# HW 5

- หา node ที่ root
  กากสูตร

$$\text{Info}(D) = -\sum_{i=1}^{n} P_i \log_2(P_i)$$

จะได้

$$\text{Info}(D) = I(9,5) \quad ; \quad 9 = \text{yes}, \quad 5 = \text{no}$$

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

หาค่า Info ในแต่ละ feature

- $\text{Info}_{age}(D) = \frac{5}{14} \underbrace{I(2,3)}_{<=30} + \frac{4}{14} \underbrace{I(4,0)}_{31.40} + \frac{5}{14} \underbrace{I(3,2)}_{>40}$

$$= \frac{5}{14}\left[-\frac{2}{5}\log_2\left(\frac{2}{5}\right) - \frac{3}{5}\log_2\left(\frac{3}{5}\right)\right] + \frac{4}{14}\left[-\frac{4}{4}\log_2\left(\frac{4}{4}\right) - \frac{0}{4}\log_2\left(\frac{0}{4}\right)\right] + \frac{5}{14}\left[-\frac{3}{5}\log_2\left(\frac{3}{5}\right) - \frac{2}{5}\log_2\left(\frac{2}{5}\right)\right]$$

$$= 0.694$$

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

$$= \frac{4}{14}\left[-\frac{2}{4}\log_2\left(\frac{2}{4}\right) - \frac{2}{4}\log_2\left(\frac{2}{4}\right)\right] + \frac{6}{14}\left[-\frac{4}{6}\log_2\left(\frac{4}{6}\right) - \frac{2}{6}\log_2\left(\frac{2}{6}\right)\right] + \frac{4}{14}\left[-\frac{3}{4}\log_2\left(\frac{3}{4}\right) - \frac{1}{4}\log_2\left(\frac{1}{4}\right)\right]$$

$$= 0.911$$

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

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

$$= 0.789$$

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

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

$$= 0.892$$

$$\text{Gain}(age) = \text{Info}(D) - \text{Info}_{age}(D) = 0.940 - 0.694 = 0.246 \quad ※$$

$$\text{Gain}(income) = \text{Info}(D) - \text{Info}_{income}(D) = 0.940 - 0.911 = 0.029$$

$$\text{Gain}(student) = \text{Info}(D) - \text{Info}_{student}(D) = 0.940 - 0.789 = 0.151$$

$$\text{Gain}(credit\_) = \text{Info}(D) - \text{Info}_{credit\_}(D) = 0.940 - 0.892 = 0.048$$

○ หาค่า Info ในแต่ละกิ่ง
- $\le 30$

| age | income | student | credit_rating | buy |
|---|---|---|---|---|
| | high | no | fair | no |
| | high | no | excellent | no |
| $\le 30$ | medium | no | fair | no |
| | low | yes | fair | yes |
| | medium | yes | excellent | yes |

$$Info(D) = I(2,3) \quad ; \quad 2 = yes, \ 3 = no$$
$$= -\frac{2}{5}\log_2\left(\frac{2}{5}\right) - \frac{3}{5}\log_2\left(\frac{3}{5}\right) = 0.971$$

$$Info_{income}(D) = \overbrace{\frac{2}{5}I(0,2)}^{high} + \overbrace{\frac{2}{5}I(1,1)}^{medium} + \overbrace{\frac{1}{5}I(1,0)}^{low}$$

$$= \frac{2}{5}\left[-\frac{0}{2}\log_2\left(\frac{0}{2}\right) - \frac{2}{2}\log_2\left(\frac{2}{2}\right)\right] + \frac{2}{5}\left[-\frac{1}{2}\log_2\left(\frac{1}{2}\right) - \frac{1}{2}\log_2\left(\frac{1}{2}\right)\right] + \frac{1}{5}\left[-\frac{1}{1}\log_2\left(\frac{1}{1}\right) - \frac{0}{1}\log_2\left(\frac{0}{1}\right)\right]$$

$$= 0.4$$

$$Info_{student}(D) = \overbrace{\frac{2}{5}I(2,0)}^{yes} + \overbrace{\frac{3}{5}I(0,3)}^{no}$$

$$= \frac{2}{5}\left[-\frac{2}{2}\log_2\left(\frac{2}{2}\right) - \frac{0}{2}\log_2\left(\frac{0}{2}\right)\right] + \frac{3}{5}\left[-\frac{0}{3}\log_2\left(\frac{0}{3}\right) - \frac{3}{3}\log_2\left(\frac{3}{3}\right)\right]$$

$$= 0$$

$$Info_{credit}(D) = \overbrace{\frac{3}{5}I(1,2)}^{fair} + \overbrace{\frac{2}{5}I(1,1)}^{excellent}$$

$$= \frac{3}{5}\left[-\frac{1}{3}\log_2\left(\frac{1}{3}\right) - \frac{2}{3}\log_2\left(\frac{2}{3}\right)\right] + \frac{2}{5}\left[-\frac{1}{2}\log_2\left(\frac{1}{2}\right) - \frac{1}{2}\log_2\left(\frac{1}{2}\right)\right]$$

$$= 0.951$$

$$Gain(income) = Info(D) - Info_{income}(D) = 0.971 - 0.4 = 0.571$$
$$Gain(student) = Info(D) - Info_{student}(D) = 0.971 - 0 = 0.971 \ ※$$
$$Gain(credit\_) = Info(D) - Info_{credit\_}(D) = 0.971 - 0.951 = 0.02$$

- $31...40$

| age | income | student | credit | buys |
|---|---|---|---|---|
| | high | no | fair | yes |
| $31...40$ | low | yes | excellent | yes |
| | medium | no | excellent | yes |
| | high | yes | fair | yes |

※ จากตาราง สามารถสรุปได้ว่า ช่วงอายุ $31...40$ มีข้อมูล buys_computer เป็น yes ทั้งหมด

— $> 40$

$$\text{Info}(D) = I(3,2) \quad ; \quad 3 = yes, \quad 2 = no$$
$$= -\frac{3}{5}\log_2\left(\frac{3}{5}\right) - \frac{2}{5}\log_2\left(\frac{2}{5}\right) = 0.971$$

| age | income | student | credit_rating | buy |
|-----|--------|---------|---------------|-----|
|     | medium | no      | fair          | yes |
|     | low    | yes     | fair          | yes |
| >40 | low    | yes     | excellent     | no  |
|     | medium | yes     | fair          | yes |
|     | medium | no      | excellent     | no  |

$$\text{Info}_{income}(D) = \overbrace{\frac{3}{5}I(2,1)}^{medium} + \overbrace{\frac{2}{5}I(1,1)}^{low}$$

$$= \frac{3}{5}\left[-\frac{2}{3}\log_2\left(\frac{2}{3}\right) - \frac{1}{3}\log_2\left(\frac{1}{3}\right)\right] + \frac{2}{5}\left[-\frac{1}{2}\log_2\left(\frac{1}{2}\right) - \frac{1}{2}\log_2\left(\frac{1}{2}\right)\right]$$

$$= 0.951$$

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

$$= \frac{3}{5}\left[-\frac{2}{3}\log_2\left(\frac{2}{3}\right) - \frac{1}{3}\log_2\left(\frac{1}{3}\right)\right] + \frac{2}{5}\left[-\frac{1}{2}\log_2\left(\frac{1}{2}\right) - \frac{1}{2}\log_2\left(\frac{1}{2}\right)\right]$$

$$= 0.951$$

$$\text{Info}_{credit}(D) = \overbrace{\frac{3}{5}I(3,0)}^{fair} + \overbrace{\frac{2}{5}I(2,0)}^{excellent}$$

$$= \frac{3}{5}\left[-\frac{3}{3}\log_2\left(\frac{3}{3}\right) - \frac{0}{3}\log_2\left(\frac{0}{3}\right)\right] + \frac{2}{5}\left[-\frac{2}{2}\log_2\left(\frac{2}{2}\right) - \frac{0}{2}\log_2\left(\frac{0}{2}\right)\right]$$

$$= 0$$

$$\text{Gain}(income) = \text{Info}(D) - \text{Info}_{income}(D) = 0.971 - 0.951 = 0.02$$
$$\text{Gain}(student) = \text{Info}(D) - \text{Info}_{student}(D) = 0.971 - 0.951 = 0.02$$
$$\text{Gain}(credit\_) = \text{Info}(D) - \text{Info}_{credit\_}(D) = 0.971 - 0 = 0.971 \quad ※$$

## Descision Tree Induction

- **Resulting tree:**