

Agenda

- ① Decision Tree Classification
- ② Decision Tree Regression
- ③ Practical Implementation
- ④ Ensemble Technique

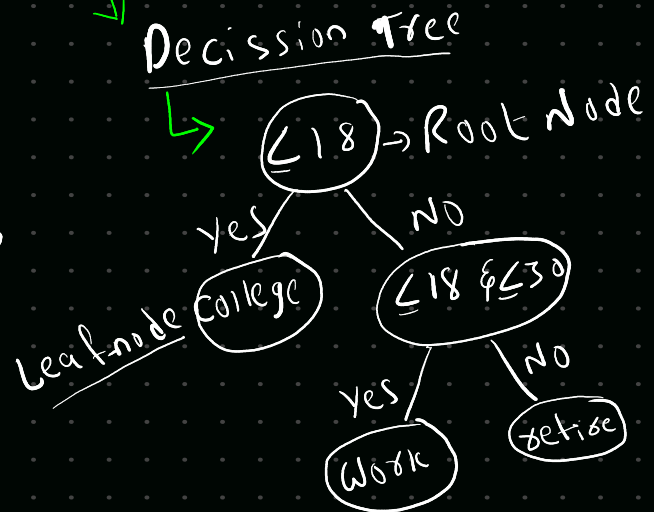
Agenda

Decision tree {Solving many usecase}

- ↳ Regression
- ↳ Classification

```

if (age < 18):
    print("college")
elif (age > 18 and age < 35):
    print("work")
else:
    print("retire")
    
```



Decision Trees

Nest if else \Rightarrow Decision tree

Sunny — 2 Yes
3 NO

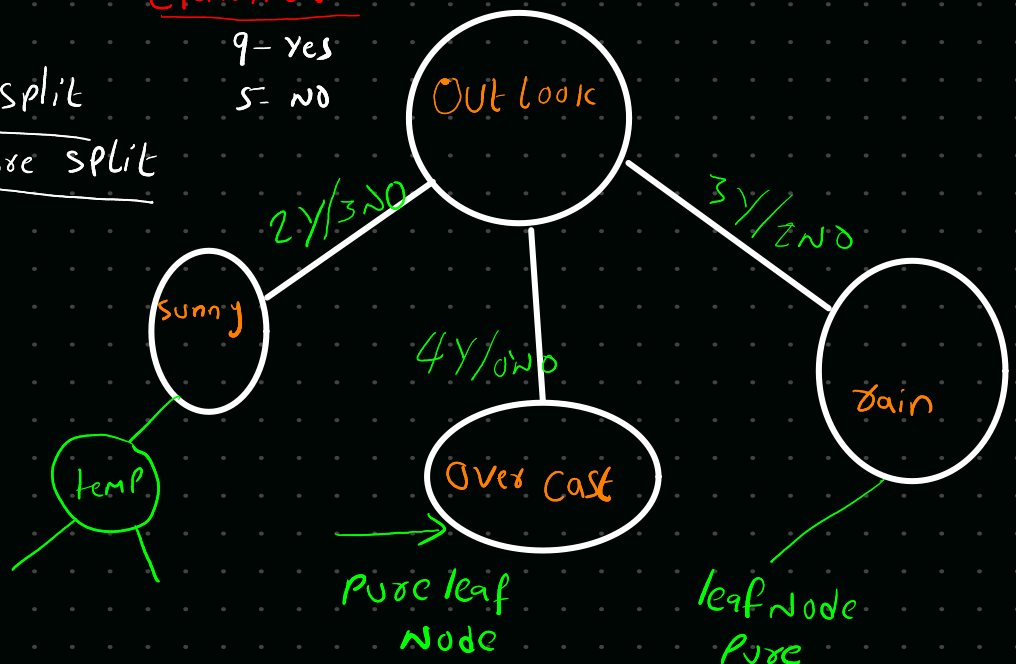
Over cast — 4 Yes
0 NO

Rain — 3 Yes
2 NO

Pure split
impose split

Classification

9- yes
5- no



① Purity \rightarrow Pure Split !!

\hookrightarrow Entropy

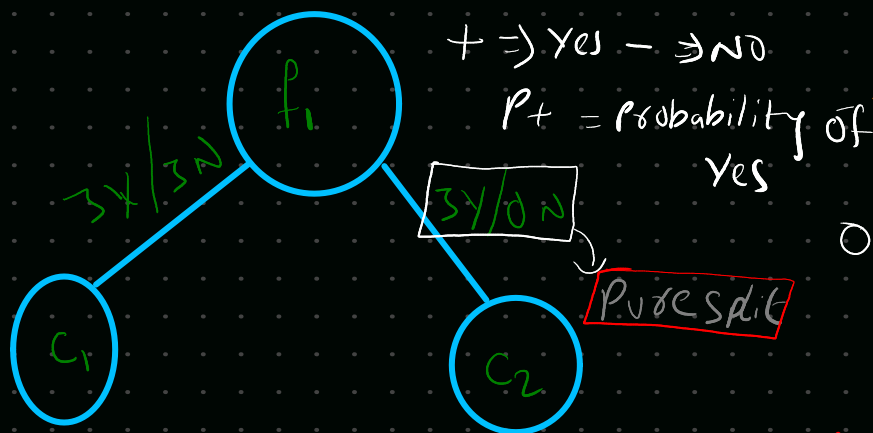
\hookrightarrow Gini Impurity

② How the features are selected

\hookrightarrow Information Gain!

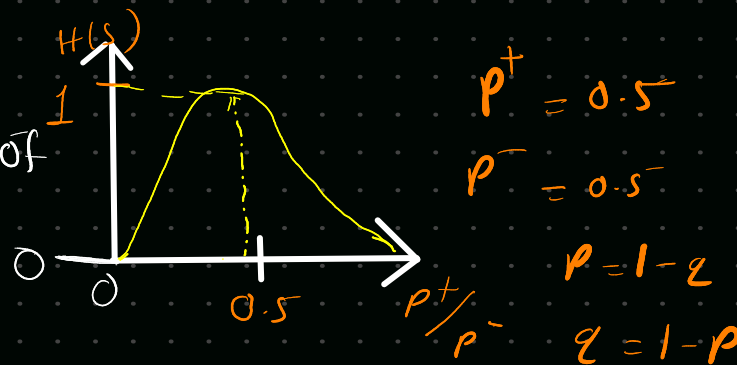
① Entropy

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



① Gini impurity

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



Entropy $H(S) = -\frac{3}{5} \log_2 \frac{3}{5} - \frac{0}{5} \log_2 \frac{0}{5}$ $H(S) = -\frac{3}{6} \log_2 \frac{3}{6} - \frac{3}{6} \log_2 \frac{3}{6}$

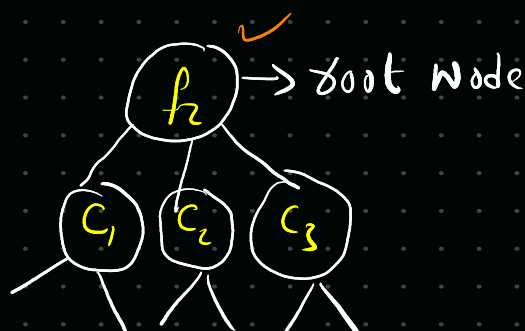
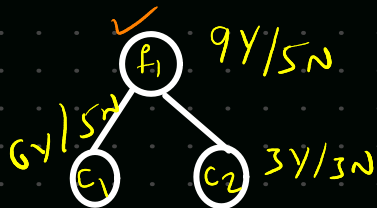
$= -1 \log_2 1$

$= \boxed{0} \rightarrow$ Pure test

\hookrightarrow Entropy

$\boxed{1}$ \rightarrow impure split

② Which feature to take to split?



Information Gain :-

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

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

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

$$\approx \boxed{0.94}$$

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

$$= 0.049$$

$$\text{Gain}(S, f_1) = 0.049$$

$$\text{Gain}(S, f_2) = 0.051$$

Using which feature should I start splitting first

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

* Gini Impurity :-

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

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

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

$$= 1 - \left[\frac{1}{2} \right] = 0.5$$

$n=2$ output { Yes
No

$2Y/2N \Rightarrow$ Impure split



Entropy = 1

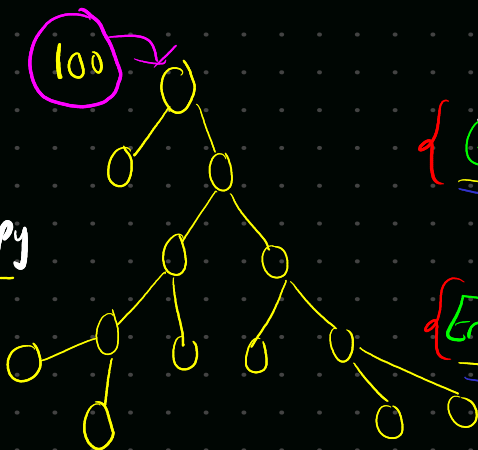
Gini Impurity = 0.5

Entropy \Rightarrow log

Gini impurity \Rightarrow Simple Math

Fast

Gini \gg Entropy



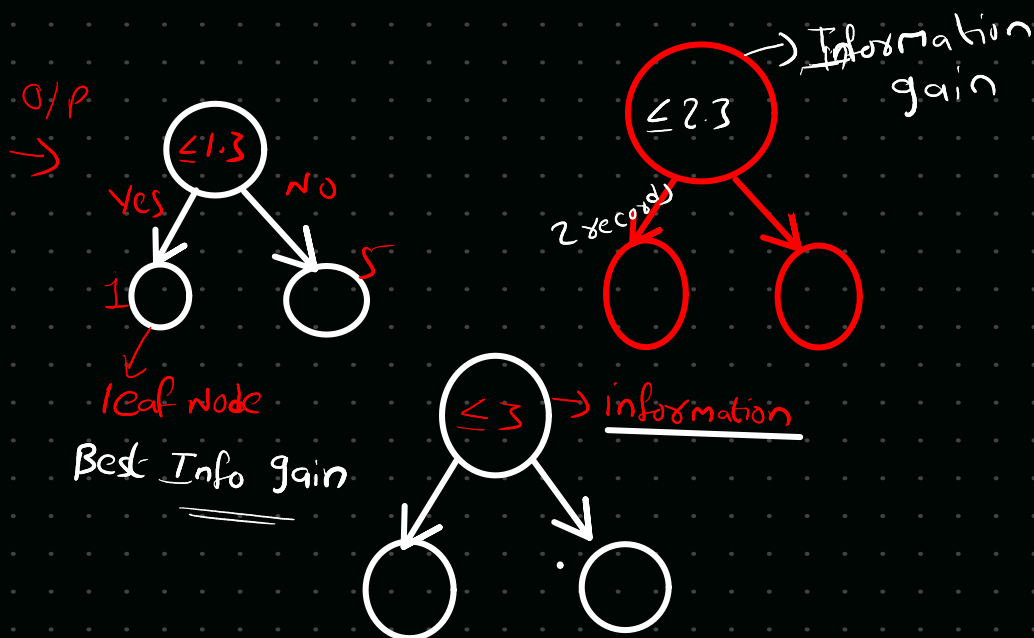
Gini

\rightarrow Should be used when huge numbers are here

Entropy

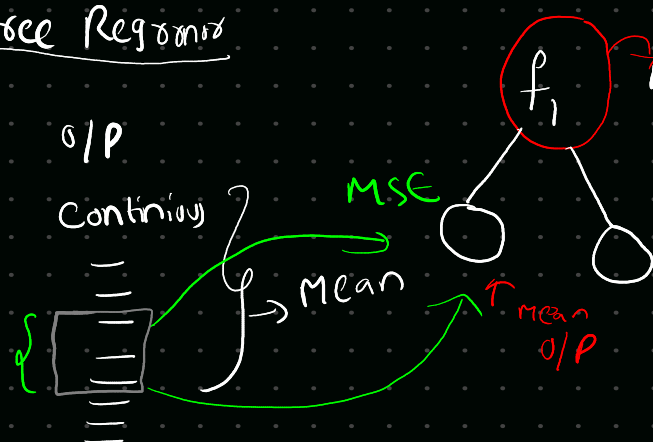
\rightarrow Should be used less no. are here

f_1 o/p $\Rightarrow f_1$
 23 \Rightarrow 1.3
 13 \Rightarrow 2.3
 9 \Rightarrow 3
 5 \Rightarrow 4
 7 \Rightarrow 5
 3 \Rightarrow 7



Decision tree Regions

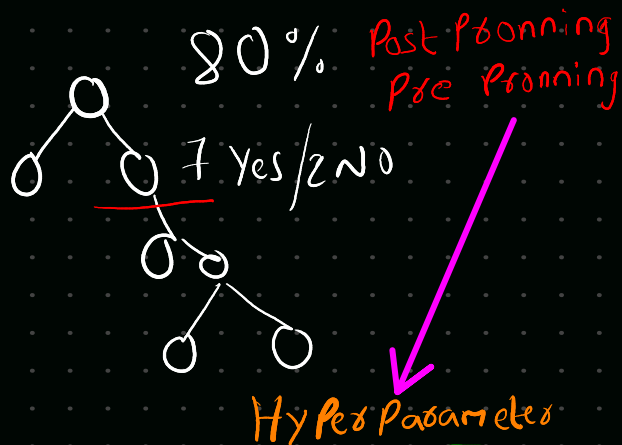
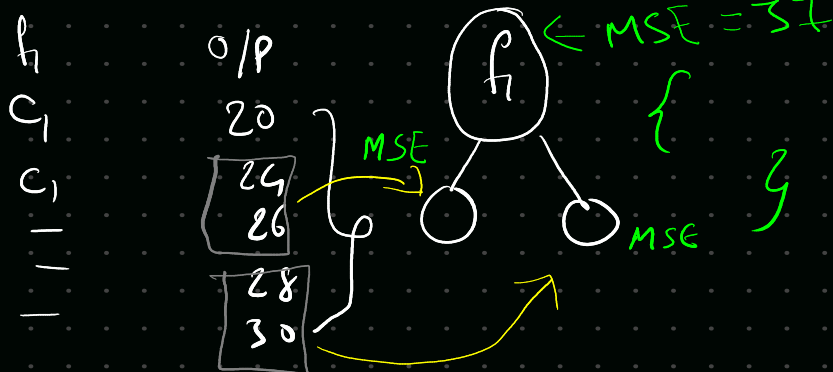
f_1 f_2 o/p
 Continuous



Mean \rightarrow MSE or MAC
 MSE - Mean squared error

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

Decision tree regions

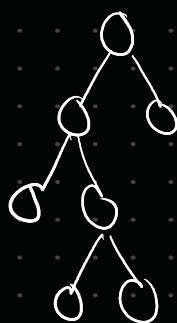


Hyper Parameters
 Max - depth
 Max - leaf \rightarrow Grid Search CV

Hyper Parameters:-

Decision \rightarrow over fitting

- ① Post Pruning
- ② Pre Pruning



Over fitting

