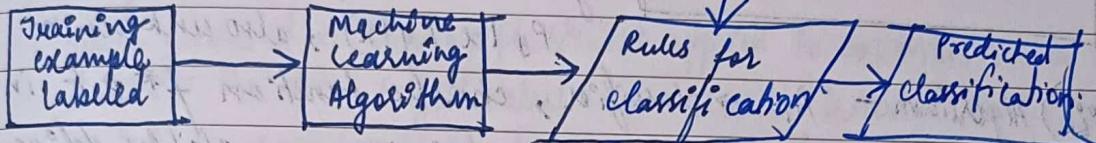


Impact of ML in BI process.

* Classification Problem -

TEL, MLA, HFL, NC, PC

TMR, P
N



→ Input - Training set of examples with class labels.

→ Aim - To predict categorical class labels for new labels.

→ Output - classifier is based on training set & class labels.

① Classification predicts categorical tables (classes).

Prediction models continuous values functions.

② Preprocessing data in preparation for classification & prediction.

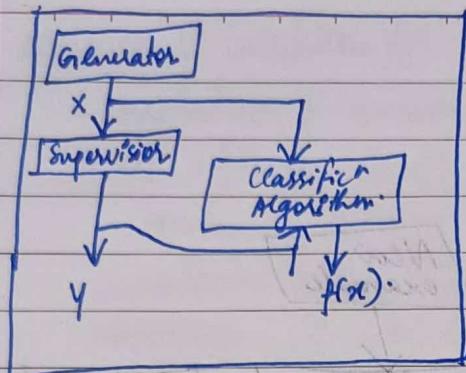
i) Data cleaning - reduces noise, handles missing values.

ii) Relevance analysis - removes irrelevant or redundant attributes.

iii) Data Transformation - generalizing data to higher level concepts
normalizing data.

③ Classification → process of finding a model that describes sets to distinguish data classes or concepts, for purpose of being able to use model to predict the class of objects whose class label is unknown. Survived model is based on analysis of a set of training data.

⇒ Probabilistic structure of learning process for classification.



- i) Generator - extracts random vector x of examples acc to an unknown probability distribution -
- ii) supervisor - Returns for each vector x of examples the value of target class acc conditional distribution $P_y | x (y | x)$, also unknown.

iii) Algorithm - $A_F \rightarrow$ classifier, chooses function $f^* \in F$ in the hypothesis space so as to minimize a suitably defined loss function (for example to fit points - least squares method)

\Rightarrow Classification Problem -

- i) suppose database D is given as $D = \{t_1, t_2, \dots, t_n\}$ & set of desired classes. $C = \{C_1, C_2, \dots, C_n\}$
- the classification problem is to define mapping $m : D \rightarrow C$ such a way that triple of database t belongs to class of C . Actually divides into equivalence classes.

2) Classification in 2 step process -

- i) Model construction -

 - sample object has a predefined class label assigned.
 - sample data is training data.
 - constructed model based on training data set is represented as classification rules, decision trees or mathematical formulae.

ii) Model usage -

- For classifying unknown objects, use construction model.
- compare class label of resultant sample with test sample.
- estimate accuracy

3) Phases of classification Model -

i) Training phase

- classification is applied to subset T for dataset D , called Training set.
- derives classification rules that allow corresponding target class y to be attached to each observation x .

ii) Test phase -

- . rules generated in training phase are applied to observation not in database D , for target value are already known.
- . to test accuracy actual target values ~~are already known~~ class of test set is compared with the predict class.

iii) Prediction Phase -

- . Actual use of classification model.
- . Assigns target class to new observation to be recorded infiltration

Categories of Classification model -

(Heuristic, separation, Probabilistic, Regression).

* Heuristic Model \rightarrow make use of classification procedures based on simple & intuitive algorithms.

- Example - Nearest Neighbour method

Separation Model -

- Divide attribute space R^n into n disjoint regions $\{S_1, S_2, \dots, S_n\}$
- Separating observation based on target class.
- Example - support vector machine, perceptron methods.

Probabilistic model -

- Hypothesis is formed regarding functional form of conditional probability, $P(y|x)$ of observations, given the target class.
- Conditional class probability.

