

Day	Outlook	Temp	Humidity	Wind	Play Tennis
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Weak	Yes
D6	Rain	Cool	Normal	Strong	No
D7	Overcast	Cool	Normal	Strong	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D10	Rain	Mild	Normal	Weak	Yes
D11	Sunny	Mild	Normal	Strong	Yes
D12	Overcast	Mild	High	Strong	Yes
D13	Overcast	Hot	Normal	Weak	Yes
D14	Rain	Mild	High	Strong	No

- Công thức:

- Select the attribute with the highest information gain
- Let  $p_i$  be the probability that an arbitrary tuple in D belongs to class  $C_i$ , estimated by  $|C_i \cap D|/|D|$
- **Expected information** (entropy) needed to classify a tuple in D:

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

- **Information** needed (after using A to split D into v partitions) to classify D:

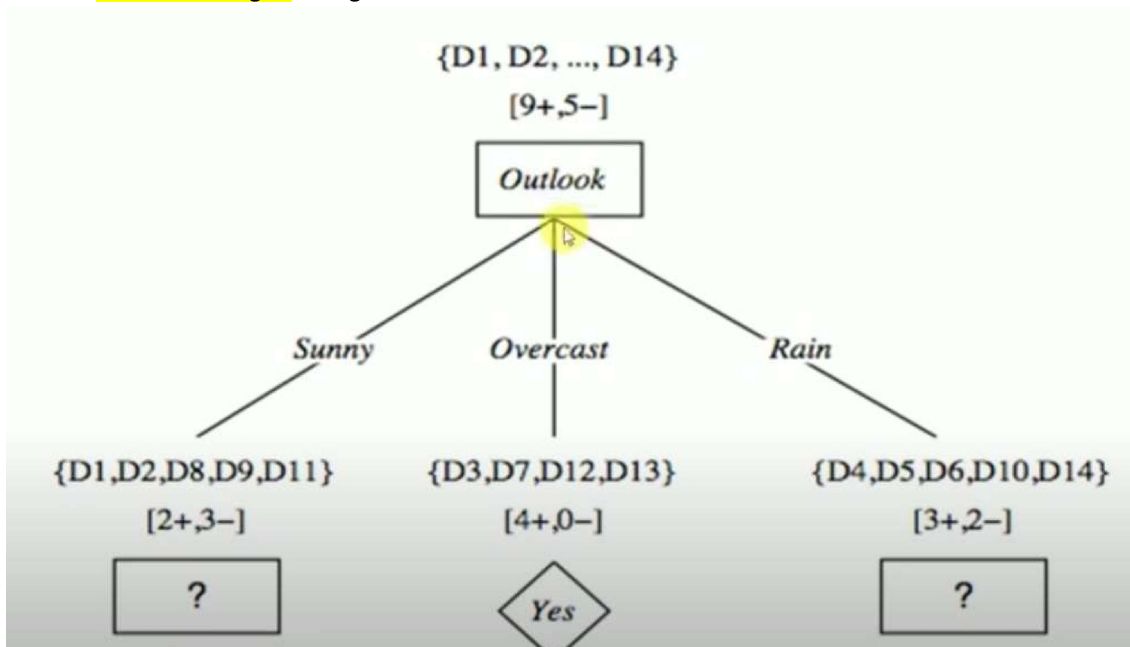
$$Info_A(D) = \sum_{j=1}^v \frac{|D_j|}{|D|} \times I(D_j)$$

- **Information gained** by branching on attribute A

$$Gain(A) = Info(D) - Info_A(D)$$

B1: Cần tìm Gain của 4 yếu tố: Outlook, Temp, Humidity, Wind để so sánh và xác định gốc trong decision tree.

- Thuộc tính “play tennis” là thuộc tính chính để tính Gain.
- “Play tennis” có 9 YES, 5 NO => D = [ 9+;5-].
- Tính Info (D) và tìm Gain của từng thuộc tính như trong ảnh. **Thuộc tính có Gain cao nhất sẽ làm gốc** trong decision tree.



