

Types

- ① Supervised Learning: In this, we have to predict the output variable, which is fully structured.
- Where we have input variables & output variables & we use an algorithm to map the input function to the output function.
 - Goal is to approximate the mapping function so well that when we have new inputs, we could predict the output correctly.

Popular Algorithms: linear and non-linear regression, decision trees

- Linear Regression: To predict a continuous value.
- Random Forest: To predict a discrete value.
- Support Vector Machine: To classify the data.

② Unsupervised Learning.

- Where we only have input data x & no corresponding output variable, i.e., we know x but not y .
- Goal is to try to model the underlying structure or distribution in the data in order to learn more about data.
- Thus, we have an input but not output. The algorithm ~~predicts~~ detects patterns based on innate characteristics of input data.

Popular Algorithms

- Apriori Algorithm
- K-Means
- Hierarchical Clustering.

③ Reinforcement Learning

- It allows software agents and machine to automatically determine the ideal behaviour within a specific context to optimize its performance.
- It is about interaction b/w environment & learning agent. Learning agent leverages two mechanism, exploration & exploitation.

Exploration = when learning agent acts on trial & error basis.

exploitation = when it acts based on the knowledge gained from the environment. It is referred to as exploitation.

Environment rewards the learning agent for correct actions, which is reinforcement signal & based on the reinforcement signal, agent improves its environment knowledge to select the next action.

Algorithms

- Classification Algorithm.

A supervised learning approach in which the computer program learns from the input given to it & uses the learn this learning to classify new observations.

- Anomaly Detection Algorithm.

It is the identification of unusual patterns in the dataset.

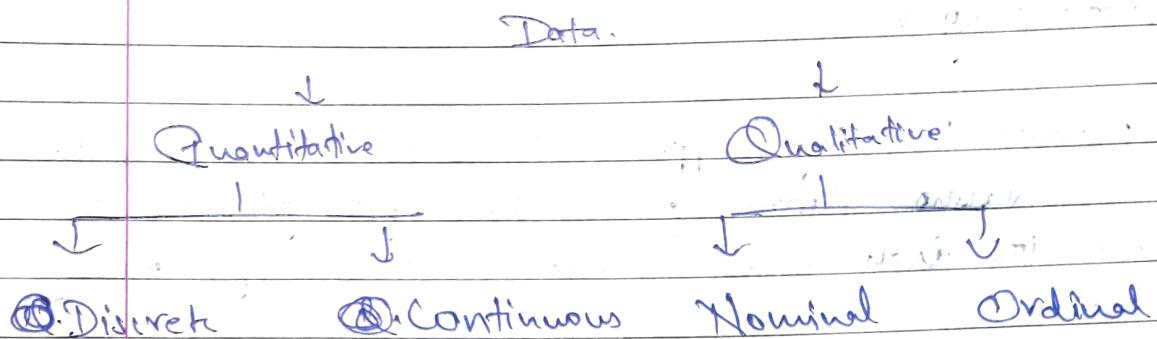
- Clustering Algorithm

Grouping the data sets based on some conditions.

- Regression Algorithm.

Data itself is predicted. Allows making prediction by learning the relationship between the features of your data & observed continuous valued response.

Statistics.



- Deals with characteristics & descriptors that cannot be easily measured but can be observed effectively.

Nominal → Data with no inherent order
 Ordinal → Data with an ordered series.

Quantitative Data

- Deals with nos. of things that could be measured objectively.

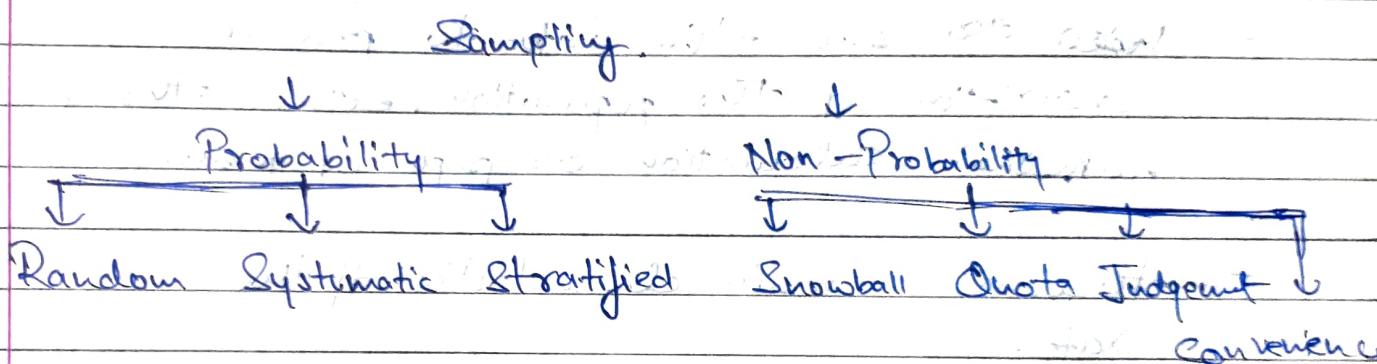
Discrete = AKA Categorical data, holds finite no. of possible values. (No. of students in class)

Continuous = Data can hold infinite no. of possible values.
 (Weighting scale.)

Population - Collection of or set of individuals or ~~and~~ objects or events whose properties are to be analyzed.

Sample - A subset of population

Sampling - Is a technique that deals with the selection of individual observations within a population. Used to infer statistical knowledge about population.



Random Sampling : A sampling technique in which samples from a large population ^{is a} are chosen using the theory of probability.

Random Sampling : Each member of the population has equal ~~chance~~ chance to be selected in the sample.

Systematic Sampling : Every n th record is chosen to be a part of the sample.

Stratified Sampling : A stratum is used to form sample. A stratum is a subset of the population that shares at least one common characteristic.

Random Sampling is used to select a sufficient no. of subjects from each stratum.

Types of Statistics

① Descriptive Statistics

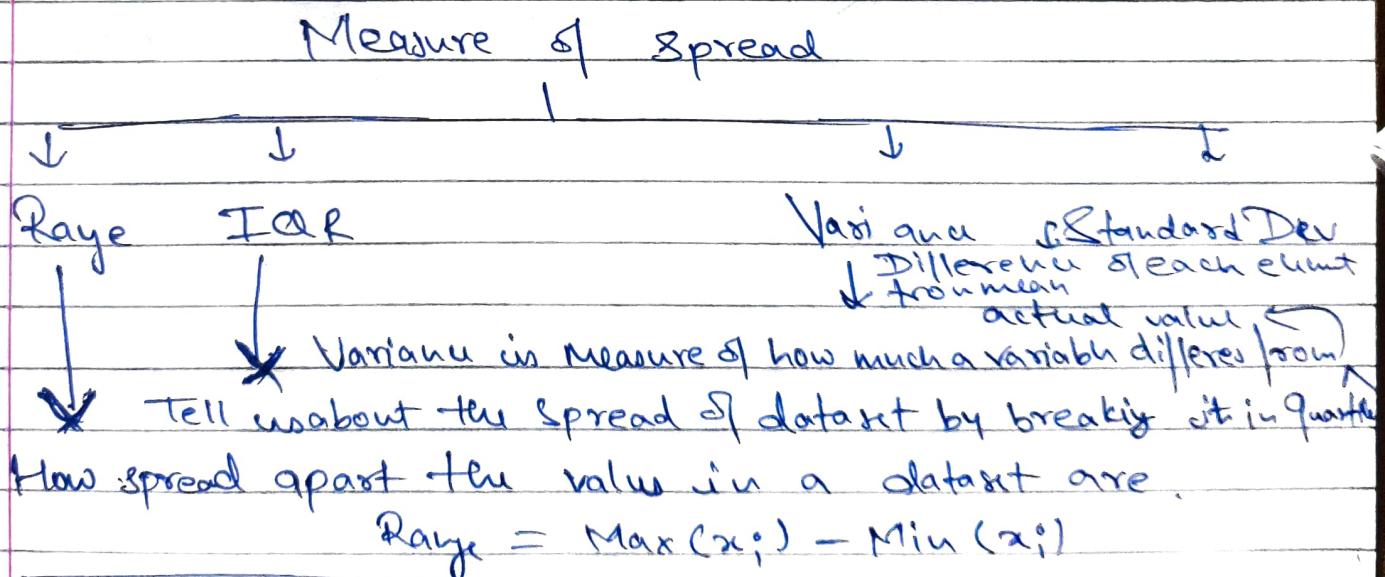
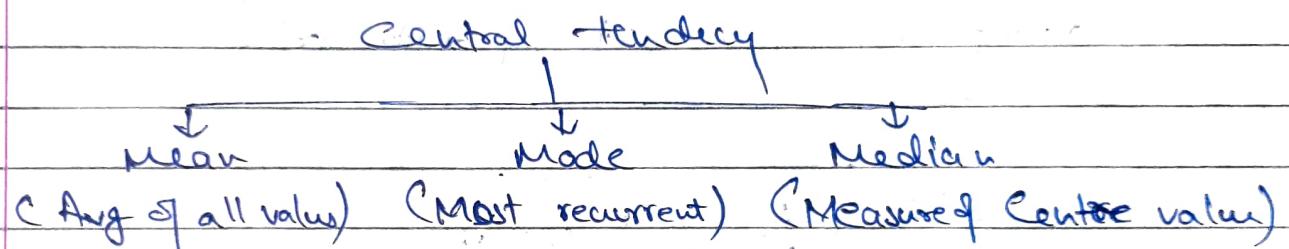
- used to use the data to provide description of the population, either through numerical calculation or graph or table.
- focuses on main characteristics of the data.

② Inferential Statistics

- It makes inferences & predictions about a population, based on a sample of data items taken from the population in question. It is not unique.
- It generalizes the large data set & applies probability to draw conclusion.

Descriptive Statistics

- It is a method to describe & understand the features of a specific dataset by giving short summaries about the sample & measure the data.
- Descriptive statistics are broken down into two categories
 - Measure of central tendency (rep. Summary of Data set)
 - Measure of variability (spread)



Quartiles

$Q_1 \Rightarrow 25\text{th}$ value

$Q_2 \Rightarrow 50\text{th}$ value

$Q_3 \Rightarrow 75\text{th}$ value

$IQR = Q_3 - Q_1$

Standard Deviation = $(x_i - \bar{x})^2$

\bar{x} = individual data points

n = No of data points

\bar{x} = mean.

$$\text{Variance. } S^2 = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n}$$

Kash

Page No.:	
Date:	

$\mu = \text{Mean of Population}$ $N = \text{no of DP in Population}$

$\bar{x} = \text{Mean of Sample}$ $n = \text{no of DP in Sample}$

\Rightarrow Population Variance is avg of Squared Deviation

$$\sigma^2 = \frac{1}{N} \sum_{i=1}^N (x_i - \mu)^2$$

\Rightarrow Sample Variance is also avg of squared differences from the mean.

$$s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$$

Standard Deviation is measure of dispersion of a set of data from its mean.

$$\text{standard deviation } \sigma = \sqrt{\frac{1}{N} \sum_{i=1}^N (x_i - \mu)^2}$$

Information Gain & Entropy.

Entropy : Measure of uncertainty in data.

$$H(S) = - \sum_{i=1}^N p_i \log_2 p_i$$

S = Set of dataset all instances in the dataset.

N = Number of distinct class values

p_i = Probability event. different class value

Information Gain : How much information a particular feature / variable gives us about final outcome.

$$\text{Gain}(A, S) = H(S) - \sum_{j=1}^{|S_j|} \frac{|S_j|}{|S|} \cdot H(S_j)$$

$H(S)$ - Entropy of whole dataset

$|S_j|$ - No. of instances with j value of a attribute A

$|S|$ - Total no. of instances in dataset A.

\vee - set of distinct value of an attribute A.

$H(S_j)$ - Entropy of subset of instance for an attribute A.

$H(A, S)$ - Entropy of an attribute A.

Use case - Decision Tree.

\Rightarrow Confusion Matrix

- It is a matrix used to describe the performance of a classification model or classifier. It uses test data for which the true values are known to calculate various additional metrics.
- Represents a tabular representation of Actual vs Predicted values

$$\text{Accuracy} = \frac{TP + TN}{TP + TN + FP + FN}$$

Actual & Predicted

		Actual		Predicted	
		No	Yes	No	Yes
Actual	No	TP	FP	TN	FN
	Yes	FN	TP	FP	TN

Predicted	0	1
0	TP	FN
1	FP	TN

Probability

$$P(x) = \frac{\text{Desired outcome}}{\text{Total}}$$

Terminologies:

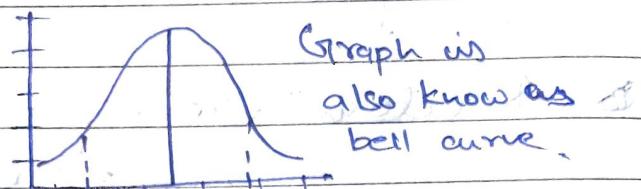
- **Random Experiment:** Experiment for which the outcome cannot be predicted with certainty.
 - **Sample Space:** The entire possible set of outcomes of a random experiment is called as the sample space (S) of that experiment.
 - **Event:** Outcomes of an experiment
- Event Type:**
- Disjoint: Don't have a common outcome.
 - Non-Disjoint: Can have common outcome

Probability Distribution

① Probability Density Function.

