→ Outlook এর gain বেশি, তাই তা হবে decision tree এর root.

☆ Info gain dataset কে more appropriately ভাগ করে,

☆ If all the instances belongs to a Particular class then we would have a leaf node

## C 4.5 Algorithm.

☆ C4.5 algorithm টা ID3 এর improvement.

☆ C4.5 -এর তথ্য প্রয়োজন হয় ⟶ Split Info_A (D)
⟶ Gain Ratio (A)

☆ C4.5

☆ ID3 $\xrightarrow{\text{Successor}}$ C4.5

## C4.5 এর সূত্রাবলী:

$$→ Split\ Info_A(D) = - n \sum_{j=1}^{n} \frac{|D_j|}{|D|} \times \log_2\left(\frac{|D_j|}{|D|}\right)$$

$$→ Gain\ Ratio\ (A) = \frac{Gain\ (A)}{Split\ Info_A\ (D)}$$

☆ Gain Ratio তে Gain (A) দরকার হয়, Gain (A) -এর সূত্র

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

☆ অর্থাৎ ID3 এর সূত্রাবলীও কাজে লাগবে,

Table: The playing tennis dataset

| Day | Outlook | Temperature | Humidity | Wind | Play |
|---|---|---|---|---|---|
| $D_1$ | Sunny | Hot | High | Weak | No |
| $D_2$ | Sunny | Hot | High | Strong | No |
| $D_3$ | Overcast | Hot | High | Weak | Yes |
| $D_4$ | Rain | Mild | High | Weak | Yes |
| $D_5$ | Rain | Cool | Normal | Weak | Yes |
| $D_6$ | Rain | Cool | Normal | Strong | No |
| $D_7$ | Overcast | Cool | Normal | Strong | Yes |
| $D_8$ | Sunny | Mild | High | Weak | No |
| $D_9$ | Sunny | Cool | Normal | Weak | Yes |
| $D_{10}$ | Rain | Mild | Normal | Weak | Yes |
| $D_{11}$ | Sunny | Mild | Normal | Strong | Yes |
| $D_{12}$ | Overcast | Mild | High | Strong | Yes |
| $D_{13}$ | Overcast | Hot | Normal | Weak | Yes |
| $D_{14}$ | Rain | Mild | High | Strong | No |

**Find the root using C4.5 algorithm from the given dataset**

$$\text{Info}(D) = -\frac{9}{14} \log_2 (9/14) - \frac{5}{14} \log_2 (5/14)$$

$$= -0.64 \times (-0.64) - 0.35 \times (-1.51)$$

$$= 0.40 + 0.52$$

$$= 0.92$$

Attribute → Outlook

$$\text{Info}_{outlook}(D) = \frac{5/14 \times (-2/5 \log_2 2/5 - 3/5 \log_2 3/5)}{\text{Sunny}} +$$

$$\frac{4/14 \times (-4/4 \log_2 4/4 - 0/4 \log_2 0/4)}{\text{Overcast}} +$$

$$\frac{5/14 \times (-3/5 \log 3/5 - 2/5 \log 2/5)}{\text{Rain}}$$

$$\left(\frac{|c_i|}{|D|}\right) = 0.35 \times \{-0.4 \times (-1.32) - 0.6 \times (-0.73)\} +$$

$$0.28 \times \{-1 \times 0 - \log_2 (0/4)\} +$$

$$0.35 \times \{-0.6 \times (-0.73) - 0.4 \times (-1.32)\}$$

$$= 0.6762$$

$$\text{Gain}(Outlook) = \text{Info}(D) - \text{Info}_{outlook}(D)$$

$$= 0.92 - 0.6762$$

$$= 0.2438$$

Split Info (Outlook) $=$ $\dfrac{-5/14 \ \log_2 5/14}{\text{Sunny}}$ —

$\dfrac{-4/14 \ \log_2 4/14}{\text{Overcast}}$ — $\dfrac{5/14 \ \log_2 5/14}{\text{Rain}}$

$= -0.35 \ \log_2 (0.35) - 0.28 \ \log_2 (0.28) - 0.35 \times \log_2 (0.35)$

$= -0.35 \times (-1.51) - 0.28 \times (-1.83) - 0.35 \times (-1.51)$

$= 0.52 + 0.51 + 0.52$

$= 1.55$

Gain Ratio (Outlook) $= \dfrac{0.2438}{1.55}$

$= 0.157.$

## Attribute ⇒ Temperature

Info $_{Temperature}$ (D) $= \dfrac{4/14 \times (-2/4 \ \log 2/4 - 2/4 \ \log_2 2/4)}{\text{Hot}}$ +

$\dfrac{6/14 \times (-4/6 \ \log_2 4/6 - 2/6 \ \log_2 2/6)}{\text{Mild}}$ +

$\dfrac{4/14 \times (-3/4 \ \log_2 3/4 - 1/4 \ \log_2 1/4)}{\text{Cool}}$

$= 0.28 \times \{ -0.5 \times (-1) - 0.5 \times (-1) \} + 0.42 \times \{ -0.66 \times (-0.599)$

