

Unit-3

* Classification - It is the process of finding good model that describe the data class or concept and its purpose to predict the class of object whose label is unknown. (categorical value) - Two step - Model making & model usage.

* Prediction - Prediction is about predicting a missing/unknown element (continuous value) of a dataset.

* Issues regarding classification and prediction

1) Data preparation

→ Data cleaning

→ Relevance analysis (feature extraction)

→ Data transformation

2) Evaluating classification method

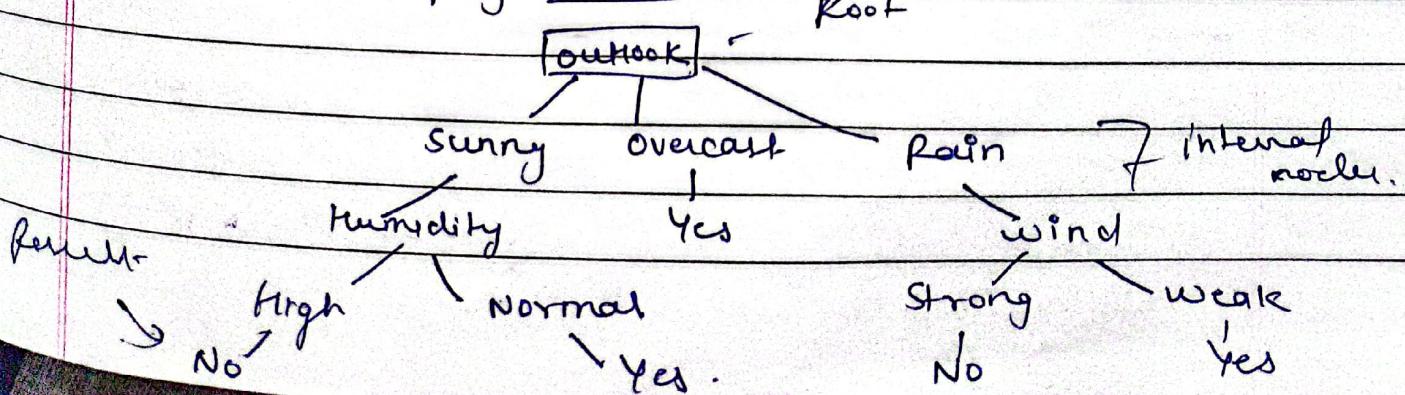
goodness of rule

- Predictive accuracy
- speed scalability
Robustness (handle noise)
Interpretability (Easy understand)

* Decision tree classification - supervised machine learning technique used for classifying problem.

- It's a tree structured classifier
- internal node represent feature of dataset
- branch represent rule
- each leaf node represent the outcome.
- Contain root node where decision tree starts.

Decision tree to play Tennis



- Bayesian classification - It is a collection of classification algorithm based on Bayes's theorem. It uses Bayes's theorem to predict occurrence of any event. These are statistical classifiers with probability understanding.

Bayes's theorem is expressed as -

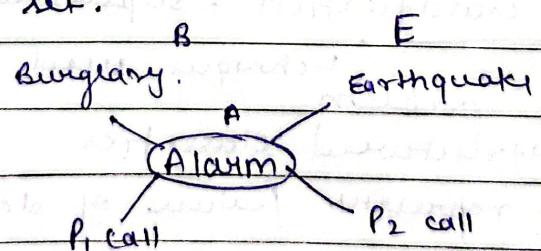
$$P(A/B) = \frac{P(B/A) P(A)}{P(B)}$$

$P(A/B)$ conditional probability that of given that B happened. is true.

$P(A)$ & $P(B)$ are probability of observing A & B independently

- Assumptions -
- No features are dependent i.e. they are independent
 - Each feature will give same importance.

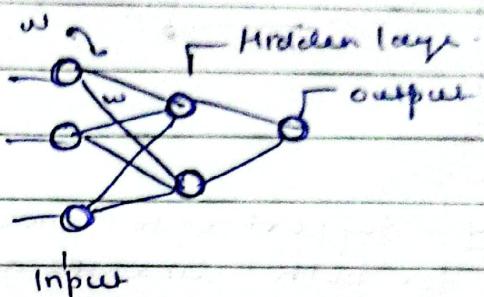
- Bayesian belief networks - It is graphical representation of different probabilistic relationships among random variables in a set.



