

# Decision Tree (DT)

## Agenda

① Build DT with numeric values

② DT regressor

③ Pre pruning and post pruning.

① Build DT with numeric values.

| I/P              | O/P           |
|------------------|---------------|
| weight (numeric) | heart disease |
| 220              | Y             |
| 180              | Y             |
| 226              | Y             |
| 190              | N             |
| 155              | N             |

After ↓ sort

| weight | heart disease |
|--------|---------------|
| 155    | N             |
| 180    | Y             |
| 190    | N             |
| 220    | Y             |
| 226    | Y             |

Avg  
or  
adjacent  
value

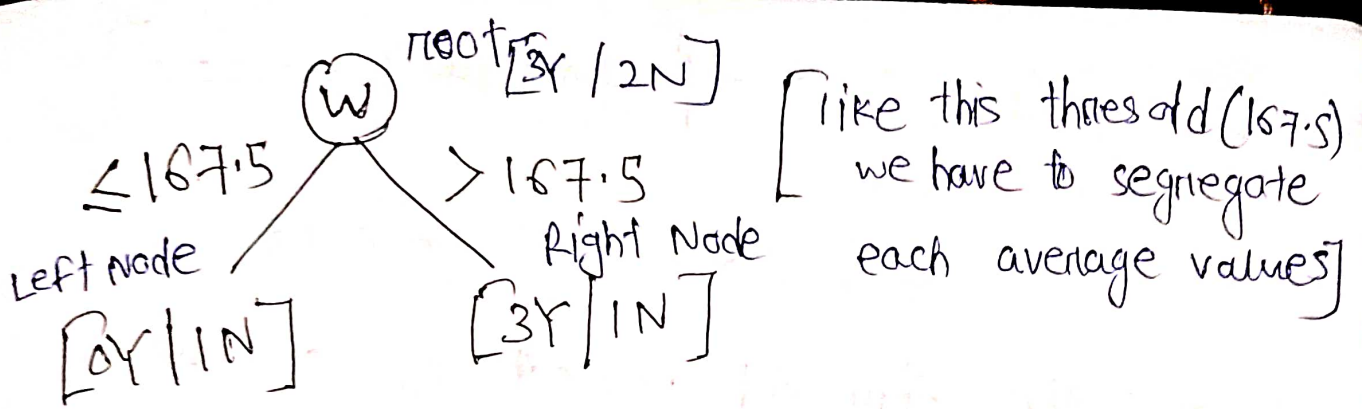
167.5  
185  
205  
223

| weight | heart disease |
|--------|---------------|
| 155    | N             |
| 180    | Y             |
| 190    | N             |
| 220    | Y             |
| 226    | Y             |

① Sort the values.

② Average of the adjacent values.

③ w.r.t every avg. value we need to find out gini impurity/entropy then Info. Gain



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

$$\text{(left node)} = 1 - \left[ \left( \frac{0}{1} \right)^2 + \left( \frac{1}{1} \right)^2 \right]$$

leaf node  $= 1 - 1 = 0$

$$\begin{aligned} \text{Gini impurity (GI)} \\ \text{(Right Node)} &= 1 - \left[ \left( \frac{3}{4} \right)^2 + \left( \frac{1}{4} \right)^2 \right] \\ &= 1 - \frac{10}{16} = 1 - \frac{5}{8} = \frac{3}{8} \\ &= 0.375 \end{aligned}$$

$$\text{Info. Gain} = \text{GI}[\text{Root}] - \sum \frac{|S_v|}{|S|} \times \text{GI}[\text{child}]$$

$$\begin{aligned} \text{GI}(\text{Root}) &= 1 - \left[ \left( \frac{3}{5} \right)^2 + \left( \frac{2}{5} \right)^2 \right] \\ &= 1 - \left( \frac{13}{25} \right) = \frac{12}{25} = 0.48 \end{aligned}$$

$3Y/2N$

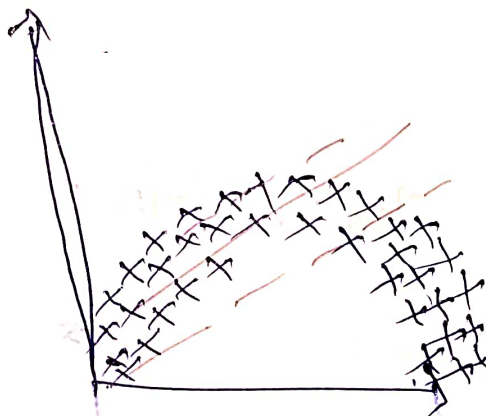
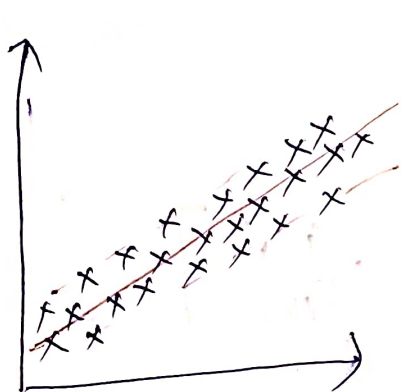
$$\text{Info. Gain} = 0.48 - \left[ \frac{1}{5} \times 0 + \frac{4}{5} \times 0.375 \right]$$

(167.5)

$$= 0.48 - 0.30 = 0.18$$

→ Find the threshold value and segregate data w.r.t each threshold.

