

16/01/25

Prev problemExample

$$\textcircled{1} \quad h_{\theta}(x^{(1)}) = 0; \quad h_{\theta}(x^{(2)}) = 0, \quad h_{\theta}(x^{(3)}) = 0$$

Initialization:  $\theta_0 = 0, \theta_1 = 20, \theta_2 = 0; J = \frac{1}{2} \sum_{i=1}^n (h_{\theta}(x^{(i)}) - y^{(i)})^2$

$$J = \frac{1}{2} [(0-3)^2 + (0-4)^2 + (0-5)^2] = \frac{1}{2} [50] = 25$$

$x_1$	$x_2$	$y$
1	2	3
2	1	4
3	3	5

$$h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 \dots$$

Initialization:  $\theta_0 = 0, \dots = 0$ .

$\alpha = 0.1$  learning rate

### ① Compute the predictions

$$\begin{aligned}
 h_{\theta}(x^{(1)}) &= \theta_0 x_0^{(1)} + \theta_1 x_1^{(1)} + \theta_2 x_2^{(1)} \\
 &= 0*1 + 0*1 + 0*2 \\
 &= 0
 \end{aligned}$$

superscript  
 | initial  
 $x_1, x_2 \dots$   
 ↓  
 input feat

$$\begin{aligned}
 h_{\theta}(x^{(2)}) &= 0*1 + 0*2 + 0*1 \\
 &= 0
 \end{aligned}$$

$$\begin{aligned}
 h_{\theta}(x^{(3)}) &= 0*1 + 0*2 + 0*1 \\
 &= 0
 \end{aligned}$$

## Compute gradients

### i) compute errors / residuals

$$1^{\text{st}} \text{ sample} \quad e^{(1)} = h_\theta(x^{(1)}) - y^{(1)}$$

$$e^{(1)} = 0 - 3 = -3$$

$$\begin{aligned} 2^{\text{nd}} \text{ sample} \quad e^{(2)} &= 0 - 4 \\ &= -4 \end{aligned}$$

$$\begin{aligned} 3^{\text{rd}} \text{ sample} \quad e^{(3)} &= 0 - 5 \\ &= -5 \end{aligned}$$

### ii) Calculate gradients

$$\begin{aligned} \frac{\partial J}{\partial \theta_0} &= \sum_{i=1}^n h_\theta(x^{(i)}) - y^{(i)} \underbrace{x_0^{(i)}}_e \\ &= -3*1 + (-4)*1 + (-5)*1 \end{aligned}$$

$$\frac{\partial J}{\partial \theta_0} = -12$$

$$\begin{aligned} \frac{\partial J}{\partial \theta_1} &= \sum_{i=1}^n h_\theta(x^{(i)}) - y^{(i)} \underbrace{x_1^{(i)}}_e \\ &= -3*1 + -4*2 + -5*3 \\ &= -26 \end{aligned}$$

$$\begin{aligned} \frac{\partial J}{\partial \theta_2} &= \sum_{i=1}^n h_\theta(x^{(i)}) - y^{(i)} \underbrace{x_2^{(i)}}_e \\ &= -3*2 + -4*1 + -5*3 \Rightarrow -6 - 4 - 15 \\ &= -25 \end{aligned}$$

### iii) Update parameters

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

prev step  
(init to 0) = 0 - 0.1 + 1(2)  
= 0 + 1.2  
= 1.2

$$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}$$

at beginning

$J = 25$  after 1st iteration,

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

$$J = 94.95$$

$$\theta_1 = \theta_1 - \alpha \frac{\partial J}{\partial \theta_1}$$

$$= 0 - 0.1 * 2.6$$

$$= 2.6$$

$$\theta_2 = \theta_2 - \alpha \frac{\partial J}{\partial \theta_2}$$

$$= 0 - 0.1 * -2.5$$

$$= 0 + 2.5$$

$$= 2.5$$

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

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

$$J(1) = \frac{1}{2} (0 - 3)^2 \rightarrow (0-3)^2$$

$$= \frac{1}{2} \times 9$$

$$J(2) = (0-4)^2 = 16$$

$$J(3) = (0-5)^2 = 25$$

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

$$= \frac{1}{2} \times 50 = 25$$

$$J(\theta) = \cancel{25}$$

after 1<sup>st</sup> iteration

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

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

After 2<sup>nd</sup> iteration

$$\theta_0 = -1.02 ; \theta_1 = -2.41 ; \theta_2 = -2.6$$

$$J = 25678.56$$

$\alpha = 0.01$  change (as we take large steps)

$$\theta_0 = 0, \theta_1 = 0, \theta_2 = 0$$

$$J = 10.8$$



After 1<sup>st</sup> iteration

$$\theta_0 = 0.19, \theta_1 = 0.43, \theta_2 = 0.41$$

$$J = 5.4$$

After 2<sup>nd</sup> iteration

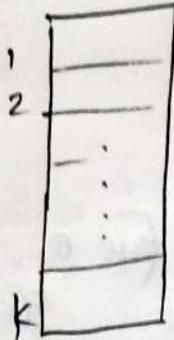
$$\theta_0 = 0.25, \theta_1 = 0.55, \theta_2 = 0.55$$

$$J = 3$$

If datasets are large we needn't use cross-val.  
Train-test split is sufficient.

24/01/26

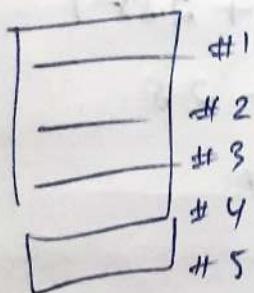
### K-fold cross validn



split as test  
other as train

	$x_1$	$x_2$	$y$	
#1	2	-1	1	
#2	0.5	1.2	0	
#3	1	2	1	
#4	-3	-2	1	
#5	4	0.1	0	

3 fold CV  
 ↳ point Acc  
 ↳ "std dev"  
 ignore  $O_0$  (bias)  
Fold 1:  
 $\{O_1, O_2\} = \{-1.8, 2.8\}$

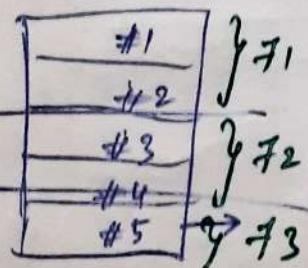


5 samples  
divide into 3 parts

$$f_1 \{O_1, O_2\} = \{2.1, 3.1\}$$

$$f_2 \{O_1, O_2\} = \{1.9, 4\}$$

compute Acc, std dev



8 parts

### Fold 1

Train set  $\{#1, #2, \cancel{\#3}, \cancel{\#4}, \cancel{\#5}\}$

Test set  $\{#5\}$

we can perform  
diff splits  
↓  
sklearn

### Fold 2

Train set  $\{#3, #4, \cancel{\#5}, \cancel{\#1}, \cancel{\#2}\}$

Test set  $\{#1, #2\}$



Train set  $\{#1, #2, #3\}$

Test set  $\{#4, #5\}$

X or  
Train  $\{#1, #3, #4, #5\}$   
Test  $\{#2\}$

Fold 2

Train set  $\{\#3, \#4, \cancel{\#5}\}$

Test set  $\{\#1, \#2\}$

$$T_1 \{ -1.8, 2.8 \}$$

$$T_2 \{ 2.1, 3.1 \}$$

$$T_3 \{ 1.9, 4 \}$$

Fold 3

Train set  $\{\#1, \#2, \cancel{\#3}\}$

Test set  $\{\#3, \#4\}$

$$\frac{\partial}{\partial \theta_j} L^{(0)} = [y - h_{\theta}(x)]x_j$$

$\left[ T_1 \Rightarrow \textcircled{1} \textcircled{2} \textcircled{3} \textcircled{4} \right] \quad (\text{use } \theta_1, \theta_2 \text{ of } T_1)$ 
 $\textcircled{1} \quad h_{\theta}(x) = \theta_0 x_0 + \theta_1 x_1 + \theta_2 x_2$   
 $= -1.8 \times 2 + 2.8 \times 1$   
 $= -3.6 - 2.8$

$$h_{\theta} = -6.4$$

$\textcircled{2} \quad h_{\theta}(2) = \theta_0 x_0 + \theta_1 x_1 + \theta_2 x_2$

$$= \cancel{-1.8} \times 0.5 + \cancel{2.8} \times 1.2$$

$$= \cancel{-0.9} + 3.36 = 2.46$$

$\textcircled{3} \quad h_{\theta}(3) = \theta_0 x_0 + \theta_1 x_1 + \theta_2 x_2$

$$= -1.8 \times 1 + 2.8 \times 2$$

$$= -1.8 + 5.6 = 3.8$$

$\textcircled{4} \quad h_{\theta}(4) = \theta_0 x_0 + \theta_1 x_1 + \theta_2 x_2$

$$= -3 \times -1.8 + 2.8 \times (-2)$$

$$= 5.4 \cancel{- 8} = \cancel{5.6} - 0.2$$

$$x_i \rightarrow 1, 2, 3, 4 \quad \theta_1 = -1.8; \theta_2 = 2.8$$

$$g(\theta_1^T x_1) = \frac{1}{1 + e^{-(-1.8)x_1}} = \frac{1}{1 + e^{-(1.8)}} = \frac{1}{1 + 601.8450};$$

$$g(\theta_2^T x_2) = \frac{1}{1 + e^{-(2.8)}} = \frac{1}{1 + 0.0854}$$

threshold  
≥ 0.5

$$g(\theta_3^T x_3) = \frac{1}{1 + e^{-(3.8)}} = \frac{1}{1 + 0.223707}$$

$$g(\theta_4^T x_4) = \frac{1}{1 + e^{-(0.2)}} = \frac{1}{1 + 0.8187} = 0.5498$$

$$g(\theta^T x_1) = \cancel{0.0016} = \hat{y} = 0 \quad (0)$$

$$g(\theta^T x_2) = 0.92131 = \hat{y} = 1 \quad (1)$$

$$g(\theta^T x_3) = \cancel{0.8171890} = \hat{y} = 1 \quad (1)$$

$$g(\theta^T x_4) = 0.54984 = \hat{y} = 1 \quad (1)$$

For test:-

$$h_0 = \theta_1 x_1 + \theta_2 x_2$$

$$= -1.8 \times 4 + 2.8 \times 0.1$$

$$\text{Accuracy} = 100\% = -7.2 + 0.28$$

$$= -6.92 \quad \rightarrow \quad g(\theta^T x_1) = \frac{1}{1 + e^{6.92}}$$

$$= \frac{1}{1013.319} = 0.000986 = \hat{y} = 1$$

Fold 2 :- Test set is #1, #2

$$\begin{aligned}h_0(x_1) &= \theta_0 x_1 + \theta_1 x_2 \\&= 2.1 \times (2) + 3.1 \times (-1) \\&= 4.2 - 3.1\end{aligned}$$

$$\#1 \ h_0(x_1) = 1.1$$

$$\begin{aligned}\#2 \ h_0(x_2) &= 2.1 \times (0.5) + 3.1 \times (1.2) \\&= 1.05 + 3.72 \\&= 4.77\end{aligned}$$

$$g_0(x_2) = \frac{1}{1 + e^{-(4.77)}} = \frac{1}{1 + 0.008480}$$

$$= 0.9915 \quad \hat{y}_2 = 0$$

$$g_0(x_1) = \frac{1}{1 + e^{-(1.1)}} = \frac{1}{1 + 0.33287}$$

$$= 0.75026$$

Accuracy = 50%

Label  $\hat{y}_1 = 1$

Fold 3 :- Test set is #3, #4

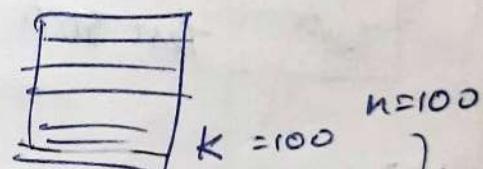
$$\begin{aligned}h_0(x_3) &= \theta_0 x_1 + \theta_1 x_2 \\&= 1.9 \times 1 + 4 \times 2 = 1.9 + 8 \\&= 9.9\end{aligned}$$

$$\begin{aligned}
 h_0(x_{44}) &= 1.9 \times (-3) + 4 \times (-2) \\
 &= -5.7 - 8 \\
 &= -13.7
 \end{aligned}$$

$$\begin{aligned}
 g(\theta)_3 &= \frac{1}{1+e^{-4.9}} = \frac{1}{1+0.00005} = \frac{1}{1.00005} = 0.9995 \\
 g(\theta)_4 &= \frac{1}{1+e^{-(13.7)}} = \frac{1}{1+890911.16597} = \frac{1}{890912.165} \\
 &\quad \text{Label } \hat{y}_3 = 1 \\
 &\quad \text{Label } \hat{y}_4 = 0 \\
 \text{Acc} &= 50\%.
 \end{aligned}$$

$k=n$

$n$  - no. of samples



$k$ -folds  $\rightarrow$  if  $k=n \rightarrow$  Leave One Out Cross Validation

computationally LOOCV is expensive

99 Train  
1 Test

\* When shd we use LOOCV?

\* Should we use?

whenever dataset < 100,  
use LOOCV

Guideline:-

$\sim$  extremely small dataset < 100

[ fold valid technique  
is not to improve accuracy  
just to increase confidence ]

To train a ML Model

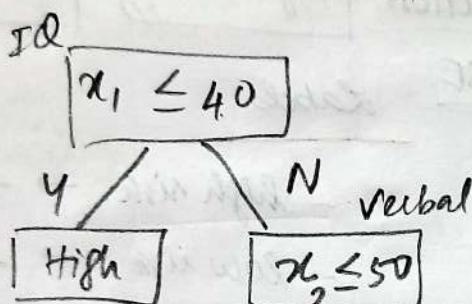
1) Cross valid  
-  $k$ -fold  $> 100$ .

- LOOCV  $< 100$

2) Train test split

# Dataset

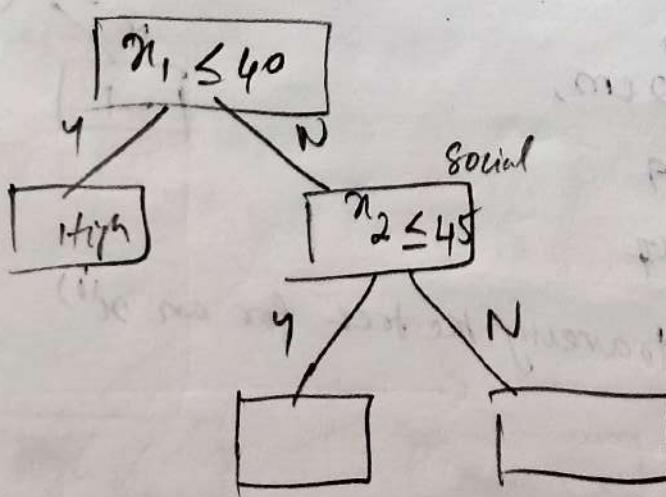
$x_1$	$x_2$	$x_3$	Risk for Autism
IQ	Social	Verbal	
100	79	86	Low
30	40	70	High
80	85	65	Low
90	82	45	High



Predict for

$x^{(i)} = \text{IQ} : 120$   
 $\text{Social} = 60$   
 $\text{Verbal} = 40$   
 Predict

High



Leaf - Label

Internal nodes

↓  
 feature, threshold/  
 value.

Randomly init a tree (theta 100)

↓  
Sum of sq. errors for all samples

↓  
Optimizing (at every <sup>interval</sup> node what  
should be the feature & what is the parameter)

We cannot  
apply grad. descent  
as the tree is not

How do we constraint trees using  
the algo?  
(my constraint?)

Tree  
↓  
Fix-cost function → average  
↓  
Optimizing → keep constraint more trees  
↓  
Constraint trees by raising param  
↓  
Take that tree whose error is ↓

### Example

$x_1$ : Temperature (degree)  
 $x_2$ : Precipitation (cm/hr)

$y$ : 1 Km

2 Km

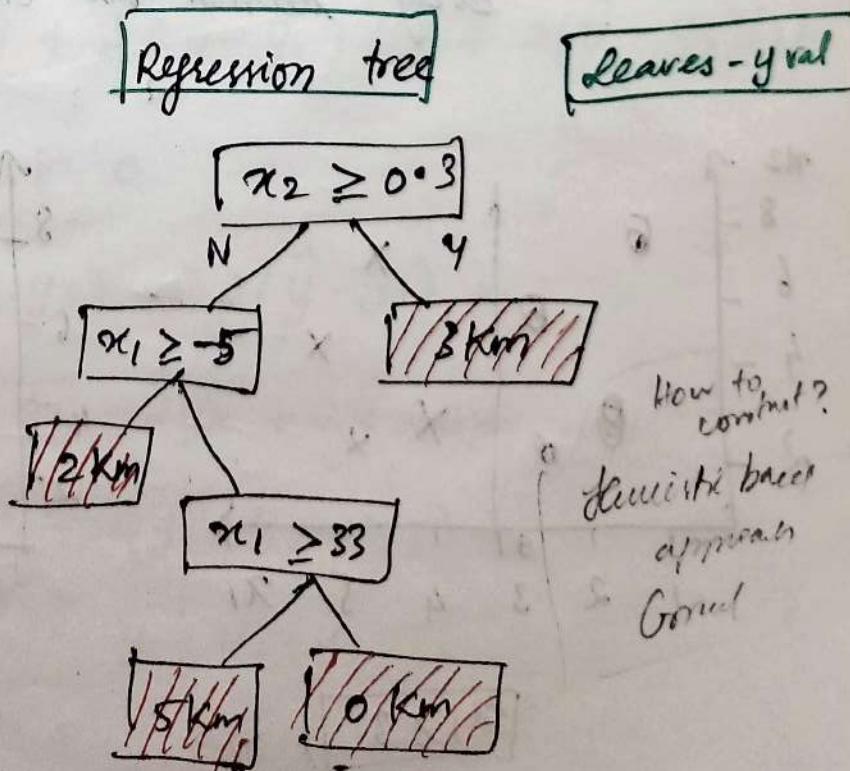
4 Km

⋮

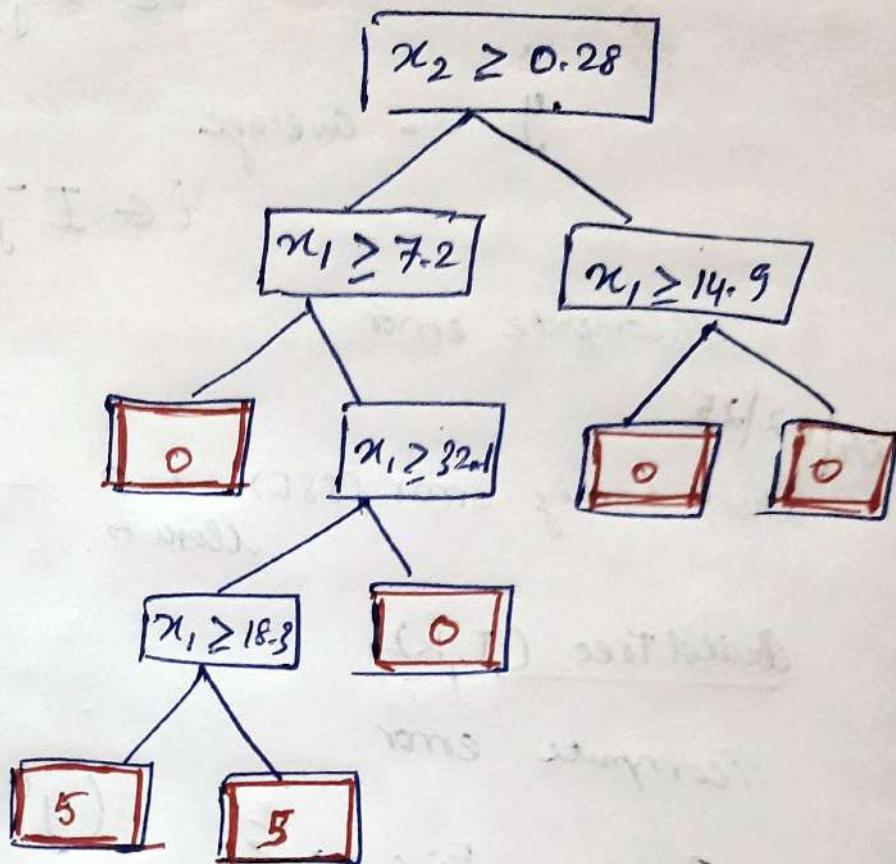
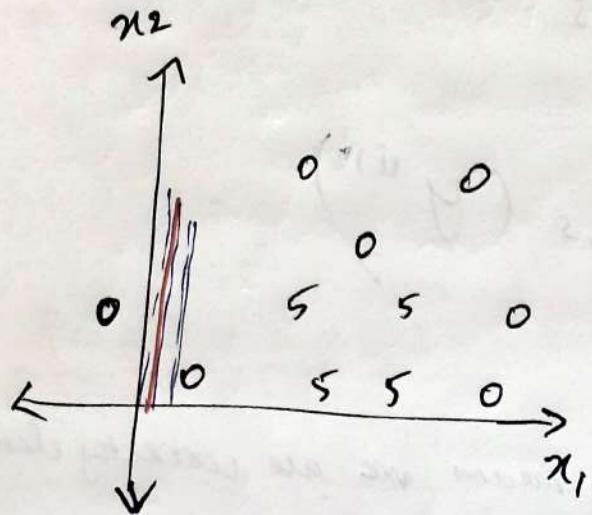
\* It needn't be  
a binary tree..

\* No imp/ order in the depth  
of the node.

\* Pick the tree which has min train error



10 samples



The feature - finite set (can't be infinity)

The split - can have real values (continuous) (- $\infty$  to  $\infty$ )

↓ problematic

app \* split value - 1 <sup>value</sup> point b/w 2 datapoints

⇒ infinite possibilities split val. → trying to eliminate

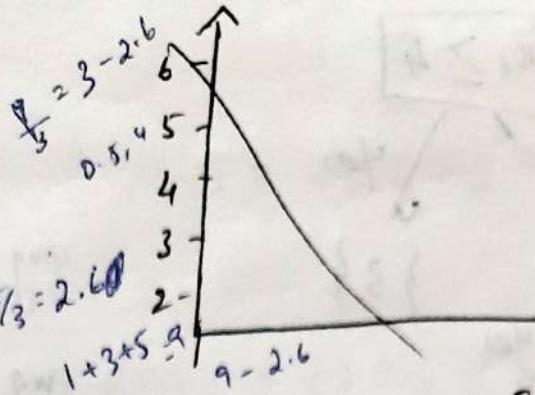
$$I = \{1, 2, 3, \dots, 11\}$$

the other possibilities

$|I| \Rightarrow$  element in the index

06/02/25

#	$x_1$	$x_2$	$y$
1	1	2	1
2	3	6	5
3	5	1	2

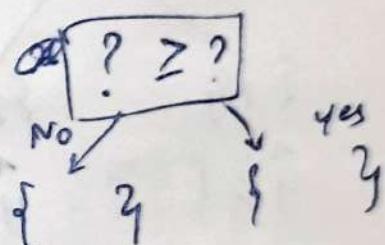


$$I = \{1, 2, 3\}$$

$$K = 2$$

construct 1 node

$$\hat{y} = \text{Avg } y^{(i)}$$



midpoints  
for  $x_1$ ,  $j = x_1, x_2$

for  $x_1$ ,

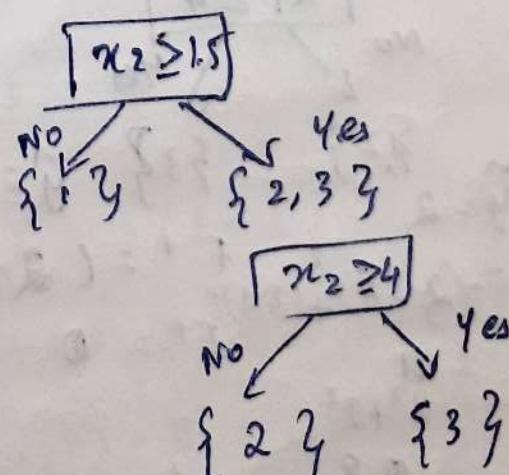
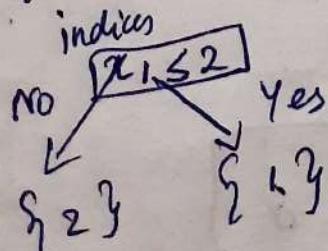
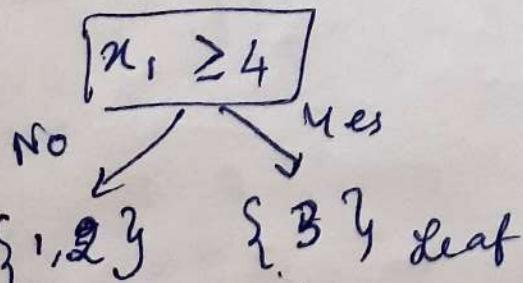
$$S = \{2, 4\}$$

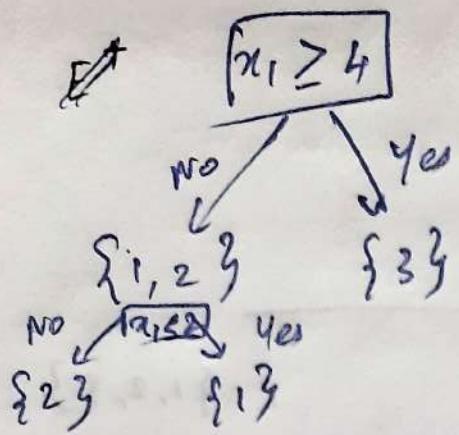
for  $x_2$

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

$$S(1-5, 4)$$

$$S(4, 2)$$

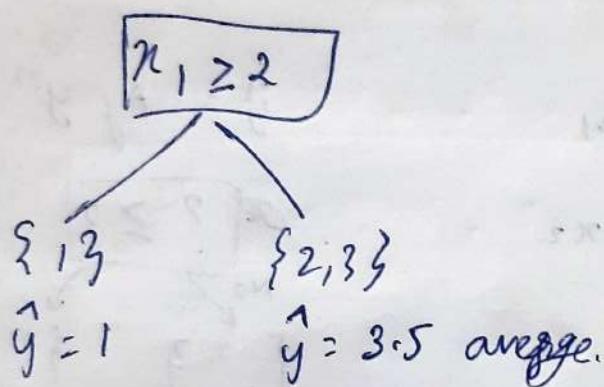




$$E_{x_1, 2}^+ = 5 -$$

$$\text{avg } \hat{y}_{\text{No}} = 1 + 5 = \frac{6}{2} = 3$$

$$\text{avg } \hat{y}_{\text{Yes}} = 1$$



$$E^- = (1-1)^2 \quad E^+ = (3.5-5)^2 + (3.5-2)^2 \\ = 0. \quad E^+ = (2.5)^2 + (2.5)^2$$

$$\therefore E^+ = 2 \cdot 2.5 + 2 \cdot 2.5 = 4.5$$

$$E_{x_1, 2} = 4.5$$

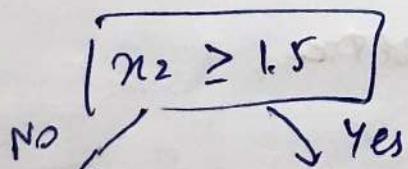
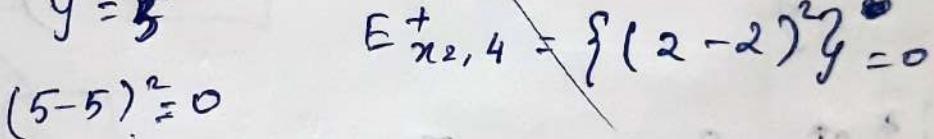
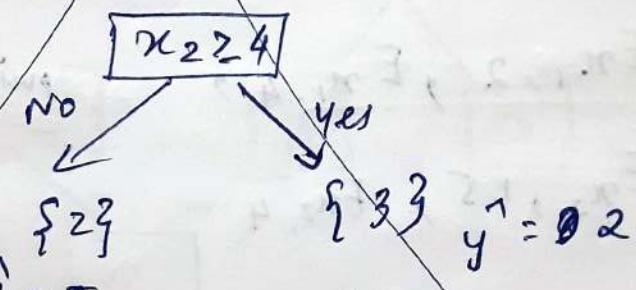
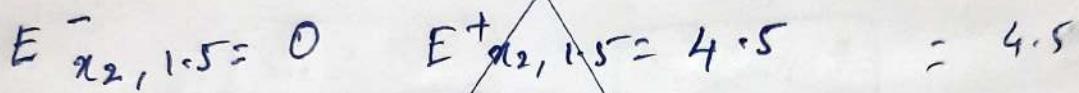
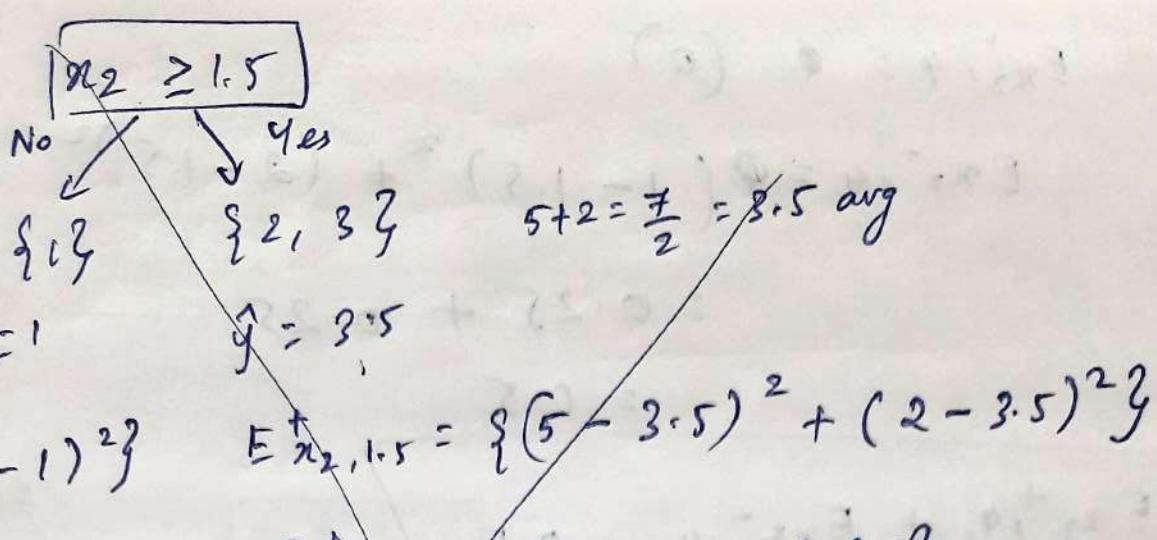
```

graph TD
    A[x1 ≥ 4] -- No --> B["{1, 2}"]
    A -- Yes --> C["{3}"]
    B -- No --> D["x1 ≤ 2"]
    D -- Yes --> E["{1}"]
    D -- Yes --> F["{2}"]
    E -- No --> G["{1}"]
    E -- Yes --> H["{2}"]
  
