

①

Design Tree - play Tennis Gain Calculation.

DAY	outlook	temperature	humidity	wind	Decision
1	Sunny	hot	high	Weak	No
2	Sunny	hot	high	Strong	No
3	overcast	hot	high	Weak	Yes
4	rainfall	mild	high	Weak	Yes
5	rainfall	cool	normal	Weak	Yes
6	rainfall	cool	normal	Strong	No
7	overcast	cool	normal	Strong	Yes
8	Sunny	mild	high	Weak	No
9	Sunny	cool	normal	Weak	Yes
10	rainfall	mild	normal	Weak	Yes
11	Sunny	mild	normal	Strong	Yes
12	overcast	mild	high	Strong	Yes
13	overcast	hot	normal	Weak	Yes
14	rainfall	mild	high	Strong	No

$$\text{Gain}(S, f_i) = H(S) - \sum \frac{|S_v|}{|S|} H(S_v)$$

~~CART~~ → Purity — Entropy (from ID3)
 — Gini Impurity (from CART)

$$\text{ID3} \rightarrow \text{Entropy} = - \sum_{i=1}^n p_i \log_2(p_i) \quad \text{Gini Coeff} = 1 - \sum_{i=1}^n p_i^2$$

$$H(S) = \text{Entropy} \rightarrow 0 \leq 1 \Rightarrow 0 \rightarrow \text{Good Entropy (pure split)}$$

$$\text{Gini Coeff} \rightarrow 0 \leq 0.5$$

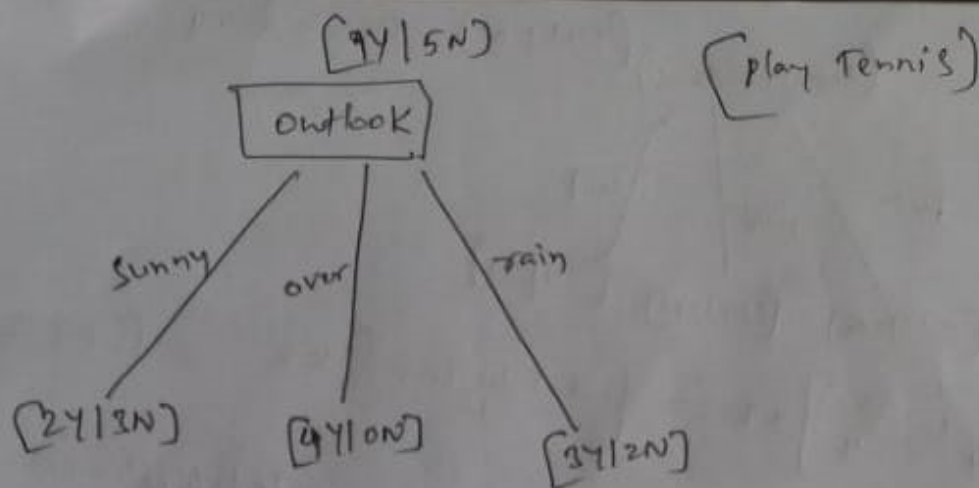
Conclusion:

∴
 Calculated Information Gain individually for all 4 features.

outlook has highest Gain ≈ 0.24 ✓

$$\left. \begin{array}{l} \text{Tempe Gain} = 0.06 \\ \text{Windy} = 0.04 \\ \text{Humidity} = 0.16 \end{array} \right\}$$

②



$$E_s = -\frac{2}{5} \log\left(\frac{2}{5}\right) - \frac{3}{5} \log\left(\frac{3}{5}\right) = 0.97$$

$$E_o = -\frac{4}{4} \log\left(\frac{4}{4}\right) - \frac{0}{4} \log\left(\frac{0}{4}\right) = 0$$

$$E_{rain} = -\frac{3}{5} \log\left(\frac{3}{5}\right) - \frac{2}{5} \log\left(\frac{2}{5}\right) = -0.6 \log(0.6) - 0.4 \log(0.4) = 0.97$$

$$H(S) \text{ (Gain (outlook))} = E_{outlook} = -\frac{9}{14} \log\left(\frac{9}{14}\right) - \frac{5}{14} \log\left(\frac{5}{14}\right) = -0.64 \log(0.64) - 0.257 \log(0.257)$$

$$H(S) = 0.94$$

$$\text{Gain} = H(S) - \sum \frac{|S_v|}{|S|} H(S_v)$$

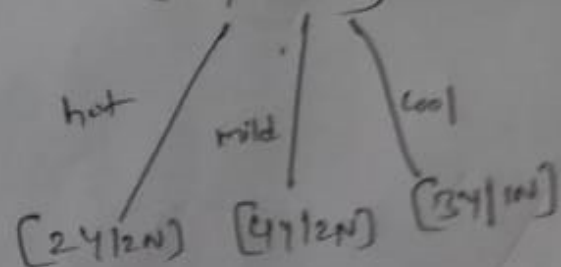
$$= 0.94 - \left[\frac{5}{14} (0.97) + 0 + \frac{5}{14} (0.97) \right]$$