* Regression Model -

- Regression is data mining function that predicts numbers -
- Good choice when all predictor variables are continuous valued -

* Evaluation of Classification Models -

- ① Various methods for estimating a classifier's accuracy are -
 - A] Holdout Method
 - B] Random sampling
 - C] cross validation
 - D] Bootstrap

All are based on randomly subsampled partitions of data.

- ② Comparison of classifier ID choose the best -
- A] confidence intervals
 B] cost benefit analysis & receiver operator curve (ROC, ROC).

* Accuracy & Error Measures -

- ① accuracy of classifier M = acc (M) is percentage of testset tuples that are correctly classified

$$\text{Acc}(M) = \frac{TP + TN}{TP + TN + FP + FN}$$

- ② success - instance class is predicted correctly

errors - instance class is predicted incorrectly

- ③ confusion Matrix

		Predicted values		Total	
		G ₁	G ₂		
A	G ₁	TP	FN	P	
	G ₂	FP	TN		
		P'	N'		

TP = class members

classified as class members

FP = non members classified as class members

FN = non class members classified as non class members

FN = class members classified as non class members

P = No. of +ve tuples P' = No. of tuples labelled +ve
 N = No. of +ve tuples N' = No. of tuples labelled -ve

③ sensitivity -

→ True positive recognition rate.

→ Proportion of positive tuples correctly identified.

$$\therefore SN = \frac{TP}{P} = \frac{TP}{TP+FN}$$

④ specificity

→ True negative recognition rate or of negative tuples correctly identified.

$$SP = \frac{TN}{N} = \frac{TN}{FP+TN}$$

Accuracy is function of SN & SP = $\boxed{ACC(M) = SN + SP}$

⑤ error rate : % of error made.

$$\text{error rate} = 1 - ACC(M) = \frac{1 - TP + FN}{TP + FP + FN + TN} = \frac{FP + FN}{TP + FP + FN + TN}$$

⑥ Precision - Measure of exactness

No. of positives truly classified as me

⑦ Recall

same as sensitivity

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

$$R = \frac{TP}{P} = \frac{TP}{TP+FN}$$

⑧ F1 score = Harmonic mean
of P & R.

$$\boxed{F1 = \frac{2 \times P \times R}{P + R}}$$

⑨ $F_B = \frac{(B+1)^2 \times P \times R}{B \times P + R}$
= weighted measure of
P & R assigns B times as much
weight to recall as precision

$$F_B = \frac{(B+1)^2 \times P \times R}{B \times P + R}$$

B is non negative
real number.

10) Classifiers can be compared w.r.t -

i) Robustness

→ Accuracy of rules doesn't change significantly when changing of test set variable.

→ Ability to handle missing data & outliers.

ii) Scalability.

→ Ability to learn from large datasets.

iii) Speed

iv) Interpretability

* Hold out Method =

① Data is split into training dataset & testing dataset.

② Training - $2/3$, Testing = $1/3$.

③ Train classifier, training dataset is used, & test dataset is used to estimate true error rate.

Training set:

Test set:

④ More training leads to better model construction more testing means error estimations are more accurate.

⑤ P - samples might not be representative. some classes will have few instances. S - stratification to ensure both T & V have same no. of instances of same class.

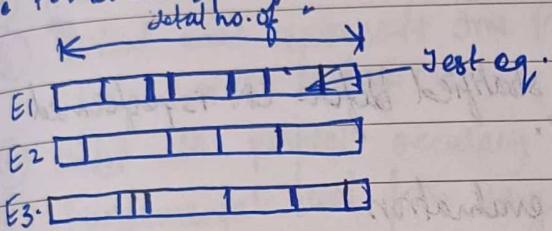
⑥ Drawbacks -

i) Requires extra dataset.

ii) Single train & test experiment. misleading error estimate.

Random Subsampling =

- Variation of hold out method.
- Hold out method is repeated k times.
- each split randomly selects fixed number examples w/o replacements.
- For each split, we retain the classifier from scratch.



• overall accuracy is calculated by averaging accuracies from each iteration-

$$E = \frac{1}{k} \sum_{i=1}^k E_i$$

⇒ Cross validation (k fold, Leave one out)

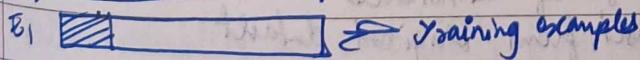
- ① ensures non overlapping, disjoint sets.
- ② Techniques for evaluating estimating performance by training several ML models on subsets of available of input data & evaluating them on complementary subset of data.
- ③ used to detect overfitting.

