

## Decision Tree

05 June 2023 20:38

- \* It comes under Supervised Learning.
- \* It is used to solve both Classification and Regression problems.

### Type of Decision Tree

#### ① Decision Tree Regression

IT is used when target is continuous.

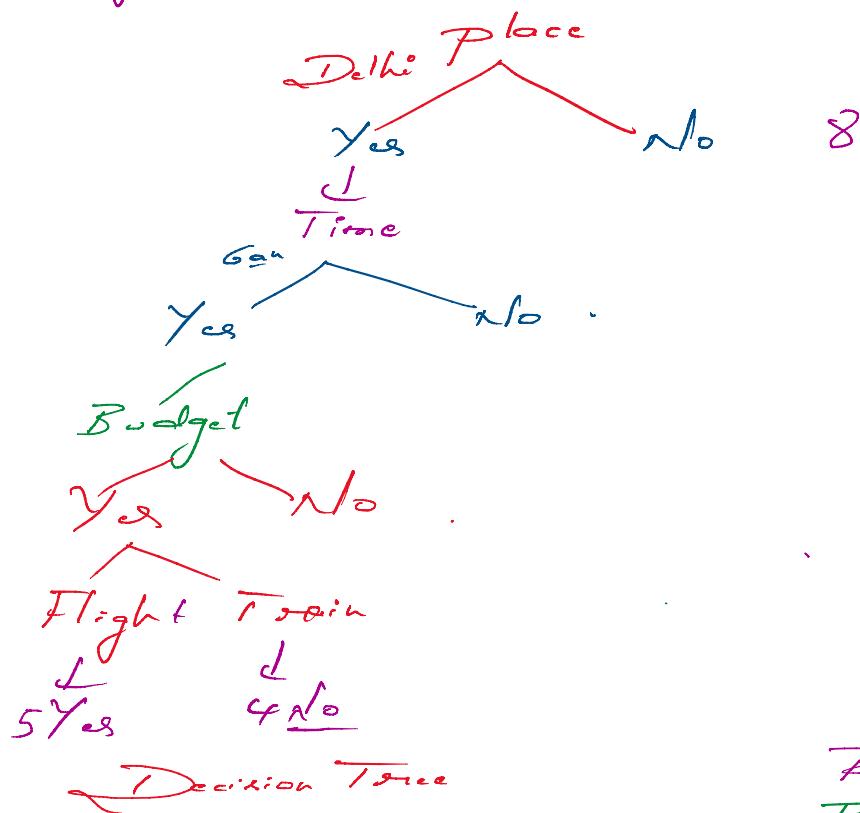
#### ② Decision Tree Classifier

IT is used when target is categorical / discrete

Decision Trees use logic to understand data and to make predictions.

Logic → if-else conditions

Eg:- place, Time, Budget → Transportation



Decision Tree

### Root Node

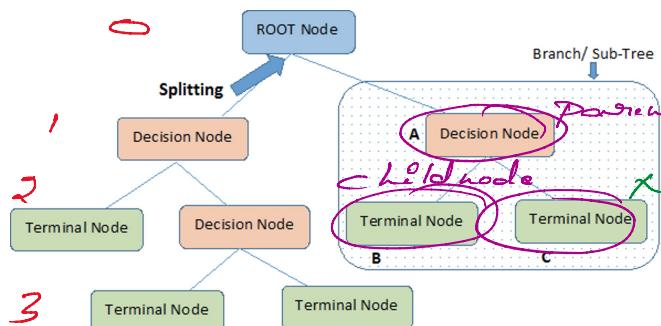
\* Root node represents entire population and this further gets splitted into two or more homogeneous sets.

### Splitting

IT is a process of dividing node into two or more sub nodes.

### Decision Node

... - P - other



Note:- A is parent node of B and C.

Depth/height of tree

node into two -

Decision node

\* When subnode is split out of further subnodes, we call those Decision node subnodes.

\* Decision node will have child nodes.

Terminal node/leaf node.

Node which has no children [no further splits]

Branch/Subtree

A subsection of decision tree

Pruning

The process of removing nodes in order to reduce the size of tree.

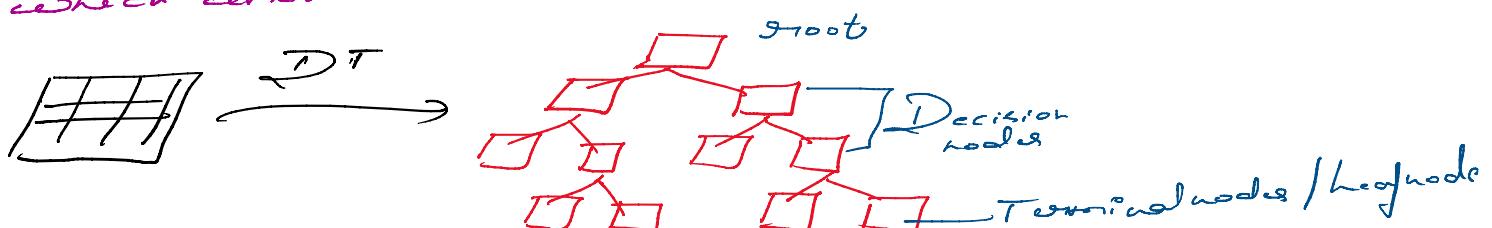
Parent node/child node

A node which is divided into subnodes is called as parent node.

Subnode of parent node is called as child node.

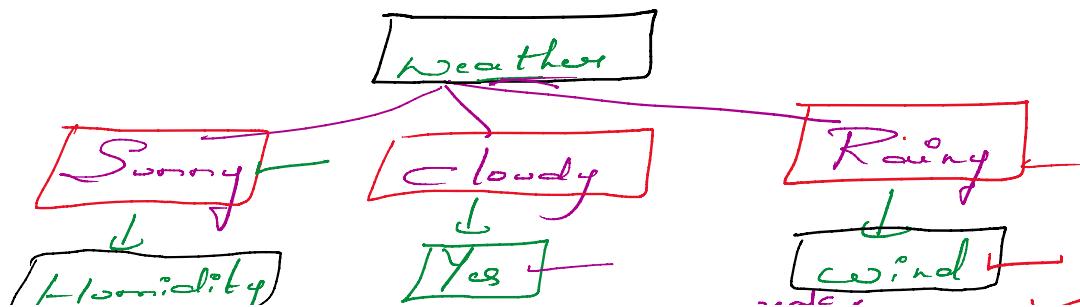
How Decision works?

Decision tree builds classification or regression model in the form of tree structure. It break down entire data into smaller and smaller subsets while at the same time an associated decision tree is incrementally developed. The final result is a tree which concludes decision nodes and leaf nodes.

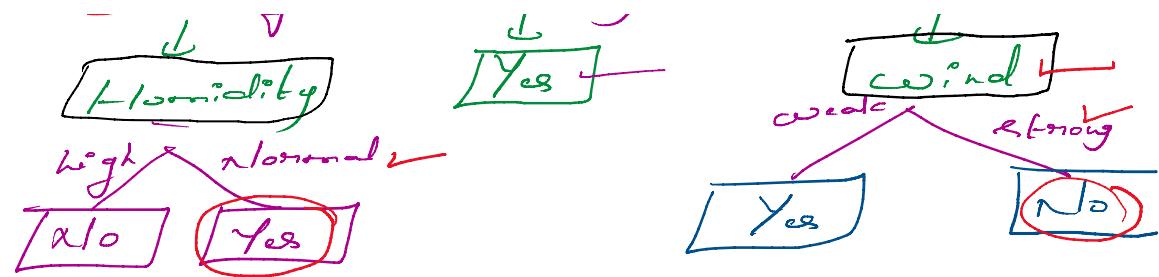


Weather	Humidity	Wind	play
sunny	high	weak	no
sunny	high	strong	no
cloudy	high	weak	yes
rainy	high	weak	yes
rainy	normal	weak	yes
rainy	normal	strong	no
cloudy	normal	strong	yes
sunny	high	weak	no

Decision Tree



cloudy	normal	strong	yes	
sunny	high	weak	no	
sunny	normal	weak	yes	
rainy	normal	weak	yes	
sunny	normal	strong	yes	
cloudy	high	strong	yes	
cloudy	normal	weak	yes	
rainy	high	strong	no	



Rainy / strong wind  
cloudy, high, weak

### Decision Tree working procedure

\* It will search for root node among input variables.  
In this example, we need to select root node among weather, windy and Humidity.

How to Select root node?

- Information gain.
- \* It measures how much information is given by all features about data.
- \* We select the feature which has maximum information gain.

Entropy.  
Entropy is nothing but the uncertainty in a dataset.  
Or measure of randomness.

Eg:- Group of friends  $\rightarrow 10$

Avenger  
6 votes  
5 votes  
100000

Mario  
2 votes  
5 votes  
 $\Rightarrow$   $\rightarrow$  Positivity

5 Yes onto split.

5  
4 Yes 2 No  
Entropy

Entropy  $\rightarrow$  pure split

Entropy

$$D_{\text{Info}} P = P \log P$$

Break  $\rightarrow 10.08$   
 $\rightarrow 10.15$   
 $P_n$

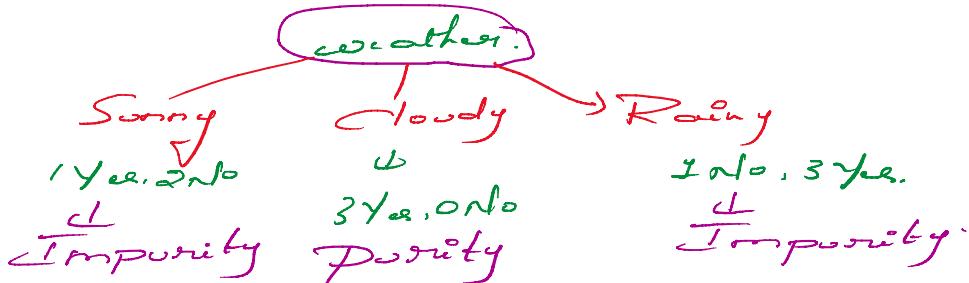
## Entropy

$$H(S) = -P_+ \log_2 P_+ - P_- \log_2 P_-$$

$P_+$  → probability of +ve class

$P_-$  → probability of -ve class.

\* Entropy measures impurity of node



Tip

\* Entropy will be 0 for pure split.

\* Entropy will be high when impurity is high.

## Information Gain

$$\text{Gain} = H(S) - \sum_{v=1}^n \frac{|S_v|}{|S|} H(S_v)$$

$S$  → Sample before split

$S_v$  → Sample after split

$H(S)$  → Entropy before split

$H(S_v)$  → Entropy after split

## Steps

- ① Calculate Entropy of data,  $H(S)$
- ② Find entropy of each category in a column
- ③ Find information gain from that column.
- ④ Repeat ② and ③ for all input variables.
- ⑤ Select one column which has maximum Information Gain.
- ⑥ Repeat steps till we get decision tree.



① Find  $H(S)$



$$P(\text{Yes}) = \frac{9}{14}$$

$$P(\text{No}) = \frac{5}{14}$$

$$H(S) = -P_+ \log_2 P_+ - P_- \log_2 P_-$$

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

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

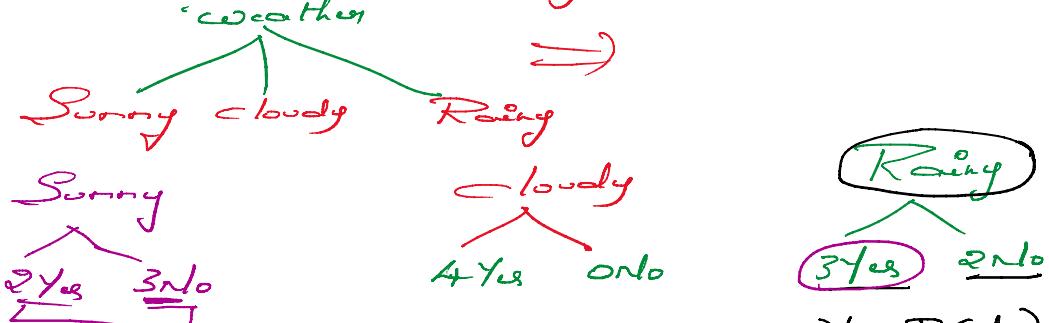
$$H(S) = -P(Y_{es}) \log_2 P(Y_{es}) - P(N_{no}) \log_2 P(N_{no})$$

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

$$H(S) = 0.92$$

Weather	Humidity	Wind	play
sunny	high	weak	no
sunny	high	strong	no
cloudy	high	weak	yes
rainy	high	weak	yes
rainy	normal	weak	yes
rainy	normal	strong	no
cloudy	normal	strong	yes
sunny	high	weak	no
sunny	normal	weak	yes
rainy	normal	weak	yes
sunny	normal	strong	yes
cloudy	high	strong	yes
cloudy	normal	weak	yes
rainy	high	strong	no

② Find Entropy of categories in weather



$$H(S) = -P(Y_{es}) \log_2 P(Y_{es}) - P(N_{no}) \log_2 P(N_{no})$$

$$H(S) = -\left(\frac{9}{14}\right) \log_2 \left(\frac{9}{14}\right) - \left(\frac{5}{14}\right) \log_2 \left(\frac{5}{14}\right) = 0.92$$

$$H(S_{sunny}) = -P(Y_{es}) \log_2 P(Y_{es}) - P(N_{no}) \log_2 P(N_{no})$$

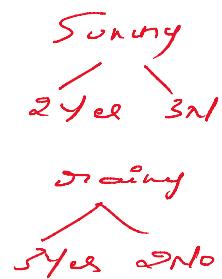
$$H(S_{sunny}) = -\left(\frac{2}{5}\right) \log_2 \left(\frac{2}{5}\right) - \left(\frac{3}{5}\right) \log_2 \left(\frac{3}{5}\right) = 0.97$$

$$H(S_{rainy}) = -P(Y_{es}) \log_2 P(Y_{es}) - P(N_{no}) \log_2 P(N_{no})$$

$$H(S_{rainy}) = -\left(\frac{3}{5}\right) \log_2 \left(\frac{3}{5}\right) - \left(\frac{2}{5}\right) \log_2 \left(\frac{2}{5}\right) = 0.97$$

$$H(S_{cloudy}) = -P(Y_{es}) \log_2 P(Y_{es}) - P(N_{no}) \log_2 P(N_{no})$$

$$H(S_{cloudy}) = -\left(\frac{4}{4}\right) \log_2 \left(\frac{4}{4}\right) - 0 \log_2 0 = 0$$



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

$$= H(S) - \left(\frac{|S_{sunny}|}{|S|}\right) H(S_{sunny}) - \left(\frac{|S_{rainy}|}{|S|}\right) H(S_{rainy}) - \left(\frac{|S_{cloudy}|}{|S|}\right) H(S_{cloudy})$$

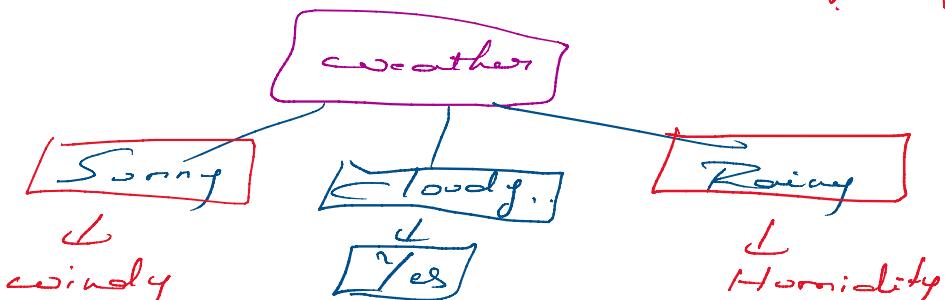
$$= 0.92 - \left(\frac{5}{14}\right) (0.97) - \left(\frac{5}{14}\right) (0.97) - \left(\frac{4}{14}\right) (0)$$

$$\boxed{\text{Gain}(weather) = 0.247}$$

$$Gain(Windy) = 0.12$$

$$Gain(Humidity) = 0.00$$

weather is root node which is giving maximum gain.



windy

Yes No

Yes

Humidity

Yes No