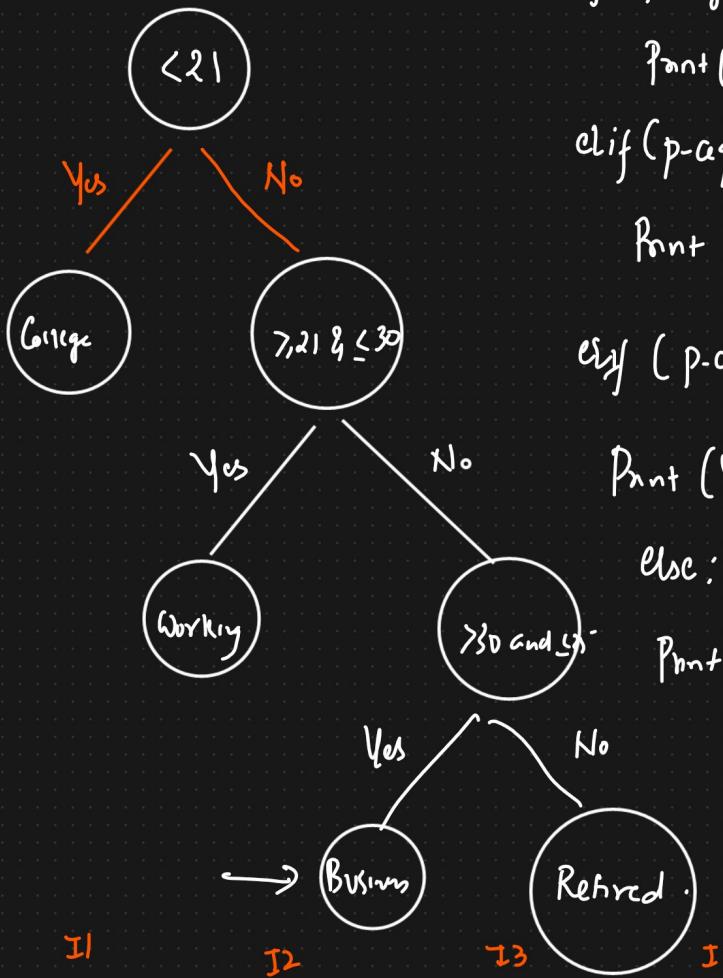


Decision Tree Classifier And Decision Tree Regressor

① Classification

② Regression Problem

$$age = 34$$



Nested if-else clause

if (P-age < 21):

Print ("He should be in college")

elif (p-age >= 21 and p-age <= 30):

Print ("He should be working")

else (p-age > 30 and p-age <= 35):

Print ("He was a business")

else:

Print ("He has retired")

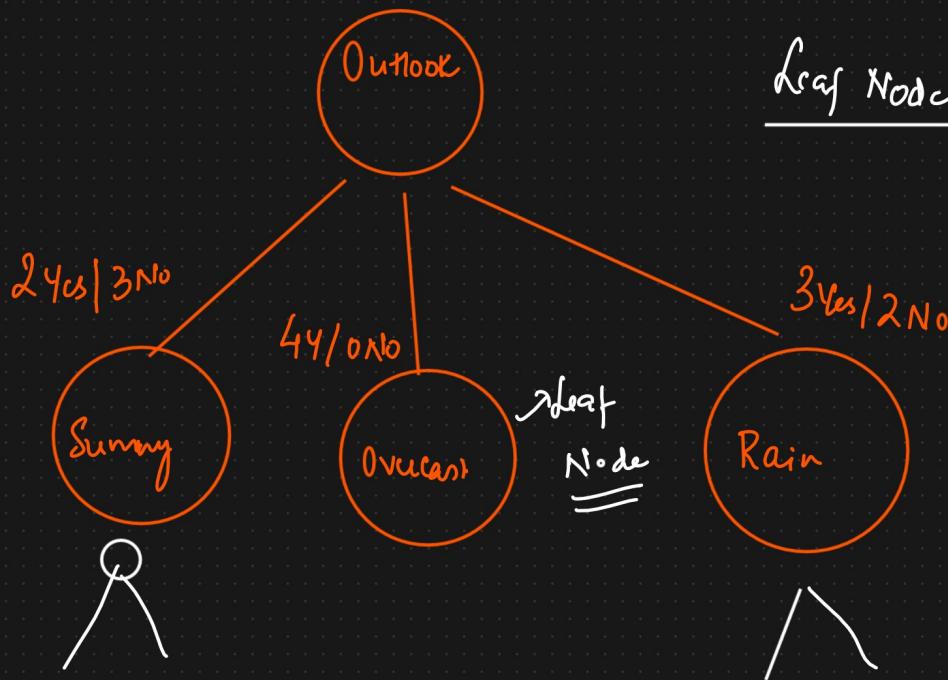
I1	I2	I3	I4
Outlook	Temperature	Humidity	Wind
Sunny ✓	Hot	High	Weak
Sunny	Hot	High	Strong
Overcast ✓	Hot	High	Weak
Rain ✓	Mild	High	Weak
Rain	Cool	Normal	Weak
Rain	Cool	Normal	Strong
Overcast	Cool	Normal	Strong
Sunny	Mild	High	Weak
Sunny	Cool	Normal	Weak
Rain	Mild	Normal	Weak
Sunny	Mild	Normal	Strong
Overcast	Mild	High	Strong
Overcast	Hot	Normal	Weak
Rain	Mild	High	Strong

