

Problem #1

Part (a)

CART

Occupation	Gender	Age	Salary
Service	Female	45	L3
Service	Male	25	L1
Service	Male	33	L2
Management	Male	25	L3
Management	Female	35	L4
Management	Male	26	L3
Management	Female	45	L4
Sales	Female	40	L3
Sales	Male	30	L2
Staff	Female	50	L2
Staff	Male	25	L1

Target variable = Salary

Impurity Before Splitting

	salary
L1	2
L2	3
L3	4
L4	2
	11

$$1 - \sum_{i=1}^{c-1} p_i(t)^2$$

L1 L2 L3 L4
S M S ST

$$P = 1 - \sum_{i=0}^{c-1} P_i(t)^2$$

$$P = 1 - \left(\frac{2}{11}\right)^2 - \left(\frac{3}{11}\right)^2 - \left(\frac{4}{11}\right)^2 - \left(\frac{2}{11}\right)^2 = 0.7272$$

Gini Index for Occupation

Occupation	Service, Management	Sales, Staff
L1	1	1
L2	1	2
L3	3	1
L4	2	0
	7	4

$$Gini(N1) = 1 - \left(\frac{1}{7}\right)^2 - \left(\frac{1}{7}\right)^2 - \left(\frac{3}{7}\right)^2 - \left(\frac{2}{7}\right)^2 = 0.6938$$

$$Gini(N2) = 1 - \left(\frac{1}{4}\right)^2 - \left(\frac{2}{4}\right)^2 - \left(\frac{1}{4}\right)^2 - \left(\frac{0}{4}\right)^2 = 0.625$$

$$\text{Weighted } Gini = \sum_{i=1}^k \frac{n_i}{n} Gini(i)$$

$$M = \frac{7}{11}(0.6938) + \frac{4}{11}(0.625) = 0.6687$$

$$Gain = P - M = 0.7272 - 0.6687 = \boxed{0.0585}$$

(B)

	Service	Sales	Management
L1	1		staff
L2	2		1
L3	2		1
L4	0	2	
	5	6	

$$Gini(N1) = 1 - \left(\frac{1}{5}\right)^2 - \left(\frac{2}{5}\right)^2 - \left(\frac{2}{5}\right)^2 = 0.64$$

$$Gini(N2) = 1 - \left(\frac{1}{6}\right)^2 - \left(\frac{1}{6}\right)^2 - \left(\frac{2}{6}\right)^2 - \left(\frac{2}{6}\right)^2 = 0.7222$$

$$M = \frac{5}{11}(0.64) + \frac{6}{11}(0.7222) = 0.6848$$

$$Gain = P - M = 0.7272 - 0.6848 = \boxed{0.0424}$$

(C)

	Service	Sales	Management
L1	2		staff
L2	2		0
L3	1		1
L4	0	3	
	5	6	

$$Gini(N1) = 1 - \left(\frac{2}{5}\right)^2 - \left(\frac{2}{5}\right)^2 - \left(\frac{1}{5}\right)^2 = 0.64$$

$$Gini(N2) = 1 - \left(\frac{0}{6}\right)^2 - \left(\frac{1}{6}\right)^2 - \left(\frac{3}{6}\right)^2 - \left(\frac{2}{6}\right)^2 = 0.6111$$

$$M = \frac{5}{11}(0.64) + \frac{6}{11}(0.6111) = 0.6242$$

$$Gain = P - M = 0.7272 - 0.6242 = \boxed{0.103}$$

	N1 Sales	N2 Service, staff Management
occupation		
L1	0	2
L2	1	2
L3	1	3
L4	$\frac{0}{2}$	$\frac{2}{8}$

$$Gini(N_1) = 1 - \left(\frac{1}{2}\right)^2 - \left(\frac{1}{2}\right)^2 = 0.5$$

$$\begin{aligned} Gini(N_2) &= 1 - \left(\frac{2}{9}\right)^2 - \left(\frac{2}{9}\right)^2 \\ &\quad - \left(\frac{3}{9}\right)^2 - \left(\frac{2}{9}\right)^2 = 0.7407 \end{aligned}$$

$$\begin{aligned} M &= \frac{2}{11}(0.5) + \frac{9}{11}(0.7407) \\ &= 0.6969 \end{aligned}$$

$$\begin{aligned} Gain &= P - M = 0.7272 - 0.6969 \\ &= 0.0303 \end{aligned}$$

	N1 Staff	N2 Service, sales Management
occupation		
L1	1	1
L2	1	2
L3	0	4
L4	$\frac{0}{2}$	$\frac{2}{9}$

$$Gini(N_1) = 1 - \left(\frac{1}{2}\right)^2 - \left(\frac{1}{2}\right)^2 = 0.5$$

$$\begin{aligned} Gini(N_2) &= 1 - \left(\frac{1}{9}\right)^2 - \left(\frac{2}{9}\right)^2 - \left(\frac{4}{9}\right)^2 - \left(\frac{2}{9}\right)^2 \\ &= 0.6913 \end{aligned}$$

$$\begin{aligned} M &= \frac{2}{11}(0.5) + \frac{9}{11}(0.6913) \\ &= 0.6565 \end{aligned}$$

$$\begin{aligned} Gain &= P - M = 0.7272 - 0.6565 \\ &= 0.0707 \end{aligned}$$

	N1 Service	N2 Sales, staff Management
occupation		
L1	1	1
L2	1	2
L3	1	3
L4	$\frac{0}{3}$	$\frac{2}{8}$

$$\begin{aligned} Gini(N_1) &= 1 - \left(\frac{1}{3}\right)^2 - \left(\frac{1}{3}\right)^2 - \left(\frac{1}{3}\right)^2 \\ &= 0.666 \end{aligned}$$

$$\begin{aligned} Gini(N_2) &= 1 - \left(\frac{2}{8}\right)^2 - \left(\frac{1}{8}\right)^2 - \left(\frac{3}{8}\right)^2 \\ &\quad - \left(\frac{2}{8}\right)^2 \\ &= 0.7187 \end{aligned}$$

$$\begin{aligned} M &= \frac{3}{11}(0.666) + \frac{8}{11}(0.7187) \\ &= 0.7043 \end{aligned}$$

$$\begin{aligned} Gain &= P - M = 0.7272 - 0.7043 \\ &= 0.0229 \end{aligned}$$

	N1 Management	N2 Sales, staff Service
occupation		
L1	0	2
L2	0	3
L3	2	2
L4	$\frac{2}{4}$	$\frac{0}{7}$

$$Gini(N_1) = 1 - \left(\frac{2}{4}\right)^2 - \left(\frac{2}{4}\right)^2 = 0.5$$

$$\begin{aligned} Gini(N_2) &= 1 - \left(\frac{2}{7}\right)^2 - \left(\frac{3}{7}\right)^2 - \left(\frac{2}{7}\right)^2 \\ &= 0.6530 \end{aligned}$$

$$\begin{aligned} M &= \frac{4}{11}(0.5) + \frac{7}{11}(0.6530) \\ &= 0.5973 \end{aligned}$$

$$\begin{aligned} Gain &= P - M = 0.7272 - 0.5973 \\ &= 0.1299 \end{aligned}$$

(2)

Gini Index for Gender

	N_1 Male	N_2 Female
L1	2	0
L2	2	1
L3	2	2
L4	0 6	2 5

$$Gini(N1) = 1 - \left(\frac{2}{6}\right)^2 - \left(\frac{2}{6}\right)^2 - \left(\frac{2}{6}\right)^2 = 0.666$$

$$Gini(N2) = 1 - \left(\frac{1}{5}\right)^2 - \left(\frac{2}{5}\right)^2 - \left(\frac{2}{5}\right)^2 = 0.64$$

$$M = \frac{6}{11}(0.666) + \frac{5}{11}(0.64) = 0.6541$$

$$Gain = P - M = 0.7272 - 0.6541 = 0.0731$$

Gini Index for Age

Salary levels	L1	L1	L3	L3	L2	L2	L4	L3	L3	L4	L2
Age											
→	25	25	25	26	30	33	35	40	45	45	50
	24	25	25	25	28	31	34	37	42	45	47
	<=	>	<=	>	<=	>	<=	>	<=	>	<=
L1	0	2	2	0	2	0	2	0	2	0	2
L2	0	3	0	3	0	3	1	2	2	1	2
L3	0	4	1	3	1	3	2	2	2	2	3
L4	0	2	0	2	0	2	0	2	1	1	1
Gini	0.7272	0.5983	0.5983	0.5983	0.5973	0.6541	0.6541	0.7077	0.7043	0.6545	0.6545
											0.7272

	N_1 \leq	N_2 $>$
L1	0	2
L2	0	3
L3	0	4
L4	0	2

$$Gini(N1) = 1 - 0 - 0 - 0 = 0 = 1$$

$$Gini(N2) = 1 - \left(\frac{2}{11}\right)^2 - \left(\frac{3}{11}\right)^2 - \left(\frac{4}{11}\right)^2 - \left(\frac{2}{11}\right)^2 = 0.7272$$

$$\text{weighted}_M = \frac{0}{11}(1) + \frac{11}{11}(0.7272) = 0.7272$$

(ii) For 25	
L1	2
L2	0
L3	1
L4	0

$$Gini(N1) = 1 - \left(\frac{2}{3}\right)^2 - \left(\frac{1}{3}\right)^2 = 0.444$$

$$Gini(N2) = 1 - \left(\frac{3}{8}\right)^2 - \left(\frac{3}{8}\right)^2 - \left(\frac{2}{8}\right)^2$$

$$= 0.6562$$

$$\text{weighted} = \frac{3}{11}(0.444) + \frac{8}{11}(0.6562)$$

$$= 0.5983$$

<u>For 28</u>		
	\leq	$>$
L1	2	0
L2	0	3
L3	2	2
L4	0	2
	4	7

$$Gini(N1) = 1 - \left(\frac{2}{4}\right)^2 - \left(\frac{2}{4}\right)^2 = 0.5$$

$$Gini(N2) = 1 - \left(\frac{3}{7}\right)^2 - \left(\frac{2}{7}\right)^2 - \left(\frac{2}{7}\right)^2 = 0.6530$$

$$M = \frac{4}{11}(0.5) + \frac{7}{11}(0.6530) \\ = 0.5973$$

<u>For 31</u>		
	\leq	$>$
L1	2	0
L2	1	2
L3	2	2
L4	0	2
	5	6

$$Gini(N1) = 1 - \left(\frac{2}{5}\right)^2 - \left(\frac{1}{5}\right)^2 - \left(\frac{2}{5}\right)^2 \\ = 0.64$$

$$Gini(N2) = 1 - \left(\frac{2}{6}\right)^2 - \left(\frac{2}{6}\right)^2 - \left(\frac{2}{6}\right)^2 \\ = 0.666$$

$$M = \frac{5}{11}(0.64) + \frac{6}{11}(0.666) \\ = 0.6541$$

(v) For 34

	\leq	$>$
L1	2	0
L2	2	1
L3	2	2
L4	0	2
	6	5

$$Gini(N1) = 1 - \left(\frac{2}{6}\right)^2 - \left(\frac{2}{6}\right)^2 - \left(\frac{2}{6}\right)^2 = 0.666$$

$$Gini(N2) = 1 - \left(\frac{1}{5}\right)^2 - \left(\frac{2}{5}\right)^2 - \left(\frac{2}{5}\right)^2 = 0.64$$

$$M = \frac{6}{11}(0.666) + \frac{5}{11}(0.64) \\ = 0.6541$$

(vi) For 37

	\leq	$>$
L1	2	0
L2	2	1
L3	2	2
L4	1	1
	7	4

$$Gini(N1) = 1 - \left(\frac{2}{7}\right)^2 - \left(\frac{2}{7}\right)^2 - \left(\frac{2}{7}\right)^2 - \left(\frac{1}{7}\right)^2 \\ = 0.7551$$

$$Gini(N2) = 1 - \left(\frac{1}{4}\right)^2 - \left(\frac{2}{4}\right)^2 - \left(\frac{1}{4}\right)^2 \\ = 0.625$$

$$M = \frac{7}{11}(0.7551) + \frac{4}{11}(0.625) \\ = 0.7077$$

(vii) For 42

	\leq	$>$
L1	2	0
L2	2	1
L3	3	1
L4	1	1
	8	3

$$Gini(N1) = 1 - \left(\frac{2}{8}\right)^2 - \left(\frac{2}{8}\right)^2 - \left(\frac{3}{8}\right)^2 - \left(\frac{1}{8}\right)^2 = 0.7187$$

$$Gini(N2) = 1 - \left(\frac{1}{3}\right)^2 - \left(\frac{1}{3}\right)^2 - \left(\frac{1}{3}\right)^2 = 0.666$$

$$M = \frac{8}{11}(0.7187) + \frac{3}{11}(0.666) \\ = 0.7043$$

(viii) For 45

	\leq	$>$
L1	2	0
L2	2	1
L3	4	0
L4	2	0
	10	1

$$Gini(N1) = 1 - \left(\frac{2}{10}\right)^2 - \left(\frac{2}{10}\right)^2 - \left(\frac{4}{10}\right)^2 - \left(\frac{2}{10}\right)^2 \\ = 0.72$$

$$Gini(N2) = 1 - \left(\frac{1}{1}\right)^2 = 0$$

$$M = \frac{10}{11}(0.72) + \frac{1}{11}(0) \quad (3) \\ = 0.6545$$

(ix) For 47

	<u><u>L</u></u>	<u><u>=</u></u>	<u><u>></u></u>
L1	2	0	
L2	2	1	
L3	4	0	
L4	2	0	
	10	1	

(x) For 50

	<u><u>L</u></u>	<u><u>=</u></u>	<u><u>></u></u>
L1	2	0	
L2	3	0	
L3	4	0	
L4	2	0	
	11	0	

$$Gini(N1) = 1 - \left(\frac{2}{10}\right)^2 - \left(\frac{2}{10}\right)^2 - \left(\frac{4}{10}\right)^2 - \left(\frac{2}{10}\right)^2 = 0.72$$

$$Gini(N2) = 1 - \left(\frac{1}{1}\right)^2 = 0$$

$$M = \frac{10}{11}(0.72) + \frac{1}{11}(0) \\ = 0.6545$$

$$Gini(N1) = 1 - \left(\frac{2}{11}\right)^2 - \left(\frac{3}{11}\right)^2 - \left(\frac{4}{11}\right)^2 - \left(\frac{2}{11}\right)^2 \\ = 0.7272$$