→ It is the equation describing a continuous probability distribution.

Properties:



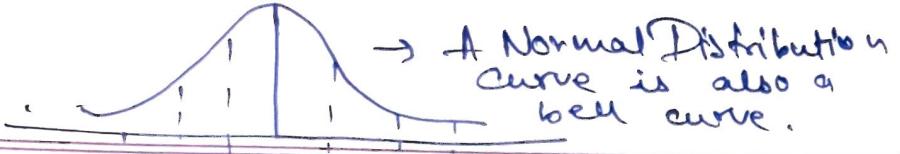
Graph is also known as bell curve.

- ① Graph of PDF is continuous in range
- ② Area bounded by the curve of density function & the x axis is equal to 1 (area below the curve)
- ③ Probability that a random variable assumes a value between a & b is equal to the area under (PDF) bounded by a & b.

② Normal Distribution (Gaussian Distribution)

- Normal Distribution is a probability distribution that associates the ^{normal} random variable x with cumulative properties
- It denotes the symmetric properties of the mean. Idea behind is that the data near the mean occurs more frequently than the data away from the mean, or data near the mean represents the entire dataset
- Graph of Normal Distribution depends on two factors: the Mean & the Standard Deviation.

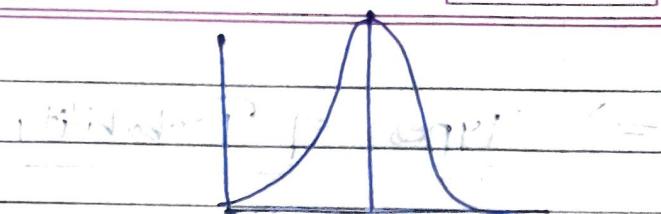
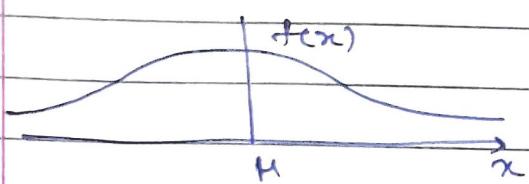
Mean: determines location of center of the graph
 Standard deviation: determines height of the graph.



Page No.:
Date:

KOBY

$$\mu - 3\sigma \quad \mu - 2\sigma \quad \mu - \sigma \quad \mu \quad \mu + \sigma \quad \mu + 2\sigma \quad \mu + 3\sigma$$

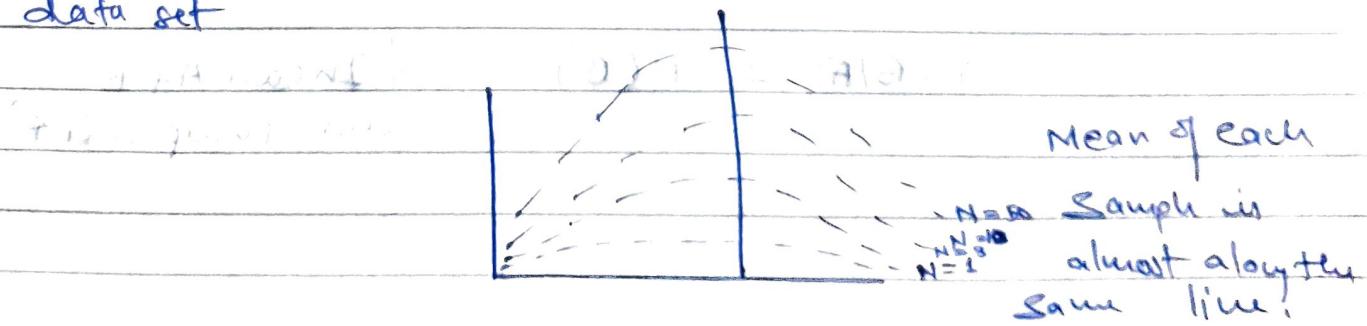


If Standard Deviation is large, curve is short & narrow if Standard Deviation is small, curve is tall & wide

∴ Standard Deviation is the spread of the data

⇒ Central Limit Theorem (only true for large dataset as small dataset has a lot of variations due to sample size factors)

- It states that the sampling distribution of the mean of any independent random variable will be normal or nearly normal if the sample size is large enough.
- If we have a large population & divided it into many samples, then the mean of all the samples of the population will be almost equal to the mean of the entire population. Hence, mean is normally distributed. Holds true only for a large data set



The accuracy or resemblance depends on two factors:

- ① No. of sample points considered
- ② Shape of underlying population

→ shape depends on

- ① Standard Deviation
- ② Mean of sample.

⇒ Types of Probability

① Marginal Probability.

→ Probability of occurrence of single event.

$$\text{No. of hearts in a deck of card} = \frac{13}{52}$$

② Joint Probability.

→ Measure of two events happening at same time.

$$\text{Probability of getting Ace of hearts} = \frac{1}{52}$$

③ Conditional Probability, happens if one event has already occurred.

→ Probability of an event based on occurrence of a previous event.

An event B will occur given that A has

already occurred with some probability.

(Increasing A & B are dependent)

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

$$P(B|A) = P(B)$$

(Increasing A & B are independent)

⇒ Baye's Theorem

- Shows the relationship between one conditional probability & its inverse.
- It is the probability of an event occurring based on prior knowledge of conditions that might be related to the same event.

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

P(A|B) is referred to as posterior probability which means the probability of occurrence of A given B.
 P(B|A) is referred to as likelihood ratio which measures the probability of occurrence of B given A.
 P(A), is referred as Prior probability which represent the actual probability distribution of A.

Cq: we have 3 bowls.

Bowl A = 2 blue + 4 red

Bowl B = 8 blue + 4 red

Bowl C = 1 blue + 3 red.

1 ball is drawn from each. Probability of drawing a blue from A given 2 blue is drawn from B & C, are drawn in total.

Sol

Event 1 = Pickly a blue ball from bowl A. $\Rightarrow A$

Event 2 = Pickly 2 blue balls. $\Rightarrow X$.

