

Machine Learning !!

Definition !!

"Study of Computer Algorithms that allows the Computer Program to automatically improve through experience" [Efficient - learning]
[Training]

By Tom Mitchell (founder of ML Department)
School of Computer Science at Carnegie Mellon University.

* "ML is teaching the machine about something"

How !!

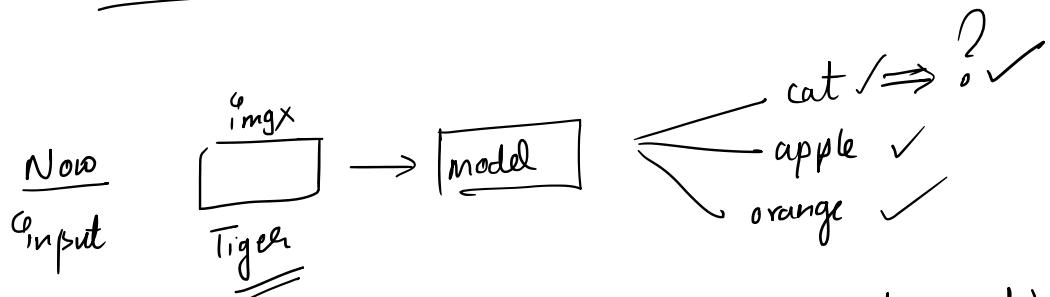
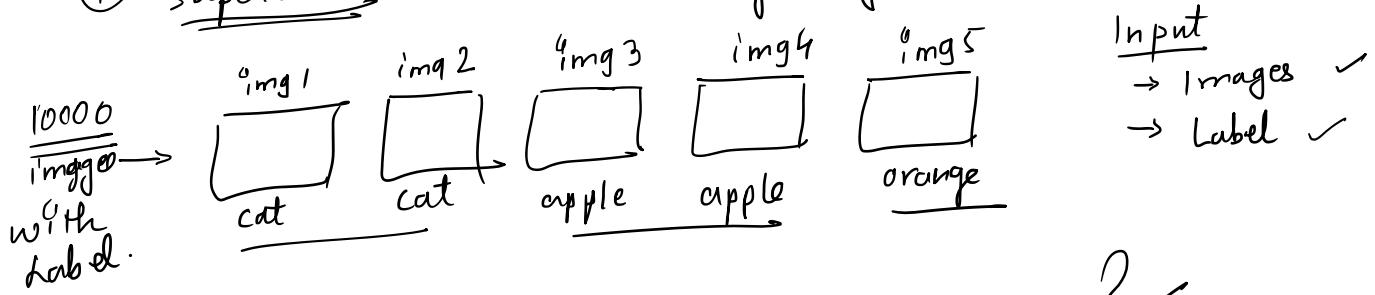
- ① Collect and clean the data
- ② Algorithm (model)
 - ↳ Selected (Readymade)
 - ↳ Built.
- ③ Teach the model essential pattern from data
(Training).
- ④ Export the model to give helpful answer.

Ex To Design a System that determine from MRI Scan, whether Tumor is present or not.

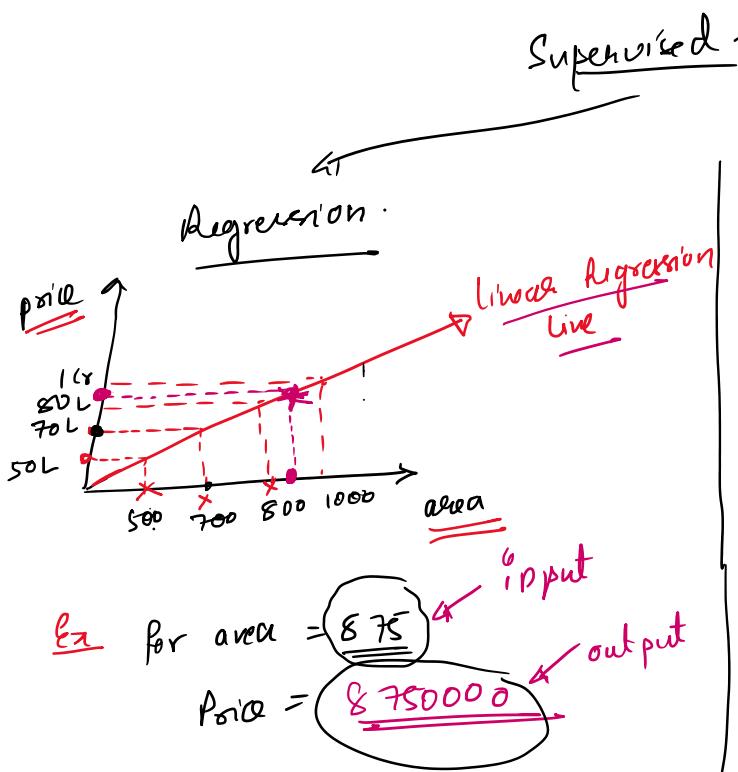
- ① Collect large No of MRI reports
Ex 10000 MRI report
6000 has Tumor
4000 do not have Tumor
Labeled Data
- ② Build an efficient algorithm that detects presence or absence of Tumor in an MRI Scan.
[Expert Consultation - Radiologist].
- ③ To the Algorithm feed the 1000 MRI scans and allow the model to learn (train).
- ✓ ④ Use around 3000 Images (MRI scans) for testing
- ✓ ⑤ Use this model to determine presence or absence of Tumor from a New Image (MRI Scan)

Types of machine learning Algorithms.

① Supervised Machine learning Algorithms .



- * With lot of images as input, model will be able to identify pattern and will be able to predict.



These off is not continuous value
but could be Boolean or
Some class / category as output.

* Span Selection

→ MRI Scan
For Tumor Detection

It also carries the slip belongs to

area \Rightarrow Independent variable

price \Rightarrow Dependent variable

Here the dependent and independent Variable can have continuous value.

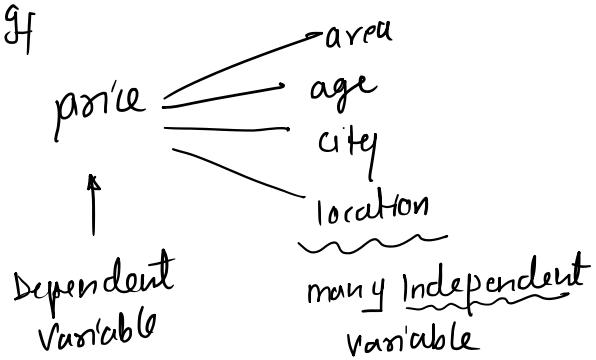
Types \Rightarrow

If price \rightarrow area.

Single Independent variable

\rightarrow Simple Regression

If

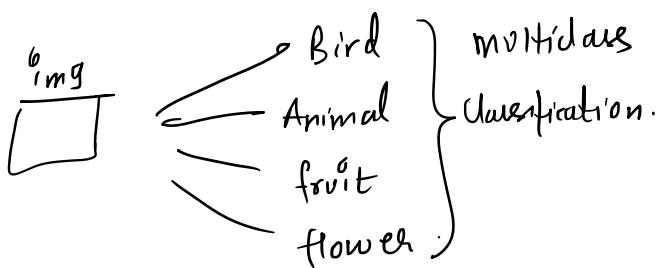


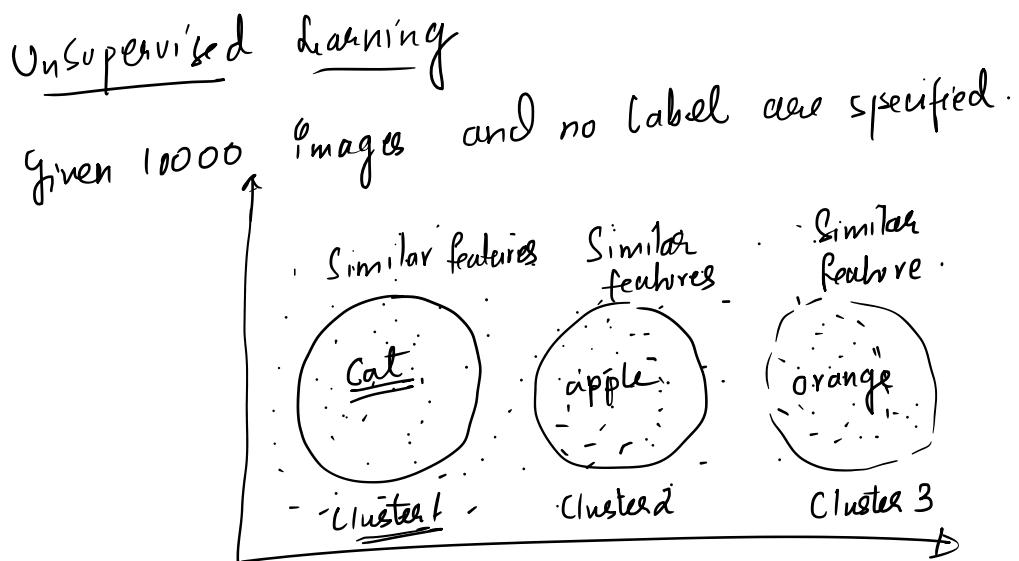
When we have more than one Independent Variable \Rightarrow multiple regression

In above cases the op belongs to one of the two classes.

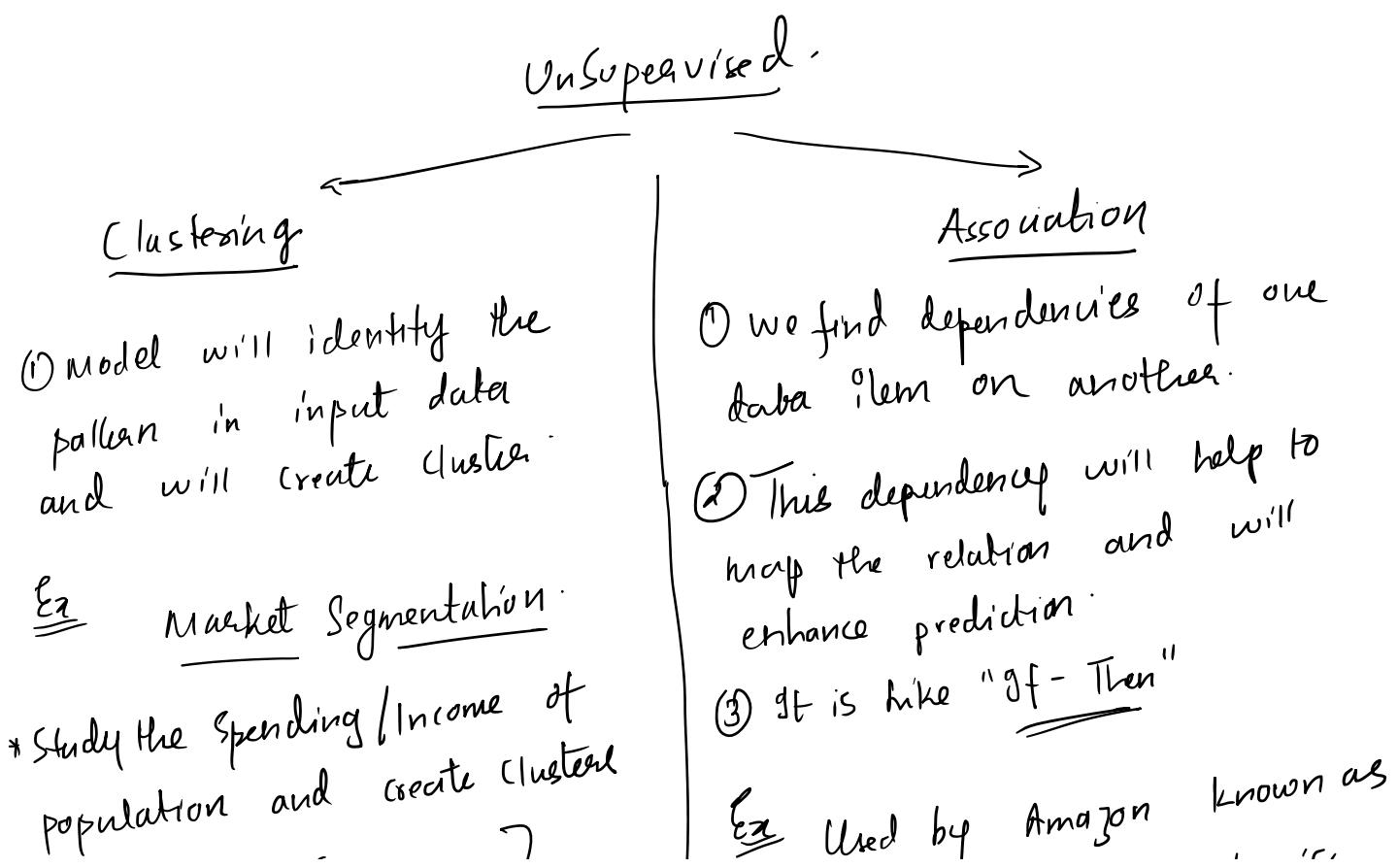
\rightarrow Binary classification.

If op can be one of many classes (more than two) Then it is called as multiclass classification.





- * Here the input is only data and no labels are associated with data.
- * Not sure about type of output
- * Unsupervised Algo will work on 10000 images and will create clusters of images based on the similarities.
- * It is "SELF LEARNING"



population and create illusion

- (1) High Earning }
(2) Medium Earning }
(3) Low Earning }

Ex Used by Amazon known as
Market Basket Analysis'

"if a person purchases cellphone
than the person has tendency
to purchase screen guard &
backcover"

Semi Supervised

Text Document Classification

Eg

1000000 articles and need to classify them
into
News
Literature
Research Paper
Medical Report

Label is not provided.

Manually labelling the 1000000 articles not possible.

→ We will label 10000 articles [Supervised Learning]
→ My model will be trained on these 10000 articles and
will use the pattern identified to classify the remaining
990000 articles [Unsupervised]

- * Uses small amount of labelled data
and large amount of unlabelled data.
- * Benefits of Both labelled and unlabelled data.
- * Overcome the challenge of finding large amount of
labelled data.

Reinforcement Learning [Experiential learning]

- Here the agent/model learns how to behave in an Environment by performing action and experiencing the result.

Type

Episodic learning

- Here we have start and end state, thus an episode is created.
 - Thus the further action is based on feedback of result of earlier action.
- Ex Fear of dog after being bitten is Episodic learning.

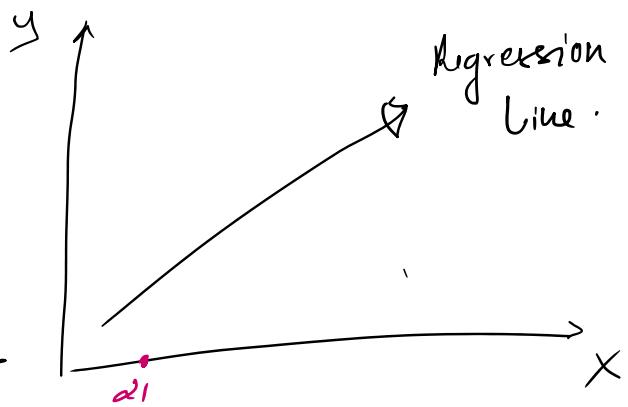
Continuous Task learning

- * There is no terminal state here
- * Agent/model that does Automated Stock Trading goes for Continuous learning

SR.NO	INDEPENDENT X	DEPENDENT Y	Actual value - $n=6$	
			$\sum x^2$	$\sum xy$
1	43	99	1849	4257
2	21	65	441	1365
3	25	79	625	1975
4	42	75	1764	3150
5	57	87	3249	4959
6	59	81	3481	4379

$$\sum x = 247 \quad \sum y = 466 \quad \sum xy = 11409 \quad \sum x^2 = 20485$$

NOW predict is $X = 55$, what is $\underline{\underline{y}}$?



$$y = a + bx + e$$

↑ ↑
 Dependent Intercept
 Variable } decides impact of x on y .

a = intercept

b = slope (coefficient of Independent Variable)

e = error

[if $x \uparrow y \uparrow$ +ve Impact]

[if $x \uparrow y \downarrow$ -ve Impact].

Let us assume $e = 0$

$$y = a + bx \quad \text{--- (1)}$$

To find a and b .

Take \sum on both side of 1

$$\sum y = \sum a + \sum bx$$

$$\sum y = a \sum 1 + b \sum x$$

$$\sum y = an + b \sum x \quad \text{--- (2)}$$

Multiply Eq 2 with Independent Variable X .

$$\sum \dots \sum x^2 \quad \text{--- (3)}$$

Multiply Eq 2 with Σx^2

$$\underline{\underline{\Sigma xy}} = a \underline{\Sigma x} + b \underline{\underline{\Sigma x^2}} \quad - (3)$$

Here $\frac{\text{Eq 2}}{\text{Eq 3}}$ $486 = a6 + b247 \quad \checkmark$
 $20485 = a247 + b11409 \quad \checkmark$

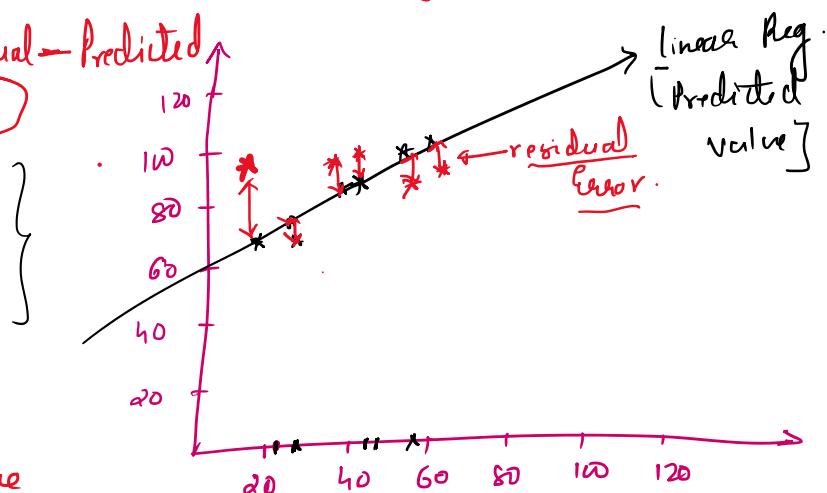
$$a = 65.14$$

$$b = 0.385$$

$$\boxed{y = 65.14 + 0.385x}$$

Using this let us calculate Predicted Y for every X

SR.NO	INDEPENDENT	DEPENDENT(ACTUAL) Y	PREDICTED Y	DIFFERENCE
	X			
1	43	99	81.695	-17.305
2	21	65	73.225	8.225
3	25	79	74.765	-4.235
4	42	75	81.31	6.31
5	57	87	87.085	0.085
6	59	81	87.855	6.855



Residual Error = Error between the Actual Value and Predicted Value.

Jaad Rakho

$$\underline{\underline{y = a + bx}}$$

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

$$b = \frac{n(\Sigma xy) - (\Sigma x)(\Sigma y)}{n(\Sigma x^2) - (\Sigma x)^2}$$

Quantitative Analysis

$$b = \frac{n(\bar{xy}) - (\bar{x})(\bar{y})}{n(\bar{x}^2) - (\bar{x})^2}$$

Given $X \& Y$ the n

Find $\underline{a}, \underline{b}$.

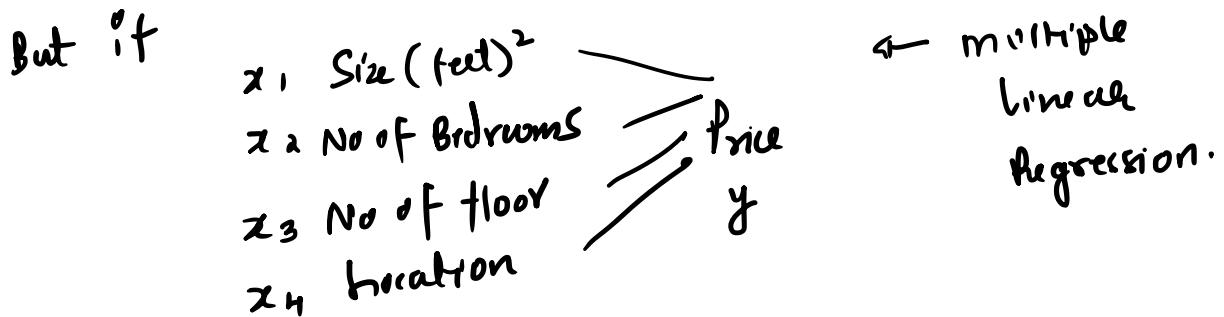
If more than one Independent Variable \Rightarrow

$$y = a + b_1 x_1 + b_2 x_2 + b_3 x_3 \dots$$

↑ impact of x_1 on y ↑ impact of x_2 on y ↑ impact of x_3 on y .

In earlier Example we had one Dependent variable ie y and one Independent variable ie x
 So it was Case of Simple linear regression.

$$\text{Ex } \underline{\text{sqftarea}} \rightarrow \underline{\text{price}} \Rightarrow \hat{y} = \underline{\beta_0} + \underline{\beta_1 x}$$



The linear regression to express the above dependency is

$$\hat{y} = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 + \beta_4 x_4$$

General form of multiple linear Regression.

$$\hat{y} = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_n x_n$$

Here $\beta_0, \beta_1, \beta_2, \dots, \beta_n$ $+ve$
OR
 $-ve$

$$h_0(x) = \beta_0 \underset{x_0=1}{\overset{\uparrow}{x_0}} + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_n x_n$$

Here $x_0, x_1, x_2, \dots, x_n$ = features

$\beta_0, \beta_1, \beta_2, \dots, \beta_n$ \Rightarrow Parameters / Coefficients.

Let X = Feature Vector =

$$\begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_n \end{bmatrix}$$

$n+1 \times 1$
 $\underline{n = \text{no of features}}$

Θ = Parameter Vector =

$$\begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_n \end{bmatrix}$$

$n+1 \times 1$

$$\therefore h_0(x) = (\Theta^T \cdot X)$$

$$\begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} \quad X$$

$$\underline{\underline{h_0(x) = \Theta^T X = \beta_0 x_0 + \beta_1 x_1 + \dots + \beta_n x_n.}}$$

x_0	x_1	x_2	y
1	7	2.6	78.5
1	1	2.9	74.3
1	11	5.6	104.3
1	11	3.1	87.6
1	7	5.2	95.9
1	11	5.5	109.2
1	3	7.1	102.7

$$\checkmark \quad y = \beta_0 + \beta_1 x_1 + \beta_2 x_2$$

To find β_0
 β_1
 β_2

Solve $\xrightarrow{3 \times 3}$

Feature Vector $= X =$

$$X = \begin{bmatrix} 1 & 7 & 2.6 \\ 1 & 1 & 2.9 \\ 1 & 11 & 5.6 \\ 1 & 11 & 3.1 \\ 1 & 7 & 5.2 \\ 1 & 11 & 5.5 \\ 1 & 3 & 7.1 \end{bmatrix}_{7 \times 3} \quad \& \quad Y = \begin{bmatrix} 78.5 \\ 74.3 \\ 104.3 \\ 87.6 \\ 95.9 \\ 109.2 \\ 102.7 \end{bmatrix}_{7 \times 1}$$

$$C = \begin{bmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \end{bmatrix}_{3 \times 1} = \frac{(X^T \cdot X)^{-1} \cdot X^T \cdot Y}{\text{Coefficient Vector}}$$

Solution

$$X = \begin{bmatrix} 1 & 7 & 2.6 \\ 1 & 1 & 2.9 \\ 1 & 11 & 5.6 \\ 1 & 11 & 3.1 \\ 1 & 7 & 5.2 \\ 1 & 11 & 5.5 \\ 1 & 3 & 7.1 \end{bmatrix}$$

$$(1) \text{ find } X^T = \begin{bmatrix} 1 & 1 & 1 & 1 & 1 & 1 \\ 7 & 1 & 11 & 11 & 7 & 11 & 3 \\ 2.6 & 2.9 & 5.6 & 3.1 & 5.2 & 5.5 & 7.1 \end{bmatrix}$$

$$(2) \text{ find } X^T \cdot X \Rightarrow \begin{bmatrix} 7 \\ 51 \\ 32 \end{bmatrix} \begin{bmatrix} 57 & 32 \\ 471 & 235 \\ 235 & 163.84 \end{bmatrix} = \begin{bmatrix} 57 & 32 \\ 471 & 235 \\ 235 & 163.84 \end{bmatrix}_{3 \times 3}$$

$$= \begin{bmatrix} 1.79 & -0.06 & -0.25 \\ -0.06 & 0.01 & -0.0011 \end{bmatrix}$$

$$= \begin{pmatrix} 1.79 & -0.00 & -0.11 \\ -0.06 & 0.01 & -0.0011 \\ -0.25 & -0.0011 & 0.0571 \end{pmatrix}$$

$$= \underline{(X^T X)^{-1}} \cdot X^T \cdot Y$$

$\hat{Y} =$

$$\begin{bmatrix} 51.6 \\ 1.5 \\ 6.72 \end{bmatrix} \quad \begin{array}{l} \beta_0 = 51.6 \\ \beta_1 = 1.5 \\ \beta_2 = 6.72 \end{array}$$

$\underline{\underline{3 \times 1}}$

$$\hat{Y} = 51.6 + 1.5x_1 + 6.72x_2$$

what is y when $x_1 = 3$ $x_2 = 2$

$$y = 51.6 + (1.5 \times 3) + 6.72 \times 2 = \underline{\underline{69.54}}$$

Consider Formulae that will be used, when only 2 independent variables specified.

[>2 features, use Matrix Algebra].

Consider

x_1	x_2	y
3	8	-3.7
4	5	3.5
5	7	2.5
6	3	11.5
2	1	5.7
3	2	2

$$y = Q_0 + Q_1 x_1 + Q_2 x_2$$

$$Q_0 = ?$$

$$Q_1 = ?$$

$$Q_2 = ?$$

$$Q_0 = \bar{y} - Q_1 \bar{x}_1 - Q_2 \bar{x}_2$$

$$Q_1 = \frac{(\sum x_2^2)(\sum x_1 y) - (\sum x_1 x_2)(\sum x_2 y)}{(\sum x_1^2)(\sum x_2^2) - (\sum x_1 x_2)^2}$$

$$Q_2 = \frac{(\sum x_1^2)(\sum x_2 y) - (\sum x_1 x_2)(\sum x_1 y)}{(\sum x_1^2)(\sum x_2^2) - (\sum x_1 x_2)^2}$$

$$\begin{aligned} \bar{y} &= \\ \bar{x}_1 &= \\ \bar{x}_2 &= \end{aligned} \} \text{ mean}$$

where

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

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

$$\sum x_2 = \sum x_2 - \frac{(\sum x_2) \bar{x}}{N}$$

$$\sum x_1 x_2 = \sum x_1 x_2 - \frac{(\sum x_1)(\sum x_2)}{N}$$

$$\sum x_1 y = \sum x_1 y - \frac{(\sum x_1)(\sum y)}{N}$$

$$\sum x_2 y = \sum x_2 y - \frac{(\sum x_2)(\sum y)}{N}$$

Sol =
$$y = 2.796 + 2.28x_1 - 1.67x_2$$

HW

- * To evaluate Performance of m.l models
- * To find how good the model fits on given data

① Karl Pearson's Coefficient of Correlation (r)

→ To calculate relationship bet two variable

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

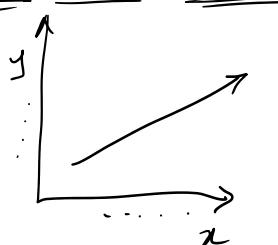
r quantifies the strength of relationship betⁿ two variables.

The value of r be between ± 1 and -1

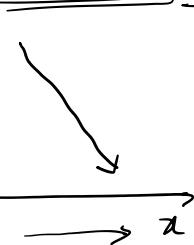
If $r = 1$ ⇒ Total +ve correlation ⇒ If $x \uparrow$ then $y \uparrow$

If $r = -1$ ⇒ Total -ve correlation ⇒ If $x \downarrow$ then $y \uparrow$

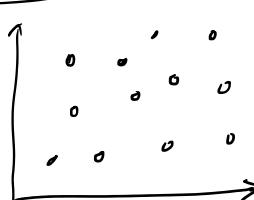
→ gives strength (degree) & direction of correlation: or $x \uparrow$ then $y \downarrow$



+ve
correlation



-ve
correlation



Zero
correlation

Q.
81.10

x	y
151	63
174	81
138	56
186	91
128	47
136	157

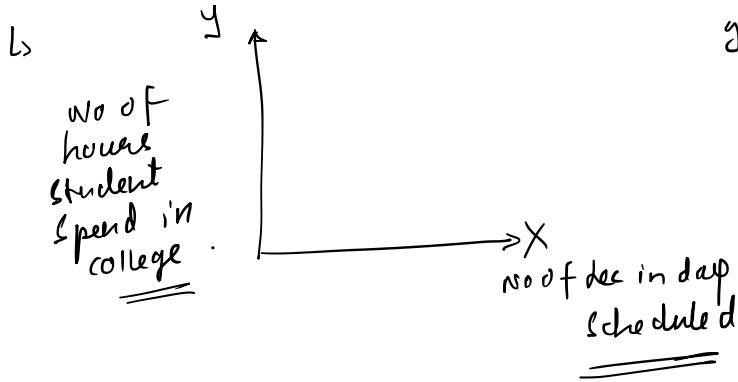
Find $r = ?$

0.9884

186	91
128	47
136	57
179	76
163	72
152	62
131	48

② R^2 method (R Square)

↳ It gives information about good of fit feature of the model.



$$\text{If } R^2 = 0.85$$

Variation in no of hours that students spend in college is 85% dependent on no of hrs scheduled.

↳ Indicate percentage of variance in dependent and independent variable pair

↳ Value varies from 0 to 1

If $R^2 = 1 \Rightarrow$ no diff betⁿ actual & predicted value.

$R^2 = 0$ means the model does not learn any relationship betⁿ variables.

$$\textcircled{1} \quad SST \text{ (Sum of squares of Total)} = \sum_{j=1}^n (y_j - \bar{y})^2$$

↑
actual
y ↑
mean of y.

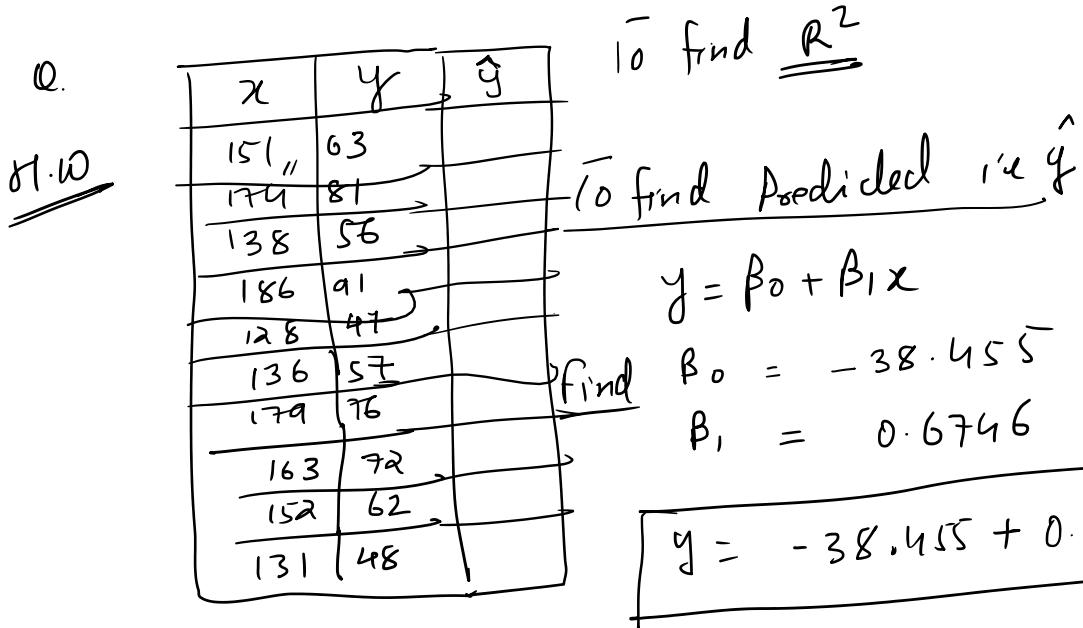
$$\textcircled{2} \quad SSR \text{ (Sum of Squares due to regression)} = \sum_{j=1}^n (\hat{y}_j - \bar{y})^2$$

↑
predicted
y ↑
mean of y.

$$\textcircled{3} \quad SSE \text{ (Sum of Squares of Error)} = \sum_{j=1}^n (y_j - \hat{y}_j)^2$$

↑
actual
y ↑
predicted
y

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



Now $x = 151$ find $\hat{y} = 63.4$
 $x = 174$ $\hat{y} = 78.92$
 $x = 138$ $\hat{y} = 54.63$

Note $R^2 = \frac{\sum (\hat{y} - \bar{y})^2}{\sum (y - \bar{y})^2}$ if line fit properly
 Then $\hat{y} \approx y$
 $\therefore R^2 \approx 1$

③ Standard Error of Estimate \rightarrow

* Measures the accuracy of prediction.

$$* \underline{b_{\text{est}}} = \sqrt{\frac{\sum (y - \hat{y})^2}{N}} = \sqrt{mSE} = \sqrt{\frac{SSE}{N}}$$

(mean
square
error)

* It reflects how well the regression model fits the dataset.

- * Smaller the value better it is
- * Larger the value worst it is

y = actual value
 \hat{y} = predicted value

Q.

H.10

x	y
151	63
174	81
138	56
186	91
128	47
136	57
179	76
163	72
152	62
131	48

$$\boxed{b_{\text{est}} = 2.909}$$

Classification →

Ex Mail → $\begin{matrix} +ve \\ \text{Spam} \end{matrix} / \begin{matrix} -ve \\ \text{Not Spam} \end{matrix}$

Online Transaction → $\begin{matrix} +ve \\ \text{Fraudulent} \end{matrix} / \begin{matrix} -ve \\ \text{Non Fraudulent} \end{matrix}$

Tumor → $\begin{matrix} \text{malignant} \\ +ve \end{matrix} / \begin{matrix} \text{Benign} \\ -ve \end{matrix}$

where $y \in \{0, 1\}$ $\begin{cases} 0 \Rightarrow -ve \text{ class} \\ 1 \Rightarrow +ve \text{ class} \end{cases}$

Here y has Discrete Value.

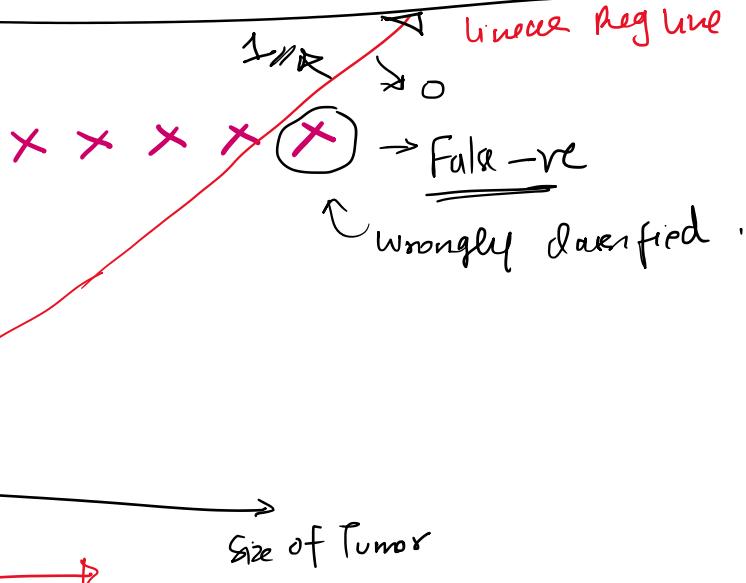
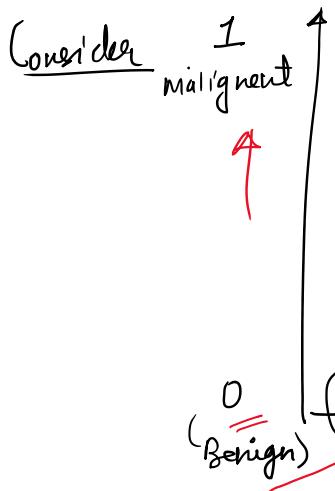
In above case OIP has only 2 class $\{0, 1\}$ } Binary classification.

To check Weather.

Possible OIP $\begin{cases} \text{Windy} \\ \text{Sunny} \\ \text{Cloudy} \\ \text{Rainy} \end{cases}$ } OIP has more than one class [multi class classification]

* Classification bothers about label and not the Exact value.

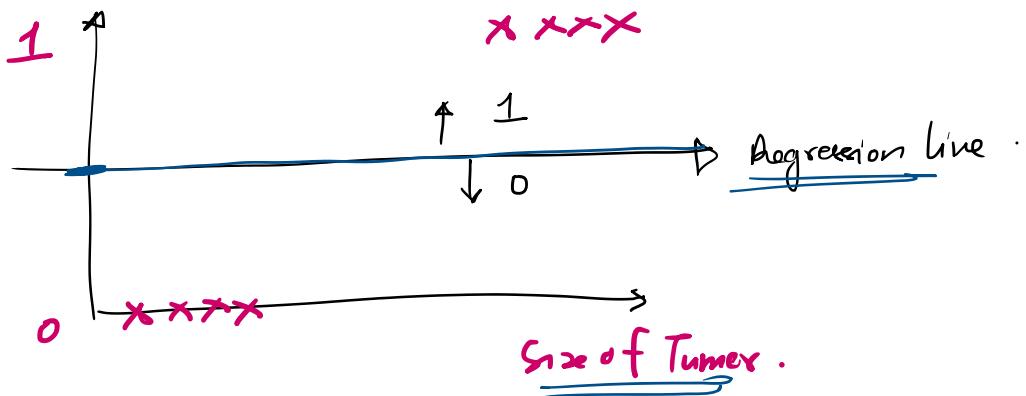
Consider (Why linear regression is not used for classification)



Suppose

Malignancy:

Not Correct



Consider

Malignancy:



here we can observe for above regression line

If tumor size $\leq 5\text{cm}$ \rightarrow Yes (1) Malignant

If tumor size $\leq 5\text{cm}$ \rightarrow Yes (1) Malignant
 tumor size $> 5\text{cm}$ \rightarrow No (0) Benign

We can say let 'P' denote Probability that $y=1$ when $\underline{\underline{X=x}}$.

$$y = \underbrace{P}_{=} (y=1 | \underline{\underline{X=x}}) = \frac{[\beta_0 + \beta_1 x]}{\text{In linear Reg}}$$

P = probability lies bet" $\underline{\underline{0 \text{ to } 1}}$

But linear function are unbounded.

and Expected o/p here is 0 or 1

So we cannot use regression to build classifier

\therefore linear regression is not suitable for classification.

For classification we will use logistic regression.

* In logistic regression we get probability score.

* It predicts the probability of occurrence of event

$$\text{Odd} = \frac{\text{No of time the Event happens}}{\text{No of time the Event will not happen}}$$

Odd = Represents chances that the event will occur

Ex If the odd of India winning against W.I. is $\underline{4:1} = \frac{\text{No of India win}}{\text{No of India not win}} = \frac{4}{1}$

Best Case \Rightarrow The odd of India winning against W.I. = $\underline{\infty}$

If odd of W.I. winning against India is 1:4

$$= \frac{\text{No of W.I. win}}{\text{No of W.I. not win}} = \frac{1}{4}$$

Worst Case \Rightarrow W.I. is winning 0 match = Odd = $\underline{0}$

Range of value that odd can take = 0 to ∞

Relationship between Odd & Probability = Odd = $\frac{P}{1-P}$

$$\therefore \text{odd} = \frac{P}{1-P}$$

Here we know that odd has range 0 to ∞

\rightarrow There is no upper bound for odd.

\rightarrow But odd has lower bound.