Decision Tree : ① ID3 Algorithm ② CART
 ↓ ↓
 (Sklarom)

Classification and Regression Tree

Algorithm

9 Yes / 5 No



① Purity → [Pure Split ??]

② How the features are selected
 ⇒ Information Gain?

Entropy
 Gini Index

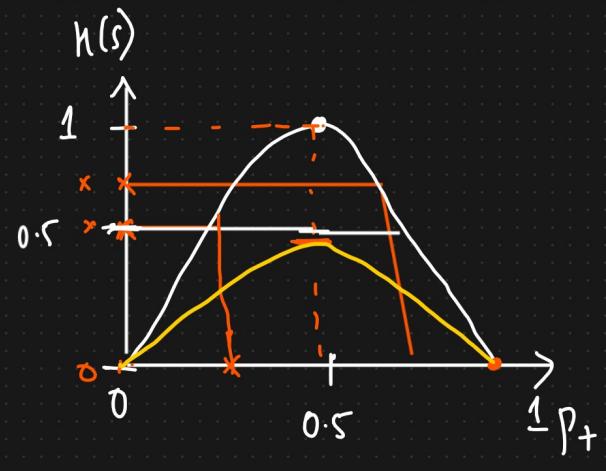
① Entropy {Multiclassification}
 {Binary Classification}

