



DS503 Machine Learning for Data Science

Assignment 02

Name: Muhammad Qasim
Roll No: 19k-1612

Refund	Status	Tax Income	Cheat				
		Mon	Tue	Wed	Thu	Fri	Sat
Yes	Single	115	125	135	145	155	NO
Yes	Single	60	70	80	90	100	NO
Yes	Single (married)	40	50	60	70	80	Yes
Yes	Married	130	140	150	160	170	Yes
Yes	Married	80	90	100	110	120	NO
Yes	Divorced	130	140	150	160	170	Yes
Yes	Divorced	80	90	100	110	120	NO
No	Single	25	35	45	55	65	NO
No	Single	65	75	85	95	105	Yes
No	Married	125	135	145	155	165	Yes
No	Married	80	90	100	110	120	NO
No	Divorced	90	100	110	120	130	Yes
No	Divorced	75	85	95	105	115	Yes

C4.S:-

Discrete Tax Income

$$\text{range} = 40 - 130$$

Buckets

low = 40 - 70 (Tax Income)

medium = 70 - 100 (Tax Income)

high = 100 - 130 (Tax Income)

Refund	Status	Tax Income	Cheat
Yes	Single	high	NO
Yes	Single	low	NO
Yes	Single	low	Yes
Yes	Married	high	Yes
Yes	Married	medium	NO
Yes	Divorced	high	Yes

Yes	Divorced	medium	No
NO	Single	medium	NO
NO	Single	low	Yes
NO	Married	high	Yes
NO	Married	medium	No
NO	Divorced	medium	Yes
NO	Divorced	medium	Yes

A) Entropy (~~CART~~):-

$$E(S) = \sum_{i=1}^c -p_i \log_2 p_i$$

$$I(\text{train } T, X) = \text{Entropy}(T) - \text{Entropy}(T, X)$$

$$E(\text{cheat}) = -\frac{6}{13} \log_2 \left(\frac{6}{13} \right) - \frac{7}{13} \log_2 \left(\frac{7}{13} \right)$$

$$E(\text{cheat}) = -0.46 \log_2 (0.46) - 0.53 \log_2 (0.53)$$

$$E(\text{cheat}) = 1.0007$$

1st Iteration:-

(A) Refund :- cheat | Total

T	R	Yes	N	Total
	Yes	3	4	7
Refund	NO	4	2	6
				13

$$E(\text{cheat}, \text{refund}) = P(\text{Yes}) \times E(3, 4) + P(\text{No}) \times E(4, 2) - ①$$

$$E(3,4) = -\frac{3}{7} \times \log_2 \left(\frac{3}{7}\right) - \frac{4}{7} \times \log_2 \left(\frac{4}{7}\right)$$

$$= -0.42 \times \log_2 (0.42) - 0.57 \times \log_2 (0.57)$$

$$\boxed{E(3,4) = 0.987} \quad - (a)$$

$$E(4,2) = -\frac{4}{6} \times \log_2 \left(\frac{4}{6}\right) - \frac{2}{6} \times \log_2 \left(\frac{2}{6}\right)$$

$$E(4,2) = -0.66 \times \log_2 (0.66) - 0.33 \times \log_2 (0.33)$$

$$\boxed{E(4,2) = 0.923} \quad - (b)$$

Substituting (a) and (b) in (i)

$$(i) \Rightarrow E(\text{cheat, refund}) = \frac{7}{13} \times 0.987 + \frac{6}{13} \times 0.923$$

$$= 0.53 \times 0.987 + 0.46 \times 0.923$$

$$\boxed{E(\text{cheat, refund}) = 0.947}$$

(B) Status			cheat		Total
			Yes	No	
status	Single		2	3	5
	Divorced		3	1	4
	Married		2	2	4
					13

$$(ii) E(\text{cheat, status}) = P(\text{Single}) \times E(2,3) + P(\text{Divorced}) \times E(3,1) + P(\text{Married}) \times E(2,2)$$