```

$$\begin{aligned} \hat{y} &= 2 & \hat{y} &= 3 \\ E^- &= (2-5)^2 + (3-1)^2 & E^+ &= (2-2)^2 \\ E^- &= 20 & E^+ &= 0 \\ E^- &= 2^2 + 2^2 \end{aligned}$$

$$E^- = 8 \quad E^+ + E^- = 8$$

$$x_{1,4} = 8$$



{3} {1, 2}

$$\bar{y} = 3$$

$$(x - \bar{x})^2$$

$$\sum x_{2,1.5} = 0 \text{ D}$$

$$\hat{y} = \frac{1+5}{2} = 3$$

$$E_{x_2, 1.5}^+ = \{(1-3)^2 + (5-3)^2\}$$

$$E_{x_2, 1.5}^+ = \{ (-2)^2 + (2)^2 \}$$

$$E_{x_2, 1.5}^+ + E_{x_2, 1.5}^- = 8$$

$$E^+_{\gamma_2, 1.5} = \delta$$

$$x_2 \geq 4$$

$$\hat{y} = 1.5$$

5

W

○  
△  
△  
△

37

2

2

ej

۲۷

5

$$E_{x_2,4}^+ = 0 \quad (0)$$

$$E_{x_2^-,4} = 0(1-1.5)^2 + (2-1.5)^2 \\ = 0.25 + 0.25 \\ = 0.5$$

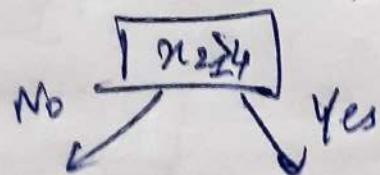
$$\boxed{E_{x_2^+,4} + E_{x_2^-,4} = 0.5}$$

$j^*, s^*$

$$\min \{ E_{x_1,2}, E_{x_1,4}, \underline{\min} \{ 4.5, 8 \\ E_{x_2,1.5}, E_{x_2,4} \} \} = 8, 0.5$$

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

Build tree



$$\bar{y} = 1.5$$

$$\{1, 3\}$$

Yes

$$\{2\}: \bar{y} = 5$$

$$E^- y - \bar{y}$$

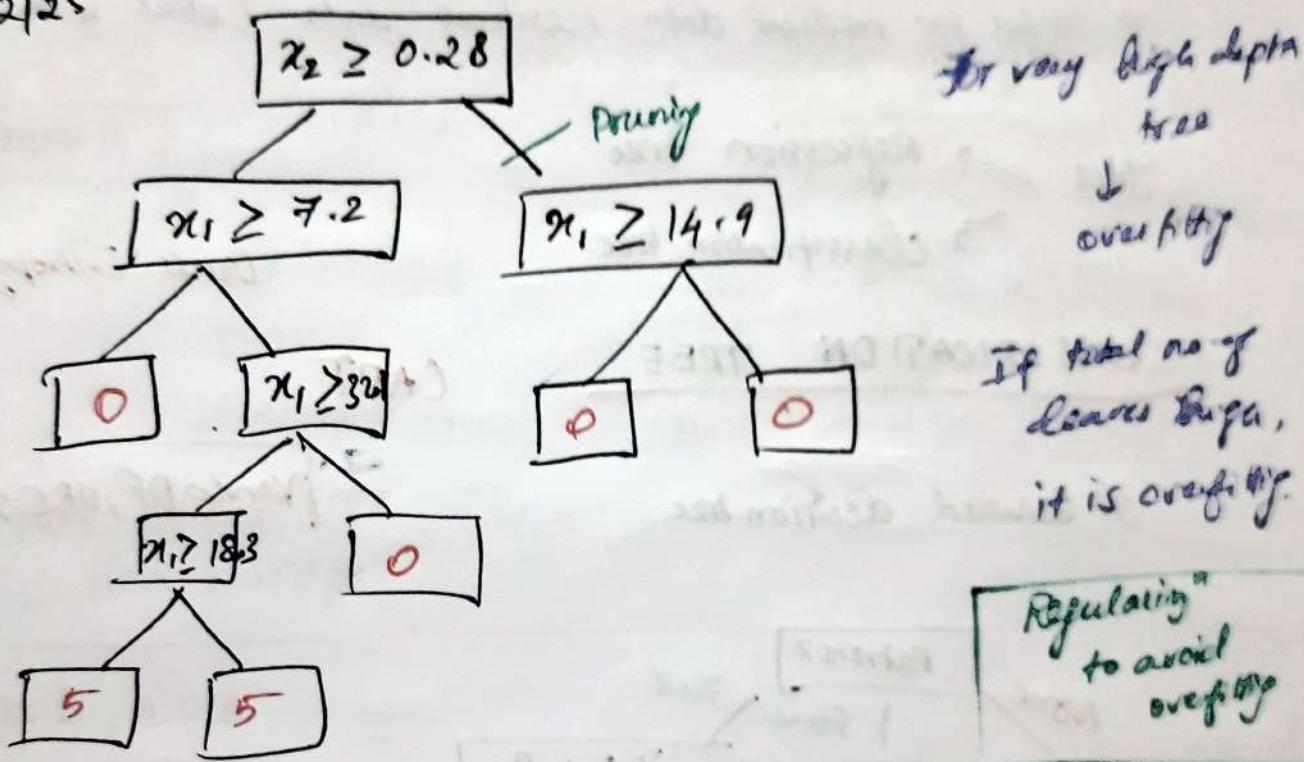
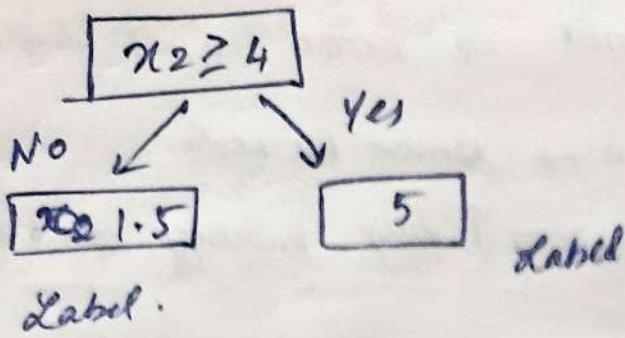
$$(1-1.5)^2 + (2-1.5)^2$$

$$E^+ y - \bar{y}$$

$$E^+ = 0$$

$$0.25 + 0.25$$

$$= 0.5$$



we need to add penalty term for the cost function.

(i)  $C_\alpha(T) = \sum_{i=1}^n L(T(x^{(i)}, y^{(i)})) + \alpha |T|$  → NO. of leaves

cost function:  $\downarrow$  Training loss

of tree

6 leaves in the present

now we reconstruct the tree by add the penalty term.

$\alpha$  - controls how much the tree grows / shrinks.

ii) Pruning → let the tree grow, and then remove the leaves.  
(eliminate the nodes) → remove 1 node → check val accuracy

compute val. error  $\rightarrow$  error (0,0). ideal.

if error is 0,0  $\rightarrow$  remove the node.

if val. error  $\uparrow$  don't knock off. don't touch

Both i) & ii) are complementary.

$\Rightarrow$  useful for medical data, critical data (start with this)

Tree  $\rightarrow$  Regression tree

$\rightarrow$  Classification tree

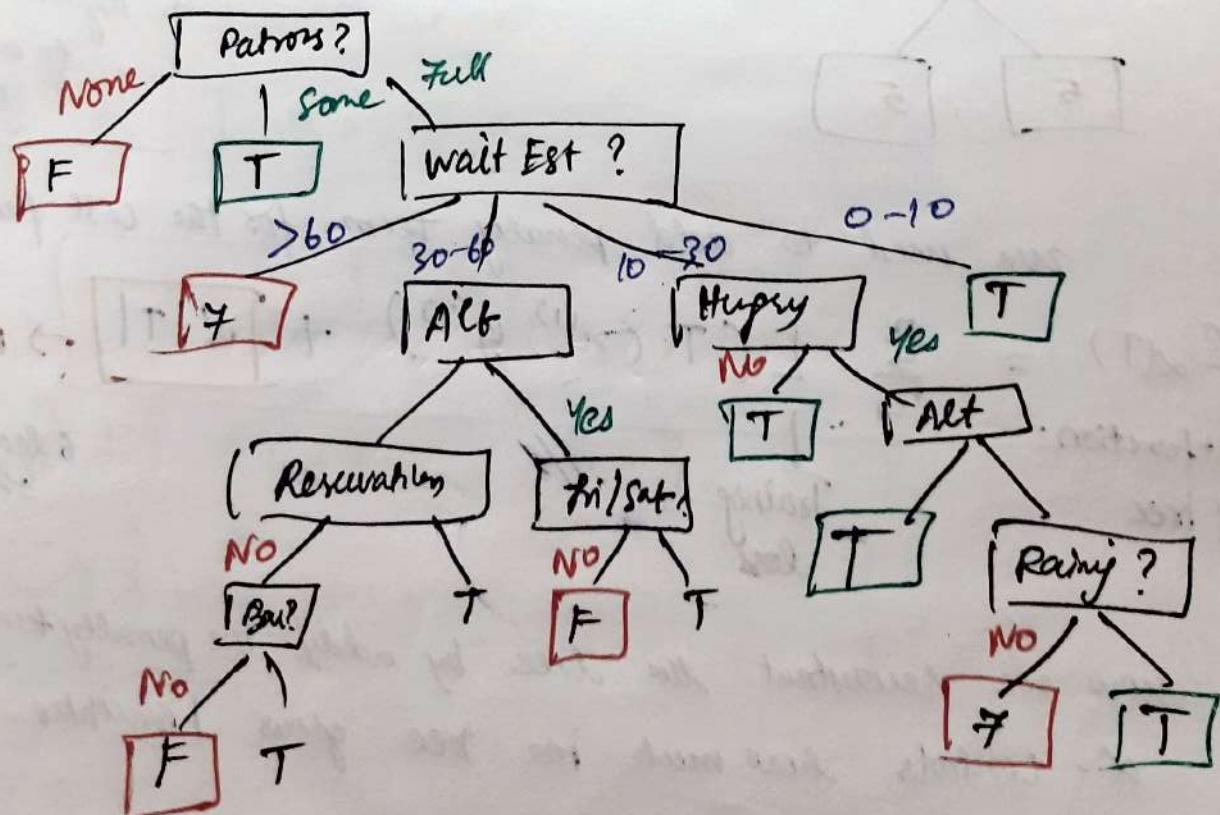
Cross Entropy loss

## CLASSIFICATION TREE

A learned decision tree

CART

Wando DF, UBC 2012

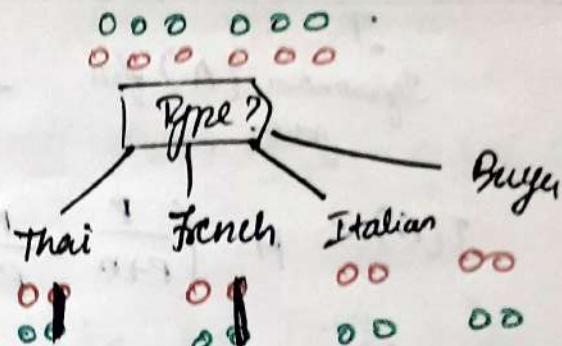
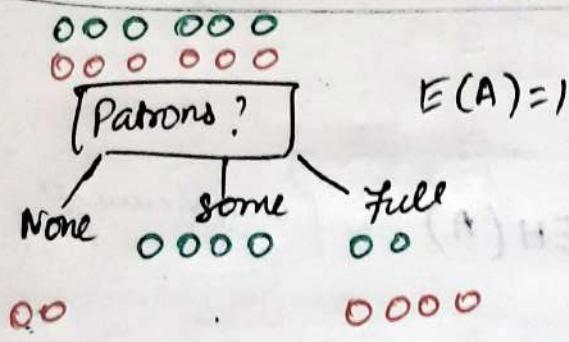


Why didn't we have SSE in logistic  $\rightarrow \hat{y} \rightarrow$  label

Shannon's Info Theory

ENTROPY - randomness / impurity / chaos in data.

$$\Rightarrow - \sum_{i=1}^K p_i \log_2(p_i)$$



↑  
More homogenous,  
less entropy.

$$\text{Eg: } p=6, n=6$$

$$H\left(\frac{p}{p+n}, \frac{n}{p+n}\right) = -\frac{p}{p+n} \log_2 \frac{p}{p+n} - \frac{n}{p+n} \log_2 \frac{n}{p+n}$$

↑  
entropy  
rep by

$$H = - \sum_{i=1}^K -p_i \log_2 p_i$$

(k-labels?)

$$= -\frac{1}{12} \log_2 \frac{1}{12} - \frac{1}{12} \log_2 \frac{1}{12}$$

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

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

$$= \frac{1}{2} + \frac{1}{2} = 0.1$$

How do we find homogeneity?

$$E(A) =$$
$$EH(A) = \sum_{i=1}^k - \left( \frac{p_i + n_i}{p+n} \right) H \left( \frac{p_i}{p_i + n_i}, \frac{n_i}{p_i + n_i} \right)$$

Expectation

$I$   
Information (A) for  $I(A)$   
gain.

$$I(A) = H \left( \frac{p}{p+n}, \frac{n}{p+n} \right) - EH(A) \quad [ \text{entropy across all form} ]$$

We will be max  $I(A)$ , make  $EH(A)$

homogeneous, entropy ↓, info gain ↑

11/02/25

1. Entropy

2. Information gain

$$IG = H(\text{Patrons}) - EH(\text{Patrons})$$

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

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

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

$$= \alpha \left[ \frac{P_1 + n_1}{12} H\left(\frac{0}{2}, \frac{2}{2}\right) + \frac{P_2 + n_2}{12} H\left(\frac{4}{4}, \frac{0}{4}\right) + \frac{P_3 + n_3}{12} H\left(\frac{2}{6}, \frac{4}{6}\right) \right]$$

for ratios  
none, some & full

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

$$= \frac{2}{12} (\log_2 0(-1 \log_2 1)) + \frac{4}{12} (-1 \log_2 1(-0 \log_2 0)) \\ + \frac{6}{12} \left( -\frac{2}{6} \log_2 \frac{2}{6} \left( -\frac{4}{6} \log_2 \left( -\frac{4}{6} \right) \right) \right)$$

$$= \frac{1}{6} (x \infty (0)) + \frac{1}{3} (-1(0) \cdot (\infty)) + \frac{1}{2} \\ + \frac{1}{2} \left( -\frac{1}{3} \log_2 \frac{1}{3} \left( -\frac{2}{3} \log_2 \left( -\frac{2}{3} \right) \right) \right)$$

$$\Rightarrow \frac{1}{2} (-\frac{1}{3} (\log 1 - \log 3)) \times (-\frac{2}{3} (\log 2 - \log 3))$$

$$= -\frac{1}{6} (0 - 1.585) \times (-\frac{2}{3} (1 - 1.585))$$

$$= 0.2641 \times \frac{1.17}{3}$$

$$= 0.2641 \times 0.39$$

expected entropy  
0.455

$$\text{IG}_1(\text{Patrons}) = H(\text{Patrons}) - EH(\text{Patrons}) \\ = 1 - 0.455$$

$$\boxed{\text{IG}_1(\text{Patrons}) = 0.541}$$

## Random forest (numerical example)

Age	Salary	Spend pattern	Buy Product	
16	0	High	Yes	
30	90,000	Low	No	
31	90,000	High	No	
55	5,00,000	Medium	No	

$$B = 2$$

1. Create 2 bags  $B_1$  &  $B_2$

$B_1$ =	Age	Buy Product	Salary	Spend pattern	Buy Pdt
	55			Med	
	30		90,000	Low	
	31		90,000	High	
	55		5,00,000	Med	No

$B_2$ =	Age	Salary	Spend pattern	Buy Pdt
	16	0	High	Yes
	16	0	High	Yes
	31	90,000	High	No
	16	0	High	Yes

Tree

$T_1$

~~for~~

if we pick 2 feature

① for features in (Age, Salary)

~~if~~

for feature in (Age, Salary)

{  
  for each  
  ^ value in features

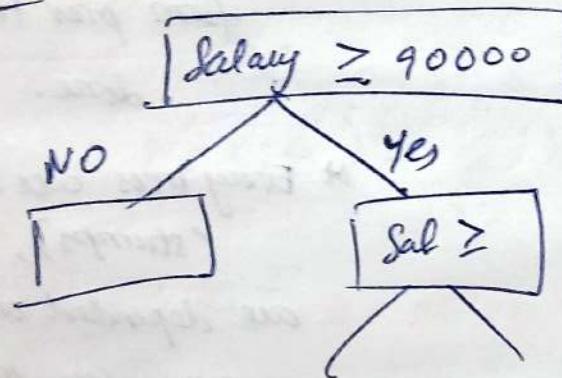
IG = (IGr.value) . append

IG for age as feature

// assume  $IG_r = 1.2$  IG of Age = 1.2

IG of Salary = 1.8 → IG of Salary is high

T1



→ select another 2 random features  $m=2$

Salary, Spend pattern

↓

$IG_2 = 1$

arrows

↓  
 $IG_2 = -0.2$

Then construct another Tree  $T_2$ ,

then aggregate the features

↓

pass key test val → 60 .

At each step instead of considering all features, we select a random subset of features usually  $m = \sqrt{p}$  features or  $m = p$ .

$x_1$	$x_2$	$y$	$r$	$\hat{y}$	new_r
1	2	7	7	7	0
2	3	9	4	3.5	0.5
5	8	1	1	10	-10

example:

Dataset

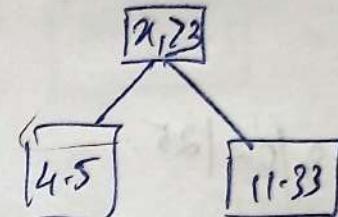
$x_1$	$y$	$r$
1	3	3
2	6	6
3	8	8
4	11	11
5	15	15

$$\hat{y} = r + \lambda(\text{value}) \times (\text{label})$$

i)

$$① \quad r = y$$

② a) fit a model



b, c) Compute  $\hat{y}$

#1, 2 -

#3, 4, 5

assume  $\lambda = 0.1$

$x_1$	$y$	$\hat{y}$
1	3	3.45
2	6	6.45
3	8	9.13
4	11	12.13
5	15	16.13

$$3 + 0.1 \times 4.5 = 3.45$$

$$F(x; \theta) + \lambda f(x; \theta)$$

update  $\hat{y}$

$$6 + 0.1 \times 4.5 = 6.45$$

$$8 * 0.1 =$$

$x_1$	$y$	$\hat{y}$	new_r
1	3	3.45	-0.45
2	6	6.45	-0.45
3	8	9.13	-1.13
4	11	12.13	-1.13
5	15	16.13	-1.13

$$\boxed{\hat{y} - \hat{y}_{\text{pred}} = \text{new}_r}$$

~~8R~~[

$$3 - (3 + 0.1 \times$$

ii) Next

$$\boxed{x_1 \geq 2.5}$$

all samples on left side  $= 0.858$

Influence

$$n_1 = 5$$

$$\sum_{i=1}^B \lambda f_b(x)$$

$$\lambda_0 f_1(5) + \lambda_1 f_2(5) + \lambda_2 f_3(5) + \dots + \lambda_B f_B(5) \\ = 0.8$$

<https://www.kaggle.com/datasets/414/notebooks/17-ensemble-methods-step-by-step>

18/02/25.

Gradient Boosting

Boosted Regression Tree (AdaBoost Regression)

Boosted Classification Tree (AdaBoost Classification)

velocity  
of boostig - Interleave  
com

$$F = f_0(x) + \lambda f_1(x) + \lambda f_2(x) + \dots + f_B(x)$$

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 subsequent version

$$F(x; \theta) = f(x; \theta) + \lambda_1 f_b(x; \theta)$$

c) update the residuals

$$r^{(1)} = r^{(0)} - \lambda f_b(x; \theta)$$

3. Final model

$$f(x; \theta) = \sum_{i=1}^3 \lambda f_i(x; \theta)$$

Numerical ex:-

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

$$f(x) = 0$$

$$f_0(x) = 0$$

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

Step c:- Update residuals

$$\delta^{(1)} = \gamma^{(1)} - \lambda f_1(x^{(1)})$$

$$\gamma^{(2)} = \gamma^{(1)} - \lambda f_1(x^{(1)})$$

$$\delta^{(3)} = \gamma^{(2)} - \lambda f_1(x^{(2)})$$

$$\gamma^{(4)} = \gamma^{(3)} - \lambda f_1(x^{(3)})$$

$$\delta^{(1)} = 4 - [0.1 \times 4.5] = 3.55$$

$$\delta^{(2)} = 5 - [0.1 \times 4.5] = 4.55$$

$$\delta^{(3)} = 6 - [0.1 \times 6.5] = 5.35$$

$$\delta^{(4)} = 7 - [0.1 \times 6.5] = 6.35$$

new best split

residuals are reduced

Iteration 2

Step a:-

Fit a model  $f_2(x) \rightarrow \gamma$

$x_1$	$x_2$	$\gamma$
1	2	3.55
2	3	4.55
3	4	5.35
4	5	6.35

Step c:- Update residuals

$$\delta_0^{(1)} = \gamma^{(1)} - \lambda f_2(x^{(1)})$$

$$\delta^{(2)} = \gamma^{(2)} - \lambda f_2(x^{(2)})$$

$$\delta^{(3)} = \gamma^{(3)} - \lambda f_2(x^{(3)})$$

$$\delta^{(4)} = \gamma^{(4)} - \lambda f_2(x^{(4)})$$

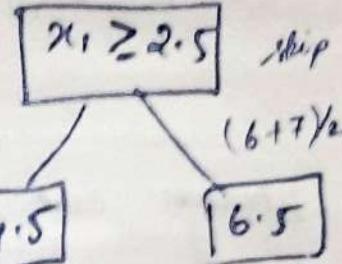
Iteration 2

$$\delta = y$$

Step a:-

a) Fit a model  $f_1(x) \rightarrow \gamma$

$$f_1(x) \rightarrow \gamma$$



Step b:- Update the model ( $\lambda = 0.1$  assume)

$\gamma^{(1,2)}$   $\gamma^{(3,4)}$

$\gamma^{(1,2)}$   $\gamma^{(3,4)}$

A

$$r_1 = 3.55 - (0.1 \times 4.05) = 3.145$$

$$r_2 = 4.55 - (0.1 \times 4.05) = 4.145$$

$$r_3 = 5.35 - (0.1 \times 5.85) = 4.765$$

$$r_4 = 6.35 - (0.1 \times 5.85) = 5.765$$

residuals have  
reduced even more.

$$f(x; \theta) = f_0(x) + \lambda f_1(x) + \lambda f_2(x) + \lambda f_3(x) \dots + \lambda f_b(x)$$

Whenever cont data on  $y \rightarrow$  to measure effectiveness  $\rightarrow$  MSE, SPE  
residuals

but in class " it doesn't work well (residual update approach)

↓  
ADABOOST class".

ADABOOST Classif":

Setup  $\Rightarrow$  Binary class" settig

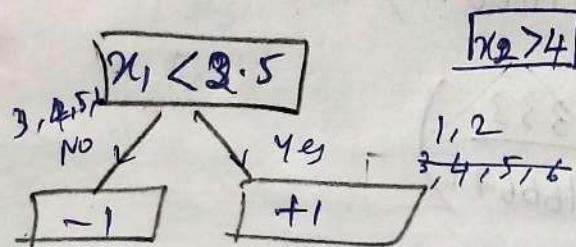
$\Rightarrow$  Multiple weak learners / models with target  $\in \{-1, +1\}$

## Adaboost class

example:

$$w^{(n)} = 1$$

If 1



#	$x_1$	$x_2$	$y$	$w$
1	1	2	+1	1
2	2	3	+1	1
3	3	3	+1	1
4	4	5	-1	1
5	5	5	-1	1
6	6	6	-1	1

$$\epsilon^{(1)} = \frac{\sum_{i=1}^N w^{(1)} I\{y^{(1)} \neq f_b(x^{(1)}; \theta_b)\}}{\sum_{i=1}^N w^{(1)}}$$

$$\epsilon_1^{(1)} = \frac{w^{(1)} \times 0}{6} = \frac{0}{6} = 0$$

$$\epsilon^{(2)} = \frac{\sum_{i=1}^N w^{(2)} I\{y^{(2)} \neq f_b(x^{(2)}; \theta_b)\}}{\sum_{i=1}^N w^{(2)}}$$

$$= 1 \times 0 = 0/6 = 0 \quad 2.$$

$$\begin{aligned} \epsilon^{(3)} &= \frac{w^{(3)} I\{y^{(3)} \neq f_b(x^{(3)}; \theta_b)\}}{\sum_{i=1}^N w^{(i)}} \\ &= 0 \quad w^{(3)}, \downarrow \quad \sum_{i=1}^N (w^{(i)}) \\ &= 1 \times 1 = 1/6 = 0.16667 \end{aligned}$$

$$\epsilon^{(4)} = \frac{w^{(4)} I\{y^{(4)} \neq f_b(x^{(4)}; \theta_b)\}}{6}$$

$$\epsilon^{(5)} = 0/6 = 0$$

$$\epsilon^{(6)} = 0$$

$$\lambda_b = \frac{1}{2} \log \left( \frac{1 - 0.16667}{0.16667} \right)$$

$$= \frac{1}{2} \log \left( \frac{0.8333}{0.16667} \right)$$

$$\lambda_b = \frac{1}{2} \log (0.8333) - \log (0.16667)$$

$$\lambda_b = 0.8047$$

Oftaining weights

$$w^{(1)} = w^{(0)} \times e^{-0.8047} = 0.4472$$

$$w^{(2)} = w^{(1)} \times e^{-0.8047} = 0.4472$$

$$w^{(3)} = w^{(2)} \times e^{-0.8047} = 0.4472$$

$$w^{(4)} = w^{(3)} \times e^{-0.8047} = 0.4472$$

$$w^{(5)} = w^{(4)} \times e^{-0.8047} = 0.4472$$

$$w^{(6)} = w^{(5)} \times e^{-0.8047} = 0.4472$$

$x_1 \quad x_2 \parallel y \parallel w$

21/02/25 absent

Gradient Boosting Algorithm (Regression)

$$h = \gamma_2 (y - \hat{y})^2$$

$$\frac{\partial L}{\partial \hat{y}} = \hat{y} - y$$

$$\left[ \frac{-\partial h}{\partial \hat{y}} = y - \hat{y} \right]$$

To minimize the loss  
we need to move  
in the dir of -ve grad

residual / Pseudo-residuals

$$* h(y, p) = -[y_i \log(p_i) + (1-y_i) \log(1-p_i)]$$

$\downarrow$   
Cross-Entropy Loss

$$\frac{\partial L}{\partial p} = -\left[\frac{y}{p} + \left(\frac{1-y}{1-p}\right)(-1)\right]$$

$$\frac{\partial L}{\partial p} = \frac{p-y}{p(1-p)} \rightarrow \text{1st derivative}$$

$\rightarrow$  its weighty factor  
we can remove it in  
2nd derivative

$$\left. \begin{array}{l} p - \text{sigmoid func} \\ p = \frac{1}{1+e^{-z}} \end{array} \right\}$$

$$\boxed{\frac{-\partial L}{\partial p} = \frac{y-p}{p(1-p)} \approx y-p}$$

$$\frac{\partial^2 L}{\partial p^2} = \frac{\partial}{\partial p^2}(p-y)$$

$$\boxed{\frac{\partial^2 L}{\partial p^2} = p(1-p)}$$

$\rightarrow$  since 2nd derivative of sigmoid  
func' =  $g(z)(1-g(z))$

Gradient Boosting for Class

Step 1: Initialization

$F_0(X) = \log(-\text{odds}) \text{ of +ve class prob}$

$$F_0 = \log \left( \frac{P}{1-P} \right)$$

$$P(+ve) = \frac{3}{6} = 0.5$$

$$F_0 = \log \left( \frac{0.5}{1-0.5} \right) = \log(1) = 0 \rightarrow \log \text{odds}$$

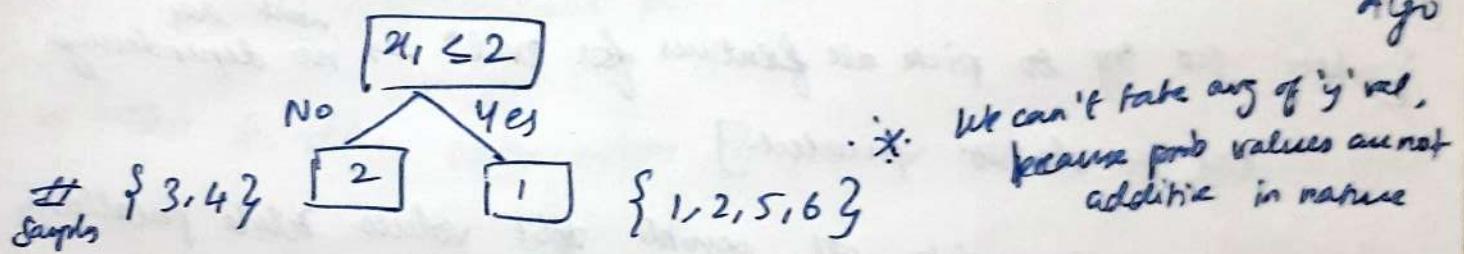
sample	x <sub>1</sub>	x <sub>2</sub>	y
1	2	3	1
2	1	2	0
3	3	4	1
4	4	5	1
5	1	1	0
6	2	2	0

Step 2: compute residuals

$$r_i = y_i - p_i$$

Sample	True - y	Initial $p=0.5$	Residual $r = y - p$
1	1		
2	0	0.5	-0.5
3	1	:	-0.5
4	1	:	0.5
5	0	:	0.5
6	0	:	-0.5

Step 3:- fit a decision stump (model - classifier)  $\Rightarrow$  Apply D. Tree Algo



$$W = \frac{\sum (y - p)}{\sum p(1-p) + \lambda} \rightarrow \text{Log(Odds) adjustmt}$$

$$W = \frac{\text{sum of residuals}}{\sum [prob \times (1-\text{Prob})]} \quad \text{where } \lambda = \text{regularizing constant.}$$

2nd term of L

$\therefore$  we use 'w' to fill the leaves

$$W(\text{right side}) = \frac{-0.5 - 0.5 + 0.5 - 0.5}{(0.5 \times 0.5) + (0.5 \times 0.5) + (0.5 \times 0.5) + (0.5 \times 0.5)}$$

$$W(\text{right}) = \frac{-1}{4} = -0.25$$

Step 4:  $W(\text{left}) = 2 \rightarrow \text{Leaves}$

Element of log odds!  $\rightarrow$  convert it into prob. value  $\rightarrow$  if its  $> 0.5$  then +ve class

$$F = F_0 + \lambda F_1 + \lambda F_2 + \dots$$

$0.1 \rightarrow \text{log(Odds) value.}$

We have constructed a linear boundary in high dim space.

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

monomial of  $x$  with degree 3.

We are capturing the relation b/w  $x_1, x_2, \dots$

if  $d=1000$

$$P \approx O(d^3)$$

$\hookrightarrow$  1 billion features

$\downarrow$

$$1 \text{ billion } O^T$$

Problem:

meaning 1M times slower than what

we do ~~in 1000 dim space~~.

for every step update in grad descent, we need to compute 1B times.

## Kernel

$$K(x, y) = \langle \phi x, \phi y \rangle \rightarrow \text{angle (dot prod)}$$

$\downarrow$   
Kernel func<sup>n</sup> transforms  $x \rightarrow \phi x$  into high dim space  
 $y \rightarrow \phi y$

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

$$y = (y_1, y_2, y_3) \rightarrow \phi(y) = \begin{bmatrix} 1 \\ y_1 \\ y_2 \\ y_3 \\ y_1^2 \\ y_1 y_2 \\ y_1 y_3 \\ y_2^2 \\ y_2 y_1 \\ y_2 y_3 \\ y_3^2 \\ y_3 y_1 \\ y_3 y_2 \end{bmatrix}$$

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

$$\text{Let } x = [1, 2, 3] \\ y = [4, 5, 6]$$

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

$\downarrow$  dot prod b/w  $\phi x \cdot \phi y$

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

$\langle \phi(x), \phi(y) \rangle = 16 + 20 + 72 + 140$   
 $100 + 180 + 72 + 180$   
 $+ 324$   
 $= 1024$

$$[1 \ 2 \ 3 \ 2 \ 4 \ 6 \ 3 \ 6 \ 9] \quad \begin{bmatrix} 16 \\ 20 \\ 24 \\ 20 \\ 25 \\ 30 \\ 24 \\ 30 \\ 36 \end{bmatrix}$$

$$K(x, y) = \langle x, y \rangle^2$$

$$[1 \ 2 \ 3] \quad \begin{bmatrix} 4 \\ 5 \\ 6 \end{bmatrix} \quad 4 + 10 + 18 = (32)^2 \\ = 1024$$

### Types of Kernels

1. Linear Kernel  $K(x, y) = x \cdot y$
2. Polynomial Kernel  $K(x, y)$
3. Radical Basis function (RBF)
4. Laplacian kernel
5. Graph kernel
6. Fisher Kernel

$$= \beta_i + \alpha \left( y^{(i)} - \sum_{j=1}^n K(x^{(j)}, x^{(i)}) \right)$$

## Testing Phase

$$\hat{y}_0(x^{(t)}) = \Theta^T \phi(x^{(t)})$$

One test instance

$$= \sum_{i=1}^n \beta_i \phi(x^{(i)}) \phi(x^{(t)}) \quad \text{- dot prod in transform space.}$$

$$y_{\text{pred}} = \sum_{i=1}^n \beta_i K(x^{(i)}, x^{(t)}) \rightarrow \text{lower dim space.}$$

↓ training      ↓ test sample

## Example

$x_1$	$x_2$	$y$
1	2	0
3	4	0

$n=2$

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

$$x = [x_1, x_2]^T \quad x = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}_{2 \times 1}$$

$\phi(x) \in \mathbb{R}^3$

$$\Theta \in \mathbb{R}^3$$

$$\Theta = [\theta_1, \theta_2, \theta_3]^T$$

$$\Theta = \sum_{i=1}^n \beta_i (\phi(x^{(i)}))$$

3dim      2dim      3dim

$$\Theta = \begin{bmatrix} ? \\ ? \\ ? \end{bmatrix}_{3 \times 1}$$

try by me

$$\beta_i = [\beta_1 \ \beta_2]_{1 \times 2} \quad \phi(x) = \begin{bmatrix} x_1^2 \\ x_2^2 \\ x_1 x_2 \end{bmatrix}_{3 \times 1}$$

~~$\theta = \beta_i \cdot \phi(x)$~~

$$\phi(x) \neq \beta_i$$

$$(\beta_i)^T \begin{bmatrix} \beta_1 \\ \beta_2 \end{bmatrix}_{2 \times 1} \quad (\phi(x))^T = [x_1^2 \ x_2^2 \ x_1 x_2]$$

$$\begin{bmatrix} \beta_1 x_1^2 + \beta_2 x_2^2 + \beta_1 (x_1 x_2) \\ \beta_2 x_1^2 + \beta_2 x_2^2 + \beta_2 (x_1 x_2) \end{bmatrix}$$

~~dim is wrong~~

$$\phi(x^{(1)}) = \begin{bmatrix} 1 \\ 4 \\ 2 \end{bmatrix}_{3 \times 1} \quad \text{kernel trick}$$

$$\theta = \beta_1 \phi(x^{(1)}) + \beta_2 \phi(x^{(2)})$$

$$\phi(x^{(2)}) = \begin{bmatrix} 9 \\ 16 \\ 12 \end{bmatrix}_{3 \times 1} \quad \theta = \begin{bmatrix} \beta_1 + 9\beta_2 \\ 4\beta_1 + 16\beta_2 \\ 2\beta_1 + 12\beta_2 \end{bmatrix}$$

$$\theta = \sum \phi(x^{(i)}) \beta_i + \phi(x^{(t)}) \beta_2$$

$$y_{\text{pred}} = \sum_{i=1}^n \beta_i k(x^{(i)}, x^{(t)}) \rightarrow \text{here we}$$

need all the training data, unlike we need only  $\theta$ s

for LSS. (Computationally expensive  $\sim 1$  billion samples)

Three lines / hyperplanes

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

$$-x_1 + 4x_2 + 7 = 0$$

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

1st line, 1st point

$$d_1 = \frac{|B_0 + B_1 x_1 + B_2 x_2|}{\sqrt{B_1^2 + B_2^2}}$$

$$= \frac{|0 + 2 \cdot 3 + 3 \cdot 4|}{\sqrt{2^2 + 3^2}} = \frac{6 + 12 - 5}{\sqrt{4 + 9}} = \frac{13}{\sqrt{13}} = \cancel{\frac{13}{\sqrt{13}}} = 3.606$$

2nd line - 2nd point

$$d_2 = \frac{|7 + (-1) \cdot 2 + 4 \cdot 3|}{\sqrt{(-1)^2 + 4^2}} = \frac{|7 - 2 + 12|}{\sqrt{1^2 + 4^2}} = \frac{17}{\sqrt{17}} = \cancel{\frac{17}{\sqrt{17}}} = 4.123$$

3rd line

~~$$d_3 = \frac{|10 + (5)(1) + (-12)(1)|}{\sqrt{5^2 + (-12)^2}} = \frac{10 + 5 - 12}{\sqrt{25 + 144}}$$~~

~~3rd line - min value~~

Dataset

$x_1$	$x_2$	$y$
3	4	+1
2	3	+1
1	-1	-1
-2	+1	-1

1st line 2nd point

$$(1) d_2 \frac{|\beta_0 + \beta_1 x_1 + \beta_2 x_2|}{\sqrt{\beta_1^2 + \beta_2^2}} = \frac{-5 + 2 \times 2 + 3 \times 2}{\sqrt{2^2 + 3^2}}$$

$$= \frac{|-5 + 4 + 6|}{\sqrt{13}} = \frac{5}{\sqrt{13}} = 1.25$$

2.21

$$(1) d_3 \frac{-5 + 1 \times (-2) + (-1) \times 3}{\sqrt{13}} = \frac{-5 + 2 - 3}{\sqrt{4}} = \frac{-6}{4} = -1.5$$

= 2.21

$$(1) d_4 \frac{-5 + 2 \times (-2) + 3 \times 1}{\sqrt{13}} = \frac{|-5 - 4 + 3|}{\sqrt{13}} = \frac{6}{\sqrt{13}} = 1.86$$

1st hyperplane margin  
Min = 1.664

$$\min(3.25, 1.25, -1.5, -1.5) \Rightarrow 1.25 \text{ d}_2$$

2nd line  
~~d<sub>1</sub>~~

$$\boxed{-x_1 + 4x_2 + 7 = 0}$$

$$\frac{|7 + (-1) \times 3 + 4 \times 4|}{\sqrt{(-1)^2 + (4)^2}} = \frac{7 - 3 + 16}{\sqrt{16+1}} = \frac{20}{\sqrt{17}} = 4.8$$

$$d_2 = \frac{|7 + (-1) \times (2) + 4 \times (3)|}{\sqrt{15}} = \frac{\cancel{-}7 - 2 + 12}{\sqrt{15}} \\ = \frac{17}{\sqrt{15}} \\ = \frac{4 \cdot 12}{\cancel{-}4 \cdot \cancel{38} \cancel{9}}$$

$$d_3 = \frac{|7 + (-1) \times (1) + (4) \times (-1)|}{\sqrt{15}} = \cancel{-} \\ = \frac{7 + -4 - 1}{\sqrt{15}} \\ = \frac{2}{\sqrt{15}}$$

$$d_4 = \frac{|7 + (-1) \times -2 + 4 \times 1|}{\sqrt{15}} = \frac{7 + 2 + 4}{\sqrt{15}} = \frac{13}{\sqrt{15}} \\ = 3 \cdot \cancel{35}6 \\ = 3 \cdot 15$$

$$\min(d_1, d_2, d_3, d_4) = 0.48$$

2nd hyperplane margin. = 0.48

- Revise:
1. Projection
  2. Norm of a vector
  3. Vector addn
  4. Orthogonality
  5. Subspaces

Line 3:

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

$$d_1 = \frac{|10 + (5 \times 3) + (-12 \times 4)|}{\sqrt{5^2 + 12^2}}$$

$$= \frac{|10 + 15 - 48|}{\sqrt{169}} = \frac{23}{\sqrt{169}} = \cancel{2.3} \quad 1.76$$

$x_1$	$x_2$	$y$
3	4	+1
2	3	
		-1
1		
-2	+1	

$$d_2 = \frac{|10 + (5 \times 2) + (-12 \times 3)|}{\sqrt{169}} = \frac{|10 + 10 - 36|}{\sqrt{169}} = \frac{16}{\sqrt{169}} = \cancel{1.23}$$

$$d_3 = \frac{|10 + 5 \times 1 + (-12 \times -1)|}{\sqrt{169}} = 2.07$$

$$d_4 = \frac{|10 + 5 \times (-2) + -12 \times 1|}{\sqrt{169}} = 0.92$$

3<sup>rd</sup> hyperplane margin = 0.92

Of all 3 margins  $M_1, M_2, M_3$ ,  $M_3$  is max, so the 1<sup>st</sup> hyperplane is optimum.

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

$$RMSE = \sqrt{MSE}$$

$$MAE \rightarrow \text{absol } \epsilon = \frac{1}{N} \sum_{i=1}^N |y_i - \hat{y}_i|$$

21/03/25

Threshold assumed = 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	0	0.60	1
10	1	0.20	0

Calculate 1) Accuracy

2) Precision

3) Recall / sensitivity / True pos Rate

4) Specificity

5) False neg rate (1 - specificity)

6) F1

⇒ ROC curve

Threshold

0.7

0.5

0.4

0.2

Build confusion matrix

	Org Pos	Org Neg	
pred pos	4	3	Org pos : 5
pred neg	1	2	Org neg : 5

1) Accuracy

$$(TP+TN)/\text{Total}$$

$$= (4+2)/10 = 6/10 = 0.60$$

2) Precision

$$TP/(TP+FP)$$

$$= 4/(4+3) = 4/7 = 0.57$$

3) Recall

$$TP/(TP+FN)$$

$$= 4/(4+1) = 4/5 = 0.80$$

4) Specificity

$$TN/(TN+FP)$$

$$2/(2+3) = 2/5 = 0.40$$

5) false positive rate

$$= 1 - \text{specificity} = 1 - 0.40 \\ = 0.60$$

6) F1 score

$$\bar{F}_1 = 2 \cdot \frac{0.57 \times 0.80}{0.57 + 0.80} = \frac{0.912}{1.37} = 0.665$$

Obs	Act	Pred score	Pred Label ( <u>O.7</u> )	O.5	O.4	O.2
1	1	0.85	1	10 FP	0	1
0	0	0.60	<del>0</del> FP 0	0 FP	0 FP	1 FP
1	1	0.70	1	1	0	1
1	1	0.40	<del>0</del> FN	0 FN	1	1
0	0	0.55	0	1 FP	1 FP	1 FP
1	0	0.50	0 FN	1	1	1
0	0	0.65	0	1 FP	1 FP	1 FP
0	0	0.35	0	0	0	1 FP
1	0	0.60	0 FN	1	1	1
0	0	0.20	0	0	0	1 FP

Confusion matrix  
for O.7

	OP	ON
PP	20	80
PN	3	25

Confusion matrix  
for O.4

	OP	ON
PP	5	3
PN	0	2

Confusion matrix  
for O.2

	OP	ON
PP	5	5
PN	0	0

Conf matrix  
for 0.5

4	3
1	2

TPR = Sensitivity

FPR = 1 - specificity

$$TPR = \frac{TP}{TP+FN}$$

$$FPR = 1 - \left( \frac{TN}{TN+FP} \right)$$

for 0.7

$$TPR = \frac{2}{2+3}$$

$$= \frac{2}{5} = 0.4$$

$$FPR = 1 - \left( \frac{5}{5+0} \right)$$

$$= 1 - 1$$

$$= 0$$

for 0.5

$$TPR = 0.8$$

$$\begin{aligned} 1 - 0.40 \\ = 0.60 \end{aligned}$$

Ideal

for 0.4

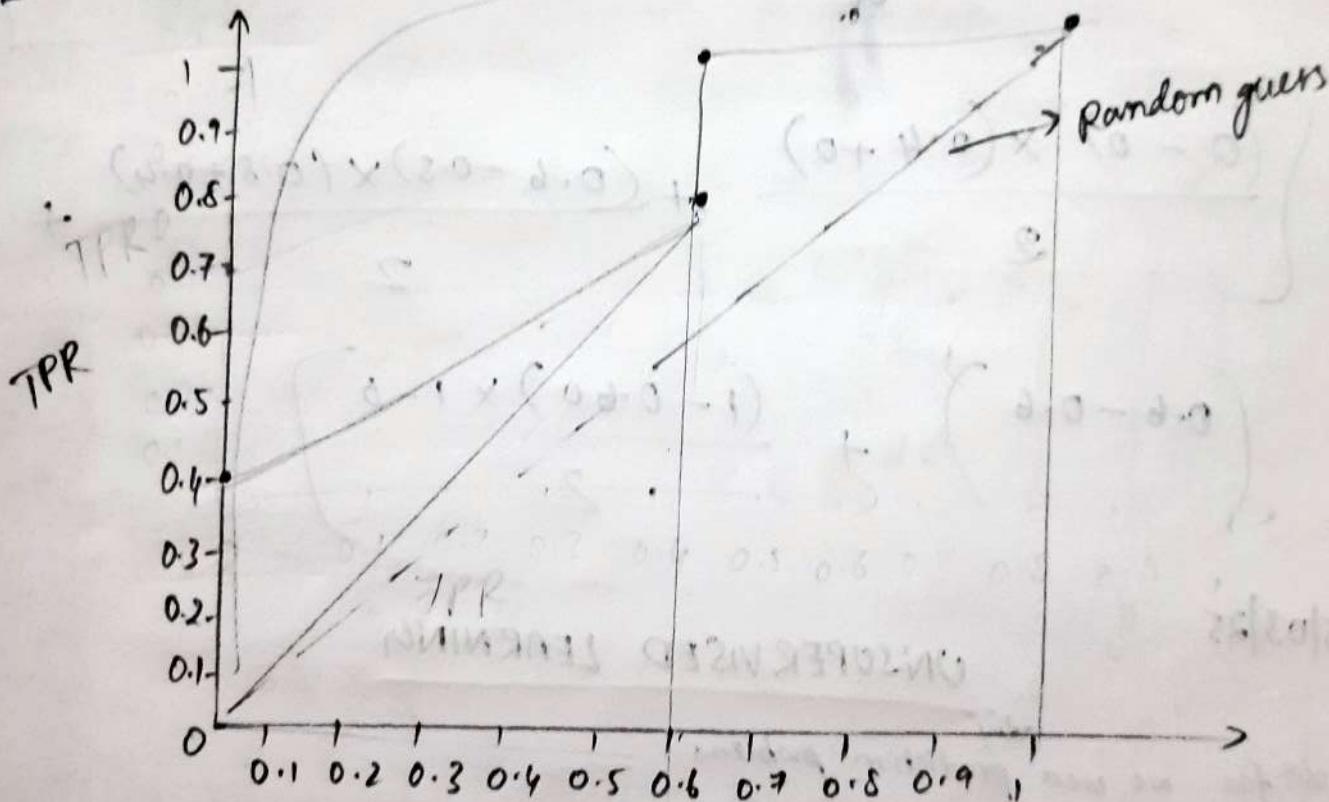
$$TPR = \frac{5}{5+0} = 1$$

$$= 0.60$$

for 0.2

$$TPR = \frac{0}{5+0} = 0$$

$$FPR = \frac{0}{0+5} = 0$$



FPR

$$\frac{1}{2}(a+b)h$$

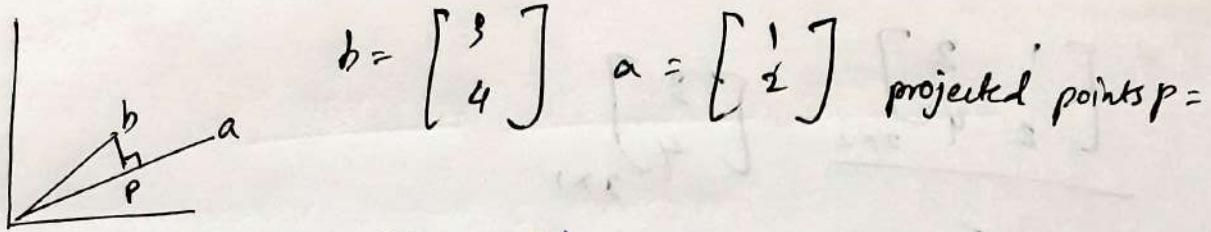
$$\frac{1}{2}(0.4 + 0.8) \times 0.6$$

$$\begin{aligned} 1 - 0.6 &= 0.4 \times 1 \\ &= 0.4 \text{ square units} \end{aligned}$$

$$AUC = \sum_{i=1}^4 \frac{(FPR_i - FPR_{i-1}) \times (TPR_i + TPR_{i-1})}{2}$$

$$= \left[ \frac{(FPR_1 - FPR_0) \times (TPR_1 + TPR_0)}{2} + \frac{(FPR_2 - FPR_1) \times (TPR_2 + TPR_1)}{2} + \right. \\ \left. \frac{(FPR_3 - FPR_2) \times (TPR_3 + TPR_2)}{2} + \frac{(FPR_4 - FPR_3) \times (TPR_4 + TPR_3)}{2} \right]$$

$$= \left[ \frac{(0 - 0) \times (0.4 + 0)}{2} + \frac{(0.6 - 0.5) \times (0.8 + 0.4)}{2} + \right. \\ \left. (0.6 - 0.6) + \frac{(1 - 0.60) \times 1 - 0}{2} \right]$$



$$b = \begin{bmatrix} 3 \\ 4 \end{bmatrix} \quad a = \begin{bmatrix} 1 \\ 2 \end{bmatrix} \quad \text{projected points } p =$$

$$X = \frac{a^T b}{a a^T}$$

$$a^T = \begin{bmatrix} 1 & 2 \end{bmatrix} \begin{matrix} 3 \\ 4 \end{matrix} = \begin{bmatrix} 3+8 \\ 1+4 \end{bmatrix} = \begin{bmatrix} 11 \\ 5 \end{bmatrix}$$

$$a a^T = \begin{bmatrix} 1 & 2 \end{bmatrix} \begin{bmatrix} 1 & 2 \end{bmatrix} = \begin{bmatrix} 1+4 \\ 2+4 \end{bmatrix} = \begin{bmatrix} 5 \\ 6 \end{bmatrix}$$

$$X = \frac{24}{5}$$

(a)  $p = Xa$

$$= \frac{24}{5} \times \begin{bmatrix} 1 \\ 2 \end{bmatrix} = \begin{bmatrix} 24/5 \\ 48/5 \end{bmatrix}$$

~~$[A^T A]^{-1} \rightarrow (24)^{-1}$~~

$$P =$$

$$P = \frac{a a^T \cdot b}{a^T b}$$

~~$= \frac{\begin{bmatrix} 1 \\ 2 \end{bmatrix} \begin{bmatrix} 1 & 2 \end{bmatrix}}{\begin{bmatrix} 1 & 2 \end{bmatrix} \begin{bmatrix} 1 \\ 2 \end{bmatrix}}$~~

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

~~$\therefore \begin{bmatrix} 1+4 \\ 5 \end{bmatrix}$~~

~~$\begin{bmatrix} 8 \\ 4 \end{math}$~~

~~$\therefore \begin{bmatrix} 5 \\ 5 \end{math}$~~ 
  
 ~~$\begin{bmatrix} 3 \\ 4 \end{math}$~~ 
  
 ~~$\therefore \begin{bmatrix} 5 \\ 5 \end{math}$~~

~~$\neq \begin{bmatrix} 1+2 \\ 2+4 \end{bmatrix} = 5$~~

$$= \frac{[\begin{matrix} 1 & 2 \\ 2 & 4 \end{matrix}]_{2 \times 2}}{2 \times 2} \cdot [\begin{matrix} 3 \\ 4 \end{matrix}]_{2 \times 1}$$

5

$$= \frac{[\begin{matrix} 3+8 \\ 6+16 \end{matrix}]}{5}$$

$$= \frac{[\begin{matrix} 11 \\ 22 \end{matrix}]}{5} \Rightarrow [\begin{matrix} 2.2 \\ 4.4 \end{matrix}]$$

### PCA

u to be a subspace of unit length (a)

$\vec{x}$

projected point of  $\vec{x}$  = proj matrix of (u)  $\cdot \vec{x}$

$$= \frac{u u^T}{\underbrace{u^T u}_{\text{unit vector}}} \cdot \vec{x}$$

$$= (u u^T) \vec{x}$$

$$= (\vec{x} \cdot u) u^T$$

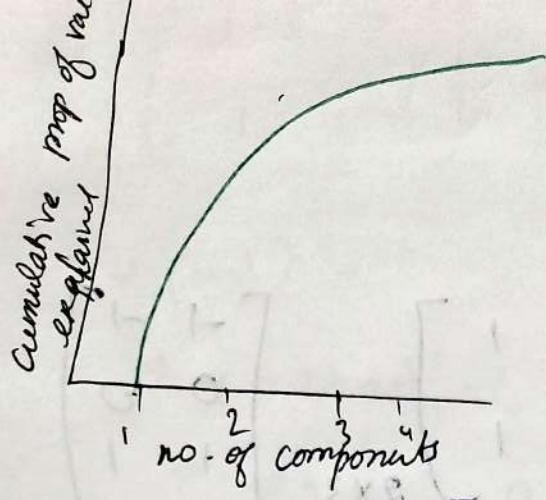
Find u such that

$$\hat{u} = \underset{u}{\operatorname{argmax}} \frac{1}{n} \sum_{i=1}^n \| \operatorname{proj}(u) \cdot x^{(i)} \|_2^2 = (\vec{x} \cdot u)^2$$

Find u which maximizes the variance

$$= \boxed{\underset{u}{\operatorname{argmax}} \frac{1}{n} \sum_{i=1}^n (\vec{x}^{(i)} \cdot u)^2}$$

The norm of a scalar multiplied by a unit vector is equal to the square of its



$$\sigma_2^2 \Rightarrow (100-100)^2 + (200-100)^2 + (300-100)^2$$

$$= \frac{1}{2} (100 + 100)^2 + (200)^2$$

$$\sigma^2 = \frac{1}{2}$$

$$\sigma_2 = 10$$

1st column std

$$\phi_1 = \frac{2-4}{2} = -\frac{2}{2} = -1$$

$$\frac{4-4}{2} = 0$$

$$\frac{6-4}{2} = 1$$

(k=1)

2nd col std

$$\Rightarrow \frac{100-200}{100} = \frac{-100}{100} \Rightarrow -1$$

$$\Rightarrow \frac{200-200}{100} \Rightarrow 0$$

$$\Rightarrow \frac{300-200}{100} \Rightarrow 1.$$

Start step A

$$[n \times d] \times PC = [n \times k]$$

this dataset

is used to train ML Model.

eg:

$$X = \begin{bmatrix} 2 & 100 \\ 4 & 200 \\ 6 & 300 \end{bmatrix}_{3 \times 2} \quad \text{apply PCA} \rightarrow Z = \begin{bmatrix} z_1? \\ z_2? \\ z_3? \end{bmatrix}_{3 \times 1}$$

$R^2 \rightarrow R^1$

① std of each column

$$\frac{x - \bar{x}}{\sigma} \text{ or } \frac{x - \mu}{\sigma} \quad X_{std} =$$

mean of 1st col = 4  $\mu_1$

mean of 2nd col = 200  $\mu_2$

$$\sigma^2 = \frac{1}{n} \sum (x_i - \mu)^2$$

$$= \frac{1}{3} \sum (2-4)^2 + (4-4)^2 + (6-4)^2$$

$$= \frac{1}{3} (2^2 + 0 + 2^2) \Rightarrow 4+4 = \frac{8}{3} = \frac{8}{2} = 4$$

$$\sigma^2 = 2.67 \quad [\sigma_1^2 = 2]$$

$$X_{std} = \begin{bmatrix} -1 & -1 \\ 0 & 0 \\ 1 & 1 \end{bmatrix}$$

②  $\frac{1}{n-1} X_{std}^T X_{std}$  Compute Covariance

$$= \begin{bmatrix} -1 & 0 & 1 \\ -1 & 0 & 1 \end{bmatrix}_{2 \times 3} \begin{bmatrix} -1 & -1 \\ 0 & 0 \\ 1 & 1 \end{bmatrix}_{3 \times 2}$$

$$\Rightarrow \begin{bmatrix} (1+0+1) & (1+0+1) \\ (1+0+1) & (1+0+1) \end{bmatrix}_{2 \times 2}$$

$$\Rightarrow \frac{1}{n-1} \begin{bmatrix} 2 & 2 \\ 2 & 2 \end{bmatrix}_{2 \times 2} \Rightarrow \frac{1}{2} \begin{bmatrix} 2 & 2 \\ 2 & 2 \end{bmatrix} \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$

③ Compute eigen values & eigenvectors of the covariance matrix

$$\begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}$$

④ Decide the needed % of variance explained & choose 'k' principal components.

$$\det(A - \lambda I) = 0$$

$$A - \lambda I = \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix} - \begin{bmatrix} \lambda & 0 \\ 0 & \lambda \end{bmatrix}$$

$$\Rightarrow \begin{bmatrix} 1-\lambda & 1 \\ 1 & 1-\lambda \end{bmatrix} \rightarrow ①$$

⑤ Transform  $X$  into new space of  $A \times PC$  principal components

⑥ Use this  $Z$  to train a model.

$$\det \begin{bmatrix} 1-\lambda & 1 \\ 1 & 1-\lambda \end{bmatrix} \Rightarrow (1-\lambda)(1-\lambda) - (1 \times 1)$$

