

12-Nov-2022

Decision Tree Regression

Build Decision Trees with numeric values

continuous

<u>Weight</u>	<u>Heart Disease (Y/N)</u>
220	Y
180	Y
225	Y
190	N
155	N

Steps

① Sort the data in ascending order

<u>weight</u>	<u>Heart disease (Y/N)</u>
155	N
180	Y
190	N
220	Y
225	Y

② find the correct threshold.

Q How can we find threshold?

Ans ③ find the average of the adjacent value

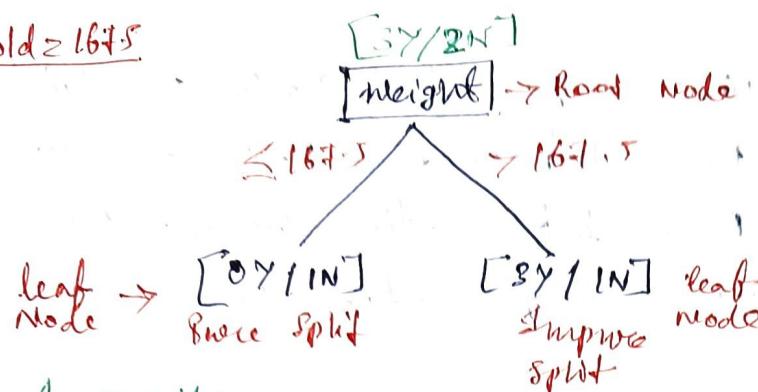
<u>Average</u>	<u>weight</u>	<u>Heart disease (Y/N)</u>
167.5	{ 155	N
	{ 180	Y
185	{ 190	N
205	{ 220	Y
222.5	{ 225	N

intercept

④ find the Gini Impurity or w.r.t every average value and finally calculate Information gain.

{ note: we can take other threshold as well. }

Threshold = 167.5



Gini Impurity

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

$$G.I. [\text{left Node}] = 0$$

$$\begin{aligned} G.I. [\text{Right Node}] &= 1 - (P_Y^2 + P_N^2) \\ &= 1 - \left(\left(\frac{3}{4}\right)^2 + \left(\frac{1}{4}\right)^2\right) \\ &= 1 - \left(\frac{9+1}{16}\right) \\ &\geq \frac{6}{16} \geq \frac{3}{8} \end{aligned}$$

Information gain

$$= 0.375$$

$$\boxed{\text{Information gain} = G.I. [\text{Root}] - \sum_{\text{values}} \frac{|S_i|}{|S|} G.I. [\text{child}]} \rightarrow \text{Weighted Impurity}$$

$$\begin{aligned} G.I. [\text{Root Node}] &= 1 - (P_Y^2 + P_N^2) \\ &= 1 - \left(\left(\frac{3}{5}\right)^2 + \left(\frac{2}{5}\right)^2\right) \\ &= 1 - \left(\frac{9}{25} + \frac{4}{25}\right) \\ &\geq \frac{12}{25} \geq 0.48 \end{aligned}$$

$$I.H = 0.48 - \left\{ \frac{1}{5} \times 0 + \frac{4}{5} \times 0.375 \right\} = 0.18 \rightarrow \text{w.r.t } 167.5 \text{ as threshold}$$

Threshold

Let's assume,

for threshold 16.5, information gain $= 0.18$ (calculated)

for threshold 18.5, information gain $= 0.20$ (assumed)

for threshold 20.5, information gain $= 0.16$ (assumed)

for threshold 22.5, information gain $= 0.3$ (assumed)

So we will choose threshold as 22.5 because it is having more information gain.



$$\text{Information Gain} = \frac{(1 - \frac{1}{3}) - 1}{(1 - \frac{1}{3})} \text{ mini entropy}$$