COPY
Date

$$E(2,3) = -\frac{2}{5} * \log_2\left(\frac{2}{5}\right) - \frac{3}{5} * \log_2\left(\frac{3}{5}\right)$$

$$= -0.4 * \log_2(0.4) - 0.6 * \log_2(0.6)$$

$$E(2,3) = 0.970 \quad - \textcircled{c}$$

$$E(3,1) = -\frac{3}{4} * \log_2\left(\frac{3}{4}\right) - \frac{1}{4} * \log_2\left(\frac{1}{4}\right)$$

$$= -0.75 * \log_2(0.75) - 0.25 * \log_2(0.25)$$

$$E(3,1) = 0.811 \quad - \textcircled{d}$$

$$E(2,2) = -\frac{2}{4} * \log_2\left(\frac{2}{4}\right) - \frac{2}{4} * \log_2\left(\frac{2}{4}\right)$$

$$= -0.5 * \log_2(0.5) - 0.5 * \log_2(0.5)$$

$$E(2,2) = 1 \quad - \textcircled{e}$$

Substituting \textcircled{c} , \textcircled{d} and \textcircled{e} in \textcircled{g}

$$\textcircled{g} \Rightarrow E(\text{cheat, status}) = \frac{5}{13} * 0.970 + \frac{4}{13} * 0.81 + \frac{4}{13} * 1$$

$$= 0.38 * 0.970 + 0.30 * 0.81 + 0.30 * 1$$

$$= 0.36 + 0.243 + 0.30$$

$$E(\text{cheat, status}) = 0.90$$

COPY

① TAX Income:

cheat

		cheat		Total
		Yes	No	
		low	high	
TAX Income	low	2	1	3
	medium	2	4	6
	high	3	1	4
				13

$$\text{iii) } E(\text{cheat}, \text{TAX Income}) = P(\text{low}) * E(2, 1) + P(\text{medium}) * E(2, 4) \\ + P(\text{high}) * E(3, 1)$$

$$E(2, 1) = -\frac{2}{3} * \log_2 \left(\frac{2}{3}\right) - \frac{1}{3} * \log_2 \left(\frac{1}{3}\right) \\ = -0.66 * \log_2 (0.66) - 0.33 * \log_2 (0.33)$$

$$E(2, 1) = 0.923 \quad \boxed{\text{f}}$$

Substituting ② ④ and ⑤ in ③

$$E(\text{cheat}, \text{TAX Income}) = \frac{3}{13} * 0.923 + \frac{6}{13} * 0.923 \\ + \frac{4}{13} * (0.811) \\ = 0.213 + 0.426 + 0.249$$

$$E(\text{cheat}, \text{TAX Income}) = 0.88 \quad \boxed{\text{g}}$$

$$\text{Gain (cheat, refund)} = 1 - 0.947$$

$$\boxed{\text{Gain (cheat, refund)} = 0.05}$$

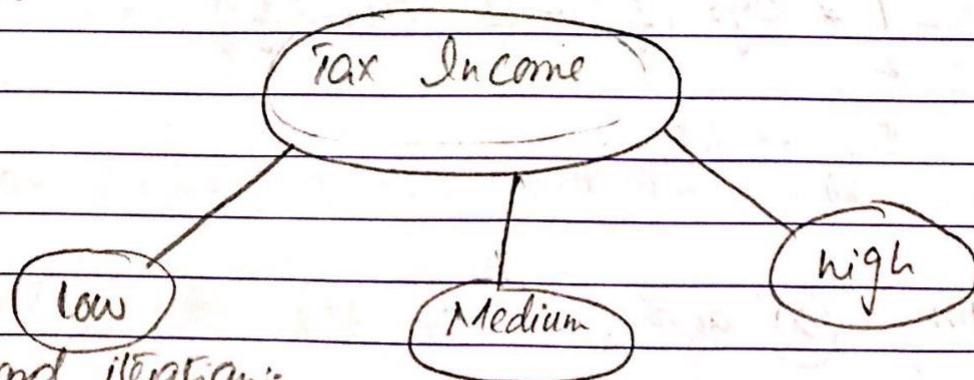
$$\text{Gain (cheat, Status)} = 1 - 0.90$$

$$\boxed{\text{Gain (cheat, Status)} = 0.10}$$

$$\text{Gain (cheat, Tax Income)} = 1 - 0.88$$

$$\boxed{\text{Gain (cheat, Tax Income)} = 0.12}$$

"Tax Income will be root node because of maximum information gain"



