## 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 |

ນນ Info (D)

$$Info(D) = I(9,5) = \left( -\frac{9}{14} \log_{(2)}\left(\frac{9}{14}\right) \right) - \left(\frac{5}{14} \log_{(2)}\left(\frac{5}{14}\right)\right)$$

$$= 0.94$$

วิธี Info$_{age}$ (D)

<= 30    31 - 40    > 40

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

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

$$I(4,0) = -\frac{4}{4}\log_{(2)}\left(\frac{4}{4}\right) - \frac{0}{4}\log_{(2)}\left(\frac{0}{4}\right) = 0$$

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

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

หา Gain (age)

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

---

วิธี Info$_{income}$ (D)

high    medium    low

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

$$I(2,2) = -\frac{2}{4}\log_{(2)}\left(\frac{2}{4}\right) - \frac{2}{4}\log_{(2)}\left(\frac{2}{4}\right) = 1$$

$$I(4,2) = -\frac{4}{6}\log_{(2)}\left(\frac{4}{6}\right) - \frac{2}{6}\log_{(2)}\left(\frac{2}{6}\right) = 0.918$$

$$I(3,1) = -\frac{3}{4}\log_{(2)}\left(\frac{3}{4}\right) - \frac{1}{4}\log_{(2)}\left(\frac{1}{4}\right) = 0.811$$

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

หา Gain (income)

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

---

วิธี Info$_{student}$ (D)

yes    No

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

$$I(6,1) = -\frac{6}{7}\log_{(2)}\left(\frac{6}{7}\right) - \frac{1}{7}\log_{(2)}\left(\frac{1}{7}\right) = 0.592$$

$$I(3,4) = -\frac{3}{7}\log_{(2)}\left(\frac{3}{7}\right) - \frac{4}{7}\log_{(2)}\left(\frac{4}{7}\right) = 0.985$$

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

หา Gain (student)

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

หา Info$_{credit\_rating}$ (D)

$$Info_{credit\_rating} (D) = \boxed{\frac{6}{14} I(\overset{Y}{6},\overset{N}{2})}^{fair} + \boxed{\frac{6}{14} I(\overset{Y}{3},\overset{N}{3})}^{excellent}$$

$$I(\overset{Y}{6},\overset{N}{2}) = -\frac{6}{8} log_{(2)}\left(\frac{6}{8}\right) - \frac{2}{8} log_{(2)}\left(\frac{2}{8}\right) = 0.8111$$

$$I(\overset{Y}{3},\overset{N}{3}) = -\frac{3}{6} log_{(2)}\left(\frac{3}{6}\right) - \frac{3}{6} log_{(2)}\left(\frac{3}{6}\right) = 1$$

แทนค่า Info$_{credit\_rating}$ (D) $= \frac{6}{14}(0.8111) + \frac{6}{14}(1) = 0.892$

หา Gain (credit_rating)

$$Gain (credit\_rating) = 0.94 - 0.892 = 0.048$$

หา Gain

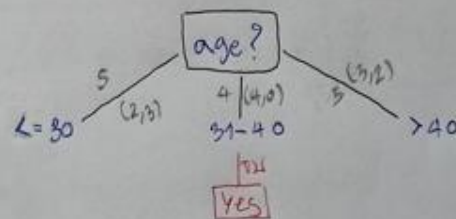| | |
|---|---|
| Gain (age) | = 0.246 |
| Gain (income) | = 0.029 |
| Gain (student) | = 0.151 |
| Gain (credit_rating) | = 0.048 |

เลือก Gain ที่มีมากที่สุด มาทำรากราก เป็นต้นไม้ แรก ซึ่งในที่นี้ คือ Gain (age)



age(<=30)

หา Info (D) ของ age(<=30)

$$Info (D) = I(\overset{Y}{2},\overset{N}{3}) = 0.971 \quad * คนมาเว็บบ้าง$$

## หา Info_income (D) ของ age

$$Info_{income}(D) \text{ ของ } age(<=30) = \boxed{\frac{2}{5}I(\overset{Y}{0},\overset{N}{2})}^{high} + \boxed{\frac{2}{5}I(\overset{Y}{1},\overset{N}{1})}^{medium} + \boxed{\frac{1}{5}I(\overset{Y}{1},\overset{N}{0})}^{low}$$

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

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

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

แทนค่า $info_{income}(D) \text{ ของ } age(<=30) = \frac{2}{5}(0) + \frac{2}{5}(1) + \frac{1}{5}(0) = 0.4$

## หา Gain (income) ของ age(<=30)

$$Gain(income) \text{ ของ } age(<=30) = 0.971 - 0.4 = 0.571$$

---

## หา Info_student (D) ของ age(<=30)

$$Info_{student}(D) \text{ ของ } age(<=30) = \boxed{\frac{2}{5}I(\overset{Y}{2},\overset{N}{0})}^{yes} + \boxed{\frac{3}{5}I(\overset{Y}{0},\overset{N}{3})}^{No}$$

จะได้ Yes → Yes (buy_computer) , No → no (buy_computer)

เลือก แนว ด้วย student เพราะ. สาขาที่ แนว ที่ จุด ได้ แม่นยำ มากที่สุด

---

## age(>40)

### หา Info (D) ของ age (>40)

$$Info(D) \text{ ของ } age(>40) = I(\overset{Y}{3},\overset{N}{2}) = 0.971 \quad * \text{ เพราะไม่บริสุทธิ์}$$

### หา Info_income (D) ของ age (>40)

$$Info_{income}(D) \text{ ของ } age(>40) = \boxed{\frac{3}{5}I(\overset{Y}{2},\overset{N}{1})}^{medium} + \boxed{\frac{2}{5}I(\overset{Y}{1},\overset{N}{1})}^{low}$$

$$I(\overset{Y}{2},\overset{N}{1}) = -\frac{2}{3}\log_{(2)}\left(\frac{2}{3}\right) - \frac{1}{3}\log\left(\frac{1}{3}\right) = 0.918$$

$$I(\overset{Y}{1},\overset{N}{1}) = 1$$

แทนค่า $Info_{income}(D) \text{ ของ } age(>40) = \frac{3}{5}(0.918) + \frac{2}{5}(1) = 0.951$

### หา Gain (income) ของ age (>40)

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

หา Info$_{student}$ (D) ของ age (>40)

$$\text{Info}_{student}(D) \text{ ของ } age(>40) = \boxed{\frac{3}{5} I(\overset{y}{2},\overset{N}{1})}^{yes} + \boxed{\frac{2}{5} I(\overset{y}{1},\overset{N}{1})}^{No}$$

$$I(\overset{y}{2},\overset{N}{1}) = -\frac{2}{3} \log_{(2)}\left(\frac{2}{3}\right) - \frac{1}{3}\log_{(2)}\left(\frac{1}{3}\right) = 0.916$$

$$I(\overset{y}{1},\overset{N}{1}) = 1$$

แทนค่า Info$_{student}$ (D) ของ age(>40) $= \frac{3}{5}(0.916) + \frac{2}{5}(1) = 0.951$

หา Gain (student) ของ age (>40)

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

---

หา Info$_{credit\_rating}$ (D) ของ age (>40)

$$\text{Info}_{credit\_rating}(D) \text{ ของ } age(>40) = \boxed{\frac{3}{5} I(3,0)}^{fair} + \boxed{\frac{2}{5} I(0,8)}^{excellent}$$

สังเกต fair → Yes (buy_Computer , excellent → No (buy_Computer)

เลือก แขนงด้วย  Credit_rating  เพราะ จำแนก แบ่งข้อมูลได้ ชัดเจน

---

สรุป