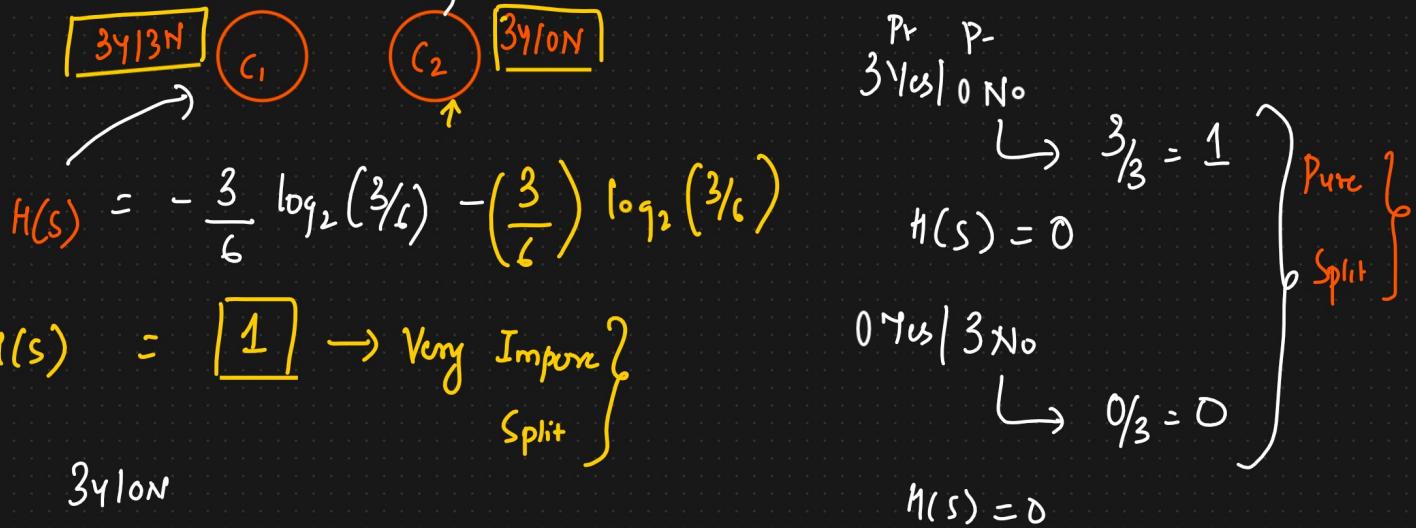
② Gini Index

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

$$H(S) = -P_{C_1} \log_2 P_{C_1} - P_{C_2} \log_2 C_2 - P_{C_3} \log_2 C_3$$

f_1
 64 / 3 N
 Leaf Node
 50% - 50% ↑

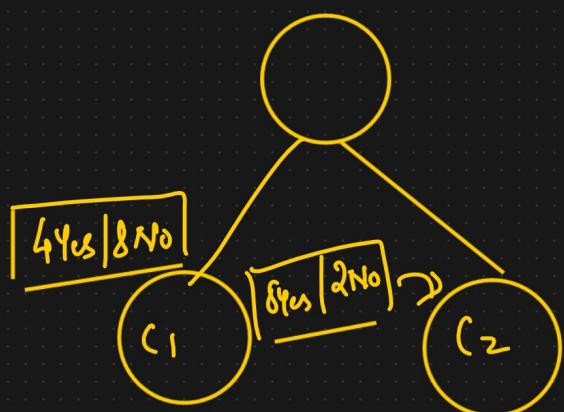




$$\begin{aligned}
 H(S) &= -\frac{3}{3} \log_2 \left(\frac{3}{3}\right) - \frac{0}{3} \log_2 \left(\frac{0}{3}\right) \\
 &= 0 - 0 \\
 &= [0] \rightarrow \text{Pure Split.}
 \end{aligned}$$

$[2\text{Yes}|3\text{No}] \rightarrow \text{Impure Split.}$

$$\begin{aligned}
 H(S) &= -\frac{2}{5} \log_2 \left(\frac{2}{5}\right) - \left(\frac{3}{5}\right) \log_2 \left(\frac{3}{5}\right) \\
 &= 0.97
 \end{aligned}$$



Gini Index

$$G.I. = 1 - \sum_{i=1}^n (p_i)^2 = 1 - \left[(p_+)^2 + (p_-)^2 \right]$$

$$\begin{aligned}
 \frac{3/6}{34/8N} &= 0.5 \\
 \frac{34/0N}{34/8N} &\Rightarrow [0.5]
 \end{aligned}$$

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

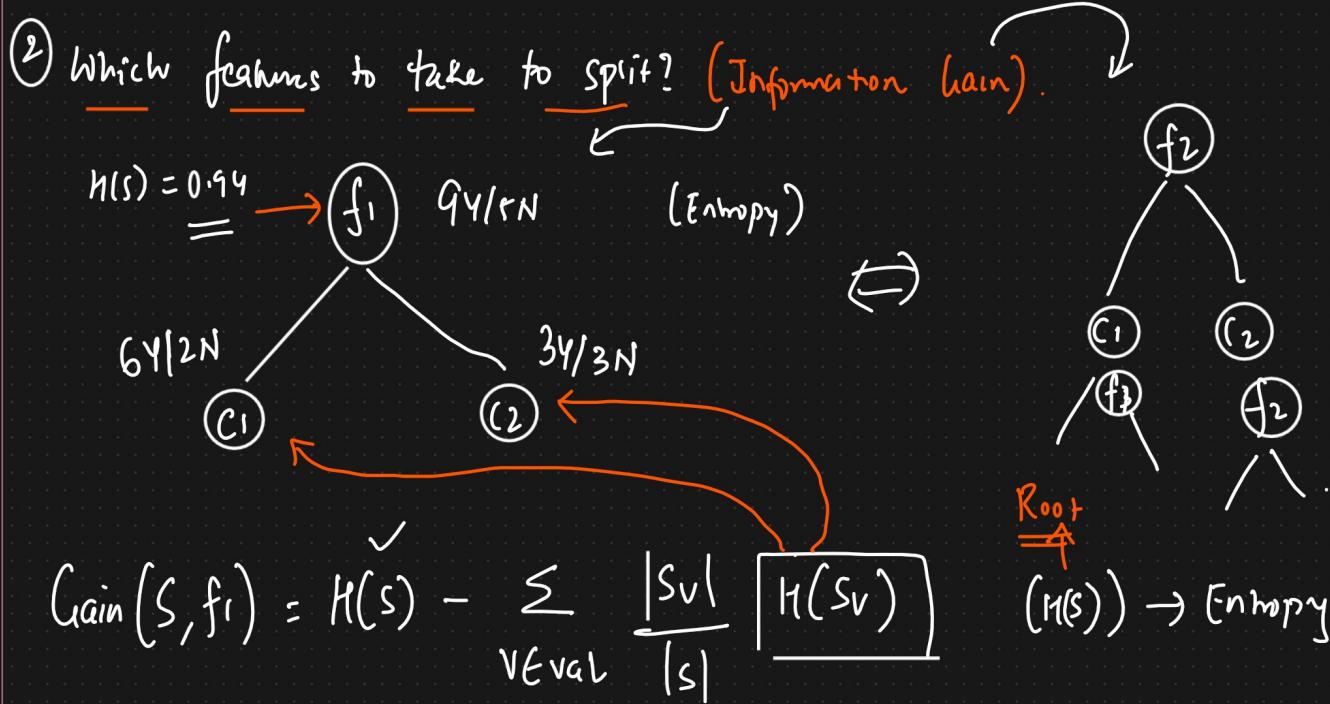
$$G.I. = 1 - \frac{1}{2} = 0.5 //$$

$$\left[\frac{4Y_{14}}{8N_0} \right] = 1 - \sum_{i=1}^n (p_i)^2 = 1 - \left[\left(\frac{4}{12} \right)^2 + \left(\frac{8}{12} \right)^2 \right] \\ = 1 - \left[\frac{1}{9} + \frac{4}{9} \right]$$

