

## Decision Tree

การสร้างต้นปัญญาตัดสินใจ (Decision Tree) การที่เลือก node ที่ดีที่สุดจาก Feature นำมาแยก class

จะหาจาก "Information Gain" ที่สูงที่สุดจากคลัตเตอร์ของ "Entropy" ซึ่งมีสูตรค่า

$$Info(D) = - \sum_{i=1}^m p_i \log_2(p_i) \quad (\text{Gain w.r.t Class}) ; p_i \text{ คือ ความน่าจะเป็น}$$

$$Info_A(D) = \sum_{j=1}^{|D|} |D_j| \times Info(D_j) \quad (\text{Gain w.r.t Feature})$$

$$Gain(A) = Info(D) - Info_A(D) \quad (\text{Information Gain ของ Feature ที่เลือก})$$

### Example

Data

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

↓  
Feature

↓  
Class

- คำนวณหา Gain w.r.t Class จากค่าจำนวน buys\_computer ฝั่งก่อนไปถัดไป Class "Yes" กับ "No"

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

$$Info(D) = I(9, 5) = \left[ \underbrace{-\frac{9}{14} \log_2(\frac{9}{14})}_{\text{Yes}} \right] + \left[ \underbrace{-\frac{5}{14} \log_2(\frac{5}{14})}_{\text{No}} \right]$$

$$= 0.940$$

Prob of Class (from values in feature)

Prob of Feature (in values)

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

2. คำนวณ Gain ของ Feature แต่ละตัว ปัจจุบันไปด้วย ( $\text{Info}_A(D) = \sum_{j=1}^v \frac{|D_j|}{|D|} \times \text{Info}(D_j)$ )

ซึ่งจะแบ่งตามความน่าจะเป็นของแต่ละ Feature ไป Class

2.1) Gain ของ age

$$\text{Info}_{\text{age}}(D) = \frac{7}{14} I(2,3) + \frac{4}{14} I(4,0) + \frac{5}{14} I(3,2)$$

$$= \frac{7}{14} \left[ -\frac{2}{5} \log_2 \left( \frac{2}{5} \right) - \frac{3}{5} \log_2 \left( \frac{3}{5} \right) \right] + \frac{4}{14} \left[ -\frac{4}{4} \log_2 \left( \frac{4}{4} \right) \right] + \frac{5}{14} \left[ -\frac{3}{5} \log_2 \left( \frac{3}{5} \right) - \frac{2}{5} \log_2 \left( \frac{2}{5} \right) \right]$$

$$= 0.694$$

Feature

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

Class

2.2) Gain ของ income

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

$$= \frac{4}{14} \left[ -\frac{1}{4} \log_2 \left( \frac{1}{4} \right) - \frac{1}{4} \log_2 \left( \frac{1}{4} \right) \right] + \frac{6}{14} \left[ -\frac{4}{6} \log_2 \left( \frac{4}{6} \right) - \frac{1}{6} \log_2 \left( \frac{1}{6} \right) \right] + \frac{4}{14} \left[ -\frac{3}{4} \log_2 \left( \frac{3}{4} \right) - \frac{1}{4} \log_2 \left( \frac{1}{4} \right) \right]$$

$$= 0.911$$

2.3 Gain ของ student

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

$$= \frac{7}{14} \left[ -\frac{3}{7} \log_2 \left( \frac{3}{7} \right) - \frac{4}{7} \log_2 \left( \frac{4}{7} \right) \right] + \frac{7}{14} \left[ -\frac{6}{7} \log_2 \left( \frac{6}{7} \right) - \frac{1}{7} \log_2 \left( \frac{1}{7} \right) \right]$$

$$= 0.489$$

2.4 Gain ของ credit\_rating

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

$$= \frac{8}{14} \left[ -\frac{6}{8} \log_2 \left( \frac{6}{8} \right) - \frac{2}{8} \log_2 \left( \frac{2}{8} \right) \right] + \frac{6}{14} \left[ -\frac{3}{6} \log_2 \left( \frac{3}{6} \right) - \frac{3}{6} \log_2 \left( \frac{3}{6} \right) \right]$$

$$= 0.892$$

Gain of Class  
↓  
Gain of Feature  
↓

3. คำนวณ Information Gain ของต特徵: Feature ( $\text{Gain}(A) = \text{Info}(D) - \text{Info}_A(D)$ )

ได้ยัง: เลือก Gain ที่สูงที่สุด ให้เป็น Root Node

3.1  $\text{Gain}(\text{age}) = 0.940 - 0.911 = 0.246$  Gain สูงที่สุดให้เป็น Root Node

