

## Bayes Theorem

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

↓  
Event A when B  
has occurred

$$= P(B|A)P(A)$$

Eg.  $n = 14$

Eg.  $\langle$  outlook = sunny, Temp. = cool,  
Humidity = high, Wind = strong  $\rangle$

$$P(\text{yes}) = \frac{5}{14}, P(\text{no}) = \frac{9}{14}$$

$$P(\text{sunny}|\text{yes}) = \frac{5}{9} \times \frac{2}{5} = \frac{2}{9}$$

$$P(\text{sunny}|\text{no}) = \frac{3}{5} \times \frac{14}{5} = \frac{42}{25}$$

$$P(\text{sunny}|\text{yes}) = \frac{5}{14} \times \frac{9}{14}$$

$$P(\text{sunny}|\text{yes}) = \frac{3}{5} \times \frac{14}{9} = \frac{42}{45} = \frac{2}{5}$$

$$P(\text{sunny}|\text{no}) = \frac{3}{5} \times \frac{14}{5} = \frac{3}{5}$$

$$P(\text{cool}|\text{yes}) = \frac{3}{9} \quad P(\text{high}|\text{yes}) = \frac{3}{9}$$

$$P(\text{cool}|\text{no}) = \frac{1}{45} \quad P(\text{high}|\text{no}) = \frac{4}{5}$$

$$P(\text{strong}|\text{yes}) = \frac{3}{9} \quad P(\text{strong}|\text{no}) = \frac{3}{1}$$

Date \_\_\_\_\_  
 Page No. \_\_\_\_\_

$$P(E) = \frac{2}{9} \times \frac{3}{9} \times \frac{3}{9} \times \frac{3}{9} = \frac{54}{6561}$$

$$P(\text{yes}) = P(yes) P(sw|y) P(c|y) P(h|y) P(alt|y) \\ = 0.0053$$

$$P(\text{no}) = P(\text{no}) P(sw|n) P(c|n) P(h|n) P(alt|n) \\ = 0.206 \quad 0.0206$$

$$V_{NB}(\text{yes}) = \frac{V_{NB}(\text{yes})}{V_{NB}(\text{yes}) + V_{NB}(n)} = 0.205$$

$$\boxed{V_{NB}(\text{no}) = 0.795}$$

Eg  $X = \{\text{Real, SUV, Domestic}\}$

$$P(\text{real } R|y) = \frac{P(R|y) P(R)}{P(R)} = \frac{\frac{3}{5} \times \frac{5}{10}}{\frac{5}{10}} = \frac{3}{5}$$

$$P(R|n) = \frac{\frac{2}{5} \times \frac{5}{10}}{\frac{5}{10}} = \frac{2}{5}$$

$$P(\text{SUV}|y) = \frac{\frac{1}{5} \times \frac{4}{10}}{\frac{5}{10}} = \frac{4}{25} \quad \frac{1}{5} \quad \frac{1}{45}$$

$$P(\text{SUV}|n) = \frac{\frac{3}{5} \times \frac{5}{10}}{\frac{5}{10}} = \frac{3}{5}$$

$$P(D|y) = \frac{\frac{2}{5} \times \frac{5}{10}}{\frac{5}{10}} = \frac{2}{5}$$

$$P(D|n) = \frac{\frac{3}{5} \times \frac{5}{10}}{\frac{5}{10}} = \frac{3}{5}$$

$$V_{NB}(y) = P(y)P(R|y)P(S|y)P(D|y)$$

$$= \frac{1}{2} \times \frac{3}{5} \times \frac{1}{5} \times \frac{2}{5}$$

$$= \frac{3}{125}$$

$$V_{NB}(n) = \frac{12}{125}$$

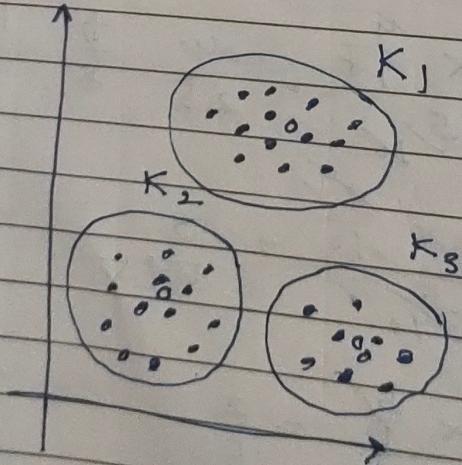
$$P(x|n) > P(x|yes)$$

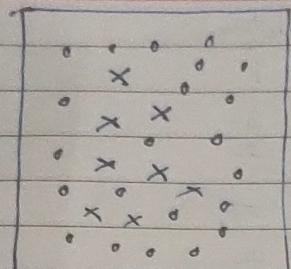
$\Rightarrow$  Not stolen

Clustering

→ no labelled  
data

- \* A list. Based [unsupervised] MLalg where D.P. close to each other are grouped in a given no. of clusters / groups





