

P1

$$\text{OUTLOOK}$$

S	/	O
Y	2	Y 4
N	3	N 0

$$\text{Gini}(S) = 1 - \left(\frac{2}{5}\right)^2 - \left(\frac{3}{5}\right)^2 = 0.48$$

$$\text{Gini}(O) = 0$$

$$\text{Gini}(R) = 1 - \left(\frac{2}{5}\right)^2 - \left(\frac{3}{5}\right)^2 = 0.48$$

$$\text{Gini}(\text{OUTLOOK}) = \frac{5}{14} \cdot 0.48 + \frac{5}{14} \cdot 0.48 = 0.342$$

$$\text{TEMPER}$$

H	/	M
Y	2	Y 4
N	2	N 2

$$\text{Gini}(H) = 1 - \left(\frac{1}{2}\right)^2 - \left(\frac{1}{2}\right)^2 = \frac{1}{2}$$

$$\text{Gini}(M) = 1 - \left(\frac{4}{6}\right)^2 - \left(\frac{2}{6}\right)^2 = \frac{4}{9}$$

$$\text{Gini}(C) = 1 - \left(\frac{1}{4}\right)^2 - \left(\frac{1}{4}\right)^2 = 0.375$$

$$\text{Gini}(\text{TEMPER}) = \frac{4}{14} \cdot \frac{1}{2} + \frac{6}{14} \cdot \frac{4}{9} + \frac{4}{14} \cdot 0.375 = 0.427$$

$$\text{Humidity}$$

H	/	N
Y	3	Y 6
N	4	1

$$\text{Gini}(H) = 1 - \left(\frac{2}{7}\right)^2 - \left(\frac{5}{7}\right)^2 = 0.489$$

$$\text{Gini}(N) = 1 - \left(\frac{1}{7}\right)^2 - \left(\frac{6}{7}\right)^2 = 0.244$$

$$\text{Gini}(\text{Humidity}) = \frac{7}{14} \cdot 0.489 + \frac{7}{14} \cdot 0.244 = 0.366$$

$$\text{Wind}$$

W	/	S
Y	6	Y 3
N	2	N 3

$$\text{Gini}(W) = 1 - \left(\frac{6}{8}\right)^2 - \left(\frac{2}{8}\right)^2 = 0.375$$

$$\text{Gini}(S) = 1 - \left(\frac{1}{2}\right)^2 - \left(\frac{1}{2}\right)^2 = \frac{1}{2}$$

$$\text{Gini}(\text{Wind}) = \frac{8}{14} \cdot 0.375 + \frac{6}{14} \cdot \frac{1}{2} = 0.428$$

$\therefore \text{Gini}(\text{Outlook}) < \text{Gini}(\text{Humidity}) < \text{Gini}(\text{Temp}) < \text{Gini}(\text{wind})$

Root node

Sunny =

$$\text{TEMPER}$$

H	/	M
Y	0	Y 1
N	2	N 1

$$\text{Gini}(H) = 0$$

$$\text{Gini}(M) = \frac{1}{2}$$

$$\text{Gini}(C) = 0$$

$$\text{Gini}(\text{Sunny-Temp}) = \frac{2}{5} \cdot \frac{1}{2} = \frac{1}{5} = 0.2$$

$$\text{Humidity}$$

H	/	N
Y	0	Y 2
N	3	N 0

$$\text{Gini}(H) = 0$$

$$\text{Gini}(N) = 0$$

$$\text{Gini}(\text{Sunny-Humidity}) = 0$$

$$\text{Wind}$$

W	/	S
Y	1	Y 1
N	2	N 1

$$\text{Gini}(W) = 1 - \frac{1}{9} - \frac{4}{9} = \frac{4}{9}$$

$$\text{Gini}(S) = 1 - \frac{1}{4} - \frac{1}{4} = 0.5$$

$$\text{Gini}(\text{Wind}) = \frac{3}{5} \cdot \frac{4}{9} + \frac{2}{5} \cdot \frac{1}{2} = 0.466$$

Rainy.

$$\text{TEMPER}$$

Mild	/	Wet
Y	1	Y 1
N	1	N 1

$$\text{Gini}(M) = 1 - \frac{4}{9} - \frac{1}{9} = 0.444$$

$$\text{Gini}(W) = 1 - \frac{1}{9} - \frac{1}{9} = 0.5$$

$$\text{Gini}(R-T) = \frac{3}{5} \cdot \frac{2}{5} + \frac{2}{5} \cdot 0.5 = 0.466$$

$$\text{Humidity}$$

H	/	N
Y	1	Y 2
N	1	N 1

$$\text{Gini}(H) = 1 - \frac{1}{4} - \frac{1}{4} = 0.5$$

$$\text{Gini}(N) = 1 - \frac{4}{9} - \frac{1}{9} = 0.444$$

P2:

Name = Yanan Liu

NetID = yl2248

2.1:

$$H(x) = -\sum p_i \log_2 p_i$$

$$H(x) = -\frac{4}{15} \log_2 \frac{4}{15} - \frac{6}{15} \log_2 \frac{6}{15} = 0.97$$

2.2:

$$H(SZ=S) = -\frac{3}{4} \log_2 \frac{3}{4} - \frac{1}{4} \log_2 \frac{1}{4} = 0.81$$

$$H(SZ=N) = -\frac{2}{5} \log_2 \frac{2}{5} - \frac{3}{5} \log_2 \frac{3}{5} = 0.917$$

$$H(SZ=L) = -\frac{1}{3} \log_2 \frac{1}{3} - \frac{2}{3} \log_2 \frac{2}{3} = 0.917$$

$$H(SZ=XL) = 0$$

$$H(SZ) = \frac{4}{15} \cdot 0.811 + \frac{6}{15} \cdot 1 + \frac{3}{15} \cdot 0.917 + 0 = 0.799$$

$$IG_1 = 0.97 - 0.799 = 0.170$$

$$H(SZF=S) = -\frac{3}{4} \log_2 \frac{3}{4} - \frac{1}{4} \log_2 \frac{1}{4} = 0.811$$

$$H(SZF=N) = -\frac{2}{5} \log_2 \frac{2}{5} - \frac{3}{5} \log_2 \frac{3}{5} = 0.917$$

$$H(SZ=L) = -\frac{1}{3} \log_2 \frac{1}{3} - \frac{2}{3} \log_2 \frac{2}{3} = 0.917$$

$$H(SZ=XL) = 0$$

$$H(SZ=LM) = -\frac{1}{3} \log_2 \frac{1}{3} - \frac{2}{3} \log_2 \frac{2}{3} = 0.917$$

$$H(SZF) = \frac{4}{15} \cdot 0.811 + \frac{3}{15} \cdot 0.917 + \frac{3}{15} \cdot 0.917 + \frac{2}{15} \cdot 0.917 = 0.766$$

$$IG_1(SZF) = 0.203$$

$$IG_1(SZ) < IG_1(SZF)$$

2.3:

$$\text{Since the equation } "x^2+y^2 \geq 2xy" \Rightarrow \frac{x+y}{2} \geq \sqrt{xy}$$

the logarithm is monotone increasing

$$\therefore \log\left(\frac{x+y}{2}\right) \geq \log\sqrt{xy} = \log\sqrt{x} + \log\sqrt{y} = \frac{1}{2}\log x + \frac{1}{2}\log y$$

$$\therefore \log\left(\frac{x}{2} + \frac{y}{2}\right) \geq \frac{1}{2}\log x + \frac{1}{2}\log y.$$

$\therefore$  when  $\lambda = \frac{1}{2}$ , the equation holds.

2.4:

$$\text{let set } f(x) = x \log x$$

$$f'(x) = \log x + \frac{1}{x}$$

$$f''(x) = \frac{1}{x^2}$$

$$\because x > 0 \quad \therefore f''(x) > 0$$

$\therefore f(x)$  is convex

$$\begin{aligned} x_1 \log \frac{x_1}{x_1+x_2} + x_2 \log \frac{x_2}{x_1+x_2} &= x_1 \cdot \frac{1}{x_1+x_2} \log \frac{x_1}{x_1+x_2} + x_2 \cdot \frac{1}{x_1+x_2} \log \frac{x_2}{x_1+x_2} \\ &= y_1 \cdot f\left(\frac{x_1}{x_1+x_2}\right) + y_2 \cdot f\left(\frac{x_2}{x_1+x_2}\right) \\ &= (y_1+y_2) \cdot \left[ \frac{y_1}{y_1+y_2} f\left(\frac{x_1}{x_1+x_2}\right) + \frac{y_2}{y_1+y_2} f\left(\frac{x_2}{x_1+x_2}\right) \right] \end{aligned}$$

