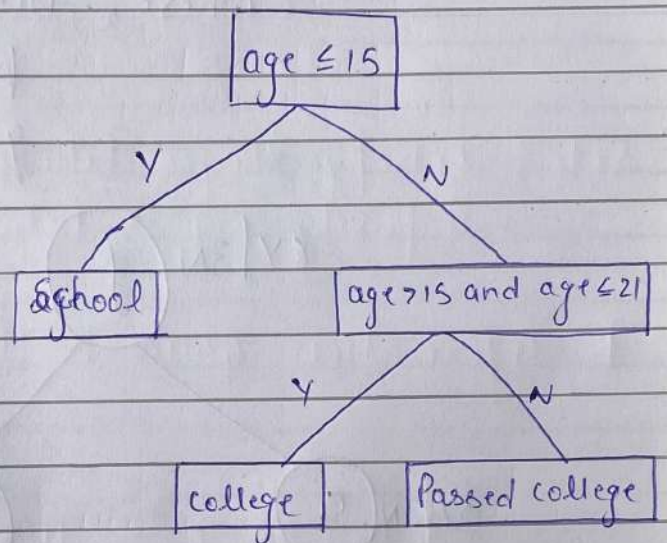


## Decision Tree

→ It is a supervised machine learning algorithm which is used to solve both classification and regression problem by splitting the dataset repeatedly till leaf node achieve.

```
age = 14
if age ≤ 15:
    print('school')
elif age > 15 and age ≤ 21:
    print('college')
else:
    print('Passed college')
```



Q But how the splitting is happening and on which condition/feature?

→

① Purity → Pure Split ??  
    └→ Entropy  
    └→ Gini Impurity

② Gain/Information Gain → How the feature are selected?

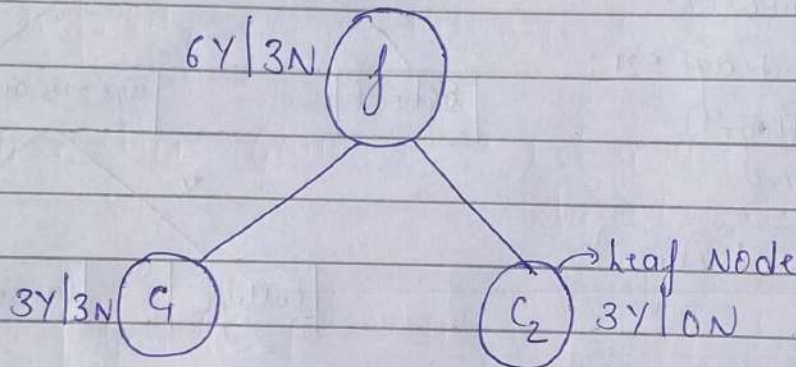
\* Entropy :  $\rightarrow$  Degree of disorganization or randomness

$$H(s) = -p \log_2 p - q \log_2 q$$

where,  $H(s) \rightarrow$  Entropy

$p \rightarrow$  probability of success

$q \rightarrow$  probability of failure.



• Entropy of  $c_1$  :

$$H(c_1) = -\frac{3}{6} \log_2 \frac{3}{6} - \frac{3}{6} \log_2 \frac{3}{6}$$

$$= 1 \rightarrow \text{Impure split}$$

$\rightarrow$  Need more splitting

• Entropy of  $c_2$  :

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

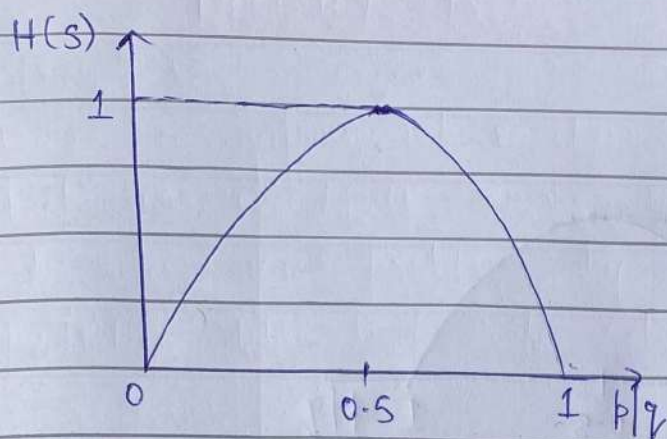
$$= -1 \log 1$$

$$= 0 \rightarrow \text{Pure split}$$

$\rightarrow$  leaf node

$\rightarrow$  No more splitting



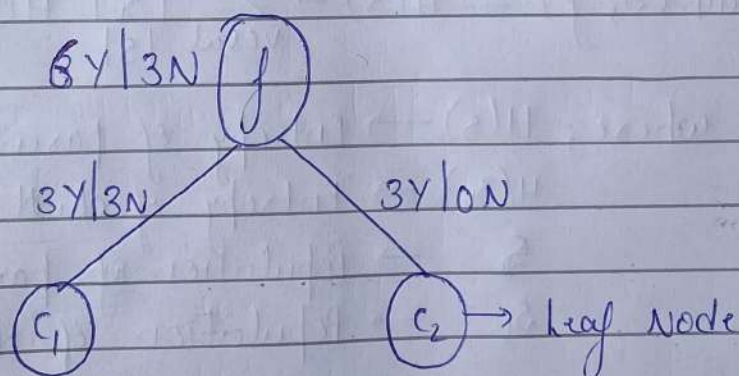


Entropy Graph

\* Gini Impurity :  $\rightarrow$  purity of split

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

$$\Rightarrow G.I = 1 - [p^2 + q^2]$$



• Gini Impurity of  $C_1$  :