Second iteration:

Tax table:

Refund	Status	Cheat
Yes	Single	No
Yes	Single	Yes
No	Single	Yes

$$E(\text{cheat}) = -\frac{2}{3} \log_2 \left(\frac{2}{3}\right) - \frac{1}{3} \log_2 \left(\frac{1}{3}\right)$$

$$E(\text{cheat}) = 0.970$$

(A) refund:

		cheat		Total
		Yes	No	
Refund	Yes	1	1	2
	No	1	0	1
				3

$$E(\text{cheat, refund}) = P(\text{Yes}) \times E(1,1) + P(\text{No}) \times E(1,0) \quad \text{--- (i)}$$

$$[E(1,0) = 0] \quad \text{--- (a)}$$

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

$$[E(1,1) = 1] \quad \text{--- (b)}$$

Substituting (a) and (b) in (i)

$$(i) \Rightarrow E(\text{cheat, refund}) = \frac{2}{3} \times 1 + \frac{1}{3} \times 0$$

$$[E(\text{cheat, refund}) = 0.66]$$

(B) Status

Date: _____

cheat
Sun Mon Tue Wed Thu Fri Sat

		Yes	No	Total	
Status	Married	0	0	0	
	Single	2	1	3	
	Divorced	0	0	0	

$$E(\text{cheat}, \text{status}) = P(\text{Single}) * E(2, 1) + P(\text{Married}) * 0 \\ + P(\text{Divorced}) * 0 - ii$$

using Previous knowledge from iteration #01 in (i)

$$E(2, 1) = 0.923$$

$$ii) E(\text{cheat}, \text{status}) = \frac{3}{3} * 0.923$$

$$E(\text{cheat}, \text{status}) = 0.923$$

$$\text{Brain}(\text{cheat}, \text{refund}) = 0.970 - 0.66$$

$$\boxed{\text{Brain}(\text{cheat}, \text{refund}) = 0.31}$$

$$\text{Brain}(\text{cheat}, \text{status}) = 0.970 - 0.923$$

$$\boxed{\text{Brain}(\text{cheat}, \text{status}) = 0.047}$$

Second node will be refund in low branch.

(Tax Income)

Low

Medium

High

refund

No

Yes

cheat = Yes

cheat = No

COPY

High table :-

refund	status	cheat
Yes	Single	No
Yes	Married	Yes
No	Married	Yes
Yes	Divorced	Yes

$$E(\text{cheat}) = -\frac{1}{4} \log_2 \left(\frac{1}{4}\right) - \frac{1}{3} \log_2 \left(\frac{1}{3}\right) - \frac{2}{4} \log_2 \left(\frac{2}{4}\right)$$

$$= 0.5 + 0.5 + 0.5$$

$$\boxed{E(\text{cheat}) = 1.5}$$

(A) Refund :-

		Yes	No	Total
	Yes	2	1	3
refund	No	1	0	1

$$(i) E(\text{cheat}, \text{refund}) = P(\text{Yes}) * E(2, 1) + P(\text{No}) * E(1, 0)$$

$$E(1, 0) = 0$$

$$E(2, 1) = 0.923 \text{ (using prior knowledge)}$$

$$(ii) E(\text{cheat}, \text{refund}) = \frac{3}{4} * 0.923 + \frac{1}{4} * 0$$

$$E(\text{cheat, refund}) = 0.692$$

(ii) Status:-

		Yes	No	Total
Status	Single	0	1	1
	Married	2	0	2
	Divorced	1	0	1

$$E(\text{cheat, status}) = P(\text{Married}) \times E(2, 0) + P(\text{Single}) \times E(0, 1)$$

$$+ P(\text{Divorced}) \times E(1, 0) \quad - (ii)$$

$$E(2, 0) = E(0, 1) = E(1, 0) = 0$$

(i)

