

2.4.1.4 Information Gain

ID3 Algorithm as attribute selection ^{measure}

minimizes information entropy / uncertainty

achieve the least randomness or impurity

definition

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

$m =$ classes number

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

sample number

$p_i = \frac{|C_i \cap D|}{|D|}$ sample belongs to class C_i

Attribute A split D into v subset

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

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

Q : calculate the $\text{Gain}(A)$

PID	Fever	Cough	Tiredness	COVID-19
1	no	mild	no	-
2	yes	no	no	-
3	yes	mild	yes	+
4	mild	yes	no	+
5	mild	mild	yes	+
6	no	mild	yes	-

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

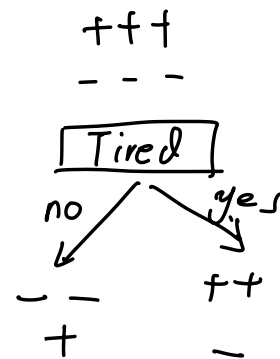
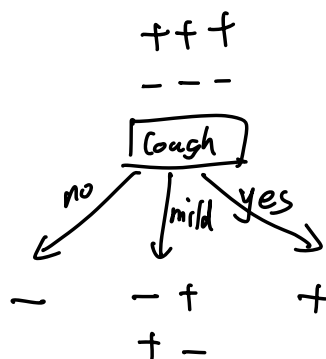
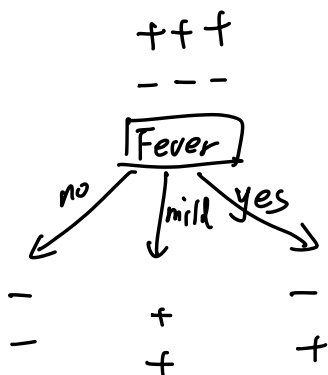
$$\text{Info}(D) = - \left[p_1 \log_2(p_1) + p_2 \log_2(p_2) \right]$$

$$p_1 = \frac{3}{6} = 0.5$$

$$p_2 = \frac{3}{6} = 0.5$$

$$\begin{aligned} \text{Info}(D) &= - \left[0.5 \log_2 0.5 + 0.5 \log_2 0.5 \right] \\ &= 1 \end{aligned}$$

②



③ Fever

$$Info_{Fever}(D) = \sum_{j=1}^3 \frac{|D_j|}{|D|} \times Info(D_j)$$

$$= \frac{2}{6} Info(D_{no}) + \frac{2}{6} Info(D_{mild}) + \frac{2}{6} Info(D_{yes})$$

where $Info(D_{no}) = -1 \log_2(1) = 0$

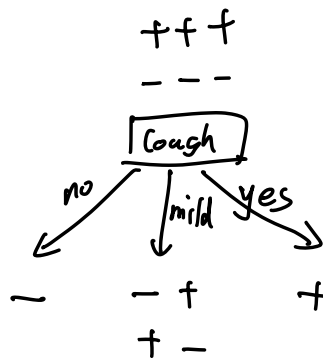
$$Info(D_{mild}) = 0$$

$$Info(D_{yes}) = -[0.5 \log_2 0.5 + 0.5 \log_2 0.5] = 1$$

So $Info_{Fever}(D) = \frac{1}{3}$

④ Cough

$$\begin{aligned}
 \text{Info}_{\text{Cough}}(D) &= \sum_{j=1}^3 \frac{|D_j|}{|D|} \times \text{Info}(D_j) \\
 &= \frac{2}{6} \text{Info}(D_{\text{no}}) + \frac{2}{6} \text{Info}(D_{\text{mild}}) + \frac{2}{6} \text{Info}(D_{\text{yes}})
 \end{aligned}$$



where $\text{Info}(D_{\text{no}}) = -1 \log_2(1) = 0$

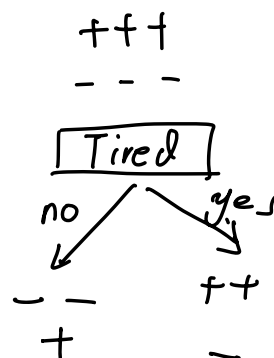
$$\text{Info}(D_{\text{mild}}) = -[0.5 \log_2 0.5 + 0.5 \log_2 0.5] = 1$$

$$\text{Info}(D_{\text{yes}}) = 0$$

So $\text{Info}_{\text{Cough}}(D) = \frac{1}{3}$

⑤ Tired

$$\begin{aligned}
 \text{Info}_{\text{Tired}}(D) &= \sum_{j=1}^2 \frac{|D_j|}{|D|} \times \text{Info}(D_j) \\
 &= \frac{3}{6} \text{Info}(D_{\text{no}}) + \frac{3}{6} \text{Info}(D_{\text{yes}})
 \end{aligned}$$



where $\text{Info}(D_{\text{no}}) = -\left[\frac{2}{3} \log_2 \frac{2}{3} + \frac{1}{3} \log_2 \frac{1}{3}\right] = 0.5850$

$$\text{Info}(D_{\text{yes}}) = 0.5850$$

So $\text{Info}_{\text{Tired}}(D) = 0.5850$

$$\textcircled{6} \text{Gain(Fever)} = 1 - \frac{1}{3} = 0.6667$$

$$\text{Gain(Cough)} = 1 - \frac{1}{3} = 0.6667$$

$$\text{Gain(Tired)} = 1 - 0.5850 = 0.4150$$

So, we can choose Fever or cough as the first attribute.

$\textcircled{7}$ Assume we choose Fever as first attribute, we can choose cough or Tiredness as the second attribute due to they have same Gain. And the classification task is over.

PID	Fever	Cough	Tiredness	COVID-19
1	no	mild	no	-
2	yes	no	no	-
3	yes	mild	yes	+
4	mild	yes	no	+
5	mild	mild	yes	+
6	no	mild	yes	-