From Jensen's equation:-

$$t f(x_1) + (1-t) f(x_2) \geq f(tx_1 + (1-t)x_2), 0 \leq t \leq 1 \quad [\text{From Wikipedia}]$$

$$\therefore x_1 \log \frac{x_1}{x_1+x_2} + x_2 \log \frac{x_2}{x_1+x_2} = (y_1+y_2) \left[ \frac{y_1}{y_1+y_2} f\left(\frac{x_1}{x_1+x_2}\right) + \frac{y_2}{y_1+y_2} f\left(\frac{x_2}{x_1+x_2}\right) \right]$$

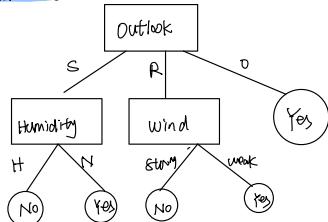
$$\geq (y_1+y_2) f\left(\frac{y_1}{y_1+y_2} \cdot \frac{x_1}{x_1+x_2} + \frac{y_2}{y_1+y_2} \cdot \frac{x_2}{x_1+x_2}\right)$$

$$= y_1+y_2 f\left(\frac{x_1+x_2}{y_1+y_2}\right) = (x_1+x_2) \log \frac{x_1+x_2}{y_1+y_2}$$

$$\therefore x_1 \log \frac{x_1}{x_1+x_2} + x_2 \log \frac{x_2}{x_1+x_2} \geq (x_1+x_2) \log \frac{x_1+x_2}{y_1+y_2} \quad \text{Proved.}$$

$Gini(R-W) = 1 \cdot \frac{2}{3} + 0.4 \cdot \frac{3}{5} = 0.66$

Decision tree:



P3:

3.1 Entropy:

$$H(X) = -\frac{2}{3} \log_2 \frac{2}{3} - \frac{1}{3} \log_2 \frac{1}{3} = 0.99$$

Feature 1:

$$H(S) = -\frac{3}{4} \log_2 \frac{3}{4} - \frac{1}{4} \log_2 \frac{1}{4} = 0.811$$

$$H(F) = -\frac{1}{5} \log_2 \frac{1}{5} - \frac{4}{5} \log_2 \frac{4}{5} = 0.72$$

$$H(F_1) = \frac{4}{5} \cdot 0.811 + \frac{1}{5} \cdot 0.72 = 0.76$$

$$IG(F_1) = 0.23$$

Feature 2:

$$H(T) = -\frac{2}{5} \log_2 \frac{2}{5} - \frac{3}{5} \log_2 \frac{3}{5} = 0.97$$

$$H(F_2) = -\frac{1}{2} \log_2 \frac{1}{2} - \frac{1}{2} \log_2 \frac{1}{2} = 1$$

$$H(F_2) = 0.971 \times \frac{5}{9} + 1 \times \frac{4}{9} = 0.982$$

$$IG(F_2) = 0.99 - 0.982 = 0.008$$

3.2:

$$H(\text{label}) = 0.99$$

threshold = 25

$$H(S) = 0$$

$$H(G_1) = -\frac{3}{8} \log_2 \frac{3}{8} - \frac{5}{8} \log_2 \frac{5}{8} = 0.953$$

$$H(25) = 0.953 \times \frac{8}{9} = 0.847$$

$$IG(TE25) = 0.143$$

threshold = 35

$$H(U) = 1$$

$$H(G_2) = -\frac{3}{7} \log_2 \frac{3}{7} - \frac{4}{7} \log_2 \frac{4}{7} = 0.985$$

$$IG(35) = 0.003$$

threshold = 45

$$H(V) = -\frac{1}{3} \log_2 \frac{1}{3} - \frac{2}{3} \log_2 \frac{2}{3} = 0.917$$

$$H(G_3) = -\frac{2}{6} \log_2 \frac{2}{6} - \frac{4}{6} \log_2 \frac{4}{6} = 0.917$$

$$H(45) = 0.917 \cdot \frac{3}{9} + \frac{6}{9} = 0.917$$

$$IG(45) = 0.073$$

2.5:

from the equation ② & ①, we have:

$$(x_1 + x_2) \log_2 \frac{x_1 + x_2}{n_m n_s} \geq x_1 \log_2 \frac{x_1}{n_m} + x_2 \log_2 \frac{x_2}{n_s}$$