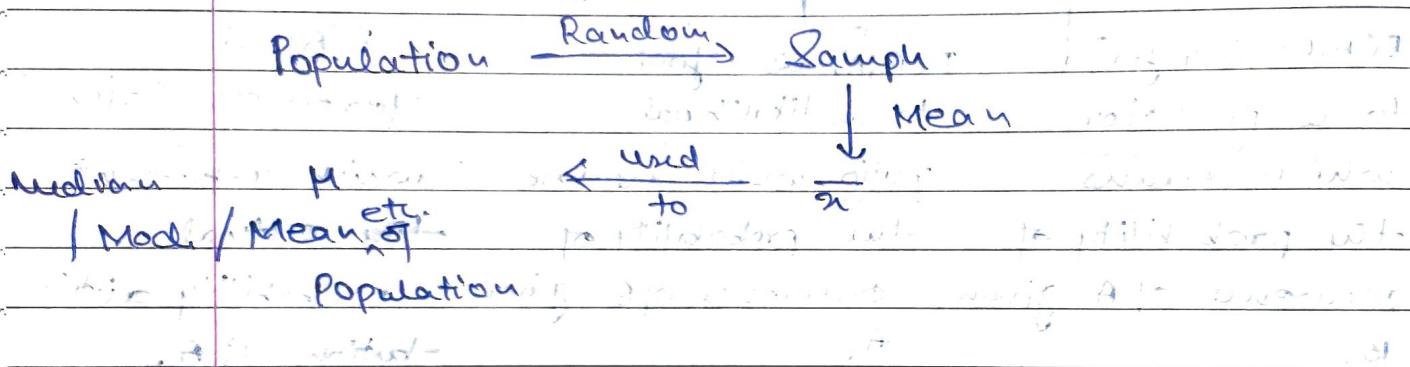
$$P(A/X) = \frac{P(A \cap X)}{P(X)}$$

$$P(A \cap X) = \frac{{}^6C_2 + {}^{12}C_2}{{}^{18}C_2}$$

$$P(X) = \frac{{}^6C_1 \times {}^{12}C_1 + {}^{12}C_1 \times {}^4C_1 + {}^4C_1 \times {}^6C_1}{{}^{22}C_2}$$

\Rightarrow Inferential Statistics & Point Estimation

Point Estimation is concerned with the use of the sample data to measure the single value that serves as an approximate value or the best estimate of an unknown population parameter.



Estimator = function of the sample used to find out the estimate.

$$\text{Ex: } \bar{x} \text{ = Sample mean}$$

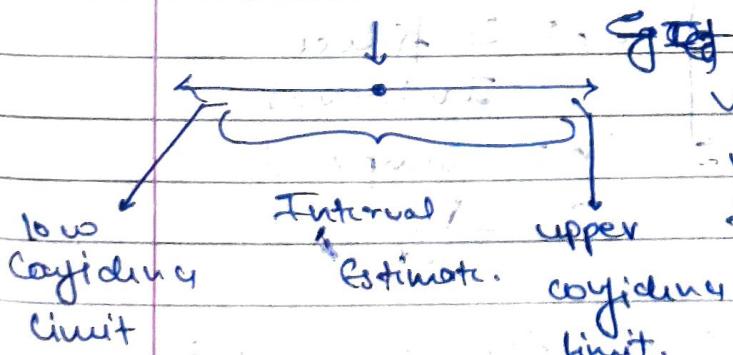
Estimate = pre-calculated value of using the function
 $\text{Ex: } \mu = \text{population mean.}$

Method to find Estimator.

- ① **Method of Moments:** form an equation in the sample statistic by comparing it with the population parameter. Then solve the equation to get the estimator.
- ② To find out we take the known facts about the population. The population parameters are extended to the sample of size n. Then predict the known parameters such as mean, mode etc.
- ③ Maximum of Likelihood: estimating a parameter.
- ④ Baye's Estimator
- ⑤ Best Unbiased Estimators.
- ⑥ Interval Estimation Method (Imp)

Interval Estimate: An interval estimate is a range of values used to estimate a population parameter with some degree of confidence.

Point Estimate



Example: If we want to find a value (Mean, mode etc) of a population, we set a range & at that range the value lies.

So, instead of predicting a point we predict a range in which a value may occur.

- Difference b/w point estimate & actual population parameter value is Sampling Error
- When μ is estimated, Sampling error = $\bar{x} - \mu$

Confidence Interval

- It is the measure of ~~your confidence about~~ ^{the probability that} that the interval estimated contains the population parameter.
- It is used to describe the amount of uncertainty associated with a sample estimate of a population parameter.
- A range of values so constructed that there is a specified probability of including the true value of a parameter within it to various levels of certainty.

Margin of error - E

- For a given level of confidence it is the greatest possible distance b/w the point estimate & the value of parameter it is estimating.

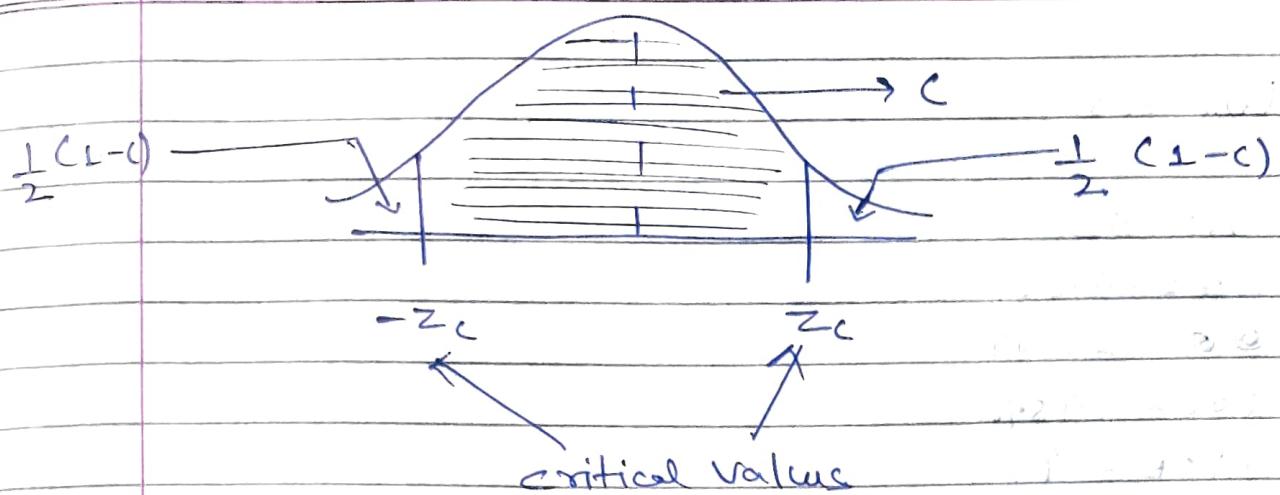
$$\text{Margin of error } E = Z_c \frac{\sigma}{\sqrt{n}}$$

