

Date.....

## Information Gain

Q What is Information Gain?

Ans Information Gain is a metric used to train Decision Trees. Specifically, this metric measures the quality of the split. The Information Gain is based on the decrease in entropy after a data set is split on an attribute. Constructing a decision tree is all about finding out the attribute that returns the highest information gain.

Q What is the formula to calculate the Information Gain?

Ans 
$$\text{Information Gain} = E(\text{Parent}) - \text{Weighted Average} * E(\text{child})$$

Example:-

Outlook	Temperature	Humidity	Windy	Play Tennis
Sunny	Hot	High	False	No
Sunny	Hot	High	True	No
Overcast	Hot	High	False	Yes
Rainy	Mild	High	False	Yes
Rainy	Cool	Normal	False	Yes
Rainy	Cool	Normal	True	No
Overcast	Cool	Normal	True	Yes
Sunny	Mild	High	False	No
Sunny	Cool	Normal	False	Yes
Rainy	Mild	Normal	False	Yes
Sunny	Mild	Normal	True	Yes
Overcast	Mild	High	True	Yes
Overcast	Hot	Normal	False	Yes
Rainy	Mild	High	True	No

Play-Tennis Data Set

Date.....

Step-1:- Calculating the E (Parent)

$$\Rightarrow E(P) = -P_{Yes} \log_2 (P_{Yes}) - P_{No} \log_2 (P_{No})$$

$$\Rightarrow -\frac{9}{14} \log_2 (9/14) - \frac{5}{14} \log_2 (5/14)$$

$$E(P) = 0.94 \quad \boxed{-1}$$

Step-2:- Calculate Entropy for Children :-

Outlook														
Sunny					Overcast					Rainy				
Day	Temperature	Humidity	Windy	Play Tennis	Day	Temperature	Humidity	Windy	Play Tennis	Day	Temperature	Humidity	Windy	Play Tennis
D1	Hot	High	Weak	No	D4	Mild	High	Weak	Yes	D10	Mild	Normal	Weak	Yes
D2	Hot	High	Strong	No	D5	Cool	Normal	Weak	Yes	D6	Cool	Normal	Strong	No
D8	Mild	High	Weak	No	D9	Cool	Normal	Weak	Yes	D11	Mild	Normal	Weak	Yes
D9	Cool	Normal	Weak	Yes	D12	Mild	High	Strong	No	D13	Mild	High	Strong	No
D11	Mild	Normal	Strong	Yes	D14	Mild	High	Strong	No					

DataSet - Sunny

DataSet - Rainy

Day	Temperature	Humidity	Windy	Play Tennis
D3	Hot	High	Weak	Yes
D7	Cool	Normal	Strong	Yes
D12	Mild	High	Strong	Yes
D13	Hot	Normal	Weak	Yes

DataSet - Overcast

$$E(D_{-Sunny}) = -P_{Yes} \log_2 (P_{Yes}) - P_{No} \log_2 (P_{No}) \quad | \quad E(D_{-Overcast}) = -P_{Yes} \log_2 (P_{Yes}) -$$

$$\Rightarrow -\frac{2}{5} \log_2 (2/5) - \frac{3}{5} \log_2 (3/5)$$

$$\Rightarrow \boxed{0.97} \quad \boxed{-2}$$

$$= \frac{4}{4} \log_2 (4/4) - 0$$

$$\Rightarrow \boxed{0} \quad \boxed{-3}$$

$$E(D_{-Rainy}) = -P_{Yes} \log_2 (P_{Yes}) - P_{No} \log_2 (P_{No})$$

$$\Rightarrow -\frac{3}{5} \log_2 (3/5) - \frac{2}{5} \log_2 (2/5)$$

$$\Rightarrow \boxed{0.97} \quad \boxed{-4}$$

Date.....

Step-3:- Calculate the Weighted Entropy of Children...

$$\begin{aligned}\text{Weighted Entropy} &\Rightarrow \frac{5}{14} * 0.97 + \frac{4}{14} * 0 + \frac{5}{14} * 0.97 \\ &\Rightarrow 0.3464 + 0.3464\end{aligned}$$

$$\boxed{\text{W.E (Children)} \Rightarrow 0.69} - \textcircled{5}$$

Taking ①, ②, ③, ④ and ⑤

$$\begin{aligned}\text{Information Gain} &= E(\text{Parent}) - \text{Weighted Average} * E(\text{Children}) \\ &\Rightarrow 0.94 - 0.69\end{aligned}$$

$$\boxed{\text{Information Gain} \Rightarrow 0.25}$$