

x_0	x_1	x_2	y	Sampling w.r.t. distribution (ii)		
1	1	2	3	\rightarrow Sample 1 i j	2	2
1	2	1	4		3	3
1	3	3	5		4	4

① compute the predictions.

for each sample

$$h_0(x^{(i)}) = \theta_0 x_0^{(i)} + \theta_1 x_1^{(i)} + \theta_2 x_2^{(i)}$$

$$h(x^{(1)}) = \theta_0 x_0^{(1)} + \theta_1 x_1^{(1)} + \theta_2 x_2^{(1)}$$

$$= 0*1 + 0*1 + 0*2 = 0$$

$$h(x^{(2)}) = \underline{\underline{0}} = 0$$

$$h(x^{(3)}) = \underline{\underline{0}} = 0$$

② compute gradient.

(i) compute errors or residuals.

$$e^{(1)} = (h_0(x^{(1)}) - y^{(1)}) = 0$$

$$e^{(2)} = 0 - 3 = -3$$

for each sample

$$e^{(1)} = -3$$

$$e^{(2)} = 0 - 4 = -4$$

$$e^{(3)} = 0 - 5 = -5$$

(ii) Calculate the gradients:-

$$\frac{\partial J}{\partial \theta_0} = \sum_{i=1}^n (h_\theta(x^{(i)}) - y^{(i)}) x_0^{(i)}$$

$$= -3*1 + (-4)*1 + (-5)*1$$

= -12

for each parameter

$$\frac{\partial J}{\partial \theta_1} = -3*1 + (-4)*2 + (-5)*3$$

= -26

$$\frac{\partial J}{\partial \theta_2} = -3*2 + (-4)*1 + (-5)*3$$

$$= -6 - 4 - 15$$

③ update parameters :- ($\alpha \rightarrow 0.1$)

$$\theta_0 = \theta_0 - \alpha \frac{\partial J}{\partial \theta_0}$$

$$\theta_0 = 0 - (0.1)(-12)$$

$$\theta_1 = \underline{\underline{0}} - (0.1)(-26)$$

$$\theta_2 = \underline{\underline{0}} - (0.1)(-25)$$

$$= \underline{\underline{2.5}}$$

$$\begin{bmatrix} 1.2 \\ 2.6 \\ 2.5 \end{bmatrix}$$

$$J = \frac{1}{2} \sum_{i=1}^n (h_\theta(x) - y)^2$$

$$= \frac{1}{2} ((-3)^2 + (-4)^2 + (-5)^2)$$

$$= \frac{1}{2} (9 + 16 + 25)$$

$$J = \underline{\underline{25}}$$

→ 9th iteration :-
→ calculating the prediction.

$$\begin{aligned} h(x^1) &= 1.2 * 1 + 2.6 * 1 + 2.5 * 2 \\ &= 1.2 + 2.6 + 5 \\ &= 8.8 \end{aligned}$$

$$\begin{aligned} h(x^2) &= 1.2 * 1 + 2.6 * 2 + 2.5 * 1 \\ &= 1.2 + 5.2 + 2.5 \\ &= \underline{\underline{8.9}} \end{aligned}$$

$$\begin{aligned} h(x^3) &= 1.2 * 1 + 2.6 * 3 + 2.5 * 3 \\ &= 1.2 + 7.8 + 7.5 \\ &= \underline{\underline{17.5}} \end{aligned}$$

→ compute gradient.

i) error:-

$$\begin{aligned} e^{(1)} &= 8.8 - 3 \\ &= \underline{\underline{5.8}} \end{aligned}$$

$$\begin{aligned} e^{(2)} &= 8.9 - 4 \\ &= \underline{\underline{4.9}} \\ e^{(3)} &= 17.5 - 5 \\ &= \underline{\underline{12.5}} \end{aligned}$$

→ Calculate the gradient:-

$$\begin{aligned} \frac{\partial J}{\partial \theta_0} &= \sum_{i=1}^n (h_\theta(x^{(i)}) - y^{(i)}) x_0^{(i)} \\ &= 5.8 * 1 + 4.9 * 1 + 12.5 * 1 \\ &= 5.8 + 4.9 + 12.5 \\ &= \underline{\underline{23.2}} \end{aligned}$$

In soft learning

- One hot encoder } for X
- Ordinal encoder } for X
- Label encoder — for Y

x_1	x_2	y
red	38.8	Y
blue	24.4	N

$$\hat{y} = \theta_0 x_0 + \theta_1 x_1 + \theta_2 x_2$$

$x_{1,0}$	$x_{1,1}$	x_2	y
0	1	38.8	Y
1	0	24.4	N

$$\hat{y} = \theta_0 x_0 + \theta_1 (x_{1,0}) + \theta_2 (x_{1,1})$$

② Data Standardization & Normalization

x_1, x_2, x_3 fact, known categorical variable

\rightarrow Normalized

$$-\infty \leq x \leq +\infty$$

2u/1

K-fold cross validation:

$\frac{1}{k}$	Sample	x_1	x_2	y
1	1	2	-1	1
2	2	5	1.2	0
3	3	1	2	1
4	4	-3	-2	1
5	5	4	0.1	0

3-fold cv: Acc = ? Std = ?

$$1 \quad [\theta_1, \theta_2] = [1.8, 2.8] \quad \textcircled{1}$$

$$2 \quad [\theta_1, \theta_2] = [2.1, 3.1]$$

$$3 \quad [\theta_1, \theta_2] = [1.9, 4.3]$$

\Rightarrow 1st fold:

Train :- I II $\rightarrow \{1, 2\} \{3, 4\}$

Test :- III $\rightarrow \{5\}$

1.8 x 5
 Acc. 1) ~~0.8 + 2.5~~ Std 1
 2) 2
 3) 3

2nd

$\begin{cases} 1.8 \\ 2.8 \end{cases} \rightarrow$ Train :- I III $\rightarrow \{1, 2\} \{5\}$

Test :- II $\rightarrow \{3, 4\}$

3rd

Train :- II III $\rightarrow \{3, 4\} \{8\}$

Test :- I $\rightarrow \{1, 2\}$

Acc 1)
 2)
 3)

Std 1)
 2)
 3)

fold 3 $[2 -1]$ $[5, 1.2]$

$$[4.802.8] \begin{bmatrix} 2 \\ -1 \end{bmatrix} :$$

$$[1.94] \begin{bmatrix} 2 \\ -1 \end{bmatrix}.$$

$$\Rightarrow 0.45 \rightarrow 0$$

$$5^{\text{old}} \Rightarrow [1.94] : \begin{bmatrix} 5 \\ 1.2 \end{bmatrix}$$

$$\Rightarrow 0.99 \rightarrow 1$$

\Rightarrow fold 2:-
 $[1 -3]$

$$[2.1, 3.1] \begin{bmatrix} 2 \\ -1 \end{bmatrix}$$

$$5^{\text{old}} \Rightarrow 4.2 + (-3.1)$$

$$k = n \rightarrow 100 \text{ CV}$$

Same one out cross validation.
 done when dataset < 100

10 \rightarrow Sample

10 \rightarrow fold.

at each point 1 sample \rightarrow test

10 \rightarrow EP. 9 \rightarrow training

features:

x_1 : date

x_2 : age

x_3 : height

x_4 : weight

x_5 : Sinus tachycardia

x_6 : min Systolic bp 24(hr)

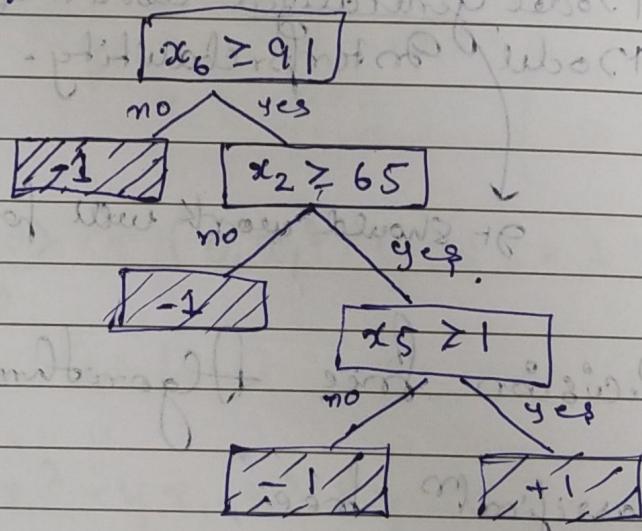
x_7 : latest diastolic bp.

Labels:

- high risk = +1

- low risk = -1

→ Graph:



→ testing / inference / Predict phase

Given $x^{(i)} = 20\ 20/11/17, 49, 172\text{cm}, 70.5\text{Kg}, 0, 115, 7.9$

→ low risk

I.Q.	Social	verbal	Risk for Autism
100	79	86	Low
30	40	70	High
80	85	68	Coco
90	82	98	High.

features:

x_1 : I.Q

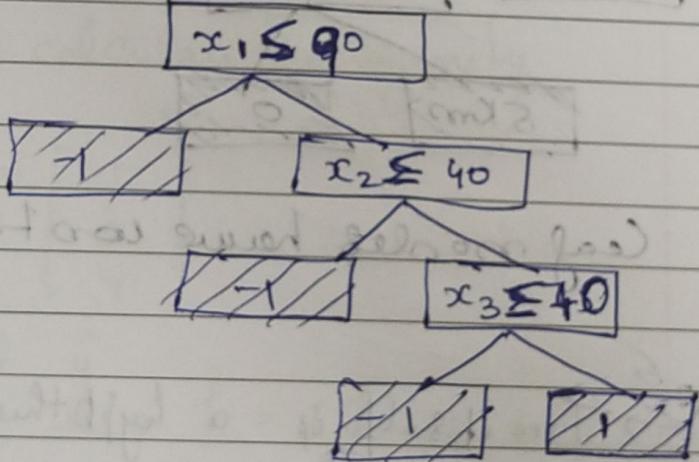
x_2 : Social

x_3 : verbal.

