

## Decision Tree Induction

### Information Gain: An Attribute Selection Measure

- ❑ Select the attribute with the highest information gain (used in typical decision tree induction algorithm: ID3/C4.5)
- ❑ Let  $p_i$  be the probability that an arbitrary tuple in  $D$  belongs to class  $C_i$ , estimated by  $|C_{i,D}|/|D|$
- ❑ 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  $A$

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

11

Example:

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

๑๗ Info (D)

$$\text{Info}(D) = I(\overset{y}{9}, \overset{n}{5}) = -\frac{9}{14} \log_2\left(\frac{9}{14}\right) - \frac{5}{14} \log_2\left(\frac{5}{14}\right) = 0.94$$

๑๗ Info<sub>age</sub> (D)

$$\text{Info}_{\text{age}}(D) = \overset{\leq 30}{\frac{5}{14} I(\overset{y}{2}, \overset{n}{3})} + \overset{31-40}{\frac{4}{14} I(\overset{y}{4}, \overset{n}{0})} + \overset{> 40}{\frac{5}{14} I(\overset{y}{3}, \overset{n}{2})}$$

$$I(2,3) = -\frac{2}{5} \log_2\left(\frac{2}{5}\right) - \frac{3}{5} \log_2\left(\frac{3}{5}\right) = 0.971$$

$$I(4,0) = -\frac{4}{4} \log_2\left(\frac{4}{4}\right) - \frac{0}{4} \log_2\left(\frac{0}{4}\right) = 0$$

$$I(3,2) = -\frac{3}{5} \log_2\left(\frac{3}{5}\right) - \frac{2}{5} \log_2\left(\frac{2}{5}\right) = 0.971$$

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

๑๗ Gain (age)

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

## u1 Info income (D)

$$\text{Info income (D)} = \frac{4}{14} I(\overset{\text{high}}{2}, \overset{\text{h}}{2}) + \frac{6}{14} I(\overset{\text{medium}}{4}, \overset{\text{h}}{2}) + \frac{4}{14} I(\overset{\text{low}}{3}, \overset{\text{h}}{1})$$

$$I(\overset{\text{h}}{2}, \overset{\text{h}}{2}) = -\frac{2}{4} \log_2\left(\frac{2}{4}\right) - \frac{2}{4} \log_2\left(\frac{2}{4}\right) = 1$$

$$I(\overset{\text{h}}{4}, \overset{\text{h}}{2}) = -\frac{4}{6} \log_2\left(\frac{4}{6}\right) - \frac{2}{6} \log_2\left(\frac{2}{6}\right) = 0.918$$

$$I(\overset{\text{h}}{3}, \overset{\text{h}}{1}) = -\frac{3}{4} \log_2\left(\frac{3}{4}\right) - \frac{1}{4} \log_2\left(\frac{1}{4}\right) = 0.811$$

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

## u1 Gain (income)

$$\text{Gain income} = 0.94 - 0.911 = 0.029$$

## u1 Info student (D)

$$\text{Info student (D)} = \frac{7}{14} I(\overset{\text{yes}}{6}, \overset{\text{h}}{1}) + \frac{7}{14} I(\overset{\text{no}}{3}, \overset{\text{h}}{4})$$

$$I(\overset{\text{h}}{6}, \overset{\text{h}}{1}) = -\frac{6}{7} \log_2\left(\frac{6}{7}\right) - \frac{1}{7} \log_2\left(\frac{1}{7}\right) = 0.592$$

$$I(\overset{\text{h}}{3}, \overset{\text{h}}{4}) = -\frac{3}{7} \log_2\left(\frac{3}{7}\right) - \frac{4}{7} \log_2\left(\frac{4}{7}\right) = 0.985$$

$$\text{Info student (D)} = \frac{7}{14} (0.592) + \frac{7}{14} (0.985) = 0.789$$

## u1 Gain (student)

$$\text{Gain student} = 0.94 - 0.789 = 0.151$$

## ๖๗ Info credit-rating (D)

$$\text{Info credit-rating (D)} = \frac{9}{14} I(\overset{\text{fair}}{6, 2}) + \frac{1}{14} I(\overset{\text{excellent}}{3, 3})$$

$$I(\overset{7}{6}, \overset{7}{2}) = -\frac{6}{8} \log_2 \left(\frac{6}{8}\right) - \frac{2}{8} \log_2 \left(\frac{2}{8}\right) = 0.811$$

$$I(\overset{7}{3}, \overset{7}{3}) = -\frac{3}{6} \log_2 \left(\frac{3}{6}\right) - \frac{3}{6} \log_2 \left(\frac{3}{6}\right) = 1$$

$$\text{Info credit-rating (D)} = \frac{9}{14} (0.811) + \frac{1}{14} (1) = 0.872$$

## ๖๘ Gain (Credit-rating)

$$\text{Gain (credit-rating)} = 0.94 - 0.872 = 0.068$$

๗๗ Gain

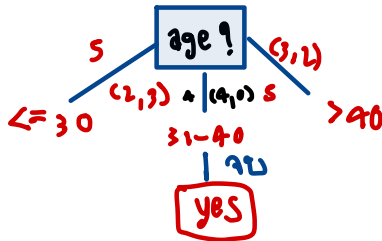
$$\text{Gain (age)} = 0.746$$

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

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

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

∴ เลือก Gain ที่มากที่สุด, การเลือกมาเป็นค่าแรก คือ Gain (age)



ex info (0) and age ( $\leq 30$ )

$\text{Info}(D) = I(2,3) = 0.977$  ↑ high ↑ medium ↑ low

$$x_1 \ln f_{income}(0) \text{ vs } age = \frac{2}{5} I^{(y, n)}(0, 2) + \frac{2}{5} I^{(y, n)}(1, 1) + \frac{1}{5} I^{(y, n)}(1, 0)$$

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

$$I(1,1) = -\frac{1}{2} \log_2\left(\frac{1}{2}\right) - \frac{1}{2} \log_2\left(\frac{1}{2}\right) = 1$$

$$I(1,0) = -\frac{1}{7} \log(2) \left(\frac{1}{7}\right) - \frac{0}{7} \log(2) \left(\frac{0}{7}\right) = 0$$

Now info income (D) vs age ( $\leq 50$ ) =  $\frac{2}{5}(0) + \frac{1}{5}(1) + \frac{1}{5}(0) = 0.4$

vi Gain (income) w/ age ( $\leq 30$ ) =  $0.971 - 0.4 = 0.571$  #

ex info student (D) vs age ( $\leq 30$ ) =  $\frac{2}{5} I(\frac{y}{2}, 0) + \frac{3}{5} I(\frac{y}{0}, 3)$

సమస్య, yes  $\rightarrow$  yes (buy-computer), No  $\rightarrow$  No (buy-computer)

เลือกแบ่งด้วย student เนื่องจากสามารถแบ่งข้อมูลได้ สมบูรณ์

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

$$\text{Info}_{\text{income}}(D) \text{ vs } \text{age} (>40) = \frac{3}{5} I(2,1) + \frac{2}{5} I(1,1)$$

$$I(2,1) = -\frac{2}{3} \log_2(\frac{2}{3}) - \frac{1}{3} \log_2(\frac{1}{3}) = 0.918$$

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

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

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

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

fair → yes (buy-computer), excellent → no (buy-computer)

credit-rating

“yes”

