

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 :

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

Handwritten notes: A red arrow points from the word "class" to the index i in the summation. Another red arrow points from the text "# v in A" to the summation symbol.

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

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

Handwritten note: A red arrow points from the text "# v in A" to the index j in the summation.

- ❑ Information gained by branching on attribute A

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

Example: Attribute Selection with Information Gain

□ Class P: buys_computer = “yes”

□ Class N: buys_computer = “no”

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

age	p _i	n _i	I(p _i , n _i)
<=30	2	3	0.971
31...40	4	0	0
>40	3	2	0.971

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

$$Info_{age}(D) = \frac{5}{14} I(2,3) + \frac{4}{14} I(4,0) + \frac{5}{14} I(3,2) = 0.694$$

$\frac{5}{14} I(2,3)$ means “age <=30” has 5 out of 14 samples, with 2 yes’es and 3 no’s.

Hence **4/14**

$$Gain(age) = Info(D) - Info_{age}(D) = 0.246$$

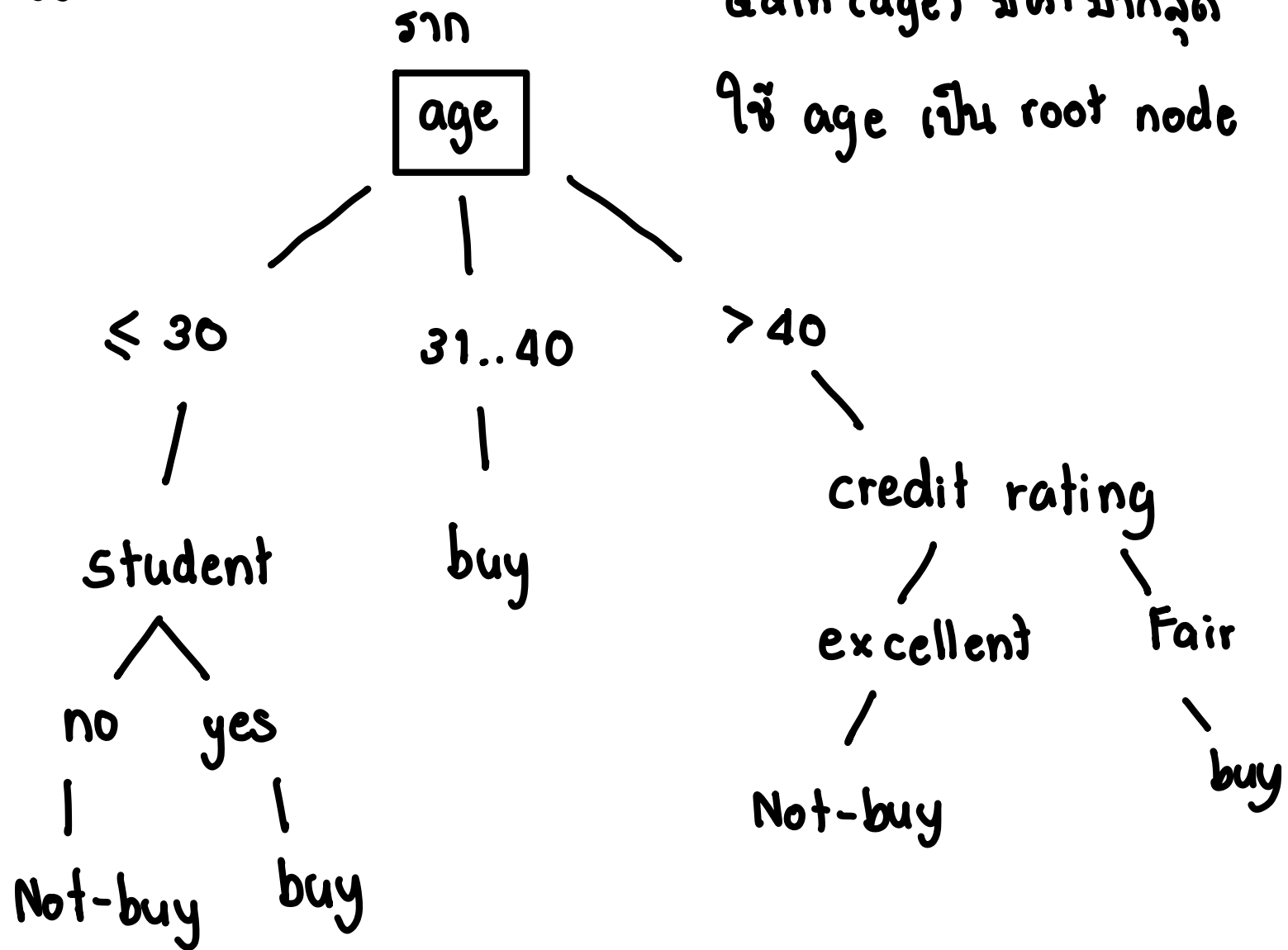
Similarly, we can get

$$Gain(income) = 0.029$$

$$Gain(student) = 0.151$$

$$Gain(credit_rating) = 0.048$$

Resulting tree :



Gain (age) มีค่ามากที่สุด
ใช้ age เป็น root node

$$\text{Info}(D) = - \sum_{i=1}^m p_i \log_2(p_i) = - \left(\frac{3}{5} \log_2 \frac{3}{5} + \frac{2}{5} \log_2 \frac{2}{5} \right) \\ \approx 0.971$$

Gain(age)

$$\text{Info}(D) - \text{Info}_A(D)$$

$$= 0.940 - 0.789$$

$$= 0.151$$

$$I(2,2) =$$

\approx

$$I(3,0) =$$

\approx

$$\leq 30$$

$$I(3,2) =$$

\approx

$$31 - 40$$

$$> 40$$

$$I(2,2)$$

4/14 NO 2 YES 2

5/14 NO

Gain (income)

$I() =$

$I() =$

$I() =$

$\text{Info}_{\text{income}} =$

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