

Bài tập về cây quyết định

Theo bảng dữ liệu ta có:

$$- |D| = 8$$

- C = 1 Bị râm, không

$$- |C_{\text{Bị râm}}| = 3, |C_{\text{không}}| = 5$$

Tính Entropy

$$- \text{Info}(D) = - \frac{3}{8} \times \left(\log_2 \frac{3}{8} \right) - \frac{5}{8} \times \left(\log_2 \frac{5}{8} \right) = 0,954 \text{ bits}$$

Xét các thuộc tính màu tóc

Đen	Bị râm	2
	không	2
Râm	Bị râm	0
	không	3
Bạc	Bị râm	1
	không	0

$$\begin{aligned} \text{Info màu tóc}(D) &= \frac{4}{8} \times \left(- \frac{2}{4} \log_2 \frac{2}{4} - \frac{2}{4} \times \log_2 \frac{2}{4} \right) \\ &\quad + \frac{3}{8} \times \left(- \frac{3}{8} \log_2 \frac{3}{8} \right) + \frac{1}{8} \times \left(- 1 \log_2 1 \right) \\ &= 0,5 \text{ bits} \end{aligned}$$

$$\rightarrow \text{Gain (màu tóc)} = 0,954 \text{ bits} - 0,5 \text{ bits} = 0,454 \text{ bits}$$

Xét thuộc tính chiều cao

Cao	Bị rấm	0
	Không	2
Thấp	Bị rấm	2
	Không	1
Tầm thường	Bị rấm	1
	Không	2

$$\begin{aligned} \text{Inf chiều cao (D)} &= \frac{2}{8} \times \left(-\frac{2}{2} \times \log_2 \frac{2}{2} \right) + \frac{3}{8} \left(-\frac{2}{3} \log_2 \frac{2}{3} - \frac{1}{3} \log_2 \frac{1}{3} \right) \\ &\quad + \frac{3}{8} \times \left(-\frac{1}{3} \times \log_2 \frac{1}{3} - \frac{2}{3} \times \log_2 \frac{2}{3} \right) \\ &= 0,689 \text{ bits} \end{aligned}$$

$$\rightarrow \text{Gion (chiều cao)} = 0,954 \text{ bits} - 0,689 \text{ bit} = 0,265 \text{ bits}$$

Xét thuộc tính cân nặng

Nhẹ	Bị râm	1
	Không	1
Nặng	Bị râm	1
	Không	2
Vừa phải	Bị râm	1
	Không	2

$$\begin{aligned}
 \text{Info cân nặng (D)} &= \frac{2}{8} \times \left(-\frac{1}{2} \times \log_2 \frac{1}{2} - \frac{1}{2} \times \log_2 \frac{1}{2} \right) \\
 &\quad + \frac{3}{8} \times \left(-\frac{1}{3} \log_2 \frac{1}{3} - \frac{2}{3} \log_2 \frac{2}{3} \right) \\
 &\quad + \frac{3}{8} \times \left(-\frac{1}{3} \log_2 \frac{1}{3} - \frac{2}{3} \log_2 \frac{2}{3} \right) \\
 &= 0,939 \text{ bits}
 \end{aligned}$$

$$\rightarrow \text{Giảm (cân nặng)} = 0,954 - 0,939 = 0,015 \text{ bits}$$

Xét thuộc tính dùng thuốc

Không	Bị râm	3
	Không	2
Có	Bị râm	0
	Không	3

$$\begin{aligned} \text{Info dùng thuốc (D)} &= \frac{5}{8} \times \left(-\frac{3}{5} \times \log_2 \frac{3}{5} - \frac{2}{5} \times \log_2 \frac{2}{5} \right) \\ &\quad + \frac{3}{8} \times \left(-\frac{3}{3} \times \log_2 \frac{3}{3} \right) \\ &= 0.607 \text{ bits} \end{aligned}$$

$$\rightarrow \text{Gain(dùng thuốc)} = 0.954 - 0.607 = 0.347 \text{ bits}$$

Vì thuộc tính màu tóc có độ thông tin lớn nhất (0.454 bits) nên màu tóc là thuộc tính được chọn để phân tách (rẽ nhánh)

* Ta có cây quyết định:

