

Created a decision tree using information gain as a measure

1. Expected information

มี class 2 class $\begin{cases} c_0 \\ c_1 \end{cases}$

class P : class = "c0" $\rightarrow 10$ คน
 class N : class = "c1" $\rightarrow 10$ คน } มีทั้งหมด 20 คน

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

$$p_i = \frac{|c_{i,0}|}{|D|}$$

$$\text{Info}(D) = I(10, 10) \quad ; \quad m = 2$$

$$= -p_{c_0} \log_2(p_{c_0}) - p_{c_1} \log_2(p_{c_1})$$

$$= -\frac{10}{20} \log_2\left(\frac{10}{20}\right) - \frac{10}{20} \log_2\left(\frac{10}{20}\right) = 1$$

2. Information

มี 3 ตัว \rightarrow Gender, Car Type และ Shirt Size

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

Gender:

gd	p_i	n_i
F	4	6
M	6	4

$$\begin{aligned} \text{Info}_{\text{Gender}}(D) &= \frac{10}{20} I(4, 6) + \frac{10}{20} I(6, 4) \\ &= \frac{10}{20} \left[-\frac{4}{10} \log_2 \frac{4}{10} - \frac{6}{10} \log_2 \frac{6}{10} \right] \rightarrow \text{Female} \end{aligned}$$

$$+ \frac{10}{20} \left[-\frac{6}{10} \log_2 \frac{6}{10} - \frac{4}{10} \log_2 \frac{4}{10} \right] \rightarrow \text{Male} \quad \approx 0.97995$$

Car Type:

ct	p_i	n_i
Fam	1	3
Sport	8	0
Lux	1	7

$$\begin{aligned} \text{Info}_{\text{Car Type}}(D) &= \frac{4}{20} I(1, 3) + \frac{8}{20} I(8, 0) + \frac{8}{20} I(1, 7) \\ &= \frac{4}{20} \left[-\frac{1}{4} \log_2 \frac{1}{4} - \frac{3}{4} \log_2 \frac{3}{4} \right] \rightarrow \text{Family} \end{aligned}$$

$$+ \frac{8}{20} \left[-\frac{8}{8} \log_2 \frac{8}{8} \right] \rightarrow \text{Sports}$$

$$+ \frac{8}{20} \left[-\frac{1}{8} \log_2 \frac{1}{8} - \frac{7}{8} \log_2 \frac{7}{8} \right] \rightarrow \text{Luxury} \quad \approx 0.37968$$

Shirt Size:

ss	p_i	n_i
S	3	2
M	3	4
L	2	2
XL	2	2

$$\begin{aligned} \text{Info}_{\text{Shirt size}}(D) &= \frac{5}{20} I(3, 2) + \frac{7}{20} I(3, 4) + \frac{4}{20} I(2, 2) + \frac{4}{20} I(2, 2) \\ &= \frac{5}{20} \left[-\frac{3}{5} \log_2 \frac{3}{5} - \frac{2}{5} \log_2 \frac{2}{5} \right] \rightarrow \text{small} \end{aligned}$$

$$+ \frac{7}{20} \left[-\frac{3}{7} \log_2 \frac{3}{7} - \frac{4}{7} \log_2 \frac{4}{7} \right] \rightarrow \text{Medium}$$

$$+ \frac{4}{20} \left[-\frac{2}{4} \log_2 \frac{2}{4} - \frac{2}{4} \log_2 \frac{2}{4} \right] \rightarrow \text{Large}$$

$$+ \frac{4}{20} \left[-\frac{2}{4} \log_2 \frac{2}{4} - \frac{2}{4} \log_2 \frac{2}{4} \right] \rightarrow \text{Extra Large} \quad \approx 0.98757$$

3. Information gained

$$\text{Gain}(A) = \text{Info}(D) - \text{Info}_A(D)$$

$$\text{Gain}(\text{Gender}) = 1 - 0.97995 = 0.02005$$

$$\text{Gain}(\text{Car Type}) = 1 - 0.37968 = 0.62032 \rightarrow \text{splitting attribute}$$

$$\text{Gain}(\text{Shirt Size}) = 1 - 0.98757 = 0.01243$$

Car Type?

Family

Sport

Luxury

Gender	Shirt Size	class
M	Small	c0
M	Large	c1
M	Extra Large	c1
M	Medium	c1

Gender	Shirt Size	class
M	Medium	c0
M	Medium	c0
M	Large	c0
M	Extra Large	c0
M	Extra Large	c0
F	Small	c0
F	Small	c0
F	Medium	c0

pure

Gender	Shirt Size	class
F	Large	c0
M	Extra Large	c1
F	Small	c1
F	Small	c1
F	Medium	c1
F	Medium	c1
F	Medium	c1
F	Large	c1

Recursive

Gender	Shirt Size	class
M	Small	c0
M	Large	c1
M	Extra Large	c1
M	Medium	c1

1. Expected information

$$\text{Info}(D) = -\frac{1}{4} \log_2 \frac{1}{4} - \frac{3}{4} \log_2 \frac{3}{4} = 0.81128$$

2. Information

2 บิต → Gender และ Shirt size

Gender:

gd	P _i	n _i
F	0	0
M	1	3

$$\text{Info}_{\text{Gender}}(D) = \frac{4}{4} \left[-\frac{1}{4} \log_2 \frac{1}{4} - \frac{3}{4} \log_2 \frac{3}{4} \right] = 0.81128$$

class P: class = "c0" → 1 บิต

class N: class = "c1" → 3 บิต

Shirt Size:

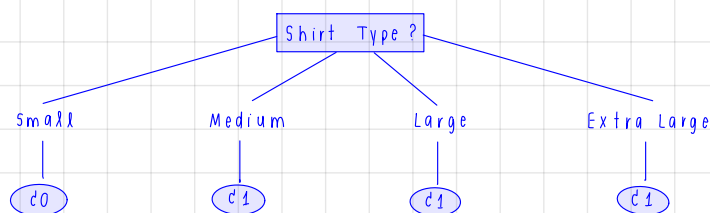
ss	P _i	n _i
S	1	0
M	0	1
L	0	1
xL	0	1

$$\begin{aligned} \text{Info}_{\text{Shirt size}}(D) &= \frac{1}{4} (-1 \log_2 1) \rightarrow \text{small} \\ &+ \frac{1}{4} (-1 \log_2 1) \rightarrow \text{Medium} \\ &+ \frac{1}{4} (-1 \log_2 1) \rightarrow \text{Large} \\ &+ \frac{1}{4} (-1 \log_2 1) \rightarrow \text{Extra Large} \\ &= 0 \end{aligned}$$

3. Information gained

$$\text{Gain}(\text{Gender}) = 0.81128 - 0.81128 = 0$$

$$\text{Gain}(\text{Shirt Size}) = 0.81128 - 0 = 0.81128 \rightarrow \text{splitting attribute}$$



Gender	shirt size	class
F	Large	c0
M	Extra Large	c1
F	small	c1
F	small	c1
F	Medium	c1
F	Medium	c1
F	Medium	c1
F	Large	c1

1. Expected information

$$\text{Info}(D) = -\frac{1}{8} \log_2 \frac{1}{8} - \frac{7}{8} \log_2 \frac{7}{8} \approx 0.54356$$

2. Information

มี 2 attribute → Gender and shirt size

Gender:

gd	P _i	n _i
F	1	6
M	0	1

$$\text{Info}_{\text{Gender}}(D) = \frac{7}{8} \left[-\frac{1}{7} \log_2 \frac{1}{7} - \frac{6}{7} \log_2 \frac{6}{7} \right] \rightarrow \text{Female}$$

$$+ \frac{1}{8} \left[-1 \log_2 1 \right] \rightarrow \text{Male} \approx 0.51771$$

class P: class = "c0" → 1 ข้อมูล

class N: class = "c1" → 7 ข้อมูล

Shirt Size:

ss	P _i	n _i
S	0	2
M	0	3
L	1	1
XL	0	1

$$\text{Info}_{\text{Shirt Size}}(D) = \frac{2}{8} (-1 \log_2 1)$$

→ small

$$+ \frac{3}{8} (-1 \log_2 1)$$

→ Medium

$$+ \frac{2}{8} \left[-\frac{1}{2} \log_2 \frac{1}{2} - \frac{1}{2} \log_2 \frac{1}{2} \right] \rightarrow \text{Large}$$

$$+ \frac{1}{8} (-1 \log_2 1)$$

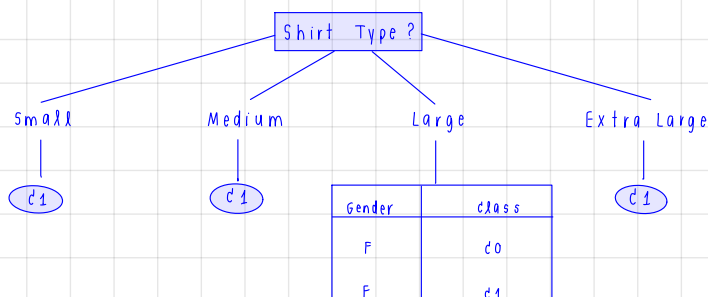
→ Extra Large

$$\approx 0.25$$

3. Information gained

$$\text{Gain}(\text{Gender}) = 0.54356 - 0.51771 = 0.02585$$

$$\text{Gain}(\text{Shirt Size}) = 0.54356 - 0.25 = 0.29356 \rightarrow \text{splitting attribute}$$



aiman

