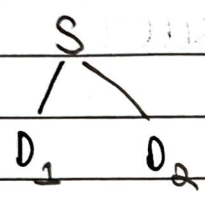


MA5710 - Assignment 03 (Devansh - B0190002)

Case Study 03: Supervised Learning - Classification

Question One



classification on main training set
can be based on either age
or salary.

Data set:

Age	Salary	Class	Index
30	65	G	1
23	15	B	2
40	75	G	3
55	40	B	4
85	100	G	5
45	60	G	6

① Age and salary as attributes.

Re-sorting (Age)

23	2
30	1
40	3
45	6
55	5
58	4

Re-sorting (Salary)

15	2
40	4
60	6
65	1
75	3
100	5

②
$$Gini\ Index\ (AD) = \frac{|D_1|}{|D|} Gini(D_1) + \frac{|D_2|}{|D|} Gini(D_2)$$

~~$$Gini(A) = \frac{|D_1|}{|D|} Gini(D_1) + \frac{|D_2|}{|D|} Gini(D_2)$$~~

$$Gini(A) = 1 - \sum_{j=1}^N p_j^2$$

N classes

p_j = frequency of a class

here N=2

(I) Age as attribute

a) Age ≤ 23 : $ID_1 = 1$, $ID_2 = 5$

$$\therefore Gini(D) = \frac{1}{5} Gini(D_1) + \frac{5}{65} \left(1 - \left(\left(\frac{1}{5} \right)^2 + \left(\frac{4}{5} \right)^2 \right) \right)$$

$$= \frac{1}{5} \cdot 0 + \frac{5}{65} \left(1 - \left(\frac{1}{25} + \frac{16}{25} \right) \right)$$

$$= \frac{5}{65} \cdot \frac{18}{25} = 0.26667$$

b) Age ≤ 30 : $ID_1 = 2$, $ID_2 = 4$

$$Gini(D) = \frac{2}{6} \left(1 - \left(\left(\frac{1}{2} \right)^2 + \left(\frac{1}{2} \right)^2 \right) \right) + \frac{4}{6} \left(1 - \left(\left(\frac{2}{4} \right)^2 + \left(\frac{2}{4} \right)^2 \right) \right)$$

$$= \frac{1}{3} + \frac{1}{3} = \frac{2}{3}$$

$$= \frac{5}{6} = 0.41666$$

c) Age ≤ 40 : $ID_1 = 3$, $ID_2 = 3$

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

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

$$= 0.4444$$

d) Age ≤ 45 : $ID_1 = 4$, $ID_2 = 2$

$$Gini(D) = \frac{2}{3} \left(1 - \left(\left(\frac{2}{4} \right)^2 + \left(\frac{2}{4} \right)^2 \right) \right) + \frac{1}{3} \left(1 - \left(\left(\frac{1}{3} \right)^2 + \left(\frac{2}{3} \right)^2 \right) \right)$$

$$= \frac{1}{3} + \frac{1}{3} = \frac{2}{3}$$

$$= 0.41666$$

e) Age ≤ 55 : $ID_1 = 6$, $ID_2 = 0$

$$Gini(D) = \frac{6}{6} \left(1 - \left(\left(\frac{2}{6} \right)^2 + \left(\frac{4}{6} \right)^2 \right) \right) + 0$$

$$= \frac{4}{9} = 0.4444$$

II Salary as Attribute

a) salary ≤ 15 : $|D_1| = 1$, $|D_2| = 5$

$$\text{Gini}(D) = \frac{1}{6} (0) + \frac{5}{6} \left(1 - \left(\left(\frac{1}{5} \right)^2 + \left(\frac{4}{5} \right)^2 \right) \right)$$

$$= 0.2667$$

b) salary ≤ 40 : $|D_1| = 2$, $|D_2| = 4$

$$\text{Gini}(D) = \frac{2}{6} (1 - 1) + \frac{4}{6} (1 - 1)$$

$$= 0$$

c) salary ≤ 60 : $|D_1| = 3$, $|D_2| = 3$

$$\text{Gini}(D) = \frac{1}{2} \left(1 - \left(\left(\frac{2}{3} \right)^2 + \left(\frac{1}{3} \right)^2 \right) \right) + \frac{1}{2} \left(1 - 1 \right)$$

$$= \frac{2}{9} = 0.222$$

d) salary ≤ 65 : $|D_1| = 4$, $|D_2| = 2$

$$\text{Gini}(D) = \frac{2}{3} \left(1 - \left(\left(\frac{2}{4} \right)^2 + \left(\frac{2}{4} \right)^2 \right) \right) + 0$$

$$= \frac{1}{3} = 0.333$$

e) salary ≤ 75 : $|D_1| = 5$, $|D_2| = 1$

$$\text{Gini}(D) = \frac{5}{6} \left(1 - \left(\left(\frac{4}{5} \right)^2 + \left(\frac{1}{5} \right)^2 \right) \right) + 0$$

$$= \frac{2}{5} = 0.4$$

f) salary ≤ 100 : $|D_1| = 6$, $|D_2| = 0$

$$\text{Gini}(D) = 1 \left(1 - \left(\left(\frac{1}{3} \right)^2 + \left(\frac{2}{3} \right)^2 \right) \right)$$

$$= \frac{4}{9} = 0.444$$

Best split among age and salary is salary
 - with salary $\leq (40+60)$ which has
 a Gini index of 0.2

salary ≤ 80

ii Root partition keeping Age ≤ 35

→ Age ≤ 35 : D = 2 samples.

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

$$\text{Gini (Age & salary } \leq 15) = \frac{1}{2} (1 - (1)^2) + \frac{1}{2} (1 - (1)^2)$$

$$= 0$$

∴ Best split = salary ≤ 15 .

→ Age ≥ 35

Age	Salary	Class	Index	Salary	Index
40	75	G	3	40	4
48	60	G	6	60	6
55	100	G	8	75	3
55	40	B	4	100	5

$$(1) \text{ Gini (Age & salary } \leq 40) = \frac{1}{4} (1 - (1)^2) + \frac{3}{4} (1 - (\frac{3}{3})^2)$$

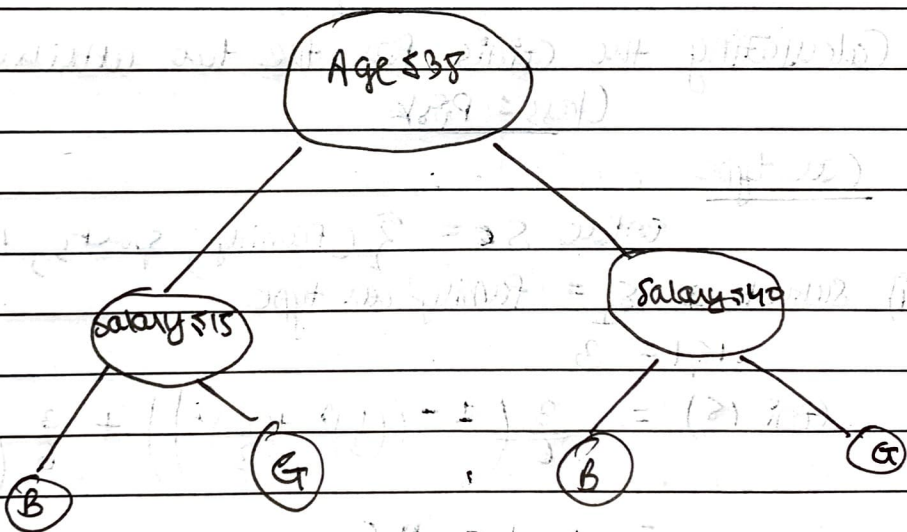
$$= 0$$

$$(i) \text{ Gini}(\text{salary} \leq 60) = \frac{2}{4} (1 - ((\frac{1}{2})^2 + (\frac{1}{2})^2)) + \frac{2}{4} (1 - (1)^2) \\ = \frac{1}{4} = 0.25$$

$$(iii) \text{ Gini}(\text{salary} \leq 75) = \frac{3}{4} (1 - ((\frac{2}{3})^2 + (\frac{1}{3})^2)) + \frac{1}{4} (1 - (1)^2) \\ = \frac{1}{3} = 0.333$$

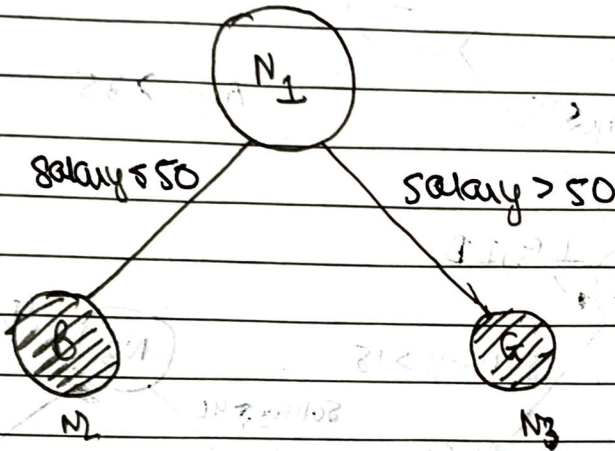
$$(iv) \text{ Gini}(\text{salary} \leq 100) = \frac{4}{4} (1 - ((\frac{2}{4})^2 + (\frac{2}{4})^2)) + 0 \\ = \frac{3}{8} = 0.375$$

\therefore Best split = salary ≤ 60 (split point $\neq 100$)

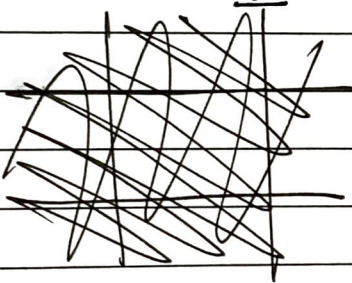


Q3) Accuracy of classifier at every node using histograms.

Ideal Classification Tree



Histogram for N_1 before updating:



	B	G	
L	0	0	→ N_1
R	2	4	

After classification (salary ≤ 50)

	B	G
L	2	0
R	0	4

N_2 :

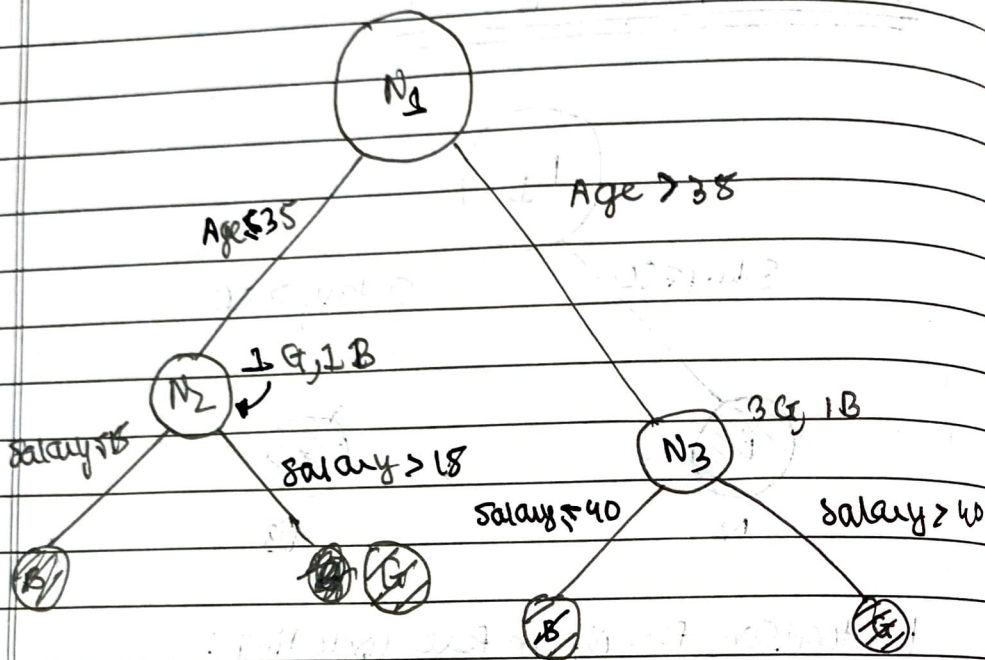
	B	G
L	-	-
R	2	0

N_3 :

	B	G
L	-	-
R	0	4

Classification Tree using root node

Age ≤ 35



class histograms

① N_2 after classification:

	B	G
$N_2 \leftarrow L$	1	1
$N_3 \leftarrow R$	1	3

$$\begin{aligned}
 \text{Accuracy to get } N_2 &= \frac{\text{number of B's in } N_2 \times 100}{\text{expected no of Bs}} \\
 &= \frac{1 \times 100}{2} = 50\%
 \end{aligned}$$

$$\begin{aligned}
 \text{Accuracy to get } N_3 &= \frac{\text{number of G's in } N_3 \times 100}{\text{expected number of G}} \\
 &= \frac{3 \times 100}{4} \\
 &= 75\%
 \end{aligned}$$

Accuracy of Nodes N_2, N_3 are 100% since they give the perfect split i.e. one class type on the left and other on the right.

~~When~~ When classifying using root node as N_1 , we do not get the perfect split as expected using the SLTD model. The accuracy is calculated for both L and R (i.e. ~~the~~ how accurately does it identify B in left Node N_2 and G in right node N_3 .)