

## DECISION TREE:

⇒ Decision tree is a classification technique which is used to classify discrete value, target attributes.

⇒ The decision tree consist of root nodes, internal nodes, and leaf nodes.

⇒ The root node will be the best attribute.

The best attribute can be derived based on the attribute selection measure like gain information or information gain based on Gini Index.

consider the dataset

Day	weather	Temp	Windy	play
1	Rain	Mild	Weak	No
2	Normal	Hot	Weak	Yes
3	wind	Mild	Strong	Yes
4	Normal	Cool	Weak	No
5	rain	Hot	Strong	No

construct or build the decision tree by finding the best attribute to decide whether to play or not



FORMULA:

(c is the number of class labels)

$$\text{Entropy } E(S) = \sum_{i=1}^c -p_i \log_2 p_i$$

$$\text{Gain}(S_A) = \text{Entropy}(S) - \sum_{i=1}^V \frac{|S_v|}{|S|} \text{Entropy}(S_v)$$

STEP-1: For building a decision tree we need to find the entropy which is used to measure the impurity in the data attribute (Entropy  $E(S)$ )

$\Rightarrow$  c is the number of class labels  
we need to find the best attribute using information gain.

FORMULA FOR INFORMATION GAIN: (given already)  
the attribute with high information gain will be considered as the best attribute.

$$\text{Entropy}(S) = \sum_{i=1}^2 -p_i \log_2 p_i$$

(2, 3)  
2 yes 3 No

$$= \left[ -\frac{2}{5} \times \log_2 \frac{2}{5} \right] + \left[ -\frac{3}{5} \times \log_2 \frac{3}{5} \right]$$

$$= -\frac{2}{5} \times \log_2 \frac{2}{5} - \frac{3}{5} \times \log_2 \frac{3}{5}$$

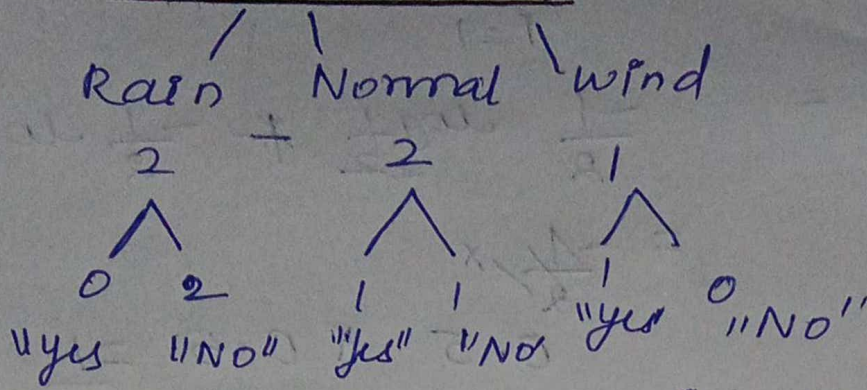
$$= \frac{-2}{5} \times (-1.3219) - \frac{3}{5} \times (-0.736)$$

$$\text{STEP-2:-} \quad = 0.53 + 0.44 = \text{Entropy}(S) = 0.97$$

Find <sup>information</sup> Gain(weather) for 3 attributes (weather, Temp, windy). The attribute with high information gain is best.



## Gain (Weather)



Entropy (Weather = Rain)

$$\sum_{i=1}^C -p_i \log_2 p_i$$

$$= -\frac{0}{2} \log_2 \frac{0}{2} - \frac{2}{2} \log_2 \frac{2}{2} = 0$$

Entropy (Weather = Normal)

$$= -\frac{1}{2} \log_2 \frac{1}{2} - \frac{1}{2} \log_2 \frac{1}{2} = 1$$

Entropy (Weather = Wind)

$$= -\frac{1}{1} \log_2 \frac{1}{1} - \frac{0}{1} \log_2 \frac{0}{1} = 0$$

$$\text{Gain (Weather)} = \text{Entropy (S)} - \sum_{i=1}^3 \frac{|S_i|}{|S|} \text{Entropy (S}_i\text{)}$$

$$= 0.97 - \frac{2}{5} \times 0 - \frac{2}{5} \times 1 - \frac{1}{5} \times 0$$

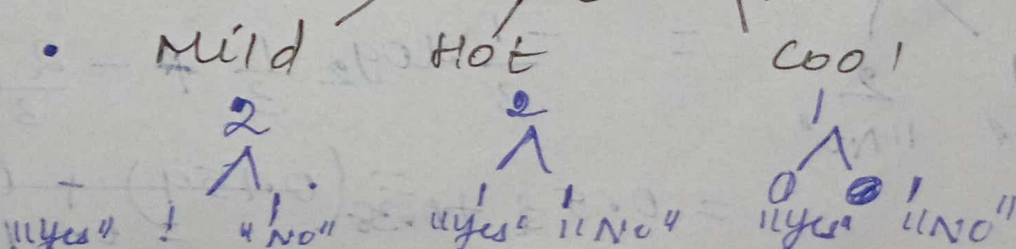
$$= 0.97 - \frac{2}{5}$$

$$\text{Gain (Weather)} = 0.57$$

STEP-3:

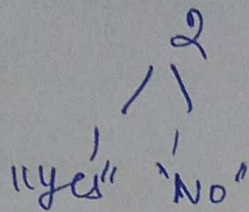
## Gain (Temp)

## Gain (Temp)





Entropy (Temp = Mild)



=

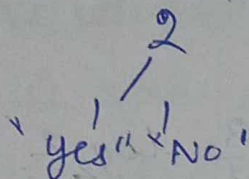
$$= \sum_{i=1}^L -p_i \log p_i$$

$$= -\frac{1}{2} \log_2 \frac{1}{2} - \frac{1}{2} \log_2 \frac{1}{2}$$

$$= -\frac{1}{2} \times 1 - \frac{1}{2} \times 1$$

$$= 0.5 + 0.5 = 1$$

Entropy (Temp = Hot)

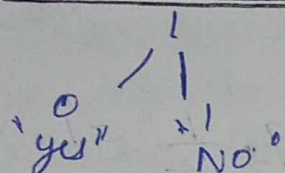


=

$$= -\frac{1}{2} \log_2 \frac{1}{2} - \frac{1}{2} \log_2 \frac{1}{2}$$

$$= 0.5 + 0.5 = 1$$

Entropy (Temp = cool)



$$= -\frac{0}{1} \log_2 \frac{0}{1} - \frac{1}{1} \log_2 \frac{1}{1}$$

$$= 0$$

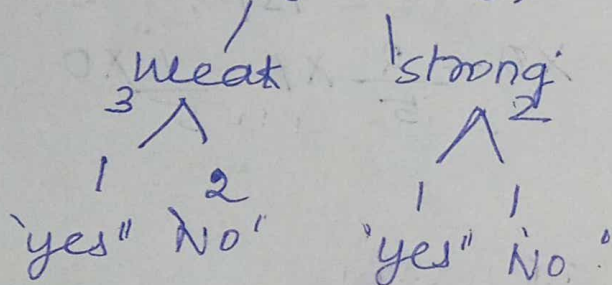
Gain (Temp) = Entropy(S) - \sum\_{i=1}^V \frac{|S\_v|}{|S|} Entropy(S\_v)

$$= 0.97 - \frac{2}{5} \times 1 - \frac{2}{5} \times 1 - \frac{1}{5} \times 0$$

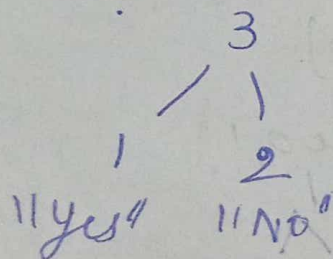
Gain (Temp) = 0.17

STEP-4:-

Gain (Windy)



Entropy (Pet Windy = weak)



=

$$= \sum_{i=1}^L -p_i \log p_i$$

$$= -\frac{1}{3} \log_2 \frac{1}{3} - \frac{2}{3} \log_2 \frac{2}{3}$$

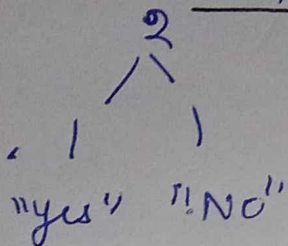
$$= 0.53 + 0.38 = 0.91$$

Entropy

Gain (windy) = strong

$$-\frac{1}{2} \times \log_2 \frac{1}{2} - \frac{1}{2} \log_2 \frac{1}{2}$$

$$= 0.5 + 0.5 = \underline{1}$$



$$\underline{\text{Gain (windy)}} = 0.97 - \frac{2}{5} \times 0.91 - \frac{2}{5} \times 1$$

$$0.97 - 0.546 - 0.4 = \underline{0.024}$$



$$\text{Gain Info (Weather)} = 0.57$$

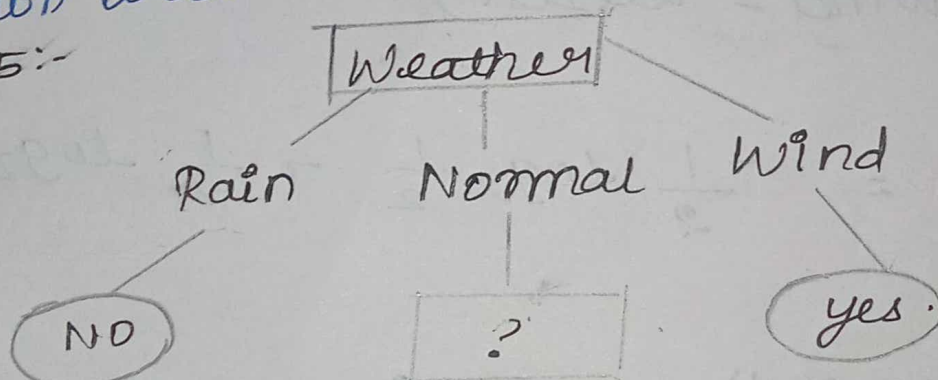
$$\text{Gain Info (Temperature)} = 0.17$$

$$\text{Gain Info (Wind)} = 0.024$$

"Weather"  $\Rightarrow$  The highest gain information value is "0.57" which represents the attribute weather.

The weather attribute is the root of the decision tree.

STEP-5:-



(yes and no  
we for Normal  
cannot  
conclude  
attribute)

To find the next attribute

New dataset  
S:-

Temperature	Wind	play
Hot	Weak	Yes
Cool	Weak	No

Entropy(S)

$\swarrow \searrow$   
 "yes" "no"

$$= -\frac{1}{2} \log_2 \frac{1}{2} - \frac{1}{2} \log_2 \frac{1}{2} = 1$$

$$\text{Gain Info (Temperature)} = \{ \text{Hot, Cool} \}$$

$$\text{Entropy (Temp = "Hot")}$$

$$\begin{array}{c} 1 \\ \swarrow \searrow \\ \text{"yes"} \quad \text{"No"} \end{array} = -\frac{1}{1} \log_2 \frac{1}{1} - \frac{0}{1} \log_2 \frac{0}{1} = 0$$

$$\text{Entropy (Temp = "Cool")}$$

$$\begin{array}{c} 1 \\ \swarrow \searrow \\ 0 \quad 1 \\ \text{"yes"} \quad \text{"No"} \end{array} = -\frac{0}{1} \log_2 \frac{0}{1} - \frac{1}{1} \log_2 \frac{1}{1} = 0$$

$$\begin{aligned} \text{Gain Info (Temperature)} &= \text{Entropy (S)} - \sum_{i=1}^V \frac{|S_i|}{|S|} \text{Entropy (S)} \\ &= 1 - \frac{1}{2} \times 0 - \frac{1}{2} \times 0. \end{aligned}$$

$$\text{Gain Info (wind)} = \{ \text{Weak} \}$$

$$\text{Entropy (wind = "Weak")}$$

$$\begin{array}{c} 2 \\ \swarrow \searrow \\ 1 \quad 1 \\ \text{"yes"} \quad \text{"No"} \end{array} = -\frac{1}{2} \log_2 \frac{1}{2} - \frac{1}{2} \log_2 \frac{1}{2} = 1.$$

$$\text{Gain Info (wind)} = 1 - \frac{2}{2} \times 1 = 0.$$

Now,

$$\text{Gain Info (Temp)} = 1$$

$$\text{Gain Info (wind)} = 0$$

Here, Gain Info of 'temp' is higher, so we take temperature as the next attribute.

+ Decision tree:

