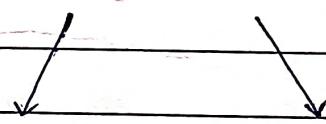


- There are 3 types of Learning they are :-
- Supervised Learning
- Unsupervised Learning
- Reinforcement Learning

1. SUPERVISED LEARNING

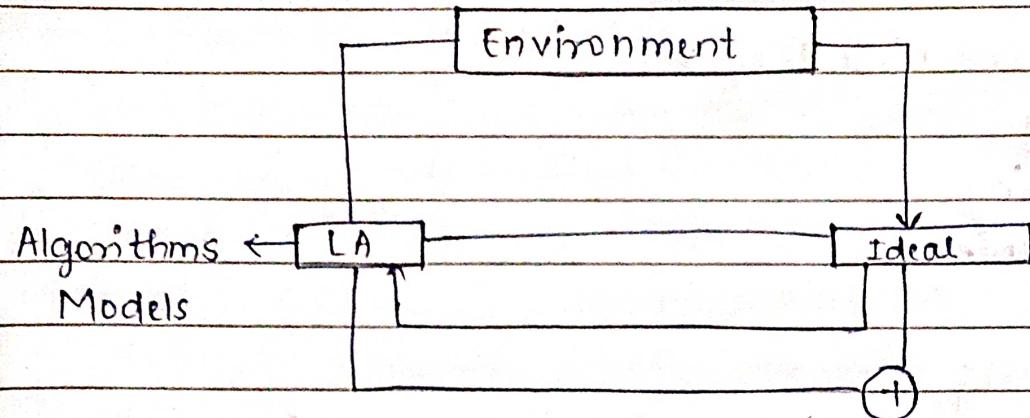
- The Learning Agent (LA) learns from the environment as well as his/her previous output. The supervisor and LA both links to the environment.
 - The supervisor is having the ideal solution.
 - The LA receives the feedback from the supervisor as well as environment.
 - \oplus - it is called as an adder - it provides the difference between Actual and predicted.
- In this way the machine trains / learns and predict the outputs.
- There are different terminology available
 - feature \rightarrow attributes
 - predictor or variable or independent variable.
 - Target variable - Dependent variable / Response variable.

Types of SUPERVISED LEARNING



Regression and Classification

- Linear Regression - Support Vector Machine
- MLR
- Lasso Regression
- Ridge Regression
- polynomial Regression
- KNN - K Nearest Neighbors
- LDA - Linear Discriminant Analysis



$$\text{Adder} = \text{Actual} - \text{Predicted}$$

2. UNSUPERVISED LEARNING

- In unsupervised Learning the model has to learn the unlabelled data.
- The IA learns the input from the environment; it has the actual output where he tries to learn from own.
- IA learns by himself not from any supervisor everytime the IA tries to improve the accuracy by comparing his previous responses.
- It learns its own pattern again and again and it generates the O/P without the corresponding input output (labelled data).

Examples of unsupervised learning are clustering, Association rules etc.

- Why we use unsupervised Learning ?

The model extract the hidden features from the dataset. It provides the insightful information about a data set.
For example - Hidden pattern

Algorithms — K-mean, K-medoid, DBSCAN, HC → clustering
 Apriori, FP-Growth, CLAT,
 Assocation

labelled data - Model is trained
 x - study time
 y - Grade.

x	y
5	D
6	C
7	B

Page No.:

Date:

Dataset

Association

T_1 : Bread, cheese, butter.

$Bread \rightarrow cheese$

T_2 : chicken pakoda, Beer

$(Bread, Butter) \rightarrow$

T_3 : Book, Pencil

$cheese \rightarrow$

MBA - Market Basket Analysis.

That hidden patterns provided to the model. The model learn and predict it.

Models of unsupervised learning

let us consider a model which contains huge number of models of and our task is to categorise the different objects.

1. Raw data

2. Interpretation - (finding hidden data from the data set)

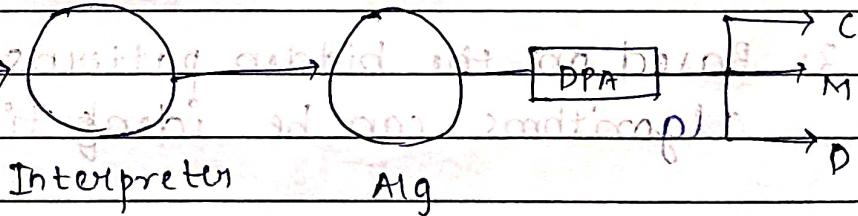
3. Algorithm

Analysis

4. DPA (Data pre-processing Algorithm)

5. Result

cat	go dog
mammal	dog
cat	cat
	dog



Features

$C \rightarrow 111, 222, 333, 444, \dots$

$M \rightarrow 111, 222, 333, 444, \dots$

$H \rightarrow 111, 222, 333, 444, \dots$

Outline of DPA - (for cleaning the data)

- (i) Removing Noise
- (ii) Handling Missing (unknown value, unassigned value)
- (iii) Identify Outlier,

- Initial we will collect a data from different sources and prepare the raw data. The raw data passed into the model.
- In this step, the model generates the features which are available in the raw dataset. It finds the hidden patterns and establish a strong association rules.

Ex -

Features	$C \rightarrow 1, 2, 3, 4$
	$M \rightarrow 11, 22, 33, 44$
Transactional	$H \rightarrow 111, 222, 333, 444$

$R_1 : \text{Bread} \rightarrow \text{cheese}$

$R_2 : \text{Pen} \rightarrow \text{Book}$

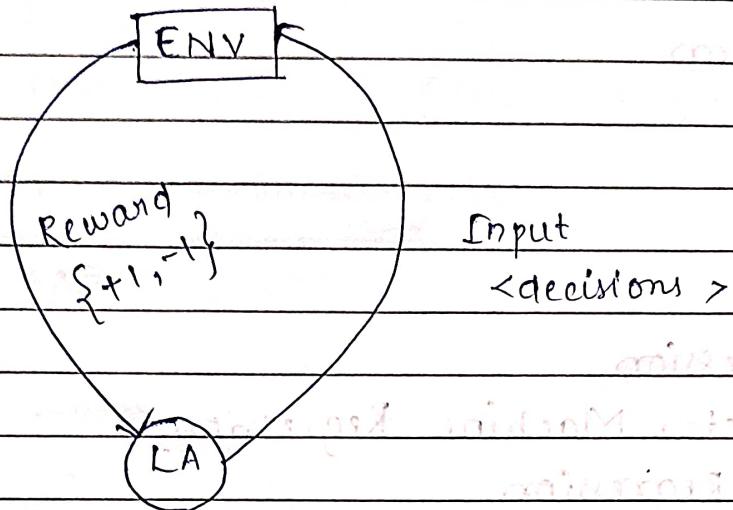
The R_1 and R_2 are called as strong association rules that help us to identify the hidden patterns.

- Based on the hidden patterns, the corresponding algorithms can be identified.
- DPA helps us to remove the Noise and identify the outliers from the dataset. It will clean the dataset. This process is called as data cleaning.

- Finally, the algorithm divides the data objects into groups according to the similarities and find the difference b/w them.

3. REINFORCEMENT LEARNING

(Reward-based Learning)



In Reinforcement learning the learning agent learns from the environment and provide a series of actions each action's are called as the decisions statements. It will be now provide the rewards each reward is further classified in two types +positive and -negative. The +ve reward is treated as +1 and the -ve reward is treated as -1. The model learns all the +ve instances and build a model. It enforces to learn hence it is called as Reinforcement learning or Reward based learning.

MBA - Market basket analysis is a strategic data mining technique used by retailers to enhance sales by gaining a deeper understanding of customer purchasing patterns.

SUPERVISED LEARNING

It provides the different algorithms to predict the future value and these are

- Linear Regression
- Multiple
- polynomial
- Ridge
- Lasso
- Logistic Regression
- Support Vector Machine Regression
- Decision Tree Regression

In all cases the machine trains with the training data and validate with the testing data.

1. LINEAR REGRESSION

The term linear Regression means how one variable is linearly dependent into another variable. The

Objective of LR is to make predictions like cells prediction, Salary prediction, Age prediction, price of production prediction.

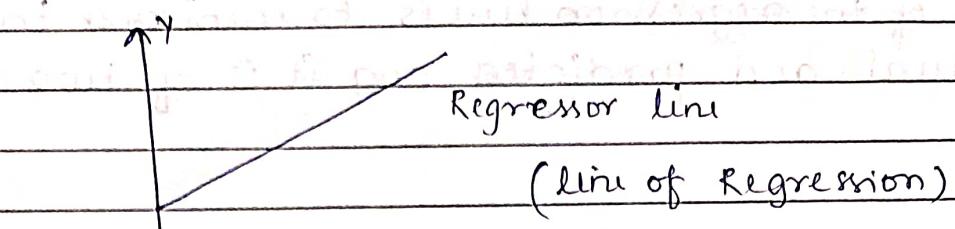
LR algorithm shows a linear relationship b/w a dependent variable(y) and one or more independent variable and Hence it is called as LRs it means find how the value variable according to the value of independent variable

Let us take an example —

Number of hours

x	$y \longrightarrow$ grades
5	D
6	C
7	B
8	A

Let us plot the graph and find the regressor line



The Regressor line

$$y = mx + c$$

OR

$$y = b_0 + b_1 x_1$$

The objective of the LR is to find the best fit line from equation where x is independent $y =$

$b=0$ intercept of the line

b_1 coefficients of the line

Types of linear regression

(i) Linear Regression → One independent variable and one dependent variable

If a single independent variable is nearly equal to the value of a numerical dependent variable then it is called simple linear regression.

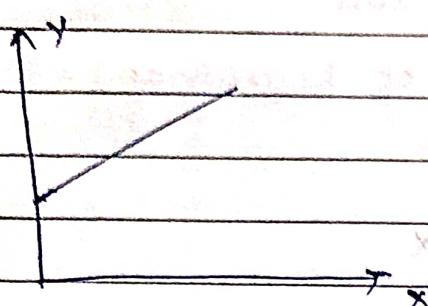
ii) Multiple Linear Regression

The equation is $y = b_0 + b_1 x_1 + b_2 x_2 + \dots + b_n x_n$ when more than 1 independent variable is used to predict the value of y then it is called MLR.

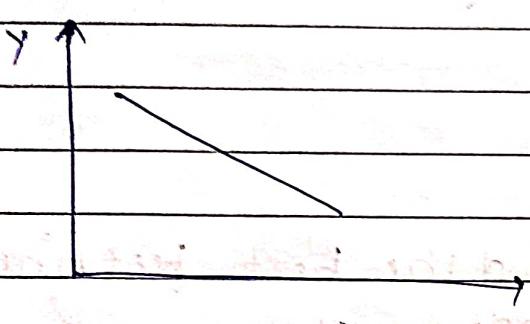
x_1	x_2	y

Linear Regression Lines

The Objective of the regression line is to minimize the error between actual and predicted and it is of two types



Positive Linear Regression
+ve linear relation.



-ve linear relation

COST FUNCTION

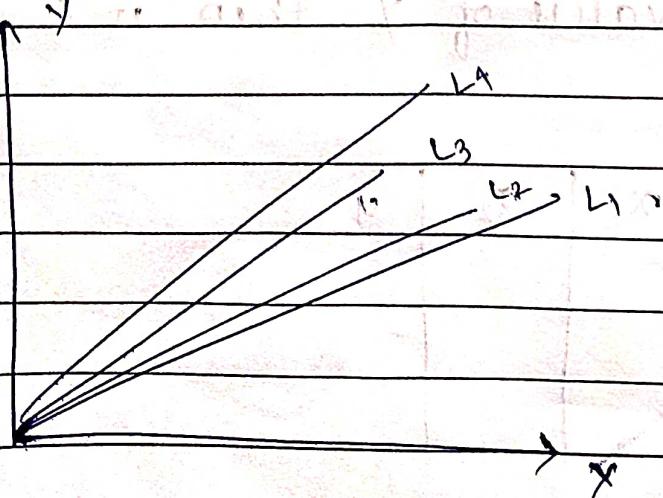
To find the best fit line we required to find the difference b/w the predicted value and the actual value the error is called as loss.

Error can be defined as Actual - Predicted.

$$\text{Loss} [E = A - P]$$

The different values for weight or coefficient of lines b_0 and b_1 gives the diff and the cost function is used to estimate the value of the coefficient for the best fit line.

In this case there are 4 lines that can be obtained and can be represented in 2 dimensional coordinate.



Out of these 4 line which one is the best line?

Those lines which obtain the less error that is the best line that means it is the less loss between actual and predicted.

I/N	Act	Pred	D
L ₁ : 5	E 90	85	5%
L ₂ : 6	E 90	86	4%
L ₃ : 7	E 90	88	2%
L ₄ : 8	E 90	89	1%

Loss

$$\text{Avg} = 3\%$$

$$\text{Mean Square Error} = \frac{1}{N} \sum_{i=1}^N (y_i - (b_0 + b_1 x_i))^2$$

LINEAR REGRESSION

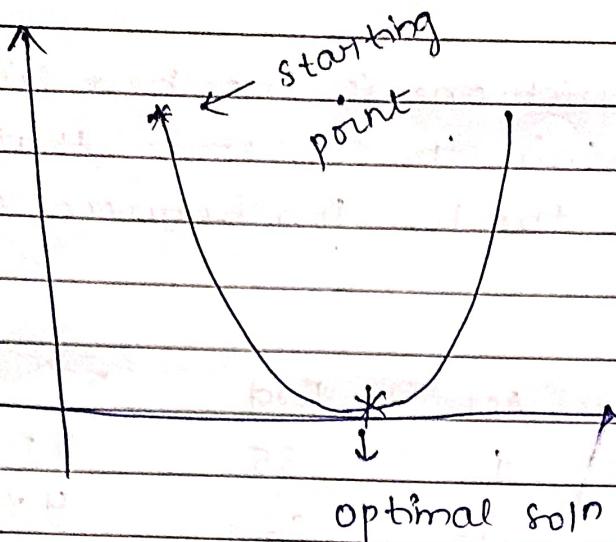
Through Linear Regression we identify the regression line by using the least square Method.

Let us consider the 4 data points if we plot then the 4 lines can be obtained by using the linear regression equation

but our objective is to find the best regression line or the optimal line

To achieve the optimization ; gradiention algorithm provides the ideal soln to minimize MSE (Mean square error). Gradientation is used to minimize the MSE by calculating the gradient of the cost function it will update the co-efficient of the line by reducing the cost function.

A gradientation can be represented as :

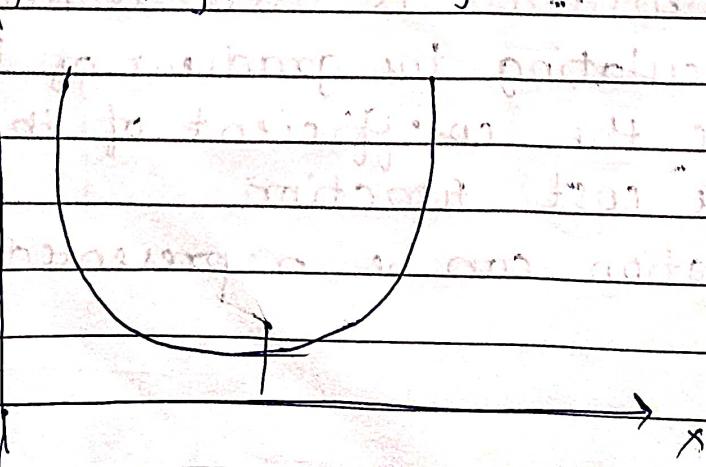
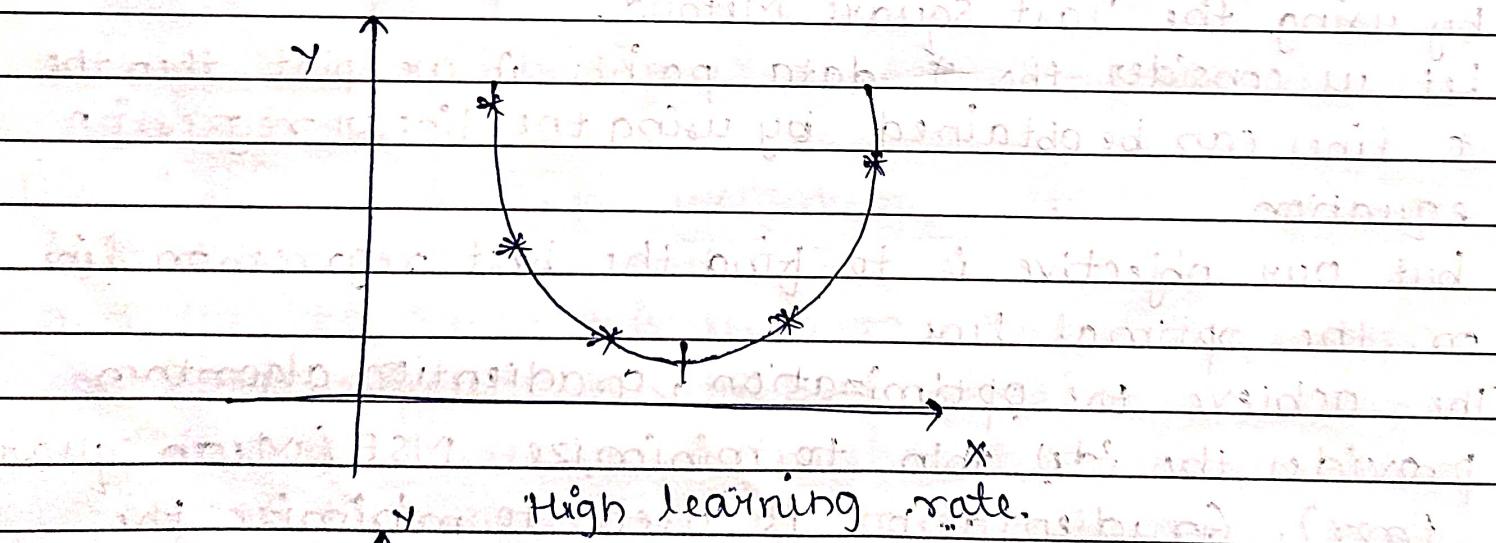


then there are 2 different scenario can be possible

- (i) If you choose to take longer steps then u may get sooner but there is a chance that you can overshoot the bottom of the pit.
- (ii) If you take lesser steps then u will reach ~~opt~~ the optimal soln. But it will take much time each step can be ideally represented as

$$\alpha = 0.0001 \text{ which is called a}$$

learning rate.



Each time the value of x is updated and checked with the optimal solution and this way gradient descent works.

