

CSC4008 homework4: Decision Tree

Name: 费祥 Student ID: 120090414

T1,

Outlook case:

$$H_p = -\frac{9}{14} \log_2 \frac{9}{14} - \frac{5}{14} \log_2 \frac{5}{14} = 0.9403$$

$$H_{c1} = -\frac{2}{5} \log_2 \frac{2}{5} - \frac{3}{5} \log_2 \frac{3}{5} = 0.971$$

$$H_{c2} = 0$$

$$H_{c3} = -\frac{2}{5} \log_2 \frac{2}{5} - \frac{3}{5} \log_2 \frac{3}{5} = 0.971$$

$$H_c = \frac{5}{14} \times 0.971 + \frac{5}{14} \times 0.971 = 0.6936$$

$$\therefore \text{information gain} = H_p - H_c = 0.2467$$

$$GINI_p = 1 - \left(\left(\frac{9}{14} \right)^2 + \left(\frac{5}{14} \right)^2 \right) = 0.4592$$

$$GINI_{c1} = 1 - \left(\left(\frac{2}{5} \right)^2 + \left(\frac{3}{5} \right)^2 \right) = 0.48$$

$$GINI_{c2} = 1 - 1^2 = 0$$

$$GINI_{c3} = 1 - \left(\left(\frac{2}{5} \right)^2 + \left(\frac{3}{5} \right)^2 \right) = 0.48$$

$$GINI_c = \frac{5}{14} \times 0.48 + \frac{5}{14} \times 0.48 = 0.3429$$

$$\therefore \text{GINI index} = 0.4592 - 0.3429 = 0.1163$$

Humidity case:

$$H_p = -\frac{9}{14} \log_2 \frac{9}{14} - \frac{5}{14} \log_2 \frac{5}{14} = 0.9403$$

$$H_{c1} = -\frac{3}{7} \log_2 \frac{3}{7} - \frac{4}{7} \log_2 \frac{4}{7} = 0.9852$$

$$H_{c2} = -\frac{1}{7} \log_2 \frac{1}{7} - \frac{6}{7} \log_2 \frac{6}{7} = 0.5917$$

$$H_c = \frac{1}{2} \times 0.9852 + \frac{1}{2} \times 0.5917 = 0.7885$$

$$\therefore \text{information gain} = H_p - H_c = 0.1518$$

$$GINI_p = 1 - \left(\left(\frac{9}{14} \right)^2 + \left(\frac{5}{14} \right)^2 \right) = 0.4592$$

$$GINI_{c1} = 1 - \left(\left(\frac{3}{7} \right)^2 + \left(\frac{4}{7} \right)^2 \right) = 0.4898$$

$$GINI_{C_2} = 1 - \left(\left(\frac{1}{7} \right)^2 + \left(\frac{6}{7} \right)^2 \right) = 0.2449$$

$$GINI_C = 0.4898 \times \frac{1}{2} + 0.2449 \times \frac{1}{2} = 0.3674$$

$$\therefore GINI \text{ index} = 0.4592 - 0.3674 = 0.0918$$

Wind case:

$$H_P = -\frac{9}{14} \log_2 \frac{9}{14} - \frac{5}{14} \log_2 \frac{5}{14} = 0.9403$$

$$H_{C_1} = -\frac{3}{6} \log_2 \frac{3}{6} - \frac{3}{6} \log_2 \frac{3}{6} = 1$$

$$H_{C_2} = -\frac{6}{8} \log_2 \frac{6}{8} - \frac{2}{8} \log_2 \frac{2}{8} = 0.8113$$

$$H_C = \frac{6}{14} \times 1 + \frac{8}{14} \times 0.8113 = 0.8922$$

$$\therefore \text{Information gain} = 0.9403 - 0.8922 = 0.0481$$

$$GINI_P = 1 - \left(\left(\frac{9}{14} \right)^2 + \left(\frac{5}{14} \right)^2 \right) = 0.4592$$

$$GINI_{C_1} = 1 - \left(\left(\frac{3}{6} \right)^2 + \left(\frac{3}{6} \right)^2 \right) = 0.5$$

$$GINI_{C_2} = 1 - \left(\left(\frac{6}{8} \right)^2 + \left(\frac{2}{8} \right)^2 \right) = 0.375$$

$$GINI_C = \frac{6}{14} \times 0.5 + \frac{8}{14} \times 0.375 = 0.4286$$

$$\therefore GINI \text{ index} = 0.4592 - 0.4286 = 0.0306$$

T₂,

$$H_P = -\frac{3}{10} \log_2 \frac{3}{10} - \frac{7}{10} \log_2 \frac{7}{10} = 0.8813$$

$$GINI_P = 1 - \left(\left(\frac{3}{10} \right)^2 + \left(\frac{7}{10} \right)^2 \right) = 0.42$$

$$\textcircled{1} H_{C_1} = -\frac{2}{4} \log_2 \frac{2}{4} - \frac{2}{4} \log_2 \frac{2}{4} = 1$$

$$H_{C_2} = -\frac{1}{6} \log_2 \frac{1}{6} - \frac{5}{6} \log_2 \frac{5}{6} = 0.65$$

$$H_C = \frac{4}{10} \times 1 + \frac{6}{10} \times 0.65 = 0.79$$

$$\therefore \text{Information gain} = 0.8813 - 0.79 = 0.0913$$

$$GINI_{C_1} = 1 - \left(\left(\frac{2}{4} \right)^2 + \left(\frac{2}{4} \right)^2 \right) = 0.5$$

$$GINI_{C_2} = 1 - \left(\left(\frac{1}{6} \right)^2 + \left(\frac{5}{6} \right)^2 \right) = 0.2778$$

$$GINI_C = \frac{4}{10} \times 0.5 + \frac{6}{10} \times 0.2778 = 0.3667$$

$$\therefore GINI \text{ index} = 0.42 - 0.3667 = 0.0533$$

$$\textcircled{2} H_A = 0$$

$$H_{C_2} = -\frac{3}{6} \log_2 \frac{3}{6} - \frac{3}{6} \log_2 \frac{3}{6} = 1$$

$$H_C = \frac{4}{10} \times 0 + \frac{6}{10} \times 1 = 0.6$$

$$\therefore \text{information gain} = 0.8813 - 0.6 = 0.2813$$

$$GINI_{C_1} = 1 - 1^2 = 0$$

$$GINI_{C_2} = 1 - \left(\left(\frac{3}{6} \right)^2 + \left(\frac{3}{6} \right)^2 \right) = 0.5$$

$$GINI_C = \frac{4}{10} \times 0 + \frac{6}{10} \times 0.5 = 0.3$$

$$\therefore GINI \text{ Index} = 0.42 - 0.3 = 0.12$$

$$\textcircled{3} H_{C_1} = -\frac{1}{2} \log_2 \frac{1}{2} - \frac{1}{2} \log_2 \frac{1}{2} = 1$$

$$H_{C_2} = -\frac{2}{8} \log_2 \frac{2}{8} - \frac{6}{8} \log_2 \frac{6}{8} = 0.8113$$

$$H_C = \frac{2}{10} \times 1 + \frac{8}{10} \times 0.8113 = 0.849$$

$$\therefore \text{information gain} = 0.8813 - 0.849 = 0.0323$$

$$GINI_{C_1} = 1 - \left(\left(\frac{1}{2} \right)^2 + \left(\frac{1}{2} \right)^2 \right) = 0.5$$

$$GINI_{C_2} = 1 - \left(\left(\frac{2}{8} \right)^2 + \left(\frac{6}{8} \right)^2 \right) = 0.375$$

$$GINI_C = \frac{2}{10} \times 0.5 + \frac{8}{10} \times 0.375 = 0.4$$

$$\therefore GINI \text{ index} = 0.42 - 0.4 = 0.02$$

since the information gain and gini index in case $\textcircled{2}$ is the largest, so the second way is the best.

T₃.

