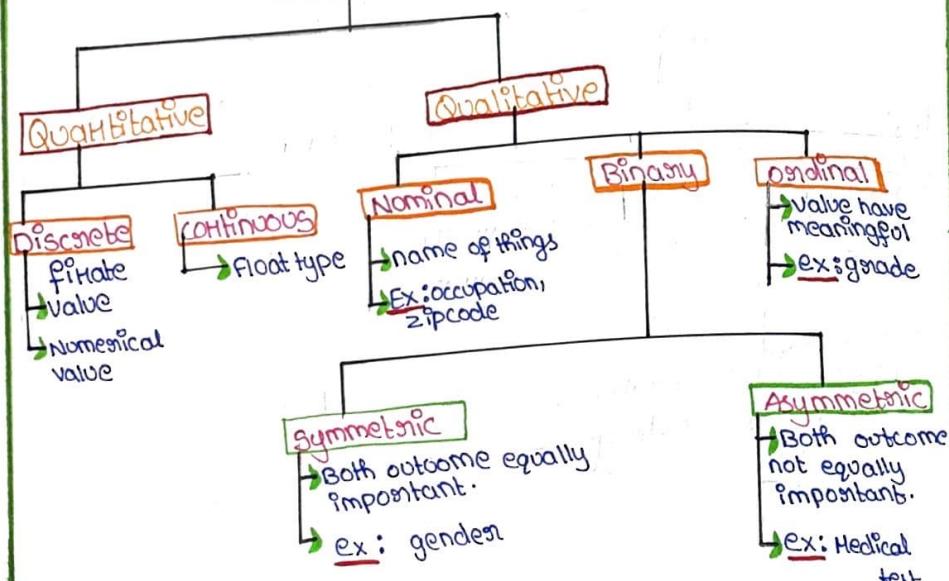
GETTING TO KNOW YOUR DATA

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

Attributes



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.

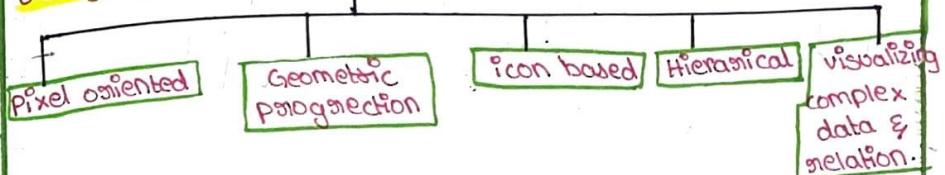
Dispersion analysis on computed measures -

- * 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,∞]

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):

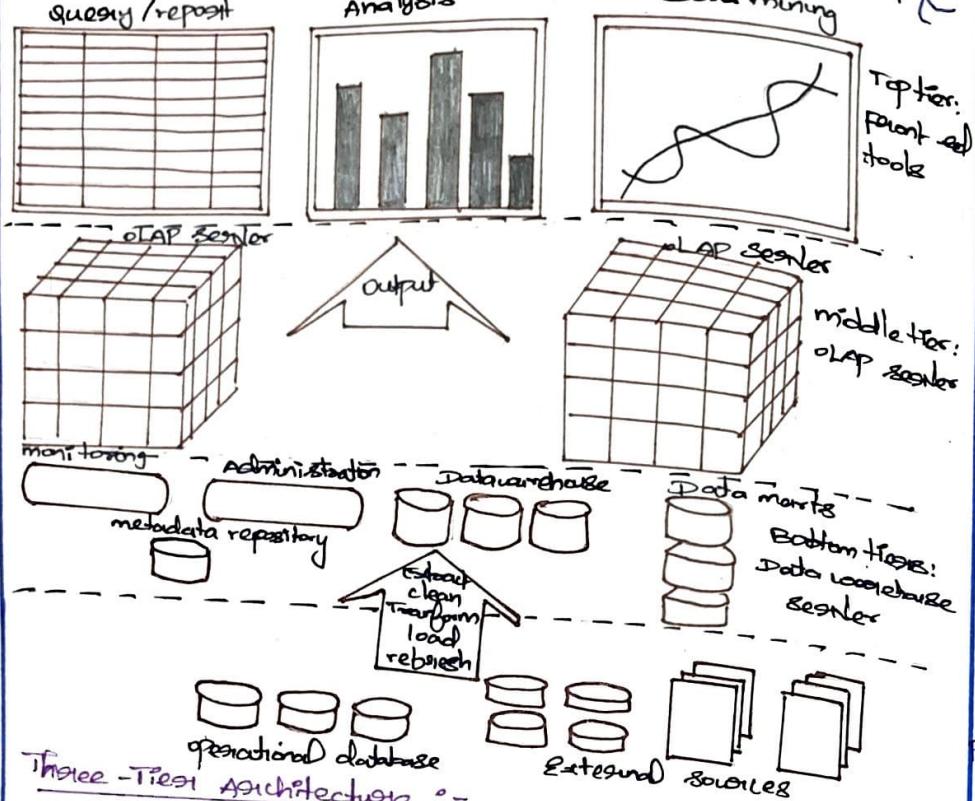
- * It is computerized program used to support determinations, judgements and course of actions in organization.

Ex:

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

Characteristic	Operation data	DSS
Data currency, granularity	Real time data	Historical data
	Atomic-detail 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



Three-Tier Architecture :-

1. Bottom-Tier :

→ DB - Server.

2. Back-end Tools:

→ Extract info/... load & clean.

3. Middle-Tier :

→ Extention of RDBMS

4. Top-Tier:

→ client layer (e.g., query tools, reporting tools, mining tools etc.)

Types of Data warehouse:-

Enterprise DW

operational data store

Data mart

↓
centralised
DSS

↓
OLTP support

↓
Subset of DW
Eg: Sales, Finance

Extraction, Transformation & Loading (ETL) (4)

During data extraction, raw data is copied from unstructured sources of data include but are not limited to:

- SQL or NoSQL databases → CRM & ERP systems → Flat files.
- Email → webpages.

Transform:

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

- Filtering, cleansing, de-duplicating, validating, & authenticating the data.
- conducting audits to ensure data quality and compliances/government regulations.
- Remaking, encrypting, or protecting data governed by industry

/formatting the data into tables or joined tables to match the schema of the target data warehouse.

Load:

The transformed data is moved from the staging area into a target data warehouse. Initial loading of all data, followed by periodic loading of incremental data changes & less often, full refreshes to ensure & replace data in the warehouse.
→ Comparison of OLTP & OLAP systems:

Feature	OLTP	OLAP
• characteristic	operational processing	informational processing
• orientation	transactions	analysis
• user	clerk, DBA, database processor	knowledge worker (e.g., manager, executive, analyst)
• Function	day-to-day operations	long-term informational requirements decision support
• DB design	ER-based, applications-oriented	star/snowflake, Sub-oriented
• Data	consistent, guaranteed up-to-date	historical, accuracy maintained
• Summarization	partitive, highly detailed	over time
• View	detailed, flat relational	Summarized, long/dated
• unit of work	short, simple transaction	summarized, multidimensional
• Access	read/write	complex query
• Focus	data in index/hash on primary key	mostly read
• Operations	transaction	info. out
• no. of records accessed	lots of scans	millions/hundreds
• no. of users	tens	↑TB
• DB size	GB to high-order GB	high performance, high availability
• Priority	high priority	highly flexible, end-user autonomy
• latency	transaction throughput	query throughput, response time

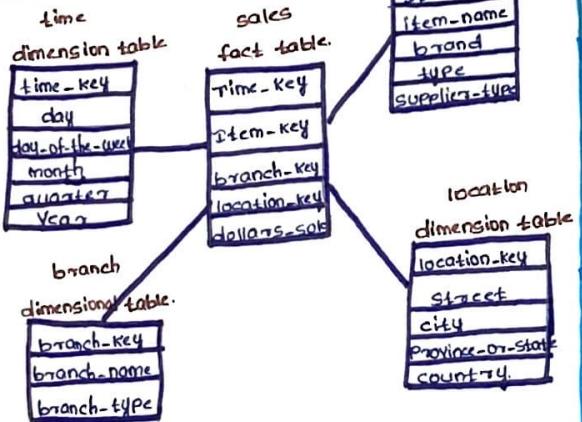
Dimensional Data Modelling

(5)

Chicago				
New York				
Toronto	818	858	16	707
Vancouver				
Q1	605	825	10	400
Q2	680	952	12	512
Q3	812	1023	31	573
Q4	927	700	21	545
H	C	P	S	

star, snowflake and fact constellation :-
* schemas for multidimensional data model

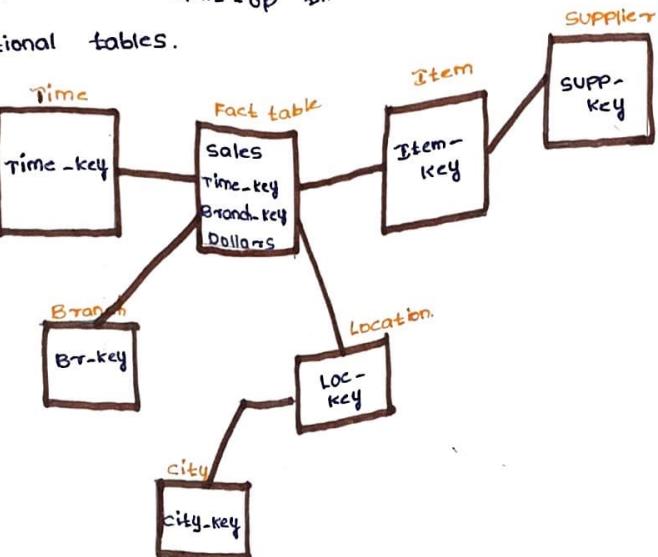
Star - schema :-



- Bulk of data with no-redundancy.
- A set of smaller attendant table [1-for each dimension]
- * It contains both dimensional table and fact table.

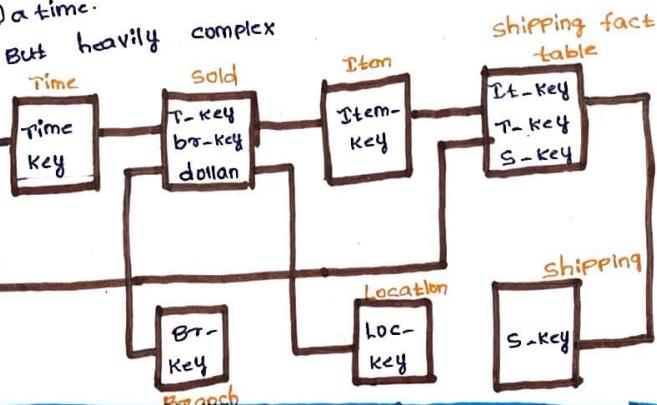
Snow-flake schema in [top-down model type]

- * Some of them are normalized
- * Normalization split-up the data into additional tables.



Fact constellation schema :-

- * It consists of dim-table that are shared by several fact tables
- * It consists of more than one star-schema @ a time.
- * But heavily complex



→ 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 ,quarter ,year)

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

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

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

snowflake Schema definition :-

define cube sales -snowflake [time, item, branch, location]:

dollars -sold = sum (sales -in-dollars), unit -Sold = count (*)

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

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

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

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

Fact constellation schema definition :-

define cube sales [item, time, 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 ,item-name ,brand ,type)

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

define dimension location as (loc-key ,street ,city ,country)

define cube shipping [time, item, shipper, from-location]:

dollars -cost = sum (cost -in-dollars), = count (*)

define dimension from-location as location in cube sales

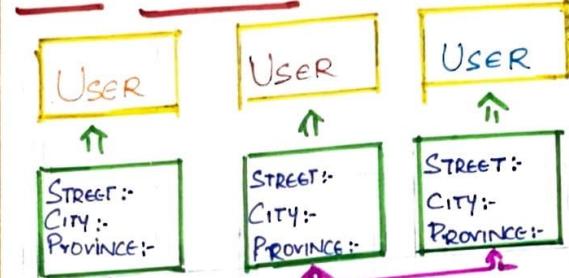
define dimension to-location as location in cube sales.

Online Analytical Processing

OLAP - Multidimensional

Information - Manage, Analysis, Interactive

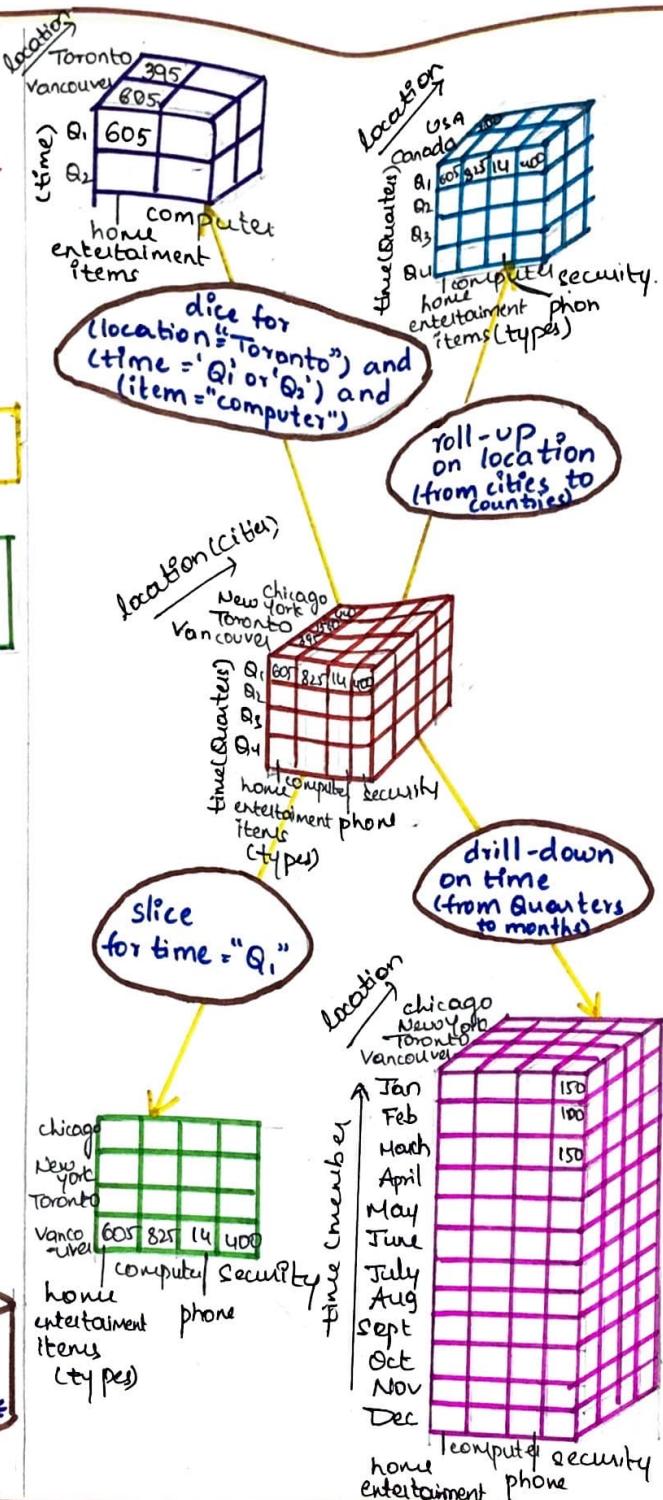
OLAP PREPARATION



App

OLAP CUBE

DATA WAREHOUSE



OLAP - OPERATIONS

ROLL UP : Aggregation on data

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

(2) Dimension Reduction

"street < city < province < country"

Aggregated hierarchy location from "city" to "country"

ROLL-DOWN : stepping down a hierarchy

* New dimension

Work :- hierarchy for dimension time

"day < Month < Quarter < year"

SLICE :- Create a new sub-cube from one-particular cube

way :- Dimensions = "time", - time = "Q1"

DICE :- Three dimension

Location = "Toronto" or "Vancouver"

Time = "Q1" or "Q2"

Item = "Mobile" or "TV"

