

Q1

- a) We calculate the gains until reaching the conditions to stop the calculation. The highest gives us the top node in the tree.

$$C4.5: \text{Gain}(A, T) = (-\text{Info}(T) - \text{Info}(A, T)) / \text{SplitInfo}(A)$$

$$\text{SplitInfo}(A) = -\sum_{v \in A} p(v) \log p(v)$$

Step 1

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

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

$$\text{SplitInfo}(CS, T) = \frac{1}{2} \log \frac{1}{2} - \frac{1}{2} \log \frac{1}{2} = 1 \quad \text{Gain}(CS, T) = \frac{1-1}{1} = 0$$

$$\text{Info}(CS, T) = \frac{1}{2} \times 1 + \frac{1}{2} \times 1 = 1$$

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

$$\text{Info}(Age, T) = \frac{2}{8} \times 1 + \frac{2}{8} \times 1 + \frac{4}{8} \times 1 = 1 \quad \text{Gain}(Age, T) = 0$$

$$\text{SplitInfo}(Age, T) = -\frac{1}{4} \log \frac{1}{4} - \frac{1}{4} \log \frac{1}{4} - \frac{1}{2} \log \frac{1}{2} = \frac{3}{2} = 1.5$$

$$\text{Info}(T_{fair}) = -\frac{3}{4} \log \frac{3}{4} - \frac{1}{4} \log \frac{1}{4} = 0.812$$

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

$$\text{Info}(T_{low}) = -\frac{3}{4} \log \frac{3}{4} - 0 \log 0 = 0$$

$$\text{SplitInfo}(Income, T) = \frac{1}{2} \log \frac{1}{2} - \frac{1}{8} \log \frac{1}{8} - \frac{3}{8} \log \frac{3}{8} = 1.4$$

$$\text{Info}(Income, T) = \frac{1}{2} \text{Info}(T_{fair}) = 0.4615$$

$$\text{Gain}(Income, T) = \frac{1 - 0.46}{1.4} = 0.39$$

Step 2

No	CS-Major	Age	Buy Bitcoin
1	Yes	old	Yes
2	Yes	middle	Yes
3	no	Young	Yes
7	no	Young	no

$$\text{Info}(T) = -\frac{1}{4} \log \frac{1}{4} - \frac{3}{4} \log \frac{3}{4} = 0.812$$

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

$$\text{Info}(T_{young}) = -\frac{1}{3} \log \frac{1}{3} - 1 \log 1 = 0.53$$

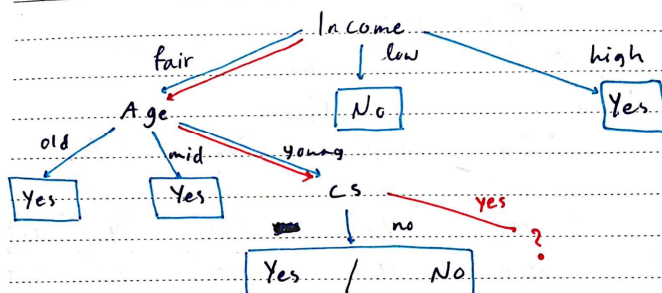
$$\text{SplitInfo}(Age, T) = -\frac{1}{2} \log \frac{1}{2} - \frac{1}{4} \log \frac{1}{4} - \frac{1}{4} \log \frac{1}{4} = 1.5$$

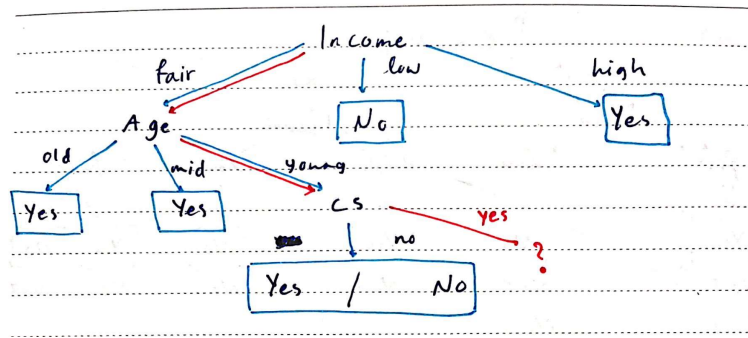
$$\text{Info}(Age, T) = -\frac{1}{2} \times 0.53 + 0 = 0.265$$

$$\text{Gain}(Age, T) = \frac{0.812 - 0.265}{1.5} = 0.364$$

$$\text{Info}(T_{CS, yes}) = -1 \log 1 - 0 \log 0 = 0$$

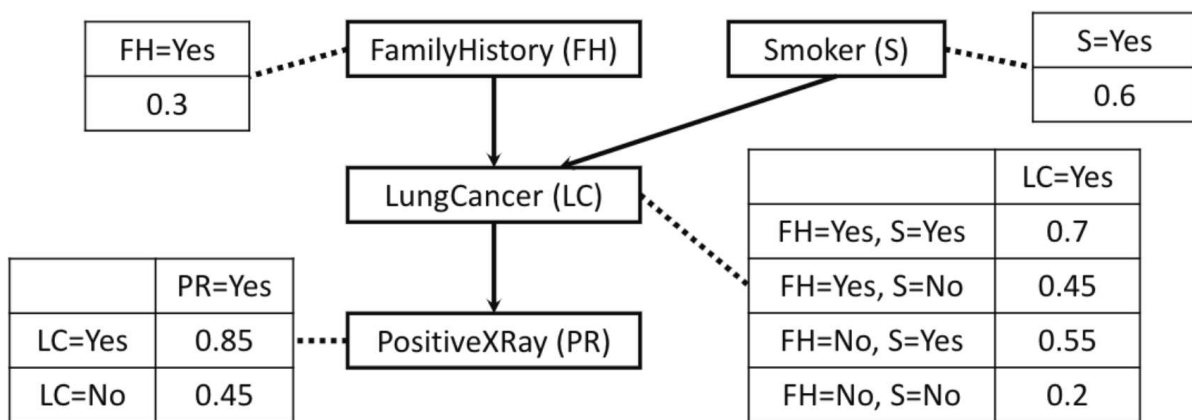
$$\text{Info}(T_{CS, no}) = -\frac{1}{2} \log \frac{1}{2} - \frac{1}{2} \log \frac{1}{2} = 1$$





- b) Looking at the tree, the arrows colored in red show the information about a young CS computer student with fair income. There is no previous example like this in the table, resulting in difficulty in predicting whether the answer is yes or no. But, according to the gain values and the probability of the young people who have a fair income, we can assume the prediction to be "Yes"

Q2



- a)
- $$\begin{aligned}
 P(LC = \text{Yes}) &= P(LC = \text{Yes} \mid FH = \text{Yes}, S = \text{Yes}) * P(FH = \text{Yes}, S = \text{Yes}) \\
 &+ P(LC = \text{Yes} \mid FH = \text{Yes}, S = \text{No}) * P(FH = \text{Yes}, S = \text{No}) \\
 &+ P(LC = \text{Yes} \mid FH = \text{No}, S = \text{Yes}) * P(FH = \text{No}, S = \text{Yes}) \\
 &+ P(LC = \text{Yes} \mid FH = \text{No}, S = \text{No}) * P(FH = \text{No}, S = \text{No}) \\
 &= 0.3 * 0.6 * 0.7 + 0.3 * 0.4 * 0.45 + 0.7 * 0.6 * 0.55 + 0.7 * 0.4 * 0.2 = 0.467
 \end{aligned}$$
- b)
- $$\begin{aligned}
 P(PR = \text{Yes} \mid FH = \text{Yes}, S = \text{Yes}) &= P(PR = \text{Yes} \mid LC = \text{Yes}) P(LC = \text{Yes} \mid FH = \text{Yes}, S = \text{Yes}) \\
 &+ P(PR = \text{Yes} \mid LC = \text{No}) P(LC = \text{No} \mid FH = \text{Yes}, S = \text{Yes}) \\
 &= 0.85 * 0.7 + 0.45 * 0.3 = 0.73
 \end{aligned}$$
- c)
- $$\begin{aligned}
 P(LC = \text{Yes} \mid PR = \text{Yes}, FH = \text{Yes}, S = \text{Yes}) &= \\
 &= \frac{P(PR = \text{Yes} \mid LC = \text{Yes}, FH = \text{Yes}, S = \text{Yes}) P(LC = \text{Yes} \mid FH = \text{Yes}, S = \text{Yes})}{P(PR = \text{Yes} \mid FH = \text{Yes}, S = \text{Yes})} \\
 &= \frac{P(PR = \text{Yes} \mid LC = \text{Yes}) P(LC = \text{Yes} \mid FH = \text{Yes}, S = \text{Yes})}{P(PR = \text{Yes} \mid LC = \text{Yes}) P(LC = \text{Yes} \mid FH = \text{Yes}, S = \text{Yes}) + P(PR = \text{Yes} \mid LC = \text{No}) P(LC = \text{No} \mid FH = \text{Yes}, S = \text{Yes})} \\
 &= 0.85 * 0.7 / (0.85 * 0.7 + (0.45 * 0.3)) = 0.8151
 \end{aligned}$$