

Training data set: Who buys computer?

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

— N H N F N
— N H N E N

— N M N F N
— Y L Y F Y
— Y M Y E Y

u1 $\text{Info}(D) = I(9,5) = -\frac{9}{14} \log_2\left(\frac{9}{14}\right) - \left(\frac{5}{14}\right) \log_2\left(\frac{5}{14}\right) = 0.940$

u1 $\text{Info}_{\text{Age}}(D) = \frac{5}{14} I(2,3) + \frac{4}{14} I(4,0) + \frac{5}{14} I(3,2) = 0.694$

$\text{Info}_{\text{income}}(D) = \frac{4}{14} I(3,1) + \frac{6}{14} I(4,2) + \frac{4}{14} I(2,2)$
 $= \frac{4}{14} \left(-\frac{3}{4} \log_2\left(\frac{3}{4}\right) - \frac{1}{4} \log_2\left(\frac{1}{4}\right) \right) + \frac{6}{14} \left(-\frac{4}{6} \log_2\left(\frac{4}{6}\right) - \frac{2}{6} \log_2\left(\frac{2}{6}\right) \right) + \frac{4}{14} \left(-\frac{2}{4} \log_2\left(\frac{2}{4}\right) - \frac{2}{4} \log_2\left(\frac{2}{4}\right) \right)$
 $= 0.911$

$\text{Info}_{\text{student}}(D) = \frac{7}{14} I(4,3) + \frac{7}{14} I(6,1)$
 $= \frac{7}{14} \left(-\frac{4}{7} \log_2\left(\frac{4}{7}\right) - \frac{3}{7} \log_2\left(\frac{3}{7}\right) \right) + \frac{7}{14} \left(-\frac{6}{7} \log_2\left(\frac{6}{7}\right) - \frac{1}{7} \log_2\left(\frac{1}{7}\right) \right)$
 $= 0.788$

$\text{Info}_{\text{credit-rating}}(D) = \frac{8}{14} I(6,2) + \frac{6}{14} I(3,3)$
 $= \frac{8}{14} \left(-\frac{6}{8} \log_2\left(\frac{6}{8}\right) - \frac{2}{8} \log_2\left(\frac{2}{8}\right) \right) + \frac{6}{14} \left(-\frac{3}{6} \log_2\left(\frac{3}{6}\right) - \frac{3}{6} \log_2\left(\frac{3}{6}\right) \right)$
 $= 0.892$

u1 Gain

$\text{Gain}(\text{age}) = \text{info}(D) - \text{info}_{\text{age}}(D) = 0.940 - 0.694 = 0.246$

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

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

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

ดังนั้น เลือก feature Gain age มีค่า 0.246

age : ≤ 30

$$\text{Info}_{\text{age: } \leq 30} (D) = I\left(\begin{smallmatrix} Y & N \\ 2 & 3 \end{smallmatrix}\right) = -\frac{2}{5} \log_2\left(\frac{2}{5}\right) - \frac{3}{5} \log_2\left(\frac{3}{5}\right) = 0.971$$

$$\begin{aligned} \text{Info}_{\text{income}} (D) &= \frac{1}{5} I\left(\begin{smallmatrix} \text{low} & Y & N \\ 1 & 0 \end{smallmatrix}\right) + \frac{2}{5} I\left(\begin{smallmatrix} M & Y & N \\ 1 & 1 \end{smallmatrix}\right) + \frac{2}{5} I\left(\begin{smallmatrix} H & Y & N \\ 0 & 2 \end{smallmatrix}\right) \\ &= \frac{1}{5} \left(-\frac{1}{1} \log_2\left(\frac{1}{1}\right) - 0\right) + \frac{2}{5} \left(-\frac{1}{2} \log_2\left(\frac{1}{2}\right) - \frac{1}{2} \log_2\left(\frac{1}{2}\right)\right) + \frac{2}{5} \left(-\frac{2}{2} \log_2\left(\frac{2}{2}\right) - 0\right) \\ &= 0.4 \end{aligned}$$

