

age	income	student	credit_rating	buys_computer
<=30	high	no	fair	no
<=30	high	no	excellent	no
31...40	high	no	fair	yes
>40	medium	no	fair	yes
>40	low	yes	fair	yes
>40	low	yes	excellent	no
31...40	low	yes	excellent	yes
<=30	medium	no	fair	no
<=30	low	yes	fair	yes
>40	medium	yes	fair	yes
<=30	medium	yes	excellent	yes
31...40	medium	no	excellent	yes
31...40	high	yes	fair	yes
>40	medium	no	excellent	no

features = x

class = y

- Expected information (entropy) needed to classify a tuple in D:

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

- Information needed (after using A to split D into v partitions) to classify D:

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

- Information gained by branching on attribute A

$$Gain(A) = Info(D) - Info_A(D)$$

● คำว่า class

$$\begin{aligned} \text{Info}(D) &= \sum_{i=1}^n P_i \log_2(P_i) \\ &= -\frac{9}{14} \log_2\left(\frac{9}{14}\right) - \frac{5}{14} \log_2\left(\frac{5}{14}\right) \\ &= 0.41 + 0.53 \end{aligned}$$

$$\text{Info}(D) = 0.940 \quad \#$$

● คำว่า Feature Info_A(D)

$$\text{สูตร Info}_A(D) = \sum_{j=1}^{\#D} \left| \frac{D_j}{D} \right| \times \text{Info}(D_j)$$

$$\square \text{Info}_{\text{age}}(D) = \frac{5}{14} I(2,3) + \frac{4}{14} I(4,0) + \frac{5}{14} I(3,2)$$

$$\begin{aligned} &= \frac{5}{14} \left[-\frac{2}{5} \log_2\left(\frac{2}{5}\right) - \frac{3}{5} \log_2\left(\frac{3}{5}\right) \right] + \frac{4}{14} \left[-\frac{4}{4} \log_2\left(\frac{4}{4}\right) - \left(\frac{0}{4}\right) \log_2\left(\frac{0}{4}\right) \right] + \\ &\quad \frac{5}{14} \left[-\frac{3}{5} \log_2\left(\frac{3}{5}\right) - \frac{2}{5} \log_2\left(\frac{2}{5}\right) \right] \end{aligned}$$

$$= 0.34676 + 0.34676$$

$$\text{Info}_{\text{age}}(D) = 0.694 \quad \#$$

$$\square \text{Info}_{\text{income}}(D) = \frac{4}{14} I(2,2) + \frac{6}{14} I(4,2) + \frac{4}{14} I(3,1)$$

$$\begin{aligned} &= \frac{4}{14} \left[\frac{2}{4} \log_2\left(\frac{2}{4}\right) - \frac{2}{4} \log_2\left(\frac{2}{4}\right) \right] + \frac{6}{14} \left[-\frac{4}{6} \log_2\left(\frac{4}{6}\right) - \frac{2}{6} \log_2\left(\frac{2}{6}\right) \right] + \\ &\quad \frac{4}{14} \left[-\frac{3}{4} \log_2\left(\frac{3}{4}\right) - \frac{1}{4} \log_2\left(\frac{1}{4}\right) \right] \end{aligned}$$

$$= 0.2857 + 0.3935 + 0.2317$$

$$\text{Info}_{\text{income}}(D) = 0.911 \quad \#$$

$$\begin{aligned}
 \square \text{Info}_{\text{credit}}(D) &= \frac{6}{14} I(3,3) + \frac{8}{14} I(6,2) \\
 &= \frac{6}{14} \left[-\frac{3}{6} \log_2 \left(\frac{3}{6} \right) - \frac{3}{6} \log_2 \left(\frac{3}{6} \right) \right] + \frac{8}{14} \left[-\frac{6}{8} \log_2 \left(\frac{6}{8} \right) - \frac{2}{8} \log_2 \left(\frac{2}{8} \right) \right] \\
 &= \frac{6}{14} [0.5 + 0.5] + \frac{8}{14} [0.3113 + 0.5] \\
 &= \left(\frac{6}{14} \times 1 \right) + \left(\frac{8}{14} \times 0.8113 \right) \\
 &= 0.4285 + 0.4636
 \end{aligned}$$