\rightarrow To remove lower bound, to have symmetrical analysis

Ex $\text{odd} = 1:6 \quad 1/6 = 0.167 \quad \left. \begin{array}{l} \\ \end{array} \right\} \text{not symmetrical}$

$$\text{odd} = 6:1 \quad 6/1 = 6 \quad \left. \begin{array}{l} \\ \end{array} \right\} \text{symmetrical}$$

But $\ln(1/6) = -1.79 \quad \left. \begin{array}{l} \\ \end{array} \right\} \text{symmetrical}$

$$\ln(6/1) = 1.79 \quad \left. \begin{array}{l} \\ \end{array} \right\} \text{symmetrical}$$

By taking log we have overcomed the lower bound

also

$\log(\text{odd})$ $\begin{cases} \text{will not have upper bound} \\ \text{will not have lower bound.} \end{cases}$

$$\underline{\underline{\log(\text{odd})}} = y = \underline{\underline{\beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_n x_n}}$$

$$\log\left(\frac{P}{1-P}\right) = \beta_0 + \beta_1 x_1 + \dots + \beta_n x_n.$$

$$\text{let } z = \beta_0 + \beta_1 x_1 + \dots + \beta_n x_n.$$

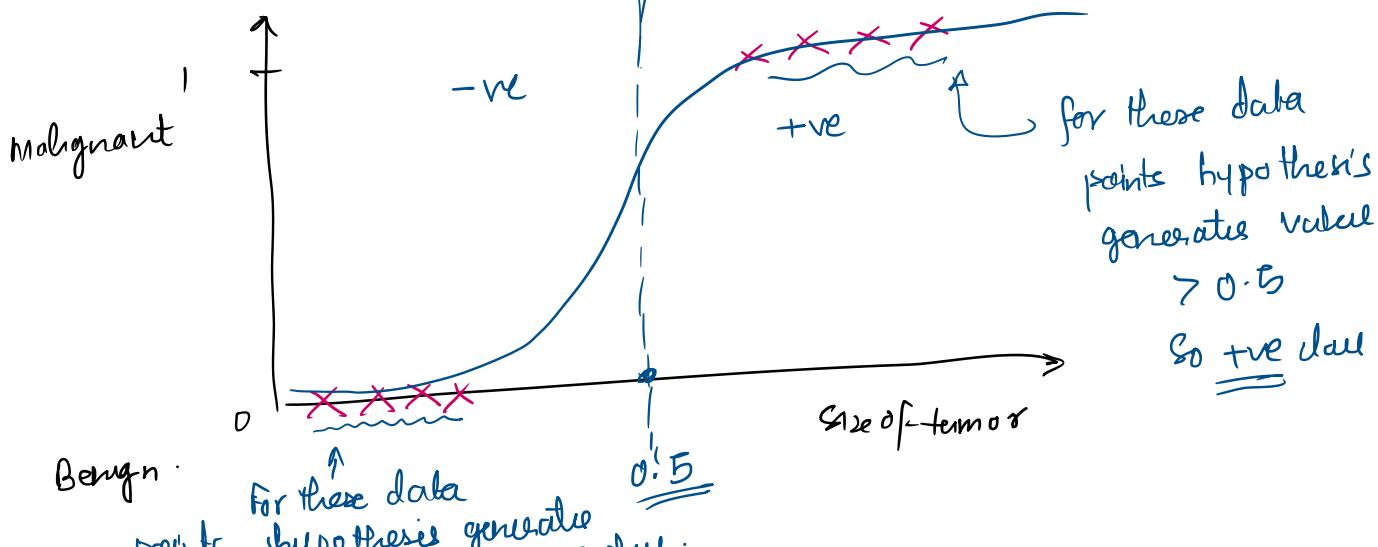
$$\therefore \log_e \left(\frac{P}{1-P} \right) = Z$$

$$\therefore \frac{P}{1-P} = e^Z$$

$$P = -e^Z P + e^Z$$

$$P(1+e^Z) = e^Z$$

$$P = \frac{e^Z}{1+e^Z} = \frac{1}{1+e^{-Z}} = \frac{1}{1 + \frac{1}{-(\beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_n x_n)}} = \text{Sigmoid } f^n$$



logistic regression Model

in linear reg $h_{\theta}(x) = \theta^T x$ → range $-\infty$ to $+\infty$

↑ feature vector
↑ parameter vector
hypothesis (prediction).

In logistic regression

$$0 \leq h_{\theta}(x) \leq 1$$

→ predicted value

$$h_{\theta}(x) = g(\theta^T x)$$

↑ Sigmoid

sigmoid

Let $z = \mathbf{Q}^T \mathbf{x}$

$$h_{\theta}(x) = g(z) = \frac{1}{1+e^{-z}}$$

↑
prediction.

Estimated probability that $y=1$ on given input x .

If $\underline{h_{\theta}(x) = 0.7}$

This means There is 70% chance that tumor is Malignant.

Since $h_{\theta}(x) = 0.7 > 0.5$

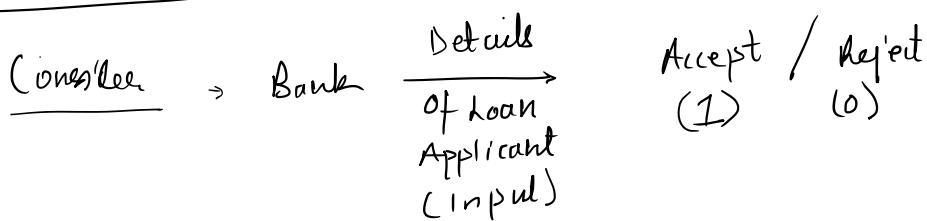
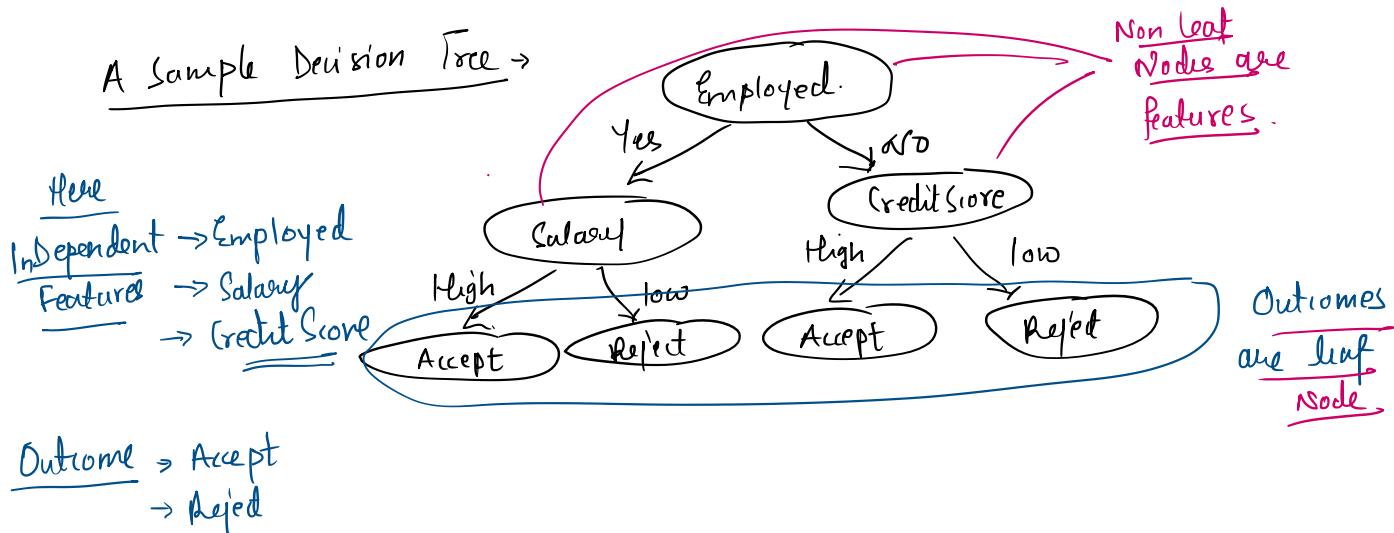
So $\boxed{y=1}$

If $\underline{h_{\theta}(x) = 0.3}$

This means There is 30% chance of tumor being malignant

Since $h_{\theta}(x) = 0.3 < 0.5$

So $\boxed{y=0}$

Decision Tree :A Sample Decision Tree →

Outcome → Accept
→ Rejected

Q) Create Decision Tree for following using Gini Index
(Classification & Regression Tree) → CART. ✓

Weekend	Weather	Parent	Money	Decision	(y)	P
w1	Sunny	Yes	Rich	Cinema	-	
w2	Sunny	No	Rich	Tennis	-	
w3	Windy	Yes	Rich	Cinema	-	
w4	Rainy	Yes	Poor	Cinema	-	
w5	Rainy	No	Rich	Stay In	-	
w6	Rainy	Yes	Poor	Cinema	-	
w7	Windy	No	Poor	Cinema	-	
w8	Windy	No	Rich	Shopping	-	
w9	Windy	Yes	Rich	Cinema	-	
w10	Sunny	No	Rich	Tennis	-	

Solution → Independent features: Weather
Parent
Money

Decision / Outcome
Cinema
Tennis
Stay In
Shopping

↳ Shopping.

Step 1

We will calculate Gini Index for Overall collection of Outcomes of Training Examples.

These are 4 possible outcomes for decision

Cinema — 6 instances
Tennis — 2 instances
StayIn — 1 instance
Shopping — 1 instance.

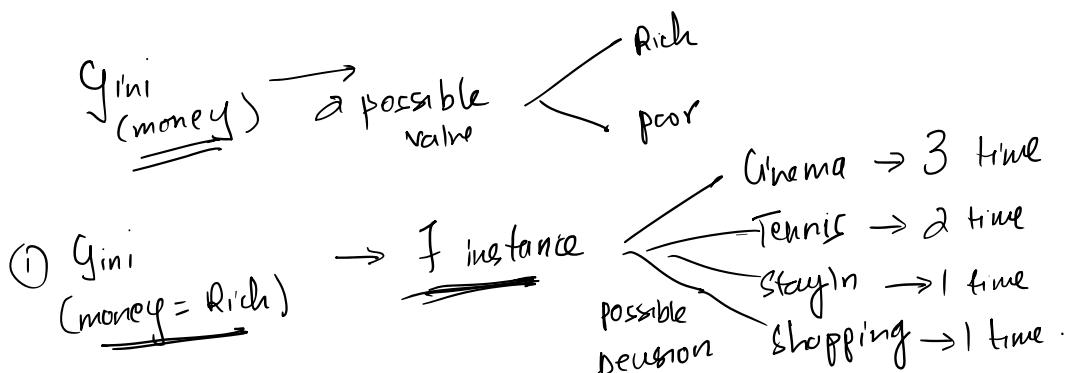
$$G_{ini} = 1 - \left(\left(\frac{6}{10}\right)^2 + \left(\frac{2}{10}\right)^2 + \left(\frac{1}{10}\right)^2 + \left(\frac{1}{10}\right)^2 \right)$$

$$(deusion) = 1 - \left(\frac{42}{100} \right) = 0.58$$

Note → In Machine Learning, Gini index/coefficient is utilized as an Impurity measure in decision tree for Classification.

$$G_{ini} = 1 - \sum_{i=1}^n (P_i)^2 \text{ where } P_i \text{ probability of outcome of specific data}$$

Step 2 To find Gini Index for Money



$$G_{ini} (money = Rich) = 1 - \left(\left(\frac{3}{7}\right)^2 + \left(\frac{2}{7}\right)^2 + \left(\frac{1}{7}\right)^2 + \left(\frac{1}{7}\right)^2 \right)$$

$$= 0.694$$

(2) G_{ini} → Cinema

(2) $Gini_{(money=poor)}$ \rightarrow 3 instance $\xrightarrow{\text{Decision}}$ Cinema.

$$= 1 - ((3/3)^2) = 0_{//}$$

Weighted Average $Gini_{(money)} = (Gini_{(money=Rich)} * \text{proportion of Rich}) + (Gini_{(money=poor)} * \text{proportion of poor})$

$$= (0.694 * 7/10) + (0 * 3/10)$$

$Gini_{(money)}$
 $= 0.485$

Step 3 Gini Index on Parent

For Parent feature $\xrightarrow[\text{values}]{\text{possible values}}$ Yes
No

$Gini_{(parent=Yes)} = 5 \text{ instances} \xrightarrow[\text{possible decision}]{\text{Cinema}}$

$$= 1 - ((5/5)^2) = 0_{//}$$

$Gini_{(parent=No)} = 5 \text{ instances} \xrightarrow[\text{possible decision}]{\text{Tennis} \rightarrow 2 \text{ times}, \text{StayIn} \rightarrow 1 \text{ times}, \text{Shopping} \rightarrow 1 \text{ time}, \text{Cinema} \rightarrow 1 \text{ time}}$

$$= 1 - ((2/5)^2 + (1/5)^2 + (1/5)^2 + (1/5)^2)$$

$$= 1 - (9/25) = 0.72$$

Weighted Average of $Gini_{(parent)} = (0 * 5/10) + (0.72 * 5/10) = 0.36$

$Gini_{(parent)} = 0.36$

$$\boxed{Gini_{(parent)} = 0.36}$$

Step 4) Gini Index for Weather.

Weather $\xrightarrow[\text{values}]{\text{Possible}} \begin{array}{l} \text{Sunny} = 3 \text{ instances} \\ \text{Windy} = 4 \text{ instances} \\ \text{Rainy} \Rightarrow 3 \text{ instances} \end{array}$

$Gini_{(\text{weather} = \text{Sunny})} \Rightarrow \begin{array}{l} 3 \text{ instances} \\ \xrightarrow{\text{possible outcomes}} \begin{array}{l} \text{Cinema} \rightarrow 1 \text{ time} \\ \text{Tennis} \rightarrow 2 \text{ time} \end{array} \end{array}$

$$= 1 - \left(\left(\frac{1}{3}\right)^2 + \left(\frac{2}{3}\right)^2 \right) = \underline{0.444}$$

$Gini_{(\text{weather} = \text{Windy})} \xrightarrow{4 \text{ instances}} \begin{array}{l} \xrightarrow{\text{possible outcomes}} \begin{array}{l} 3 \text{ time Cinema} \\ 1 \text{ time Shopping} \end{array} \end{array}$

$$= 1 - \left(\left(\frac{3}{4}\right)^2 + \left(\frac{1}{4}\right)^2 \right) = \underline{0.375}$$

$Gini_{(\text{weather} = \text{Rainy})} \xrightarrow{3 \text{ instances}} \begin{array}{l} \xrightarrow{\text{possible outcomes}} \begin{array}{l} 2 \text{ cinema} \\ 1 \text{ Stay In} \end{array} \end{array}$

$$= 1 - \left(\left(\frac{2}{3}\right)^2 + \left(\frac{1}{3}\right)^2 \right) = \underline{0.444}$$

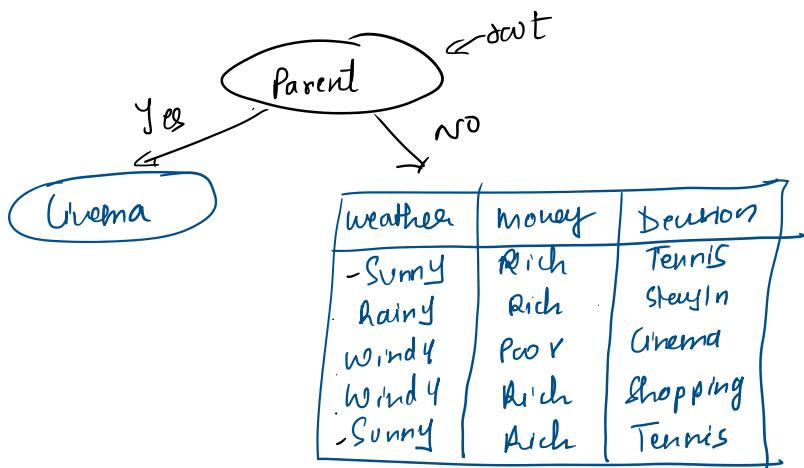
Weighted Average $Gini_{(\text{Weather})} = \frac{(0.44 \times \frac{3}{10}) + (0.375 \times \frac{4}{10}) + (0.44 \times \frac{3}{10})}{3} = \underline{0.414}$

$$\boxed{1 - r_{\text{parent}} = 0.486}$$

$$\boxed{\begin{aligned} Gini(\text{money}) &= 0.486 \\ Gini(\text{parent}) &= 0.36 \\ Gini(\text{weather}) &= 0.416 \end{aligned}}$$

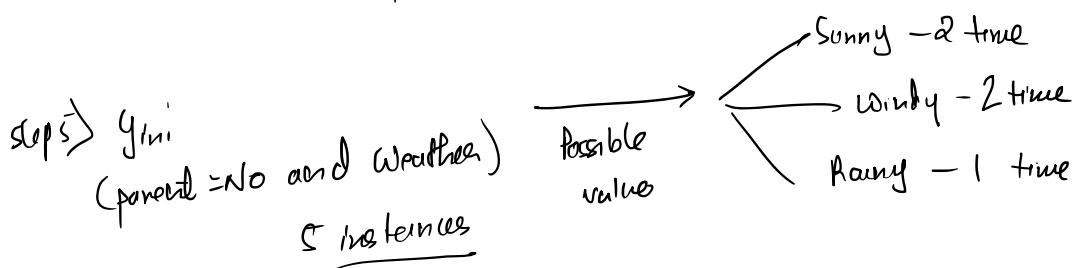
Minimum Gini \rightarrow Minimum Impurity in Decision.
Here Minimum Gini Value = $Gini(\text{parent}) = 0.36$

So the root of Decision is Parent



We need to find $Gini(\text{parent}=\text{No} \text{ and } \text{weather})$

Also $Gini(\text{parent}=\text{No} \text{ and } \text{money})$



$Gini(\text{parent}=\text{No} \text{ and } \text{weather}=\text{Sunny})$

$$\text{2 instance} \rightarrow \text{Tennis} = 1 - \left(\left(\frac{1}{2} \right)^2 \right) = 0 //$$

$Gini(\text{parent}=\text{No} \text{ and } \text{weather}=\text{Windy})$

$$\text{2 instance} \rightarrow \begin{aligned} \text{Cinema} &= 1 - \left(\left(\frac{1}{2} \right)^2 + \left(\frac{1}{2} \right)^2 \right) = 0.5 \\ \text{Shopping} &= \end{aligned}$$

$$\text{Gini}'_{(\text{parent} = \text{No} \text{ and } \text{Weather} = \text{Rainy})} \xrightarrow[1 \text{ instance}]{\text{possible outcome}} \text{StayIn} = 1 - ((Y_1)^2) = 0$$

$$\boxed{\text{Weighted Average Gini}_{(\text{parent} = \text{No} \text{ and } \text{Weather})} = 0.5 \times \frac{4}{5} = \underline{\underline{0.2}}}.$$

Step 6) $\text{Gini}'_{(\text{parent} = \text{No} \text{ and } \text{Money})}$

Rich = 4 time
poor = 1 time
possible values.

5 instance

$$\text{Gini}'_{(\text{parent} = \text{No} \text{ and } \text{Money} = \text{Rich})} \xrightarrow[\text{(4P)}]{\text{possible outcome}} \begin{array}{l} \text{Tennis} - 2 \\ \text{StayIn} - 1 \\ \text{Shopping} - 1 \end{array}$$

$$= 1 - \left((2/4)^2 + (1/4)^2 + (1/4)^2 \right) = \underline{\underline{0.625}}$$

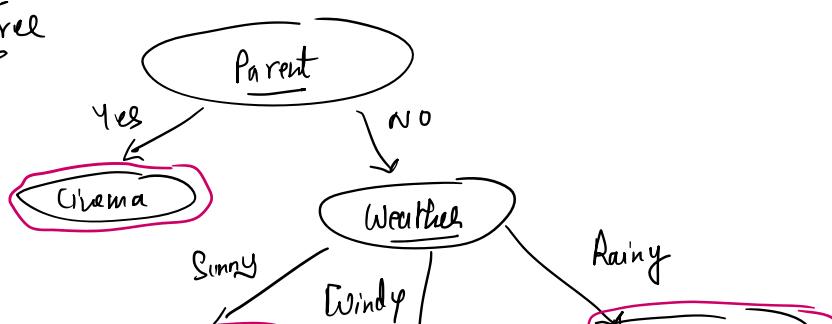
$$\text{Gini}'_{(\text{parent} = \text{No} \text{ and } \text{Money} = \text{Poor})} \xrightarrow{\text{possible outcome}} \text{Cinema} = 1 - ((Y_1)^2) = 0$$

$$\boxed{\text{Weighted Average Gini}'_{(\text{parent} = \text{No} \text{ and } \text{Money})} = 0.625 \times \frac{4}{5} = 0.5}$$

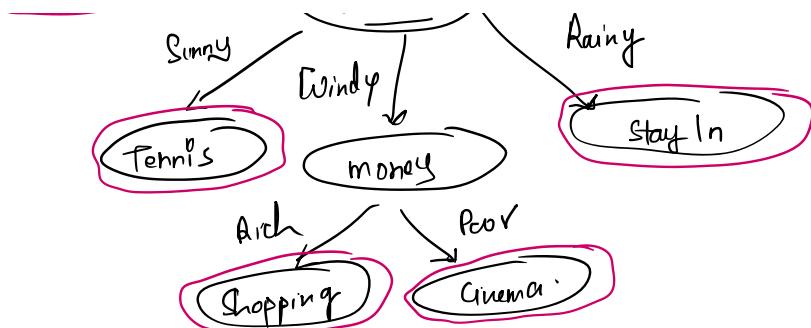
Here $\text{Gini}'_{(\text{parent} = \text{No} \text{ and } \text{Weather})}$ has smallest value

So now Next Node = Weather

* Updated Decision Tree



Aho



HW

Construct an optimal Decision Tree for following

outlook	Temperature	Humidity	windy	Play (Decision)
Sunny	Hot	High	False	NO
Sunny	Hot	High	True	NO
Overcast	Hot	High	F	Y
Rainy	Mild	High	F	Y
Rainy	Cool	Normal	F	Y
Rainy	Cool	Normal	T	NO
Overcast	Cool	Normal	T	Y
Sunny	Mild	High	F	NO
Sunny	Cool	Normal	F	NO
Rainy	Mild	Normal	F	Y
Sunny	Mild	Normal	T	Y
Overcast	Mild	High	T	Y
Overcast	Hot	Normal	F	Y
Rainy	Mild	High	T	NO

Total given Data points

Consider \rightarrow

height



- ① \rightarrow Test Set
- ② \rightarrow Train Set

weight

bias ↑
variance ↓

(Underfitting)

model be
longer straightline

weight

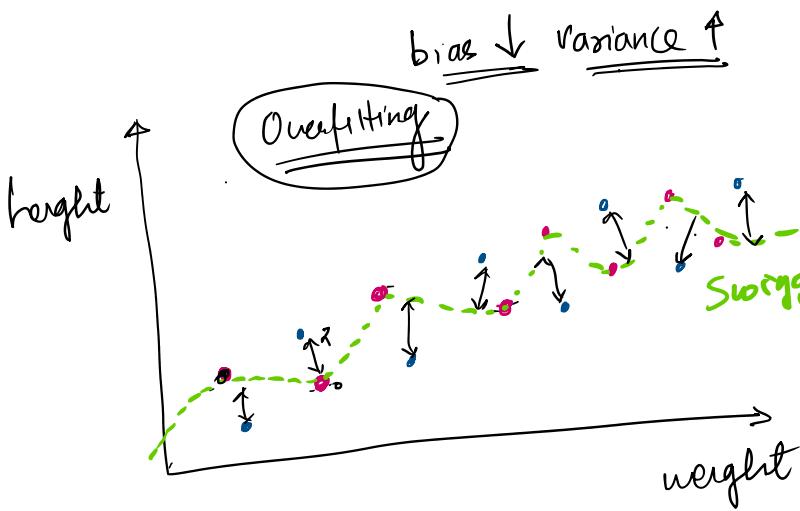
height

Train set

Here any simple straight line cannot match all data points in

Training Set

Bias \rightarrow Inability of ML methods
(which reg) to capture true
relationship



How the Swiggly line
handles the true
relationship betⁿ
weight & height
So since the line
almost passes through
every training data
points we say \Rightarrow little/
zero

bias ↓ variance ↑

Overfitting

↓

Swiggly line

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

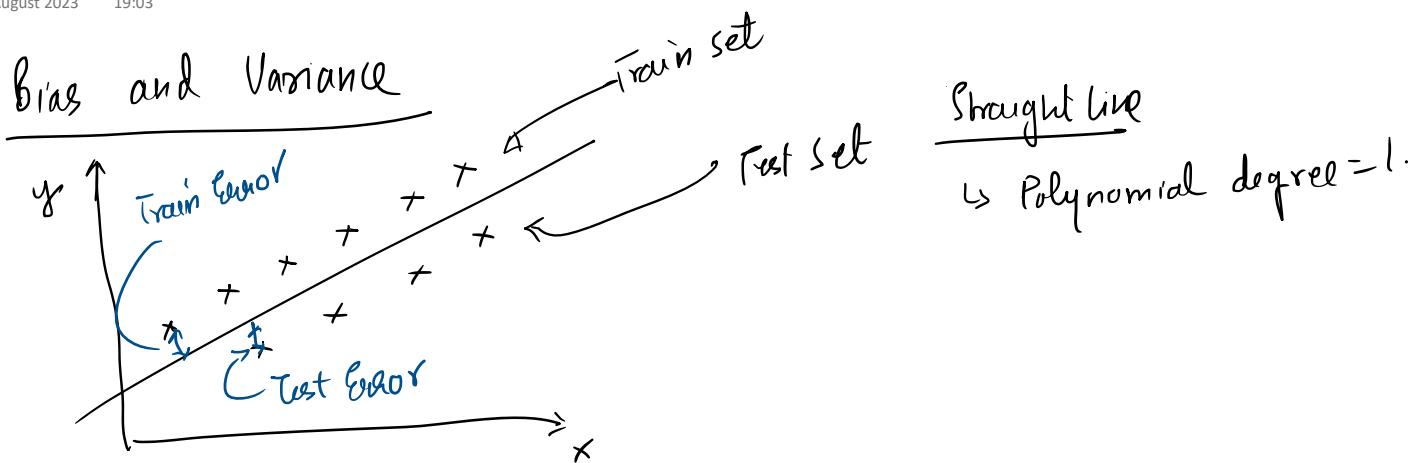
↓

↓

↓

points we say \Rightarrow Zero
Bias.

Variance \Rightarrow Difference in fits betⁿ
the Datasets (Train & Test Data Sets).



- * Underfitting :- When the model is simple (lower degree polynomial), the model might not fit the train set data points.

So here the difference betⁿ Actual and predicted value for Train set is higher

This is bias.

So In Underfitting bias ↑

There will error for predicted value and actual value for Test set

So the diff in error for Test and Train set is low

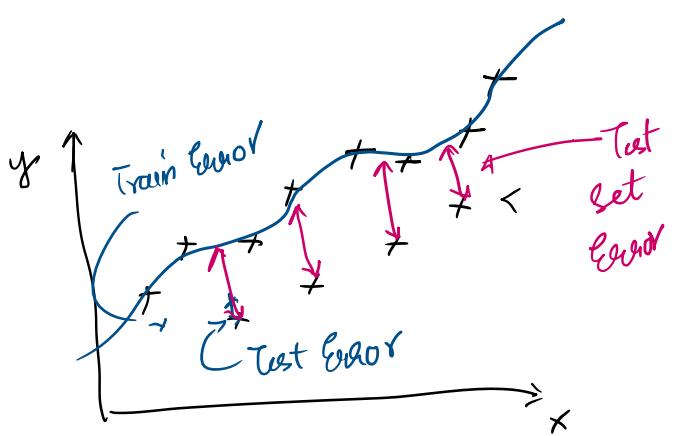
This is Variance

So In Underfitting Variance ↓

* Underfitting

- For Simple Model (lower order Polynomial)
- ◦ Bias ↑ Variance ↓

Bias ↑ Variance ↑



Here model is complex
(higher order polynomial)

→ So Error with Train Set is very less
∴ Bias is ↓

Since the model perfectly fits the train set
it is case of Overfitting

and there is significant diff (error) with Test set

∴ Variance is ↑

* Overfitting →

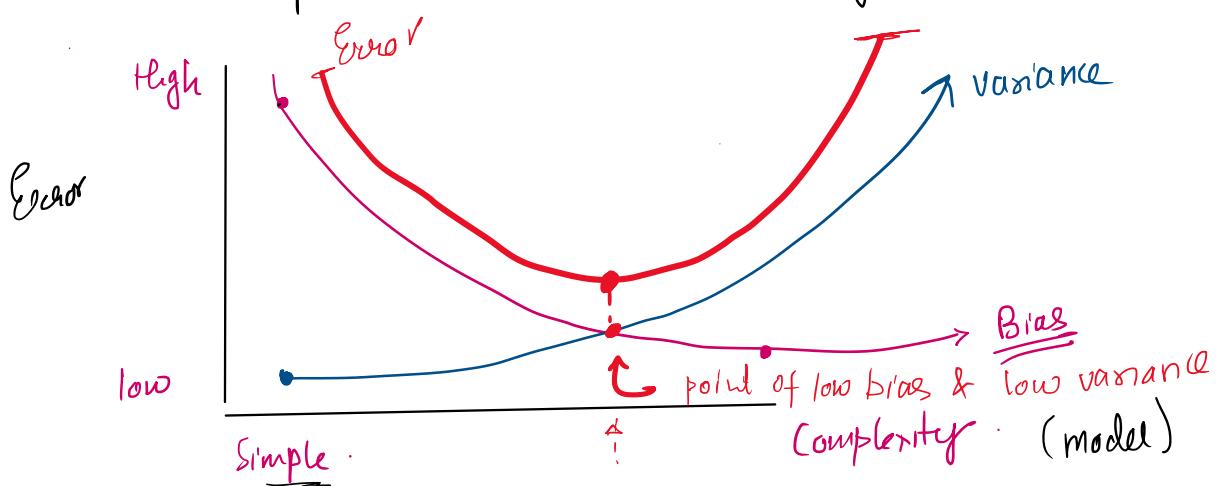
For Complex model (higher order polynomial)

bias ↓ variance ↑

GMP Bias - Variance Tradeoff

(1) If model is too simple and has very few parameters → Underfitting
 → High bias
 → Low variance

(2) If model is complex and has large no of parameters → Overfitting
 → Low bias
 → High variance.



Bias Variance Tradeoff says

We need a model that gives

(1) low bias

(2) low variance.

i.e we need to find point of low bias and low variance.

The methods to achieve this →

* Regularization (Penalize ' θ ' parameters)

* Boosting }

* Bagging }

* ... to Minimize the Total Error: $\text{Bias}^2 + \text{Variance} + \text{Irreducible Error}$

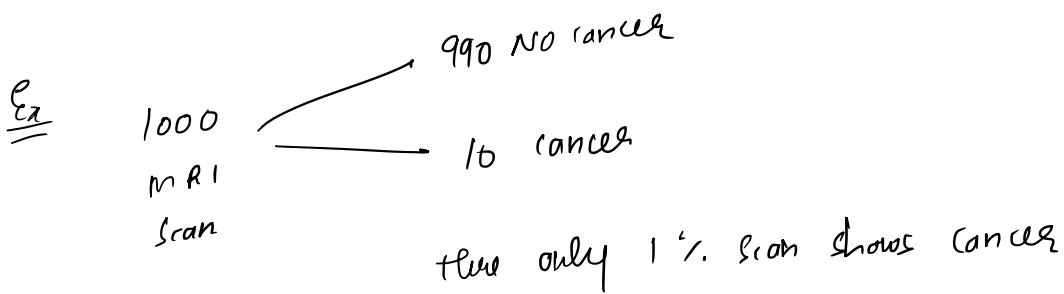
* ~~we need to minimize the Total Error: Bias + Variance + irreducible error.~~

We need to Minimize the Total Error: Bias + Variance + irreducible error.

$$\hat{y} = \theta_0 + \theta_1 x_1 + \theta_2 x_2$$

* Performance Metric :-

Consider a Skewed Data Set (One Sided).



- it is very less -
- Normally Allowed
- might be Ignored
- But this is incorrect.

In Such Case to Evaluate Performance We may use

Confusion Matrix →

		<u>Actual</u>	
		<u>True</u>	<u>False</u>
<u>Predicted</u>	<u>True</u>	True +ve (TP)	False +ve
	<u>False</u>	False -ve	True -ve

The above Confusion Matrix is OK for 2 classes

Question →

Cancer ?

Predicted

		<u>Actual</u>	
		Has Cancer	Do not have Cancer
<u>Predicted</u>	<u>Has Cancer</u>	TP	FP
	<u>Do not have Cancer</u>	FN	TN

If we have More than Two classes

↓ , Actual , ↓ , +

Question →
Person Watched
Chakde ?

if person
watched +ve
chakde
else -ve

Actual

		Chakde		KGF	DDLJ
		TP	FP	FP	FP
Predicted	Chakde	TP	FP	FP	FP
	Not Chakde	FN	TN	TN	TN
		FN	TN	TN	TN

Predicted that these people have not watched Chakde so N;

Person has not watched Chakde & Predicted also Not watching Chakde.

Question → Watched Chakde?

For 4 classes

Actual

		Chakde	KGF	DDLJ	Gadar
		TP	FP	FP	FP
Predicted	Chakde	TP	FP	FP	FP
	Not Chakde	FN	TN	TN	TN
		FN	TN	TN	TN

Chakde Nahi dikha phir bhi predicted Dekha

Chakde Nahi Dekha But predicted Nahi Dekha

Chakde Nahi Dekha And predicted Nahi Dekha.

* TP (True +ve) = Correctly Identified +ve (True HAI → +ve Predicted)

FP (False +ve) = Incorrectly Identified +ve (True Nahi → +ve Predicted)

TN (True -ve) = Correctly Identified -ve ⇒ (-ve HAI → -ve Predicted)

FN (False -ve) ⇒ Incorrectly Identified -ve ⇒ (-ve Nahi HAI → -ve Predicted)

$$* 1) \text{Accuracy} = \frac{\text{Total Correct Prediction}}{\text{Total Prediction}} = \frac{TP + TN}{TP + TN + FP + FN}$$

$$2) \text{Error Rate} = 1 - \text{Accuracy}$$

$$3) \text{True Positive Rate} = \frac{\text{Correctly Identified } +ve}{\text{Total Actual } +ve} = \frac{TP}{TP + FN}$$

(TPR)
[Sensitivity / Recall]

$$4) \text{False Negative Rate} = \frac{\text{Incorrectly Identified } +ve}{\text{Total Actual } +ve} = \frac{FN}{TP + FN}$$

(Also +ve false-negative prediction)

(FNR).

$$5) \text{True Negative Rate} = \frac{\text{Correctly Identified } -ve}{\text{Total Actual } -ve} = \frac{TN}{TN + FP}$$

(TNR)
[Specificity]

$$6) \text{False Negative Rate} = \frac{\text{Incorrect Identified } -ve}{\text{Total Actual } -ve} = \frac{FP}{TN + FP}$$

(FNR)

$\overline{TPR} = \frac{\text{Correct } +ve}{\text{Total } +ve}$

 $\overline{FNR} = \frac{\text{Incorrect } +ve}{\text{Total } +ve}$

$\overline{TNR} = \frac{\text{Correct } -ve}{\text{Total } -ve}$

 $\overline{FP} = \frac{\text{Incorrect } -ve}{\text{Total } -ve}$

Note To Identify Heart Disease

$$\text{Sensitivity} = TPR = \frac{\text{Correct Identified } +ve}{\text{Total } +ve}$$

= What percentage of patient with heart disease are correctly

$$\text{Sensitivity} = \frac{\text{True Positives}}{\text{Total Positive}}$$

Total true

patients with disease are correctly identified

$$\text{Specificity} = TNR = \frac{\text{Correct Identified -ve}}{\text{Total -ve}} = \text{What percentage of patients without heart disease are correctly identified.}$$

$$FPR = 1 - \text{Specificity}$$

$$= 1 - \frac{\text{TN}}{\text{TN} + \text{FP}} = \frac{\text{TN} + \text{FP} - \text{TN}}{\text{TN} + \text{FP}} = \frac{\text{FP}}{\text{TN} + \text{FP}}$$



Training Error →

If \hat{y} is prediction error that we get when we apply the model on same data from where it is trained.

$$\mathcal{E}_{\text{Train}} = \frac{1}{n} \sum_{i=1}^n \text{Error} \left(\hat{y}_i, y_i \right)$$

↑
for all the samples
↑
prediction of x_i

↑
actual value -

(2) Test Error: If \hat{y} is prediction error we get when we apply model on altogether different data set (test set) and not on the data on which it is trained.

$$\mathcal{E}_{\text{Test}} = \frac{1}{n} \sum_{i=1}^n \text{Error} \left(\hat{y}_i, y_i \right)$$

↑
for Test Set //

(3) Generalization Error ⇒ also known as Out of Sample Error

→ Measure of how accurately an algorithm

is able to predict outcome values for previously unseen data -

→ We want to know how the model will perform

→ We want to know how the model will perform on future data (we do not have today)

→ For Future we do not have x_i (input)
 y_i (output)

$$E_{\text{gen}} = \int_{\substack{\text{overall} \\ \text{possible value} \\ \text{of } x \& y}} \ell_{\text{error}} \left(\underbrace{f(x_i)}_{\substack{\uparrow \\ \text{Predicted}}}, \underbrace{y_i^*}_{\substack{\uparrow \\ \text{actual}}} \right) \underbrace{P(y, x)}_{\substack{\uparrow \\ \text{How often} \\ \text{we expect} \\ \text{such } x \& y}} dx.$$

Usually

$$\underbrace{E_{\text{train}} \leq E_{\text{gen}}}$$

as we do not have value of future $P(y, x)$

so we do not compute generalizⁿ error,

We approximate it with Testing Error.

- * To decide which one should we work for between Sensitivity & Specificity
- * If identifying +ve is more important to us, then we will select algo that has high sensitivity
- * If correctly identifying -ve is more important then we will select algo that has high specificity

Precision & Recall \Rightarrow

\rightarrow Used for Information Retrieval.

\rightarrow Google Search Engine \Rightarrow

\rightarrow query fired

\rightarrow Have millions of related records

\rightarrow From these top 10-100 records are returned.

$$\text{Precision} = \frac{\text{Correctly predicted +ve}}{\text{Total +ve Predicted}} = \frac{TP}{TP + FP}$$

Range (0 - 1)

$$\text{Recall} = \frac{\text{Correctly Identified +ve}}{\text{Total Actual +ve}} = \frac{TP}{TP + FN}$$

(TPR)
(range 0 - 1)

* Ideally we want precision to be high (≈ 1) for a good classifier

(range 0 - 1)

* Ideally we want Precision to be high (*i.e.* 1) for a good classifier

$$\text{Precision} = 1 = \frac{\text{TP}}{\text{TP} + \text{FP}} = 1 \Rightarrow \text{when } \boxed{\text{FP} = 0}$$

* Ideally we want Recall to be very high (*i.e.* 1) for a good classifier

$$\therefore \text{Recall} = 1 = \frac{\text{TP}}{\text{TP} + \text{FN}} \Rightarrow \text{when } \boxed{\text{FN} = 0}$$

So ideally a good classifier has High Precision & Recall

But in reality there is trade-off

* When we tweak our model to increase one, then the other decreases.
(update)

Q) Explain

F1-Score

- * In reality we need a metric that takes into account both precision and recall.

- * F1 score is a metric that takes into account both precision & recall.

- * F1 Score is harmonic mean of Precision & Recall.

$$\text{F1-Score} = \frac{\frac{2}{\text{Precision}} + \frac{1}{\text{Recall}}}{\text{Precision} + \text{Recall}} = \frac{2 * \text{Precision} * \text{Recall}}{\text{Precision} + \text{Recall}}$$

Harmonic mean of two variables $\frac{2}{\frac{1}{a} + \frac{1}{b}} = \frac{2ab}{a+b}$

for n variables

$$H = \frac{n}{\frac{1}{n_1} + \frac{1}{n_2} + \dots + \frac{1}{n_n}} = \frac{n}{\sum_{i=1}^n \frac{1}{n_i}}$$

If F1-Score = 1 \Rightarrow when Precision = 1
Recall = 1

- * When Precision and Recall both are high then F1-Score is high

When to Use F1-Score \rightarrow

\rightarrow Accuracy is not a good metric to use when we have class imbalance.

