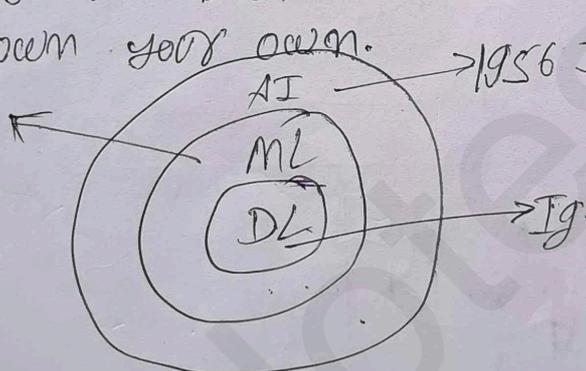


IML

*Artificial Intelligence AI is man made thinking power, it is a branch of computer science by which we can create intelligent machine which can behave, think like human and ~~the~~ are able to make decision. AI exists when a machine can have human based skills, like learning, reasoning, problem solving. It is a process of building intelligence machine from vast volume of data. AI uses complex algo and methods to build machine that can make decision on your own.

1959
Arthur
Samuel



→ 1956 John Mc Carthy

→ Fig.

Vast
volume
of data

During ~~world war~~ second world war the first computer was developed to break the government communication. in 1950 Turing published a paper in Journal

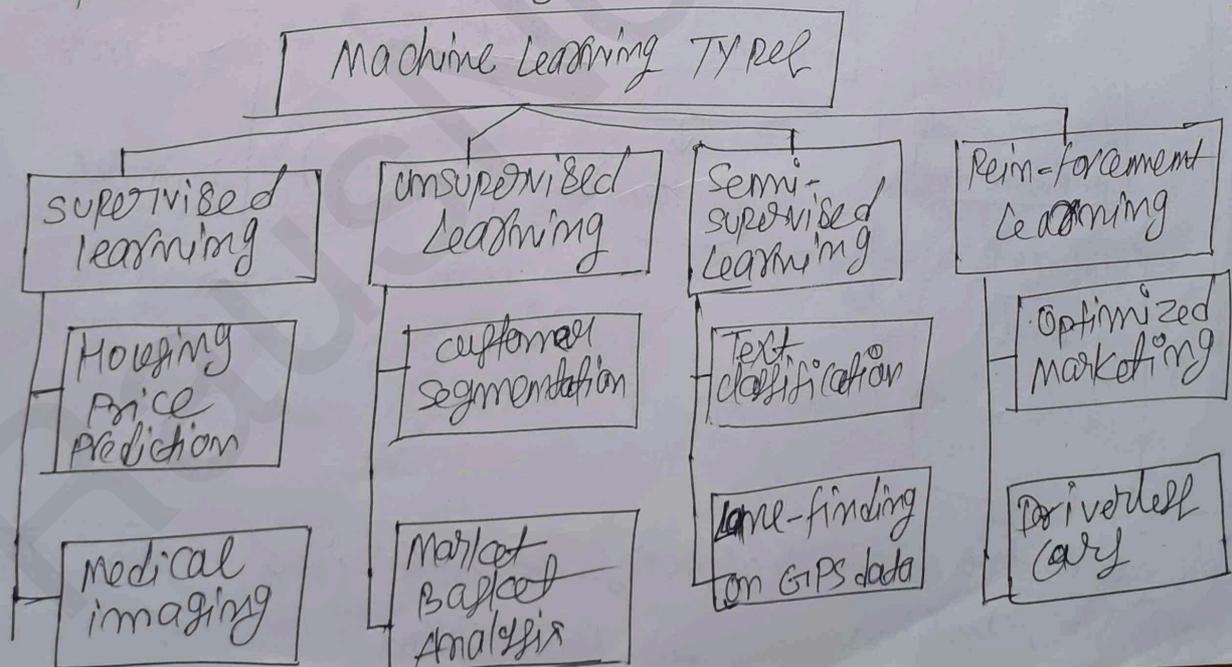
PSB 19/05/2021

Machine Learning: Machine learning is a subset of Artificial intelligence (AI) that enables systems to learn from data, identify patterns and make decisions with minimal human intervention.

Machine learning is used today for a wide range of commercial purposes, including suggesting products to consumers based on their past purchases, predicting stock market and translating text from one language to another. Machine learning contains a set of algorithms that work on a huge amount of data. Data is fed to these algorithms to train them, and on the basis of training, they build the model & perform a specific task.

Machine learning is divided into mainly four types

- (i) supervised Machine Learning
- (ii) unsupervised Machine Learning
- (iii) Semi-supervised Machine Learning
- (iv) Reinforcement Learning



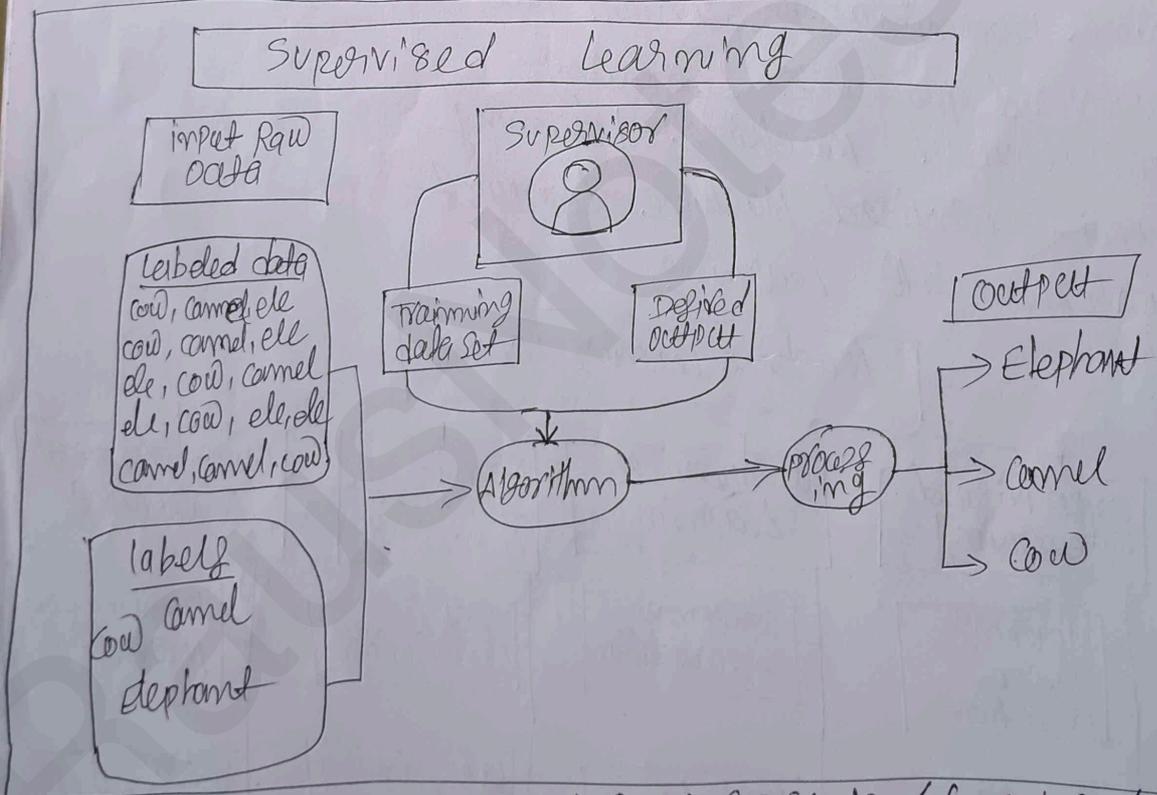
① Supervised Machine Learning: Supervised machine learning is based on supervision. It means in the supervised learning technique, we train the machines using the "labelled" dataset, and based on the training, the machine predicts the output.

labeled data: Data that includes both input features and their output labels

Ex ⇒ input: features of a house (size, location, no of rooms)
 output: price of the house

supervised learning maps the input variable (x) with the output variable (y).

Application: (i) spam detection (ii) Recommendation system (iii) fraud detection (iv) stock market prediction (v) weather forecasting (vi) sports analysis



There are two main categories of supervised learning that are (i) classification (ii) Regression

Classification: Classification algorithms are used to solve the classification problems in which the output variable is categorical, such as "yes" or "no", male or female, red or blue, etc. The classification algorithm predicts the categories present in the dataset. Ex → spam detection
Here are some classification algorithms.

- Logistic Regression
- Random forest
- Decision Tree

Regression: Regression predicting continuous target variables, which represent numerical values. Regression algorithms are used to solve regression problems in which there is a linear relationship b/w input and output variables. These are used to predict continuous output variables, such as market trends, weather prediction. Here are some regression algorithms.

- Linear Regression
- Polynomial Regression
- Random forest
- Decision Tree

Advantages of Supervised Machine Learning

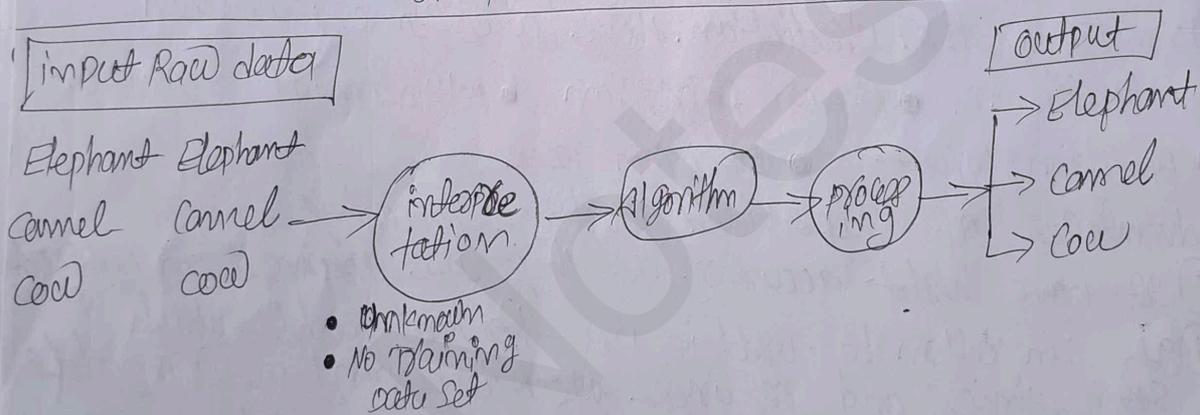
- (i) They have high accuracy as they are trained on labeled data.
- (ii) They can often be used in pre-trained models which save time and resources when developing new models from scratch.
- (iii) In supervised learning, we can have ~~an~~ exact idea about the class of objects.
- (iv) The process of decision making in supervised learning models is often interpretable.

Disadvantages: (i) These algorithms are not able to solve complex tasks.

- (ii) They may predict the wrong output if the test data is different from the training data.
- (iii) They can be time-consuming and costly as they rely on labeled data only.
- (iv) They require lots of computational time to train them.

② Unsupervised Machine Learning: Unsupervised learning is a type of machine learning technique in which an algorithm discovers patterns and relationships using unlabeled data. Unlike supervised learning, it is different from the supervised learning technique; as its name suggests, there is no need for supervision. In unsupervised machine learning, the machine is trained using the unlabeled dataset, and the machine predicts the output without any supervision. The primary goal of unsupervised learning is often to discover hidden patterns, similarities, or clusters within the data, which can then be used for various purposes.

Unsupervised Learning



Application of unsupervised learning

- (i) clustering
- (ii) Anomaly detection
- (iii) Recommendation system
- (iv) Image and video compression
- (v) Data preprocessing

* There are two main categories of unsupervised learning that are mentioned below.

- (i) clustering
 - (ii) Association
- ① Clustering: Clustering is the process of grouping data points into clusters based on their similarity. This technique is useful for identifying patterns and relationships.

in data without the need for labeled examples. The clustering technique is used when we want to find the inherent groups from the data. It is a way to group the objects into a cluster such that the objects with the most similarities remain in one group and have fewer or no similarities with the objects of other groups. Ex \Rightarrow Grouping the customers by their purchasing behaviour.

Here some clustering algorithm.

- k-means clustering algorithm
 - DBSCAN algorithm
 - Mean-shift algorithm
- Association: Association rule learning is an unsupervised learning technique, which finds interesting relations among variables within a large dataset. The main aim of this learning algorithm is to find the dependency of one data item on another data item and map those variables accordingly so that it can generate maximum profit. This algorithm is mainly applied in market basket analysis, web usage mining, continuous production.
- Some popular algorithms of Association rule learning are
- Apriori Algorithm
 - Eclat
 - FP-growth algorithm

* Advantage of unsupervised Machine learning

- (i) It helps to discover hidden patterns and various relationships b/w the data.
- (ii) used for tasks such as customer segmentation, anomaly detection, and data exploration.
- (iii) It does not require labeled data and reduces the effort of data labeling.

Disadvantages of unsupervised machine learning

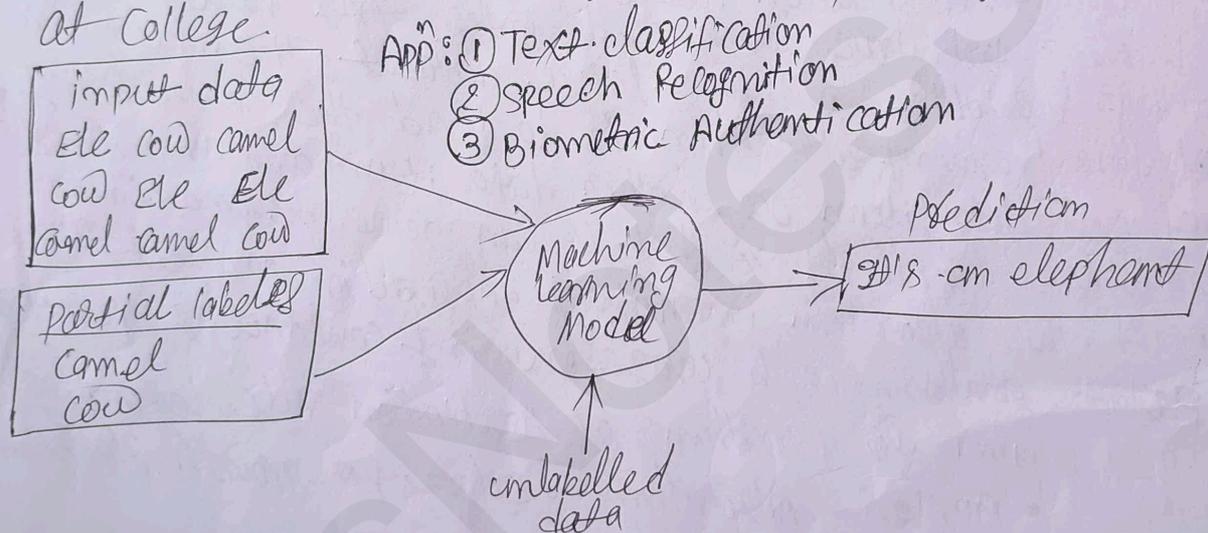
- Without using labels, it may be difficult to predict the

Quality of the model's output.

- suffer interpretability may not be clear and may not have meaningful interpretations.

③ Semi-supervised learning: Semi-supervised learning is a type of machine learning algorithm that lies b/w supervised and unsupervised machine learning. It represents the intermediate ground b/w supervised (with labelled data) and unsupervised ~~learning~~ learning (without labelled data) algorithms, and uses the combination of labelled and unlabeled datasets during the training period.

Example: Student has to revise himself after analyzing the same concept under the guidance of an instructor at college.



Advantage: ① It is simple and easy to understand the algorithm.

② It is highly efficient.

③ It is used to solve drawbacks of supervised and unsupervised algorithms.

Disadvantages:

① Iterations results may not be stable.

② we cannot apply these algorithms to network-level data.

③ Accuracy is low.

④ Reinforcement learning: Reinforcement learning works on a feedback-based process, in which an AI agent automatically explore its surrounding by hitting & trial, acting, learning from experience, and improving its performance. Agent gets rewarded for each good action and get punished for each bad action. In this technique, the model keeps on increasing its performance using reward feedback to learn the behavior or patterns.

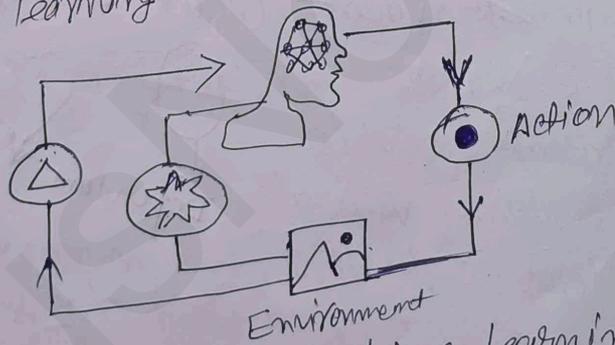
Ex ⇒ Google self driving car,

Applications of reinforcement machine learning

- ① Game playing
- ② Robotics
- ③ Automation vehicles
- ④ Autonomous vehicle
- ⑤ Recommendation system
- ⑥ Healthcare
- ⑦ Natural language processing
- ⑧ Agriculture

Here are some of most common reinforcement learning algorithms -

- ① Q-learning
- ② State-Action-Reward-State-Action
- ③ Deep Q-learning



Advantages of Reinforcement Machine learning:

- ① It has autonomous decision-making that is well-suited for tasks and that can learn to make a sequence of decisions, like robotics and game-playing
- ② It is used to solve complex problems that cannot

be solved by conventional technique.

Disadvantages of Reinforcement Machine Learning

- (i) Training Reinforcement Learning agents can be computationally expensive and time consuming.
- (ii) It is not preferable to solving simple problems.
- (iv) It needs a lot of data, a lot of computation that makes it costly.

Unit - 2

* Simple Linear Regression: Simple Linear Regression is a statistical technique used to model the relationship between a dependent variable y and a single independent variable x . It assumes a linear relationship between these two variables, represented by the equation of a straight line. It is also a type of supervised machine learning algorithm that learns from the labelled datasets. It predicts the continuous output variables based on the independent input variable.

Equation of linear regression $y = ax + b$

x → independent variable

y → dependent variable

a → is the slope of the line (how much y changes for a unit change in x).

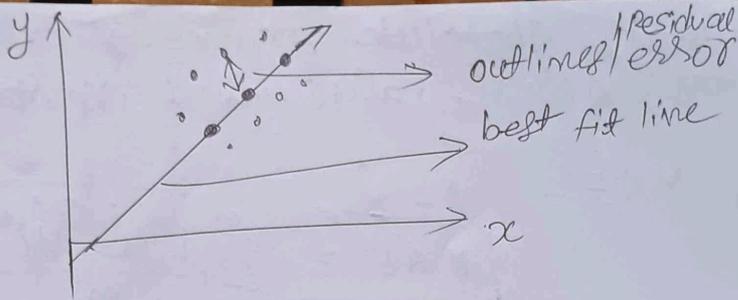
b → is the intercept

$$y = ax + b$$

where

$$a = \frac{n(\sum xy) - \sum x \cdot \sum y}{n(\sum x^2) - (\sum x)^2}$$

$$b = \frac{\sum y - a(\sum x)}{n}$$



$$\boxed{\text{Residual error} = \text{Actual} - \text{Predicted}}$$

Q Expenditure of business (in thousands) for every year is shown in table below.

x (year)	1	2	3	4	5
y (Expenditure)	12	19	29	37	45

Solⁿ:

Sr. N	x	y	xy	x ²
1	1	12	12	1
2	2	19	38	4
3	3	29	87	9
4	4	37	148	16
5	5	45	225	25

$$\Sigma x = 15, \Sigma y = 142 \Rightarrow \Sigma xy = 510 \quad \Sigma x^2 = 55$$

$$b = \frac{\Sigma y - a \Sigma x}{n}$$

$$y = ax + b$$

$$a = \frac{n \Sigma xy - \Sigma x \Sigma y}{n \Sigma x^2 - (\Sigma x)^2} = \frac{5 \times 510 - 15 \times 142}{5 \times 55 - (15)^2} = \frac{5 \times 510 - 15 \times 142}{5 \times 55 - 225} = \frac{2550 - 2130}{275 - 225} = \frac{420}{50} = 8.4$$

$$b = \frac{142 - 8.4 \times 15}{5} \quad \therefore b = 3.2$$

$$\text{Equation of line} \Rightarrow \boxed{y = 8.4x + 3.2}$$

$$\text{for } x = 6 \Rightarrow y = 8.4 \times 6 + 3.2$$

$$y = 53.6 \text{ thousands}$$

Q10] Perform comparison b/w simple linear reg and polynomial regression with example.

$$\sum (y - \hat{y})$$

* Evaluate Metrics for linear regression

(i) Mean Absolute Error (MAE): Mean absolute error is the average difference between actual value and predicted value.

$$MAE = \frac{\sum_{i=1}^n |y_i - \hat{y}_i|}{n}$$

$$\text{Residual error} = y - \hat{y}$$

X (item)	y (Sales)
11	80
12	90
13	100
14	110
15	120

x	y _i	\hat{y}_i
16	80	75
17	75	85

$$MAE = \frac{|80 - 75| + |75 - 85|}{2} = \frac{15}{2} = 7.5$$

(ii) Mean square Error (MSE): Mean square error is the square of average difference b/w actual value and predicted value.

$$MSE = \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{n}$$

$$MSE = \frac{(80 - 75)^2 + (75 - 85)^2}{2} \Rightarrow \frac{25 + 100}{2} = \frac{125}{2} = 62.5$$

(RMSE)

(iii) Root Mean Square Error = $\sqrt{MSE} \Rightarrow \sqrt{62.5} \Rightarrow 7.90$

Root Mean Square Error is the root of ~~square~~ Mean Square Error.