Outlook làm gốc vì có chỉ số Gain cao nhất

- Các bước tính các Gain

Handwritten calculations on graph paper:

$$D = [9+; 5-]$$

$$Info(D) = -\frac{9}{14} \cdot \log_2 \frac{9}{14} - \frac{5}{14} \cdot \log_2 \frac{5}{14} = 0,94.$$

t) Outlook.

$$D_{Sunny} \leftarrow [2+; 3-] \Rightarrow Info(D_{Sunny}) = -\frac{2}{5} \cdot \log_2 \frac{2}{5} - \frac{3}{5} \cdot \log_2 \frac{3}{5} = 0,971.$$

$$D_{Overcast} \leftarrow [4+; 0-] \Rightarrow Info(D_{Overcast}) = -\frac{4}{4} \cdot \log_2 \frac{4}{4} - \frac{0}{4} \cdot \log_2 \frac{0}{4} = 0.$$

$$D_{Rain} \leftarrow [3+; 2-] \Rightarrow Info(D_{Rain}) = -\frac{3}{5} \cdot \log_2 \frac{3}{5} - \frac{2}{5} \cdot \log_2 \frac{2}{5} = 0,971.$$

$$\begin{aligned} \text{Gain}(D, \text{Outlook}) &= \text{Info}(D) - \frac{5}{14} \cdot \text{Info}(D_{\text{sunny}}) \\ &\quad - \frac{4}{14} \cdot \text{Info}(D_{\text{overcast}}) - \frac{5}{14} \cdot \text{Info}(D_{\text{rain}}) \\ &= 0,2464. \end{aligned}$$

+ ) Temp.

$$\begin{aligned} D_{\text{Hot}} &\leftarrow [9+, 5-] \Rightarrow \text{Info}(D_{\text{Hot}}) = \frac{-2}{9} \cdot \log_2 \frac{2}{9} - \frac{2}{9} \cdot \log_2 \frac{2}{9} \\ &= 1,0. \end{aligned}$$

$$\begin{aligned} D_{\text{Mild}} &\leftarrow [4+, 2-] \Rightarrow \text{Info}(D_{\text{Mild}}) = \frac{-4}{6} \cdot \log_2 \frac{4}{6} - \frac{2}{6} \cdot \log_2 \frac{2}{6} \\ &= 0,9183. \end{aligned}$$

$$\begin{aligned} D_{\text{Cool}} &\leftarrow [3+, 1-] \Rightarrow \text{Info}(D_{\text{Cool}}) = \frac{-3}{4} \cdot \log_2 \frac{3}{4} - \frac{1}{4} \cdot \log_2 \frac{1}{4} \\ &= 0,8113. \end{aligned}$$

$$\begin{aligned} \text{Gain}(D, \text{Temp}) &= \text{Info}(D) - \frac{4}{14} \cdot \text{Info}(D_{\text{Hot}}) - \frac{6}{14} \cdot \text{Info}(D_{\text{Mild}}) \\ &\quad - \frac{4}{14} \cdot \text{Info}(D_{\text{Cool}}) = 0,0289. \end{aligned}$$

HÔNG HÀ



+) Humidity

$$D_{\text{High}} [3+; 4-] \Rightarrow \text{Info}_{\text{High}} = -\frac{3}{7} \log_2 \frac{3}{7} - \frac{4}{7} \log_2 \frac{4}{7} = 0,9852$$

$$D_{\text{Normal}} [6+; 1-] \Rightarrow \text{Info}_{\text{Normal}} = -\frac{6}{7} \log_2 \frac{6}{7} - \frac{1}{7} \log_2 \frac{1}{7} = 0,5916$$

$$\text{Gain}(D, \text{Humidity}) = \text{Info}(D) - \left( \frac{7}{14} \times \text{Info}(D_{\text{High}}) \right) - \left( \frac{7}{14} \times \text{Info}(D_{\text{Normal}}) \right)$$

$$= 0,94 - \frac{7}{14} \times 0,9852 - \frac{7}{14} \times 0,5916 = 0,1516$$

+) Wind

$$D_{\text{Strong}} [3; 3-] \Rightarrow \text{Info}(D_{\text{Strong}}) = 1$$

$$D_{\text{Weak}} [6+; 2-] \Rightarrow \text{Info}(D_{\text{Weak}}) = -\frac{6}{8} \log_2 \frac{6}{8} - \frac{2}{8} \log_2 \frac{2}{8} = 0,8113$$

$$\text{Gain}(D, \text{Wind}) = \text{Info}(D) - \left( \frac{6}{14} \times \text{Info}(D_{\text{Strong}}) \right) - \left( \frac{8}{14} \times \text{Info}(D_{\text{Weak}}) \right)$$

$$= 0,94 - \frac{6}{14} \times 1 - \frac{8}{14} \times 0,8113 = 0,0478$$

B2: Tiếp tục tìm gốc của nhánh Sunny từ Outlook. Lọc những ngày Outlook = "Sunny", ta được bảng như bên dưới:

Day	Temp	Humidity	Wind	Play Tennis
D1	Hot	High	Weak	No
D2	Hot	High	Strong	No
D8	Mild	High	Weak	No
D9	Cool	Normal	Weak	Yes
D11	Mild	Normal	Strong	Yes

Các bước làm tương tự như B1 để tìm gốc tiếp theo

+ Temp:

$$D_{\text{sunny}} = [2+; 3-]$$

$$\text{Injo}(D_{\text{sunny}}) = -\frac{2}{5} \cdot \log_2 \frac{2}{5} - \frac{3}{5} \cdot \log_2 \frac{3}{5} = 0,97.$$

$$D_{\text{Hot}} \leftarrow [0+; 2-] \Rightarrow \text{Injo}(D_{\text{Hot}}) = 0,0.$$

$$D_{\text{Mild}} \leftarrow [1+; 1-] \Rightarrow \text{Injo}(D_{\text{Mild}}) = 1,0.$$

$$D_{\text{Cool}} \leftarrow [1+; 0-] \Rightarrow \text{Injo}(D_{\text{Cool}}) = 0,0.$$

$$\text{Gain}(D_{\text{sunny}}, \text{Temp}) = \text{Injo}(D) - \frac{2}{5} \cdot \text{Injo}(D_{\text{Hot}})$$

$$- \frac{2}{5} \cdot \text{Injo}(D_{\text{Mild}}) - \frac{1}{5} \cdot \text{Injo}(D_{\text{Cool}})$$

$$= 0,57.$$

+ Humidity:

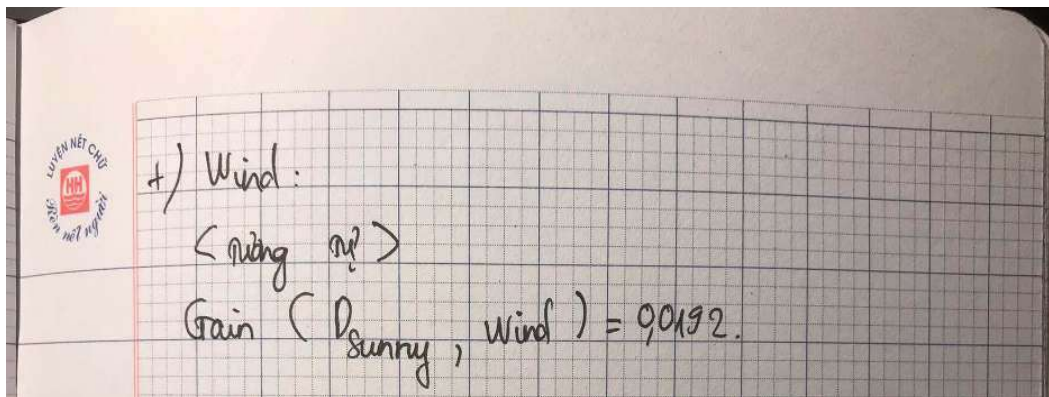
$$D_{\text{High}} \leftarrow [0+; 3-] \Rightarrow \text{Injo}(D_{\text{High}}) = 0,0.$$

