

Reading Assignment

Question 1

$$(a) P_+ = \frac{120}{200} = \frac{3}{5}$$

$$P_- = \frac{80}{200} = \frac{2}{5}$$

$$\text{Gini Index} = 1 - \sum_{i=1}^n (P_i)^2$$

$$= 1 - \frac{9}{25} - \frac{4}{25}$$

$$= \frac{12}{25} = \boxed{0.48}$$

(b) Left (50 positive, 10 negative) \rightarrow 60 samples

Right (70 positive, 70 negative) \rightarrow 140 samples

$$\text{Gini}_{\text{left}} = 1 - \frac{25}{36} - \frac{1}{36} = \frac{10}{36} = 0.2778$$

$$\text{Gini}_{\text{right}} = 1 - \frac{1}{4} - \frac{1}{4} = \frac{1}{2} = 0.5$$

$$\text{Gini}_{\text{weighted}} = \frac{\overset{+50}{60} \times \frac{10}{3636} + \overset{+70}{140} \times \frac{1}{21}}{200} = \frac{13}{30} = 0.4333$$

Since Gini index decreased after the split, purity is improved

Question 2

(a) Possible splits: 1.5, 2.5, 3.5, 4.5, 6.5, 7.5

Split Point	Left set	Right set	SSE
1.5	10	12, 15, 18, ..., 30	$0 + 271.45 = 271.45$
2.5	10, 12	15, 18, ..., 30	$2 + 170.83 = 172.83$
3.5	10, 12, 15	18, 21, 25, 28, 30	$12.67 + 97.2 = 109.87$
4.5	10, 12, 15, 18	21, 25, 28, 30	$36.75 + 46 = 82.75$
5.5	10, 12, 15, 18, 21	25, 28, 30	$78.8 + 12.67 = 91.47$
6.5	10, 12, 15, 18, 21, 25	28, 30	$158.83 + 2 = 160.83$
7.5	10, 12, 15, 18, 21, 25, 28	30	265.94

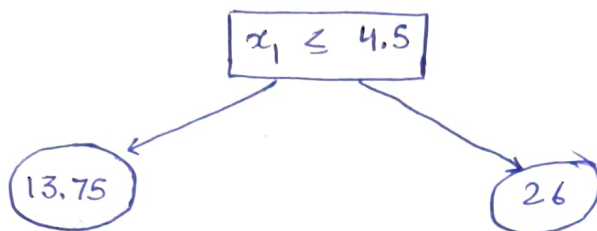
Best split: $x_1 \leq 4.5 \Rightarrow \text{SSE} = 82.75$

(b) Left node : $x_1 \leq 4.5$ Mean, $\bar{y}_{LHS} = \frac{10+12+15+18}{4} = 13.75$

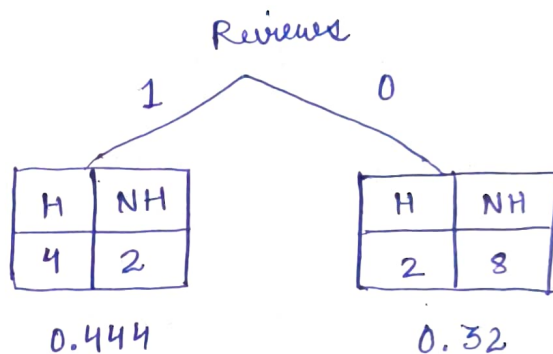
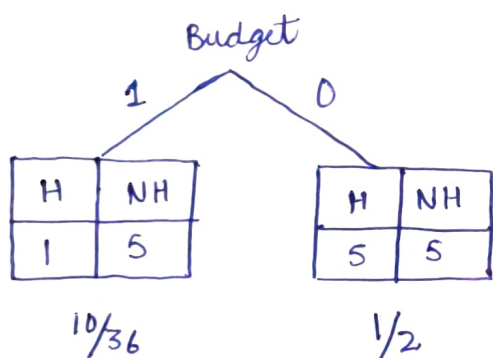
SSE = 36.75

Right node : $x_1 > 4.5$ Mean, $\bar{y}_{RHS} = \frac{21+25+28+30}{4} = 26$

SSE = 46



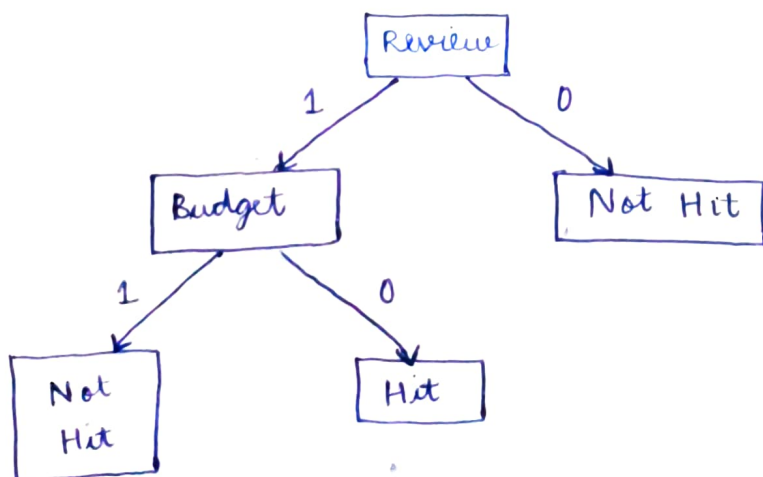
Question 3



Gini_{weighted} = 0.4687

Gini_{weighted} = 0.366 < Gini_{weighted} (Budget)

Final Tree



Question 4

(a) For whole dataset, $H(x) = - \left(\frac{6}{10} \log_2 \frac{6}{10} + \frac{4}{10} \log_2 \frac{4}{10} \right) = 0.971$

(i) Attribute : Price

Low

4 samples \rightarrow 3 yes, 1 no

$$\text{Entropy} = - \left(\frac{3}{4} \log_2 \frac{3}{4} + \frac{1}{4} \log_2 \frac{1}{4} \right) \\ = 0.811$$

Med

4 \rightarrow 1 yes, 3 no

$$\text{Entropy} = 0.811$$

High

2 \rightarrow 2 yes

$$\text{Entropy} = 0$$

$$\text{Weighted entropy, } E_{\text{price}} = \frac{4}{10} \times 0.811 + \frac{4}{10} \times 0.811 + \frac{2}{10} \times 0 = 0.6488$$

$$\text{Information Gain, } IG_{\text{price}} = 0.971 - 0.6488 = 0.3222$$

(ii) Attribute : Maintenance

Low

2 \rightarrow 2 yes

$$\text{Entropy} = 0$$

Med

4 \rightarrow 2 yes, 2 no

$$\text{Entropy} = 1$$

High

4 \rightarrow 2 yes, 2 no

$$\text{Entropy} = 1$$

$$E_{\text{maintenance}} = \frac{2}{10} \times 0 + \frac{4}{10} \times 1 + \frac{4}{10} \times 1 = 0.8$$

$$IG_{\text{maintenance}} = 0.971 - 0.8 = 0.171$$

(iii) Attribute : Capacity

Low 2

2 \rightarrow 2 yes

$$\text{Entropy} = 0$$

Med 4

6 \rightarrow 3 yes, 3 no

$$\text{Entropy} = 1$$

High 5

2 \rightarrow 1 yes, 1 no

$$\text{Entropy} = 1$$

$$E_{\text{capacity}} = \frac{2}{10} \times 0 + \frac{6}{10} \times 1 + \frac{2}{10} \times 1 = 0.8$$

$$IG_{\text{capacity}} = 0.971 - 0.8 = 0.171$$

(iv) Attribute : Airbag

Yes

5 \rightarrow 2 yes, 3 no

$$\text{Entropy} = 0.971$$

No

5 \rightarrow 4 yes, 1 no

$$\text{Entropy} = 0.722$$

$$E_{\text{airbag}} = \frac{5}{10} \times 0.971 + \frac{5}{10} \times 0.722 = 0.8465$$

$$IG_{\text{airbag}} = 0.971 - 0.8465 = 0.1245$$

Best attribute to split on $\rightarrow \text{Max}^m IG$

\Rightarrow Price is the root node

(b) To decide next split, we will

- recalculate entropy with each child node
- calculate IG with remaining attributes
- choose the one with highest IG as next split

Group 1 : Price = Low

Profitable : 3 yes, 1 no \rightarrow Entropy = 0.811

Maintenance \rightarrow Low : 2 yes \rightarrow Entropy = 0

\rightarrow Med : 0 yes \rightarrow Entropy = 0

\rightarrow High : 1 yes \rightarrow Entropy = 0

$$IG = 0.811$$

Best split \rightarrow Maintenance

Likewise,

Group 2 : Price = Med

Profitable : $H(X) = 0.811$

$$IG_{\text{airbag}} = 0.311$$

$$IG_{\text{maintenance}} = 0.311$$

Best split \rightarrow Either airbag or maintenance

Group 3 : Price = High

All are profitable , Entropy = 0

No split required

