

Homework 5

Decision Tree

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 0	yes	fair \	yes ✓
>40	low 0	yes	excellent ✓	no \
31...40	low 0	yes	excellent ✓	yes ✓
<=30	medium \	no ✓	fair \	no \
<=30	low 0	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 \

бадд root

Class

$$\begin{aligned}
 \text{Info}(D) &= - \sum_{i=1}^M p_i \log_2(p_i) \\
 &= I(\overset{N_{\text{yes}}}{9}, \overset{N_{\text{no}}}{5}) \\
 &= - \left(\frac{9}{14} \log_2 \frac{9}{14} \right) + \left(- \frac{5}{14} \log_2 \frac{5}{14} \right) \\
 &= - \frac{9}{14} \log_2 \frac{9}{14} - \frac{5}{14} \log_2 \frac{5}{14} \\
 &= - \frac{9}{14} (-0.637) - \frac{5}{14} (-1.485) \\
 &= 0.940 \#
 \end{aligned}$$

Feature

$$\begin{aligned}
 \text{Info}_{\text{age}}(D) &= \sum_{j=1}^v \left| \frac{D_j}{D} \right| \times \text{Info}(D_j) \\
 &= \overset{<= 30}{\frac{5}{14} I(2, 3)} + \overset{31 \dots 40}{\frac{4}{14} I(4, 0)} + \overset{> 40}{\frac{5}{14} I(3, 2)} \\
 &= \frac{5}{14} \left[-\frac{2}{5} \log_2 \left(\frac{2}{5} \right) - \frac{3}{5} \log_2 \left(\frac{3}{5} \right) \right] + \frac{4}{14} \left[-\frac{4}{4} \log_2 \left(\frac{4}{4} \right) - \frac{0}{4} \log_2 \left(\frac{0}{4} \right) \right] + \frac{5}{14} \left[-\frac{3}{5} \log_2 \left(\frac{3}{5} \right) - \frac{2}{5} \log_2 \left(\frac{2}{5} \right) \right] \\
 &= \frac{5}{14} (0.529 + 0.442) + \frac{4}{14} (0 + 0.0000000000) + \frac{5}{14} (0.442 + 0.529) \\
 &= \frac{5}{14} (0.971) + 0.0000000000 + \frac{5}{14} (0.971) \\
 &= 0.347 + 0.347 \\
 &= 0.694 \quad \#
 \end{aligned}$$

$$\begin{aligned}
 \text{Info}_{\text{income}}(D) &= \sum_{j=1}^v \left| \frac{D_j}{D} \right| \times \text{Info}(D_j) \\
 &= \overset{\text{high}}{\frac{4}{14} I(2, 2)} + \overset{\text{medium}}{\frac{6}{14} I(4, 2)} + \overset{\text{low}}{\frac{4}{14} I(3, 1)} \\
 &= \frac{4}{14} \left[-\frac{2}{4} \log_2 \left(\frac{2}{4} \right) - \frac{2}{4} \log_2 \left(\frac{2}{4} \right) \right] + \frac{6}{14} \left[-\frac{4}{6} \log_2 \left(\frac{4}{6} \right) - \frac{2}{6} \log_2 \left(\frac{2}{6} \right) \right] + \frac{4}{14} \left[-\frac{3}{4} \log_2 \left(\frac{3}{4} \right) - \frac{1}{4} \log_2 \left(\frac{1}{4} \right) \right] \\
 &= \frac{4}{14} (0.5 + 0.5) + \frac{6}{14} (0.390 + 0.528) + \frac{4}{14} (0.311 + 0.5) \\
 &= \frac{4}{14} + \frac{6}{14} (0.918) + \frac{4}{14} (0.811) \\
 &= 0.286 + 0.394 + 0.232 \\
 &= 0.912 \quad \#
 \end{aligned}$$

$$\begin{aligned}
 \text{Info}_{\text{student}}(D) &= \sum_{j=1}^v \left| \frac{D_j}{D} \right| \times \text{Info}(D_j) \\
 &= \overset{\text{No}}{\frac{7}{14} I(3, 4)} + \overset{\text{Yes}}{\frac{7}{14} I(6, 1)} \\
 &= \frac{7}{14} \left[-\frac{3}{7} \log_2 \left(\frac{3}{7} \right) - \frac{4}{7} \log_2 \left(\frac{4}{7} \right) \right] + \frac{7}{14} \left[-\frac{6}{7} \log_2 \left(\frac{6}{7} \right) - \frac{1}{7} \log_2 \left(\frac{1}{7} \right) \right] \\
 &= \frac{7}{14} (0.524 + 0.461) + \frac{7}{14} (0.191 + 0.401) \\
 &= \frac{7}{14} (0.985) + \frac{7}{14} (0.592) \\
 &= 0.493 + 0.296 \\
 &= 0.789 \quad \#
 \end{aligned}$$