Ex \rightarrow let say 99% of people visiting site are onlookers and not purchasing anything.

\rightarrow Suppose we have a model that predicts that 1% people visiting site are onlookers.

" 1% error is acceptable

people visiting site are onlookers.

→ The model is 1% wrong, Generally 1% error is acceptable

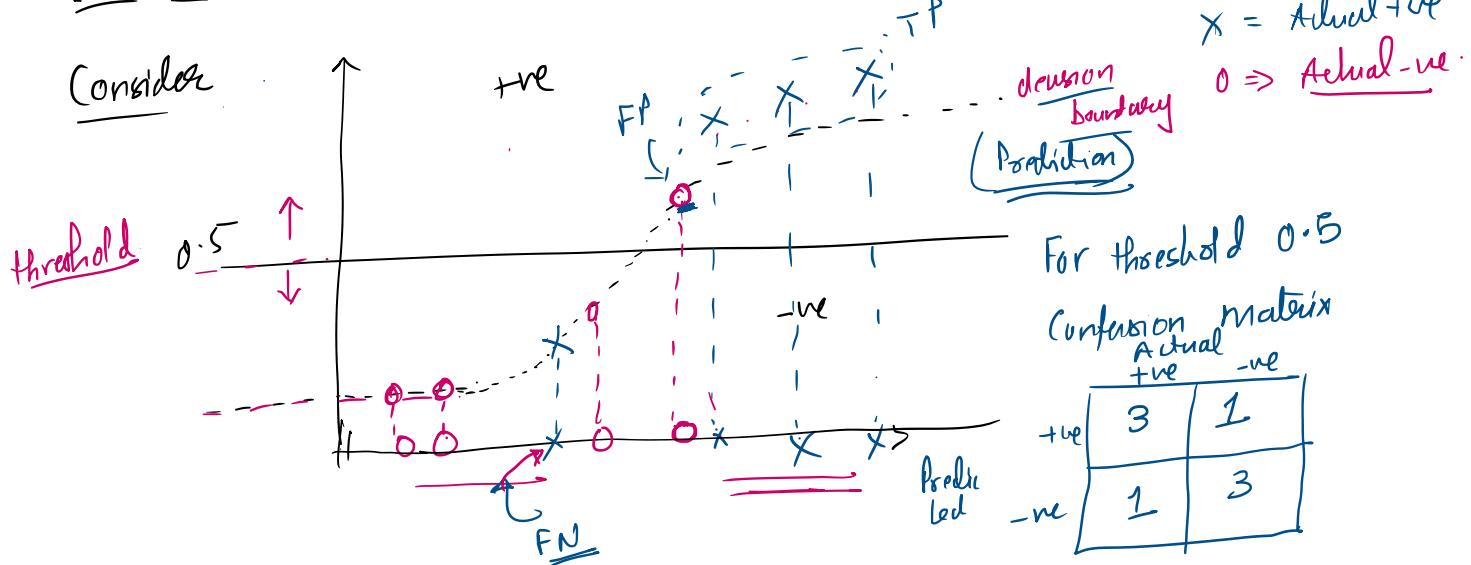
→ But such model in this case is useless

→ In such case instead of accuracy, we will prefer

F1-Score.

ROC [Receiver Operator characteristic]

Consider

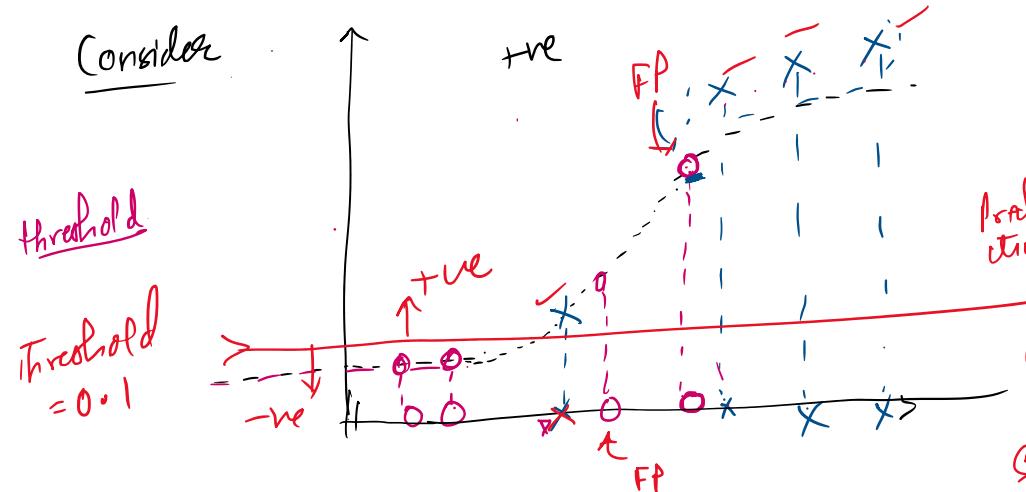


$$\text{Sensitivity} = \frac{3}{4} \quad \text{TPR}$$

$$\text{Specificity} = \frac{3}{4} \quad \text{TNR}$$

$$\boxed{\text{FP.R} = 1 - \text{Specificity}}$$

Consider



Now threshold = 0.1

Confusion matrix

		Actual +ve	Actual -ve
Actual +ve	4	2	
Actual -ve	0	2	

$$\text{Sensitivity} = \frac{4}{4} = 1$$

$$\text{Specificity} = \frac{2}{4} = \frac{1}{2}$$

- * Consider for logistic regression, where we identify a threshold point and prepare confusion matrix and calculate Sensitivity & Specificity.

- and calculate Sensitivity & Specificity
- * If threshold changes then the confusion matrix and accordingly the Sensitivity and Specificity changes.
 - * We can have many such thresholds betⁿ 0 → 1
 - * We want to analyze the performance at diff threshold and want to identify the best of it.
 - * For this we will plot TPR and FPR for diff threshold.
-

→ From the above ROC curve for Random Forest model, it will help us to find the best threshold.

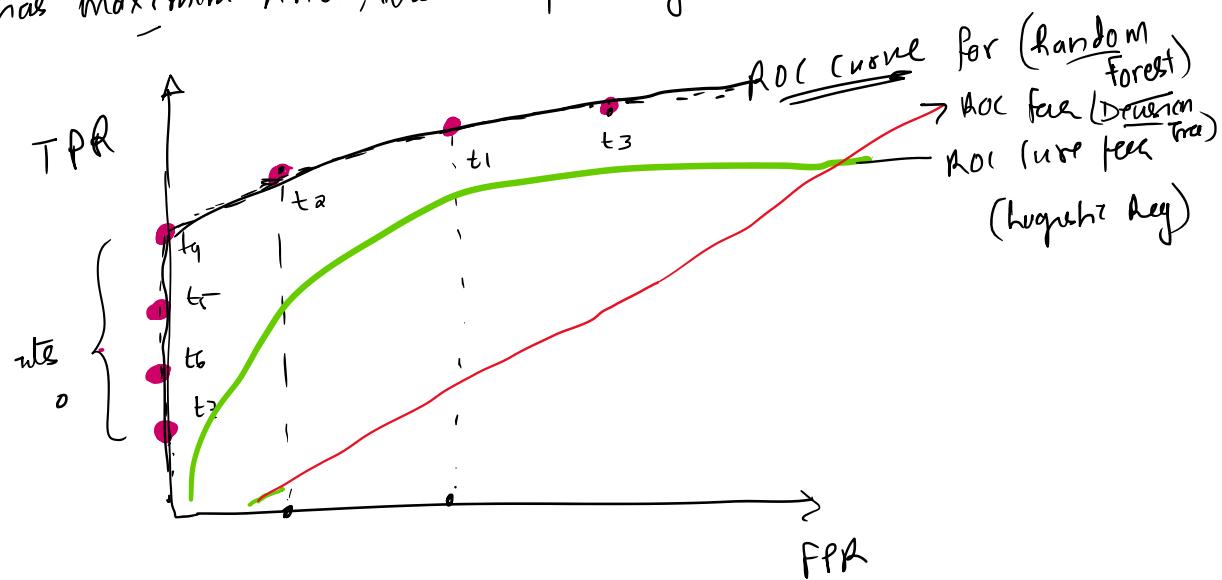
Q) Explain

AUC [Area Under Curve] → It is a method to compare ROC for more than one method and will help to judge which one is better.

→ Here will find Area Under ROC for each method.
→ ROC that has maximum Area, the corresponding method will be best.

ROC curve for (Random Forest)

→ ROC that has maximum Area, the better it is.



From above ROC Curves, we find that there is maximum area under ROC Curve for Random forest, hence the method Random forest will be the best method.

Q) Explain Kappa Statistic?

- * Kappa Statistic or Cohen's Kappa is statistical measure of inter-rater reliability for categorical variable.
- * It is used when two/more raters apply a criteria based on a tool to assess whether or not some condition occurs.
- * Ex: let say Two doctors rates whether or not each of 20 patients has diabetes based on symptoms.
- * If two raters uses same criteria on same target to evaluate and then their agreement is very high then we will have evidence of reliable rating.
- * If their agreement is not very high then →
 - either criterion tool is not useful
 - or raters are not trained enough.
- * Kappa Statistics correct for chance agreement and not percent agreement.

Evaluator A

	Yes	No
Yes	40	20
No	15	40

Evaluator B

- 35 times both agreed - said Yes
- 60 times both agreed → said No
- 20 time A said NO but B said Yes
- 15 times A said Yes but B said No.

11

\Rightarrow do time A said NO but B said YES
 \Rightarrow is time A said YES but B said NO.

Cohen suggested following statistics. \Rightarrow

value ≤ 0	\Rightarrow No agreement
0.01 \rightarrow 0.20	\Rightarrow as <u>none</u> to <u>slight</u>
0.21 \rightarrow 0.40	\Rightarrow as <u>fair</u>
0.41 \rightarrow 0.60	\Rightarrow as <u>moderate</u>
0.60 \rightarrow 0.80	\Rightarrow as <u>substantial</u>
0.81 \rightarrow 1.00	\Rightarrow <u>perfect agreement</u>

* Rather than calculating the percentage of items, the raters agreed on Cohen's Kappa attempts to account the fact that rater may happen to agree on some items purely by chance

Ex Two curators asked to rate 70 paintings.

		Curator C2	
		Yes	No
Curator C1	Yes	25	10
	No	15	20

Step 1) Calculate Relative Agreement betⁿ Curators.

$$P_o = \frac{\text{Both Said Yes} + \text{Both Said No}}{\text{Total}} = \frac{25 + 10}{70} = 0.6429$$

Step 2) Calculate hypothetical probabilities of chance Agreement betⁿ Curators.

$$P(\text{Yes}) = \underline{C_1(\text{Yes})} \times \underline{C_2(\text{Yes})} = \underline{\frac{25}{70}} \times \underline{\frac{25}{70}} = 0.2857$$

$$\checkmark P(Y_{10}) = \frac{C_1(Y_{10})}{\text{Total Ans}} * \frac{C_2(Y_{10})}{\text{Total Ans}} = \frac{(25+10)}{70} * \frac{(25+15)}{70} = \underline{\underline{0.2857}}$$

$$\checkmark P(N_0) = \frac{C_1(N_0)}{\text{Total Ans}} * \frac{C_2(N_0)}{\text{Total Ans}} = \frac{(15+20)}{70} * \frac{(10+20)}{70} = \underline{\underline{0.214285}}$$

$$\underline{\underline{P_e}} = P(Y_{10}) + P(N_0) = 0.2857 + 0.214285 = \underline{\underline{0.5}}$$

$$\text{Calculate Cohen's Kappa} = K = \frac{P_o - P_e}{1 - P_e} = \frac{(0.6429 - 0.5)}{1 - 0.5} = \underline{\underline{0.2857}}$$

gt is in range $0.21 \rightarrow 0.40$ so the agreement betⁿ two
 Cuthakor is fair

Module - 3

Ensemble learning →

3.1 Understand Ensembles

K-fold Cross Validation

⇒ Boosting

Stumping (Adaboost)

XGBoost

3.2 ⇒ Bagging

Subagging

Random Forest

Comparison with Boosting

Different ways to
combine classifiers

ML VCP LEC 5 16 AUGUST 2023 Page 59

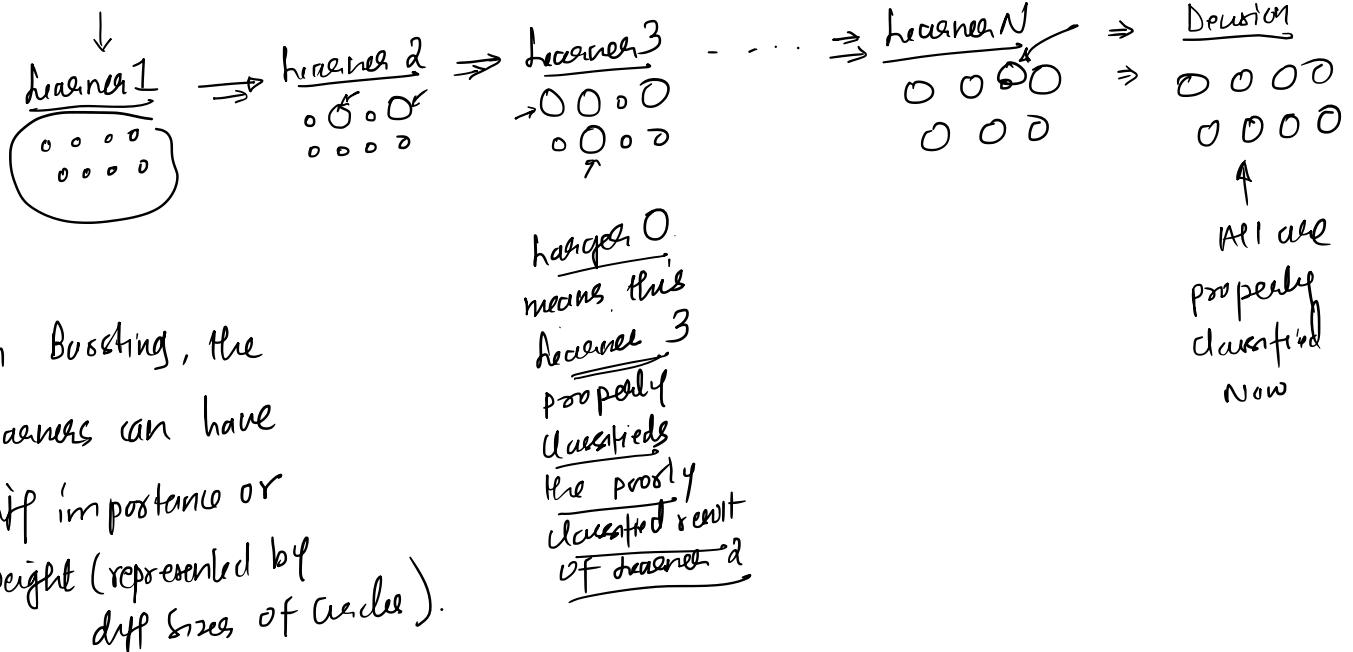
* Ensemble ->

- * In ML, ensemble is a model that combines the prediction from two or more models.
- * The models that contributes to Ensemble are known as ensemble members.
- * The members may or may not be trained on same training data and they may be of same type or different type.
- * It's very powerful method to improve the performance of the model.
- * It's technique that uses group of weak learners in order to create a strong and aggregated learner.

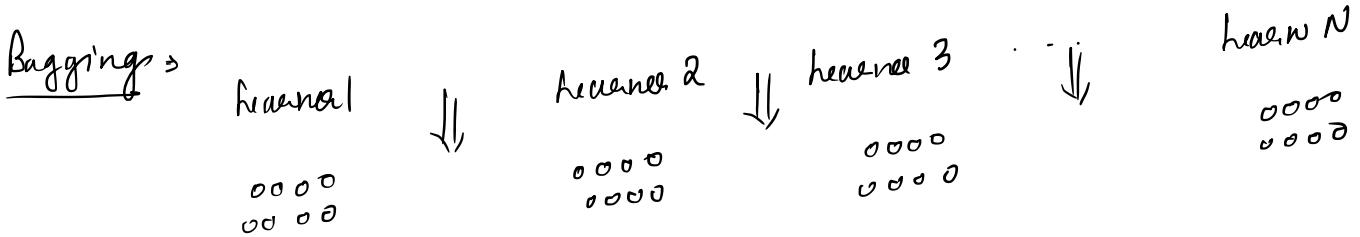
→ The Ensemble technique helps to reduce the Variance (By Bagging) and Bias (By Boosting) and thus helps in improving the predictions.

- Boosting model :-
- * It falls inside family of Ensemble method.
 - * It consists of filtering or weighting the data that is used to train team of Weak Learners, so that the new learner can give more weight on sample that is poorly classified by previous learner.
 - In Boosting the learners are trained

Sequentially

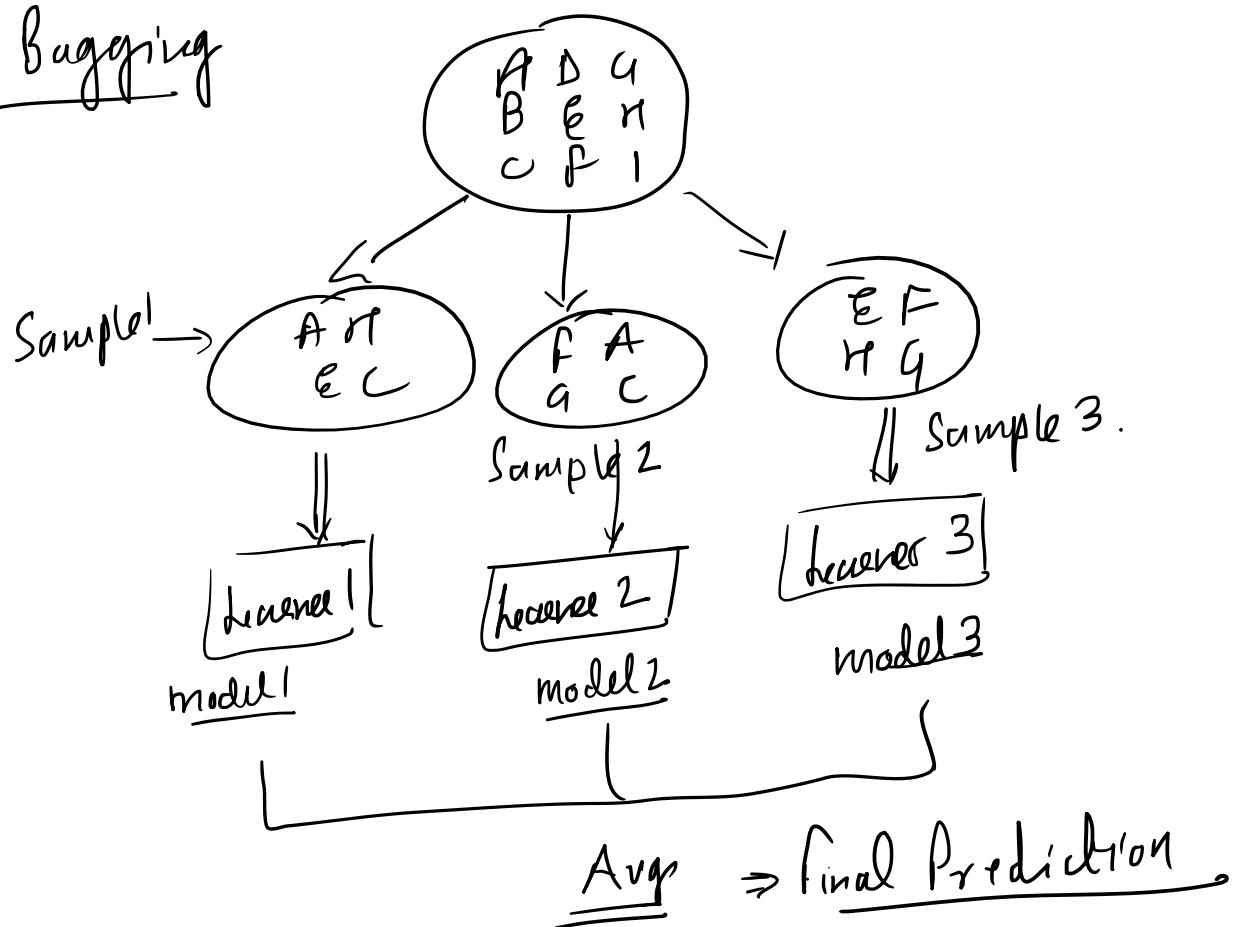


* In Boosting, the learners can have diff importance or weight (represented by diff sizes of circles).



- ⇒ In Bagging the weak learners are trained in parallel using randomness
- ⇒ All learners have same weights.

Note Bagging

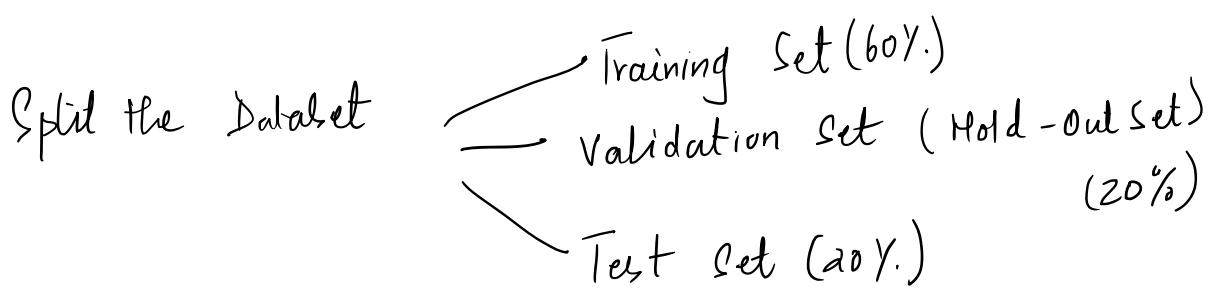


- * Bagging helps in reducing variance
- * Boosting helps in reducing bias.

Cross Validation →

- In Supervised ML
- Train a Model on a Dataset
- Trained model is used to predict the target given new sample.
- How to know if model we have trained will produce effective and accurate result on new input

Cross Validation → It is process that ensures the model will perform well on new Data.



Training Set = part of data on which model is trained.
(This dataset will help to * build model)

* Validation Set ⇒ * Evaluate the Model

- will help to check if model overfits or Underfits -
- update the parameters and again train the model
- Repeat this until the model performs best on Validation set

→ Repeat until ...
Validation set

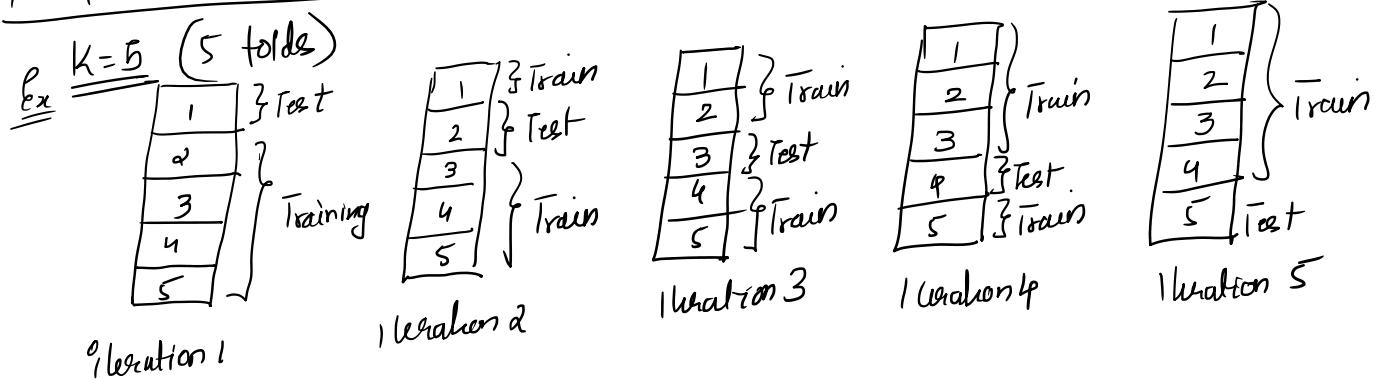
Test Set \Rightarrow^* Prediction

→ The fully trained model after being evaluated on validation set can be used on test set to generate Estimation.

Q) Types of Cross Validation →

- ① The standard validation set Approach
- ② Leave one out cross validation.
- ③ K-fold Cross Validation

K-fold Cross Validation →

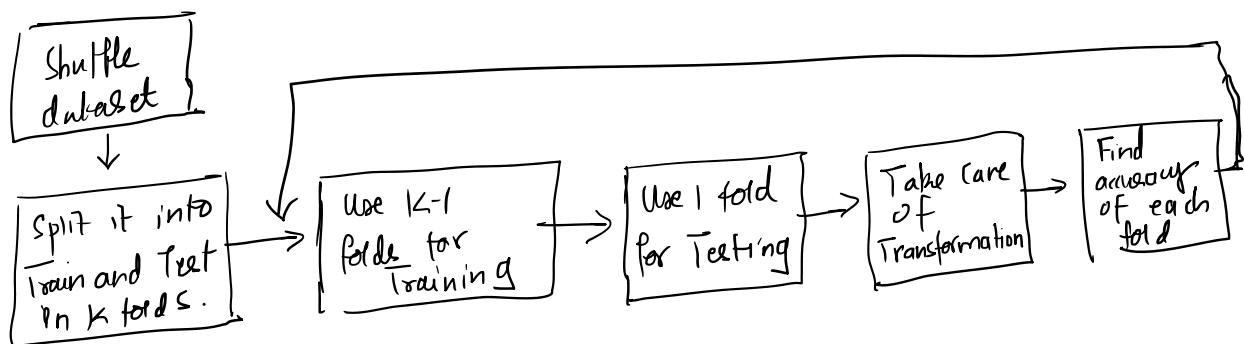


K fold → K fold helps us to build the model in generalized form.

→ To achieve this K fold Cross Validation splits the dataset into Training Testing & Validation.

→ Here Test and Train data will support building the model.

→ Life cycle of K-fold Cross Validation.



- * The No of iterations ideally is K time
- * Finding mean of accuracy score of each iteration will give the consistency of the Trained model.

Rules

$$\textcircled{1} \quad \underline{k \geq 2}$$

if $k=2 \rightarrow$ just 2 iteration.

if $k=n \geq 2 \rightarrow$ $n-1$ for Training
1 for Testing

\textcircled{2} most commonly used value of $\underline{k=10}$

\textcircled{3} If k is very large then the running time of process will increase.

\textcircled{4} The value of k is inversely proportional to size of data i.e if dataset size is small then number of folds can increase.

Bagging \rightarrow Random Forest

* Random Forest is example of Bagging.

* You must know Decision Tree Construction [Gini Index]

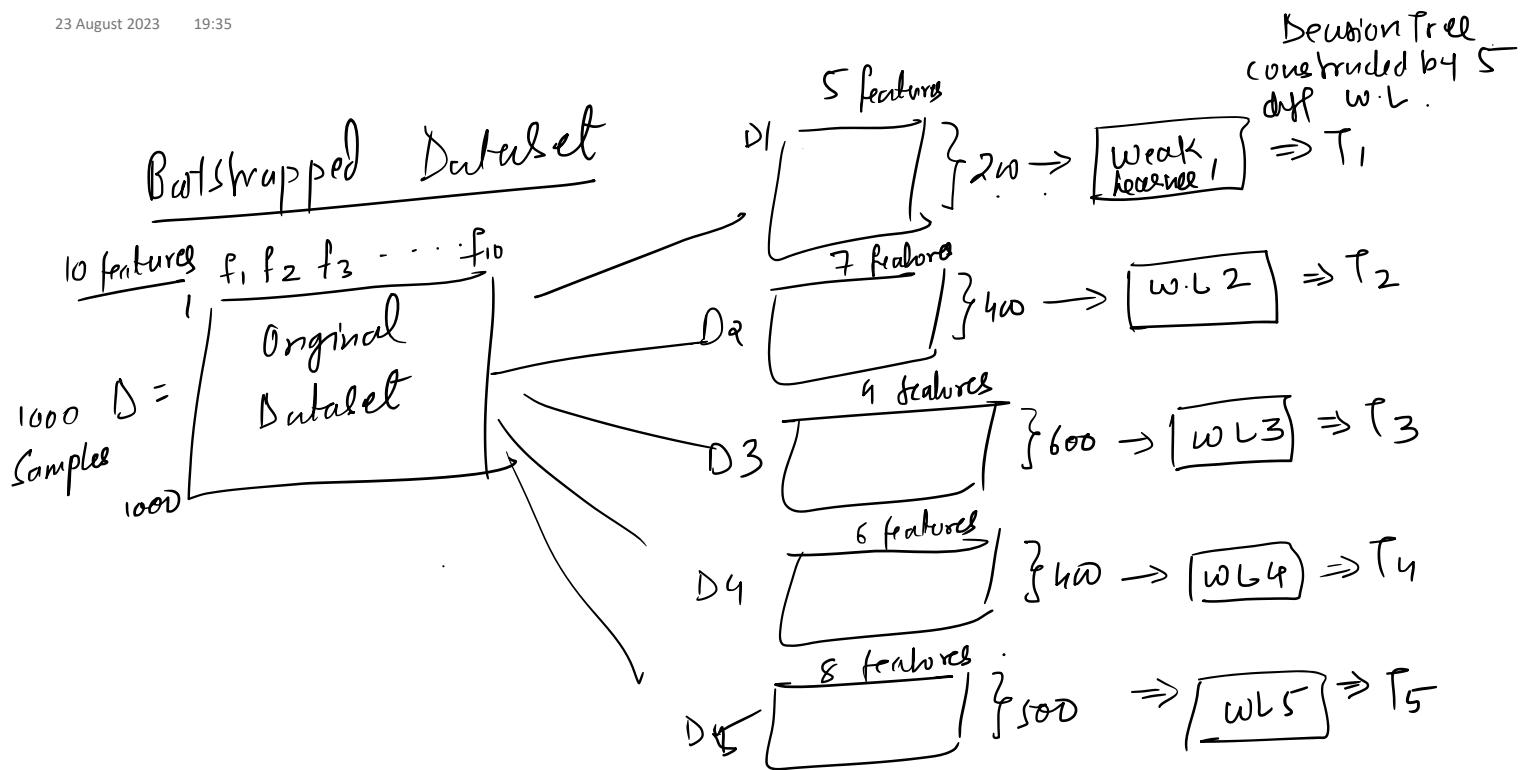
+ Decision Tree

- |
 | Used for Regression \Rightarrow Regression Tree
- |
 | Used for Classification / Categorization \Rightarrow Classification Tree.

* Decision Tree is easy to Build and interpret in practical.

* The Decision Tree is not very accurate on new test sample

* Random Forest Combines Simplicity of Decision Tree with flexibility resulting into vast improvement in accuracy.



We have 5 decision Tree Constructed one each by a weak learner.

Now Test Sample = S_{test}

Now S_{test} is subjected to each of the 5 Decision Tree \Rightarrow

Suppose

$$T_1(S_{test}) = \text{Yes}$$

$$T_2(S_{test}) = \text{No}$$

$$T_3(S_{test}) = \text{Yes}$$

$$T_4(S_{test}) = \text{Yes}$$

$$T_5(S_{test}) = \text{No}$$

In 5 cases 3 cases are Yes
2 cases are No

Final Prediction Yes.

* Construct Random Forest

Original Set

	Chest Pain	Good Blood Circulation	Blocked Arteries	Weight	Heart Disease
S1	No	No	No	125	No
S2	Yes	Yes	Yes	180	Yes
S3	Yes	Yes	No	210	No
S4	Yes	No	Yes	167	Yes

No Yes Yes No ?

Step 1 → Create Bootstrapped Dataset

1. could be / not be of same size
2. Samples are randomly selected
3. Allowed to pick same sample more than once.

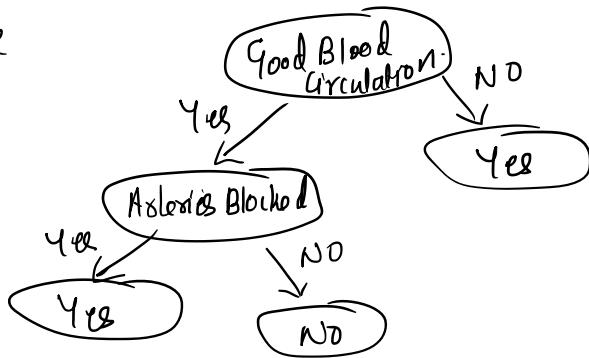
Bootstrapped Dataset

	Chest pain	Good Blood Circulation	Blocked Arteries	Weight	Heart Disease
	Yes	Yes	Yes	150	Yes
	No	No	No	125	No
	Yes	No	Yes	167	Yes
	Yes	No	Yes	167	Yes -

Step 2 → Create a Decision Tree Using Bootstrapped dataset but
Use random subset of Variables (features/columns)

Let us consider we use only 2 features
(Good Blood Circulation and Blocked Arteries)

Let the Tree be
(Gini Index)



Step 3 → go to step 1 and Repeat

- * Ideally we repeat it for 100 times
- * So we have T_1, T_2, \dots, T_{100} (large NO of Trees)
- * Each time the Tree Constructed is a weak learner.
[As all the features and samples are not considered while Tree construction].
- * Now we have Random Forest of 100 Trees and will be more effective than Individual Decision Tree.

⑥ How To Use the Random Forest (Here 100 Trees) →

Consider a Test Sample

Chest pain	Good Blood Circulation	Blocked Arteries	Weight	Heart Disease
Yes	No	No	168	?

Here we will run it through each of 100 Trees and

will note the Prediction

Let say of 100 Trees

80 Trees Predicted Yes

20 Trees Predicted No

Answer = Yes

Answer = Yes

Q) To find Random Forest is effective or not ?

- * There might be some Sample not considered by any of Bootstrapped Dataset.
- * We will create a New Dataset with such samples. This is known as "Out of Bag" Dataset.
- * Now we will chk how many samples from Out of Bag Dataset are predicted correctly.
- * Numbers of Incorrectly Predicted Out of Bag Samples = Out of Bag Error.

Q) How to decide on Number of Columns to use for Building Trees in Random Forest ?

① In our Case DataSet has 4 features

→ Create a Random Forest of Trees using 4 Features = F_1
→ Create a Random Forest of Trees using 3 Features = F_2

100 Tree

100 Tree

② Now chk the Accuracy of Out of Bag Data on each of the Random Forest

③ Use the one that gives More Accurate Answer.

Summarize

Step 1: Create a Bootstrapped Data Set

Step 2: Create a Decision Tree using Subset of features

Step 3: Repeat Step 1 and 2 approx 100 times.

Step 4: Use Out of Bag Data Set to determine Out of Bag Error

Step 5: If Out of Bag Error is high then Repeat Step 1 to Step 4 until Out of Bag Error is considerably low.

Note → To compare the Performance for diff Random Forest Create diff Random Forest using different number of features

Use the one that gives Highest Accuracy.

ADABOOST (ADAPTIVE BOOST)

→ (Tree Stumping)

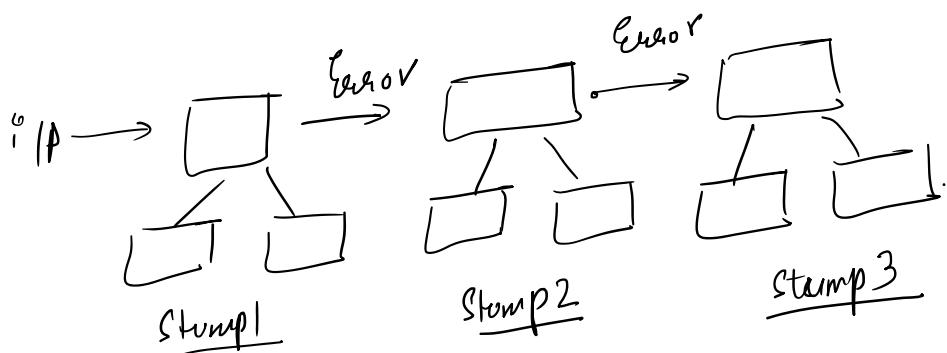
- * Adaboost creates forest of Trees from scratch and then it uses it to make classification.
- * [In Random Forest everytime we make a Tree with different depth (depending on number of features considered).]
- * There is no predetermined maximum depth of Tree.]
- * In ADABOOST, the forest of Trees are usually just node and Two leaves

```

graph TD
    Node[ ] --> Leaf1[ ]
    Node --> Leaf2[ ]
    style Node fill:none,stroke:none
    style Leaf1 fill:none,stroke:none
    style Leaf2 fill:none,stroke:none
  
```
- * A Tree with one node and two leaves only is known as Tree Stump.
- * In ADABOOST we have forest of Stump instead of Tree.
- * Stump alone is not great for classification as it takes only one parameter to make decision.
- * Stumps are weak learners.
- * ADABOOST combines many stumps -
- * In Random forest each tree contributes equally in final decision.

final decision

- * In ADABOOST, some stump may get more say in final decision than others.
- * In Random Forest the decision made by Trees are independent to other, so the order of tree generation is not important
- * In ADABOOST, we have forest of Stump and hence order is Important.
- * Here the error of first stump is corrected by next stump and so on -



→ Here error of one stump influences the other stump
So the stumps are generated in sequence.

Idea Behind ADABOOST. (Boosting Technique \rightarrow reduce Bias)

1. Adaboost combines lot of weak learners to make classification
2. The weak learners are almost always Stump.
- ... Some will have more say in classification

in the run-

3. Some stump will have more say in classification than other stump.
4. Each stump is made by taking previous stump mistakes into account

Working of AdaBoost with Example.

Consider

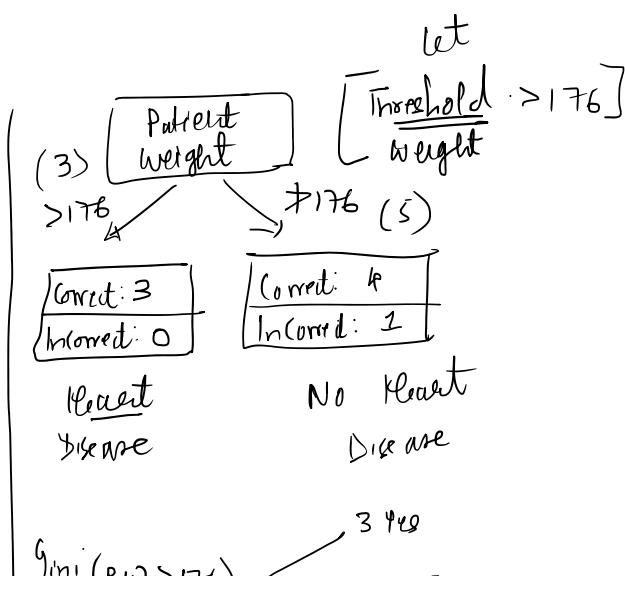
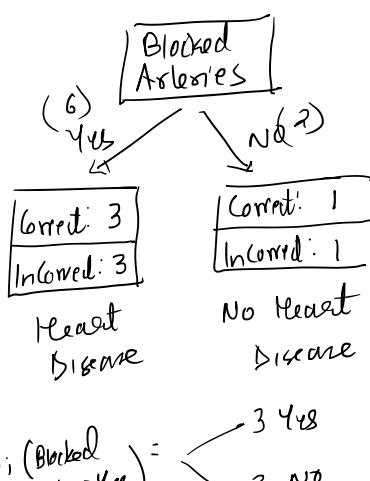
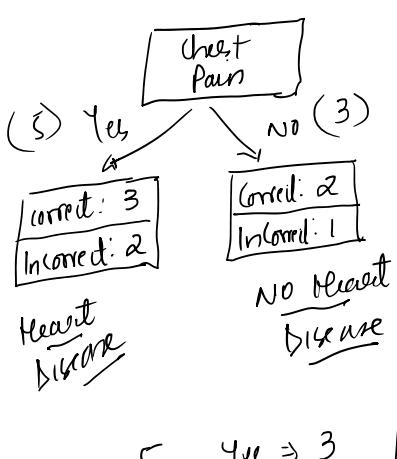
Chest Pain	Blocked Arteries	Patient Weight	Heart Disease
Yes	Yes	205	Yes
No	Yes	180	Yes
Yes	No	210	Yes
Yes	Yes	167	Yes
No	Yes	156	No
No	Yes	125	No
Yes	No	168	No
Yes	Yes	172	No

Step 1: Assign weight to every sample = $\frac{1}{\text{Total No of Sample}} = \frac{1}{8}$
(Initial)

* (Initially equal weights are assigned to each sample)

Chest Pain	Blocked Arteries	Patient Weight	Heart Disease	Sample weight
Yes	Yes	205	Yes	$\frac{1}{8}$
No	Yes	180	Yes	$\frac{1}{8}$
Yes	No	210	Yes	$\frac{1}{8}$
Yes	Yes	167	Yes	$\frac{1}{8}$
No	Yes	156	No	$\frac{1}{8}$
No	Yes	125	No	$\frac{1}{8}$
Yes	No	168	No	$\frac{1}{8}$
Yes	Yes	172	No	$\frac{1}{8}$

Step 2: Create stamp on each feature



D¹⁴

$Gini'(chest pain = \underline{Yes}) = \begin{cases} 4/8 \Rightarrow 3 \\ \underline{No} \Rightarrow 2 \end{cases}$ <p style="text-align: center;">\uparrow Recent Disease</p> $= 1 - ((3/8)^2 + (5/8)^2)$ $= \underline{0.48}$ $Gini'(chest pain = \underline{No}) = \begin{cases} 4/8 = 2 \\ 3 \underline{No} = 1 \end{cases}$ $= 1 - ((2/3)^2 + (1/3)^2)$ $= \underline{0.444}$	$Gini'(\text{Blocked Arteries} = \underline{Yes}) = \begin{cases} 3 \underline{Yes} \\ 3 \underline{No} \end{cases}$ $= 1 - ((3/6)^2 + (3/6)^2)$ $= \underline{0.5}$ $Gini'(\text{Blocked Arteries} = \underline{No}) = \begin{cases} 1 \underline{Yes} \\ 1 \underline{No} \end{cases}$ $= 1 - ((1/2)^2 + (1/2)^2)$ $= \underline{0.5}$ <div style="border: 1px solid black; padding: 5px; display: inline-block;"> $Gini'(\text{Blocked Arteries}) = 0.5$ </div>	$Gini'(\text{P.W} > 176) = \begin{cases} 3 \underline{Yes} \\ 0 \underline{No} \end{cases}$ $= 1 - ((3/3)^2)$ $= \underline{0}$ $Gini'(\text{P.W} \neq 176) = \begin{cases} 4 \underline{Yes} \\ 1 \underline{No} \end{cases}$ $= 1 - ((4/5)^2 + (1/5)^2)$ $= \underline{0.32}$ <div style="border: 1px solid black; padding: 5px; display: inline-block;"> $Gini'(\text{Weight}) = 0.2$ </div>
--	---	---

weighted $Gini^0(\text{chest pain})$

$$= 0.48 * 5/8 + 0.44 * 3/8$$

$$= \underline{0.4665} \approx \underline{\underline{0.47}}$$

* The smallest Gini Index is $\underline{0.2}$ for $\underline{\underline{\text{weight}}} > 176$ feature

So first Stump will be $\underline{\underline{\text{weight}}} > 176$

Step 3 → To calculate Amount of say of $\underline{\underline{\text{weight}}} > 176$ [feature].

$$\text{Amount of Say} = \frac{1}{2} \ln \left(\frac{1 - \text{Total Error}}{\text{Total Error}} \right)$$

Total Error: for a stamp the sum of weight associated with incorrectly classified sample.

Here only one sample is incorrectly classified.

Yes	Yes	167	Yes	1/8
-----	-----	-----	-----	-----

$$\therefore \underline{\underline{\text{Total Error}}} = \underline{\underline{1/8}}$$

$$* \quad \text{Amt of Say} = \frac{1}{2} \ln \left(\frac{1 - \frac{1}{8}}{\frac{1}{8}} \right) = \frac{0.97}{(\text{High})}$$

[Note : If Amt of say is less then the stump is not doing good job]

- * now we have \rightarrow first stump (weight > 176)
 - \rightarrow Amt of say.

Now initially every sample had equal weight \rightarrow

Note Before second stump we must modify weights of samples

- * Weights of correctly classified samples \Rightarrow must be decrease
- " " incorrectly classified samples \Rightarrow must be increase
- [So that next stump will focus more on incorrectly identified samples].

$$\text{New Sample weight (correctly classified)} = \text{Sample weight (old)} \times e^{-\text{amount of say}}$$

$$\text{New Sample weight (incorrectly classified)} = \text{Sample weight (old)} \times e^{+\text{amount of say}}$$

$$\therefore \text{for Incorrectly Classified Sample} = \frac{1}{8} \times e^{+0.97} = 0.33$$

$$\text{For Correctly Classified Sample} = \frac{1}{8} \times e^{-0.97} = 0.05$$

Chest Pain	Blocked Arteries	Patient weight	Heart Disease	Sample weight	New Sample weight	Normalized Sample weight
Yes.	Yes	205	Yes	$\frac{1}{8}$	0.05	$0.05/0.68 = 0.07$
No.	Yes	180	Yes	$\frac{1}{8}$	0.05	0.07
Yes.	No	210	Yes	$\frac{1}{8}$	0.05	0.07
Yes.	Yes	167	Yes	$\frac{1}{8}$	0.33	$0.33/0.68 = 0.49$

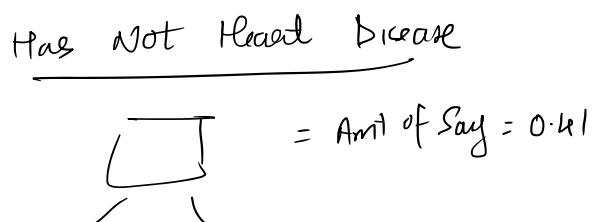
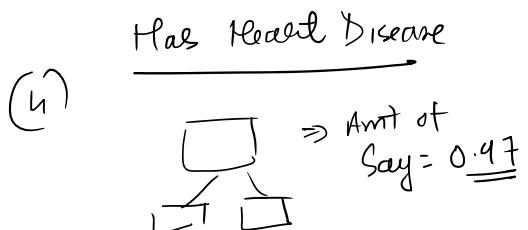
	No	210	Yes	$\frac{1}{8}$	0.05	
✓	Yes	48	(167)	<u>Yes</u>	<u>0.33</u>	$0.33 / 0.68 = 0.49$
✓	No	48	156	No	0.05	0.07
✓	No	48	125	No	0.05	0.07
.	Yes.	No	168	No	0.05	0.07
.	Yes	48	172	No	0.05	0.07
						$\uparrow \text{Sum} = 0.68$

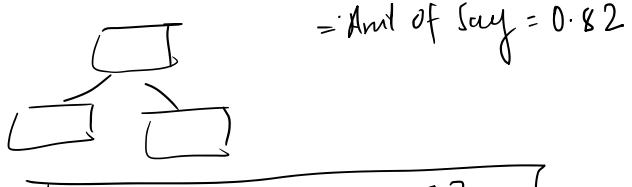
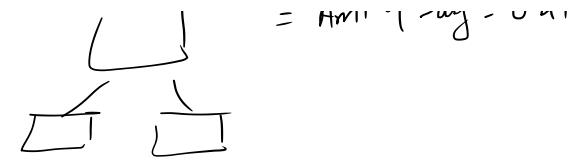
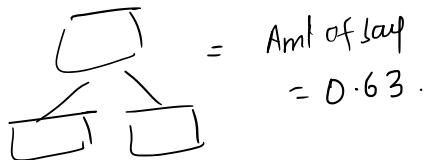
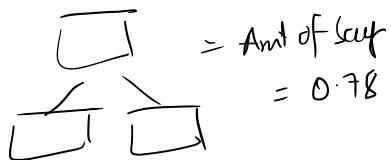
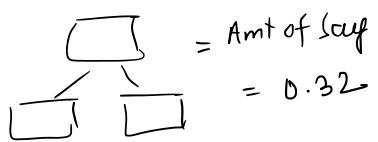
* Now we have New Normalized Sample weights:

- * The Incorrectly classified Sample has Higher New Sample weight
→ which means Stump 2 will focus more on it.
- * The Correctly classified Sample has Comparatively lower Sample weight
→ which means Stump 2 will focus less on it.
- * Similarly we can have Many such Stumps
- * This is how Adaboost creates and uses Stumps.

* Let us chk how forest of stumps made by Adaboost does Classification →

- ① let us consider a test sample
- ② we will pass the test sample through each stump
- ③ let there be few stumps that classifies Sample as Yes (^{Has Heart Disease}) and few stumps classifies Sample as No (^{Do not have heart disease}).





Total Amt of Say (No) = 1.23

Total Amt of Say for stumps that classifies as Heart Disease
(Yes)

is greater

∴ The Test Sample is classified as Yes (Has Heart Disease).

Summary of Adaboost

① Assign sample weight (Initial equal)

② Create stump for each feature.

③ Use Gini Index to identify first stump.

④ For first stump:

① Calculate Total Error

② Amt of Say

④ Use "j" to train the stump.

⑤ For first stump:

① Calculate Total Error

② Amt of Say

⑥ Update the sample weights for each sample

$$\text{For Incorrectly classified weight} = \frac{\text{updated weight}}{\text{old weight}} + \text{Amt of Say}$$

$$\text{For correctly classified weight} = \frac{\text{updated weight}}{\text{old weight}} - \text{Amt of Say}$$

⑦ Now Specify updated weight for each sample & Normalize the weight to generate New Sample weight.

⑧ Identity new stump based on New Sample Weight & Repeat step 3 to 7

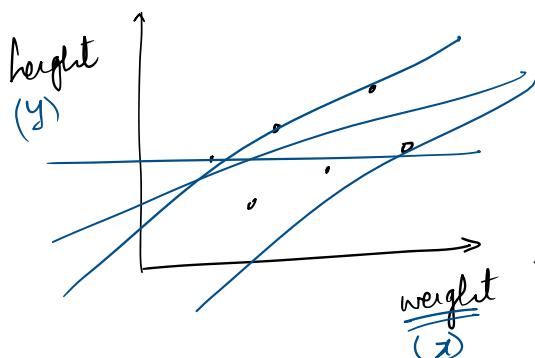
⑨ For Classification (Testing).

8.1 \Rightarrow Run the test sample through all the stumps.

8.2 \Rightarrow Calculate Amt of Say for sample classifying Yes & No

8.3 \Rightarrow Classify the test sample based on largest sum of Amount of Say

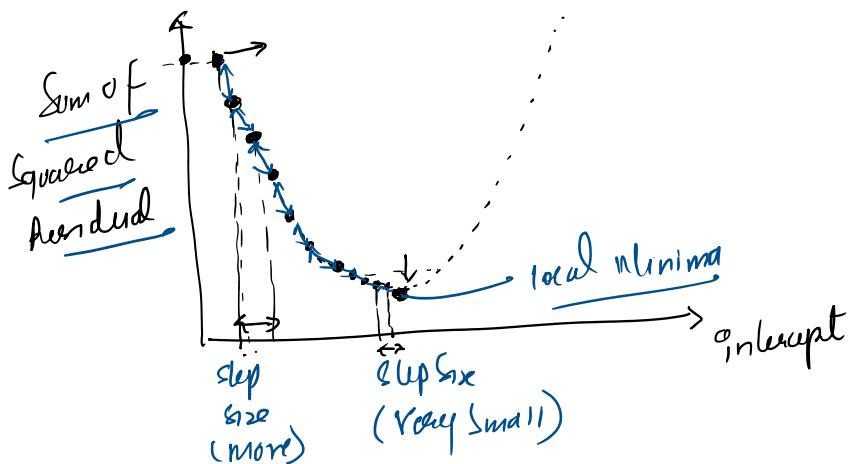
Gradient Descent



$$\text{predicted height} = \underset{\text{parameters}}{\underset{\text{intercept}}{c}} + \underset{\text{slope}}{m} \underset{x}{\underset{\text{weight}}{x}}$$

* Here we can have many line with different intercept and slope values

- * We want a line with intercept and slope such that the Cost function should give Minimum value.
- * Let the Cost function be Sum of Squared Error (Residual error)
- * Let us Consider only intercept
- * Let's plot a graph of Intercept v/s Sum of Squared Residual Error.



To find Minimum value for Sum of Squared Residual Error, we will start with very small value of intercept and will increase the value by very small amount.

x	y
0.5	1.4
2.3	1.9
	2.2

Let the Cost = Sum of Squared Error
F^n be (Residual).

Ez

0.5	1.4	Let the loss = sum of squares F^n be (Residual).
2.3	1.9	
2.9	3.2	
		↑ Actual

$$\text{Sum of Squared Residual Error} = \left(1.4 - (\text{Intercept} + \text{slope} * 0.5) \right)^2 \\ + \left(1.9 - (\text{Intercept} + \text{slope} * 2.3) \right)^2 \\ + \left(3.2 - (\text{Intercept} + \text{slope} * 2.9) \right)^2$$

Here we want to find value of Intercept and Slope that give minimum Sum of Squared Error.

Step 1 → Find $\frac{d}{d\text{Intercept}}$ (Sum of Squared Residual Error)

$\frac{d}{d\text{slope}}$ (Sum of Squared Residual Error).

$$\text{Sum of Squared Residual Error} = \left(1.4 - (\text{Intercept} + \text{slope} * 0.5) \right)^2 \\ + \left(1.9 - (\text{Intercept} + \text{slope} * 2.3) \right)^2 \\ + \left(3.2 - (\text{Intercept} + \text{slope} * 2.9) \right)^2$$

Using chain Rule $\left[\frac{d}{dx} x^2 = 2x \frac{d}{dx} x \right]$

$$\frac{d(\text{Sum of Squared Error})}{d\text{Intercept}} = 2 * \left(1.4 - (\text{Intercept} + \text{slope} * 0.5) \right) \\ * \frac{d}{d\text{Intercept}} \left(1.4 - (\text{Intercept} + \text{slope} * 0.5) \right) \quad (-1) \\ + 2 * \left(1.9 - (\text{Intercept} + \text{slope} * 2.3) \right) \\ \frac{d}{d\text{Intercept}} \left(1.9 - (\text{Intercept} + \text{slope} * 2.3) \right) \quad (-1) \\ + 2 * \left(3.2 - (\text{Intercept} + \text{slope} * 2.9) \right) \\ \frac{d}{d\text{Intercept}} \left(3.2 - (\text{Intercept} + \text{slope} * 2.9) \right). \quad (-1)$$

$$\frac{d}{d \text{Intercept}} (3.2 - (\text{Intercept} + \text{slope} \times 2.9)) = (-1)$$

Eq 1.

$$\begin{aligned} \frac{d(\text{SS } \epsilon)}{d \text{Intercept}} &= (-2)(1.4 - (\text{Intercept} + \text{slope} \times 0.5)) + \\ &(-2)(1.9 - (\text{Intercept} + \text{slope} \times 2.3)) + \\ &(-2)(3.2 - (\text{Intercept} + \text{slope} \times 2.9)). \end{aligned}$$

$$\begin{aligned} \frac{\text{Sum of Squared Residual Error}}{(\text{SS } \epsilon)} &= \left(1.4 - (\text{Intercept} + \text{slope} \times 0.5)\right)^2 \\ &+ \left(1.9 - (\text{Intercept} + \text{slope} \times 2.3)\right)^2 \\ &+ \left(3.2 - (\text{Intercept} + \text{slope} \times 2.9)\right)^2 \end{aligned}$$

Note

$$\begin{aligned} \frac{d}{d \text{slope}} (\text{SS } \epsilon) &= \frac{d}{d \text{slope}} (1.4 - (\text{Intercept} + \text{slope} \times 0.5))^2 \\ &+ \frac{d}{d \text{slope}} (1.9 - (\text{Intercept} + \text{slope} \times 2.3))^2 \\ &+ \frac{d}{d \text{slope}} (3.2 - (\text{Intercept} + \text{slope} \times 2.9))^2. \end{aligned}$$

Eq 2.

$$\begin{aligned} \frac{d}{d \text{slope}} (\text{SS } \epsilon) &= \frac{(-2)(0.5)(1.4 - (\text{Intercept} + \text{slope} \times 0.5))}{+ (-2)(2.3)(1.9 - (\text{Intercept} + \text{slope} \times 2.3))} \\ &+ (-2)(2.9)(3.2 - (\text{Intercept} + \text{slope} \times 2.9)). \end{aligned}$$

We got the above 2 equations.
Let start with random values of Intercept & slope

Let

$\text{Intercept} = 0$
$\text{slope} = 1$

and $\text{Intercept} = 0$ & $\text{slope} = 1$ in Eq 1 & Eq 2

Substitute $\hat{y}_{\text{intercept}} = 0$ & $\underline{\text{slope}} = 1$ in Eq 1 & Eq 2

$$\frac{d(\text{SSE})}{d \text{slope}} = -0.8 \Rightarrow \text{slope of line w.r.t parameter slope}$$

$$\frac{d(\text{SSE})}{d \text{intercept}} = -1.6 \Rightarrow \text{slope of line w.r.t parameter } \underline{\text{intercept}}$$

For step size

$$\text{New Step Size} = \text{slope of line at point} * \text{learning Rate}$$

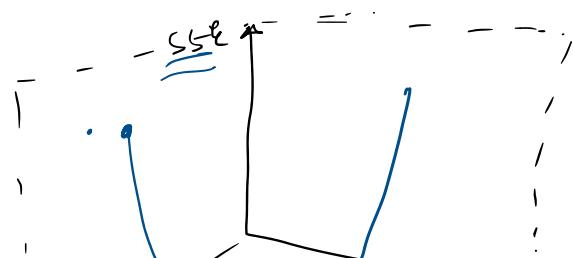
* Learning Rate \Rightarrow takes smaller value [$\underline{10^{-6}}$ to 1.0]

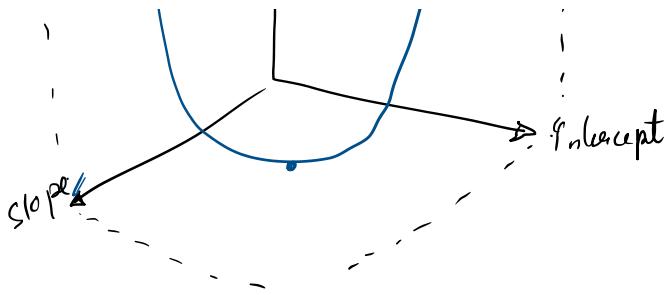
Let us take learning rate = 0.01

$$\begin{aligned} \text{Step Size (Intercept)} &= \frac{d(\text{SSE})}{d \text{intercept}} * 0.01 = -1.6 * 0.01 = -0.016 \\ \text{Step Size (slope)} &= \frac{d(\text{SSE})}{d \text{slope}} * 0.01 = -0.8 * 0.01 = -0.008 \end{aligned}$$

$$\text{New Intercept} = \text{old } \hat{y}_{\text{intercept}} - \frac{\text{Step Size (Intercept)}}{0.016} = 0 - (-0.016) = \underline{0.016}$$

$$\text{New Slope} = \text{old Slope} - \frac{\text{Step Size (slope)}}{0.008} = 1 - (-0.008) = \underline{1.008}$$





Substitute the new $y_{\text{intercept}}$ and slope in $\frac{d(\text{SSE})}{d y_{\text{intercept}}}$ & $\frac{d(\text{SSE})}{d \text{slope}}$

and Calculate new step size for $y_{\text{intercept}}$ & slope

And again Calculate New $y_{\text{intercept}}$ & New slope]

* We will repeat until the step size becomes very small
 $\Rightarrow \underline{\underline{0.001}}$

or Some maximum Number of steps is reached Ex: 1000

[For above Dataset Best fitting line will have
 $y_{\text{intercept}} = 0.95$ & slope = 0.64]

[This is How Gradient Descent Optimizes Parameters]

Summary → Gradient Descent

① get eqn of line [Identify Simple linear or Polynomial]

② Want $y_{\text{intercept}}$ & slope value that minimizes cost J^*

③ Identify cost (loss) function. [We took SSE]

④ Find Eq for $\frac{d(\text{SSE})}{d y_{\text{intercept}}}$ & $\frac{d(\text{SSE})}{d \text{slope}}$

→ ⑤ Start with some initial value of slope & $y_{\text{intercept}}$.

⑥ Find the value of $\frac{d(\text{SSE})}{d y_{\text{intercept}}}$ & $\frac{d(\text{SSE})}{d \text{slope}}$

(6) Find the value of $\frac{d(\text{SSE})}{d\text{intercept}} \times \frac{d(\text{SSE})}{d\text{slope}}$

(7) Calculate $\text{StepSize}_{(\text{intercept})} = \frac{d(\text{SSE})}{d\text{intercept}} \times \text{learning rate}$

$$\text{StepSize}_{(\text{slope})} = \frac{d(\text{SSE})}{d\text{slope}} \times \text{learning rate}$$

(8) Update $\underline{\text{intercept}} = \text{old intercept} - \text{stepsize}_{(\text{intercept})}$

update $\underline{\text{slope}} = \text{old slope} - \text{stepsize}_{(\text{slope})}$

(9) Repeat step 6 - 8 until stepsize is very small
or Number of iterations $\leq \underline{1000}$

Gradient Boost for Regression

- * Gradient Boost start with single leaf instead of Tree/stump.
- * Leaf represents initial guess for the weights of the samples -

Consider

Height	Favorite color	Gender	<u>Weight</u> (Initial)
1.6	Blue	male	88
1.6	Green	Female	76
1.5	Blue	Female	56
1.8	Red	male	73
1.5	Green	male	77
1.4	Blue	Female	57

Step 1 → First let make initial guess = weight (Average) = 71.2

Step 2 → Like Adaboost, even Gradient Boost builds tree based on previous tree

But unlike Adaboost, the gradient boost tree are larger than stump

[However the height is still restricted [No of leaf 8 → 32] in each tree]

* Gradient Boost will build Another Tree based on Error of previous Tree.

+ [Gradient Boost continuous to build Tree in this fashion until it has made number of Trees you asked for or the Additional tree fails to improve the fit]

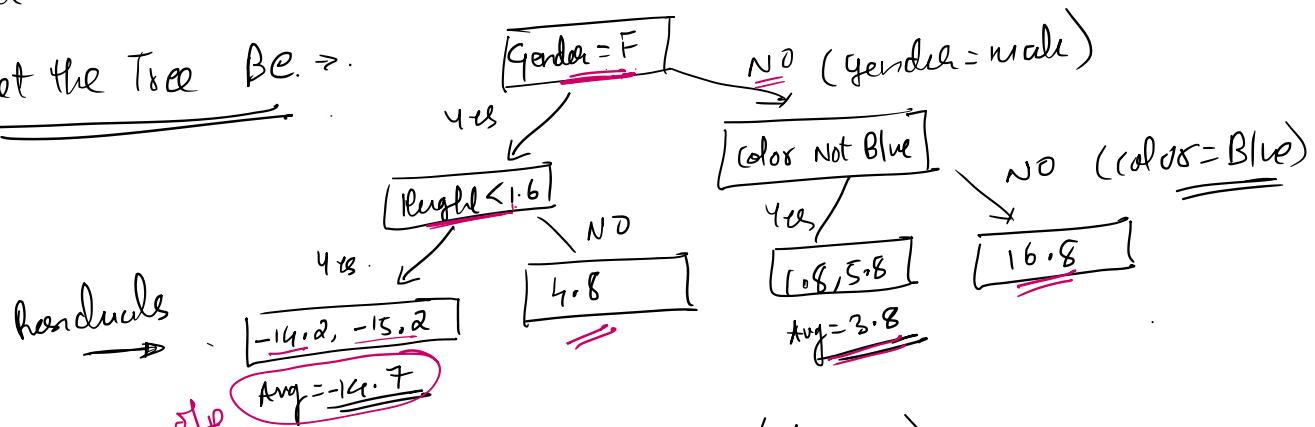
\Rightarrow Step 3 : Initial Guess = weight (Average) = 71.2

$$\text{Predicted Value} = \boxed{71.2} + \frac{\text{Tree1}}{\uparrow \text{Initial guess}} + \frac{\text{Tree2}}{\uparrow \text{Boosting}}$$

Height	Favorite Color	Gender	Weight	Residual	(Observed - Average)
1.6	Blue	male	88	16.8	
1.6	Green	Female	76	-4.8	
1.5	Blue	Female	58	-15.2	
1.8	Red	male	73	1.8	
1.5	Green	male	77	5.8	
1.4	Blue	Female	57	-14.2	

Let Build a Tree using Residual and Not weight

Let the Tree Be. \Rightarrow



If there will be only one Tree (Assume).

Then:

$$\text{Predicted value} = \text{Initial Guess} + \text{Learning Rate} * \text{Residual}$$

$\text{Initial Guess} = \boxed{71.2}$
 $\text{Learning Rate} = \boxed{0.1}$
 $\text{Residual} = \boxed{16.8}$

Decision tree structure:

```

graph TD
    Root["Gender = F"] -- NO --> Node1["Gender = male"]
    Root -- YES --> Node2["color not Blue"]
    Node1 -- NO --> Node3["Height < 1.6"]
    Node1 -- YES --> Node4["4.8"]
    Node2 -- NO --> Node5["1.8, 5.8"]
    Node2 -- YES --> Node6["Avg = 3.8"]
    Node3 -- NO --> Node7["-14.2, -15.2"]
    Node3 -- YES --> Node8["Avg = -14.7"]
  
```

Now first Sample = $\boxed{1.6 \text{ | Blue | male | ? }}$

$$0.1 * 16.8 - 71.2 + 0.1 * 16.8 = \boxed{72.88}$$

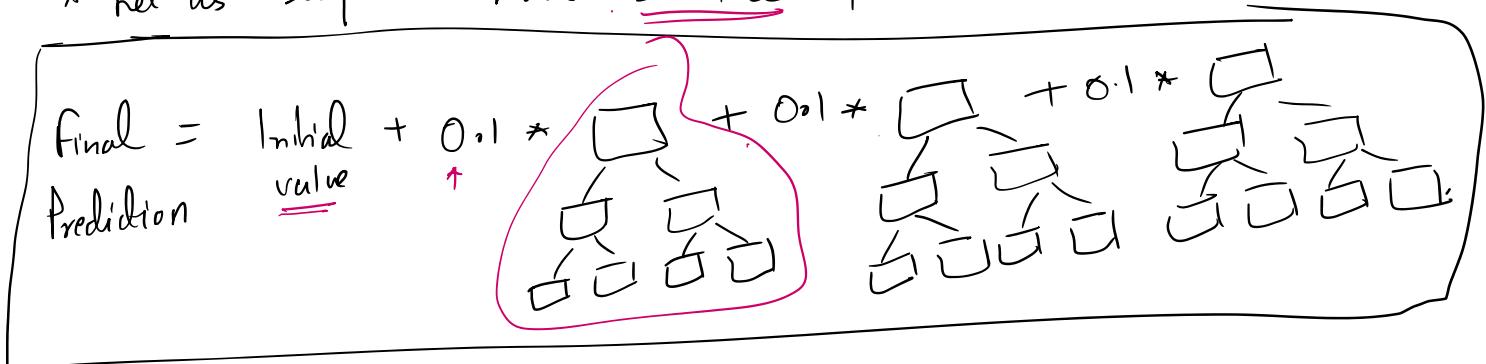
$$\text{Predicted value} = 71.2 + 0.1 * 16.8 = \underline{\underline{72.88}}$$

* Repeat this for Every Sample & Calculate Predicted Value for Each Sample
 → Call it as New Prediction.

Height	Favorite color	Gender	Weight	Residual	New prediction	New Residual
1.6	Blue	male	58	16.8	72.88	R ₁
1.6.	Green	Female	76	-48	P ₁	R ₂
1.5	Blue	Female	58	-15.2	P ₂	R ₃
1.8	Red	male	73	1.8	P ₃	R ₄
1.5	Green	male	77	5.6	P ₄	R ₅
1.4	Blue	Female?	57	-14.2	P ₅	R ₆

- * Now find New Residual for New Prediction for all the Samples
- * Using this New Residual Construct New Tree and So on
- Until the New Tree does not make significant contribution to the New Prediction.

* Let us say we have 3 Tree for above Dataset →



[Gradient Boost Principle → Taking lot of small steps in Right Direction will result into better predictions]

Gradient Boost for Classification →

Consider

Likes Popcorn	Age	Favorite color	loves Javan movie
Yes	12	Blue	Yes (1)
Yes	87	Green	Yes (1)
No	44	Blue	No (0)
Yes	19	Red	No (0)
No	32	Green	Yes (1)
No	14	Blue	Yes (1)

↑
Actual values

↓ Probability

Step 1: Start with Initial Prediction →

[note: When we use gradient boost for classification, the initial prediction for every individual sample → log(odd)]

$$\text{odd} = \frac{\text{No of Yes}}{\text{No of No}} = \frac{\Phi}{\alpha} = 2$$

$$\log_e(\text{odd}) = \log_e(2) = 0.7$$

To use this we have to convert it into Probability =

$$\text{Probability of loving Javan movie} = \frac{1}{1 + e^{-\log(\text{odd})}} = \frac{1}{1 + e^{-0.7}} \approx \frac{0.7}{\Phi}$$

So Initial Prediction for Probability of all the samples

for loving Javan movie = 0.7

Initial Prediction 0.7

Likes Popcorn	Age	Favorite color	loves Javan movie	Residual Error
Yes	12	Blue	Yes (1)	0.3
Yes	87	Green	Yes (1)	0.3
No	44	Blue	No (0)	-0.7

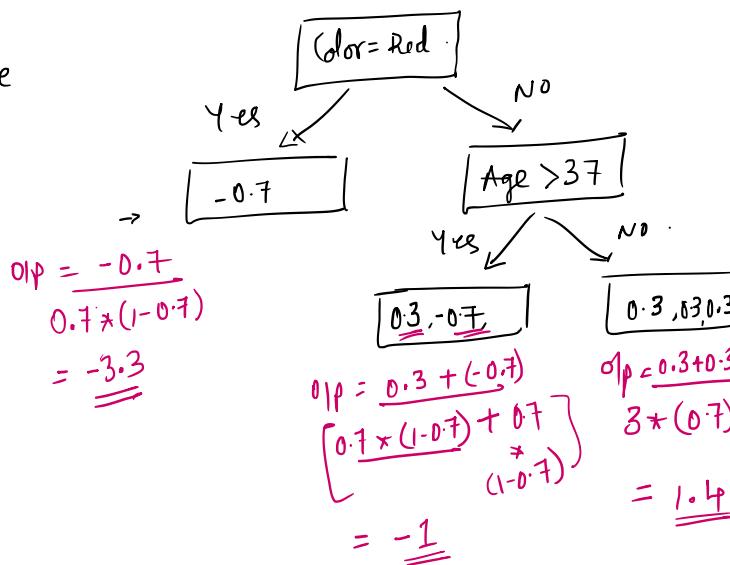
Yes	87	Green	Yes (1)	<u>0.7</u>
No	418	Blue	No (0)	<u>-0.7</u>
Yes	19	Red	No (0)	<u>-0.7</u>
No	32	Green	Yes (1)	<u>0.3</u>
No	14	Blue	Yes (1)	<u>0.3</u>

Atrial

(Initial residual error)

Now let us Build Tree Color And Age

Let the Tree be



For Classification:

0lp value at every leaf

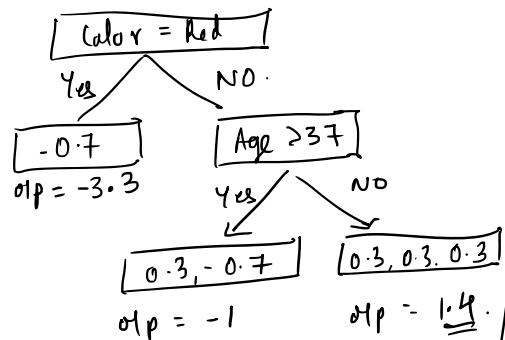
$$= \frac{\sum \text{Residual}}{\sum}$$

$$\sum (\text{Previous Prob} * (1 - \text{Previous Prob}))$$

Here this time Previous Prob =

Initial Prediction

$$\text{New Prediction} = \text{Initial Prediction} + 0.8 * (\text{Learning Rate})$$



Consider first sample

the popcorn			Initial
Age	Color	Initial	
Yes 12 Blue		0.7	

$$\text{New Prediction} = 0.7 + 0.8 * 1.4 = 1.8$$

$$\text{Convert into Probability} = \frac{1}{1+e^{-1.8}} = 0.9$$

Similarly calculate for all the other samples -

Now the New Prediction be : \rightarrow Initial Prediction + Residual Factor

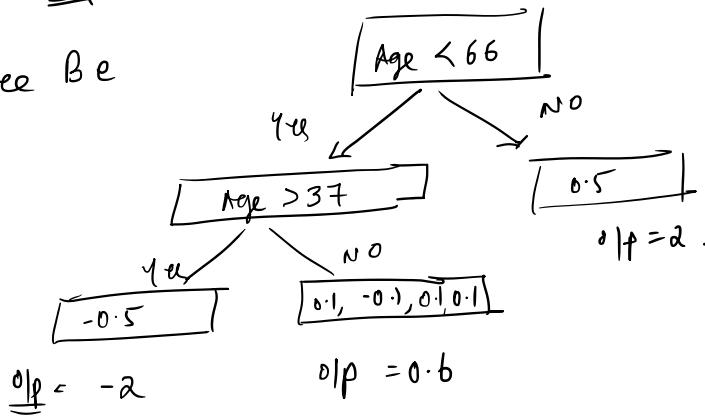
Likes Popcorn	Age	Favourite color	Does Javan movie	<u>Residual Factor</u>	New Prediction	New Residual
Yes	12	Blue	Yes (1)	0.3	0.9	0.1
Yes	87	Green	Yes (1)	0.3	0.5	-0.5
No	44	Blue	No (0)	-0.7	0.5	-0.1
Yes	19	Red	No (0)	-0.7	0.1	-0.1
No	32	Green	Yes (1)	0.3	0.9	0.1
No	14	Blue	Yes (1)	0.3	0.9	0.1

Actual values

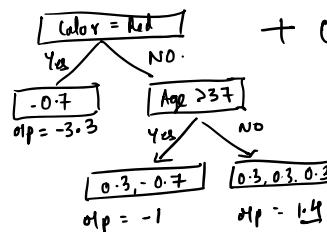
Here New prediction is not same for all the samples -

- * Here New prediction is not same for all the samples
- * Each sample has its own prediction and calculated New residual for each sample.
- * Again on the basis of New Residual let us Build a Tree Using only Age only

Let the Tree Be

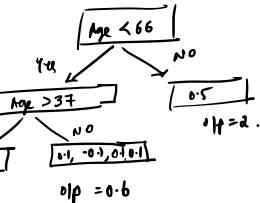


$$\therefore \text{New Prediction} = \text{Initial Prediction} + 0.8 * (\text{Learning Rate})$$



$$+ 0.8 *$$

$$+ 0.8 *$$



* Apply each sample in database on the above model and calculate New Predictions

Apply ... any ... now ...

New Predictions

* Again calculate New Residual and so on -

* This process repeats until we have made the maximum specified number of Trees or the residual gets very very small.

Summary : Gradient Boost for Classification

① Consider the dataset

② find log(odd) on the target

③ find Initial Probability = $\frac{1}{1 + e^{-\text{log(odd)}}}$.

④ Now calculate Initial Residual based on Initial Probability

⑤ Build a Tree using some of the features and Initial Residual (-Tree1)

⑥ Calculate New Prediction for = Initial + learning * Tree1
Every Sample Prediction Rate

⑦ Using New Predictive Calculate New Residual

⑧ Construct a new Tree based on some feature and New Residual and so on.

⑨ Repeat until max no of Trees are constructed or the residuals are insignificant.

Summary : Gradient Boost for Regression

① Consider the dataset

② find Initial Prediction = Avg of observed value.

③ Now calculate Initial Residual based on Initial Prediction

④ Build a Tree using some of the features and Initial Residual (-Tree1)

⑤ Calculate New Prediction for = Initial + learning * Tree1
Every Sample Prediction Rate

- ⑤ Calculate New Prediction tree = Initial + accuracy * tree
Every Sample Prediction Rate
- ⑥ Using New Predictions Calculate New Residual.
- ⑦ Construct a new Tree based on some feature and New Residual and so on.
- ⑧ Repeat until max no of Trees are constructed or the residuals are insignificant.