$\therefore n_m^+ = n_m^m + n_m^+$

$$\therefore n_m^+ \log_2 \frac{n_m^+}{n_m n_s} + n_m^- \log_2 \frac{n_m^-}{n_m n_s} \geq n_m^+ \log_2 \frac{n_m^+}{n_m} \quad ①$$

$$n_m^+ \log_2 \frac{n_m^+}{n_m n_s} + n_m^- \log_2 \frac{n_m^-}{n_m n_s} \geq n_m^- \log_2 \frac{n_m^-}{n_m} \quad ②$$

When ① + ②, we will get:

$$n_m^+ \log_2 \frac{n_m^+}{n_m n_s} + n_m^- \log_2 \frac{n_m^-}{n_m n_s} + n_m^+ \log_2 \frac{n_m^+}{n_m} + n_m^- \log_2 \frac{n_m^-}{n_m}$$

$$\geq n_m^+ \log_2 \frac{n_m^+}{n_m} + n_m^- \log_2 \frac{n_m^-}{n_m}$$

$\therefore$  by divide  $n$  on both sides and made some slight change:

$$-\frac{n_m^+ \log_2 \frac{n_m^+}{n_m n_s}}{n} + \left(-\frac{n_m^- \log_2 \frac{n_m^-}{n_m n_s}}{n}\right) + \left(-\frac{n_m^+ \log_2 \frac{n_m^+}{n_m}}{n}\right) =$$

$$-\frac{n_m^+ \log_2 \frac{n_m^+}{n_m}}{n} + \left(-\frac{n_m^- \log_2 \frac{n_m^-}{n_m}}{n}\right) + \left(-\frac{n_m^+ \log_2 \frac{n_m^+}{n_m}}{n}\right)$$

$$\therefore -\frac{n_m^+ \log_2 \frac{n_m^+}{n_m}}{n} - \frac{n_m^- \log_2 \frac{n_m^-}{n_m}}{n} - \frac{n_m^+ \log_2 \frac{n_m^+}{n_m}}{n} =$$

$$-\frac{n_m^+ \log_2 \frac{n_m^+}{n_m}}{n} \geq -\frac{n_m^+ \log_2 \frac{n_m^+}{n_m^+ n_m^-}}{n} - \frac{n_m^- \log_2 \frac{n_m^-}{n_m^+ n_m^-}}{n}$$

$$-\frac{n_m^+ \log_2 \frac{n_m^+}{n_m^+ n_m^-}}{n} + \frac{n_m^- \log_2 \frac{n_m^-}{n_m^+ n_m^-}}{n} = \frac{n_m^+ \log_2 \frac{n_m^+}{n_m^+ n_m^-}}{n} + \frac{n_m^- \log_2 \frac{n_m^-}{n_m^+ n_m^-}}{n}$$

$$\therefore \frac{n_m^+ \log_2 \frac{n_m^+}{n_m^+ n_m^-}}{n} \geq -\frac{n_m^+ \log_2 \frac{n_m^+}{n_m^+ n_m^-}}{n} + \frac{n_m^- \log_2 \frac{n_m^-}{n_m^+ n_m^-}}{n}$$

Therefore:

$$\therefore H(S) = \sum_{j=1}^k \frac{|S_j|}{|S|} \times H(S_j) \geq H(S) - \frac{1}{|S|} \sum_{j=1}^k |S_j| \cdot H(S_{Uj})$$

$$\therefore IG(LSSF) \geq IG(LSS)$$

$$\text{threshold} = 5.5$$

$$H(\text{L}) = -\frac{2}{5}\log_2 \frac{2}{5} - \frac{3}{5}\log_2 \frac{3}{5} = 0.971$$

$$H(\text{R}) = -\frac{1}{2}\log_2 \frac{1}{2} - \frac{1}{2}\log_2 \frac{1}{2} = 1$$

$$H(\text{LR}) = \frac{5}{8} \cdot 0.971 + \frac{3}{8} = 0.983$$

$$IG(\text{threshold} = 5.5) = 0.057$$

threshold = 6.5

$$H(\text{L}) = 1$$

$$H(\text{R}) = -\frac{1}{3}\log_2 \frac{1}{3} - \frac{2}{3}\log_2 \frac{2}{3} = 0.917$$

$$IG(\text{LBS}) = 0.99 - 1 \times \frac{6}{9} - 0.917 \times \frac{3}{9} = 0.057$$

threshold = 7.5

$$H(\text{L}) = 1$$

$$H(\text{R}) = 0$$

$$IG = 0.99 - \frac{8}{9} = 0.102$$

Therefore, when threshold = 7.5, the IG is highest

3.3.

Feature No. 1 will be my root node

since it has highest IG

3.4.