Info_{credit}(D)

$$= \sum_{j=1}^3 \left| \frac{D_j}{D} \right| \times \text{Info}(D_j)$$

$$= \frac{9}{14} \overset{\text{excellent}}{I(6,2)} + \frac{6}{14} \overset{\text{Low}}{I(3,3)}$$

$$= \frac{9}{14} \left[-\frac{6}{9} \log_2 \left(\frac{6}{9} \right) - \frac{2}{9} \log_2 \left(\frac{2}{9} \right) \right] + \frac{6}{14} \left[-\frac{3}{6} \log_2 \left(\frac{3}{6} \right) - \frac{3}{6} \log_2 \left(\frac{3}{6} \right) \right]$$

$$= \frac{9}{14} (0.311 + 0.5) + \frac{6}{14} (0.5 + 0.5)$$

$$= \frac{9}{14} (0.811) + \frac{6}{14}$$

$$= 0.464 + 0.429$$

$$= 0.893 \#$$

Gain

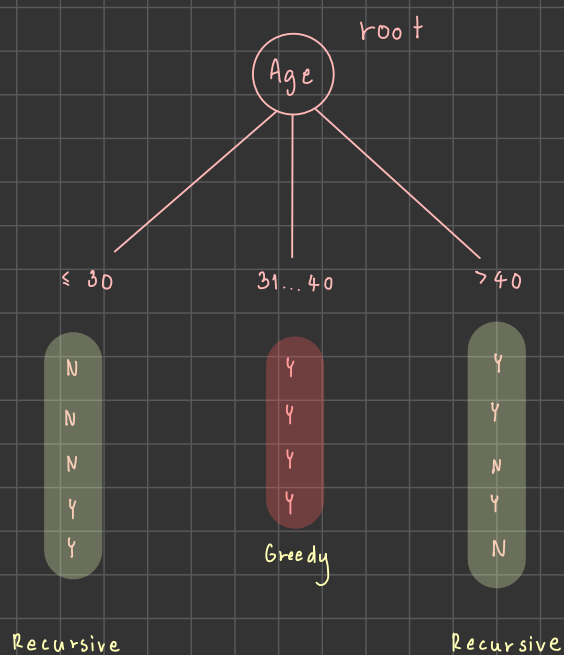
$$\text{age} = \text{Info}(D) - \text{Info}_{\text{age}}(D) = 0.940 - 0.649 = 0.291$$

$$\text{income} = \text{Info}(D) - \text{Info}_{\text{income}}(D) = 0.940 - 0.912 = 0.028$$

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

$$\text{credit_rating} = \text{Info}(D) - \text{Info}_{\text{credit}}(D) = 0.940 - 0.893 = 0.047$$

∴ เลือก age เป็น root เพราะ มีค่ามากที่สุด ซึ่งเป็นทางเลือกที่ดีที่สุด #



" Feature 1 Age ≤ 30 "

Class

$$\begin{aligned}
 \text{Info}(D) &= - \sum_{i=1}^M p_i \log_2(p_i) \\
 &= I(\overset{\text{Yes}}{2}, \overset{\text{No}}{3}) \\
 &= - \frac{2}{5} \log_2\left(\frac{2}{5}\right) - \frac{3}{5} \log_2\left(\frac{3}{5}\right) \\
 &= 0.971 \#
 \end{aligned}$$

Features

$$\begin{aligned}
 \text{Info}_{\text{income}}(D) &= \overset{\text{high}}{\frac{2}{5} I(0, 2)} + \overset{\text{medium}}{\frac{2}{5} I(1, 1)} + \overset{\text{low}}{\frac{1}{5} I(1, 0)} \\
 &= \frac{2}{5} \left[-\frac{0}{2} \log_2\left(\frac{0}{2}\right) - \frac{2}{2} \log_2\left(\frac{2}{2}\right) \right] + \frac{2}{5} \left[-\frac{1}{2} \log_2\left(\frac{1}{2}\right) - \frac{1}{2} \log_2\left(\frac{1}{2}\right) \right] + \frac{1}{5} \left[-\frac{1}{1} \log_2\left(\frac{1}{1}\right) - 2 \log_2(0) \right] \\
 &= 0.4 \#
 \end{aligned}$$

$$\begin{aligned}
 \text{Info}_{\text{student}}(D) &= \overset{\text{Yes}}{\frac{3}{5} I(0, 3)} + \overset{\text{No}}{\frac{2}{5} I(2, 0)} \\
 &= \frac{3}{5} \left[-\frac{0}{3} \log_2\left(\frac{0}{3}\right) - \frac{3}{3} \log_2\left(\frac{3}{3}\right) \right] + \frac{2}{5} \left[-\frac{2}{2} \log_2\left(\frac{2}{2}\right) - \frac{0}{2} \log_2\left(\frac{0}{2}\right) \right] \\
 &= 0 \#
 \end{aligned}$$

$$\begin{aligned}
 \text{Info}_{\text{credit}}(D) &= \overset{\text{fair}}{\frac{3}{5} I(1, 2)} + \overset{\text{excellent}}{\frac{2}{5} I(1, 1)} \\
 &= \frac{3}{5} \left[-\frac{1}{3} \log_2\left(\frac{1}{3}\right) - \frac{2}{3} \log_2\left(\frac{2}{3}\right) \right] + \frac{2}{5} \left[-\frac{1}{2} \log_2\left(\frac{1}{2}\right) - \frac{1}{2} \log_2\left(\frac{1}{2}\right) \right] \\
 &= 0.551 + 0.4 \\
 &= 0.951 \#
 \end{aligned}$$