Two centers

$C_1$   $C_2$

Clustering

Hard clustering

Soft clustering

Each datapoint

assigned to only belongs to a single cluster cluster within

(K-means, K-medoid) a certain prob.

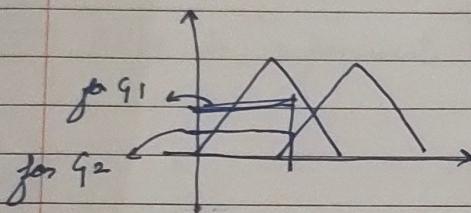
KNN

also known as  
Membership values

(Fuzzy C-Mean)

$$A = \{a, u_a\}$$

element membership value



Suppose we have following pts -

$A_1(2, 10), A_2(2, 5), A_3(8, 4), B_1(5, 8), B_2(7, 5)$   
 $B_3(6, 4), C_1(1, 2), C_2(4, 9)$

$\bullet A_1$

$\bullet B_1$

$\bullet C_1$

$$d_{ave} = \sqrt{(x_i - x_j)^2 + (y_i - y_j)^2}$$

$$d_{A_1 A_1} = 0$$

$$d_{A_1 A_2} = \sqrt{5}$$

$$d_{A_1 B_1} = 6\sqrt{2}$$

$$d_{B_1 B_1} = 0$$

$$d_{B_1 B_2} = \sqrt{13}$$

$$d_{B_1 B_3} = \sqrt{17}$$

$$d_{C_1 C_1} = 0 \quad d_{C_1 B_2} = \sqrt{58}$$

			410	58	1,2	Cluster
A <sub>1</sub>	2	10	0	5.11	8.05 2.65	1
A <sub>2</sub>	2	5	5	4.24	3.15	3
A <sub>3</sub>	8	4	8.49	5	7.28	2
B <sub>1</sub>	5	8	3.61	0	7.21	2
B <sub>2</sub>	7	5	7.07	3.61	6.71	2
B <sub>3</sub>	1	4	7.21	4.12	5.39	2
C <sub>1</sub>	1	2	8.06	7.21	0	3
C <sub>2</sub>	4	9	2.24	1.41	7.62	2

$$\text{New centroid} = \frac{A_3 + B_1 + B_2 + B_3}{5}$$

$$\text{New cen. } A_1 = (2, 10)$$

$$B_1 = (6, 6)$$

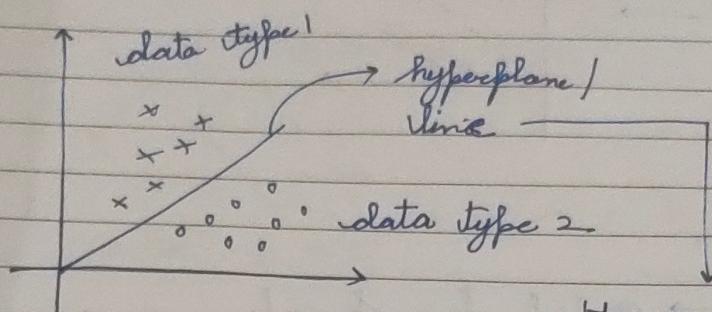
$$C_1 = (1.5, 3, 5)$$

Repeat the process with new centroids until the centroids are stable.

Points	(185, 72)	(170, 56)	Cluster
185, 72	0	21.93	1
170, 56	24.515	0	2
168, 60	18.78	4.47	2
179, 68	7.21	15	1
182, 72	3	20	1
188, 77	5.83	27.65	1