Let us consider a dataset which consists of 2 features x and y and find the followings

- (i) find the b_0
- (ii) find b_1
- (iii) Generate the Regressor line.
- (iv) find the Error term
- (v) Find R^2 .

x	y
1	2
2	4
3	5
4	4
5	5

Apply the linear regression algorithm for the following

The Formula

$$y = b_0 + b_1 x \rightarrow ①$$

calculate the b_1 from the formula.

$$b_1 = \frac{(x - \bar{x})(y - \bar{y})}{(x - \bar{x})^2}$$

Given $\bar{x} = 3$
 $\bar{y} = 4$.

calculate b_0 from the formula.

x	y	$x - \bar{x}$	$y - \bar{y}$	$(x - \bar{x})^2$	$(y - \bar{y})^2$	$(x - \bar{x})(y - \bar{y})$	\hat{y}
1	2	-2	-2	4	4	4	2.8
2	4	-1	0	1	0	0	3.4
3	5	0	1	0	1	0	4
4	4	1	0	1	0	0	4.6
5	5	2	1	4	1	2	5.2

$$b_1 = \frac{\sum (y - \bar{y})(\hat{y} - \bar{y})}{\sum (x - \bar{x})^2}$$

$$b_1 = \frac{0.8 - 0.6}{0.8 + 0.36} = \frac{0.16}{1.16} = 0.14$$

$$\boxed{b_1 = 0.6}$$

Calculate b_0 from equation ①

$$y = b_0 + b_1 x$$

$$\hat{y} = 0.6 + b_1 x \quad y = b_0 + 0.6(3)$$

$$y = b_0 + 1.8$$

$$4 - 1.8 = b_0$$

$$2.2 = b_0$$

$$y = b_0 + b_1 x$$

$$y = 2.2 + 0.6x$$

Regressor line

Best fitted line

Optimal line for the prediction.

Calculate \hat{y} by substituting the values of x

$$\hat{y} = 2.2 + 0.6(1)$$

$$\hat{y} = 2.8$$

Subtract y from \hat{y} .

$$MSE = \frac{\sum (\hat{y} - y)^2}{N-2}$$

$$MSE = 0.8944$$

MSE - Mean square error and it can obtain by using thi formula

$$MSE = \sqrt{\frac{\sum (\hat{y} - y)^2}{N-2}}$$

This MSE is also called as standard error of estimate (SEE)

The total loss of the model is 0.89 which is less than 1 i.e. the model behaves well for training data as well as unseen data.

R^2 - coefficient of determination.

If $x = 10$ what would be y

$$y = 2.2 + 0.6 \times 10$$

$$y = 2.2 + 6.0$$

$$y = 62.2$$

* Calculate total sum of squares
(SST).

The formula for SST is $\sum (y_i - \bar{y})^2$

$$\sum (y_i - \bar{y})^2$$

* Calculate the coefficient of determination

$$R^2 = 1 - \frac{\sum (y_i - \hat{y})^2}{\sum (y_i - \bar{y})^2}$$

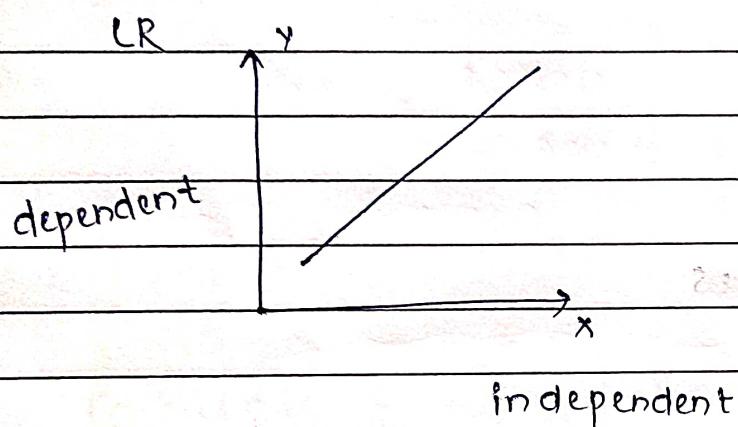
$$R^2 = 0.6$$

$$R = \{-1, 0, 1\}$$

* Multiple Linear Regression

The term Multiple Linear Regression finds the regression line b/w the dependent variable and independent variable.

It is a statistical method used to model the relationship b/w multiple independent variable and a single dependent variable.



MLR = Y

OLS = S.E.

$y | x_1 \ x_2 \ \dots$

$$y = b_0 + b_1 x_1 + \dots + b_n x_n + \epsilon \quad \text{OR} \quad y = f(x)$$

$$y = b_0 + b_1 x_1 + \dots + b_n x_n + \epsilon$$

$$y = f(x)$$

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_n x_n + \epsilon$$

$$\beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_n x_n + \epsilon$$

$$y = b_0 + b_1 x_1 + b_2 x_2 \quad \text{OR}$$

$$y = b_0 + b_1 x_1 + b_2 x_2 + \epsilon$$

In Multiple linear regression we need to find b_0, b_1 , as well as b_2 from the equation $y = b_0 + b_1 x_1 + b_2 x_2$

The formula for $b_0 = \bar{y} - b_1 \bar{x}_1 - b_2 \bar{x}_2$

$$b_1 = \frac{(\sum x_2^2)(\sum x_1 y) - (\sum x_1 x_2)(\sum x_2 y)}{(\sum x_1^2)(\sum x_2^2) - (\sum x_1 x_2)^2}$$

$$b_2 = \frac{(\sum x_1^2)(\sum x_2 y) - (\sum x_1 x_2)(\sum x_1 y)}{(\sum x_1^2)(\sum x_2^2) - (\sum x_1 x_2)^2}$$

$$= \frac{(1000)(1300) - (20)(10)}{(1000)(1300) - (20)(10)} = \frac{1300000 - 200}{1300000 - 400} = \frac{1299800}{1299600} = 1.00008261402$$

Q. Find the regressor line of the following tabular data.

QD (y)	Price (x_1)	(x_2) Income	$y[\bar{y}]$	$x_1[x_1 - \bar{x}_1]$	$x_2[x_2 - \bar{x}_2]$
100	5	1000	1000	20	200
75	7	600	-5	1	-200
80	6	1200	0	0	400
70	6	500	-10	0	-300
50	8	300	-30	2	-500
65	7	400	-15	1	-400
90	5	1300	10	10	500
100	4	1100	20	-2	300
110	3	1300	30	-3	500
60	9	300	-20	3	-500

$$\bar{y} = \frac{\sum x_1 d + \sum x_2 d + \bar{x}}{n} = \frac{\sum x_1 d}{n} + \bar{x} = \bar{x} + ad$$

y^2	x_1^2	$y x_1$	$y x_2$	$x_1 x_2$	$x_1^2 x_2$	$\bar{y} = 800$
400	1	-20	4000	-200	4000	100
25	1	-5	1000	-200	4000	$\bar{x}_1 = \frac{60}{10} = 6$
0	0	0	0	0	160000	100
100	0	0	3000	0	90000	$\bar{x}_2 = \frac{8000}{10} = 800$
900	4	-60	15000	-1000	250000	100
225	1	-15	6000	-400	160000	100
100	1	-10	5000	-500	250000	,
400	4	-40	6000	-600	90000	100
900	9	-90	15000	-1500	250000	,
100	9	-60	10000	-1500	250000	,
		-300	65000	-5900	1580000	,

$$b_1 = \frac{(\sum x_2)^2 (\sum x_1 y) - (\sum x_1 x_2)(\sum x_2 y)}{(\sum x_1^2)(\sum x_2^2) - (\sum x_1 x_2)^2} = \frac{(1580000)(-300)}{(30)(1580000)} = -(-5900)(6500)$$

$$= \frac{(-5900)(6500)}{(30)(1580000)}$$

$$b_1 = -7.1882 - (-5900)^2$$

$$b_2 = (30)(65000) - (-5900)(650 - 300)$$

$$(30)(1580000) - (-5900)^2$$

$$1950000 - 1770000$$

$$= \frac{47400000 - 34810000}{12590000}$$

$$= \underline{180000}$$

$$12590000$$

$$b_2 = 0.01429$$

$$y = b_0 + b_1 x_1 + b_2 x_2$$

$$b_0 = \bar{y} - b_1 \bar{x}_1 - b_2 \bar{x}_2$$

$$= 80 - (-7.188)(6) - (0.014)(800)$$

$$= b_0 = 111.928$$

$$y = b_0 + b_1 x_1 + b_2 x_2$$

$$Y = 111.928 + (-7.188)(6) + (0.014)(800)$$

$$y = 80$$

The Regression coefficient for x on y is written as b_{xy} .

$$y \text{ on } x = b_{yx}$$

The coefficient of correlation can be represented as r and it can be written to as

$$r = \sqrt{b_{xy} * b_{yx}}$$

where b_{xy} = coefficient of x on y .

b_{yx} = coefficient of y on x .

The value of r lies between $\{-1 \text{ to } +1\}$

if $r = +1$ then it is highly correlated

- $r = -1$ there is no correlation among x and y .

$r = 0$ there is partial correlation among x and y

Let us consider two regression line $x + 6y = 6$, and $3x + 2y = 10$. (1)

find the correlation coefficient (r). (2)

treat eqn ① as y on x .

$$x + 6y = 6.$$

$$6y = 6 - x$$

$$y = \frac{6-x}{6}$$

$$\text{in the } y = \frac{1-x}{6}$$

$$b_{xy} = -\frac{1}{3}.$$

the coefficient of correlation can be determined as

$$r = \sqrt{b_{xy} * b_{yx}}$$

$$r = \sqrt{\frac{-2}{3} * \frac{-1}{6}}$$

$$r = 0.33$$

$$r \text{ is}$$

From the two given equations there is a partial correlation among two variables.

of two lines

Q. The two equations of the two regressions or the obtained in a correlation analysis are the following

$$2x = 8 - 3y \rightarrow ①$$

$$2y = 5 - x \rightarrow ②$$

Obtain the value of the correlation coefficient

$$2x = 8 - 3y$$

$\frac{1}{2}$

$$2y = 5 - x$$

~~2y~~

$$x = 4 - \frac{3}{2}y$$

$$y = \frac{5-x}{2}$$

$$x = 4 - \frac{3}{2} \left(\frac{5-x}{2} \right)$$

$$b_{xy} = -\frac{3}{2}$$

$$x = 4 - \frac{15+3x}{4}$$

$$b_{yx} = -\frac{1}{2}$$

$$4x = 16 - 15 + 3x$$

$$x = 1$$

$$r = \sqrt{\frac{3}{2} \times \frac{1}{2}}$$

$$r = 0.86$$

* A student obtained the following two regression equations
Do you agree with him.

$$6x = 15y + 21$$

$$21x + 14y = 56.$$

$$x = \frac{15}{6}y + \frac{21}{6}$$

$$\frac{21x}{56} + \frac{14y}{56} = 0.$$

$$x = \frac{15}{6}y + \frac{7}{2}$$

$$\frac{14y}{56} = \frac{21x}{56}$$

$$b_{xy} = \frac{15}{6}$$

$$\frac{1}{4}y = \frac{21x}{56}$$

$$r = \sqrt{\frac{15}{6} \times \frac{3}{2}} = 1.93$$

$$y = \frac{84x}{56}$$

$$y = 3\frac{1}{2}x$$

- If the obtained line is correct
- any regression coefficient is -ve and other are +ve
- if both the regression coefficients are > 1
- It is not possible.

* The height of the fathers and sons are given below,
find the height of the son when the height of the
father is 70 inches.

Father	Son		
71	69	65	59
68	64	66	62
66	65		
67	63	$\bar{x} = 759/11 = 69$	
70	65	$\bar{y} = 704/11 = 64$	
71	62		
70	65		
73	64		
72	66		

The above tabular data indicates to find the regression equation and we have to fit the regression equation for answering the derived value.

The linear regression x on y can be written as

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

In the above equation the deviations are taken from actual mean

$$r = \frac{\sum xy}{\sqrt{\sum x^2}} \text{ can be written as } \frac{\sum xy}{\sqrt{\sum x^2}}$$

deviation of actual mean.

$$\frac{\sum xy}{\sum x^2} \quad \bar{x} = 69 \quad \bar{y} = 64$$

The calculation required

\bar{x}	$(x - \bar{x})^2$	x^2	\bar{y}	$(y - \bar{y})^2$	y^2	xy
71	2	49	69	48	450	25
68	-3	9	64	0	400	0
66	-2	4	65	1	4225	-10
67	-1	1	63	-1	3969	-6
70	1	1	65	1	4225	1
71	2	4	62	-2	3844	-4
70	1	1	65	1	4225	1
73	4	16	64	0	4096	0
72	3	9	66	2	4356	4
65	-4	16	59	-5	3481	25
66	-3	9	62	-2	3844	4
					66	39

$$Y - \bar{Y} = \frac{\sum xy}{\sum x^2} (x - \bar{x})$$

$$Y - 64 = \frac{39}{74} (x - 69)$$

$$Y - 64 = 0.52(x - 69)$$

$$Y = 0.52x - 35.88 + 64$$

$$Y = 0.52x + 28.12$$

$$Y = 0.52x + 28.12 = 64.4$$

Hence the likely height son when the height of the father is 70 inches shall be 64.53 inches.

Q. Given the following data $\bar{x} = 36$ $\bar{y} = 85$ $\sigma_x = 11$ $\sigma_y = 8$.
 $r = 0.06$.

find the two regression equations

$$Y - 85 = 0.06 \left(\frac{11}{8}\right) (x - 36)$$

$$Y - 85 = \frac{0.66}{8} (x - 36)$$

$$Y - 85 = 0.08x (x - 36)$$

$$Y - 85 = 0.08x - 2.88.$$

$$\underline{Y - 82.12 = x.}$$

$$0.08$$

$$Y = 0.08x - 2.88 + 85$$

$$\times Y - 1026.5 = x.$$

$$\checkmark Y = 0.08x + 82.12$$

$$0.08$$

$$Y - 85 + 2.88 = 0.08x.$$

$$Y - 82.12 = 0.08x$$

$$b_{\text{on } y} = \frac{x - \bar{x}}{\sigma_x} = r \frac{\sigma_y}{\sigma_x} (y - \bar{y})$$

$$x - 36 = 0.06 \left(\frac{8}{11} \right) (y - 85)$$

$$x - 36 = 0.04(y - 85) \quad \checkmark x = 0.04y + 32.6$$

$$\bar{x} = 0.04y - 3.4 + 36$$

Q. Find the two regression equation for the following two series what is most likely value of x when $y = 20$ and most likely value of y when $x = 20$

x	y	$(x - \bar{x})$	x^2	$(y - \bar{y})$	y^2	xy
35	23	5.	25	-2	4	-10
25	27	-8	64	2	4	-10
29	26	-1	1	1	1	-1
31	21	1	1	-4	16	-4
27	24	-3	9	-1	1	3
24	20	-6	36	-5	25	30
33	29	3	9	(-4)(16)	81	-12
36	30	6	36	5	25	30
240	200	0	142	(-8)(82)	92	26
8	8					

$$\bar{x} = 30 \quad \bar{y} = 25$$

y on x .

$$y - \bar{y} = r \frac{\sigma_y}{\sigma_x} (x - \bar{x})$$

$$= 0.04(y - 85)$$

x on y .

$$(x - \bar{x}) = r \frac{\sigma_x}{\sigma_y} (y - \bar{y})$$

As there is no information about coefficient of correlation then we assume deviation takes from assumed mean Hence

$$r \frac{\sigma_x}{\sigma_y} = \frac{\text{sum of } \Sigma xy}{\Sigma x^2}$$

$$x - 30 = 0.35 (y - 25)$$

$$x = 0.35 y - 8.75 + 30$$

$$\hat{x} = 0.35 y + 21.25 = 28.25$$

Substitute

Similarly the regression equation y on x

$$y - 25 = 0.54 (x - 30)$$

$$y = 0.54 x - 16.2 + 25$$

$$y = 0.54 x + 8.8$$

$$y = 19.6$$

Q. You are given below the following information about advertisement expenditure and shells. The correlation coefficient 0.8.

Calculate the two

find the likely shell when advertisement expenditure 25.

What would be the advertisement budget if the company wants attained shells target of 150 cr.

$$x \text{ on } y = bxy \\ y \text{ on } x = byx$$

M	T	W	T	F	S	S
Page No.:						
Date:	20	21	22	23	24	25

The regression equation x on y

Adv. Exp. Sales y

Mean 20 120

SD 5 25

$$\bar{x} = 20 \quad \bar{y} = 120 \quad \sigma_x = 5 \quad \sigma_y = 25$$

$$x - \bar{x} = 0.8(y - \bar{y})$$

$$x - 20 = 0.8(y - 120)$$

$$x = 0.8y - 96 + 20$$

$$x = 0.8y + 16$$

$$x - 20 = 0.8 \left(\frac{\sum}{5} (y - 120) \right)$$

$$x - 20 = 0.16y - 19.2$$

$$x = 0.16y - 19.2 + 20 = 0.8y + 0.8$$

$$x = 0.16y - 0.8 \cdot 0.8 = 0.16(25) + 0.8 = 4.8$$

~~$$y - \bar{y} = 0.16(x - 20) \quad y - \bar{y} = 0.8 \left(\frac{\sum}{5} (x - \bar{x}) \right)$$~~

~~$$y = 0.16x - 3.2 + 120$$~~

~~$$y - 120 = 4x - 80$$~~

~~$$y = 0.16x + 116.8$$~~

~~$$y = 4x - 80 + 120$$~~

~~$$y = 4x + 40$$~~

$$y = 640$$

\bar{x} = is the mean of x series

\bar{y} = is the mean of y series

r_{xy} is the regression coefficient of x on y .

σ_y

$r \frac{\sigma_y}{\sigma_x}$ is the regression coefficient of y on x .

$$\text{Q. prove } r = \sqrt{b_{xy} * b_{yx}}$$

We know that

$$b_{xy} = r \frac{\sigma_x}{\sigma_y} \rightarrow ①$$

$$b_{yx} = r \frac{\sigma_y}{\sigma_x} \rightarrow ②$$

$$b_{xy} * b_{yx} = r \frac{\sigma_x}{\sigma_y} * r \frac{\sigma_y}{\sigma_x}$$

$$r = \sqrt{b_{xy} * b_{yx}}$$

Assumptions for the coefficient of correlation

(i) Both the regression coefficient will have the same sign
 either there will be a +ve or -ve it is never
 possible that one coefficient is +ve another one
 is -ve.

(ii) The coefficient of correlation are should not
 greater than 1. must be less than 1.

(iii) If $b_{xy} = 1.2$

$$r = 1.29$$

$$b_{yx} = 1.4$$

M	T	W	T	F	S	S
Page No.:						
Date:						

(iii) The coefficient of correlation will have the same sign as that of regression coefficient.
 If regression coefficient have -ve sign or r will also be -ve and if regression coefficient +ve r is also +ve.

Logistic Regression

(Predict the future data point)

The main objective of the Logistic Regression is to find a line which is perfectly separable between +ve and -ve classes. It is a binary class classification problem where the target variable (y) is a binary nature.

INPUT FEATURE OUTPUT FEATURE

Test 1

Test 2

y_{ij}

i.

2

1

+1

3

4

0

-1

5

6

1

+1

Testing

8

9

0

-1

10

18

0

-1



New feature

Test 1 and Test 2 are called as the input feature and y is the target feature.

Let us take the 10 data points out of this 80% data will be training and 20% data will be used for testing.

Robust

Model

Once the model is trained and it is ready to accept the unseen data.

	tut 1	tut 2	y
New data	11	8	?

The model predicts either 1 or 0. If it 1 then the person is having a symptom otherwise not. If the model says the robust then the difference b/w training and testing is minimized. that means a model generalises.

Training Accuracy	Testing Accuracy	Difference
92%	91%	1%

If we have the two I/P features then we can visualise with the line if we have 3 I/P features we can visualise with the plane otherwise hyperplane.