$$\text{Info}_{\text{credit}}(D) = 0.892$$

$$\begin{aligned}
 \square \text{Info}_{\text{student}}(D) &= \frac{7}{14} I(6,1) + \frac{7}{14} I(3,4) \\
 &= \frac{7}{14} \left[-\frac{6}{7} \log_2 \left(\frac{6}{7} \right) - \frac{1}{7} \log_2 \left(\frac{1}{7} \right) \right] + \frac{7}{14} \left[-\frac{3}{7} \log_2 \left(\frac{3}{7} \right) - \frac{4}{7} \log_2 \left(\frac{4}{7} \right) \right] \\
 &= \frac{7}{14} [0.191 + 0.401] + \frac{7}{14} [0.524 + 0.461] \\
 &= 0.296 + 0.493
 \end{aligned}$$

$$\text{Info}_{\text{stu}}(D) = 0.788 \quad \#$$

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

$$\text{Gain}(\text{age}) = 0.940 - 0.694 = 0.246 \quad \star \text{Root node}$$

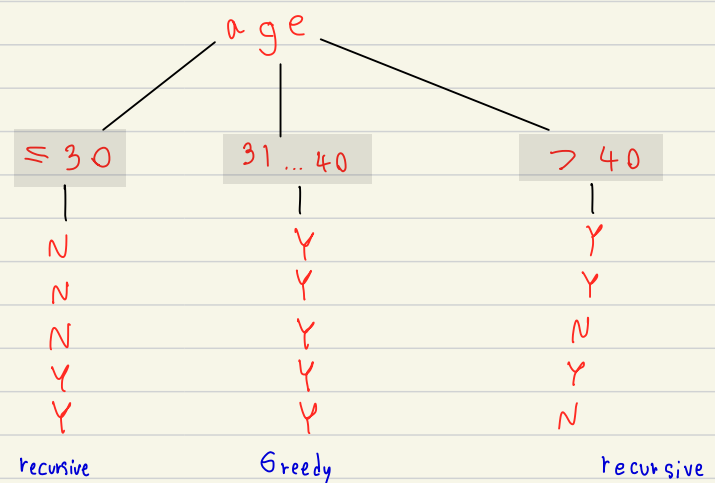
$$\text{Gain}(\text{income}) = 0.940 - 0.911 = 0.029$$

$$\text{Gain}(\text{student}) = 0.940 - 0.788 = 0.152$$

$$\text{Gain}(\text{credit-rating}) = 0.940 - 0.892 = 0.048$$

Recursive age ≤ 30

	age	income	student	credit_rating	buys_computer
0	≤ 30	high	no	fair	no
0	≤ 30	high	no	excellent	no
	31...40	high	no	fair	yes
	>40	medium	no	fair	yes
	>40	low	yes	fair	yes
	>40	low	yes	excellent	no
	31...40	low	yes	excellent	yes
0	≤ 30	medium	no	fair	no
0	≤ 30	low	yes	fair	yes
	>40	medium	yes	fair	yes
0	≤ 30	medium	yes	excellent	yes
	31...40	medium	no	excellent	yes
	31...40	high	yes	fair	yes
	>40	medium	no	excellent	no



- **Info Class**

$$\begin{aligned} \text{Info}(D) &= I(2,3) \\ &= -\frac{2}{5} \log_2 \frac{2}{5} - \frac{3}{5} \log_2 \left(\frac{3}{5} \right) \\ &= 0.5288 + 0.4422 \end{aligned}$$

$$\text{Info}(D) = 0.9710$$

- **Info Feature**

$$\begin{aligned} \text{Info}_{\text{income}}(D) &= \frac{2}{5} I(0,2) + \frac{2}{5} I(1,1) + \frac{1}{5} I(1,0) \\ &= \frac{2}{5} \left[-\frac{0}{2} \log_2 \left(\frac{0}{2} \right) - \frac{2}{2} \log_2 \left(\frac{2}{2} \right) \right] + \frac{2}{5} \left[-\frac{1}{2} \log_2 \left(\frac{1}{2} \right) - \frac{1}{2} \log_2 \left(\frac{1}{2} \right) \right] + \frac{1}{5} \left[-\frac{1}{1} \log_2 \left(\frac{1}{1} \right) - \frac{0}{1} \log_2 \left(\frac{0}{1} \right) \right] \\ &= \frac{2}{5} (0 + 0) + \frac{2}{5} (0.5 + 0.5) + \frac{1}{5} (0 + 0) \\ &= \frac{2}{5} \cdot 1 \end{aligned}$$

$$\text{Info}_{\text{income}}(D) = 0.4 \quad \#$$

$$\begin{aligned} \text{Info}_{\text{student}}(D) &= \frac{2}{5} I(2,0) + \frac{3}{5} I(0,3) \\ &= \frac{2}{5} \left[-\frac{2}{2} \log_2 \left(\frac{2}{2} \right) - \frac{0}{2} \log_2 \left(\frac{0}{2} \right) \right] + \frac{3}{5} \left[-\frac{0}{3} \log_2 \left(\frac{0}{3} \right) - \frac{3}{3} \log_2 \left(\frac{3}{3} \right) \right] \\ &= \frac{2}{5} [0 + 0] + \frac{3}{5} [0 + 0] \\ &= 0 + 0 \end{aligned}$$

$$\text{Info}_{\text{student}}(D) = 0 \quad \#$$

$$\text{Info}_{\text{credit}}(D) = 0.951 \text{ \#}$$

age	income	student	credit_rating	buys_computer
<=30	high	no	fair	no
<=30	high	no	excellent	no
31...40	high	no	fair	yes
>40	medium	no	fair	yes
>40	low	yes	fair	yes
>40	low	yes	excellent	no
31...40	low	yes	excellent	yes
<=30	medium	no	fair	no
<=30	low	yes	fair	yes
>40	medium	yes	fair	yes
<=30	medium	yes	excellent	yes
31...40	medium	no	excellent	yes
31...40	high	yes	fair	yes
>40	medium	no	excellent	no



Info (D) = 0.971 #

$$\begin{aligned}
 \text{Info}_{\text{income}}(D) &= \frac{3}{5} I(2,1) + \frac{2}{5} I(1,1) \\
 &= \frac{3}{5} \left[-\frac{2}{5} \log_2 \left(\frac{2}{5} \right) - \frac{1}{5} \log_2 \left(\frac{1}{5} \right) \right] + \frac{2}{5} \left[-\frac{1}{2} \log_2 \left(\frac{1}{2} \right) - \frac{1}{2} \log_2 \left(\frac{1}{2} \right) \right] \\
 &= \frac{3}{5} [0.399 + 0.528] + \frac{2}{5} [0.5 + 0.5] \\
 &= 0.5509 + 0.4
 \end{aligned}$$

$$\text{Info}_{\text{income}}(D) = 0.9509 \quad \#$$

$$\begin{aligned}
 \text{Info}_{\text{credit}}(D) &= \frac{2}{5} I(0,2) + \frac{3}{5} I(3,0) \\
 &= \frac{2}{5} \left[-\frac{0}{2} \log_2 \left(\frac{0}{2} \right) - \frac{2}{2} \log_2 \left(\frac{2}{2} \right) \right] + \frac{3}{5} \left[-\frac{3}{3} \log_2 \left(\frac{3}{3} \right) - \frac{0}{3} \log_2 \left(\frac{0}{3} \right) \right] \\
 &= \frac{2}{5} [0] + \frac{3}{5} [0] \\
 &= 0
 \end{aligned}$$

$$\text{Info}_{\text{credit}}(D) = 0 \quad \#$$

$$\text{Gain}(\text{income}) = 0.710 - 0.9509 = 0.0201$$

$$\text{Gain}(\text{credit_rating}) = 0.710 - 0 = 0.710 \quad \# \text{ Gain info}$$

Resulting tree:

