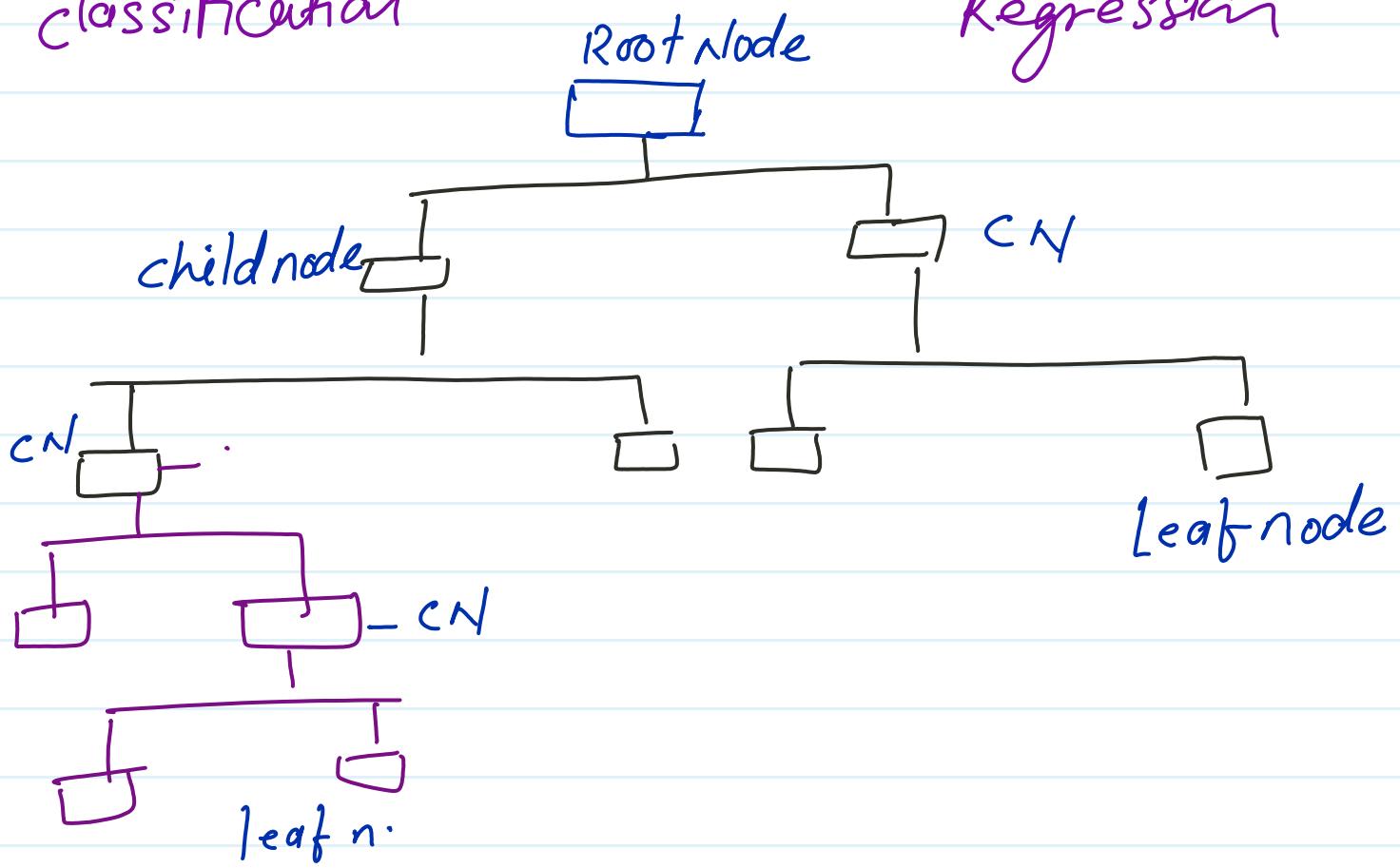


Decision tree

classification

Regression



method of solve decision Tree
algorithm -

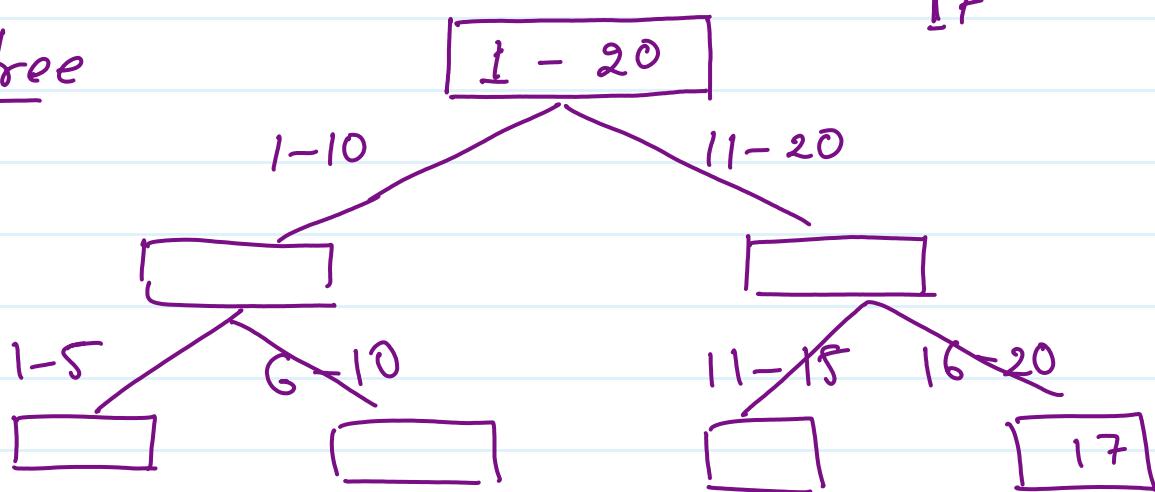
- ✓ ① ID3
- ② Cart

Gini Index

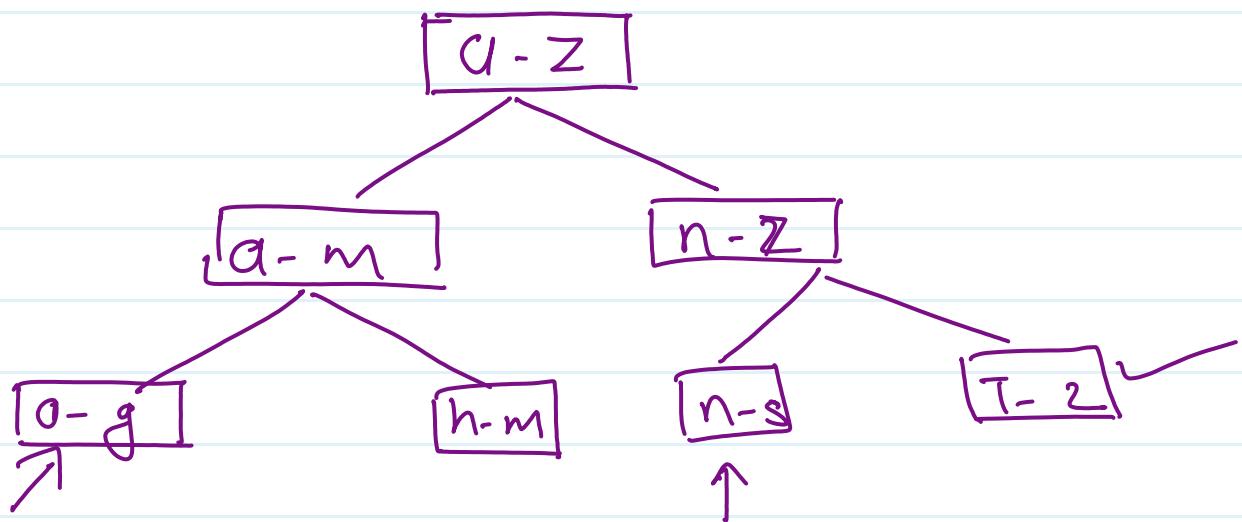
Entropy
Information Gain

Dataset [5, 6, 3, 11, 10, 2, 3, 7, 8]

Binary tree



[Red, yellow, purple, blue, white, green, black]