④ Relative MSE: (RelMSE)

$$\text{RelMSE} = \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{\sum_{i=1}^n (y_i - \bar{y})^2}$$

$\bar{y} \rightarrow$ mean

x_i	sales (y)
11	80
12	90
13	100
14	110
15	120

The avg of y is: $\frac{80+90+100+110+120}{5}$

$$\Rightarrow \frac{500}{5} = 100 \quad \boxed{\bar{y} = 100}$$

$$\text{RelMSE} = \frac{(80-75)^2 + (75-85)^2}{(80-100)^2 + (75-100)^2}$$

$$\text{RelMSE} = \frac{725}{1025} = 0.707$$

Item	A.V	P.V
16	80	75
17	75	85

* Coefficient of variation ~~$CV = \frac{\text{RelMSE}}{\bar{y}}$~~

$$\boxed{CV = \frac{RMSE}{\bar{y}}}$$

$$\Rightarrow CV = \frac{\sqrt{62.5}}{100}$$

$$\boxed{CV = 0.08}$$

⑤ Coefficient of Determination (R^2)

X	Y	X^2	Y^2
1	3	1	9
2	4	4	16
3	8	9	64
4	4	16	16
5	5	25	25
$\Sigma X = 15$	$\Sigma Y = 18$	$\Sigma X^2 = 55$	$\Sigma Y^2 = 120$

$$a = \frac{n \Sigma xy - \Sigma x \Sigma y}{n \Sigma x^2 - (\Sigma x)^2} \Rightarrow \frac{5 \times 58 - 15 \times 18}{5 \times 55 - 15 \times 15}$$

$$a = \frac{290 - 270}{275 - 225}$$

$$a = \frac{20}{50} \Rightarrow a = 0.4$$

$$\boxed{b = \frac{\Sigma y - a \Sigma x}{n}}$$

$$b = \frac{18 - (0.4)(15)}{5}$$

$$b = \frac{18 - 6}{5} \Rightarrow b = \frac{12}{5}$$

$$\boxed{b = 2.4}$$

$$\text{COD} \Rightarrow \boxed{y = 0.4x + 2.4}$$

Let's predict value of y for $x = 1, 2, 3, 4, 5$ & calculate coefficient of determination.

X	y	\hat{y}	$(y_p - \bar{y})^2$	$(y_p - \bar{y})^2$	$(\hat{y} - \bar{y})$	$(y - \bar{y})^2$
1	3	2.8	0.04			
2	4	3.2	0.04			
3	2	3.6				
4	4	4.0				
5	5	4.4				

for $x_1 = 0.4(x) + 2.4 \Rightarrow x_1 = 0.4(1) + 2.4 \Rightarrow x_1 = 0.4 + 2.4 = 2.8$

for $x_2 = 0.4(2) + 2.4 \Rightarrow 0.8 + 2.4 \Rightarrow 3.2$

and so on till x_5 .

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

$$\bar{y} = \frac{3 + 4 + 2 + 4 + 5}{5}$$

$$\bar{y} = \frac{18}{5}$$

$$\bar{y} = 3.6$$

* Multiple Linear regression: Linear Reg is a statistical approach for modeling relationship b/w dependent and one independent variable but multiple linear regression accepts two model relationship b/w two or more features to fit a linear equation to predict one dependent variable. Regression model learns a function from the dataset and use it to predict y values for unknown x_1, x_2, x_3 by adding more predictors to simple linear regression the technique ~~helps~~ helps to understand better how the predictors effects the outcome variable as a whole. for example prediction of CO2 emission based on no. of size and no. of cylinders in a car. Multiple linear reg. tries to fit a reg line to multi dimensional data point.

$$y = a_1x_1 + a_2x_2 + a_3x_3 + \dots + a_mx_m + b$$

a_1, a_2, \dots → slope
 b → intercept
 x → y

Agent Salary

X	Experiment	Salary y
35	1	20
28	1	10
40	2	35
50	4	?

$y =$

$$y = a_1x_1 + a_2x_2 + b$$

$$a_1 = \frac{(\sum x_1^2)(\sum x_2 y) - (\sum x_1 x_2)(\sum x_1 y)}{(\sum x_1^2)(\sum x_2^2) - (\sum x_1 x_2)^2}$$

$$a_2 = \frac{(\sum x_1^2)(\sum x_2 y) - \sum(x_1 x_2) \cdot \sum(x_1 y)}{(\sum x_1^2)(\sum x_2^2) - (\sum x_1 x_2)^2}$$

$$b = \bar{y} - a_1 \bar{x}_1 - a_2 \bar{x}_2$$

x_1	x_2	y	x_1^2	x_2^2	$x_1 x_2$	$x_1 y$	$x_2 y$
35	1	20	1225	1	35	700	20
25	1	10	625	1	25	250	10
40	2	35	1600	4	80	1400	70
50	4	?					

$$\begin{aligned} \sum x_1 &= 100 \\ \sum x_2 &= 4 \\ \sum y &= 65 \end{aligned}$$

$$\begin{aligned} \sum x_1^2 &= \\ \sum x_2^2 &= 6 \\ \sum x_1 x_2 &= 140 \end{aligned}$$

$$\begin{aligned} \sum x_1 y &= \\ \sum x_2 y &= 100 \end{aligned}$$

$$a_1 =$$

$$b = \bar{y} - a_1 \bar{x}_1 - a_2 \bar{x}_2$$

$$y = a_1 x_1 + a_2 x_2 + b$$

$$y = 0.09 x_1 + 14.5 x_2 + 4.5$$

$$y = 0.09 x_1 (50) + 14.5 (4) + 4.5$$

$$y = 67 \text{ K}$$

* Poly
a reg
the indepen
of a

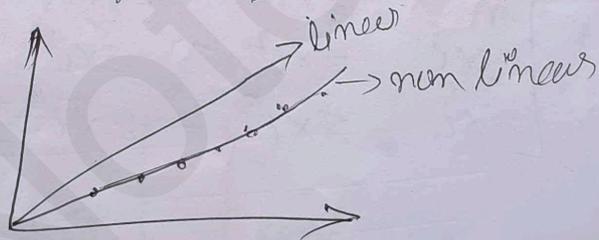
* polynomial regression: polynomial is a regression algorithm that models the relationship between an independent variable and a dependent variable as an nth degree. It is known as a special case of MLR.

$$y = a_1 x^1 + a_2 x^2 + a_3 x^3 + \dots + a_n x^n + \text{FP}$$

Polynomial function

We add some polynomial terms to the multiple linear regression equation to convert it into polynomial regression. We required it to increase accuracy. (i) Data set used of non-linear nature. (ii) Why we need polynomial regression.

If we apply a linear model on the linear dataset it provides a good result, but if we apply the same model without any modification on a non-linear dataset, it will produce a drastic output.



In ~~the~~ this figure data points are arranged in non linear fashion. We need ~~the~~ a ~~linear~~ polynomial regression.

In above we take a data set which is arranged in non linear. If we convert it with linear model it hardly covers any data point. A curve is suitable to cover most of the datapoints i.e. polynomial model. If dataset is polynomial

Regression instead of simple LR. Simple & MLR are also polynomial regression with a single degree (or). But PR is linear equation with n^{th} degree (x^n). If we add a degree to our LR then it will be converted to PR.

Use: When there is a non-linear co-relation b/w 2 variables - it can easily fit a vast range of curvatures.

$$Y = b + \sum_{i=1}^n a_i x^i + \text{FP}$$

$$a = X^{-1} B$$

$$X = \begin{bmatrix} n & \sum x_i & \sum x_i^2 \\ \sum x_i & \sum x_i^2 & \sum x_i^3 \\ \sum x_i^2 & \sum x_i^3 & \sum x_i^4 \end{bmatrix} \quad \begin{bmatrix} \sum y_i \\ \sum (x_i y_i) \\ \sum (x_i^2 y_i) \end{bmatrix} \begin{matrix} \rightarrow b \\ \rightarrow a_1 \\ \rightarrow a_2 \end{matrix}$$

Q.1

x	y
1	1
2	4
3	9
4	15

$$y = a_1 x_1 + a_2 x_1^2 + b$$

$$a = X^{-1} B$$

$$a \approx \begin{bmatrix} b \\ a_1 \\ a_2 \end{bmatrix}$$

$$a = \begin{bmatrix} n & \sum x_i & \sum x_i^2 \\ \sum x_i & \sum x_i^2 & \sum x_i^3 \\ \sum x_i^2 & \sum x_i^3 & \sum x_i^4 \end{bmatrix}^{-1} \begin{bmatrix} \sum y_i \\ \sum (x_i y_i) \\ \sum (x_i^2 y_i) \end{bmatrix}$$

x	y	$x_i y_i$	x_i^2	x_i^3	x_i^4	$x_i^2 y_i$
1	1	1	1	1	1	1
2	4	8	4	8	16	16
3	9	27	9	27	81	81
4	15	60	16	64	256	240

$$\bar{y} = \begin{bmatrix} 0.75 \\ 0.95 \\ 2.75 \end{bmatrix}$$

$$\begin{bmatrix} 6 \\ a_1 \\ a_2 \end{bmatrix} = \begin{bmatrix} 4 & 30 & 30 \\ 10 & 100 & 100 \\ 30 & 354 & 338 \end{bmatrix}^{-1} \begin{bmatrix} 2 \\ 96 \\ 338 \end{bmatrix}$$

$$X^{-1} = \frac{1}{|adj X|}$$

$$|adj X| = \begin{vmatrix} 30 & 100 & 10 & 30 \\ 100 & 354 & 20 & 100 \\ 10 & 30 & 4 & 10 \\ 30 & 100 & 10 & 30 \end{vmatrix} \quad adj = \begin{bmatrix} 620 & -540 & 100 \\ -540 & 516 & -100 \\ 100 & -100 & 20 \end{bmatrix}$$

$$\begin{bmatrix} b \\ a_1 \\ a_2 \end{bmatrix} = \begin{bmatrix} 4 & 10 & 30 \\ 10 & 30 & 100 \\ 30 & 100 & 354 \end{bmatrix}$$

$$\text{adj } A = \begin{bmatrix} 620 & -540 & 100 \\ -540 & 516 & -100 \\ 100 & -100 & 20 \end{bmatrix}$$

$$|A| = 248 \cdot 0 + (-5400) + 3000 \Rightarrow |A| = \cancel{13920}$$

$$|A| = 80$$

$$A^{-1} = \frac{1}{80} \begin{bmatrix} 620 & -540 & 100 \\ -540 & 516 & -100 \\ 100 & -100 & 20 \end{bmatrix}$$

$$A^{-1} = \begin{bmatrix} \frac{620}{80} & \frac{-540}{80} & \frac{100}{80} \\ \frac{-540}{80} & \frac{516}{80} & \frac{-100}{80} \\ \frac{100}{80} & \frac{-100}{80} & \frac{20}{80} \end{bmatrix}$$

$$A^{-1} = \begin{bmatrix} \frac{31}{4} & -\frac{27}{4} & \frac{5}{4} \\ -\frac{27}{4} & \frac{129}{20} & -\frac{5}{4} \\ \frac{5}{4} & -\frac{5}{4} & \frac{1}{4} \end{bmatrix}$$

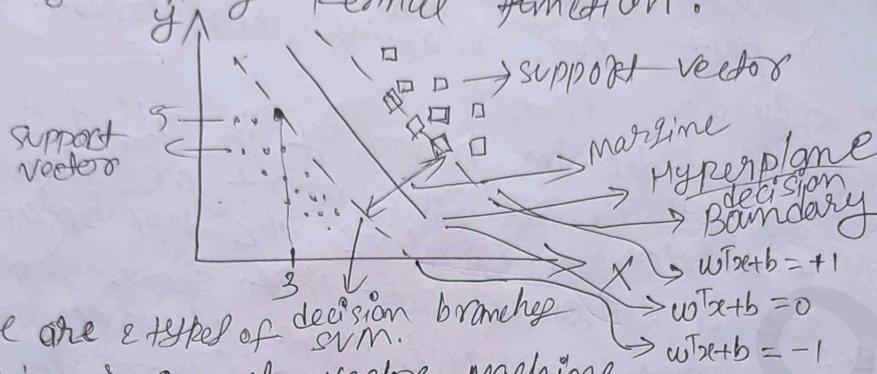
$$A^{-1} = \begin{bmatrix} 7.75 & -6.75 & 1.25 \\ -6.75 & -6.45 & -1.25 \\ 1.25 & -1.25 & 0.25 \end{bmatrix} \times \begin{bmatrix} 2 \\ 96 \\ 338 \end{bmatrix} \Rightarrow \begin{bmatrix} 183.2 \\ -183.2 \\ -210 \end{bmatrix}$$

$$\begin{bmatrix} b \\ a_1 \\ a_2 \end{bmatrix} = \begin{bmatrix} 210 \\ 183.2 \\ -33 \end{bmatrix}$$

$$y = 183.2x_1 + (-33)x_2 + (-210)$$

$$y = 183.2x_1 - 33x_2 - 210$$

* Support vector machine (SVM): SVM is a supervised machine algorithm used for both classification and regression. It can solve both linear and nonlinear problems using kernel function.



There are 2 types of decision SVM.

- (i) Linear support vector machine
- (ii) Non-Linear support vector machine

(i) Linear SVM : Linear SVM is used for linearly separable data which means if a dataset can be classified into two classes by using a single straight line then such data is linearly separable data.

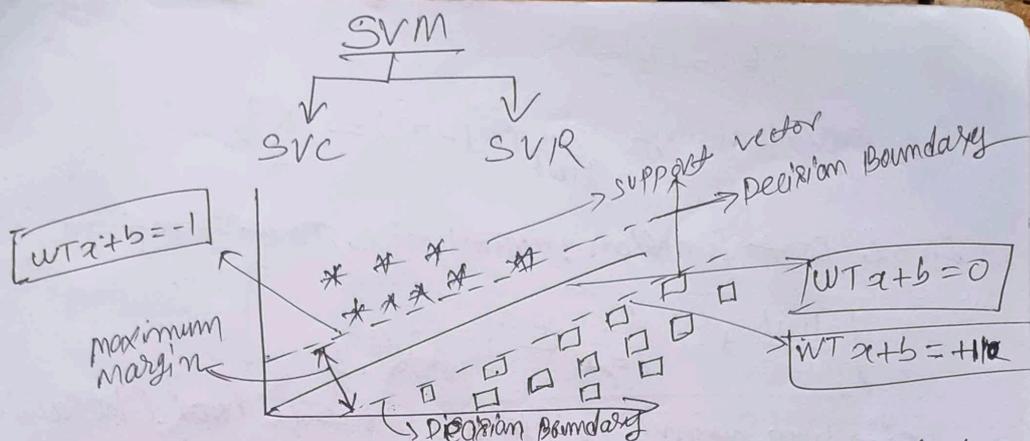
(ii) Non Linear SVM : Non-SVM used for non linear separable data, which means if a data set cannot be classified by using a straight line such data is termed as non-linear data, SVM works by finding a hyper plane in a high dimensional space that best separates data into diff classes. It aims to find a function that approximate the relationship b/w input variables and continuous

target variable by minimizing prediction error.

Goal is to separate best suitable line or decision boundary that can separate n dimensional space ~~into~~ classes so that we can easily put new data points ^{into} into correct category in future.

This best decision boundary is called hyperplane. SVM chooses the extreme point or support vector that ~~used~~ helps in creating the hyperplane we can see a strange cat that some features of dog. ^{Ex:} we want a model that can accurately identify whether it is a cat or dog. To train model with a lot of images of cat or dog so that it can learn about the different features of dog, SVM creates a decision boundary b/w these two data and choose extreme case of dog and cat. ^{Ex} \Rightarrow face detection

Hyperplane: There can be multiple lines or the decision boundary to separate the classes in n dimensional space but we need to find out the best decision boundary that help to classify the data points, this best boundary is known as a hyperplane. The dimension of the hyperplane depend upon the features present in the data set which means if there are two features then hyperplane will be a straight line. If there are three features then hyperplane will be a 2D plane.



Decision Boundaries are also known as marginal planes
 vectors helpful to make decision boundaries are known
 as support vectors.

Hard Margin / Soft Margin

$$y = wx + c$$

$$y = w^T x + c$$

$ax + by + c$ Equation of line

$$by = -ax - c \Rightarrow y = \frac{-ax}{b} - \frac{c}{b}$$

Let it a line $3x + 2y + 4 = 0$

$$3(4) + 2(4) + 4 \Rightarrow 12 + 8 + 4 = 24 \quad (4, 4)$$

$w^T x_1 + b = +1$
 $w^T x_2 + b = -1$

} how SVM used for classification of linearly separable data,

$w^T(x_1 - x_2) = +2$ dividing w by magnitude, we get a vector

$$\frac{w^T(x_1 - x_2)}{|w|} = \frac{2}{|w|} \rightarrow \text{we need to maximize it}$$

Marginal plane distance

{ maximum margin }
 give

$$y = \begin{cases} +1 & w^T x + b \geq +1 \\ -1 & w^T x + b \leq -1 \end{cases}$$

must see

$$y = w^T x + b \geq 1$$

for minimising $\frac{w^T (x_1 - x_2)}{|w|} = \frac{|w|}{2}$

final cost function: Minimise cost function

$$\left[\frac{|w|}{2} + c_i \sum_{i=1}^n \epsilon \right] \text{ (eta)}$$

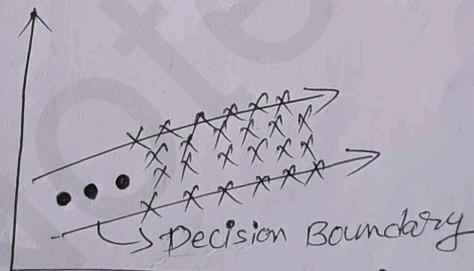
$c_i \rightarrow$ how many points we want to avoid because some errors are acceptable.

