

# Assignment - 11

PAGE No. \_\_\_\_\_  
DATE: / / 202

Q1

Observ.	x	y	xy	$x^2$	$y_p$	$(y_p - y)^2$
1	10	25	250	100	26.25	1.56
2	15	28	420	225	30.18	4.56
3	20	35	700	400	34.11	4.41
4	25	42	1050	625	38.04	0.04
5	30	48	1290	900	41.97	1.06
6	35	45	1575	1225	46.90	0.81
7	40	50	2000	1600	49.53	0.03
8	45	52	2340	2025	53.76	3.24
$\Sigma x = 220$	$\Sigma y = 320$	$\Sigma xy = 9625$	$\Sigma x^2 = 7100$			
$n = 8$						

$$\text{if } Y = mx + c$$

$$m = \frac{n \sum xy - \sum x \sum y}{\sum x^2 - (\sum x)^2}, \quad c = \frac{\sum y - m \sum x}{n}$$

$$m = \frac{8 \cdot 9625 - 220 \cdot 320}{8 \cdot 7100 - 220} = 0.7857, \quad c = 320 - 0.7857 \cdot 220 = 18.39$$

$$\therefore m = 0.7857, \quad c = 18.39.$$

$$\therefore Y = 0.7857x + 18.39.$$

② Predicted sales for 32000 spend

$$x = 32$$

$$Y_p = 0.7857 \cdot 32 + 18.39$$

$$= 25.1424 + 18.39$$

$$= 43.53 \text{ Lakhs}$$

(c)

### MSE

$$MSE = \frac{1}{n} \sum (Y_a - Y_p)^2$$

$$MSE = \frac{\sum (Y_a - Y_p)^2}{n}$$

$$= \frac{27.78}{8} = 3.47$$

for  $R^2 = 1 - \frac{\text{Sum of Squared error}}{(\text{Total variance in } Y)^2}$

$$1 - \frac{\sum (Y_a - Y_p)^2}{\sum (Y - \bar{Y})^2}$$

$$\downarrow \bar{Y} = \frac{320}{8} = 40$$

$$= 1 - \frac{27.78}{676}$$

$$= 1 - 0.041$$

$$= 0.959$$

Q2

Model	Training Acc.	Test Acc.
A	95%	70%
B	85%	83%
C	60%	58%

Q3

Model A shows high variance i.e. overfitting because it has very high training accuracy but lower test accuracy.

Model B shows well balanced behavior  
good accuracies and close accuracies  
in both training and testing.

Model E shows high bias i.e. underfitting  
i.e. both training and testing accuracies  
are low and close.

(ii) As complexity increases means from  
linear to polynomial, biases  
decrease because the model  
can fit training data more closely

while variance increase because the  
model become more sensitive to small  
changes.

(iii) Validation or regularisation can be used

→ It helps to control overfitting and keep  
small model generalised.

Q3

Outlook	Temp.	Humidity	Wind	Play
Sunny	Hot	High	Weak	N
Sunny	Hot	High	Strong	N
Overcast	Hot	High	Weak	Y
Rain	Mild	High	Weak	Y
Rain	Cool	Normal	Weak	Y

(i)

Entropy of Target (Play)

$$\text{Entropy of play} = S(3,2) = -\frac{2}{3} \log \frac{2}{3} - \frac{2}{3} \log \frac{2}{3}$$

$$= -0.6(-0.73) - 0.4(1.32)$$

$$= 0.4421 + 0.5288$$

$$= 0.9709$$

(ii)

$P_{G1}$  (Outlook)

~~$\text{Entropy of sunny } S(2,2) = -0 - \frac{2}{2} \log \frac{2}{2} = 0$~~

~~$\text{Entropy of Overcast } (1,0) = -1.0 - 0 = 0$~~

~~$\text{Entropy of Rain } (2,0) = \frac{2}{2} \cdot 0 - 0 = 0$~~

$$I_{G1} \text{ of Outlook} = S_{bs} - 0 - 0 - 0$$

$$= 0.9709$$

III

### TGr of temp

$$\text{Entropy of NOT } \Theta(1,2) = -\frac{1}{3} \log_2 \frac{1}{3} - \frac{2}{3} \log_2 \frac{2}{3}$$

$$= 0.918$$

$$\text{Entropy of cool } \Theta(0,1,0) = 0$$

$$\text{Entropy of Mid } (1,0) = 0$$

$$\rightarrow \text{TGr of temp} = 0.918 - 0.5508 = 0.420$$

### TGr of Humidity

$$\text{Entropy of High } (2,2) = -\frac{2}{4} \log_2 \frac{2}{4} - \frac{2}{4} \log_2 \frac{2}{4}$$

$$= 1$$

$$\text{Entropy of Normal } (1,0) = 0$$

$$\rightarrow \text{TGr of Humidity} = 0.918 - \frac{4}{5} = 0.171$$

### TGr of wind

$$\text{Entropy of weak } (3,1) = -\frac{3}{4} \log_2 \frac{3}{4} - \frac{1}{4} \log_2 \frac{1}{4}$$

$$= 0.811$$

$$\text{Entropy of strong } (0,1) = 0$$

$$\text{TGr of wind} = 0.918 - \frac{4}{5} = 0.081$$

$$= 0.811 - 0.322 = 0.489$$

Since TGr of outlook is greater than  
so root node is outlook

Outlook R

①

Sunny

Overcast

Rain

↓  
No

↓  
Yes

↓  
Yes.

④

outlook = Rain, Temp = Cool, Humidity = N,  
Wind = Strong.

Since outlook of rain is 0, so

we can directly say from decision  
tree, play is Yes.

④

Age	Income	Bought
22	20	0
25	30	0
28	35	1
35	50	1
40	60	1

→ For a new customer ( $\text{Age} = 30, \text{Income} = 40$ )

→ Step 1: Distance calculation

$$d_+ = (22, 20) \Delta (30, 40)$$

$$= \sqrt{8^2 + 20^2} = \sqrt{464} = 21.5$$

$$d_1 = (28, 30) \text{ & } (30, 40)$$

$$\Rightarrow \sqrt{5^2 + 10^2} = \sqrt{125} = 11.18$$

$$d_2 = (28, 35) \text{ & } (30, 40)$$

$$\Rightarrow \sqrt{2^2 + 5^2} = 5.39$$

$$d_3 = (35, 50) \text{ & } (30, 40)$$

$$\Rightarrow \sqrt{5^2 + 10^2} = 11.18$$

$$d_4 = (30, 60) \text{ & } (30, 40)$$

$$\Rightarrow \sqrt{10^2 + 20^2} = \sqrt{500} = 22.36$$

$$\therefore K = 3$$

$$\text{1st} \rightarrow (28, 35) \rightarrow \text{Buyt} = 1$$

$$\text{2nd} \rightarrow (28, 30) \rightarrow \text{Buyt} = 0$$

$$\text{3rd} \rightarrow (35, 50) \rightarrow \text{Buyt} = 1$$

Final Classification :-  $\text{Buyt} = 1$

~~Age = 30 & Fnew = 40~~

~~DS~~

Model TP FP FN TN

A 40 10

B 45 25 15 15

$M = 100$

$$\text{Accuracy} = \frac{TP + TN}{TP + FN + FP + TN}$$

$$\text{Precision} = \frac{TP}{TP + FP}$$

$$\text{Recall} = \frac{TP}{TP + FN}$$

for Model A

~~$$A_A = \frac{40 + 30}{100} = 0.7 = 70\%$$~~

~~$$P_A = 40 / (40 + 10) = 0.8 = 80\%$$~~

~~$$R_A = 40 / 60 = 0.67 = 67\%$$~~

for Model B

~~$$A_B = 45 + 15 / 100 = 0.6 = 60\%$$~~

~~$$P_B = 45 / (45 + 25) = 0.64 = 64\%$$~~

~~$$R_B = 45 / 45 + 15 = 0.75 = 75\%$$~~

(ii)

~~FN (false negative) i.e. actually positive but ~~was~~ predicted negative.~~

~~Model B will be choose~~

~~whereas FN are more costly~~

~~because it has high recall than model A.~~

(ii)

FP → False positive i.e. Actually negative but predicted positive  
Model A will be chosen when FP are more costly because of higher precision