Why 2D coordinate space

Test 1 and Test 2 consist of +ve and -ve

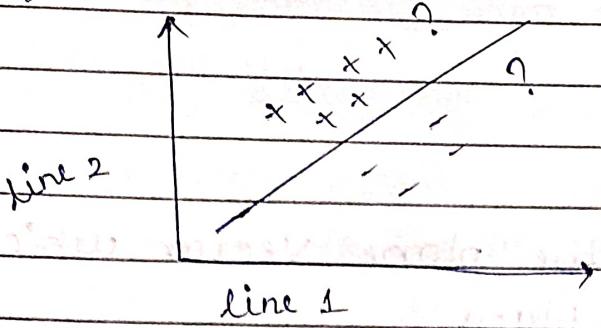
	+	-	+	-	?
Test 1	+	-	0	-1	?
Test 2	+	-	+	-	?
					Hyperplane, anti line

Test 1

* Linearly separable

The concept says that all the data points equally distributed.

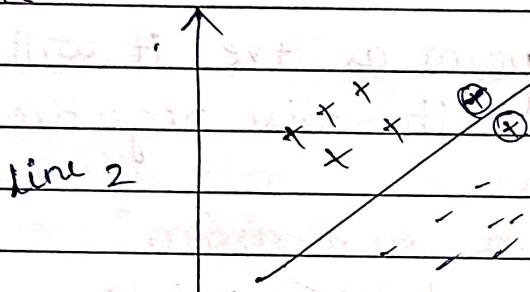
example :-



* Almost linearly Separable

In this case all most all the data points are classified and some points available both sides hence it is called as Almost linearly Separable.

example :-



* Non-linearly Separable.

You take any line in any direction you cannot separate the +ve and -ve classes in the dataset that's why we call as non-linear lines.

Logistic regression won't work with this dataset.

Q. How this line helpful for predicting future data points?

The logistic regression found as the best line it will use this equation of the line for predicting future data points. Now apply equation of the line

$$w_1x_1 + w_2x_2 + w_0 = 0 \longrightarrow ①$$

$$w: [w_1, w_2]$$

$$x: [x_1, x_2]$$

where w = Normal to the plane vector which is
perpendicular to the plane

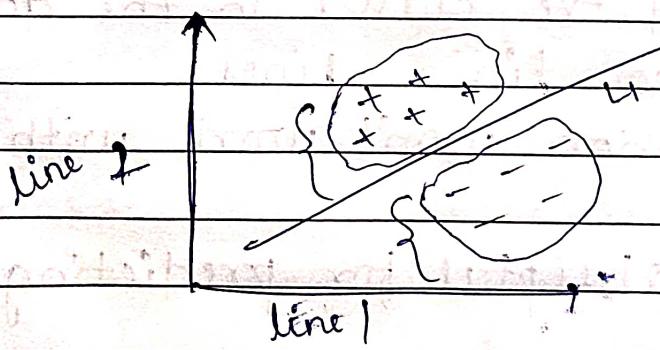
Assume, w is the perpendicular line if the point is
available on the side then it's a +ve
distance otherwise -ve distance

based on this concept our logistic regression
our future predicted data point

If you got the distance point as +ve it will
predict the +ve classes otherwise negative classes.

Q. Why it is the best line?

Out of this 8 data points 4 +ve points and
4 -ve points are perfectly classified with
each other out of these 8 lines L_i is
the best line because maximum data
points are correctly classified.



If we put all the concepts into the equations to identify the best line we use

$$z_i = y_i d_i \rightarrow ①$$

where y_i is the actual value

d_i is the predicted value.

z_i is the Resultant.

$$w_1 x_1 + w_2 x_2 + w_0 = 0$$

$$w: [w_1, w_2] \quad x: [x_1, x_2]$$

$$[w^T x_i + w_0 = 0] \rightarrow ②$$

$w_0 = \text{Origin}$

* How to calculate d_i ?

$$d_i = \frac{w^T x_i + w_0}{\|w\|} \rightarrow ③$$

From equation ③ if the line passes through the origin then d_i can be written as

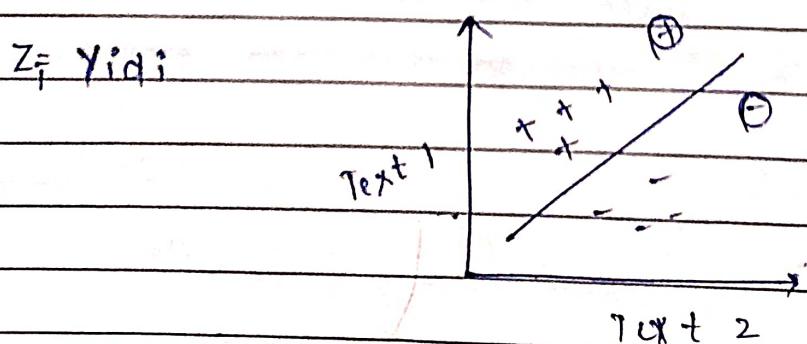
$$d_i = \frac{w^T x_i}{\|w\|} \rightarrow ④$$

| w^T - transpose of w

from equation ④

$$d_i = w^T x_i \rightarrow ⑤$$

equation ④ and ⑤ we are taking our assumptions and $y_i d_i$ can have the followings



future data points

1. If y_i is +ve

d_i is +ve

then the resultant is $z_i = y_i d_i$ is +ve.

2. If y_i is +ve

d_i is -ve

then the resultant is $z_i = y_i d_i$ is -ve.

3. If y_i is -ve

d_i is +ve

then the resultant is $z_i = y_i d_i$ is -ve.

4. If y_i is -ve

d_i is -ve

then the resultant is $z_i = y_i d_i$ is +ve.

$$\operatorname{argmax} \sum_{i=1}^N y_i d_i$$

$$w = w_1, w_2$$

Let us consider to calculate $y_i d_i$ all the datapoints and sum them

$$\sum_{i=1}^N y_i d_i$$

$$y_1 d_1 +$$

$$y_2 d_2 +$$

$$y_3 d_3 +$$

$$y_4 d_4 +$$

:

:

$$y_n d_n$$

The general term of the equation is

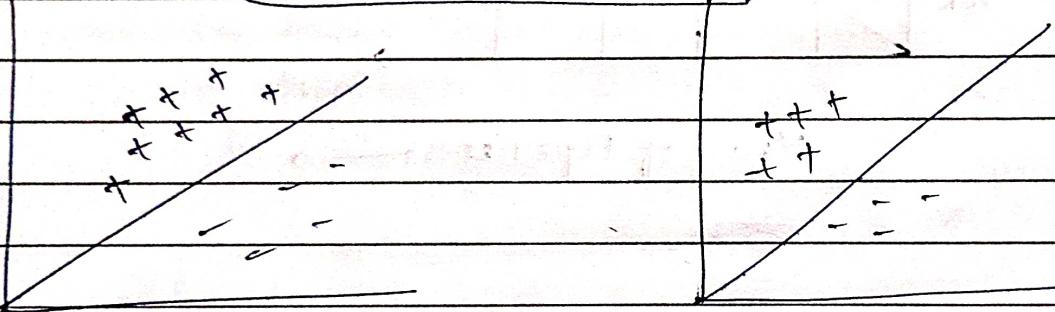
$$\sum_{i=1}^N y_i d_i \longrightarrow \textcircled{5} \quad \textcircled{6}$$

Let us take two lines L_1 and L_2 calculate the summation of both L_1 and L_2 line for which line summation value got maximum.

Let $L_1 = 8 + ve$ and $L_2 = 5 + ve$ out of these two lines L_1 is having correctly classified and the logistic regression equation can be written as for best line is

$$\operatorname{argmax} \sum_{i=1}^N y_i d_i \longrightarrow \textcircled{7}$$

$$w = w_1, w_2.$$



The above is called as optimization function the meaning of this equation is for what value of w the entire summation function will be maximum that w 's need to be calculated.

We will solve this equation and identify w_1 and w_2 in equation 7 and it can be written as $y_i d_i$

$$\text{Substitute } d_i = w^T x_i$$

Hence eq(1) can be written as

$$z_i = y_i (w^T x_i) \longrightarrow \textcircled{8}$$

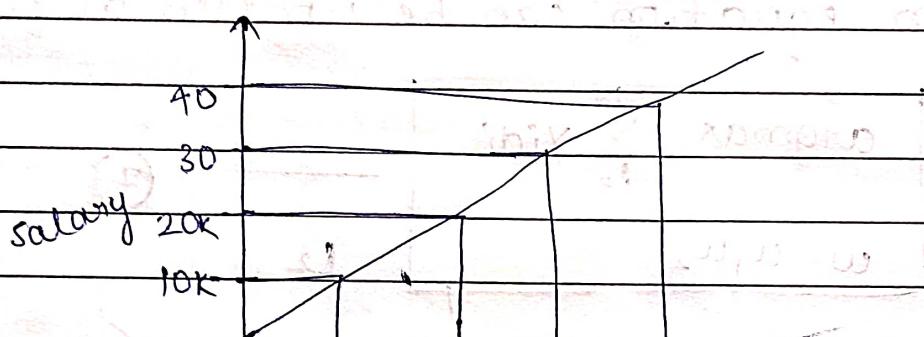
$$\text{Substitute } w^T x_i$$

$$\text{Hence } z_i = y_i (w_1 x_1 + w_2 x_2) \longrightarrow \textcircled{9}$$

In this equation ⑨ $y_0 + w_1 x_1 + w_2 x_2$ are given and w_0 and w_1 are required to calculate.

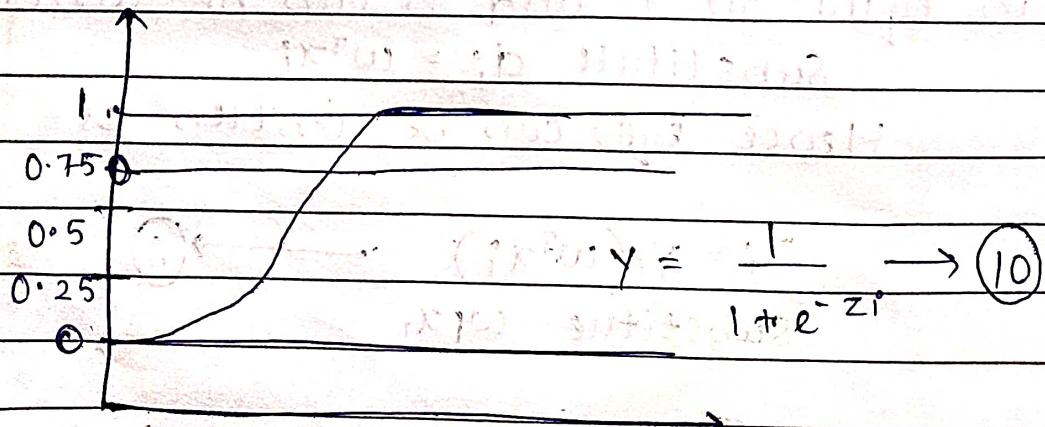
Let us take two different examples

- An Organization want to determine the employees based on the year of experience in this example we will apply the ~~linear~~ linear regression and we can predict the salary.



Years of Experience →

- In this example, if the Organization want to know whether employee would get employed based on their performance. In this case the linear regression won't help us. We will clip the line from 0 and 1 and convert it into a sigmoid curve.



$$y =$$

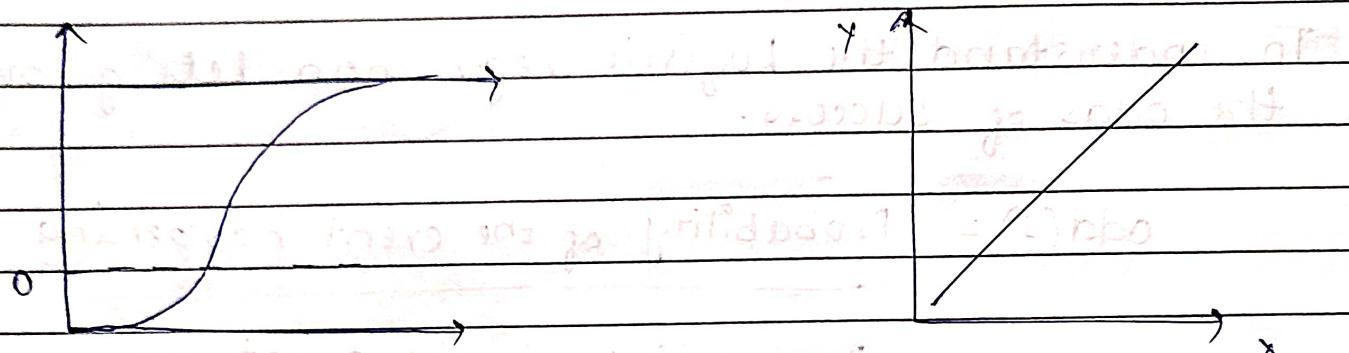
$$\frac{1}{1 + e^{-zi}}$$

→ 10

* Difference b/w Logistic Regression and Linear Regression.

classification

Prediction, if t



sigmoid function
(S-shaped)

$$y = ax + b$$

$$y = b_0 + b_1 x$$

$$y = \frac{1}{1 + e^{-xi}}$$

- (i) This approach is useful when the nature of the dataset says binary or multi class classification approach
- * It is used to predict the future events by taking dependent and independent features.
- (ii) In Logistic regression the features are categorical nature.
- * The response of linear regression is a continuous nature.
- (iii) In linear regression the best fit line obtain in the straight line.
- (iv) In Logistic regression the line is obtained as Sigmoid, shape.

* In linear regression the equation is

$$y = b + ax$$

(iv) In LR the equation is

$$y = \frac{1}{1 + e^{-zi}}$$

* How does the logistic regression algorithm works.?

To understand the Logistic regression let's go over the odds of success.

odd(s) = Probability of the event happening

Probability of the event not happening

$$\text{odd}(q) = \frac{P}{1-P} \quad (1)$$

The values of odd's range are 0 to ∞

and the values of the probability lies

Consider the equation of the straight line $y = B_0 + B_1 x$

Now to predict the odds of success, we take log on odd formula

$$y = B_0 + B_1 x$$

$$\log\left(\frac{P(x)}{1-P(x)}\right) = B_0 + B_1 x \quad (2)$$

exponential on both sides.

$$e^{\ln} \left(\frac{P(x)}{1-P(x)} \right) = e^{B_0 + B_1 x}$$

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

$$e^{\ln(20)} = 20$$

$$\frac{P(x)}{1-P(x)} = y$$

13

In the above eqn we substitute y on right hand side and eqn becomes

$$\frac{P(x)}{1-P(x)} = y \quad \rightarrow 14$$

Simplify the eqn 14

$$\therefore y = e^{B_0 + B_1 x}$$

$$P(x) = y(1-P(x))$$

$$P(x) = y - yP(x)$$

$$P(x) + yP(x) = y$$

$$P(x) = \frac{y}{1+y} \quad \rightarrow 15$$

Substitute the value of y in these equations. 15

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

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

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

The sigmoid function equation can be obtained as

$$P(x) = \frac{1}{1 + e^{-(B_0 + B_1 x)}} \rightarrow (17)$$

We know the eqn of the straightline $y = B_0 + B_1 x$
Substitute in the equation for (17)

$$P(x) = \frac{1}{1 + e^{-y}} \rightarrow (18)$$

Q. Let us consider the student dataset which are
entrance based on the histogram who are
selected and non-selected based on the
Logistic represent.

$$b_0 = 1$$

$$b_1 = 8$$

assuming the mark $K = 60$

(a) compute the resultant class

$$y = b_0 + b_1 x_1 = 1 + 8 \times 60 = 481$$

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

$$P(x) = \frac{1}{1 + e^{-481}} = 1$$

The threshold value as 0.5 compare the result with the threshold value and interpret it.

We observed that the resultant value is greater than the threshold (0.5). Therefore, the candidate mark with 60 is selected.

Q. Let us consider the dataset having two features hours study, pass and fail.

Apply the logistic regression to answer the following question.

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

(ii) Atleast how many hours student should study that max he will pass the score with the prob. max 95%.

The optimization produce the model hours fail/pas is $\log(\text{odds}) = -64 + 2 \times \text{hours}$

Assume the model suggested the best fit is

optimizer for odds of the passing score is $\log(\text{odds}) = -64 + 2 \times h$

$$\text{Sol} - P(x) = 1$$

$$(i) \quad 1 + e^{-z}$$

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

$$Z = -64 + 2 \times 33$$

$$Z = 2$$

$$P(x) = 1$$

$$1 + e^{-2}$$

$$P(x) = 0.88$$

(ii)

$$P(x) = \frac{1}{1+e^{-z}} = 0.95$$

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

$$1+e^{-z} = \frac{1}{0.95}$$

$$1+e^{-z} = 1.0526$$

$$e^{-z} = 1.0526 - 1$$

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

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

$$-z = \cancel{-2.94} \quad 2.94$$

$$2.94 = -64 + 2 * \text{hour}$$

$$66.94 = 2h$$

$$\underline{66.94 = h}$$

2

$$\boxed{33.47 = h}$$

L₁ Regularization [Lasso Regression]

M	T	W	T	F	S
Page No.:					
Date:					

Ridge Regression. (L₂ Regularization)

- The objective of the Ridge regression is to prevent the over fitting to meet the model's generalization ability.
- It is used when your dataset is having multi co-linearity features.
- When the model is having high variance / low bias and high bias with low variance as highly correlated of one attribute with other attributes then Ridge Regression will be useful.
- The Ridge Regression can be defined as cost function of linear regression plus (+) penalty term. The above concept is called as a Regularization.

Mathematically, The Ridge Regression can be defined as

$$\text{minimise } \left(\text{SSE} + \lambda \sum_{j=1}^p \beta_j^2 \right) \rightarrow ①$$

→ linear regression (sum of squares) of error
 λ term is need to find overfitting / underfitting.

Where λ is a hyperparameter of L₂ regularization
 $(\lambda$ is the regularization parameter)

β - coefficient of linear regression model.

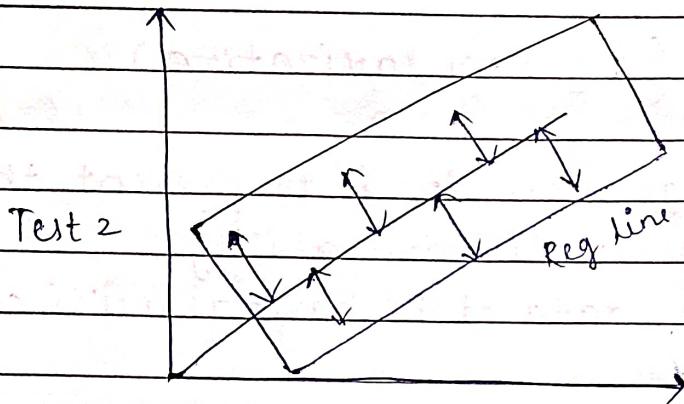
p - Number of features.

The above equation can also be written as

Ridge regression = Cost Function +

(of LR)

Hypoparameter of L₂



Q. How to avoid overfitting?

To avoid the overfitting issue we required the cross validation technique that helps in estimating the error over the test set and it will decide what parameters works best for the model.