XGBoost (Extreme Gradient Boost)

(1) XGBoost for Regression

Consider

Dosage of Drug	Effectiveness of Drug
10	-10
20	7
25	8
35	-7

y ✓ (Actual)

for Classification

Dosage	Effectiveness
10	0
20	1
25	1
35	0

✓

||

Build Prediction Model Using XGBoost

Step 1: Assume Initial threshold for effectiveness = 0.5

Step 2: Calculate Residual

$$\text{Initial Prediction} = 0.5$$

Dosage of Drug	Effectiveness of Drug	Residual
10	-10	-10.5
20	7	6.5
25	8	7.5
35	-7	-7.5

↑

Step 3: Consider all the Residuals in Root initially.

$$\checkmark [-10.5, 6.5, 7.5, -7.5]$$

$$\text{Similarity} = \frac{(-10.5 + 6.5)^2}{4 + 0} = 4$$

$$\underline{\text{Note}} \\ \underline{\text{Score}} \\ \text{Similarity} = \frac{(\text{Sum of Residual})^2}{\text{No of Residual} + \lambda}$$

[λ = Regularization Parameter]
→ prevents Overfitting

$$\text{gain} = \text{Similarity}(\text{left}) + \text{Similarity}(\text{right}) - \text{Similarity}(\text{Root})$$

Now let us see if we can cluster the residual better on Similarity by split

Determine dosage value for split →

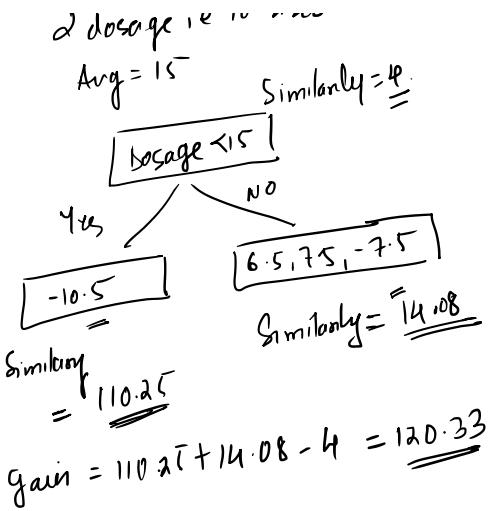
Case 1: Consider Average of first 2 dosage i.e 10 & 20

$$\text{Avg} = 15$$

$$\text{Similarity} = 10$$

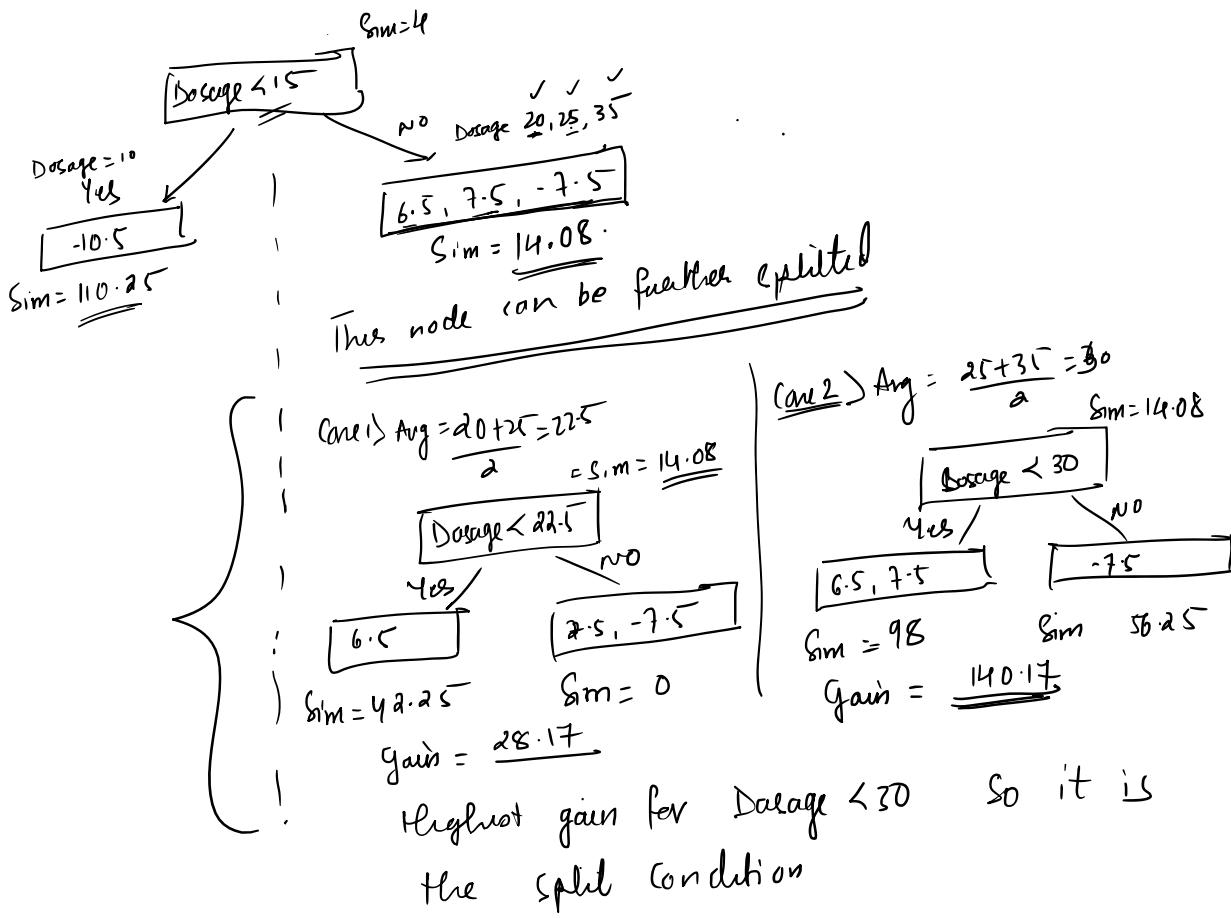
Case 2: Consider dosage as Average of 20 & 25
Average = 22.5
Similarity = 4.
↓
/ dosage < 22.5

Case 3: Dosage as Average of 25 & 35 Avg = 30
Similarity = 4
↓
/ dosage < 30
Yes ↗
↓
/ 7.5

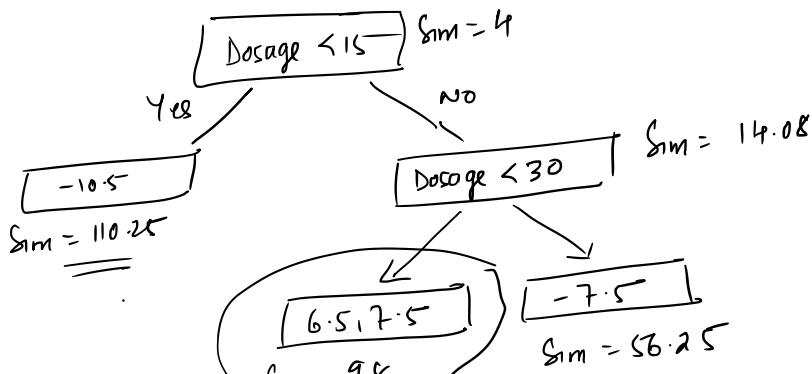


Since the Dosage < 15 gives highest gain = 120.33

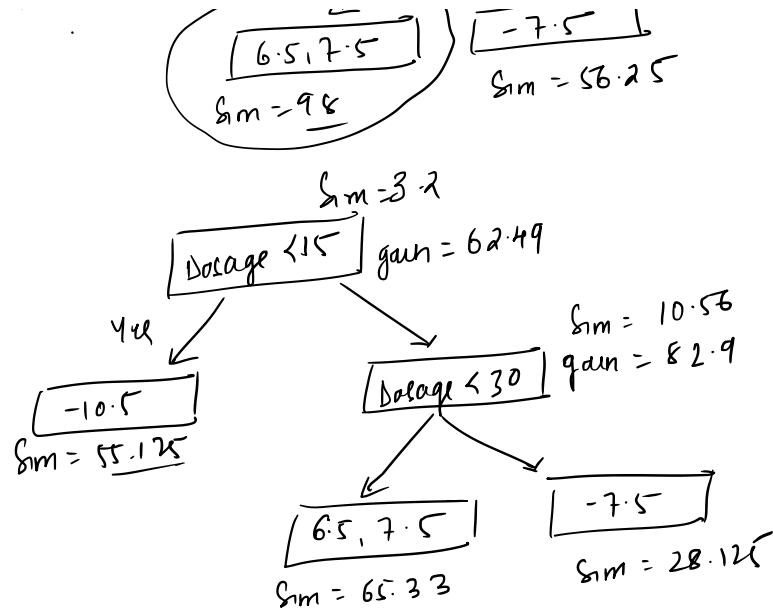
The split factor = Dosage < 15



Novo Tree $\lambda = 0$



Now for $\underline{x=1}$



Note
when $\lambda > 0$ the Similarity Score decreases (\downarrow)
Also Gain decreases (\downarrow).]

Pruning : Kyon chahiye

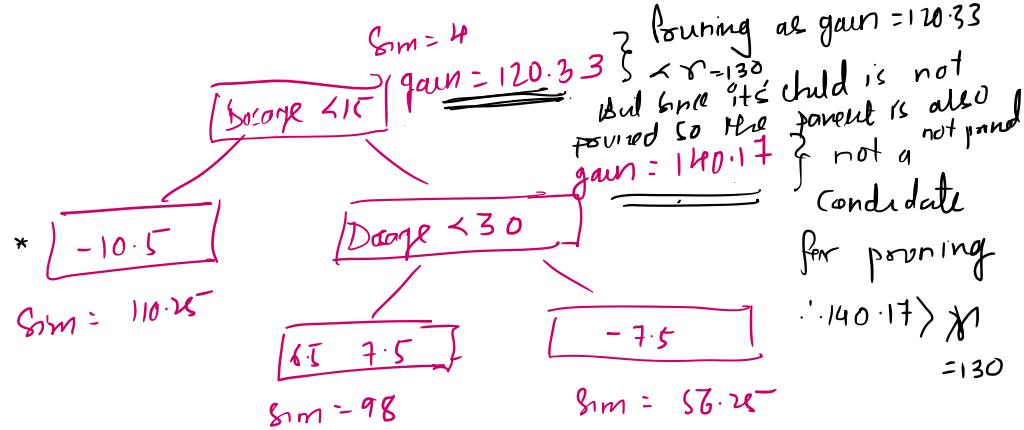
- ① Right distribution of Tree
- ② Ignoring unnecessary branches
- ③ Decrease computation
- ④ Less Time
- ⑤ Avoid Overfitting

Pruning of trees in XGBoost is purely done on gain

but we assume a threshold gain = γ [gamma].

- [Rule for Pruning] * if at branch $\underline{\text{gain} - \gamma} < 0$ then prune and go above
- * if $\underline{\text{gain} - \gamma} \geq 0$ then do not Prune
- * if child is not pruned then parent also will not be pruned]

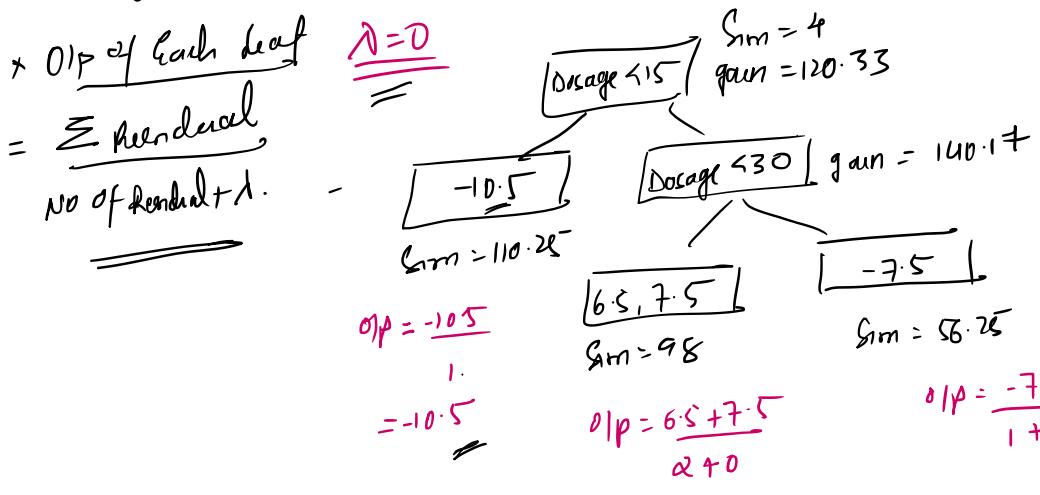
$$\text{Tree } \lambda = 0 \\ \boxed{\text{let } \gamma = 130}$$



* Note
If γ is sufficiently large then it may result into complete tree pruning till root \rightarrow Extreme Pruning

* Regardless of value of λ and γ but we assume the tree as \rightarrow

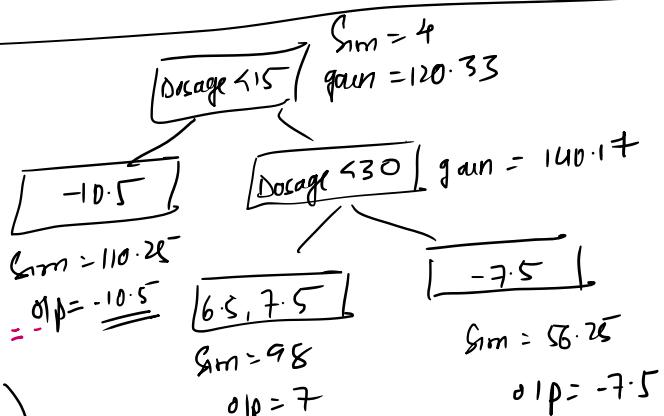
* Regardless of value of λ and η let us assume the tree as \Rightarrow



New Prediction \Rightarrow

$$0.5 + 0.3 *$$

(In XGBoost
the learning
rate is known as η
as Eta (E)
and default value =
 $= 0.3$)



Now prediction for Dosage $\underline{= 10}$

$$= 0.5 + 0.3 * (-10.5) = -\underline{2.65}$$

for Dosage $= 10$

Initial prediction = 0.5

New Prediction = $\underline{-2.65}$

* Find New Prediction for all the other 3 Dosages.

* Using the New Residual Build a New Prediction that will give Smaller Residual

* Keep Building the Tree until the Max No of Trees Constructed or Residual is Significantly Small.

Residual is significantly small.

Let consider few above case there are 4 Tree (T_1, T_2, T_3, T_4)

* final prediction = $0.5 + 0.3 * \underline{\underline{T_1}} + 0.3 * \underline{\underline{T_2}} + 0.3 * \underline{\underline{T_3}} + 0.3 * \underline{\underline{T_4}}$.

XGBoost (Extreme Gradient Boost)

- * Implementation of Gradient Boost Decision Tree.
- * Decision Tree is constructed in Sequential form.
- * Weights play important role in XGBoost
- * Weights are assigned to all independent variables which are fed into Decision Tree that predict result.
- * Weights of variables predicted wrong by previous tree is increased and that variable are fed in next Decision Tree.
- * These individual predictors/classifiers can ensemble to give strong and precise model.
- * XGBoost is faster than Gradient Boost
- * There is stop criterion for Tree splitting in XGBoost
 - XGBoost uses max depth parameter that it starts pruning the tree backward.
 - This pruning improves computational performance and helps to overcome problem of "Overfitting"

XG Boost for Classification →

- ① Convert the observed value into Probabilities (0 & 1).
- ② Decide on initial prediction (e.g. 0.5)
- ③ find Residual (observed - Initial Prediction).
- ④ Create Araf with all Residuals.
- ⑤ Calculate Similarly =
$$\frac{\sum (\text{Residual})^2}{\sum_{\text{prob}} \text{prev} * (1 - \underset{\substack{\uparrow \\ (\text{Initial prob})}}{\text{previous prob}}) + \lambda} \quad [\lambda = 0 \text{ or } \lambda = 1]$$
- ⑥ Now to decide where to Split
 [Take number of Averages case and calculate gain for each.
 Select the split with gain]
- ⑦ XGB have threshold for minimum No of Residuals in the leaf
 * To find Minimum No of Residuals in each leaf
 Cover is calculated.

$$\text{Cover} = \sum (\text{previous prob}) * (1 - \text{prev prob})$$

[In Regression \Rightarrow Cover $\Rightarrow 1$ (default)]
- ⑧ Construct a Tree.
- ⑨ Calculate o/p value at each leaf =
$$\frac{\sum (\text{Residual})}{\text{No of Residuals} + \lambda}$$
.
- ⑩ Calculate New Predicted Value =
$$\underset{\text{Initial}}{\text{Initial}} + \underset{\text{Learning Rate}}{\text{Learning}} * \text{Tree}$$

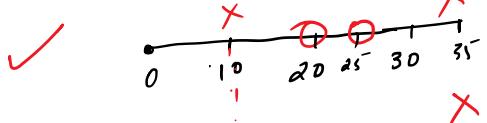
(10) Calculate New Predicted Value = Initial + Learning * Tree
Prob Rate

(11) Repeat the process until max No of Trees are reached -
or Residual becomes Insignificant.

(12) [Note] > New Prediction will be log(odd).
To convert into Probability = $\frac{1}{1+e^{-\text{log(odd)}}}$.
=> This gives new Probability]

XGB for Classification \rightarrow Example
 Q) XGBoost is used with below Dataset

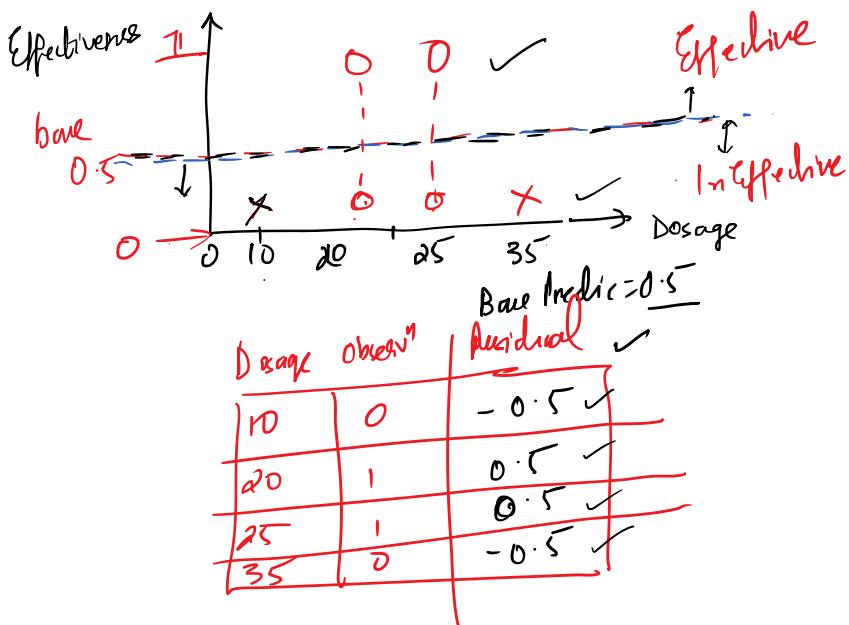
* Let us take 4 sample for our simplicity
 $X = \text{Ineffective drug}$
 $O = \text{Effective drug}$



Dosage	Effective
10	-10
20	7
25	8
35	-7

① Step \rightarrow make Initial (Base Prediction) = 0.5 [Default value]

* Regardless of dosage there is 50% chance that drug is effective.



Dosage	Effective
10	0
20	1
25	1
35	0

$$\underbrace{\text{Previous prob}}_{\text{prob}} * (1 - \text{Previous prob}) + \lambda \cdot \underbrace{\text{Residual}}_{(\text{Residual})^2}$$

Start the Tree with single leaf

(formulae)

$$[-0.5, 0.5, 0.5, -0.5] \quad \text{Similarly } \underline{0} \checkmark$$

Here Prev Probability for all samples = 0.5 (Initial base probability)

Consider

$\text{Dosage } 15$	$\text{Sum} = 0$	$\lambda = 0$
-0.5	$0.5, 0.5, -0.5$	

$\text{Gain} = \underline{1.33}$

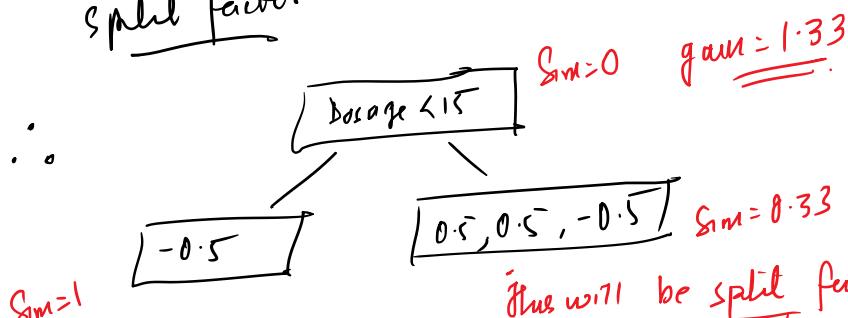
$\text{Sum} = \frac{(-0.5)^2}{(0.5 \times 0.5)} = \underline{1}$

$\text{Sum} = \frac{(0.5 + 0.5 - 0.5)^2}{(0.5 \times 0.5) + (0.5 \times 0.5) + (0.5 \times 0.5)} = \underline{0.33}$

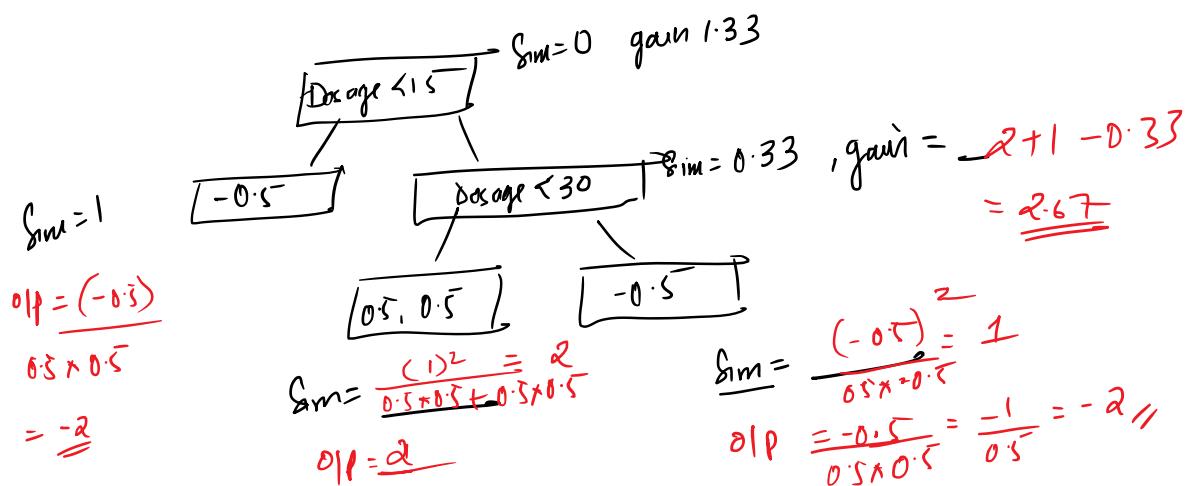
$$\frac{(-0.1)}{(0.5 \times 0.5)} = -\frac{1}{4}$$

$$(0.5 \times 0.5) \text{ gain} = -$$

Let us assume that Dosage < 15 gives best gain, then Dosage < 15 will be split factor.



This will be split further.
Let Dosage < 30 is split factor
that give highest gain



$$\text{OIP for each leaf} = \frac{\sum \text{Sum of Averaged}}{\sum \text{Previous} \times (1 - \text{Previous}) + \lambda} = \text{for } \lambda=0$$

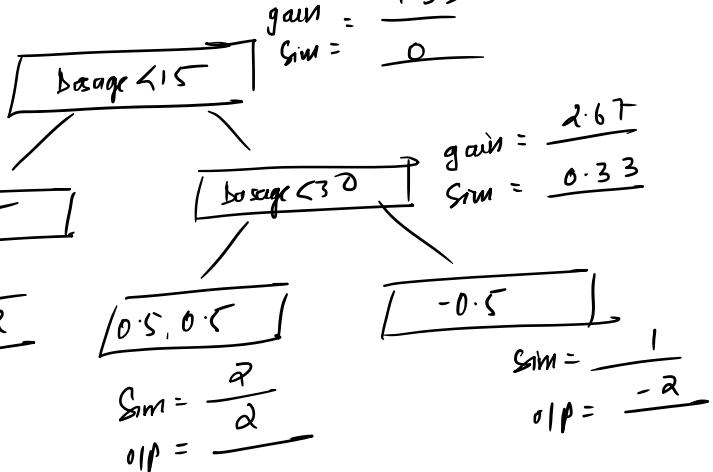
Now
Prediction

$$= 0.5 + 0.3 +$$

↑
Averaging
Rate

$$S_m = \frac{1}{-2}$$

$$\text{OIP} = \frac{2}{2}$$



for Dosage = 10 Prediction $\Rightarrow 0.5 + 0.3 * (-2) = 0.5 - 0.6 = \underline{\underline{-0.1}}$
 Thus is log(odd)

$$\text{New Probab} = \frac{1}{1 + e^{-0.1}} = \underline{\underline{0.524}}$$

For Dosage = 10 ✓ Original Predic 0.5 ✓ New Prediction 0.524.

Residual = -0.5 ✓ Residual = -0.0524 ✓
 Smaller Residual than Before.

- * Similarly do it for all the samples.
- * Find New Probabilities & Residual.
- * Again Construct New Tree
- * Continue until Max No of Trees reached or Residual generated becomes insignificant.

VIMP Bagging & Boosting

(i) Bagging (Bootstrapping Aggregate)

- ↳ Weak learners organized in parallel.
- ↳ Independent weak learners
- ↳ Used to reduce Variance

- Steps:
- ① Multiple subsets are created from Original Dataset
Selecting observations with replacement.
 - ② A base model (weak) is created for each of subset
 - ③ The weak model runs in parallel and are independent to each other.
 - ④ The final predictions are determined by combining the predictions from all the models.

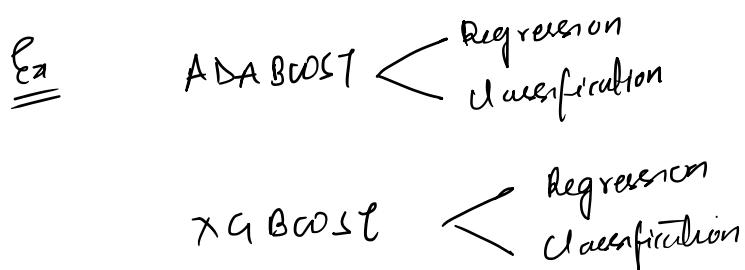
Ex: Random Forest.

Boosting \rightarrow Basic \rightarrow

if a data point is incorrectly predicted by first model
then the next model will correct the prediction thus giving better result for immediate next model.

- Steps
- ① A subset is created from Original Dataset
 - ② Initially all the data points are given equal weight.
 - ③ A base model (Initial prediction) is created on dataset
 - ④ This base model is used to make prediction (Initial) on whole dataset
- error ... \dots calculated for each sample.

- whole dataset
- (1) Residual (Initial) are calculated for each sample.
 - (2) Now the observation that are incorrectly classified are given higher weight
 - (3) Another model is created based on predictions of previous model.
 [This model tries to correct the errors from the previous model]
 - (4) Similarly we will have multiple models (each based on previous model)
 - (5) The final model (Strong learner) is weighted mean of all the model



Q) Which one is better?

Depends on

- (1) Data
- (2) Simulation
- (3) Circumstances

* If Individual single model has high Bias.
then Bagging will not improve Bias

However Boosting will improve the Bias.

* If Single Model overfits then Bagging is the best option
as Boosting will not help in overfitting.

[Note ↑ Bias → Boosting is helpful
 ↑ Variance → Bagging is helpful]

Q) Similarities of Bagging & Boosting

- ① Both are ensemble methods (we need weak learners)
- ② Both generate several training data sets by random sampling.
- ③ Both make final decisions by averaging the N learners (Regression) or majority voting (Classification)
- ④ Both are good at Reducing Variance and providing higher stability.

Q) Differences of Bagging & Boosting

Bagging

- ① Bagging is simplest way of combining predictions that belongs to same type
- ② Bagging ^{primarily} helps in reducing Variance
- ③ Each model (tree) receives equal weight
- ④ Each model is built independently.
- ⑤ Different Training Dataset are randomly built i.e. replacement from

Boosting

- ① Boosting is way of combining predictions that belongs to diff types.
- ② Boosting primarily helps in reducing Bias
- ③ Boosting models are weighted according to their performance (Amount of say).
- ④ New model is built influenced by performance of previously built model.
- ⑤ Every new subset contains Element that were misclassified

are randomly built
with replacement from
entire training dataset

⑥ g_t is better for
Overfitting.

⑦ Ex: Random Forest

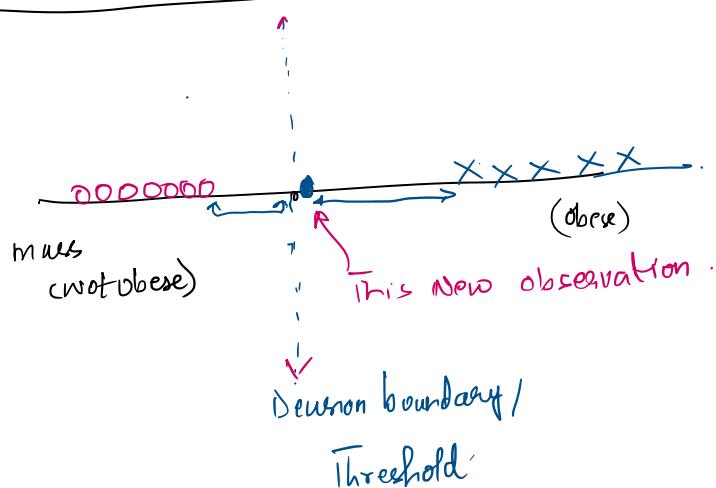
Element that were misclassified
by previous model.

⑥ g_t is better for Underfitting

⑦ Ex. Gradient Boost
ADA Boost
XG Boost

SVM [Support Vector Machine]

Consider \rightarrow



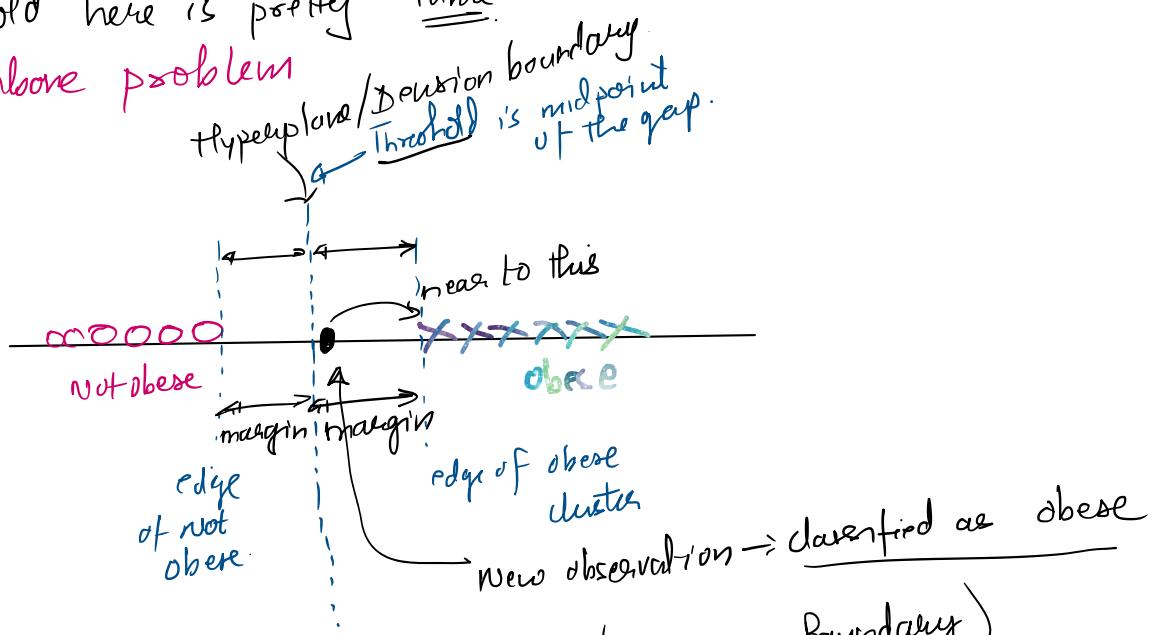
- * Here the new observation will be classified as obese. Even though the new observation is more near to Not obese as compared to obese.

* The threshold here is pretty large.

To Solve the above problem

Consider \rightarrow

We will consider threshold at mid-point of the gap betn the edges of cluster



Here the margin is gap betn threshold (Decision Boundary) and the edge of the cluster

\rightarrow then the margin is maximum

\rightarrow Known as Maximum Margin Classifier

\rightarrow Here in above case there is No Misclassification

→ Also there are no Outliers

Consider → Here we have outliers in NOT obese (one)

* If we consider Margin as the edge of the cluster considering the Outlier.

* Here the new observation will be classified as not obese, even though it is near to obese.

* Here the Maximum Margin Classifier is very sensitive to outlier.

* Here even though we have outlier, however all the given data points are correctly classified → low Bias

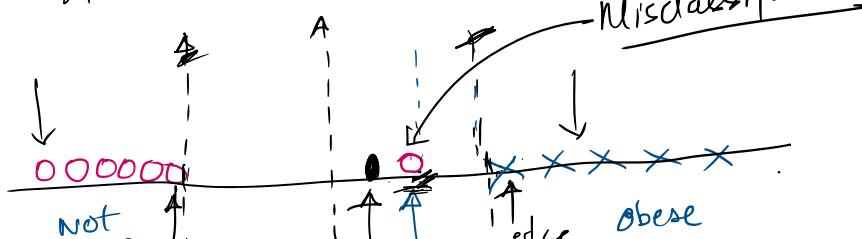
* But for the new observation. It incorrectly classified ⇒ high Variance

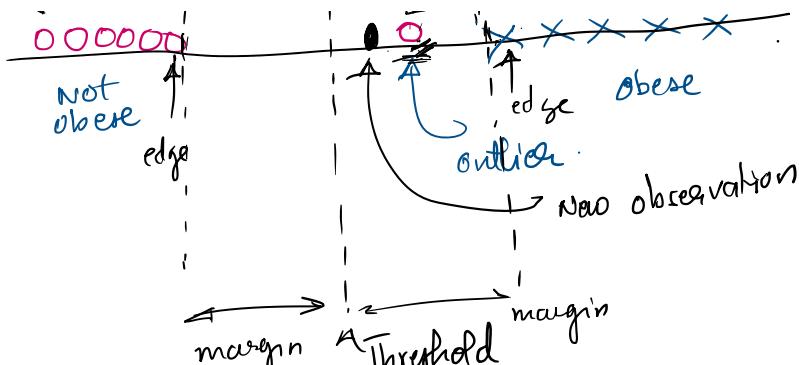
Note > With Hard Margin →

- Misclassification not allowed
- All data points correctly classified ⇒ low Bias
- New observation incorrectly classified ⇒ high Variance

Now: We can do better (To make the threshold Incentive to outlier)

→ For this to happen we must allow Misclassification.





Soft Margin

* Here since while considering edge of cluster, we have allowed misclassification, so the margin now known as soft margin

Note With Soft Margin :-

→ Misclassification is Allowed

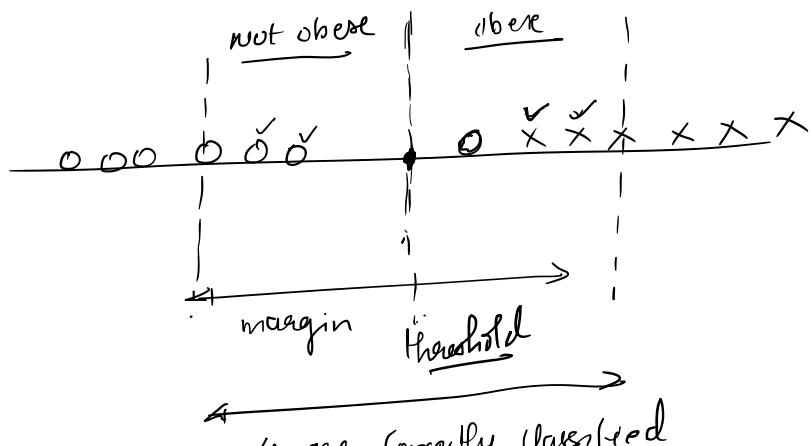
→ All given data points are not correctly classified (outliers). → High Bias

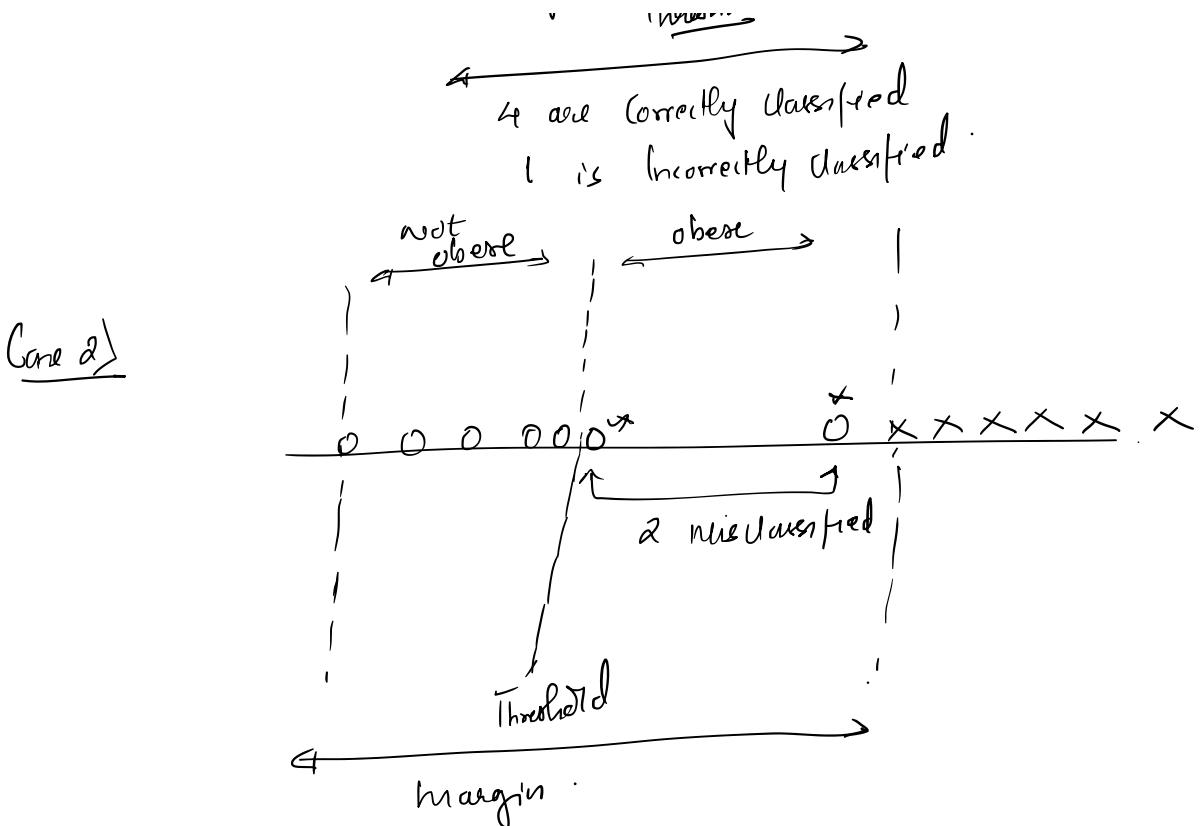
→ New observations are correctly classified \Rightarrow Low Variance

Soft Margin :- When Misclassification is allowed then the distance b/w threshold and given observation is known as Soft Margin.

* For above case we can have many soft margins.

Consider Case I :-





- * We will determine which soft margin is better by observing how many misclassification each soft margin allows.

Imp.
When we use soft margin to determine the location of threshold, then we are using "Soft Margin Classifier"
Also known as "Support Vector Classifier"

- * Support Vector Classifier → Main Aim is to create a Decision Boundary with allowed misclassification
→ It allows small amount of Overlapping

But if the Dataset is as follows →

- * Here there is very high level of overlapping
... at some

* Here there is very high level of overfitting

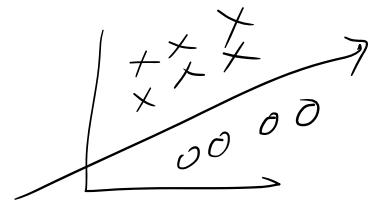
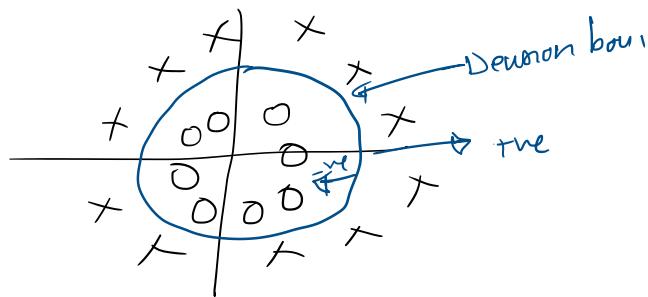
<u>not obse</u>	<u>obse</u>	<u>not obse</u>
0 0 0 0 0	X X X X X	0 0 0 0 0

Marks

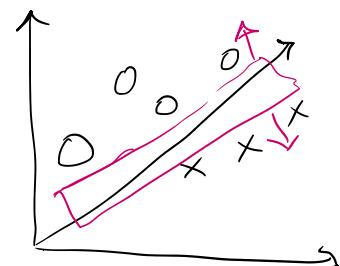
* Here we cannot have linear separator.

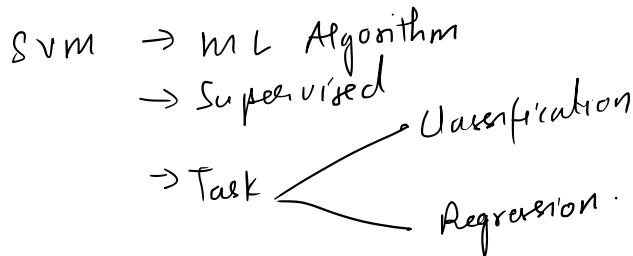
Also \rightarrow [For 1-D data Decision Boundary is a point
For 2-D Data " " is a line
For 3-D Data " " is a plane]

Also



Here also Linear Decision Boundary Not Possible.

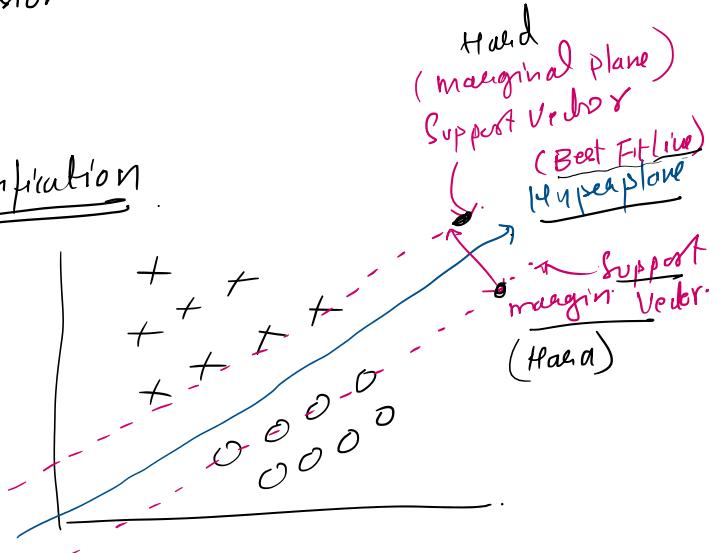


Note

* Support Vector Classifier for Classification.