$Z_c = \text{Confidence Interval}$

$n = \text{Sample size}$

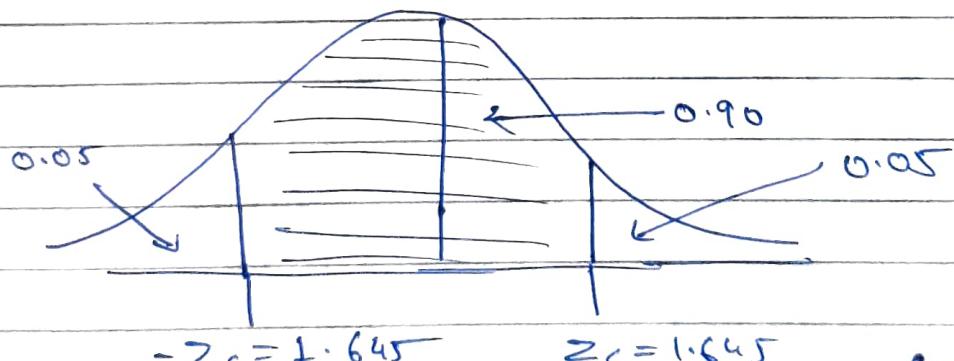
$\sigma = \text{S.D.}$

The level of confidence C is the probability that the interval estimate contains the population parameter.



c is the area beneath the normal curve b/w the critical values. Z-score can be calculated using standard normal table.

4 level of confidence = 90%. i.e. 90% confident that interval contains population parameter.



correspondingly
Z score = ± 1.645

Sequence :

- ① Sample statistic (used to calculate population parameter)
- ② Select confidence level.
- ③ Margin of error
- ④ Specify confidence interval

Use Case

$$n = 32$$

$$\bar{x} = 74.22$$

$$\sigma = 23.44$$

$$(Z) L.O.C = 95\%$$

$$(E) M.O.E = ?$$

$$Q \Rightarrow E = Z \cdot \frac{\sigma}{\sqrt{n}}$$

$$E = 1.96 \cdot \frac{23.44}{\sqrt{32}} \approx 8.12$$

for 2.7 estimation L.O.C $\pm \sqrt{32} \cdot 1.96 = 10.27$

Hypothesis Testing

- Used to check whether hypothesis is accepted or rejected.

 Steps :

- ① State the hypothesis : Stating null & alternative hypothesis.
- ② Formulate an analysis plan: This stage involve construction of analysis plan.
- ③ Analyse Sample Data: Calculation & interpretation of the test statistic as described in the analysis plan.
- ④ Interpret Result: The application of the decision rule described in the analysis plan.

NULL HYPOTHESIS : Result is no different from (H_0) assumption

ALTERNATIVE HYPOTHESIS : Result disproves the assumption.

Eg: 4 boys Nick, John, Bob, Harry are picked randomly each day. What is the probability that John is not cheating.

Sol

$$P(\text{John picked for a day}) = \frac{1}{4} = 0.25$$

$$P(\text{John picked for 3 days}) = \frac{1}{4} \times \frac{3}{4} \times \frac{3}{4} = 0.42$$

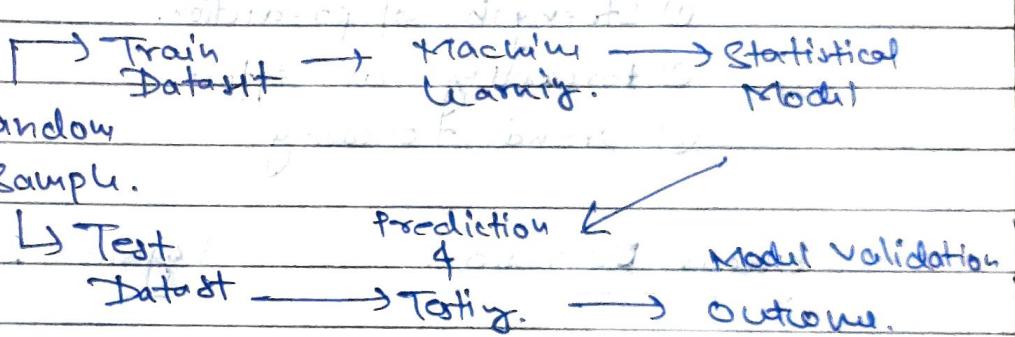
$$P(\text{John ————— 12 days}) = \left(\frac{1}{4}\right)^{12} = 0.032.$$

As, $P(\text{Event}) < 5\% \text{ (threshold)}$; John is cheating.

Super Vised Learning

→ Where we have input variable (x) & output variable (y)
 & we use the algorithm to map the function
 from the input to output.

Historic → Random
 Data. Sample.



New: → Model → Prediction.

Algorithms

- ① Linear Regression (Estimate real values)
- ② Logistic Regression (Estimate discrete/binary etc values)
- ③ Decision Tree (Classification)
- ④ Random forest (Ensemble of DT)
- ⑤ Naive Bayes Classification (Classification)

Regression

→ It is a form of predictive modelling technique which investigates relationship between a dependent & independent variable.

Uses:

- ① Strength of prediction.
- ② forecasting an effect.
- ③ Trend forecasting.

Linear Regression

→ Data is modelled using straight line

→ $y = mx + c$, we find correlation b/w x &

y.

→ Data is modelled via sigmoid function.

→ for continuous Variables → Categorical Variable.

→ Value of variable as an output

→ Accuracy is measured by r-squared, Adjusted R squared etc.

→ Maps continuous x to continuous y

Logistic Regression

→ The probability of some obtained event is represented as a linear function of a combination of predictor variables.

→ Data is modelled via sigmoid function.

→ Categorical Variable.

→ Probability of occurrence of an event as an output.

→ Accuracy, precision, Recall, f1 score, ROC curve, confusion Matrix, etc.

→ Maps continuous x to categorical y

Linear Regression is used in: sales, with

cost, profit, with different factors.

- Evaluating Trends & sales forecast.
- Analyzing effect of price changes
- Assessment of risk in financial services & insurance domain.

y-axis = Dependent Variable

x-axis = independent variable.

