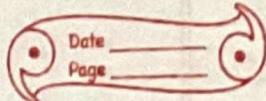


Name: Siddharth Betala  
 Roll: BE19B032  
 MAS710 MMI Assignment -3



## ① Training Data:

Index	Age	Salary	Class
1	30	65	G
2	23	15	B
3	40	75	G
4	55	40	B
5	55	100	G
6	45	60	G

Thid Column

has been added  
by me for  
use ahead.

Usually we decide the root node condition by attribute based on given Gini impurity value.

(i) Age & Salary as Attributes

Since 'Age' & 'Salary' are both numeric attributes, we need to presort the data.

Presorting based on age.		Presorting based on salary.	
Age	Index	Age	Index
23	2	15	2
30	1	40	4
40	3	60	6
45	6	65	1
55	5	75	3
55	4	100	5

For a ~~root~~ potential node,

$$\text{Gini Impurity } (G) = \frac{|D|}{|D|} G_{\text{ini}}(D_1) + \frac{|D_a|}{|D|} G_{\text{ini}}(D_a)$$

$$G_{\text{ini}}(D_K) = 1 - \sum_{j=1}^N p_j^2$$

~~for K=1, 2~~

N = No. of Classes

$p_j$  = Frequency / Probability  
of a class.

(A) Age as Attribute

a)  $\text{Age} \leq 23$ :

$\text{Age} \leq 23$			
	True	False	
$G_1$	B	G	B
0	1	4	1

$|D_1| = 1$        $|D_a| = 2$ .

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

$$= 1 - 1$$

$$= 0$$

$$= 1 - \left( \frac{4}{4+1} \right)^2 - \left( \frac{1}{4+1} \right)^2$$

$$= 1 - \frac{16}{25} - \frac{1}{25}$$

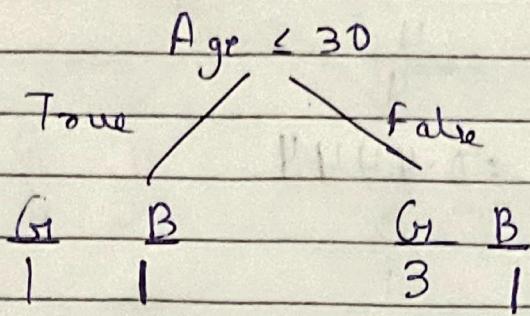
$$= \frac{8}{25}$$

$$\text{Gini}(\text{Age} \leq 23) = \frac{|D_1|}{|D|} G_{\text{ini}}(D_1) + \frac{|D_a|}{|D|} G_{\text{ini}}(D_a)$$

$$= \frac{1}{6}(0) + \frac{5}{6}\left(\frac{8}{25}\right)$$

$$= \frac{4}{15} = 0.2667$$

(b) Age  $\leq 30$ :



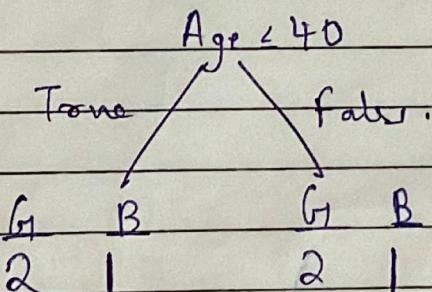
$$\begin{aligned} \text{Gini}(D_1) &= 1 - \left(\frac{1}{1+1}\right)^2 - \left(\frac{1}{1+1}\right)^2 \\ &= 1 - \frac{1}{2} - \frac{1}{2} \\ &= 0.5 \end{aligned}$$

$$\begin{aligned} \text{Gini}(D_2) &= 1 - \left(\frac{3}{1+3}\right)^2 - \left(\frac{1}{1+3}\right)^2 \\ &= 1 - \frac{9}{16} - \frac{1}{16} \\ &= \cancel{\frac{6}{16}} - \frac{3}{8} \end{aligned}$$

$$\begin{aligned} \text{Gini}(\text{Age} \leq 30) &= \frac{2}{6} (0.5) + \frac{4}{6} \cancel{\frac{3}{8}} \\ &= \frac{6}{24} + \frac{24}{24} \\ &= \frac{9}{24} = \frac{3}{8} = 0.375 \end{aligned}$$

$$\begin{aligned} \text{Gini}(\text{Age} \leq 30) &= \frac{2}{6} (0.5) + \frac{4}{6} \times \frac{3}{8} \\ &= \frac{1}{6} + \frac{1}{4} \\ &= \frac{5}{12} = 0.4167 \end{aligned}$$

(c) Age  $\leq 40$ :