⇒ k cross validation:

Step 1 - Data is split into k subsets of equal size usually by FS.

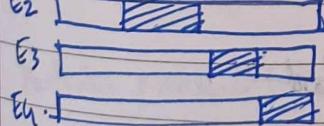
Step 2 - each subset is in turn used for both training & testing.

\downarrow Total no. of examples



• Advantage is all examples-

• are used for both training & testing.



• ERROR ESTIMATE

$$\bar{E} = \frac{1}{k} \sum_{i=1}^k E_i$$

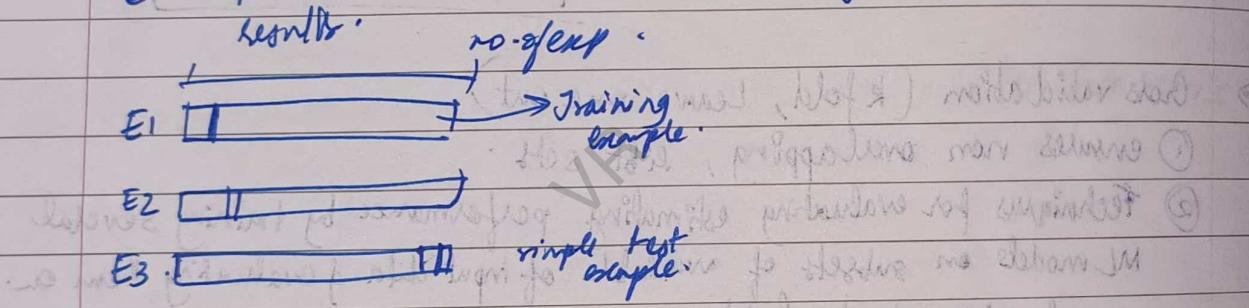
⇒ Leave One Out CV -

- If O has N examples, N experiments need to be performed.
- For every experiment, training uses $N-1$ examples remaining for testing.
- $E = \frac{1}{N} \sum_{i=1}^N E_i$ is avg error rate of test examples, giving true error.

⇒ Stratified CV - subsets are stratified before CV is performed

stratified 10-fold CV.

- ① gives accurate estimate of evaluation.
- ② estimates variance gets variance gets reduced due to stratified.
- ③ repeated 10 times & result will all averaged based on proportions



⇒ Bootstrapping -

- ① CV uses sampling set without replacement i.e., once samples or instance is selected, it can't be selected again for training or testing.
- ② bootstrap uses sampling with testing w.r.t. replacement to get training set.
- ③ Training set - Dataset of k instances is sampled using replacement k times to form training set of k instances.
- ④ Test set - separate dataset from original dataset (not part of training set).
- ⑤ Best error setting - estimator for small dataset

ROC Curve.

- Receiver Operating Characteristics curve.
- ① used to visually compare classification methods.
- ② originating roots from signal detection theory.
- ③ Tradeoff b/w TP rate & FP rate.
- ④ Accuracy of model can be measured by AUC.
- ⑤ vertical axis represents True Positive Rates & horizontal axis represents False Positive Rates.
- ⑥ Model with perfect accuracy will have AUC = 1.
- ⑦ Fundamental tool for diagnostic test evaluation.

* Bayesian Methods -

- belongs to family of probabilistic classification model.
 - once prior & class conditional probability are known it can explicitly calculate posterior probability.
 - Bayesian classifier are statistical classifiers.
 - predict class membership probabilities.
 - given observation belongs to particular class.
 - exhibit high accuracy & speed when it comes to large databases.
- Prior probability = Initial probability.

* Baye's Theorem -

Prior Probability = probability of event based on established knowledge, before any data is collected.

Posterior probability = revised probability of an event, after taking into consideration new data.