Entropy OR Gini Index

Information Gain



In DT we do not required to label categorical data into numerical

$x_1 \ x_2 \ x_3$

* Entropy and Gini Index →
purity split in dataset

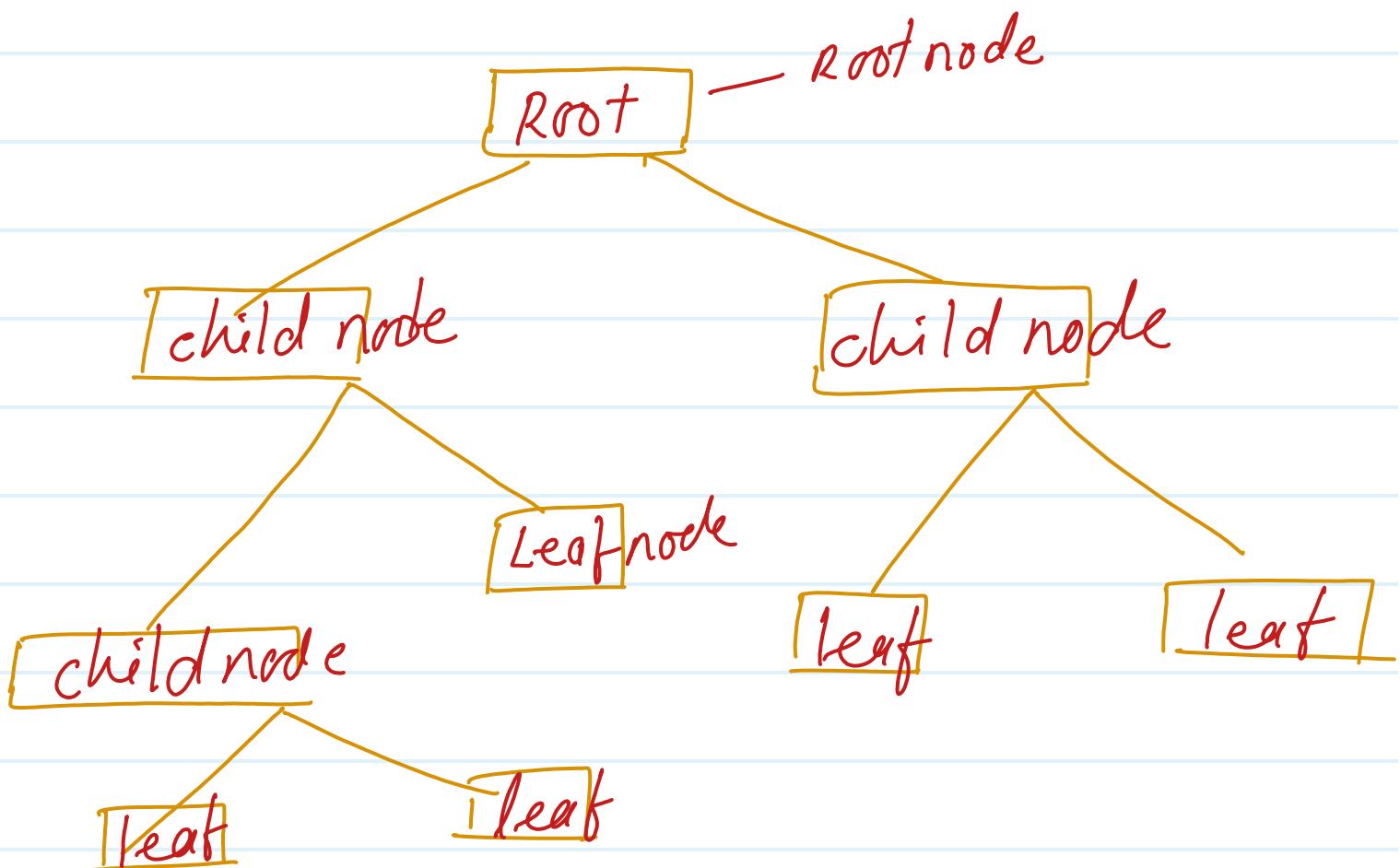
* Information Gain → DT feature

<u>class.</u>	weight	height	O/P	split
	60	160	ob	obese/nooben.
	70	170	no	
	80	180	ob	
	90	190	no	
	100	200	no	

⑪ DT Regressor

Regression we use standard deviation / MSE / MAE

weight	height	BMI
60	160	21
70	170	22
80	180	20
85	190	23
90	195	24



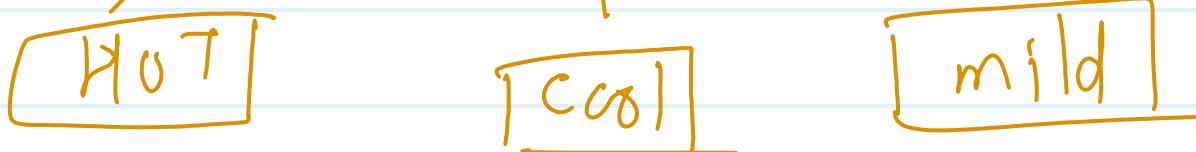
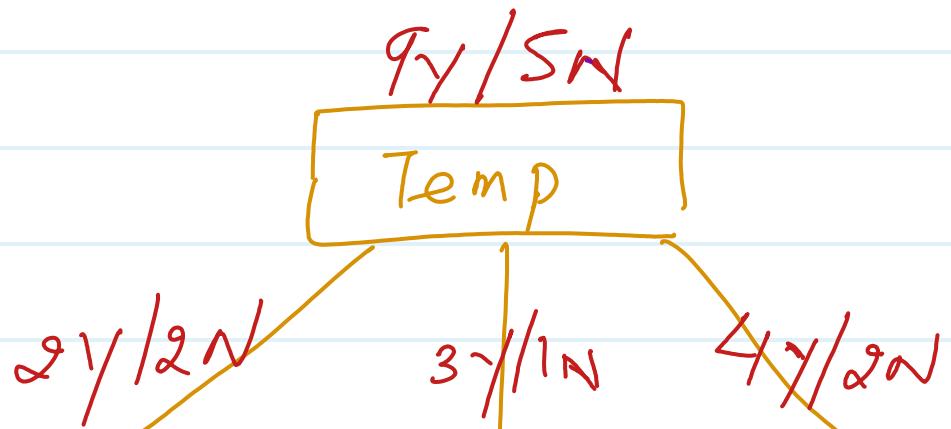
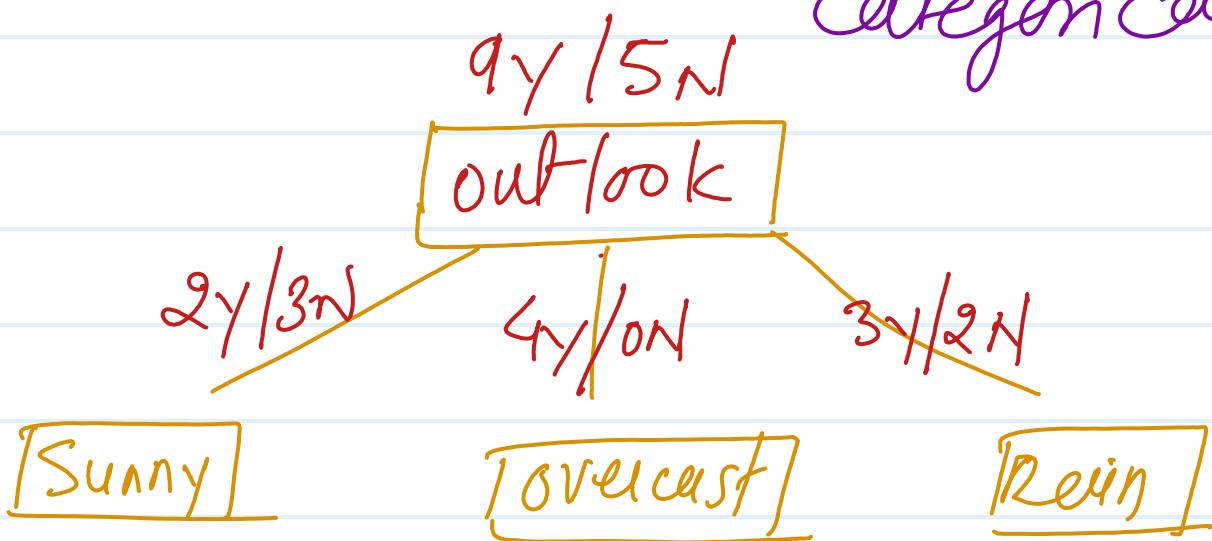
Decision Tree Classifier

Tennis

outlook	Temp	humidity	wind	play
sunny	H	High	weak	N
Sunny	H	H	strong	N
overcast	H	H	W	Y
rain	M	H	W	Y
rain	C	Normal	W	Y
rain	C	N	S	N
overcast	C	N	W	Y
sunny	M	H	W	N
Sunny	C	N	W	Y
rain	M	N	W	Y
sunny	M	N	S	Y
overcast	M	high	S	Y
overcast	H	N	W	Y
rain	M	H	S	N

= feature can be numeric and categorical

= output can be numeric and categorical



① Entropy $H_{(S)}$

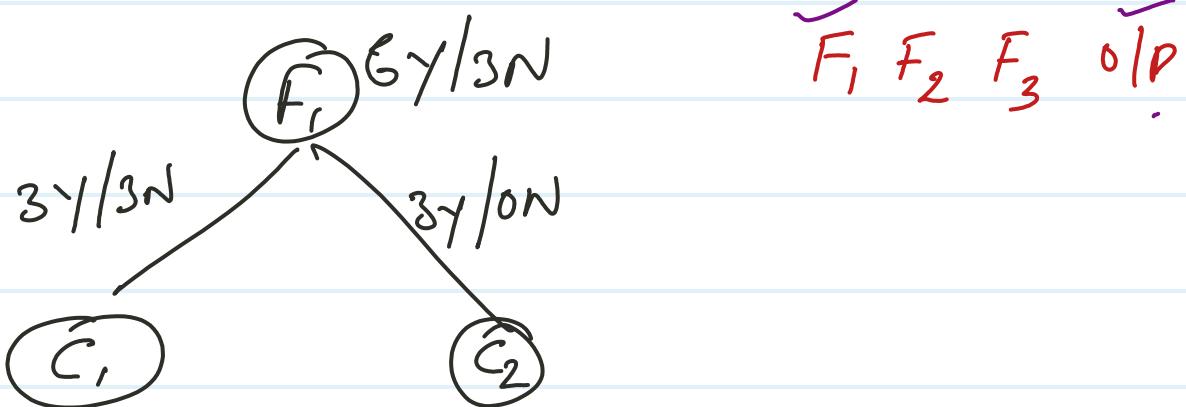
⇒ Formula (Binary class)

$$H_{(S)} = -P_{\text{yes}} \log_2(P_{\text{yes}}) - P_{\text{no}} \log_2(P_{\text{no}})$$

Multiclass

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

Example



$$C_1 \Rightarrow H_{(S)} = -\frac{3}{6} \log\left(\frac{3}{6}\right) - \frac{3}{6} \log\frac{3}{6}$$

\Rightarrow 1 impure split

$$C_2 \Rightarrow H_{(S)} = -\frac{3}{3} \log \frac{3}{3} - \frac{0}{3} \log \frac{0}{3}$$

\Rightarrow 0 pure split

For the pure split of feature

pure entropy should be zero (0)

for impure split = 1

⑤ Gini index (Impurity) -

main formula -

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

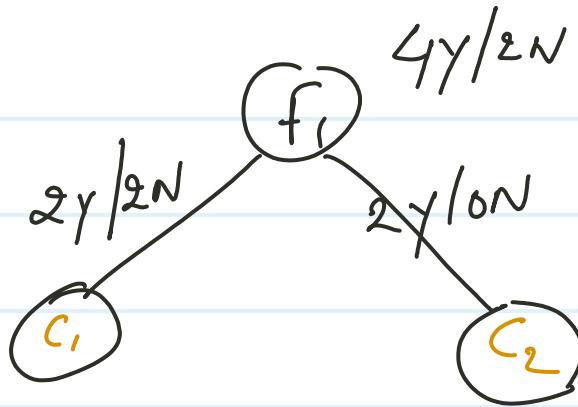
binary class.

$$G.I. = 1 - \sum_{i=1}^n [(P_{C_1})^2 + (P_{C_2})^2]$$

multiclass

$$G.I. = 1 - \sum_{i=1}^n [(P_{C_1})^2 + (P_{C_2})^2 + (P_{C_3})^2 + \dots]$$

Example



$$\underline{C_1} \Rightarrow G.I. = 1 - \left[\left(\frac{2}{4}\right)^2 + \left(\frac{2}{4}\right)^2 \right]$$

$$= 0.5 \quad \checkmark$$

$$C_2 \Rightarrow G.I. = 1 - \left[\left(\frac{0}{2}\right)^2 + \left(\frac{0}{2}\right)^2 \right]$$