Labels

- high risk $\rightarrow +1$
- low risk $\rightarrow -1$

Graph.



- 1) feature
2) threshold.

Parameters - for every interval mod

feature threshold }

→ Regression setting :-

x_1 temperature - $^{\circ}\text{C}$.

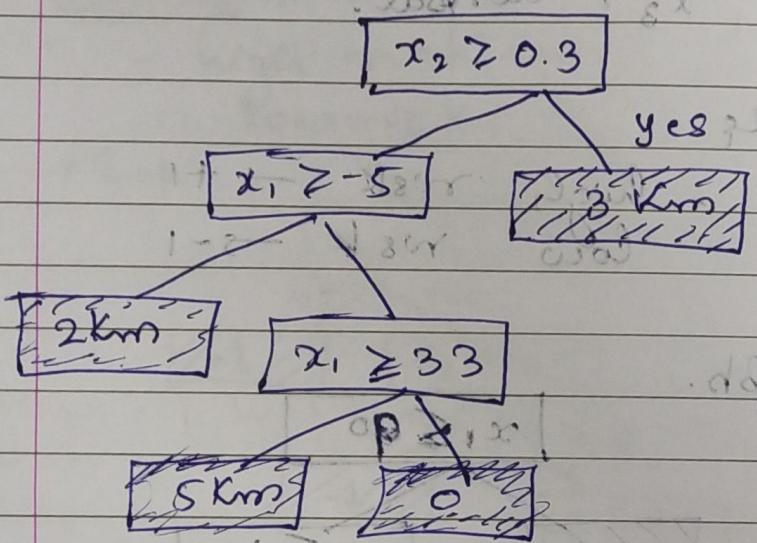
x_2 precipitation (cm/hr).

y : Walking - 1 Km

2 Km

3 Km

Regression tree :-



Leaf nodes have continuous values

0, 2, 3, . . .

→ Tree itself is a hypothesis function.

1) Split dim (feature)

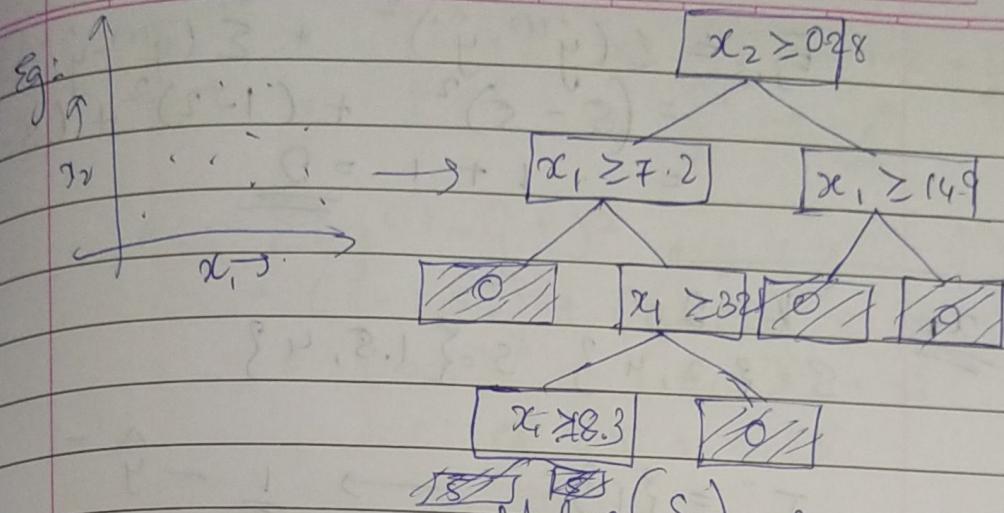
2) Split value

Learning Partition of Space

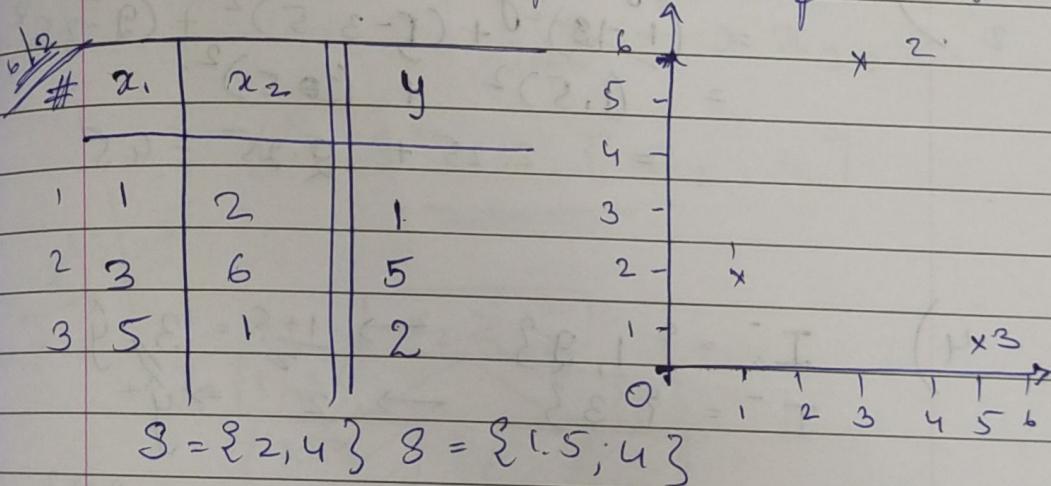
Decision boundary

feature
 x_1^*, x_2^* value

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



→ To decide split size → people will usually take the middle of each node or the data point itself.



$$\frac{x_1}{2} \quad K=2.$$

$$x \rightarrow I^+ = \{1\} \quad \boxed{\frac{1+0}{2}} \rightarrow \hat{y}^-$$

$$x \rightarrow I^+ = \{3, 5\} \quad \boxed{\frac{3+5}{2}} \rightarrow \hat{y}^+$$

$$E_{J,S} = \sum (y^{(i)} - \hat{y}^+)^2 + \sum (y^{(i)} - \hat{y}^-)^2$$

$$= (1-1)^2 + (5-\underline{4})^2 + (2-4)^2$$

$$= 1 + 4 = 5 //$$

$$4) \quad I^- = \{1, 3\} \quad I^+ = \{5\}$$

$$\hat{y}^- = \underline{2} \quad \hat{y}^+ = 5$$

$$\begin{aligned}
 E_{j,s} &= \sum (y^{(i)} - \hat{y}^+)^2 + \sum (y^{(i)} - \hat{y}^-)^2 \\
 &= (5-5)^2 + (1-2)^2 + (3-2)^2 \\
 &= \underline{\underline{1+1=2}}
 \end{aligned}$$

x_2 $k = 2$.

$$S = \{2, 4\} \quad S = \{1.5, 4\}$$

$$\begin{aligned}
 I^- &= \{1\} \rightarrow \underline{\underline{1}} - \hat{y}^- \\
 I^+ &= \{2, 3\} \rightarrow \underline{\underline{2.5}} \cdot \frac{5+2}{2} = 3.5
 \end{aligned}$$

$$\begin{aligned}
 1) \quad E &= \sum (y^{(i)} - \hat{y}^+)^2 + \sum (y^{(i)} - \hat{y}^-)^2 \\
 &= (1-1.5)^2 + (5-3.5)^2 + (2-3.5)^2 \\
 &= (1.5)^2 + (10.5)^2 \\
 &= \underline{\underline{2.25+2.25=4.5}}
 \end{aligned}$$

$$\begin{aligned}
 4) \quad I^- &= \{1, 2\} \rightarrow \frac{1+2}{2} = 1.5 \hat{y}^- \\
 I^- &= \{3\} \rightarrow \underline{\underline{3}} \hat{y}^+
 \end{aligned}$$

$$\begin{aligned}
 E &= \sum (y^{(i)} - \hat{y}^+)^2 + \sum (y^{(i)} - \hat{y}^-)^2 \\
 &= (1-3)^2 + (5-3)^2 + (2-2)^2 \\
 &= 4+4 \\
 &= \underline{\underline{8}}
 \end{aligned}$$

x_2 $S = \{1.5, 4\}$

$$I^- = \{3\} \rightarrow \underline{\underline{3}} \hat{y}^- \Rightarrow 2 \hat{y}^-$$

$$I^+ = \{1, 2\} \rightarrow \underline{\underline{1.5}} \hat{y}^+ \Rightarrow \frac{1+2}{2} = 1.5 \hat{y}^+$$

$$\begin{aligned}
 1.5) \quad E &= \sum (y^{(i)} - \hat{y}^+)^2 + \sum (y^{(i)} - \hat{y}^-)^2 \\
 &= (2-2)^2 + (1-3)^2 + (5-3)^2 \\
 &= \underline{\underline{4+4}} \\
 &= \underline{\underline{8}}
 \end{aligned}$$

$$q) I^- = \{1, 3\} \Rightarrow \hat{y}^- \rightarrow \frac{1+2}{2} = 1.5$$

$$I^+ = \{2\} \Rightarrow \hat{y}^+ \rightarrow \underline{\underline{\frac{2}{5}}}$$