Gini Index

$$= 1 - \frac{5}{9} = \boxed{\frac{4}{9}} \approx 0.444$$

$$\left[\frac{8Y_{14}}{2N_0} \right] = 1 - \left[\left(\frac{8}{10} \right)^2 + \left(\frac{2}{10} \right)^2 \right] \\ = 1 - \left[\frac{16}{25} + \frac{1}{25} \right] = 1 - \frac{17}{25} = \frac{8}{25} \approx 0.32$$



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

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

$\sqrt{64/2N}$

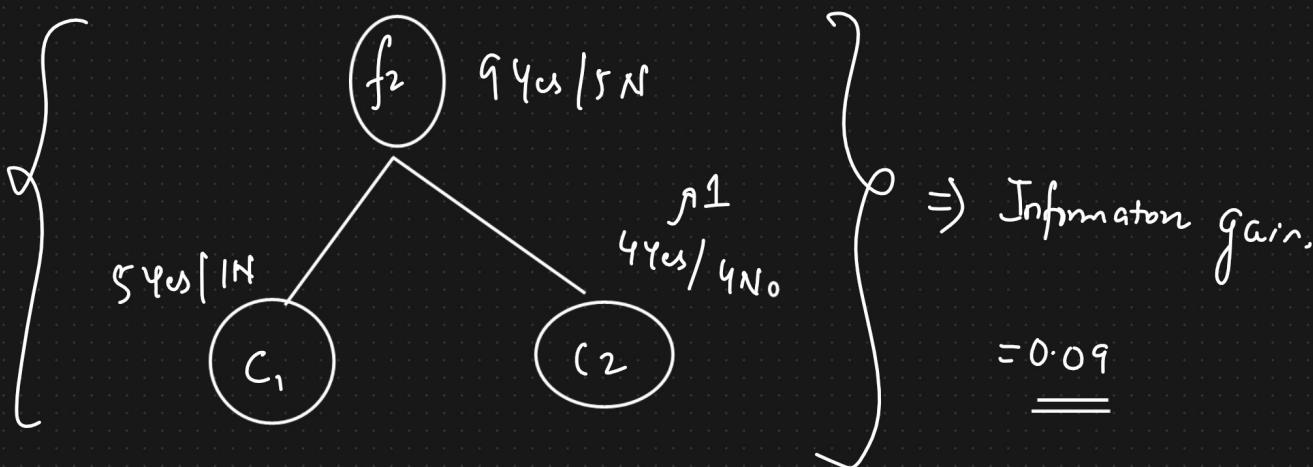
$$H(C_1) = -\frac{6}{8} \log_2 \left(\frac{6}{8} \right) - \frac{2}{8} \log_2 \left(\frac{2}{8} \right) = 0$$

$$H(C_2) = 1$$

↗

$$(34/3N) \quad \text{Gain } (S, f_1) = 0.94 - \left[\frac{8}{14} (0.81) + \frac{6}{14} \times 1 \right]$$

