# Decision Tree

(i) DT classifier
(ii) DT Regressor

① DT classifier — Type of Technique

(i) ID3
(ii) CART

✳ Entropy and Gini Index →
purity split in dataset
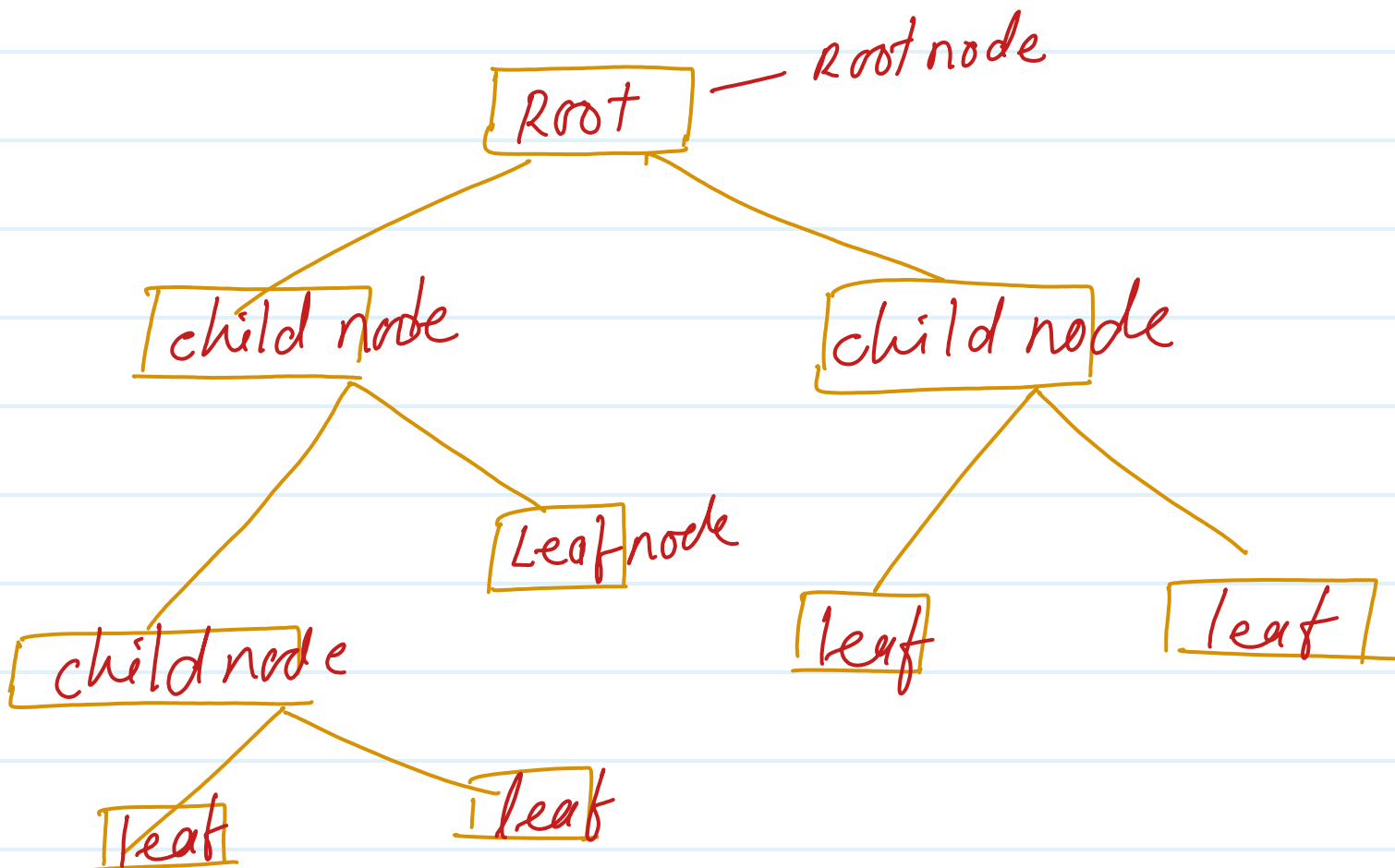
✳ Information Gain → DT feature
split

| weight | height | o/p obese/Noobese. |
|--------|--------|--------------------|
| 60 | 160 | ob |
| 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 |