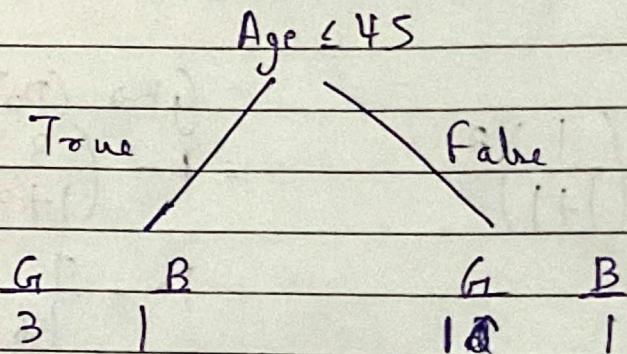
$$\begin{aligned} \text{Gini}(D_1) &= \text{Gini}(D_2) = 1 - \left(\frac{1}{3}\right)^2 - \left(\frac{2}{3}\right)^2 \\ &\quad \leftarrow \frac{4}{9} \end{aligned}$$

$$Gini(\text{Age} \leq 40) = \frac{3}{6} \times \frac{4}{9} + \frac{3}{6} \times \frac{4}{9}$$

$$= \frac{4}{9}$$

$$= 0.4444$$

(d) Age  $\leq 45$ :



$$Gini(\text{Age} \leq 45) =$$

$$\begin{aligned} & Gini(D_1) \\ & = 1 - \left(\frac{3}{4}\right)^2 - \left(\frac{1}{4}\right)^2 \\ & = \frac{3}{8} \end{aligned}$$

~~$$\begin{aligned} & Gini(D_2) \\ & = 1 - \left(\frac{2}{3}\right)^2 - \left(\frac{1}{3}\right)^2 \\ & = \frac{1}{9} \end{aligned}$$~~

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

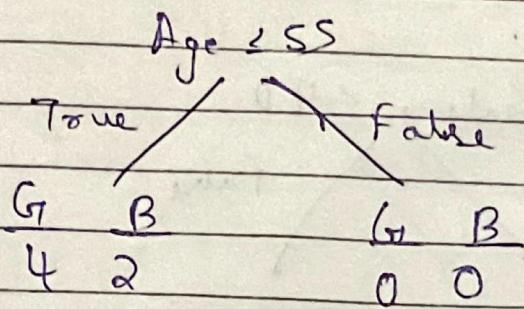
$$Gini(\text{Age} \leq 45)$$

$$= \frac{4}{6} \times \frac{3}{8} + \frac{2}{6} \times \frac{1}{2}$$

$$= \frac{1}{4} + \frac{1}{6}$$

$$= \frac{5}{12} = 0.4167$$

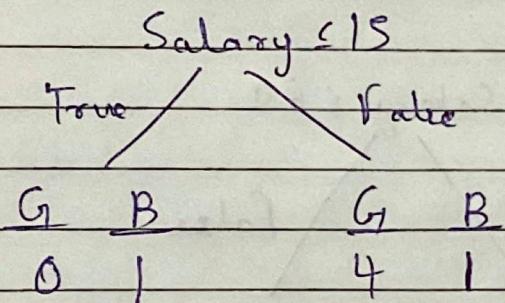
(e) Age < 55:



$$\begin{aligned}
 \text{Gini}(\text{Age} < 55) &= 1 - \left(\frac{4}{6}\right)^2 - \left(\frac{2}{6}\right)^2 \\
 &= \frac{4}{9} \\
 &= 0.4444
 \end{aligned}$$

(B) Salary as Attribute

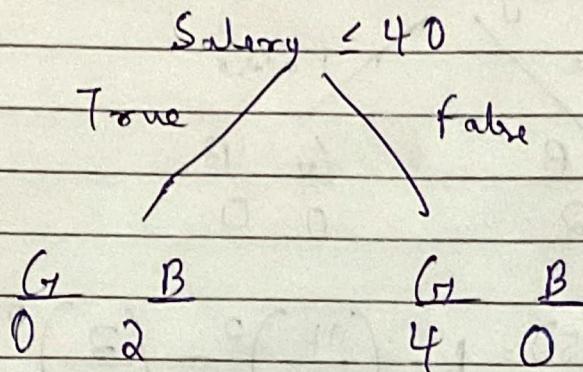
(a) Salary < 15:



$$\begin{aligned}
 \text{Gini}(D_1) &= 1 - \left(\frac{0}{1}\right)^2 - \left(\frac{1}{1}\right)^2 \\
 &= 0
 \end{aligned}
 \quad
 \begin{aligned}
 \text{Gini}(D_2) &= 1 - \left(\frac{4}{5}\right)^2 - \left(\frac{1}{5}\right)^2 \\
 &= \frac{8}{25}
 \end{aligned}$$

$$\begin{aligned}
 \text{Gini}(\text{Salary} < 15) &= \frac{1}{6} \times 0 + \frac{5}{6} \times \frac{8}{25} \\
 &= \frac{4}{15} = 0.26667
 \end{aligned}$$

(b) Salary  $\leq 40$ :

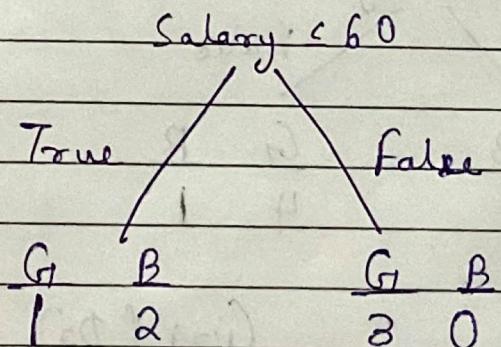


$$\text{Gini}(P_1) = 1 - \left(\frac{0}{2}\right)^2 - \left(\frac{2}{2}\right)^2 = 0$$

$$\text{Gini}(D_2) = 1 - \left(\frac{0}{4}\right)^2 - \left(\frac{4}{4}\right)^2 = 1 - 1 = 0$$

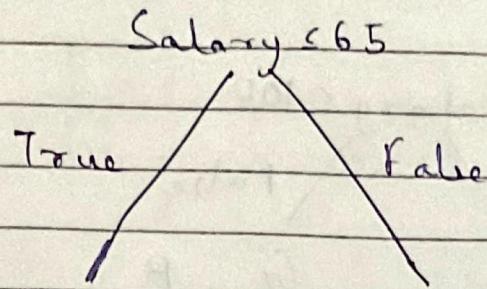
$$\text{Gini}(\text{Salary} \leq 40) = \frac{2}{6} \times 0 + \frac{4}{6} \times 0 = 0.$$

(c) Salary  $\leq 60$ :



$$\begin{aligned} \text{Gini}(P_1) &= 1 - \left(\frac{1}{3}\right)^2 - \left(\frac{2}{3}\right)^2 = 0 \\ \text{Gini}(D_2) &= 1 - \left(\frac{3}{3}\right)^2 - \left(\frac{0}{3}\right)^2 = 0 \\ \text{Gini}(\text{Salary} \leq 60) &= \frac{3}{6} \times \frac{1}{3} \times \frac{4}{9} + \frac{3}{6} \times 0 = \frac{2}{9} = 0.2222 \end{aligned}$$

(d)  $\text{Salary} \leq 65$ :

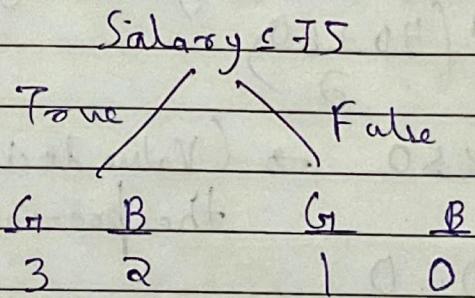


$$\begin{aligned} \text{Gini}(D_1) &= 1 - \left(\frac{2}{4}\right)^2 - \left(\frac{2}{4}\right)^2 \\ &= 0.5 \end{aligned}$$

$$\begin{aligned} \text{Gini}(D_2) &= 1 - \left(\frac{2}{2}\right)^2 \\ &= 0 \end{aligned}$$

$$\begin{aligned} \text{Gini}(\text{Salary} \leq 65) &= \frac{4}{6} \times \frac{1}{2} + \frac{2}{6} \times 0 \\ &= \frac{1}{3} = 0.3333 \end{aligned}$$

(e)  $\text{Salary} \leq 75$ :

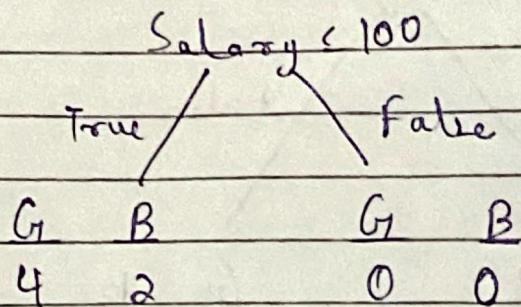


$$\begin{aligned} \text{Gini}(D_1) &= 1 - \left(\frac{3}{5}\right)^2 - \left(\frac{2}{5}\right)^2 \\ &= \frac{12}{25} \end{aligned}$$

$$\begin{aligned} \text{Gini}(D_2) &= 1 - \left(\frac{1}{1}\right)^2 \\ &= 0 \end{aligned}$$

$$\text{Gini}(\text{Salary} \leq 75) = \frac{5}{6} \times \frac{12}{25} + \frac{1}{6} \times 0 = 0.4$$

(f)  $\text{Salary} \leq 100$ :



$$\begin{aligned} \text{Gini}(D_1) &= 1 - \left(\frac{4}{6}\right)^2 - \left(\frac{2}{6}\right)^2 \\ &= \frac{4}{9} \end{aligned}$$

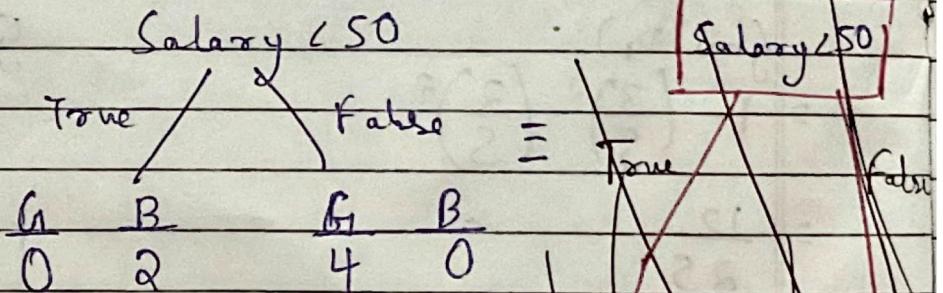
$$\text{Gini}(\text{Salary} \leq 100) = \frac{4}{9} = 0.4444$$

Best Split is given by the smallest Gini Impurity value.

Smallest Gini Impurity is for the split:  $\text{Salary} \leq 40$ .  
The ~~second~~ same split can be given by the strict condition:  $\text{Salary} < \frac{40+60}{2}$

$\Rightarrow \text{Salary} < 50 \rightarrow$  (Value decided based on the pre-sorting order)

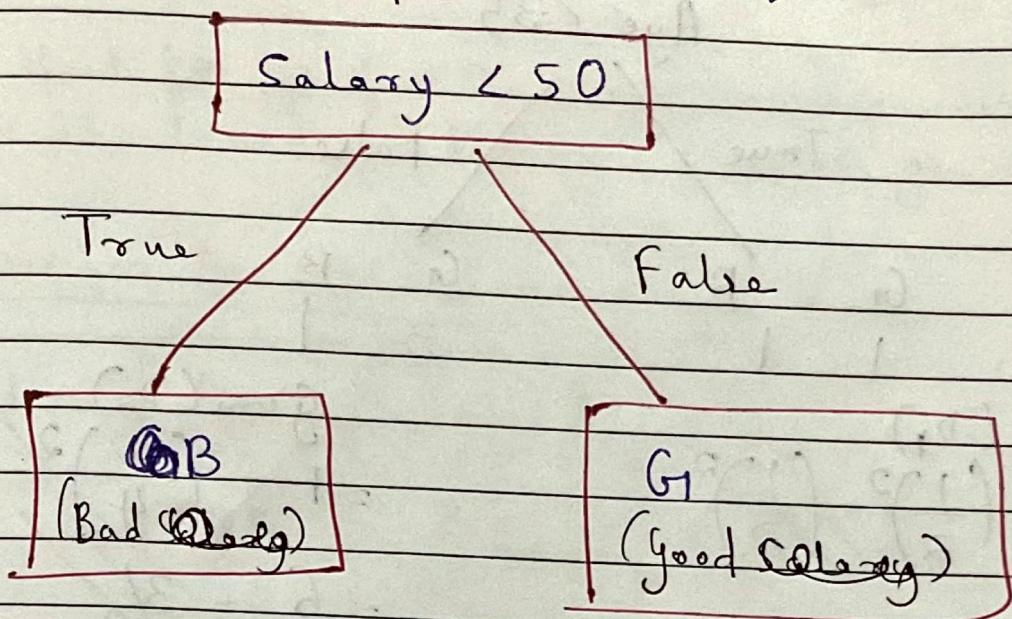
$$\text{Gini}(\text{Salary} < 50) = 0$$



(We don't need to classify further as these nodes are pure. Hence, these are leaf nodes)

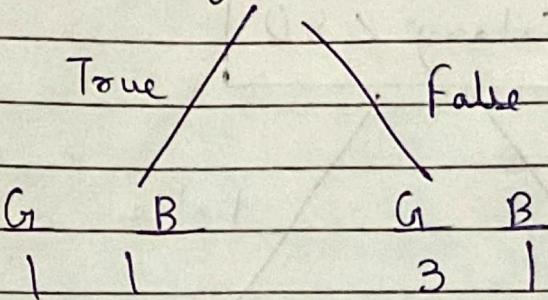
Good  
Salary/  
G1

Based on the best split the classifier tree is given as:



## (ii) Root Partition Keeping ( $\text{Age} \leq 35$ )

$\text{Age} \leq 35$



$$\begin{aligned} \text{Gini}(D_1) \\ = 1 - \left(\frac{1}{2}\right)^2 - \left(\frac{1}{2}\right)^2 \\ = 0.5 \end{aligned}$$

$$\begin{aligned} \text{Gini}(D_2) \\ = 1 - \left(\frac{3}{4}\right)^2 - \left(\frac{1}{4}\right)^2 \\ = \frac{6}{16} = \frac{3}{8} \\ = 0.375 \end{aligned}$$

We see that the nodes  $D_1$  &  $D_2$  are both impure.

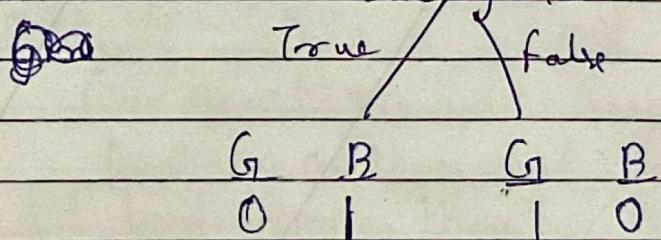
We can try classifying the tree further by calculating the gini impurity value using Salary as an attribute. Note that we'll proceed only if the subsequent gini values are smaller than  $\text{Gini}(D_1)$  &  $\text{Gini}(D_2)$ . i.e., the any further classification should lead to purer classes than the current ones.

Classes for  $\text{Age} \leq 35 \wedge \text{True}$  are.

(Pre-Sorted based on Salary)

Index	Age	Salary	Class
2	23	15	B
1	30	65	G1

$\text{Salary} \leq 15$



$$\begin{aligned} \text{Gini}(D_1) &= 1 - (1)^2 = 0 \\ \text{Gini}(D_2) &= 1 - (1)^2 = 0 \end{aligned}$$

$$Gini(\text{Salary} \leq 15) = 0 \times \frac{1}{2} + 0 \times \frac{1}{2} = 0.$$

$\therefore$  Best Split for the node corresponding to  $(\text{Age} \leq 35: \text{True})$  is  $\text{Salary} \leq 15$  with gini impurity of 0.

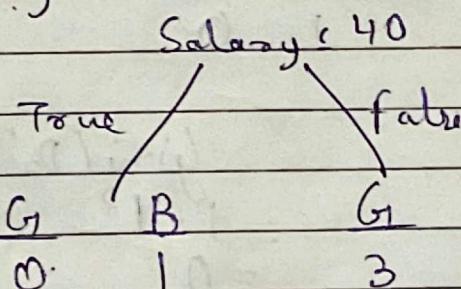
Class for  $\text{Age} \leq 35: \text{False}$  are

Index	Age	Salary	Class
40	75		
45	60		
55			
55			

(Pre-sorted based on salary)

Index	Age	Salary	Class
4	55	40	B
6	45	60	G
3	40	75	G
5	55	100	G

(a) ~~Given~~  $\text{Salary} \leq 40:$

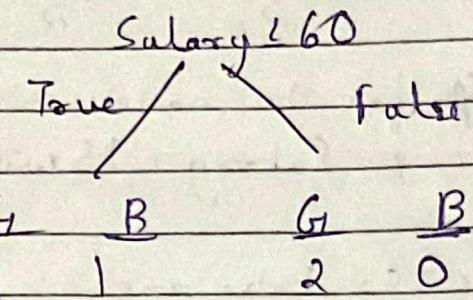


$$\begin{aligned} Gini(D_1) \\ = 1 - 1^2 \\ = 0 \end{aligned}$$

$$\begin{aligned} Gini(D_2) \\ = 1 - \left(\frac{3}{3}\right)^2 \\ = 1 - 1 = 0. \end{aligned}$$

$$Gini(\text{Salary} \leq 40) = \frac{1}{4} \times 0 + \frac{3}{4} \times 0 = 0.$$

(b) Salary  $\leq 60$ :

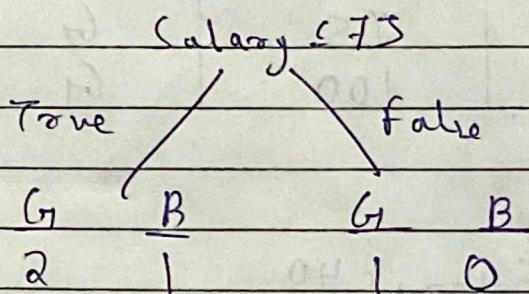


$$\begin{aligned}
 & \text{Gini}(D_1) \\
 &= 1 - \left(\frac{1}{2}\right)^2 - \left(\frac{1}{2}\right)^2 \\
 &= 0.5
 \end{aligned}
 \quad
 \begin{aligned}
 & \text{Gini}(D_2) \\
 &= 1 - \left(\frac{2}{2}\right)^2 \\
 &= 0
 \end{aligned}$$

Gini(Salary  $\leq 60$ )

$$\begin{aligned}
 &= \frac{2}{4} \times 0.5 + \frac{2}{4} \times 0 \\
 &= 0.25
 \end{aligned}$$

(c) Salary  $\leq 75$ :



$$\begin{aligned}
 & \text{Gini}(D_1) \\
 &= 1 - \left(\frac{2}{3}\right)^2 - \left(\frac{1}{3}\right)^2 \\
 &= \frac{4}{9}
 \end{aligned}
 \quad
 \begin{aligned}
 & \text{Gini}(D_2) \\
 &= 1 - 1^2 \\
 &= 0
 \end{aligned}$$

Gini(Salary  $\leq 75$ )

$$\begin{aligned}
 &= \frac{3}{4} \times \frac{4}{9} = \frac{1}{3} \\
 &= 0.3333
 \end{aligned}$$

(d)  $\text{Salary} \leq 100$

$\text{Salary} \leq 100$

True / False

G	B	G	B
3	1	0	0

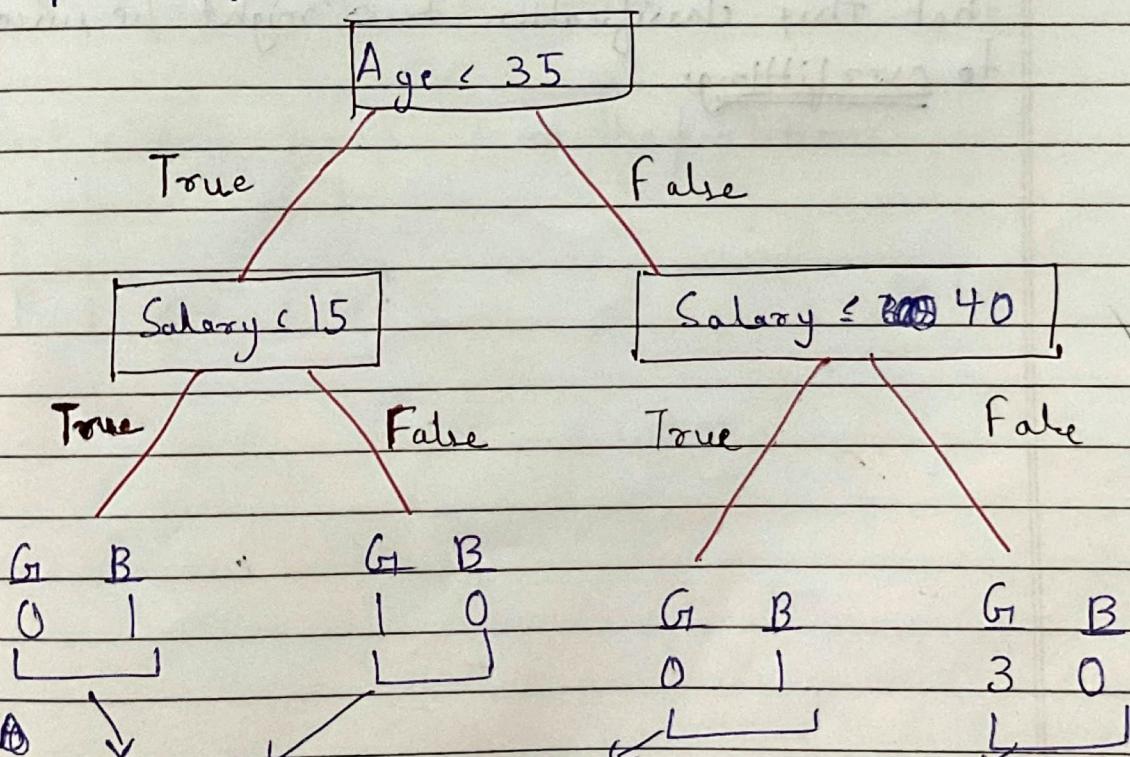
$$\text{Gini}(D_1) = 1 - \left(\frac{3}{4}\right)^2 - \left(\frac{1}{4}\right)^2 \\ = \frac{3}{8}$$

$\text{Gini}(\text{Salary} \leq 100)$

$$= \frac{4}{4} \times \frac{3}{8}$$

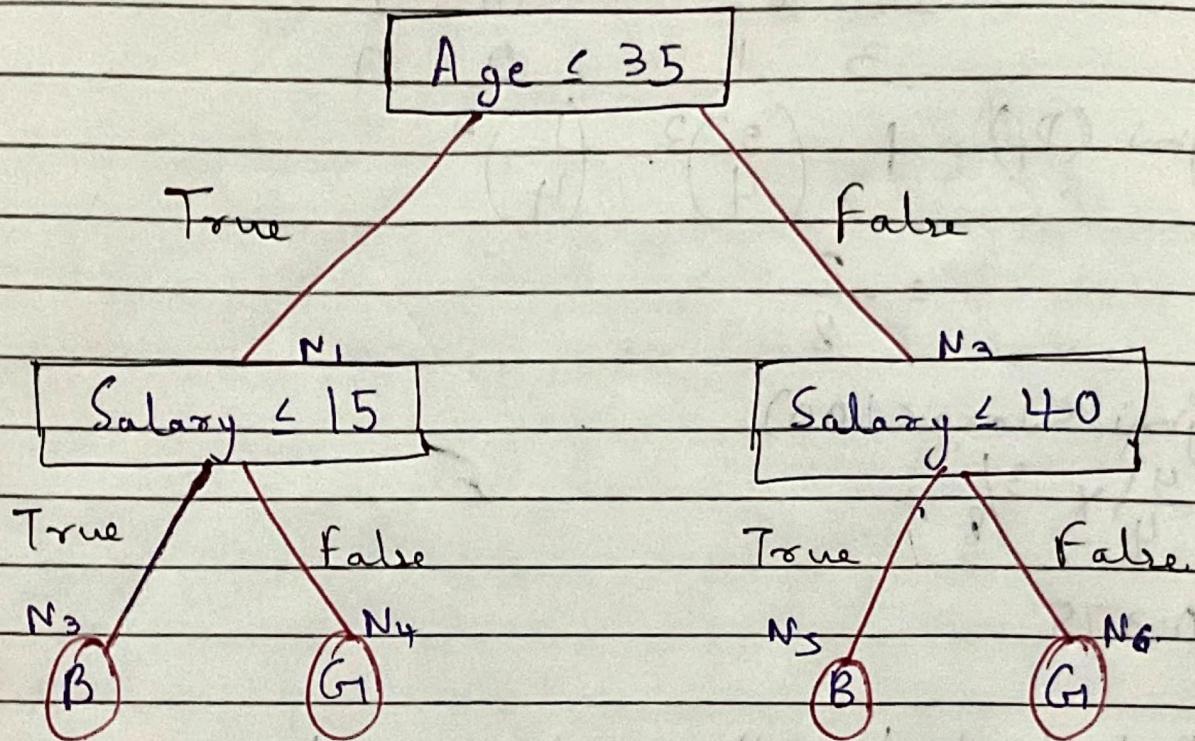
$$= 0.375$$

∴ Best Split for the node corresponding to  $(\text{Age} < 35 : \text{False})$  is  $\text{Salary} \leq 40$  with Gini impurity of 0.



All these nodes are pure, & give an exact classification

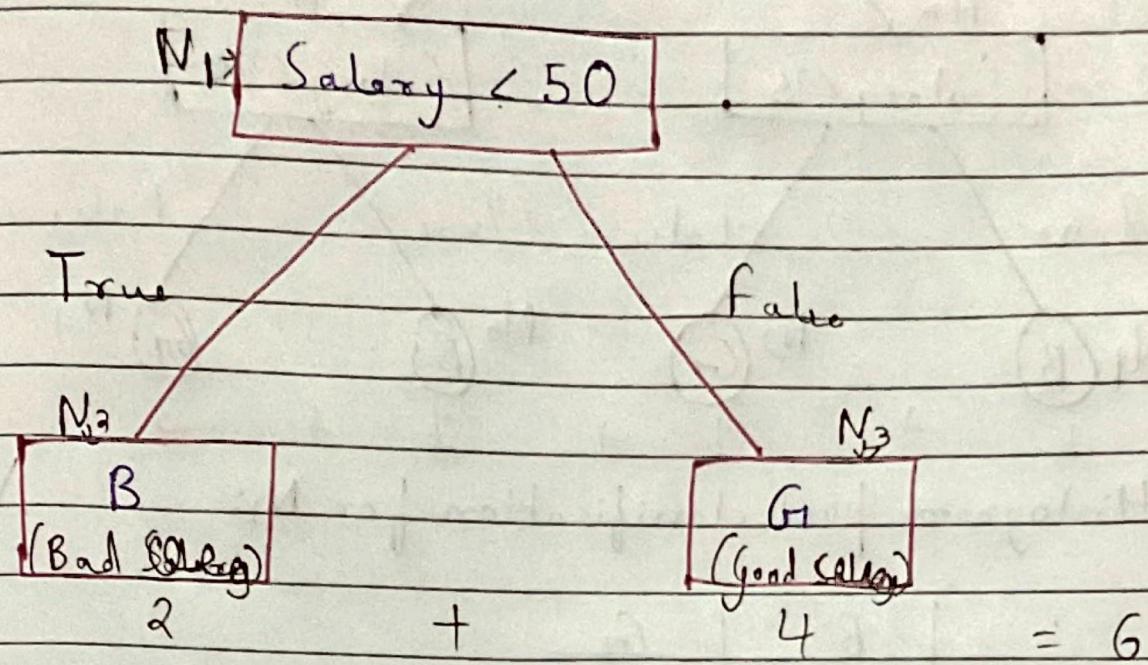
So there's no need for further classification.  
 These 4 can be taken as leaf nodes.  
 So the decision tree with  $\text{Age} \leq 35$  as root partition is given as:



Note: Only one training example make it to the nodes:  $N_3, N_4, N_5$ . So it is important to note that this classification tree might be susceptible to overfitting.

(2) Here, we have to compute the accuracy of the classifier from Q2 at every node using class histogram.

Ideal Classification tree (Refer ①(i))



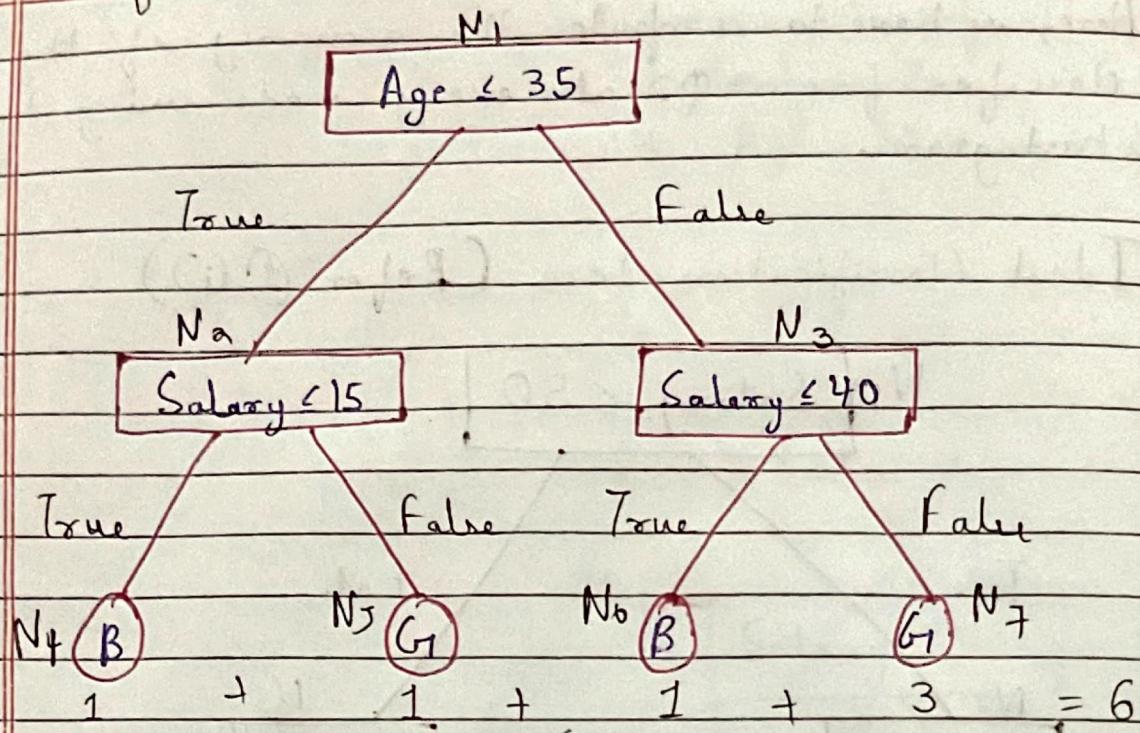
Initial histogram for  $N_1$ :

	B	G
L	0	0
R	2	4

Post classification based on the partition:

	B	G
L	2	0
R	0	4

Classification Tree using root node Age < 35



Histogram post classification for N<sub>1</sub>:

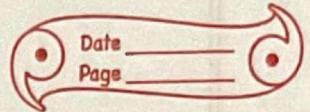
	B	G <sub>1</sub>
(N <sub>a</sub> ) L	1	1
(N <sub>3</sub> ) R	1	3

Accuracy for left Branch / to get N<sub>a</sub> node

$$\begin{aligned}
 &= \frac{\text{No. of } B^? \text{ s in } N_a}{\text{Expected No. of } B^? \text{ s in } N_1} \times 100 \\
 &= \frac{1}{2} \times 100 \\
 &= 50.1.
 \end{aligned}$$

Accuracy for Right Branch / to get N<sub>3</sub> node

$$\begin{aligned}
 &= \frac{\text{No. of } G_1^? \text{ s in } N_3}{\text{Expected no. of } G_1^? \text{ s in } N_1} \times 100 \\
 &= \frac{3}{4} \times 100 \\
 &= 75.1.
 \end{aligned}$$



Accuracy of getting nodes  $N_4$ ,  $N_5$ ,  $N_6$  &  $N_7$  is 100%. because these are leaf nodes, i.e., they offer a perfect split.

When we the root partition as (Age  $< 35$ ), we do not get the perfect split as given by the SLIP model. Hence, we calculate the accuracy for both the left and right branches (how accurately does it identify the label 'B' in the left branch & the label 'G' in the right branch)