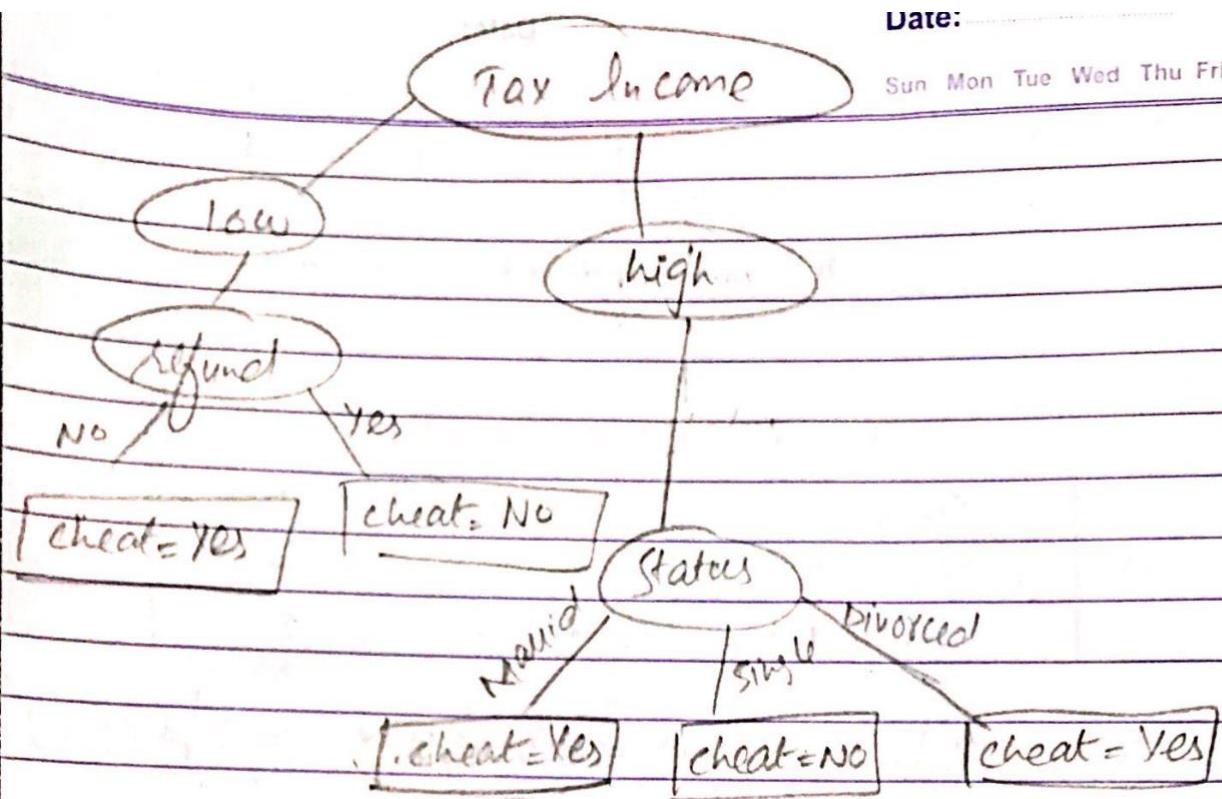
$$E(\text{cheat, status}) = 0$$

$$\text{Gain}(\text{cheat, refund}) = 1.5 - 0.692$$

$$\text{Gain}(\text{cheat, refund}) = 0.808$$

$$\text{Gain}(\text{cheat, status}) = 1.5$$

Status will be second node in medium branch due to maximum IG.



medium table :-

refund	Status	cheat
Yes	Married	NO
Yes	Divorced	NO
No	Single	NO
No	Married	NO
No	Divorced	Yes
No	Divorced	Yes

$$E(\text{cheat}) = \frac{4}{6} \times \log_2 \left(\frac{4}{6} \right) - \frac{2}{6} \log_2 \left(\frac{2}{6} \right) = 0.923$$

(A) refund:

		Yes	No	Total
refund	Yes	0	2	2
refund	NO	2	2	4
				6

$$(i) \cdot E(\text{cheat, refund}) = P(\text{Yes}) \cdot E(0,2) + P(\text{No}) \cdot E(2,2)$$

using Prior knowledge from previous iteration

$$E(0,2) = 0$$

$$E(2,2) = 1$$

$$(i) \Rightarrow E(\text{cheat, refund}) = \frac{2}{6} \times 0 + \frac{4}{6} \times 1$$

$$\boxed{E(\text{cheat, refund}) = 0.66}$$

(ii) Status:-

		Yes	No	Total
Status	Single	0	1	1
	Married	0	2	2
	Divorced	2	1	3
				6

$$(ii) \Rightarrow E(\text{cheat, refund}) = P(\text{single}) \cdot E(0,1) + P(\text{married}) \cdot E(0,2) \\ + P(\text{divorced}) \cdot E(2,1)$$

using prior knowledge from previous iterations

$$E(0,2) = E(0,1) = 0$$

$$E(2,1) = 0.923$$

$$(ii) \Rightarrow E(\text{cheat, refund}) = 0 + 0 + 0.46 = 0.46$$

COPY

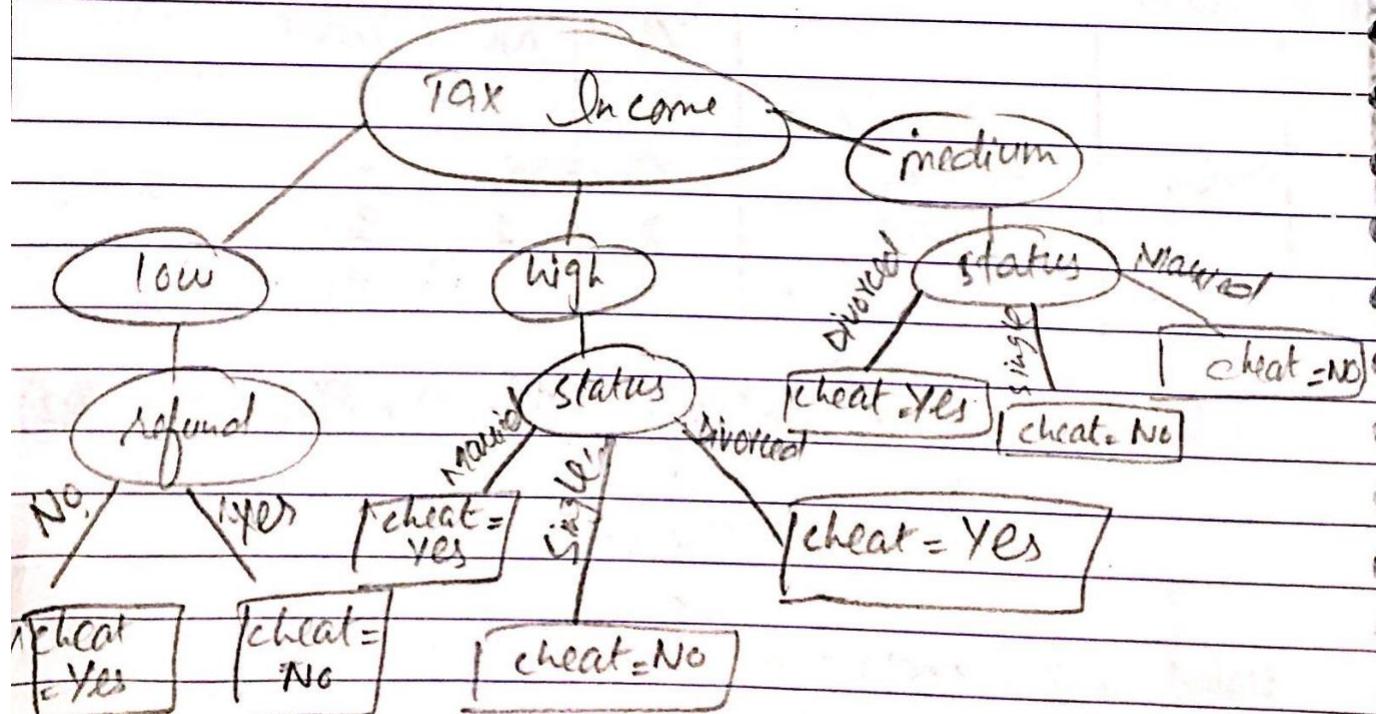
$$\text{Gain}(\text{cheat}, \text{refund}) = 0.923 - 0.66$$

$$\boxed{\text{Gain}(\text{cheat}, \text{refund}) = 0.263}$$

$$\text{Gain}(\text{cheat}, \text{status}) = 0.923 - 0.46$$

$$\boxed{\text{Gain}(\text{cheat}, \text{status}) = 0.463}$$

Status will be second node in medium branch due to maximum IG.



