

# Decision Tree

① Decision Tree classifier [classification]

② Decision Tree Regressor [Regression]

Decision Tree Classifier  $\Rightarrow$

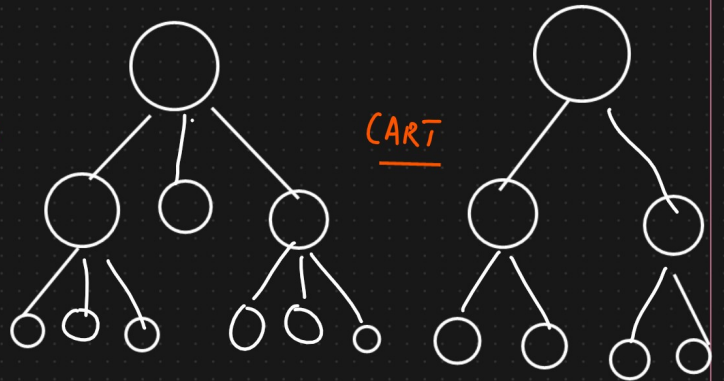
Two types

① ID3 [Iterative Dichotomiser 3]

② CART ✓ [Classification And Regression Tree]

ID3

CART



Age = 14

if (age  $\leq 15$ ):

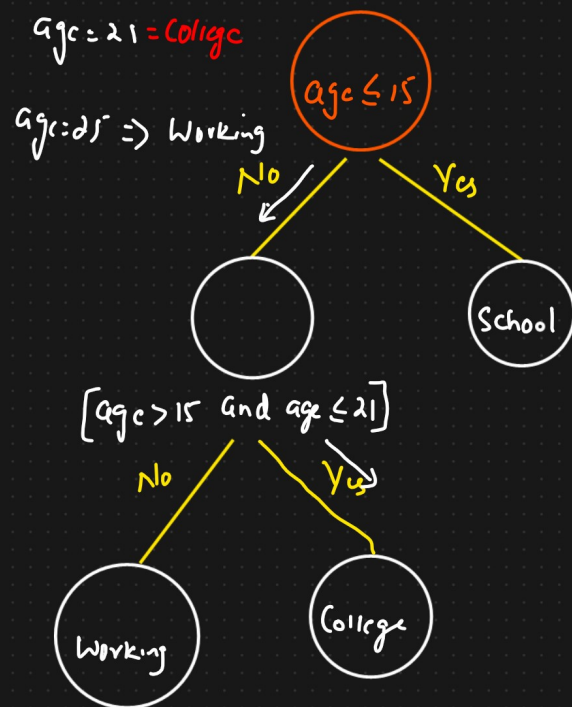
Print ("School")

elif (age > 15 and age  $\leq 21$ ):

Print ("College")

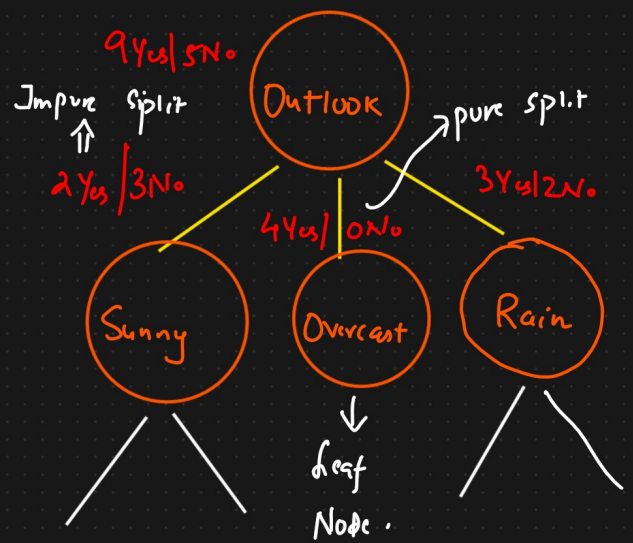
else:

Print ("Working")



Dataset → Predict Play Tennis or Not

Day	Outlook	Temperature	Humidity	Wind	Play Tennis
1	Sunny	Hot	High	Weak	No
2	Sunny	Hot	High	Strong	No
3	Overcast	Hot	High	Weak	Yes
4	Rain	Mild	High	Weak	Yes
5	Rain	Cool	Normal	Weak	Yes
6	Rain	Cool	Normal	Strong	No
7	Overcast	Cool	Normal	Strong	Yes
8	Sunny	Mild	High	Weak	No
9	Sunny	Cool	Normal	Weak	Yes
10	Rain	Mild	Normal	Weak	Yes
11	Sunny	Mild	Normal	Strong	Yes
12	Overcast	Mild	High	Strong	Yes
13	Overcast	Hot	Normal	Weak	Yes
14	Rain	Mild	High	Strong	No



① Purity check : Pure Split or Impure Split

{ Entropy  
 Gini Impurity } Measure of purity. ✓

② What feature you need to select to start the split? → Information Gain

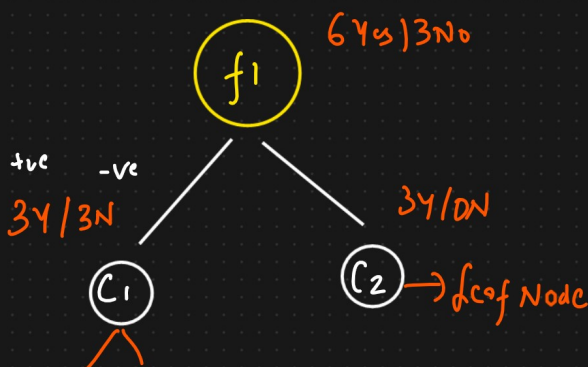
Binary classification

① Entropy

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

$P_+$  = probability of positive category

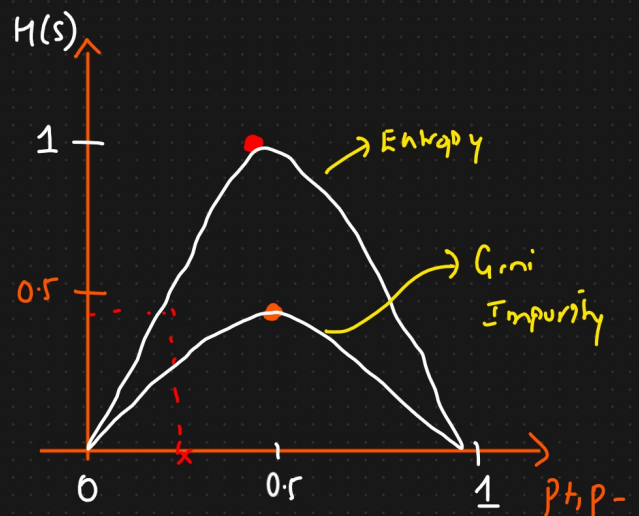
$P_-$  = probability of negative category



$$\begin{aligned}
 H(C_1) &= -P_+ \log_2 P_+ - P_- \log_2 P_- \\
 &= -\frac{3}{6} \log_2 \left(\frac{3}{6}\right) - \left(\frac{3}{6}\right) \log_2 \left(\frac{3}{6}\right)
 \end{aligned}$$

② Gini Impurity

$$G.I = 1 - \sum_{i=1}^n (p_i)^2$$



$$= -\frac{1}{2} \log_2(1/2) - (1/2) \log_2(1/2).$$

$$= 1 \Rightarrow \text{Impure Split}$$

$$H(c_2) = -\frac{3}{3} \log_2(3/3) - 0/3 \log_2(0/3)$$

$$= 0 \Rightarrow \text{Pure Split}$$

Multiclass

$c_1 \quad c_2 \quad c_3$   
Yes / No / May Be

$$H(s) = -p_{c_1} \log_2 p_{c_1} - p_{c_2} \log_2 p_{c_2} - p_{c_3} \log_2 p_{c_3}$$

② Gini Impurity

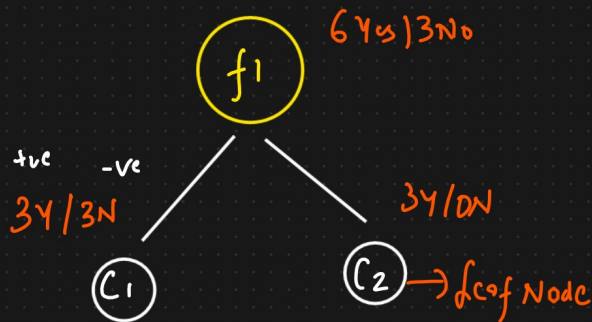
$$G.I = 1 - \sum_{i=1}^n (p_i)^2$$

$$= 1 - [(p_+)^2 + (p_-)^2]$$

$$= 1 - \left[\left(\frac{3}{6}\right)^2 + \left(\frac{3}{6}\right)^2\right]$$

$$= 1 - \left[\frac{1}{4} + \frac{1}{4}\right]$$

$$= 1/2 = 0.5 \Rightarrow \text{Impure Split}$$



$$1 - \left[\left(\frac{3}{3}\right)^2 + \left(\frac{0}{3}\right)^2\right]$$

$$= 1 - 1$$

$$= 0 // \Rightarrow \text{Pure Split}$$

② What feature you need to select to

start the split?  $\rightarrow$  Information Gain

$f_1 \quad f_2 \quad f_3 \quad \text{q/p}$



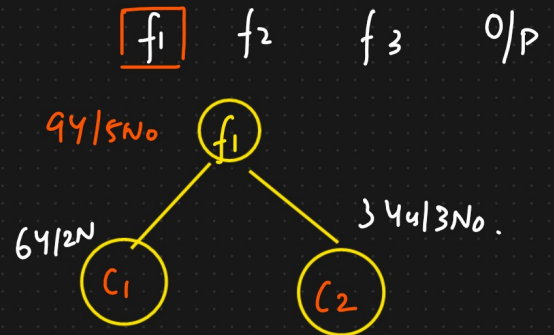


$f_1$   $f_2$

Information Gain

Entropy of the root node

$$Gain(S, f_1) = \boxed{H(S)} - \sum_{v \in Val} \frac{|S_v|}{|S|} H(S_v)$$



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

$$= -\frac{9}{14} \log_2 \frac{9}{14} - \frac{5}{14} \log_2 \frac{5}{14}$$

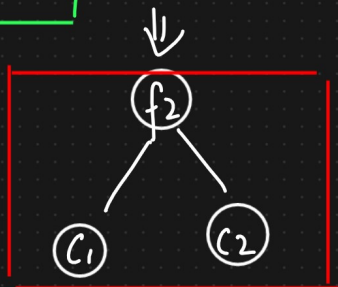
$$\approx 0.94 //$$

$$H(c_1) = -\frac{6}{8} \log_2 \left(\frac{6}{8}\right) - \frac{2}{8} \log_2 \frac{2}{8} \approx \boxed{0.81}$$

$$H(c_2) = \underline{1} //$$

$$Gain(S, f_1) = 0.94 - \left[ \frac{8}{14} * 0.81 + \frac{6}{14} * 1 \right]$$

$$Gain(S, f_1) = 0.049$$



$$\boxed{Gain(S, f_2) = 0.051} > \boxed{Gain(S, f_1) = 0.049}$$

Information gain is more when we split using  $f_2$ .

## Entropy Vs Gini Impurity

When dataset is small  $\rightarrow$  Entropy  
dataset is huge  $\rightarrow$  Gini Impurity