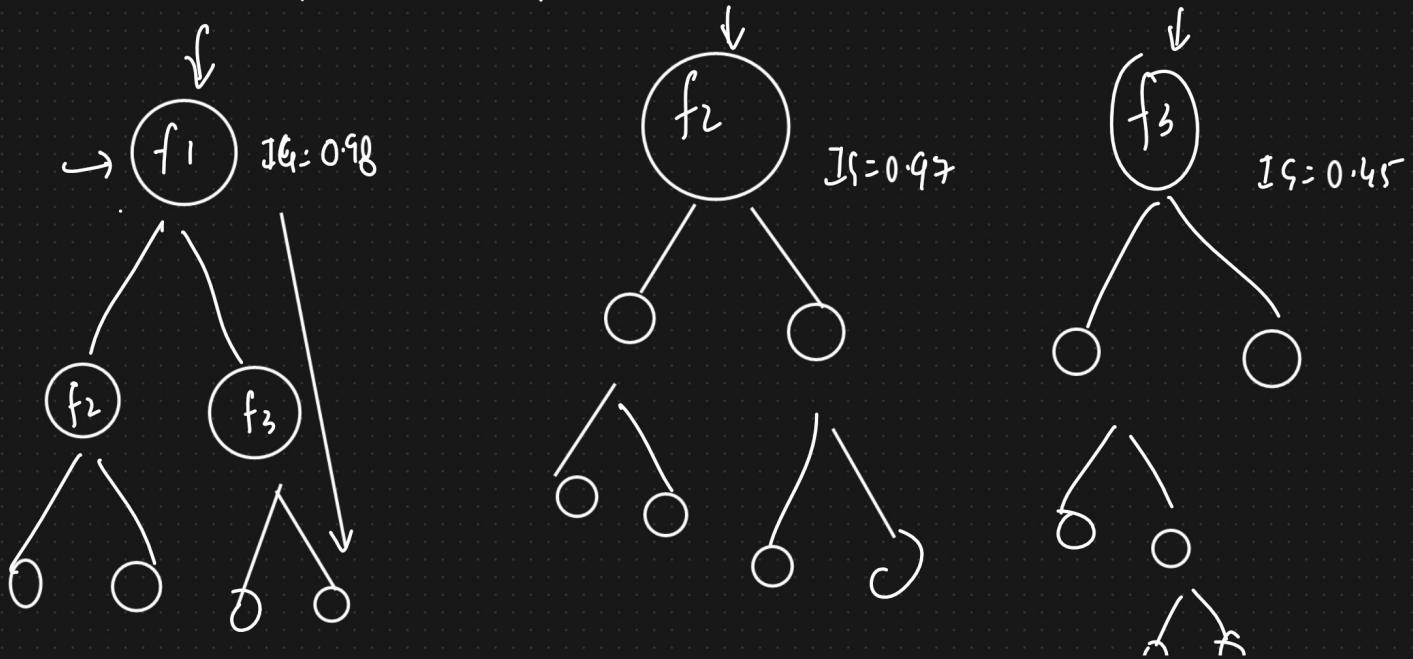
$$\frac{|SV|}{|T|} = \frac{6}{14} = 0.94 - \left[0.462 + 0.42 \right] = 0.049$$

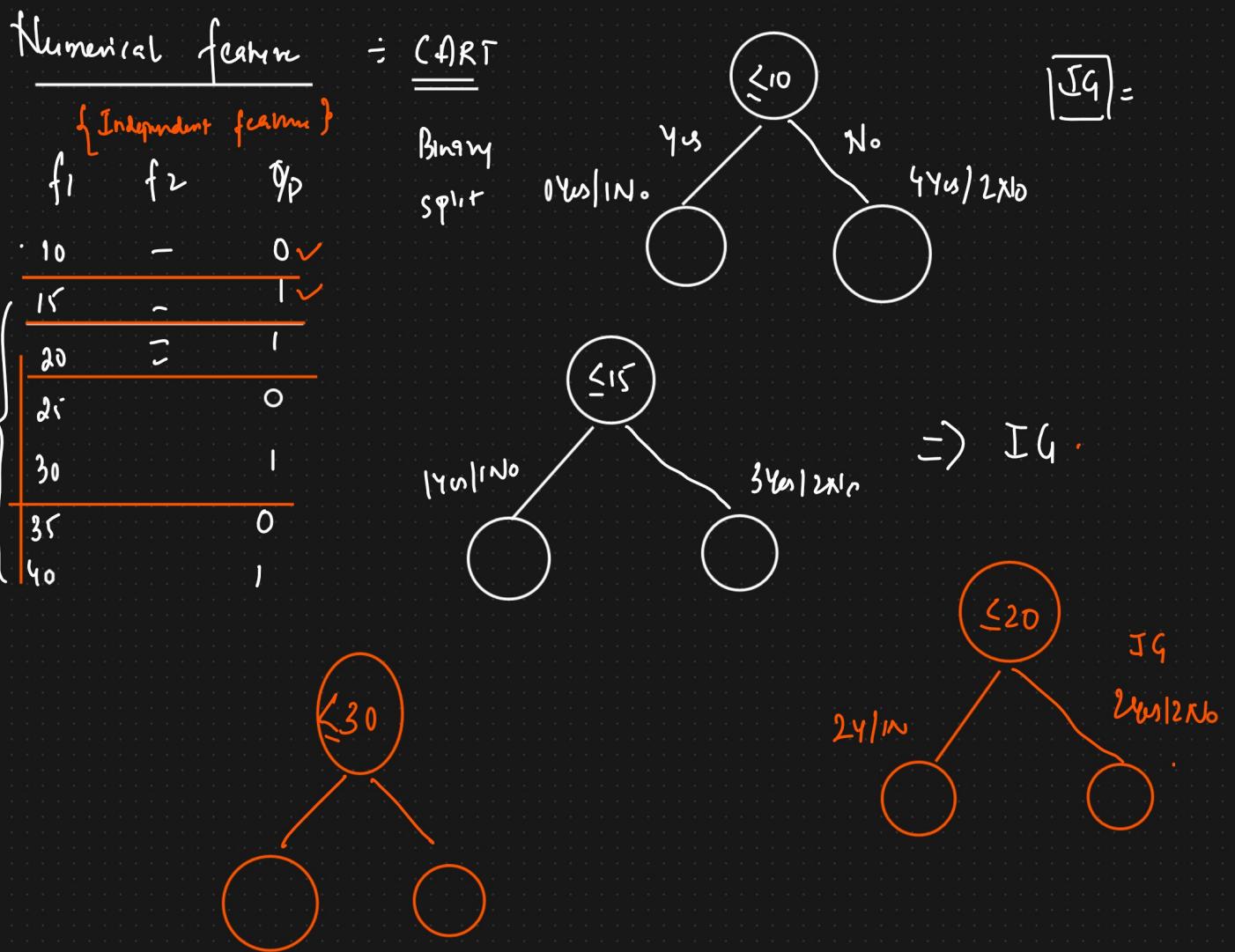


$$\text{Gain} = 0.94 - (6/14) * 0.65 - (8/14) * 1 = 0.09$$

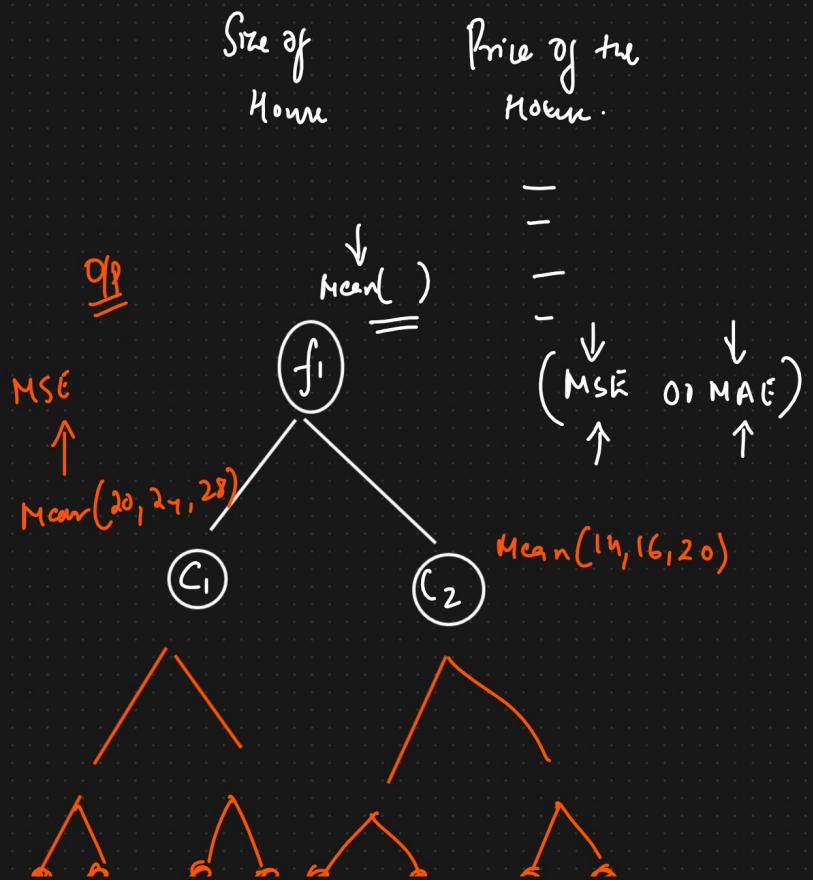
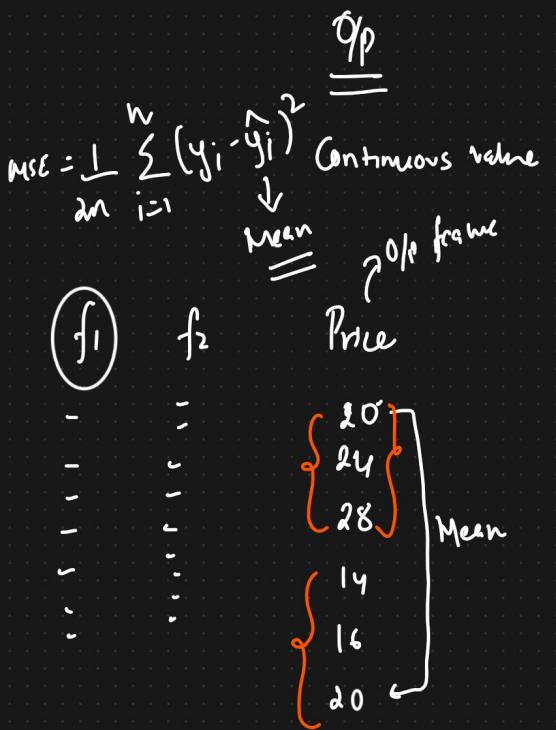
$$IG(f_2) > IG(f_1)$$