$\epsilon \rightarrow$ eta is the submission of distance of misclassification of points from marginal plane.

kernel:
 $\begin{cases} \rightarrow \text{sigmoid kernel} \\ \rightarrow \text{linear kernel} \\ \rightarrow \text{RBF kernel} \end{cases}$

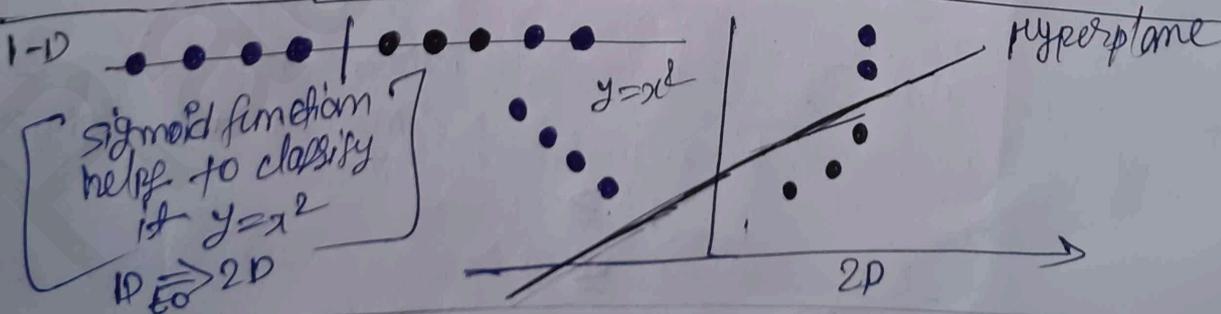
$y = (x^2)$ kernel Trick is used in SVM

* SVM Regression:



Regression to find a hyperplane such that most of the data points are inside marginal line:

same for $\left[\frac{|w|}{2} + c_i \sum_{i=1}^n \epsilon \right] \rightarrow \text{(eta)}$

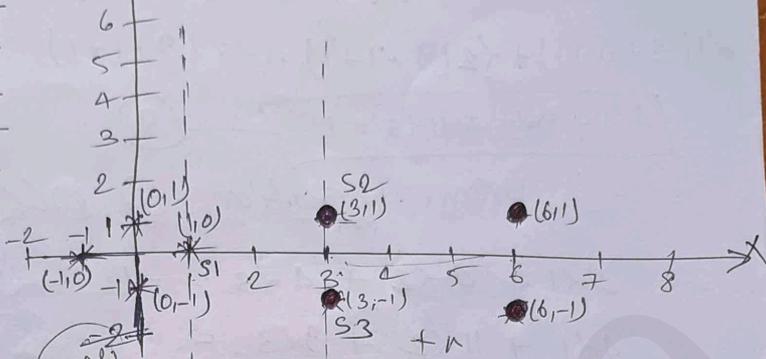


Q Generate a hyperplane, with the following points

(3,1)	+1	✓
(3,-1)	+1	✓
(6,1)	+1	-
(6,-1)	+1	-
(1,0)	-1	-
(0,1)	-1	-
(0,-1)	-1	-
(-1,0)	-1	-

Solⁿ: y (How sum used for classifying linearly separable data)

Solⁿs



$$s_1 = \begin{pmatrix} 1 \\ 0 \end{pmatrix}, s_2 = \begin{pmatrix} 3 \\ 1 \end{pmatrix}, s_3 = \begin{pmatrix} 3 \\ -1 \end{pmatrix}$$

Adding bias vectors to convert these vectors to augmented vectors.

$$\bar{s}_1 = \begin{pmatrix} 1 \\ 0 \\ 1 \end{pmatrix}, \bar{s}_2 = \begin{pmatrix} 3 \\ 1 \\ 1 \end{pmatrix}, \bar{s}_3 = \begin{pmatrix} 3 \\ -1 \\ 1 \end{pmatrix}$$

$\alpha_1, \alpha_2, \alpha_3$ are slopes

$$\alpha_1 \bar{s}_1 \bar{s}_1 + \alpha_2 \bar{s}_1 \bar{s}_2 + \alpha_3 \bar{s}_1 \bar{s}_3 = -1$$

$$\alpha_1 \bar{s}_1 \bar{s}_2 + \alpha_2 \bar{s}_2 \bar{s}_2 + \alpha_3 \bar{s}_2 \bar{s}_3 = +1$$

$$\alpha_1 \bar{s}_1 \bar{s}_3 + \alpha_2 \bar{s}_2 \bar{s}_3 + \alpha_3 \bar{s}_3 \bar{s}_3 = +1$$

$$\alpha_1 \begin{pmatrix} 1 \\ 0 \\ 1 \end{pmatrix} \begin{pmatrix} 1 \\ 0 \\ 1 \end{pmatrix} + \alpha_2 \begin{pmatrix} 1 \\ 0 \\ 1 \end{pmatrix} \begin{pmatrix} 3 \\ 1 \\ 1 \end{pmatrix} + \alpha_3 \begin{pmatrix} 1 \\ 0 \\ 1 \end{pmatrix} \begin{pmatrix} 3 \\ -1 \\ 1 \end{pmatrix} = -1$$

$$\alpha_1 (1+0+1) + \alpha_2 (3+0+1) + \alpha_3 (3+0+1) = -1$$

$$\boxed{2\alpha_1 + 4\alpha_2 + 4\alpha_3 = -1} \quad \text{--- (1)}$$

$$\alpha_1 \begin{pmatrix} 1 \\ 0 \\ 1 \end{pmatrix} \begin{pmatrix} 3 \\ 1 \\ 1 \end{pmatrix} + \alpha_2 \begin{pmatrix} 3 \\ 1 \\ 1 \end{pmatrix} \begin{pmatrix} 3 \\ 1 \\ 1 \end{pmatrix} + \alpha_3 \begin{pmatrix} 3 \\ 1 \\ 1 \end{pmatrix} \begin{pmatrix} 3 \\ -1 \\ 1 \end{pmatrix} = 1$$

$$\alpha_1 (3+0+1) + \alpha_2 (9+1+1) + \alpha_3 (9-1+1) = 1$$

$$\boxed{4\alpha_1 + 11\alpha_2 + 9\alpha_3 = 1} \quad \text{--- (2)}$$

$$\alpha_1 \begin{pmatrix} 3 \\ -1 \\ 1 \end{pmatrix} + \alpha_2 \begin{pmatrix} 3 \\ -1 \\ 1 \end{pmatrix} + \alpha_3 \begin{pmatrix} 3 \\ -1 \\ 1 \end{pmatrix}$$

$$\alpha_1(3+0+1) + \alpha_2(9-1+1) + \alpha_3(9+1+1)$$

$$\boxed{4\alpha_1 + 9\alpha_2 + 11\alpha_3 = 1} \quad \text{--- (3)}$$

From equation (1), (2) & (3)

$$2\alpha_1 + 4\alpha_2 + 4\alpha_3 = -1 \quad \text{--- (1)}$$

$$4\alpha_1 + 11\alpha_2 + 9\alpha_3 = 1 \quad \text{--- (2)}$$

$$4\alpha_1 + 9\alpha_2 + 11\alpha_3 = 1 \quad \text{--- (3)}$$

From eqn (2) & (3)

$$\cancel{4\alpha_1} + 11\alpha_2 + 9\alpha_3 = 1$$

$$\cancel{4\alpha_1} + 9\alpha_2 + 11\alpha_3 = 1$$

$$2\alpha_2 + (-2\alpha_3) = 0$$

$$2\alpha_2 - 2\alpha_3 = 0$$

$$\boxed{\alpha_2 = \alpha_3}$$

multiply (2) in eq (1)

$$\cancel{4\alpha_1} + 8\alpha_2 + 8\alpha_3 = -2 \quad \text{taking eq}^n \text{ (3)}$$

$$\cancel{4\alpha_1} + 9\alpha_2 + 11\alpha_3 = 1$$

$$-\alpha_2 - 3\alpha_3 = -3$$

$$\alpha_2 + 3\alpha_3 = 3 \quad \text{put } \alpha_3 = \alpha_2$$

$$\alpha_2 + 3\alpha_2 = 3$$

$$4\alpha_2 = 3$$

$$\boxed{\alpha_2 = \frac{3}{4}}$$

$$\therefore \alpha_2 = \alpha_3$$

$$\boxed{\alpha_3 = \frac{3}{4}}$$

put α_2 and α_3 in eq (1)

$$2x_1 + 4x_2 + 4x_3 = -1$$

$$2x_1 + 4 \times \frac{3}{4} + 4 \times \frac{3}{4} = -1$$

$$2x_1 + 3 + 3 = -1$$

$$2x_1 + 6 = -1 \rightarrow 2x_1 = -1 - 6$$

$$2x_1 = -7$$

$$x_1 = -\frac{7}{2}$$

$$K1 = -3.5$$

$$\vec{w} = \sum_{i=1}^n \alpha_i \vec{s}_i$$

$$\vec{w} = -3.5 \begin{pmatrix} 1 \\ 0 \\ 1 \end{pmatrix} + 0.75 \begin{pmatrix} 3 \\ 1 \\ 1 \end{pmatrix} + 0.75 \begin{pmatrix} 3 \\ -1 \\ 1 \end{pmatrix}$$

$$\vec{w} = \begin{pmatrix} -3.5 \\ 0 \\ -3.5 \end{pmatrix} + \begin{pmatrix} 2.25 \\ 0.75 \\ 0.75 \end{pmatrix} + \begin{pmatrix} 2.25 \\ -0.75 \\ 0.75 \end{pmatrix}$$

$$\vec{w} = \begin{pmatrix} 1 \\ 0 \\ -2 \end{pmatrix} \left[\omega \text{ slope?} \right]$$

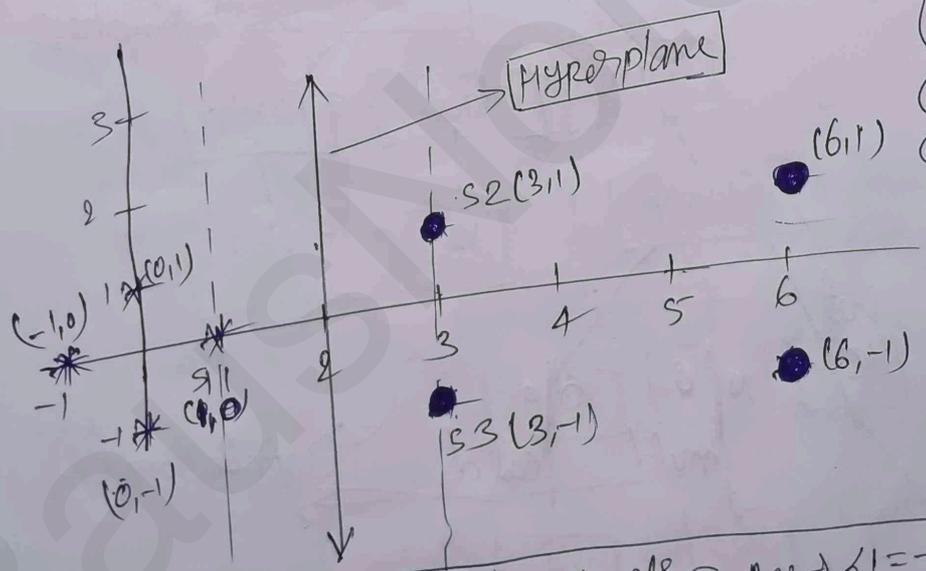
→ bias & intercept

$$w = \begin{pmatrix} x \\ 0 \\ y \end{pmatrix}$$

$$b = -2$$

$$b + 2 = 0$$

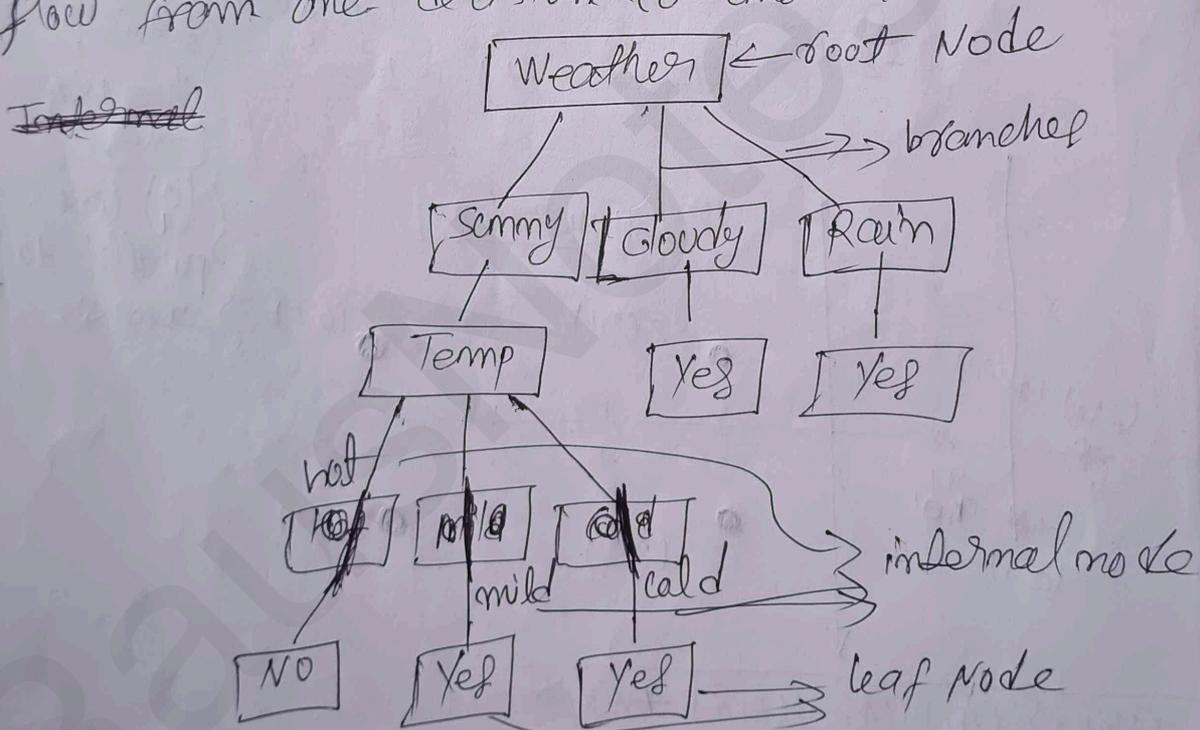
hyperplane will be parallel to x & y (best if $\begin{pmatrix} 1 \\ 1 \end{pmatrix}$ them as to x & y)
 $\begin{pmatrix} 1 \\ 0 \end{pmatrix}$ line parallel to x & y 's



* How \Rightarrow $(4,1)$ $(4,-1)$ $(6,0)$ +ve pts } Any $\rightarrow x_1 = -3$
 $(1,0)$ $(0,1)$ $(0,-1)$ -ve pts } $x_2 = 0$
 $x_3 = 0$

Decision Tree (Classification)

Decision Tree is a supervised algorithm used for both classification and regression. It models decision of tree like structure where internal node represent attribute test, branches represent root node attribute value and leaf node represent final decision. It is a graphical representation of different options for solving a problem. It has a hierarchical tree structure having root node, it is starting points that represent entire dataset. Branches are the lines that connect nodes. It shows the flow from one decision to another.



* Root Node: Root Node is the starting point that represents the entire dataset.

* Branches: These are the lines that connect nodes. It shows the flow from one decision to another.

* Internal Nodes: Internal nodes are points where decisions are made based on the input features.

* Leaf Node: These are the terminal nodes at the end of branches that represent final outcomes or predictions.

* Gini Index:

* Entropy (H): Entropy measures the impurity or uncertainty in the dataset. The formula for entropy is:

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

$$H(S) = - [P_+ \log_2 P_+] - [P_- \log_2 P_-]$$

where $H(S) \Rightarrow$ entropy of set S.
 $c \Rightarrow$ c is the no of classes in dataset.
 $p_i \Rightarrow p_i$ is the proportion of samples belonging to class i.

$\log_2 p_i \Rightarrow$ is the logarithm (base 2) of p_i .

$P_- \Rightarrow$ Probability of No

$P_+ \Rightarrow$ Probability of Yes

* Information Gain (IG): Information Gain measures the reduction in entropy after a dataset is split on an attribute. It is given by:

$$IG(S, A) = H(S) - \sum_{v \in \text{values}(A)} \frac{|S_v|}{|S|} H(S_v)$$

$H(S_v)$: is the entropy of subset S_v .

- $IG(S, A)$ is the information gain for splitting on attribute A
- $H(S)$ is the entropy of the original dataset.
- S_v is the subset of S where attribute A takes value v.
- $\frac{|S_v|}{|S|}$ is the proportion of S that falls into subset

Q: Decision Tree:

Day	Weather	Temperature	Humidity	wind	play
Day 1	Sunny	Hot	High	weak	NO
Day 2	Sunny	Hot	High	strong	NO
Day 3	cloudy	Hot	High	weak	Yes
Day 4	Rain	Mild	High	weak	Yes
Day 5	Rain	Cool	Normal	weak	Yes
Day 6	Rain	cool	Normal	strong	NO
Day 7	cloudy	cool	Normal	strong	Yes
Day 8	Sunny	Mild	High	weak	NO
Day 9	Sunny	Cool	Normal	weak	Yes
Day 10	Rain	Mild	Normal	weak	Yes
Day 11	Sunny	Mild	Normal	strong	Yes
Day 12	cloudy	Mild	High	strong	Yes
Day 13	cloudy	Hot	Normal	weak	Yes
Day 14	Rain	Mild	High	strong	NO

Now calculate IG (Information Gain) of Weather

Step 1: Entropy of Entire dataset

$$H(W) = - \sum_{i=1}^n p_i \log_2 p_i \quad \left\{ \begin{array}{l} \text{whole} \\ \text{or} \end{array} \right. \quad H(W) = -P_+ \log_2 P_+ - P_- \log_2 P_-$$

$$H(W) = H(+, -) = -\frac{9}{14} \log_2 \frac{9}{14} - \frac{5}{14} \log_2 \frac{5}{14} = 0.94$$

Step 2: Entropy of all attributes (weather)

• Entropy of Sunny $\{+2, -3\} = -\frac{2}{5} \log_2 \frac{2}{5} - \frac{3}{5} \log_2 \frac{3}{5} = 0.97$

• Entropy of cloudy $\{+4, -0\} = -\frac{4}{4} \log_2 \frac{4}{4} - \frac{0}{4} \log_2 \frac{0}{4} = 0$

• Entropy of Rain $\{+3, -2\} = -\frac{3}{5} \log_2 \frac{3}{5} - \frac{2}{5} \log_2 \frac{2}{5} = 0.97$

$$\text{Information Gain} = \text{Entropy(whole data)} - \frac{5}{14} \text{Ent}(S) - \frac{4}{14} \text{Ent}(C) - \frac{5}{14} \text{Ent}(R)$$

$$= 0.246$$

* Calculate IGI of Temperature

Step 1: Entropy of entire dataset

$$H(+9, -5) = 0.94$$

Step 2: Entropy of Hot $\{+2, -2\} = -\frac{2}{4} \log_2 \frac{2}{4} - \frac{2}{4} \log_2 \frac{2}{4} = 1.0$

~~Entropy~~ Entropy of mild $\{+4, -2\} = -\frac{4}{6} \log_2 \frac{4}{6} - \frac{2}{6} \log_2 \frac{2}{6} = 0.91$

" " cool $\{+3, -1\} = -\frac{3}{7} \log_2 \frac{3}{7} - \frac{1}{7} \log_2 \frac{1}{7} = 0.81$

$$\text{Information Gain} = \text{Entropy}(\text{whole dataset}) - \frac{4}{14} \text{Ent}(H) - \frac{6}{14} \text{Ent}(M) - \frac{4}{14} \text{Ent}(C) \Rightarrow 0.029$$

* Calculate IGI of wind

Step 1: Entropy of entire dataset

$$H(+9, -5) = 0.94$$

Step 2: Entropy of all attributes

Entropy of Strong $\{+3, -3\} = -\frac{3}{6} \log_2 \frac{3}{6} - \frac{3}{6} \log_2 \frac{3}{6} = 1.0$

Entropy of ~~Strong~~ Weak $\{+6, -2\} = -\frac{6}{8} \log_2 \frac{6}{8} - \frac{2}{8} \log_2 \frac{2}{8} = 0.81$

$$\text{IGI} = \text{Entropy}(\text{whole data}) - \frac{6}{14} \text{Ent}(S) - \frac{8}{14} \text{Ent}(W) = 0.0478$$

* Calculate IGI of Humidity

Step 1: Entropy of entire dataset $H(+9, -5) = -\frac{9}{14} \log_2 \frac{9}{14} - \frac{5}{14} \log_2 \frac{5}{14} = 0.94$

Step 2: Entropy of all attributes:

Entropy of High $\{+3, -4\} = -\frac{3}{7} \log_2 \frac{3}{7} - \frac{4}{7} \log_2 \frac{4}{7} = 0.98$

" of Normal $\{+6, -1\} = -\frac{6}{7} \log_2 \frac{6}{7} - \frac{1}{7} \log_2 \frac{1}{7} = 0.59$

$$IG = \text{Entropy}(\text{whole data}) - \frac{7}{14} \text{Ent}(H) - \frac{7}{14} \text{Ent}(N)$$

$$= 0.15$$

$$IG(\text{weather}) = 0.246$$

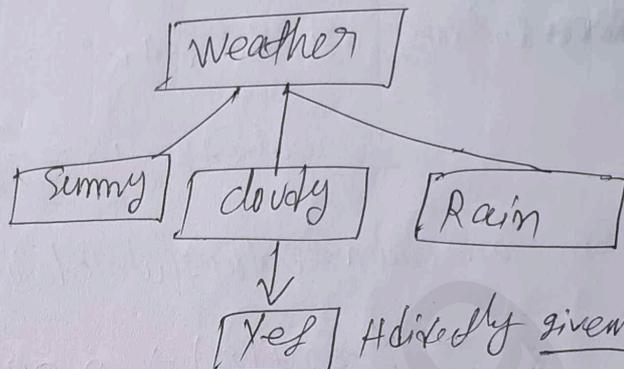
$$IG(\text{Temp}) = 0.029$$

$$IG(\text{Humidity}) = 0.15$$

$$IG(\text{wind}) = 0.0478$$

Root Node = Max of all

Root Node = Weather (0.246)



Now we will calculate IG from sunny point of view.