$$= 0.94 - \left[2 \times \frac{5}{14} (0.97) \right] = 0.94 - \left[2 \times 0.35 \times 0.97 \right]$$

$$= 0.94 - 0.69$$

$$\text{Gain}_{outlook} = \underline{0.24}$$

(Temperature) (9415N) $H(S) = 0.94$



$$E_{\text{hot}} = \left[-\frac{2}{4} \log\left(\frac{2}{4}\right) \right] \times 2 = \left[-0.5 \log(0.5) \right] \times 2 = [-0.5(-1)] \times 2 = 1$$

$$E_{\text{mild}} = \frac{4}{6} = 0.66 \Rightarrow \frac{E}{6} = 0.33 = -0.66 \log(0.66) - 0.33 \log(0.33)$$

$$= (-0.66 \times -0.59) - (0.33 \times -1.5)$$

$$= 0.38 + 0.49 = 0.87$$

$$E_{\text{cool}} = \frac{3}{4} = 0.75, \frac{1}{4} = 0.25$$

$$= -0.75 \log(0.75) - 0.25 \log(0.25)$$

$$= -0.41 - (-2)$$

$$= +0.307 + 0.5 = 0.807$$

$$\text{Gain} = 0.94 - \left[\frac{4}{14} (1) + \frac{6}{14} (0.87) + \frac{4}{14} (0.807) \right]$$

$$= 0.94 - [0.28 + 0.372 + 0.22]$$

$$= 0.94 - [0.87] = 0.06$$

$$\text{Gain}_{\text{temp}} = 0.06 \quad (0.06)$$

(3)

(94/5N) humidity $H(3) = 0.94$.

high

normal.

(34/4N)

(64/11N)

$$E_{\text{high}} \Rightarrow \frac{3}{7} = 0.42 \quad \frac{4}{7} = 0.571$$

$$\begin{aligned}
 &= -0.42 \log(0.42) - 0.571 \log(0.571) \\
 &= -0.42(-1.25) - 0.571(-0.801) \\
 &= 0.525 + 0.461 = 0.98
 \end{aligned}$$

$$\begin{aligned}
 E_{\text{normal}} &= \frac{-6 \log(6) - 0.16 \log(0.16)}{7} \\
 &= \frac{-6(2.58) - 0.16(-2.6)}{7} \\
 &= \frac{-15.44}{7}
 \end{aligned}$$

$$\begin{aligned}
 E_{\text{normal}} &= -\frac{6}{7} \log\left(\frac{6}{7}\right) - \frac{1}{7} \log\left(\frac{1}{7}\right) \\
 &= -0.85 \log(0.85) - 0.14 \log(0.14) \\
 &= -0.85(-0.23) - 0.14(-2.83) \\
 &= 0.19 + 0.39 = 0.58
 \end{aligned}$$

$$\begin{aligned}
 \text{Gain (humidity)} &= 0.94 - \left[\underbrace{0.98 \times \frac{7}{14} + 0.58 \times \frac{7}{14}}_{0.49 + 0.29} \right] \\
 &= 0.94 - 0.78 \\
 &= 0.16
 \end{aligned}$$

(94/5N) Windy $H(S) = 0.94$.

Weak

Strong

(64/2N)

(34/84)

$H(S) = 1$ very impure.
strong

$$E_{\text{weak}} = \frac{6}{8} = 0.75 \quad \frac{2}{8} = 0.25$$

$$= -0.75 \log(0.75) - 0.25 \log(0.25)$$

$$= -0.75(-0.415) - 0.25(-2)$$

$$= 0.311 + 0.5 = 0.811$$

$$\text{Gain} = 0.94 - \left[0.811 \times \frac{8}{14} + 0.11 \times \frac{6}{14} \right]$$

$$= 0.94 - [0.463 + 0.42]$$

$$\text{Windy Gain} = 0.94 - 0.891 = 0.04$$

④

Entropy

① Algorithm used for calculation
ID3

② Value max purity = 0 \Rightarrow [0, 1]
max impurity 1

③ Formula: $-\sum_{i=1}^n p_i \log_2(p_i)$

④ measure of Information that
Indicates disorder of the features
with target.

⑤ Use logarithms more complex.

Gini Impurity

CART Algorithm,

[0, 0.5] \rightarrow 0.5 is
max purity.

$$\Rightarrow 1 - \sum p_i^2$$

measures the frequency at which
any element of the data set will be
misclassified when it is randomly labeled

calculation is faster