* Consider Binary Classification →

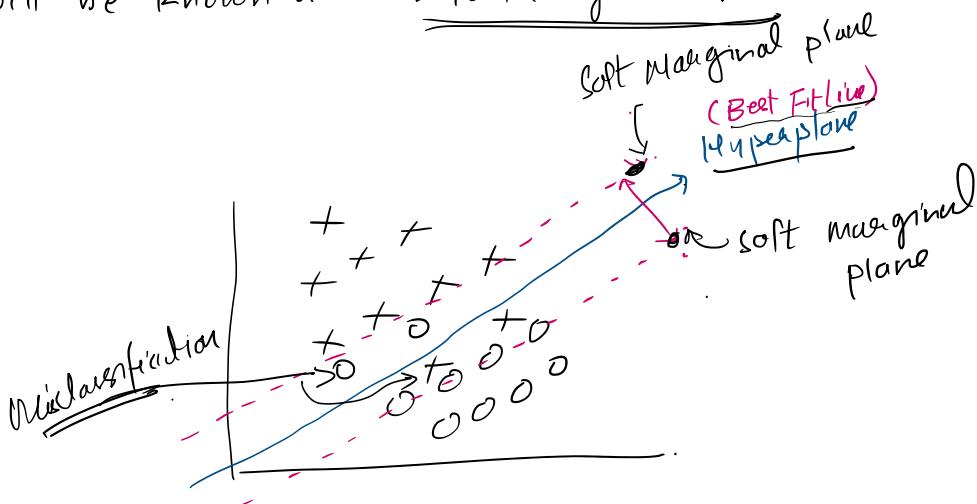
Aim →
 Create a Hyperplane and
Marginal plane such that
the Margin is Maximum →



* If we get perfect Marginal plane & Hyperplane [No Misclassification]
 Then we can linearly separate the Data points.

This Margin will be known as Hard Marginal Plane

* If there are Misclassification allowed then the plane will be known as "Soft Marginal Plane"



Note Consider Hyperplane \rightarrow Here 2-D Data
hyperplane will be line.

$$\text{Eq: } y = mx + c$$

$$\text{or } y = \theta_1 x_1 + \theta_0$$

$$\text{or } y = \underline{\theta_1 x_1 + b}$$

For multilinear regression \Rightarrow Many features $x_1 x_2 x_3 \dots x_n$
Each feature will have weight $w_1 w_2 w_3 \dots w_n$

$$y = \underline{w_1 x_1 + w_2 x_2 + w_3 x_3 + \dots + w_n x_n} + b.$$

$$y = \underline{\underline{w^T X}} + b$$

\uparrow $t_{\text{predicted}}$

actual value

$$w^T = [w_1 \ w_2 \ \dots \ w_n]$$

$$X = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$$

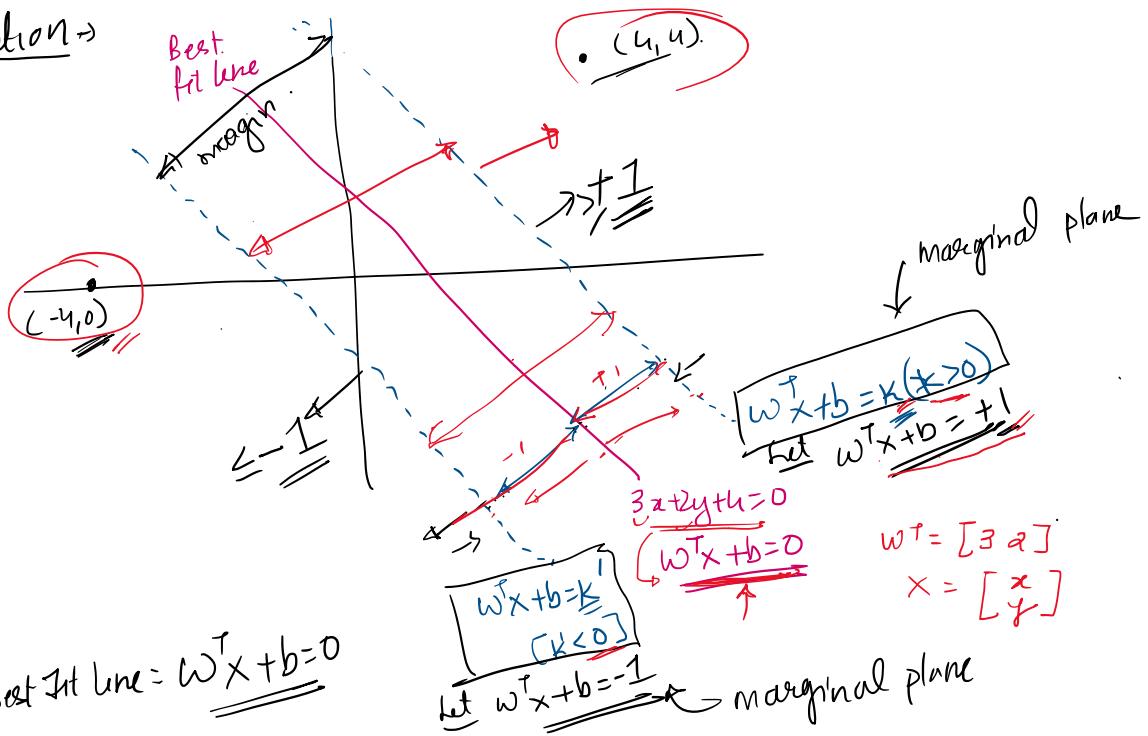
Consider \Rightarrow

SVM For Classification →
(consider)

$$\text{Eq}^1 \Rightarrow 3x + 2y + b = 0$$

put $(-1, 0) = 3(-1) + 2(0) + b < 0$

$$(1, 1) = 3(1) + 2(1) + b \\ = 2b > 0$$



Our aim is to draw two Marginal planes (+ve & -ve side)
and need to ensure the distance (Margin) is maximum.

* We want to find distance bet'n the Marginal Planes.

Let's find difference:

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

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

$$w^T(x_1 - x_2) = 2$$

$$\therefore w^T(x_1 - x_2) = 2$$

Here w = slope (coefficient) magnitude
 direction } Vector.

To convert w^T into Vector i.e. $\vec{w^T}$

To convert w^T into Vector ie \vec{w}^T

divide by $\|w\|$

\therefore

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

difference
(margin)

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

\therefore To Maximize the Margin \rightarrow

$$\checkmark \boxed{\text{Maximize } (w, b) = \frac{2}{\|w\|}}$$

Constraints \rightarrow

y_i
Actual/
observed

+1

-1

when $\vec{w}^T x + b \geq 1$
(point lies outside of $\vec{w}^T x + b = 1$)

when $\vec{w}^T x + b \leq -1$
(point lies below of $\vec{w}^T x + b = -1$)

The constraints are for Correctly Classified Data point

\therefore Final constraints for S.V Classifier.

$$\boxed{y_i * (\vec{w}^T x + b) \geq 1}$$

↑
observed
(actual)

↑
predicted
value

For correctly
classified Data point

\therefore To Maximize $= \frac{2}{\|w\|}$, subject to

$$\boxed{y_i * (\vec{w}^T x + b) \geq 1}$$

$y_i = +1$
 $\vec{w}^T x + b \geq +1$

To Maximize $= \frac{\alpha}{\|w\|}$, Subject to \leq

(Also can be written as)

* To Minimize $= \frac{\|w\|}{\alpha}$, Subjected to $y_i * (w^T x + b) \geq 1$
only for Correctly Classified Data points.

here we have not considered for Misclassification
But in real world there will be always Misclassification.

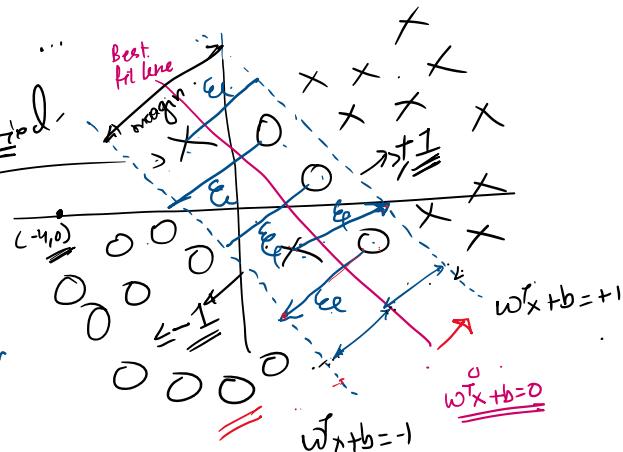
To Consider for Misclassification we have to use
"Hyper parameters"

We allow Misclassification

Now we will use Misclassified
two Hyperparameters

$\epsilon \Rightarrow \text{distance bet'}$
the Misclassified
point and correct
Marginal plane

$C_i \Rightarrow \text{No of Allowed Misclassified}$
points



Now Final Cost F^n (For All data points) [Correctly &
Incorrectly Classified]

* Objective's

$$\text{To Minimize}_{(w,b)} \frac{\|w\|}{\alpha} + C_i \sum_{i=1}^{C_i} \epsilon_i$$

Constraint $y_i * (w^T x + b) \geq 1$

for Hyperparameters

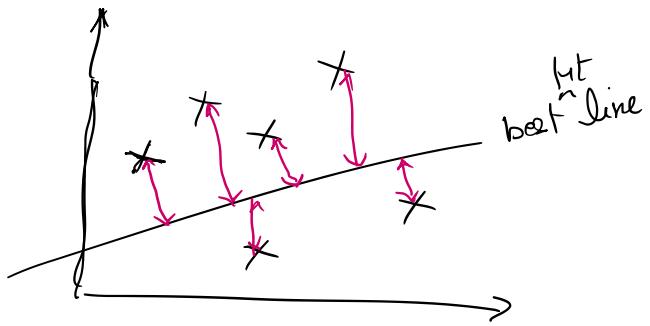
ii. Labeled Cost F^n for SVC used for Classification

~~Non-penalized~~ Identified Misclassification ·
The label cost Γ^n for SVC used for classification
[For Allowed Misclassification].

Support Vector For Regression \Rightarrow

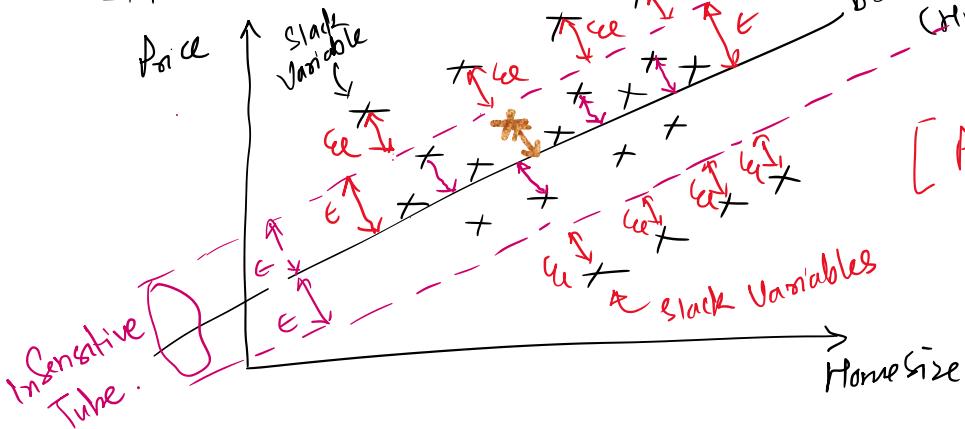
- * In Linear Regression

Objective \rightarrow Minimize MSE
[Mean Squared Error]



Best Fit line \Rightarrow It is the line that gives Minimum MSE

Support Vector Regression \Rightarrow



$w^T x + e$ (margin plane)
 $w^T x - e$ (margin plane)
best fit line (Hyperplane)
 $w^T x$

[Data Points Inside ϵ Insensitive Tube \rightarrow Support Vector]

\rightarrow Here we are calculating e (Insensitive Tube) [Allowed Margin of Error]

- * We can ignore the error if difference between observed and predicted value is $\leq \epsilon$

All the points inside ϵ insensitive tube are not considered for calculating error.

- * The data points outside the ϵ insensitive tube are known as Slack Variables.

- * For a point x_i within ϵ tube margin value

- * For a point x_i within Insensitive (ϵ) tube

$$|y_i - w^T x_i^*| \leq \epsilon$$

y_i^* = observed value
 x_i^* = predicted value

constraint

$$\text{if } |y_i^* - w^T x_i^*| \leq \epsilon$$

We can say that prediction is good and we will not consider the difference in error calculation.

- * For Slack Variable [points outside ϵ Insensitive tube]

we need to calculate distance of slack variable from Hyperplane

$$\begin{aligned} &= \text{distance betn slack variable} + \epsilon \\ &\quad \& \text{Marginal plane} \\ &= \epsilon_e + \epsilon \end{aligned}$$

Objective \Rightarrow

$$\text{Cost } F^n \Rightarrow \text{Minimize } (w, b)$$

$$\frac{\|w\|}{\alpha} + C_i^* \sum_{e_i=1}^{l_i} \epsilon_e^*$$

Constraints

$$|y_i^* - w^T x_i^*| \leq \epsilon + \epsilon_e^*$$

Hinge Loss

Hyperparameter

for points outside ϵ Insensitive tube

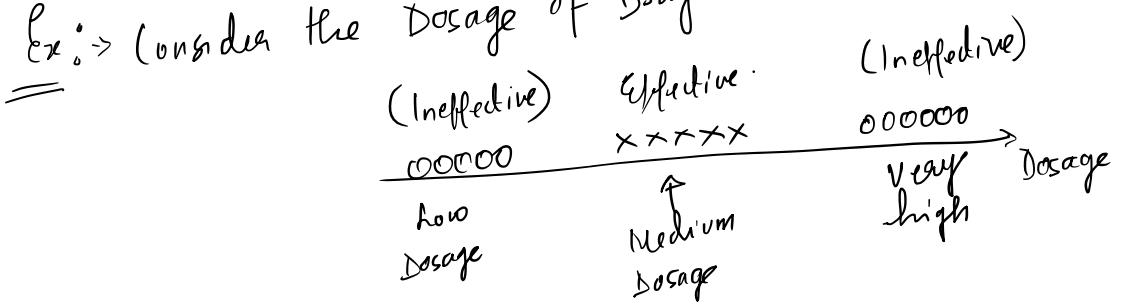
We have considered up till now Support Vector Classifier.

(Soft Margin)

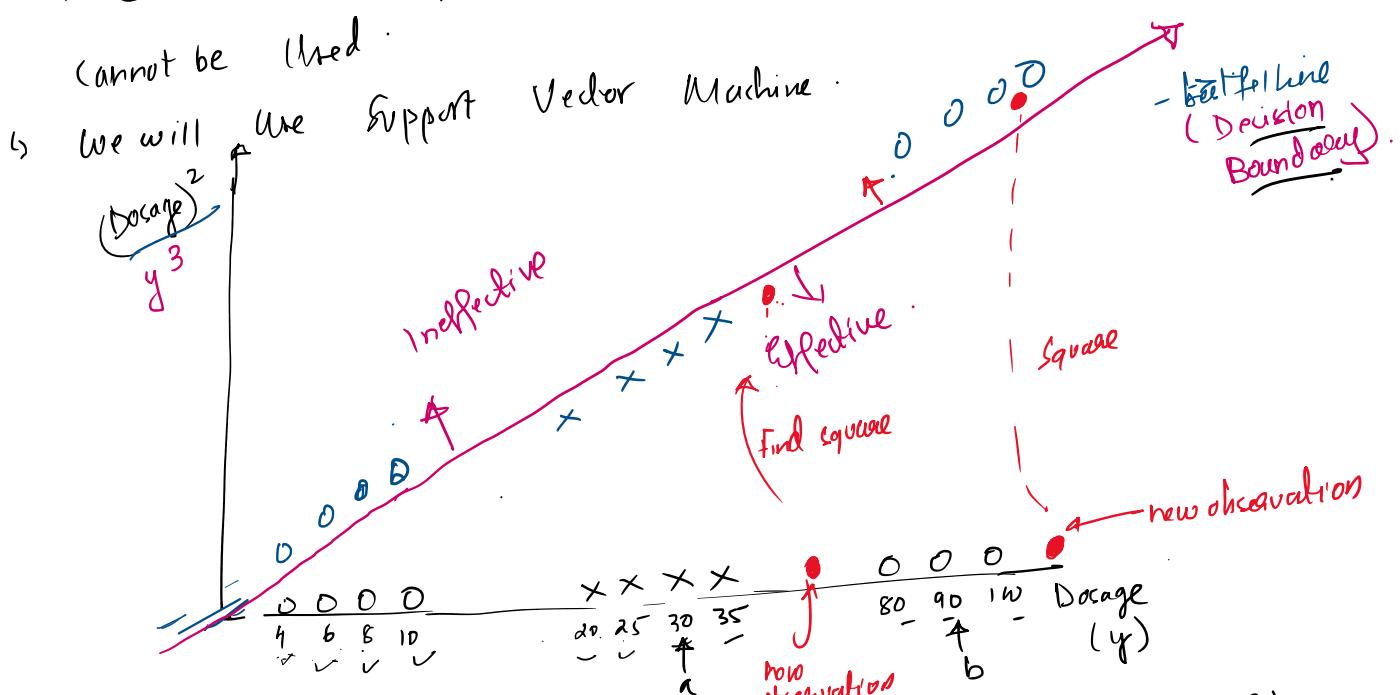
[We were interested in finding Max Margin Classifier]

- * If there are high amount of Overlapping than Soft Margin classifier will not be efficient
- * In Such Case we will use Support Vector Machine.

Ex: Consider the Dosage of Drug



→ Here we have lot of Overlapping so Maximum Margin classifier.



→ From above By increasing dimension of data from 1D to 2D and by using Square (as function) we are able to separate the data

* From above we are able to separate the unseparable points and by using Square (as function) we are able to separate the unseparable points.

* We can also use it to predict for new observation
→ For new observation take square of the observation and check if the square lies above (ineffective) or below (effective) the decision boundary.

SVM Main Idea →

- * Start data with low dimension (given).
- * Now move data to higher dimension [Using some function]
- * Find Support Vector Classifier to classify the data into different classes.

Question Arises → Why Square, why not Cube or some other function?

→ How to decide, how to transform.

- * The Function Used be known as Kernel
- * Types of Kernel function in SVM.
 - (1) Polynomial Kernel
 - (2) Radial Basis Kernel.

Note)
if data points in $\underline{1D}$ \rightarrow convert in $\underline{2D} \Rightarrow$ line as Decision boundary
 $\underline{2D} \rightarrow$ convert in $\underline{3D} \Rightarrow$ Plane as Decision Boundary

2)
if Data points in $1D \rightarrow$ Point is Decision Boundary
 $2D \rightarrow$ Line is Decision Boundary
 $3D \rightarrow$ Plane is Decision Boundary

Polynomial Kernel

- * It calculates higher Dimensional Relationship between observations (data points).
 - * The Kernel that transforms 1D to dD is Polynomial Kernel.
 - * It may look like $(a \cdot b + r)^d$
 where a and b are two different observations in dataset
 (any two data points).
 - r = coefficient of Polynomial.
 - d = degree of Polynomial.
 - * Here we use SVM with Polynomial Kernel to compute relationship betⁿ observations in higher dimension and then find good classifier.
 - * Let a and b be two observations.
 - * Let $r = 1/2$ & $d = 2$ } Note ⇒ value of r and d is determined by cross validation.
- $$\Rightarrow (a \cdot b + r)^d = (a \cdot b + 1/2)^2$$
- $$= (a \cdot b + 1/2) \cdot (a \cdot b + 1/2)$$
- $$= a^2 b^2 + \frac{1}{2} ab + \frac{1}{2} ab + \frac{1}{4}$$
- $$= ab + a^2 b^2 + \frac{1}{4} \Rightarrow \text{Polynomial.}$$
- = can be written as dot product of

= can be written as dot product of

$$\Rightarrow = \left(\underline{\underline{a}}, \underline{\underline{a}}^2, \frac{1}{2} \right) \cdot \left(\underline{\underline{b}}, \underline{\underline{b}}^2, \frac{1}{2} \right)$$

For data point a

Original value of a
Higher dimension value of a

Original value of b
Value of b in higher dimension

[Dot product is sum of 1st term multiplied, 2nd term multiplied and so on]

Note: $\frac{1}{2}$ is third axis but value same so ignore.

So $(a+b)^d$ is used to get higher dimension "rel" between two data points a & b .

Ex: $a=9$ $b=14$
 $\underline{\underline{a}}=\underline{\underline{b}}$ $\underline{\underline{b}}=\underline{\underline{a}}$

\Rightarrow Higher Dimension Relationship betⁿ
 $a=9, b=14$

 $\Rightarrow (\underline{\underline{9}} + \underline{\underline{14}} + \underline{\underline{\frac{1}{2}}})^2 \Rightarrow 126.5^2 \Rightarrow \underline{\underline{16002.25}}$

↑
Relationship representation of
 $a=9, b=14$ in
Higher Dimension

* We can find this Higher Dimension Relationship betⁿ every pair of data points.

VVIMP \rightarrow Kernel Trick

- ① Kernel function actually never does any transformation in higher dimension.
- ② Instead it calculates relation betⁿ observations and visualizes them in higher dimension.

visualizes them in higher Dimension.

③ Support Vector Classifier uses this Visualization for Classification.

④ This is known as Kernel Trick.

$d=1$ \Rightarrow The Polynomial Kernel compute the relationship betⁿ each pair of observation in 1D.

$d=2$ \Rightarrow " " " " " " " " $\in \mathbb{R}^D$.

Best value of d can be found by Cross Validation.

Note Since Polynomial Kernel is $\underbrace{(a \cdot b + 1)^2}_{\text{↑}} = \underbrace{(a, a^2, 1)}_{\text{↑}} \cdot \underbrace{(b, b^2, 1)}_{\text{↑}}$

\rightarrow we actually do not have to Transform to understand higher dimension relationship.

\rightarrow All we need to do is calculate dot product betⁿ each pair of point.

[not important for Exam].

not important for exam].
let us consider infinite Dimension Concept \rightarrow .

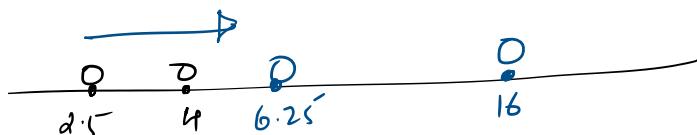
Let a & b be two observations. Let $a=0$, $d = (a \times b + 0)^d = a \cdot b$

① If $d = \alpha$

$$(a \times b + 0)^2 = a^2 + b^2$$

$$a = a \cdot 5 \quad (a+b+0)^2 = (a \cdot 5)^2 + (4)^2$$

6. 25 16



the dimension is not changed instead the kernel shrinks the data down the original axis.

Summarize \Rightarrow Let a and b be two observations.

$$x=0 \quad d=1$$

$$r=0 \quad \underline{d = \infty} \quad a^2 \cdot b^2$$

$$r=0 \quad d=3 \quad a^3 \cdot b^2$$

$$r=0 \quad d=\underline{\infty} \quad a^\infty \cdot b^\infty$$

let us add all

$$= a^1 b^1 + a^2 b^2 + a^3 b^3 - \dots + a^{\infty} \cdot b^{\infty}$$

$$= \begin{pmatrix} a, a^2, a^3, \dots & a^\infty \\ \uparrow & \uparrow & \uparrow & \ddots & \uparrow \end{pmatrix} \cdot \begin{pmatrix} b, b^2, b^3, \dots & b^\infty \\ \uparrow & \uparrow & \uparrow & \ddots & \uparrow \end{pmatrix}$$

This gives a set provided with coordinates for infinite number of dimension.

Radial Basis Kernel [VImp]

- * Works in Infinite Dimension.
- * If there are lot of Overlapping in data points classification then Support Vector Classifier cannot be used to linearly separate.
- * One way to deal with overlapping data is to use SVM with Radial Kernel.
- * Radial Kernel uses Radial Basis function!

$$-\gamma \cdot \frac{(a-b)^2}{\epsilon}$$

Here a and b are two observations.

$(a-b)^2$ \Rightarrow diff betⁿ the measurement is squared giving us the squared distance betⁿ two observation.

$\gamma \Rightarrow$ scales the squared distance and thus scales the influence.

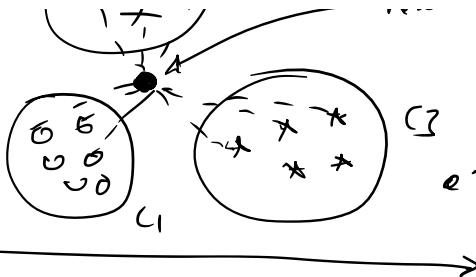
- * Since Radial Basis function finds Support Vector classifier in Infinite Dimension it is not possible to visualise it.

- * When applied to new observation, it behaves like Weighted Nearest Model.

- * Radial Kernel checks the influence of the Nearest



Influence of the Nearest Neighbour on New observation and accordingly makes classification.



- * Let see how radial kernel determines how much influence each observation in training dataset has on classifying new observation.

$$\text{Let } \begin{cases} a = 2.5 \\ b = 4 \end{cases} \Rightarrow \frac{-1}{e} (2.5 - 4)^2 \Rightarrow 0.01 \rightarrow \textcircled{1}$$

$$\gamma = 1$$

$$\begin{aligned} & \begin{cases} a = 2.5 \\ b = 4 \\ \gamma = 2 \end{cases} \Rightarrow \frac{-2}{e} (2.5 - 4)^2 \Rightarrow 0.01 \end{aligned} \quad \left. \begin{array}{l} \text{The } \gamma \text{ scales} \\ \text{the influence.} \end{array} \right\}$$

$$\begin{aligned} & \begin{cases} a = 2.5 \\ b = 16 \\ \gamma = 1 \end{cases} \Rightarrow \frac{-1}{e} (2.5 - 16)^2 \Rightarrow \frac{-1}{e^{13.5}} \approx \text{close to zero} \end{aligned} \quad \hookrightarrow \textcircled{2}$$

From ① & ② If points are close enough then we have high influence.

and if points are far off then we have low influence.

- * Thus to calculate influence b/w two data points, we can plug in the values in Radial Basis fn with appropriate γ .



... in on infinite Dimension.

* We get relationship in infinite dimension.

Not from Exam \Rightarrow Let a and b are two observations.

$$\text{Let } \frac{a+b}{2} = \frac{-1}{e^{\frac{a+b}{2}}} (a^2 + b^2 - 2ab) = \frac{-1}{e^{\frac{a+b}{2}}} (a^2 + b^2). e^{\frac{ab}{2}} \quad (1)$$

Taylor Series Expansion \Rightarrow [Allows any f^n to split in infinite sum]

$$f(x) = f(a) + \frac{f'(a)}{1!}(x-a) + \frac{f''(a)}{2!}(x-a)^2 + \dots + \frac{f^{(\infty)}(a)}{\infty!}(x-a)^\infty$$

Let $f(x) = e^x$.

$$e^x = e^a + \frac{e^a}{1!}(x-a) + \frac{e^a}{2!}(x-a)^2 + \dots + \frac{e^a}{\infty!}(x-a)^\infty$$

Let $a=0$

$$e^x = 1 + \frac{x}{1!} + \frac{x^2}{2!} + \dots + \frac{x^\infty}{\infty!}$$

$$e^{ab} = 1 + \frac{ab}{1!} + \frac{ab^2}{2!} + \dots + \frac{a^\infty b^\infty}{\infty!}$$

$$e^{ab} = \left(1, \frac{1}{\sqrt{1!}}, \frac{1}{\sqrt{2!}}, \dots, \frac{1}{\sqrt{\infty!}} \right) \cdot \left(1, \frac{1}{\sqrt{1!}}, \frac{1}{\sqrt{2!}}, \dots, \frac{1}{\sqrt{\infty!}} \right)$$

from (1) & (2)

$$e^{ab} = \frac{-1}{e^{\frac{a+b}{2}}} \left(1, \frac{a}{\sqrt{1!}}, \frac{a^2}{\sqrt{2!}}, \dots, \frac{a^\infty}{\sqrt{\infty!}} \right) \cdot \left(1, \frac{b}{\sqrt{1!}}, \frac{b^2}{\sqrt{2!}}, \dots, \frac{b^\infty}{\sqrt{\infty!}} \right)$$

$$\text{hel } S = \frac{-1}{e^2} (a^2 + b^2)$$

$$= \left(S, 1, \frac{a}{\sqrt{1}}, \frac{a^2}{\sqrt{2!}}, \dots, \frac{a^\infty}{\sqrt{\infty!}} \right) \cdot \left(S, 1, \frac{b}{\sqrt{1}}, \frac{b^2}{\sqrt{2!}}, \dots, \frac{b^\infty}{\sqrt{\infty!}} \right)$$

We can see Radial Kernel is equal to Dot product that has coordinate for infinite No of Dimension.

→ Dimensionality Reduction →

① PCA [Principal Component Analysis]

Consider

→

weight	height	DBP	SBP	Health?
		↓ Diastolic B.P.	↓ Systolic B.P.	

Say health of person depends on 4 features



To represent this data → We need 4 Dimension

For more features → We need more Dimension.

- * Visualizing data in More than 3 dimension is difficult.
- * Computation on 4 Features is also Complex

Possible Solution →

Height	weight	DBP	SBP	Health -
		↑ B.P.	↑ B.P.	?

[Body Mass Index] BMI

[Blood Pressure]

- ① look for Strong correlated features: Here Height & weight are strongly correlated
Also DBP and SBP are strongly correlated

∴ There is a strong correlation of Height & weight → BMI

→ We can combine the effective correlation of Height & Weight \Rightarrow BMI
 " " " " " " " " " " DBP and SBP \Rightarrow BP.

Now instead of 4 columns we have 2 columns.

BMI	BP

→ This is 2-D data. Easy to Compute
 Easy to Visualize.

We have Reduced the dimension of data from 4-D to 2-D.

Note: Consider d-columns :-
 Let \rightarrow $x_1 \quad x_2$

DBP	SBP
78	126
80	128
81	127
82	130
84	130
86	132

$$(BP) = \alpha_1 \underline{x_1} + \alpha_2 \underline{x_2}$$

α_1 & α_2 are weights of features
 x_1 and x_2 are features

Ex → ① $BP = 0.8 \underline{DBP} + 0.6 \underline{SBP}$

we are giving more weights to DBP as compared to SBP.

Ex ② Let $BP = \text{mean of } DBP \text{ & } SBP$.

for mean

DBP	SBP	BP
78	126	102
80	128	104
81	127	104
82	130	106
84	130	107
86	132	109

$$BP = \frac{DBP + SBP}{2} = 0.5 DBP + 0.5 SBP$$

$$\begin{aligned} \alpha_1 &= 0.5 \\ \alpha_2 &= 0.5 \end{aligned}$$

∴ $x_1 = DBP$ and $x_2 = SBP$.

82	150	107
84	130	107
86	132	109

E₁(3)

Let BP = Sum of DBP and SBP.

$$BP = DBP + SBP \Rightarrow \alpha_1 = 1 \\ \alpha_2 = 1$$

For Sum

DBP	SBP	Sum
78	126	204
80	128	208
81	127	208
82	130	212
84	130	214
86	132	218

↑

Now PCA \Rightarrow Principal Component Analysis ✓
 → It is a method to find the linear combination that
accounts for as much variability as possible
in combined Variable \hookrightarrow (Maximum variance)

To understand $y = \alpha_1 x_1 + \alpha_2 x_2$
 we want value of α_1 and α_2 such that the variance in
 the value of y should be maximum
 B'coz more the variance \rightarrow more is the information]

Note: for $y = \alpha_1 x_1 + \alpha_2 x_2$
 (combined variable)
 → Here to maximize variance large value of α_1 and α_2 can be proposed
 → We will place restriction on value of weights i.e. $(\alpha_1, \alpha_2 \dots)$
 such that $\alpha_1^2 + \alpha_2^2 + \dots + \alpha_n^2 = 1$

↗ If there are n features, we will need ' n ' weights

group ex: we had 2 features: we need α_1 and α_2

$$\therefore BP = \alpha_1 DBP + \alpha_2 SBP$$

constraint $\alpha_1^2 + \alpha_2^2 = 1$

but $\alpha_1 = 0.8$ $(0.8)^2 + (0.6)^2 = 0.64 + 0.36 = \underline{1.00}$ ✓

Let $\alpha_1 = 0.8$ $\alpha_2 = 0.6$

$$(0.8)^2 + (0.6)^2 = 0.64 + 0.36 = \underline{\underline{1.00}} \quad \checkmark$$

Let $\alpha_1 = 0.8$ and $\alpha_2 = 0.6$

DBP	SBP	BP = $\alpha_1 \text{DBP} + \alpha_2 \text{SBP}$
78	126	138
80	128	140.8
81	127	141.0
82	130	143.6
84	130	145.2
86	132	148.0
Mean		142.8

$$\text{Variance} = \frac{1}{n-1} \sum_{i=1}^n (\text{BP} - \text{mean BP})^2$$

$$= \frac{1}{5} \left((138 - 142.8)^2 + (140.8 - 142.8)^2 + (141.0 - 142.8)^2 + \dots \right)$$

$$= \underline{\underline{12.74}}$$

Let for above 6 input Sample

* We can take many values of α_1 and α_2 such that $\alpha_1^2 + \alpha_2^2 = 1$

and calculate variance for each α_1 and α_2 of the above 4 cases.

We get largest variance = 12.74

α_1	α_2	$\text{var}(y)$
0.8	0.6	12.74
0.6	0.8	11.8
0.98	0.2	10.4
0.2	0.98	7.4

Tanika 1 : $\therefore \boxed{BP = 0.8 \text{DBP} + 0.6 \text{SBP}}$

To find α_1 and α_2

We got our α_1 and α_2

We got our x_1 and x_2
 Taika 2 \Rightarrow
Now \Rightarrow Consider again \Rightarrow

DBP	SBP
78	126
80	128
81	127
82	130
84	130
86	132

construct
Covariance = Matrix

$$\begin{bmatrix} \text{cov}(x, x) & \text{cov}(x, y) \\ \text{cov}(y, x) & \text{cov}(y, y) \end{bmatrix}$$

$$\underline{\text{cov}}(x, y) = \frac{1}{N-1} \sum (x_i - \bar{x})(y_i - \bar{y}) = \underline{\underline{\text{cov}}(y, x)}$$

For above Sample data set

Covariance Matrix \Rightarrow

DBP	SBP
8.17	5.97
5.97	4.97

Eigen values = $\begin{bmatrix} -0.8 \\ -0.6 \end{bmatrix}$ ✓
 from above Covariance Matrix.

$$\begin{aligned} \therefore BP &= -0.8 DBP - 0.6 SBP \\ \text{in } BP &= 0.8 DBP + 0.6 SBP \end{aligned} \quad \left. \right\} \checkmark$$

Applications \Rightarrow

(1) Consider

DBP	SBP	Weight	Height

BP	BMI

Reduced Dimension by
 combining features.

BP

BMI

Reduced Dimension "I"
Combining features -

(2)

	DBP	SBP	Weight	Height	Pulse
Pat 1	78	126	82	5.1	90
Pat 2	80	128	74	5.4	80
	81	127	64	5.11	70
	82	130	82	6.2	90
	84	130	92	6.3	95
	86	132	85	6.10	94

Here if we have many no of features and we need to compare the health chart/parameters of samples, it is difficult

↓ ↓
But if we combine the above 5 col in a col -

C1	C2
v1	v2
v3	v4
v5	v6

Now with less no of features comparing the health chart/parameters of samples is relatively easy:

- * PCA \rightarrow It is Dimensionality Reduction Technique.

Consider: For House Price Prediction

Size	Location	Yoc	OC Y/N	Builder	Color	Garden	Swimming	Metro Distance	P'ice

- * Here Price of House depends on above features.
- + There Many No of features (Problem of Plenty).
- * Issues:
 - Represent "Problem" ✓
 - Overfitting Problem ✓
- * We need to Reduce Dimension
- * We need to Identify Important Components.

Dimensionality Reduction:

- + Reduces the dimension of feature Space
- Ex if there are 100 features/col in dataset and you want to get only 10 features then with dimensionality reduction technique we can achieve this.

- * It transforms dataset which is in n dimension Space to n' dimension Space where $n' < n$.

to n dimension space where $n \sim \dots$

Why Dimensionality Reduction?

↳ Normally it is argued that many features gives

More accurate result

But after some point the model performance decreases

(Overfitting) with increase in number of features.

* This is Curse of Dimensionality

* So Dimension Reduction is crucial.

* PCA enables us to identify the correlation and patterns in the dataset so that it can be transformed into new dataset with lower dimension without loss of important information.

PCA (Principal Component Analysis)

LDA
SVD

House Prediction

1	2	3	4	5	6	7	8	9	10
Size	Locah	YOC	Oc/N	Buidle	Color	Garden	Swimmi	Distance from Metro	Price

"Problem of"

Plenty "

- Represent "problem"
- Overfitting problem

So we need to
Reduce Dimension.

- * We need to identify
Important component.

* PCA is dimension Reduction Technique ✓

* Dimensionality Reduction →

* Reduces the dimension of feature space. [Reduces no of dimensions/features for Analysis.]

Ex If there are 100 features / col in dataset and you want to get only 10 features then with dimensionality reduction technique we can & achieve this:

* It transforms dataset which is in [n dimension Space] to n' dimension space where $n' \leq n$

Why Dimensionality Reduction →

↳ Normally it is argued that many features gives more accurate result.

↳ However after some point the performance of model decreases (Overfitting) with increase in no of features.

↳ It increases (overfitting) with increase in no of features.

↳ This is known as "Curse of Dimensionality"

So Dimension Reduction is crucial.

PCA enables us to identify the correlation and pattern in a dataset so that it can be transformed into new dataset of significantly lower dimension without loss of important information.

Eg.

Height	Weight	DBP	SBP	Cholesterol
90	90	90	110	?
90	90	Highly Correlated	BP	?

BMI	BP
x^2	y^2

Steps in Performing PCA

Step 1 → Get the data

Step 2 → Subtract the mean and produce New dataset (Row Data Adjust)

Step 3 → Calculate the Covariance Matrix

Step 4 → Calculate the Eigen Vectors and Eigen values of the Covariance Matrix

Step 5 → Choosing Components and forming a feature Vector.
Feature Vector = (eig1, eig2, eig3 ... eign)

Step 6 → Derive the new data set

Final Dataset = Row Feature Vector * Row Data Adjust

* Row Feature Vector → is matrix with the Eigen Vectors in the column transposed so that Eigen Vectors are now in the rows with most significant Eigen Vectors on top.

* Final Dataset is set with data items in col and dimensions along rows.

Step 7 → Getting the old data back.

Step 7

→ Getting the old data back.

$$\underline{\text{Final Data}} = \underline{\text{Row Feature Vector}} \times \underline{\text{Row Data Adjust}}$$

$$\underline{\text{Row Data Adjust}} = \underline{\text{Row Feature Vector}^{-1}} \times \text{Final Data}$$

$$= \underline{\text{Row Feature Vector}^T \times \text{Final Data}}$$

$$\text{Now Original Data} = \text{Row Data Adjust} + \text{Original Mean}$$

Perform PCA on following Dataset

✓	2	2.5	0.5	2.2	1.9	3.1	2.3	2.0	1.0	1.5	1.1
✓	y	2.4	0.7	2.9	2.2	3.0	2.7	1.6	1.1	1.6	0.9

Step 1: Given dataset \rightarrow

$$\text{for } \lambda_2 \rightarrow$$

(2.5, 2.0)	(0.5, 0.7)	(3.1, 2.9)	
(0.827)	-1.777	-0.492	0.274
		1.875	0.912
			-0.099
			0.046
			0.017
			1.00

ID (converted Data) \uparrow

(1.9, 2.2)

2	2.5	0.5	2.2	1.9	3.1	2.3	2.0	1.0	1.5	1.1
y	2.4	0.7	2.9	2.2	3.0	2.7	1.6	1.1	1.6	0.9

Step 2: calculate mean

$$\bar{x} = \frac{1.81}{11}$$

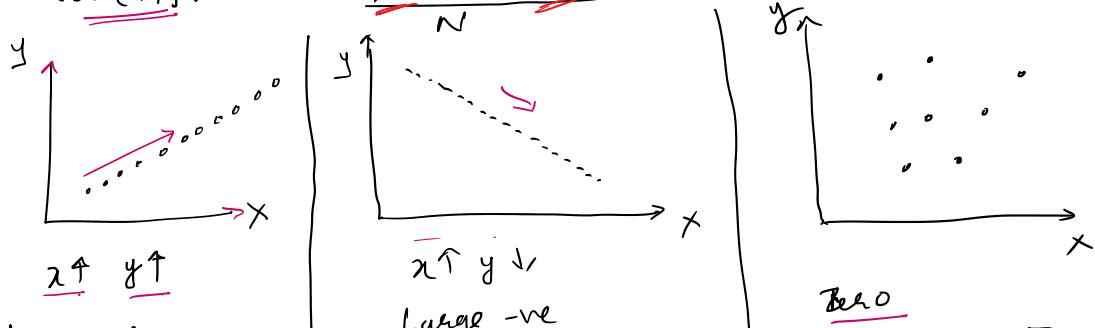
$$\bar{y} = \frac{1.91}{11}$$

Step 3 Adjusted Dataset (So the Dataset Mean is Zero).

