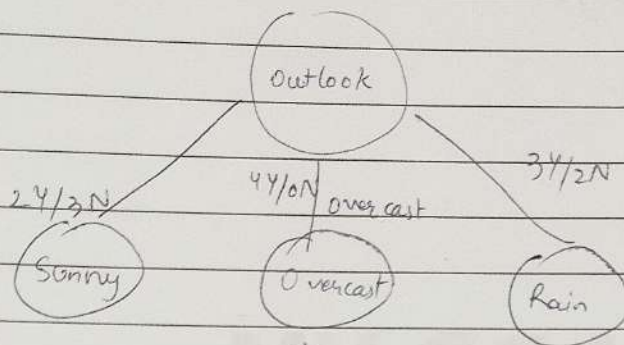


DECISION TREE



Pure leaf Node. \rightarrow Either with 0N/0Y.

(1) How can we say if it is a pure split \rightarrow
 \rightarrow Entropy
 \rightarrow Gini Coeff.

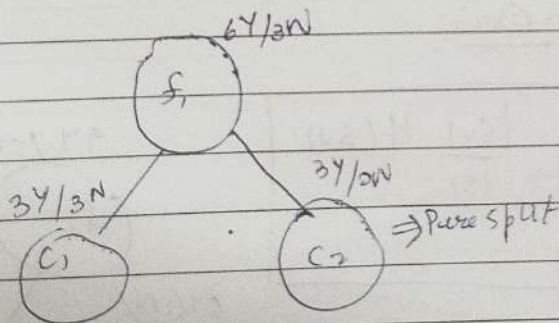
(2) How features are selected
 \rightarrow Information Gain

ENTROPY

$$H(S) = -\sum_{i=1}^n P_i \log_2 P_i \rightarrow H_S = -P \log_2 P - (1-P) \log_2 (1-P)$$

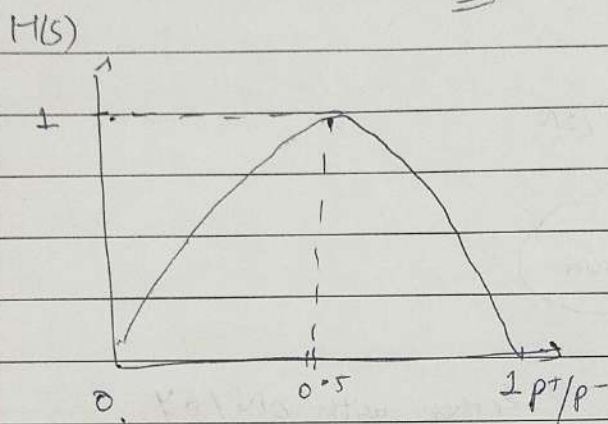
GINI IMPURITY

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



Entropy right node $H(S) = -\frac{3}{3} \log_2 \frac{3}{3} - \frac{0}{3} \log_2 \frac{0}{3}$
 $= 0$

Entropy of left node = $-0.5 \log_2 0.5 - 0.5 \log_2 0.5$
 $= -1 \log_2 0.5$
 $= \underline{\underline{1}}$