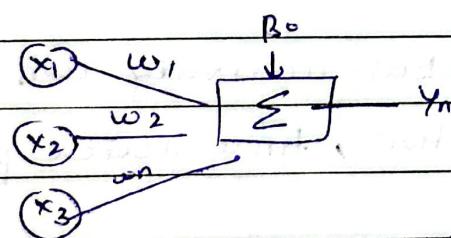
It is probabilistic graphical model, which represents a set of variables and their conditional dependencies using a directed acyclic graph.

- Multilayer feed forward Neural Network - It is an interconnected artificial neural network with multiple layers that has neuron with weights associated with them, and they compute result using activation function.
- . In this flow is from input to output, do not contain loop, feedback, no signal move in backward direction.

↳ Input layer
↳ Hidden layer
↳ Output layer.



- Back propagation algorithm - It is an algorithm that back propagates the error from output nodes to the input nodes. Simply called back propagation of error.



Steps - Initialise weights \rightarrow Propagate input forward \rightarrow
Backpropagate error \rightarrow terminating condition.

weights are modified to minimise the error.

Adv - High tolerance to noisy data
classify untrained patterns

Disadv - long training

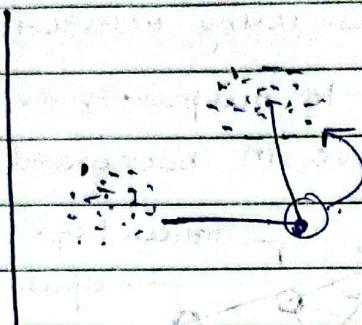
- Poor interpretability
- large parameters required

- k-nearest neighbour - KNN algorithm assumes the similarity between new data and available data and assign class based on similarity. It is also called a lazy learner algorithm because it does not learn from training set immediately instead it stores and at time of classification it performs action.

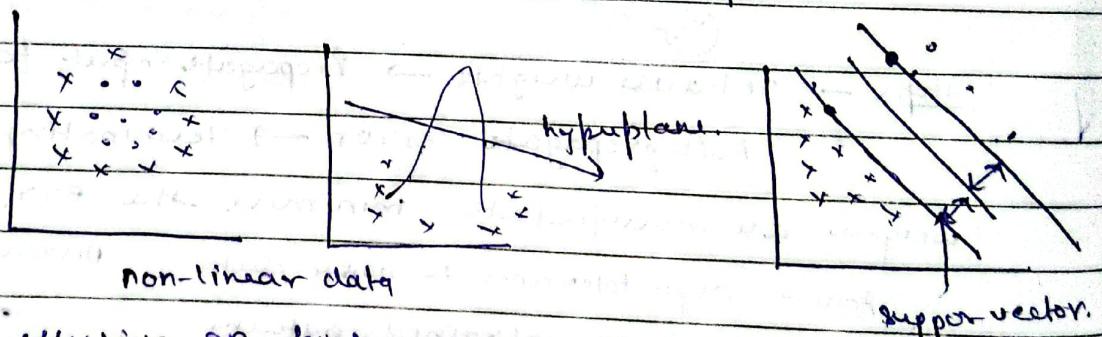
- It stores data and when it gets new data, then it classifies data into new category similar to new data.

steps - ~~seed~~ - k nearest neighbour

- Among these count data point each category
- Assign new point to category which the no. of neighbour is maximum.



- * SVM - Support vector machine - classification algorithm for both linear and non-linear data.
 - It uses non-linear mapping to transform data into higher dimension.
 - With this now data can always be separated by hyperplane.
 - Training can be slow, but accuracy can be high.
 - Application - digit recognition, time series prediction.



- very effective on high dimensional data.

- genetic algorithm - It is based on analogy to biological evolution.
 - An initial population is created consisting randomly generated rules.
 - Based on survival of fittest a new population is formed.
 - Offspring are generated by crossover and mutation.
 - Process continues until a population P evolves where each rule in P satisfies a pre-specified threshold.

- Cluster analysis - cluster - a collection of data objects
 - similar data are in one cluster
 - dissimilar to others lie in other clusters.
- cluster analysis - grouping a set of data object into cluster.
 - it is unsupervised technique
- Application - to get insights from data
 - as it is a preprocessing step for some algorithm.

general application - pattern recognition, image processing.

examples - Marketing - discover groups in customer base

Land use - identification of similar land

Insurance - identify high claim cost groups

City Planning - identify group of house based on

house type
Earthquake - observe earth quake epicentres

- good clustering
 - high quality clusters
 - quality depends on similarity measure
 - ability to find hidden pattern.
 - high inter class similarity.