$$\text{New cen. } = [(183.5, 72.25) \rightarrow (169, 58)]$$

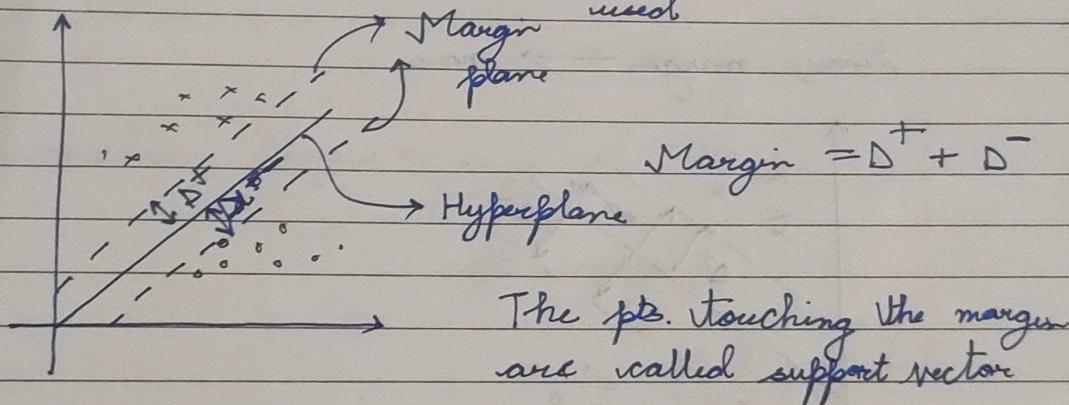
- \* SVM (Support Vector Machine)
- \* Classifier
- \* Regression



However, low margin poses a problem

→ SVM is a supervised MLA

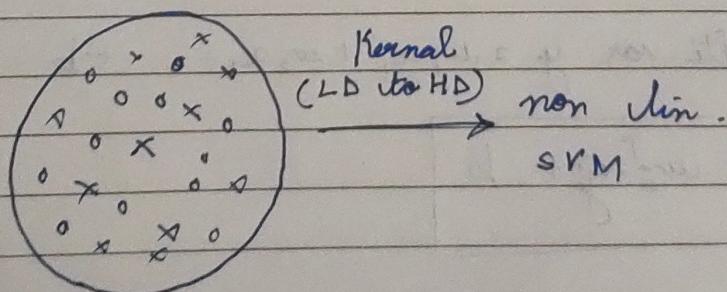
Hence SVM is useful

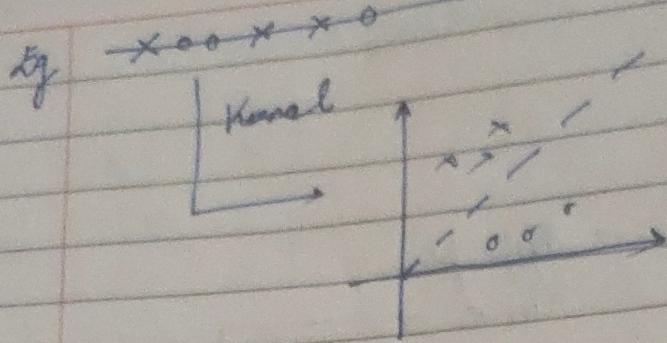


The pts. touching the margin are called support vector

The best classifier has the max. margin

If data is non-linearly separable



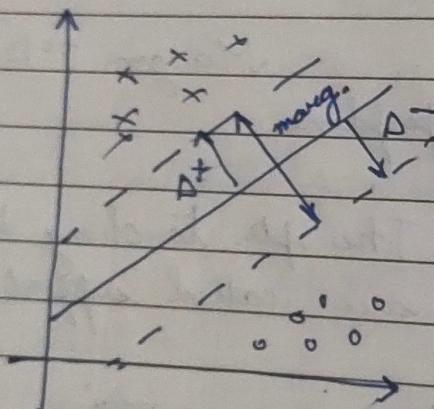


→ Kernel

- Radial basis kernel
- Quadratic kernel

Linear SVM

Since, margin → max.