$$\frac{2}{3} - \frac{1}{3} = \frac{1}{3}$$

$$1 - \frac{1}{3} = \frac{2}{3}$$

Decision Tree \rightarrow Feature A = most informative

Feature A = Coldest Part

$$(1 - \frac{1}{3}) - 1 =$$

$$(1 - \frac{1}{3}) - 1 =$$

$$\frac{2}{3} - \frac{1}{3} = \frac{1}{3}$$

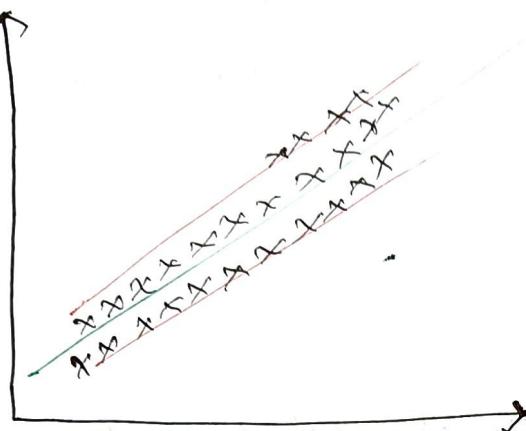
Feature A = Freezing + Warm \rightarrow 0.33 = 0.33

Most informative feature

Decision-Trees - Regression

Q why do we need DTR?

A:

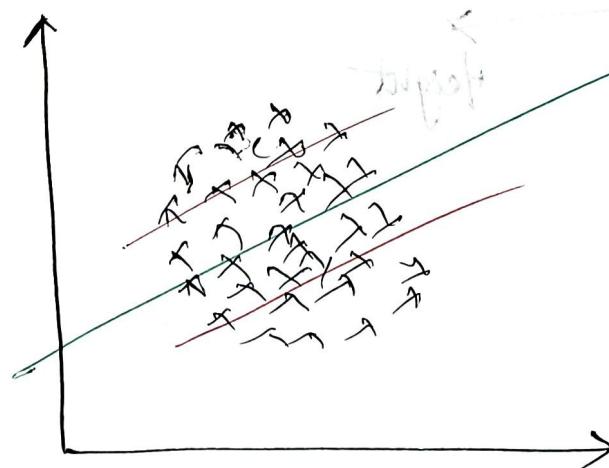


Best fit line

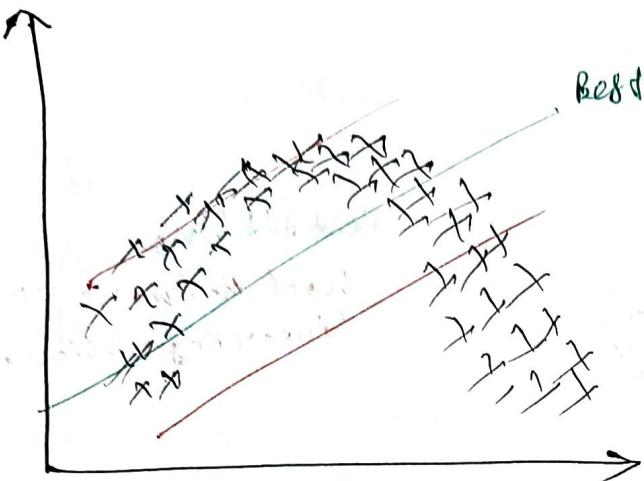
Worst R-Square w.r.t linear Regression, SVM



