

(*) Construct Decision tree for the following transactions using ID3

TID	Refund	Marital Status	Taxable income	Cheat
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes

Let Consider attribute Taxable income

Sort The ^{Continuous} Values of attribute Taxable income

60K \rightarrow No

70K \rightarrow No

75K \rightarrow No } Change in the
85K \rightarrow Yes } Target Concept
Value

90K \rightarrow Yes

95K \rightarrow Yes } Change in the
100K \rightarrow No } Target Concept
Value

120K \rightarrow No

125K \rightarrow No

220K \rightarrow No

In this,

We should pick a Threshold, c , that produces the greatest Information gain.

\therefore Identifying adjacent examples that differ in their target classification.

From this, ~~Good~~ Thresholds are

$$C_1 = \frac{75K + 85K}{2} = 80$$

$$C_2 = \frac{95K + 100K}{2} = 97.5$$

Now,

We have to Calculate Information

Gain on attribute Taxable Income

Let these Threshold

$$C_1 = 80$$

$$C_2 = 97.5$$

Let's consider C_1

Before that, Calculate Entropy(D).

$$\text{Entropy(D)} = - \sum_{i=1}^2 p_i \log_2 p_i$$

(Info(D))

$$= - \left[\left(\frac{3}{10} \right) \log_2 \left(\frac{3}{10} \right) + \left(\frac{7}{10} \right) \log_2 \left(\frac{7}{10} \right) \right]$$

$$= \underline{\underline{0.88 \text{ bits}}}$$

Info (D)

Taxable income

wrt $C_1, >8015$

$$= \sum_{j=1}^2 \frac{|D_j|}{|D|} \text{Info}(D_j)$$

$$= \left(\frac{7}{10}\right) \text{Info}(D_{C_1, >8015}) + \left(\frac{3}{10}\right) \text{Info}(D_{C_1, \leq 8015})$$

~~$$\left(\frac{7}{10}\right) \left[-\left(\frac{3}{7}\right) \log_2 \left(\frac{3}{7}\right) \right]$$~~

$$= \frac{7}{10} \left[-\frac{3}{7} \log_2 \frac{3}{7} - \frac{4}{7} \log_2 \frac{4}{7} \right]$$

$$+ \frac{3}{10} \left[-\frac{0}{3} \log_2 \left(\frac{0}{3}\right) - \frac{3}{3} \log_2 \left(\frac{3}{3}\right) \right]$$

$$= \underline{\underline{0.68}}$$

$$\text{Gain} \left(\begin{array}{l} \text{Taxable income} \\ \text{wrt } C_1, >8015 \end{array} \right) = \begin{array}{l} \text{Info}(D) \\ - \text{Info}(D) \\ \text{Taxable income} \\ \text{wrt } C_1, >8015 \end{array}$$

$$= 0.88 - 0.68$$

$$= \underline{\underline{0.20}}$$

Info(D)

Taxable income

w.r.t $C_2 > 97.5$

$$= \phi \sum_{j=1}^2 \frac{|D_j|}{|D|} \text{Info}(D_j)$$

$$= \left(\frac{4}{10}\right) \left[-\left(\frac{0}{4}\right) \log_2\left(\frac{0}{4}\right) - \left(\frac{4}{4}\right) \log_2\left(\frac{4}{4}\right) \right]$$

$$+ \left(\frac{6}{10}\right) \left[-\left(\frac{3}{6}\right) \log_2\left(\frac{3}{6}\right) - \left(\frac{3}{6}\right) \log_2\left(\frac{3}{6}\right) \right]$$

$$= \underline{\underline{0.6}}$$

Gain (Taxable income)

w.r.t $C_2 > 97.5$

$$= 0.88 - 0.6$$

$$= \underline{\underline{0.28}}$$

Next Consider attribute Marital Status

$$\text{Info}(D)_{\text{Marital Status}} = \sum_{j=1}^3 \frac{|D_j|}{|D|} \text{Info}(D_j)$$

