

## Decision tree

(q1) A company wants to predict whether a customer will subscribe to a premium membership based on their demographic and browsing behaviour data. Use the decision tree to predict the given example

Age	Gender	Income	Browsing time	Subscription
20-30	Male	High	10am - 12pm	Yes
20-30	Female	Medium	2pm - 4pm	Yes
30-40	Male	Low	8am - 10am	No
30-40	Female	High	4pm - 6pm	Yes
>40	Male	Medium	6pm - 8pm	Yes
>40	Female	Medium	8am - 10am	No
>40	Male	High	12pm - 2pm	Yes
20-30	Female	Low	10am - 12pm	No
20-30	Male	Medium	2 pm - 4 pm	Yes
30-40	Female	High	8 am - 10 am	Yes

$$\Rightarrow S = [7+, 3-]$$

$$\text{Entropy } (S) = - \frac{\text{No of Yes}}{\text{Total}} \log_2 \frac{\text{No of yes}}{\text{total}} - \frac{\text{No of no}}{\text{Total}} \log_2 \frac{\text{No of no}}{\text{total}}$$

$$= - \frac{7}{10} \log_2 \frac{7}{10} - \frac{3}{10} \log_2 \frac{3}{10} = 0.881$$

Attribute: Gender

$$S_{\text{Male}} = [4+, 1-]$$

$$\text{Entropy } (S_{\text{Male}}) = - \frac{4}{5} \log_2 \frac{4}{5} - \frac{1}{5} \log_2 \frac{1}{5} = 0.721$$

$$S_{\text{female}} = [3+, 2-]$$

$$\text{Entropy}(S_{\text{female}}) = -\frac{3}{5} \log_2 \frac{3}{5} - \frac{2}{5} \log_2 \frac{2}{5}$$

$$= 0.971$$

$$\text{Gain}(S, \text{Gender}) = \text{Entropy}(S) - \sum \frac{|S_v|}{|S|} \text{Entropy}(S_v)$$

$$= 0.881 - \frac{5}{10} \text{Entropy}(S_{\text{male}}) - \frac{5}{10} \text{Entropy}(S_{\text{female}})$$

$$= 0.881 - \frac{5}{10} \times 0.721 - \frac{5}{10} \times 0.971 = 0.035$$

Attribute: age

$$S_{20-30} = [3+, 1-] \quad S_{30-40} = [2+, 1-]$$

$$\text{Entropy}(S_{20-30}) = -\frac{3}{4} \log_2 \frac{3}{4} - \frac{1}{4} \log_2 \frac{1}{4} = 0.811$$

$$\text{Entropy}(S_{30-40}) = -\frac{2}{3} \log_2 \frac{2}{3} - \frac{1}{3} \log_2 \frac{1}{3} = 0.918$$

$$S_{>40} = [1+, 1-]$$

$$\text{Entropy}(S_{>40}) = -\frac{1}{2} \log_2 \frac{1}{2} - \frac{1}{2} \log_2 \frac{1}{2} = 1$$

$$\text{Gain} = \text{Entropy} - \sum \frac{|S_v|}{|S|} \text{Entropy}(S_v)$$

$$= 0.881 - \frac{4}{10} \text{Entropy}(S_{20-30}) - \frac{3}{10} \text{Entropy}(S_{30-40})$$

$$- \frac{2}{10} \text{Entropy}(S_{>40})$$

$$= 0.881 - \frac{4}{10} \times 0.811 - \frac{3}{10} \times 0.918 - \frac{2}{10} \times 1 = 0.0812$$

Attribute : Income

$$\text{Gain } S_{\text{High}} = [4+, 0-]$$

$$S_{\text{medium}} = [3+, 1-]$$

$$S_{\text{low}} = [0+, 2-]$$

$$\text{Entropy } (S_{\text{High}}) = 0$$

$$\text{Entropy } (S_{\text{medium}}) = 0.811$$

$$\text{Entropy } (S_{\text{low}}) = 0$$

$$\text{Gain} = 0.881 - 0.811 \times \frac{3}{10} = 0.6377$$

Attribute : Browsing time

$$S_{10\text{am}-12\text{pm}} = [1+, 1-]$$

$$S_{2\text{pm}-4\text{pm}} = [2+, 0-]$$

$$S_{8\text{am}-10\text{am}} = [1+, 2-]$$

$$S_{4\text{pm}-6\text{pm}} = [1+, 0-]$$

$$S_{6\text{pm}-8\text{pm}} = [1+, 0-]$$

$$S_{12\text{pm}-2\text{pm}} = [1+, 0-]$$

$$\text{Entropy } (S_{10\text{am}-12\text{pm}}) = 0$$

$$\text{Entropy } (S_{2\text{pm}-4\text{pm}}) = 0$$

$$\text{Entropy } (S_{8\text{am}-10\text{am}}) = 0.918$$

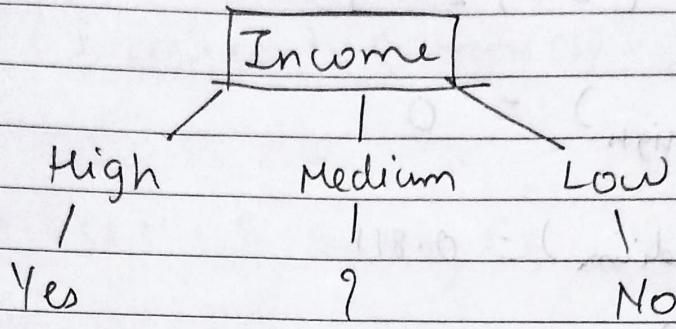
$$\text{Entropy } (S_{4\text{pm}-6\text{pm}}) = 0$$

$$\text{Entropy } (S_{6\text{pm}-8\text{pm}}) = 0$$

$$\text{Entropy } (S_{12\text{pm}-2\text{pm}}) = 0$$

$$\text{Gain} = 0.881 - \frac{3}{10} \times 0.918 = 0.6056$$

We will consider income as root node



Age	Gender	Browsing time	Subscription
20-30	Female	2pm - 4pm	Yes
>40	Male	6pm - 8pm	Yes
>40	Female	8am - 10am	No
20-30	Male	2pm - 4pm	Yes

$$S = [3+, 1-]$$

$$S = -\frac{3}{4} \log_2 \frac{3}{4} - \frac{1}{4} \log_2 \frac{1}{4} = 0.811$$

Attribute : Age

$$S_{20-30} = [2+, 0-]$$

$$S_{>40} = [1+, 1-]$$

$$\text{Entropy } (S_{20-30}) = 0$$

$$\text{Entropy } (S_{>40}) = 1$$

$$\text{Gain} = 0.811 - 1 \times \frac{2}{5} = 0.411$$

Attribute : Gender

$$S_{\text{male}} = [2+, 0-]$$

$$S_{\text{female}} = [1+, 1-]$$

$$\text{Entropy } (S_{\text{male}}) = 0$$

$$\text{Entropy } (S_{\text{female}}) = 1$$

$$\text{Gain} = 0.811 - 1 \times \frac{2}{5} = 0.411$$

Attribute : Browsing time

$$S_{2 \text{pm} - 4 \text{pm}} = [2+, 0-]$$

$$S_{6 \text{pm} - 8 \text{pm}} = [1+, 0-]$$

$$S_{8 \text{am} - 10 \text{am}} = [0+, 1-]$$

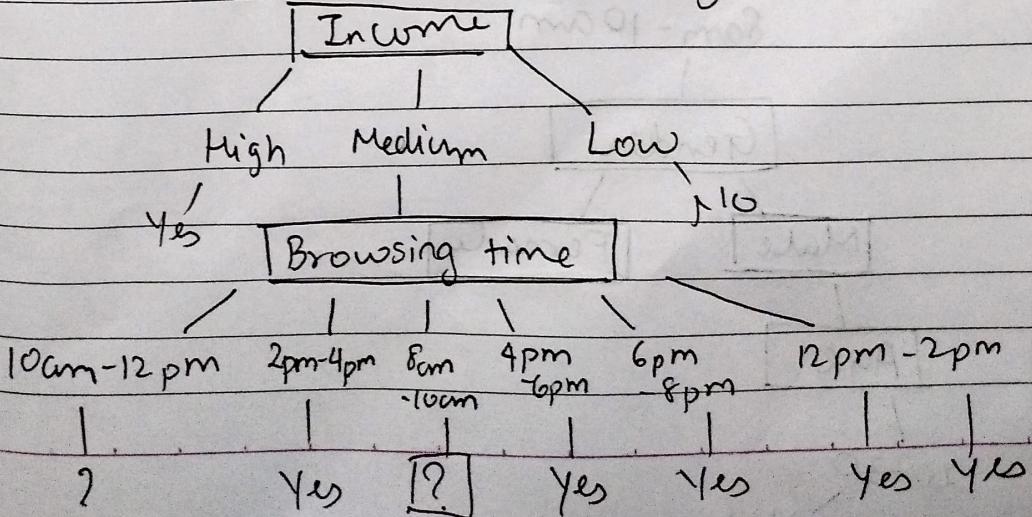
$$\text{Entropy } (S_{2 \text{pm} - 4 \text{pm}}) = 0$$

$$\text{Entropy } (S_{6 \text{pm} - 8 \text{pm}}) = 0$$

$$\text{Entropy } (S_{8 \text{am} - 10 \text{am}}) = 0$$

$$\text{Gain} = 0.811 - 0 = 0.811$$

The root node becomes browsing time



For 10am - 12pm

Age	Gender	Subscription
20-30	Male	Yes
20-30	Female	No

Attribute : Age

$$E(S_{20-30}) = 1$$

$$\text{Gain} = 1 - 1 \times 1 = 1 - 1 = 0$$

$$\text{Gain} = 0$$

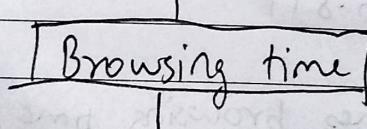
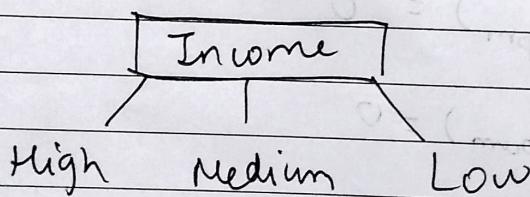
Attribute : Gender

$$E(S_{\text{Male}}) = 0$$

$$E(S_{\text{Female}}) = 0$$

$$\text{Gain} = 1$$

Considering gender as root node



Gender

Male

Female

Age

20-30

yes