$$E = \sum (y^{(i)} - \hat{y}^+)^2 + \sum (y^{(i)} - \hat{y}^-)^2$$

$$= (1 - 1.5)^2 + (2 - 1.5)^2 + (5 - 1.5)^2$$

$$= (0.5)^2 + (0.5)^2 + 0$$

$$= 0.25 + 0.25$$

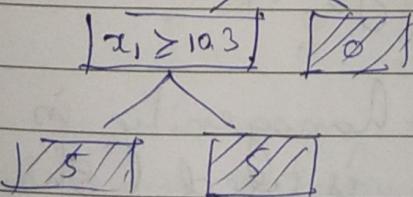
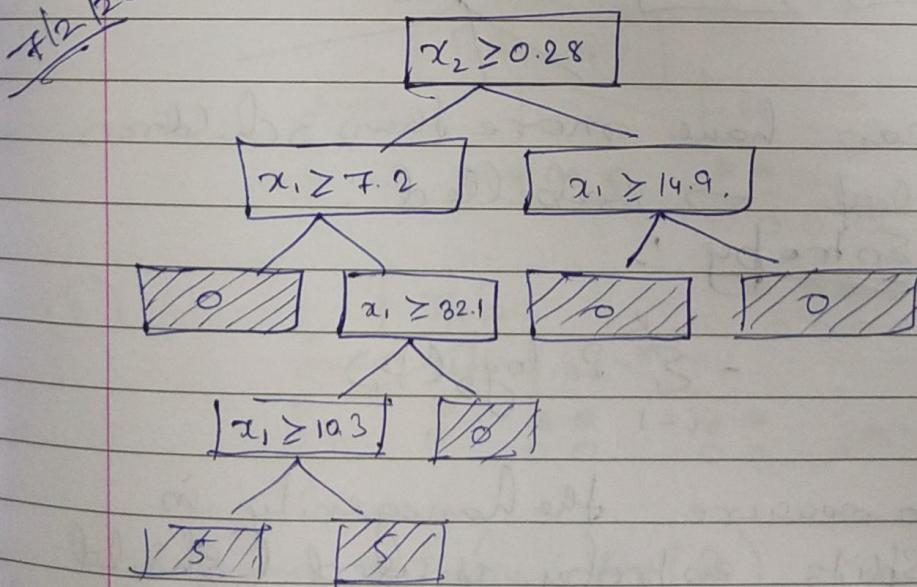
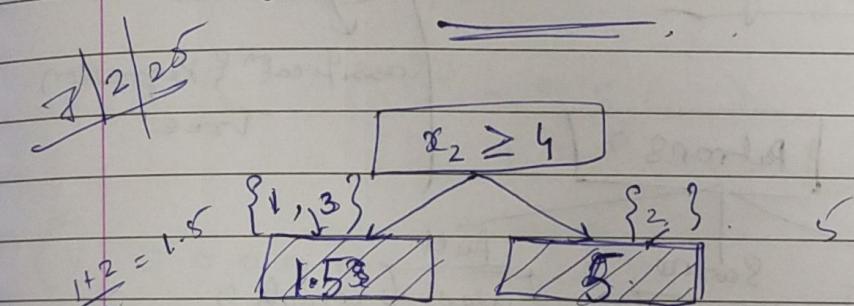
$$\underline{\underline{= 0.5}}$$

$$j^*, s^*$$

$$\min = \{4.5, 8, 8, 0.5\}$$

$$= \underline{\underline{0.5}} \rightarrow x_2, 4.$$

$$j^*, s^* = [x_2, 4]$$



$IG = \text{Information Gain}$

$\frac{\text{no. of true samples}}{\text{no. of false samples}} \rightarrow \text{negligible samples}$
 $P = 6, N = 6.$

bases algos

$$H\left(\frac{P}{P+N} \rightarrow \frac{n}{n+P}\right) = -\frac{P \log_2 \frac{P}{P+N}}{P+N}$$

$$\Rightarrow -\frac{P}{P+N} \log_2 \frac{P}{P+N} - \frac{N}{P+N} \log_2 \frac{N}{P+N}$$

$$= -\sum_{i=1}^k -P_i \log_2 P_i$$

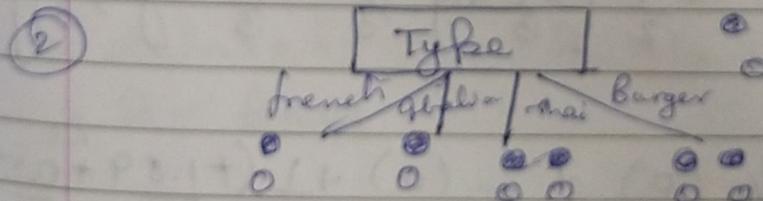
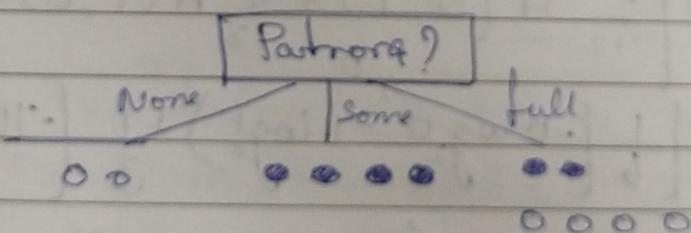
$$H\left(\frac{6}{12}\right) = \frac{6}{12} \log_2 \frac{6}{12} - \frac{6}{12} \log_2 \frac{6}{12}$$

$$= -\frac{1}{2} \log_2 0.5 - 0.5 \log_2 0.5$$

$$= -\frac{1}{2}(-1) - 0.5(-1)$$

$$= \frac{1}{2} + 0.5 = \underline{\underline{1}}$$

① 0 0 0 0 0 0 → red.
 0 0 0 0 0 0 → green



Expected entropy: EH(A)

$$EH(A) = \sum_{i=1}^K - \left(\frac{P_i + n_i}{P+n} \right) H \left(\frac{P_i}{P_i + n_i}, \frac{n_i}{P_i + n_i} \right)$$

maximize.

~~C~~ I.G. :

$$IG(A) = H \left(\frac{P}{n+P}, \frac{n}{P+n} \right) - EH(A).$$

" / 2 "

$$H(\text{Patrons}) = H \left(\frac{P}{P+n}, \frac{n}{P+n} \right)$$

$$H \left(\frac{6}{12}, \frac{6}{12} \right) = \sum_{i=1}^n P_i \log_2 P_i - \frac{1}{P+n} \log_2 \frac{1}{P+n}$$

$$EH(\text{Patrons}) = \sum_{i=1}^3 - \left(\frac{P_i + n_i}{P+n} \right) H \left(\frac{P_i}{P_i + n_i}, \frac{n_i}{P_i + n_i} \right)$$

$$\Rightarrow + \left[\frac{P_1 + n_1}{12} H \left(\frac{0}{2}, \frac{2}{2} \right) + \frac{P_2 + n_2}{12} H \left(\frac{4}{4}, \frac{0}{4} \right) + \right]$$

$$\Rightarrow + \left[\frac{2}{12} H(0, 1) + \frac{4}{12} H(1, 0) + \frac{6}{12} H \left(\frac{2}{3}, \frac{1}{3} \right) \right]$$

$$\Rightarrow + \left[\frac{1}{6} \left[-0 \log_2 0 - 1 \log_2 1 \right] + \frac{1}{3} \left[-1 \log_2 1 - 0 \log_2 0 \right] \right]$$

$$+ \frac{1}{2} \left(-\frac{1}{3} \log_2 \frac{1}{3} - \frac{2}{3} \log_2 \frac{2}{3} \right)$$

$$\Rightarrow + \left[\frac{1}{6} (0) + \frac{1}{3} (0) + \frac{1}{2} \left(+1.89 + 0.599 \right) \right]$$

$$= + \left(\frac{1}{2} (0.53 + 0.399) \right)$$

$$= \underline{\underline{0.464}}$$

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

$$TG(\text{Patrons}) = \frac{H(\text{Patrons})}{EH(\text{Patrons})}$$

$$= 1 - 0.48$$

$$= \underline{\underline{0.54}}$$

→ Same calculation for $TG(\text{Type}) = 0$.

So Patron Split is more homogeneous.

Example 6:

x_1	y	r
1	3	3
2	6	6
3	8	8
4	11	11
5	15	15

a) $x_1 \geq 3$.

$f(x_1, x_2)$ $\boxed{x_1 \geq 3}$

$\boxed{u \cdot 5} \quad \boxed{\frac{11+15+8}{2}}$

$\frac{3+6+8}{3} \quad \frac{11+15+8}{2}$

x_1	y	\hat{y}	r_{new}	$\hat{y} = r + \lambda(\text{Value})$
1	3	3.45 - 0.45	3 + 0.1 * 4.5 = 3.45	
2	6	6.45 - 0.45	6 + 0.1 * 4.5 = 6.45	
3	8	9.13 - 1.13	8 + 0.1 * 11.3 = 9.13	
4	11	12.13 - 1.13	:	
5	15	16.13 - 1.13		

λ, r_{new}

\Rightarrow Difference

$x_1 = 5$

$x f_1(5) + \lambda f_2(5) + f_3(5) \dots$

$= 0.8$

$\therefore x f_b(5)$

Boosted Regression Tree:- (Ada Boost Regression)

Boosted Classification Tree (Ada Boost classifier)

1) Initialize $f(x) = 0$
 $r = y$

2) for iteration 1 to B

- a) fit a model, $f_b(x; \theta) \rightarrow r$
- b) update the final model by adding

Stump function

$$F(x, \theta) = f(x, \theta) + \lambda$$

- c) update the residuals.