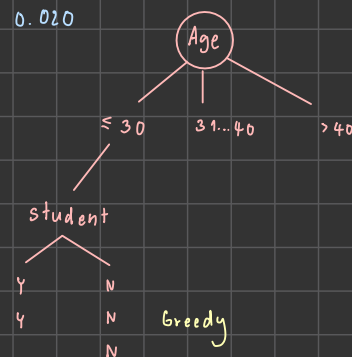
Gain

$$\text{income} \quad = \quad \text{Info}(D) - \text{Info}_{\text{income}}(D) \quad = \quad 0.971 - \quad = \quad$$

$$\text{student} \quad = \quad \text{Info}(D) - \text{Info}_{\text{student}}(D) \quad = \quad 0.971 - 0 \quad = \quad 0.971$$

$$\text{credit_rating} \quad = \quad \text{Info}(D) - \text{Info}_{\text{credit}}(D) \quad = \quad 0.971 - 0.951 \quad = \quad 0.020$$

\therefore เลือก student เพราะ มีค่า Gain สูง และเป็นหมวดเลือกที่เล็กที่สุด #



" Feature 3 Age > 40 "

Class

$$\begin{aligned} \text{Info}(D) &= - \sum_{i=1}^M p_i \log_2 p_i \\ &= I(3, 2) \\ &= - \frac{3}{5} \log_2 \left(\frac{3}{5} \right) - \frac{2}{5} \log_2 \left(\frac{2}{5} \right) \\ &= 0.971 \# \end{aligned}$$

Features

$$\begin{aligned} \text{Info}_{\text{income}}(D) &= \frac{3}{5} I(2, 1) + \frac{2}{5} I(1, 1) \\ &= \frac{3}{5} \left[- \frac{2}{3} \log_2 \left(\frac{2}{3} \right) - \frac{1}{3} \log_2 \left(\frac{1}{3} \right) \right] + \frac{2}{5} \left[- \frac{1}{2} \log_2 \left(\frac{1}{2} \right) - \frac{1}{2} \log_2 \left(\frac{1}{2} \right) \right] \\ &= 0.551 + 0.4 \\ &= 0.951 \# \end{aligned}$$

$$\begin{aligned} \text{Info}_{\text{student}}(D) &= \frac{2}{5} I(1, 1) + \frac{3}{5} I(2, 1) \\ &= \frac{2}{5} \left[- \frac{1}{2} \log_2 \left(\frac{1}{2} \right) - \frac{1}{2} \log_2 \left(\frac{1}{2} \right) \right] + \frac{3}{5} \left[- \frac{2}{3} \log_2 \left(\frac{2}{3} \right) - \frac{1}{3} \log_2 \left(\frac{1}{3} \right) \right] \\ &= 0.4 + 0.551 \\ &= 0.951 \# \end{aligned}$$

$$\begin{aligned} \text{Info}_{\text{credit}}(D) &= \frac{3}{5} I(3, 0) + \frac{2}{5} I(1, 1) \\ &= \frac{3}{5} \left[- \frac{3}{3} \log_2 \left(\frac{3}{3} \right) - \frac{0}{3} \log_2 \left(\frac{0}{3} \right) \right] + \frac{2}{5} \left[- \frac{1}{2} \log_2 \left(\frac{1}{2} \right) - \frac{1}{2} \log_2 \left(\frac{1}{2} \right) \right] \\ &= 0.4 \# \end{aligned}$$

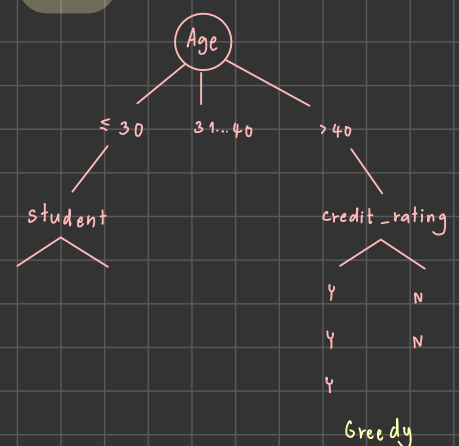
Gain

$$\text{income} \quad = \quad \text{Info}(D) - \text{Info}_{\text{income}}(D) \quad = \quad 0.971 - 0.951 \quad = \quad 0.2$$

$$\text{student} \quad = \quad \text{Info}(D) - \text{Info}_{\text{student}}(D) \quad = \quad 0.971 - 0.951 \quad = \quad 0.2$$

$$\text{credit_rating} \quad = \quad \text{Info}(D) - \text{Info}_{\text{credit}}(D) \quad = \quad 0.971 - 0.4 \quad = \quad 0.571$$

∴ เลือก credit เพราะ มีค่า Gain สูง ช่วยแบ่งกลุ่มข้อมูลที่ชัดเจน #



"Decision Tree Induction"