Due to $\text{max_depth} = 2$ we adjust some records into branches based on maximum occurrences.

$$\text{Gini index (Attribute = value)} = 1 - \sum_{i=1}^n (p_i)^2$$

$$\text{Gini index (Attribute)} = \sum_{V=\text{value}} P_V \times G_I(V)$$

1st Iteration:

(A) refund :-

		Yes	No	Total
	Yes	3	4	7
refund	No	4	2	6
				13

$$\text{Gini (refund = Yes)} = 1 - \left[\left(\frac{3}{7} \right)^2 + \left(\frac{4}{7} \right)^2 \right] = 1 - 0.18 - 0.32 = 0.5$$

$$\text{Gini (refund = No)} = 1 - \left[\left(\frac{4}{6} \right)^2 + \left(\frac{2}{6} \right)^2 \right] = 1 - 0.44 - 0.11 = 0.45$$

$$\text{Gini (refund)} = \frac{7}{13} \times 0.5 + \frac{6}{13} \times 0.45$$

$$= 0.26 + 0.20$$

$$\boxed{\text{Gini (refund)} = 0.46}$$

(B) Status

		Yes	No	Total
Status	Single	2	3	5
	Married	3	1	4
	Divorced	2	2	4
				13

$$\text{Gini}(\text{Status} = \text{Single}) = 1 - \left[\left(\frac{2}{5} \right)^2 + \left(\frac{3}{5} \right)^2 \right] = 1 - 0.16 - 0.36 \\ = 0.48$$

$$\text{Gini}(\text{Status} = \text{Married}) = 1 - \left[\left(\frac{3}{4} \right)^2 + \left(\frac{1}{4} \right)^2 \right] = 1 - 0.56 - 0.06 \\ = 0.378$$

$$\text{Gini}(\text{Status} = \text{Divorced}) = 1 - \left[\left(\frac{2}{4} \right)^2 + \left(\frac{2}{4} \right)^2 \right] = 1 - 0.25 - 0.25 \\ = 0.5$$

$$\text{Gini}(\text{status}) = \frac{5}{13} \times 0.48 + \frac{4}{13} \times 0.37 + \frac{4}{13} \times 0.5 \\ = 0.18 + 0.113 + 0.153$$

$$\boxed{\text{Gini}(\text{status}) = 0.446}$$

(e) Tax Income.

		Yes	No	Total
	Low	2	1	3
Tax Income	Medium	2	4	6
	High	3	1	4
				13

$$\text{Gini}(\text{Tax Income} = \text{Low}) = 1 - \left[\left(\frac{2}{3} \right)^2 + \left(\frac{1}{3} \right)^2 \right] = 1 - 0.44 - 0.11 = 0.49$$

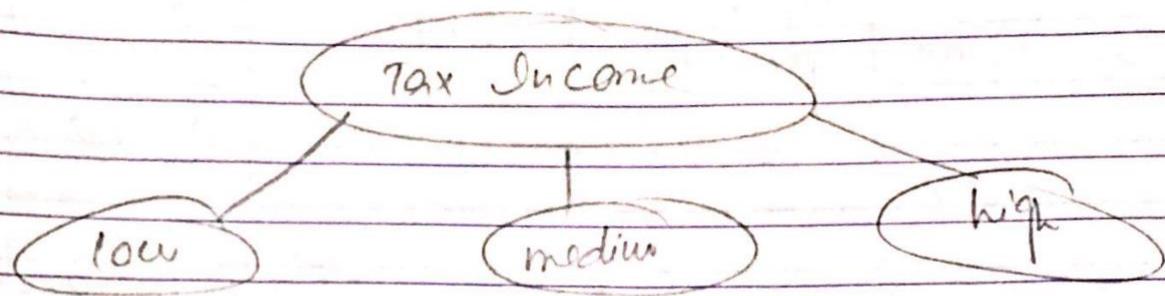
$$\begin{aligned} \text{Gini}(\text{Tax Income} = \text{medium}) &= 1 - \left[\left(\frac{2}{6} \right)^2 + \left(\frac{4}{6} \right)^2 \right] \\ &= 1 - 0.11 - 0.44 \\ &= 0.45 \end{aligned}$$