Worst Decision Tree



Worst Decision Tree

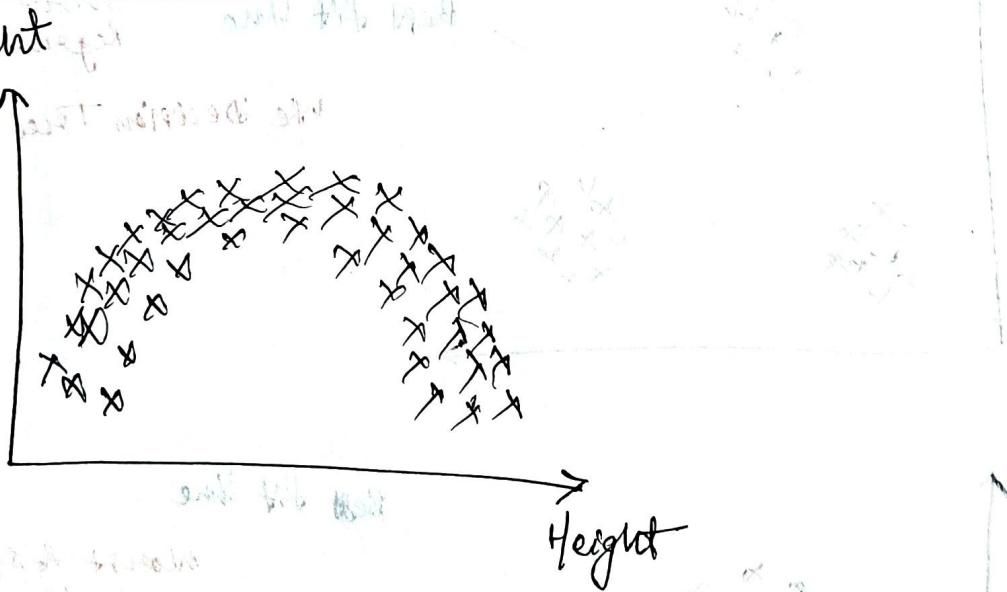


Best fit line

Worst fit square with
Linear Regression SVM
Use Decision Tree

Note: - We can use SVM Kernel and can plot data in higher dimension but that works well for classification mode. SVM Kernel will not work properly in regressor model.

Non-linear data



~~Height~~

~~Height~~

65

Height	Weight
165	50
160	35
180	90
170	85
175	70



- ① Sort the values in ascending order and calculate adj. avg. value.

<u>Average</u>	<u>Height</u>	<u>Weight</u>
162.5	{ 160	50
	{ 165	65
167.5	{ 170	85
172.5	{ 175	70
177.5	{ 180	90

Classification Problem steps in DT

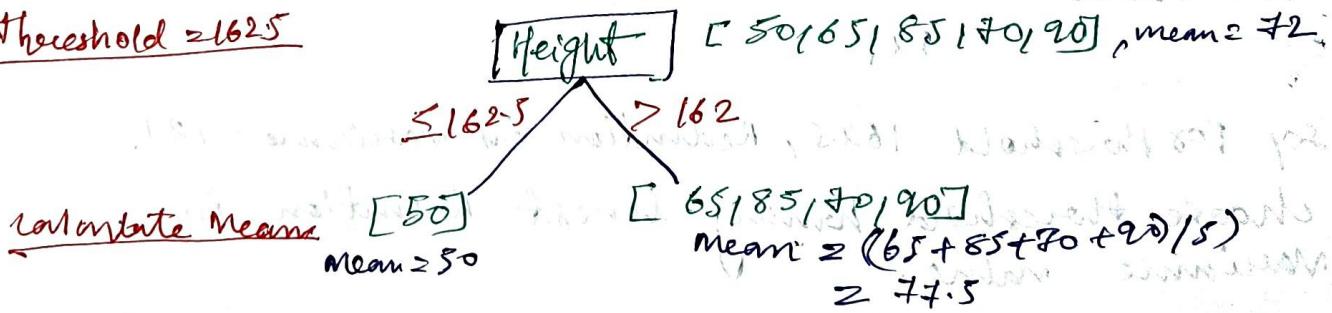
- ① Entropy
- ② Gini Impurity
- ③ Information Gain

Regression Problem steps in DT

- ① Mean
- ② MSE / MAE / RMSE
- ③ Reduction of Variance

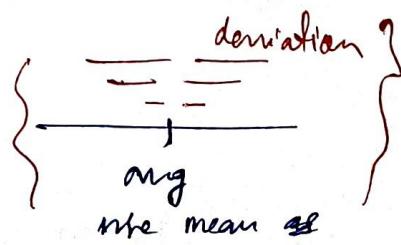
} we will select the feature to threshold having less Mean, MSE, RMSE, MAE values.

Threshold = 162.5



calculate MSE / MAE / RMSE

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



Root Node

$$\text{Height}[\text{Variance}] \geq \frac{(72-50)^2 + (72-65)^2 + (72-85)^2 + (72-90)^2 + (72-90)^2}{5}$$

$$= 206$$

With the standard deviation of the variance node.

Height [50]

$$\text{Variance}[\text{left}] = \frac{(50-50)^2}{1} = 0$$

[65, 85, 70, 90]

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

$$= 106.25$$

Note:

Variance and MSE are same.