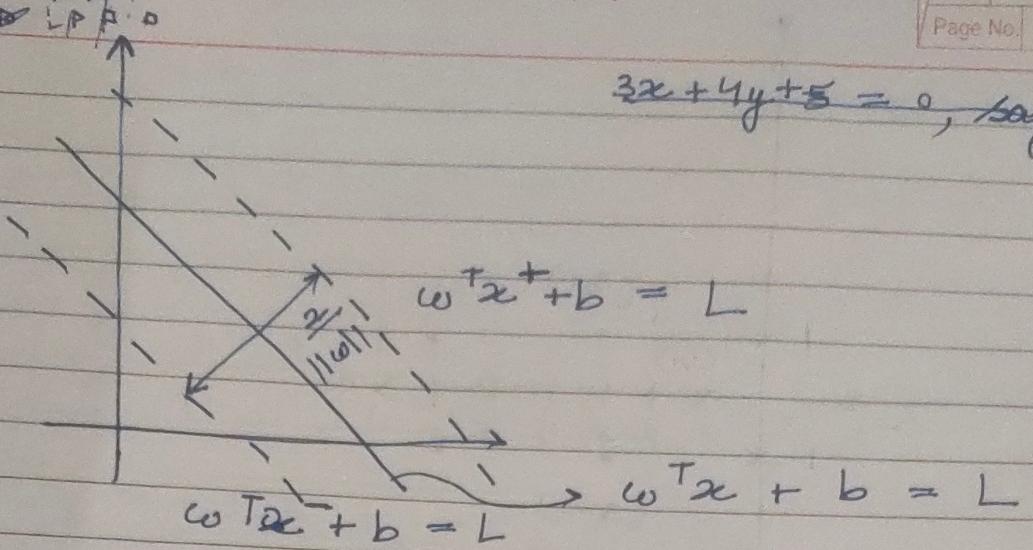
$$y = mx + c, \quad y = b_1 x + b_0$$

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

For multi var.  $y = w_1 x_1 + w_2 x_2 + \dots + b$

$$y = w_i^T x_i + b$$

$$3x + 4y + 5 = 0, \text{ say}$$



$$w^T(x^+ - x^-) = 2$$

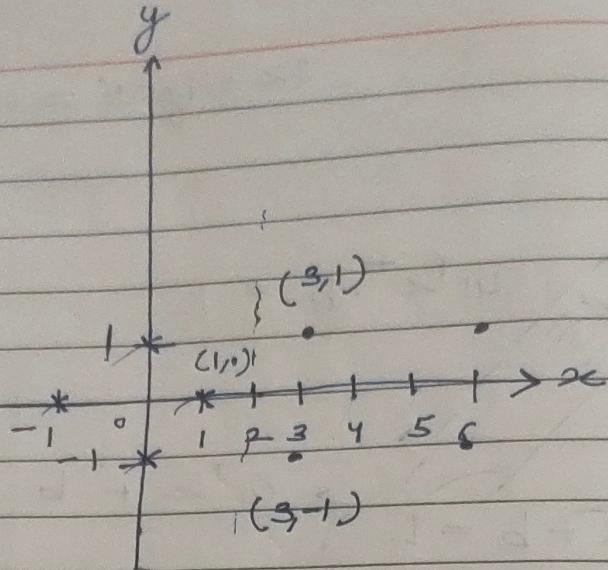
$$x^+ - x^- = \frac{2}{\|w\|} = \text{Margin}$$

→ Loss function

For misclass., Margin =  $\frac{2}{\|w\|} + C_i \sum_{i=1}^n C_i \epsilon_i$   
 (soft margin)

dist. of data  
 pts. from margin  
 plane

$$(3, 1), (3, -1), (6, 1), (6, -1) \rightarrow +1 + 1 \\ (1, 0), (0, 1), (-1, 0), (0, -1) \rightarrow -1$$



$$S_1 \rightarrow (1, 0), S_2 = (3, 1), S_3 = (3, -1)$$

(To find the line  $y = mx + c$ )

Riass

$$\tilde{S}_1 (1, 0, 1), \tilde{S}_2 (3, 1, 1), \tilde{S}_3 (3, -1, 1)$$

$$\alpha_1 \tilde{S}_1 + \alpha_2 \tilde{S}_2 + \alpha_3 \tilde{S}_3 = -1$$

$$\alpha_1 \tilde{S}_2 + \alpha_2 \tilde{S}_1 + \alpha_3 \tilde{S}_3 = +1$$

$$\alpha_1 \tilde{S}_3 + \alpha_2 \tilde{S}_2 + \alpha_3 \tilde{S}_1 = +1$$

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

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

$$2\alpha_1 + 4\alpha_2 + 4\alpha_3 = -1$$

$$\text{Given, } 4\alpha_1 + 11\alpha_2 + 9\alpha_3 = 1$$

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

$$\begin{bmatrix} 2 & 4 & 4 \\ 4 & 11 & 9 \\ 4 & 9 & 11 \end{bmatrix} \begin{bmatrix} \alpha_1 \\ \alpha_2 \\ \alpha_3 \end{bmatrix} = \begin{bmatrix} -1 \\ 1 \\ 1 \end{bmatrix}$$