$$= 0 \quad \checkmark$$

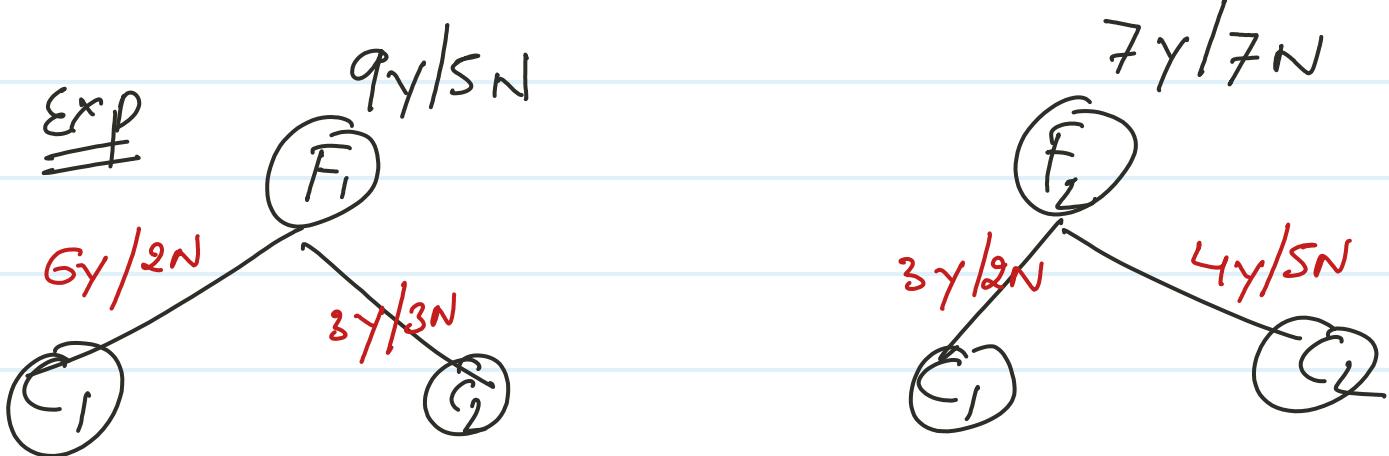
Range of entropy = 0 to 1

Gini impurity (index) = 0 to 0.5

③ Information Gain

formula -

$$\text{gain}(S, f_i) = H(S) - \sum_{i=1}^n \frac{|S_v|}{|S|} H(S_v)$$



\Rightarrow For $F_1 \Rightarrow$

$$H(S) = -\frac{9}{14} \log \frac{9}{14} - \frac{5}{14} \log \frac{5}{14}$$

$$\boxed{H(S) = 0.94}$$

$$C_1 \Rightarrow H(S) = -\frac{6}{8} \log \frac{6}{8} - \frac{2}{8} \log \frac{2}{8}$$

$$\boxed{H(S) = 0.81}$$

$C_2 \Rightarrow$

$$H(s) = -\frac{3}{3} \log \frac{3}{3} - \frac{0}{3} \log \frac{0}{3}$$

$$\boxed{H(s) = 1}$$

gain of f_1 \Rightarrow

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

$$\Leftarrow \text{gain}(s, f_1) = \underline{0.049}$$

$$f_2 \rightarrow H(s) = -\frac{7}{7} \log \frac{7}{7} - \frac{7}{7} \log \frac{7}{7}$$

$$= 0$$

$$C_1 \rightarrow H(s) = -\frac{3}{5} \log \frac{3}{5} - \frac{2}{5} \log \frac{2}{5}$$

$$\begin{aligned} &= 0.133 + 0.159 \\ &= \boxed{0.29} \end{aligned}$$

$$C_2 \Rightarrow H(S) = -\frac{4}{9} \log \frac{4}{9} - \frac{5}{9} \log \frac{5}{9}$$

$= 0.019$

$$f_2 \text{ gain}(S, F_2) = 0 - \left[\frac{5}{14} \times 0.29 + \frac{9}{14} \times 0.019 \right]$$

$$= 0 - [0.10 + 0.009]$$

$= -0.10$

F_1

0.49

F_2

0.56

F_3

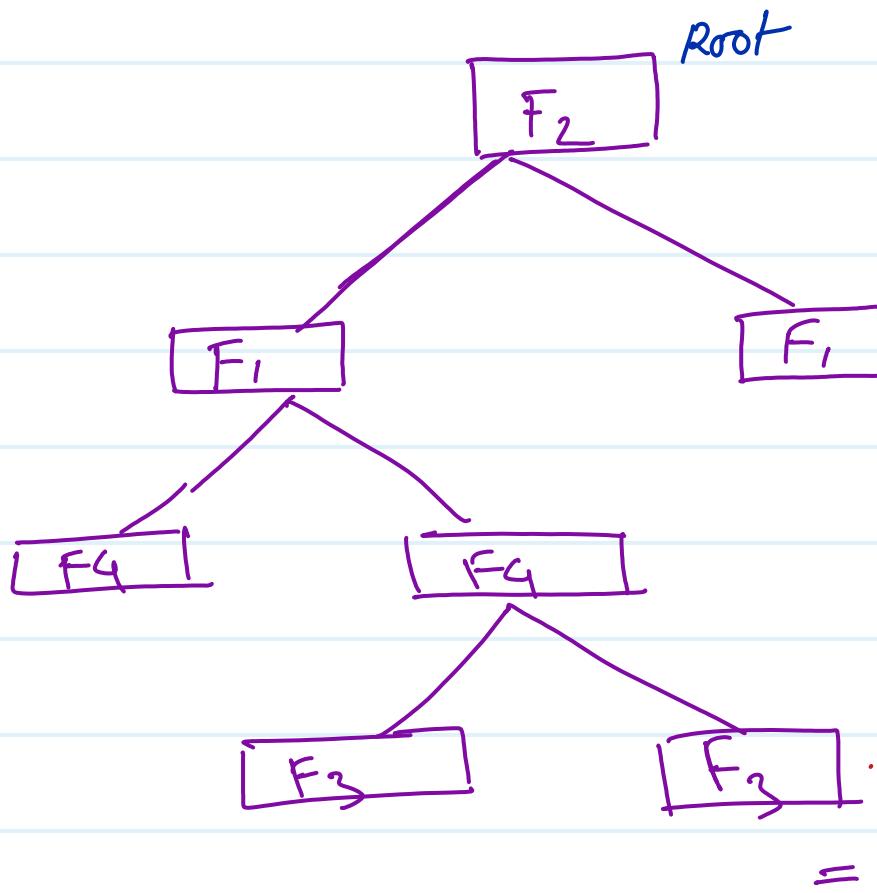
0.025

F_4

0.10

Since F_2 has higher value of information gain's among the all feature so that it will be our root node.