Baye's theorem statement

Let A_1, A_2, \dots, A_n be events representing partition of sample space. If $P(A_i) \neq 0; i=1, 2, n$
 Let B be any other defined event on S . If $P(B) \neq 0$ then

$$P(A_i|B) = \frac{P(A_i) \cdot P(B|A_i)}{\sum_{i=1}^n P(A_i) \cdot P(B|A_i)} = \frac{P(A_i) \cdot P(B|A_i)}{P(B)}$$

Proof -

$$\begin{aligned} S &= A_1 \cup A_2 \cup \dots \cup A_n \text{ & } B \in S \\ \therefore B &= B \cap S \\ &= B \cap (A_1 \cup A_2 \cup \dots \cup A_n) \\ &= (B \cap A_1) \cup (B \cap A_2) \cup \dots \cup (B \cap A_n) \end{aligned}$$

(Distributive Law)

$\therefore (B \cap A_i)$ are mutually disjoint for $i=1$ to n ;

By addition theorem of probability .

$$\begin{aligned} P(B) &= P(B \cap A_1) + P(B \cap A_2) + \dots + P(B \cap A_n) \\ &= \sum_{i=1}^n P(B \cap A_i) = \sum_{i=1}^n P(A_i) \cdot P(B|A_i) \end{aligned}$$

Also we have ,

$$P(B \cap A_i) = P(B) \cdot P(A_i|B)$$

$$P(A_i|B) = \frac{P(B \cap A_i)}{P(B)} = \frac{P(A_i) \cdot P(B|A_i)}{\sum_{i=1}^n P(A_i) \cdot P(B|A_i)}$$

Formula for Baye's Theorem,

$$P(A|B) = \frac{P(A \cap B)}{P(B)} = \frac{P(A) \cdot P(B|A)}{\sum_{i=1}^n P(A_i) \cdot P(B|A_i)}$$

{ Red | Domestic | SUV }

$$P(R|Y) \quad P(D|Y) \quad P(SUV|Y) \quad P(R|N) \quad P(D|N) \quad P(SUV|N)$$

Yes .

$$\begin{aligned} n &= 5 \\ n(R) &= 3 \end{aligned}$$

$$P(R|Y) = 3/5$$

$$P(B|Y) = ?$$

$$(R|S)9 \cdot (S|A)4 = (R|A)4$$

① Bayesian Belief Networks -

- * conditional Independence is defined b/w subsets of variables.
- * N/w provides graphical model of causal reln of which learning can be performed.
- * Trained N/w can be used for classification.
- * Two components.
 - i) Directed Acyclic Graph
 - ii) set of conditional probability tables.

* Directed Acyclic Graph -

- each node represents random variable, discrete or continuous.
- discrete or continuous variables which may correspond to actual attribute or hidden variable
- probabilistic dependence is represented by each graph.
- $A \rightarrow B$: A is immediate predecessor or parent of B.
- each variable is conditionally independent of its non-descendents.
- eg - ① Alarm for Burglary & Earthquake -

$$\begin{aligned}
 p(AC, VC, A \wedge E, JE) &= p(AC|A) \times p(VC|A) \times p(A|B \wedge E) p(JE) p(JB) \\
 &= 0.90 \times 0.70 \times 0.001 \times 0.99 \times 0.99 \\
 &= 0.00062
 \end{aligned}$$

② Appl'n of BN -

ML, statistics, computer vision, NLP, speech recognition, email
categorization.

* Logistic Regression

- regression analysis technique where's outcome variable is binary or dichotomous.
- supervised learning algorithm.
- doesn't involve decision tree more like linear regression.
- can be used with 2 types of target variables-
 - i) categorical TV (binary/dichotomous).
 - ii) continuous TV (representing probability values or proportion).

* Binary logistic Regression model -

$$P(Y) = \frac{e^z}{1+e^z}$$

where $z_0 = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_k x_k$

* Logistic regression function is written as $\frac{p(Y=1)}{1-p(Y=1)}$

$$\therefore \ln \left(\frac{p(Y=1)}{1-p(Y=1)} \right) = \beta_0 + \beta_1 x_1 + \dots + \beta_k x_k$$

- Multiple linear regression allows response variable y to be modelled as linear function of 2+ predictor variables.

- Logistic function is similar to MLR such models are called.

Generalized Linear Model (GLM).

How, errors do not follow normal distribution if there exists a function of outcome variable that takes linear variable.

Clustering (Classification).

Unsupervised

training's ample not provided.
data is not labelled.

How can I group their set of items?