$\quad - 0.33 \times (-1.59) + 0.28 \times \{ -0.75 \times (-0.41) -$

$\quad 0.25 \times (-2) \}$

$= 0.892$

Gain( Temparature) $=$ Info (D) $-$ Info $_{Temparature}$ (D)

$\qquad = 0.92 - 0.892$

$\qquad = 0.028$

Split Info ( Temparature) $= - \dfrac{4/14 \ \log_2 \ 4/14}{Hot} = -$

$- \dfrac{6/14 \ \log_2 \ 6/14}{Mild} - \dfrac{4/14 \ \log_2 \ 4/14}{Cool}$

$= - 0.28 \ \log_2 (0.28) - 0.42 \ \log_2 (0.42) - 0.28 \ \log_2 (0.28)$

$= - 0.28 \times (-1.83) - 0.42 \times (-1.25) - 0.28 \times (-1.83)$

$= 0.51 + 0.525 + 0.51$

$= 1.545$

Gain Ratio ( Temparature) $= \dfrac{0.028}{1.545}$

$\qquad = 0.018$

Attribute ⟹ Humidity

$$\text{Info}_{\text{Humidity}}(D) = \underbrace{\frac{7/14 \times \{-3/7 \log_2 3/7 - 4/7 \log_2 4/7)}{\text{High}}}_{} +$$

$$\underbrace{\frac{7/14 \times (-6/7 \log_2 6/7 - 1/7 \log_2 1/7)}{\text{Normal}}}_{}$$

$$= 0.5 \quad \times \{-0.42 \times (-1.25) - 0.57 \times (-0.81)\} +$$

$$0.5 \quad \times \{-0.85 \times (-0.23) - 0.14 \times (-2.83)\}$$

$$= 0.78$$

$$\text{Gain (Humidity)} = \text{Info}(D) - \text{Info}_{\text{Humidity}}(D)$$

$$= 0.92 - 0.785$$

$$= 0.135$$

$$\text{Split Info (Humidity)} = \underbrace{\frac{-7/14 \log_2 7/14}{\text{High}}}_{} -$$

$$\underbrace{\frac{-7/14 \log_2 7/14}{\text{Normal}}}_{}$$

$$= -0.5 \log_2 (0.5) - 0.5 \log_2 (0.5)$$

$$= -0.5 \times (-1) - 0.5 \times (-1)$$

$$= 0.5 + 0.5$$

$$= 1$$

$$\text{Gain Ratio (Humidity)} = \frac{0.135}{1}$$

$$= 0.135$$

Attribute → Wind

$$\text{Info}_{wind} (D) = \frac{8/14 \ (- 6/8 \ \log_2 6/8 - 2/8 \ \log_2 2/8)}{\text{Weak}} +$$

$$\frac{6/14 \ (- 3/6 \ \log_2 3/6 - 3/6 \ \log_2 3/6)}{\text{Strong}}$$

$$= 0.57 \times \{ - 0.75 \times (-0.41) - 0.25 \times (-2)\} +$$
$$0.42 \times \{ - 0.5 \times (-1) - 0.5 \times (-1)\}$$

$$= 0.88$$

$$\text{Gain (Wind)} = \text{Info}_D - \text{Info}_{wind} (D)$$
$$= 0.92 - 0.88$$
$$= 0.04$$

$$\text{Split Info (Wind)} = \frac{- 8/14 \ \log_2 8/14}{\text{Weak}} -$$

$$\frac{6/14 \ \log_2 6/14}{\text{Strong}}$$

$$z = - 0.57 \times \log_2 (0.57) - 0.42 \times \log_2 (0.42)$$

$$= - 0.57 \times (-0.81) - 0.42 \times (-1.25)$$

$$= 0.46 + 0.525$$

$$= 0.985$$

Gain Ratio (Wind) $= \dfrac{0.09}{0.985}$

$= 0.09$

Here

Gain Ratio (Outlook) $= 0.2438$

Gain Ratio (Temperature) $= 0.0181$

Gain Ratio (Humidity) $= 0.135$

Gain Ratio (Wind) $= 0.09$

Gain Ratio of Outlook is higher, so it would be the splitting feature.

## ঐ Wrote the usage of Gini Index

__Ans.__

The Gini index is used in classification & Regression Trees (CART) algorithm.

## Gini (CART) Formula

a. $Gini(D) = 1 - \sum_{i=1}^{N} P_i^2$

b. $Gini_A(D) = \frac{|D_1|}{|D|} Gini(D_1) + \frac{|D_2|}{|D|} Gini(D_2)$  [এই slide এ আছে]

c. $\Delta Gini(A) = Gini(D) + Gini_A(D)$

d. $Gini_A(D) = \sum_{J=1}^{n} \left( \frac{|D_j|}{|D|} \times Gini\ |D_j| \right.$  [এই নিয়ম— $Gini_A(D)$ easily বের হয়]

Attribute    Total Instance.

Table: The playing tennis dataset

| Day | Outlook | Temperature | Humidity | Wind | Play |
|-----|---------|-------------|----------|------|------|
| $D_1$ | Sunny | Hot | High | Weak | No |
| $D_2$ | Sunny | Hot | High | Strong | No |
| $D_3$ | Overcast | Hot | High | Weak | Yes |
| $D_4$ | Rain | Mild | High | Weak | Yes |
| $D_5$ | Rain | Cool | Normal | Weak | Yes |
| $D_6$ | Rain | Cool | Normal | Strong | No |
| $D_7$ | Overcast | Cool | Normal | Strong | Yes |
| $D_8$ | Sunny | Mild | High | Weak | No |
| $D_9$ | Sunny | Cool | Normal | Weak | Yes |
| $D_{10}$ | Rain | Mild | Normal | Weak | Yes |
| $D_{11}$ | Sunny | Mild | Normal | Strong | Yes |
| $D_{12}$ | Overcast | Mild | High | Strong | Yes |
| $D_{13}$ | Overcast | Hot | Normal | Weak | Yes |
| $D_{14}$ | Rain | Mild | High | Strong | No |

**Find the root using CART(Gini Index) algorithm from the given dataset**

Solve:

$$Gini \cdot (D) = 1 - \sum_{i=1}^{n} P_i^2$$

$$= 1 - \{ (9/14)^2 + (5/14)^2 \}$$

$$= 1 - \{ (0.64)^2 + (0.35)^2 \}$$

$$= 1 - (0.40 + 0.12)$$

$$= 1 - 0.52$$

$$= 0.48$$

Attribute $\Rightarrow$ Outlook

$$Gini_{Sunny} (D) =$$

$$Gini_{Outlook} (D) = \frac{5/14 \{ 1 - ((2/5)^2 + (3/5)^2) \}}{Sunny} +$$

$$\frac{4/14 \{ 1 - ((4/9)^2 + (0/4)^2) \}}{Overcast} +$$

$$\frac{5/14 \{ 1 - ((3/5)^2 + (2/5)^2) \}}{Rainy}$$

$$= 0.35 \times (1 - 0.16 - 0.36) + 0.28 \times (1 - 1) +$$
$$0.35 \times (1 - 0.36 - 0.16)$$

$$= 0.35 \times 0.48 + 0.35 \times 0.48$$

$$= 0.998$$

$$\Delta \text{ Gini (Outlook)} = 0.48 - 0.998$$
$$= -0.518$$

Attribute ⇒ Temperature

$$\text{Gini}_{\text{Temperature}} \; (D) = \dfrac{4/14 \{ 1 - ((2/4)^2 + (2/4)^2) \}}{\text{Hot}} +$$

$$\dfrac{6/14 \{ 1 - ((4/6)^2 + (2/6)^2) \}}{\text{Mild}} +$$

$$\dfrac{4/14 \{ 1 - ((3/4)^2 + (1/4)^2) \}}{\text{Cool}}$$

$$= 0.28 \times 0.5 + 0.42 \times 0.44 + 0.28 \times 0.375$$

$$= 0.4298$$

$$\Delta \text{ Gini ( Temperature)} = 0.48 - 0.4298$$
$$= 0.050$$

Attribute → Humidity

$$\text{Gini}_{\text{Humidity}}^{(D)} = \frac{7/14 \{ 1 - ((3/7)^2 + (4/7)^2) \}}{\text{High}} +$$

$$\frac{7/14 \{ 1 - ((6/7)^2 + (1/7)^2) \}}{\text{Normal}}$$

$$= 0.5 \times 0.48 + 0.5 \times 0.244$$

$$= 0.362$$

$$\Delta \text{Gini (Humidity)} = 0.48 - 0.362$$

$$= 0.118$$

Attribute ⇒ Wind

$$\text{Gini}_{\text{Wind}}^{(D)} = \frac{8/14 \{ 1 - ((6/8)^2 + (2/8^2)) \}}{\text{Weak}} +$$

$$\frac{6/14 \{ 1 - ((3/6)^2 + (3/6)^2) \}}{\text{Strong}}$$

$$= 0.57 \times 0.375 + 0.42 \times 0.5$$

$$= 0.42$$

$$\Delta \text{Gini (Wind)} = 0.48 - 0.42$$

$$= 0.06$$

Here

$\Delta Gini (Outlook) = -0.518$

$\Delta Gini (Temperature) = 0.050$

$\Delta Gini (Humidity) = 0.118$

$\Delta Gini (Wind) = 0.06$

Here highest Delta gini comes from Humidity. So, it would be the root.

Ay —