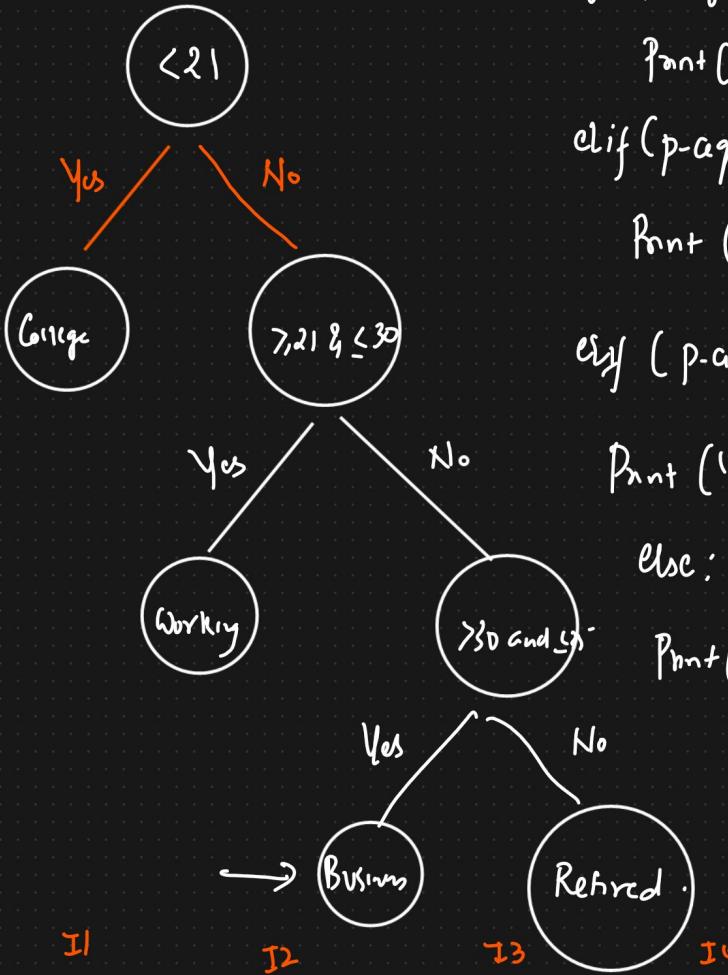


# Decision Tree Classifier And Decision Tree Regressor

- ① Classification
- ② Regression Problem

age = 34



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

Nested if-else clause

if (P-age < 21):

Print ("Age should be in college")

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

Print ("Age should be working")

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

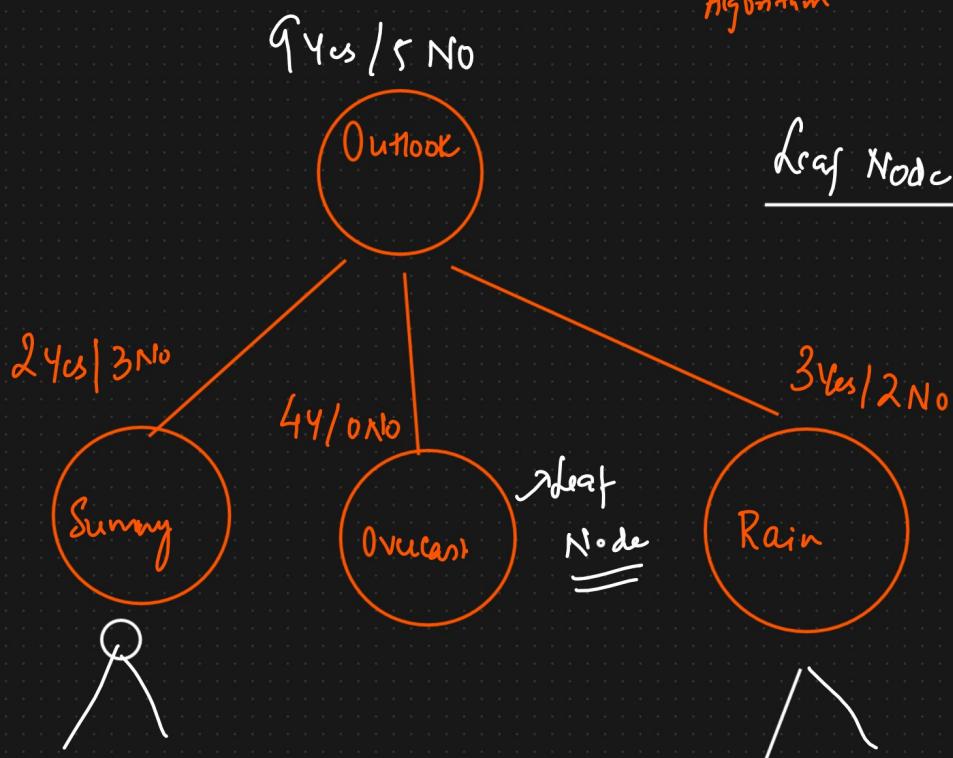
Print ("Age was a business")

else:

Print ("Age has retired")

Decision Tree : ① ID3 Algorithm      ② CART (Sklarn)

Classification and Regression Tree  
Algorithm



① Purity  $\rightarrow$  Pure Split ??

Entropy  
Gini Index

② How the features are selected

$\Rightarrow$  Information Gain?

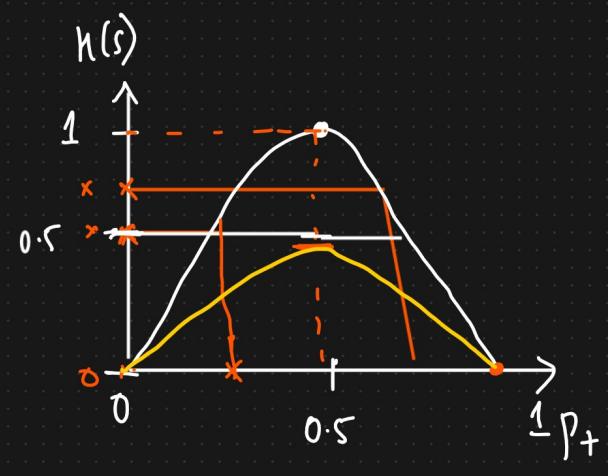


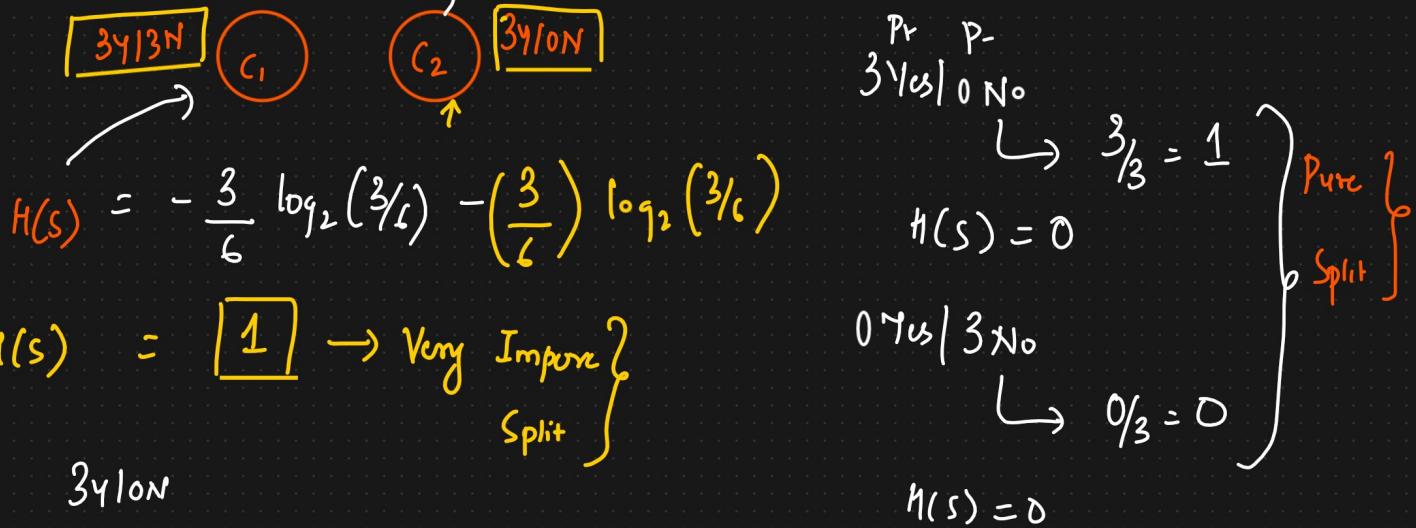
① Entropy {Multiclassification}  
② Gini Index {Binary Classification}

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

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

$f_1$   
64 / 3 N  
leaf Node  
50% - 50% ↑





$$H(S) = -\frac{2}{3} \log_2 \left(\frac{2}{3}\right) - \frac{1}{3} \log_2 \left(\frac{1}{3}\right)$$

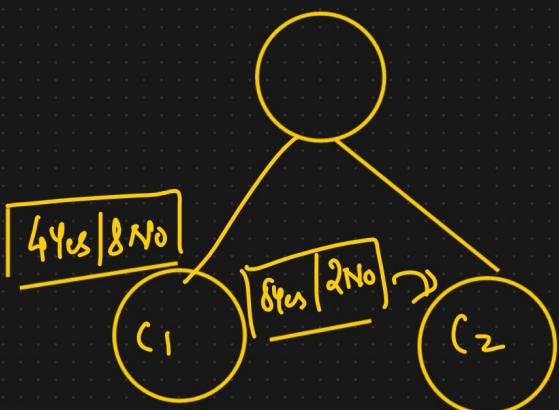
$$= 0 - 0$$

$$= \boxed{0} \rightarrow \text{Pure Split.}$$

$\boxed{2Y|3N}$  → Impure Split.

$$H(S) = -\frac{2}{5} \log_2 \left(\frac{2}{5}\right) - \frac{3}{5} \log_2 \left(\frac{3}{5}\right)$$

$$= 0.97$$



### Gini Index

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

$$\frac{3}{6} = 0.5 \quad = 1 - \left[ \left(\frac{1}{2}\right)^2 + \left(\frac{1}{2}\right)^2 \right]$$

$$\frac{3Y|3N}{\underline{\underline{0.5}}} \Rightarrow \boxed{0.5}$$

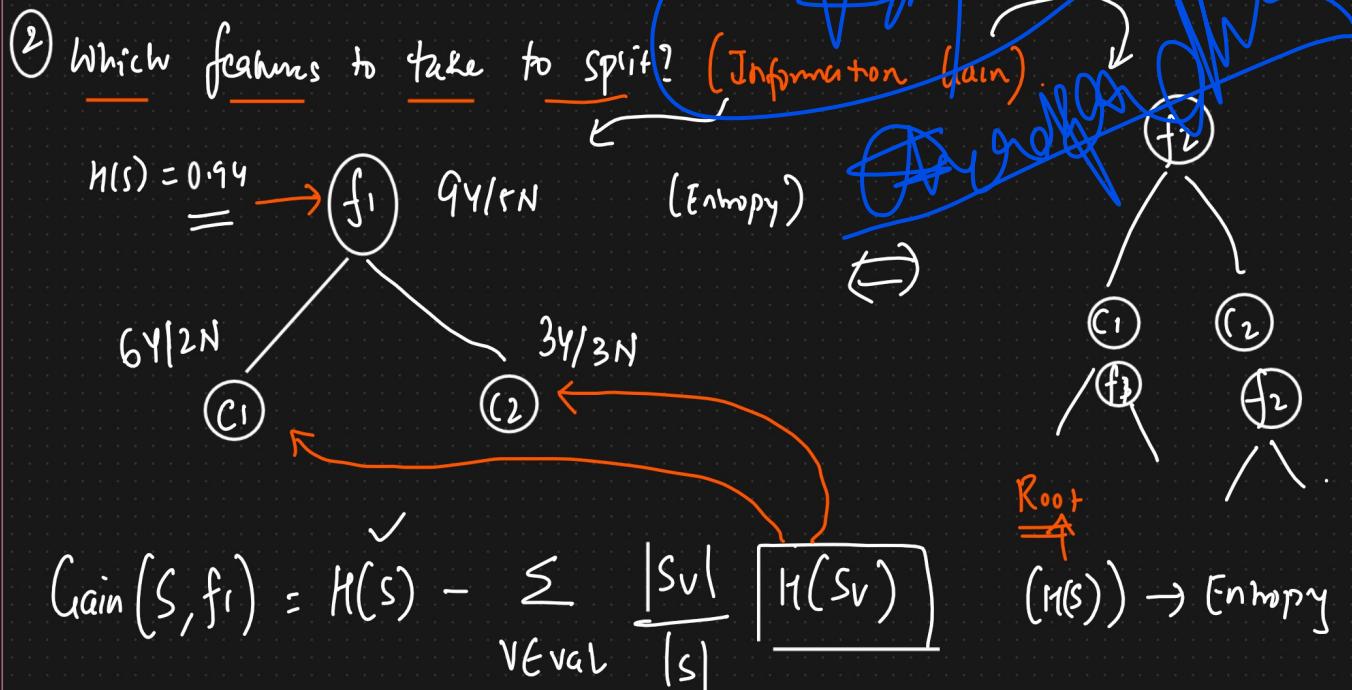
$$G.I. = 1 - \frac{1}{2} = 0.5 //$$

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

Gini Index

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

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



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

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

$\sqrt{6Y/2N}$

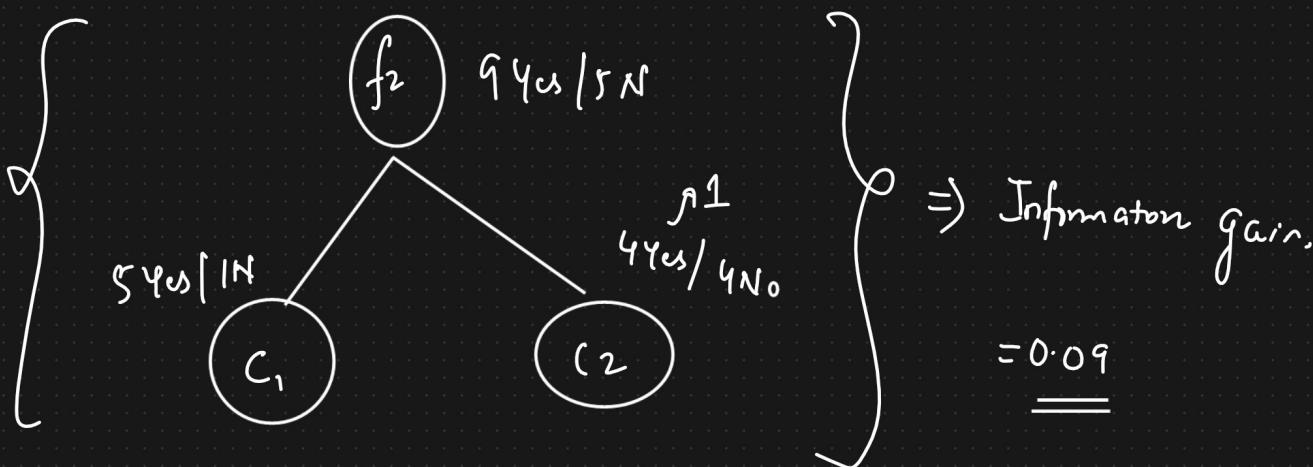
$$H(C_1) = -\frac{6}{8} \log_2 \left( \frac{6}{8} \right) - \frac{2}{8} \log_2 \left( \frac{2}{8} \right) = 0$$

$$H(C_2) = 1$$

↗

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

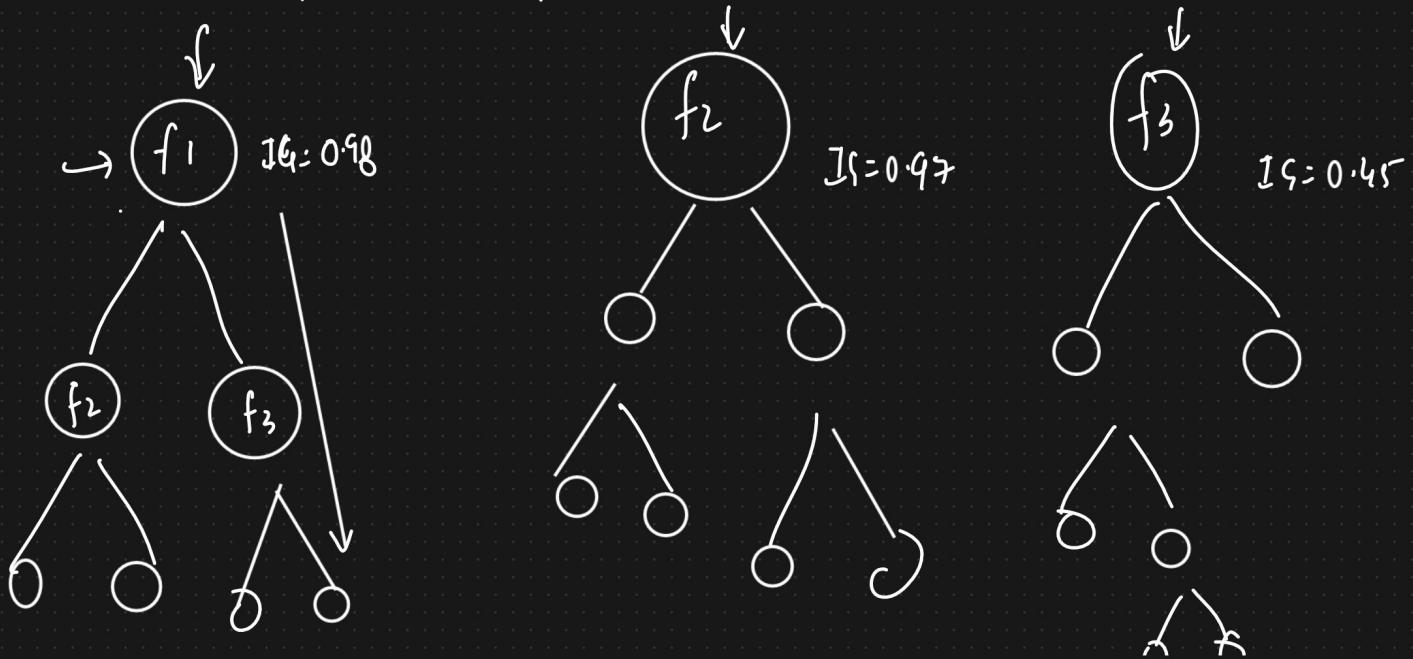
$$\frac{|SV|}{|T|} = \frac{6}{14} = 0.94 - \left[ 0.462 + 0.42 \right] = 0.049$$

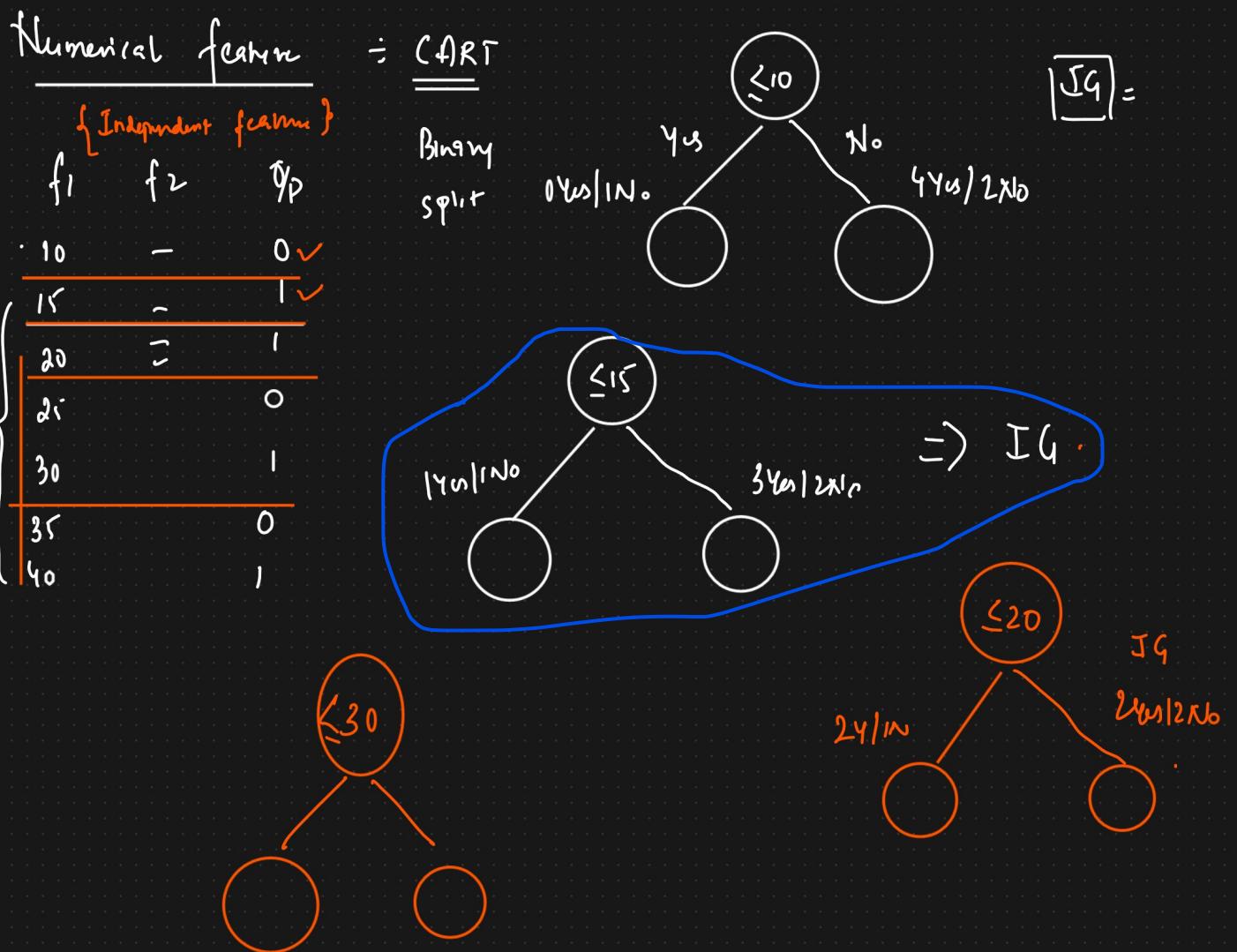


$$\text{Gain} = 0.94 - (6/14) * 0.65 - (8/14) * 1 = 0.09$$

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

feature 2





## Decision Tree Regressor

$$y_p = \frac{\sum_{i=1}^n (y_i - \bar{y}_i)^2}{\text{Mean}}$$

Continuous value

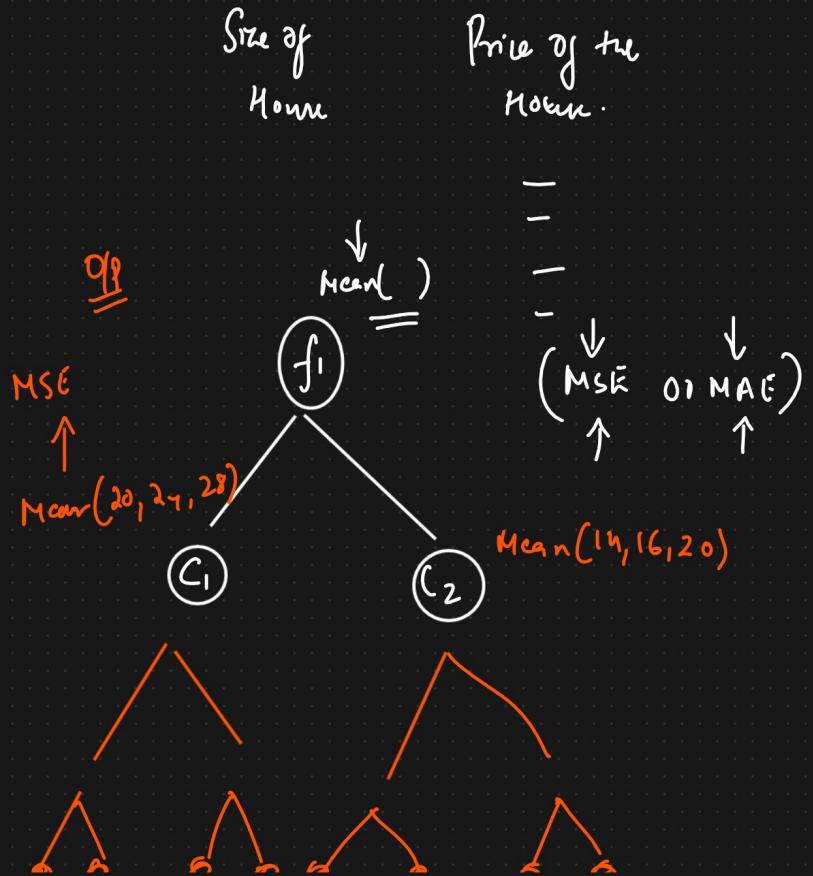
$\bar{y}_i$  = 10, 14, 16, 20, 24, 28

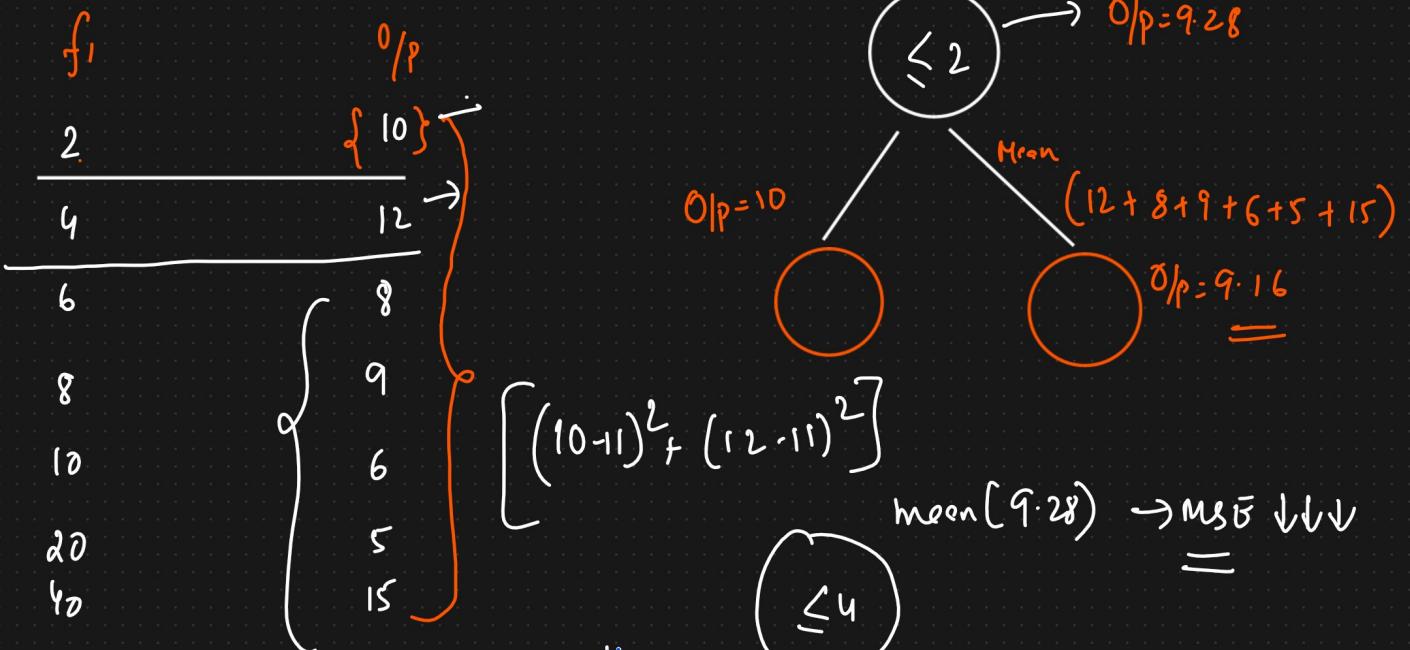
$f_1$  = Size of House

$f_2$  = Price

Mean = {20, 24, 28} = 24

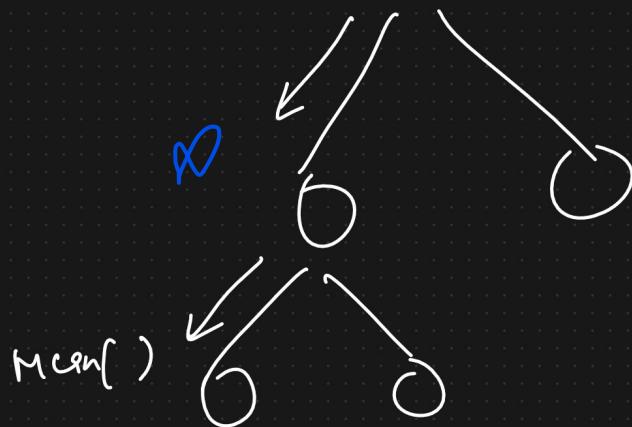
Mean = {14, 16, 20} = 16





Mean Square Error

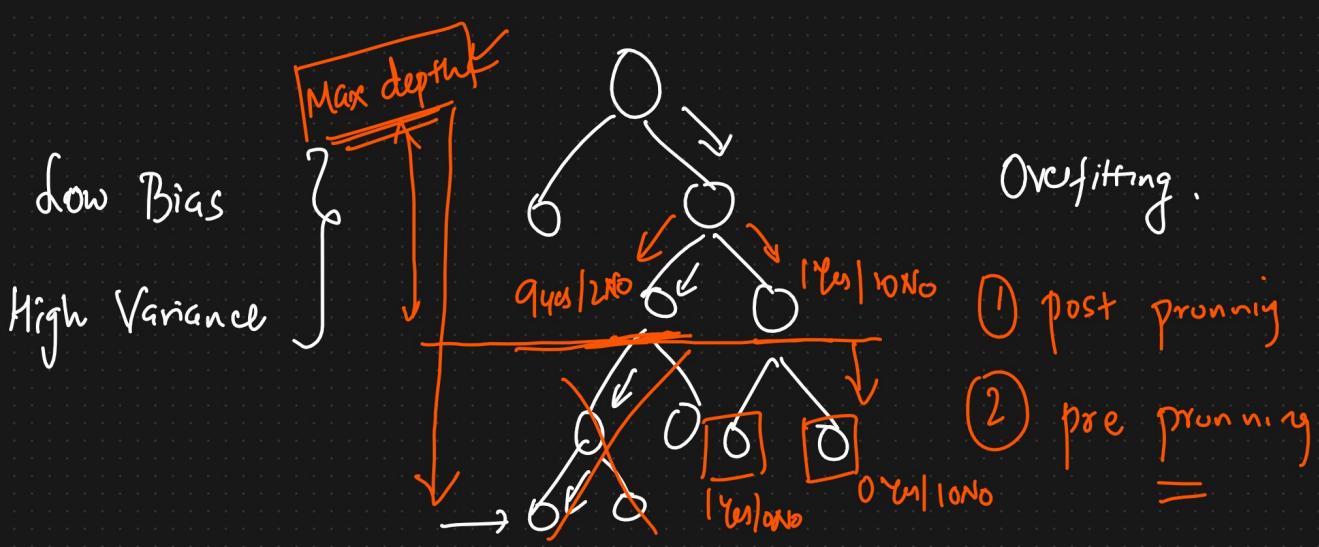
Mean Absolute Error



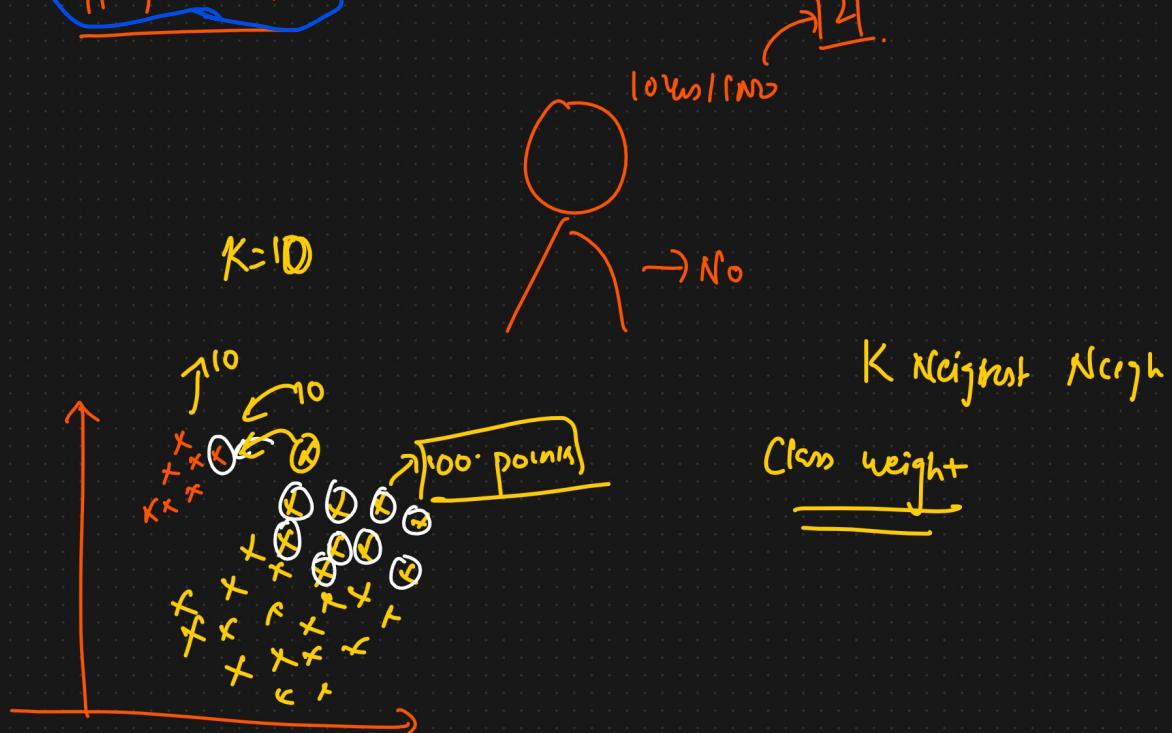
Decision Tree

Overfitting

Training Accuracy ↑  
 Test Accuracy Low  
 $\begin{cases} \text{Low bias} \\ \text{High variance} \end{cases}$



## Hypoparameter



$f_1$      $f_2$      $f_3$     O/P    Regression     $\{ \quad \leftarrow \quad \}$   
 -       -       -       -       { what is the approach to fill these  
 -       -       -       -       NAN values? }.