~~Calculate IG of~~

Day	weather	Temp	Humidity	wind	play
Day 1	Sunny	Hot	High	weak	No
Day 2	Sunny	Hot	High	strong	No
Day 3	Sunny	Mild	High	weak	No
Day 4	Sunny	cool	normal	weak	Yes
Day 5	Sunny	Mild	Normal	strong	Yes

* Calculate IG of Temperature

Step 1: Entropy of Sunny $\{+2, -3\} = -\frac{2}{5} \log_2 \frac{2}{5} - \frac{3}{5} \log_2 \frac{3}{5} = 0.97$

Step 2: Entropy of all attributes: (Temp)

Entropy of Hot $\{+2, -2\} = -\frac{2}{4} \log_2 \frac{2}{4} - \frac{2}{4} \log_2 \frac{2}{4} = 0$

" " Mild $\{+1, -1\} = -\frac{1}{2} \log_2 \frac{1}{2} - \frac{1}{2} \log_2 \frac{1}{2} = 1.0$

" " Cool $\{+1, -0\} = -\frac{1}{1} \log_2 \frac{1}{1} - \frac{0}{1} \log_2 \frac{0}{1} = 0$

$$IG = \text{Entropy}(\text{sunny}) - \frac{2}{5} \text{Ent}(H) - \frac{2}{5} \text{Ent}(M) - \frac{1}{5} \text{Ent}(C)$$

$$= 0.57$$

Calculate IG of Humidity

Step 1: Entropy of Sunny $\Rightarrow H(S+2-3) = 0.97$
 Step 2: Entropy of all attributes:

Entropy of High $S+0-3 = -\frac{0}{3} \log_2 \frac{0}{3} - \frac{3}{3} \log_2 \frac{3}{3} = 0$

" " Normal $S+2-0 = -\frac{2}{2} \log_2 \frac{2}{2} - \frac{0}{2} \log_2 \frac{0}{2} = 0$

IG = Entropy (Sunny) - $\frac{3}{5}$ Ent (H) - $\frac{2}{5}$ Ent (M) = 0.97

Calculate IG of Wind

Step 1: Entropy of Sunny $S+2-3 = 0.97$

Step 2: Entropy of all attributes

• Entropy of Strong $S+1-1 = -\frac{1}{2} \log_2 \frac{1}{2} - \frac{1}{2} \log_2 \frac{1}{2} = 1$

" " Weak $S+1-2 = -\frac{1}{3} \log_2 \frac{1}{3} - \frac{2}{3} \log_2 \frac{2}{3} = 0.918$

~~MI~~

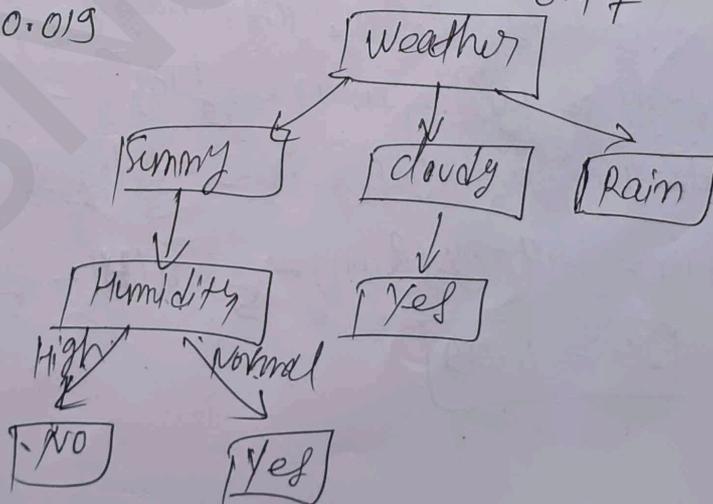
IG = Entropy (Sunny) - $\frac{2}{5}$ Ent (S) - $\frac{3}{5}$ Ent (W)
 = 0.019

IG (Temp) = 0.57

IG (Humidity) = 0.97

- IG (wind) = 0.019

Next Node (Humidity) 0.97



Day	weather	Temperature	Humidity	play/wind	play
Day 4	Rain	Mild	High	Weak	Yes
Day 5	Rain	cool	Normal	weak	Yes
Day 6	Rain	cool	Normal	strong	No
Day 10	Rain	Mild	Normal	strong	Yes
Day 14	Rain	Mild	High	Strong	No

* Calculate IG of Temperature

Step 1: Entropy of Rain $S+3, -2$ $\{ = -\frac{3}{5} \log_2 \frac{3}{5} - \frac{2}{5} \log_2 \frac{2}{5} = 0.97$

Step 2: Entropy of all attributes

- Entropy of Hot $S+0, -0$ $\{ = -\frac{0}{2} \log_2 \frac{0}{2} - \frac{0}{2} \log_2 \frac{0}{2} = 0$
- Entropy of Mild $S+2, -1$ $\{ = -\frac{2}{3} \log_2 \frac{2}{3} - \frac{1}{3} \log_2 \frac{1}{3} = 0.918$
- Entropy of cool $S+1, -1$ $\{ = -\frac{1}{2} \log_2 \frac{1}{2} - \frac{1}{2} \log_2 \frac{1}{2} = 1.0$

$$IG = \text{Entropy}(\text{Rain}) - \frac{0}{5} \text{Ent}(H) - \frac{3}{5} \text{Ent}(M) - \frac{2}{5} \text{Ent}(C)$$

$$IG = 0.019$$

* Calculate IG of Humidity

Step 1: Entropy of Rain $S+3, -2$ $\{ = -\frac{3}{5} \log_2 \frac{3}{5} - \frac{2}{5} \log_2 \frac{2}{5} = 0.97$

Step 2: Entropy of all attributes:

- Entropy of High $S+1, -1$ $\{ = -\frac{1}{2} \log_2 \frac{1}{2} - \frac{1}{2} \log_2 \frac{1}{2} = 1$
- Entropy of Normal $S+2, -1$ $\{ = -\frac{2}{3} \log_2 \frac{2}{3} - \frac{1}{3} \log_2 \frac{1}{3} = 0.918$

$$IG = \text{Entropy}(\text{Rain}) - \frac{2}{5} \text{Ent}(H) - \frac{3}{5} \text{Ent}(N)$$

$$IG = 0.019$$

* Calculate IGI of wind

Entropy of Rain $\{+3, -2\} = -\frac{3}{5} \log_2 \frac{3}{5} - \frac{2}{5} \log_2 \frac{2}{5} = 0.97$

Entropy of all attributes:

Entropy of strong $\{+0, -2\} = -\frac{0}{2} \log_2 \frac{0}{2} - \frac{2}{2} \log_2 \frac{2}{2} = 0$

Entropy of weak $\{+3, -0\} = -\frac{3}{3} \log_2 \frac{3}{3} - \frac{0}{3} \log_2 \frac{0}{3} = 0$

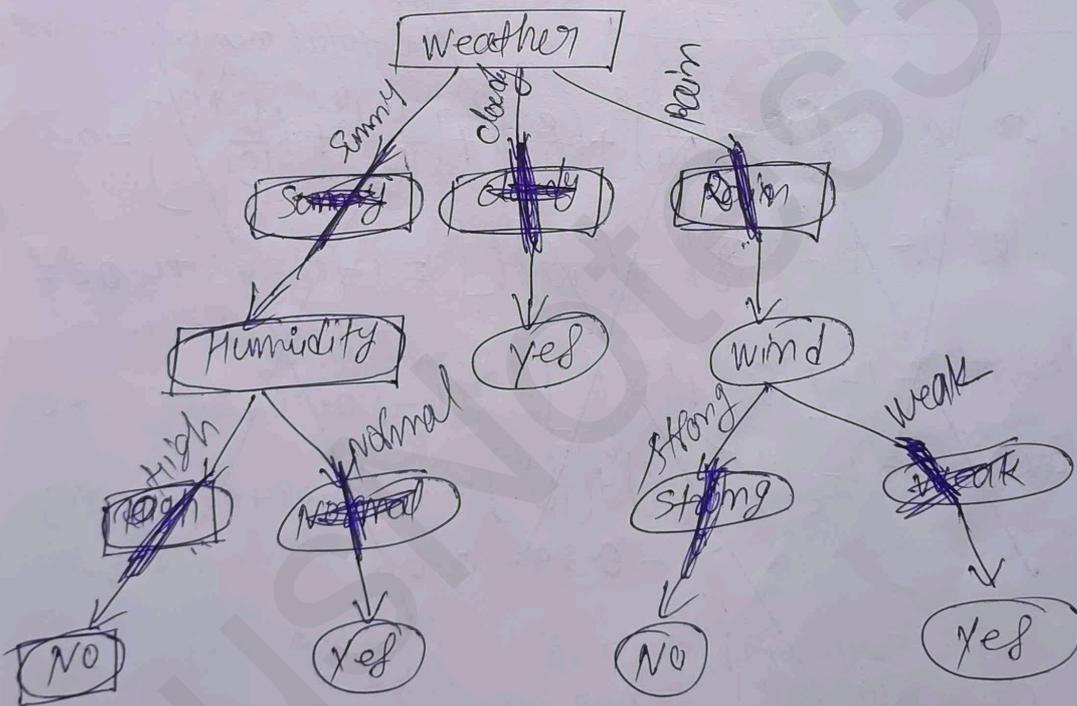
$IGI = Ent(Rain) - \frac{2}{5} Ent(S) - \frac{2}{5} Ent(W)$

$IGI = 0.97$

Now, $IGI(\text{temp}) = 0.019$

$IGI(\text{Hum}) = 0.019$

$IGI(\text{wind}) = 0.97 \checkmark$ (Leaf Node)



Weekend	weather	parents	Money	Decision
W1	Sunny	Yes	Rich	cinema
W2	Sunny	NO	Rich	Tennis
W3	windy	Yes	Rich	cinema
W4	Rainy	Yes	Poor	cinema
W5	Rainy	NO	Rich	stay in
W6	Rainy	Yes	Poor	cinema
W7	windy	NO	Poor	cinema
W8	windy	NO	Rich	shopping
W9	windy	Yes	Rich	cinema
W10	Sunny	NO	Rich	Tennis

There are 4 possible outcome variables ~~cinema~~, ~~Tennis~~, ~~stay in~~, ~~shopping~~, ~~weather~~, ~~parents~~, ~~Money~~, ~~Decision~~

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

p_i is the probability of class i .
 n is the total number of classes.

$$Gini\ (Decision) = 1 - \left[\left(\frac{6}{10}\right)^2 + \left(\frac{2}{10}\right)^2 + \left(\frac{1}{10}\right)^2 + \left(\frac{1}{10}\right)^2 \right] = 0.58$$

$$Gini\ (Money) = 1 - \left[\left(\frac{7}{10}\right)^2 + \left(\frac{3}{10}\right)^2 \right] = 1 - [0.49 + 0.09] = 1 - 0.58 = 0.42$$

$$Gini\ (Parents) = 1 - \left[\left(\frac{5}{10}\right)^2 + \left(\frac{5}{10}\right)^2 \right] = 1 - [0.25 + 0.25] = 0.5$$

$$Gini\ (Weather) = 1 - \left[\left(\frac{3}{10}\right)^2 + \left(\frac{4}{10}\right)^2 + \left(\frac{3}{10}\right)^2 \right] = 1 - [0.09 + 0.16 + 0.09] = 1 - 0.34 = 0.66$$

$$Gini\ (Decision) = 0.58$$

$$Gini\ (Money) = 0.42$$

$$Gini\ (Parents) = 0.5$$

$$Gini\ (Weather) = 0.66$$

In decision tree (classification) variables there are 4 possible out-comes
cinema, Tennis, stay-in, shopping

$$\text{Gini (Decision)} = \left[1 - \sum_{i=1}^n p_i^2 \right] \left\{ \begin{array}{l} p_i = \text{is the probability} \\ \text{of class } i. \\ n \text{ is the total number} \\ \text{of classes} \end{array} \right.$$

Gini for entire dataset

$$\text{Gini (Decision)} = \left[1 - \left(\frac{6}{10} \right)^2 + \left(\frac{2}{10} \right)^2 + \left(\frac{1}{10} \right)^2 + \left(\frac{1}{10} \right)^2 \right] = 0.58$$

there are three elements that affect Decision
So we will calculate Gini for money, parents,
weather then we will get our root Node.

① Gini (Money) \rightarrow 2 variables Rich \rightarrow 7 instances
Poor \rightarrow 3 instances

Gini (Rich) \Rightarrow 7 instances \rightarrow 3 cinema
 \rightarrow 2 Tennis
 \rightarrow 1 stay-in
 \rightarrow 1 shopping

$$\text{Gini (Rich)} = \left[1 - \left(\frac{3}{7} \right)^2 + \left(\frac{2}{7} \right)^2 + \left(\frac{1}{7} \right)^2 + \left(\frac{1}{7} \right)^2 \right] = 0.694$$

Gini (Poor) \Rightarrow 3 instances \rightarrow 3 cinema

$$\text{Gini (Poor)} = \left[1 - \left(\frac{3}{3} \right)^2 \right] = 0$$

weighted Avg Gini for Money = ~~Gini Idx~~ \times (proportion)

$$\text{Avg} = \left(0.694 \times \frac{7}{10} \right) + \left(0 \times \frac{3}{10} \right) = 0.486$$

* Gini (parents) \Rightarrow 2 variables \rightarrow Yes \rightarrow 5 Inst
 \rightarrow No \rightarrow 5 Inst

Gini (Yes) \Rightarrow 5 instances \rightarrow 5 Cinema

$$\boxed{\text{Gini(Yes)} = 0}$$

Gini (No) \Rightarrow 5 Inst \rightarrow

- \rightarrow 2 Tennis
- \rightarrow 1 shopping
- \rightarrow 1 stay-in
- \rightarrow 1 cinema

$$\text{Gini(No)} = \left[1 - \left(\frac{2}{5}\right)^2 - \left(\frac{1}{5}\right)^2 - \left(\frac{1}{5}\right)^2 - \left(\frac{1}{5}\right)^2 \right] =$$

$$\boxed{\text{Gini(No)} = 0.72}$$

Weight Avg Gini for parents $= \left(0 \times \frac{5}{10}\right) + \left(0.72 \times \frac{5}{10}\right)$

$$\text{Avg} = 0.36$$

* Gini (Weather) \Rightarrow 3 variables \rightarrow Sunny \rightarrow 3 inst
 \rightarrow Windy \rightarrow 4 inst
 \rightarrow Rainy \rightarrow 3 inst

Gini (Sunny) \Rightarrow 3 inst \rightarrow

- 1 Cinema
- 2 Tennis

$$\text{Gini(Sunny)} = \left[1 - \left(\frac{1}{3}\right)^2 - \left(\frac{2}{3}\right)^2 \right] = 0.444$$

Gini (Windy) \Rightarrow 4 inst \rightarrow

- 3 Cinema
- 1 shopping

$$\text{Gini(Windy)} = \left[1 - \left(\frac{3}{4}\right)^2 - \left(\frac{1}{4}\right)^2 \right] = 0.375$$

Gini (Rainy) \Rightarrow 3 inst \rightarrow

- 2 Cinema
- 1 stay-in

$$\text{Gini(Rainy)} = \left[1 - \left(\frac{2}{3}\right)^2 - \left(\frac{1}{3}\right)^2 \right] = 0.444$$

Weighted Avg Gini for weather $= \left(0.444 \times \frac{3}{10}\right) + \left(0.375 \times \frac{4}{10}\right) + \left(0.444 \times \frac{3}{10}\right)$

$$\text{Avg} = 0.416$$

Gini for Money = 0.486

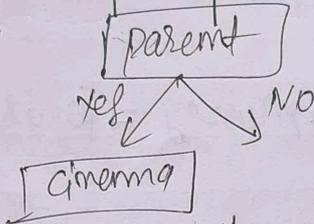
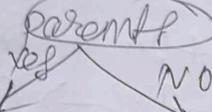
" " parent = 0.36

weather = 0.416

choose \rightarrow minimum

Gini for parent

Week	Weather	Parent	Money	Decision	Week	Weather	Parent	Money	Decision
W1	Sunny	Yes	Rich	cin	W2	Sunny	NO	Rich	Tena
W3	wind	Yes	Rich	cin	W5	Rain	NO	Rich	stay in
W4	Rain	Yes	poor	cin	W4	wind	NO	poor	cin
W8	Rain	Yes	poor	cin	W8	wind	NO	Rich	shop
W9	wind	Yes	Rich	cin	W10	summ	NO	Rich	Tena



* Gini (parent = No & weather)

Gini(Weather) \rightarrow 3 variable \rightarrow sunny \rightarrow 2 inf
 \rightarrow Rain \rightarrow 1 inf
 \rightarrow windy \rightarrow 2 inf

Gini(sunny) \Rightarrow 2 inf \rightarrow 2 ~~inf~~ Fermis

$$\text{Gini(sunny)} = 1 - \left[\frac{2}{3}\right]^2 = 0$$

Gini(Rainy) \Rightarrow 1 inf \rightarrow 1 stay-in

$$\text{Gini(Rainy)} = \left[1 - \frac{1}{3}\right]^2 = 0$$

Gini(Windy) = 2 inf \rightarrow 1 shopp
 \rightarrow 1 chem

$$\text{Gini(windy)} = \left[1 - \left(\frac{1}{2}\right)^2 + \left(\frac{1}{2}\right)^2\right] = 0.5$$

$$\text{Weighted Avg Gini for weather} = \left(0 \times \frac{2}{5}\right) + \left(0 \times \frac{1}{5}\right) + \left(0.5 \times \frac{2}{5}\right)$$

$$\text{Avg} = 0.2$$

Decision

Gini (parent = No & money)
 Gini (Money) \Rightarrow vari \rightarrow 4 Ing \rightarrow Rich
 \rightarrow 1 Ing \rightarrow Poor

Gini (Rich) \Rightarrow 4 ing \rightarrow 2 Term
 \rightarrow 1 shop
 \rightarrow 1 stay-in

$$Gini(Rich) = \sqrt{1 - \left(\frac{2}{4}\right)^2 - \left(\frac{1}{4}\right)^2 - \left(\frac{1}{4}\right)^2} = 0.625$$