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

F N E N
F N E Y
F Y

$$\begin{aligned} \text{Info}_{\text{credit_rating}} (D) &= \frac{2}{5} I\left(\begin{smallmatrix} F \\ 1, 2 \end{smallmatrix}\right) + \frac{2}{5} I\left(\begin{smallmatrix} E \\ 1, 1 \end{smallmatrix}\right) \\ &= \frac{2}{5} \left(-\frac{1}{3} \log_2\left(\frac{1}{3}\right) - \frac{2}{3} \log_2\left(\frac{2}{3}\right)\right) + \frac{2}{5} \left(-\frac{1}{2} \log_2\left(\frac{1}{2}\right) - \frac{1}{2} \log_2\left(\frac{1}{2}\right)\right) \\ &= 0.951 \end{aligned}$$

$$\text{Gain}(\text{income}) = \text{info}_{\text{age: } \leq 30} (D) - \text{info}_{\text{income}} (D) = 0.971 - 0.400 = 0.571$$

$$\text{Gain}(\text{student}) = \text{info}_{\text{age: } \leq 30} (D) - \text{info}_{\text{student}} (D) = 0.971 - 0 = 0.971$$

$$\text{Gain}(\text{credit_rating}) = \text{info}_{\text{age: } \leq 30} (D) - \text{info}_{\text{credit_rating}} (D) = 0.971 - 0.951 = 0.020$$

ดังนั้นเลือก Node student เพราะได้ Gain สูงสุด

age: 31 ... 40

$$\begin{aligned} \text{Info}_{\text{age: } 31 \dots 40} (D) &= I(4, 0) \\ &= -\frac{4}{4} \log_2\left(\frac{4}{4}\right) - \frac{0}{4} \log_2\left(\frac{0}{4}\right) \\ &= 0 \end{aligned}$$

ไม่มีตัวอย่าง หรือตัวอย่างน้อยเกินไป Yes 4 อัน

age : 40

$$\begin{aligned} \text{Info}_{\text{age: } \geq 40} (D) &= I(3, 2) \\ &= -\frac{3}{5} \log_2\left(\frac{3}{5}\right) - \frac{2}{5} \log_2\left(\frac{2}{5}\right) \\ &= 0.971 \end{aligned}$$

$$\begin{aligned} \text{Info}_{\text{income}} (D) &= \frac{2}{5} I(1, 1) + \frac{3}{5} I(2, 1) \\ &= \frac{2}{5} \left(-\frac{1}{2} \log_2\left(\frac{1}{2}\right) - \frac{1}{2} \log_2\left(\frac{1}{2}\right)\right) + \frac{3}{5} \left(-\frac{2}{3} \log_2\left(\frac{2}{3}\right) - \frac{1}{3} \log_2\left(\frac{1}{3}\right)\right) \\ &= 0.951 \end{aligned}$$

$$\begin{aligned} \text{Info}_{\text{student}} (D) &= \frac{3}{5} I(2, 1) + \frac{2}{5} I(1, 1) \\ &= \frac{3}{5} \left(-\frac{2}{3} \log_2\left(\frac{2}{3}\right) - \frac{1}{3} \log_2\left(\frac{1}{3}\right)\right) + \frac{2}{5} \left(-\frac{1}{2} \log_2\left(\frac{1}{2}\right) - \frac{1}{2} \log_2\left(\frac{1}{2}\right)\right) \end{aligned}$$

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

$$\begin{aligned} \text{Gain}(\text{income}) &= \text{info}_{\text{age} > 40}(\text{CD}) - \text{info}_{\text{income}}(\text{CD}) = 0.971 - 0.951 = 0.02 \\ \text{Gain}(\text{student}) &= \text{info}_{\text{age} > 40}(\text{CD}) - \text{info}_{\text{student}}(\text{CD}) = 0.971 - 0.951 = 0.02 \\ \text{Gain}(\text{credit_rating}) &= \text{info}_{\text{age} > 40}(\text{CD}) - \text{info}_{\text{credit_rating}}(\text{CD}) = 0.971 - 0 = 0.971 \end{aligned}$$

∴ Credit_rating ကို Gain အရ ရွေး