$$\begin{aligned} \text{Gini}(\text{Tax Income} = \text{high}) &= 1 - \left[\left(\frac{3}{4} \right)^2 + \left(\frac{1}{4} \right)^2 \right] \\ &= 1 - 0.562 - 0.0625 \\ &= 0.375 \end{aligned}$$

$$\begin{aligned} \text{Gini}(\text{Tax Income}) &= \frac{3}{13} \times 0.49 + \frac{6}{13} \times 0.45 + \frac{4}{13} \times 0.375 \\ &= 0.11 + 0.20 + 0.115 \end{aligned}$$

$$\text{Gini}(\text{Tax Income}) = 0.425$$

Since Tax Income has lowest gini index it will be root node.



Second Iteration :-

① For low branch :- (using table from previous pages).

(A) Refund.

		Yes	No	Total
refund	Yes	1	1	2
	No	1	0	1

$$\text{gini}(\text{refund} = \text{Yes}) = 1 - \left[\left(\frac{1}{3} \right)^2 - \left(\frac{1}{3} \right)^2 \right] = 1 - 0.44 - 0.01 \\ = 0.52$$

$$\text{gini}(\text{refund} = \text{No}) = 1 - \left[\left(\frac{1}{1} \right)^2 - \left(\frac{0}{1} \right)^2 \right] \\ = 0$$

$$\text{gini}(\text{refund}) = \frac{2}{3} * (0.52) + 0$$

$$\boxed{\text{gini}(\text{refund}) = 0.52}$$

B) Status

		Yes	No	Total
Status	Single	2	1	3
	Married	0	0	0
	Divorced	0	0	0

3

$$\text{gini}(\text{status} = \text{Single}) = 1 - \left[\left(\frac{2}{3} \right)^2 - \left(\frac{1}{3} \right)^2 \right] \\ = 0.45$$

$$\text{gini}(\text{status} = \text{Married}) = 1 - \left[\left(\frac{0}{3} \right)^2 - \left(\frac{0}{3} \right)^2 \right] \\ = 1$$

$$\text{gini}(\text{status} = \text{Divorced}) = 1 - \left[\left(\frac{0}{3} \right)^2 - \left(\frac{0}{3} \right)^2 \right] \\ = 1$$

$$\text{gini}(\text{status}) = \frac{3}{3} \times 0.45 + 0 + 0$$

$$\text{gini}(\text{status}) = 0.45$$

Since refund has lowest gini index it be the second node in low branch.

Tax Income

19
Date:

Sun Mon Tue Wed Thu Fri Sat

low

medium

high

refund

No

yes

[cheat. Yes] [cheat. No]

For high branch - (using table from previous pages).

A Refund

		Yes	No	Total
Refund	Yes	2	1	3
	No	1	0	1
				4

$$\text{gini}(\text{refund} = \text{Yes}) = 1 - \left[\left(\frac{2}{3} \right)^2 - \left(\frac{1}{3} \right)^2 \right]$$

$$= 0.45$$

$$\text{gini}(\text{refund} = \text{No}) = 1 - \left[\left(\frac{1}{1} \right)^2 - \left(\frac{0}{1} \right)^2 \right]$$

$$= 1 - 1 = 0$$

$$\text{gini}(\text{refund}) = \frac{3}{4} \times 0.45 + 0 = 0.33$$

(B) Status

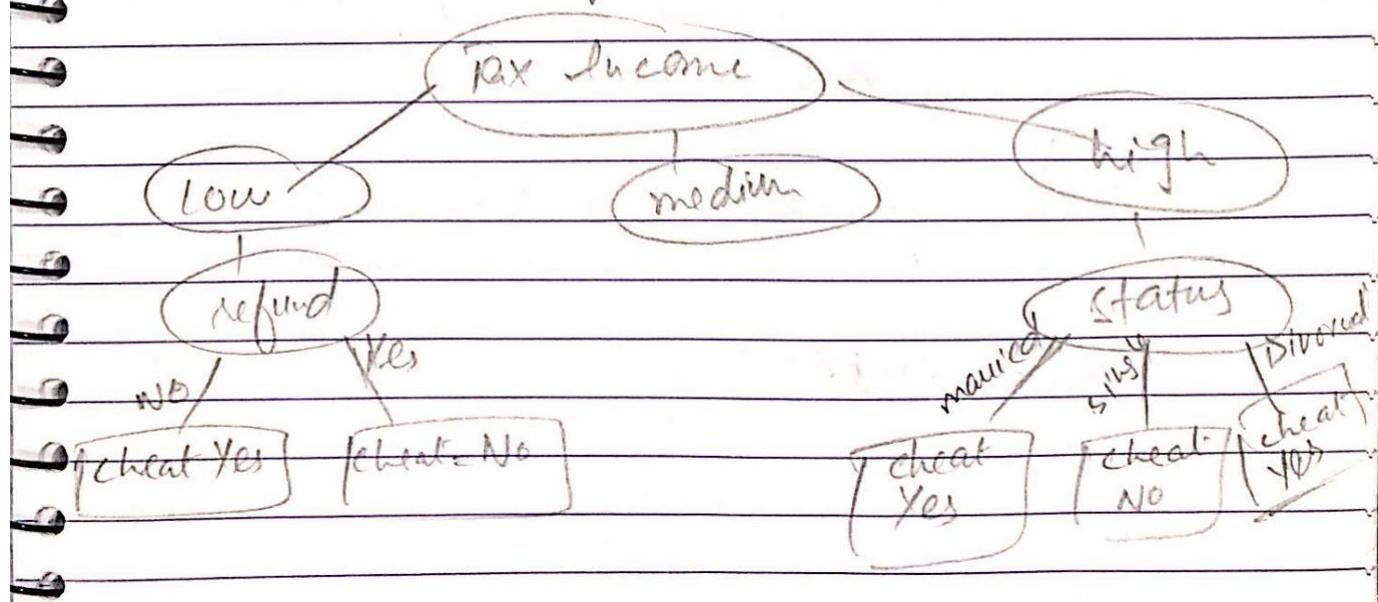
		Yes	No	Total
status	Singl	0	1	1
	Married	2	0	2
	Divorced	1	0	1
				4

$$\text{gini}(\text{status} = \text{Singl}) = 1 - \left[\left(\frac{0}{1} \right)^2 + \left(\frac{1}{1} \right)^2 \right] \\ = 0$$

$$\text{gini}(\text{status} = \text{Married}) = \text{gini}(\text{Divorced}) = 0$$

$$\text{gini}(\text{status}) = \frac{1}{4} * 0 + \frac{2}{4} * 0 + \frac{1}{4} * 0 = 0$$

Since status has lowest gini index it will be second node in high branch.



For medium branch

using table from previous pages.

(A) refund:

	Yes	Yes	No	Total
Yes	0	2	2	
refund	No	2	2	4
			6	

$$\text{gini}(\text{refund} = \text{Yes}) = 1 - \left[\left(\frac{0}{2} \right)^2 + \left(\frac{2}{2} \right)^2 \right] \\ = 0$$

$$\text{gini}(\text{refund} = \text{No}) = 1 - \left[\left(\frac{2}{4} \right)^2 + \left(\frac{2}{4} \right)^2 \right] \\ = 0.5$$

$$\text{gini}(\text{refund}) = \frac{2}{6} * 0 + \frac{4}{6} * 0.5$$

$$\boxed{\text{gini}(\text{refund}) = 0.33}$$

(B) Status

		Yes	No	Total
Status	Single	0	1	1
	Married	0	2	2
	Divorced	2	1	3
				6

$$\text{gini}(\text{status} = \text{Single}) = \text{gini}(\text{status} = \text{Married}) = 0$$

$$\text{gini}(\text{status} = \text{Divorced}) = 1 - \left[\left(\frac{2}{3} \right)^2 + \left(\frac{1}{3} \right)^2 \right]$$

$$= 0.45$$

$$\text{gini}(\text{status}) = \frac{1}{6} \times 0 + \frac{2}{6} \times 0 + \frac{3}{6} \times 0.45$$

$$\boxed{\text{gini}(\text{status}) = 0.225}$$

Since status has lowest gini index it be second node in medium branch.

Some instance are adjusted based on maximum occurrence criteria due to max depth constraint, can't grow tree further.

