

# Decision Tree Induction

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

- wt Info (D)

$$\text{Info (D)} = I(9,5) = \frac{9}{14} \log_{(2)} \left( \frac{9}{14} \right) - \frac{5}{14} \log_{(2)} \left( \frac{5}{14} \right) = 0.94$$

- wt Info<sub>age</sub> (D)

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

$$I(2,3) = -\frac{2}{5} \log_{(2)} \left( \frac{2}{5} \right) - \frac{3}{5} \log_{(2)} \left( \frac{3}{5} \right) = 0.971$$

$$I(4,0) = -\frac{4}{4} \log_{(2)} \left( \frac{4}{4} \right) - \frac{0}{4} \log_{(2)} \left( \frac{0}{4} \right) = 0$$

$$I(3,2) = -\frac{3}{5} \log_{(2)} \left( \frac{3}{5} \right) - \frac{2}{5} \log_{(2)} \left( \frac{2}{5} \right) = 0.971$$

$$\text{min Info}_{\text{age}}(D) = \frac{5}{14} (0.971) + \frac{4}{14} (0) + \frac{5}{14} (0.971) = 0.694$$

- wt Gain (age)

$$\text{Gain (age)} = 0.94 - 0.694 = 0.246$$

- u1 Info<sub>income</sub> (D)

$$\text{Info}_{\text{income}}(D) = \frac{4}{14} I_{\text{high}}^{\text{Yes}}(2,2) + \frac{6}{14} I_{\text{medium}}^{\text{Yes}}(4,2) + \frac{4}{14} I_{\text{low}}^{\text{Yes}}(3,1)$$

$$I(2,2) = -\frac{2}{4} \log_{(2)}\left(\frac{2}{4}\right) - \frac{2}{4} \log_{(2)}\left(\frac{2}{4}\right) = 1$$

$$I(4,2) = -\frac{4}{6} \log_{(2)}\left(\frac{4}{6}\right) - \frac{2}{6} \log_{(2)}\left(\frac{2}{6}\right) = 0.918$$

$$I(3,1) = -\frac{3}{4} \log_{(2)}\left(\frac{3}{4}\right) - \frac{1}{4} \log_{(2)}\left(\frac{1}{4}\right) = 0.811$$

$$\text{minim Info}_{\text{income}}(D) = \frac{4}{14}(1) + \frac{6}{14}(0.918) + \frac{4}{14}(0.811) = 0.911$$

- u1 Gain(income)

$$\begin{aligned} \text{Gain(income)} &= 0.94 - 0.911 \\ &= 0.029 \end{aligned}$$

- u1 Info<sub>student</sub> (D)

$$\text{Info}_{\text{student}}(D) = \frac{7}{14} I_{\text{Yes}}^{\text{Yes}}(6,1) + \frac{7}{14} I_{\text{No}}^{\text{Yes}}(3,4)$$

$$I(6,1) = -\frac{6}{7} \log_{(2)}\left(\frac{6}{7}\right) - \frac{1}{7} \log_{(2)}\left(\frac{1}{7}\right) = 0.592$$

$$I(3,4) = -\frac{3}{7} \log_{(2)}\left(\frac{3}{7}\right) - \frac{4}{7} \log_{(2)}\left(\frac{4}{7}\right) = 0.985$$

$$\text{minim Info}_{\text{student}}(D) = \frac{7}{14}(0.592) + \frac{7}{14}(0.985)$$

$$= 0.789$$

- u1 Gain(student)

$$\begin{aligned} \text{Gain(student)} &= 0.94 - 0.789 \\ &= 0.151 \end{aligned}$$

- u1 Info<sub>credit\_rating</sub>(D)

$$\text{Info}_{\text{credit\_rating}}(D) = \frac{8}{14} I(\overset{\text{fair}}{6}, \overset{\text{N}}{2}) + \frac{6}{14} I(\overset{\text{N}}{3}, \overset{\text{Y}}{3})$$

$$I(6,2) = -\frac{6}{8} \log_{(2)}\left(\frac{6}{8}\right) - \frac{2}{8} \log_{(2)}\left(\frac{2}{8}\right) = 0.911$$

$$I(3,3) = -\frac{3}{6} \log_{(2)}\left(\frac{3}{6}\right) - \frac{3}{6} \log_{(2)}\left(\frac{3}{6}\right) = 1$$

$$\text{nnwn Info}_{\text{credit\_rating}}(D) = \frac{8}{14} (0.911) + \frac{6}{14} (1) = 0.992$$

- u1 Gain(credit\_rating)

$$\begin{aligned} \text{Gain}(\text{credit\_rating}) &= 0.99 - 0.992 \\ &= 0.048 \end{aligned}$$

Gain

$$\text{Gain}(\text{age}) = 0.246$$

$$\text{Gain}(\text{income}) = 0.029$$

$$\text{Gain}(\text{student}) = 0.151$$

$$\text{Gain}(\text{credit\_rating}) = 0.048$$

เลือก Gain ที่มากที่สุดคือ age (0.246) > Gain (age)

age (<= 30)

- u1 Info(D) ของ age (<= 30)

$$\text{Info}(D) = I(2,3) = 0.911$$

- u1 Info<sub>income</sub>(D) ของ age (<= 30)

$$\text{Info}_{\text{income}}(D) \text{ ของ age } (<= 30) = \frac{2}{3} I(\overset{\text{high}}{0}, \overset{\text{N}}{2}) + \frac{2}{5} I(\overset{\text{medium}}{1}, \overset{\text{N}}{1}) + \frac{1}{5} I(\overset{\text{low}}{1}, \overset{\text{N}}{0})$$

$$I(0,2) = -\frac{0}{2} \log_{(2)}\left(\frac{0}{2}\right) - \frac{2}{2} \log_{(2)}\left(\frac{2}{2}\right) = 0$$

$$I(1,1) = -\frac{1}{2} \log_{(2)}\left(\frac{1}{2}\right) - \frac{1}{2} \log_{(2)}\left(\frac{1}{2}\right) = 1$$

$$I(1,0) = -\frac{1}{1} \log_{(2)}\left(\frac{1}{1}\right) - \frac{0}{1} \log_{(2)}\left(\frac{0}{1}\right) = 0$$

$$\text{Inform Info}_{\text{income}}(D) \text{ for } \text{age} (<=30) = \frac{2}{5}(0) + \frac{2}{5}(1) + \frac{1}{5}(0) = 0.4$$

- or  $\text{Gain}(\text{income}) \text{ for } \text{age} (<=30)$

$$\text{Gain}(\text{income}) \text{ for } \text{age} (<=30) = 0.971 - 0.4 = 0.571$$

- or  $\text{Info}_{\text{student}}(D) \text{ for } \text{age} (<=30)$

$$\text{Info}_{\text{student}}(D) \text{ for } \text{age} (<=30) = \frac{2}{5} I(2,0)^{\text{yes}} + \frac{3}{5} I(0,3)^{\text{no}}$$

สังเกต Yes  $\rightarrow$  Yes (buy Computer)  
No  $\rightarrow$  No (buy Computer)

เลือกด้วย Student เพราะสามารถแบ่งข้อมูลได้แบบสมบูรณ์

$\text{age} (>40)$

$$\text{Info}(D) \text{ for } \text{age} (>40) = I(3,2) = 0.971$$

- or  $\text{Info}_{\text{income}}(D) \text{ for } \text{age} (>40) = \frac{3}{5} I(2,1)^{\text{medium}} + \frac{2}{5} I(1,1)^{\text{low}}$

$$I(2,1) = -\frac{2}{3} \log_{(2)} \left( \frac{2}{3} \right) - \frac{1}{3} \log_{(2)} \left( \frac{1}{3} \right) = 0.918$$

$$I(1,1) = 1$$

$$\text{Inform Info}_{\text{income}}(D) \text{ for } \text{age} (>40) = \frac{3}{5}(0.918) + \frac{2}{5}(1) = 0.951$$

- or  $\text{Gain}(\text{income}) \text{ for } \text{age} (>40)$

$$\text{Gain}(\text{income}) \text{ for } \text{age} (>40) = 0.971 - 0.951 = 0.02$$

- or  $\text{Info}_{\text{student}}(D) \text{ for } \text{age} (>40)$

$$\text{Info}_{\text{student}}(D) \text{ for } \text{age} (>40) = \frac{3}{5} I(2,1)^{\text{yes}} + \frac{2}{5} I(1,1)^{\text{no}}$$

$$I(2,1) = -\frac{2}{3} \log_{(2)} \left( \frac{2}{3} \right) - \frac{1}{3} \log_{(2)} \left( \frac{1}{3} \right) = 0.918$$

$$I(1,1) = 1$$

$$\text{Inform Info}_{\text{student}}(D) \text{ for } \text{age} (>40) = \frac{3}{5}(0.918) + \frac{2}{5}(1) = 0.951$$

- or  $\text{Gain}(\text{student}) \text{ for } \text{age} (>40)$

$$\text{Gain}(\text{student}) \text{ for } \text{age} (>40) = 0.971 - 0.951 = 0.02$$

- Info  $\text{credit\_rating} (0)$  vs  $\text{age} (>40)$

$$\text{Info}_{\text{credit\_rating}} (0) \text{ vs } \text{age} (>40) = \frac{3}{5} I(3,0) + \frac{2}{5} I(0,2)$$

ถ้า  $\text{fair} \rightarrow \text{Yes (buy\_computer)}$   
 $\text{excellent} \rightarrow \text{No (buy\_computer)}$

เลือกแบ่งด้วย  $\text{credit\_rating}$  เพราะสามารถแบ่งข้อมูลได้สมบูรณ์

สรุป