$$r^{(i)} = y^{(i)} - \lambda f_b(x; \theta)$$

3) final model

$$R(x, \theta) = \sum_{b=1}^B \lambda f_b(x, \theta)$$

→ Eg :-

	x_1	x_2	y
1	2	4	
2	3	5	
3	4	6	
4	5	7	

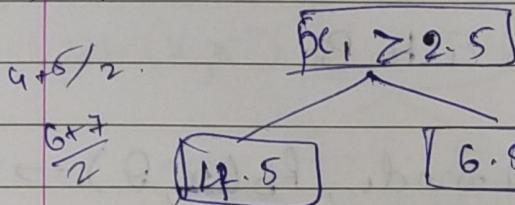
$$f(x) = 0$$

$$f_0(x) = 0$$

Iterations 1 :-

x_1	x_2	y	r	r
1	2	4	4	3.55
2	3	5	8	4.55
3	4	6	6	5.35
4	5	7	7	6.35

a) fit model, $f_1(x) \rightarrow r$.



b) update the model $\lambda = 0.1$

$$F(x) = f_0(x) + \lambda f_1(x)$$

c) update the residuals.

$$r_1 = r_1 - \lambda f_1(x) = 4 - 0.45 = 3.55$$

$$r_2 = r_2 - \lambda f_1(x) = 5 - 0.45 = 4.55$$

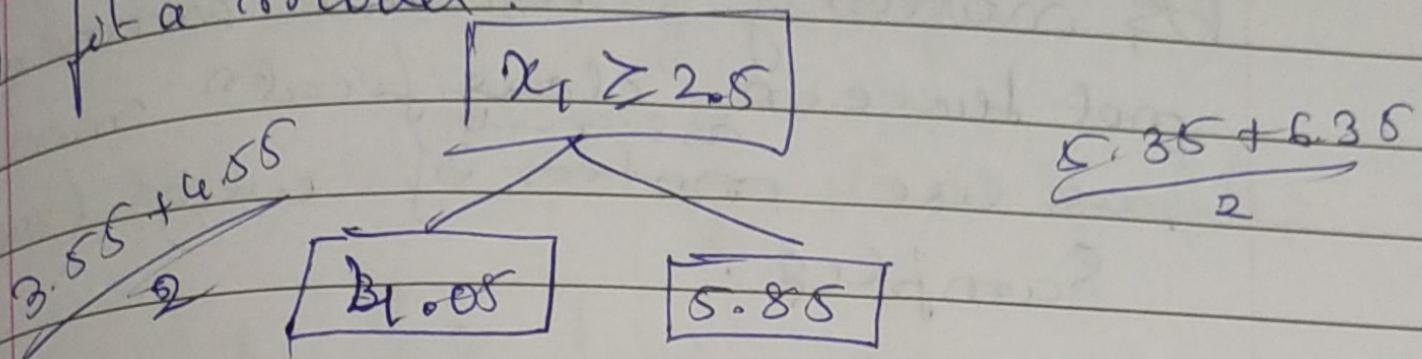
$$r_3 = r_3 - \lambda f_1(x) = 6 - 0.45 = 5.35$$

$$r_4 = r_4 - \lambda f_1(x) = 7 - 0.45 = 6.35$$

Iterations 2 :-

x_1	x_2	r
1	2	3.55
2	3	4.55
3	4	5.35
4	5	6.35

a) fit a model.



b) update the model

$$f(x) = f_0(x) + \lambda f_1(x) + \lambda f_2(x)$$

c) update the residuals

$$r_1 = y^{(1)} - \lambda f^{(1)}(x) = 3.55 - 0.405 = 3.145$$

$$r_2 = y^{(2)} - \lambda f^{(2)}(x) = 4.145$$

$$r_3 = y^{(3)} - \lambda f^{(3)}(x) = 4.765$$

$$r_4 = y^{(4)} - \lambda f^{(4)}(x) = 5.765$$

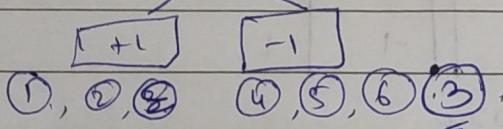
$$\lambda \propto \frac{1}{\omega}$$

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

YOUVA

Eg:-	α_1	α_2	y	ω
	1	2	+1	1
	2	3	+1	1
	3	3	+1	1
	4	8	-1	1
	5	5	-1	1
	6	6	-1	1

4) fit a model $x_i \in \{+1, -1\}$



$$\rightarrow E = \frac{1}{6} = 0.1667$$

$$\rightarrow \lambda = \frac{1}{2} \log_e \left(\frac{1 - 0.1667}{0.1667} \right)$$

$$= \frac{1}{2} \log_e (4.998).$$

$$= \frac{1}{2} (1.6092)$$

$$= \underline{0.808}$$

$$\omega^n = \omega^m e^{-\lambda b}$$

$$\omega^1 = 1 \cdot e^{-0.805}$$

$$= 0.447 = \omega^2 = \omega^4 = \omega^5 = \omega^6$$

$$\omega^n = \underline{\omega^m e^{-\lambda b}}$$

$$\omega^3 = 1 \cdot e^{\lambda b} = e^{0.805} = 2.236$$

M T W T F S S

Page No.:

Date:

YOUVA

x_1	x_2	y	w
1	2	+1	0.447
2	3	+1	0.447
3	3	+1	2.236
4	5	-1	0.447
5	5	-1	0.447
6	6	-1	0.447

→ Gradient Boosting for classification.

Sample	x_1	x_2	y
1	2	3	1
2	1	2	0
3	3	4	1
4	4	5	1
5	1	1	0
6	2	2	0

Step 1 :-

Initialization:-

$P_0(x) = \log \text{ odds of positive class probability}$

$$\left| P = \frac{3}{6} = 0.5 \right|$$

$$P_0 = \log \left(\frac{P}{1-P} \right) = \log \left(\frac{0.5}{0.5} \right)$$

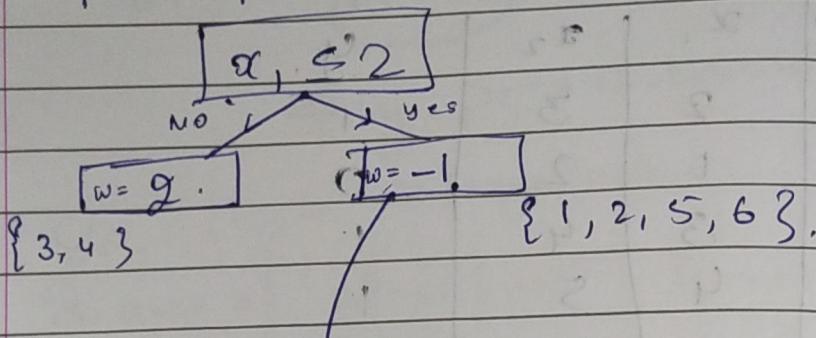
20,

⇒ Step 2 :- compute residuals.

$$r = y - P$$

Sample	Tree Y	Initial P=0.5	Residual $r = y - P$
1	1	0.5	0.5
2	0	0.5	-0.5
3	1	0.5	0.5
4	1	0.5	0.5
5	0	0.5	-0.5
6	0	0.5	-0.5

Step 3: fit a decision stump.



here for regression we fit avg of ys but, in classif. → we can't use that because $0.9 + 0.3 \rightarrow > 1$

(Probability is not linearly additive, hence we use concept of log odds):

hence,

$$w = \sum (y - p) \quad \text{regularization const}$$

~~and dir. $\lambda \sum (P \ln(p) + (1-p) \ln(1-p))$~~

= sum of residuals

$$\sum [P \ln(p) + (1-p) \ln(1-p)]$$

log odds adjustment

for yes (1, 2, 5, 6)

$$w = 0.5 + (-0.5) + (-0.5) + (-0.5)$$

$$(4)(0.5(-0.5))$$

$$0.25$$

using avg

$$= -1$$

$$= -1$$

for left NO

$$w = \frac{1}{2(0.5(0.5))} \leq \frac{1}{0.5} = \underline{\underline{2}}$$

Log odds w/baseline.

- for residuals also convert to P & Subtract.
- $F = F_0(x) + F_1(x) + F_2(x_{12})$.
- for P from log(odd)

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

This is for leaf nodes are having
log(odd) → convert into probability
& then give labels.

if $P \geq 0.5 \rightarrow \text{label} = 1$
else $\rightarrow 0$.

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

↓

$$\theta = \theta + \alpha \sum_{i=1}^n (y^{(i)} - \theta^T x^{(i)}) x^{(i)}$$

$\rightarrow x \in \mathbb{R}^r$

$$\theta = \theta + \alpha \sum_{i=1}^n (y^{(i)} - \theta^T \phi(x^{(i)})) \phi(x^{(i)})$$

monomial of
x with degree
3 -

$$\phi(x) = \begin{bmatrix} 1 \\ x_1 \\ x_2 \\ \vdots \\ x_1^2 \\ x_1 x_2 \\ x_1 x_3 \\ \vdots \\ x_1^3 \\ x_1^2 x_2 \\ x_1 x_2^2 \end{bmatrix}$$

$d = 1000$

$P \approx O(d^3)$

$1000^3 = 1,000,000,000$

$P = 1 \text{ billion.}$

for update theta
we should compute
1 billion computation
(1 million times lower
than normal)

Kernel function :-

$$k(x, y) = \langle \phi(x), \phi(y) \rangle$$

1st vector

dot product

transformed
 x, y .

Eg : $x = (x_1, x_2, x_3)$

$y = (y_1, y_2, y_3)$

attribute 2nd vector (not target value).

$$\phi(x) = \begin{bmatrix} x_1 x_1 \\ x_1 x_2 \\ x_1 x_3 \\ x_2 x_1 \\ x_2 x_2 \\ x_2 x_3 \\ x_3 x_1 \\ x_3 x_2 \\ x_3 x_3 \end{bmatrix}$$

$x = [1, 2, 3]$

$y = [4, 5, 6]$

→ compute $\rightarrow \langle \phi(x), \phi(y) \rangle = ?$

$$\phi(x) = \begin{bmatrix} 1 \\ 2 \\ 3 \\ 2 \\ 4 \\ 6 \\ 3 \\ 6 \\ 9 \end{bmatrix} \quad \phi(y) = \begin{bmatrix} 16 \\ 20 \\ 24 \\ 20 \\ 25 \\ 30 \\ 24 \\ 30 \\ 36 \end{bmatrix}$$

$$\langle \phi(x), \phi(y) \rangle =$$

$$16 + 40 + 72 + 40 + 100 + 180 + 72 + 180 + 324$$

$$= \underline{\underline{1024}}$$

But it is computationally expensive so, we can do $K(x, y) = (\langle x, y \rangle)^2$.

$$\begin{aligned} \text{↓ in feature space} &= [1, 2, 3] [4, 5, 6] \\ &= 4 + 10 + 18 \\ &= (32)^2 = \underline{\underline{1024}} \end{aligned}$$

⇒ Kernel methods.

1. Linear kernel, $K(x, y) = xy$
2. Polynomial.
3. Radial Basis Function
4. Laplacian Kernel
5. Graph Kernel.
6. Fisher Kernel.

Testing phase :-

$$h_0(x^{(t)}) = \theta^T \phi(x^{(t)}) \quad \begin{matrix} \text{train} \\ \text{sample} \end{matrix}$$

$$= \sum_{i=1}^n \beta_i \phi(x^{(i)}) \phi(x^{(t)}) \quad \begin{matrix} \text{inst} \\ \text{of test} \\ \text{sample} \end{matrix}$$

$$= \sum_{i=1}^n \beta_i K(x^{(i)}, x^{(t)}) \quad \begin{matrix} \text{of test} \\ \text{sample} \end{matrix}$$

dimension $\rightarrow 1 \times n$

Eg:-

	x_1	x_2	y	
1	1	2	0	$\phi(x) = \begin{bmatrix} x_1^2 \\ x_2^2 \\ x_1 x_2 \end{bmatrix}$
3	4	0		

$$x = [x_1, x_2]^T \leftarrow s_{2x_1} \quad \beta \rightarrow \beta_1, \beta_2$$

$x \in \mathbb{R}^2$.

$$\phi(x) \in \mathbb{R}^3 \quad \theta \in \mathbb{R}^3 \rightarrow [\theta_1, \theta_2, \theta_3]$$

$$\theta = \sum_{i=1}^3 \beta_i \phi(x^{(i)})$$

$$= \beta_1 \left[\begin{array}{l} \end{array} \right]$$

$$\phi(x_1) = \begin{bmatrix} 1 \\ 4 \\ 2 \end{bmatrix}$$

$$\phi(x_2) = \begin{bmatrix} 9 \\ 16 \\ 12 \end{bmatrix}$$

$$= \beta_1 \begin{bmatrix} 1 \\ 4 \\ 2 \end{bmatrix} + \beta_2 \begin{bmatrix} 9 \\ 16 \\ 12 \end{bmatrix}$$

$$= \begin{bmatrix} \beta_1 \\ 4\beta_1 \\ 2\beta_1 \end{bmatrix} + \begin{bmatrix} 9\beta_2 \\ 16\beta_2 \\ 12\beta_2 \end{bmatrix}$$

$$\begin{bmatrix} \theta_1 \\ \theta_2 \\ \theta_3 \end{bmatrix} = \begin{bmatrix} \beta_1 + 9\beta_2 \\ 4\beta_1 + 16\beta_2 \\ 2\beta_1 + 12\beta_2 \end{bmatrix}$$

* But in Kernel method we need training set even when we are in testing phase but in Linear or any other model we just had θ 's.

$$d = \frac{|\beta_0 + \beta_1 x_1 + \beta_2 x_2|}{\sqrt{\beta_1^2 + \beta_2^2}}$$

3 lines :- (Hyperplanes)

	Dataset		
	x_1	x_2	y
$2x_1 + 3x_2 - 5 = 0$	3	4	+1
$-x_1 + 4x_2 + 7 = 0$	2	3	+1
$5x_1 - 12x_2 + 10 = 0$	1	-1	-1
	-2	+1	-1

→ distance b/w each points.

i) hyperplane.

$$2x_1 + 3x_2 - 5 = 0$$

$$d = \frac{|2x_1 + 3x_2 - 5|}{\sqrt{2^2 + 3^2}} = \frac{|2x_1 + 3x_2 - 5|}{\sqrt{13}}$$

$$(3, 4) d_1 = \frac{|2(3) + 3(4) - 5|}{\sqrt{13}} = \frac{13}{\sqrt{13}} = \frac{\sqrt{13}}{3.6}$$

$$(2, 3) d_2 = \frac{|2(2) + 3(3) - 5|}{\sqrt{13}} = \frac{8}{\sqrt{13}} = 2.218$$

$$(1, -1) d_3 = \frac{|2(1) + 3(-1) - 5|}{\sqrt{13}} = \frac{6}{\sqrt{13}} = 1.664$$

$$(-2, 1) d_4 = \frac{|2(-2) + 3(1) - 5|}{\sqrt{13}} = \frac{6}{\sqrt{13}} = \frac{6}{\sqrt{13}} = 1.664$$

2) hyperplane $-x_1 + 4x_2 + 7 = 0$

$$d_1 = \frac{|-1x_1 + 4x_2 + 7|}{\sqrt{(-1)^2 + (4)^2}} = \frac{|-x_1 + 4x_2 + 7|}{\sqrt{17}}$$

$$(3, 4) d_1 = \frac{|-3 + 4(4) + 7|}{\sqrt{17}} = \frac{20}{\sqrt{17}} = 4.85$$

$$(2, 3) d_2 = \frac{|-2 + 4(3) + 7|}{\sqrt{17}} = \frac{17}{\sqrt{17}} = 4.123$$

$$(1, -1) d_3 = \frac{|-1 + 4(-1) + 7|}{\sqrt{17}} = \frac{2}{\sqrt{17}} = 0.485$$

$$(-2, 1) d_4 = \frac{|2 + 4(1) + 7|}{\sqrt{17}} = \frac{13}{\sqrt{17}} = 3.152$$

3) hyperplane

$$5x_1 - 12x_2 + 10 = 0$$

$$d = \frac{|5x_1 - 12x_2 + 10|}{\sqrt{5^2 + (-12)^2}} = \frac{|5x_1 - 12x_2 + 10|}{\sqrt{169}} = \frac{13}{13} = 13$$

$$(3, 4) d_1 = \frac{|5(3) - 12(4) + 10|}{13} = \frac{+23}{13} = 1.769$$

$$(2, 3) d_2 = \frac{|5(2) - 12(3) + 10|}{13} = \frac{-16}{13} = 1.230$$

$$(1, -1) d_3 = \frac{|5(1) - 12(-1) + 10|}{13} = \frac{2.076}{13}$$

$$(-2, 1) d_4 = \frac{|5(-2) - 12(1) + 10|}{13} = \frac{0.923}{13}$$

$$\Rightarrow \max(M_1, M_2, M_3) = \max \{ 1.664, 0.485, 0.923 \}$$

$$M_2 = \underline{0.485} \quad 1.664$$

Refresh linear algebra concepts:-

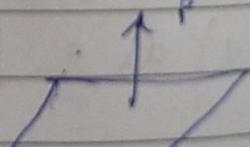
- ① projection
- ② Norm of a vector
- ③ vector addition
- ④ orthogonality
- ⑤ Subspaces

to prove that a point β_0 is always \perp to the plane

Take two pts. x_1, x_2

$$\beta \cdot x_1 = 0$$

$$\beta \cdot x_2 = 0$$



21/3/25

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

⇒ Eg.: Evaluation matrix. Threshold → 0.5

Observation	Actual Label	Predicted Score	Predicted Label
1	1	0.85	1
2	0	0.60	1 * (FP)
3	1	0.70	1
4	1	0.40	0 * (FN)
5	0	0.55	1 * (FP)
6	1	0.50	1
7	0	0.65	1 * (FP)
8	0	0.35	0
9	1	0.60	1
10	0	0.20	0

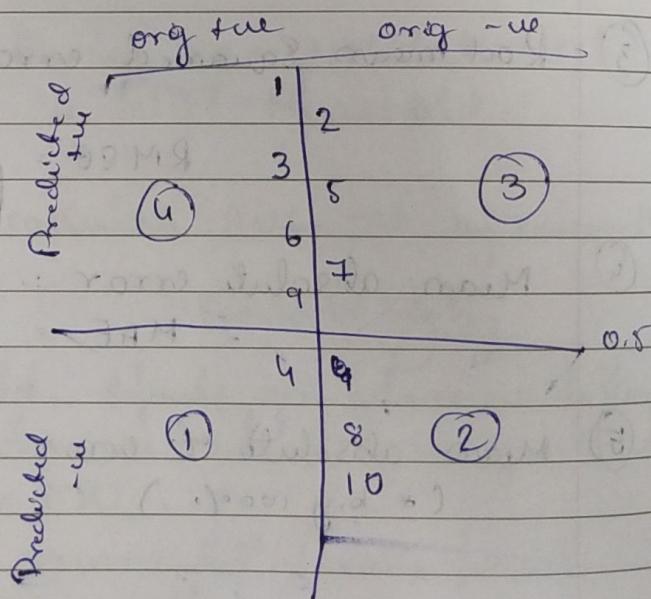
- 1) Accuracy ?
 - 2) Precision ?
 - 3) Recall or Sensitivity or True Positive Rate
 - 4) Specificity
 - 5) False +ve rate ($1 - \text{Specificity}$)
 - 6) F_1 - Score
 - 7) ROC Plot
 - 8) AUC
- ⇒ Accuracy

$$\cdot TP = 4$$

$$\cdot TN = 2$$

$$\cdot FP = 3$$

$$\cdot FN = 1$$



$$Accuracy = \frac{TP + TN}{Total} = \frac{4+5}{10} = \frac{9}{10} = 0.6$$

$$2) Precision := \frac{TP}{TP+FP} = \frac{4}{4+3} = \frac{4}{7} = 0.571$$

$$3) Sensitivity / Recall = \frac{TP}{TP+FN} = \frac{4}{4+1} = \frac{4}{5} = 0.8$$

$$4) Specificity = \frac{TN}{TN+FP} = \frac{2}{2+3} = \frac{2}{5} = 0.4$$

$$5) False +ve ratio = 1 - Specificity = 1 - 0.4 = 0.6$$

$$6) F_1 \text{ Score} = \frac{2 \times \frac{\text{Precision} \times \text{recall}}{\text{Precision} + \text{recall}}}{\frac{8 \times 0.571 \times 0.8}{0.571 + 0.8}} = \frac{2 \times 0.333.18}{1.333.18} = 0.666$$