$x = x - \bar{x}$	$x_i^0 - \bar{x}$	0.69	-1.31	0.39	0.09	1.29	0.49	0.19	-0.81	-0.31	-0.71
$y = y - \bar{y}$	$y_i^0 - \bar{y}$	0.49	-1.21	0.99	0.29	1.09	0.79	-0.31	-0.81	-0.31	-1.01

Note \rightarrow To find Covariance Matrix \rightarrow it is measure of the extent to which the corresponding elements from two set of ordered data move in same direction.

$$\text{cov}(x, y) = \frac{1}{N} \sum (x_i^0 - \bar{x})(y_i^0 - \bar{y})$$



$$\text{Covariance Matrix} \rightarrow C = \frac{1}{N} (x_{i=1} - \bar{x})(x_{i=1} - \bar{x})$$

$$C = \begin{bmatrix} \text{cov}(x, x) & \text{cov}(x, y) \\ \text{cov}(y, x) & \text{cov}(y, y) \end{bmatrix}$$

$$\begin{bmatrix} 0.61655555 & 0.61544444 \\ 0.61544444 & 0.71655556 \end{bmatrix}$$

$$\text{Covariance Matrix} \rightarrow C = \begin{bmatrix} \text{cov}(x, x) & \text{cov}(x, y) \\ \text{cov}(y, x) & \text{cov}(y, y) \end{bmatrix}$$

$$\text{cov}(x, x) = \frac{1}{N-1} \sum_{k=1}^N (x_{ik} - \bar{x})(x_{jk} - \bar{x})$$

Note $\text{cov}(x, y) = \text{cov}(y, x)$

Step 5: Compute Eigen Values & Vector.

$$C - \lambda I = 0$$

$$\begin{bmatrix} 0.61655555 - \lambda & 0.61544444 \\ 0.61544444 & 0.71655556 - \lambda \end{bmatrix} = 0$$

$$\lambda^2 - 1.3332\lambda + 0.0630244 = 0$$

$$\boxed{\begin{array}{l} \lambda_1 = 0.0490834 \\ \lambda_2 = 1.284627712 \end{array}} \quad \text{Eigen Values}$$

* The largest Eigen Value gives first PCA

Eigen Vector \rightarrow Let $u = \begin{bmatrix} v_1 \\ v_2 \end{bmatrix}$ and $v = \begin{bmatrix} v_1 \\ v_2 \end{bmatrix}$ be Eigen Vectors for λ_1

$$\begin{bmatrix} 0.56747215 & 0.61544444 \\ 0.61544444 & 0.66747216 \end{bmatrix} \begin{bmatrix} v_1 \\ v_2 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

$$0.56747215v_1 + 0.61544444v_2 = 0 \quad \text{--- (1)}$$

$$0.61544444v_1 + 0.66747216v_2 = 0 \quad \text{--- (2)}$$

Let $v_1 = 1$ in eq (1)

$$v_2 = -0.92205266$$

$$\therefore u = \begin{bmatrix} 1 \\ -0.92205266 \end{bmatrix} \quad \boxed{u_1 \quad u_2}$$

for λ_2

$$\begin{bmatrix} 0.61544444 & 0.61544444 \\ -0.66747216 & -0.56747215 \end{bmatrix} \begin{bmatrix} v_1 \\ v_2 \end{bmatrix} = 0$$

$$\begin{bmatrix} -0.66747216 \\ 0.61544444 \end{bmatrix} \quad \begin{bmatrix} 0.61544444 \\ -0.567472152 \end{bmatrix} \quad \left| \begin{bmatrix} v_1 \\ v_2 \end{bmatrix} = 0 \right.$$

$$-0.66747216v_1 + 0.61544444v_2 = 0 \quad \text{--- (3)}$$

$$0.61544444v_1 + -0.567472152v_2 = 0 \quad \text{--- (4)}$$

Consider Eq 3 put $v_2 = 1$

$$v_1 = 0.92205266$$

$$v = \begin{bmatrix} 0.92205266 \\ 1 \end{bmatrix} \quad \begin{array}{l} v_1 \\ v_2 \end{array}$$

(6) Normalize the vector \rightarrow to convert it into vector of unit length
How \rightarrow divide it by length of vector.

$$u = \frac{1}{\sqrt{v_1^2 + v_2^2}} \begin{bmatrix} v_1 \\ v_2 \end{bmatrix} = \frac{1}{1.360213596} \begin{bmatrix} 1 \\ -0.922052616 \end{bmatrix}$$

$$u = \begin{bmatrix} 0.7351786533 \\ -0.6778734 \end{bmatrix} \quad \checkmark$$

$$u = \frac{1}{\sqrt{v_1^2 + v_2^2}} \begin{bmatrix} v_1 \\ v_2 \end{bmatrix} = \frac{1}{1.360213596} \begin{bmatrix} 0.922052616 \\ 1 \end{bmatrix}$$

$$u = \begin{bmatrix} 0.6778734 \\ 0.7351786533 \end{bmatrix} \quad \checkmark$$

Since we got both the Eigen Vector, let us construct Feature Vector

$$\text{FeatureVector} = \begin{bmatrix} u \\ 0 \\ 0.6778734008 \\ 0.735178635 \end{bmatrix}$$

$$\text{Final Data} = \text{FeatureVector}^T \times \text{Data adjust}$$

$$= \begin{bmatrix} u \\ 0.7351786533 \\ 0.6778734008 \end{bmatrix} \times \begin{bmatrix} -0.6678734008 \\ 0.735178635 \end{bmatrix}$$

$$U = \begin{bmatrix} 0.7357786535 \\ 0.6778734008 \end{bmatrix} \quad \begin{bmatrix} -0.6678 + 5.9000 \\ 0.735778635 \end{bmatrix}$$

$$\begin{bmatrix} 0.69 & -1.31 & 0.39 & 0.09 & 1.29 & 0.49 & 0.19 & -0.21 & -0.31 & -0.71 \\ 0.49 & -1.21 & 0.99 & 0.29 & 1.09 & 0.79 & -0.31 & -0.51 & -0.31 & -0.71 \end{bmatrix}_{2 \times 10}$$

depended on λ_2

$U = \begin{bmatrix} 0.1751 \\ 0.827 \end{bmatrix}$ \Rightarrow $1D$ (converted) Data.

$U^T = \begin{bmatrix} 1 & -0.162 & -0.384 & -0.130 & 0.209 & -0.175 & 0.349 & -0.046 & -0.017 & -1.00 \\ 0.827 & -1.777 & -0.992 & 0.274 & 1.875 & 0.912 & -0.099 & 0.046 & 0.017 & 1.00 \end{bmatrix}_{10 \times 2}$

from λ_2 find value of (x_1, x_2) \Rightarrow will have combination of both

Note: Here only if λ_2 is used as it is greater than given first PCA

\therefore Calculating feature Vector for λ_2 let it be θ as calculated above

$$\theta = \begin{bmatrix} 0.92205266 \\ 1 \end{bmatrix} \quad \text{feature Vector} = \begin{bmatrix} 0.92205266 \\ 1 \end{bmatrix}$$

only θ

Final Subset = Feature Vector $^T \times$ Adjusted Data.

$$= [0.92205266 \ 1] ^T \begin{bmatrix} \dots \\ \dots \end{bmatrix}_{2 \times 10}$$

$$= \begin{bmatrix} \dots \\ \dots \end{bmatrix}_{1 \times 10}$$

Single Row

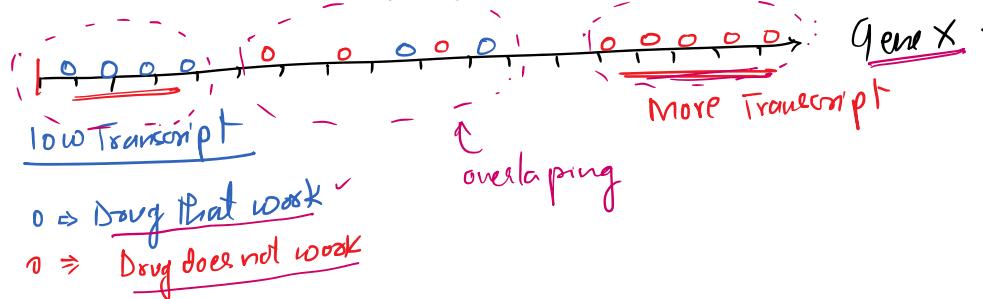
Original Dataset = 2×10 \Rightarrow Reduced Dimension
 Final Dataset = 1×10 \leftarrow from $2D$ to $1D$

$\text{PCA} \rightarrow$ LDA (Linear Discriminant Analysis)
SVD (Singular Value Decomposition)

LDA [Linear Discriminant Analysis].

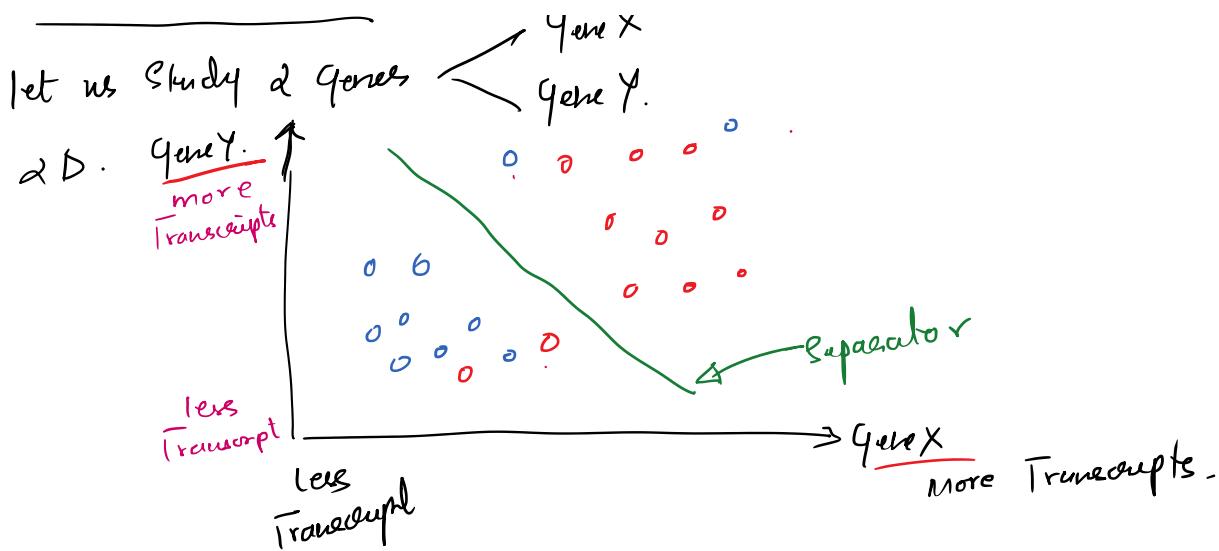
- * Suppose \Rightarrow We got a cancer drug ✓
 ↳ It works great for some people ✓
 ↳ But it makes it worse for other people ✓
- * How do we decide whom to give the drug? ?
- * May be gene study of patients will be of some help.

(1) Consider only one gene \Rightarrow Gene X. [I-D]



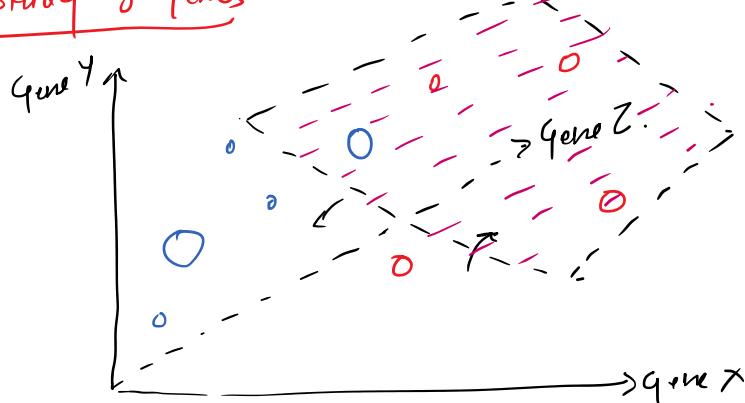
- * From above we can see In cases with More Transcript drug is not working.
- * Whereas In cases of low Transcript the drug tends to work.
- * Now there are few exceptions (there are overlaps).
- * Summary \Rightarrow Gene X does OK job but it has few overlaps.

* Can we do better?
 Let us Study of Genes $\begin{cases} \text{Gene X} \\ \text{Gene Y} \end{cases}$



It is better than studying only one gene but still there are overlapping.

Let us Study 3 Genes



In case of 3 genes we will need to represent the information in 3-D.

→ Here the separator is a plane.

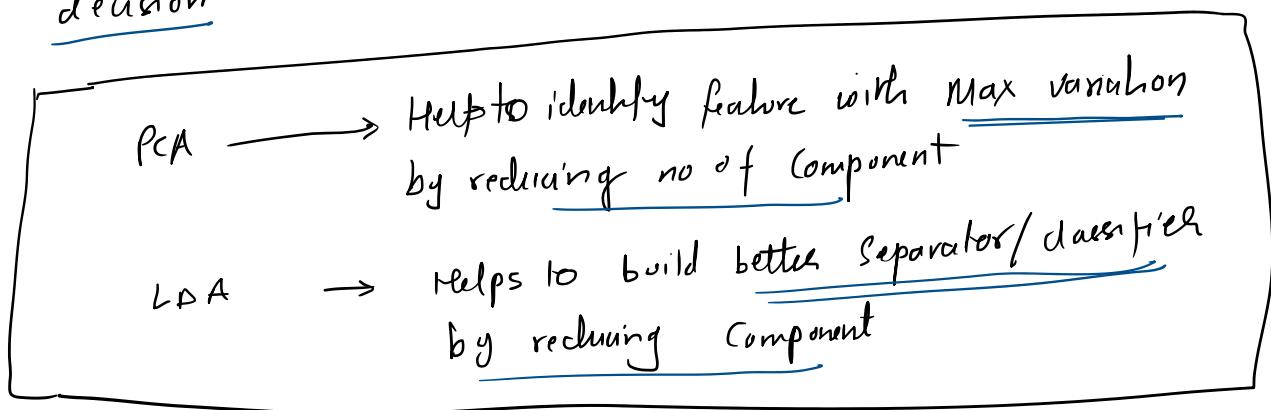
Suppose we want to study 4 genes $\xrightarrow{\text{Need}} \underline{\text{4-D Represent}}$
 $\rightarrow \text{Can't Draw 4D Graph}$.

* Here more the Dimension, More will be Complication in graph representation.

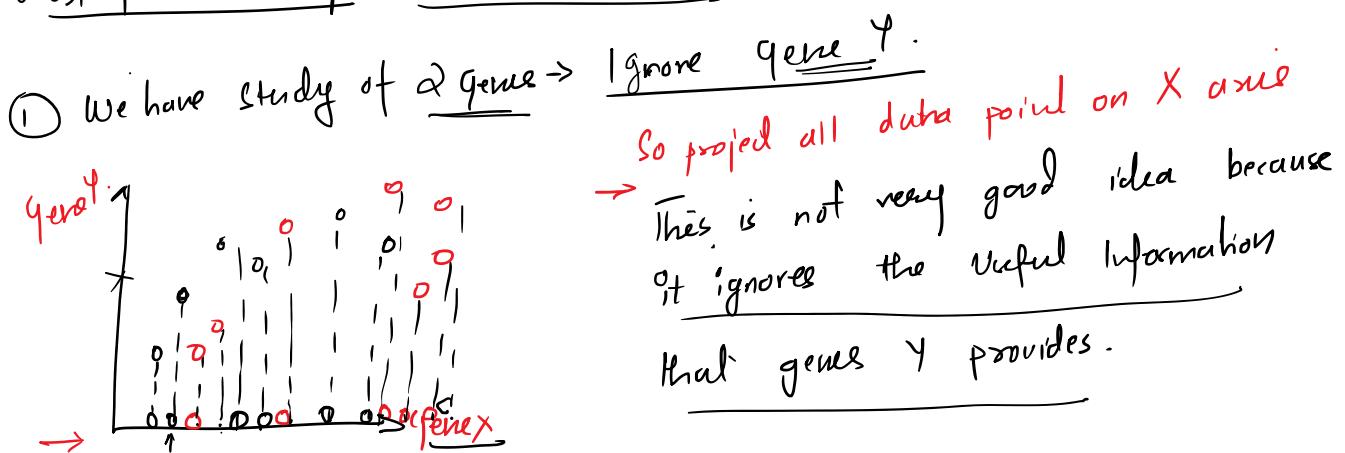
+ PCA reduces dimensions by focusing on the genes with most variance (More the variance, More the information).

- * Normally PCA is useful in plotting data with lot of dimension (or genes) onto simple XY plane \rightarrow
- * Here we are not interested much in identifying genes with most variance
Instead (LDA) we are interested in Maximizing Separability betn the two groups so we can make best classification decision.

Note



- * Worst possible ways to Reduce Dimension \Rightarrow



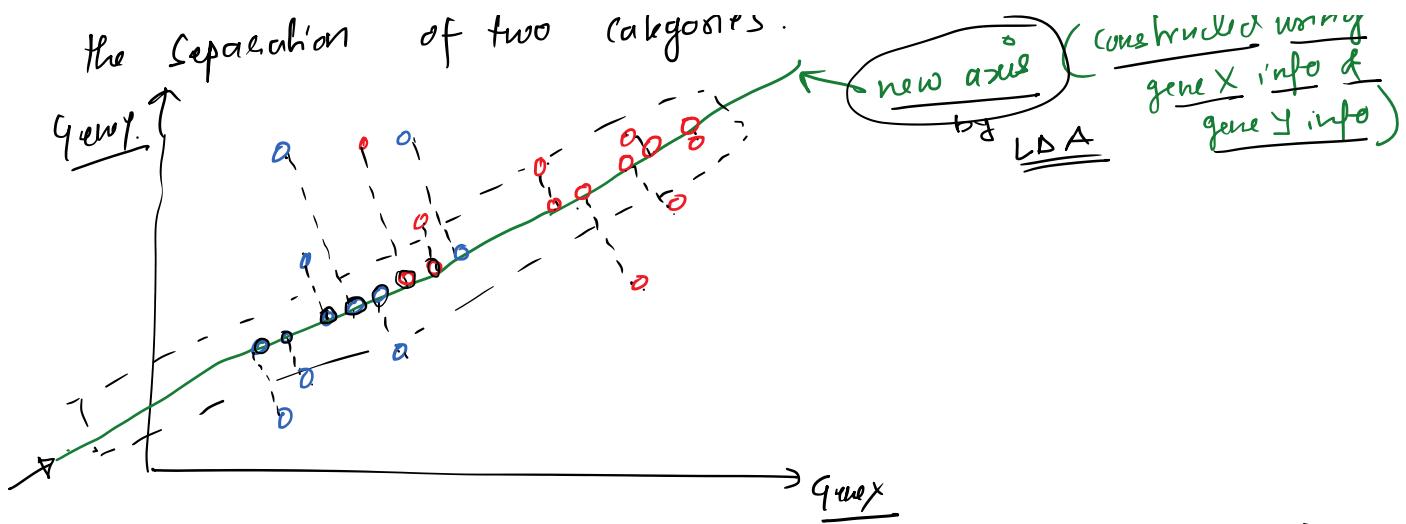
- * LDA provides better way

\rightarrow LDA reduces an α -D graph to 1-D graph.

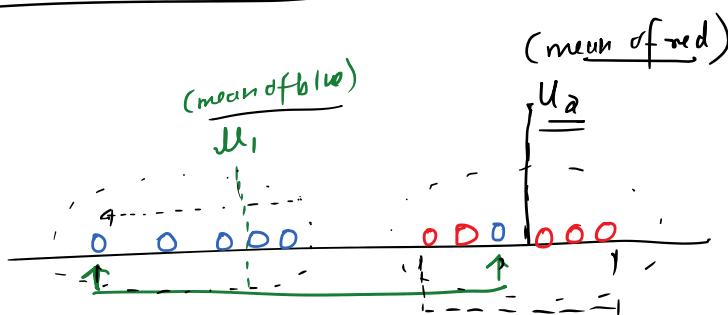
\rightarrow LDA uses both the genes to create a new axis and project the data on new axis in a way to maximize the separation of two categories.

\rightarrow new axis

constructed wrong gene X info &



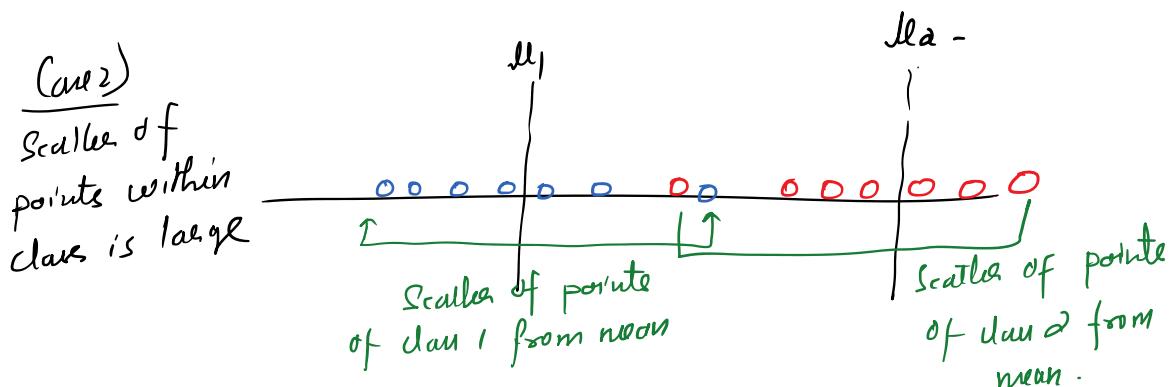
How LDA creates a New Axis's \rightarrow \bar{u}_1 and \bar{u}_2 are mean of respective class.



Condition) means of two classes are near.
 \bar{u}_1 \bar{u}_2
 \rightarrow the points in each class are widely scattered

The new axis is created according to two criteria
(considered simultaneously)

- ① Maximize the distance between mean of two classes.)
- ② Minimize the Variation (LDA calls it scatter) and is represented by S^2) for each category.

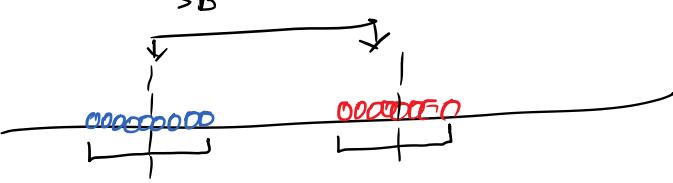


$S_B \rightarrow$ Between class scatter.

⇒ points of points

(con 3) Scatter of points
within class is small.

$S_B \rightarrow$ Between class scatter.



$\underline{S_W}$
within class scatter
Scatter of points is small within class.

Goal $\rightarrow \sqrt{S_B} \Rightarrow$ Maximize } Must be focused
 $\sqrt{S_W} \Rightarrow$ Minimize } on both simultaneously.

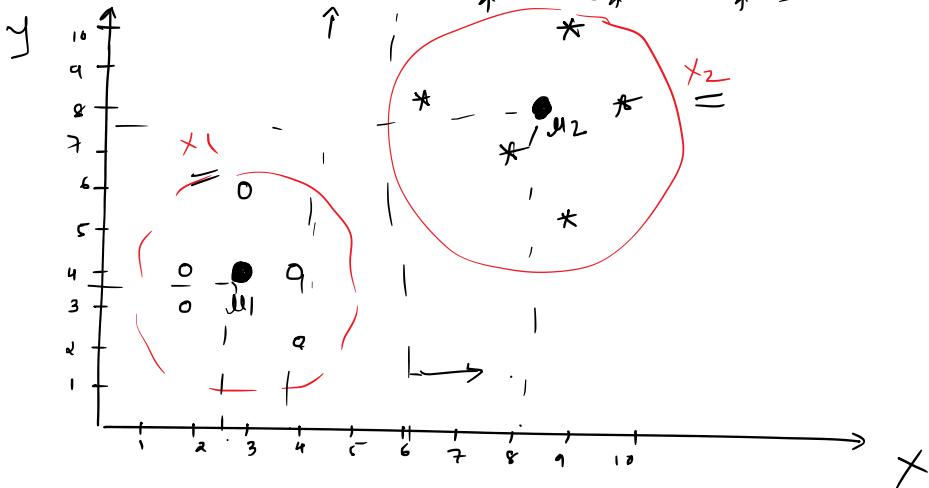
Mathematical Representn \Rightarrow
$$\frac{(m_1 - m_2)^2}{S_{\text{B}}^2 + S_{\text{W}}^2} = \frac{\text{Square of diff bet two mean}}{\text{Sum of scattered within each category/class}}$$

Ideally
$$\frac{(m_1 - m_2)^2}{S_{\text{W}}^2 + S_{\text{W}}^2} \Rightarrow \begin{cases} \text{Large} \\ \text{Small} \end{cases} \checkmark$$

Consider \Rightarrow Compute the (linear Discriminant projection) for the following two Dimensional Dataset -

$$\text{Class 1} \rightarrow \underline{\underline{X}_1} = (x_1, y_1) = \{ (4, 2), (2, 4), (2, 3), (3, 6), (4, 4) \} \Rightarrow \circ$$

$$\text{Class 2} \rightarrow \underline{\underline{X}_2} = (x_2, y_2) = \{ (9, 10), (6, 8), (9, 5), (8, 7), (10, 8) \} \Rightarrow *$$



Step 1 Find (two means)

$$\begin{aligned} \underline{\underline{\mu}}_1 &= \frac{1}{N_1} \sum x \in \underline{\underline{X}}_1 \\ &= \frac{1}{5} \left[\left(\frac{4}{2} \right) + \left(\frac{2}{4} \right) + \left(\frac{2}{3} \right) + \left(\frac{3}{6} \right) + \left(\frac{4}{4} \right) \right] = \frac{1}{5} \begin{pmatrix} 15 \\ 19 \end{pmatrix} \\ &= \underline{\underline{\underline{\underline{\mu}}}}_1 \end{aligned}$$

$$\begin{aligned} \underline{\underline{\mu}}_2 &= \frac{1}{N_2} \sum x \in \underline{\underline{X}}_2 \\ &= \frac{1}{5} \left[\left(\frac{9}{10} \right) + \left(\frac{6}{8} \right) + \left(\frac{9}{5} \right) + \left(\frac{8}{7} \right) + \left(\frac{10}{8} \right) \right] = \frac{1}{5} \begin{pmatrix} 42 \\ 38 \end{pmatrix} \\ &= \underline{\underline{\underline{\underline{\mu}}}}_2 \end{aligned}$$

Step 2 Covariance Matrix of Both class.

Covariance Matrix of Class 1. ($\underline{\underline{X}}_1$)

$$\underline{\underline{S}}_1 = \frac{1}{N-1} \sum_{x \in \underline{\underline{X}}_1} (x - \underline{\underline{\mu}}_1)(x - \underline{\underline{\mu}}_1)^T$$

$$\begin{aligned}
 S_1 &= \frac{1}{N-1} \sum_{x \in X_1} (x - \mu_1)(x - \mu_1)^T \\
 &= \frac{1}{4} \left[\left[\begin{pmatrix} 4 \\ 2 \end{pmatrix} - \begin{pmatrix} 3 \\ 3.8 \end{pmatrix} \right]^2 + \left[\begin{pmatrix} 2 \\ 4 \end{pmatrix} - \begin{pmatrix} 3 \\ 3.8 \end{pmatrix} \right]^2 + \left[\begin{pmatrix} 2 \\ 3 \end{pmatrix} - \begin{pmatrix} 3 \\ 3.8 \end{pmatrix} \right]^2 \right. \\
 &\quad \left. + \left[\begin{pmatrix} 3 \\ 6 \end{pmatrix} - \begin{pmatrix} 3 \\ 3.8 \end{pmatrix} \right]^2 + \left[\begin{pmatrix} 4 \\ 4 \end{pmatrix} - \begin{pmatrix} 3 \\ 3.8 \end{pmatrix} \right]^2 \right] \\
 &= \begin{bmatrix} 1 & -0.25 \\ -0.25 & 2.2 \end{bmatrix}_{//}
 \end{aligned}$$

Covariance Matrix of Second class

$$\begin{aligned}
 S_2 &= \frac{1}{N-1} \sum_{x \in X_2} (x - \mu_2)(x - \mu_2)^T \\
 &= \frac{1}{4} \left[\left[\begin{pmatrix} 9 \\ 10 \end{pmatrix} - \begin{pmatrix} 8.4 \\ 7.6 \end{pmatrix} \right]^2 + \left[\begin{pmatrix} 6 \\ 8 \end{pmatrix} - \begin{pmatrix} 8.4 \\ 7.6 \end{pmatrix} \right]^2 + \left[\begin{pmatrix} 9 \\ 5 \end{pmatrix} - \begin{pmatrix} 8.4 \\ 7.6 \end{pmatrix} \right]^2 \right. \\
 &\quad \left. + \left[\begin{pmatrix} 8 \\ 7 \end{pmatrix} - \begin{pmatrix} 8.4 \\ 7.6 \end{pmatrix} \right]^2 + \left[\begin{pmatrix} 10 \\ 8 \end{pmatrix} - \begin{pmatrix} 8.4 \\ 7.6 \end{pmatrix} \right]^2 \right] \\
 &= \begin{bmatrix} 2.3 & -0.05 \\ -0.05 & 3.3 \end{bmatrix}
 \end{aligned}$$

Step 3 Within Class Scatter Matrix (Minimize)

$$\begin{aligned}
 S_W &= S_1 + S_2 \\
 &= \begin{bmatrix} 1 & -0.25 \\ -0.25 & 2.2 \end{bmatrix} + \begin{bmatrix} 2.3 & -0.05 \\ -0.05 & 3.3 \end{bmatrix} \\
 &= \begin{bmatrix} 3.3 & -0.3 \\ -0.3 & 5.5 \end{bmatrix}_{//}
 \end{aligned}$$

Step 4 Between Class Scatter Matrix (Maximize).

$$\begin{aligned}
 S_B &= (\mu_1 - \mu_2)(\mu_1 - \mu_2)^T \\
 &= \left[\begin{pmatrix} 3 \\ 3.8 \end{pmatrix} - \begin{pmatrix} 8.4 \\ 7.6 \end{pmatrix} \right] \left[\begin{pmatrix} 3 \\ 3.8 \end{pmatrix} - \begin{pmatrix} 8.4 \\ 7.6 \end{pmatrix} \right]^T
 \end{aligned}$$

$$= \begin{bmatrix} -5.4 \\ -3.8 \end{bmatrix} \begin{bmatrix} -5.4 & -3.8 \end{bmatrix}$$

$$S_B = \begin{bmatrix} 29.16 & 20.52 \\ 20.52 & 14.44 \end{bmatrix}$$

Step 5 The LDA projection is then obtained as solution of the generalized Eigen Value Problem

$$S_{\omega}^{-1} S_B \omega = \lambda \omega$$

$$\Rightarrow |S_{\omega}^{-1} S_B - \lambda I| = 0$$

$$\Rightarrow \left| \left(\begin{bmatrix} 3.3 & -0.3 \\ -0.3 & 5.5 \end{bmatrix} \right)^{-1} \begin{pmatrix} 29.16 & 20.52 \\ 20.52 & 14.44 \end{pmatrix} - \lambda \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \right| = 0$$

$$\Rightarrow \left| \begin{pmatrix} 0.3045 & 0.0166 \\ 0.0166 & 0.1827 \end{pmatrix} \begin{pmatrix} 29.16 & 20.52 \\ 20.52 & 14.44 \end{pmatrix} - \lambda \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \right| = 0$$

$$= \left| \begin{pmatrix} 9.2213 - \lambda & 6.489 \\ 4.2339 & 2.9794 - \lambda \end{pmatrix} \right| \quad \text{--- } I$$

$$\Rightarrow (9.2213 - \lambda)(2.9794 - \lambda) - 6.489 \times 4.2339 = 0$$

$$\lambda^2 - 12.2007 \lambda = 0 \Rightarrow \lambda(\lambda - 12.2007) = 0$$

$$\lambda \Rightarrow \boxed{\begin{array}{l} \lambda_1 = 0 \\ \lambda_2 = 12.2007 \end{array}} \quad \text{2 Eigen Values.}$$

Now By Substituting $\lambda = \lambda_1$ in I we get

$$\boxed{\text{Eigen Vector } I = \underline{\omega_1} = \begin{pmatrix} -0.5755 \\ 0.8178 \end{pmatrix}}$$

Eigen Vector $I = \underline{\omega_1} = \begin{pmatrix} -0.5755 \\ 0.8178 \end{pmatrix}$

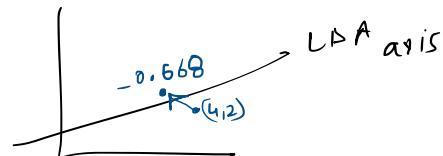
Eigen vector $\underline{w}_1 = \underline{\underline{w}}_1 = \begin{pmatrix} 0.8178 \\ 0.9088 \end{pmatrix}$

By

Substituting $\lambda = \lambda_1$ in I we get

Eigen Vector $\underline{w}_2 = \underline{\underline{w}}_2 = \begin{pmatrix} 0.9088 \\ 0.4173 \end{pmatrix}$

Slip 6 $\text{Final } \underline{y} = \underline{\underline{w}}^T \underline{\underline{x}}$
 Data \uparrow Input Data
 $\text{Projection Vector} \quad \downarrow$



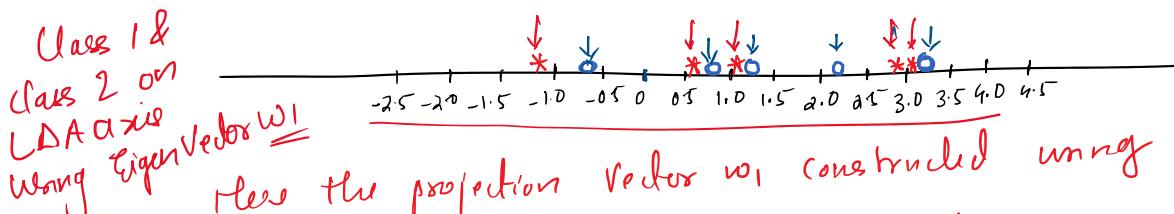
For Projection Vector \underline{w}_1

Final Data of $\underline{x}_1 = \underline{\underline{w}}_1^T \underline{\underline{x}}_1 = \begin{bmatrix} -0.5755 & 0.8178 \end{bmatrix} \begin{bmatrix} \underline{x}_1 \\ \underline{x}_2 \end{bmatrix}$

$= \begin{bmatrix} -0.668 \\ 2.12 \end{bmatrix} \quad 1.3022 \quad 3.18 \quad 0.9688$

Final Data of $\underline{x}_2 = \underline{\underline{w}}_1^T \underline{\underline{x}}_2 = \begin{bmatrix} -0.5755 & 0.8178 \end{bmatrix} \begin{bmatrix} \underline{x}_1 \\ \underline{x}_2 \end{bmatrix}$

Red $= \begin{bmatrix} 2.9985 & 3.0894 & -1.0905 & 1.1206 & 0.7874 \end{bmatrix}$



Smaller Eigen value leads to bad separability.
 $(\lambda_1=0)$

For Projection Vector \underline{w}_2

$$\text{final Data } X_1 = w_2^T X_1$$

$$= \begin{bmatrix} 0.9088 & 0.4173 \end{bmatrix} \begin{bmatrix} (1) \\ (2) \\ (3) \\ (4) \\ (5) \end{bmatrix}$$

Class 1 on $= \begin{bmatrix} 4.4698 & 3.4868 & 3.0645 & 5.2302 & 5.3044 \end{bmatrix}$

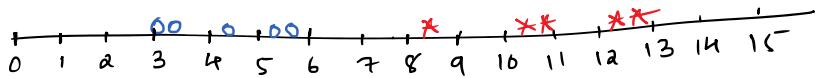
LDA using w_2 vector ✓

$$\text{Final Data } X_2 = w_2^T X_2$$

$$* = \begin{bmatrix} 0.9088 & 0.4173 \end{bmatrix} \begin{bmatrix} (9) \\ (8) \\ (5) \\ (7) \\ (6) \end{bmatrix}$$

Class 2 $\rightarrow = \begin{bmatrix} 12.3522 & 8.7912 & 10.2657 & 10.1951 & 12.4264 \end{bmatrix}$

on LDA
using w_2 vector



Here the projection vector corresponding to larger eigen value λ_2
leads to good separability.

This is LDA \rightarrow when the 2D data points reduced
to 1-D data points

\rightarrow and also using projection vector were
able to predict a good separator
of classes.

SVD [Singular Value Decomposition]

- * We normally use 2D matrix to represent Data Values where column represents features and row represents samples data points
- * Matrix computation with all the values in matrix sometime become redundant or computationally expensive.
- * We need to represent matrix in a form such that the most important part of matrix which is needed for further computation could be extracted easily.
- * This can be done by SVD

SVD Theorem \Rightarrow A rectangular matrix A_{mn} can be decomposed it into product of 3 matrices.

$$A_{mn} = U \sum_{m \times m}^{m \times n} V^T$$

$$\begin{matrix} m \times m & m \times n \\ \downarrow m \times n & \downarrow m \times n \end{matrix}$$

where $U_{m \times m} \Rightarrow$ Orthogonal matrix

$\sum_{m \times n} \Rightarrow$ Diagonal matrix of Eigen Values

$V^T_{n \times n} \Rightarrow$ Transpose of orthogonal matrix $V_{n \times n}$

Columns of U are the orthonormal Eigen Vectors of AA^T

& columns of V are the orthonormal Eigen Vectors of AA^T & \sum is diagonal matrix.

The elements of \sum are the square roots of Eigen values of

U & V in decreasing order.

Ex Find the SVD of $\underline{A} = \begin{bmatrix} 3 & 1 & 1 \\ -1 & 3 & 1 \end{bmatrix}_{2 \times 3}$

Note we have to Decompose the matrix A into $\underline{\underline{U}} \underline{\Sigma} \underline{V^T}$.

To find $\underline{\Sigma}$
 \underline{U}
 $\underline{V^T}$.

Solution \rightarrow To find $\underline{U} := \underline{\underline{A}} \underline{A^T} = \begin{bmatrix} 3 & 1 & 1 \\ -1 & 3 & 1 \end{bmatrix} \begin{bmatrix} 3 & -1 \\ 1 & 3 \\ 1 & 1 \end{bmatrix} = \begin{bmatrix} 11 & 1 \\ 1 & 1 \end{bmatrix}$

for Eigen Values $\Rightarrow \begin{vmatrix} 11-\lambda & 1 \\ 1 & 11-\lambda \end{vmatrix} = 0 \quad \text{--- (I)}$

$$\Rightarrow \lambda^2 - 22\lambda + 120 = 0$$

$$\boxed{\lambda_1 = 10} \quad \boxed{\lambda_2 = 12}$$

For Eigen Vector $\rightarrow \begin{bmatrix} 11-\lambda & 1 \\ 1 & 11-\lambda \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$

(one) \Rightarrow Substitute $\lambda = \lambda_1 = 10$

$$\begin{bmatrix} 1 & 1 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

$$\frac{x_1}{1} = -\frac{x_2}{1} = 1$$

$$\therefore x_1 = 1 \quad x_2 = -1$$

Eigen Vector $x_1 = \underline{\underline{\begin{bmatrix} 1 \\ -1 \end{bmatrix}}}$

(contd) Substitute $\lambda = \lambda_2 = 10$

$$\begin{bmatrix} 11-\lambda & 1 \\ 1 & 11-\lambda \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

$$\begin{bmatrix} -1 & 1 \\ 1 & -1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

Eigen Vector $X_2 = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$

* In U the Eigen Vector generated by larger Eigen value will be the first column.

$$\therefore U = \begin{bmatrix} 1 & 1 \\ 1 & -1 \end{bmatrix}$$

\uparrow \uparrow
Eigen vector Eigen vector
of $\lambda=12$ of $\lambda=10$

Now we need to Normalize the matrix \Rightarrow divide by length of respective Vector.

*
$$U = \begin{bmatrix} 1/\sqrt{2} & 1/\sqrt{2} \\ 1/\sqrt{2} & -1/\sqrt{2} \end{bmatrix}$$

\uparrow \uparrow
length of length of Vector
Vector $\lambda=12$ $\lambda=10$
 $\sqrt{1^2+1^2}$
 $= \sqrt{2}$ $\sqrt{1^2+(-1)^2} = \frac{1}{\sqrt{2}}$

Slope 2 To find V

$$A^T A = \begin{bmatrix} 3 & -1 \\ 1 & 1 \end{bmatrix} \begin{bmatrix} 3 & 1 & 1 \\ -1 & 3 & 1 \end{bmatrix} = \begin{bmatrix} 10 & 0 & 2 \\ 0 & 10 & 4 \\ 2 & 4 & 2 \end{bmatrix}$$

These are Eigen Values

$$\lambda^3 - 22\lambda^2 + 132\lambda = 0$$

$$\begin{array}{l} \therefore \lambda_1 = 0 \\ \lambda_2 = 10 \\ \lambda_3 = 12 \end{array}$$

We need to find Eigen Vectors for the 3 Eigen values.

(case 1) Eigen Vector for $\lambda_1=0$

$$\rightarrow \begin{bmatrix} 10 & 0 & 2 \\ 0 & 10 & 4 \\ 2 & 4 & 2 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = 0$$

By Cramers Rule

$$\frac{x_1}{|10 \ 0 \ 2|} = \frac{x_2}{|0 \ 10 \ 4|} = \frac{x_3}{|10 \ 0 \ 10|}$$

$$= \frac{x_1}{-20} = -\frac{x_2}{40} = \frac{x_3}{10} = \frac{-1}{20}$$

$$x_1 = 1 \quad x_2 = 2 \quad x_3 = -5$$

$$\text{Eigen Vector } X_1 = \begin{bmatrix} 1 \\ 2 \\ -5 \end{bmatrix}$$

(case 2) Eigen Vector for $\lambda = 10$

$$\begin{bmatrix} 0 & 0 & 2 \\ 0 & 0 & 4 \\ 2 & 4 & -8 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}$$

By Cramers Rule

$$\frac{x_1}{|0 \ 0 \ 2|} = -\frac{x_2}{|0 \ 0 \ 4|} = \frac{x_3}{|2 \ 4 \ -8|} \Rightarrow \frac{x_1}{-16} = -\frac{x_2}{-8} = \frac{x_3}{0} = \frac{-1}{8}$$

$$\therefore x_1 = 2 \quad x_2 = -1 \quad x_3 = 0$$

$$\text{Eigen Vector } x_2 = \begin{bmatrix} 2 \\ -1 \\ 0 \end{bmatrix}$$

Case 3) Eigen Vector for $\lambda_3 = 12$

$$\begin{bmatrix} -2 & 0 & 2 \\ 0 & -2 & 4 \\ 2 & 4 & -10 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}$$

By Cramer's Rule

$$\frac{x_1}{\begin{vmatrix} 0 & 2 \\ -2 & 4 \end{vmatrix}} = \frac{-x_2}{\begin{vmatrix} -2 & 2 \\ 0 & 4 \end{vmatrix}} = \frac{x_3}{\begin{vmatrix} -2 & 0 \\ 0 & -2 \end{vmatrix}} \Rightarrow \frac{x_1}{4} = \frac{-x_2}{-8} = \frac{x_3}{4} = \frac{1}{4}$$

$$\therefore x_1 = 1 \quad x_2 = -2 \quad x_3 = 1$$

$$\text{Eigen Vector } x_3 = \begin{bmatrix} 1 \\ -2 \\ 1 \end{bmatrix}$$

$$\text{Now } V = \begin{bmatrix} 1 & 2 & 1 \\ 2 & -1 & 2 \\ 1 & 0 & 5 \end{bmatrix}$$

\uparrow \uparrow \uparrow
 Eigen Vector for $\lambda = 12$ Eigen Vector for $\lambda = 10$ Eigen Vector for $\lambda = 0$.

Now Normalizing the matrix \rightarrow Divide by length of Vector.

$$V = \begin{bmatrix} \sqrt{56} & 2/\sqrt{56} & 1/\sqrt{30} \\ 2/\sqrt{56} & -1/\sqrt{5} & 2/\sqrt{30} \\ 1/\sqrt{56} & 0 & -5/\sqrt{30} \end{bmatrix}$$

\uparrow \uparrow \uparrow
 length of Vector length of Vector length of Vector

$$= \begin{bmatrix} \sqrt{11}/\sqrt{56} & 2/\sqrt{56} & 1/\sqrt{30} \end{bmatrix}$$

$$V^T = \begin{bmatrix} 1/\sqrt{6} & 2/\sqrt{6} & 1/\sqrt{6} \\ 2/\sqrt{5} & -1/\sqrt{5} & 0 \\ 1/\sqrt{30} & 2/\sqrt{30} & -5/\sqrt{30} \end{bmatrix}$$

Step 3 To find Σ (or D)

$$\Sigma = \begin{vmatrix} \sqrt{12} & 0 & 0 \\ 0 & \sqrt{10} & 0 \\ 0 & 0 & \sqrt{6} \end{vmatrix}_{3 \times 3}$$

\uparrow
diag matrix \Rightarrow the diag elements are square root of eigen values in decreasing order.

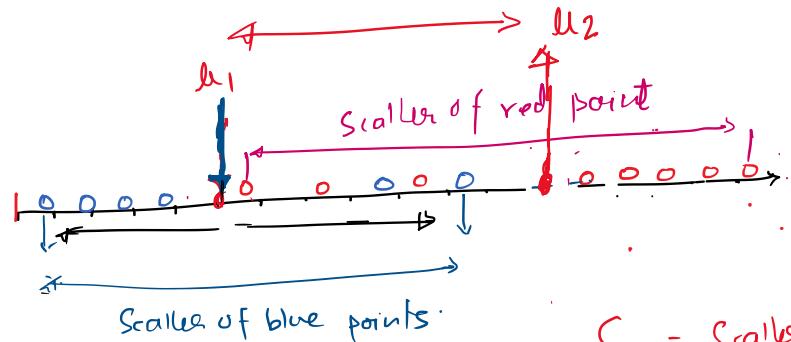
$$A = U \Sigma V^T$$

$$A = \begin{bmatrix} 1/\sqrt{2} & 1/\sqrt{2} \\ 1/\sqrt{2} & -1/\sqrt{2} \end{bmatrix} \begin{bmatrix} \sqrt{12} & 0 & 0 \\ 0 & \sqrt{10} & 0 \\ 0 & 0 & \sqrt{6} \end{bmatrix} \begin{bmatrix} 1/\sqrt{6} & 2/\sqrt{6} & 1/\sqrt{6} \\ 2/\sqrt{5} & -1/\sqrt{5} & 0 \\ 1/\sqrt{30} & 2/\sqrt{30} & -5/\sqrt{30} \end{bmatrix}$$

$U \quad \Sigma \quad V^T$

$=$

- Consider
2 groups
* Group of Blue points
* Group of Red points



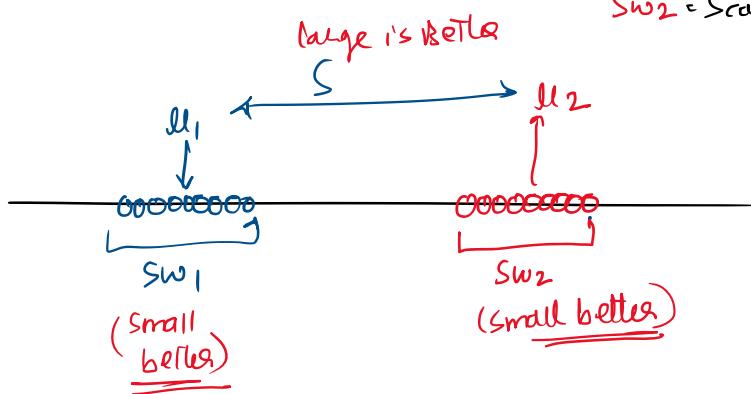
$$\begin{aligned} u_1 &= \text{Mean of blue point} \\ u_2 &= \text{mean of red point} \end{aligned}$$

S_B = Scatter between two groups
= Should be large

S_W → Scatter within a group.
= Should be small

S_{W1} = Scatter within group 1
(Blue points)

S_{W2} = Scatter within group 2
(Red points)



Note

In math Any Square Matrix A can be represented as

3 component matrix

$$A = U D U^T$$

(Decomposed)

→ $A = \underbrace{U}_{\text{Eigen Vector}} \underbrace{D}_{\text{Diagonal matrix of Eigen Value}} \underbrace{U^T}_{\text{Transpose of Eigen Vector}}$

This makes majority of operations on matrix to be easy as operation on Diagonal matrix is easy.

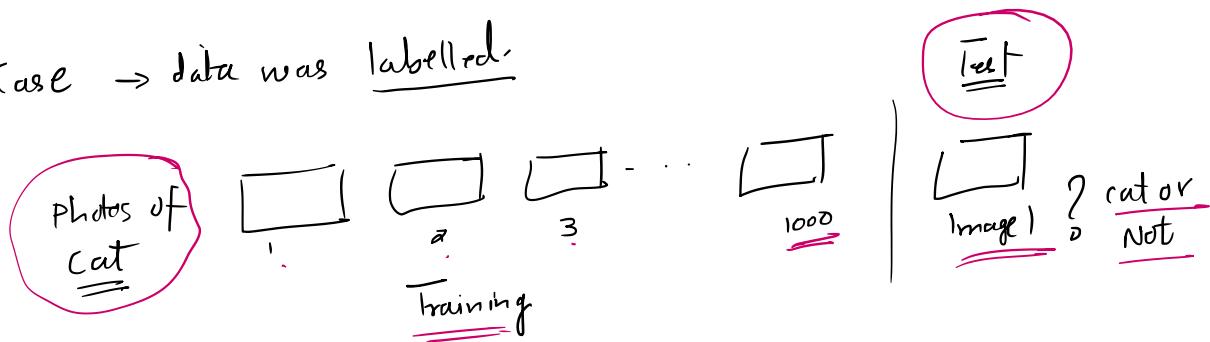
as operation on Diagonal matrix is easy .

In Real life \rightarrow Database \Rightarrow Matrix \Rightarrow Rectangular Matrix

Module 5

Clustering →

In earlier case → data was labelled.

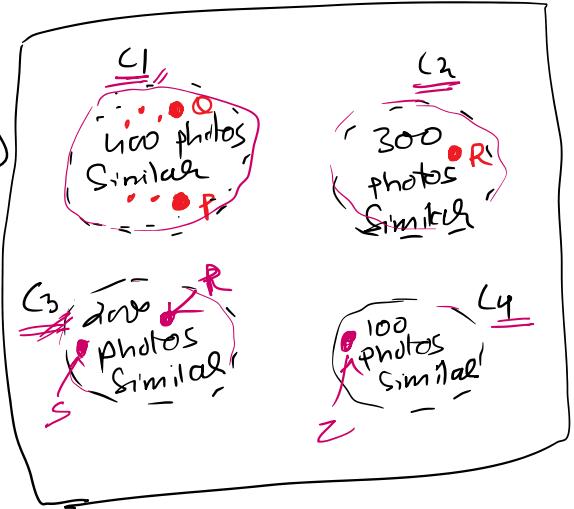


Supervised learning

Now consider → 1000 photos (not labelled)

- ✓ 400 → flowers ✓
- ✓ 300 → cars ✓
- ✓ 200 → houses ✓
- ✓ 100 → animal ✓

(Not Available)



Clustering is the task of dividing the population or data points into no of groups such that the data points in the same group are more similar to other data points in same group than those in other group.

Clustering →
 * Aim is to Segregate groups with Similar traits/features
and assign them into clusters."

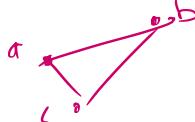
Distance Metric : It is a parameter that tells how close two points are.

Given element a, b, c in a set, a distance metric is defined as a function with following properties:

(1) Non-negativity $\rightarrow d(a, b) \geq 0$

(2) Symmetry $\rightarrow d(a, b) = d(b, a)$

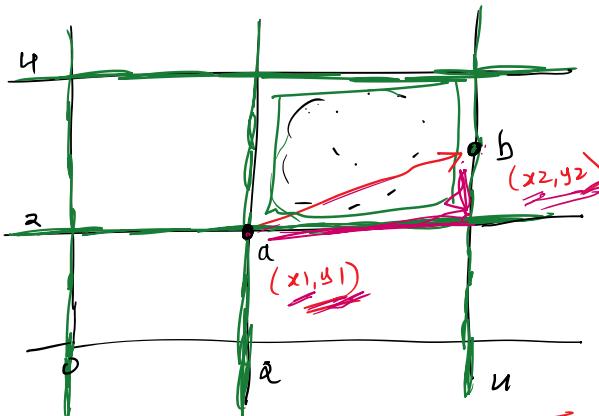
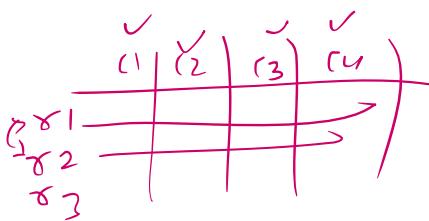
(3) Triangle Inequality $d(a, c) \leq d(a, b) + d(b, c)$. $a, b, c \in S$



Different Distance Metric

(1) Euclidean Distance \rightarrow Most commonly used

\rightarrow When data is dense or continuous, this is the best proximity measure.



$$\text{Euclidean Distance (vector)} = d(a, b) = \sqrt{\sum_{i=1}^n (a_i - b_i)^2}$$

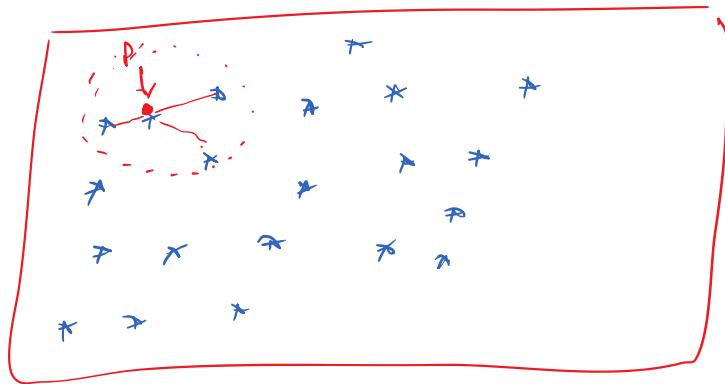
$$= \sqrt{(x_1 - x_2)^2 + (y_1 - y_2)^2}$$

If for a point $P = (p_1, p_2, \dots, p_n)$ and $q = (q_1, q_2, \dots, q_n)$.

$$d(p, q) = d(q, p) = \sqrt{(p_1 - q_1)^2 + (p_2 - q_2)^2 + \dots + (p_n - q_n)^2}$$

$$= \sqrt{\sum_{i=1}^n (p_i - q_i)^2}$$

$$= \sqrt{\sum_{i=1}^n (p_i - q_i)^2}$$

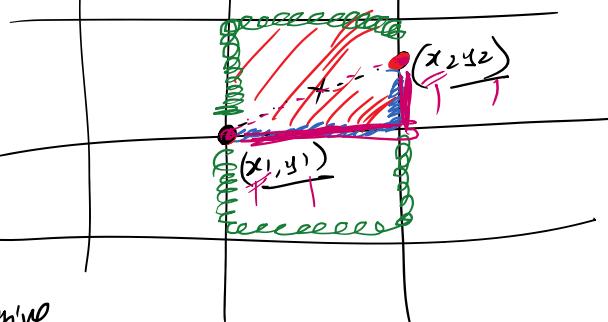


② Manhattan Distance

* If we are trying to find the distance betⁿ two buildings that are separated by several blocks. We can only walk through Sidewalk that are parallel and cannot walk diagonally through buildings.

In such case
we cannot use
Euclidean
distance.

Here Euclidean
distance will not give
realistic estimate for the distance



Manhattan distance -

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

* Manhattan Distance is a metric in which the distance between two points is the sum of absolute difference of their Cartesian coordinates.
* Total sum of difference betⁿ the x coordinate & y coordinate.

* Manhattan distance metric is also known as Manhattan length, Rectilinear distance, L1 distance, L1 Norm, City block distance

taxi cab metric.,

$$\begin{aligned} a &= (a_1, a_2, \dots, a_n) \\ b &= (b_1, b_2, \dots, b_n) \end{aligned}$$

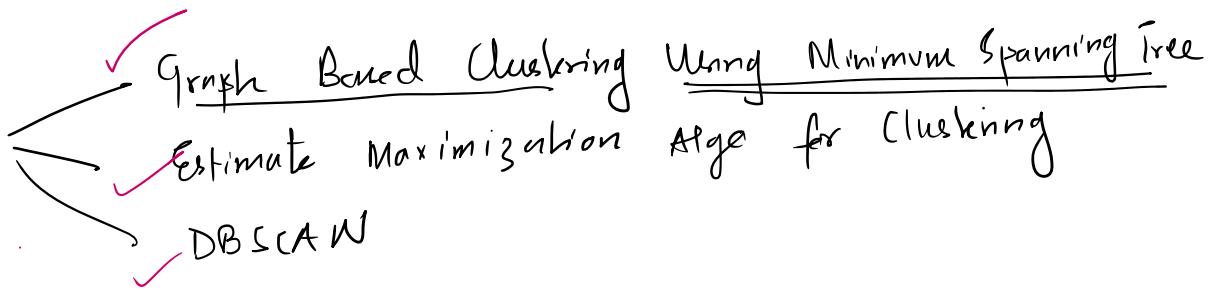
... has n features -

* If the data points has n features -

Manhattan Distance

$$d(a, b) = |a_1 - b_1| + |a_2 - b_2| + |a_3 - b_3| + \dots + |a_n - b_n|$$
$$= \sum_{i=1}^n |a_i - b_i|$$

Syllabus



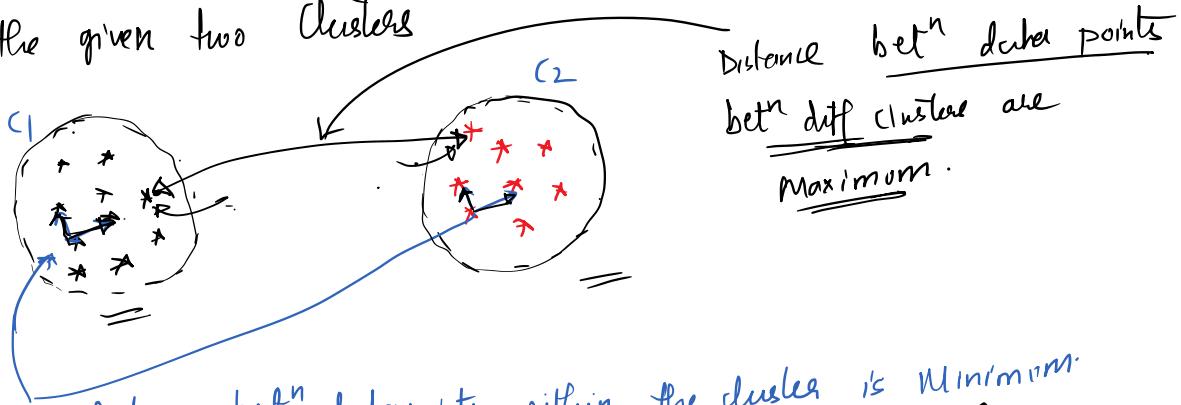
① Graph Based Clustering →

(1) First consider How to construct Graph of given data.

→ K Neighbourhood graph ✓

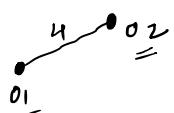
→ Epsilon Neighbourhood graph ✓

Consider the given two clusters

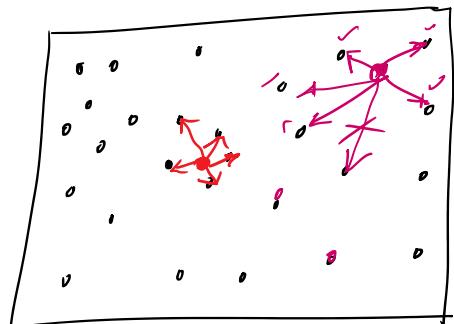


* In Graph Based Clustering the objects are represented as nodes in graph.

* The weight of edge/branch bet^n the objects
defines the distance bet^n the objects



① K Neighbourhood Approach



* To find distance of a point with all other points will be computationally challenging.

* We will consider nearest K points

Let K=5. ✓

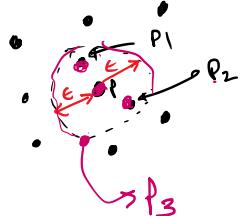
So we will calculate distance of
nearest 5 data points are will be .. .

$n = 1000$

So we will minimum -
nearest 5 data points are will be
labelled as Nearest Neighbour of node

Challenge → choosing K value is challenging task
in graph construction.

(a) Epsilon Neighborhood Graph → Objects within a radius of epsilon from
a given object are considered nearest
neighbours.



Here P_1 & P_2 are neighbours of P_0

Challenge → Selecting a proper epsilon value.

* Graph Based Clustering Using Minimum Spanning Tree →

Algo

1. Determine MST of given graph. ✓
2. Delete (edges) branch iteratively. ✓
3. Each connected component is a cluster. ✓

→ Subgraph of given graph
→ Some No of vertex
→ and no cycle
→ Subgraph is connected

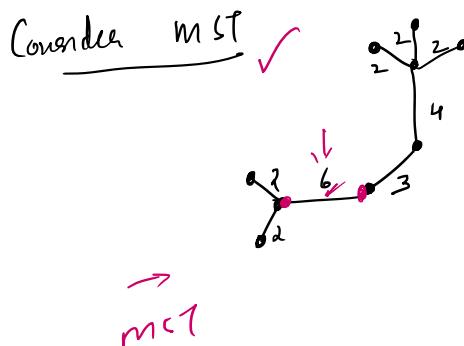
Different Strategies to Delete Branches.

- ① Delete Branch with max weight → ✓
- ② Delete Inconsistent Branch → ✓
- ③ Delete by Analysis of weight. ✓

① Deletion of Branch with max weight

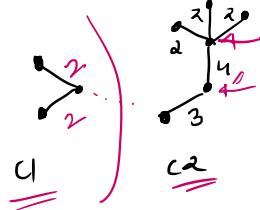
1. In each step create two clusters by deleting the branch with

- Y
1. In each step create two clusters by deleting the branch with max weight
 2. Repeat until the desired no of clusters is reached

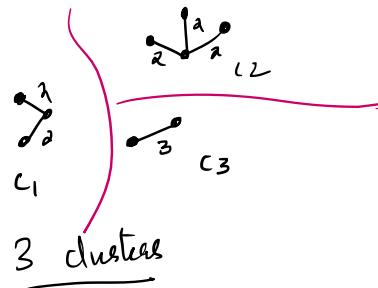


Given no of cluster desired is 3

Step 1) delete edge with weight 6



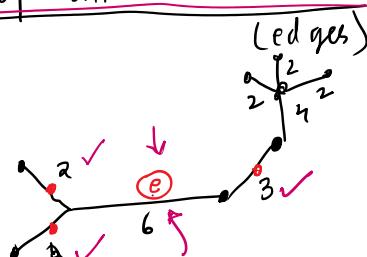
Step 2) delete edge with weight 4.



(2) Delete Inconsistent Branches -

- * A branch is inconsistent if the corresponding weight (de) is much larger than the reference value \bar{d}_e
- * The reference value \bar{d}_e can be defined by the average weight of all the branches adjacent to edge e

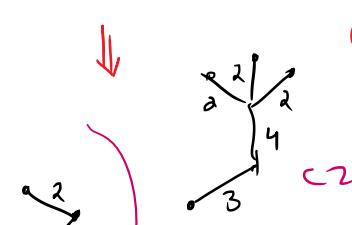
Consider MST:

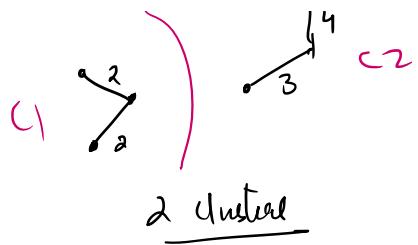


$$\text{Here } \bar{d}_e = \frac{2+2+3}{3} = \frac{7}{3} = 2.3$$

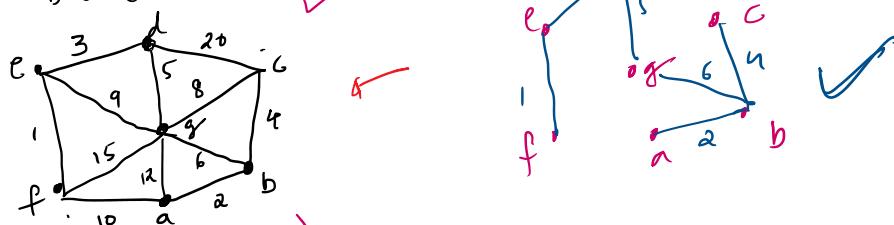
Here $d_e = 6 > \bar{d}_e$

So delete the edge e,

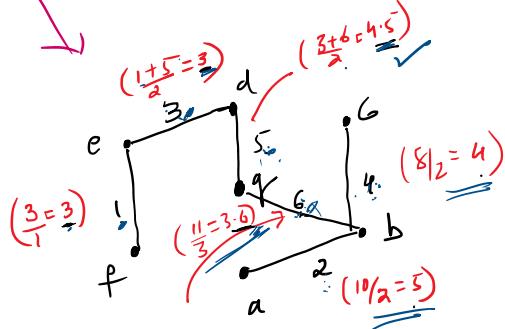




Q) Construct MST and provide clustering of graph Using
Inconsistent branches.



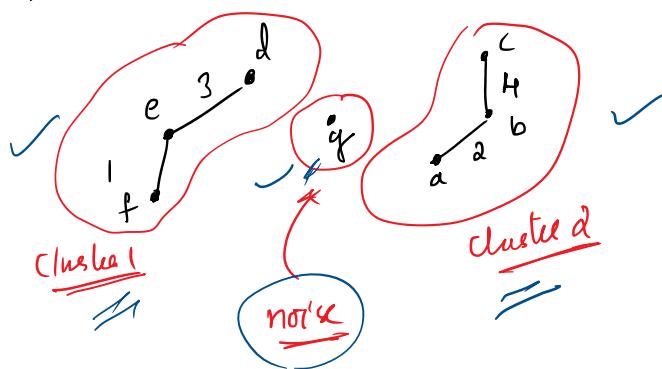
Step 1) MST of the Graph



Step 2) For deletion we are using Inconsistent branch approach.

Here edge with 5 > reference value (4.5) so delete edge 5

Also edge with 6 > reference value (3.6) so delete edge 6.



② Estimation Maximization Algo

- This Algo is for Maximum likelihood Estimation in presence of latent variable (missing).
- The most discussed application of EM Algorithm is used for Clustering with Mixture Model. (where data points have missing value)

- * (E-step) Estimation → Estimate are expectations for machine & then classify the data into some classes.
- * (M-step) Maximization → Whatever we estimated should now be maximized.
- * Iteratively repeat E-step and M-step until the values converges (no difference)

Imagine There are 2 coins A & B (Biased)
 Here One is likely to get Head more no of times and other is likely to get tail more no of tail

- * You pick one at random and toss it
which one was it

- Sol:
- Repeat this 1 time:
 - * Pick a coin randomly (known whether A or B)
 - * Toss 10 times
 - * Read No of Head & Tail

Consider the result

	Coin A	Coin B
First Round		5H, 5T
Second Round	9H, 1T	

Consider the

First Round	5H, 5T
Second Round	9H, 1T
Third Round	8H, 2T
Fourth Round	4H, 6T
Fifth Round	7H, 3T
Total	24H 6T [9H, 11T]

Probability of getting Head with Coin A = $\frac{14}{30} = \underline{\underline{0.80}}$

Coin A yields head 80% of time and tail 45% of time.

- * What if only results are given.
we need to guess the percentage of heads that each coin yields & Part I
Also need to guess which coin was picked at each round of toss. Part B

→ (on Not given)

First Round	(5H, 5T)	?
Second Round	9H, 1T	?
Third Round	8H, 2T	?
Fourth Round	4H, 6T	?
Fifth Round	7H, 3T	?

Revert

Initial Assumption

Assume \rightarrow the probability of Head for Coin A = 60%. $\Rightarrow P_A = \underline{\underline{0.6}}$
Coin B = 50%. $\Rightarrow P_B = \underline{\underline{0.5}}$

First Round \Rightarrow 

Compute the Likelihood that it was coin A & coin B Using

Binomial Distribution i.e.

$$P(k) = \binom{n}{k} p^k (1-p)^{n-k}$$

n = no of toss
k = successful (head)

Q => is probability on n trials with k success

Likelihood of A = $P_A(h)^k (1 - P_A(h))^{n-k}$

Given n=10, k=5
 $= (0.6)^5 \times (1 - 0.6)^{10-5} = 0.0007962624$
Tail

Likelihood of B = $P_B(h)^k (1 - P_B(h))^{n-k}$

n=10, k=5
 $= (0.5)^5 \times (1 - 0.5)^{10-5} = 0.0009765625$

Normalize by Using $\frac{A}{A+B}$ [To convert likelihood into Probability]

for coin A = $\frac{0.0007962624}{(0.0007962624 + 0.0009765625)}$

= 0.45

coin B = $\frac{0.0009765625}{(0.0007962624 + 0.0009765625)}$

= 0.55

Estimating Likely No of head & tails for coin A & B for
First Round.

for coin A = $0.45 \times 5 = 2.2 H$
 $= 0.45 \times 5 = 2.2 T$

for coin B = $0.55 \times 5 = 2.8 H$
 $= 0.55 \times 5 = 2.8 T$

Consider Round 2 (9H, 1T)

Consider Round 2 (9H, 1T)

$$\text{Likelihood of } A = \frac{0.6^9 (1-0.6)^{10-9}}{0.6^9 (1-0.5)^{10-9}} = \underline{\underline{0.0040310784}}$$

$$\text{Likelihood of } B = \underline{\underline{0.009765625}}$$

Normalize A & B

$$P(A) = 0.80$$

$$\frac{A}{A+B}$$

$$P(B) = 0.20$$

$$\frac{B}{A+B}$$

Estimate likely no. of Head & Tail for coin A & B in Round 2

$$\begin{array}{l} \text{No. of H for coin A} = \frac{0.80 \times 9}{0.80 + 0.20} = 7.2H \\ \text{T " " A} = \frac{0.20 \times 1}{0.80 + 0.20} = 0.8T \end{array}$$

$$\begin{array}{l} \text{No. of H for coin B} = \frac{0.20 \times 9}{0.20 + 0.80} = 1.8H \\ \text{T " " B} = \frac{0.80 \times 1}{0.20 + 0.80} = 0.8T \end{array}$$

After 5 Rounds

	Coin A	Coin B
Round 1	2.2H, 2.2T	2.8H, 2.8T.
Round 2	7.2H, 0.8T	1.8H, 0.2T.
Round 3	5.9H 1.5T	2.1H 0.5T
Round 4	1.4H 2.1T	2.6H 3.9T.
Round 5	4.5H 1.9T	2.5H 1.1T
Total	<u>21.3H</u> <u>8.6T</u>	<u>11.7H</u> <u>8.4T</u>

Probability of getting Head

$$Q_A = \frac{21.3}{21.3 + 8.6} = \underline{\underline{0.71}}$$

$$Q_B = \frac{11.7}{11.7 + 8.4} = \underline{\underline{0.58}}$$

After 1 step of E & M

$$\begin{array}{l} Q_A = 0.71 \\ Q_B = 0.58 \end{array}$$

... in ... Similar

→

Repeat until New Value of $\overbrace{\theta_A \& \theta_B}^x$ is more similar to Equal.

In this Example

The Value Converge to

$$\left. \begin{array}{l} \theta_A = 0.8 \\ \theta_B = 0.52 \end{array} \right\} \checkmark \checkmark \checkmark$$

Part 1 Answered } Percentage of Head Coin A yields = 80%
 Coin B yields = 52%

$$\left. \begin{array}{l} \theta_A = 0.8 \\ \theta_B = 0.52 \end{array} \right\}$$

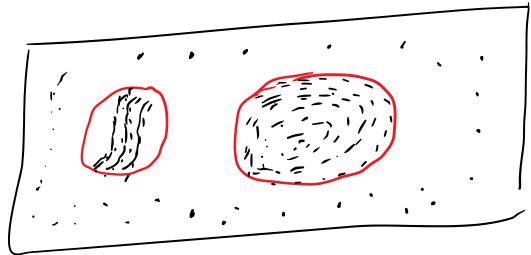
Now Consider the Result →

	H	T	Which coin?
Round 1	5 ✓	5 ✓	$P(H) = \frac{5}{10} = 0.5$ Coin B
Round 2	9	1	$P(H) = \frac{9}{10} = 0.9$ Coin A
Round 3	8	2	$P(H) = \frac{8}{10} = 0.8$ Coin A
Round 4	4	6	$P(H) = \frac{4}{10} = 0.4$ Coin B
Round 5	7	3	$P(H) = \frac{7}{10} = 0.7$ Coin A

Not to Part 2

DBSCAN (Density Based Spatial Clustering of Application with Noise)

- * In Dense Region there is possibility of cluster than in Sparse Region
- * Density based approach is better if the cluster is of arbitrary shape.



Key Features

To understand DBSCAN Algo we need to understand.

ϵ Maximum radius of neighbourhood.

MinPts : Minimum No of points in an ϵ neighbourhood of that point.



ϵ neighbourhood of $q = \{4\}$
 ϵ ————— — $P = \{3\}$

Density of q is high.
Density of p is low.

Terminologies \Rightarrow Given ϵ and Minpts then categorize the objects into 3 group

① Core Point \Rightarrow if for a point the neighbourhood has the points \geq Minpoint. Then it is Core point.

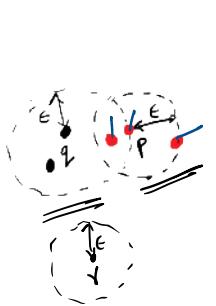
q is core point



\sqcup Minpts = 4

No of points in neighbourhood of $q = 4 \rightarrow$
So q is core point

② Border Point \rightarrow If for a point the neighbourhood has points $< \text{Minpt}$ but it should be neighbours of core point



$$\text{Minpt} = 4$$

No. of points in neighbourhood of $P = 3$

$$\text{Minpt} (k)$$

& P is neighbour of Q (core point)

So ' P ' is border point

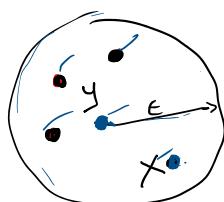
③ Outlier \rightarrow If a point is neither a core point nor a border point, it is called outlier.

"It is outlier"

* Directly Density Reachable \rightarrow A point X is directly density reachable from point Y w.r.t. epsilon & Minpt if

1. X belongs to neighbourhood of Y ($\text{dist}(X, Y) \leq \epsilon$)

2. Y is core point -



$$\text{Minpt} = 4$$

① as Y a core point \Rightarrow

no. of points in neighbourhood of $Y = 5 > \text{Minpt}$

$\therefore Y$ is core point

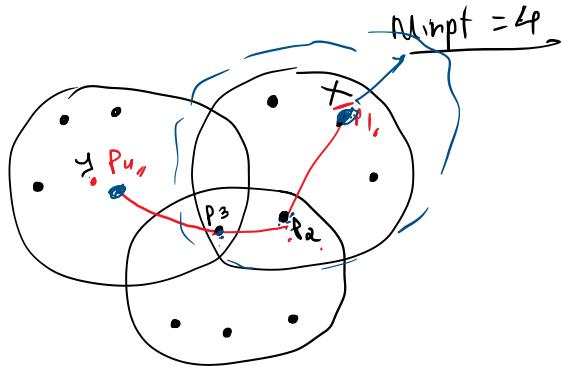
② X belongs to neighbourhood of Y .

$\therefore X$ is Directly Density Reachable from Y .

② Density Reachable \rightarrow

A point X is density reachable from Point Y w.r.t ϵ & Minpt if there is a chain of points P_1, P_2, \dots, P_n and $P_1 = X$

$\& P_n = Y$ such that P_{i+1} is directly density reachable from P_i



here $X \wedge Y$ are core points

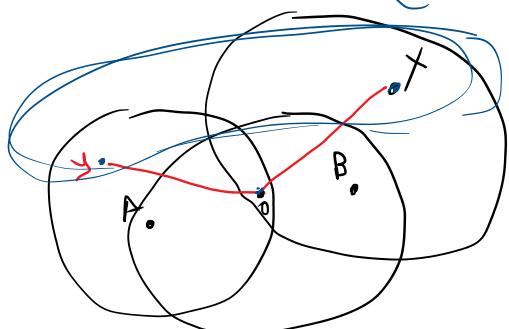
Here P_2 is directly density reachable from X

P_3 is directly density reachable from P_2

Y is directly density reachable from P_3

$\therefore X$ is density reachable from Y

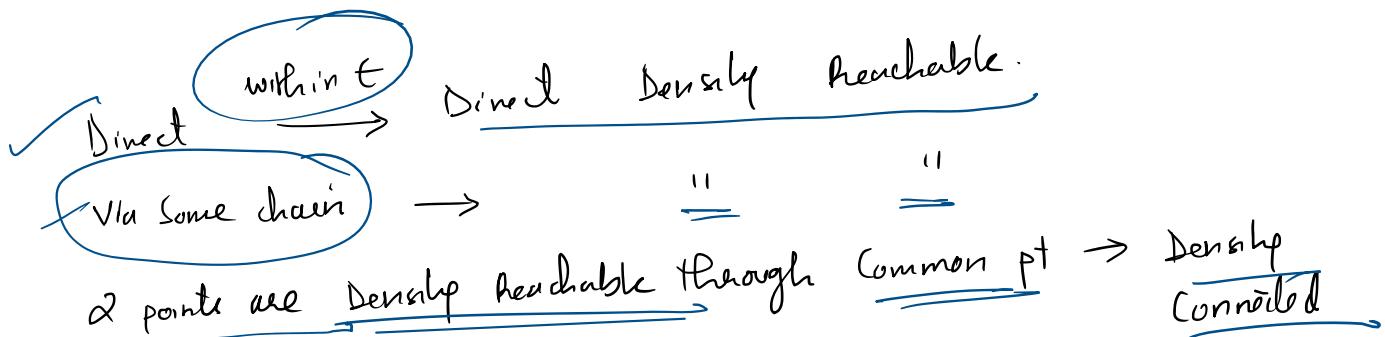
Density Connected \rightarrow A point X is density connected from point Y w.r.t ϵ and $Minpt$ if there exists a point O such that both $X \wedge Y$ are density reachable from O w.r.t ϵ & $Minpt$



$Y \wedge O$ are direct density reachable

$X \wedge O$ are direct density reachable

then $X \wedge Y$ are density connected



Algo

1. Arbitrary Select a Point P

2. Retain all points density reachable from P wst ϵ & Min Point

3. If P is core point then cluster is formed

4. If P is border point then DBSCAN visits next point of dataset

5. Continue the process until all the points are processed.

Q1. a) Elaborate Bagging Ensemble Learning.