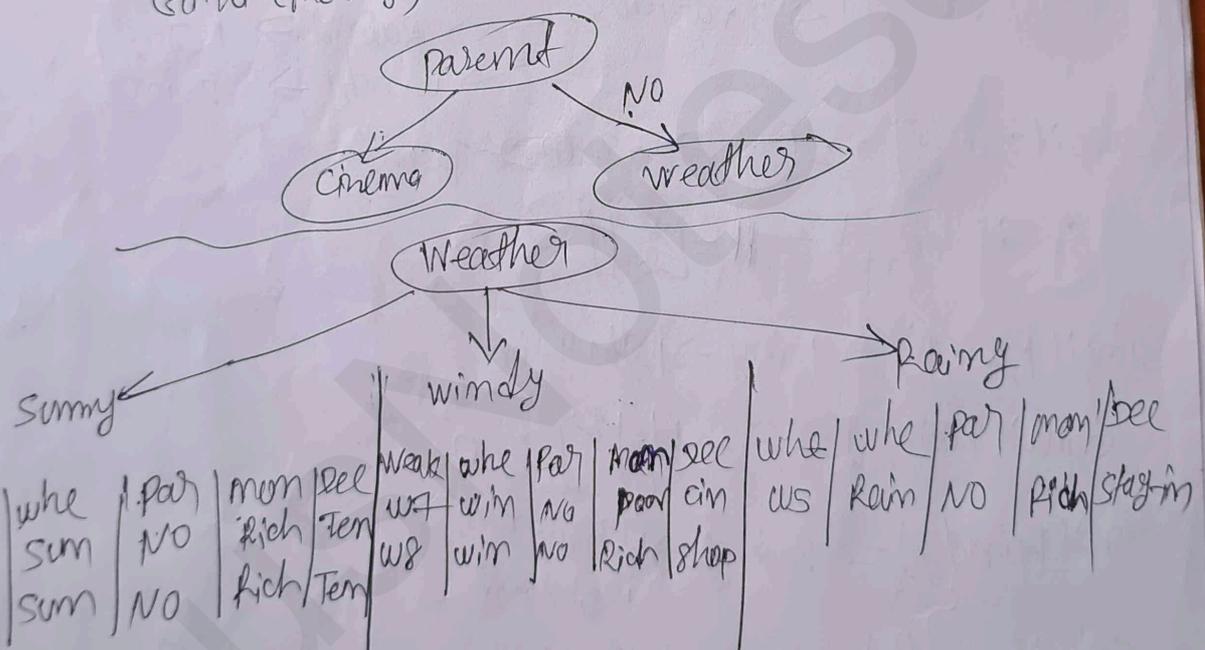
Gini (Poor) \Rightarrow 1 ing \rightarrow cinema

$$Gini(Poor) = 0$$

Weighted Avg Gini for Money = $(0.625 \times \frac{4}{5}) + (0 \times \frac{1}{5})$

$$Avg = 0.5$$

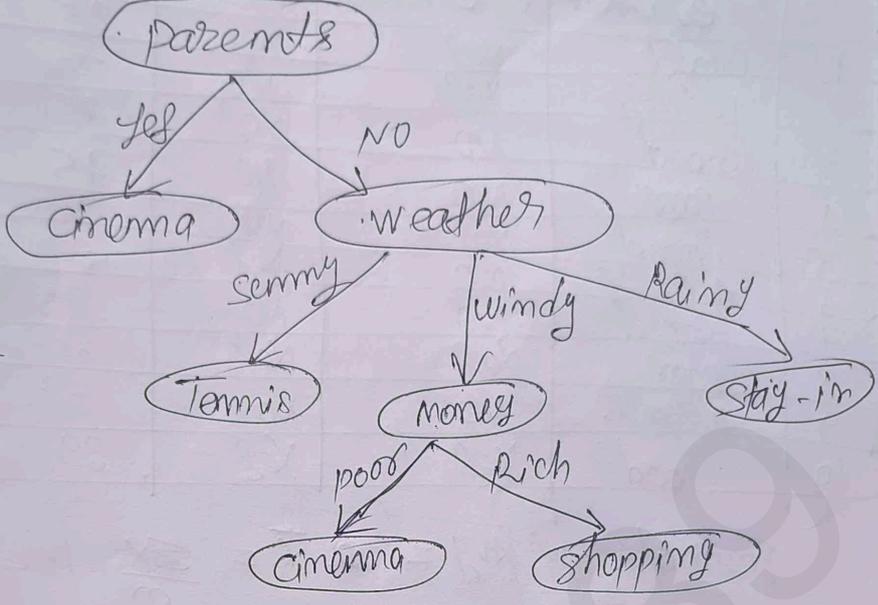
Gini (Weather) = 0.2 \checkmark Selected
 Gini (Money) = 0.5



Now decision tree is ready

~~Regression tree~~

Decision Tree:



see previous.

(Regression Tree)

SNO	Assessment	Assignment	Project	Result (%)
1	Good	Yes	Yes	95
2	Average	Yes	No	70
3	Good	No	Yes	75
4	Poor	No	No	45
5	Good	Yes	Yes	98
6	Average	No	Yes	80
7	Good	No	No	75
8	Poor	Yes	Yes	65
9	Average	No	No	58
10	Good	Yes	Yes	89

$$\text{Formula of Avg} = \bar{x} = \frac{\sum x}{n}$$

$$\text{Standard deviation } S = \sqrt{\frac{\sum (x - \bar{x})^2}{n}}$$

Computing standard deviation for each attribute with respect to the target attribute which is result.

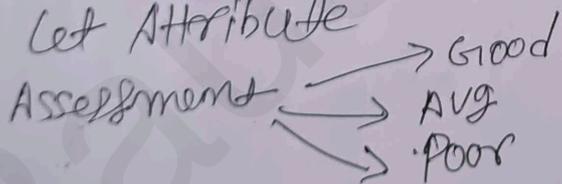
$$\text{Avg} = (95 + 70 + 75 + 45 + 98 + 80 + 75 + 65 + 58 + 89) / 10 = 75$$

Standard deviation for entire data set = 16.55

$$S = \sqrt{\frac{(95-75)^2 + (70-75)^2 + (75-75)^2 + (45-75)^2 + (98-75)^2 + (80-75)^2 + (75-75)^2 + (65-75)^2 + (58-75)^2 + (89-75)^2}{10}}$$

✓ $S = 16.55$ total std deviation

Let Attribute



Assessment = Good

$$\text{Avg} = (95 + 78 + 98 + 75 + 89) / 5 = 86.4$$

$$S = \sqrt{\frac{(95 - 86.4)^2 + (78 - 86.4)^2 + (98 - 86.4)^2 + (75 - 86.4)^2 + (89 - 86.4)^2}{5}}$$

$$S = 10.9$$

* Assessment = Average

$$\text{Avg} = (70 + 80 + 58) / 3 = 69.3$$

$$S = \sqrt{\frac{(70 - 69.3)^2 + (80 - 69.3)^2 + (58 - 69.3)^2}{3}}$$

$$S = 11.01$$

Assessment = poor

$$\text{Avg} = (45 + 65) / 2 = 55$$

$$S = \sqrt{\frac{(45 - 55)^2 + (65 - 55)^2}{2}} \quad S = 14.14$$

Assessment	S	instances
(Good)	10.9	5
(Avg)	11.01	3
(poor)	14.14	2

~~Sum of std dev~~
= ~~16.55~~

weighted standard deviation for Assessment:

$$\left(\frac{5}{10}\right) \times 10.9 + \left(\frac{3}{10}\right) \times 11.01 + \left(\frac{2}{10}\right) \times 14.14 = 11.58$$

standard deviation reduction for Assessment

$$\Rightarrow 16.55 - 11.58 = 4.97$$

* ~~standard~~ Assignment = Xef

$$\text{Avg} = (95 + 70 + 98 + 65 + 89) / 5 = 83.4$$

avg see previous.

$$\text{Standard dev} = \sqrt{\frac{(95-83.4)^2 + (70-83.4)^2 + (98-83.4)^2 + (65-83.4)^2 + (89-83.4)^2}{5}}$$

$$S = 14.98$$

Assignment = NO

$$\text{Avg} = (75 + 45 + 80 + 75 + 58) / 5 = 66.6$$

$$S = \sqrt{\frac{(75-66.6)^2 + (45-66.6)^2 + (80-66.6)^2 + (75-66.6)^2 + (58-66.6)^2}{5}}$$

$$S = 14.7$$

Standard deviation for Assignment

Assignment	S	impct of
Yes	14.98	.5
NO	14.7	.5

$$\text{Weighted standard deviation for App} \Rightarrow \left(\frac{.5}{1.0}\right) \times 14.98 + \left(\frac{.5}{1.0}\right) \times 14.7 = 14.84$$

Standard deviation reduction for ~~App~~ Assignment

$$\Rightarrow .1655 - 14.84 = 1.71$$

Project = Yes

$$\text{Avg} = 95 + 75 + 98 + 80 + 65 + 89 / 6 = 83.7$$

$$\text{Std dev} = \sqrt{\frac{(95-83.7)^2 + (75-83.7)^2 + (98-83.7)^2 + (80-83.7)^2 + (65-83.7)^2 + (89-83.7)^2}{6}}$$

$$S = 12.6$$

$$\text{Project} = \text{NO} \quad \text{Avg} = 70 + 45 + 75 + 58 / 4 = 62$$

$$\text{std dev} = \sqrt{\frac{(70-62)^2 + (48-62)^2 + (75-62)^2 + (98-62)^2}{4}}$$

$$\boxed{\text{std dev} = 13.39}$$

std dev for project

project	std deviation	data instances
Yes	12.6	6
No	13.39	4

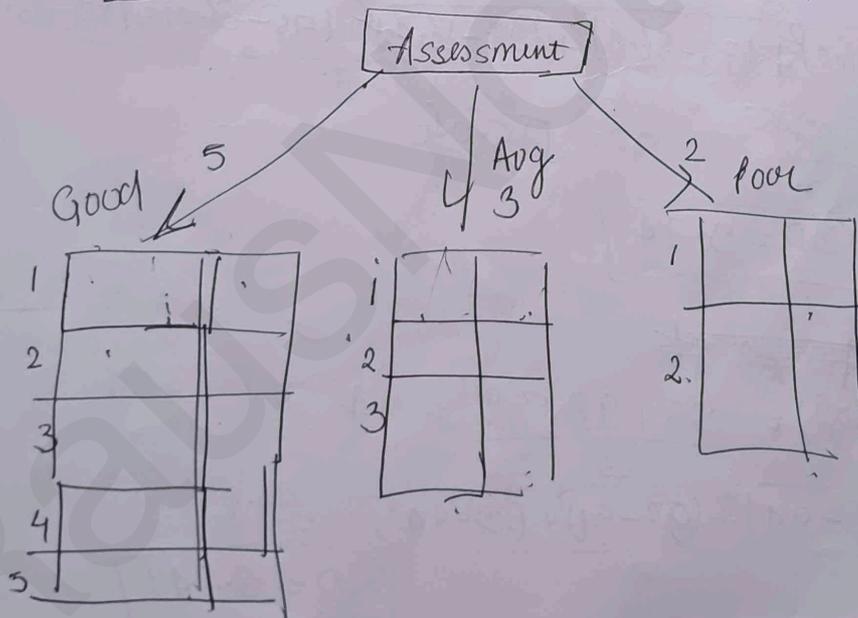
Weighted standard deviation for project

$$\Rightarrow \left(\frac{6}{10}\right) \times 12.6 + \left(\frac{4}{10}\right) \times 13.39 = 12.92$$

$$\text{std dev reduction for project} = 16.55 - 12.92 = 3.63$$

Standard deviation Reduction for each Att

Attribute	std dev Reduction	max = Assessment
Assessment	4.97	
Assignment	1.71	
Project	3.63	



see previous

name of

Assignment (yes) = 70 SD = 0

Assignment (No) = $\frac{80+58}{2} = 69$

SD = $\sqrt{\frac{(80-69)^2 + (58-69)^2}{2}}$ SD = 11

weighted SD for Assignment $\Rightarrow \frac{1}{3} \times 0 + (\frac{2}{3}) \times 11 =$

Reduction SD for Assignment $\Rightarrow 8.99 - 7.33 =$ 1.65

Project (yes) \Rightarrow Avg = 80 SD = 0

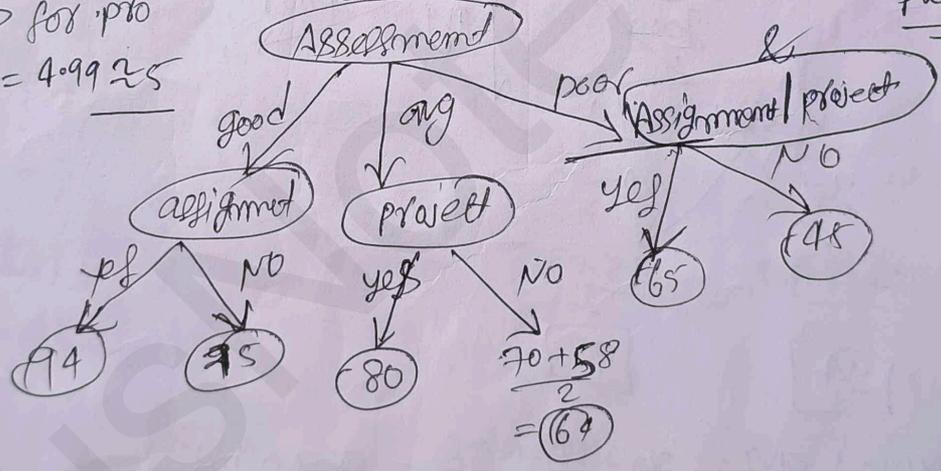
Project (No) = $\frac{70+58}{2} = 64$

SD = $\sqrt{\frac{(70-64)^2 + (58-64)^2}{2}} = \sqrt{\frac{72}{2}} = 6$

weighted SD for p80 = $(\frac{1}{3}) \times 0 + (\frac{2}{3}) \times 6 =$ 4

Reduction SD for p80 = $8.99 - 4 = 4.99 \approx 5$

Maximum project



$$\text{Assignment (No)} \Rightarrow \frac{75+75}{2} = 75$$

$$SD = \sqrt{\frac{(75-75)^2 + (75-75)^2}{2}} = 0$$

$$\text{Weighted SD for appi} \Rightarrow \left(\frac{3}{5}\right) \times 3.74 + 0 \times \left(\frac{2}{5}\right) = \cancel{2.748}$$

$$= 2.244 \quad \text{Reduced SD for Appi} = 9.74 - 2.244 = \boxed{7.496}$$

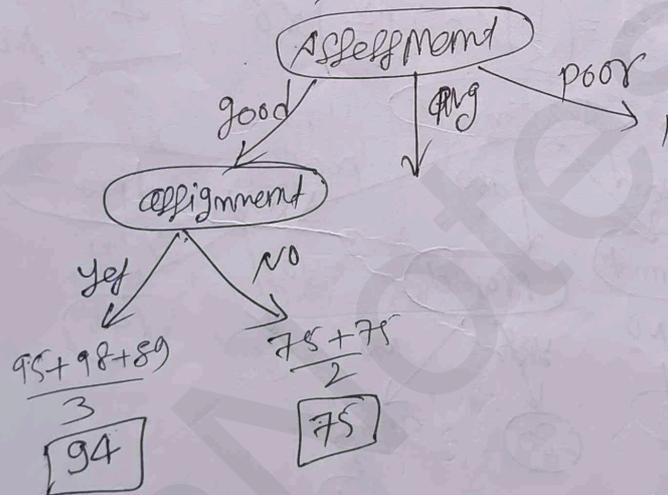
$$\text{Project (Yes)} \Rightarrow \text{Avg} = \frac{95+75+98+89}{4} = 89$$

$$SD = \sqrt{\frac{(95-89)^2 + (75-89)^2 + (98-89)^2 + (89-89)^2}{4}} = \sqrt{\frac{313}{4}} = 8.84$$

$$\text{Project (No)} = 75 \quad \boxed{SD=0}$$

$$\text{Weighted SD for Pro} \Rightarrow \left(\frac{4}{5}\right) \times 8.84 + \left(\frac{1}{5}\right) \times 0 = 7.072$$

$$\text{Reduction SD for Pro} \Rightarrow 9.74 - 7.072 = \boxed{2.668}$$



$$\textcircled{2} \text{ Average: } \text{avg} = \frac{70+80+58}{3} = \cancel{69.33}$$

$$SD = \sqrt{\frac{(70-69.3)^2 + (80-69.3)^2 + (58-69.3)^2}{3}}$$

$$SD = \sqrt{\frac{242.67}{3}} \quad \boxed{SD = 8.99}$$

... the previous.

Assessment

↓ Good

SN	Appel	Assigmn	Project	Refuel
1	good	Yes	Yes	95
3	good	NO	Yes	75
5	good	Yes	Yes	98
7	good	NO	NO	75
10	good	Yes	Yes	89

↓ Avg

SN	Appel	Assigmn	Project	Refuel
2	Avg	Yes	NO	70
6	Avg	NO	Yes	80
9	Avg	NO	NO	58

Poor

SN	Appel	Assigmn	Project	Refuel
4	Poor	NO	NO	45
8	Poor	Yes	Yes	65

Need to calculate

① ~~NO~~ Good

$$\text{Avg} = \frac{95 + 75 + 98 + 75 + 89}{5} = \frac{432}{5} = 86.4$$

$$\text{S.D. } \sigma = \sqrt{\frac{(95-86.4)^2 + (75-86.4)^2 + (98-86.4)^2 + (75-86.4)^2 + (89-86.4)^2}{5}}$$

$$\text{SD} = \sqrt{\frac{475.2}{5}} \quad \text{SD} = \sqrt{95.04}$$

$$\boxed{\text{SD} = 9.74}$$

~~Assessment~~ Assignment

$$\text{Assignment (Yes)} = \frac{95 + 98 + 89}{3} = 94$$

$$\text{SD} = \sqrt{\frac{(95-94)^2 + (98-94)^2 + (89-94)^2}{3}} = \dots$$

$$\boxed{\text{SD} = 3.741}$$

(P20)

Random forest Describe RF algo to improve accuracy. improve classifier

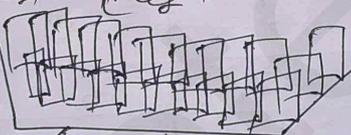
Random forest: Random forest is a supervised machine learning algorithm that is used for both classification and regression. It is an ensemble learning technique that creates multiple learning decision trees and merges them together to ~~predict~~ produce a more accurate and stable prediction.

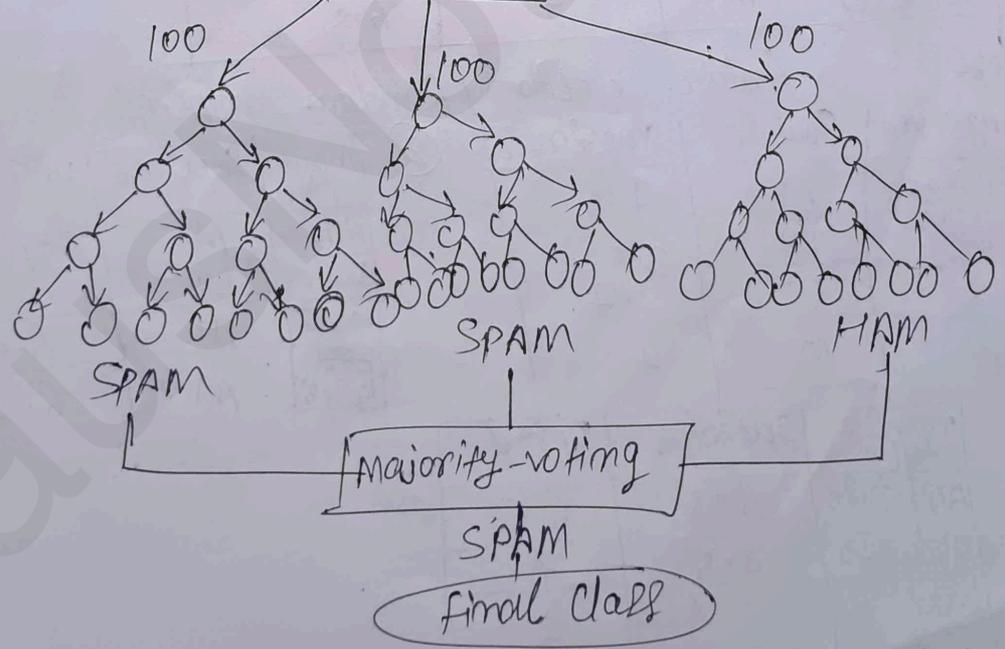
steps of Algorithm

step 1: Create Bootstrap dataset from original data by randomly choosing data (repetition is allowed)

step 2: Create Randomized Decision Tree from Bootstrap Dataset.

step 3: Finally output of the random forest is the class selected by most trees.