Root — Root node

child node      child node

Leaf node      leaf    leaf

child node

leaf    leaf

# DT classifier

| 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

9y/5N

**outlook**

2y/3N     4y/0N     3y/2N

Sunny     overcast     Rain

9y/5N

**Temp**
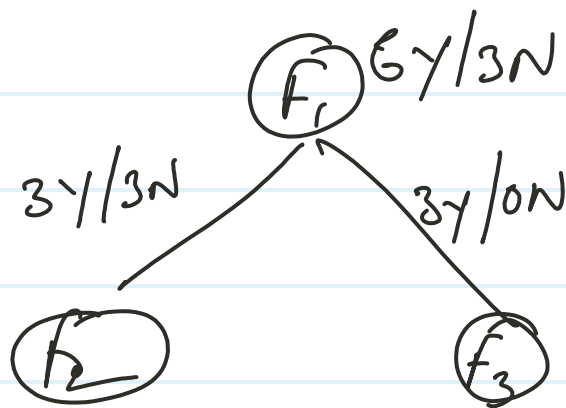
2y/2N     3y/1N     4y/2N

HOT     Cool     mild

# ⓐ Entropy  H(s)

## Formula (Binary)

$$H_{(s)} = -P_{yes} \log_2(P_{yes}) - P_{No} \log_e(P_{No})$$

## multiclass

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

Exm



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

$$\Rightarrow \underline{1}$$

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

$$\Rightarrow 0$$

For the pure split of feature

pure entropy should be zero (0)

for impure split $= 1$

② Gini index (Impurities) .
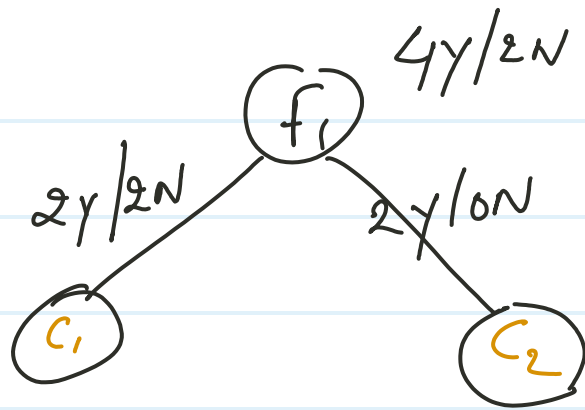
main formula —

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

binary class

$$G.I. = 1 - \sum_{i=1}^{n} \left[ (P_{c_1})^2 + (P_{c_2})^2 \right]$$

multiclass

$$G.I. = 1 - \sum_{i=1}^{n} \left[ (P_{c_1})^2 + (P_{c_2})^2 + (P_{c_3})^2 + \cdots \right]$$

__Exm__



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

$$\boxed{= 0.5}$$

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

$$\boxed{= 0}$$

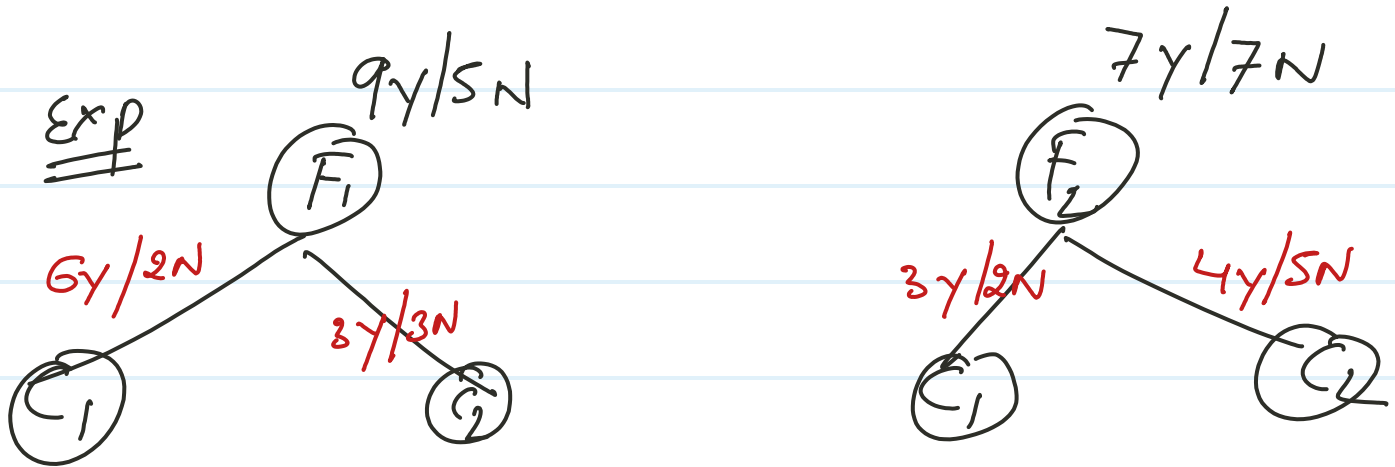Range of entropy = 0 to 1

Gini impurity (index) = 0 to 0.5

# Ⓒ Information Gain

formula —

$$gain(S, f_1) = H(s) - \Sigma \frac{|S_v|}{|S|} H(S_v)$$

Exp #



for $F_1$ ⇒

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

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

$C_1$ ⇒ $H(s) = -\frac{6}{8} \log \frac{6}{8} - \frac{2}{8} \log \frac{2}{8}$

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

$C_2 \Rightarrow$

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

$$\boxed{H_{(S)} = 1}$$

gain of $f_1 \Rightarrow$

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

$$gain\ (S, f_1) = 0.049$$

$f_2 \rightarrow$ $H_{(S)} = -\frac{7}{7} \log \frac{7}{7} - \frac{7}{7} \log \frac{7}{7}$

$$= 0$$

$C_1 \rightarrow H_{(S)} = -\frac{3}{5} \log \frac{3}{5} - \frac{2}{5} \log \frac{2}{5}$

$$= 0.133 + 0.159$$

$$\boxed{= 0.29}$$

$$C_2 \Rightarrow H_{(s)} = -\frac{4}{9} \log \frac{4}{9} - \frac{5}{9} \log \frac{5}{9}$$

$$\boxed{= 0.014}$$

$$f_2 \quad \text{gain}(S, f_2) = 0 - \left[ \frac{5}{14} \times 0.29 + \frac{9}{14} \times 0.014 \right]$$

$$= 0 - \left[ 0.10 + 0.009 \right]$$

$$= -0.10$$

| $F_1$ | $F_2$ | $F_3$ | $F_4$ |
|-------|-------|-------|-------|
| 0.49  | 0.56  | 0.025 | 0.10  |

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

⭐ Independent analysis befor making DT

build DT with numerical feature

| weight | heart De. |
|--------|-----------|
| 220 | Y |
| 180 | Y |
| 225 | Y |
| 190 | N |
| 155 | N |

| weight | | Heart |
|--------|--------|-------|
| 155 | | N |
| 180 | > 167.5 | Y |
| 190 | > 185 | N |
| 220 | > 204 | Y |
| 225 | > 222.4 | Y |

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

$$\boxed{\text{Weight}} \quad 3Y/2N$$

Left $\quad \leq 167.5$ $\qquad\qquad$ $> 167.5$ $\quad$ Right

$$\boxed{0Y \mid 1N} \qquad\qquad \boxed{3Y \mid 1N}$$

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

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

$$\text{gini (Right)} \quad = \quad 1 - \left[ \left(\frac{3}{4}\right)^2 + \left(\frac{1}{4}\right)^2 \right]$$

$$= \quad 0.375$$

$$\text{Information gain} \quad = \quad G.I.[\text{Root}] - \sum_{\text{value}} \frac{|Sv|}{|S|} \, G.I.$$

$$[\text{child}]$$

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

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

$$I.G.[167.5] \quad = \quad 0.18 \, \underline{=}$$

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

# ✳ DT Regression :-

| 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 : - Shoot the value of
height column (x feature)

② step - Find Adjcent Avg. value
b/w data point

③ step - Find Information gain with
help of entropy and Gini Index.

height

$< 162.5$       $> 162.5$

(50)

$(65, 85, 70, 90)$

mean $= 77.5$

Regression —

① mean $= 77.5$

② MSE, RMSE, MAE

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

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

$$height \ (variance) = \frac{(72 - 50)^2 + (72 - 65)^2 + (72 - 85)^2}{5} + \frac{(72 - 70)^2 + (72 - 90)^2}{5}$$

$height_{(variance)} = 206$

$$var_{(Right)} = \frac{(77.5-65)^2 + (77.5-85)^2 + (77.5-70)^2 + (77.5-90)^2}{4}$$

$var_{(Right)} = 106.25$

$var_{(left)} = 50$

## Reduction in variance

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

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

Reduction variance $= 121$

we calculate MSE for all the datapoint whichever is less will be thresold.

| height | gender | 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.

