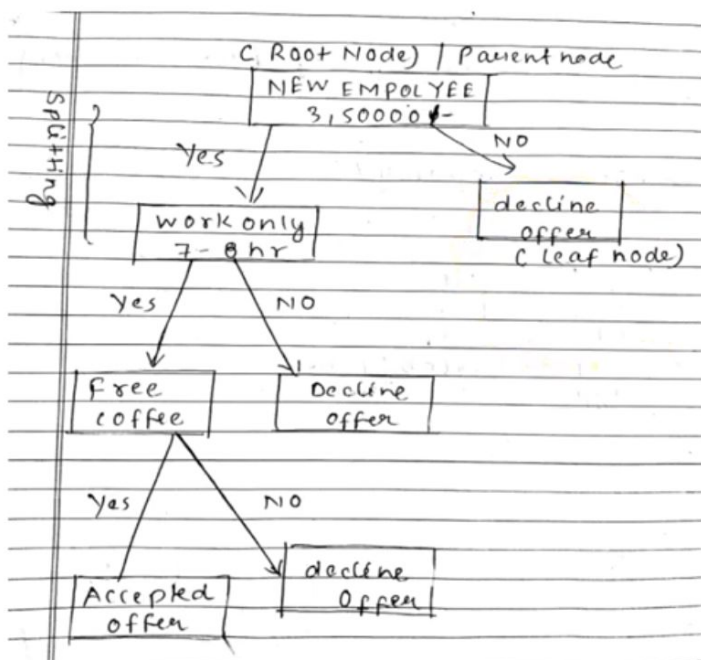


What is Decision tree

A decision tree is a graphical representation of possible solution to a decision based on certain condition it's called a decision tree.



How does a tree decide where to split

1. GINI INDEX

Measure the impurity used to build a decision tree.

2. Information Gain

Select the node will be highest information gain.

3. Reduction in variance

If your data is pure then less variance in data.

Lowest variance is good for creating a variance

entropy

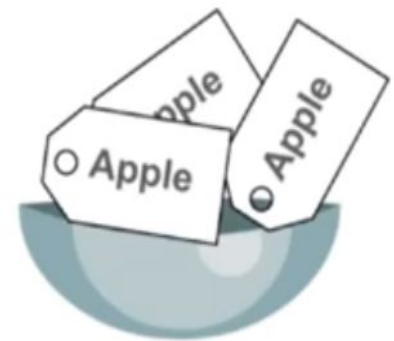
- Measure the impurity of the substances.

- What is impurity

- Impurity=0

- Impurity \neq 0

$$\text{Entropy}(s) = -P(\text{yes}) \log_2 P(\text{yes}) - P(\text{no}) \log_2 P(\text{no})$$



Entropy(features)= -P(yes)log₂P(Yes) -P(No)log₂P(No)
Entropy(outlook=sunny)=
P(yes)

Information gain

- ☒ Decide which attribute should be selected as the decision tree node.

Information Gain = Entropy(S) – [(Weighted Avg) x Entropy(each feature)]

Problematic 1000

Day	outlook	Temperature	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

Step ①) Out of 14 Instances we have 9 Yes
And 5 NO

compute the Entropy of entire data-set

$$E(S) = -P(\text{Yes}) \log_2 P(\text{Yes}) - P(\text{No}) \log_2 P(\text{No})$$

$$= -\left(\frac{9}{14}\right) \times \log_2 \left(\frac{9}{14}\right) - \left(\frac{5}{14}\right) \log_2 \frac{5}{14}$$

$$E(S) = 0.41 + 0.53$$

$$= 0.94$$

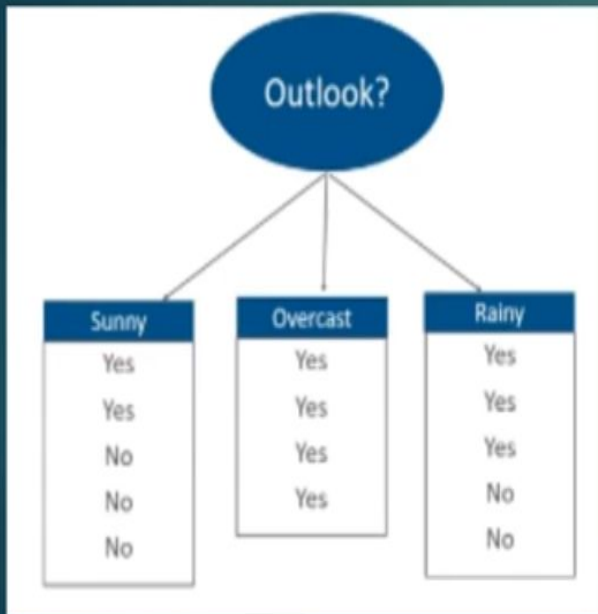
Which node select as root node

Outlook

temperature

humidity

windy



$$E(\text{Outlook} = \text{Sunny}) \\ = -\frac{2}{5} \log_2 \frac{2}{5} - \frac{3}{5} \log_2 \left(\frac{3}{5} \right) \\ = 0.971$$