$$D_{\text{Normal}} \leftarrow [2+; 0-] \Rightarrow \text{Injo}(D_{\text{Normal}}) = 0,0.$$

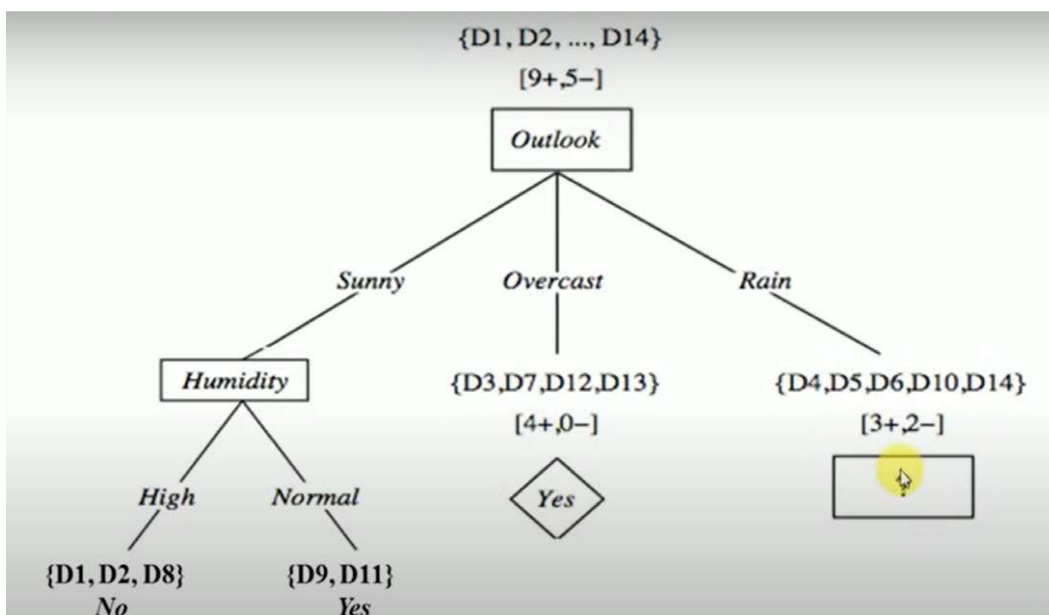
$$\text{Gain}(D_{\text{sunny}}, \text{Humidity}) = \text{Injo}(D) - \frac{3}{5} \cdot \text{Injo}(D_{\text{High}})$$

$$- \frac{2}{5} \cdot \text{Injo}(D_{\text{Normal}})$$

$$= 0,97.$$



- Gain của Humidity là cao nhất nên Humidity tiếp tục làm gốc của nhánh tiếp theo.
- Dựa vào trong phần tính của Humidity (ở bước 2 này) , D high và D normal đều bằng không => Kết thúc (như ảnh dưới)



B3: Tiếp tục tương tự B2, tìm gốc của nhánh Rain từ Outlook.

Day	Temp	Humidity	Wind	Play Tennis
D4	Mild	High	Weak	Yes
D5	Cool	Normal	Weak	Yes
D6	Cool	Normal	Strong	No
D10	Mild	Normal	Weak	Yes
D14	Mild	High	Strong	No

$$Gain(S_{Rain}, Temp) = 0.0192$$

$$Gain(S_{Rain}, Humidity) = 0.0192$$

$$Gain(S_{Rain}, Wind) = 0.97$$

Ta có kết quả cuối cùng.