7) ROC

Thresholds	TP	TN	FP	FN	TPR	FPR
0.7	2	5	0	3	0.4	0
0.5	4	2	3	1	0.8	0.6
0.2	5	2	3	0	1	0.6

$$\rightarrow 0.7 \rightarrow TP \rightarrow 2 \quad TPR \Rightarrow \frac{TP}{TP+FN} = \frac{2}{2+3} = \frac{2}{5} = 0.4$$

$$TN \rightarrow 5 \quad FPR \Rightarrow 1 - \frac{TN}{TN+FP}$$

$$FN \rightarrow 3$$

$$FP \rightarrow 0$$

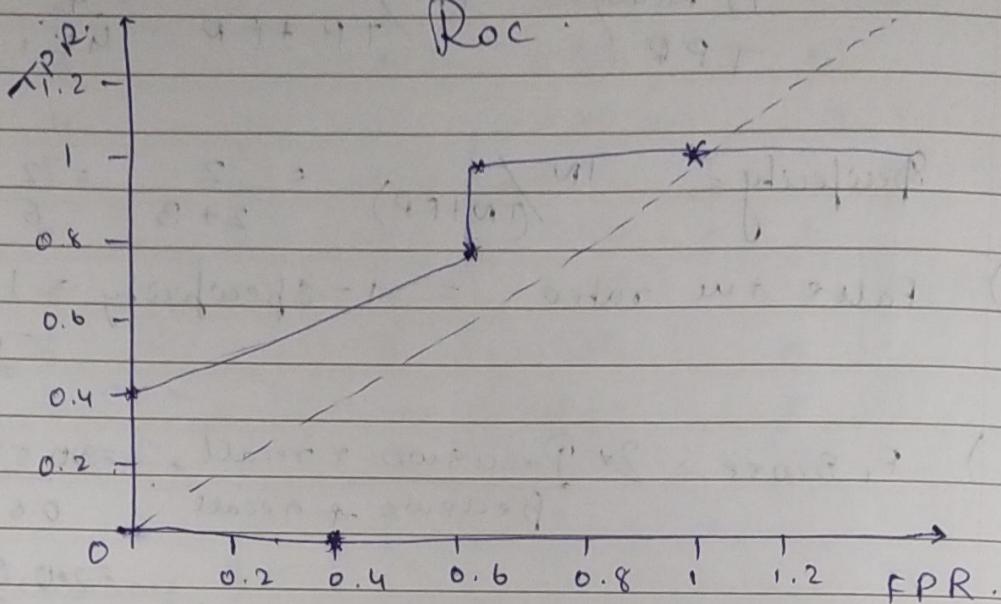
$$\rightarrow 0.4 \rightarrow TP \rightarrow 5 \quad TPR = \frac{5}{5} = 1 \quad FPR = 1 - \frac{2}{2+3} = 1 - \frac{2}{5} = 0.6$$

$$TN \rightarrow 2 \quad FN \rightarrow 0 \quad FP \rightarrow 3$$

$$0.2 \rightarrow TP = 0 \quad TPR = \frac{TP}{TP+FN} = \frac{0}{0+1} = 0$$

$$TN = 4 \quad FPR = 1 - TNR = 1 - \frac{TN}{TN+FP}$$

$$FP = 5 \quad FN = 1 \quad = 0.555 = 1 - TPR$$



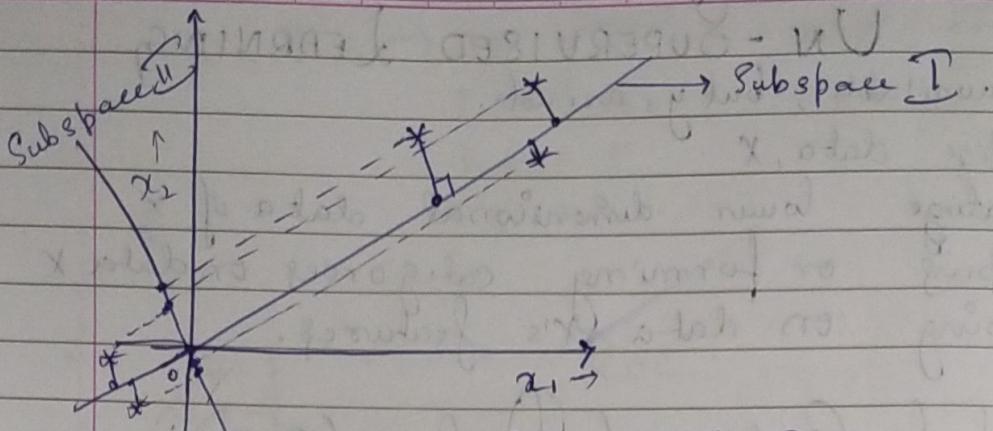
$$AUC = \sum_{i=1}^4 (FPR_{i+1} - FPR_i) \cdot (TPR_i + TPR_{i+1})$$

$$= (0 - 0.6) \cdot (0.4 - 0.8) + (0.6 - 1) \cdot (0.8 - 1)$$

$$= \frac{0.6 \cdot 0.4}{2} - 0.6 = 0.12 - 0.6 = -0.48$$

$$\underline{\underline{0.76}}$$

$$AUC = \sum_{i=1}^{n-1} (FPR_{i+1} - FPR_i) \cdot (TPR_{i+1} + TPR_i)$$



$$x = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \in \mathbb{R}^2 \Rightarrow z = [z_1] \in \mathbb{R}$$

We are going to linearly project each point on that subspace (linear line). Orthogonal projection

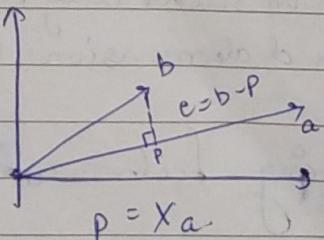
Eg:- 2 Subspace I & II.

A goal is to minimize variance in Subspace I \rightarrow variance is more b/w the projected points as distance is more but in Subspace II \rightarrow it's less.

So, Subspace I is more better than II.
 \because It captures more info.

Find flat Subspace among all possible spaces where my projected points have longer variance.

Linear alg. bsa.



We are trying to project b on a & to find $p = ?$

a \perp or to c.

$$a^T \cdot c = 0 \quad (\text{Dot product})$$

$$a^T \cdot (b - p) = 0$$

$$a^T \cdot (b - x_a) = 0$$

$$a^T b = x a^T a$$

$$x = \frac{a^T b}{a^T a} \rightarrow \text{Scalar}$$

$$p = x a.$$

$$P = \frac{a a^T b}{a^T a} \rightarrow \text{projection matrix } \frac{a a^T}{a^T a}$$

$$p = A (A^T A)^{-1} A^T b$$

$$\text{Eq: } b = \begin{bmatrix} 3 \\ 4 \end{bmatrix} \quad a = \begin{bmatrix} 1 \\ 2 \end{bmatrix} \quad p = ?$$

$$P = A (A^T A)^{-1} A^T b$$

$$= \begin{bmatrix} 1 \\ 2 \end{bmatrix} \cdot \left[\begin{bmatrix} 1 & 2 \end{bmatrix} \begin{bmatrix} 1 \\ 2 \end{bmatrix} \right]^{-1} \cdot b$$

$$x = \frac{a^T b}{a^T a} \rightarrow$$

$$= \begin{bmatrix} 1 \\ 2 \end{bmatrix} \cdot \begin{bmatrix} 1+4 \\ 1+4 \end{bmatrix}^{-1} \cdot b$$

$$= \begin{bmatrix} \frac{1}{2} \\ \frac{1}{2} \end{bmatrix} \cdot \begin{bmatrix} 5 \\ 5 \end{bmatrix}^{-1} \cdot b = \frac{8}{5} \begin{bmatrix} 3 \\ 4 \end{bmatrix}$$

$$= \begin{bmatrix} 3 \\ 4 \end{bmatrix}$$

$$= \underbrace{\begin{bmatrix} 1 & 2 \\ 2 & 4 \end{bmatrix}}_{2 \times 2} \underbrace{\begin{bmatrix} 3 \\ 4 \end{bmatrix}}_{2 \times 1} \downarrow$$

$$= \begin{bmatrix} 5 \\ 3+8 \\ 6+16 \end{bmatrix} = \begin{bmatrix} 5 \\ 11 \\ 22 \end{bmatrix} = \begin{bmatrix} 1 \\ 2.2 \\ 4.4 \end{bmatrix}$$

→ continue PCA ...

\vec{u} to be subspace of unit length
 \vec{x} → point to be projected.

p = projected pt of \vec{x} = proj(a) $\cdot \vec{x}$

Algorithm : K-means (K, dataset) {

- 1) Randomly assigns a number from 1 to K to each observation. These serve as initial clusters assignments.
- 2) Iterate 3) until the cluster assignments stop changing.
 - a) for each of the K-clusters, compute the cluster centroid.
 - b) Assign each observations to the cluster whose centroid is closest (based on Euclidean distance)

Sample.

	x_1	x_2	$(x_1 - \bar{x})^2 + (x_2 - \bar{x})^2$	Given:
1	1	1	$(1-5)^2 + (1-5)^2 = 16$	$C_1 = \{2, 3, 4, 6\}$
2	1	2	$(1-5)^2 + (2-5)^2 = 18$	$C_2 = \{1, 5\}$
3	2	2	$(2-5)^2 + (2-5)^2 = 18$	$C_3 = \{3, 4\}$
4	8	8	$(8-5)^2 + (8-5)^2 = 28$	$C_4 = \{5, 6\}$
5	8	9	$(8-5)^2 + (9-5)^2 = 29$	
6	9	8	$(9-5)^2 + (8-5)^2 = 34$	

Step 1:- $C_1 = \frac{1}{n} \begin{bmatrix} 5 \\ 5 \\ 5 \end{bmatrix} \rightarrow x_1 = \frac{20}{3} = 6.67$ & $x_2 = \frac{20}{3} = 6.67$

$C_2 = \frac{1}{n} \begin{bmatrix} 4.5 \\ 5 \end{bmatrix} \quad x_1 = \frac{9}{2} = 4.5$ & $x_2 = \frac{10}{2} = 5$

Step 2:- $d_{11} = ?$

$d_{12} = ?$

$C_1 \quad d_{11} = \sqrt{(5-1)^2 + (5-1)^2} = \sqrt{32} = 5.65$

$d_{12} = \sqrt{(4.5-1)^2 + (5-1)^2} = 5.31$

1 2

$$d_{21} = \sqrt{(5-1)^2 + (8-2)^2} \\ = \sqrt{16+9} = 5$$

$$\text{Or } d_{22} = \sqrt{(4.5-1)^2 + (4.5-2)^2} \\ = \sqrt{18} = 4.6$$

$$\rightarrow d_{31} = \sqrt{(5-2)^2 + (8-2)^2} \\ = \sqrt{18} = 4.24$$

$$C_2 \quad d_{32} = \sqrt{(4.5-2)^2 + (5-2)^2} \\ = \sqrt{18} = 4.24$$

$$\rightarrow d_{41} = \sqrt{(8-8)^2 + (5-8)^2} \\ = \sqrt{9+9} = 4.24$$

$$C_2 \quad d_{42} = \sqrt{(4.5-8)^2 + (8-8)^2} \\ = \sqrt{4.609}$$

$$\rightarrow d_{51} = \sqrt{(5-8)^2 + (5-9)^2} = 5$$

$$d_{52} = \sqrt{(4.5-8)^2 + (8-9)^2} = 5.31$$

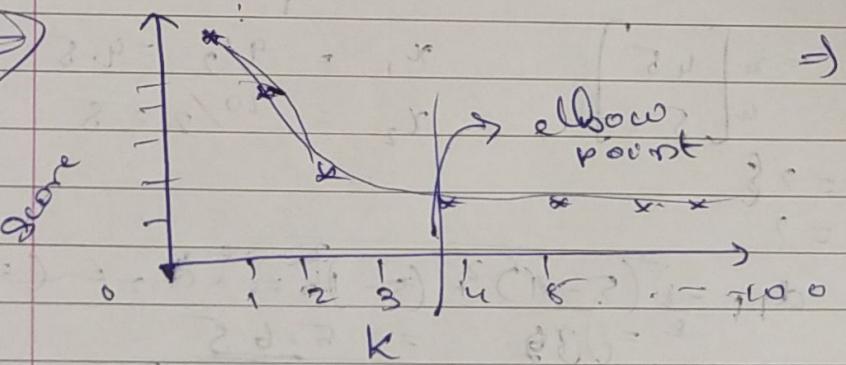
After 18th Iteration.

$$C_1 = \{4, 5, 6\}$$

$$C_2 = \{1, 2, 3\}$$

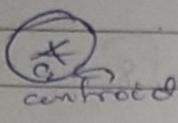
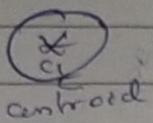
Cal centroid.

2/1/20
→



⇒ when to stop adding
of K. is not making
much difference
in C₂.

→ Centroid linkage.



Possimilarity b/w Dot product;

• Dot product \rightarrow similarity

• 1 - Dot product \rightarrow dissimilarity.

	x_1	x_2
A	1	1
B	2	1.1
C	4	3
D	5	4
E	6	5

Single linkage.
Euclidean distance.

Step 1: compute the distance matrix

	A	B	C	D	E
A	0	1.41	3.61	4.2	5.6
B	1.41	0	2.8	4.2	5.6
C	3.61	2.8	0	1.41	2.8
D	4.2	4.2	1.41	0	1.41
E	5.6	5.6	2.8	1.41	0

$$d_{AB} = \sqrt{(x_2^* - x_1^*)^2 + (y_2 - y_1)^2}$$

$$= \sqrt{1^2 + 0^2} = 1$$

$$d_{AC} = \sqrt{(3)^2 + (2)^2} = \sqrt{13} = 3.61$$

$$d_{AD} = \sqrt{4^2 + 3^2} = \sqrt{25} = 5$$

$$d_{AE} = \sqrt{5^2 + 4^2} = \sqrt{41} = 6.4$$

$$d_{BC} = \sqrt{4^2 + 4^2} = \sqrt{32} = 2.8$$

$$d_{BD} = \sqrt{9^2 + 9^2} = \sqrt{18} = 4.2$$

$$d_{CD} = \sqrt{1+1} = \sqrt{2} = \underline{\underline{1.41}}$$

$$d_{CE} = \sqrt{4+4} = \sqrt{8} = \underline{\underline{2.8}}$$

$$d_{DE} = \sqrt{1+1} = \underline{\underline{1.41}}$$

→ Minimum → AB.

merge (AB)	C	D	E	
(AB)	0	2.82	4.2	5.6
C	2.82	0	1.41	2.8
D	4.2	1.41	0	1.41
E	5.6	2.8	1.41	0

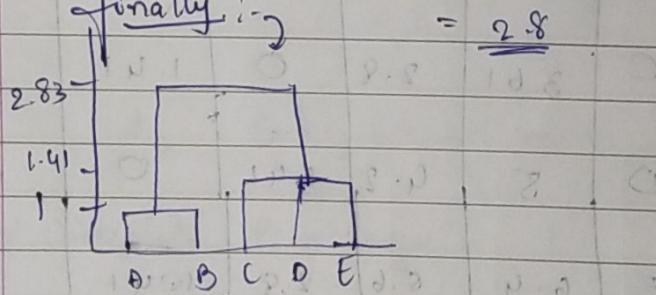
$$AB \rightarrow C = \min(d_{AC}, d_{BC})$$

$$= \min(5.6, 2.8)$$

finally :-

$$= \underline{\underline{2.8}}$$

repeat ↓



~~4/4/25~~

Maine Bayes Classifier :- (usually used in spam detection)

→ Input → Text modality. ↗ by looking @ context

Eg :- "To buy this product" \Rightarrow

Is by looking @ constant

Some algo to convert
so, that we can
do further
processes

* It should capture file

Semantic meaning of the sentence

Eg :- King ,  empire

~~rephrase~~ should be closer

dictionary

a		0
able		0
allow		0
buy	← take index →	1
move		0
Play		0
Product	→ ↑	0
this		1
to		1
	↓	0
		50,000

→ One way to preprocess.

* take index from dictionary.

- * Part 1 Corresponding index of the words of Sentences & rest all are zeros.

Eg:-	x_1	x_2	x_3	Class
1	1	0	1	Spam
2	1	1	0	Spam
3	1	1	0	Spam
4	0	0	1	No Spam
5	0	0	0	No Spam
6	1	1	1	No Spam
7	0	1	0	Spam
8	1	0	0	Spam
9	0	1	0	Spam
10	1	0	0	Spam
11	1	1	0	No Spam

x^3	\rightarrow	distancia
x_1	x_2	x_3
0	0	0
0	0	1
0	1	0
1	0	0
0	1	1
1	1	0
1	0	1
0	1	0
0	0	0

x_1	x_2	x_3	Prob Spams	Prob not Spams
0	0	0	0	$\frac{1}{4}$
0	0	1	0	$\frac{1}{4}$
0	1	0	$\frac{2}{7}$	0
0	1	1	$\frac{2}{7}$ 0	0
0	0	0	$\frac{2}{7}$ 0	0
1	0	1	$\frac{1}{7}$	0
1	1	0	$\frac{1}{7}$	$\frac{1}{4}$
1	1	1	0	$\frac{1}{4}$

$$2^3 - 1 = 7$$

Problems with this approach is if features \rightarrow 1 million possible combinations \rightarrow $2^{1\text{million}}$

It is a working model, but not practical.