Reduction in Variance

$\Delta \text{Var} = \frac{w_i \cdot w_j}{\sum w_i \cdot w_j}$, weight = $\frac{1}{5}$

$$\frac{\text{Reduction in Variance}}{\text{Variance}} = \text{Var}(\text{root}) - \sum_{i=1}^n w_i \cdot \text{var}[\text{child}]$$

wi's weighted Average

$$\text{Reduction in Variance} = 206 - \left(\frac{1}{5} \times 0 + \frac{4}{5} \times 106.25 \right)$$

$$\therefore \text{Reduction in Variance} = 121$$

so, for threshold 162.5, Reduction in Variance = 121.

→ choose threshold having lowest Reduction in Variance value.

Let's assume,

for threshold 162.5, $RV = 121$ (calculated)

for threshold 167.5, $RV = 142$ (assumed)

for threshold 172.5, $RV = 106$ (assumed)

for threshold 177.5, $RV = 131$ (Assumed)

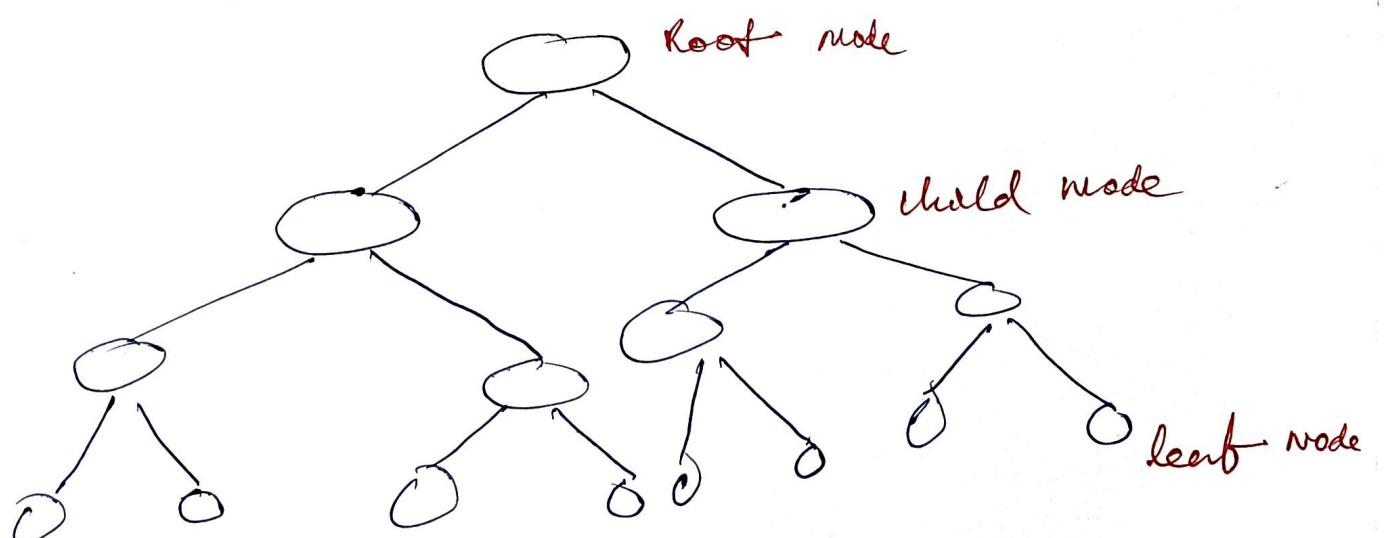
hence we will choose threshold as 162.5 because it has lowest reduction in variance value.

→ In multiple features take the 1st features and build the tree using least RV value.

→ Gini impurity will always perform binary classification.

{ Blackbox Model → we cannot visualize the mathematical calculation }
whitebox Model → we can visualize the mathematical calculation behind the algorithm.

Pre-Pruning and Post Pruning



→ Pruning means cutting.

→ We do pruning to avoid overfitting in the model.

Pre-Pruning

→ While creating decision tree, decide whether to prune or not.