The cost function of the simple LR is

$$\sum_{i=1}^M (y_i - \hat{y})^2 = \sum_{i=1}^M \left(y_i - \sum_{j=0}^p w_j x_{ij} \right)^2 \rightarrow (3)$$

where

M is the number of instances

p is the number of features

$$\hat{y} = w_0 x_0 + w_1 x_1 + \dots + w_N x_N + b$$

b is the constant

This can be written as the linear model based on

n number of features and considering a

single feature w_0 is the slope

b is the intercept

The LR looks for optimizing w and p such that it minimizes the cost functions

In Ridge regression the cost function is altered by adding the penalty term which is equivalent to square of the magnitude of the coefficient and the equation can be written as.

$$\sum_{i=1}^M (y_i - \hat{y})^2 + \lambda \sum_{j=0}^P w_j^2 \rightarrow ④$$

Which is called as cost function of linear regression. Where w_j is the square of magnitude of the coefficients.

The ideal value of the λ should be 0 but not exactly zero. The RR shrinks the coefficient and it allows to reduce the complexity of the model. So the lower the λ value the model provides good results if $\lambda = 0$ then RR becomes LR.

* Let us take one example $y = \beta_0 + 1.2x_1 + 20x_2 + 39x_3$. In the above case x_3 plays the important feature the magnitude of the coefficient is high and it will give less importance to x_1 and x_2 .

Our objective is to minimise the high coefficient variables otherwise it requires more expansive cost and the model becomes complex.

To avoid the above problem (Reducing the models complexity) Regularization is used.

* The ridge regression can be defined

$$\text{Ridge} = \text{Loss} + \lambda \|w\|^2 \rightarrow ⑤$$

$$\text{where } \|w\| = w_1^2 + w_2^2 + \dots + w_n^2.$$

$$\text{Hence, Ridge} = \text{loss} + \text{penalty.} \rightarrow ⑥$$

M	T	W	T	F	S
---	---	---	---	---	---

 Page No.: | Date: |

The penalty is to reduce the losses where γ is the constant. If w is the vector of coefficients.

Lasso Regression → Feature selection technique

The main objective of Lasso regression is, to minimise the difference b/w observed value and predicted value as well as penalty term λ and absolute size of the coefficients.

As the coefficients are determined by the absolute size it help us to solve the overfitting problems.

LASSO - Least absolute shrinkage and selection operation.

It is a kind of LR model that use a penalty term and called L_1 regularisation.

* OBJECTIVE FUNCTION OF LASSO REGRESSION

In Lasso regression the objective function is to define the as +

between actual and predicted.

Mathematically,

$$\text{minimize}_{\beta} \frac{1}{2m} \sum_{i=1}^n (y_i - \hat{y})^2 + \lambda \sum_{j=1}^p \|\beta_j\|$$

$\left. \begin{array}{l} \text{Sum of squares of the residual (RSS)} \\ \text{Cost function of linear regression} \end{array} \right\}$

Where,

Lasso = RSS + penalty.

n = number of observations / instances

p = number of features / predictors.

β = Coefficients.

y_i = Actual

\hat{y}_i = predicted

λ = hyperparameter / a constant for Lasso.

The above function is called as a objective function of Lasso regression.

The constants controls the strength of the penalty term the higher value of the λ gives more regularization.

The cost function of Lasso regression can be written as

$$\sum_{i=1}^M (y_i - \hat{y})^2 = \sum_{i=1}^M (y_i - \sum_{j=1}^p w_j x_{ij})^2 + \lambda \sum_{j=1}^p |w_j|$$

Let us take one example where,

$$y = \pi_0 + 1.2x_1 + 20x_2 + 39x_3$$

After some iteration the y becomes

$$y = 0.9 + 0x_1 + 0x_2 + 5x_3$$

In the above example we understood that the higher coefficient impact to the model and it requires more cost so that the model becomes complex. In order to avoid this situation we use the concept of regularisation.

The Lasso regression not only helps in reducing overfitting but also help us to reduce in the feature selection purpose.

- * The main objective of Lasso is to use a regularization as well as feature selection. This type of regularization can leads to zero coefficients.
- * Some of the features are completely neglected (x_1 and x_2)

In the above example x_1 and x_2 variables are completely neglected. This values are not important for determining the value of y . Lasso is called as a feature selection technique.

$$Y = 0.9 + 5x_3$$

The complexity reduced because unnecessary features are eliminated overfitting problem solved because we have done feature selection. When we increase the value of λ some of the variable become zero.

* ~~ORDER~~ ~~REGRESSION~~ II

* ORDER LEAST SQUARE REGRESSION

The main objective of OLS regression is to minimize the sum of square difference between the actual and the predicted values by the linear model.

The main role of OLS is to estimate the coefficients of a linear regression model.

At the end \hat{y} can be estimated.

When we perform linear regression we use OLS method to estimate the coefficients of the linear equation that best fits the data.

$$y = \beta_0 + \beta_1 x$$

$$y = 0.2 + 6(2)$$

$$y = 0.2 + 12$$

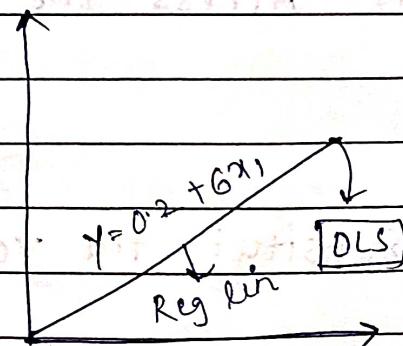
$$y = 12.2$$

$$x$$

$$2$$

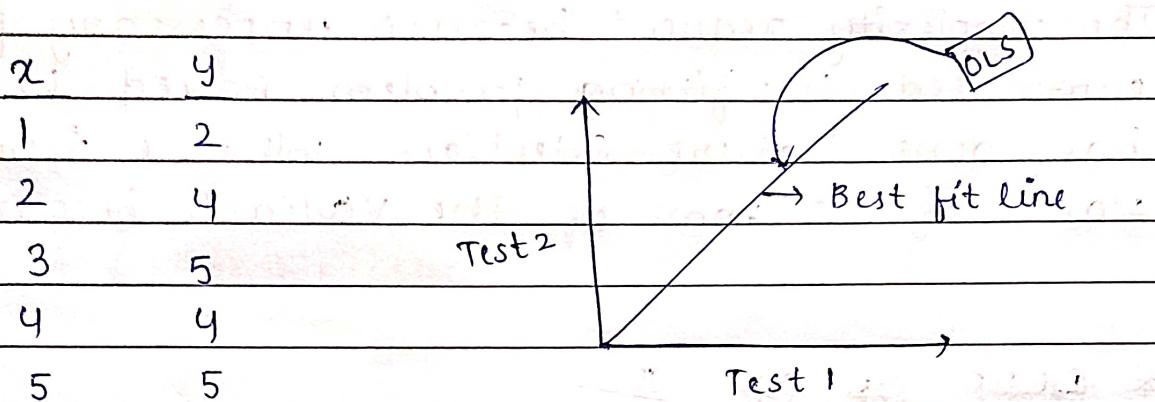
$$3$$

$$4$$



Steps.

- (i) Let us take the dataset x and y apply the OLS regression technique to estimate the regression equation. Use the hypothesis testing to determine whether the estimated coefficients are statistically significant or not.



The estimated regressor line

$$y = 2.2 + 0.6x \rightarrow \text{regressor line}$$

- (ii) \hat{y} value
- 2.8
 - 3.4
 - 4
 - 4.6
 - 5.2

- (iii) Access the goodness of fit

$$\hat{\epsilon}^2 = \sum \hat{\epsilon}_i^2$$

substitute the value of Numerator and denominator

$$R^2 = 1 - \frac{\sum (\hat{y}_i - \bar{y})^2}{\sum (y_i - \bar{y})^2}$$

$$R^2 = 0.60$$

The R^2 value indicates that 60% of the variance y is explained with the linear regression equations.

(iv) TheOLS In OLS we will use the hypothesis testing to determine whether statistically significant or not.

* The Null hypothesis

$$H_0 : \beta_0 = 0$$

β_0 = intercept.

similarly, $H_0 : \beta_1 = 0$

β_1 = Slope

* Alternate hypothesis

$$H_0 : \beta_0 \neq 0$$

(Intercept)

$$H_0 : \beta_1 \neq 0$$

(slope)

* performe the hypothesis test for β_0 .

(i) estimated coefficient ($\hat{\beta}_0$) = 2.2

(ii) SE (standard error of $\hat{\beta}_0$) [From linear regression output]

(iii) DF - Degree of freedom

$$DF = n - 2$$

where, n is the number of observation.

(iv) calculate the t value $t = \frac{\hat{\beta}_0}{SE(\hat{\beta}_0)}$ t-distributions

defn $t = \frac{\hat{\beta}_0 - \beta_0}{SE(\hat{\beta}_0)}$ distribution of statistics

$SE(\hat{\beta}_0)$ student

(v) Compare the calculate t value with a critical value from the t distribution table.

(vi)

$SE(\hat{\beta}_0)$

By calculating the SE of the intercept ($\hat{\beta}_0$) the formula is

$$SE = \sqrt{\frac{MSE}{n} + \frac{(\bar{x})^2}{SS_x}}$$

where MSE = Mean square error / residual error

n - number of observation

ss_x - sum of square of independent variable

\bar{x} - Mean of x .

In this equation Known values are

$$n = 5$$

$$\bar{x} = 3$$

$$MSE = \frac{\sum (\hat{y} - y)^2}{n-2}$$

$$MSE = (1 - R^2) * var(y)$$

$$R^2 = 0.4 * \frac{\sum_{i=1}^n (y_i - \bar{y})^2}{n-1}$$

$$var(y) = 0.4 * 1.5$$

$$MSE = 0.6$$

* To calculate ss_x we have the formula

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

$$ss_x = 10$$

Plugged in the value of

$$SE = \sqrt{\frac{SSE}{N} + \frac{(X)^2}{SSX}} = \sqrt{\frac{0.6}{5} + \frac{9}{10}} = 1.009$$

The degree of freedom = 5 - 2

$$DF = 3$$

calculate the t value

$$\beta_0 = 2.2$$

$$t = \frac{\beta_0}{SE(\beta_0)} = \frac{2.2}{1.00} = 2.2$$

→ The comparison

The calculate t value is = 2.2

and our DF is 3

we can find the critical t value from t distribution

table at level of significance is 0.05 is 3.182.

(*)

Since 2.2 falls below 3.182 We do not have the evidence to reject the null hypothesis

$$H_0: \beta_0 = 0$$

Therefore, we conclude that intercept is not statistically significant different from 0 at 0.05 level of significance.

UNIT-2

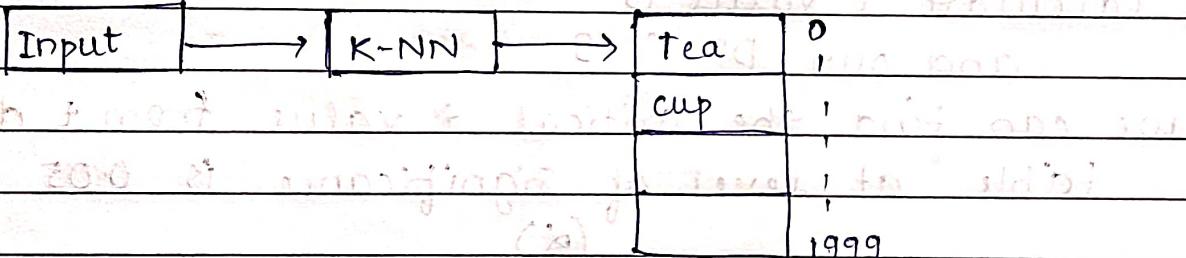
M T W T F S S

Page No.:

Date:

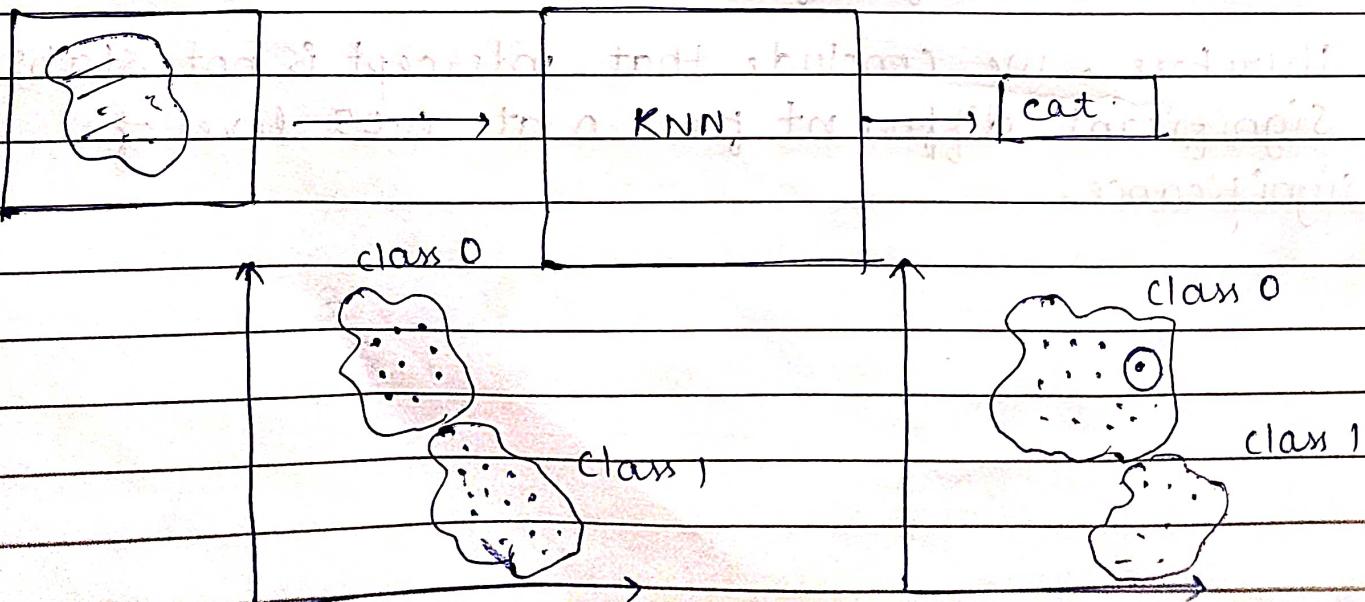
K-NN Algorithm :- (K-Nearest Neighbour)

- (i) KNN stands for K-Nearest Neighbour
- (ii) It is a supervised Machine learning algorithm. That can be used for both regression and classification.
- (iii) A classification problem has a discrete value at its output. Each row in the dataset typical called an example, observation or a datapoint while each column is called a predictor, dimension, independent variable or feature.
- (iv) KNN algorithm captures the idea of similarity, sometimes it is called as proximity or closeness.
- (v) KNN algorithm assumes the similarity b/w the new data and available data and put the new data into the category that is most similar to available categories.

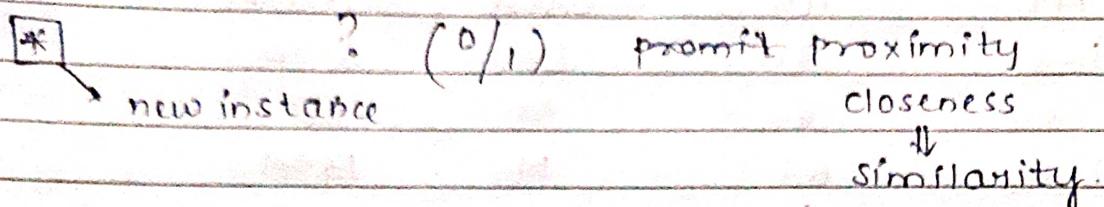


80% → Training

20% → Testing



KNN algorithm can also used for regression as well as classification but mostly used for classification algorithm. It is a non-parametric algorithm, because it does not make any assumptions to the undivided data.



The new instance can be decided by using the distance measure formula

① Euclidean Euclidian distance.

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

② Manathan

$$|x_2 - x_1| + |y_2 - y_1|$$

It is also called as Lazy learner algorithm. It does not learn from the training set immediately instead it stores the dataset and at the time of classification it performs an action ^{from} on the dataset.

* How does K-NN works?

- (i) Initialise the number of K neighbours
- (ii) Calculate their ~~first~~ distance using euclidian distance
- (iii) Take the ~~nearest~~ K nearest neighbour as per ~~egs~~ euclidian distance.
- (iv) Count the number of points in each category.
- (v) Assign the new datapoints to the category for which the number of neighbour is maximum.

Q. Let us take one dataset which consist of the following information apply KNN classifier to predict their class level.

HT	wt	class	Distance
167	51	UN	6.70
182	62	N	13
176	69	N	13.41
173	64	N	7.61
172	65	UN	8.2
174	56	UN	4.12
169	58	N	1.41
173	57	N	3
170	55	N	2
170.	57	?	0

N

Find the distance b/w new instance to first instance. to estimate using distance formula.

$$\sqrt{(170-167)^2 + (57-51)^2} = 6.70$$

Ideally $K=3$ or
 $K=5$

M	T	W	T	F	S	S
Page No.:						
Date:						

- ④ Restructure the table according to the minimum distance.

HT	WT	Class	Distance	Rank
169	58	N	1.41	1
170	55	N	2	2
173	57	N	3	3
174	56	UN	4.1	4
167	51	UN	6.70	5
173	64	N	7.6	6
172	65	N	8.2	7
182	62	N	13	8
176	69	N	13.41	9
170	57	N		

- ⑤ Provide the rank according to the minimum distance.

Step-3

Assume $K=3$

$R=1 \rightarrow (N)$

$R=2 \rightarrow (N)$

$R=3 \rightarrow (N)$

underweight (UN)

Normal

170, 57, N

As 170 (Height) and 57 (Weight) is closer to the three nearest values so the new instance will be of normal class.

SVM - Support Vector Machine

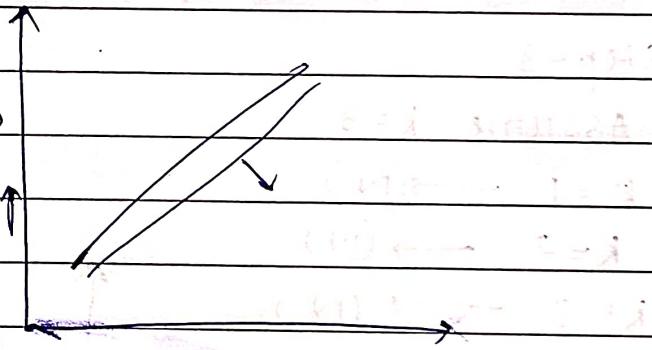
SVM is a supervisor learning. It can be possible both classification as well as regression. If we predict the class level then it is called as Support Vector for Regression (SVR). If it classify the class level then it is called as support vector for classification (SVC).

The main objective of SVM is to predict the level class level by using based the concept of hyperplane.

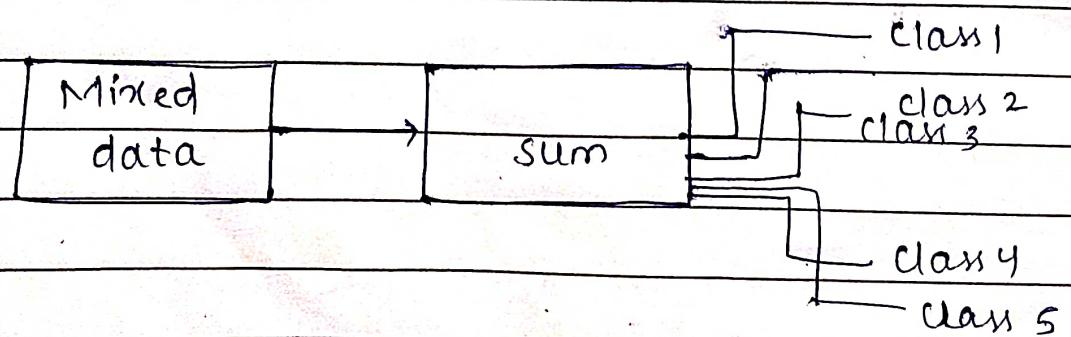
Representation of hyperplane.

