

Data Mining
Assignment A4

Sindhuja Yerramalla
U00839259.

Question 1

(a) Entropy of Target Attribute "Buy Bitcoin"

$$Info(T) = P_+ \log P_+ - P_- \log P_-$$

$$= -\frac{4}{8} \log \frac{4}{8} - \frac{4}{8} \log \frac{4}{8}$$

$$\boxed{Info(T) = 1}$$

for attribute CS-Major

Entropy of "yes" given Target attribute (T)

$$Info(Yes) = -\frac{3}{4} \log \frac{3}{4} - \frac{1}{4} \log \frac{1}{4}$$

$$\boxed{Info(Yes) = 1}$$

Entropy of "no" given Target attribute (T)

$$Info(No) = -\frac{3}{4} \log \frac{3}{4} - \frac{1}{4} \log \frac{1}{4}$$

$$\boxed{Info(No) = 1}$$

Entropy of CS-Major given target as attribute (T)

$$\Rightarrow Info(CS-Major, T) = P_{Yes} \times (Info(Yes)) + P_{No} (Info(No))$$

$$= \frac{4}{8}(1) + \frac{4}{8}(1) = 1$$

$$\boxed{Info(CS-Major, T) = 1}$$

Info Gain of 'cs-major' given Target attribute (T)

$$\text{Gain}(\text{cs-major}, T) = \frac{(\text{Info}(T) - \text{Info}(\text{cs-major}, T))}{\text{Split Info}(\text{cs-major})}$$

$$\text{Split Info}(\text{cs-major}) = -\frac{4}{8} \log \frac{4}{8} - \frac{4}{8} \log \frac{4}{8}$$

$$\boxed{\text{Split Info}(\text{cs-major}) = 1}$$

$$\text{Gain}(\text{cs-major}, T) = \frac{1-1}{1} = 0 \%$$

For attribute Age

Entropy of "young" given Target attribute (T)

$$\text{Info}(T_{\text{young}}) = -\frac{3}{4} \log \frac{3}{4} - \frac{3}{4} \log \frac{2}{4}$$

$$\boxed{\text{Info}(T_{\text{young}}) = 1}$$

Entropy of "old" given Target attribute (T)

$$\text{Info}(T_{\text{old}}) = -\frac{1}{2} \log \frac{1}{2} - \frac{1}{2} \log \frac{1}{2}$$

$$\boxed{\text{Info}(T_{\text{old}}) = 1}$$

Entropy of "middle" given Target attribute (T)

$$\text{Info}(T_{\text{middle}}) = -\frac{1}{2} \log \frac{1}{2} - \frac{1}{2} \log \frac{1}{2}$$

$$\boxed{\text{Info}(T_{\text{middle}}) = 1}$$

Entropy of "Age" given Target attribute (T)

$$\text{Info}(\text{Age}) = P_{\text{old}}(\text{Info}(T_{\text{old}})) + P_{\text{middle}}(\text{Info}(T_{\text{middle}})) + P_{\text{young}}(\text{Info}(T_{\text{young}}))$$

$$= \frac{2}{8} + \frac{2}{8} + \frac{4}{8}$$

$$= \frac{1}{4} + \frac{1}{4} + \frac{1}{2} = 1$$

$$\therefore \text{Info}(\text{Age}) = 1$$

$$\text{Split Info}(\text{Age}) = -\frac{2}{8} \log \frac{2}{8} - \frac{2}{8} \log \frac{2}{8} - \frac{4}{8} \log \frac{4}{8}$$

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

$$\boxed{\text{Split Info}(\text{Age}) = 1.5}$$

$$\therefore \text{Gain}(\text{Age}, T) = \frac{(\text{Info}(T) - \text{Info}(\text{Age}, T))}{\text{Split Info}(\text{Age})}$$

$$= \frac{1-1}{1.5} = 0$$

$$\boxed{\text{Gain}(\text{Age}, T) = 0}$$

for attribute "Income"