$$\alpha_1 = -3.5, \alpha_2 = \alpha_3 = 0.75$$

$$y = \omega^T x + c$$

$$\omega^T x = \sum_i \alpha_i s_i$$

$$= -3.5 \begin{bmatrix} 1 \\ 0 \\ 1 \end{bmatrix} + 0.75 \begin{bmatrix} 3 \\ 1 \\ 1 \end{bmatrix} + 0.75 \begin{bmatrix} 3 \\ -1 \\ 1 \end{bmatrix}$$

$$= \begin{bmatrix} -3.5 \\ 0 \\ 1 \end{bmatrix} + \begin{bmatrix} 0 \\ 0 \\ 1.5 \end{bmatrix}$$

$$= \begin{bmatrix} ① \\ ② \\ -2 \end{bmatrix}$$

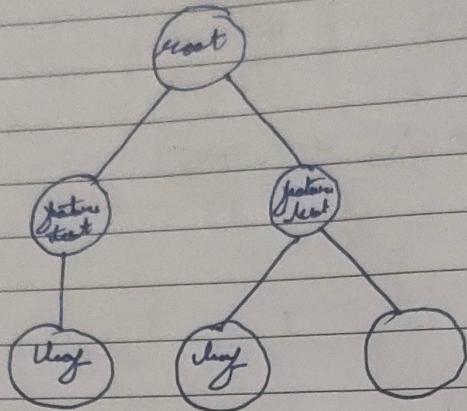
$$y = \begin{bmatrix} 1 \\ 0 \\ -2 \end{bmatrix} \begin{bmatrix} x_1 & x_2 & x_3 \end{bmatrix} + c$$

$$\Rightarrow y = 2$$

## Decision Tree

Supervised algo

→ regression + classification



if (it is raining)

print (Buy umbrella)

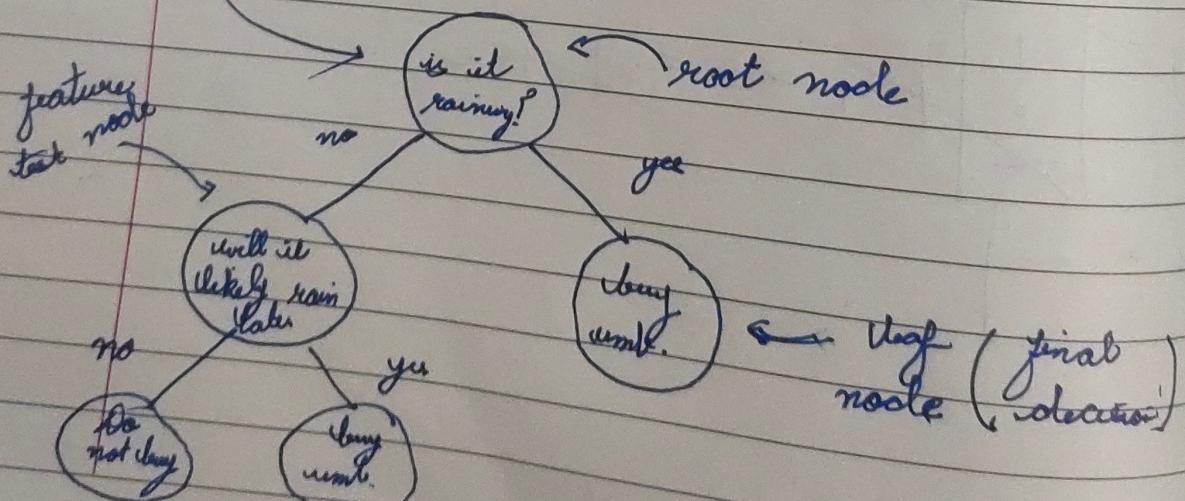
else if (it will likely rain later)

print (Buy umbrella)

} Fig. of how  
dec. tree works

else

(print (Do not buy))



- Information Gain
- Gini Index

$$\text{Gain}(S, \text{outlook}) = \text{Entropy}(S) - \sum \frac{|S_i|}{|S|} \text{Entropy}(S_i)$$

Eg. Values (outlook) = {sunny, overcast, rain}

$$S = [9+, 5-]$$

$$\text{Entropy}(S) = - \sum_{i=1}^n P_i \log_2 P_i$$

$$\Rightarrow \text{Entropy}(S) = - \frac{9}{14} \log_2 \frac{9}{14} - \frac{5}{14} \log_2 \frac{5}{14} = 0.94$$

$$S_{\text{sunny}} = [2+, 3-] \Rightarrow - \frac{2}{5} \log_2 \frac{2}{5} - \frac{3}{5} \log_2 \frac{3}{5} = 0.971$$

$$S_{\text{overcast}} = [4+, 0-] \Rightarrow - \frac{4}{4} \log_2 \frac{4}{4} + \frac{0}{4} \log_2 \frac{0}{4} = 0$$

$$S_{\text{rain}} = [3+, 2-] \Rightarrow - \frac{3}{5} \log_2 \frac{3}{5} - \frac{2}{5} \log_2 \frac{2}{5} \\ = 0.941$$

$$\text{Gain}(S, \text{outlook}) = 0.94 - \frac{5}{14}(0.971) - \frac{4}{14}(0) \\ - \frac{5}{14}(0.971)$$

$\text{Info. Gain} = 0.2404$

Now, you attribute temp.

$$S = [9+, 5-], E(S) = -\frac{9}{14} \log_2 \frac{9}{14} - \frac{5}{14} \log_2 \frac{5}{14} = 0.94$$

$$S_{hot} = [2+, 2-] \rightarrow E(S) = 1$$

$$S_{cold}, E(S) = 0.9183$$

$$S_{cool}, E(S) = 0.8113$$

$$\text{Gain} = 0.94 - \frac{4}{14} - \frac{6}{14} 0.9183 - \frac{4}{14} 0.8113 = 0.0289$$

Sim. for humidity  $\rightarrow$  high, low

$$S = [9+, 5-], E(S) = 0.94$$

$$S_{high} = [3+, 4-], E(S) = -\frac{3}{7} \log_2 \frac{3}{7}$$

$$-\frac{4}{7} \log_2 \frac{4}{7} = 0.9852$$

$$S_{low} = [6+, 1-], E(S) = -\frac{6}{7} \log_2 \frac{6}{7}$$

$$-\frac{1}{7} \log_2 \frac{1}{7} = 0.591$$

$$\text{Gain} = 0.1516$$

wind

Sim. for wind

$$S = [9+, 5-], E(S) = 0.94$$

$$S_{\text{leaves}} = [3+, 3-], E(S) = 1$$

$$\text{Gain} = 0.0478$$

Info. gain of outlook is highest  $\leftarrow$  root node  
 And so on

Day	Temp	Hum.	Wind	Play Tennis
D1	H	H	W	N
D2	H	H	S	N
D8	M	H	W	Y
D9	C	N	W	Y
D1	H	H	S	Y

Table  
 afterwards  
 go sunny

$$\text{Gain}(S_S, \text{temp}) = 0.570$$

$$\text{Gain}(S_S, \text{hum.}) = 0.97 \rightarrow \text{next root node}$$

$$\text{Gain}(S_S, W) = 0.0192$$

Now, go sunny

$$\text{Gain}(S_R, H) = 0.0192$$

$$\text{Gain}(S_R, T) = 0.0192$$

$$\text{Gain}(S_R, H) = 0.0192$$

$$\text{Gain}(S_R, H) = 0.97 \rightarrow \text{next node}$$

Drawback of dec. tree  $\rightarrow$  overfitting

Gini index

$$G(S) = 1 - \sum_{i=1}^n p_i^2$$

\* Root has min. gini index

e.g. Finding  $G(S)$  of decision

$$G(S) = 1 - \left[ \left(\frac{2}{10}\right)^2 + \left(\frac{7}{10}\right)^2 + \left(\frac{1}{10}\right)^2 + \left(\frac{0}{10}\right)^2 \right]$$

$$= 0.58$$

Wear, for attribute money ; case - poor

$$G(S) = 1 - \left(\frac{3}{3}\right)^2 = 0 \quad \left( \frac{3}{3} \text{ as all 3 have some off p & dec.} \right)$$

Fornich

$$G(S) = 1 - \left[ \left(\frac{2}{7}\right)^2 + \left(\frac{3}{7}\right)^2 + \left(\frac{1}{7}\right)^2 + \left(\frac{1}{7}\right)^2 \right]$$

$$= 0.694$$

Weighted avg (money)

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

No for parent :-

→ Yes

→ No

$$q(s) = 1 - \left(\frac{5}{5}\right)^2 = 0 \quad q(s) = 1 - \left[\left(\frac{2}{5}\right)^2 + \left(\frac{1}{5}\right)^2 + \left(\frac{1}{5}\right)^2 + \left(\frac{1}{5}\right)^2\right]$$
$$= 1 - 0.075$$
$$= 0.72$$