Unknown no. of classes

Used to understand data

The function maps the data
into one of several clusters
which is one of several clusters.

which is grouping of data items

based on similarities b/w them.

No. of clusters is not known

before clustering these are

identified after completion.

① Supervised -

② Training sample provided.

③ Data is labelled.

④ What class does object belong to?

⑤ Known number of classes -

⑥ Used to classify further observations -

⑦ Used to classify features or even
data into one of several
predefined categorical clusters -

⑧ No. of classes is known before

classification as there is

predefined output based on input -

Clustering -

- Clusters are groups such that objects in one group are similar to each other & objects from separate groups are dissimilar.

- Clustering analysis is powerful data mining tool. Cluster is group of objects that belong to same class.

- Clustering is process of partitioning a set of data into sets of meaningful sub-clusters.

e.g. → Hierarchical, Fuzzy, Density based on model based -

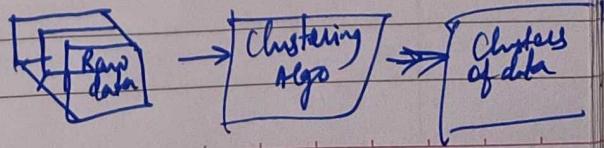
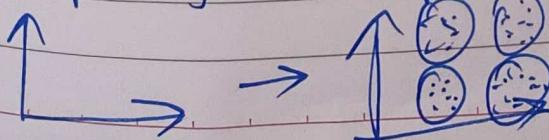
- Desirable properties of cluster algorithms -

- i) Scalability (w.r.t time & space)

- ii) Ability to deal with different types of data

- iii) Minimal domain knowledge requirement to determine parameters.

- iv) Interpretability & Usability



- 13/19
- good clustering technique will produce high quality clusters.
with high intra-class similarities & low inter-class similarity
 - i) Cluster centroid - point whose parameters values are mean of parameters values of all points in cluster.
 - ii) Distance - common metric to see if similarity in population.

(7) Types of data in cluster analysis -

- i) Interval Scaled iii) Nominal, ordinal & ratio variables.
- ii) Binary variables iv) Variable of mixed type

(8) Types of clustering algorithm -

- i) Overlapping
- ii) Exclusive
- iii) Grid based
- iv) Probabilistic
- v) Hierarchical
- vi) Density Based

⇒ Appln of Clustering -

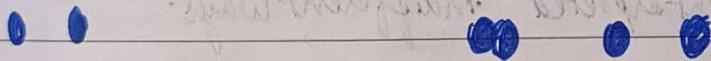
- 1) Marketing → identify distinct group of customer.
- 2) Land Use → similar land identification for easier database.
- 3) Insurance - moto insurance policy holders with high claim.
- 4) Urban planning - house groups based on type, value & geography.
- 5) Seismology - observed earth quake epicenter.
- 6) Biology - classify species.
- 7) Libraries -

→ typical requirements) Properties of Clustering in Data Mining.

- ① Scalability (w.r.t time & space)
- ② Ability to deal with different types of attributes.
- ③ Ability to deal with noisy data.
- ④ Minimal requirement of domain knowledge to determine i/p parameters.
- ⑤ Interpretability & usability.
- ⑥ constraint based clustering
- ⑦ High dimensionality
- ⑧ Discovery of clusters with arbitrary shape.
- ⑨ Incremental clustering & insensitive to order of input records.

→ Types of clusters -

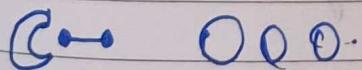
- ① Well separated
 - points in nodes are closer than points in separate ones.
- ② Prototype based
 - objects in cluster is closer to prototype center
 - numerical data prototype center
 - categorical medoid center.



③ Contiguity based →

point increases 2 to

if points in same cluster than in others.



④ Density based

low density high density.

→

high density low density.

→

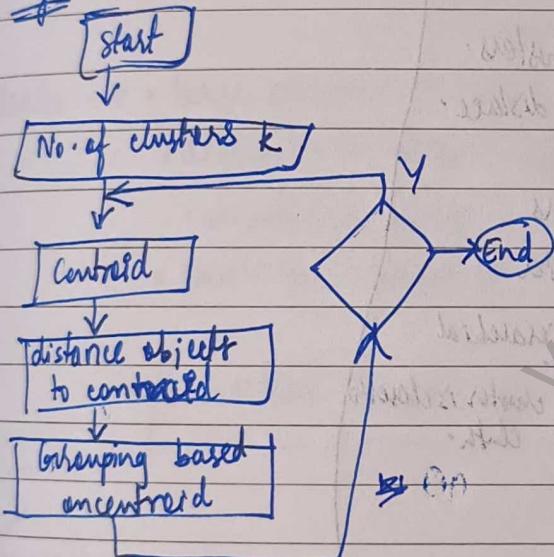
- Desired features of Cluster Analysis -
- ① Scalability
 - ② Minimal i/p parameters.
 - ③ only one scan of dataset
 - ④ Ability to stop & resume.
 - ⑤ Robustness
 - ⑥ Ability to discover diff cluster shapes.
 - ⑦ Different data types.
 - ⑧ Result independent of data i/p order.

- Problems with clustering
- ① Do not address all requirements adequately
 - ② Complexity due to large dataset + high dimensionality
 - ③ effectiveness depends on distance.
 - ④ If obvious distance measures doesn't as most of which isn't always possible, especially, multi-dimensional space.
 - ⑤ Result may be interpreted in different ways.

- * Partition Methods -
- ① Partition Method constructs a partition of a database D of n objects into set of K clusters.
 - ② these are clustering methods used to classify observations into multiple groups based on their similarity.
 - ③ Two types
 - ① Global optimal method. exhaustively enumerates all partitions.
 - ② Heuristic - k-means & k-medoids.

K-Means

Attempts to minimize total squared errors.
worst off $O(n^2)$
sensitive to outliers
convex shape required.
efficient for separate clusters
Centroid Based.

Algo -K-Medoids

- ① Minimizes sum of dissimilarities b/w point labelled to be in classical cluster & points designated as center of cluster.
- ② Less efficient
- ③ not sensitive to outliers
- ④ convex shape not required.
- ⑤ efficient for separate clusters & small datasets.

Representative Object-Basedalgo -

- Select k -points as initial medoids
- Assign all points to closest medoid.
- see if any point is better medoid.
- Repeat $2 \frac{1}{2} \times k$ until medoids don't change.

Strength \rightarrow 1) efficiency $\rightarrow O(nk)$.

- 2) often terminates local optimum
- 3) deterministic annealing to get global solution

Weakness of k-means -

- ① Applicable only when mean is defined
- ② cannot handle noisy data.
- ③ needs to be told in advance
- ④ Not suitable w.r.t non-convex shape.

K-Means



- Heuristic Method
- supervised learning algo.
- solves clustering problems
- groups unlabeled data in diff clusters
- $k \rightarrow$ define no. of predefined clusters needs to be created
- $k=2 \rightarrow$ 2 clusters
- each dataset belongs to one group has similar project
- mini sum of b/w data points clusters
- distance calculation \rightarrow Euclidean distance

Properties \rightarrow

- Always k -clusters
- At least one item cluster
- clusters are non-hierarchical & don't overlap
- every member of cluster is closer to cluster

Adv \rightarrow

- efficient in computation.
- easy to implement.

Disadv \rightarrow

- only when mean is defined
- need to specify k , no. of clusters in advance.
- trouble with noisy data & outliers.
- not suitable to discover cluster with non-convex shapes.

Hierarchical Methods -

- ① Uses distance matrix as clustering criteria.
- ② Widely used data analysis tool.
- ③ Idea is to build binary tree of data that successively merges similar group of points. Visualizing this tree provides useful summary of data.

Adv → • simple to implement • easy & results in hierarchy (more info)
• doesn't need to pre-specify no. of clusters.

Disadv → • large clusters break.

- difficult to handle diff sized clusters & shapes.
- sensitive to noise & outliers
- can't be changed or deleted once done.

Hierarchical clustering

Agglomerative Clust.

- bottom up approach
- each item → own cluster
- identify small cluster
- relatively clustered are merged together
- also known as ABNEES

Divisive clustering

- top down approach
- all item → one cluster
- large cluster
- large clusters are successively divided
- also known as DIANA

4 diff methods -

i) single linkage -

→ minimum method, where shortest distance from any object of one cluster to object of another cluster is considered
 ↳ susceptible to noise / outliers.

$$D(A, B) = \min \{d(i, j) \mid i \in A, j \in B\}$$

ii) complete linkage -

→ maximum method, largest distance is considered.

$$D(A, B) = \max \{d(i, j) \mid i \in A, j \in B\}$$

iii) Average linkage -

→ distance b/w clusters A & B is average of all distances b/w pairs of $i \in A$ & $j \in B$.

→ Mean distance between elements of cluster.

$$D(A, B) = \text{Mean } \{d(i, j) \mid i \in A, j \in B\}$$

→ can handle noise

→ biased towards global

clusters.

iv) Centroid linkage -

→ distance b/w cluster centroid / means.

Evaluation of clustering Value -

- 1) internal validation methods.
 - Quality of clustering method can be evaluated w/o using external info.
 - Two types -
 - 1) cohesion → metric to evaluate how closely the elements are clustered.
 - 2) separation → metric to evaluate level of separation b/w clusters.
 - cohesion & separation do not perform well for density based algorithm.

2) external validation -

- Associated with supervised learning problem.
- additional info required -

(A) silhouette score -

- ① Measures how elements in one cluster are close to points in neighbouring clusters.
- ② Principle - Maximum internal cohesion & Max^m cluster separation.
 - how similar objects is to its own cluster compared to other clusters.
- ③ Finds optimal value of k, during clustering.
- ④ For each element, silhouette value is calculated.
 - well clustered elements = 1
 - poorly clustered elements = -1

* Association rule -

- unsupervised learning method.
- descriptive method \rightarrow discover reln b/w hidden in large dataset.
- used to mining transactions in databases.
- FREQUENTLY BROUGHT TOGETHER
- goal \rightarrow discover interesting relationships among items
- support \rightarrow popularity of product of all product transactions.

$$\text{support (A)} = \frac{\text{No. of transaction in which A appears}}{\text{Total no. of transactions}}$$
- confidence \rightarrow likelihood of purchasing both products A & B.

$$\text{confidence (A} \rightarrow \text{B)} = \frac{\text{No. of transaction includes both A \& B}}{\text{No. of transaction include only product A}}$$
- Itemset \rightarrow set of one or more items.
- Frequent itemset \rightarrow an itemset whose support is greater than or equal to minsup threshold.

* Market Basket Analysis -

- determine what products customers purchase together.
- name \rightarrow idea of customers following all purchases in shopping (market basket) during grocery shopping.
- Association analysis & Frequent itemset mining.
- uses if-then scenario rules.

If item A is purchased then item B is likely to be purchased.
 Rule \rightarrow If {A} Then {B}.

Algorithm \rightarrow Association rule & Apriori

Applicn \rightarrow

- Retail
- Telecommunications
- Banks
- Insurance
- Medical

Lift Ratio -

Ratio of confidence to expected confidence.
expected confidence is confidence divided by frequency of tells us how much better the rule is at predicting result.

$$\text{lift} = \frac{(A \cap B)/A}{(B/\text{total})} = \frac{\text{confidence}}{(B/\text{total})}$$

Leverage = measures diff in $P(x)$ & $P(y)$ appearing together.

Leverage ($x \rightarrow y$) = support ($x \wedge y$) - Support (x) * support (y)

* Appls of MBA

- Retails,

Telecommunications By defn all these items will occur as freq as Banks predict same support count.

Health Medical

* two step process-

i) Find all frequent item sets.

ii) Generate strong association rule from frequent sets.

$X \rightarrow Y$ means transaction containing items from set X tend to contain items from set Y .

* Apriori Algorithm

- classic ML algo

- designed to work on databases having transactions

- aimed to find subsets which are common to atleast a min. number

- solves frequent items problem

- used for mining frequent itemsets & deriving association rules from frequent

- major components are support, confidence, lift

- i) calculates support of item size $k=1$

- ii) Prune candidate set by eliminating items w support less than threshold

- iii) Do same for $k \neq 1$ by joining frequent items

Adv → easy
→ Time Pyme
→ easy on large dataset

Limitation → needs several iterations
→ needs several iterations → education
→ uses uniform minimum support → medical field
→ difficult to find rarely occurring items → forestry
→ some competing approaches focus on
particularity & sampling.