### Decision Tree Regressor



By using Linear Regression or SVM we will not get high accuracy ( $R^2$ ) in regression.

So we will use Decision tree regressor (DTR)

| <u>1/P<br/>(numeric)<br/>height</u> | <u>q/P<br/>(numeric)<br/>weight</u> |                | <u>height<br/>1/P</u> | <u>o/P<br/>weight</u> |
|-------------------------------------|-------------------------------------|----------------|-----------------------|-----------------------|
| 165                                 | 65                                  |                | 160                   | 50                    |
| 160                                 | 50                                  |                | 162.5 <               | 65                    |
| 180                                 | 90                                  | Sort<br>values | 165                   | 85                    |
| 170                                 | 85                                  |                | 167.5 <               | 70                    |
| 175                                 | 70                                  |                | 170                   | 90                    |
|                                     |                                     |                | 172.5 <               |                       |
|                                     |                                     |                | 175                   |                       |
|                                     |                                     |                | 177.5 <               |                       |
|                                     |                                     |                | 180                   |                       |

(DT)  
Classification

→ entropy / Gini Impurity

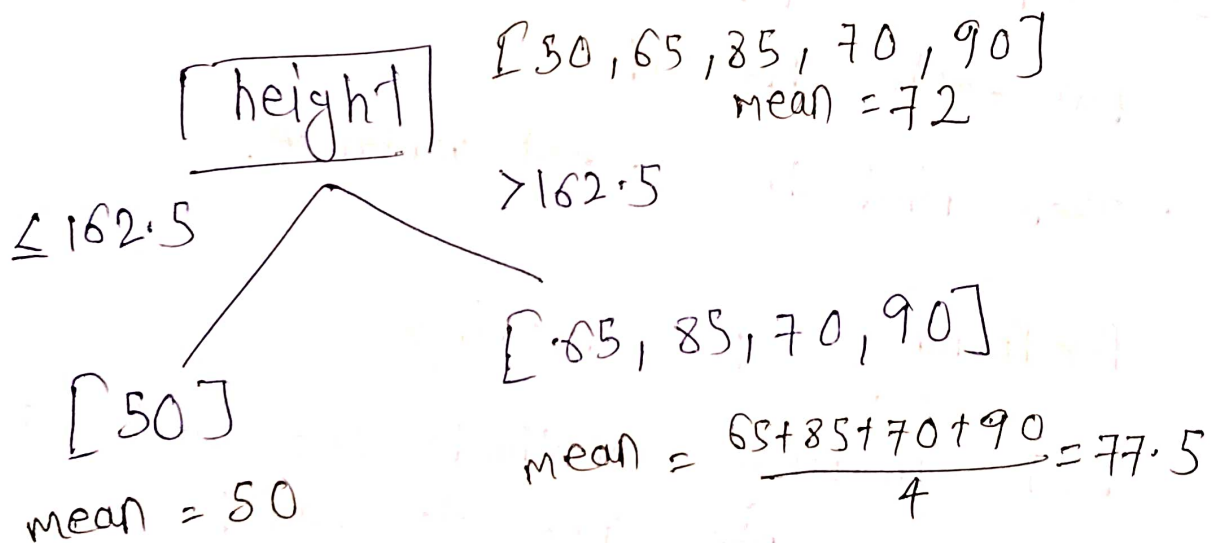
→ Information Gain

(DT)  
vs Regression

→ Mean

→ MSE / MAE / RMSE

→ Reduction in Variance



$$\text{MSE / variance} = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2$$

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

$$= 206$$

for left variance = 0



$$\text{Var [Right]} = \frac{(77.5 - 65)^2 + (77.5 - 85)^2 + (77.5 - 70)^2 + (77.5 - 90)^2}{4}$$

$$= 106.25$$

$$\text{Reduction in variance} = \text{Var [root]} - \sum_{i=1}^n w_i \times \text{Var [child]}$$