- Types of data in clustering analysis

- Interval-valued variable - interval scaled variable are continuous measurement of roughly linear scale

• standardise data

- calculate mean absolute deviation.

$$m_f = \frac{1}{n} (n_{1f} + n_{2f} + \dots + n_{nf})$$

- calculate standardised measurement (z-score)

$$z = \frac{n_{if} - m_f}{S_f}$$

Similarity and dissimilarity are measured on the basis of distance.

$$d(i, j) = \sqrt{(x_{i1} - x_{j1})^2 + (x_{i2} - x_{j2})^2 + \dots + |x_{ip} - x_{jp}|^2}$$

$q=1$ = Manhattan distance $q=2$ Euclidean distance

Binary variable - Binary variable can take only two values. Example - gender - male or female.

↳ Symmetric - like male-female may family

↳ Asymmetric - like disease.

- Nominal Variable A generalization of binary variable it can take more than 2 states e.g. red, green, yellow

Method-1 Simple matching - $d(i, j) = \frac{p-m}{p}$

p total variables m: m no of match when i, j have same state.

- Ordinal Variable - It can be discrete or continuous.
order is important

- can be treated like interval scaled.

- Ratio scaled variable - A positive measurement on a non-linear scale approximately at exponential scale.

such as - $A e^{Bt}$ or $\$$

method - - first treat them like scaled variable

- apply logarithmic

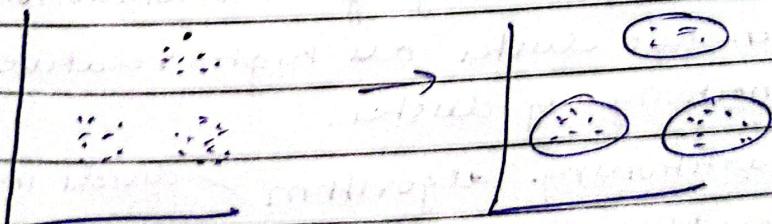
- finally treat them as continuous ordinal data.

- Variable of mixed type - A database may contain all the six types of variables.

all these combined called mix-type variable.

clustering method .

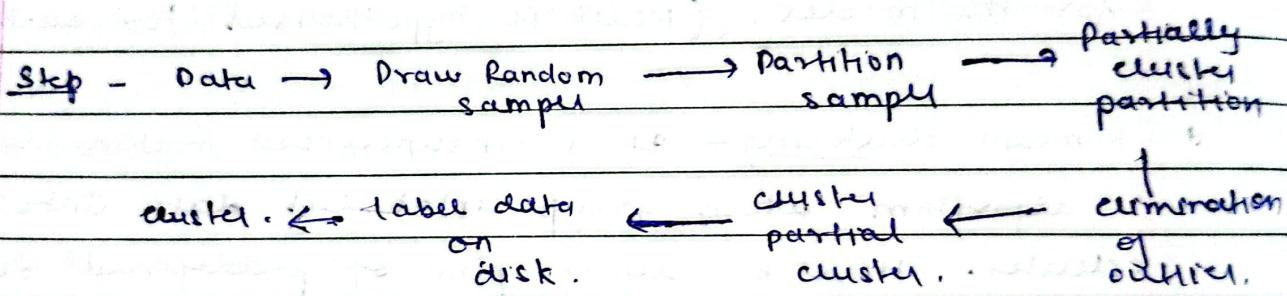
- 1) Partitioning algorithms - construct various partitions and evaluate them by some criteria. (ACP)
 - 2) Hierarchy algorithm - create a hierarchical decomposition of set of data using some criteria.
 - 3) Density based - based on connectivity and density function
 - 4) Grid based - based on multilevel granularity structure
In this method grid is formed i.e object space is quantized into finite no. of cells. It is fast
 - 5) Model based - In this all cluster are hypothesized. In order to find the data which is best suited for the model. (Model is hypothesized for each cluster)
- * K-means clustering - It is unsupervised machine learning algorithm which group unlabelled data into different cluster. Here k - defines. p no. of predefined cluster.
 It is a centroid based algorithm where each cluster is associated with a centroid. The main aim of algo is to minimise sum of distance b/w data point.