$$= \frac{|D_{\text{single}}|}{|D|} \text{Info}(D_{\text{single}})$$

$$+ \frac{|D_{\text{married}}|}{|D|} \text{Info}(D_{\text{married}})$$

$$+ \frac{|D_{\text{divorced}}|}{|D|} \text{Info}(D_{\text{divorced}})$$

$$= \left(\frac{4}{10}\right) \left[-\left(\frac{2}{4}\right) \log_2\left(\frac{2}{4}\right) - \frac{2}{4} \log_2\left(\frac{2}{4}\right) \right]$$

$$+ \left(\frac{4}{10}\right) \left[-\left(\frac{0}{4}\right) \log_2\left(\frac{0}{4}\right) - \frac{4}{4} \log_2\left(\frac{4}{4}\right) \right]$$

$$+ \left(\frac{2}{10}\right) \left[-\left(\frac{1}{2}\right) \log_2\left(\frac{1}{2}\right) - \frac{1}{2} \log_2\left(\frac{1}{2}\right) \right]$$

$$= 0.4 + 0.2$$

$$= \underline{\underline{0.6}}$$

$$\begin{aligned} \text{En Gain (Marital Status)} \\ &= 0.88 - 0.6 \\ &= \underline{\underline{0.28}} \end{aligned}$$

Next Consider attribute Refund

Info(D)

Refund

$$= \sum_{j=1}^2 \frac{|D_j|}{|D|} \text{Info}(D_j)$$

$$= \frac{|D_{\text{yes}}|}{|D|} \text{Info}(D_{\text{yes}})$$

$$+ \frac{|D_{\text{no}}|}{|D|} \text{Info}(D_{\text{no}})$$

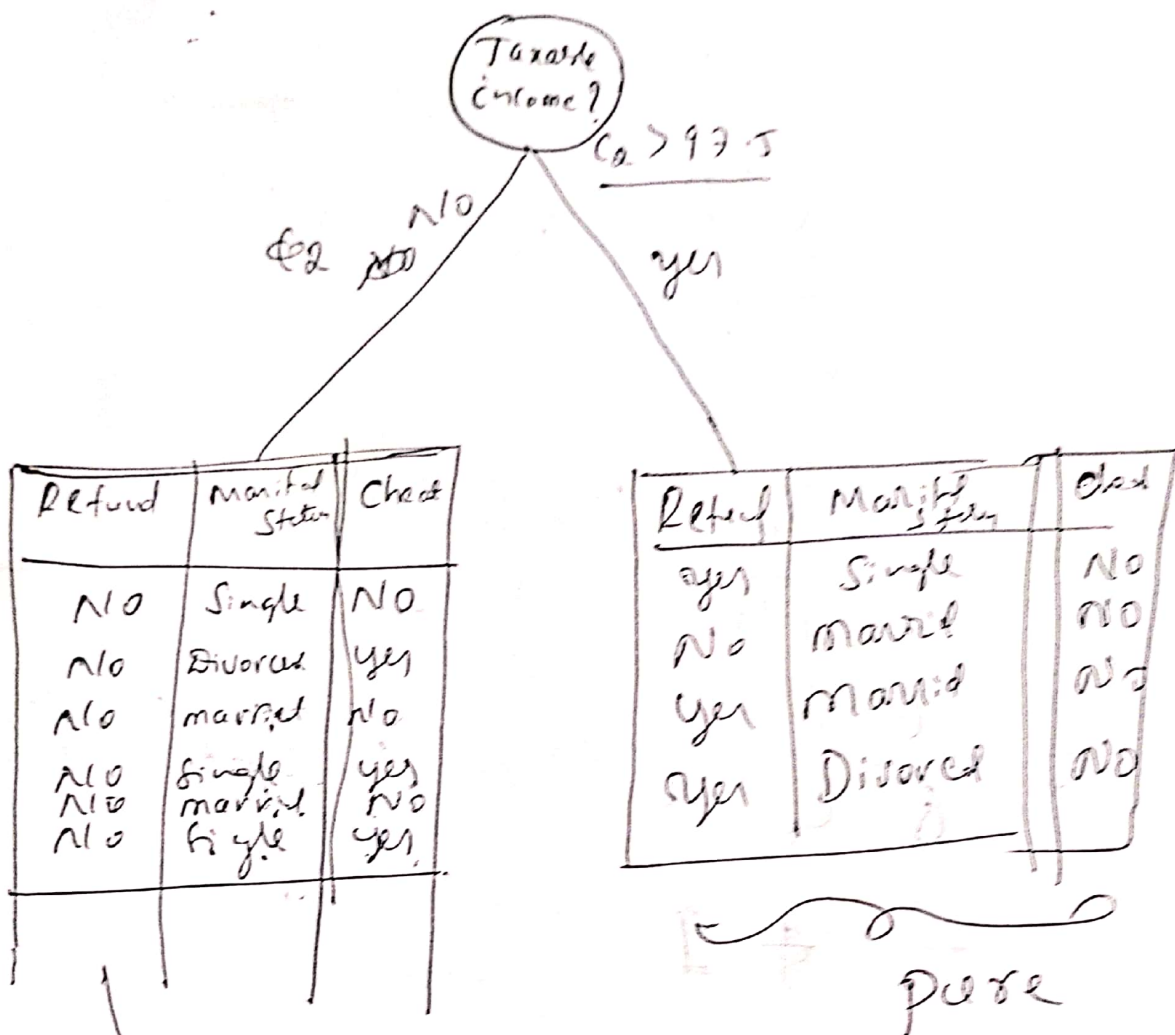
$$= \left(\frac{3}{10}\right) \left[-\left(\frac{0}{3}\right) \log_2\left(\frac{0}{3}\right) - \left(\frac{3}{3}\right) \log_2\left(\frac{3}{3}\right) \right]$$

$$+ \left(\frac{7}{10}\right) \left[-\left(\frac{3}{7}\right) \log_2\left(\frac{3}{7}\right) - \left(\frac{4}{7}\right) \log_2\left(\frac{4}{7}\right) \right]$$

$$= \underline{\underline{0.68}}$$

$$\text{Gain(Refund)} = 0.88 - 0.68$$

$$= \underline{\underline{0.20}}$$



Since it is impure, Split

$$b \text{ Info}(D)$$

$$= - \sum_{i=1}^2 p_i \log_2 p_i$$

$$= - \left[\frac{3}{6} \log_2 \left(\frac{3}{6} \right) + \left(\frac{3}{6} \right) \log_2 \left(\frac{3}{6} \right) \right]$$

$$= 1$$

29/11/20

$$\text{Info}(D) \\ \text{Refund} = \sum_{j=1}^v \frac{|D_j|}{|D|} \text{Info}(D_j)$$

$$= \frac{|D_{\text{No}}|}{|D|} \text{Info}(D_{\text{No}})$$

$$= \frac{6}{6} \left[-\frac{3}{6} \log_2\left(\frac{3}{6}\right) - \frac{3}{6} \log_2\left(\frac{3}{6}\right) \right]$$

$$= 1$$

$$\text{Gain}(\text{Refund}) = 1 - 1 \\ = 0$$

$$\text{Info}(D)_{\text{market status}} = \sum_{j=1}^v \frac{|D_j|}{|D|} \text{Info}(D_j)$$

$$= \frac{3}{6} \left[-\frac{2}{3} \log_2\left(\frac{2}{3}\right) - \left(\frac{1}{3}\right) \log_2\left(\frac{1}{3}\right) \right]$$

$$+ \frac{2}{6} \left[-\frac{0}{2} \log_2\left(\frac{0}{2}\right) - \left(\frac{2}{2}\right) \log_2\left(\frac{2}{2}\right) \right]$$

$$+ \frac{1}{6} \left[-\frac{1}{1} \log_2\left(\frac{1}{1}\right) - \frac{0}{1} \log_2\left(\frac{0}{1}\right) \right]$$

$$= 0.4591$$

Ex Gain (marital status)

$$= 1 - 0.4591$$

$$= 0.540$$