Weighted avg = 0.36

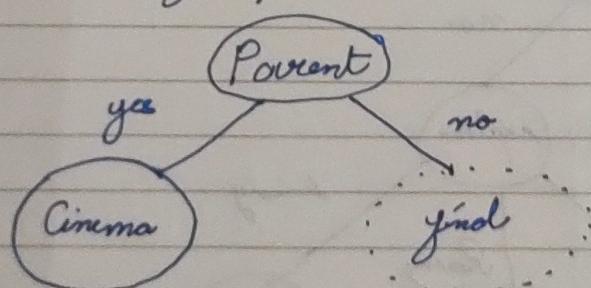
Now for weather :-

→ Sunny      → Windy      → Rainy

$$q(s) = 1 - \left[\left(\frac{1}{3}\right)^2 + \left(\frac{2}{3}\right)^2\right] = 0.44$$
$$q(s) = 1 - \left[\left(\frac{3}{4}\right)^2 + \left(\frac{1}{4}\right)^2\right] = 0.375$$
$$q(s) = 1 - \left[\left(\frac{2}{3}\right)^2 + \left(\frac{1}{3}\right)^2\right] = 0.44$$

Weighted avg =  $0.44 \cdot \frac{3}{10} + 0.44 \cdot \frac{3}{10} + 0.375 \cdot \frac{4}{10} = 0.414$

⇒ I nodes is for parent



Now, finding  $q(s)$  for all those cases when parent said no.

weather —  
 Money —  
 Sunny —  
 $\rightarrow q(s) = 1 - \left(\frac{1}{2}\right)^2 = 0$

Rainy —  
 $q(s) = 1 - \left(\frac{1}{2}\right)^2 = 0$

Windy —  
 $1 - \left(\frac{1}{2}\right)^2 = 0.5$

weighted avg =  $0.5 \times \frac{2}{5} = 0.2$

Money —

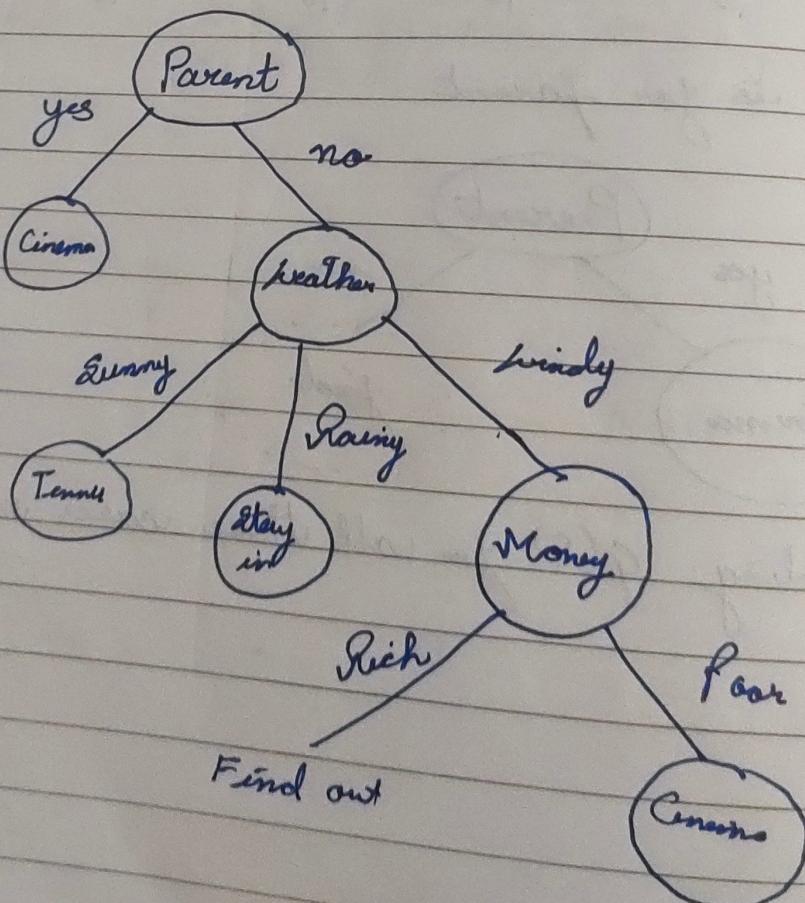
Rich —  
 $q(s) = 1 - \left[\left(\frac{1}{4}\right)^2 + \left(\frac{1}{4}\right)^2 + \left(\frac{1}{2}\right)^2\right]$

