# Decision Tree Induction

## 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

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

Exaple:

| 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 |

## บ1 Info (D)

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

Y        N

$$= 0.94$$

## บ1 Info age (D)

               <= 30               31-40             > 40

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

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

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

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

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

## บ Gain (age)

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

## ผ Info_income (D)

$$\text{Info}_{income}(D) = \boxed{\frac{4}{14} \; I\left(\overset{Y}{2},\overset{N}{2}\right)} + \boxed{\frac{6}{14} \; I\left(\overset{Y}{4},\overset{N}{2}\right)} + \boxed{\frac{4}{14} \; I\left(\overset{Y}{3},\overset{N}{1}\right)}$$

(high, medium, low)

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

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

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

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

## และ Gain (income)

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

---

## ผ Info_student (D)

$$\text{Info}_{student}(D) = \boxed{\frac{7}{14} \; I\left(\overset{Y}{6},\overset{N}{1}\right)} + \boxed{\frac{7}{14}\left(\overset{Y}{3},\overset{N}{4}\right)}$$

(yes, No)

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

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

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

## ผ Gain (Student)

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

ค Info credit_rating (D)
_____

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

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

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

แทนค่า $\text{Info}_{credit\_rating}(D) = \frac{8}{14}(0.8111) + \frac{6}{14}(1) = 0.892$

ค Gain (credit_rating)

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

จาก Gain

Gain (age) $= 0.246$

Gain (income) $= 0.29$

Gain (student) $= 0.151$

Gain (credit_rating) $= 0.048$

เลือก Gain ที่ค่ามากที่สุดมาพิจารณาเป็นตัวแรก ซึ่งในที่นี้คือ Gain (age)



ทดสอบบทราย พรจินไพส์ 623021051-4 ฯ