Naive Bayes Assumption

x_i 's are independent given y

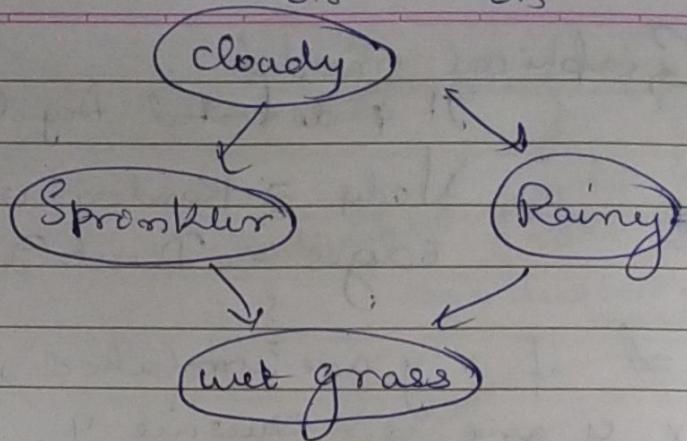
↳ conditional independence.

$$\text{Eg: } P(x_{2087}/y) = P(x_{2087}/y, x_{38193})$$

(for given y , x_{38193} is not influencing x_{2087})

Eg for given $y=1$ (Spam), "this" is not influencing y "buy"

$$P(C=F) = 0.5 \quad P(C=T) = 0.5$$



		$P(R C)$	
C		$P(R=F)$	$P(R=T)$
F	F	0.8	0.2
	T	0.2	0.8

		$P(S C)$	
C		$P(S=F)$	$P(S=T)$
F	F	0.5	0.5
	T	0.9	0.1

		$P(w S, R)$	
S	R	$P(w=F)$	$P(w=T)$
F	F	1.0	0.0
T	F	0.1	0.9
F	T	0.1	0.9
T	T	0.01	0.99

~~(2) (a) 25
Inference, & Update :-~~

$$\Rightarrow P(S=1) = \sum_{C=0} \sum_{R=0} \sum_{w=0} P(C, R, w, S=1)$$

$$= \sum_C \sum_R \sum_w P(C) P(R|C) P(S=1|R) P(w|R)$$

$$= (0.5)(0.2 + 0.8) =$$

$$= P(C=0) P(S=1|C=0) P(R=0|C=0) P(w=0|S=1, R=0)$$

$$+ P(C=0) P(S=0|C=0) P(R=0|C=0) P(w=0|S=0, R=0)$$

CRW
000

001

010

011

111

Continued 2 pages later

as $C \uparrow \rightarrow M \uparrow \rightarrow \downarrow$ Sensitivity.

$M \uparrow \rightarrow$ \downarrow bias \uparrow Variance

$M \downarrow \rightarrow$ \uparrow bias \downarrow Variance.

$$\therefore P(c=0) P(S=1/c=0) P(R=0/c=0) P(w=0/S=1, R=0)$$

+

$$P(c=0) P(S=1/c=0) P(R=0/c=0) P(w=1/S=1, R=0)$$

+

$$P(c=0) P(S=1/c=0) P(R=1/c=0) P(w=0/S=1, R=1)$$

+

$$P(c=1) P(S=1/c=1) P(R=0/c=1) P(w=0/S=1, R=0)$$

+

$$P(c=1) P(S=1/c=1) P(R=1/c=1) P(w=0/S=1, R=1)$$

+

$$P(c=1) P(S=1/c=1) P(R=0/c=1) P(w=1/S=1, R=0)$$

+

$$P(c=0) P(S=1/c=0) P(R=1/c=0) P(w=1/S=1, R=1)$$

+

$$P(c=1) P(S=1/c=1) P(R=1/c=1) P(w=1/S=1, R=1)$$

$$\Rightarrow (0.8 * 0.5 * 0.8 * 0.1) + (0.5 * 0.5 * 0.8 * 0.9) +$$

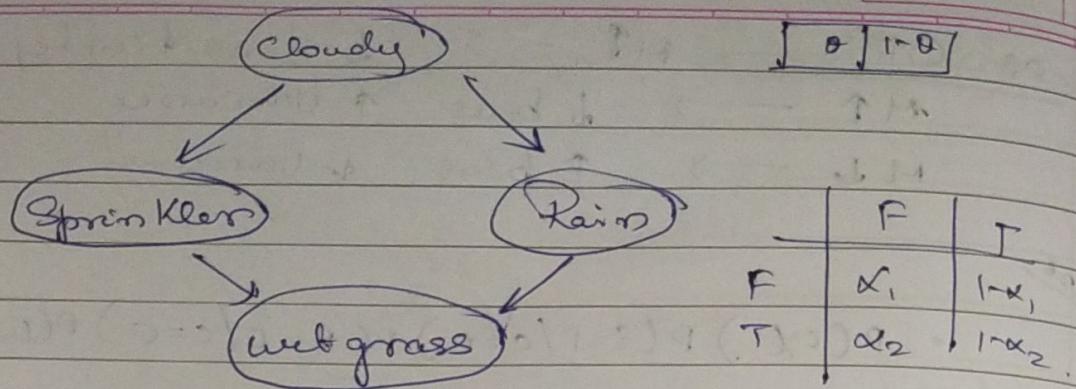
$$(0.5 * 0.8 * 0.2 * 0.01) + (1 - 1 - 1)$$

(2) 9(2, 2, 1, 7) a (1, 1, 1) a (1, 1, 1) a (1, 1, 1) a (1, 1, 1)

and 2, 2, 2 (cont'd part) below will be
plotted for more clarity and analysis.

Learning phase

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

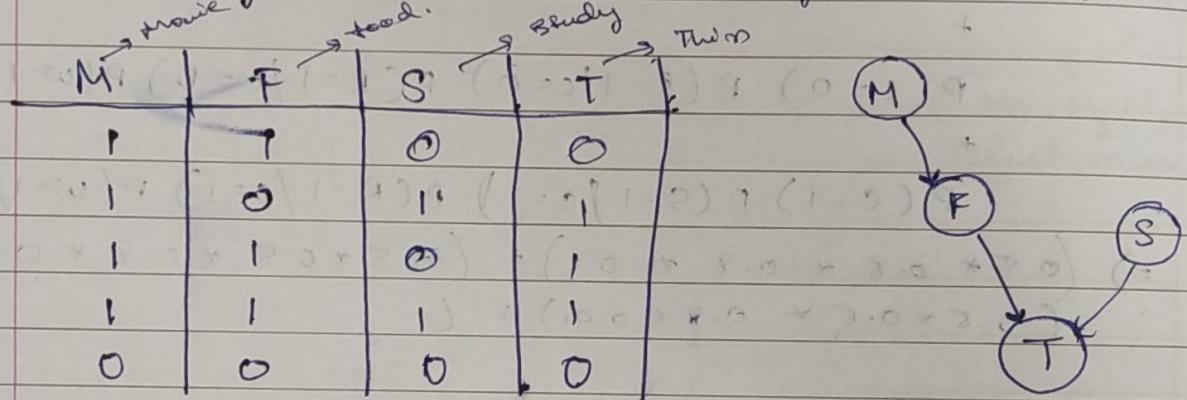


	F	T			F	T	
F	γ_1	$1-\gamma_1$			β_1	$1-\beta_1$	
T	$1-\gamma_2$	$1-\gamma_2$			β_2	$1-\beta_2$	
			F T		β_3	$1-\beta_3$	
				T T	β_4	$1-\beta_4$	

~~↓ July 25~~

→ t-SNE - non linear approaches like PCA
 \downarrow
 linear.

→ write pseudocode == PCA + kfold + kmeans.



$$P(M, F, S, T) = P(M) P(F|M) P(T|F, S) P(S)$$

build multiple model (graph/tree). Select the model which gave highest joint probability.

→ Ways to learn parameters (α, β, θ 's)

- 1) MLE appro
- 2) Bayesian Learning approach

i) MLE (Maximum Likelihood Estimation) - eqn.

$$P(M|\theta) = \theta^4(1-\theta)^1$$

$$P(S|\alpha) = \alpha^2(1-\alpha)^3$$

$$P(F|M=0) = \beta_1^0(1-\beta_1)^1$$

$$P(F|M=1) = \beta_2^3(1-\beta_2)^1$$

$$P(T|F=0, S=0) = \beta_1^0(1-\beta_1)^1$$

$$P(T|F=1, S=1) = \beta_2^3(1-\beta_2)^1$$

↓
write joint probability as equal to 0
differentiate
likelihood function
 \log maximize.

$$\frac{m}{n} \hat{\theta}_{MLE} = \frac{x_1}{5}$$

$$\hat{\alpha}_{MLE} = \frac{2}{5}$$

$$\hat{\beta}_{MLE} = \frac{4}{5}$$

Bayesian Learning approach :-

→ we use prior knowledge.

→ MLE → frequentist approach

Step 1:- Given $x_{1-n} = \{x_1; x_2; x_3; \dots; x_n\}$
write down the expression likelihood.

Step 2:- Specify a prior, $p(\theta)$

Step 3:- Compute posterior, $p(\theta|x_1, \dots, x_n) = \frac{p(x_1, \dots, x_n|\theta)}{p(x_1, \dots, x_n)}$

→ Beta distribution is used to choose the parameter

for prior knowledge we can use beta distribution to capture it

$$\text{eg: } \text{Beta}(1, 1) \text{ or } \text{Beta}(2, 2).$$

$$\text{Let } p(\alpha) = \text{Beta}(10, 1)$$

$$p(\theta) = \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} \theta^{\alpha-1} (1-\theta)^{\beta-1}$$

↓
parameters of
Beta

$$\text{Prior rule: } p(\alpha) = \text{Beta}(10, 1) = \alpha^{10-1} (1-\alpha)^{1-1}$$

$$= \underline{\alpha^9} (1-\alpha)^0$$

$$P(\alpha|S) \propto P(S|\alpha) p(\alpha)$$

$$\approx \alpha^2 (1-\alpha)^3 [\alpha^9 (1-\alpha)^0]$$

$$P(\alpha|S) \approx \underline{\alpha^{11} (1-\alpha)^3} \Rightarrow E[\alpha|S] = \frac{11}{11+3} = \frac{11}{14}$$

So, the α value obtained by MLE was 0.4, ≈ 0.7
 Bayesian $\rightarrow 0.7$ (for $\text{Beta}(10, 1)$)