PIVOT :- PIVOT-Rotation Operation

TOTAL NO OF CUBOIDS :-

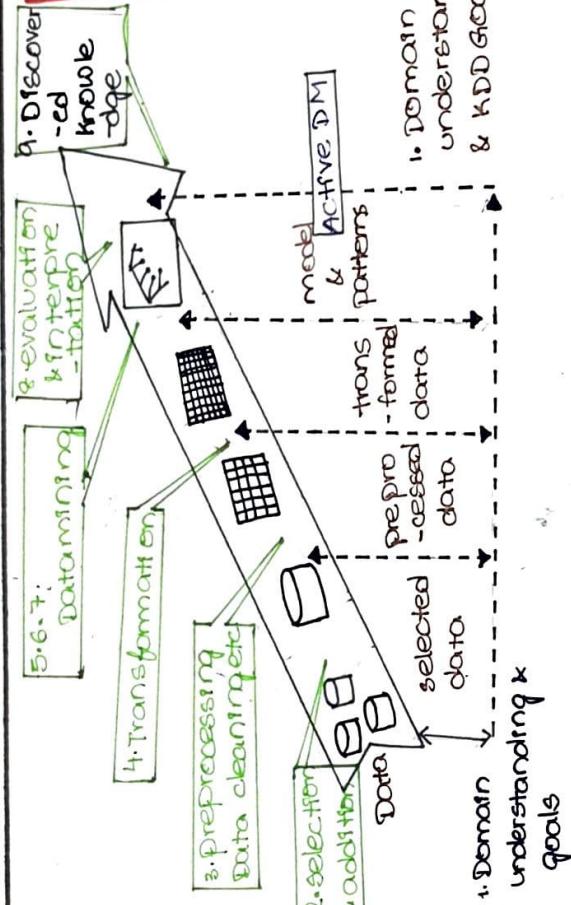
$$\prod_{i=1}^n (L_i + 1)$$

L_i = Levels associated with dimension

1 = Virtual top level

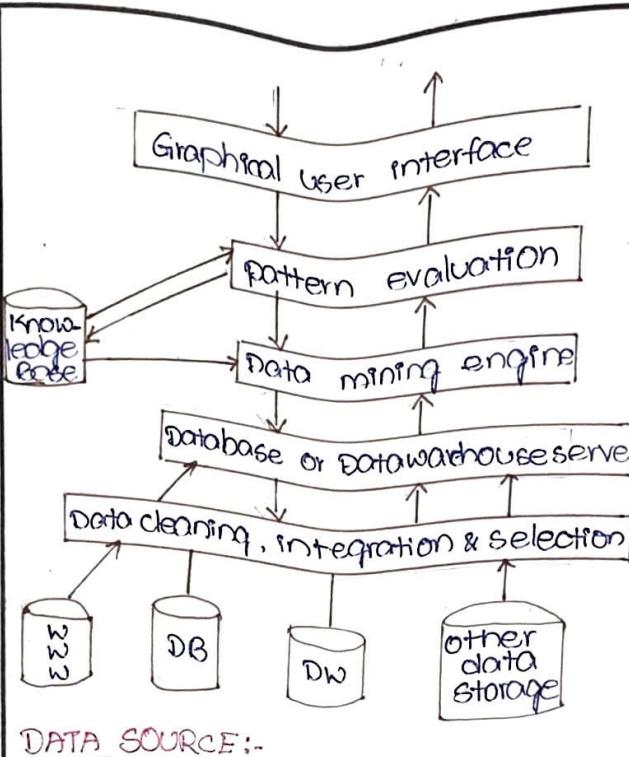
DATA CUBE COMPUTATION

KNOWLEDGE DISCOVERY IN DATABASE (KDD):-



PROCESSING:-

1. Data cleaning
2. Data integral
3. Data reduction
4. Data transfer
5. Data discretization
6. Data mining
7. pattern evaluation
8. Representing knowledge.



DATA SOURCE:-

Database, www, text files and other documents.

organization store data in datawarehouse different process:-

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

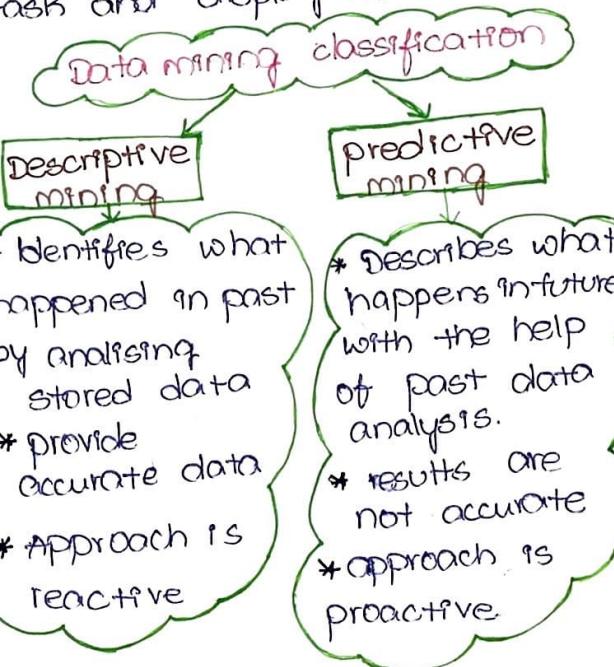
Data mining engine: contains modules like association, characterization, classification, clustering, prediction, time-series analysis.

pattern evaluation engine:-

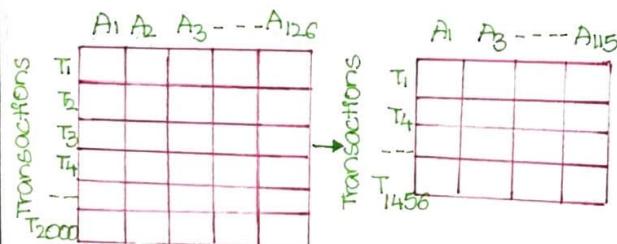
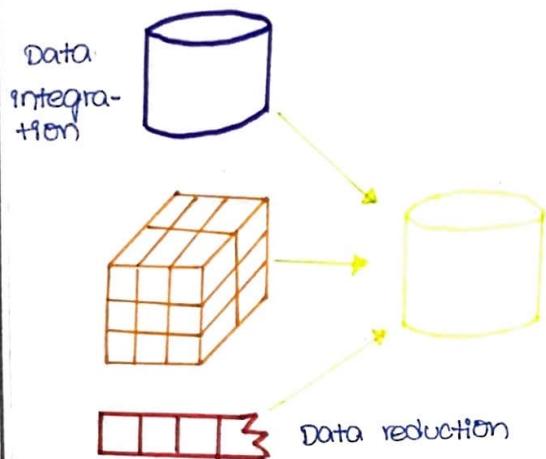
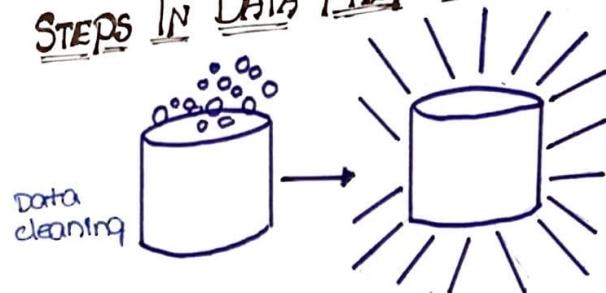
- * Measure of investigation of the pattern using a threshold value.
- * collaborates with data mining engine to focus the search on exciting patterns.

Graphical user interface:-

- * communicates between the data mining system and the user.
- * This module cooperates with the data mining system.
- * user specifies a query or a task and displays the results.



STEPS IN DATA PREPROCESSING:



Data transformation: $-2, 32, 100, 59, 48 \rightarrow -0.02, 0.32, 1.00, 0.59, 0.49$

Preprocessing:-

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

* 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 - lacking attribute values
eq: age = "NOISY"

* contain noise or error salary = "10"

INCONSISTENT: contain discrepancies, age = 20

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:-

* 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

* use commercial tools :-

data scrubbing or data auditing.

DATA MIGRATION & INTEGRATION:

* Allows transformation to be specified

ETL:- allows users to specify transformations through qui.

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.

* object identification

* Derivable data.

* Redundant attributes detect using correlation analysis and covariance analysis.