Example:  300 emails



the previous.

the same of

Features of Random forest

- ① It takes less training time as compared to other algorithms.
- ② It predicts output with high accuracy, even for the large dataset it runs efficiently.
- ③ It can also maintain accuracy when a large proportion of data is missing.
- ④ It can handle both classification and Regression.
- ⑤ It reduces overfitting.
- ⑥ It combines multiple decision trees to improve accuracy.
- ⑦ Low variance
- ⑧ High accuracy, no need of normalization.

Ex →

MangoType	Source	Sweetness	Diameter
Alphanso	Maharashtra (MH)	4.5	7
Alphanso	MH	4.5	6.5
Desheri	UP	3.8	9
Desheri	UP	3.6	7.9
Kesari	MH	3.7	6
Kesari	MH	4	6.5

123, 124, 134

Table 1: Feature Selection / Random Sampling

①

MangoType	Source	Sweetness
Alphanso	M	4.5
Desheri	UP	3.8
Kesari	M	3.7
Kesari	M	4

②

MangoType	Source	Diameter
Alphanso	M	7
Desheri	UP	7.9
Kesari	M	6.5
Kesari	M	6.5

③

MangoType	Sweetness	Diameter
Alphanso	4.5	6.5
Alphanso	4.5	7
Kesari	3.7	6
Desheri	3.8	9

To predict India will win or loose by watching
 run of virat kohli.

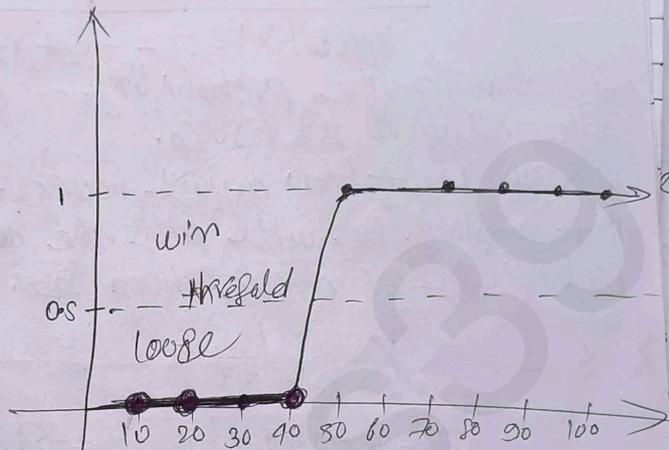
$$y = \frac{1}{1 + e^{-x}} = y = \frac{1}{1 + e^{-(B_0 + B_1 x)}}$$

B_0 = Intercept

B_1 = slope

Q.

Run (x)	Win or loose (y)
10	0
20	0
40	0
50	1
70	1
80	1
90	1
100	1



In our example the run above threshold tends to win
 the run below threshold tends india will loose.

Advantages of Logistic Regression.

- (i) performance better when data is linearly separable.
- (ii) it does not required computational cost
- (iii) does not require fine tuning tuning.
- (iv) Easy to implement and train the model.
- (v) Give the size of how relevant.

one the previous. ... must some of

Logistic Regression (classification alg)

~~Linear Regression~~

Logistic Regression: Logistic Regression is a supervised learning algorithm, it is a classification algorithm, used for predicting the output of a categorical dependent variable. The output can be either 'yes' or 'no', it gives the probability b/w 0 and 1.

Ex \Rightarrow disease, fraud detection, spam detection.

Types of Logistic Regression

- (i) Binary (0 or 1)
- (ii) Multinomial (three or more ordered type) [A, B, C]

Sigmoid function $y = \frac{1}{1 + e^{-x}}$ y is the dependent variable

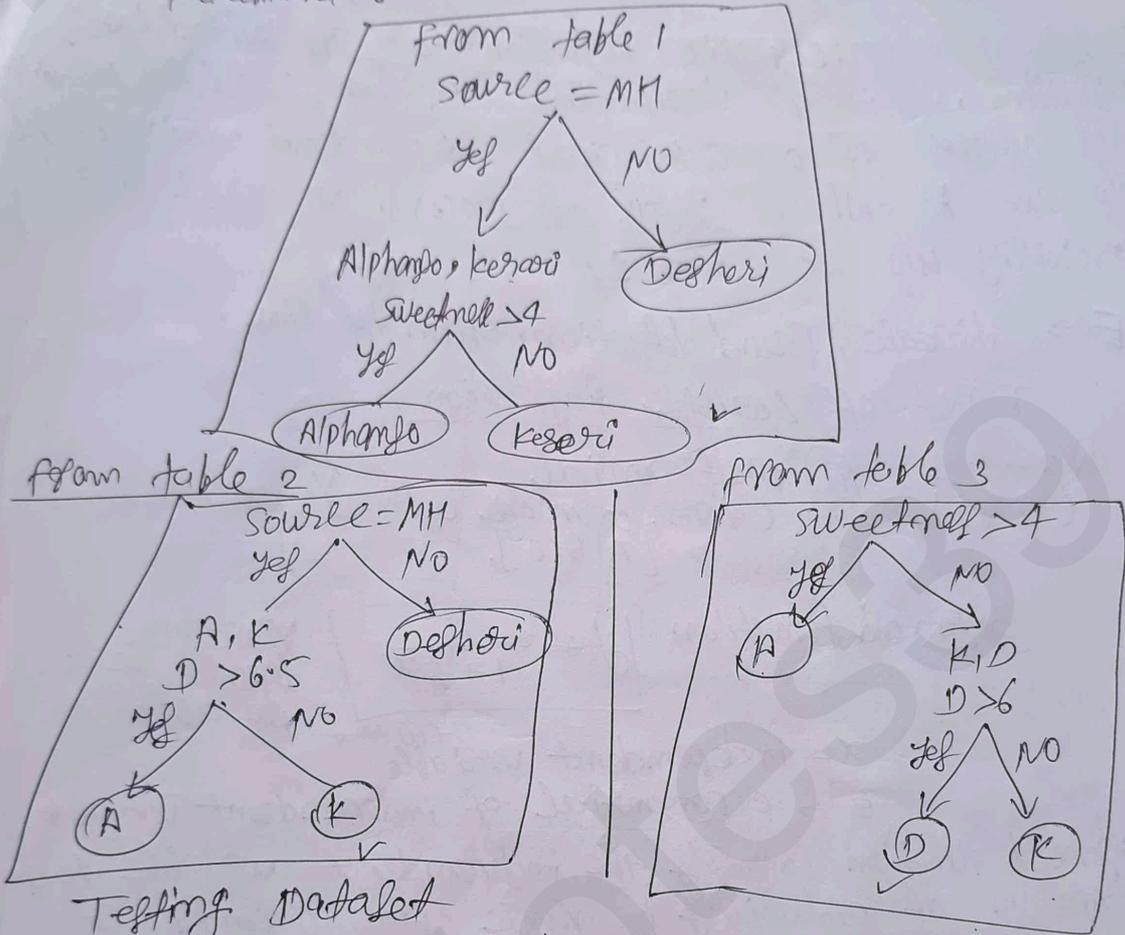
$x \rightarrow$ independent variable

$e^x \rightarrow$ exponential of independent variable

Sigmoid function converts the independent variable to an expression of probability in range b/w 0 and 1. It maps real value b/w 0 and 1. Target variable is must be binary and desired outcome is represented by factor level 1. No multi-collinearity in the model which means independent variable of each must be independent of each other.

Lg transform its output using logistic sigmoid function where y is a independent with the help of sigmoid function we are able to reduce a loss during time of training because it eliminates the gradient problems in the ml model during the training.

This particular process of creating Bootstrap training dataset is known as random sampling with replacement.



Testing Dataset

Source	Sweetness	Diameter
MH	3.9	6.4

According to first ~~tree~~ tree type = kesari
 " " second " " = kesari
 " " third " " = Daghari

Ans = kesari

one more previous.

The new point is classified as category 2 because most of its closest neighbours are blue squares. KNN assigns the category based on the majority of nearby points.

- The black diamonds represent category 1 and blue squares represent category 2.
- The new data point checks its closest neighbours.
- Since the majority of its closest neighbours are blue squares (category 2) KNN predicts the new data point belongs to category 2.

Working of KNN algorithm

Step 1: Selecting the optimal value of k (nearest neighbours)

Step 2: Calculating distance: Euclidean distance

$$\sqrt{|x_1 - y_1|^2 + |x_2 - y_2|^2}$$

Step 3: Finding Nearest Neighbours.

Step 4: Voting for classification or Taking Average for Regression.

Example:

IMDb Rating	Duration	Genre
8.0 (Mission Impo)	160	Action
8.2 (Gadar 2)	170	Action
7.2 (Rocky & Rani)	168	Comedy
8.2 (OMG 2)	155	Comedy

Now predict the genre of Barbie movie with IMDb rating ~~7.4~~ 7.4 and duration 114 minutes.

Distance between cluster centers and data-points

Data points	Distance from $K_1 (2, 1.75)$	Distance from $K_2 (4.5, 4)$	Assigned center
$q_1 (1, 1)$	1.25	4.61	K_1
$q_2 (2, 1)$	0.75	3.9	K_1
$q_3 (2, 3)$	1.25	2.69	K_1
$q_4 (3, 2)$	1.03	2.5	K_1
$q_5 (4, 3)$	2.36	1.12	K_2
$q_6 (5, 5)$	4.42	1.12	K_2

Euclidean distance

$$\sqrt{(x_1 - y_1)^2 + (x_2 - y_2)^2}$$

cluster 1 of $K_1 = \{q_1, q_2, q_3, q_4\}$

cluster 2 of $K_2 = \{q_5, q_6\}$

cluster elements are same as in the previous iteration

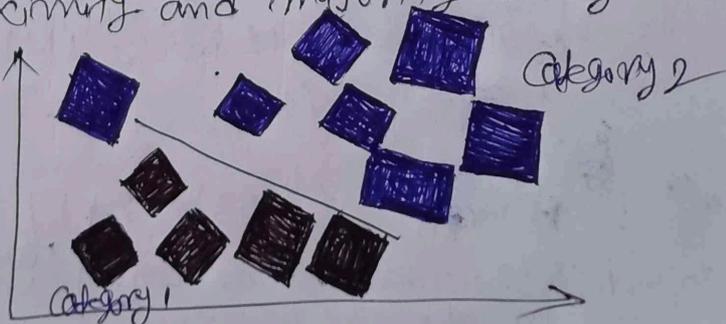
final cluster is:

cluster 1: $\{(1, 1), (2, 1), (2, 3), (3, 2)\}$

cluster 2: $\{(4, 3), (5, 5)\}$

(2) **k-Nearest Neighbour (k-NN)**: k-Nearest Neighbour (k-NN) is a supervised machine learning algorithm used for classification and regression. It is a lazy algorithm, meaning it does not explicitly learn a model during training. Instead, it memorizes the training dataset and makes decisions only during prediction. k-NN works by using proximity and majority voting to make predictions.

Ex:



Step 5: Repeat from step 2 until we get same cluster center or same cluster elements as in the previous iteration.

Distance table :-

Data point	Distance from $k_1 (2, 1.33)$	Distance from $k_2 (3.67, 3.67)$	Assigned center
$q_1 (1, 1)$	1.05	3.78	k_1
$q_2 (2, 1)$	0.33	3.15	k_1
$q_3 (2, 3)$	1.67	1.8	k_1
$q_4 (3, 2)$	1.204	1.8	k_1
$q_5 (4, 3)$	2.605	0.75	k_2
$q_6 (5, 5)$	4.74	1.88	k_2

cluster 1 of $k_1 = \{q_1, q_2, q_3, q_4\}$ } calculate ~~but not~~ Euclidean distance for cluster 1
 cluster 2 of $k_2 = \{q_5, q_6\}$ } $\sqrt{(x_1 - x_2)^2 + (y_1 - y_2)^2}$
 $\sqrt{(4-2)^2 + (1-1.33)^2} = 1.05$
 cluster elements are ^{not} same as in the previous iteration

cluster 1: $\{(1, 1), (2, 1), (2, 3), (3, 2)\}$ } for cluster 2
 $E_d = \sqrt{(1-3.67)^2 + (1-3.67)^2}$
 $E_d = 3.78$
 cluster 2: $\{(4, 3), (5, 5)\}$

cluster 1: $\{q_1, q_2, q_3, q_4\}$

cluster 2: $\{q_5, q_6\}$

Recalculating the cluster centers

$$k_1 = \frac{1}{4} [q_1 + q_2 + q_3 + q_4] \Rightarrow \frac{1}{4} [(1, 1) + (2, 1) + (2, 3) + (3, 2)]$$

$$= \frac{1}{4} [8, 7] = (2, 1.75)$$

$$k_2 = \frac{1}{2} [q_5 + q_6] = \frac{1}{2} [(4, 3) + (5, 5)] \Rightarrow \frac{1}{2} (9, 8) = (4.5, 4)$$

So cluster elements and centers are not same as in the previous.

3. We repeat the process for a given number of iterations and at the end, we have our clusters.

Q. Use k-means clustering algorithm to divide the following data into two clusters.

x_1	1	2	2	3	4	5
x_2	1	1	3	2	3	5

ans \Rightarrow step 1: choosing randomly 2 cluster centers.

say $K_1 = (2, 1)$ $K_2 = (2, 3)$

step 2: finding the distance b/w the cluster centers and each data points. Euclidean dist = $\sqrt{(x_1 - y_1)^2 + (x_2 - y_2)^2}$

Data point	Distance from $K_1(2, 1)$	Distance from $K_2(2, 3)$	Assigned cluster	We have to find distance b/w each center and given data point $(x_1, y_1) = (1, 1)$ $(y_1, y_2) = (2, 1)$ for cluster
$q_1 (1, 1)$	1	2.24	K_1	
$q_2 (2, 1)$	0	2	K_1	
$q_3 (2, 3)$	2	0	K_2	
$q_4 (3, 2)$	1.41	1.41	K_1	
$q_5 (4, 3)$	2.83	2	K_2	
$q_6 (5, 5)$	5	3.61	K_2	

step 3: cluster 1 of $K_1 = \{q_1, q_2, q_4\}$

cluster 2 of $K_2 = \{q_3, q_5, q_6\}$

step 4: Recalculate the cluster centers.

$$K_1 = \frac{1}{3} [(1, 1) + (2, 1) + (3, 2)] = \frac{1}{3} (6, 4) = (2, 1.33)$$

$$K_2 = \frac{1}{3} [(2, 3) + (4, 3) + (5, 5)] = \frac{1}{3} (11, 11)$$

$$K_2 = \frac{1}{3} (11, 11) = (3.67, 3.67)$$

Euclidean distance = $\sqrt{(x_1 - y_1)^2 + (x_2 - y_2)^2}$
 $= \sqrt{(4-2)^2 + (1-1)^2}$
 $= \sqrt{1} = 1$ and so on

Step 5: repeat the process

K-Means clustering

* K - no of cluster/groups

Clustering: dividing the dataset into diff groups.

Clustering: clustering is the task of dividing a dataset into groups (clusters) such that data points in the same group are more similar to each other than to those in other groups. It is used when we don't have labeled data.

What is K-Means clustering?

as \Rightarrow K-Means clustering is an unsupervised machine learning algorithm which groups the unlabeled dataset into different clusters. K-Means clustering is a technique used to organize data into groups based on their similarity.

Example: Online store uses K-Means to group customers based on purchase frequency and spending, creating segments like Budget shoppers, frequent buyers and Big spenders for personalized marketing.

The algorithm works by first randomly picking some central points called 'centroids' and each data point is then assigned to the closest centroid forming a cluster.

The algorithm will categorize the items into K-group or clusters of similarity. To calculate that similarity, we will use the Euclidean distance of a measurement.

Working of K-Means clustering

1. First, we randomly initialize K points, called mean or cluster centroid.
 2. We categorize each item to its closest mean, and we update the mean's coordinate, which are the averages of the items categorized in that cluster so far.
- one the previous. \Rightarrow are not same of

Q. Calculate the probability of pass for the student who studied 33 hr.

Hours study	pass (1) / Fail (0)
29	0
15	0
33	1
28	1
39	1

$$\log(\text{odds}) = z = -64 + 2 \times \text{hours}$$

(i) Calculate the probability of pass for the student who studied 33 hours.

(ii) At least how many hours student should study that makes he will pass the course with the probability of more than 95%.

solⁿ: (i) $y = \frac{1}{1 + e^{-z}} = y = \frac{1}{1 + e^{-z}}$

$$z = -64 + 2 \times \text{hr} = -64 + 2 \times 33$$

$$\boxed{z = 2}$$

$$p = \frac{1}{1 + e^{-z}} \quad \boxed{p = 0.88}$$

Conclusion: If student studied 33 hours; then there is 88% chance that the student will pass the exam.

(ii) $p = \frac{1}{1 + e^{-z}} = 0.95$

$$0.95 * (1 + e^{-z}) = 1$$

$$0.95 * e^{-z} = 1 - 0.95$$

$$e^{-z} = \frac{0.05}{0.95}$$

$$e^{-z} = 0.0526$$

$$\ln(e^{-z}) = \ln(0.0526)$$

$$\ln(e^z) = x$$

$$-z = \ln(0.0526) = -2.94$$

$$\boxed{z = 2.94}$$

$$\log(\text{odds}) =$$

$$z = -64 + 2 \times \text{hr}$$

$$z = 2.94 = -64 + 2 \times \text{hr}$$

$$\boxed{\text{hr} = 33.87 \text{ hr}}$$

* Naive Bayes IML-2

Def \Rightarrow Naive Bayes classifiers are supervised machine learning algorithms used for classification tasks, based on Bayes theorem to find probabilities. It is named as 'naive' because it assumes the presence of one feature does not affect other features. The 'Bayes' part of the name refers to ~~the~~ the basis in Bayes theorem.

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

Features of Naive Bayes

- (i) Based on Bayes Theorem: Naive Bayes theorem uses Bayes Theorem to compute the probability of each class given a set of input features.
- (ii) Naive Assumption of feature independence: It assumes all features are independent of each other, given the class label.
- (iii) Simple and fast: Naive Bayes is a simple and powerful algorithm.
- (iv) Lower number of parameters: It requires fewer parameters which makes it efficient and easy to train.
- (v) ~~App~~ Application: Spam filtering, Recommendation system.

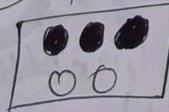
(1) Independent Event

Rolling a dice
 $\{1, 2, 3, 4, 5, 6\}$

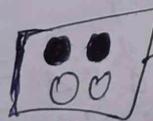
$$P(1) = \frac{1}{6} \quad P(2) = \frac{1}{6}$$
$$P(3) = \frac{1}{6}$$

Tossing a coin = $\frac{1}{2}$

(2) Dependent Event



$$\rightarrow P(B) = \frac{3}{5}$$



$$\rightarrow P(B) = \frac{2}{4}$$

$$P(B) = \frac{3}{5} \times \frac{2}{4}$$

height (cm)	weight (kg)	class	Distance
169	58	normal	1.4
170	55	normal	2
173	57	normal	3
174	56	underweight	4.1
167	51	underweight	6.7
173	64	normal	7.6
172	65	normal	8.2
182	62	normal	13
176	69	normal	13.4
170	57	?	

if $k=1$ normal

$k=2$ normal

$k=3$ = normal

$k=4$ underweight
but majority
is = Normal

$k=5$ = Normal

? = Normal

Step 1: calculate distances

$$\text{distance to } (8.0, 160) = \sqrt{(7.4 - 8)^2 + (114 - 160)^2} = 46$$

$$\text{distance to } (6.2, 160) = \sqrt{(7.4 - 6.2)^2 + (114 - 170)^2} = 56.01$$

$$\text{to } (7.2, 168) = \sqrt{(7.4 - 7.2)^2 + (114 - 168)^2} = 54.00$$

$$\text{to } (8.2, 155) = \sqrt{(7.4 - 8.2)^2 + (114 - 155)^2} = 41.00$$

Step 2: Select k Nearest Neighbours

k=1	k=3
best distance = 41	41, 46, 54
41 = dist(8.2, 155)	↓ ↓ ↓
OMG2 = Comedy	(OMG2) (MI) (R&P)
	(Comedy) (Action) (Comedy)

Step: Majority voting → Comedy

hence, Barbie movie with IMDb rating 7.4 and duration 114 minutes, its genre is "Comedy"

Q.