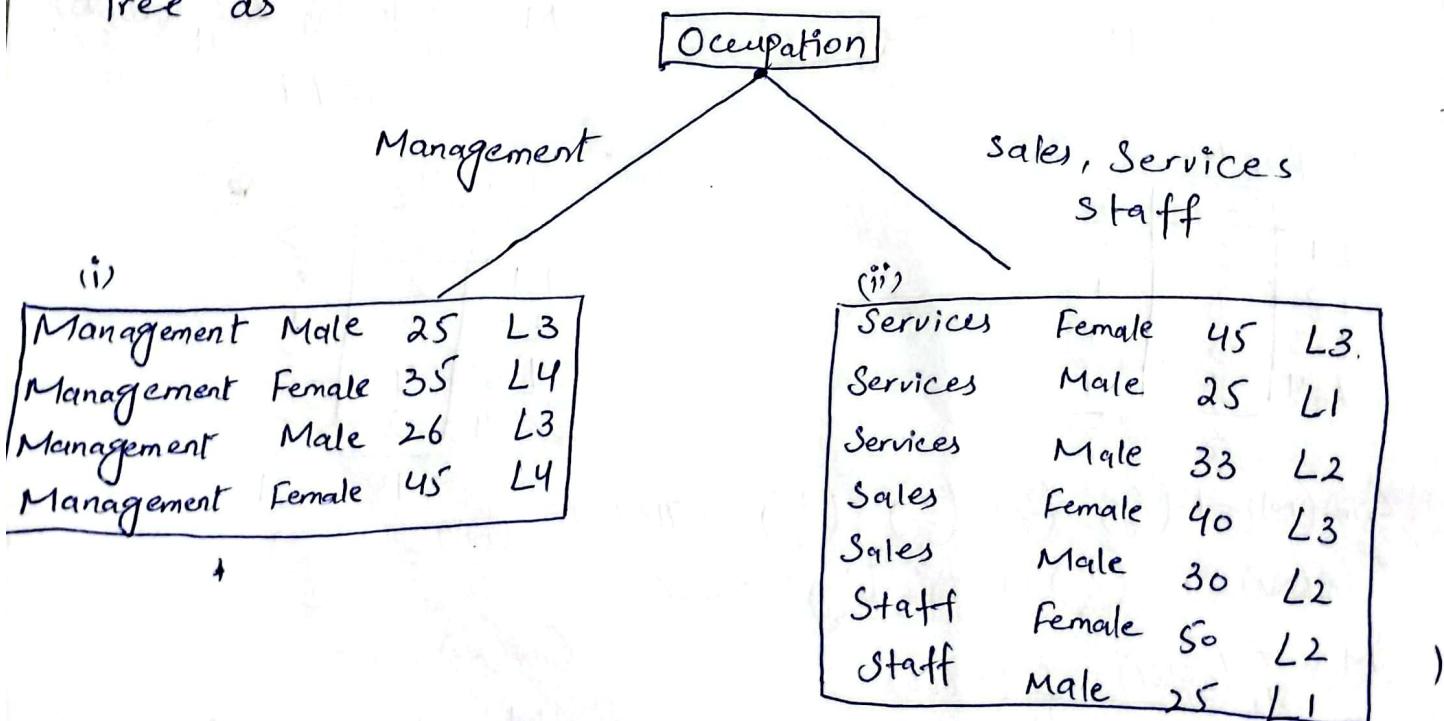
$$Gini(N2) = 1 - 0 = 1$$

$$M = \frac{11}{11}(0.7272) + \frac{0}{11}(1) \\ = 0.7272$$

As Occupation with Management / sales, services, staff split (binary) of And Age with threshold of 28 has same smallest gini index.

So we can pick anyone of these.
So, we selected Occupation as Root node.

Tree as



from Table (ii)
Parent/secondary

L1	2	$G_{ini} = 1 - \left(\frac{2}{7}\right)^2 - \left(\frac{3}{7}\right)^2 - \left(\frac{2}{7}\right)^2$
L2	3	
L3	2	
L4	0	$P = 0.6530$
		$\frac{7}{7}$

Gini Index for Gender

	Female	Male
L1	0	2
L2	1	2
L3	2	0
L4	0	0
	$\frac{3}{7}$	$\frac{4}{7}$

$$G_{ini}(N1) = 1 - \left(\frac{1}{3}\right)^2 - \left(\frac{2}{3}\right)^2$$

$$= 1 - 0.111 - 0.444$$

$$= 0.444$$

$$G_{ini}(N2) = 1 - \left(\frac{2}{4}\right)^2 - \left(\frac{2}{4}\right)^2$$

$$= 1 - 0.25 - 0.25$$

$$= 0.5$$

weighted Gini

$$M = \frac{3}{7} (0.444) + \frac{4}{7} (0.5)$$

$$= 0.476$$

Gini Index for Age

L1	L1	L2	L2	L3	L3	L2	
Age							
25	25	30	33	40	45	50	
22	25	27	31	36	42	47	50
1	0	2	0	2	0	2	0
2	0	3	0	3	1	2	1
3	0	2	0	2	0	2	0
4	0	0	0	0	0	0	0
Min	0.6530	0.3428	0.3428	0.476	0.476	0.6	0.5708
Max	0.6530	0.3428	0.3428	0.476	0.476	0.6	0.6530

For 22

$$\begin{aligned} \text{Gini}(N1) &= 1 - 0 - 0 - 0 \\ &= 1 \end{aligned}$$

$$\begin{aligned} \text{Gini}(N2) &= 1 - \left(\frac{2}{7}\right)^2 - \left(\frac{3}{7}\right)^2 - \left(\frac{2}{7}\right)^2 \\ &\approx 0.6530 \end{aligned}$$

$$\begin{aligned} M &= \frac{2}{7}(1) + \frac{5}{7}(0.65) \\ &= 0.6530 \end{aligned}$$

For 25

$$\begin{aligned} \text{Gini}(N1) &= 1 - \left(\frac{2}{2}\right)^2 = 0 \\ \text{Gini}(N2) &= 1 - \left(\frac{3}{5}\right)^2 - \left(\frac{2}{5}\right)^2 = 0.48 \\ M &= \frac{2}{7}(0) + \frac{5}{7}(0.48) \\ &= 0.3428 \end{aligned}$$

For 27

$$\text{Gini}(N1) = 1 - \left(\frac{2}{2}\right)^2 = 0$$

$$\text{Gini}(N2) = 1 - \left(\frac{3}{5}\right)^2 - \left(\frac{2}{5}\right)^2 = 0.48$$

$$M = 0.3428$$

For 31

$$\begin{aligned} \text{Gini}(N1) &= 1 - \left(\frac{2}{3}\right)^2 - \left(\frac{1}{3}\right)^2 \\ &= 0.44 \end{aligned}$$

$$\begin{aligned} \text{Gini}(N2) &= 1 - \left(\frac{2}{4}\right)^2 - \left(\frac{2}{4}\right)^2 \\ &= 0.5 \end{aligned}$$

$$\begin{aligned} M &= \frac{3}{7}(0.44) + \frac{4}{7}(0.5) \\ &= 0.476 \end{aligned}$$

For 36

$$M = 0.476$$

For 42

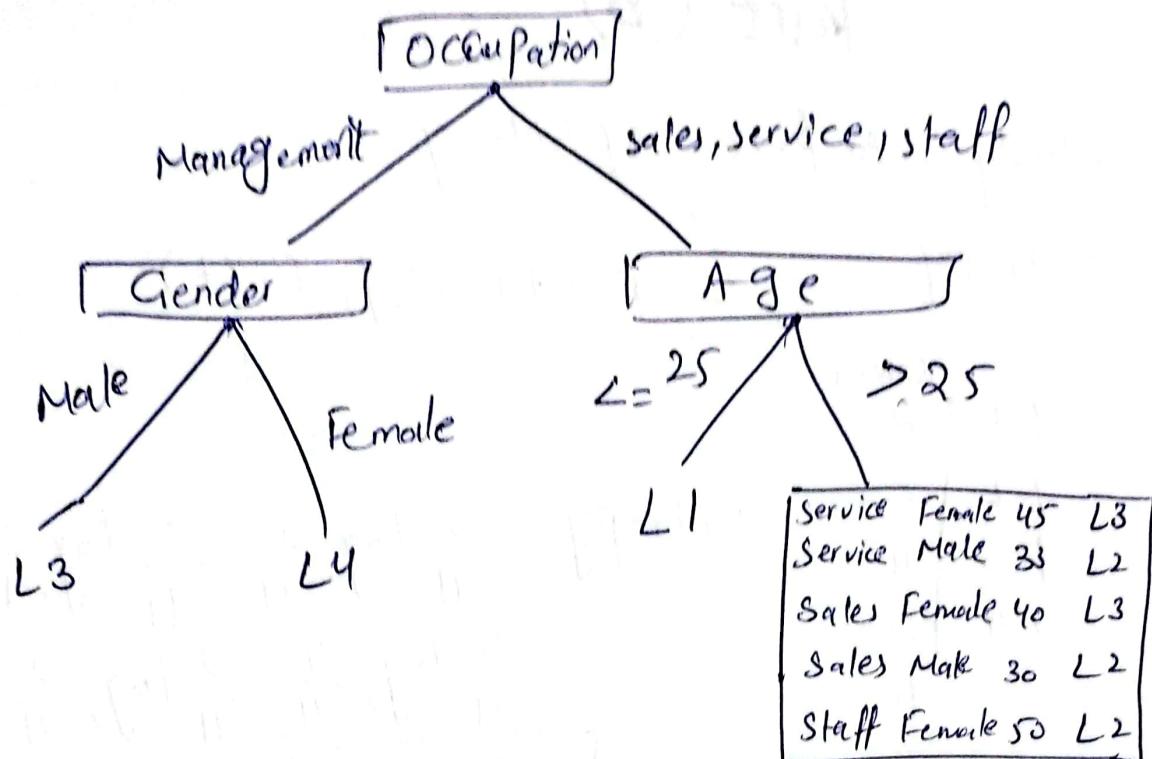
$$\begin{aligned} \text{Gini}(N1) &= 0.64 \\ \text{Gini}(N2) &= 0.5 \\ M &= 0.6 \end{aligned}$$

For 47

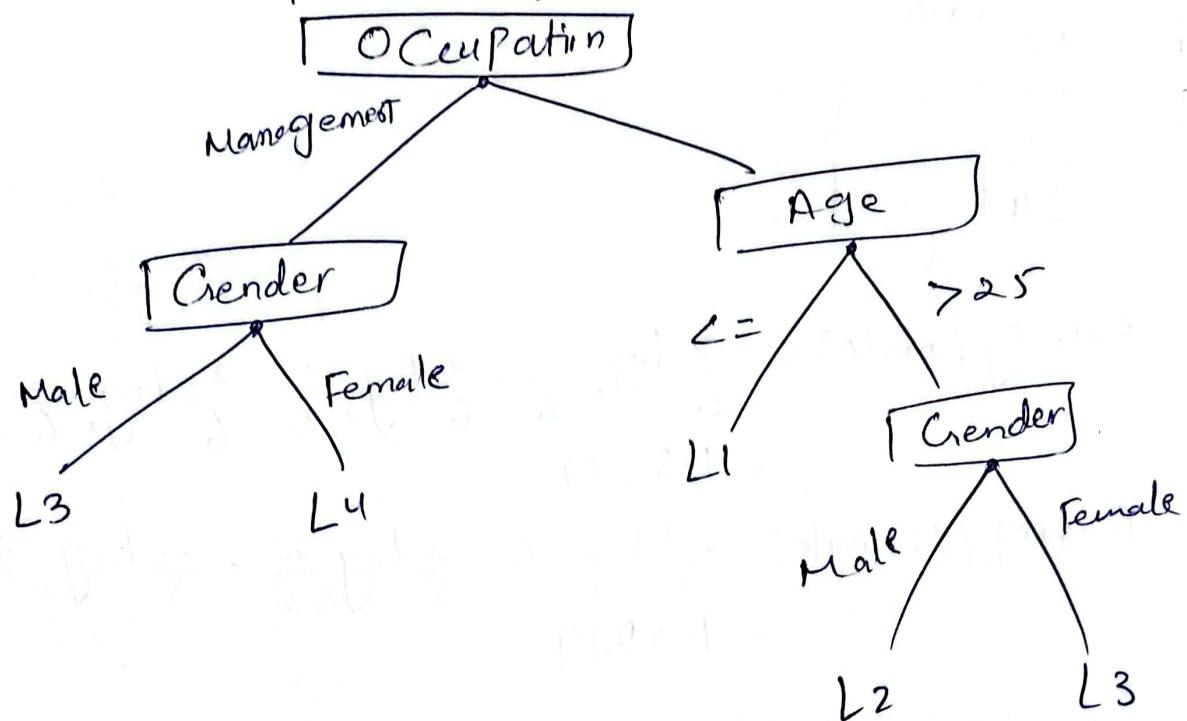
$$\begin{aligned} \text{Gini}(N1) &= 0.66 \\ \text{Gini}(N2) &= 0 \\ M &= 0.5708 \end{aligned}$$

For 50

$$\begin{aligned} \text{Gini}(N1) &= 0.6530 \\ \text{Gini}(N2) &= 0 \\ M &= 0.6530 \end{aligned}$$



As occupation, Age already been covered
 so we can split on gender only - further.



One Record is misclassified with this tree.

Part (b) ID3

Entropy before splitting.

$$\text{Entropy} = - \sum_{i=1}^{c-1} p_i(t) \log_2 p_i(t)$$

Salary	
L1	2
L2	3
L3	4
L4	2
	11

$$\begin{aligned} \text{Entropy} &= -\frac{2}{11} \log_2 \frac{2}{11} - \frac{3}{11} \log_2 \frac{3}{11} \\ &\quad - \frac{4}{11} \log_2 \frac{4}{11} - \frac{2}{11} \log_2 \frac{2}{11} \\ &= \boxed{1.9362} \end{aligned}$$

Entropy for Gender

	Male	Female
L1	2	0
L2	2	1
L3	2	2
L4	0	2
	6	5

$$\begin{aligned} \text{Entropy(Male)} &= -\frac{2}{6} \log_2 \frac{2}{6} - \frac{2}{6} \log_2 \frac{2}{6} - \frac{2}{6} \log_2 \frac{2}{6} \\ &= 1.5849 \end{aligned}$$

$$\begin{aligned} \text{Entropy(Female)} &= -\frac{1}{5} \log_2 \frac{1}{5} - \frac{2}{5} \log_2 \frac{2}{5} - \frac{2}{5} \log_2 \frac{2}{5} \\ &= 1.5219 \end{aligned}$$

Weighted Entropy

$$\begin{aligned} &= \frac{6}{11} (1.5849) + \frac{5}{11} (1.5219) \\ &= \boxed{1.5562} \end{aligned}$$

28 Copy

Q#

Entropy for Age

age ≤ 25 Total = 3

age > 25 Total = 8

	≤ 25	> 25
L1	2	0
L2	0	3
L3	1	3
L4	0	2
	$\frac{3}{11}$	$\frac{8}{11}$

$$\text{Entropy}(N1) = -\frac{2}{3} \log_2 \frac{2}{3} - \frac{1}{3} \log_2 \frac{1}{3} = 0.9182$$

$$\begin{aligned} \text{Entropy}(N2) &= -\frac{3}{8} \log_2 \frac{3}{8} - \frac{3}{8} \log_2 \frac{3}{8} - \frac{2}{8} \log_2 \frac{3}{8} \\ &= 1.5612 \end{aligned}$$

$$\begin{aligned} \text{Weighted Entropy} &= \frac{3}{11} (0.9182) + \frac{8}{11} (1.5612) \\ &= 1.3858 \end{aligned}$$

Entropy for Occupation

②

	service staff Management	sales
L1	2	0
L2	2	1
L3	3	1
L4	2	0
	$\frac{9}{11}$	$\frac{2}{11}$

$$\begin{aligned} \text{Entropy}(N1) &= -\frac{2}{9} \log_2 \frac{2}{9} - \frac{2}{9} \log_2 \frac{2}{9} - \frac{3}{9} \log_2 \frac{3}{9} - \frac{2}{9} \log_2 \frac{2}{9} \\ &= 1.9749 \end{aligned}$$

$$\text{Entropy}(N2) = -\frac{1}{2} \log_2 \frac{1}{2} - \frac{1}{2} \log_2 \frac{1}{2} = 1$$

$$\begin{aligned} \text{Weighted Entropy} &= \frac{9}{11} (1.9749) + \frac{2}{11} (1) \\ &= 1.7976 \end{aligned}$$

	staff	sales, service management
L1	1	1
L2	1	2
L3	0	4
L4	0	2
	$\frac{2}{2}$	$\frac{9}{9}$

$$\text{Entropy}(N1) = -\frac{1}{2} \log_2 \frac{1}{2} - \frac{1}{2} \log_2 \frac{1}{2} = 1$$

$$\text{Entropy}(N2) = -\frac{1}{9} \log_2 \frac{1}{9} - \frac{2}{9} \log_2 \frac{2}{9} - \frac{4}{9} \log_2 \frac{4}{9} = 1.8365$$

$$\text{Weighted Entropy} = \frac{2}{11}(1) + \frac{9}{11}(1.8365) = \boxed{1.6844}$$

	service	sales, staff management
L1	1	1
L2	1	2
L3	1	3
L4	0	2
	$\frac{0}{3}$	$\frac{8}{8}$

$$\text{Entropy}(N1) = -\frac{1}{3} \log_2 \frac{1}{3} - \frac{1}{3} \log_2 \frac{1}{3} - \frac{1}{3} \log_2 \frac{1}{3} = 1.5849$$

$$\text{Entropy}(N2) = -\frac{1}{8} \log_2 \frac{1}{8} - \frac{2}{8} \log_2 \frac{2}{8} - \frac{3}{8} \log_2 \frac{3}{8} - \frac{2}{8} \log_2 \frac{2}{8} = 1.9056$$

$$\text{Weighted Entropy} = \frac{3}{11}(1.5849) + \frac{8}{11}(1.9056) = \boxed{1.8181}$$

	Management	sales, staff services
L1	0	2
L2	0	3
L3	2	2
L4	2	0
	$\frac{4}{4}$	$\frac{7}{7}$

$$\text{Entropy}(N1) = -\frac{2}{4} \log_2 \frac{2}{4} - \frac{2}{4} \log_2 \frac{2}{4} = 1$$

$$\text{Entropy}(N2) = -\frac{2}{7} \log_2 \frac{2}{7} - \frac{3}{7} \log_2 \frac{3}{7} - \frac{2}{7} \log_2 \frac{2}{7} = 1.5566$$

$$\text{Weighted} = \frac{4}{11}(1) + \frac{7}{11}(1.5566) = \boxed{1.3542}$$

	service, management	sales, staff
L1	1	1
L2	1	2
L3	3	1
L4	2	0
	$\frac{7}{7}$	$\frac{4}{4}$

$$\text{Entropy}(N1) = 1.8423$$

$$\text{Entropy}(N2) = 1.5$$

$$\text{Weighted} = \boxed{1.7178}$$

	service, sales	staff, management
L1	1	1
L2	2	1
L3	2	2
L4	0	2
	$\frac{5}{5}$	$\frac{6}{6}$