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

$$= 206 - 85 = 121$$

→ we have to find the reduction in variance for each threshold

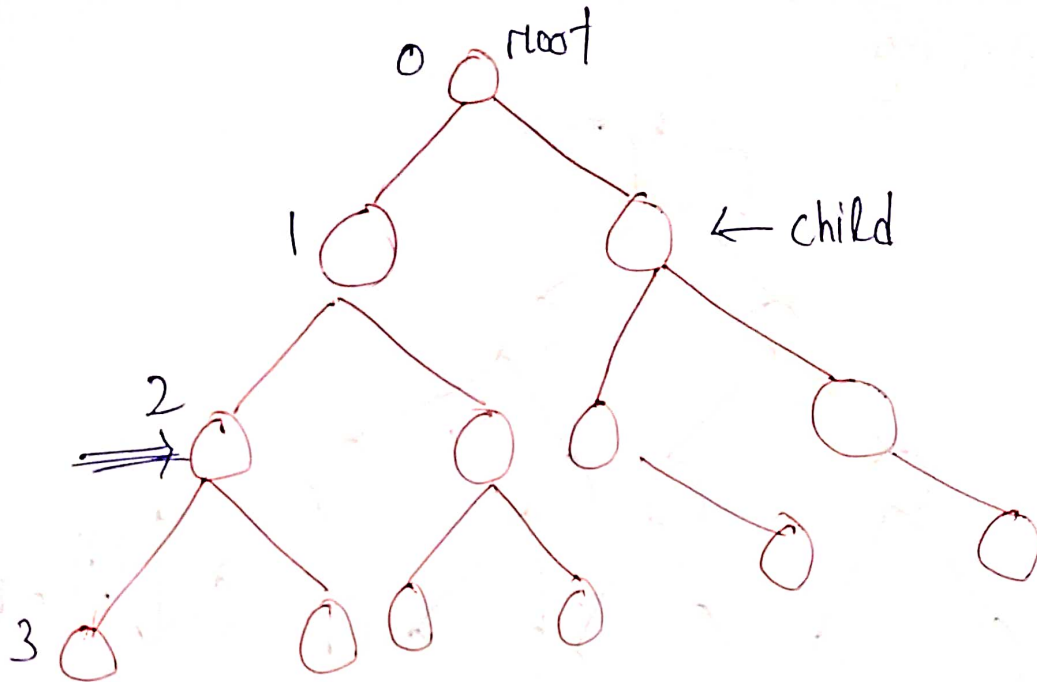
→ Then find the lowest variance.

### Pre Pruning / post Pruning

Both pruning methods are used to cut or prune the tree to avoid overfitting.

### Pre Pruning

Prepruning is done while we are creating decision tree.



if  $\text{max\_depth} = 2$

we will create tree upto 2.

if  $\text{minimum\_sample\_leaf} = 10$

we cut the tree if it should not obey

these criteria.

$\text{minimum\_sample\_split} = 5$

if that node has ~~more~~<sup>less</sup> sample then.  
 $\text{minimum\_sample\_split}$  then ~~only~~ we will not  
 split that node.

①  $\text{max\_depth}$

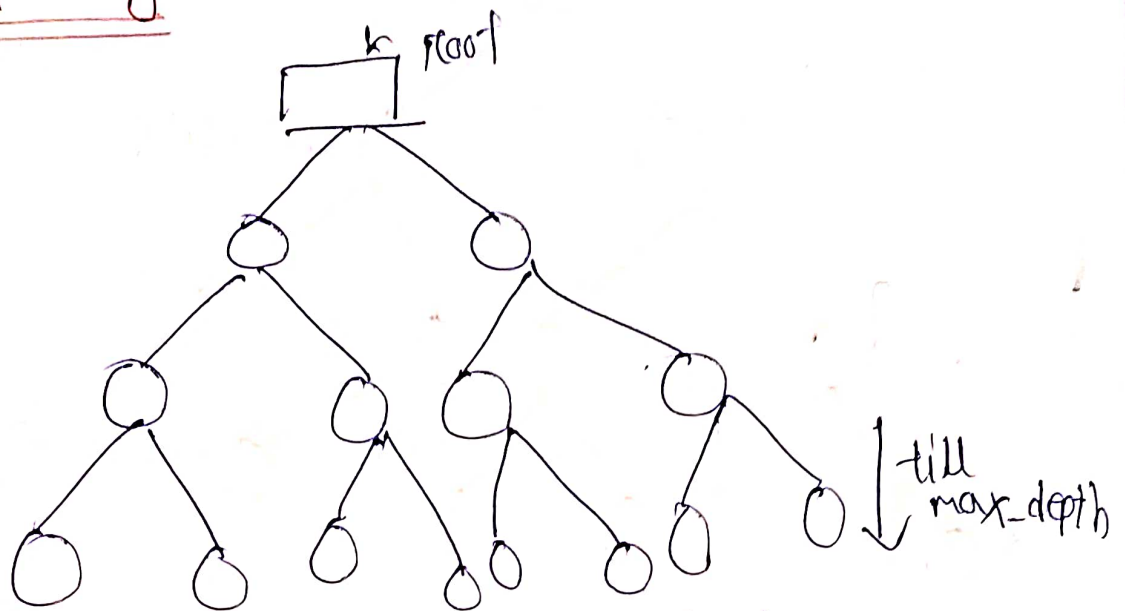
②  $\text{minimum\_sample\_split}$

③  $\text{minimum\_sample\_leaf}$

④  $\text{max\_feature}$

} hyper parameters  
 for  
 prepruning

## Post-pruning



→ Build Decision tree till last

\*) ccp\_alpha (cost complexity pruning)

→ cut the decision tree using ccp\_alpha.

ccp\_alpha: This is a threshold ~~to~~ compared with entropy / gini-impurity  
ccp\_alpha ↓ (low)      Decision tree depth ↓ (low)

ccp\_alpha ↑      DT depth ↑