Correlation analysis:

χ^2 (chi-square) test.

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

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

covariance: $\text{cov}(A, B) = E((A - \bar{A})(B - \bar{B})) =$

$$\frac{1}{n} \sum_{i=1}^n (a_i - \bar{A})(b_i - \bar{B})$$

$$\sigma_{A,B} = \sqrt{\text{cov}(A, B)} / \sigma_A \sigma_B$$

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

negative covariance:- if $\text{cov}_A, B < 0$ A is larger than its expected value, B is small than expected value.

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

Data Pre-Processing : Data Reduction & Transformation... (9)

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

* Produces the same analytical result.

Data Reduction Strategies

Dimensionality reduction

Wavelet Transformation :- 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 vectors defines space.

Numerosity reduction :- Regression and Log-linear

Linear :- Data modelled to fit a straight line.

Multiple regression :- Allow response to variable x to be modeled as linear function.

Histogram analysis :- Divide data into buckets and store average (sum) of each bucket per partition rule. 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 smaller

min-max normalization :-

$$v' = \frac{v - v_{\min A}}{v_{\max A} - v_{\min A}} \cdot (new_max_A - new_min_A) + new_min_A$$

Z Score - Normalization :-

$$v' = \frac{v - \bar{x}_A}{\sigma_A}$$

\bar{x} = Mean σ = Standard deviation

Normalization by decimal scaling :-

$$v' = \frac{v - \bar{v}}{\sigma} \quad \because \bar{v} = \text{smallest integer such that } \max(v') \leq 1$$

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

* Interval tables 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.

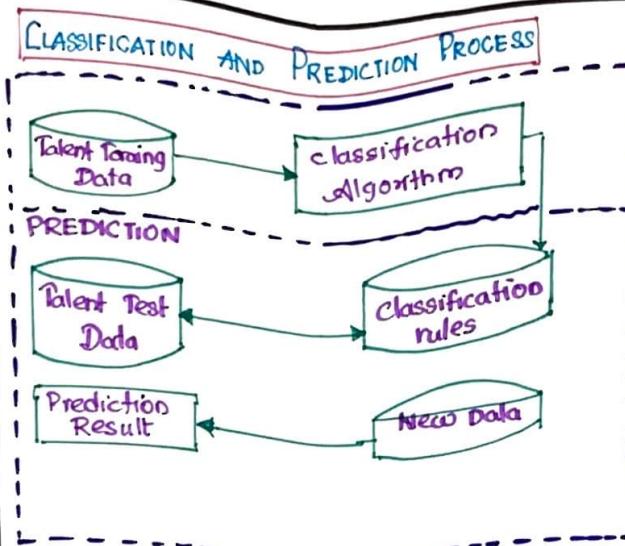
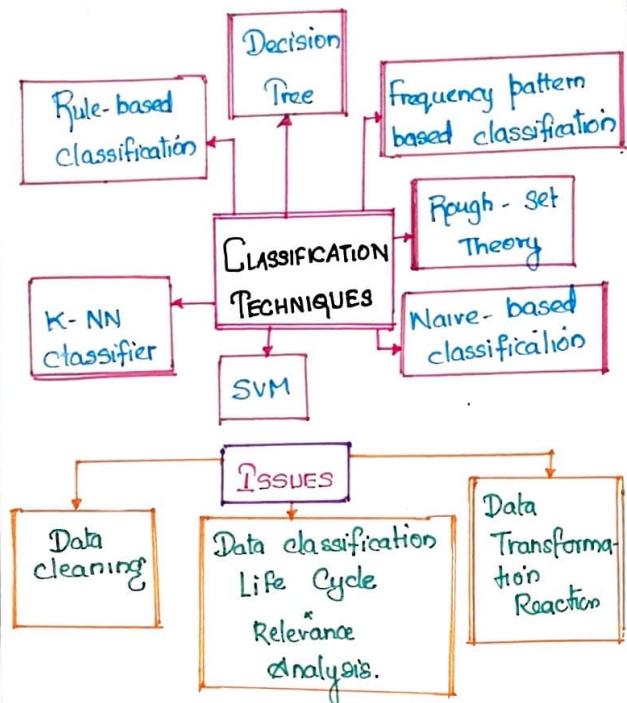
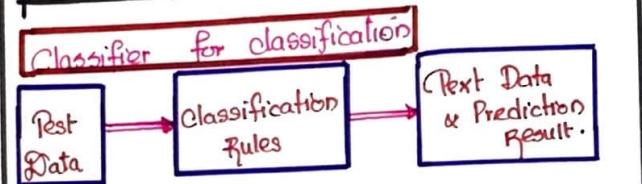
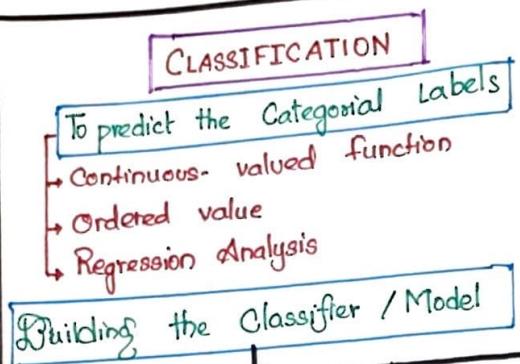
⇒ A & B are the lowest & highest values of attribute wide $w = (B - A)/n$.

Equal-depth (frequency) Partitioning :-

⇒ Divide the range into n intervals, each containing approximately same number of samples.

⇒ Good data scaling :-

concept hierarchy organizes concepts hierarchically and is usually associated with each dimension in a data warehouse.



COMPARISON OF CLASSIFICATION & PREDICTION

- SPEED → Value of Attributes
- ROBUSTNESS → Noisy Data
- SCALABILITY → Class prediction Label
- INTERPRETABILITY → Efficiency

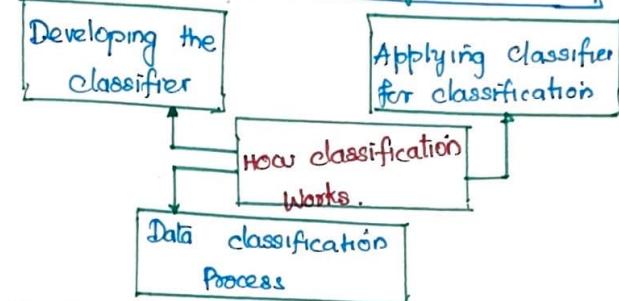
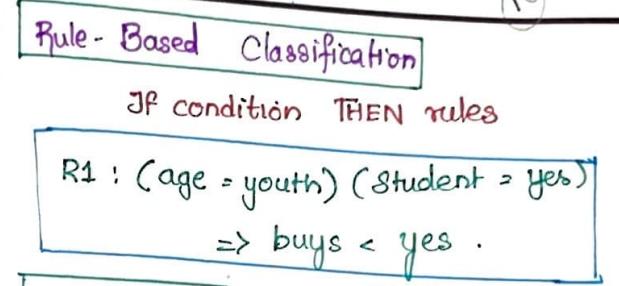
$$\text{Gain Ratio } (A) = \frac{\text{Gain}(A)}{\text{split Info}_A}$$

Decision Tree Induction

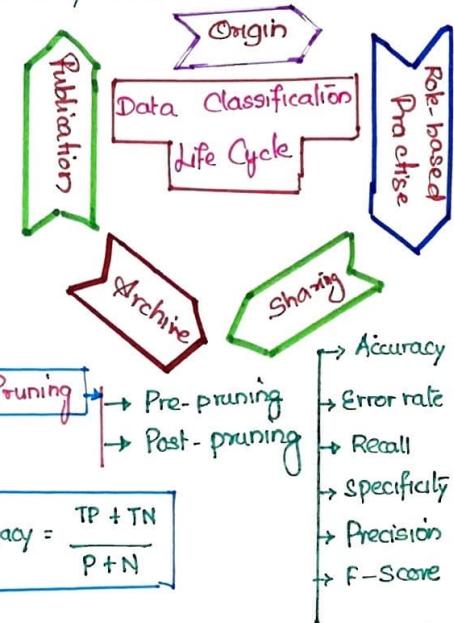
- Multi-dimensional data
- Domain-knowledge
- Top-down recursive
- Divide-conquer Manner

Bayes' Theorem

$$P(H/x) = P(x/H) P(H)$$



- Essential Concept of Measurement**
- Measurement of Internet Traffic
 - Understanding Measurement
 - Per-Packet and Per-flow Measurement
 - Evaluation Criteria for Perf. Measurement
 - Accuracy Recall.



