

1. (a) Entropy = $-\sum_i p_i \log_b(p_i)$

Class 0 = $P(0) = \frac{1}{2}$

Class 1 = $P(1) = \frac{1}{2}$

Entropy = $-p(1) \log_2(p(1)) - p(0) \log_2(p(0))$

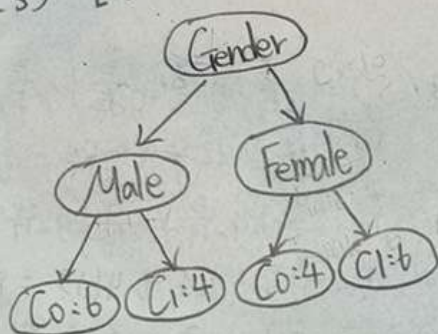
= $-\frac{1}{2} \log_2(\frac{1}{2}) - \frac{1}{2} \log_2(\frac{1}{2})$

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

= 1

(b.) Information gain (Gender)

= Entropy(S) - [(Male/S) Entropy(Male) + (Female/S) Entropy(Female)]



Entropy(Male) = $-\frac{6}{10} \log_2(\frac{6}{10}) - \frac{4}{10} \log_2(\frac{4}{10})$

= $(-0.6 \times -0.737) - (0.4 \times -1.321)$

= $0.4422 + 0.5284 = 0.9706$

Entropy(Female) = $-\frac{4}{10} \log_2(\frac{4}{10}) - \frac{6}{10} \log_2(\frac{6}{10})$

= 0.9706

=> Entropy(S) - [(Male/S) Entropy(Male) + (Female/S) Entropy(Female)]

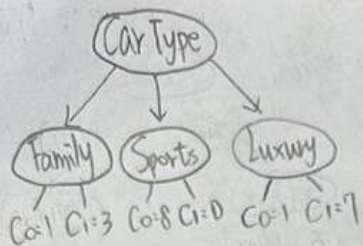
=> $1 - [\frac{10}{20} \cdot 0.9706 + \frac{10}{20} \cdot 0.9706]$

= $1 - 0.9706$

= 0.0294

\therefore Gender 對 information gain 提升不高
不太能有效地區分 data

(C.)



$$\text{Entropy (Family)} = -\frac{1}{4} \log_2 \frac{1}{4} - \frac{3}{4} \log_2 \frac{3}{4} \\ \approx 0.5 + 0.31125 = 0.81125$$

$$\text{Entropy (Sports)} = -\frac{8}{8} \log_2 \frac{8}{8} - 0 = 0$$

$$\text{Entropy (Luxury)} = -\frac{1}{8} \log_2 \frac{1}{8} - \frac{7}{8} \log_2 \frac{7}{8} \\ \approx 0.375 + 0.16852 = 0.54352$$

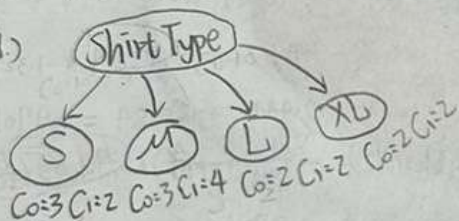
Information Gain (Car Type)

$$= 1 - \left[\frac{4}{20} \cdot 0.81125 + \frac{8}{20} \cdot 0 + \frac{2}{20} \cdot 0.54352 \right]$$

$$= 1 - [0.16225 + 0 + 0.21741]$$

$$= 0.62034$$

(d)



$$\text{Entropy (S)} = -\frac{3}{5} \log_2 \frac{3}{5} - \frac{2}{5} \log_2 \frac{2}{5} \\ \approx 0.4422 + 0.5284 = 0.9706$$

$$\text{Entropy (M)} = -\frac{3}{7} \log_2 \frac{3}{7} - \frac{4}{7} \log_2 \frac{4}{7} \\ \approx 0.5239 + 0.4613 = 0.9852$$

$$\text{Entropy (L)} = -\frac{2}{4} \log_2 \frac{2}{4} - \frac{2}{4} \log_2 \frac{2}{4} \\ = \frac{1}{2} + \frac{1}{2} = 1$$

$$\text{Entropy (XL)} = -\frac{2}{4} \log_2 \frac{2}{4} - \frac{2}{4} \log_2 \frac{2}{4} = 1$$

Information Gain (Shirt Type)

$$= 1 - \left[\frac{5}{20} \cdot 0.9706 + \frac{7}{20} \cdot 0.9852 + \frac{4}{20} \cdot 1 + \frac{4}{20} \cdot 1 \right]$$

$$= 1 - [0.24265 + 0.34482 + 0.2 + 0.2] = 0.01253$$

(e.)

比較各分法的 information gain 可得知

$$\text{Gender} = 0.0294$$

$$\text{Car Type} = 0.54352$$

$$\text{Shirt Type} = 0.01253$$

不太能有效區分 data

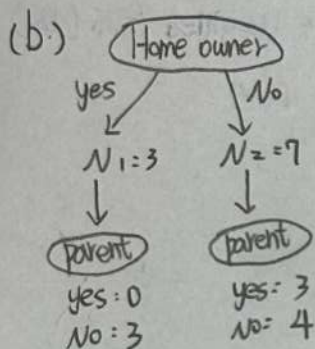
能有效區分 data

\therefore Car Type is better. ✖

2.

| Parent | |
|--------|---|
| No | 7 |
| Yes | 3 |

(a) Gini index = $1 - \left(\frac{7}{10}\right)^2 - \left(\frac{3}{10}\right)^2 = 0.42$



Node N_1 Gini index = $1 - \left(\frac{0}{3}\right)^2 - \left(\frac{3}{3}\right)^2 = 0$

Node N_2 Gini index = $1 - \left(\frac{3}{7}\right)^2 - \left(\frac{4}{7}\right)^2 = \frac{24}{49} \approx 0.49$

Weighted Gini (N_1) = $\frac{3}{10} \cdot 0 = 0$

Weighted Gini (N_2) = $\frac{7}{10} \cdot 0.49 = 0.343$

Gini (Child) = Weighted Gini (N_1) + Weighted Gini (N_2)
 $= 0 + 0.343 = 0.343$

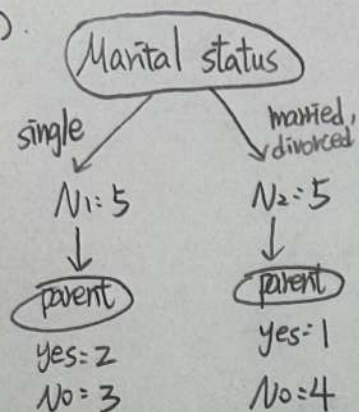
Gain in the Gini index

$= \text{Gini (Parent)} - \text{Gini (Child)}$

$= 0.42 - 0.343 = 0.077$

\therefore Home owner 不是很好的 attribute 方法,
 因為 the gain in Gini index 不高, 小於 0.2
所以沒有很好的區分 data.

(C.)



$$\text{Node } N_1 \text{ Gini index} = 1 - \left(\frac{2}{5}\right)^2 - \left(\frac{3}{5}\right)^2 = 0.48$$

$$\text{Node } N_2 \text{ Gini index} = 1 - \left(\frac{1}{5}\right)^2 - \left(\frac{4}{5}\right)^2 = 0.32$$

$$\text{Weighted Gini } N_1 = \frac{5}{10} \cdot 0.48 = 0.24$$

$$\text{Weighted Gini } N_2 = \frac{5}{10} \cdot 0.32 = 0.16$$

$$\begin{aligned} \text{Gini (Child)} &= \text{Weighted Gini } (N_1) + \text{Weighted Gini } (N_2) \\ &= 0.24 + 0.16 = 0.4 \end{aligned}$$

Gain in the Gini index

$$= \text{Gini (Parent)} - \text{Gini (Child)}$$

$$= 0.42 - 0.4$$

$$= 0.02$$

\therefore Marital status 也不是很好的 attribute 方法,
the gain in the Gini index 小於 0.2,
所以沒有很好的區分 data.