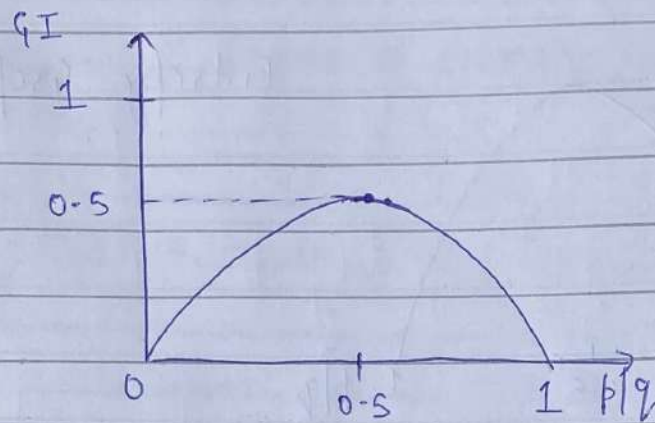
$$GI = 1 - \left[ \left(\frac{1}{2}\right)^2 + \left(\frac{1}{2}\right)^2 \right]$$

$$= 1 - \frac{1}{2} = 0.5 \rightarrow \text{Impure split} \rightarrow \text{Need more splitting}$$

• Gini Impurity of  $C_2$  :

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

$$= 0 \rightarrow \text{Pure split} \rightarrow \text{leaf Node} \rightarrow \text{no more splitting}$$



Note : Entropy ranges between 0 to 1 whereas gini impurity ranges between 0 to 0.5.

\* Information Gain :  $\rightarrow$  Determine which feature should select for splitting.

$$\text{Gain}(S, f) = H(S) - \sum_{v \in \text{val}} \frac{|S_v|}{|S|} H(S_v)$$

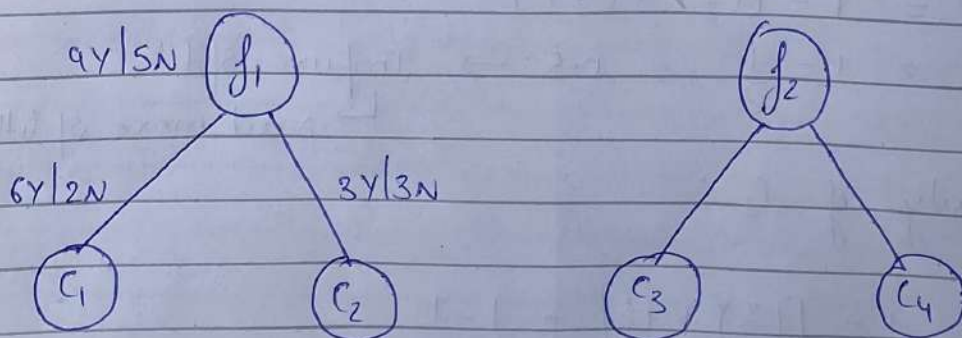
where,  $H(S) \rightarrow$  Entropy of parent node

$H(S_v) \rightarrow$  Entropy of child node

$S \rightarrow$  Population of parent node

$S_v \rightarrow$  Population of child node

Eg : we have two features ( $f_1$  and  $f_2$ ). which feature should we select for splitting.





• Entropy :

$$H(S) = -\frac{9}{14} \log_2 \frac{9}{14} - \frac{5}{14} \log_2 \frac{5}{14} \left. \vphantom{H(S)} \right\} \text{Parent Node}$$

$$= 0.94$$

$$H(C_1) = -\frac{6}{8} \log_2 \frac{6}{8} - \frac{2}{8} \log_2 \frac{2}{8} \left. \vphantom{H(C_1)} \right\} C_1$$

$$= 0.81$$

$$H(C_2) = -\frac{3}{6} \log_2 \frac{3}{6} - \frac{3}{6} \log_2 \frac{3}{6} \left. \vphantom{H(C_2)} \right\} C_2$$

$$= 1$$

Root ~~at~~  
Child Nodes

• Information Gain :

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

$$= 0.049$$

Similarly, let's suppose information gain for  $f_2$

$$\text{Gain}(S, f_2) = 0.052$$

$$\therefore \text{Gain}(S, f_2) > \text{Gain}(S, f_1)$$

So, feature  $f_2$  will be selected for splitting