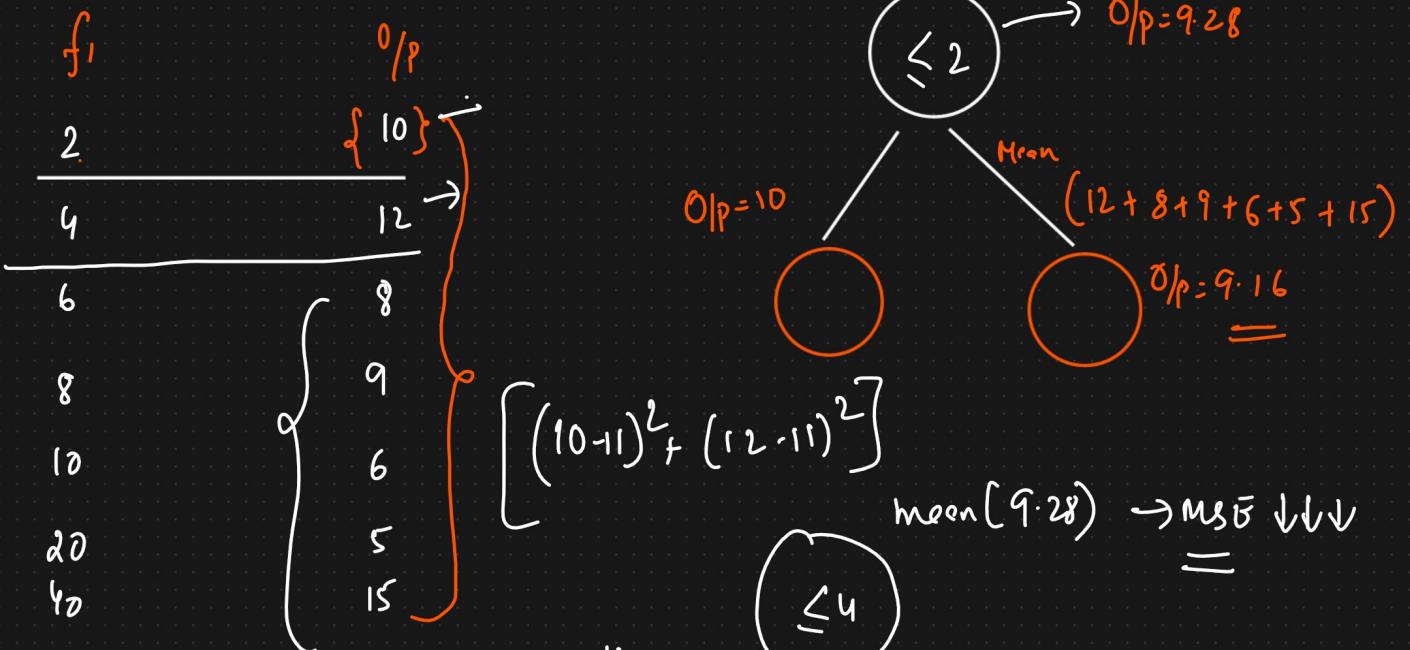
feature 2





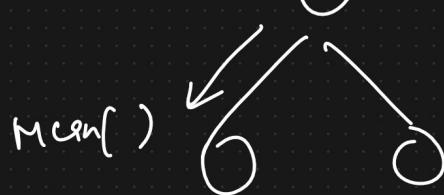
Decision Tree Regressor





Mean Square Error

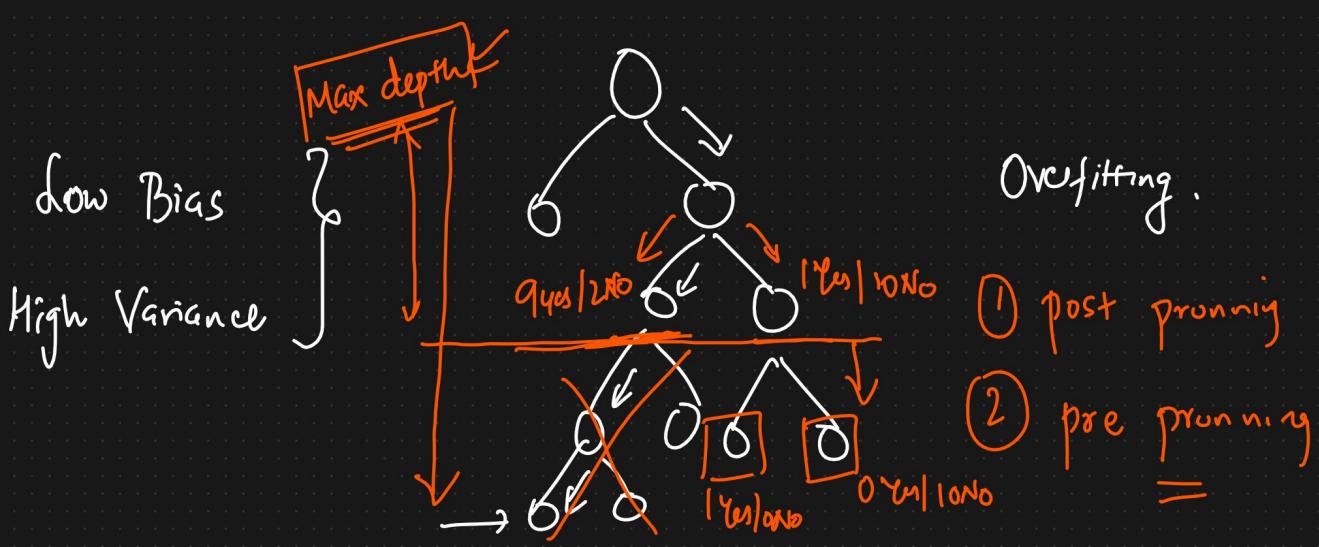
Mean Absolute Error



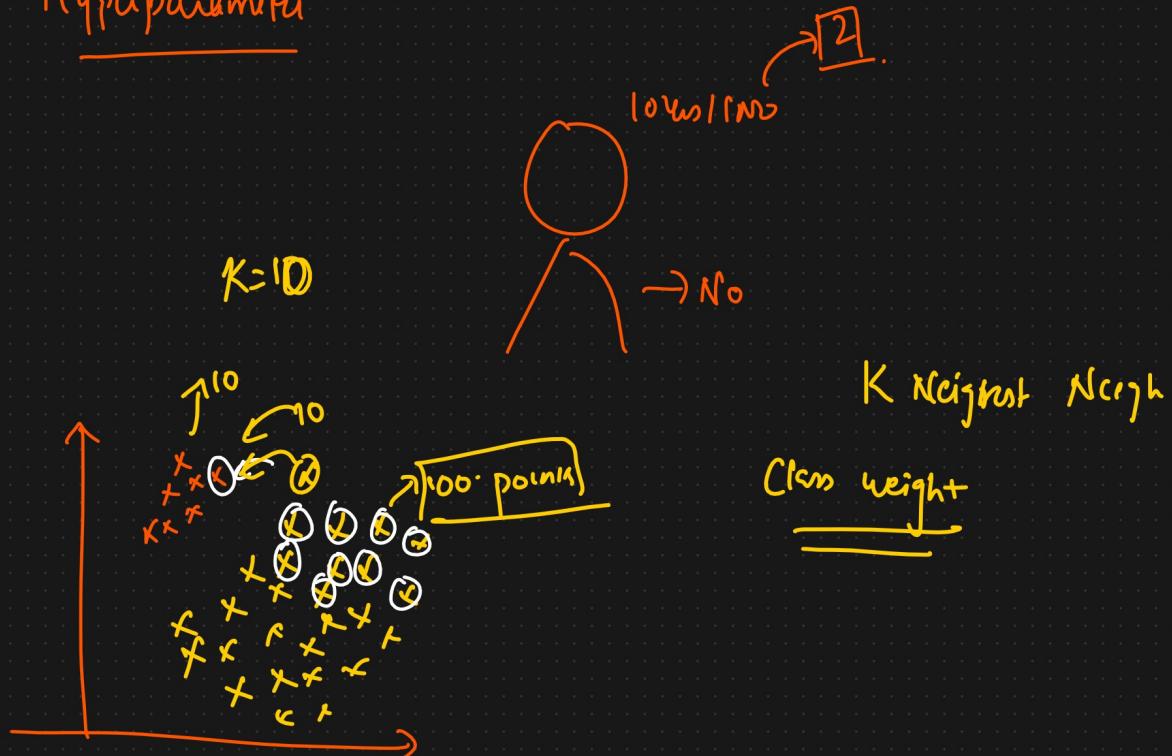
Decision Tree

Overfitting

Training Accuracy \uparrow $\left\{ \begin{array}{l} \text{Low bias} \\ \text{High Variance} \end{array} \right\} \leftarrow$
 Test Accuracy \downarrow $\left\{ \begin{array}{l} \text{High variance} \\ \text{High bias} \end{array} \right\} \uparrow$



Hypoparameter



$f_1 \quad f_2 \quad f_3 \quad O/P \quad \xrightarrow{\text{Regression}} \quad \{ \quad \}$

- - - - -

{ what is the approach to fill these NAN values? }