- Steps - Select number k and initialise
 select random k point as centroid.
 Assign data point to their closest centroid.
 calculate new centroid.
 Repeat the process.

core clustering - It is a hierarchical clustering technique that adopt a middle ground b/w centroid based and all point extremes. It start with single point cluster and move to merge until desired clusters are formed.

- Used for identifying spherical & non-spherical clusters.
- useful for discovering groups and identify distribution
- instead of one point centroid, we use a set of well defined representative point for efficiently handle cluster & outliers.



- Chameleon clustering - chameleon is a hierarchical clustering algorithm that use dynamic modeling to divide similarity among pairs of clusters.
- Two cluster are merged only if interconnectivity and closeness between two cluster are high relative to internal interconnectivity of cluster.
- It use graph partitioning algorithm - divide in sub clusters
- Use an agglomerative hierarchical clustering find cluster by combining sub-cluster.
- DBSCAN - Density based spatial clustering of application with noise.

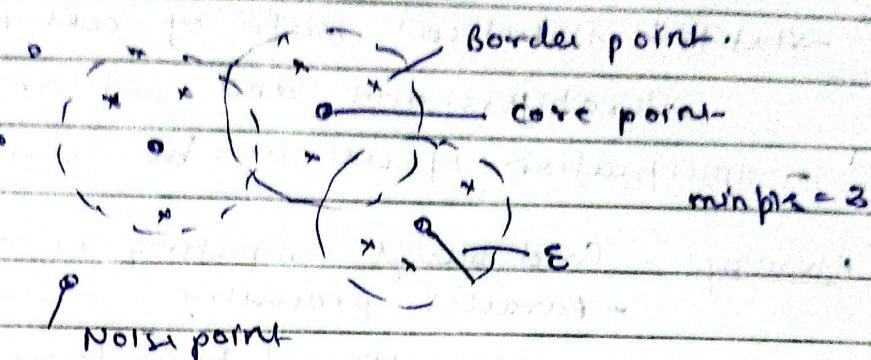
It is a density based clustering method. It can

discovered clusters of different shapes and sizes from large amount of data, which contain noise and outliers.

other two parameters -

minpts - The minimum no. of points in a region

eps - A distance measure that is used to locate point in neighbour to point.



core point - A point is core if it has more than minpoints

Border point - A point which has fewer than minpoints

Noise - A point non a core point.

* optics - ordering points to identify cluster structure.

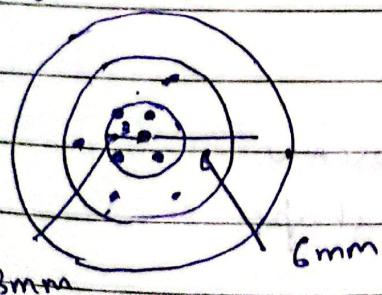
- draw inspiration from DBSCAN

- Add two more concepts

1) core distance - It is minimum value of radius required.

to classify a given point as a core point.

If given point is not core then core distance is undefined.



$\text{minpts} = 5$

core distance = 3mm.

2) Reachability distance - It is defined w.r.t. another data point q. Reachability between p & q is max of core

distance of p and Euclidean b/w p and q.

Reachability distance is not defined if q is not core.

* String clustering - (Statistical Information grid)

- Spatial area is divided into rectangular cell.

- there are several level of cell correspond to different level of resolution

- each cell at high level is separated to multiple smaller cell

- statistical data of cell is computed and stored beforehand and can answer query.

- specification of cell can be count, mean, min, max.

Advantages - Grid based computing is query independent

- parallel processing, incremental update.

disadvantages - No diagonal boundary detected.

* Clique clustering - Clique is a density-based and grid-based subspace clustering

grid-based - use grid

Density - density based

- find out cluster by taking density threshold and no of grid as input parameter.

- large dimension data

- very scalable

- use apriori algorithm.

Advantages - better than kmean, DBSCAN
find cluster of any shape
simple method

disadvantage - if size of cell is unsuitable estimation will be too much, unable to find cluster

- Model based clustering is statistical approach to data clustering. The data is considered to have created from a finite combination of component model. model is → probability distribution.
 - Statistical approach - It assign object to cluster according to weight (probability distribution) New means are computed.
 - Neural Network approach - Represent each cluster as an example acting as a prototype of cluster New object are distributed to cluster whose exemplar is most similar according to some distance measure.
- Outlier detection
- Outlier detection -
- outlier are a set of objects that are dissimilar from remainder of data.
- Problem - find top outliers.
- Application - credit card fraud detection
telecom fraud.
customer segmentation
medical analysis.