# # pre-pruning and post-pruning

$$max - Depth = 5$$

$$mini - sample - leaf = 10$$

$$min - sample - split = 8$$

$$max - feature = 6$$

These 4 hyperparameter selected for pre-pruning before build D.T. algorithms.

Post pruning =)

① make DT till end
② cut DT. using ccp value.
③ ccp value is nothing but thesold for gini / Entropy.

ccp value is responsible for depth of Tree.. If ccp is less, the depth will be less.
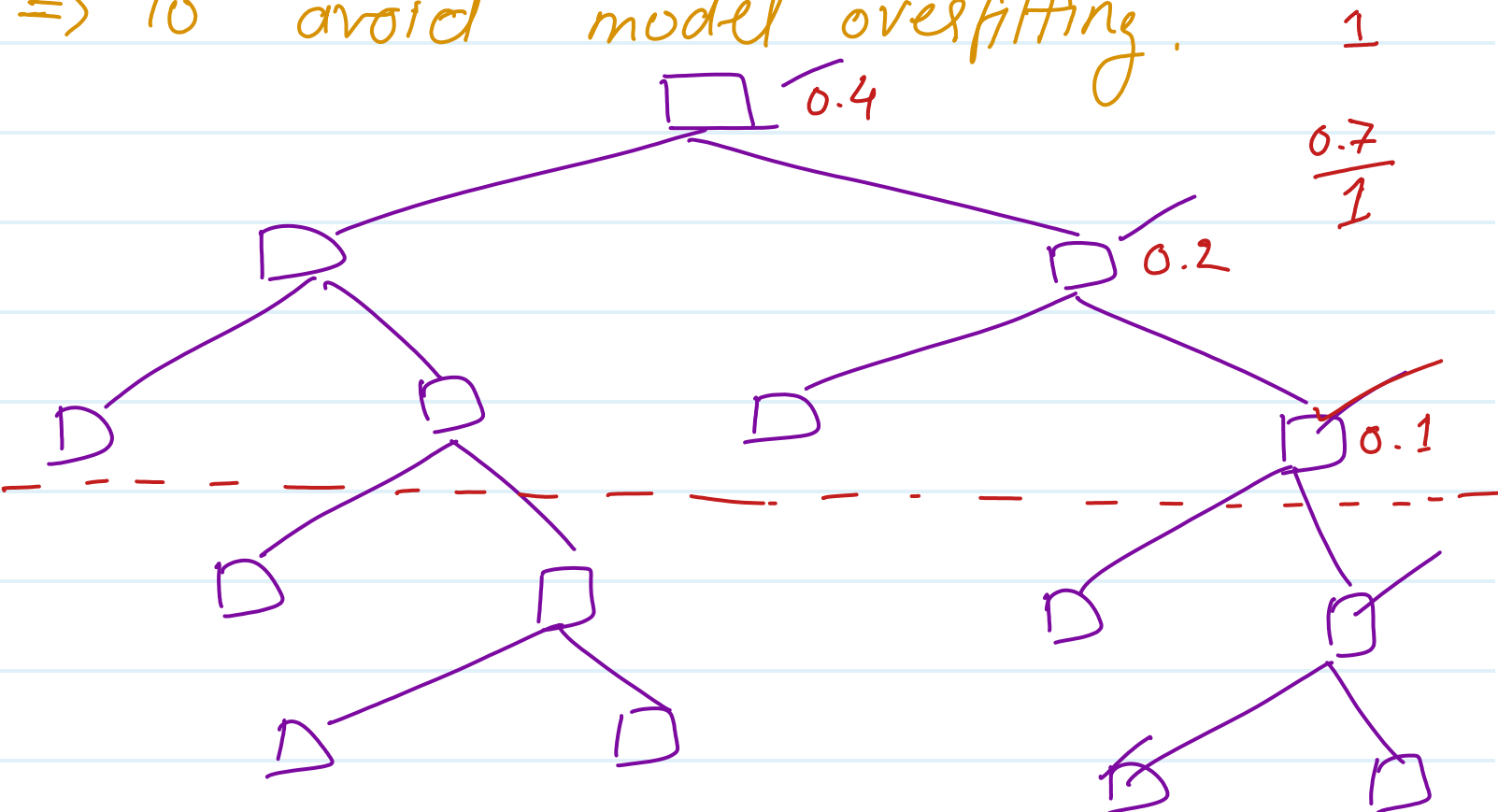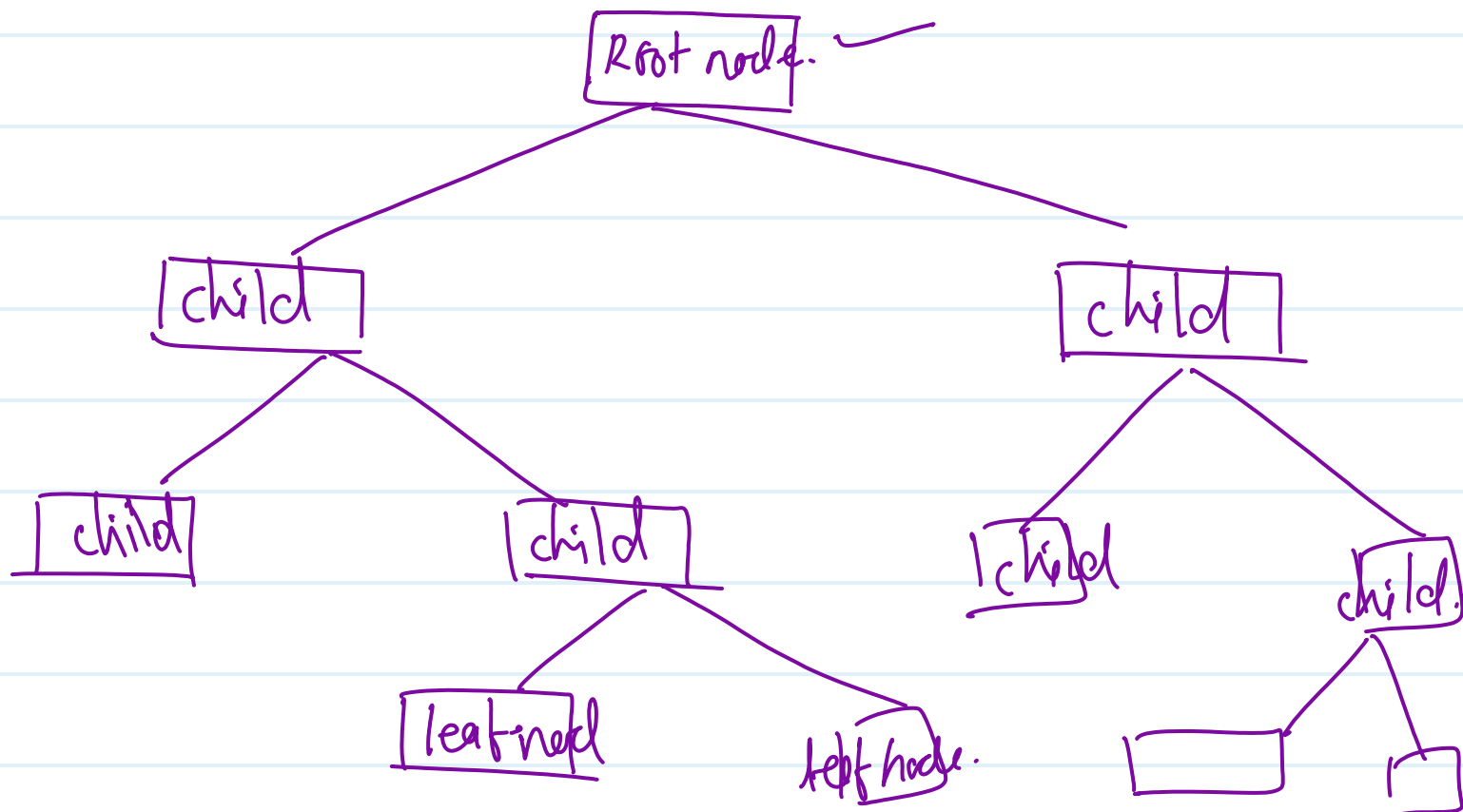High ccp value the depth will be more.

$$ccp = [0.4, 0.5, 0.6, 0.01]$$

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

① When we have large dataset
at this time we use pre-prunning.

② When we have small dataset
at this time we use postprunning.

Why we use post or pre-pruning?

⇒ To avoid model ~from~ overfitting.

A tree diagram with a root node connecting to two child nodes, each branching to further child nodes and leaf nodes.

```
                    Root node.
                   /          \
              child            child
             /     \          /     \
         child     child   child     child.
                  /    \            /      \
          leaf node  leaf node.   [  ]    [ ]
```

| $x_1$ | $x_2$ | $x_3$ | $y$ |
|-------|-------|-------|-----|
|       |       |       |     |

$x_2 =$