Entropy of "high", given target attribute (T)

$$\text{Info}(T_{\text{high}}) = -\frac{1}{1} \log \frac{1}{1} = 0$$

$$\boxed{\text{Info}(T_{\text{high}}) = 0}$$

Entropy of "fair" given target attribute (T)

$$\begin{aligned} \text{Info}(T_{\text{fair}}) &= -\frac{3}{4} \log \frac{3}{4} - \frac{1}{4} \log \frac{1}{4} \\ &= 0.311 + 0.5 \end{aligned}$$

$$\boxed{\text{Info}(T_{\text{fair}}) = 0.811}$$

Entropy of "low" given target attribute (T)

$$\text{Info}(T_{\text{low}}) = -\frac{3}{3} \log \frac{3}{3} = 0$$

$$\boxed{\text{Info}(T_{\text{low}}) = 0}$$

∴ Entropy of "Income" given target attribute (T)

$$\text{Info}(\text{Income}, T) = P_{\text{fair}} \text{Info}(T_{\text{fair}}) + P_{\text{high}} \text{Info}(T_{\text{high}}) + P_{\text{low}} \text{Info}(T_{\text{low}})$$

$$= \frac{4}{8} \times 0.811 + \frac{1}{8} \times 0 + \frac{3}{8} \times 0$$

$$\boxed{\therefore \text{Info}(\text{Income}, T) = 0.406}$$

$$\begin{aligned} \text{Split Info} &= -\frac{4}{8} \log \frac{4}{8} - \frac{1}{8} \log \frac{1}{8} - \frac{3}{8} \log \frac{3}{8} \\ &= 0.5 + 0.375 + 0.531 \end{aligned}$$

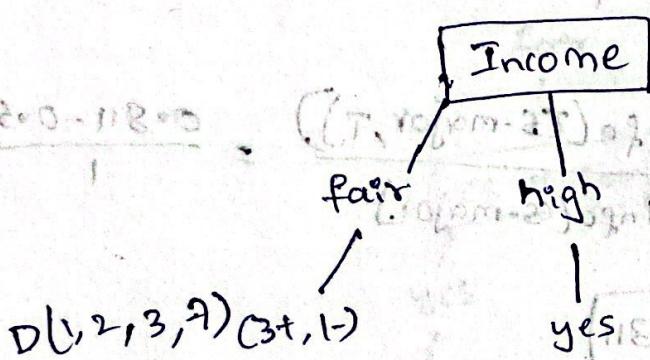
$$\boxed{\text{split Info}(\text{Income}) = 1.406}$$

$$\text{Gain}(\text{Income}, T) = \frac{\text{Info}(T) - \text{Info}(\text{Income}, T)}{\text{splitInfo}(\text{Income})}$$

$$= \frac{1 - 0.406}{1.406}$$

$$\boxed{\text{Gain}(\text{Income}, T) = 0.422}$$

\therefore Gain of (Income, T) is greater than other two attribute gains, so, Income will be taken as the root node of our decision Tree.



Now consider Records D_1, D_2, D_3, D_4

$$\begin{aligned}\text{Info}(T) &= -P_1 \log P_1 - P_2 \log P_2 \\ &= -\frac{3}{4} \log \frac{3}{4} - \frac{1}{4} \log \frac{1}{4} \\ &= 0.811\end{aligned}$$

$$\boxed{\therefore \text{Info}(T) \approx 0.811}$$

For attribute "cs-Major"

$$\text{Info}(T_{yes}) = -\frac{2}{2} \log \frac{2}{2} = 0$$

$$\text{Info}(T_{no}) = -\frac{1}{2} \log \frac{1}{2} - \frac{1}{2} \log \frac{1}{2} = 1$$

$$\text{Info}(cs-major, T) = \frac{2}{4} \times 0 + \frac{2}{4} \times 1 = 0.5$$

$$\text{split Info}(cs-major) = -\frac{2}{4} \log \frac{2}{4} - \frac{2}{4} \log \frac{2}{4} = 1$$

$$\text{split Info}(cs-major) = 1$$

$$\text{Gain}(cs-major, T) = \frac{\text{Info}(T) - \text{Info}(cs-major, T)}{\text{split Info}(cs-major)} = \frac{0.811 - 0.5}{1}$$

$$\therefore \text{Gain}(cs-major, T) = 0.311$$

for attribute Age

Entropy of "Age" given Target attribute

$$\text{Info}(T_{old}) = -\frac{1}{4} \log \frac{1}{4} = 0$$

$$\text{Info}(T_{young}) = -\frac{1}{2} \log \frac{1}{2} - \frac{1}{2} \log \frac{1}{2} = 1$$

$$\text{Info}(T_{middle}) = -\frac{1}{4} \log \frac{1}{4} = 0$$

$$\Rightarrow \text{Info}(age, T) = \frac{1}{4} \times 0 + \frac{1}{4} \times 0 + \frac{2}{4} \times 1 = 0.5$$

$$\begin{aligned} \text{split Info}(Age) &= -\frac{1}{4} \log \frac{1}{4} - \frac{1}{4} \log \frac{1}{4} - \frac{2}{4} \log \frac{2}{4} \\ &= 0.1054 \end{aligned}$$