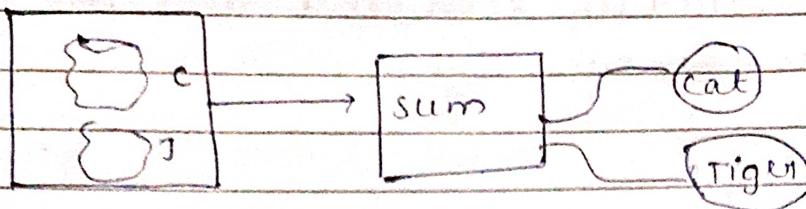
It is also called as a decision boundary which will classify both the +ve and -ve classes.

It is the goal of SVM to create the best line or decision boundary that can segregate n dimensional space into classes so that it can easily classify the new instance in the future. The best decision boundary is called as hyperplane.



There may be 'n' number of boundaries. we choose one of the best from it we choose and call it as hyperplane.

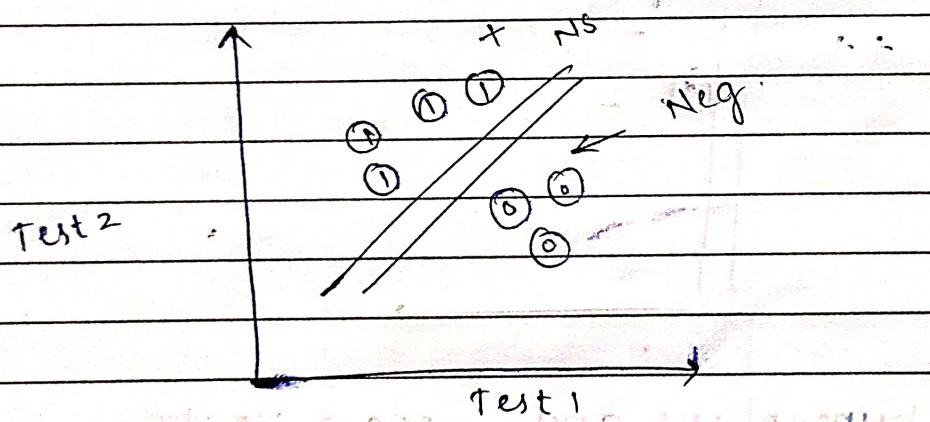




Data point	test 1	test 2	Result	
1	-	-	yes	1
2	-	-	yes	1
3	-	-	yes	1
4	-	-	No	0
5	-	-	No	0
6	-	-	No	0
New instance	?	-	yes.	?

Predict the class level.

By seeing the result of the above table we can generate / estimate identify the hyperplane which will divide the +ve and -ve classes.

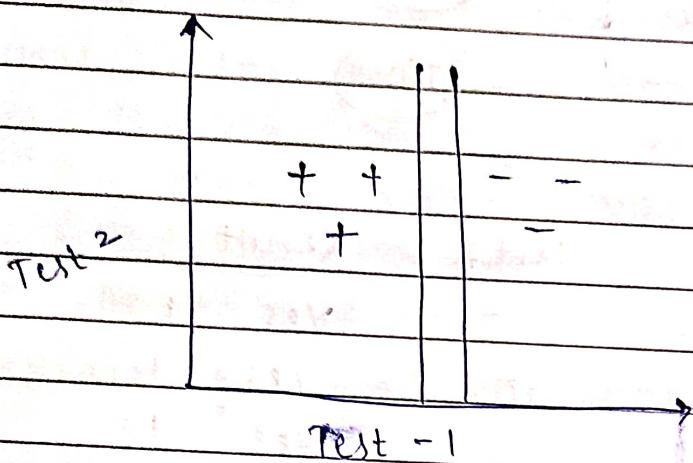


There are two types of SVM

(i) Linear SVM

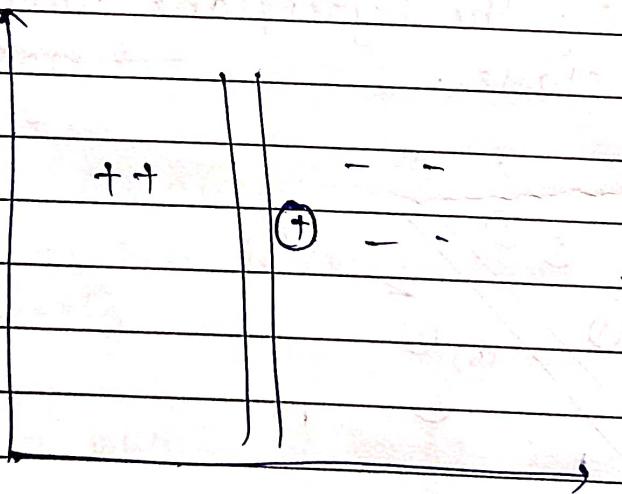
In linear SVM we use the linearly separable data which means if a dataset can be classified into two classes by using a single straight line then such

data is termed as linearly separable and the classifier is called as linear SVM.



(iii) Non-linear SVM

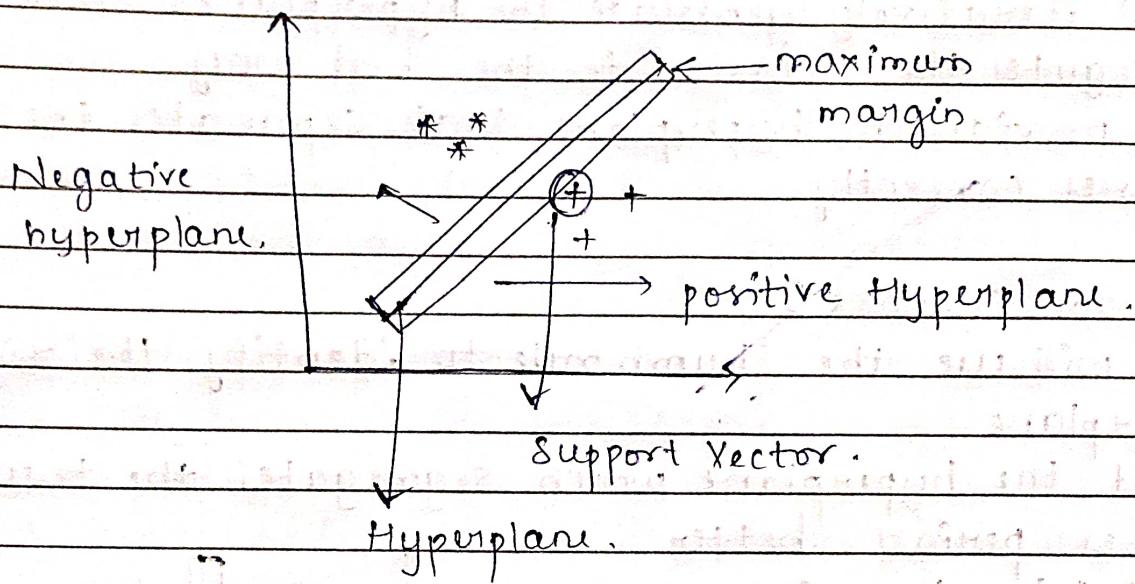
It is used for non-linearly separated data which means if a dataset cannot be classified by using a straight line then such data is termed as non-linearly separable and the classifier is called as Non-linear SVM.



* The concept of hyperplane and support vectors in SVM

Those points which touch the boundary and very closest to boundary they are called as the support vector

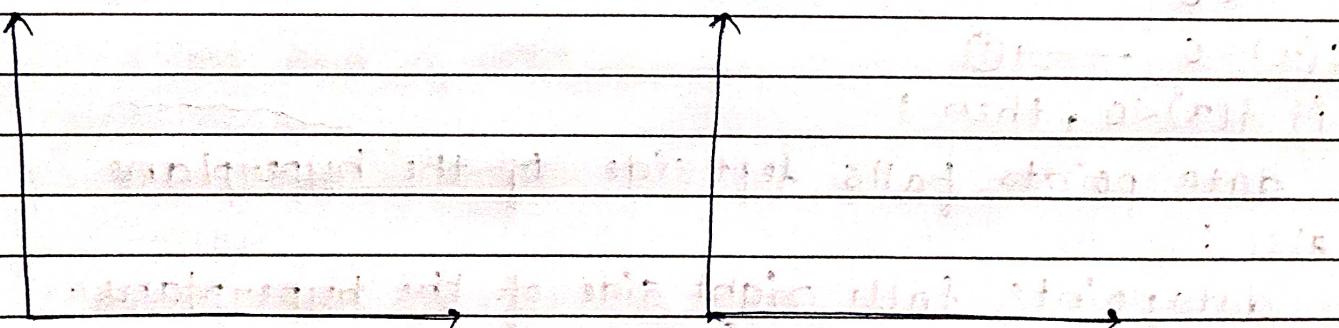
There can be multiple decision line to segregate the classes in n dimensional p but we need to find best decision boundary that helps to classify the datapoints.



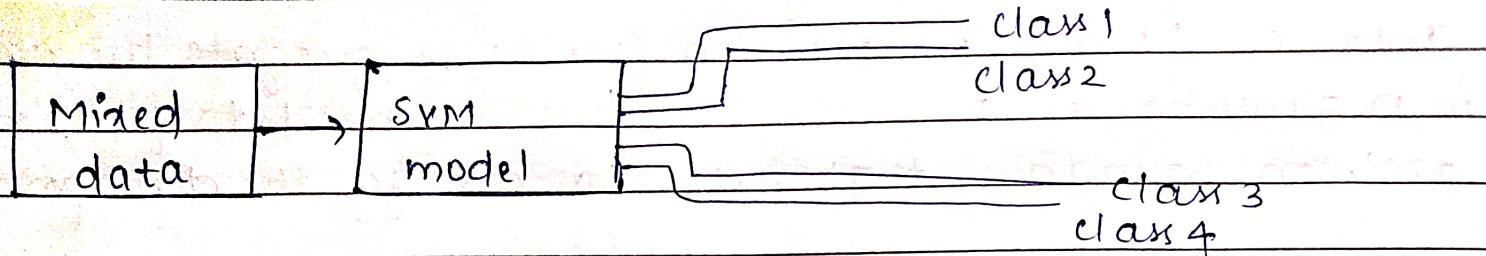
Margin :-

It may be defined as the gap between two lines on the closest data point of the different classes. It may be positive hyperplane and negative hyperplane.

It can be calculated as the perpendicular distance from the line to support vectors.



Large margin is considered as a good margin and a good hyperplane similarly small margin is considered as a bad margin not and not a accurate hyperplane between the feature.



- 1) SVM iteratively generate the hyperplanes that segregate the classes in the best way.
- 2) It provides the hyperplane that separates the classes correctly.

* How does it work?

- 1) We will use the thumb rule to identify the right hyperplane.
- 2) Select the hyperplane which segregates the two classes period better.
- 3) Maximizing the distance between nearest data points and hyperplane that distance is called as Margin.

$f(x) \rightarrow$ SVM objective function

> 0

< 0

$$f(x) = 0 \longrightarrow \textcircled{1}$$

if $f(x) < 0$, then :

data points falls left side of the hyperplane.

else :

data point falls right side of the hyperplane.

$$[f(x) = w_1x_1 + b_1] \longrightarrow \textcircled{2}$$

\therefore consider the 2-dimensional data points and the equation is $w_1x_1 + b_1$

where w is normal to the line
 x is Input vector
 b is known as constant (Weight Vector)

* Advantages of SVM :-

- It works better if the number of features are large
 Ex - spam filtering, disease prediction, crop recommendation etc.
- SVMs are conceptually easy to understand.

* Disadvantages of SVM :-

- It has the two major constraints
- It works only with real numbers
- All the data points in all the dimensions define numeric value only.
- It works only binary class classification
- Training the SVM is an efficient and time consuming process when the data is large.
- When we have a noise in the data set then it does not work well.

* Application Area of SVM :-

- Regression Analysis
- pattern Recognition
- prediction
- Mail is spam or not.

* Optimization Objective

- SVM allows to minimize the classification error while maximizing the margin.

- SVM allows to minimize the classification error while maximizing the margin.
- The optimization objective function of SVM can be defined as

$$\boxed{\text{minimize}_{w,b} \frac{1}{2} \|w\|^2}$$

③

w is the weight vector, subjective constraint that

$$y_i(wx_i + b) \geq 1 \quad \text{for all } i \rightarrow ④$$

$$\text{We know that } d_i = wx_i + b \rightarrow ⑤$$

$$\text{Substituting eqn } ⑤ \text{ in } ④$$

$$y_i(wx_i + b) \rightarrow ⑥$$

$$y_i(wx_i + b) \geq 1 \text{ for all training examples of, } (x_i, y_i)$$

* Decision function/ Boundary :-

The decision function for classifying the new data.

point x is given by

$$f(x) = wx + b$$

If $f(x) \geq 0$ the data point classify one class and if

$f(x) < 0$ the data point classify other class.

* Elasticnet Regression

The main objective of the elasticnet regression is to improve the performance of the model. Individually, the model perform well but to make the robust model we need to develop the hybridizations.

Model $\rightarrow 92\%$

Lasso Model

Ridge Model $\rightarrow 93\%$

The term ridge regression can be defined as

$$\text{Ridge} = \text{Loss} + \text{Penalty}$$

$$= \text{Loss} + \lambda \|\mathbf{w}\|^2$$

$\|\mathbf{w}\|^2$ = square of magnitudes.

The objective of ridge regression is to handle the regularization. Similarly, Lasso is used for feature selection because some of the features are eliminated. Those dataset is having multi-collinearity affect we will use lasso regression combining these two regression a new regression is formed which is called as elasticnet regression.

$$\text{Ridge} = \text{Loss} + \|\mathbf{w}\|^2 \rightarrow ①$$

$$\text{Lasso} = \text{Loss} + \lambda \|\mathbf{w}\| \rightarrow ②$$

Combining these two equations we get

$$\text{ENR} = \text{Loss} + \lambda_1 \|\mathbf{w}\|^2 + \lambda_2 \|\mathbf{w}\|$$

So the cost function of ENR is ↑

It is called as hybrid regularization.

* What is VIF - Variation inflation factor?

$$VIF_i = \frac{1}{1-R^2}$$

It detects multicollinearity in regression analysis. The term multicollinearity means a dataset is having correlation between the predictors. The VIF estimates how much variation of the regression coefficient is inflated / suffered due to multicollinearity in the models.

$$\frac{1}{1-R^2} = \text{Tolerance}$$

$$1-R^2 = \text{Tolerance}$$

Where, R^2 = Unadjusted coefficient of determination for i^{th} independent variable. Reciprocal of VIF is called as tolerance.

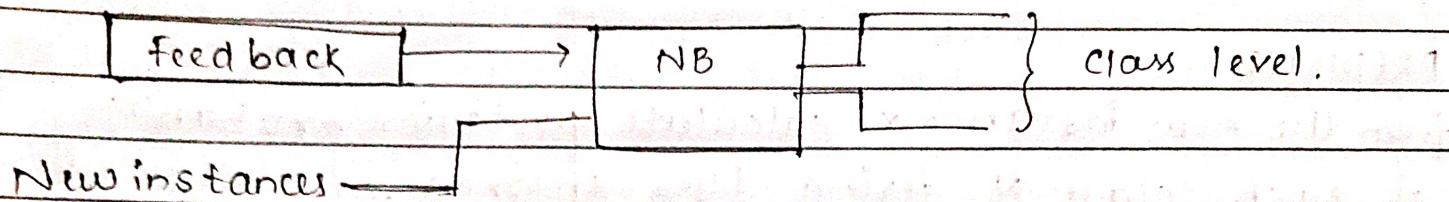
If $VIF = 1$ No correlated

If $VIF = 5$ Moderately correlated

If $VIF > 5$ Highly correlated

NB (Naive Bayes) classifier :-

It is a kind of supervised learning which main objective is for classification. It allows to predict the class levels.



St.	FB	Result
1		1
2		1
3		1
4		0
5		0

6

↑ new instances

The model is trained with a training data of 5 samples when new instance comes the model has to credit the class level. It is based on the bayes theorem with the assumption that features are conditionally independent given in the class level.

• Steps of the Naïve Bayes classifier

1. Data preparation.

split the dataset into two parts training and testing let x and y are two features

$x \rightarrow$ independent.

$y \rightarrow$ dependent.

2. Calculate the prior probability $P(Y)$ for each class in the training examples.
- For each feature x_i and each class y calculate the likelihood, $P(x_i|y)$ which is the probability distribution of feature x_i for given class y .

PREDICTIONS

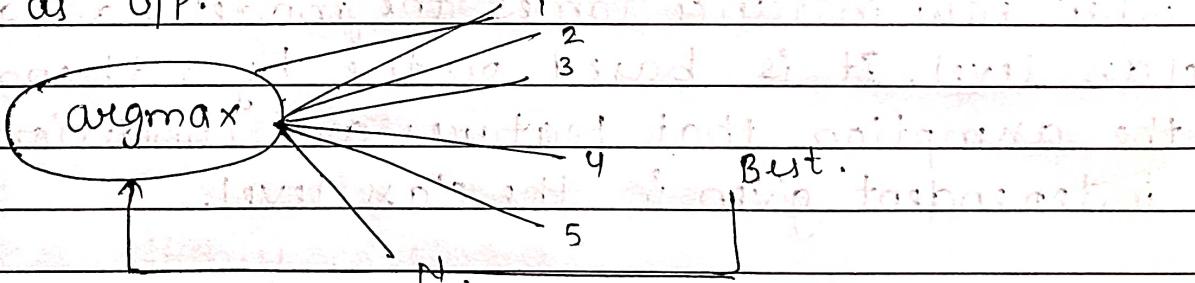
Given the new instance x calculate posterior probability of each class y using Bayes theorem.

$$P(Y|x)$$

$$P(Y|x) = \frac{P(x|Y) \cdot P(Y)}{P(x)}$$

→ Bayes theorem formula

- since $P(x)$ is constant for all classes it can be ignored during comparison.
- The class with the highest posterior probability is predicted as O/P .



optimization and which will do best.
posterior probability.

- It can be represented as

$$y = \operatorname{argmax} p(y|x)$$

4. Evaluation

Evaluate the performance of the ~~classification~~ classifier using appropriate matrix such as accuracy, precision, Recall, F1 score, Decision etc.

* BAYE'S THEOREM

It is a fundamental concept in the probability theory it describe the probability of prior knowledge that might be related to the event.

Mathematically,

The Baye's theorem is represented as

$$P(Y|X) = \frac{P(X|Y) \cdot P(Y)}{P(X)}$$

$P(Y|X)$: posterior probability.

$P(X|Y)$: likelihood

$P(X)$: prior probability of the predictor X .

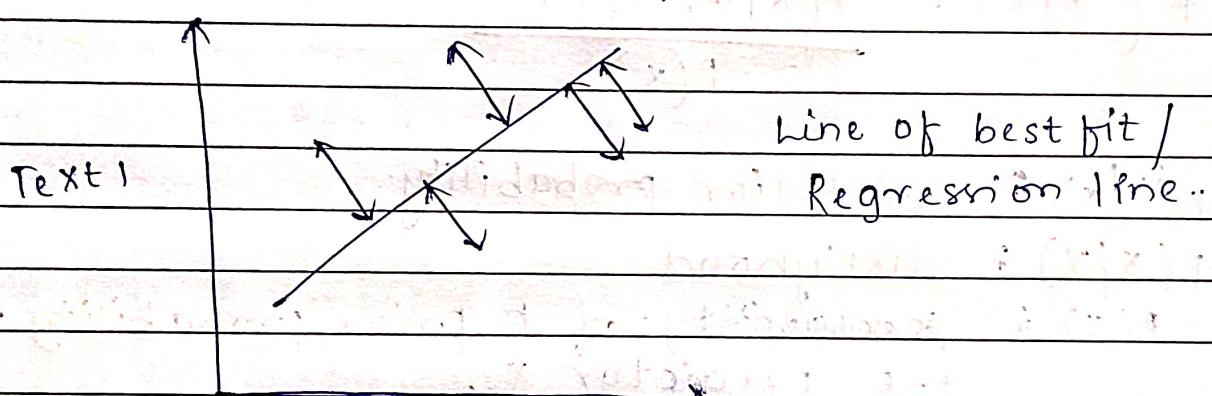
$P(Y)$: prior probability of class Y .

* Applications of Naive Bayes.

- (i) Sentiment Analysis
- (ii) Disease prediction
- (iii) Crop recommendation
- (iv) Emotional detection
- (v) protein Analysis

GRADIENT DESCENT

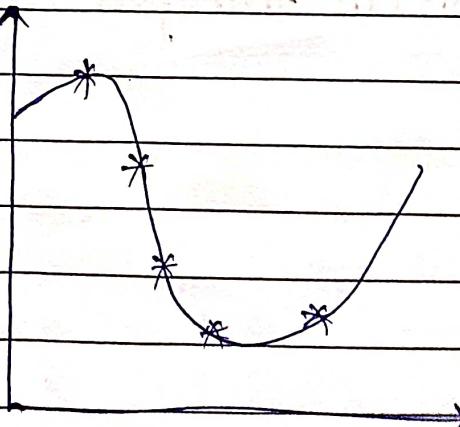
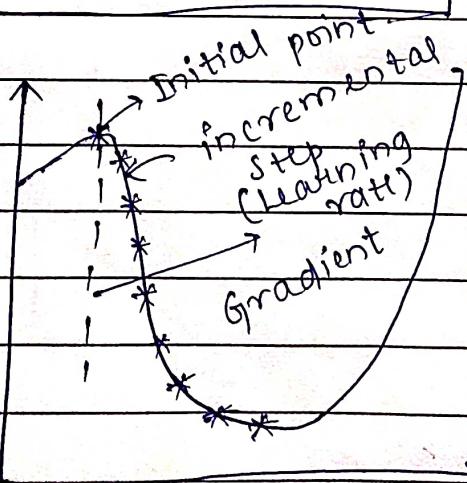
- It is a kind of optimization technique which is used to minimize the loss function of machine learning algorithm.
- In mathematical terminology optimization algorithms refers to the task of minimizing or maximizing an objective function $f(x)$.
- Similarly in ML, optimization is the task of minimizing the cost function.
- The main objective of Gradient descent is to minimize the loss or cost function of the Machine learning classifier.



$$\text{Error} = \text{Actual} - \text{predicted}$$

$$\sum_{i=0}^n (y_i - \bar{y})^2$$

$$x := x - \alpha \frac{\partial L}{\partial w}$$



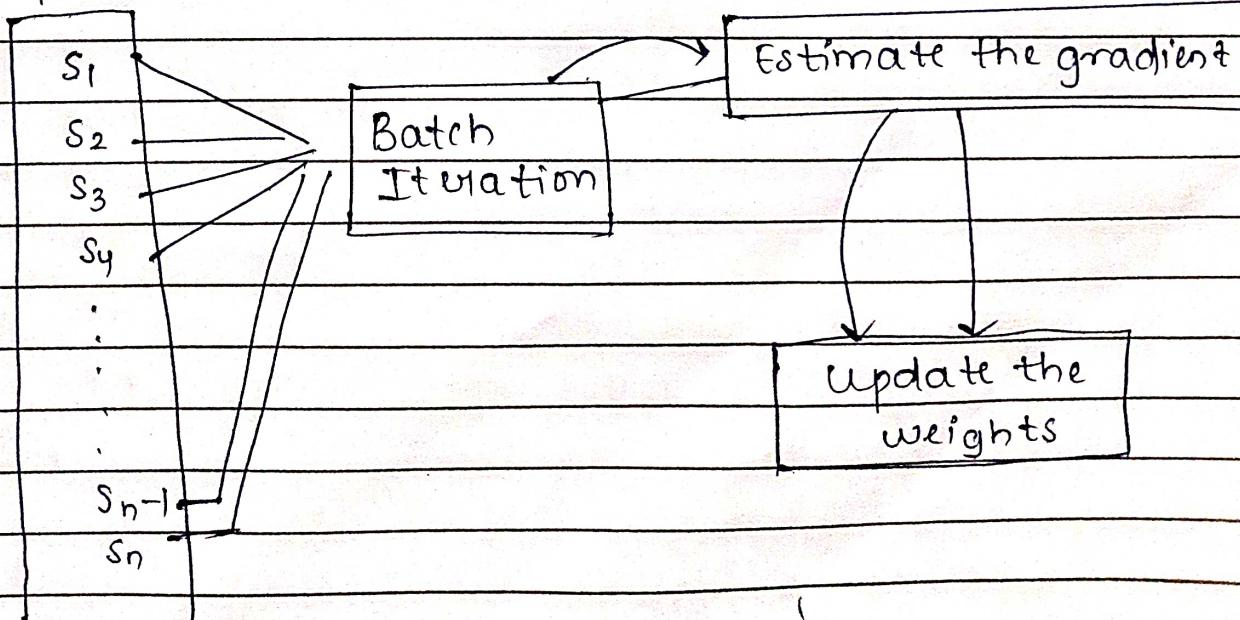
- Gradient descent is defined as one of the most commonly used iterative optimization algorithm of machine learning to train the machine learning and deep learning model.
- It helps to find local minima of a function.
- If we move towards a negative gradient of the function at the current point it will give a local minima.
- If we move towards a positive direction of the function at the current point it will get a local maxima of that function.

Gradient estimation
weights update.

* Types of Gradient Descent :-

- (*) There are 3 types of Gradient Descent.
- (i) Batch Gradient
- (ii) Mini Batch Gradient
- (iii) Stochastic Gradient.

(i) BATCH GRADIENT



- In Batch Gradient Descent we use all the training samples to prepare a single batch and pass into the network.
- It will compute the gradient loss function for each sample.
- Then update the weights.
- If the training sample size is more than it will take much time to compute the gradient otherwise it is the best approach.
- A hypothesis of the gradient can be defined

$$h_w(x) = w_0 + w_1 x$$

- The parameters that need to be updated are w_0 and w_1 because they are unknown.
- The cost function can be defined as $J(w_0, w_1)$ in term of gradient descent it will be

$$= \frac{1}{2m} \sum_{i=1}^m (h_w(x_i) - y_i)^2$$

MINI BATCH GRADIENT

In mini batch gradient some of the training samples created a iteration which is called as a batch. Each batch further responsible to estimate the gradient.

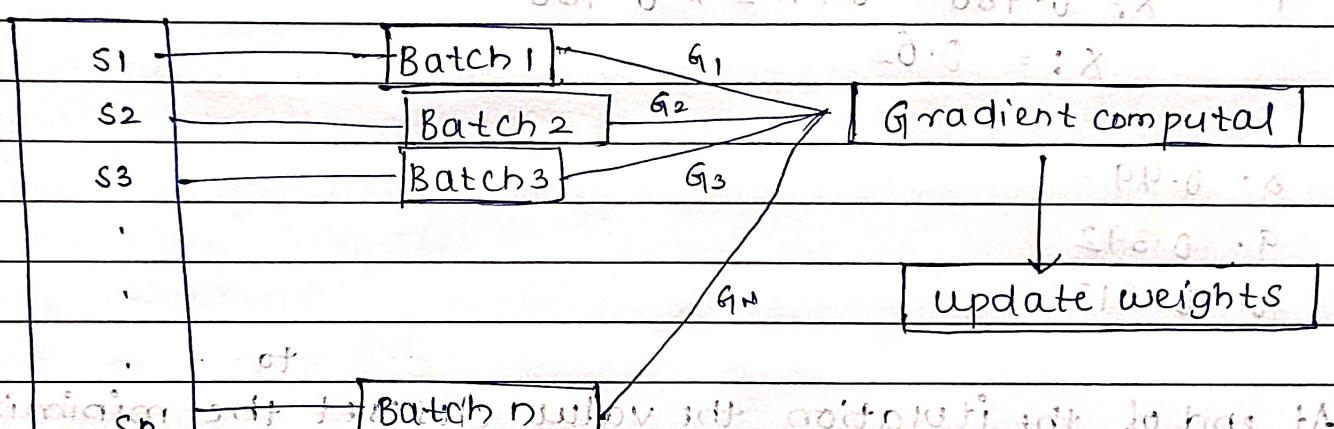
The last step of mini batch gradient is to compute the weights.

It updates the parameters using the small random subset of training data. It reduces the computational burden compared to batch gradient descent.

It combines the benefits of batch gradient descent and stochastic gradient descent.

* STOCHASTIC GRADIENT

In this gradientation each training sample prepared a batch and each batch is responsible to find a gradient.



$$x_i = \bar{x}_i \text{ of } f(x) \text{ with minimum old and new value}$$

$$\text{gradient function } \frac{\partial f}{\partial x} \text{ with respect to } x^i \text{ in gradient}$$

$$x := 3$$

$$\alpha = 0.1$$

$$f(x) = x^2$$

$$2.4 \quad 1.92$$

$$1.92 \quad 1.536$$

$$1.228 \quad 0.976$$

$$1. x := 3 - 0.1 \times 2x \quad 2. x := 2.4 - 0.1 \times 2 \times 2.4$$

$$3 - 0.1 \times 2.3$$

$$x = 1.92 \quad 1.92$$

$$x := 2.4 \text{ until not high loss function from repeat}$$

$$3. x := x - \alpha \times 2x$$

$$1.92 - 0.1 \times 2 \times 1.92$$

$$x = 1.536$$

$$4. x := 1.536 - 0.1 \times 2 \times 1.536$$

$$x := 1.228$$

$$5. x := 1.228 - 0.1 \times 2 \times 1.228$$

$$x := 0.976$$

$$6. x := 0.976 - 0.1 \times 2 \times 0.976$$

$$x := 0.786$$

$$7. x := 0.786 - 0.1 \times 2 \times 0.786$$

$$x := 0.62$$

$$8. 0.49$$

$$9. 0.392$$

$$10. 0.313$$

to

At end of the iteration the value closest to the minimization function $f(x) = 0$.

The entire procedure will terminate either the function will converge or the maximum iterations reached.

* Linear regression in terms of gradient descent approach.

In linear regression, our objective was to find the best fitted line.

Algorithm - Gradient descent for linear regression.

1. Start with the initial value for the slope m intercept b .

$$y = mx + b$$

m = slope

$m : 0$

b = intercept, y $b : 0$

2. Cost function

choose a cost function to measure the error between actual and predicted

for example in linear regression the mean square error is commonly used

20

$$MSE = \frac{1}{N} * \sum (y_i - (mx + b))^2$$

where N is the number of observations.

3. Gradient calculation

perform the iterations to update the parameters based on gradient.

the gradient indicates the direction of steepest ascent so we subtract a fraction of it from the parameters to descent towards the minimum.

Repeat until convergence or maximum number of iterations.

4. Partial derivative of the MSE

the partial derivative of the MSE with respect to slope (m)

is $\frac{\partial MSE}{\partial m}$ $\frac{\partial MSE}{\partial q}$

$$-\left(\frac{2}{n}\right) * \sum x(y - mx + b)$$

the partial derivative of MSE with respect to MSE

$$\frac{\partial MSE}{\partial b} = -\left(\frac{2}{n}\right) * \sum (y - (mx + b))$$

This equation represents the rate of change of MSE with respect to change in the intercept. It accounts for the contribution of each datapoint to the error without multiplication by the input feature since the intercept does not directly depend on x .

5. Prediction

after convergence use the learning parameter to make predictions for new input x predict y using $y = mx + b$.

let's consider

* Gradient descent with logistic regression

(i) Initialize the parameters

weight : 0

intercept : 0

(ii) Define the sigmoid function

$$f(x) = \frac{1}{1 + e^{-x}}$$

$$\text{sigmoid } z = f(z) = w_1x_1 + w_2x_2 + \dots + w_nx_n + b$$

(iii) Compute the predictions

For each data points

$$z = w_1x_1 + w_2x_2 + \dots + w_nx_n + b$$

$$\text{Sigmoid} \rightarrow \text{pred} = \text{Sigmoid}(z)$$

(iv) Compute gradients

$$\frac{\partial \text{cost}}{\partial w_1}, \frac{\partial \text{cost}}{\partial w_2}, \frac{\partial \text{cost}}{\partial b}$$

We need to find the gradient of the cost function (binary cross entropy) with respect to parameters weights using the following formula.

$$\frac{\partial \text{cost}}{\partial w_1} = \left(\frac{1}{m} \right) * \sum_{i=1}^N (y_{\text{pred}} - y_i) * x_i$$

$$\frac{\partial \text{cost}}{\partial b} = \left(\frac{1}{m} \right) * \sum_{i=1}^N (y_{\text{pred}} - y_i)$$

Here m is the number of data points

(v) update the parameters using the gradients and learning rate α

$$1. w_i - \text{new} = w_i - \alpha * \frac{\partial \text{cost}}{\partial w_i} \text{ for } i=1 \text{ to } N$$

$$2. b - \text{new} = b - \alpha * \frac{\partial \text{cost}}{\partial b}$$

(vi) Repeat iterations.

Repeat steps 3-5 until convergence or stopping criteria met.

(vii) Convergence can be determined based on changes in the cost function or parameter value.

• Compute the gradients

The gradient of the cost function with respect to the parameters weight and intercept can be written as

$$\frac{\partial \text{cost}}{\partial w_i} = \frac{1}{m} \left(\sum_{i=1}^m (\hat{y}_i - y_i) x_i \right)$$

$$\frac{\partial \text{cost}}{\partial b} = \frac{1}{m} \left(\sum_{i=1}^m (\hat{y}_i - y_i) \right)$$

The cost function / binary cross entropy function of logistic regression given by

$$\text{cost}(y, \hat{y}) = -\frac{1}{m} \sum_{i=1}^m \left[(y_i \log(\hat{y}_i)) + (1-y_i) \log(1-\hat{y}_i) \right]$$

- Q. Let us take a dataset having x and y . cost function / cross entropy and their values are given - $x = [1, 2]$ $y = [2, 4]$
- Apply gradient descent technique in linear regression cost function and find the optimal parameter and update it.

① Initialize : $b = 0$ // Intercept
 $m = 0$ // slope
 $\alpha = 0.1$ // step size

② calculate the predicted value :
 $y_{\text{predicted}} = [m + b \text{ for } i : 1 \text{ to } N]$

$$y_{\text{pred}} = [0x1 + 0, 0x2 + 0]$$

$$y_{\text{pred}} = [0, 0]$$

③ Compute the cost function.

$$\text{MSE} = \left(\frac{1}{2}\right) * \sum (y - y\text{-predicted})^2$$

$$= \frac{1}{2} * [(2-0)^2 + (4-0)^2]$$

$$= \frac{1}{2} * (4+16)$$

$$= 10$$

④ Calculate the gradients

In this step we required to calculate the partial derivative of the cost function of linear regression with required to slope is

$$\frac{\partial \text{MSE}}{\partial m} = -\left(\frac{2}{n}\right) * \sum x(y - (mx + b))$$

$$= -\left(\frac{2}{n}\right) * [1 * (2 - (0 * 10)) + 2 * (4 - (0 * 2 + 0))]$$

$$= -[1 * (2-0) * 2 * (4-0)]$$

$$= -[1 * 2 * 8]$$

$$\frac{\partial \text{MSE}}{\partial M} = -10$$

ΔM

$$\text{⑤ } \frac{\partial \text{MSE}}{\partial b} = \left(\frac{2}{n}\right) * \sum (y - (mx + b))$$

$$= \left(\frac{2}{2}\right) * [2 - (0 * 1 + 0) + 4 - (0 * 2 + 0)]$$

$$= -6$$

⑥ $x := x - \alpha * \frac{\partial \text{MSE}}{\partial b}$

calculate the updated parameter:

$$m_{\text{new}} = m - \alpha * \frac{\partial \text{MSE}}{\partial m}$$

$$= 0 - 0.1 * (-10)$$

$$m_{\text{new}} = 1$$

$$b_{\text{new}} = b - \alpha * \frac{\partial \text{MSE}}{\partial b}$$

$$= 0 - 0.1 * -6 = 0.6.$$

At the end of the first iteration m value changes to 1 and b value changes to 0.6

$$y = mx_1 + b$$

$$y = 1(1) + 0.6$$

$$y_1 = 1.6$$

$$y = mx_2 + b$$

$$y_2 = 1(2) + 0.6$$

$$y_2 = 2.6$$

② $\text{MSE} = \frac{1}{N} * \sum (y_i - (mx + b))^2$

$$= \frac{1}{2} * \left[(2 - 1.6)^2 + (4 - 1.6)^2 \right] = 1.06$$

- calculate the gradients

$$\frac{\partial \text{MSE}}{\partial m} = -\left(\frac{2}{n}\right) * \sum x(y - (mx + b))$$

$$= \left(-\frac{2}{2}\right) * \left[1(2 - (1 \times 1 + 0.6)) + 2(4 - (1 \times 2 + 0.6))\right]$$

$$\frac{\partial \text{MSE}}{\partial m} = -3.2$$

$$\frac{\partial \text{MSE}}{\partial b} = -\frac{2}{n} * \sum (y - (mx + b))$$

$$= \left(-\frac{2}{2}\right) * \left[(2 - (1 \times 1 + 0.6)) + 4 - (1 \times 2 + 0.6)\right]$$

$$\frac{\partial \text{MSE}}{\partial b} = -1.8$$

- update the gradients.

$$m_{\text{new}} = m - \alpha * \frac{\partial \text{MSE}}{\partial m}$$

$$= 1 - 0.1 * (-3.2) = 1.32$$

$$b_{\text{new}} = b - \alpha * \frac{\partial \text{MSE}}{\partial b}$$

$$= 0.6 - 0.1 * (-1.8) = 0.78$$

We can continue this process for more iterations until parameters converged to their optimal values or until a predefined number of iterations is reached.

$$\text{Q.2. } x = [1, 2, 3, 4, 5] \quad y = [2, 4, 5, 4, 5] \quad \alpha = 0.01$$

- Initialize.

$$m = 0$$

$$b = 0$$

$$\alpha = 0.01$$

② Calculate the predicted value.

$$\begin{aligned}y_{\text{pred}} &= mx + b \\&= [0 \cdot 1 + 0, \dots, 0 \cdot 5 + 0]\end{aligned}$$

$$[0, 0, 0, 0]$$

③ Cost function

$$\begin{aligned}\text{MSE} &= \frac{1}{5} \sum (y - y_{\text{pred}})^2 \\&= \frac{1}{5} [(2)^2 + (4)^2 + (5)^2 + (4)^2 + (5)^2] \\&= 17.2\end{aligned}$$

④ Calculate $\frac{\partial \text{MSE}}{\partial m}$ and $\frac{\partial \text{MSE}}{\partial b}$

$$\begin{aligned}\frac{\partial \text{MSE}}{\partial m} &= -2 * \frac{1}{n} \sum x(y - (mx + b)) \\&= -\frac{2}{5} * [1(2-0) + 2(4-0) + 3(5) + 4(4) + 5(5)] \\&= -\frac{2}{5} * 66\end{aligned}$$

$$\frac{\partial \text{MSE}}{\partial m} = -26.4$$

$$\frac{\partial \text{MSE}}{\partial b} = -2 * \frac{1}{n} \sum (y - (mx + b))$$

$$\frac{\partial \text{MSE}}{\partial b} = -\frac{2}{5} * 20 = -8$$

$$\frac{\partial \text{MSE}}{\partial b} = -8$$

$$m_{\text{new}} = m - \alpha \times \frac{\partial \text{MSE}}{\partial m}$$

$$= 0 - 0.01 \times (-26.4)$$

$$= 0.264$$

$$b_{\text{new}} = b - \alpha \frac{\partial \text{MSE}}{\partial b}$$

$$= 0 - 0.01 \times (-8)$$

$$= 0.08$$

2nd iterations.

1. Initialize $m = 0.264$ $b = 0.08$ $\alpha = 0.01$

② $y\text{-pred} = mx + b$

$$= ((0.264 \times 1 + 0.08), (0.264 \times 2 + 0.08),$$

$$(0.264 \times 3 + 0.08), (0.264 \times 4 + 0.08), \\ (0.264 \times 5 + 0.08))$$

$$= [0.344, 0.608, 0.872, 1.136, 1.4]$$

③ $MSE = \frac{1}{5} * \sum (y - y\text{-pred})^2$

$$= \frac{1}{5} * [(2 - 0.344)^2 + (4 - 0.608)^2 + (5 - 0.872)^2 + \\ (4 - 1.136)^2 + (5 - 1.4)^2]$$

$$MSE = 10.488$$

④ calculate the gradients.

$$\frac{\partial MSE}{\partial m} = \left(-\frac{2}{n}\right) * \sum x(y - (mx + b))$$

$$= \left(-\frac{2}{5}\right) * [1(2 - (0.264 \times 1 + 0.08)) + \\ 5 * (5 - (0.264 \times 5 + 0.08))]$$

$$\frac{\partial MSE}{\partial m} = -20.166$$

$$\frac{\partial MSE}{\partial b} = -\frac{2}{n} * \sum (y - (mx + b))$$

$$\frac{\partial MSE}{\partial b} = -6.25$$

⑤ update the gradients

$$m_{\text{new1}} = m - \alpha * \frac{\partial \text{MSE}}{\partial m}$$

$$= 0.264 - 0.01 * (-20.106)$$

$$= 0.465$$

$$b_{\text{new}} = b - \alpha * \frac{\partial \text{MSE}}{\partial b}$$

$$= 0.08 - 0.01 * (-6.25)$$

$$= 0.14$$

We can continue this process for more iterations until parameters converged to their optimal values or until a predefined number of iterations is reached.

* MINI BATCH GRADIENT

Let us consider the dataset contains one feature and target variable apply mini batch gradient descent to optimization linear regression cost functions.

$$\text{Q.1 } x : [1, 2, 3, 4, 5] \text{ // features}$$

$$y : [2, 4, 6, 8, 10] \text{ // Target}$$

steps to calculate the mini Batch

① Initialization

$$y = mx + b$$

$$m = 0$$

$$b = 0$$

② create a mini batch (For one mini batch)

$$x_{\text{batch}} : [1, 2]$$

$$y_{\text{batch}} : [2, 4]$$

(3) calculation) compute the prediction using correct parameters.

$$\begin{aligned}y_{\text{pred}} &= wx + b \\&= w \cdot x_{\text{batch}} + b\end{aligned}$$

(4) calculate the gradients of the cost function with respect to weights and bias

$$\frac{\partial \text{MSE}}{\partial w} = -\frac{2}{n} \sum_{i=1}^n (y_i - (wx_i + b)) x_i$$

$$\frac{\partial \text{MSE}}{\partial b} = -\frac{2}{n} \sum_{i=1}^n (y_i - (wx_i + b))$$

To calculate $\frac{\partial w}{\partial w}$ we have :-

$$\frac{\partial w}{\partial w} = \left(-\frac{2}{n}\right) * \text{sum}((y_{\text{batch}} - w * x_{\text{batch}} + b) * x_{\text{batch}})$$

$$\frac{\partial b}{\partial b} = \left(-\frac{2}{n}\right) * \text{sum}(y_{\text{batch}} - w * x_{\text{batch}} + b)$$

(5) update the parameters

$$w_{\text{new}} = w - \alpha * \frac{\partial w}{\partial w}$$

$$b_{\text{new}} = b - \alpha * \frac{\partial b}{\partial b}$$

$$\alpha = 0.01$$

(6) Repeat steps 2-4 until the convergence.

$$Q \cdot X: [1, 2, 3, 4, 5] \quad Y: [2, 4, 6, 8, 10]$$

$$\alpha = 0.01$$

(1) Initialization

$$w = 0$$

$$b = 0$$

$$x_{\text{batch}} = [1, 2] \quad y_{\text{batch}} = [2, 4]$$

$$\textcircled{2} \quad y_{\text{pred}} = w * x_{\text{batch}} + b$$

$$y_{\text{pred}} = 0 * 1 + 0, 0 * 2 + 0$$

$$[0, 0]$$

\textcircled{3} Compute.

$$\frac{\partial \text{MSE}}{\partial w} = -\frac{2}{N} \sum_{i=1}^N (y_i - (w x_i + b)) * x_i$$

$$= -\frac{2}{2} [1 * 2]$$

$$= -10$$

$$\frac{\partial \text{MSE}}{\partial b} = -\frac{2}{N} \sum_{i=1}^N (y_i - w x_i + b)$$

$$= -\frac{2}{2} \sum_{i=1}^2 [(2 - 0) + (4 - 0)]$$

$$\frac{\partial \text{MSE}}{\partial b} = -6$$

$$= 0.1$$

$$\frac{\partial \text{MSE}}{\partial b}$$

$$b_{\text{new}} = b - \alpha * \frac{\partial \text{MSE}}{\partial b} = 0.06$$

Select the mini batch 2

$$x : [3, 4]$$

$$\text{and } w_{\text{new}} = 0.1$$

$$y : [6, 8]$$

$$b_{\text{new}} = 0.06$$

\textcircled{1} Initialize.

$$w = 0.1$$

$$b = 0.06$$

$$x = [3, 4] \quad y = [6, 8]$$

$$\textcircled{2} \quad y_{\text{pred}} = w x + b$$

$$= [0.1 * 3 + 0.06, 0.1 * 4 + 0.06]$$

$$= [0.36, 0.46]$$

$$\textcircled{3} \quad \frac{\partial \text{MSE}}{\partial w} = -\frac{2}{N} \sum_{i=1}^N (y_i - (wx_i + b) * x_i)$$

$$\frac{\partial \text{MSE}}{\partial w} = -47.08$$

$$\frac{\partial \text{MSE}}{\partial b} = -\frac{2}{N} \sum_{i=1}^N (y_i - wx_i + b)$$

$$\frac{\partial \text{MSE}}{\partial b} = -13.18$$

∂b :

$$w_{\text{new}} = 0.1 - 0.01 * -47.08 = 0.5768$$

$$b_{\text{new}} = 0.06 - 0.01 * -13.18 = 0.1918$$

* Multiple Linear Regression with Gradient Descent Approach:-

- It involves predicting the target variables based on multiple input features.
- The main objective of this algorithm is to update the multiple parameters corresponding to each input features along with the intercept.
- Let us consider a dataset having two features along with the target variables our objective is to find multiple linear regression model to find the optimal values of the coefficient m_1 and m_2 and the intercept of linear regression.
- The equation we will use for multiple linear regression is

$$y = m_1x_1 + m_2x_2 + b \rightarrow ①$$

Algorithm

- ① Initialize the parameters.

Let's proceed the Gradient descent algorithm with initialization parameter.

Set

$$m_1 = 0$$

$$m_2 = 0$$

$$b = 0$$

$$\alpha = 0.01$$

- ② Find the cost function

The cost function of Multiple Linear regression is termed as MSE.

$$MSE = \left(\frac{1}{N} \right) * \sum (y - (m_1x_1 + m_2x_2 + b))^2$$

(3)

Gradient Descent Iteration

Perform the iterations to update the parameters m_1, m_2, b using the Gradient descent of the cost function with respect to each parameter.

$$(i) \frac{\partial \text{MSE}}{\partial m_1}$$

$$(ii) \frac{\partial \text{MSE}}{\partial m_2}$$

$$(iii) \frac{\partial \text{MSE}}{\partial b}$$

(4)

Update the parameters

The parameters will be updated by using the formula
 $x := x - \alpha \cdot \text{cost function}$

(5) prediction

After convergence we use the learning parameters to make prediction on new data.

$$y_{\text{pred}} = m_1 x_1 + m_2 x_2 + b$$

$$Q. x_1 = [1, 2]$$

$$x_2 = [2, 3]$$

$$y = [3, 4]$$

$$(1) m_1 = 0 \quad b = 0$$

$$m_2 = 0 \quad \alpha = 0.01$$

$$(2) y_{\text{pred}} = [0, 0]$$

$$(3) \text{MSE} = \left(\frac{1}{N} \right) * \sum (y - (m_1 x_1 + m_2 x_2 + b))^2$$

$$= \frac{1}{2} * [(3-0)^2 + (4-0)^2]$$

$$\text{MSE} = 12.5$$

$$(4) \frac{\partial \text{MSE}}{\partial m_1} = \left(\frac{-2}{N} \right) * \sum (x_1 (y - (m_1 x_1 + m_2 x_2 + b)))$$

$$\frac{\partial \text{MSE}}{\partial m_2} = -\frac{2}{N} * \sum (x_2(y - (m_1x_1 + m_2x_2 + b)))$$

$$\frac{\partial \text{MSE}}{\partial b} = -\frac{2}{N} * \sum (y - (m_1x_1 + m_2x_2 + b))$$

$$\frac{\partial \text{MSE}}{\partial m_1} = -\frac{2}{2} * [(1 * (3-0)) + (2 * (4-0))]$$

$$\text{answ} = -11 \quad \text{from the equation of line}$$

$$\frac{\partial \text{MSE}}{\partial m_2} = -\frac{2}{2} * [(2 * (3-0)) + (3 * (4-0))]$$

$$\text{answ} = -18 \quad \text{from the equation of line}$$

$$\frac{\partial \text{MSE}}{\partial b} = -2 * \frac{(3+4)}{2} = -7$$

$$m_1\text{-new} = m_1 - \alpha \frac{\partial \text{MSE}}{\partial m_1}$$

$$= 0 - 0.01 * -11$$

$$= 0.11$$

$$m_2\text{-new} = m_2 - \alpha \frac{\partial \text{MSE}}{\partial m_2}$$

$$= 0 - 0.01 * -18$$

$$= 0.18$$

$$b\text{-new} = b - \alpha \frac{\partial \text{MSE}}{\partial b}$$

$$= 0 - 0.01 * (-7)$$

$$= 0.07$$

The new updated parameters are.

- ① Initialize $m_1 = 0.11$ $b = 0.07$
- $m_2 = 0.18$

$$\begin{aligned}
 (2) \quad y_{\text{pred}} &= m_1 x_1 + m_2 x_2 + b \\
 &= [(0.11 \times 1 + 0.18 \times 2 + 0.07), (0.11 \times 2 + 0.18 \times 3 + 0.07)] \\
 &= [0.54, 0.83]
 \end{aligned}$$

$$\begin{aligned}
 (3) \quad \text{MSE} &= \frac{1}{N} * \sum (y - (m_1 x_1 + m_2 x_2 + b))^2 \\
 &= \frac{1}{2} * ((3 - 0.54)^2 + (4 - 0.83)^2) \\
 &= \frac{1}{2} (6.05 + 10.04) = 8.045
 \end{aligned}$$

$$\begin{aligned}
 (4) \quad \frac{\partial \text{MSE}}{\partial m_1} &= -2 * \frac{1}{N} \sum x_1 (y - (m_1 x_1 + m_2 x_2 + b)) \\
 &\stackrel{\text{as } \frac{-2}{2} * (1 * ((3 - 0.54) + 2 * (4 - 0.83)))}{=} -1 * 8.8 \\
 &= -8.8
 \end{aligned}$$

$$\begin{aligned}
 \frac{\partial \text{MSE}}{\partial m_2} &= -2 * \frac{1}{N} \sum x_2 (y - (m_1 x_1 + m_2 x_2 + b)) \\
 &= -2 * \frac{1}{2} [2 + (3 - 0.54) + 3 * (4 - 0.83)] \\
 &= -14.43 \\
 \frac{\partial \text{MSE}}{\partial b} &= -2 * \frac{1}{N} * [y - (m_1 x_1 + m_2 x_2 + b)] \\
 &= -2 * \frac{1}{2} * [(3 - 0.54) + (4 - 0.83)] \\
 &= -5.63
 \end{aligned}$$

$$\begin{aligned}
 m_1 - \text{new} &= m_1 - \alpha * \frac{\partial \text{MSE}}{\partial m_1} \\
 &= 0.11 - 0.01 * (-8.8) = 0.19
 \end{aligned}$$

$$m_2\text{-new} = m_2 - \alpha \cdot \frac{\partial \text{MSE}}{\partial m_2}$$

$$= 0.18 - 0.01 \times (-14.43)$$

$$= 0.32$$

$$b\text{-new} = b - \alpha \cdot \frac{\partial \text{MSE}}{\partial b}$$

$$= 0.07 - 0.01 \times (-5.63)$$

$$= 0.12$$

* Non-parametric Test :-

• Spearman correlation coefficient Rank :-

It is a non-parametric test which is used to calculate the value of $r(r)$.

The spearman Rank correlation can be written by the formula

$$\rho_R = 1 - \frac{6 \sum D^2}{N(N^2 - 1)}$$

where N is the number of observations

D is the difference b/w Rank x and Rank y .

The values of ρ_R lies between $[-1, 1]$

If $\rho_R = 1 \rightarrow$ Perfect positive correlation

$\rho_R = -1 \rightarrow$ perfect Negative correlation

$\rho_R = 0 \rightarrow$ No correlation

If $\rho_R = 0.08 \rightarrow$ Moderate correlation

Q. Let us take a dataset. Apply spearman rank correlation coefficient. Interpret the results.

x	y	R_x	R_y	D	D^2
10	7	3.5	3.5	0	0
8	4	5.5	6	-0.5	0.25
12	6	2.5	5	-3	9
15	7	1	3.5	-2.5	6.25
8	9	5.5	11	4.5	20.25
10	8	3.5	10.5	1.5	2.25

$$D = R_x - R_y$$

$$\rho_R = 1 - \frac{6 \sum D^2}{N(N^2 - 1)}$$

$$1 - \frac{6 \times 38}{6(36-1)}$$

$$= -0.085.$$

* Wilcoxon-Rank sum Test :-

• Wilcoxon Rank sum Test also known as Mann-Whitney U test

• It is a non-parametric test used to access whether two independent samples come from the same distribution

• The main objective of this test is to compare the medians of two independent groups, it is also called as alternative to two-sample t-test.

① Formulate the Hypothesis

Define the null and alternative hypothesis

H_0 : Null hypothesis

H_1 : Alternative hypothesis

The Null hypothesis tells that there is no difference b/w distribution of two groups.

The Alternative hypothesis tells that there is a difference between distribution of two groups.

The Alternative hypothesis tells that there is a difference between distribution of two groups.

② Rank the Data

Combine the data from both the groups.

Rank them in the ascending order.

Assign the ranks with the smallest value.

If there is same value then assign the average rank.

③ Calculate the sum of Rank

Calculate the sum of Rank x (T_x)

Calculate the sum of Rank y (T_y)

④ Calculate the statistic U prior to this calculate

$$U_x = T_x - \frac{n_1(n_1+1)}{2}$$

$$U_y = T_y - \frac{n_2(n_2+1)}{2}$$

$$U = \min(U_x, U_y)$$

⑤ Determine the critical value

Look the wilcoxon rank sum table where the critical value depends on sample sizes and level of significance.

For example - $n_1 = n_2 = 3$

$\alpha = 0.01$, critical value = 9

$$U = 7.8$$

- ⑥ Compare the test statistic to the normal value
 If U is less than or equal to critical value then reject the Null hypothesis else accept the alternative hypothesis.

⑦ Interpretation

If the Null hypothesis is rejected then it suggests that there is a significance between distribution of two groups.

If the Null hypothesis is not rejected then it indicates that there is not enough evidence to conclude a significant difference.

- Q. Let us consider the two different groups and their values are available apply wilcoxon Rank sum test and interpret it.

X 60 45 23 32

Y 10 25 20 54 32 65 8.

① Hypothesis

$$H_0 = \text{Both are same}$$

$$H_1 = \text{Both are not same}$$

H_0 : Both are not same

H_1 : Both are same

Level	Y	Y	Y	X	Y	X	Y	X	Y	X	Y
Data	8	10	20	23	25	32	32	45	54	60	65
Rank	1	2	3	4	5	6.5	6.5	8	9	10	11

$$T_x = 4 + 6.5 + 8 + 10 = 28.5$$

$$T_y = 1 + 2 + 3 + 5 + 6.5 + 9 + 11 = 37.5$$

$$U_x = 28.5 - \frac{\frac{2}{4}(5)}{7}$$

$$U_x = 18.5$$

$$U_y = 37.5 - \frac{7(8)4}{7}$$

$$U_y = 9.5$$

Cross Verification.

Test

$$1) U_x + U_y = n_1 n_2$$

$$18.5 + 9.5 = 4 \times 7$$

$$28.5 = 28$$

$$2) T_x + T_y = (n_1 + n_2)(n_1 + n_2 + 1)$$

$$28.5 + 37.5 = (4+7)(4+7+1)$$

$$66 = 11 \times 10 \frac{5}{2}$$

$$66 \neq 55$$

$$U_{\min} = (U_x, U_y)$$

$$U = (18.5, 9.5)$$

$$U_{\min} = 9.5$$

$$Q. A = 78, 82, 85, 88, 90$$

$$B = 72, 75, 80, 83, 86$$

critical value = 8

$$\alpha = 0.05$$

$$\underline{\text{Sol}}: \quad ① \quad n_1 = 5$$

$$n_2 = 5$$

$$\alpha = 0.05$$

$$n_1 = n_2$$

H_0 : Both are not same

H_1 : Both are same

Level	B	B	A	B	A	B	A	B	A	A
Data	72	75	78	80	82	83	85	86	88	90
Rank	1	2	3	4	5	6	7	8	9	10

$$T_x = 3 + 5 + 7 + 9 + 10 = 34$$

$$T_y = 1 + 2 + 4 + 6 + 8 = 21$$

$$U_x = \frac{T_x - n_1(n_1+1)}{2}$$

$$= 34 - \frac{5(5+1)}{2}$$

$$= 34 - \frac{5(6)}{2}$$

$$= 34 - 15$$

$$U_x = 19$$

$$U_y = \frac{T_y - n_2(n_2+1)}{2}$$

$$= 21 - \frac{5(6)}{2}$$

$$= 21 - 15$$

$$U_y = 6$$

$$U_{\min}(U_x, U_y)$$

$$U_{\min} = 6.$$

$$6 < 8$$

Reject the Null hypothesis

Mann Whitney Test for large sample.
 When your the sample size is more i.e. $n_1 \geq 10$ or $n_2 \geq 10$
 then we called it as a Mann whitney Test for large sample.

• Steps

- ① Test the hypothesis
- ② compute U_x and U_y
- ③ compute Test statistics

where

$$z = \frac{U - U_0}{\sigma}$$

$$\text{where } U = \min(U_x, U_y)$$

$$U_0 = \frac{(n_1)(n_2)}{2}$$

$$\sigma = \sqrt{\frac{n_1 n_2 (n_1 + n_2 + 1)}{12}}$$

- ④ compare the test statistics to the critical value

If compared $z >$ critical value at level of given significance,

We reject H_0 otherwise we fail to reject $[H_0]$

- Q. The below data below show the salaries at randomly selected advertisement in two different occupation [education and Marketing]

<u>education</u>	<u>Marketing</u>
22	28
40	40
18	20
25	45
15	50
23	39
16	26
19	55
21	48
30	41
42	

Use mann whitney test at 1% level of significance at the median check the salary in the education is lower than the median salary in the marketing.

① Test the Hypothesis

$$n_1 = 10$$

$$n_2 = 11$$

$$H_0 = \text{Both are not same} \checkmark$$

$$H_1 = \text{Both are same}$$

Label	E	E	E	E	M	E	(E)	E	N	M	E
Data	15	16	18	19	20	21	22	23	25	26	28
Rank	1	2	3	4	5	6	7	8	9	10	11

②

$$\text{with } E_x = 66.5 \text{ and } E_y = 116.5$$

$$\text{as } E_x < E_y \rightarrow 164.5 - 66.5 = 98 > 2.33$$

$$U_x = E_x - n_1(n_1 + 1)$$

$$U_x = 66.5 - 55$$

$$U_x = 11.5$$

$$U_y = E_y - n_2(n_2 + 1)$$

$$U_y = 165 - 66.5$$

$$U_y = 99$$

$$U = \min(U_x, U_y)$$

$$U = 11.5$$

(3)

$$U_0 = \frac{10(11)}{2} = 55$$

$$Z = \frac{11.5 - 55}{\sqrt{14.20}} = \frac{11.5 - 55}{14.20} = -3.063.$$

$$\sigma = \sqrt{\frac{10(10)(10+11+1)}{12}} = \sqrt{14.20} = 3.77$$

M	M	E	M	M	M	M	M	M	M
39	40	40	41	42	45	48	50	55	
13	14.5	14.5	16	17	18	19	20	21	

(4) The critical value at 1% level of significance is

$-3.09 < -2.33$ so we reject the

Therefore we conclude that

- Q. A 42 20 51 39 57 60 23 24 28 30
 B 30 42 25 29 35

① Test the Hypothesis

$$n_1 = 10 \quad n_2 = 5$$

Label	A	A	A	B	A	B	A	B	B	A	A	B	A	A
Data	20	23	24	25	28	29	30	30	35	39	42	42	51	52
Rank	1	2	3	4	5	6	7	7	9	10	11	11.5	11.5	14

$$T_x = 80 \quad T_y = 38.$$

$$U_x = 80 - \frac{10(11)}{2} = 80 - 55 = 25$$

$$U_y = 38 - \frac{10 \cdot 5(6)^3}{2} = 38 - 150 = -112$$

$$U_y = 23$$

$$U = \min(U_x, U_y) = 23$$

$$\boxed{U=23}$$

$$U_0 = \frac{5 \cdot 10(5)}{2} = 25$$

$$\sigma = \sqrt{\frac{10(5)(10+5+1)}{12}} = \sqrt{\frac{8 \cdot 164}{12}} = \sqrt{136} = 11.66$$

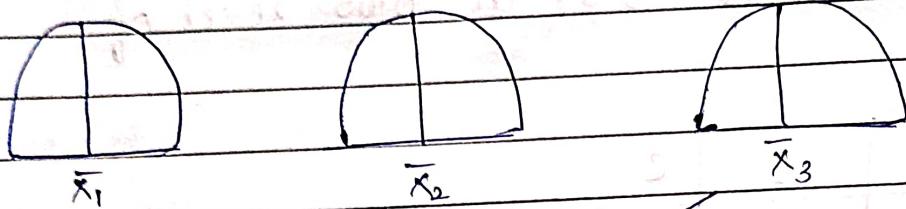
$$Z = \frac{23-25}{8.164} = -0.244$$

$$\alpha = -5.24$$

* Analysis of Variance (ANOVA)

Variance - It is defined as expectation of the square deviation of a random variable of their mean.

- For comparison of more than 2 populations or population having more than two subgroups of the same population we use the same technique?



Formula for ANOVA is

$$\text{ANOVA} = \frac{\text{Variability between the mean (x)}}{\text{Variability within the mean (y)}}$$

$$\text{Total Variance} = x + y$$

Assumption :-

- Each population is having normal distribution.
- The population from which the samples are drawn have the equal variance.
- Each sample is drawn randomly and they are independent.

$$s_1^2 = s_2^2 = \dots = s_k^2$$

$$\text{Null Hypothesis: } \mu_1^2 = \mu_2^2$$

$$\text{Alternate: } \mu_1^2 \neq \mu_2^2$$

Q. To access the significance of possible variation in performance in a certain test between the convert School of a city. A common entrance test is given to a no. of students taken randomly from 5th class of three different schools. Their concerned results are given below:- Perform the 1 way ANOVA test, the critical value of distribution is 3.89 at 0.05 level of significance.

A	B	C
9	13	14
11	12	13
13	10	17
9	15	7
8	5	9

$$\bar{x}_A = \frac{50}{5} = 10 \quad \bar{x}_B = \frac{55}{5} = 11 \quad \bar{x}_C = \frac{60}{5} = 12$$

- ① calculate the mean $\bar{x}_A, \bar{x}_B, \bar{x}_C$
- ② $\bar{x} = \bar{x}_A + \bar{x}_B + \bar{x}_C = 10 + 11 + 12 = 33$
- ③ calculate sse (variation b/w the sample)

$$sse = \sum \left[(\bar{x}_A - \bar{x})^2 + (\bar{x}_B - \bar{x})^2 + (\bar{x}_C - \bar{x})^2 \right]$$

$\bar{x}_A - \bar{x}$	$(\bar{x}_A - \bar{x})^2$	$\bar{x}_B - \bar{x}$	$(\bar{x}_B - \bar{x})^2$	$\bar{x}_C - \bar{x}$	$(\bar{x}_C - \bar{x})^2$
-1	1	0	0	1	1
-1	1	0	0	1	1
-1	1	0	0	1	1
-1	1	0	0	1	1
-1	1	0	0	1	1

$$SSC = 5 + 0 + 5 = 10$$

So the value of $SSC = 10$

④ Degree of freedom

$$V = C - 1$$

c is the number of columns.

$$V = 2$$

⑤ calculate the Mean square

$$\text{Mean square} = \frac{SSC}{DF(v)} = \frac{10}{2} = 5$$

⑥ Now calculate SSE (sum of square within the samples)

$$SSE = (A - \bar{x}_A)^2 + (B - \bar{x}_B)^2 + (C - \bar{x}_C)^2$$

$A - \bar{x}_A$	$(A - \bar{x}_A)^2$	$B - \bar{x}_B$	$(B - \bar{x}_B)^2$	$C - \bar{x}_C$	$(C - \bar{x}_C)^2$
-1	1	2	4	2	4
1	1	1	1	1	1
3	9	-1	1	5	25
-1	1	4	16	-5	25
-2	4	-6	36	-3	9
	16		58		64
				SSF = 120	

6 (7) Degree of freedom

$$V_2 = n - c \\ = 5 - 2$$

$$V_2 = 3.$$

$$V_2 = n - c \\ = 5 - 2$$

$$V_2 = 3.$$

$$V_2 = n - c \\ = (5 \times 3) - 3$$

$$\boxed{V_2 = 12}$$

(8)

$$\text{Mean square} = \frac{\text{SSE}}{V_2} = \frac{138}{12} = 11.5$$

(9)

$$F = \frac{5}{11.5} = 0.434$$

The calculated value of F value = 0.435 and
F distribution tabular value at $\alpha = 0.05$ = 3.89

(10)

Compare these two values so our tabulated value is greater than 0.435 Hence Null hypothesis is passed and no significant variation in the schools marks sheet.

Q. Assign 3 additional groups to rows) 322 Statistics with (a)

$$A \quad B \quad C \quad \text{Critical value}$$

$$2 \quad 3 \quad 4 \quad 5 \quad 6 \quad 7 \quad 8 \quad 9$$

$$4 \quad 5 \quad 6 \quad 7 \quad 8 \quad 9$$

$$(1) \bar{x}_A = \frac{12}{3} = 4$$

$$\bar{x}_B = \frac{15}{3} = 5$$

$$\bar{x}_C = \frac{18}{3} = 6$$

(2)

$$\bar{x} = \frac{4+5+6}{3} = \frac{15}{3} = 5$$

(3)

$$SSC = \sum [(x_A - \bar{x})^2 + (x_B - \bar{x})^2 + (x_C - \bar{x})^2]$$

$(\bar{x}_A - \bar{x})$	$(\bar{x}_A - \bar{x})^2$	$(\bar{x}_B - \bar{x})$	$(\bar{x}_B - \bar{x})^2$	$(\bar{x}_C - \bar{x})$	$(\bar{x}_C - \bar{x})^2$
-1	1	0	0	1	1
+1	1	0	0	1	1
-1	1	0	0	1	1
	<u>3</u>	<u>0</u>	<u>0</u>	<u>3</u>	<u>3</u>

$$SSC = 3 + 0 + 3 = 6.$$

(4)

$$V = C - 1 \text{, where } C \text{ is the number of observations}$$

$$V = 3 - 1 = 2 \text{, degrees of freedom}$$

(5)

$$\text{Mean square} = \frac{SSC}{DF(V)} = \frac{6}{2} = 3$$

(6)

$$SSE = \sum (A - \bar{x}_A)^2 + (B - \bar{x}_B)^2 + (C - \bar{x}_C)^2$$

P	S	Σ	ΣA	ΣB	ΣC	ΣD
-2	P	4	-2	4	-2	4
0	S	0	0	0	0	0
2	<u>A</u>	4	<u>2</u>	<u>4</u>	<u>2</u>	<u>4</u>
	<u>8</u>	<u>8</u>	<u>8</u>	<u>8</u>	<u>8</u>	<u>8</u>

$$SSE = 8 + 8 + 8 = 24$$

(7)

$$V_2 = (n-1) = 3 \text{, degrees of freedom}$$

$$V_2 = 6.$$

(8)

$$\text{Mean square} = \frac{SSE}{V_2} = \frac{24}{6} = 4.$$

(9)

$$f = \frac{3}{4} = 0.75.$$

* ANN (Artificial Neural Network)

* Performance Model of Regression and classification.

In regression there are different approaches available which exactly tells how is your model.

- RMSE - Root Mean Square Error
- MSE - Mean square error
- MAE - Mean Absolute error
- RMSLE - Root mean square log error.

① RMSE =

$$\sqrt{\frac{\sum (y_i - \hat{y}_i)^2}{N}}$$

Let us take one example, find RMSE of the above dataset using LR

BA	TA	(y_i)	Predicted TA	y	MEASURE
500	50	52	52	2	4
550	48	57	57	-9	81
600	60	62	62	-2	4
700	65	72	72	-7	49
800	75	82	82	-7	49
		0	0	0	0
		187	187		

$$RMSE = \sqrt{\frac{187}{85}} = 6.115$$

RMSE ↑ model is biased

RMSE ↓ accuracy ↑

Based on the bill amount we have to identify the trip amount the training data tax as a model and identify the function.

② MSE - It is one of the most preferred matrix for the regression task.

It is simply the average of square difference between actual and predicted value. As it squares the differences it penalizes even a small error which lead to over estimate how bad the model is.

The formula for MSE is

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

Let us take a dataset which consist of dependent and independent variable find performance of a model using MSE.

HT	WT (y_i)	\hat{y}	E	E^2
160	72	74.8448	-2.8448	+8.0928
171	76	77.1328	-1.1328	1.2832
182	77	79.4208	-2.4208	5.8602
180	83	79.0048	3.9952	15.9616
154	76	73.5968	2.4032	5.7753

The regressor line is $\hat{y} = 41.5648 + 0.208x_i$ 36.9731

$$\text{MSE} = \frac{36.9731}{5} = 7.3946.$$

MSE is always positive and the value closer to 0 or lower value is better.

If error is zero then model is Overfitted.

$$(3) \text{MAE} = \frac{1}{N} \left\{ \sum_{i=1}^N |y_i - \hat{y}| \right\}$$

N = Number of obs.
 y_i = actual
 \hat{y} = predicted

It is the absolute difference b/w the actual and the predicted value MAE is the most robust to outliers and does not penalize the errors as extremely as MSE. It can be used

Q. Find out the models performance by using MAE of the following dataset.

Act	predicted	Error	Error
100	130	-30	30
150	170	-20	20
200	220	-20	20
250	260	-10	10
300	325	-25	25

$$MAE = \frac{105}{5} = 21$$

(4) RMSLE (Root Mean square log error) - It is calculated at logarithmic scale while calculating

UNIT - 4

M T W T F S S

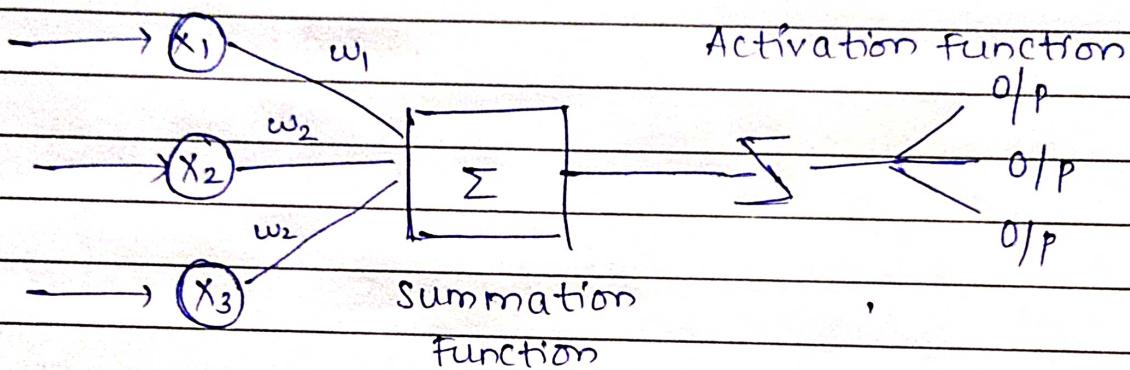
Page No.:

Date:

- ANN - Artificial Neural Network.

The ANN stands for Artificial Neural Network which is part of ML. The Human brain is a complex structure and it is interconnected network of simple processing elements called neurons.

The behaviour of neuron can be captured by using the following model.



The total input can be calculated as

$$Y_{in} = w_1x_1 + w_2x_2 + w_3x_3$$

$$Y_{in} = \sum_{i=1}^n x_i w_i$$

Where x_1, x_2, x_3 are the inputs

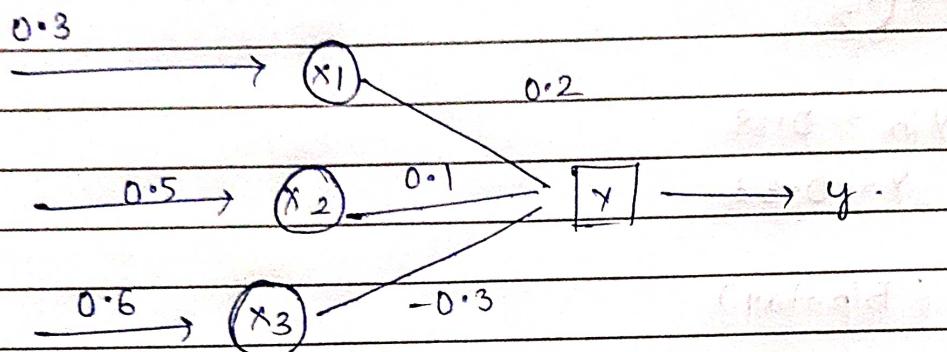
w_1, w_2, w_3 are the weights

Input signals are verified through weights to generate the final output by y .

The sum is passed to the activation function which is known as Non-linear activation function or filter or squash function.

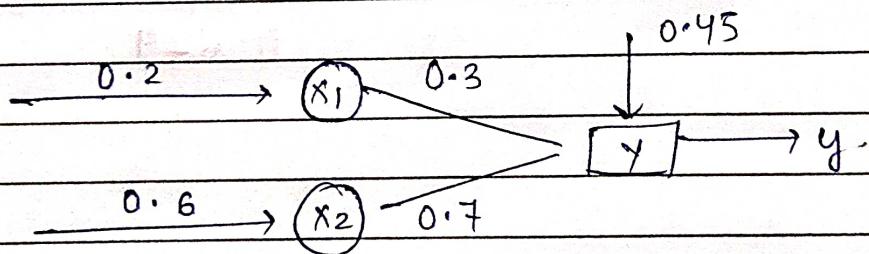
A very commonly used activation function is the thresholding function compared with a threshold value ϕ if the value of $\phi > 0$ then the o/p is 1, Otherwise 0.

Q. From the given model calculate the total input and weight.



$$Y_{in} = \sum_{i=1}^N x_i w_i$$

$$Y_{in} = -0.07$$



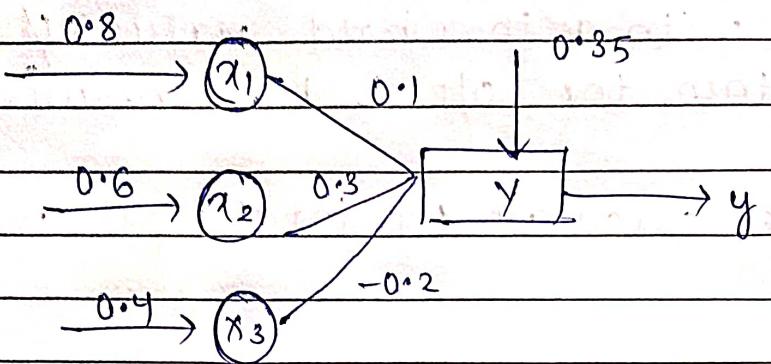
$$Y_{in} = 0.48$$

$$y = 0.48 + 0.45 = 0.93$$

Q.

Obtain the o/p of the neuron y for the network shown in figure using activation function

- ① binary sigmoidal
- ② Bipolar Sigmoidal.



$$y_{in} = 0.18$$

$$y = 0.53$$

(2) sigmoidal (Bipolar)

$$Y = f(M) = \frac{2}{1 + e^{-y_{in}}} - 1 = \frac{2}{1 + e^{-0.53}} - 1 = 1.258 - 1 = 0.259$$

$$\textcircled{1} \quad y = f(Y_M) = \frac{1}{1 + e^{-y_{in}}} = \frac{1}{1 + e^{-0.53}} = 0.629$$

(binary)