Height (cm)	weight (kg)	class
167	51	underweight
182	62	Normal
176	69	Normal
173	64	Normal
172	65	Normal
174	86	underweight
169	58	Normal
173	57	Normal
170	55	Normal
170	57	?

Q. Estimate conditional probabilities of each attribute (color, legs, height, smelly) for the species classes $\{M, H\}$ using the data given in the table. Use these probabilities estimate the probability values for the new instance (color = Green, legs = 2, Height = Tall and Smelly = No)

No	color	legs	height	Smelly	Species
1	white	3	short	yes	M
2	green	2	Tall	NO	M
3	Green	3	short	yes	M
4	white	3	short	yes	M
5	Green	2	short short	NO	M
6	white	2	Tall	NO	H
7	white	2	Tall	NO	H
8	white	2	short	yes	H

$$P(M) = \frac{4}{8} = 0.5, P(H) = \frac{4}{8} = 0.5$$

New Instance

(color = Green, legs = 2, Height = Tall, smelly = No)

Color	M	H	legs	M	H	Height	M	H	Smelly	M	H
white	2/4	3/4	2	1/4	4/4	Tall	1/4	2/4	yes	3/4	1/4
Green	2/4	1/4	3	3/4	0/4	short	3/4	2/4	NO	1/4	3/4

$$P(M | \text{New Instance}) = P(M) * P(\text{color} = \text{Green} | M) * P(\text{legs} = 2 | M) * P(\text{Height} = \text{tall} | M) * P(\text{Smelly} = \text{no} | M)$$

$$P(M | \text{New Instance}) = 0.5 * \frac{2}{4} * \frac{1}{4} * \frac{1}{4} * \frac{1}{4} = \cancel{0.017} 0.00390$$

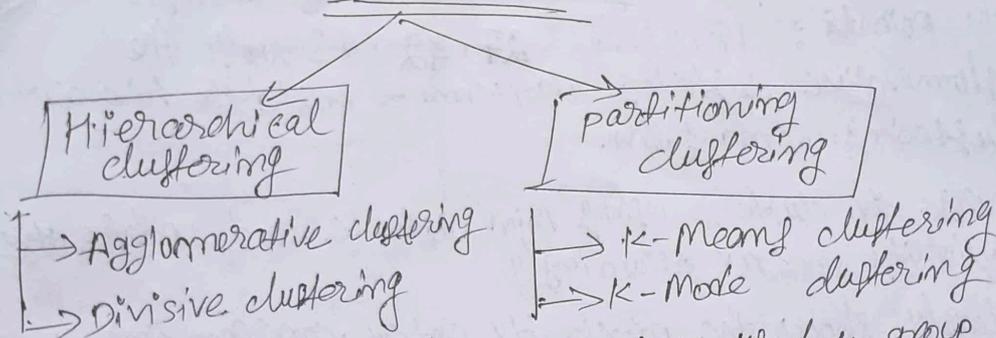
$$P(H | \text{New Instance}) = P(H) * P(\text{color} = \text{Green} | H) * P(\text{legs} = 2 | H) * P(\text{Height} = \text{tall} | H) * P(\text{Smelly} = \text{no} | H)$$

$$P(H | \text{New Instance}) = 0.5 * \frac{1}{4} * \frac{1}{4} * \frac{2}{4} * \frac{3}{4} = 0.0468$$

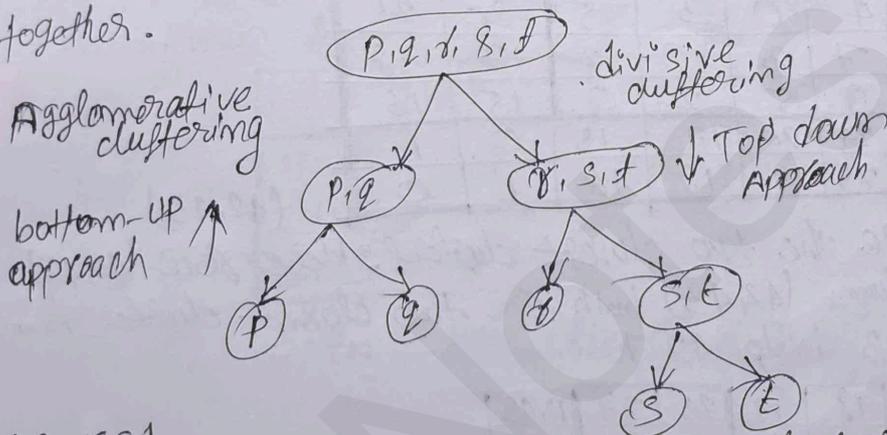
$$\therefore P(H | \text{New Instance}) > P(M | \text{New Instance})$$

Hence the new instance belongs to species H

Clustering



Hierarchical clustering is a technique used to group similar data points together based on their similarity. It is creating a hierarchy or tree-like structure. A dendrogram is like a family tree for clusters. It shows how individual data points or groups of data merge together.



(AGNES)
 Agglomerative clustering: Agglomerative clustering is a type of hierarchical clustering that builds clusters by merging them successively. It is a bottom-up approach, meaning

1. Each data points starts in its own cluster.
2. The algorithm repeatedly merges the two closest clusters

until all points are in a single cluster.
 (DIANA)
 Divisive clustering: Divisive clustering is the opposite of Agglomerative clustering - it is a top-down hierarchical clustering approach.

Iteratively compare the ^{leader} cluster data points to each of the observations. similar data points give 0; dissimilar data points give 1

Comparing leader p1 to the observation p2

blonde | amber | fair
brunette | gray | brown
+1 | +1 | +1 = 1+1+1 = 3

like will calculate all the dissimilarities and put them in a matrix...

	cluster 1 (P1)	cluster 2 (P3)	cluster 3 (P8)	cluster
P1	0 ✓ (minimum)	2	2	
P2	3 ✓	3	3	cluster 1
P3	3	1 ✓	3	cluster 2
P4	3	3	1 ✓	cluster 3
P5	1 ✓	2	2	cluster 1
P6	3	3	2 ✓	cluster 3
P7	2	0 ✓	2	cluster 2
P8	2	2	0 ✓	cluster 3

After step 2; the observations p1, p2, p5 are assigned to cluster 1, p3, p7 to cluster 2 and p4, p6, p8 are assigned to cluster 3

✓ → minimum in row
Same → any list

Step 3: Define new model for the clusters

	hair color	eye color	skin color
cluster 1	brunette	amber	fair
cluster 2	red	green	fair
cluster 3	black	hazel	brown

Repeat step steps 2-4

Q. Imagine we have a dataset that has the information about hair color, eye color and skin color of persons. We aim to group them based on the available information.

Person	hair color	eye color	skin color
P1	blonde	amber	fair
P2	brunette	gray	brown
P3	red	green	brown
P4	black	hazel	brown
P5	brunette	amber	fair
P6	black	gray	brown
P7	red	green	fair
P8	black	hazel	fair

Step 1: Pick k observations at random and use them as leaders / clusters. choosing $k=3$ P1, P7, P8 as leaders

Leaders

P1	blonde	amber	fair
P7	red	green	fair
P8	black	hazel	fair

Person	hair color	eye color	skin color
P1	blonde	amber	fair
P2	brunette	gray	brown
P3	red	green	brown
P4	black	hazel	brown
P5	brunette	amber	fair
P6	black	gray	brown
P7	red	green	fair
P8	black	hazel	fair

Step 2: calculate the dissimilarities (no. of mismatches) and assign each observation to its closest cluster.

	18	22	25, 27	42, 43
18	0	4	7	24
22	4	0	3	20
25, 27	7	3	0	15
42, 43	24	20	15	0

(42, 43), (25, 27), ~~(18, 22)~~

Step 4: Repeat until
 The distance b/w 22 and (25, 27) is minimum so merge (25, 27) into 22.

	18	22, 25, 27	42, 43
18	0	4	24
22, 25, 27	4	0	20
42, 43	24	20	0

(42, 43), ((25, 27), 22)

Step 5: The distance b/w (22, 25, 27) and 18 is minimum so merge (22, 25, 27) into 18.

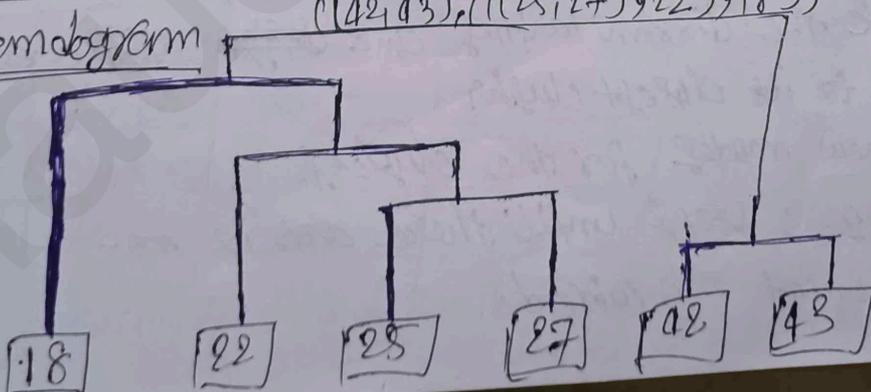
	18, 22, 25, 27	42, 43
18, 22, 25, 27	0	24
42, 43	24	0

((42, 43), ((25, 27), 22), 18)

Step 6: The distance b/w (42, 43) and (18, 22, 25, 27) is minimum so merge (42, 43) into (18, 22, 25, 27)

	18, 22, 25, 27, 42, 43
18, 22, 25, 27, 42, 43	0

• Dendrogram ((42, 43), ((25, 27), 22), 18)



Q. Consider the following set of 6 one dimensional data points: 18, 22, 25, 27, 42, 43. Apply the
 (i) agglomerative clustering algorithm to build the hierarchical clustering dendrogram.

(ii) merge the clusters using min distance and update the proximity matrix accordingly

(iii) clearly show the proximity matrix corresponding to each iteration of the algorithm.

Solⁿ:

Step 1: Compute distance between all pairs of clusters

	18	22	25	27	42	43
18	0	4	7	9	24	25
22	4	0	3	5	20	21
25	7	3	0	2	17	18
27	9	5	2	0	15	16
42	24	20	17	15	0	1
43	25	21	18	16	1	0

(42, 43)

Step 2: merge the two closest clusters: Here row (42, 43) and column (42, 43) both are two closest clusters hence merge 43 into 42 here.

	18	22	25	27	42, 43
18	0	4	7	9	24
22	4	0	3	5	20
25	7	3	0	2	17
27	9	5	2	0	15
42, 43	24	20	17	15	0

(42, 43) (~~27, 27~~)

Step 3: update the distance matrix: The distance b/w 27 and 25 is minimum so merge 25 and 27 and in a cluster, Merge 27 into 25.

Self-organizing Maps (SOMs)

Self-organizing Maps (SOMs) also known as Kohonen maps, are a type of unsupervised neural network used for dimensionality reduction, clustering and data visualization. They were invented by Finnish professor Teuvo Kohonen in the 1980s. It trained its network through a competitive learning algorithm. SOM is used for clustering and mapping techniques to map multidimensional data onto lower-dimensional which allows people to reduce complex problems for easy interpretation. SOM has two layers, one is the input layer and the other one is the output layer.

Architecture of SOM

pts. Discuss the various steps to create a self-organized map

The SOM consists of two layers

- (i) Input layer: Each neuron corresponds to a feature in the input vector.
- (ii) Output layer: Usually a 2D grid (eg. 10×10) of neurons (map grid).

Each node (neuron) in the output layer has a weight vector of the same dimension as the input vector.

Consider the network shown in figure which considered for training sample each vector of length 4 and two output units. Train the SOM network by determining the class membership of the input data

Training samples:

$$x_1 = (1, 0, 1, 0) \quad x_2 = (1, 0, 0, 0)$$

$$x_3 = (1, 1, 1, 1) \quad x_4 = (0, 1, 1, 0)$$

bring to
Dainip

New leader

	hair color	eye color	skin color
cluster 1	brunette	amber	fair
cluster 2	red	green	fair
cluster 3	black	hazel	brown

person	hair color	eye color	skin color
P1	blonde	amber	fair
P2	brunette	gray	brown
P3	red	green	brown
P4	black	hazel	brown
P5	brunette	amber	fair
P6	black	gray	brown
P7	red	green	fair
P8	black	hazel	fair

Compare ~~cluster~~ cluster 1 to the observation ^{P1} give
1 dissimilarity.
likewise, calculate all the dissimilarities and put
them in a matrix. Assign each observation to its
closest cluster.

	cluster 1	cluster 2	cluster 3	cluster
P1	1 ✓	2	3	cluster 1
P2	2 ✓	3	2	cluster 1
P3	3	1 ✓	2	cluster 2
P4	3	3	0 ✓	cluster 3
P5	0 ✓	2	3	cluster 1
P6	3	3	1 ✓	cluster 3
P7	2	0 ✓	3	cluster 2
P8	2	2	1 ✓	cluster 3

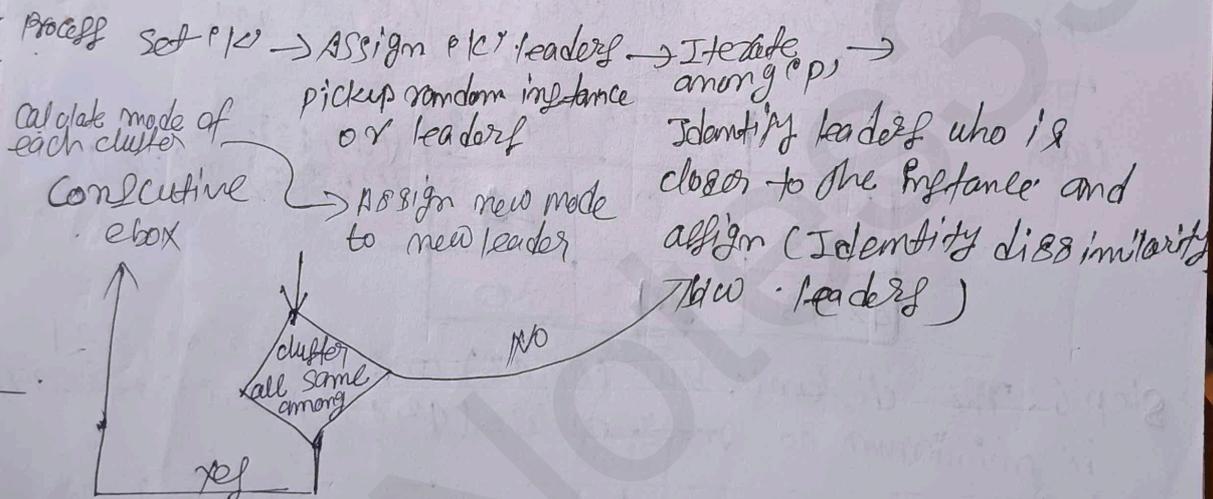
Stop because
 $1.1 = 1.2$

→ diagram

The observations P1, P2, P5 are assigned to cluster 1, P3, P7 are assigned to cluster 2, and P4, P6, P8 are assigned to cluster 3. We stop here as we see here there is no change in the assignment of observations.

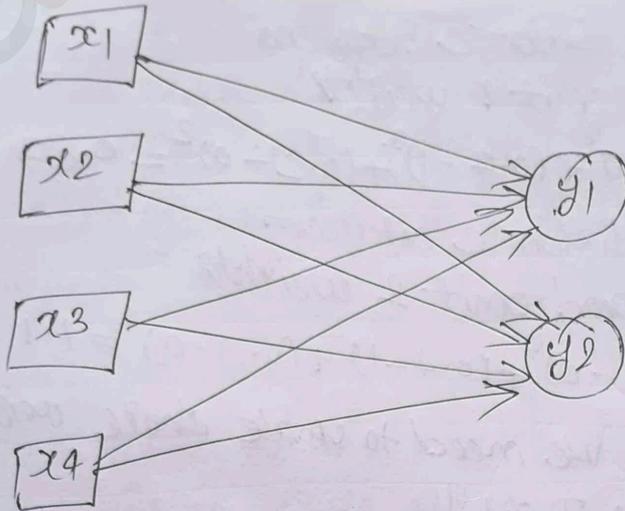
K-Mode clustering

K-Mode is one of the unsupervised algorithms used to cluster categorical variables unlike traditional clustering algorithms that use distance metrics. K-Mode works by identifying the modes or most frequent values within each cluster to determine its centroid. K-Mode is ideal for clustering categorical data such as customer demographics, market segments or survey responses. It is a powerful tool for data analysts and scientists to gain insights into their data and make informed decisions.



or working of K-Mode clustering

- ① pick k observations at random and use them as cluster leaders.
- ② calculate the dissimilarities and assign each observation to its closest cluster.
- ③ Define new modes for the clusters.
- ④ Repeat 2-3 steps until there are no reassignments required.



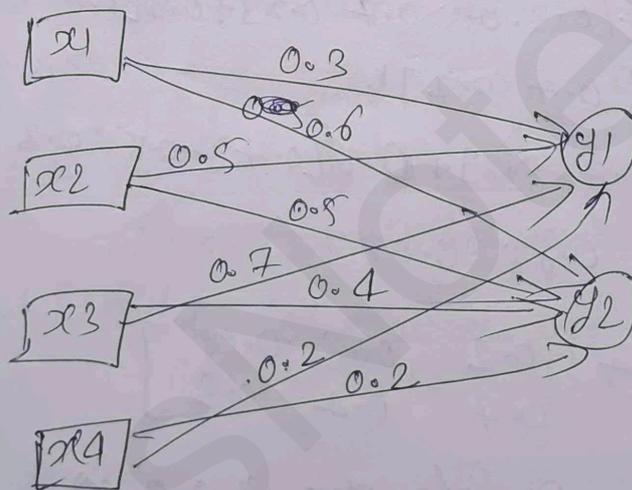
output units: unit 1, unit 2

learning rate $\alpha = 0.6$

initial weight matrix

$$\begin{bmatrix} \text{unit 1} \\ \text{unit 2} \end{bmatrix} = \begin{bmatrix} 0.3 & 0.5 & 0.7 & 0.2 \\ 0.6 & 0.5 & 0.4 & 0.2 \end{bmatrix}$$

ans \rightarrow



Iteration 1: Training sample $x_1 = (1, 0, 1, 0)$

weight matrix:

$$\begin{bmatrix} \text{unit 1} \\ \text{unit 2} \end{bmatrix} = \begin{bmatrix} 0.3 & 0.5 & 0.7 & 0.2 \\ 0.6 & 0.7 & 0.4 & 0.2 \end{bmatrix}$$

Compute Euclidean distance between

$x_1: (1, 0, 1, 0)$ and unit 1 weights

$$d^2 = (0.3 - 1)^2 + (0.5 - 0)^2 + (0.7 - 1)^2 + (0.2 - 0)^2 = 0.87$$

Compute Euclidean distance between

$x_1: (1, 0, 1, 0)$ and unit 2 weights

$$d^2 = (0.6 - 1)^2 + (0.7 - 0)^2 + (0.4 - 1)^2 + (0.3 - 0)^2 = 1.1$$

unit 1 wins! now we need to update those weights
here $[0.3, 0.5, 0.7, 0.2]$

$$w_j(t+1) = w_j(t) + \eta(t) (x_s - w_j(t))$$

New weights old weights learning rate input old weights

update the weights of the ~~win~~ winning unit

$$\text{New unit 1 weight} = [0.3, 0.5, 0.7, 0.2] + 0.6([1, 0, 1, 0] - [0.3, 0.5, 0.7, 0.2])$$

$$= [0.3, 0.5, 0.7, 0.2] + 0.6([0.7, -0.5, 0.3, -0.2])$$

$$= [0.72, 0.2, 0.88, 0.08]$$

$$\begin{bmatrix} \text{unit 1} \\ \text{unit 2} \end{bmatrix} = \begin{bmatrix} 0.72 & 0.2 & 0.88 & 0.08 \\ 0.6 & 0.7 & 0.4 & 0.3 \end{bmatrix}$$

Iteration 2: Training Sample $x_2: (1, 0, 0, 0)$

weights Matrix

$$\begin{bmatrix} \text{unit 1} \\ \text{unit 2} \end{bmatrix} = \begin{bmatrix} 0.72 & 0.2 & 0.88 & 0.08 \\ 0.6 & 0.7 & 0.4 & 0.3 \end{bmatrix}$$

Compute Euclidean distance between $x_2: (1, 0, 0, 0)$
and unit 1 weights $d^2 = (0.72 - 1)^2 + (0.2 - 0)^2 + (0.88 - 0)^2 +$

$x_1: (1, 0, 1, 0) \rightarrow \text{unit 1}$
 $x_2: (1, 0, 0, 0) \rightarrow \text{unit 1}$
 $x_3: (1, 1, 1, 1) \rightarrow \text{unit 2}$
 $x_4: (0, 1, 1, 0) \rightarrow \text{unit 2}$

} Epoch

This process is continued for many epochs until the feature map does not change.

DIANA - Divisive Analysis Clustering

Q. Given the dataset (a, b, c, d) and the distance matrix in table, Apply divisive analysis clustering algorithm (DIANA) to form clusters.

	a	b	c	d	e
a	0	9	3	6	11
b	9	0	7	5	10
c	3	7	0	9	2
d	6	5	9	0	8
e	11	10	2	8	0

Solⁿ:

→ Step 1: initially, $C_1 = \{a, b, c, d, e\}$

Step 2: $C_i = C_1$ and $C_j = \phi$

Step 3: Initial Iteration

Let us calculate the average dissimilarities of the objects in C_i with the other objects in C_i .

Average dissimilarity of a

$$a = \frac{1}{4} * (d(a, b) + d(a, c) + d(a, d) + d(a, e))$$

$$a = \frac{1}{4} (9 + 3 + 6 + 11) = 7.25$$

$$= [0.6 \cdot 0.7 \ 0.4 \ 0.3] + 0.6 [0.4 \ 0.3 \ 0.6 \ 0.7]$$

$$= [0.84 \ 0.88 \ 0.76 \ 0.72]$$

$$\begin{bmatrix} \text{unit 1} \\ \text{unit 2} \end{bmatrix} = \begin{bmatrix} 0.89 & 0.08 & 0.35 & 0.03 \\ 0.84 & 0.88 & 0.76 & 0.72 \end{bmatrix}$$

Iteration 4:

Training sample $x_4: (0, 1, 1, 0)$

weight matrix

$$\begin{bmatrix} \text{unit 1} \\ \text{unit 2} \end{bmatrix} = \begin{bmatrix} 0.89 & 0.08 & 0.35 & 0.03 \\ 0.84 & 0.88 & 0.76 & 0.72 \end{bmatrix}$$

Compute Euclidean distance b/w

$x_4: (0, 1, 1, 0)$ and unit 1 weights

$$d^2 = (0.89 - 0)^2 + (0.08 - 1)^2 + (0.35 - 1)^2 + (0.03 - 0)^2 = 2.06$$

Compute Euclidean distance b/w $x_4: (0, 1, 1, 0)$ and unit 2

weights $d^2 = (0.84 - 0)^2 + (0.88 - 1)^2 + (0.76 - 1)^2 + (0.72 - 0)^2$

$$= 1.3$$

$$w_{ij}(t+1) = w_{ij}(t) + \eta(t)(x_{is} - w_{ij}(t))$$

unit 2 wins

update the weights of the winning unit:

New unit 2 weights:

$$[0.84 \ 0.88 \ 0.76 \ 0.72] + 0.6 [0 \ 1 \ 1 \ 0] - [0.84 \ 0.88 \ 0.76 \ 0.72]$$

$$= [0.34 \ 0.95 \ 0.9 \ 0.29]$$

$$\begin{bmatrix} \text{unit 1} \\ \text{unit 2} \end{bmatrix} = \begin{bmatrix} 0.89 & 0.08 & 0.35 & 0.03 \\ 0.34 & 0.95 & 0.9 & 0.29 \end{bmatrix}$$

Best mapping units for each of the sample taken are

$$+(0.08-0)^2 = 0.74$$

Compute Euclidean distance between

$x_2: (1, 0, 0, 0)$ and unit 2 weights

$$d^2 = (0.6-1)^2 + (0.7-0)^2 + (0.4-0)^2 + (0.3-0)^2 = 0.9$$

unit 1 wins

$$w_j(t+1) = w_j(t) + \eta(t)(x_s - w_j(t))$$

update the weights of the winning unit:

$$\text{New unit 1 weights} = [0.72 \ 0.2 \ 0.88 \ 0.08] +$$

$$0.6[1 \ 0 \ 0 \ 0] - [0.72 \ 0.2 \ 0.88 \ 0.08]$$

$$= [0.89 \ 0.08 \ 0.35 \ 0.03]$$

$$\begin{bmatrix} \text{unit 1} \\ \text{unit 2} \end{bmatrix} : \begin{bmatrix} 0.89 & 0.08 & 0.35 & 0.03 \\ 0.6 & 0.7 & 0.4 & 0.3 \end{bmatrix}$$

Iteration 3:

Training sample $x_3: (1, 1, 1, 1)$

weight matrix:

$$\begin{bmatrix} \text{unit 1} \\ \text{unit 2} \end{bmatrix} : \begin{bmatrix} 0.89 & 0.08 & 0.35 & 0.03 \\ 0.6 & 0.7 & 0.4 & 0.3 \end{bmatrix}$$

Compute Euclidean distance between

$x_3: (1, 1, 1, 1)$ and unit 1 weights

$$d^2 = (0.89-1)^2 + (0.08-1)^2 + (0.35-1)^2 + (0.03-1)^2$$

$$= 2.2$$

Compute Euclidean distance btw $x_3: (1, 1, 1, 1)$ and unit 2

$$\text{weights } d^2 = (0.6-1)^2 + (0.7-1)^2 + (0.4-1)^2 + (0.3-1)^2$$

unit 2 wins. $w_j(t+1) = w_j(t) + \eta(t)(x_s - w_j(t))$

update the weights of the winning unit:

$$\text{New unit 2 weights} = [0.6 \ 0.7 \ 0.4 \ 0.3] + 0.6([1 \ 1 \ 1 \ 1] - [0.6 \ 0.7 \ 0.4 \ 0.3])$$

Similarly we have

Average dissimilarity of $a = 7.25$

" " " " $b = 7.75$

" " " " $c = 8.85$

" " " " $d = 7.00$

" " " " $e = 7.75$

The highest average distance is 7.75 and there are two corresponding objects.

We choose one of them, b , arbitrarily.

We move b to c .

We now have $c_i = \{a, c, d, e\}$ and $c_j = \phi \cup \{b\} = \{b\}$

Step 3: Remaining iteration

(i) 2nd iteration

$c_i = \{a, c, d, e\}$ and $c_j = \{b\}$

$$D_a = \frac{1}{2} (d(a,c) + d(a,d) + d(a,e)) - \frac{1}{1} (d(a,b)) = \frac{20}{2} - 9 = -2.33$$

$$D_c = \frac{1}{3} (d(c,a) + d(c,d) + d(c,e)) - \frac{1}{1} (d(c,b)) = \frac{14}{3} - 7 = -2.33$$

$$D_d = \frac{1}{3} (d(d,a) + d(d,c) + d(d,e)) - \frac{1}{1} (d(d,b)) = \frac{23}{3} - 7 = 0.67$$

$$D_e = \frac{1}{3} (d(e,a) + d(e,c) + d(e,d)) - \frac{1}{1} (d(e,b)) = \frac{21}{3} - 7 = 0$$

D_d is the largest and $D_d > 0$

So we move d to c

We now have $c_i = \{a, c, e\}$ and $c_j = \{b\} \cup \{d\} = \{b, d\}$

(ii) 3rd iteration

$$D_a = \frac{1}{2} (d(a,c) + d(a,e)) - \frac{1}{2} (d(a,b) + d(a,d))$$

$$D_a = \frac{14}{2} - \frac{15}{2} = -0.5$$

$$D_c = \frac{1}{2} [d(c,a) + d(c,e)] - \frac{1}{2} (d(c,b) + d(c,d))$$

$$D_c = \frac{5}{2} - \frac{16}{2} = -13.5$$

of a weighted sum of the inputs, which then passed through an activation function.

The weighted sum =

$$z = w_1 x_1 + w_2 x_2 + \dots + w_n x_n = x^T W$$

The step function compares this weighted sum to a threshold. If the input is larger than the threshold value, the output is 1; otherwise, it is 0.

$$h(z) = \begin{cases} 0 & \text{if } z < \text{Threshold} \\ 1 & \text{if } z \geq \text{Threshold} \end{cases}$$

A perceptron consists of a single layer of threshold logic units (TLUs) with each (TLU) fully connected to all input nodes. The output of the fully connected layer is

computed as:

$$f_w, b(x) = h(x^T W + b)$$

where x is the input w is the weight for each input neurons and b is the bias and h is the step function.

The weight update formula is

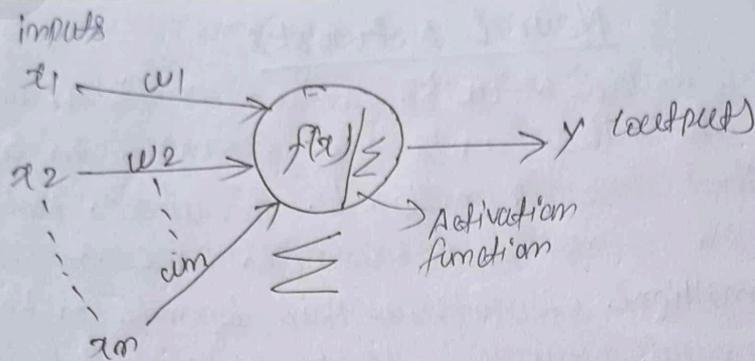
$$w_{i,j} = w_{i,j} + \eta (y_j - \hat{y}_j) x_i$$

where $w_{i,j}$ is the weight b/w i^{th} imp and j^{th} o/p neuron.

x_i is the i^{th} imp value.

y_j is the actual value, and \hat{y}_j is the predicted value.

η is the learning rate



Q10 what is perceptron? How does a perceptron work in Neural networks?

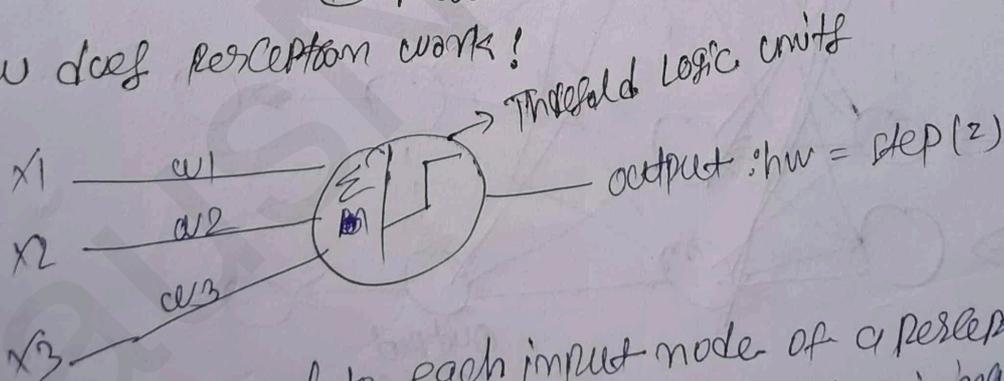
ans → perceptron is a type of neural network that performs binary classification that maps input features to an output decision, usually classifying data into one of two categories such as 0 or 1.

perceptron consists of a single layer of input nodes that are fully connected to a layer of output nodes. It is particularly good at learning linearly separable patterns. It utilizes a variation of artificial neurons called

Threshold Logic Unit (TLU).

- Types of perceptron
- (i) single layer perceptron
 - (ii) multi-layer perceptron.

How does perceptron work?



A weight is assigned to each input node of a perceptron indicating the importance of that input in determining the output. The ~~output~~ perceptron output is ~~also~~ calculated

$$D_e = \frac{1}{2}(d(e, a) + d(e, c)) - \frac{1}{2}(d(e, b) + d(e, d))$$

$$D_e = \frac{13}{2} - \frac{18}{2} = -2.5$$

None of them are +ve here

All are negative, so we stop and form the clusters c_i and c_j .

$$c_i = \{a, c, e\} \text{ and } c_j = \{b, d\}$$

Step 4:

To divide c_i and c_j , we compute their diameters.

$$\text{diameter}(c_i) = \max\{d(a, c), d(a, e), d(c, e)\}$$

$$\text{diameter}(c_i) = \max\{3, 11, 2\} = 11$$

$$\text{diameter}(c_j) = \max\{d(b, d)\} = 5$$

The cluster with the large diameter is c_i .

So we now split c_i .

We repeat the process by taking $c_i = \{a, c, e\}$

$$\therefore c_i = \{a, c, e\}, c_j = \emptyset$$

Iteration-1

$$D_a = \frac{1}{2}(d(a, c) + d(a, e)) = \frac{1}{2}(14) = \underline{\underline{7}}$$

$$D_c = \frac{1}{2}(d(c, a) + d(c, e)) = \frac{1}{2} \times 5 = 2.5$$

$$D_e = \frac{1}{2}(d(e, a) + d(e, c)) = \frac{1}{2} \times 13 = 6.5$$

D_a is the maximum so we move a to c_j .

$$\text{Now } c_i = \{c, e\}, c_j = \{a\}$$

Iteration-2

$$D_c = \frac{1}{2}(d(c, e)) - \frac{1}{2}(d(c, a)) = 2 - 3 = -1$$

$$D_e = \frac{1}{2}(d(e, c)) - \frac{1}{2}(d(e, a)) = 2 - 11 = -9$$

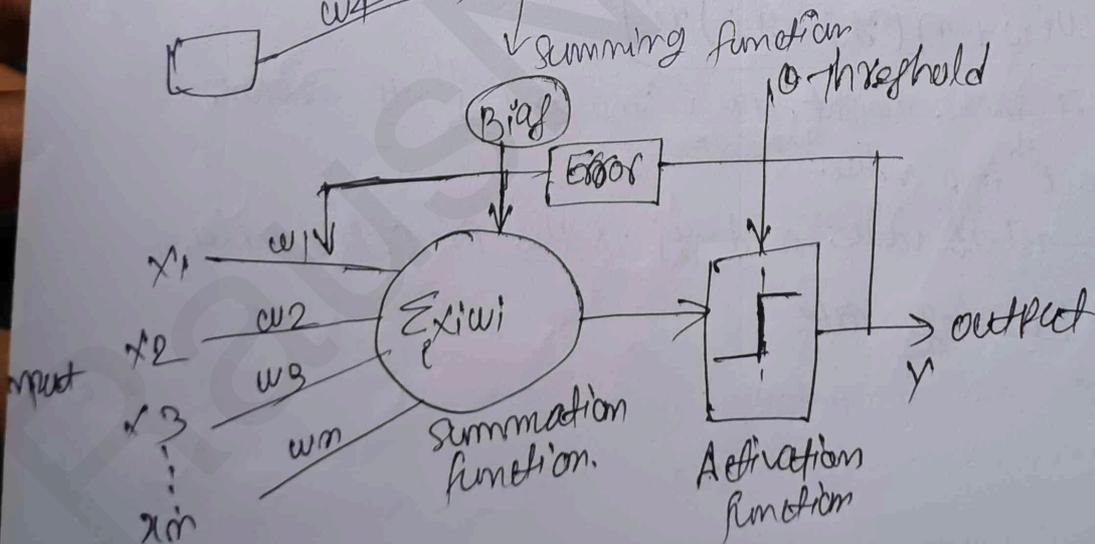
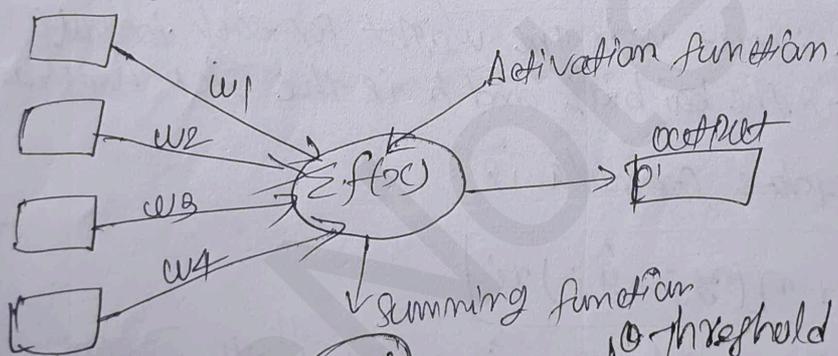
Since none of them are +ve here

All are negative, so we stop and form the clusters c_i and c_j .

* Single Layer Perceptron: Single layer perceptron is inspired by biological neurons and their ability to process information. To understand the SLP we first need to break down the workings of a single artificial neuron which is the fundamental building block of neural networks. An artificial neuron is a simplified computational model that mimics the behavior of a biological neuron. It takes inputs, processes them and produces an output.

- It receives signal from outside.
- Processes the signal and decides whether we need to send information or not.
- Communicate the signal to the target cell, which can be another neuron.

Input layer



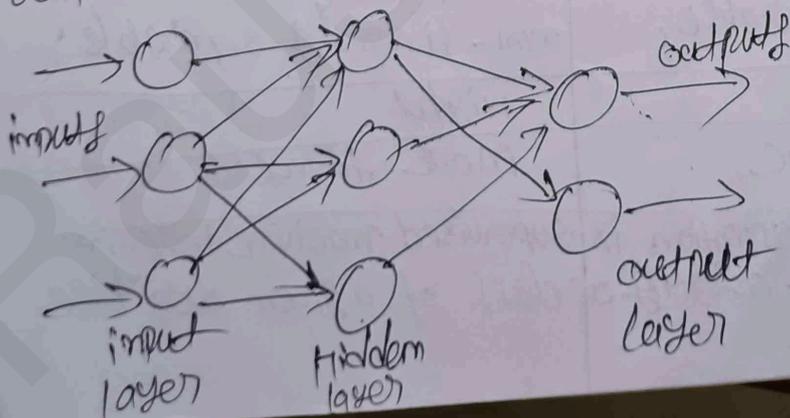
Multi-layer perceptron: Multi-layer perceptron is an artificial neural network widely used for solving classification and regression tasks. MLP transform input data from one dimension to another. It is called "multi-layer" because it contains an input layer, one or more hidden layers, and an output layer. The purpose of an MLP is to model complex relationships b/w inputs and outputs, making it a powerful tool for various machine learning tasks.

Key Components of Multi-layer perceptron

(i) Input layer: Each neuron in this layer corresponds to an input feature. ~~the input~~ for instance, if you have three input features, the input layer will have three neurons.

(ii) Hidden layers: An MLP can have any no. of hidden layers, with each layer containing any number of nodes. These layers process the information ~~received~~ received from the input layer.

(iii) Output layer: The output ~~layer~~ layer generates the final prediction or result. If there are multiple outputs, the output layer will have a corresponding no. of neurons.



* Difference b/w Multi-class and Multi-label classification

features	Multi-class classification	Multi-label classification
Definition	In multi-class classification, each instance is classified into one class out of more than two possible classes.	In multi-label classification, each instance can be assigned multiple classes.
Example	classifying an image as cat, dog	classifying a movie into action, comedy, thriller.
Exclusion	The classes are mutually exclusive.	The classes are not mutually exclusive
Label representation	one-hot-encoded vectors	multi-hot encoded binary vectors
Loss Function	Categorical cross-Entropy	Binary cross-Entropy

* Diff b/w features	Linear and Non-linear classification	Linear classification	Non-linear classification
Definition	A classification that uses a straight line to separate data points into classes.	A classification that uses a curved plane	A classification that uses a curve to separate data points into classes
Decision Boundary	Straight line	Decision Decision Trees	Curved plane
Example	Logistic Regression		
Work well	works well when data is linearly separable.		works well when data is non-linearly separable.
Complexity	Low		High
Training Time	Less Time		More time

Define: classification is supervised machine learning algorithm used to predict category or class of a given data base on past labeled data.