Poor —  
 $q(s) = 0$

~~0.375~~ 62.5

wt. avg =  $0.375 \times \frac{4}{5} = 0.300$

$\Rightarrow$  Next root is weather

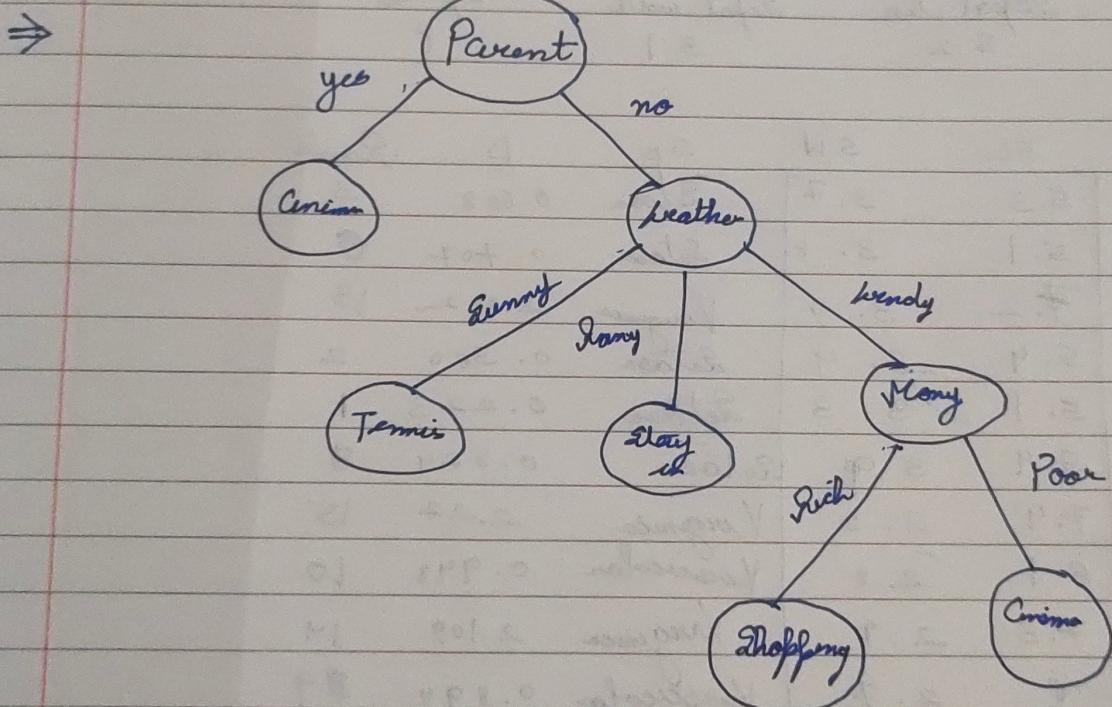


W2 S R T  
 W5 R R S1  
 [W9 W R S ✓]  
 W10 S R T

For weather →

→ Sunny → Rainy → Windy

$$g(S) = 1 - 1 = 0 \quad g(S) = 0 \quad g(S) = 0$$



KNN - K Nearest Neighbours  
Supervised

→ Class + regression

$$* \text{Distance } (x, y) = \sqrt{(x_1 - y_1)^2 + (x_2 - y_2)^2}$$

Ginn rosata —

sepal len.	sepal width	Species?
5.2	3.1	?

Liner date	SL	SW	SP	D	Rank
	5.3	3.7	Setosa	0.608	3
	5.1	3.8	Setosa	0.707	8
7.2	3.0		Virginica	2.002	13
5.4	3.4		Setosa	0.360	2
5.1	3.3		Setosa	0.223	1
5.4	3.9		Setosa	0.824	8
7.4	2.3		Virginica	2.22	15
6.1	2.8		Versicolor	0.948	10
7.3	2.9		Virginica	2.109	14
6	2.7		Versicolor	0.894	89
5.8	2.7		Virginica	0.67	45
6.3	2.3		Versicolor	1.36	12
5.1	2.5			0.6708	5
6.3	2.5	"		0.6088	5
5.5	2.4	"		1.25	11
		"		0.7015	7

Date \_\_\_\_\_  
Page No. \_\_\_\_\_

If  $k = 1 \rightarrow$  choose the class with <sup>the</sup> shortest Dist.  
i.e. ~~distances~~

If  $k = n \rightarrow$  choose the  $n^{\text{th}}$  shortest Dist.