

## 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  $A$

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

age	$p_i$	$n_i$	$I(p_i, n_i)$
$\leq 30$	2	3	0.971
31...40	4	0	0
$> 40$	3	2	0.971

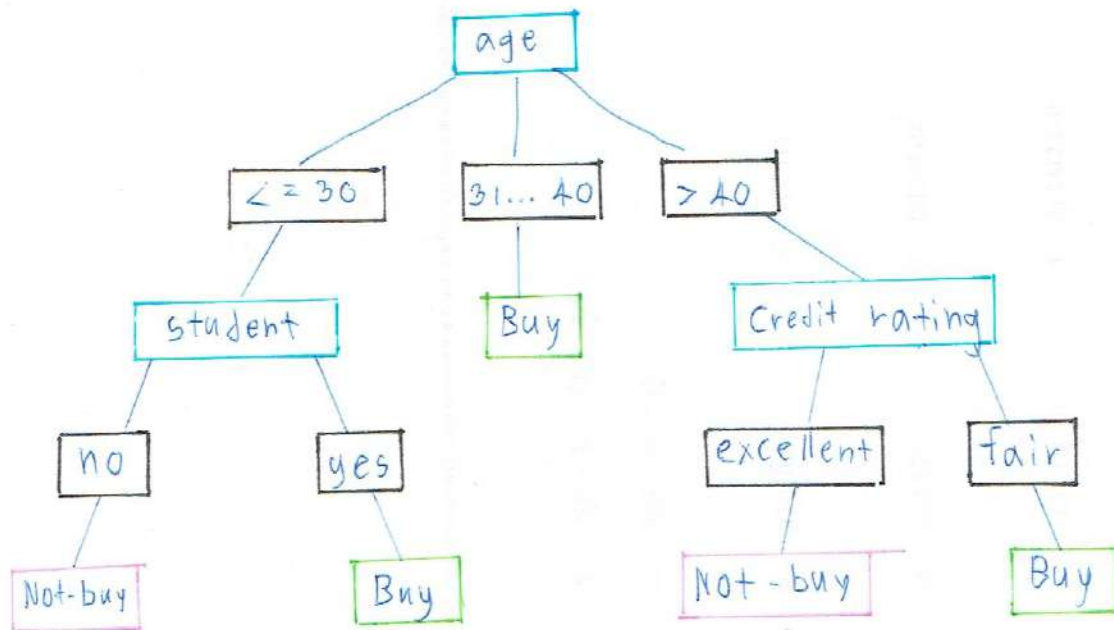
age	income	student	credit rating	buys computer
$\leq 30$	high	no	fair	no
$\leq 30$	high	no	excellent	no
31...40	high	no	fair	yes
$> 40$	medium	no	fair	yes
$> 40$	low	yes	fair	yes
$> 40$	low	yes	excellent	no
31...40	low	yes	excellent	yes
$\leq 30$	medium	no	fair	no
$\leq 30$	low	yes	fair	yes
$> 40$	medium	yes	fair	yes
$\leq 30$	medium	yes	excellent	yes
31...40	medium	no	excellent	yes
31...40	high	yes	fair	yes
$> 40$	medium	no	excellent	no

Resulting tree :

အမျိုးအမည်အရ

အသစ်ဝယ်

6230210A1-7



91 Info (D)

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

91 Info<sub>age</sub>(D)

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

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

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

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

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

91 Gain (age)

$$\text{Gain}(\text{age}) = \text{Info}(D) - \text{Info}_{\text{age}}(D)$$

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

91 Info<sub>income</sub>(D)

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

$$I(2,2) = 1$$

$$I(4,2) = 0.918$$

$$I(3,1) = 0.911$$

$$\begin{aligned} \text{Info}_{\text{income}}(D) &= \frac{4}{14} (1) + \frac{6}{14} (0.918) + \frac{4}{14} (0.911) \\ &= 0.9108 \end{aligned}$$

91 Gain (income)

$$\begin{aligned} \text{Gain}_{\text{income}} &= 0.94 - 0.9108 \\ &= 0.02918 \end{aligned}$$

997 Info student (D)

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

$$I(6,1) = 0.392$$

$$I(3,4) = 0.943$$

$$\text{bkkh} \text{ Info}_{\text{student}}(D) = \frac{7}{14} (0.392) + \frac{7}{14} (0.943) = 0.789$$

998 Gain (student)

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

999 Info credit-rating (D)

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

$$\text{bkkh} \text{ Info}_{\text{credit-rating}}(D) = \frac{8}{14} (0.611) + \frac{6}{14} (1) = 0.892$$

999 Gain credit-rating

$$\text{Gain}(\text{credit-rating}) = 0.99 - 0.892 = 0.098$$

✓ บัดนี้ Gain ดังนี้

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

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

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

$$\text{Gain}(\text{credit-rating}) = 0.049$$

✓ เปรียบเทียบ Gain ซึ่งค่ามากที่สุดจะเป็นค่าที่เหมาะสมที่สุด ได้ Gain (age)

age ( $\leq 30$ )

Info(D) vs age ( $\leq 30$ )

$$\text{Info}(D) = I(2,3) = 0.991$$

Info income (D) vs age

$$\text{Info income vs age } (\leq 30) = \frac{2}{5} \overset{\text{High}}{I(0,2)} + \frac{2}{5} \overset{\text{Median}}{I(1,1)} + \frac{1}{5} \overset{\text{Low}}{I(1,0)}$$

$$I(0,2) = 0$$

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

$$I(1,0) = 0$$

uncertain info income (D) vs age ( $\leq 30$ )

$$\begin{aligned} &= \frac{2}{5}(0) + \frac{2}{5}(1) + \frac{1}{5}(0) \\ &= 0.4 \end{aligned}$$

Gain (income vs age ( $\leq 30$ ))

$$= 0.991 - 0.4 = 0.591$$

Info student (D) vs age ( $\leq 30$ )

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

$$I(2,0) = 0$$

$$I(0,3) = 0$$

ดังนั้น Yes  $\rightarrow$  yes (buy-computer) , No  $\rightarrow$  No (buy-computer)

เนื่องจาก student เป็นนักเรียนไม่จำเป็นต้องซื้อคอมพิวเตอร์

age ( $>40$ )

241 Info(D) vs age ( $>40$ )

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

242 Info income(D) vs age ( $>40$ )

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

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

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

243 Info income(D) vs age ( $>40$ )

$$= \frac{3}{5} (0.918) + \frac{2}{5} (1) = 0.951$$

244 Gain income vs age ( $>40$ )

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

245 Info student(D) vs age ( $>40$ )

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

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

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

246 Info student(D) vs age ( $>40$ )

$$= \frac{3}{5} (0.918) + \frac{2}{5} (1) \\ = 0.951$$

247 Gain(student) vs age ( $>40$ )

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



Info credit-rating (D) and age ( $>40$ )

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

$I(3,0) = 0$        $I(0,2) = 0$

ถ้า fair  $\rightarrow$  Yes (buy computer), excellent  $\rightarrow$  No (buy computer)

ถ้า credit-rating บวกลบ แล้วดูว่าใช่หรือไม่

ใช่