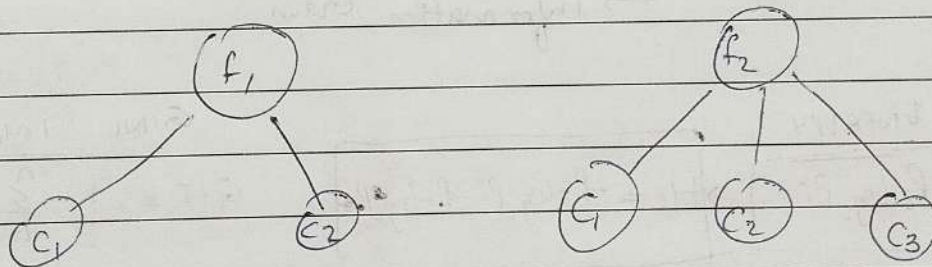
Entropy b/w $\boxed{0-1}$

Using Purity Test,

If $H(S) = 0 \rightarrow$ Pure split

$H(S) = 1 \rightarrow$ Impure split

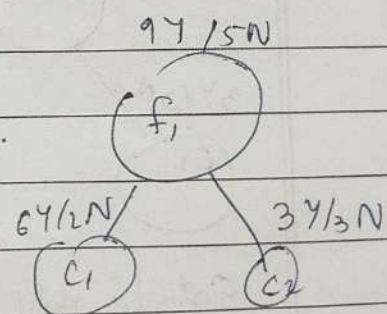
(2) which feature to take to split?



\rightarrow Here we use Information Gain.

$$\text{Gain}(S, f_i) = H_S - \sum_{\substack{\text{value} \\ \text{belong}}} \frac{|S_v|}{|S|} H(S_v)$$

$\downarrow \quad \downarrow$
 Sample Feature f_i



$H(S) \rightarrow$ Root node = $-P_+ \log_2 P_+ - P_- \log_2 P_-$

$$= -\frac{9}{14} \log_2 \frac{9}{14} - \frac{5}{14} \log_2 \left(\frac{5}{14} \right) = 0.94$$

$$H_{\text{src}}(C_1) = -\frac{6}{8} \log_2\left(\frac{6}{8}\right) - \frac{2}{8} \log_2\left(\frac{2}{8}\right) = 0.81$$

$$H_{\text{src}}(C_2) = -\frac{3}{6} \log_2(0.5) - \frac{3}{6} (\log_2(0.5)) = \underline{\underline{1}}$$

$$\begin{aligned} \text{Gain}(S, F_1) &= 0.94 - \left[\frac{8}{14} \times 0.81 + \frac{6}{14} \times 1 \right] \\ &= 0.94 - [0.574 + 0.4285] \\ &= 0.94 - [0.8913] = \underline{\underline{0.049}} \end{aligned}$$

$$\text{Gain}(S, F_2) = \underline{\underline{0.051}} \rightarrow \text{Using which feature should I start splitting first?}$$

$$\text{Gain}(S, F_2) > \text{Gain}(S, F_1)$$

\Rightarrow We will use feature ^{first} 2 to split data.

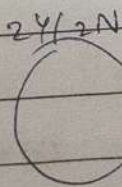
GINI IMPURITY

$$G.I = 1 - \sum_{i=1}^n (P_i)^2 \rightarrow \text{Gini Value } [0 - 0.5]$$

$n=2$ output $\begin{cases} \text{Yes} \\ \text{No} \end{cases}$

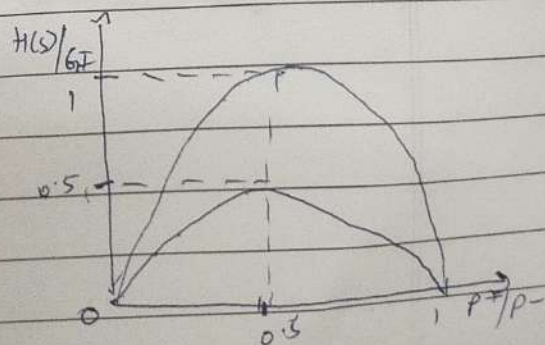
$$G.I = 1 - \left[(P_+)^2 + (P_-)^2 \right]$$

Eg:-



$$\begin{aligned} G.I. &\Rightarrow 1 - [0.5^2 + 0.5^2] \\ &= \underline{\underline{0.5}} \end{aligned}$$

→ Impure split



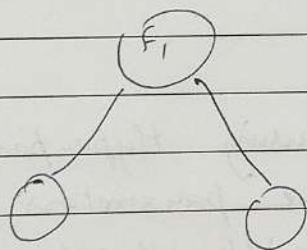
Decision Tree Regressor

f_1 f_2 O/P
Continuous

⇒ Calculate mean of O/P & we will not check Entropy / Gini coeff,
MSE or MAE will get selected calculated.

$$\frac{1}{2m} \sum_{i=1}^n (y_i^1 - y_i)^2$$

⇒ Suppose we select f_1

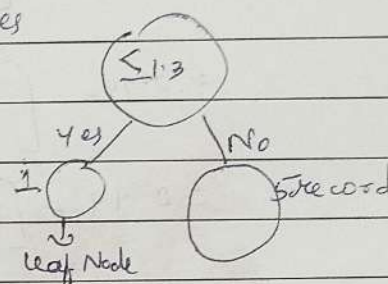


- Calculate MSE or MAE
- When we split Node certain records will go to nodes.
- O/P of nodes will be average or mean.
- Then again node will split, records will go & mean will be calculated.
- Along the way MSE will get reduced & will approach towards leaf node or 0.

⇒ Because there is log fn in case of Entropy, Gini Impurity is fast.

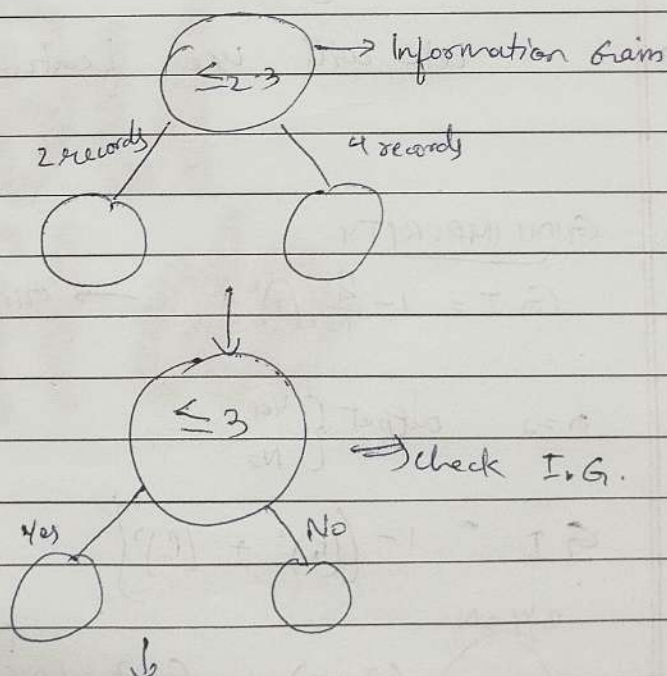
for Continuous Variables

| f_i | O/P | O/P ⇒ f_i | Sort Values |
|-------|-----|-------------|-------------|
| 2, 3 | | 1, 3 | |
| 1, 3 | | 2, 3 | |
| 4 | | 3 | |
| 5 | | 4 | |
| 7 | | 5 | |
| 3 | | 7 | |



Calculate Information gain of this

After create D.T for ~~1, 2, 3~~ 2, 3

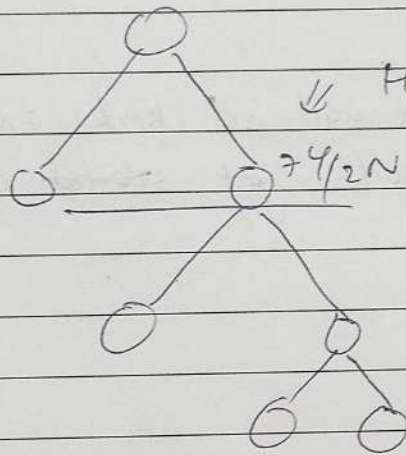


After this it will check which I.G. is larger & then it will get selected.

Hyperparameters

- Decision tree leads to overfitting
- To prevent overfitting we will do -
 - ① Post Pruning
 - ② Pre Pruning

① Post Pruning



⇓ Here More than 80% ~~for~~ is there to get $74/2N$ Yes. we will cut the tree.

- ② Pre Pruning - It is done using Hyperparameters
Max depth, max-leaf → These parameters can be ~~done~~ ^{set} by Grid Search CV.