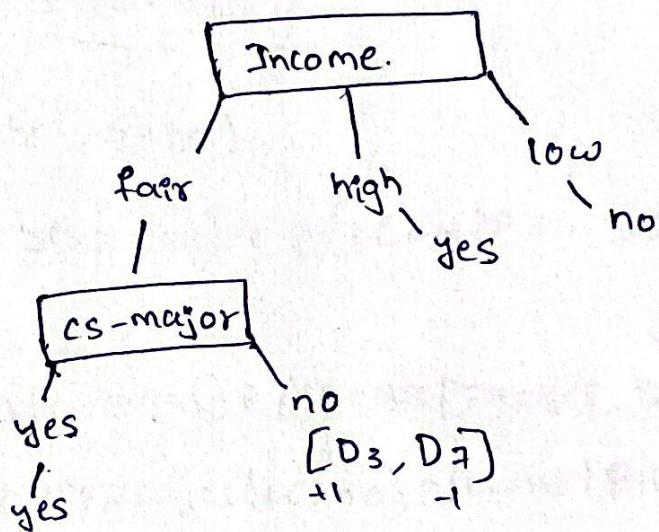
$$\text{Gain}(\text{Age}, T) = (\text{Info}(T) - \text{Info}(\text{Age}, T)) / \text{Split Info}(\text{Age})$$

$$= \frac{0.811 - 0.5}{1.5} = 0.2073$$

$$\boxed{\text{Gain}(\text{Age}, T) = 0.2073}$$

Since $\text{Gain}(\text{cs-major}, T)$ is greater than $\text{Gain}(\text{Age}, T)$, cs-major will be the next node in the tree.

Decision tree :-



b) The question is asking about the user who is studying cs-major and his income is fair. As per the decision tree we got in the previous solution, The user will buy bitcoin.

Question 2

$$a) P(LC = \text{yes}) = \sum_{x \in \{\text{yes, no}\}} \sum_{y \in \{\text{yes, no}\}} P(LC = \text{yes} | FH = x, S = y) + P(FH = x, S = y)$$

$$\therefore P(LC = \text{yes} | FH = \text{yes}, S = \text{yes}) + P(FH = \text{yes}, S = \text{yes}) + \\ P(LC = \text{yes} | FH = \text{yes}, S = \text{no}) + P(FH = \text{yes}, S = \text{no}) + \\ P(LC = \text{yes} | FH = \text{no}, S = \text{yes}) + P(FH = \text{no}, S = \text{yes}) + \\ P(LC = \text{yes} | FH = \text{no}, S = \text{no}) + P(FH = \text{no}, S = \text{no})$$

$$P(LC=yes | FH=no, S=yes) * P(FH=no, S=yes) +$$

$$P(LC=yes | FH=no, S=no) * P(FH=no, S=no)$$

$$\therefore P(LC=yes) = 0.7 \times 0.3 \times 0.6 + 0.45 \times 0.3 \times 0.4 +$$

$$0.55 \times 0.7 \times 0.6 + 0.2 \times 0.7 \times 0.4$$

$$\approx 0.467$$

$$\boxed{\therefore P(LC=yes) = 0.467}$$

b) $P(PR=yes | FH=yes, S=yes)$

$$= \sum_{x \in \{yes, no\}} P(PR=yes | LC=x) * P(LC=x | FH=yes, S=yes)$$

$$= P(PR=yes | LC=yes) P(LC=yes | FH=yes, S=yes) + P(PR=yes | LC=no) P(LC=no | FH=yes, S=yes)$$

$$= 0.85 \times 0.7 + 0.45 \times 0.3 = 0.595 + 0.135$$

$$= 0.73$$

$$\boxed{\therefore P(PR=yes | FH=yes, S=yes) = 0.73}$$

c) $P(LC=yes | PR=yes, FH=yes, S=yes)$

$$= \frac{P(PR=yes | LC=yes, FH=yes, S=yes)}{P(PR=yes | FH=yes, S=yes)} * P(LC=yes | FH=yes, S=yes)$$

$$\begin{aligned}
 &= \frac{P(PR=\text{yes} | LC=\text{yes}) P(LC=\text{yes} | FH=\text{yes}, S=\text{yes})}{\sum_{x \in \{\text{yes, no}\}} P(PR=\text{yes} | LC=x) P(LC=x | FH=\text{yes}, S=\text{yes})} \\
 &= \frac{0.85 \times 0.7}{0.85 \times 0.7 + 0.45 + 0.3} \\
 &= \frac{0.595}{0.93} = 0.815
 \end{aligned}$$

\Rightarrow The new person is more likely to have lung cancer.