~~$\rightarrow 1-\lambda-\lambda+\lambda^2 - 1$~~

$$\Rightarrow (1-\lambda)^2 - 1$$

~~$\rightarrow 1-2\lambda+\lambda^2 - 1$~~

$$\Rightarrow 1-2\lambda+\lambda^2 - 1 \Rightarrow$$

$$\Rightarrow -2\lambda + \lambda^2 = 0$$

$$\lambda(\lambda-2) = 0$$

$$\lambda_1 = 0, \lambda_2 = 2$$

To find eigenvectors for  $\lambda_2 = 2$

$$(A - \lambda I)v = 0 \quad \text{subs } \lambda \text{ in ①}$$

$$(A - 2I) = \begin{bmatrix} 1-2 & 1 \\ 1 & 1-2 \end{bmatrix} = \begin{bmatrix} -1 & 1 \\ 1 & -1 \end{bmatrix}$$

$$\begin{bmatrix} -1 & 1 \\ 1 & -1 \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

$$-1x + y = 0 \quad (\text{or}) \quad 1x - 1y = 0$$

$\boxed{x=y}$

$$v_1 = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$

$x, y$

For  $\lambda_1 = 0$

$$(A - 0I) \Rightarrow \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}$$

$$\begin{bmatrix} 1 & 1 \\ -1 & 1 \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

$$x + y = 0$$

$$x = -y$$

28/03/25

$$v_2 = \begin{bmatrix} 1 \\ -1 \end{bmatrix}$$

normalize:

$$\frac{1}{\sqrt{1^2 + 1^2}}$$

$$\begin{bmatrix} 1 \\ -1 \end{bmatrix}$$

→ vector

↓ magnitude

$$\frac{1}{\sqrt{2}} \begin{bmatrix} 1 \\ -1 \end{bmatrix}$$

→ to make it as unit vector

$$\begin{bmatrix} 0 & 0 \\ 1 & 1 \end{bmatrix}$$

④

We only pick  $\frac{1}{\sqrt{2}} \begin{bmatrix} 1 \\ -1 \end{bmatrix}$  ( $v_1$ )

$$z_{pca} = \sum x - v_1$$

$$z_{pca} = \begin{bmatrix} -1 & -1 \\ 0 & 0 \\ 1 & 1 \end{bmatrix}_{3 \times 2} \times \frac{1}{\sqrt{2}} \begin{bmatrix} 1 \\ -1 \end{bmatrix}_{2 \times 1}$$

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

$$= \begin{bmatrix} -\frac{2}{\sqrt{2}} \\ 0 \\ \frac{2}{\sqrt{2}} \end{bmatrix} = \begin{bmatrix} -\frac{\sqrt{2}}{2} \\ 0 \\ \frac{\sqrt{2}}{2} \end{bmatrix}_{3 \times 1}$$

#	$x_1$	$x_2$
1	1	1
2	1	2
3	2	2
4	8	8
5	8	9
6	9	8

K=2

$$D = C_1 = \{2, 3, 4, 6\} \rightarrow 1+2+8+4 = 20/4 = 5$$

$$C_2 = \{1, 5\} \rightarrow \begin{bmatrix} 5 \\ 5 \end{bmatrix}$$

$$\text{Centroid of } C_1 = 3.75$$

$$\text{Centroid of } C_2 = 3$$

$$\begin{bmatrix} 4.5 \\ 5 \end{bmatrix}$$

~~$$② - 3.75 = -1.75$$~~

~~$$② - 3 = 1$$~~

$\textcircled{2}$  belongs to  $C_2$ .

for

~~$$③ - 3.75 = 0.75$$~~

~~$$3 - 3 = 0$$~~

$\textcircled{3}$  belongs to  $C_2$

for

~~$$\textcircled{4}$$~~

~~$$4 - 3.75 = 0.25$$~~

~~$$4 - 3 = 1$$~~

$\textcircled{4}$  belongs to  $C_1$

for  $\textcircled{5}$

~~$$6 - 3.75 = 2.25$$~~

~~$$6 - 3 = 3$$~~

$\textcircled{5}$  belongs to  $C_1$

Centroid

$$C_1 = \begin{bmatrix} 5 \\ 5 \end{bmatrix}$$

$$C_2 = \begin{bmatrix} 4.5 \\ 5 \end{bmatrix}$$

$$d_{11} = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$

$$5 - 1$$

$$5 - 1$$

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

$$4.5 - 1$$

$$4 - 1$$

$$d_{11} \rightarrow C_2 = \begin{bmatrix} 3.5 \\ 4 \end{bmatrix}^2$$

$$d_{11} \rightarrow 4 = \begin{bmatrix} 4 \\ 4 \end{bmatrix}$$

$$\sqrt{32} = 5.65$$

$$12.25 -$$

$$d_{11} \rightarrow c_1 = 5.65$$

$$d_{11} \rightarrow c_2 = 5.31 \checkmark$$

min

$d_{11}$  shd be in cluster ②

$$d_{12} = \begin{bmatrix} 1 \\ 2 \end{bmatrix}$$

$$d_{12} \rightarrow c_1 = \begin{bmatrix} 4 \\ 3 \end{bmatrix}$$

$$d_{12} \rightarrow c_2 = \begin{bmatrix} 3.5 \\ 3 \end{bmatrix}$$

$$\sqrt{4^2 + 3^2}$$

$$= 5$$

$$\sqrt{3.5^2 + 3^2}$$

$$= 4.609 \checkmark$$

min

$d_{12}$  shd be in cluster ②

$$d_{13} = \begin{bmatrix} 2 \\ 2 \end{bmatrix}$$

$$d_{13} \rightarrow c_1 = \begin{bmatrix} 3 \\ 3 \end{bmatrix}$$

$$= \sqrt{3^2 + 3^2}$$

$$= \sqrt{18}$$

$$= 4.24$$

$$d_{13} \rightarrow c_2 = \begin{bmatrix} 2.5 \\ 3 \end{bmatrix}$$

$$= \sqrt{(2.5)^2 + 3^2}$$

$$= \sqrt{15.25}$$

$$= 3.905 \checkmark$$

min

$d_{13}$  shd be in cluster ②

$$d_{14} = \begin{bmatrix} 8 \\ 8 \end{bmatrix}$$

$$d_{14} \rightarrow c_1 = \begin{bmatrix} -3 \\ -3 \end{bmatrix}$$

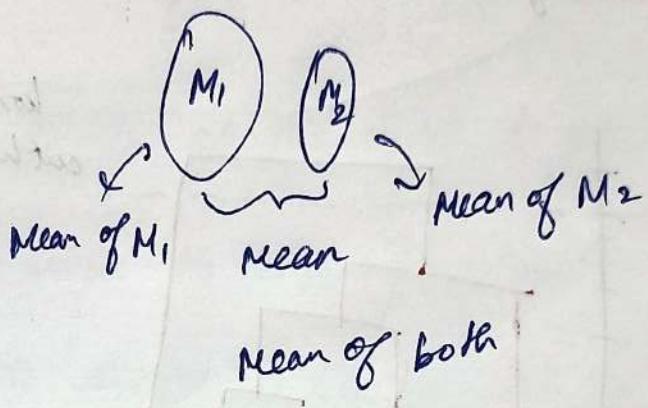
$$d_{14} \rightarrow c_2 =$$

## Complete linkage

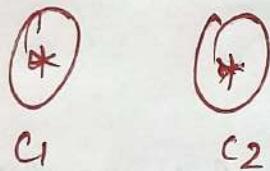
min  $\rightarrow$  max

Choosing min / max  $\rightarrow$  comes from the data.

## Average linkage



## Centroid linkage



cluster centroid  $c_1$

cluster centroid  $c_2$

Dissimilarity b/w  
clst pdt

clst pdt is high  $\rightarrow$  similar  
 $1 - \text{clst pdt}$   $\rightarrow$  will give dissimilarity

Mostly used & complete link

## Hierarchical

Agglomerative  $\rightarrow$  build branches from branches  
example:

Data

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

single linkage

Euclidean distance

$$5\text{-observations} \rightarrow 5C_2 \Rightarrow \frac{5!}{3!2!} = \frac{20}{2} = 10.$$

	A	B	C	D	E
A	0	1	3.6	5	6.4
B	1	0	2.8	4.24	5.65
C	3.6	2.8	0	1.41	2.82
D	5	4.24	1.41	0	1.41
E	6.4	5.65	2.82	1.41	0

Fill up the values.

1st row:

$$A = \begin{bmatrix} 1 \\ 1 \end{bmatrix}, B = \begin{bmatrix} 2 \\ 1 \end{bmatrix} \Rightarrow D_{BA} = \begin{bmatrix} 1 \\ 0 \end{bmatrix} \Rightarrow \sqrt{1^2 + 0^2} = 1$$

$$A = \begin{bmatrix} 1 \\ 1 \end{bmatrix}, B = \begin{bmatrix} 4 \\ 3 \end{bmatrix} = \begin{bmatrix} 3 \\ 2 \end{bmatrix} \Rightarrow \sqrt{9 + 4} \Rightarrow \sqrt{13} = 3.605$$

$$A = \begin{bmatrix} 1 \\ 1 \end{bmatrix}, B = \begin{bmatrix} 5 \\ 4 \end{bmatrix} = \begin{bmatrix} 4 \\ 3 \end{bmatrix} \Rightarrow \sqrt{16 + 9} \Rightarrow 5$$

$$A = \begin{bmatrix} 1 \\ 1 \end{bmatrix}, E = \begin{bmatrix} 6 \\ 5 \end{bmatrix} = \begin{bmatrix} 5 \\ 4 \end{bmatrix} \Rightarrow \sqrt{25 + 16} \Rightarrow$$

$$\sqrt{41} \Rightarrow 6.4$$

~~Row:~~

$$B = \begin{bmatrix} 2 \\ 1 \end{bmatrix}, C = \begin{bmatrix} 4 \\ 3 \end{bmatrix} = \begin{bmatrix} 2 \\ 2 \end{bmatrix} \Rightarrow \sqrt{4 + 4} \Rightarrow 2.82$$

$$B = \begin{bmatrix} 2 \\ 1 \end{bmatrix}, D = \begin{bmatrix} 5 \\ 4 \end{bmatrix} = \begin{bmatrix} 3 \\ 3 \end{bmatrix} \Rightarrow \sqrt{9 + 9} = 4.24$$

$$B = \begin{bmatrix} 2 \\ 1 \end{bmatrix}, E = \begin{bmatrix} 6 \\ 5 \end{bmatrix} = \begin{bmatrix} 4 \\ 4 \end{bmatrix} \Rightarrow 5.65$$

3rd row:

$$C = \begin{bmatrix} 4 \\ 3 \end{bmatrix} \quad D = \begin{bmatrix} 5 \\ 4 \end{bmatrix} \rightarrow \begin{bmatrix} 1 \\ 1 \end{bmatrix} \Rightarrow \sqrt{2} = 1.41$$

$$C = \begin{bmatrix} 4 \\ 3 \end{bmatrix} \quad E = \begin{bmatrix} 6 \\ 5 \end{bmatrix} \rightarrow \begin{bmatrix} 2 \\ 2 \end{bmatrix} \Rightarrow \sqrt{4+4} = 2.82$$

4th row:

$$D = \begin{bmatrix} 5 \\ 4 \end{bmatrix} \quad E = \begin{bmatrix} 6 \\ 5 \end{bmatrix} \rightarrow \begin{bmatrix} 1 \\ 1 \end{bmatrix} \Rightarrow \sqrt{2} = 1.41$$

Min dist  $A \rightarrow B$

Merge A & B  $\rightarrow$  Now clusters are : (AB), CDE

	AB	C	D	E
AB	0	2.82	1.41	5.65
C	2.82	0	1.41	2.82
D	1.41	1.41	0	1.41
E	5.65	2.82	1.41	0

1st row

$$C = \begin{bmatrix} 4 \\ 3 \end{bmatrix} \quad B = \begin{bmatrix} 1 \\ 1 \end{bmatrix} \rightarrow \begin{bmatrix} 2 \\ 2 \end{bmatrix} \Rightarrow \sqrt{4+4} = \sqrt{8} = 2.82$$

$$C = \begin{bmatrix} 4 \\ 3 \end{bmatrix} \quad A = \begin{bmatrix} 1 \\ 1 \end{bmatrix} \rightarrow \begin{bmatrix} 3 \\ 2 \end{bmatrix} \Rightarrow \sqrt{9+4} = \sqrt{13} = 3.602$$

$$D = \begin{bmatrix} 5 \\ 4 \end{bmatrix} \quad A = \begin{bmatrix} 1 \\ 1 \end{bmatrix} \Rightarrow \begin{bmatrix} 4 \\ 3 \end{bmatrix} \Rightarrow \sqrt{16+9} \Rightarrow 5$$

$$D = \begin{bmatrix} 5 \\ 4 \end{bmatrix} \quad B = \begin{bmatrix} 2 \\ 1 \end{bmatrix} \Rightarrow \begin{bmatrix} 3 \\ 3 \end{bmatrix} \Rightarrow \sqrt{9+9} \Rightarrow 4.24 \text{ min}$$

$$E = \begin{bmatrix} 6 \\ 5 \end{bmatrix} \quad A = \begin{bmatrix} 1 \\ 1 \end{bmatrix} \Rightarrow \begin{bmatrix} 5 \\ 4 \end{bmatrix} \Rightarrow \sqrt{25+16} = 6.4$$

$$E = \begin{bmatrix} 6 \\ 5 \end{bmatrix} \quad B = \begin{bmatrix} 2 \\ 1 \end{bmatrix} \Rightarrow \begin{bmatrix} 4 \\ 4 \end{bmatrix} \Rightarrow \sqrt{32} = 5.65 \text{ min}$$

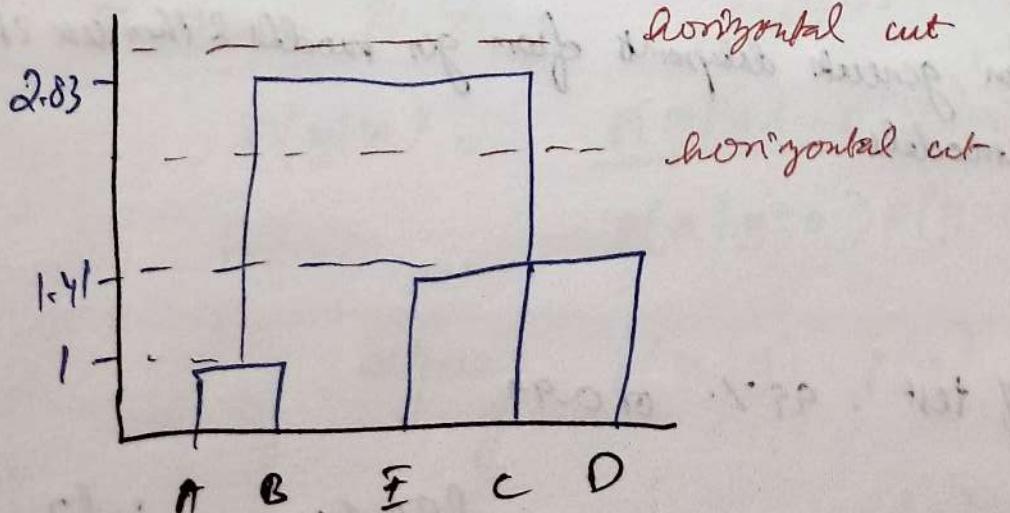
2nd row

$$C \& D \Rightarrow 1.41$$

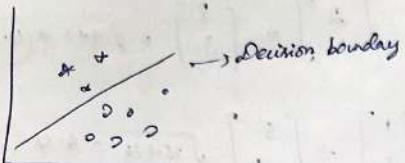
$$C \& E = 2.82$$

3rd row

$$D \& E \Rightarrow 1.41$$



03/04/25

DISCRIMINATIVE LEARNING

Discriminative model  
Generative learning - Generative models

$x_1$  ← o →  
 $x_2$  ← + →  
 $x_3$  ← ++ →  
o - computer-like  
+ - computer.

When dataset is small, we can generate datapoints from gen models & then use it for discriminative models.

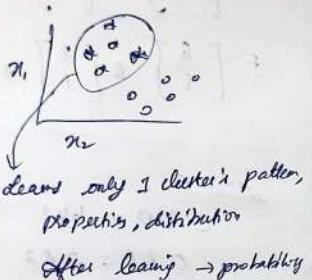
Puzzle:

1) Accuracy of test = 99% or 0.99

2) Tested +ve

3) Rare disease  $\rightarrow 1 \text{ in } 10000$ 

Should I be worried?



$$P(D|T) = \frac{P(D) \cdot P(T|D)}{P(D) P(T|D) + P(T|D') P(D')}$$

$$P(D) = 0.0001$$

$$P(D') = 0.9999$$

$$= \frac{0.0001 \times 0.99}{0.0001 \times 0.99 + 0.01 \times 0.9999} = \frac{0.000099}{0.000099 + 0.000999} = \frac{0.000099}{0.000999} = 0.0098039216$$

Discriminative model

$$P(y|x)$$

↓  
conditional

Generative model

$$P(x, y) = P(x|y) \cdot P(y)$$

Joint probability      ↓  
cond "margin"      ↓  
marginal

Prediction Problems

$$P(y|x) = \frac{\text{likelihood}}{\text{Posterior}} = \frac{P(x|y) \cdot P(y)_{\text{prior}}}{P(x|y=0)P(y=0) + P(x|y=1)P(y=1)}$$

$$\underset{y}{\operatorname{argmax}} \quad P(x|y) \cdot P(y)$$

04/04/25

GENERATIVE MODELSNaive Bayes Classifier

Applic: IRL - spam,  
nospam

email with these words  
↓

"To buy this product"

↓

$J_T \text{ or } [ ]$

Complexity :-

- context transfer is difficult.
- semantic meanings need to be captured

Convert to vector containing numbers

"kij --- enpire ---"

## SEMANTIC SIMILARITY

this rep shd be closer to enpire

relative distance b/w kij & enpire shd be closer to  
that of kij & animal.

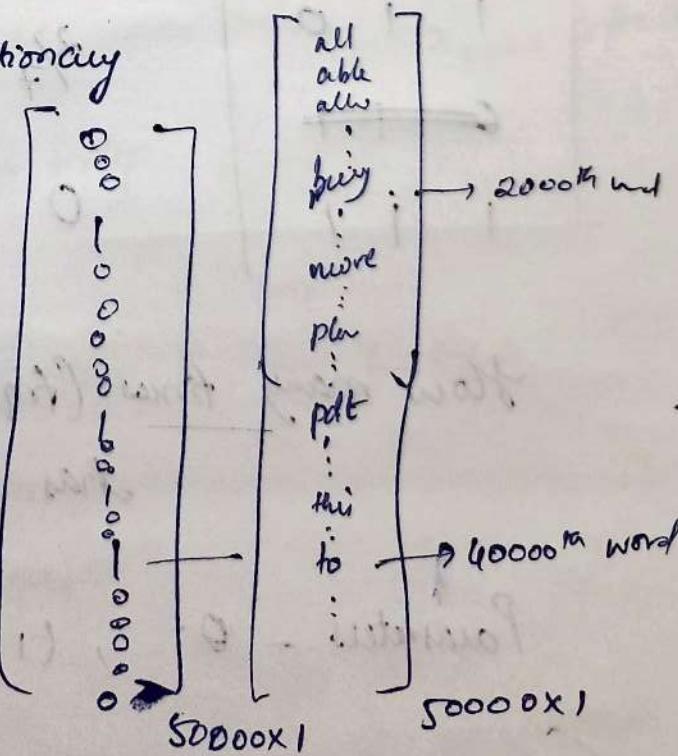
Assume that

we have 50,000 words like in a dictionary

"To buy this product"

Text modality → Vector =

	$x_1$	$x_2$	$x_3$	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	Nospam



Assign prob for every row!

#	1	0	0	Spam
2	0	1	0	Spam
3	1	0	0	Spam
4	1	1	0	NoSpam

$2^3 = 8$  combinations

All possible comb'ns - once we fill all the prob, for the new sample.

$x_1, x_2, x_3$	prob $y=1$ (spam)	Prob $y=0$ (NoSpam)
0 0 0	0.1	0.9
0 0 1	0	0.4 0.4
0 1 0	0.2/7	0
0 1 1	0	0
1 0 0	2/7	0
1 0 1	1/7	0
1 1 0	2/7	1/4
1 1 1	0	1/4

$$\frac{0.28}{7} \frac{20}{20}$$

How many times (freq) that, a combination - 100 has led to spam?

Parameters -  $\theta$ ,  $(1-\theta)$

(Multinomial distb) we have vector of possibilities  $\theta_1, \theta_2, \dots, \theta_{20}, \dots$  etc.

\* if our test sample is 110.  $\rightarrow$  we classify it as spam

$$\frac{2}{7} = 0.28$$

$\Rightarrow$  if our feature explodes (like 50,000 or more)  $\frac{1}{L_1} = 0.25$

we can't explore all possibilities,  
(exponential), computationally expensive.

Though it's a working model, it is not practical.

$$\hat{P}_y = \frac{\sum_{i=1}^n I(y^{(i)}=1)}{n}$$

[Fraction of samples in a class]

Naive Bayes assumes - boolean

follows binomial dist<sup>b</sup>

eg:

X  
1) Free win now

2) Win a prize

3) Hello how are you

4) Let's win it

5) Free lunch today

Y  
spam

spam

Not spam

Not spam

Not spam

Features:

	word <i>x<sub>1</sub></i> "free"	word <i>x<sub>2</sub></i> "win"	Label
#1	Yes/1	Yes/1	spam
#2	No/0	Yes/1	spam
#3	No/0	No/0	Not spam
#4	No/0	Yes/1	Not spam
#5	Yes/1	No/0	Not spam

Text msg: Free win

Step 1:

Calculate prior probabilities.

Prior prob  $y=\text{spam} = 2/5$ ,  $y=\text{not spam} = 3/5$

$$P(y=1 | X) = \frac{P(X|y=1) P(y=1)}{P(X|y=1) P(y=1) + P(X|y=0) P(y=0)}$$

↓    ↓  
 spam      free win

$$= P(\cancel{X \cap Y})$$

① Prior

$$P(\text{spam} = 1) = 2/5 \text{ or } 0.4$$

$$P(\text{Notspam}) = 3/5 = 0.6$$

② Likelihood

For spam class, how many times free appear? ~~Free~~

Free appears 1 of 2 spam class

from train data

$$\left[ \begin{array}{l} P(\text{Free} = \text{Yes} \mid \text{spam} = 1) = 1/2 \\ P(\text{Win} = \text{Yes} \mid \text{spam} = 1) = 2/2 = 1 \end{array} \right]$$

For not spam class,

$$P(\text{Free} = \text{Yes} \mid \underset{\text{or notspam}}{\text{spam}} = 0) = 1/3$$

$$P(\text{Win} = \text{Yes} \mid \text{notspam}) = 1/3$$

③ Bayes rule - Inference | Test phase

$P(\text{Free} = \text{yes}, \text{win} = \text{yes}) \rightarrow \text{Text}$

$$P(\text{spam} \mid X) = \frac{P(X \mid \text{spam})}{\overbrace{P(X \mid \text{spam}) + P(X \mid \text{notspam})}} P(\text{spam})$$

$$\begin{aligned} P(X \mid \text{spam}) &= P(X_1, X_2 \mid \text{spam}) \\ &= P(X_1 \mid \text{spam}), P(X_2 \mid \text{spam}) \\ &= 1/2 \times 1 \end{aligned}$$

$$P(\text{spam} | X) = \frac{1}{2} \times \frac{2}{5}$$

$$= \frac{1}{5} = 0.2$$

$$\boxed{P(\text{spam} | X) = 0.2}$$

similarly

$$P(\text{not spam} | X) = P(X | \text{not spam}) P(\text{not spam})$$

$$P(X | \text{not spam}) \times \frac{3}{5}$$

$$P(X | \text{not spam}) = P(X_1, X_2 | \text{not spam})$$

$$= P(X_1 | \text{not spam}) P(X_2 | \text{not spam})$$

$$= \frac{1}{3} \times \frac{1}{3}$$

$$P(\text{not spam} | X) = \frac{1}{3} \times \frac{1}{3} \times \frac{3}{5}$$

$$= \frac{1}{15} = 0.066$$

$$\boxed{P(\text{not spam} | X) = 0.066}$$

free Discout  $P(X_1 = 1)$

$\downarrow$  free  $\rightarrow 0$   
 $(X_1 = 0)$

09/04/25

# PROBABILISTIC GRAPHICAL MODELS (PCM)

(or)  
BAYESIAN NETWORKS

Bag of words  
Text to vector

Daphne Co-  
Stanford

HMMs  $\rightarrow$  PGM

practive tables to represent  
probabilities.

{ r, r, r, b }

$$P(B_1) = \begin{array}{c|c} B_1 \\ \hline r & \frac{3}{4} \\ b & \frac{1}{4} \end{array}$$

(Marginal)

$$P(B_2 | B_1) = \begin{array}{c|c} B_2 & B \\ \hline r & \begin{array}{c|c} B_1 \\ \hline r & \frac{2}{3} \\ b & 1 \end{array} \\ b & \begin{array}{c|c} B_1 \\ \hline r & \frac{1}{3} \\ b & 0 \end{array} \end{array} =$$

(Conditional)

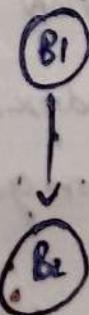
CPT - cond. prob. table

$$P(B_1, B_2) = P(B_2, B_1) = \begin{array}{c|c} B_2 & B \\ \hline r & \begin{array}{c|c} B_1 \\ \hline r & \frac{3}{4} \times \frac{2}{3} \\ b & \frac{1}{4} \times \frac{1}{3} \end{array} \\ b & \begin{array}{c|c} B_1 \\ \hline r & 0 \\ b & 0 \end{array} \end{array}$$

(Joint)

Joint table - prob  
(overall sums to 1)

Graphical representation - SIMPLE GRAPHICAL MODEL



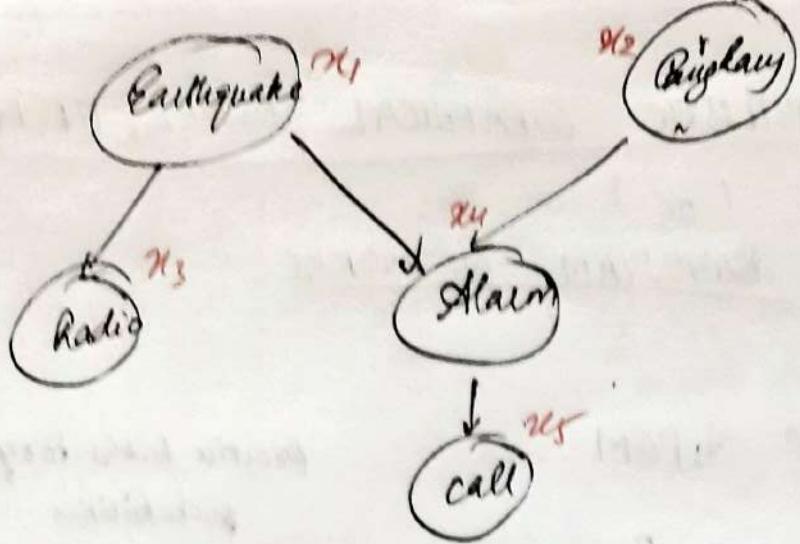
Tedra Pearl - Stanford

Tony Award

Directed Acyclic

Nodes - random variable  
Edges - Direct Influence (causation)

causatn  $\rightarrow$  difficult to compute



define Boolean

call - obs var

Earthquake & Burglary  
hidden

Joint probability

$$P(B, E, A, R, C)$$

$$= P(B) \cdot P(E|B) P(A|BE)$$

$$P(R|BEA) P(C|B, E, A, R)$$

↓

This egn is equivalent to

constructing the joint prob table  $2^5 = 32$

and assumptions

$$P(C|B, E, A, R)$$

knock off others

only depends on

the immediate parent

$$P(R|B, E, A) = P(R|E)$$

we can ignore others

$$P(A|BE) =$$

Conditional  
independence

Don't confuse w/ AND time series - these features  
are not sequential

Constructing Joint prob

E/B	R	A	C
0	0	0	0
0	-	-	-
-	-	-	-
-	-	-	-
1	1	1	1

constructing prob table

is diff  $2^5 = 32$ .

[There are Bayes - we assumed

independence to reduce space]

↓  
Now we are  
assuming differently

ASSUMPTION

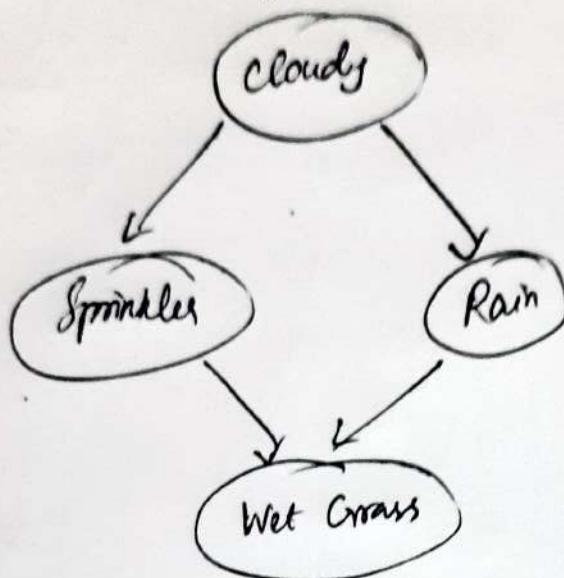
$x_i$ 's are independent  
of their ancestors given  
 $n$  parents

X

$$p(x_1, x_2, \dots, x_d) = \prod_{i=1}^d p(x_i \mid \text{parents}(x_i))$$

Other eg:

$$p(C=F) \rightarrow 0.5 \quad p(C=T)=0.5$$



$$p(S=1)?$$

$$p(R|C)$$

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

$$p(S|C)$$

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

$$p(W|R,S)$$

R	$p(W=F)$	$p(W=T)$
F	1	0
F	0.1	0.9
T	0.1	0.9
T	0.01	0.99