$$\text{Entropy}(N1) = 1.5219$$

$$\text{Entropy}(N2) = 1.9182$$

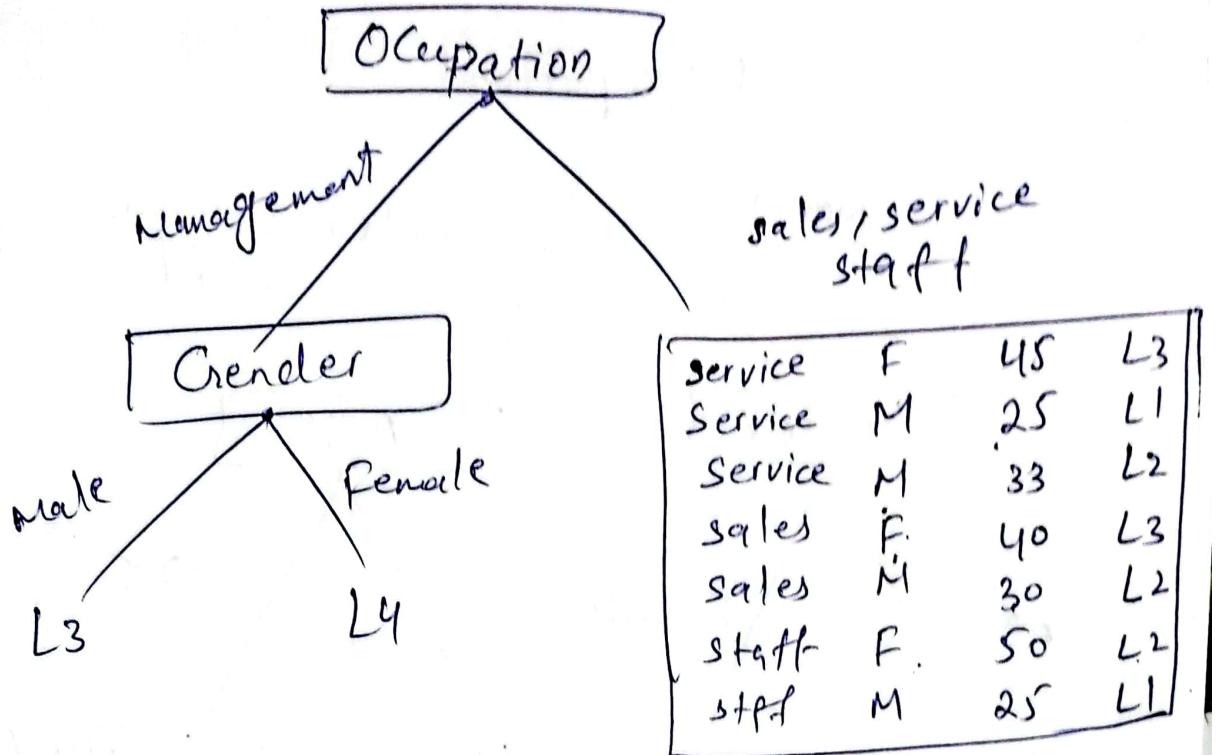
$$\text{Weighted} = \boxed{1.7380}$$

	staff, service	Management sales
L1	2	0
L2	2	1
L3	1	3
L4	0	2
	$\frac{5}{5}$	$\frac{6}{6}$

$$\text{Entropy}(N1) = 1.5219$$

$$\text{Entropy}(N2) = 1.4591$$

$$\text{Weighted} = \boxed{1.4876}$$



Now for table.

Entropy for Gender

	Male	Female
L1	2	0
L2	2	1
L3	0	2
L4	0	0
	4	3

$$\text{Entropy}(\text{male}) = 1$$

$$\text{Entropy}(\text{female}) = 0.918$$

$$\begin{aligned}\text{weighted Entropy} &= \frac{4}{7}(1) + \frac{3}{7}(0.918) \\ &= 0.5714 + 0.393 \\ &= 0.9648\end{aligned}$$

Entropy for Age

Age ≥ 25

Total = 2

Age > 25

Total = 5

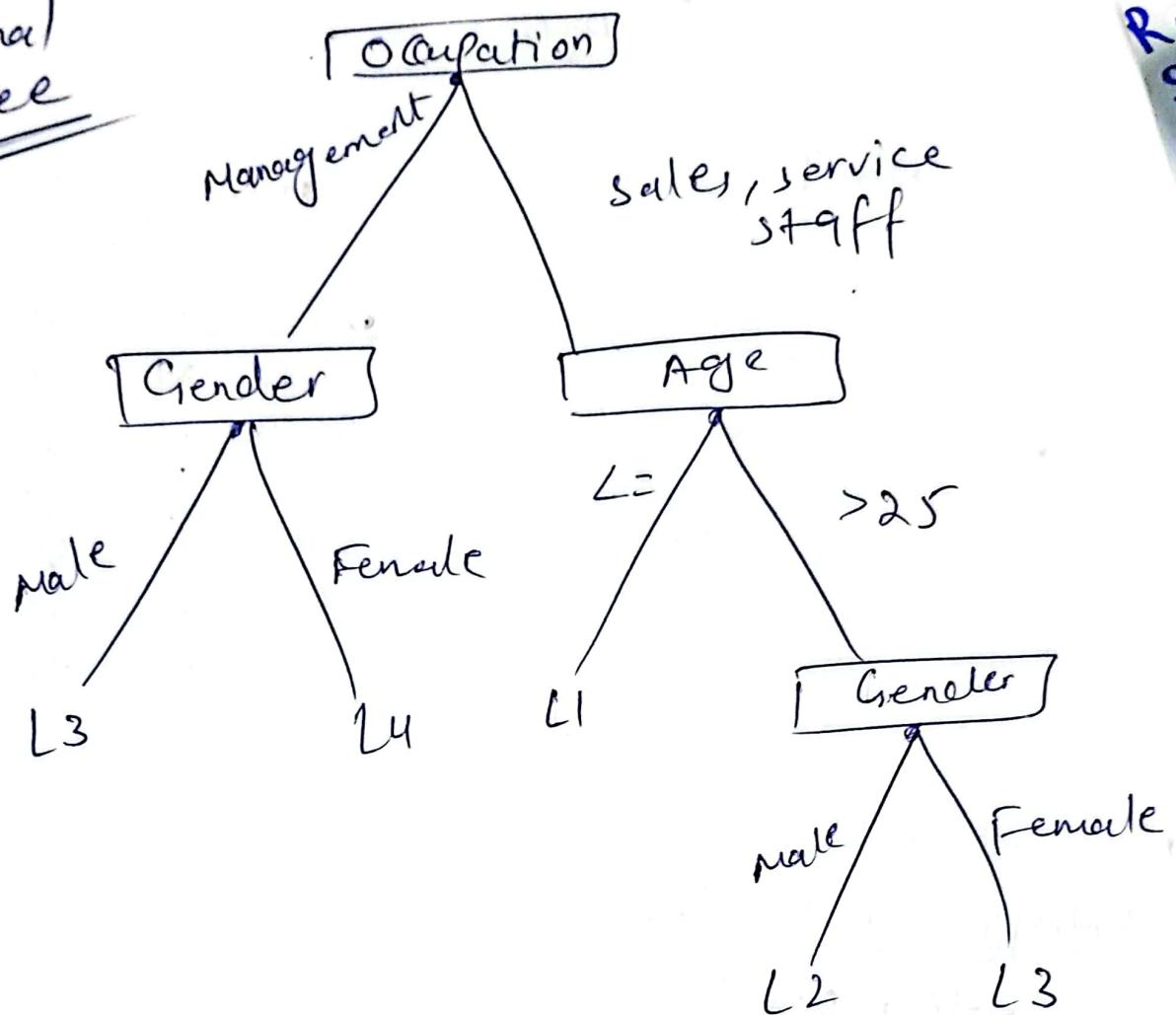
	≤ 25	> 25
L1	1	0
L2	0	3
L3	0	2
L4	0	0
	1	5

$$\text{Entropy}(N1) = 0$$

$$\text{Entropy}(N2) = 0.9709$$

$$\text{weighted} = 0.6935$$

Final
Tree



One record is misclassified in this tree as well.

- (c) Both trees are same.
(d) For Part(a) and (b)

- ① R₁ (Occupation = Management) \wedge (Gender = Male) \rightarrow L₃
- ② R₂ (Occupation = Management) \wedge (Gender = Female) \rightarrow L₄
- ③ R₃ (Occupation = sales/service/staff) \wedge (Age \leq 25) \rightarrow L₁
- ④ R₄ (Occupation = sales/service/staff) \wedge (Gender = Male) \wedge (Age $>$ 25) \rightarrow L₂
- ⑤ R₅ (Occupation = sales/service/staff) \wedge (Age $>$ 25) \wedge (Gender = Female) \rightarrow L₃.

problem #2

solution

- (a) No, the rules are not mutually exclusive as more than one rules are covering its diff combinations.
- (b) Yes , the rule set is exhaustive as it is covering the every combination of attributes.
- (c) Rule (3,4) can be put at top to get maximum coverage by putting toughest requirement first as it involves many tests.
- (d) No , we don't need any default class as rule set is exhaustive.

Problem #3

Solution :-

(ii) ~~R1~~ $\rightarrow +$ (covers 12 positive and 4 negative examples)

$$\underline{R1} \text{ Rule's Accuracy} = \frac{12}{12+4} = \frac{12}{16} = 75\%$$

$$\underline{R2} \text{ Rule's Accuracy} = \frac{45}{45+15} = \frac{45}{60} = 75\%$$

$$\underline{R3} \text{ Rule's Accuracy} = \frac{115}{115+100} = \frac{115}{215} = 53\%$$

$R1$ & $R2$ are the best rules according to Accuracy and $R3$ is the worst rule.

$$P_0 = 150 \quad N_0 = 320$$

$R1$ Foil's Information gain.

$$= P_1 * \left(\log_2 \frac{P_1}{P_1+N_1} - \log_2 \frac{P_0}{P_0+N_0} \right)$$

$$P_1 = 12 \quad N_1 = 4$$

$$= 12 * \left(\log_2 \frac{12}{12+4} - \log_2 \frac{150}{150+320} \right)$$

$$= 14\%.$$

$R2$

$$P_1 = 45 \quad N_1 = 15$$

$$= 45 * \left(\log_2 \frac{45}{45+15} - \log_2 \frac{150}{150+320} \right)$$

$$= 55\%.$$

$R3$

$$P_1 = 115 \quad N_1 = 100$$

$$= 115 * \left(\log_2 \frac{115}{115+100} - \log_2 \frac{150}{150+320} \right) = 85\%$$

$R1$ is worst and $R3$ is best rule as per Foil info

Problem #4

Solution:

Nearth neighbor

unweighted Majority Voting

$$X_1 = 4.0$$

X	0.5	3.0	4.5	4.6	4.9	5.2	5.3	5.5	7.0	9.5	label
Y	-	-	+	+	+	-	-	+	-	-	
$X_1 = 4.0$	3.5	1.0	0.5	0.6	0.9	1.2	1.3	1.5	3	5.5	
1			+								
3			+	+	+						+
5		-	+	+	+	-					+
9	-	-	+	+	+	-	-	+	-		+
											-

$$X_2 = 7.5$$

X	0.5	3.0	4.5	4.6	4.9	5.2	5.3	5.5	7.0	9.5	label
Y	-	-	+	+	+	-	-	+	-	-	
$X_2 = 7.5$	7	4.5	3	2.9	2.6	2.3	2.2	2	0.5	2	
1									-		-
3								+	-	-	-
5							-	-	+	-	-
9	-	+	+	+	-	-	-	+	-	-	-

Instance Weighted Majority Voting

$x_1 = 4.0$

$$\text{Distance weight} = \frac{1}{d^2}$$

X	0.5	3.0	4.5	4.6	4.9	5.2	5.3	5.5	7.0	9.5	Label
y	-	-	+	+	+	-	-	+	-	-	
$x_1 = 4$	3.5	1	0.5	0.6	0.9	1.2	1.3	1.5	3	5.5	
Distance weight	0.08	1	4	2.77	1.23	0.69	0.59	0.44	0.11	0.33	
1			+								+
3			+	+	+						+
5		-	+	+	+	-					+
9	-	-	+	+	+	-	-	+	-		+

For 5 $(0.5)(4) + (0.6)(2.77) + (0.9)(1.23)$ OR $(1.2)(0.69) + (1)(1)$
 $[4.769 \quad \text{OR} \quad 1.828]$ As $4.76 > 1.828$

So "+" would be the label.

For 9 $(0.5)(4) + (0.6)(2.77) + (0.9)(1.23) + (1.5)(0.44)$ OR
 $((3.5)(0.08) + (1)(1) + (1.2)(0.69) + (1.3)(0.59) + (3)(0.11)) = [5.429 \quad \text{OR} \quad 3.205]$
As $(5.429 > 3.205)$ so "+" is label

$x_2 = 7.5$

X	0.5	3.0	4.5	4.6	4.9	5.2	5.3	5.5	7.0	9.5	Label
y	-	-	+	+	+	-	-	+	-	-	
$x_2 = 7.5$	7	4.5	3	2.9	2.6	2.3	2.2	2	0.5	2	
Distance weight	0.020	0.049	0.11	0.118	0.14	0.18	0.20	0.25	4	0.25	
1									-		-
3								+	-	-	-
5						-	-	+	-	-	-
9	-	+	+	+	-	-	-	+	-	-	-

For 3

$$\left[(2)(0.25) \text{ OR } (0.5)(1) + (2)(0.25) \right] = \begin{cases} 0.5 \text{ OR } 2.5 \\ (+) \quad (-) \end{cases}$$

As $2.5 > 0.5$ so " $-$ " is the label.
 $(-) > (+)$

Na
Re
s

For 5

$$\left[(2)(0.25) \text{ OR } (0.5)(1) + (2)(0.25) + (2.2)(0.20) + (2.3)(0.18) \right]$$

$$\left[0.5 \text{ OR } 3.354 \right] \text{ As } 3.354 > 0.5 \text{ so } "-" \text{ is label.}$$

For 9

$$\left[(2)(0.25) + (2.6)(0.14) + (2.9)(0.11) + (3)(0.11) \text{ OR } (0.5)(1) + (2)(0.25) + (2.2)(0.20) + (2.3)(0.18) + (4.5)(0.04) \right]$$

$$= \begin{cases} 1.513 \text{ OR } 3.534 \\ (+) \quad (-) \end{cases} \text{ As } 3.534 > 1.513 \text{ so } (-) > (+)$$

" $-$ " is the label.

Problem # 5

Solution:-

$$(a) R_1 \text{ Accuracy} = \frac{12}{12+3} = \frac{12}{15} = 0.8 = 80\%$$

$$R_3 \text{ Accuracy} = \frac{8}{8+4} = \frac{8}{12} = 0.66 = 66\%$$

$$R_2 \text{ Accuracy} = \frac{7}{7+3} = \frac{7}{10} = 0.7 = 70\%$$

R_1 Best and R_3 worst.

$$R_1 \text{ Accuracy} = \frac{12}{12+3} = \frac{12}{15} = 80\%$$

$$R_2 \text{ Accuracy} = \frac{7}{7+3} = \frac{7}{10} = 70\%$$

$$R_3 \text{ Accuracy} = \frac{6}{6+4} = \frac{6}{10} = 60\%$$

R_1 is best and R_3 is Worst.

(C) $R_1 \text{ Accuracy} = \frac{12}{12+3} = \frac{12}{15} = 80\%$

$$R_3 \text{ Accuracy} = \frac{6}{6+2} = \frac{6}{8} = 75\%$$

$$R_2 \text{ Accuracy} = \frac{7}{7+3} = \frac{7}{10} = 70\%$$

R_1 is best and R_2 is Worst.

Problem #6

Solution:

(i) For $x-y$ dimensions.

(a) sorted for x

10	30
50	90
60	70
65	50
70	60

(b) sorted for y

10	30	50	60	70	90	95	25
65	80	85	90	95	75	85	65
50	70	60	50	90	80	75	90
75	25	50	60	70	85	90	65

Median = 70

75	25
80	85
90	65
95	75

⑥ sorted for x

$$\begin{bmatrix} 10 & 30 \\ 65 & 50 \end{bmatrix}$$

Median = 37.5

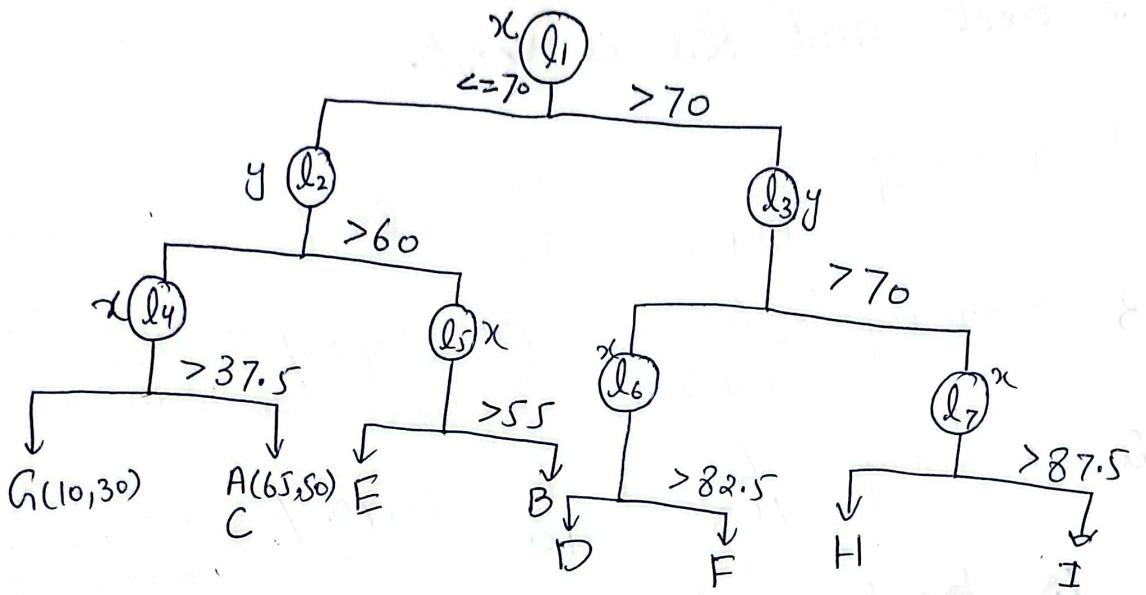
$$\begin{bmatrix} 50 & 90 \\ 60 & 70 \end{bmatrix}$$

Median = 55

$$\begin{bmatrix} 75 & 25 \\ 90 & 65 \end{bmatrix} \quad \begin{bmatrix} 80 & 85 \\ 95 & 75 \end{bmatrix}$$

Median = 82.5 Median = 87.5

KD-Tree.



$Z = (55, 60)$ is close to Point A(65, 50) traversing tree as $l_1 \rightarrow l_2 \rightarrow l_4 \rightarrow A$.

For y - x dimensions

⑦ sorted for y

$$\begin{bmatrix} 75 & 25 \\ 10 & 30 \\ 65 & 50 \\ 70 & 60 \\ 90 & 65 \end{bmatrix}$$

Median = 65

⑧ sorted for x

$$\text{Median} = \begin{bmatrix} 10 & 30 \\ 65 & 50 \\ 70 & 60 \\ 75 & 25 \end{bmatrix}$$

Median = 70

sorted for y

$$\begin{bmatrix} 10 & 30 \\ 65 & 50 \end{bmatrix} \quad \begin{bmatrix} 75 & 25 \\ 70 & 60 \end{bmatrix}$$

Median = 40 Median = 42.5

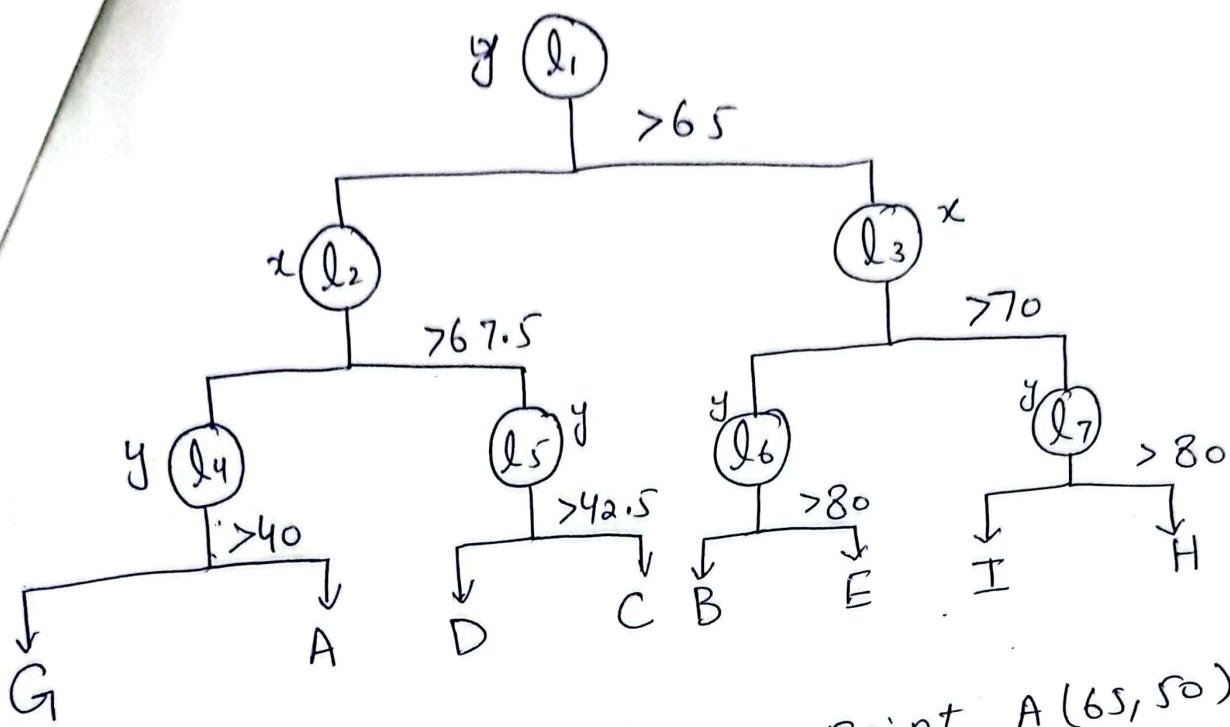
$$\begin{bmatrix} 60 & 70 \\ 50 & 90 \end{bmatrix}$$

Median = 80

$$\begin{bmatrix} 95 & 75 \\ 80 & 85 \end{bmatrix}$$

Median = 80

KD-Tree



$Z = (55, 60)$ is closest to point $A(65, 50)$ traversing tree as $l_1 \rightarrow l_2 \rightarrow l_4 \rightarrow A$.

There is no difference in the performance of both trees. Both trees classify new point in same manner and still find the same nearest neighbor for point.