F_1	F_2	F_3	F_4
0.49	0.56	0.025	0.10



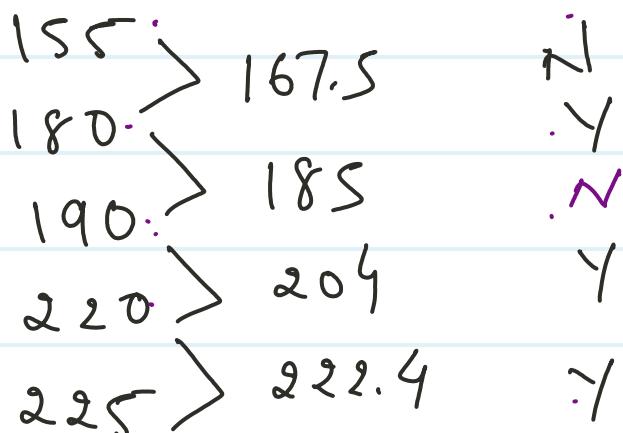
* Independent analysis before making DT

build DT with numerical feature

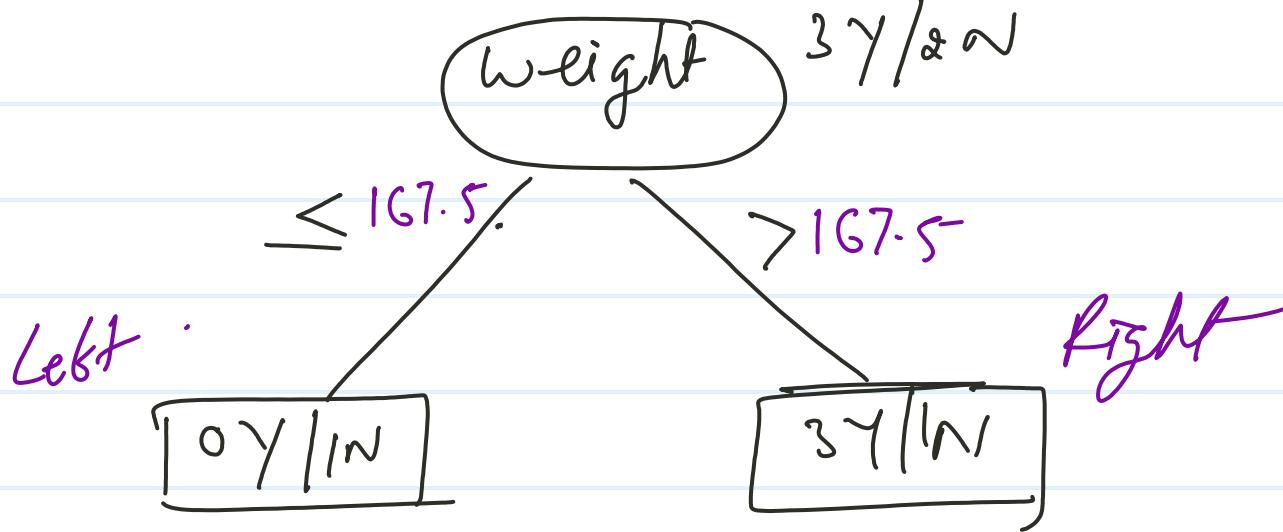
weight heart De.

220	Y
180	Y
225	Y
190	Z
155	Z

weight Heart



with respect to every point avg. value need
to find out gini index / entropy



$$\text{Gini impurity} = 1 - \sum_{i=1}^n p_i^2$$

$$\text{gini (Left)} = 0$$

$$\begin{aligned}\text{gini (Right)} &= 1 - \left[\left(\frac{3}{4}\right)^2 + \left(\frac{1}{4}\right)^2 \right] \\ &= 0.375\end{aligned}$$

$$\text{Information gain} = \text{G.I. [Root]} - \sum \frac{|\text{Sv}|}{|\text{Sv}|} \text{G.I.}_{[\text{child}]}$$

[child]

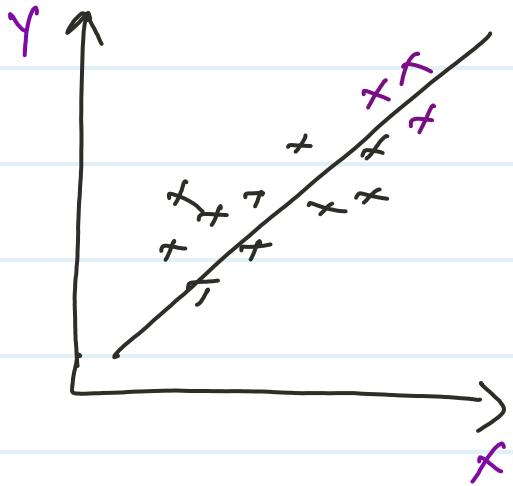
$$\text{G.I. [Root]} = 0.48$$

$$\text{I.G. } [167.5] = 0.48 \left[\frac{1}{5} \times 0 + \frac{4}{5} \times 0.375 \right]$$

$$\text{I.G. } [167.5] = 0.18 =$$

Information Gain should be high
and Gini Index should be low.

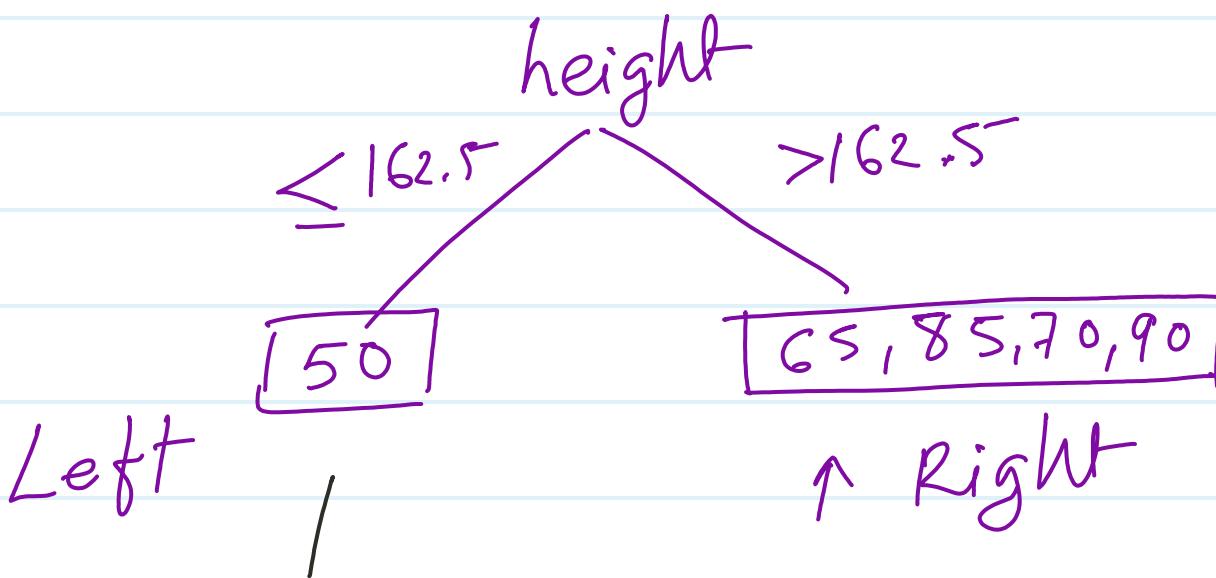
* DT Regression



height weight

160	> 162.5	50
165	> 167.5	65
170	> 172.5	85
175	> 177.5	70
180	> 177.5	90

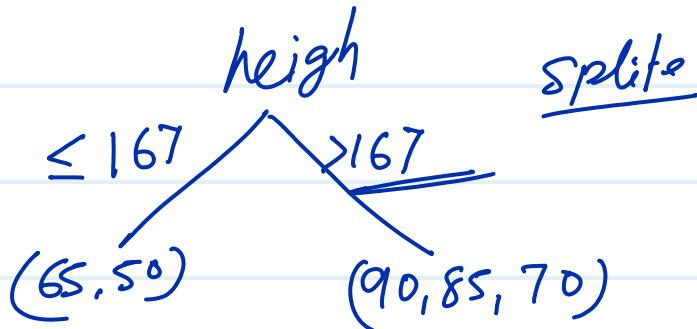
height



50

$$\text{mean} = 77.5$$

height



Ex

height weight

162.5	< 165	65	160 >
167.5	< 160	50	165 >
172.5	< 180	90	170 >
177.5	< 170	85	175 >
	175	70	180

Regression problem weight calculated
with respect to height

- ① step - Short the value of height column (X Feature)
- ② step - Find Adjacent Avg. value b/w data point
- ③ step - Find Information gain with help of entropy and Gini Index.



$$\text{mean} = 77.5$$

Regression -

$$\textcircled{1} \text{ mean} = 77.5$$

$$\textcircled{2} \text{ MSE, RMSE, MAE}$$

$$\text{MSE} = \frac{1}{n} \sum_{i=1}^n (y - \hat{y})^2$$

$$\text{overall mean} = \frac{50 + 65 + 85 + 70 + 90}{5} = 72.$$

$$\begin{aligned} \text{height (variance)} &= (72 - 50)^2 + (72 - 65)^2 + (72 - 85)^2 \\ &\quad + (72 - 70)^2 + (72 - 90)^2 \\ &\quad \hline 5 \end{aligned}$$

$$\text{height}_{\text{variance}} = 206$$

$$\begin{aligned} \text{var}(\text{right}) &= \frac{(77.5 - 65)^2 + (77.5 - 85)^2}{4} \\ &\quad + (77.5 - 20)^2 + (77.5 - 90)^2 \end{aligned}$$

$$\text{var}(\text{right}) = 106.25$$

$$\text{var}(\text{left}) = 50$$

Reduction in variance

$$= \text{var}(\text{root}) - \sum_{i=1}^n w_i \times \text{var}[\text{child}]$$

$$= 206 - \left[\frac{1}{5} \times 0 + \frac{4}{5} \times 106.25 \right]$$

$$\text{Reduction variance} = 121$$

We calculate MSE for all the datapoint whichever is less will be threshold..

x_1 height	x_2 gender	y weight
160	M	65
165	F	70
170	M	80
175	M	90
180	F	100

From height /Gender, . choose root node

for height $MSE = 55.5$

for gender $MSE = 53$

so value of gender MSE is less
It will be our root node. ✓

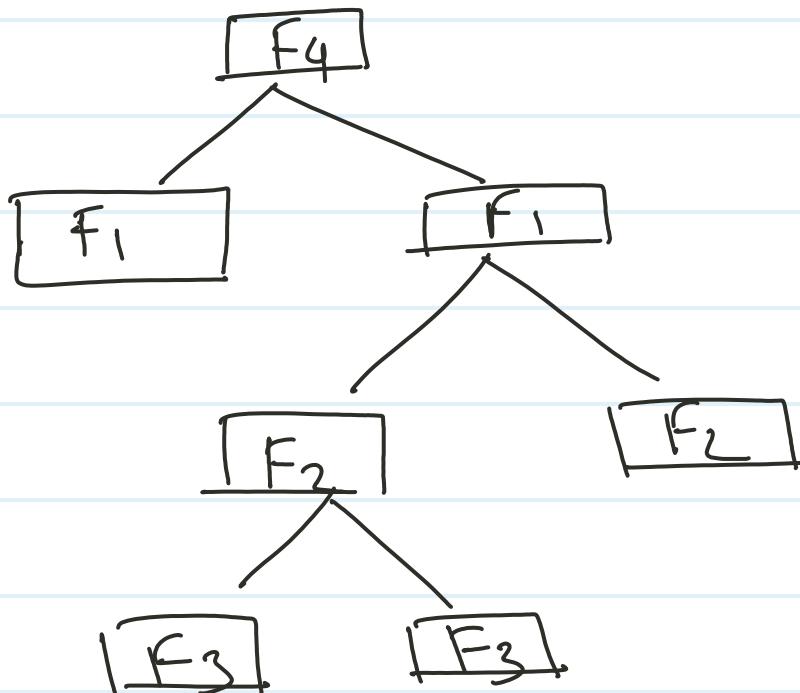
F_1	F_2	F_3	F_4	Target(y)
-	-	-	-	-
-	-	-	-	-
-	-	-	-	-
-	-	-	-	-
-	-	-	-	-

$$F_1 \Rightarrow \text{MSE} = 49$$

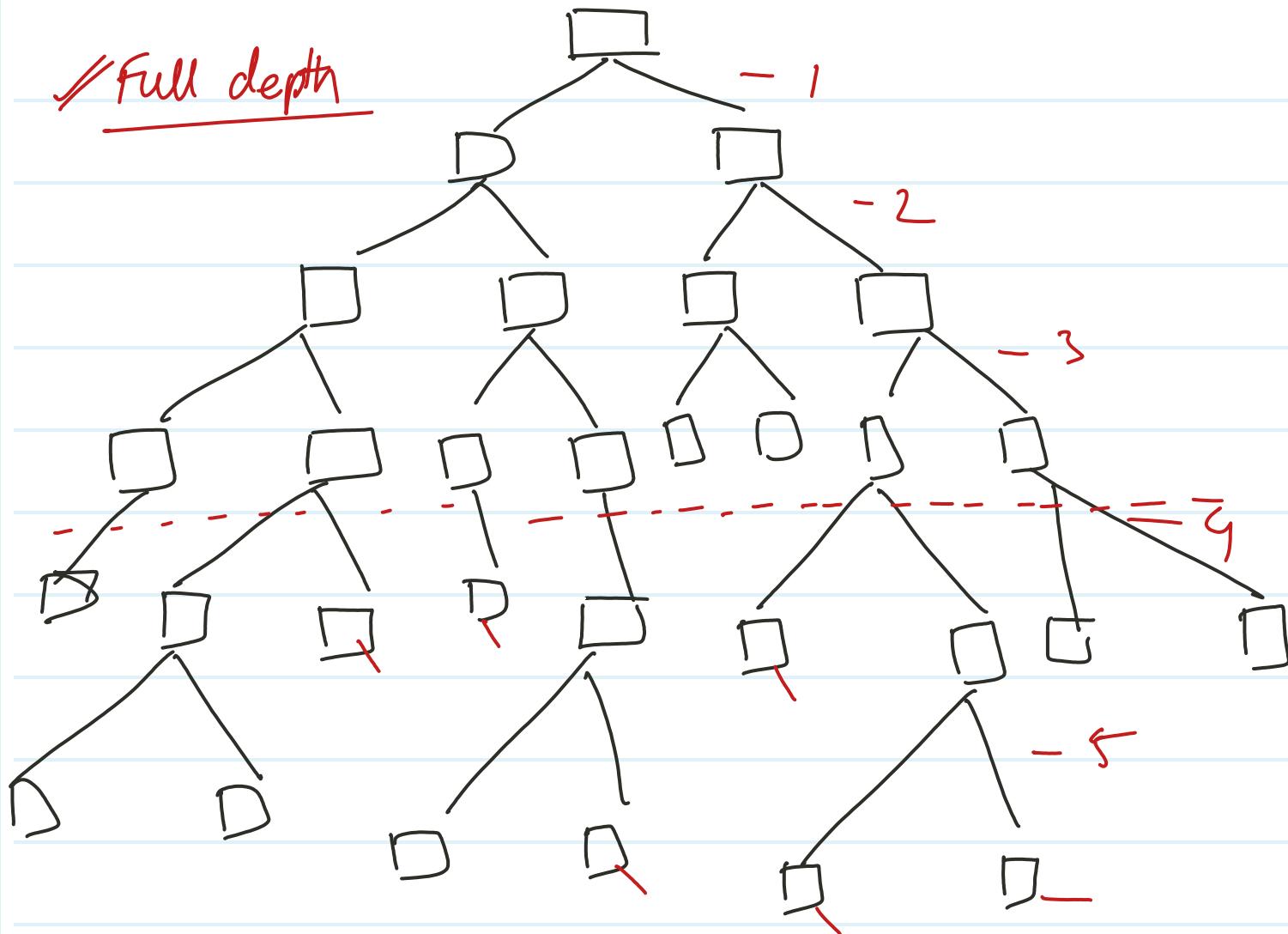
$$F_2 \Rightarrow \text{MSE} = 55$$

$$F_3 \Rightarrow \text{MSE} = 60$$

$$F_4 \Rightarrow \text{MSE} = 43$$



Full depth



Pre-pruning and post-pruning

$$\text{max-depth} = 5$$

$$\text{min-sample-leaf} = 10$$

$$\text{min-sample-split} = 8$$

$$\text{max-feature} = 6$$

[4, 5, 6, 8, 10]

[5, 6, 8, 10, 12]

[4, 6, 8, 10, 12]

[2, 4, 6, 8, 10]

These 4 hyperparameters selected for pre-pruning before build DT algorithms.

Post pruning \Rightarrow

- ① make DT full first
- ② cut DT. using CCP value.
- ③ CCP value is nothing but threshold for gini / Entropy.

CCP value is responsible for depth of Tree.. If CCP is less, the depth will be less.

High CCP value the depth will be more.

$$\text{CCP} = [0.4, 0.5, 0.6, 0.01]$$

For model training either we can use pre-pruning or post-pruning.

- ① When we have large dataset at this time we use pre-pruning.
- ② When we have small dataset at this time we use postpruning.

Why we use post or pre-pruning?

⇒ To avoid model overfitting. ^{from} 1