$$E(\text{Outlook} = \text{Overcast}) \\ = 0$$

$$E(\text{Outlook} = \text{Rainy}) = 0.971$$

Information from outlook

$$I(\text{Outlook}) = \frac{5}{14} \times 0.971 + \frac{4}{14} \times 0 + \frac{5}{14} \times 0.971 \\ = 0.693$$

Information gained from outlook

$$\text{Gain}(\text{Outlook}) = E(S) - I(\text{Outlook}) \\ = 0.94 - 0.693$$

Entropy(features) = $-P(\text{yes})\log_2 P(\text{Yes}) - P(\text{No})\log_2 P(\text{No})$

Entropy(outlook=sunny) = $-\frac{2}{5}\log_2(\frac{2}{5}) - \frac{3}{5}\log_2(\frac{3}{5}) = 0.971$

Entropy(outlook=overcast) = $-\frac{4}{4}\log_2(1) - 0 = 0$

E(outlook=Rainy) = $-\frac{3}{5}\log_2(\frac{3}{5}) - \frac{2}{5}\log_2(\frac{2}{5}) = 0.971$

Information from outlook = sum of (weight * E(each))

Information from outlook = $\frac{5}{14} \times 0.971 + \frac{4}{14} \times 0 + \frac{5}{14} \times 0.971 = 0.693$

Information gained from outlook = $E(S) - \text{Information}(\text{Outlook})$

Entropy(target)

Entropy(Play/tennis) = $-P(\text{yes})\log_2(P(\text{yes})) - P(\text{NO})\log_2(\text{No})$

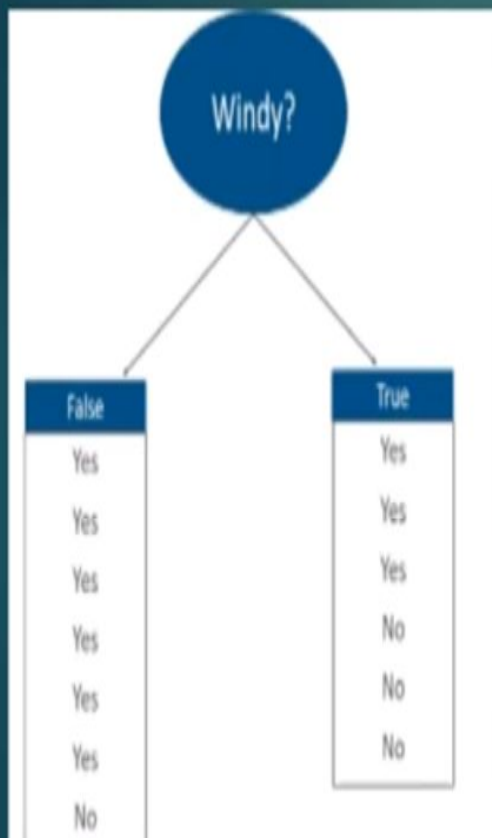
$$P(\text{yes})=9/14 \quad P(\text{no})=5/14$$

$$\text{Entropy}(\text{Target})=-9/14 \cdot \log_2(9/14) - (5/14) \cdot \log_2(5/14)=0.94$$

$$\text{Information gained from outlook} = E(S) - \text{Information}(\text{Outlook})$$

$$E(S)=\text{entropy of target}$$

$$\text{IG}=0.94-0.693=0.247(\text{outlook})$$



$$E(\text{Windy} = \text{TRUE}) = 1$$

$$E(\text{Windy} = \text{FALSE}) = 0.811$$

Information from windy

$$I(\text{Windy}) = \frac{8}{14} \times 0.811 + \frac{6}{14} \times 1 = 0.892$$

Information gained from windy

$$\begin{aligned} \text{gain}(\text{Windy}) &= E(S) - I(\text{Windy}) \\ &= 0.94 - 0.892 \\ &= 0.048 \end{aligned}$$

Find out Rest of two with same Procedure:

Outlook	Temperature
max	Information : 0.811
Information : 0.693	Gain : 0.029
Gain : 0.247	

humidity	windy
Information: 0.788	Information : 0.892
Gain: 0.152	Gain : 0.048