It can be found using distance based approach and deviation based approach.

Unit-N

Data warehouse - data warehouse is collection of data marts representing historical data from different operations in the company. Data is stored in structure optimised for querying and data analysis in data warehouse. Table design, dimension & organisation should be consistent so that query are consistent.

	Operational data	Informational data
Data content	current values	summarised, archived, derived
Organisation	By application	by subject
Stability	Dynamic	Static until refreshed
Structure	optimised for transaction	optimised for complex query
Access frequency	High	medium / low.
Response	2-3 s	several seconds.
	updateable	non-updateable

Database system	Data warehouse
• Support operational process	support analysis and performance reporting.
• Capture & maintain data.	Explore data history.
• Current data	Data update on scheduled process.
• Data is update when transaction occur	subject oriented.
• Application oriented	summarised and consolidated.
• Primitive & highly detailed	Star / Snowflake based.
• ER based	

* multidimensional data model - it is a method which is used for ordering data in database along with good arrangement and assembling of content in database, this allow customer to interrogate questions comparatively fast.

it represent data in form of data cubes which allow to view data from many dimension & perspective.

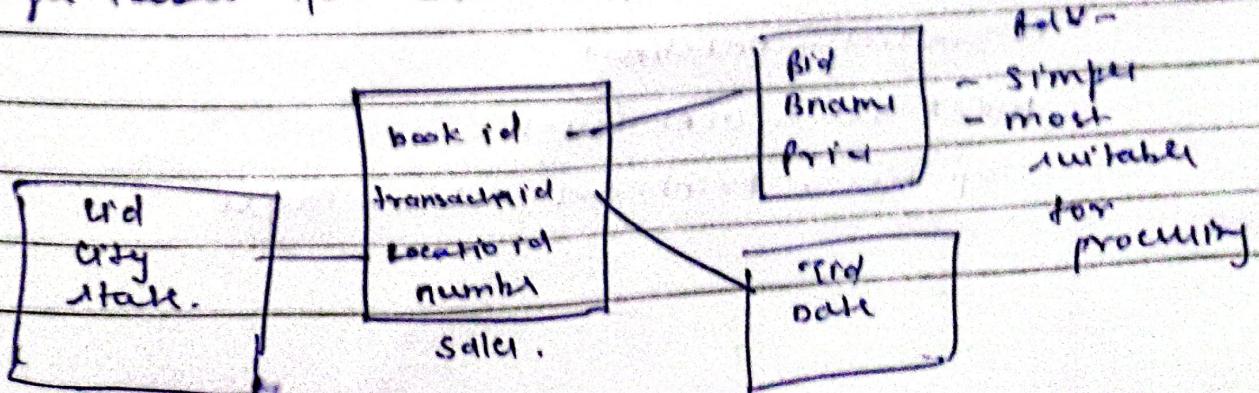
Advantages

- easy to handle
 - easy to maintain
 - better performance
 - better representation
 - work on complex system
- Schema for multidimensional database

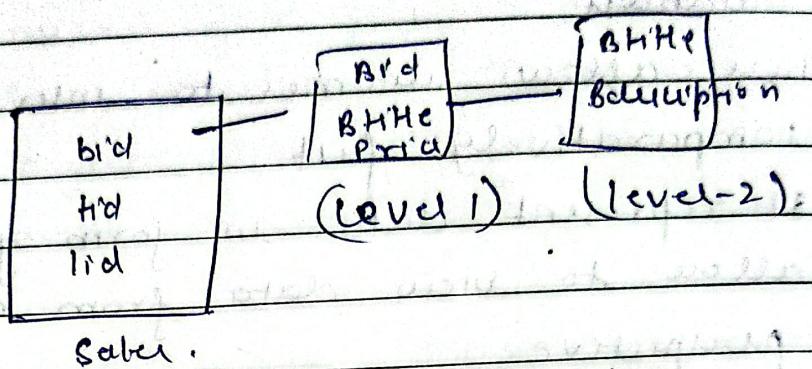
schema logical representation of entire database
it includes name and description of all record type.

Star schema - each dimension in star schema is represented with only one dimension.

- This dimension table contains the set of attributes.
- Single table for each dimension.



Snowflake schema - Some dimension table in snowflake schema are normalized. Normalization splits up data into additional tables.

