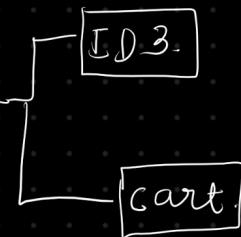


## \* Decision tree classifier

\* Decision Tree classifier



a) Entropy and gini index → purity split.

b) Information gain → features to select for  
DT construction.

age = 14

if (age ≤ 15)

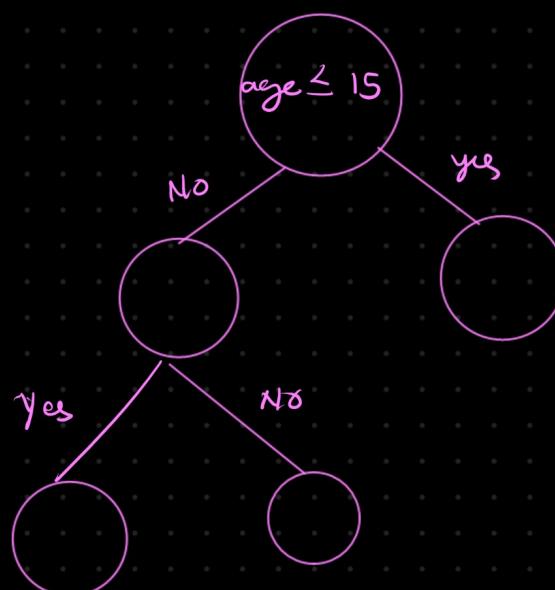
    print ("School")

elif (age < 15 and age ≥ 21):

    print ("The person may be college")

else :

    print ("The person has passed")



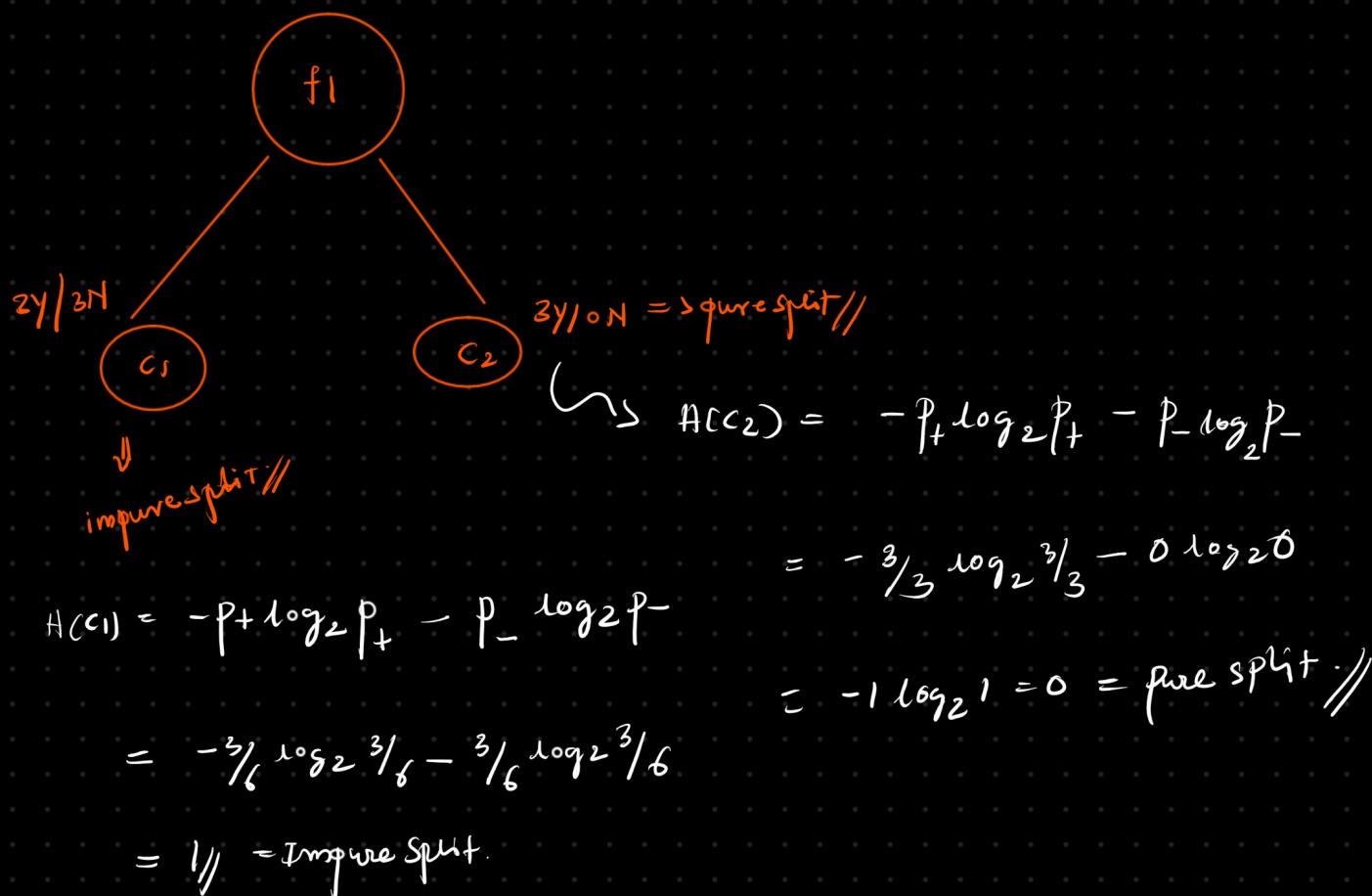
# 1) Entropy

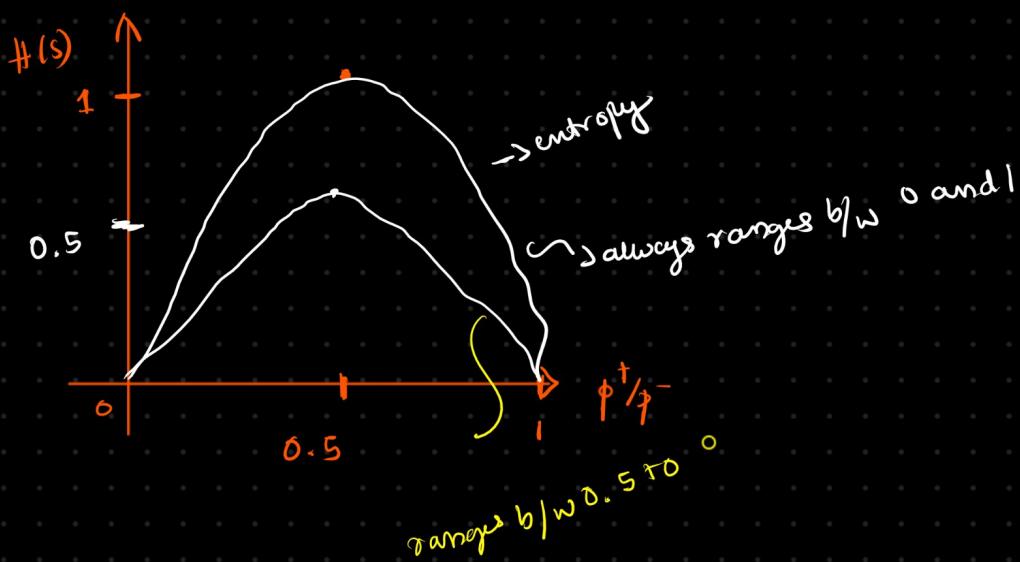
$$H(S) = -P_+ \log_2 P_+ - P_- \log_2 P_-$$

# 2) Gini Impurity

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

{ binary classification }  
 one category -  $P_-$  = negative class  
 another category  $+P_+$  = positive class





## Gini index

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

$$= 1 - ((p_1)^2 + (p_2)^2)$$

$$= 1 - ((\frac{1}{2})^2 + (\frac{1}{2})^2)$$

$$= 0.5 \Rightarrow \text{impure}$$

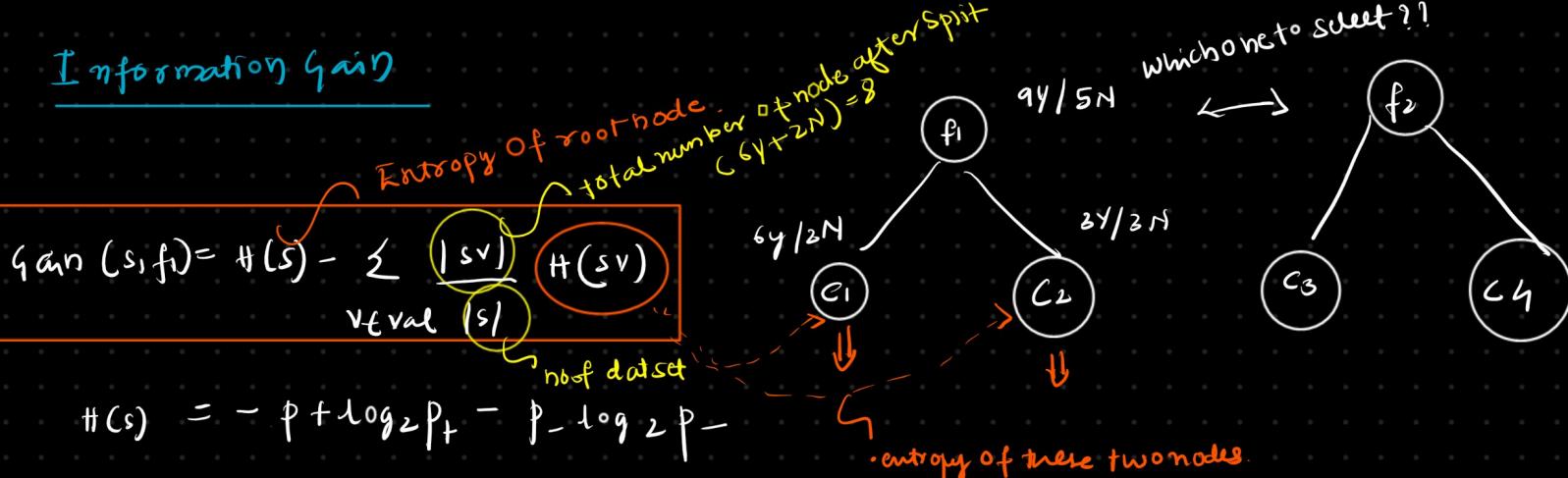
3Y/0N

$$= 1 - \left( \left( \frac{3}{5} \right)^2 \right)$$

$$= 1 - 1$$

$$= 0 // \Rightarrow \text{pure.}$$

## Information Gain



$$H(S) = -p_f \log_2 p_f - p_c \log_2 p_c$$

$$= -\frac{9}{14} \log_2 \frac{9}{14} - \frac{5}{14} \log_2 \frac{5}{14}$$

$$\approx 0.9411$$

$$\text{Gain}(S, f_1) = 0.9411 - \left[ \frac{8}{14} * 0.81 + \frac{6}{14} * 1 \right]$$

$$IG(S|f_1) = 0.4911 \quad \text{with respect to } f_1$$

$$IG(S|f_2) = 0.51 \quad \text{with respect to } f_2 \quad \text{it should start from } f_2 \text{ as root node}$$

When should we use Entropy and GINI index

$$\rightarrow H(s) = -p_+ \log_2 p_+ - p_- \log_2 p_-$$

↓

$$H(s) = -p_{c_1} \log_2 p_{c_1} - p_{c_2} \log_2 p_{c_2} - p_{c_3} \log_2 p_{c_3}$$

$$\rightarrow \text{GiniIndex} = 1 - \sum_{i=1}^n (p_i)^2 \Rightarrow$$

\* whenever the dataset is "small use entropy"

\* entropy → Small dataset

\* Gini index → Large dataset } when we use decision tree classifier

Decision tree with numerical features

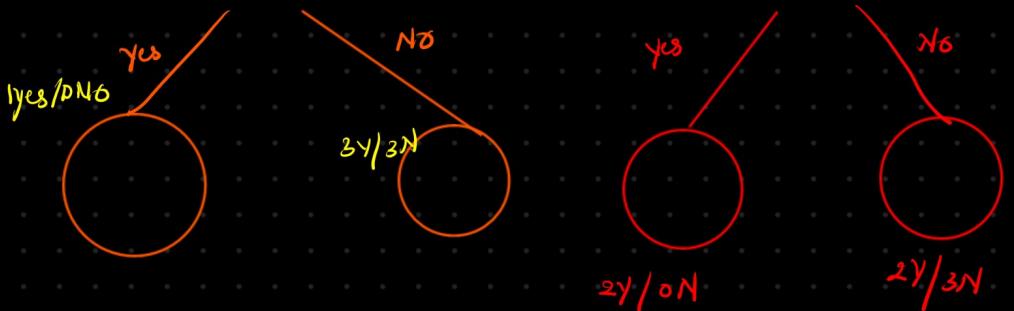
f1	%/P
Yes	
Yes	
No	

Step - sort the values

\* threshold



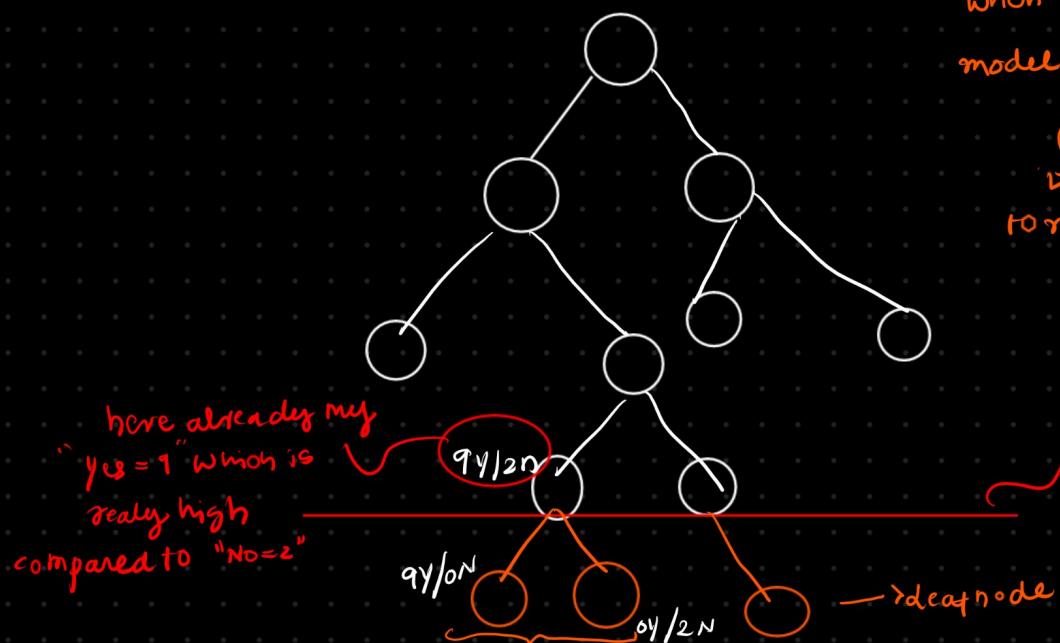
5.2	No
6.7	Yes
8.9	No
10.5	Yes



\* there will be many splits like this which split will have larger information gain then we will consider that.

## Post pruning & Pre pruning

### Training dataset



- when we train the model until the leaf node our model will be "overfitting"



To reduce this we use

- post-pruning.
- pre-pruning.

→ post pruning

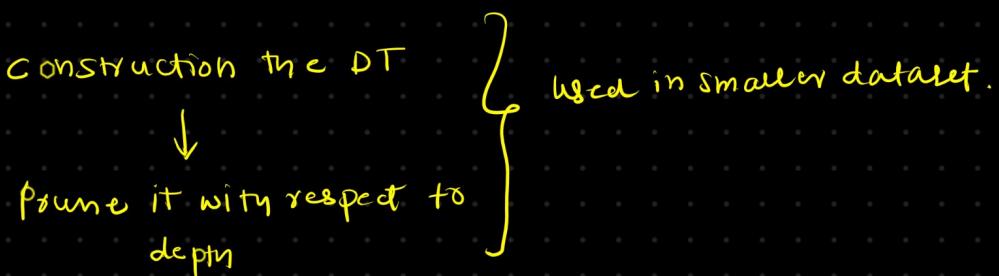
↳ construct DT

↳ prune it with respect to the

↳ for small length

- S
- even though i did the split these type of split will cause overfitting

## ① Post pruning



## ② Pre - pruning : tune with hyperparameters while constructing the decision tree

- max\_depth = 3
- max\_features
- split\_ratio
- criterion
- splitter.

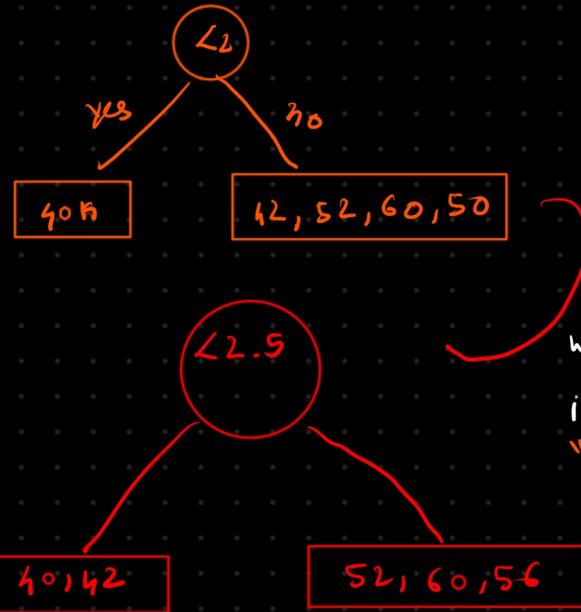
## Decision Tree Regression

Dataset

→ Exp

	gap	salary
threshold 2	yes	40K
2.5	yes	42K
3	no	52K
4	no	60K
4.5	yes	50K

continuous value



between this  
which feature should  
i select. so we use  
"variance reduction"

$$\textcircled{1} \text{ Variance} = \frac{1}{n} \sum_{i=1}^n (y_i - \bar{y})^2 \quad [\text{mean square Error}]$$

Average.

$$\text{Variance (Root)} = \frac{1}{5} \left( (40-50)^2 + (42-50)^2 + (52-50)^2 + (60-50)^2 + (56-50)^2 \right)$$

$[40K, 42K, 52K, 60K, 56K]$      $[40K, 42K, 52K, 60K, 56K]$

$$= \frac{1}{5} [100 + 64 + 4 + 100 + 36]$$

Variance reduction 6

$$= 60.8 //$$

Left side of tree:  $\sqrt{C_1} = \sqrt{100} = 10$

Right side of tree:  $\sqrt{C_2} = \sqrt{51} = 7.1$

Bottom left leaf:  $\sqrt{C_1} = \sqrt{82} = 9.1$

Bottom right leaf:  $\sqrt{C_2} = \sqrt{46.66} = 6.8$

$$\text{Variance}(c_1) = \frac{1}{4} [(40 - 50)^2] \\ = 100 //$$

$$\text{Variance}(c_2) = \frac{1}{5} \left[ (42 - 50)^2 + (52 - 50)^2 + (50 - 50)^2 + (56 - 50)^2 \right] \\ = 51 //$$

$\approx 0.04 //$

Variance reduction

$$= \text{Var}(\text{root}) - \underbrace{(\omega_i)}_{\substack{\text{ratio of that particular child node to the total} \\ \text{no.}}} \text{Var}(\text{child})$$

$$= 60.8 - \left[ \frac{1}{8} \times 100 + \frac{4}{5} \times 51 \right]$$

$$= 60.8 \left[ 20 + \frac{4}{5} \times 51 \right]$$

$$= 60.8 - 20 - 40.8$$

Variance reduced = 0

$\downarrow$   
left split

Variance reduction = 0.04

$\downarrow$   
right split