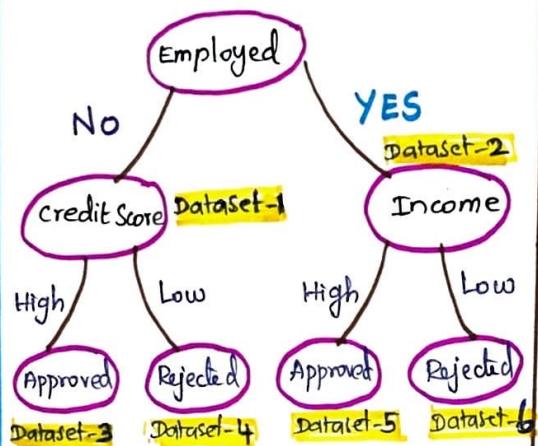
DECISION TREE

* 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)
2. GINI Index

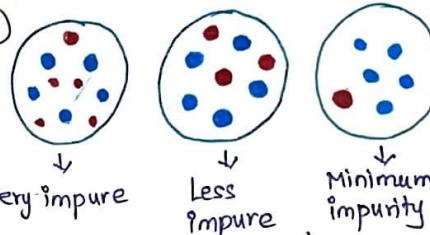
1. Information Gain :- (IG):

→ IG is calculated for a split by subtracting the weighted "entropies" of each branch from the "original entropy".

→ i.e measure the change in entropy

→ Entropy : → Measures the impurity [or] uncertainty in a group of observations.

→ e.g)



→ Formula :

$$E = - \sum_{i=1}^n p_i \log_2 p_i$$

Example:

$$\text{Set } X = \{a, a, a, b, b, b, b, b\}$$

Total instances = 8

Instance b = 5

Instance a = 3

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

(ii) GINI Index :-

→ Gini Index & entropy are the criterion for calculating information gain (IG).

→ D.T. algorithm use IG to split a node. Both GiniIndex & Entropy are the measures of impurity of a node.

→ Formula

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

Where,

p_i → The probability [or] percentage of class C_i in a node.

$$\text{IG} = (\text{Entropy of Parent node}) - (\text{Sum of weighted entropies of child node})$$

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 independence 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 $P(B) \neq 0$

* $P(A|B)$ → Conditional Probability

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

* $P(B|A)$ → Conditional Probability → Likelihood

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

* $P(A)$ & $P(B)$ → Probabilities of observing A & B independently of each other. ↳ called Marginal Probability

Example :

Fruit = {Yellow, Sweet, Long}

FRUIT	Yellow	Sweet	Long	TOTAL
Orange	350	450	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.

① Probability (Yellow | orange) = $\frac{P(\text{orange}|\text{yellow}) \cdot P(\text{yellow})}{P(\text{orange})}$

$$= \frac{\frac{350}{800} \times \frac{800}{1200}}{\frac{650}{1200}} = \frac{0.4375 \times 0.667}{0.541667} = 0.5387 //$$

∴ Probability of yellow/orange = 0.53

② Prob. (Sweet | orange) = $\frac{P(\text{orange}|\text{sweet}) \cdot P(\text{sweet})}{P(\text{orange})}$

$$= \frac{\frac{450}{850} \times \frac{850}{1200}}{\frac{650}{1200}} = \frac{0.529 \times 0.7083}{0.541667} = 0.6917 //$$

∴ Prob. (Sweet | orange) = 0.69

11
 $P(\text{Fruit}|\text{orange}) = P(\text{Yellow}|\text{orange}) \times P(\text{Sweet}|\text{orange}) \times P(\text{Long}|\text{orange})$

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

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

③ $P(\text{Fruit}|\text{Banana}) = \frac{P(\text{Yellow}) \cdot P(\text{Sweet})}{P(\text{Banana})} \times P(\text{Long})$

$$= 1 \times 0.75 \times 0.89 = 0.65$$

④ $P(\text{Fruit}|\text{others}) = P(\text{Yellow}) \times P(\text{others}) \times P(\text{Sweet}) \times P(\text{others}) \times P(\text{Long}) \times P(\text{others})$

$$= 0.33 \times 0.66 \times 0.33 = 0.072$$

From the above,

→ Fruit of orange is zero,
 \therefore eliminated

→ Fruit of Banana & } Which
 Fruit of other fruits } is
 the most probability
 of
 Yellow, Sweet & Long.
 ?

Conclusion :-

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

Rule Based classification

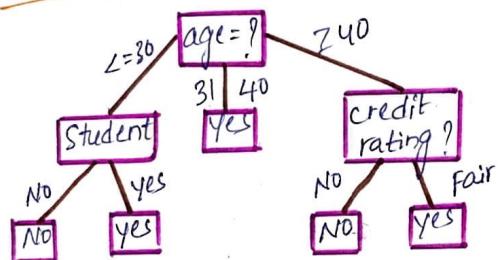
- * It uses set of IF-THEN rules for classifications

Example:

If age = "youth" and student = "yes"
then buys computer = "yes"

- * Rule Extraction from decision tree

Example:



- * chances to buy the computer
- person with age ≤ 30 & student
(more chances to buy a computer)
- person with age ≥ 40 & credit score is fair
(more chances to buy a computer)
- person with age ≤ 30 & not a student
(less chances to buy a computer)
- person with age ≥ 40 & low credit score
(less chances to buy a computer)

- * Coverage & correctly fit the target function
- * Accuracy [how much confidence?] $R \rightarrow$ Rule

$N_{\text{covers}} = \text{No. of tuples covered by } R$

$N_{\text{Correct}} = \frac{\text{No. of tuples correctly classified by } R}{N_{\text{covers}}}$

- * 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

Rule based ordering per class

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

Rule induction

- * Sequential covering alg/method used

* Seq covering → Extracts rules from dataset directly

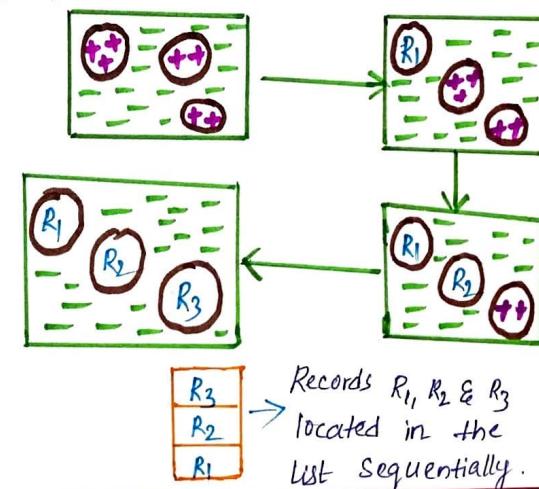
* Seq covl → Mainly based on one of the Evaluation measure
(i) Accuracy (ii) coverage

- * 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
 - ↳ It extracts best rule for y class
 - If all training \in class $y \Rightarrow +ve$ (accepted)
 - If all training \notin class $y \Rightarrow -ve$ (Rejected)
- 3) Get only desirable values (only +ve)
- 4) Eliminate records (-ve)
- 5) New rule is added to the bottom of R

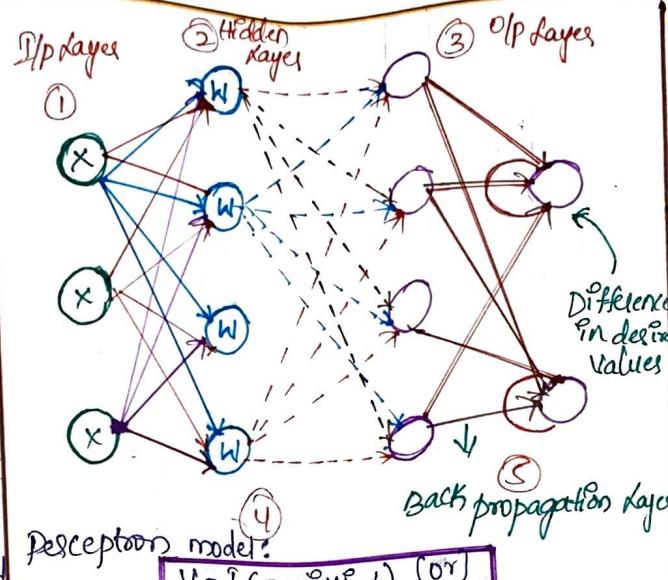
Example:



CLASSIFICATION BY BACK-PROPAGATION

Back propagation :-

- * Neurobiological to develop and list - Computation analogous of neuron
- * Artificial Neural NW(ANN) uses back-propagation as a learning algorithm
 - "to Compute a gradient descent with respect to weights."
- * Desired o/p's are compared to achieved system outputs & then the system are turned by adjusting connections 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/p's → x
 - O/p layer: → Difference in desired values
 - Hidden layer: → calculates the O/p & data is ready @ the O/p layer.
 - The gradient-loss-function Calculates the difference b/w the D/w-O/p & its probable O/p



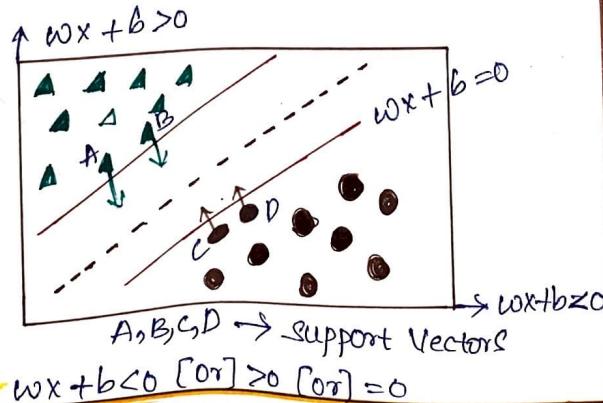
Perception model?

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

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

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



- (14) 13
- * Increase the max-Margin width, the points support vectors found well, if need.

- * Sometimes, deleting the support-vectors easily change/pick the position of the optimal hyperplane.

Association classification

- * Mine data to 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 values

- * Classification based on predictive association rule

- * Generation of predictive rules but allow covered rules to retain with reduced weight.

EVALUATING ACCURACY OF A CLASSIFIER - CROSS VALIDATION, BOOTSTRAP.

Measuring Accuracy:

- * Model evaluation is the process through which we quantify of a system's predictions.
- * The performance & usage of the classifier's model is decided in terms of accuracy measures at the end.
- * Measuring accuracy using validation set of class-labeled tuples.

confusion matrix:

True Positives (TP)	False Negatives (FN)
False Positives (FP)	True Negatives (TN)

- * **TP** - predicted +ve & its true
- * **TN** - predicted -ve & its true
- * **FP** - predicted +ve & its false
- * **FN** - predicted -ve & its false.

A confusion matrix - $N \times N$ matrix, $N \rightarrow$ No. of class being predicted.

classifier Accuracy:

- * % of test-set tuples that are correctly classified.

Measures:

Accuracy :-

The proportion of the total no. of predictions that were correct.

$$Acc = \frac{TP + TN}{TP + FP + TN + FN}$$

Error Rate:

$$ER = 1 - Accuracy$$

Sensitivity :- True +ve recognition rate [or] Recall [or] True +ve rate.

$$Recall = \frac{TP}{TP + FN}$$

specificity:

True -ve recognition rate

$$specificity = \frac{TN}{TN + FP}$$

- * A.k.a as % of -ve instances out of the total actual -ve instances.

classifier Evaluation matrix:

Precision : % of tuples classified as +ve were actually +ve.

$$Precision = \frac{TP}{TP + FP}$$

Recall:

what % positive tuples did the classifier label as positive?

$$Recall = \frac{TP}{TP + FN}$$

F1-Measure / F-Score:-

Harmonic mean of precision & recall.

Highest F1-Score is the better value.

$$F = \frac{2 \times Precision \times Recall}{Precision + Recall}$$

Evaluating classifier Accuracy:

Holdout: → Reserves 1/2 for training & 1/2 for testing

→ Reserves 2/3 for training & 1/3 for testing.

satisfied Sampling :- Each class is represented with approximation in training a test.

Cross Validation:

Step 1: Data split into 'k' subsets of equals.

Step 2: Each subset is trained used for testing & remainder for training subset satisfied before cross-validation.

Evaluating classifier accuracy:

Bootstrap:

- Works well with small data sets.
- Samples the often training tuples uniformly with replacement
- Each time a tuple is selected to be selected again & re-added to the training set.

b32 bootstrap:

- A data set with 'd' tuples is sampled 'd' times, with replacement resulting a training set of 'd' samples
- The data tuples that did not make it into the training set end up becoming the test set.

36.2 → end with bootstrap.
36.8 → test set.

Repeat the sampling procedure k-times.

$$Acc(M) = \frac{1}{k} \sum_{i=1}^k (b32 \times Acc$$

$$\in M_i) \text{ test-set} + 0.368 \times Acc(M_i) \text{ train-set})$$

CLUSTER ANALYSIS / CLUSTERING

* Clustering is the process of partitioning a set of data into subclasse[s] / subsets

Application Area:-

→ 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.

Ex:-



K-Means clustering is an example for partitioning method

Hierarchical Methods

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

Ex:-

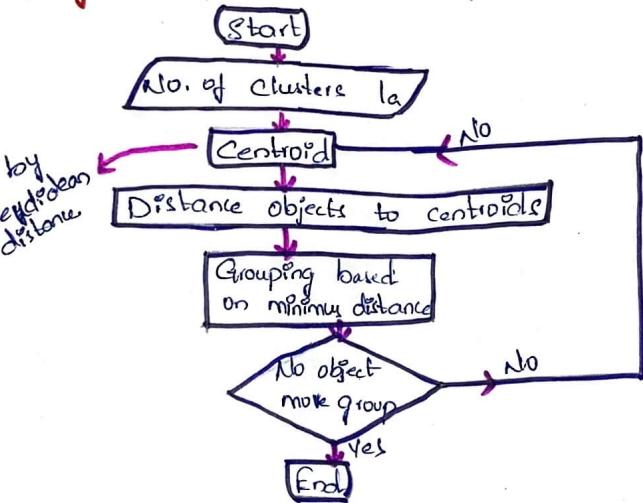


Chameleon is an example for hierarchical clustering.

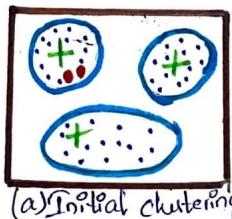
K-Means clustering:-

- * Partitioning clustering approach
- * Each cluster is associated with a "centroid" (i.e. centrepoint)
- * Each point is assigned to the cluster with the "closest Centroid."
- * No. of clusters, K must be specified.
- * Closeness is measured by "Euclidean Distance"
- * Initial centroid are often chosen "randomly"

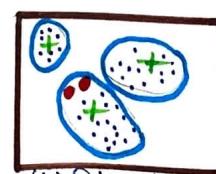
Algorithm / Flow chart Steps:



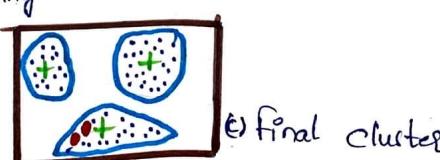
Example:-



(a) Initial clustering



(b) Iterate



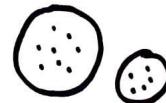
(c) Final clustering

LIMITATIONS OF K-MEANS:-

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

→ Densities &

→ Non-globular shapes



→ & the data contains outliers/noise.

K-Medoids

* K-Medoids also called "Partitio[n]ing Around Medoid"

* The cost in K-Medoids algorithm is given as

$$C = \sum_{i=1}^n \sum_{j=1}^{k-1} |P_i - C_j|$$

Algorithm:-

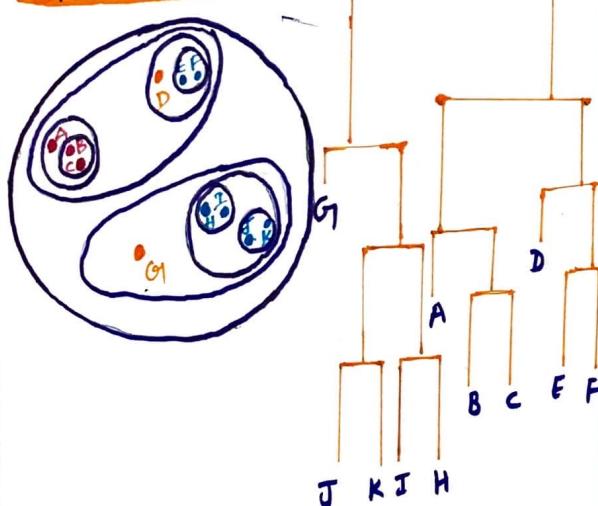
1. Initialize: Select k random points out of the n data points as the medoids.
2. Associate each data point to the closest medoid by using any common distance metric methods.
3. While the cost decreases: For each medoid m, for each data O point which is not a medoid:
 - Swap m and O, associate each data point to the closest medoid and recompute the cost
 - If the total cost is more than that in the previous step, undo the swap.

Heirarchical clustering:

- * Grouping data into a tree of clusters
- * It begins by treating every data point as a separate cluster.

Steps:

- 1) Identify the 2 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.



AGGLOMERATIVE VS DIVISE

- * Bottom up Approach
- * It successively merges the objects (or) group close to one another until all the groups.

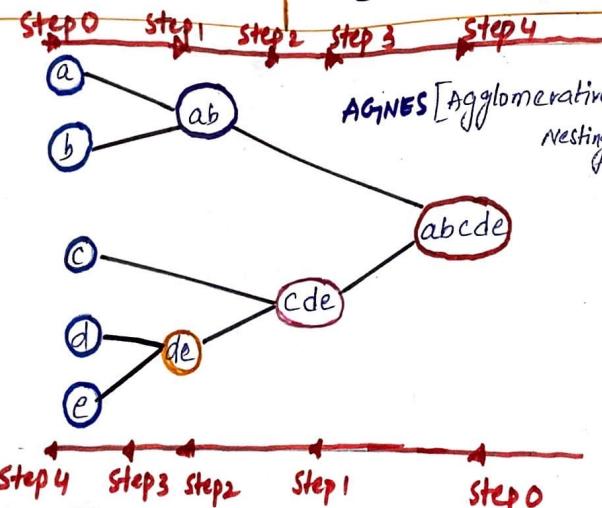
Ex: **AGNES**

- Agglomerative NESTING

- * Top down approach
- * In each successive iteration a cluster is split into smaller clusters, until eventually each object is in one cluster, or a termination condition hold

Ex: **DIANA**

Divisive Analysis.



AGNES:

- * 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.

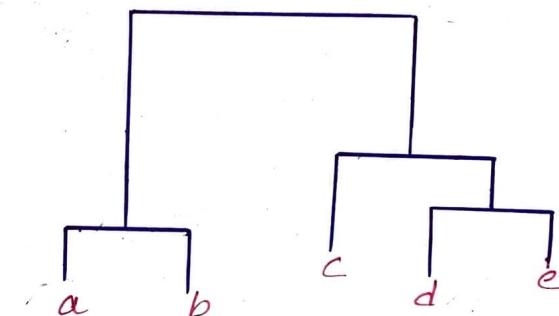
DIANA:

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

Steps

- 1) Initially all points to the dataset belong to one cluster
- 2) partition the cluster into two least similar clusters
- 3) proceed recursively to form new cluster until the desired number of cluster is obtain.

Dendogram:



- * Dendogram shows the clusters are merged

- * A tree structure called a dendrogram is commonly used to represent the process of heirarchical 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-step process :

1) Find all frequent item sets

→ Apriori Algorithm

→ FP-Growth 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. represented 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. item-set.

→ It is a compact data-structure that stores quantitative infn. 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 YUT.base with count = smallest count of nodes in Y;
4. else
 - for each i in T.header do begin
 5. Output Y = T.base U {i} with i.count;
 6. if T.FP-array is defined
 7. Construct a new header table for Y's - FP-Tree from T.FP-array
 8. else
 - Construct a new header-table from T.
 9. Construct Y's conditional FP-Tree Ty & Possibly if FP-array Ay ;
 10. If Ty ≠ φ
 11. call FP.growth * (Ty) ;
 12. End.

APRIORI ALGORITHM.

Steps in Apriori :

Step-1 : Determine the support of itemsets in the transactional database & Select the "minimum support" & "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 [or] min-confidence.

Step-4 :

Sort the rules as the decreasing order to fit.

Flow chart :

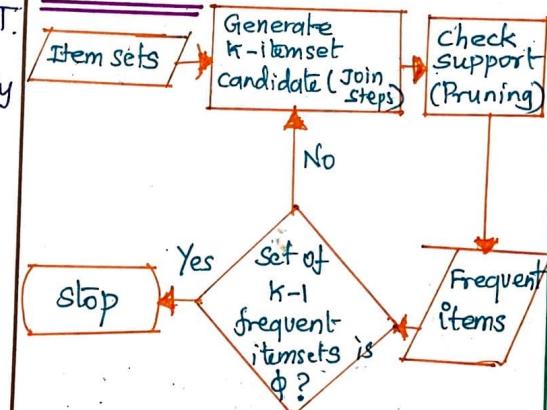
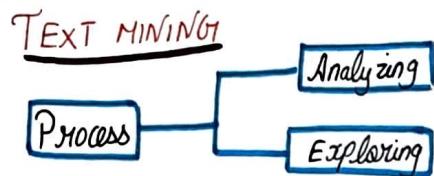


Fig. Flowchart of Apriori



Findings / Outputs \Rightarrow concepts, pattern, topic, etc....

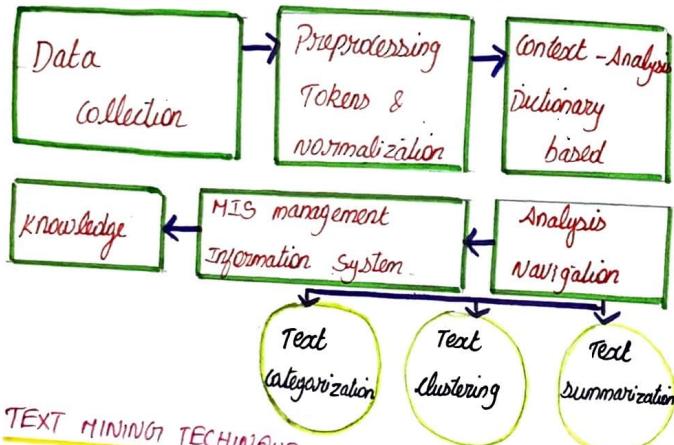
METHODS:-

STEP 1: Gathering unstructured information

STEP 2: Preprocessing - cleaning

STEP 3: Analysis - Analyze the patterns

STEP 4: Information extraction, retrieval, categorization, clustering & summarization



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

KEY-WORD BASED ASSOCIATION

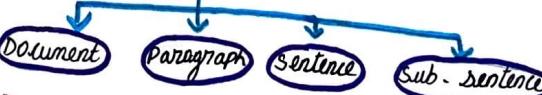
- \rightarrow collects set of keywords/items that occurred frequently
- \rightarrow Preprocess the text \rightarrow by parsing, stemming, removing the stop words
- \rightarrow Implement association mining algorithm
- \rightarrow Each document consider as transaction
- \rightarrow view a set of keywords 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

TEXT CLASSIFICATION

- \rightarrow Analyze at different levels

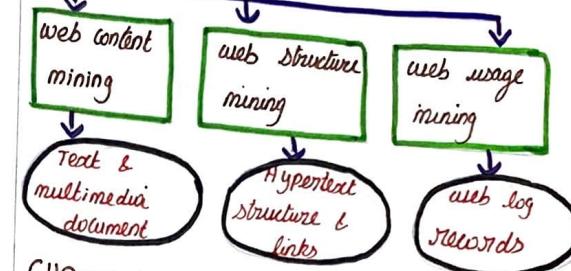


DOCUMENT CLUSTERING

- \rightarrow Group the related document based on their content

WEB MINING

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



CHALLENGES IN WEB MINING

- * The complexity of web pages
- * The web is a dynamic data source
- * Diversity of client networks
- * Relevancy of data
- * The web is too broad

ADVANTAGES

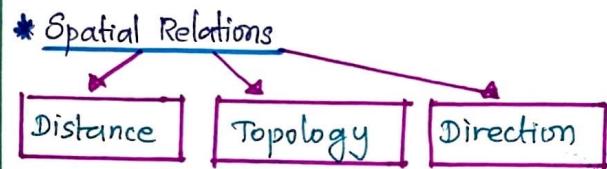
- \rightarrow Find useful information from the user interactions
- \rightarrow understanding consumers behaviour, need & buying patterns

APPLICATIONS

- * Google search engine
- * CRM (customer relationship mgt)
- * Marketing and conversion tool
- * Audience behavior analysis
- * Testing and analysis of a site
- * Data analysis on website and application

Mining Spatial Data

- Geographical [or] spatial information to produce business intelligence (or) other results.
- Data related to spatial description of the objects such as
 - co-ordinates
 - Areas
 - Latitudes
 - Perimeters
 - etc.,



* Finding & Analyze helps :-

- Earthquake points
- Climate / Weather predictions
- Google Map
- GPS (Global Positioning System)
- Trend Analysis

* Clustering Analysis Methods :-

- PAM → Partitioning Around Medoids
 - [//] also 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

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 :

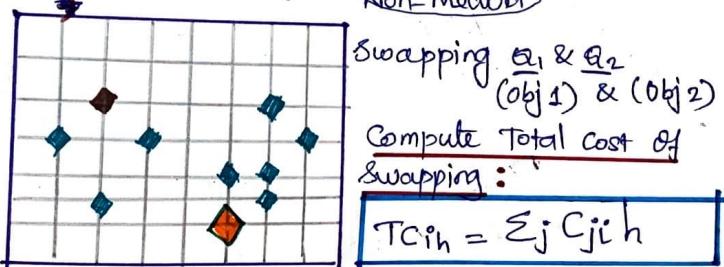
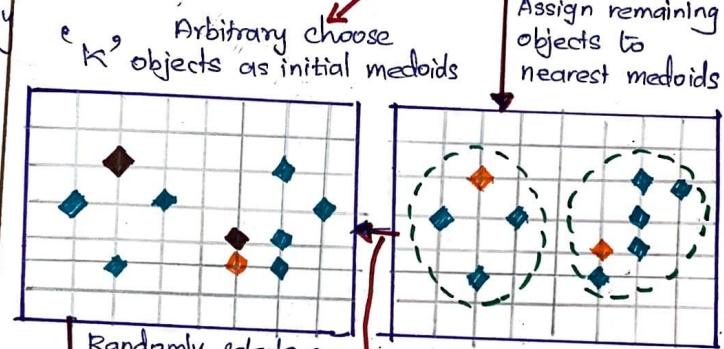
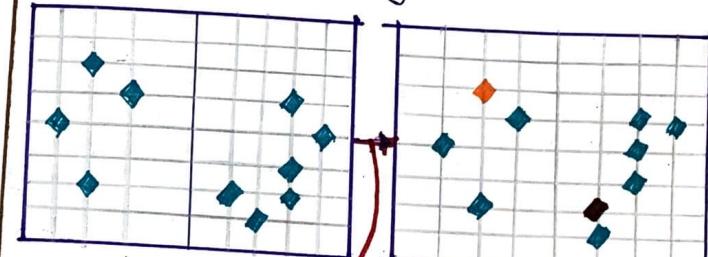
- Follows the raster data model
- Data contains the mapping of continuous data in space.
- Represented as a satellite image, digital surface model etc.,
- The visual representation of coverage data is in the form of a "Grid" [or] triangulated irregular network.

Partitioning Around Medoids (PAM)

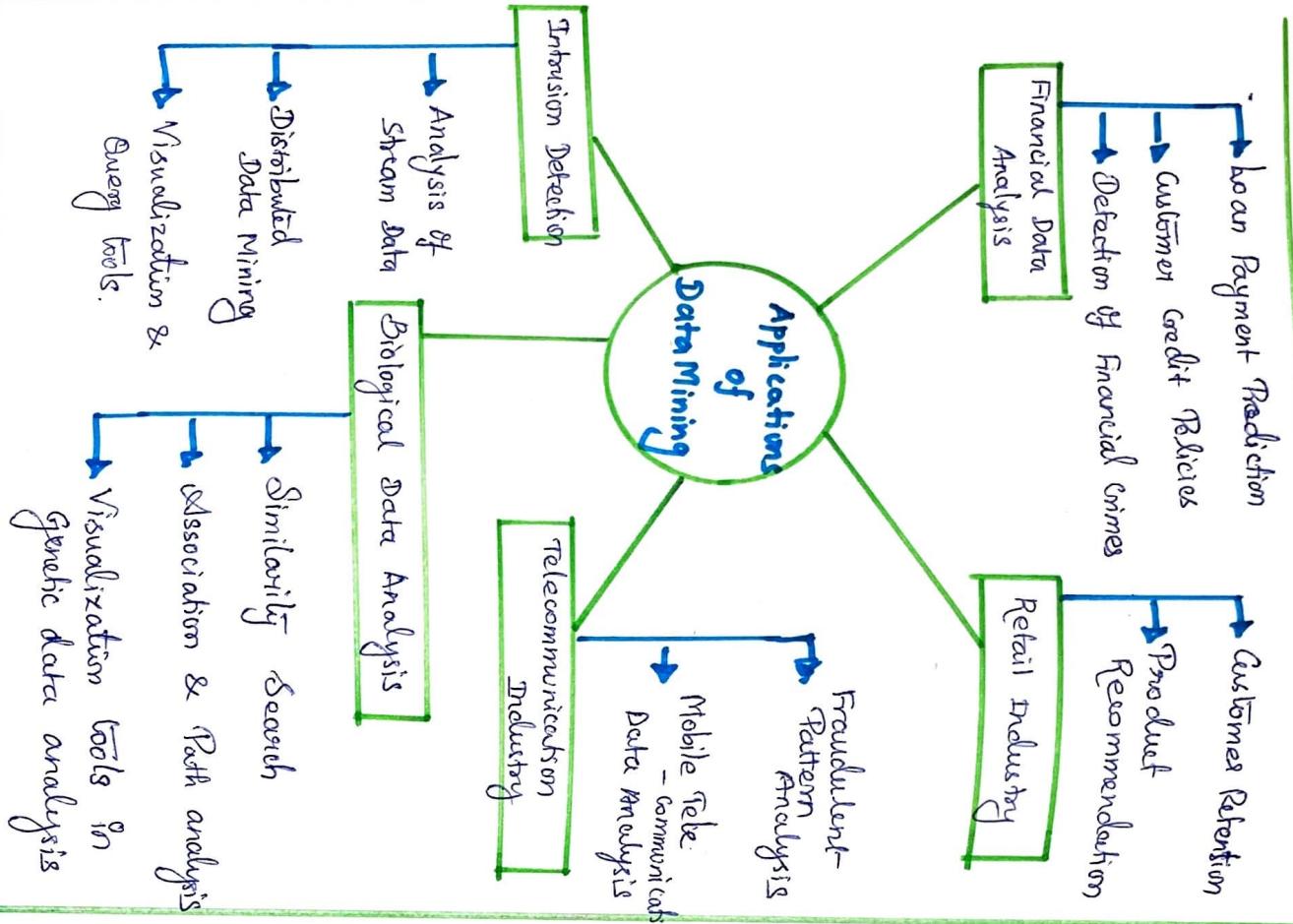
* Find a sequence of objects called "medoids" that are centrally located in clusters.

* Starts from initial set of medoids & iteratively replaces one of the medoids by one of the non-medoids

↓
if it improves the total distance of the resulting clustering.



APPLICATIONS IN DATAMINING



TRENDS IN DATA MINING