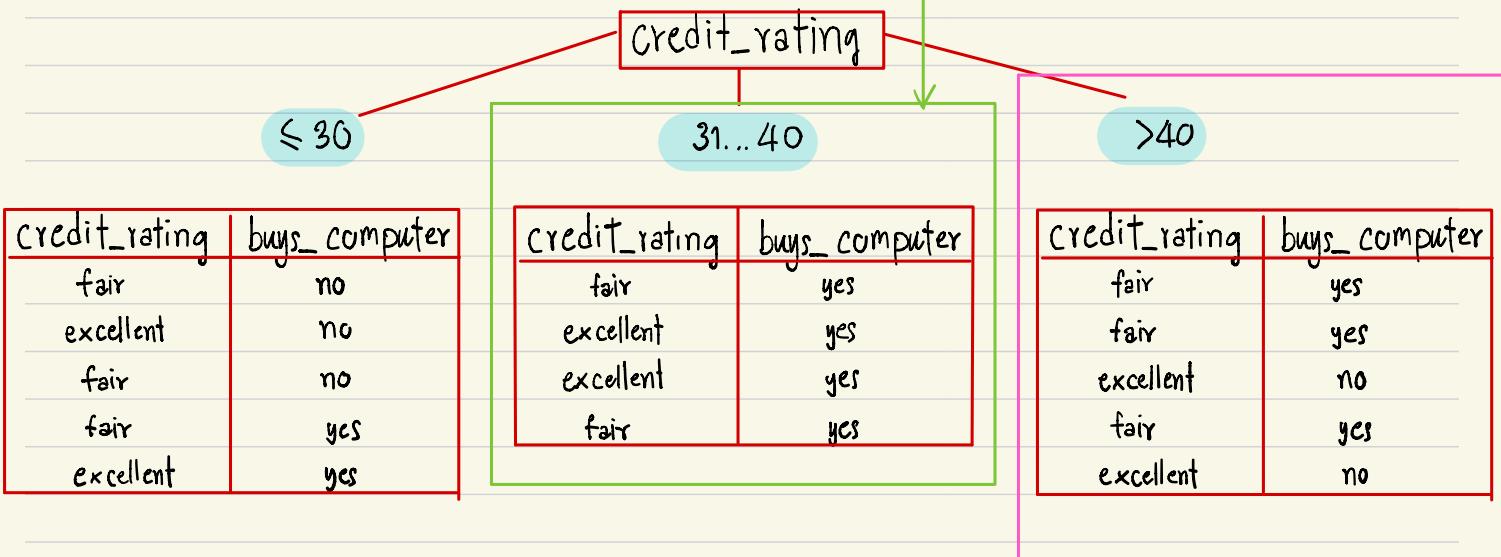
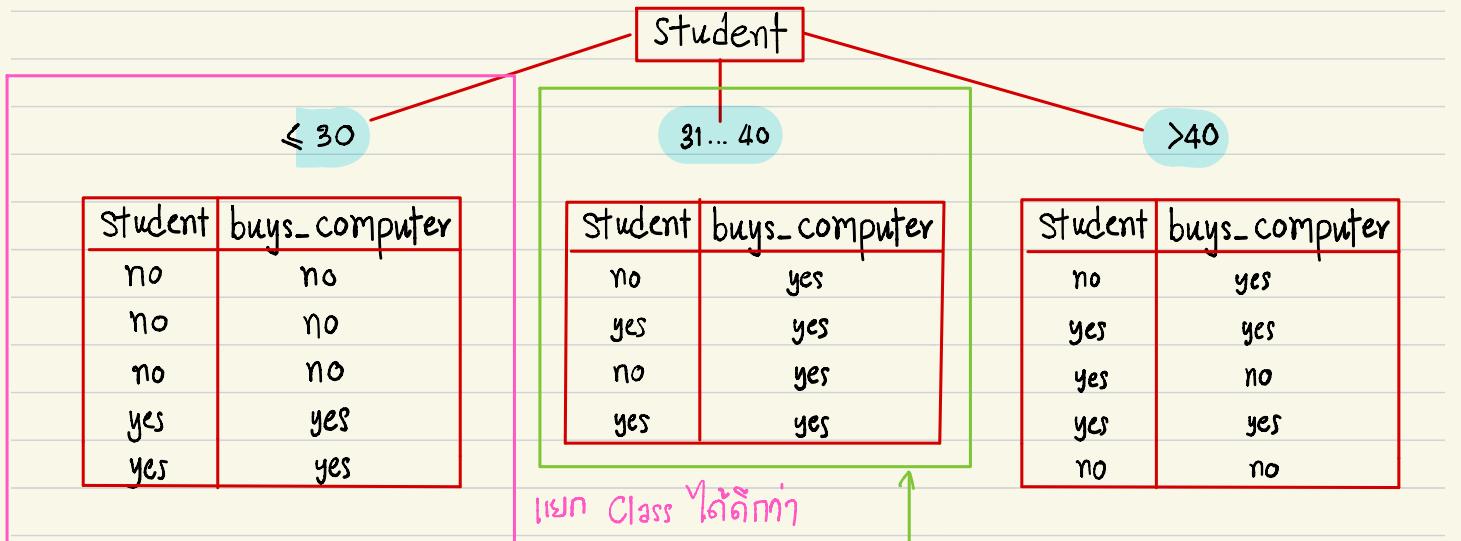
3.2  $\text{Gain}(\text{income}) = 0.940 - 0.911 = 0.029$  ทำต่อตั้ง เพราะ Decision Tree อาจไม่ไปสิ่งใดๆ

3.3  $\text{Gain}(\text{student}) = 0.940 - 0.489 = 0.451$

3.4  $\text{Gain}(\text{credit\_rating}) = 0.940 - 0.892 = 0.048$

### Decision Tree

1) แยกกลุ่มของ Feature มากที่สุดใน Root Node



2. តើនូវ Values ទាក់ទងនិង Class របស់ Feature នៅរីយៈពេល Decision Tree ។

