# Decision - Tree

| 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 |

$14$

Yes = 9
No = 5

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

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

2. $\text{Info}_A(D) = \sum_{j=1}^{v} \frac{|D_j|}{|D|} \times \text{Info}(D_j)$ **Feature**

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

$= \frac{5}{14} \left[ -\frac{2}{5} \log_2 \left(\frac{2}{5}\right) - \frac{3}{5} \log_2 \left(\frac{3}{5}\right) \right] + \frac{4}{14} \left[ -\frac{4}{4} \log_2 \left(\frac{4}{4}\right) \right] + \frac{5}{14} \left[ -\frac{3}{5} \log_2 \left(\frac{3}{5}\right) - \frac{2}{5} \log_2 \left(\frac{2}{5}\right) \right] = 0.694$

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

$= \frac{4}{14} \left[ -\frac{2}{5} \log_2 \left(\frac{2}{4}\right) - \frac{2}{4} \log_2 \left(\frac{2}{4}\right) \right] + \frac{6}{14} \left[ -\frac{4}{6} \log_2 \left(\frac{4}{6}\right) - \frac{2}{6} \log_2 \left(\frac{2}{6}\right) \right] + \frac{4}{14} \left[ -\frac{3}{4} \log_2 \left(\frac{3}{4}\right) - \frac{1}{5} \log_2 \left(\frac{1}{5}\right) \right] = 0.911$

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

$= \frac{7}{14} \left[ -\frac{3}{7} \log_2 \left(\frac{3}{7}\right) - \frac{4}{7} \log_2 \left(\frac{4}{7}\right) \right] + \frac{7}{14} \left[ -\frac{1}{7} \log_2 \left(\frac{6}{7}\right) - \frac{1}{7} \log_2 \left(\frac{1}{7}\right) \right] = 0.789$

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

$= \frac{8}{14} \left[ -\frac{6}{8} \log_2 \left(\frac{6}{8}\right) - \frac{2}{8} \log_2 \left(\frac{2}{8}\right) \right] + \frac{6}{14} \left[ -\frac{3}{6} \log_2 \left(\frac{3}{6}\right) - \frac{3}{6} \log_2 \left(\frac{3}{6}\right) \right] = 0.892$

3. $\text{Gain}(A) = \text{Info}(D) - \text{Info}_A(D)$  <span style="color:blue">ค่า Gain สูงสุดจะเป็น root node</span>

   3.1 $\text{Gain}(\text{age}) = 0.940 - 0.694 = 0.246$ //

   3.2 $\text{Gain}(\text{income}) = 0.940 - 0.911 = 0.029$

   3.3 $\text{Gain}(\text{student}) = 0.940 - 0.789 = 0.151$

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

4. แยกกลุ่ม feature ตามค่าใน root node

   4.1 $< = 30$               คำนวน Inforntion Grain

   $\text{Info}(D) = I(2,3) = 0.971$

   $\text{Info}_{\text{income}}(D) = \frac{2}{3} I(0,2) + \frac{2}{5} I(1,1) + \frac{1}{5} I(1,0) = 0.4$

   $\text{Info}_{\text{student}}(D) = \frac{2}{5} I(2,0) + \frac{3}{5} I(0,3) = 0$

   $\text{Info}_{\text{credit\_rating}}(D) = \frac{3}{5} I(1,2) + \frac{2}{5} I(1,1) = 0.951$

   $\text{Grain}(\text{Income}) = 0.971 - 0.4 = 0.571$
   $\text{Grain}(\text{Student}) = 0.971 - 0 = 0.971$ //
   $\text{Grain}(\text{credit\_rating}) = 0.971 - 0.951 = 0.02$

   เลือก $\text{Grain}(\text{Student})$ เป็น Node

   4.2 $30...40$

   | age | income | student | credit_rating | buys_computer |
   |-----|--------|---------|---------------|---------------|
   | 31...40 | high | no | fair | yes |
   | 31...40 | low | yes | excellent | yes |
   | 31...40 | medium | no | excellent | yes |
   | 31...40 | high | yes | fair | yes |

   $yes = 4$     $no = 0$

   เมื่ออายุ $31...40$ ตอบ Yes

   ใน buys_computer

   4.3 $> 40$

   $\text{Info}(D) = I(3,2) = -\frac{3}{5}\log_2\left(\frac{3}{5}\right) - \frac{2}{5}\log_2\left(\frac{2}{5}\right) = 0.971$

   $\text{Info}_{(\text{income})}(D) = \frac{3}{5} I(2,1) + \frac{2}{5} I(1,1) = 0.951$

   $\text{Info}_{\text{student}}(D) = \frac{3}{5} I(2,1) + \frac{2}{5} I(1,1) = 0.951$

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

   คำนวน Inforntion Grain

   $\text{Grain}(\text{Income}) = 0.971 - 0.951 = 0.02$
   $\text{Grain}(\text{Student}) = 0.971 - 0.951 = 0.02$
   $\text{Grain}(\text{credit\_rating}) = 0.971 - 0 = 0.971$ //

   เลือก $\text{Grain}(\text{Credit\_rating})$ เป็น Node

5.

age
- ≤30 → Student
  - Yes → Buy
  - No → Not Buy
- 31...40 → Buy
- >40 → Credit_rating
  - excellent → Not-Buy
  - fair → Buy