Exq

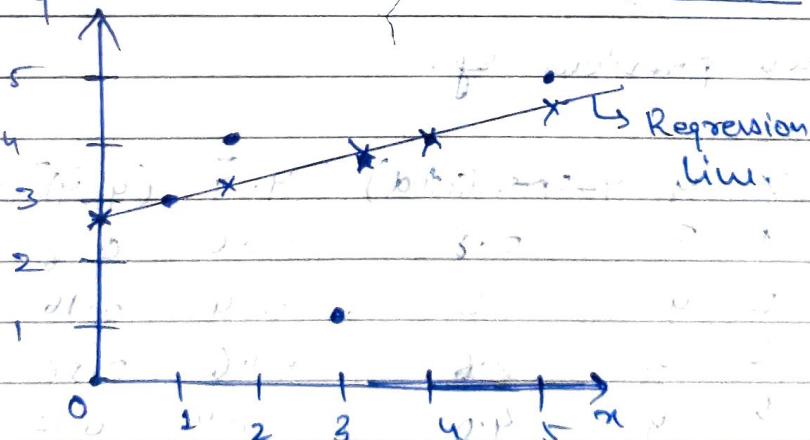
$$\begin{array}{|c|c|c|c|c|} \hline X & Y & x-\bar{x} & y-\bar{y} & (x-\bar{x})^2 & (x-\bar{x})(y-\bar{y}) \\ \hline 1 & 3 & -2 & -0.6 & 4 & -1.2 \\ \hline 2 & 4 & -1 & 0.4 & 1 & -0.4 \\ \hline 3 & 2 & 0 & -1.6 & 0 & 0 \\ \hline \end{array}$$

$$\text{Sum of } x = 6, \text{ Sum of } y = 9, \text{ Sum of } x^2 = 14, \text{ Sum of } xy = 10$$

$$\frac{5}{3}, \frac{5}{3}, \frac{14}{10}, \frac{10}{10}$$

$$\frac{5}{3} = 3.6, \text{ Slope } m = \frac{10}{14} = 0.714, \text{ Intercept } c = \frac{9 - 0.714 \times 5}{10} = 2.8$$

Mean



Linear regression eq $\Rightarrow y = mx + c$. $\therefore m = 4/10$

$$\therefore 3.6 = 0.4 \times 3 + c$$

slope, y-intercept

$$m = \frac{\sum (x-\bar{x})(y-\bar{y})}{\sum (x-\bar{x})^2}$$

$x - \bar{x}$ = difference (distance)

of point x. mean of x

$y - \bar{y}$ = similar to x.

\therefore Predicted values for when $x = \{1, 2, 3, 4, 5\}$

$$y = 0.4x + 2.4 = 2.8$$

$$y = 0.4x + 2.4 = 3.2$$

etc.

- Now, we calculate the distance b/w actual & predicted value, and reduce the error/distance.
- The line with least error will be best fit line.

To check the goodness of fit we use R squared method.

- R squared value is a statistical measure of how close the data are to the fitted regression line.
- Also called coefficient of determination or coefficient of multiple determination.

Eq: From previous Eq.

	x_i	y_i	$y - \bar{y}$ (predicted)	$(y - \bar{y})^2$	$(y_p - \bar{y})(y_p - \bar{y})$
1	3	2.8	-0.6	0.36	-0.8 0.6
2	4	3.2	0.4	0.16	-0.4 0.16
3	2	3.6	-1.6	2.56	0 0
4	4	4.0	0.4	0.16	0.4 0.16
Σ	$\frac{5}{2}$	4.4	1.4	1.96	0.8 3.64
		$\frac{3.5}{2}$		$\frac{5.2}{2}$	$\frac{1.6}{2}$

Predicted Value

$$R^2 = \frac{\sum (y_p - \bar{y})^2}{\sum (y - \bar{y})^2}$$

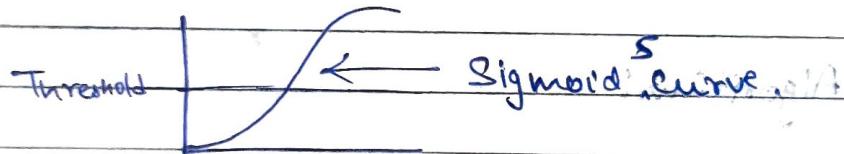
$$R^2 = \frac{1.6}{5.2}$$

$R^2 \approx 0.3$, It is not a good fit.

Logistic Regression.

- It produces results in a binary format which is used to the outcome of a categorical dependent variable.
- Outcome should be discrete / categorical such as:

0 1	High / low	f. / y N
-------	------------	------------
- Since the value of y will be b/w 0 & 1, ~~linear like has to be clipped by 0 & 1.~~
- A Sigmoid function caps all other values from $-\infty$ to ∞ to 0 & 1



- for values on curve, we use threshold value to either say they are 1 or 0.
- Eg: If curve gives value of 0.8, it is > 0.5 hence, will be given as 1 or rounded to 1.
- Logistic Reg eq is derived from straight line eq.

Eq of st line

$$Y = C + B_1 X_1 + B_2 X_2 \dots \rightarrow \text{Range } -\infty \text{ to } \infty.$$

Since, the Logistic regression has y p value b/w 0 & 1. we will $1-y$.

$$\therefore \frac{y}{1-y}$$

Use cases

- ① Weather Prediction (Rain or not)
- ② Determine illness

Classification

→ Process of identifying & dividing the dataset into different categories or group by adding labels

Algorithms:

Decision Tree

Random Forest

Naive Bayes & NB with pos. ratio & KNN.

Logistic Regression

Linear Regression

SUM

Decision Tree: graphical representation of all the possible solution to a decision.

Random forest :- Built multiple decision trees & merges them together.

- More accurate & stable
- Corrects decision trees mistakes of overfitting
- Trained with bagging method.

Naive Bayes:- Based on Bayes Theorem

- Assumes that the presence of a particular feature in a class is (is) unrelated to the presence of any other feature.

KNN - Stores all the available cases & classifies new cases based on a similarity measure.

- While K in KNN algorithm we take which to take vote from.

Decision Tree Using KART Algorithm

- Machine Learning Model

- ① Root Node - Parent node (Root Node)
- ② Leaf Node - End of the tree
- ③ Splitting - Dividing of structure
- ④ Pruning - Removing (Subnode) -
not giving better result

→ Entropy: Measure of impurity (impurity with respect to target attribute)

$$\text{Entropy}(S) = -P(\text{yes}) \log_2 P(\text{yes}) - P(\text{no}) \log_2 P(\text{no})$$

ratio

S = No of sample space.

$$P(S) = 0.5 \text{ Entropy} = 1.0 \text{ since } 0.5$$

$$P(S) = 0.0 \text{ Entropy} = 0.0 \text{ since } 0.0$$

worst stage of it

→ Information Gain: Measure the reduction in entropy

Decides which attribute should be selected at the decision node

$$\text{Information Gain} = \text{Entropy}(S) - [(\text{Weighted Avg}) \times \text{Entropy each feature}]$$

Ex → Temp, outlook dataset.

14 instances, 9 yes, 5 no.

$$E(S) = -P(\text{Yes}) \log_2 P(\text{Yes}) - P(\text{No}) \log_2 P(\text{No})$$

$$E(S) = -\frac{9}{14} \log_2 \frac{9}{14} - \frac{5}{14} \log_2 \frac{5}{14}$$

$$E(S) = 0.94$$

Selection of Root Node:

lets take column Outlook.

Sunny → 2Y 3N Overcast → 4Y 1N Rainy → 3Y 2N

$$E(\text{Sunny}) = -\frac{2}{5} \log_2 \left(\frac{2}{5}\right) - \frac{3}{5} \log_2 \left(\frac{3}{5}\right) = 0.971$$

$$E(\text{Overcast}) = 0$$

$$E(\text{Rainy}) = -\frac{3}{5} \log_2 \left(\frac{3}{5}\right) - \frac{2}{5} \log_2 \left(\frac{2}{5}\right) = 0.971$$

$$\text{Information Gain} = \frac{5}{14} * 0.971 + \frac{4}{14} * 0 + \frac{5}{14} * 0.971 = 0.693$$

from Outlook

$$= 0.693$$

$$\text{Information Gain} = 0.94 - 0.693 = 0.247$$

lets take column Rainy & Temp:

Hot → 2Y 2N Mild → 4Y 2N Cool → 3Y 1N.

$$E(\text{Hot}) = -\frac{2}{4} \log_2 \left(\frac{2}{4}\right) - \frac{2}{4} \log_2 \left(\frac{2}{4}\right) =$$

Similarly, performed for all the attributes.

⇒ Random forest

- Capable of performing regression & classification
- Randomly select n variables out of N variables in the dataset.
- Create decision trees, one per variable
- Create different decision trees with different combination of variables.
- Use majority voting to get final result

=) KNN Algorithm (Lazy Learner)

- It stores all the available cases & classifies the new data or can based on similarity measure.
- k denotes no of voting neighbors / nearest members to be checked.

User Case:

Recommendation System.

File Separation.

Distance b/w the new data point & test neighbour is calculated by:

① Euclidean Method.

② Manhattan Method.

Steps:

① Handle Data (if any) (Illustration)

② Similarity (calculating distance b/w two data instances)

③ Neighbours (locate k most similar data instances)

④ Response (Generate a response from a set of data instances)

⑤ Accuracy (summarize the accuracy of prediction.)

⑥ Main (tie all together)

⇒ Naive Bayes

Based on Bayes theorem.

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

Use cases with examples

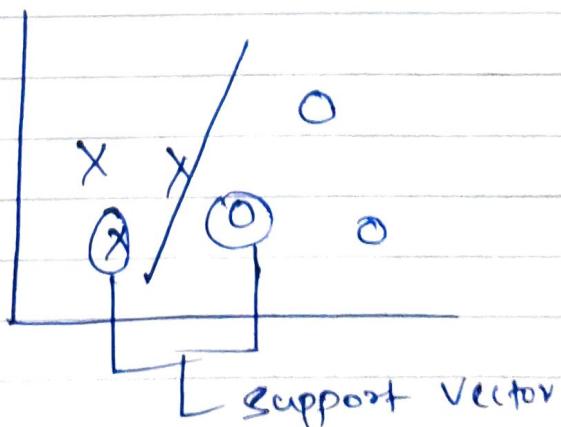
- Text Classification
- Spam Filtering
- Medical Diagnostic

3 Models on ~~Naive Bayes~~ Scikit Learn.

- ① Gaussian (used in classification, assumes feature follow Normal Distribution)
- ② Multinomial (used for discrete count)
- ③ Bernoulli (used for feature vectors are binary)

=) SVM

- Supervised classification method that separates data wby hyperplane.
- Classification & Regression Algorithm.
- Can be used to separate non-linear data with the help of kernel function
- Aim is to draw a hyper plane to separate 2 classes.
- First we draw a random ~~data~~ hyperplane.
- Then we check the distance with the nearest data points aka, support vectors.
- An optimum hyperplane will have maximum distance from each of these support vectors.
- Distance b/w hyperplane & support vector is called margin.



In case of non linear data?

- In this we ~~do~~ transform the 2D space into 3D space.

- Use Case

- Cancer Classification

Classification of cancer in 2D space is very difficult

but if we transform it into 3D space it will be easier

Classification of cancer in 3D space is very difficult

Classification of cancer in 2D space is very difficult

Classification of cancer in 3D space is very difficult

⇒ Unsupervised Learning

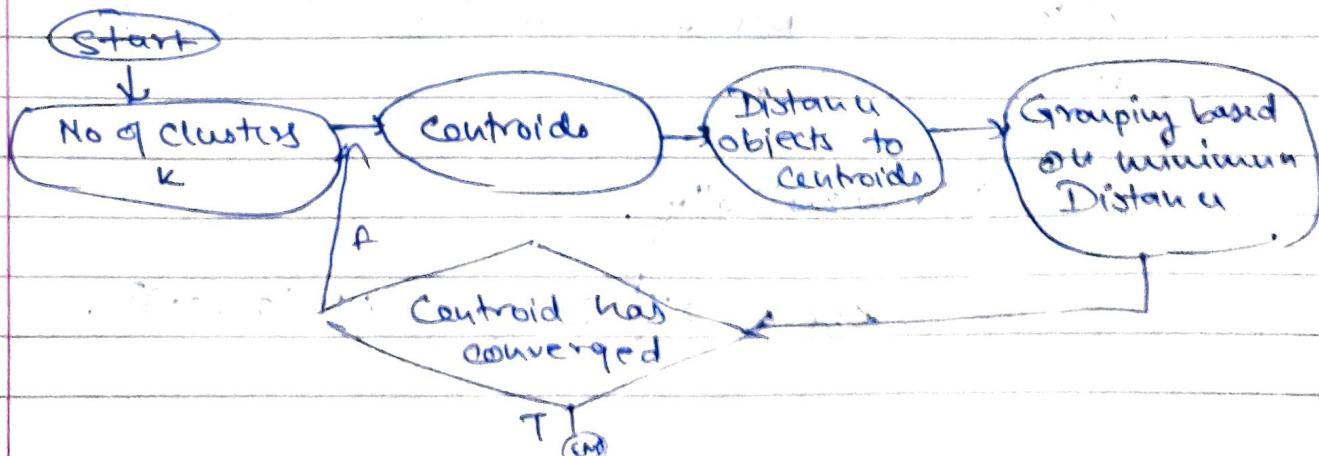
- It is a type of machine learning algorithm used to draw inferences from datasets consisting of input data without labeled responses.
- Training data is collection of information without any label (label) and we need to find out what is the nature of data.

⇒ Clustering: To group the objects with similar results with it response.

- Dividing the dataset into groups of similar datapoints without prior knowledge.
- Types: -
 • Exclusive Clustering → e.g. K-means clustering
 • Overlapping Clustering → Fuzzy means clustering.
 • Hierarchical Clustering.

⇒ K-Means Clustering: It uses with following steps:

- Algorithm that groups similar datapoints.



Steps for calculations:

- Decide k: No of clusters to be made.
- Then we provide centroids of all the clusters.
- Also calculate Euclidean Distance of the points from each centroid & assigns the point to the closest cluster.
- Next the centroids are calculated again, when we have one new cluster.
- The distance of the points from the center of the clusters are calculated again & points are assigned to the closest cluster.
- Again new centroid of this cluster is calculated.
- These steps are repeated until we have a repetition in centroid or new centroids are very close to the previous ones.

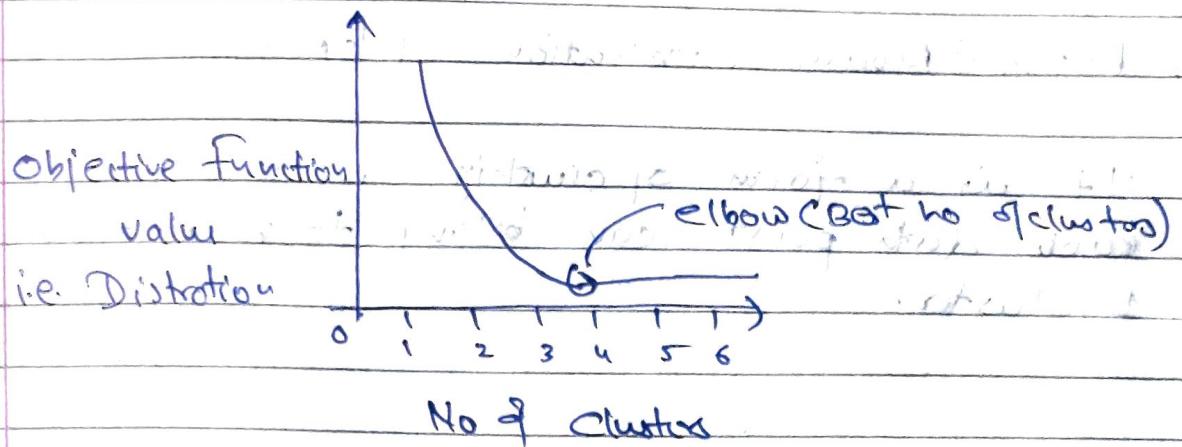
How to decide no. of Clusters (k).

① The Elbow Method.

- Compute the sum of squared error for some values of k. The SSE is defined as the sum of the squared distances of all the points of each member of the cluster from its centroid.

Mathematically,

$$SSE = \sum_{i=1}^k \sum_{x \in C_i} \text{dist}(x, c_i)^2.$$



Cons of k-means:

- Can't handle ~~see~~ & noisy data & outliers.
- All items are forced to cluster.

\Rightarrow Fuzzy C Means (extension of k-means)

- It is a form of clustering in which each data point can belong to more than 1 cluster.

Drawback

- Need to define
- Need to define membership cutoff value.
- These are non deterministic.

\Rightarrow Hierarchical Clustering

- Alternative approach which builds a hierarchy from the bottom up, & does not require to specify no of clusters.
- 2 types :

Agglomerative & Division



Builds from
bottom level

Has all points
Initially in one
cluster.

Disadvantage

- Too slow on big data.

=) Market basket Analysis

- Two algorithms:

→ ① Association Rule Mining. ~~It's~~ It's used to see
→ Shows how items are associated with each other.

3 ways to measure association:

① Support: Gives fraction of transactions which

contains both items A & B.

Support = $\frac{\text{freq}(A, B)}{\text{Total number of transactions}}$

② Confidence: Given how often (often) the items A & B

occur together, given no item

of its time, whether A occurs.

(A) joint occurrence of B given A.

Confidence = $\frac{\text{freq}(A \cap B)}{\text{freq}(A)}$

(B) joint occurrence of A given B.

③ Lift: It indicates the strength of a rule over the

random co-occurrence of A & B

(A) joint occurrence of A & B

Lift = $\frac{\text{support}(A \cap B)}{\text{supp}(A) \times \text{supp}(B)}$

→ ② Apriori Algorithm: It uses frequent item sets to generate association rules. It is based on the concept that a subset of a frequent item set must also be frequent item set.

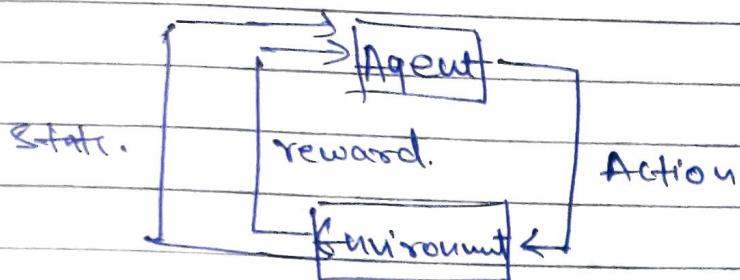
⇒ Reinforcement Learning

- Reinforcement Learning is a type of ML where a agent learns to behave in their environment by performing actions & seeing the results.
- 2 components:
 - Agent (Reinforcement Algo) - Learns from Trial & Error.
 - Environment (The world we live in) : response
- Action (A) - All possible states which Agent can take
- State (s) - Current Condition returned by the environment
- Reward (R) - An instant return from the environment to approve the last action.
- → Policy (π) - Approach the agent determines the next action based on the current stat.
- → Value (V) - Expected long term return with discount as opposed to short term return (R)
- → Action Value (\hat{Q}) - Similar to Value, except, it takes an extra parameter, the current action.

- Reward Maximization theory states that, a RL Agent must be trained in such a way that he takes the best action so that the reward is Maximum.

- Exploration & Exploitation → It is about exploring & capturing more information about an environment.
 - ↓
 - It is about using the already known information about exploited info to highlight rewards

Markov Decision Process : Approach for mapping a solution in reinforcement learning



- Uses Bellman's Equation for decision making.

$$V(s) = \max_a [R(s, a) + \gamma V(s')]$$

Value of being in a particular state (room)

Reward function, outputs a reward value.

action.

Discount factor

State at which robot goes from s

• State