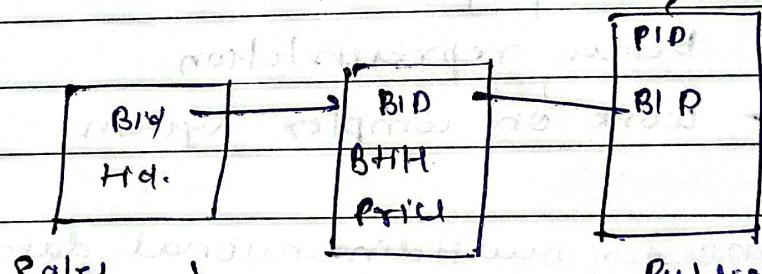


Adv - less redundancy

Easy to update

Disadv - complex.

Fact constellation schema - also known as galaxy schema. It's that multiple tables share same dimension table.



Adv

disadv

- complex

- difficult to maintain

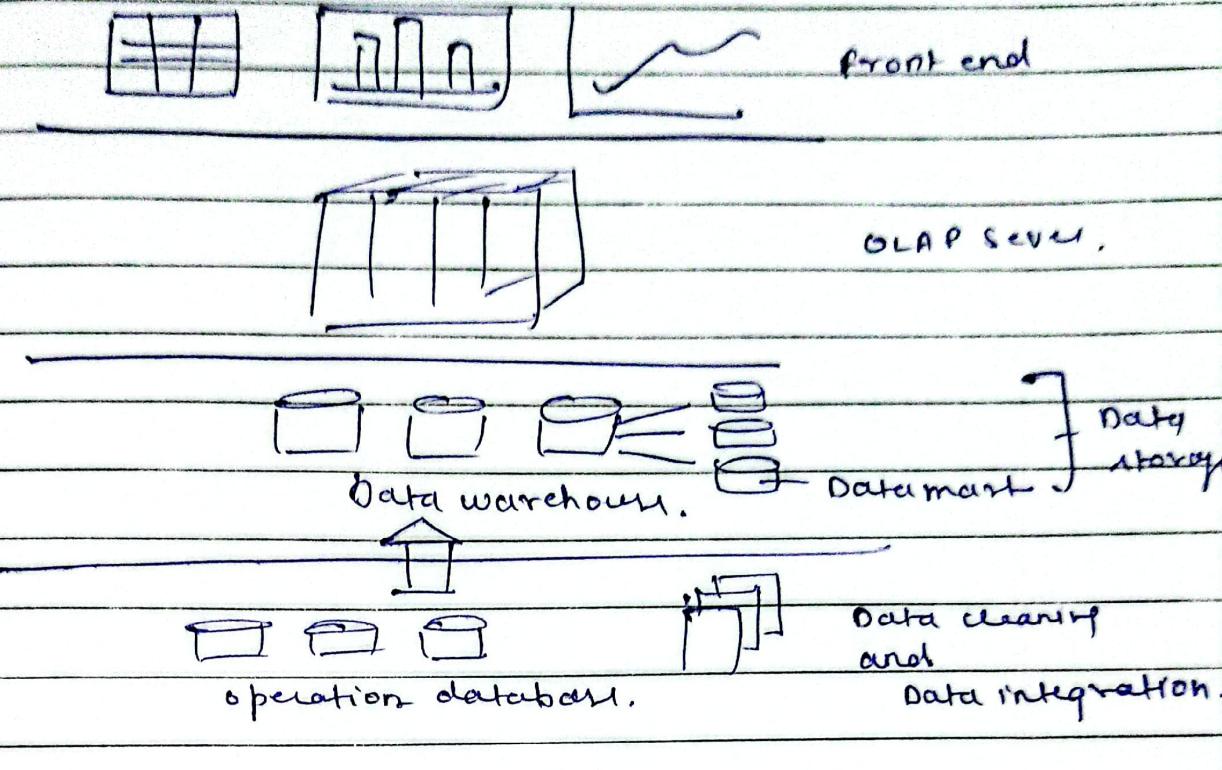
- Dimension tables are large.

Data warehouse - 3 tier architecture

Bottom tier - Database

Middle tier - OLAP server

Top tier - Front-end client layer.

architecture

Data mart - A data mart contain subset of corporate data that is of value to specific group of user.

- independent data mart are sourced from data captured from one or more operational system.
- dependent data mart are sourced directly from enterprise data warehouse.