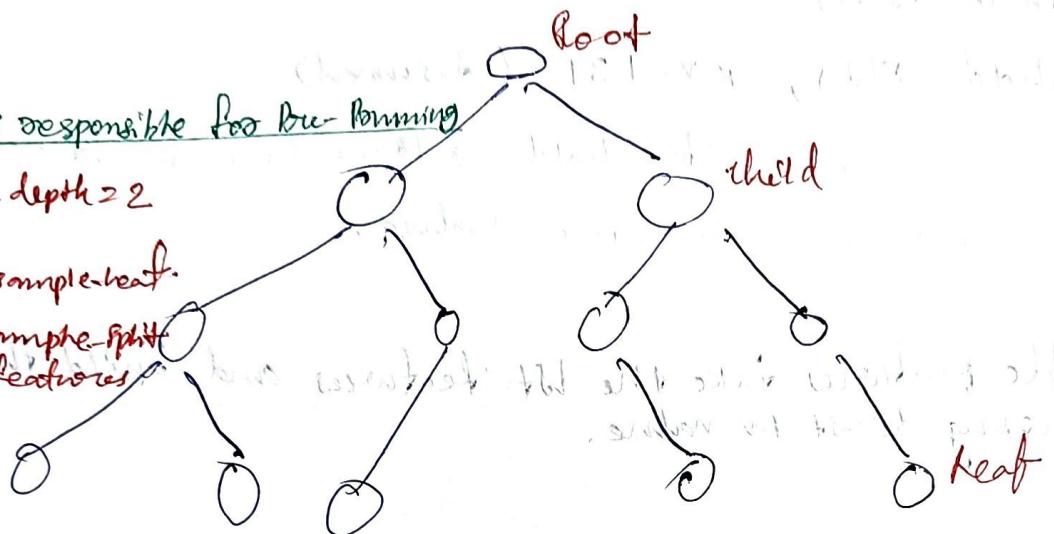
Factors responsible for Pre-Pruning

① max-depth = 2

② min-sample-leaf

③ min-sample-split

④ max-features



→ Before creating DT, we decide factors for pre-pruning.

Hyperparameters for Pre-Pruning

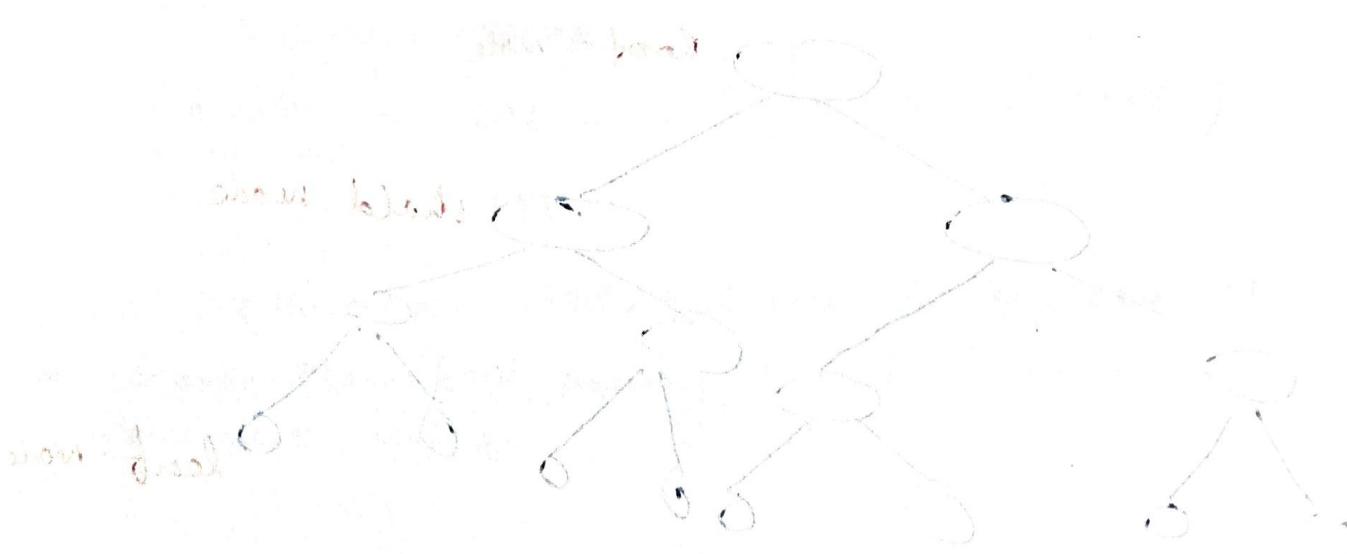
① max-depth = 2

② min-sample-leaf

③ min-sample-split

④ max-features = max-features

→ When Ratio will be high we can pruning ~~at the full that point~~ perform



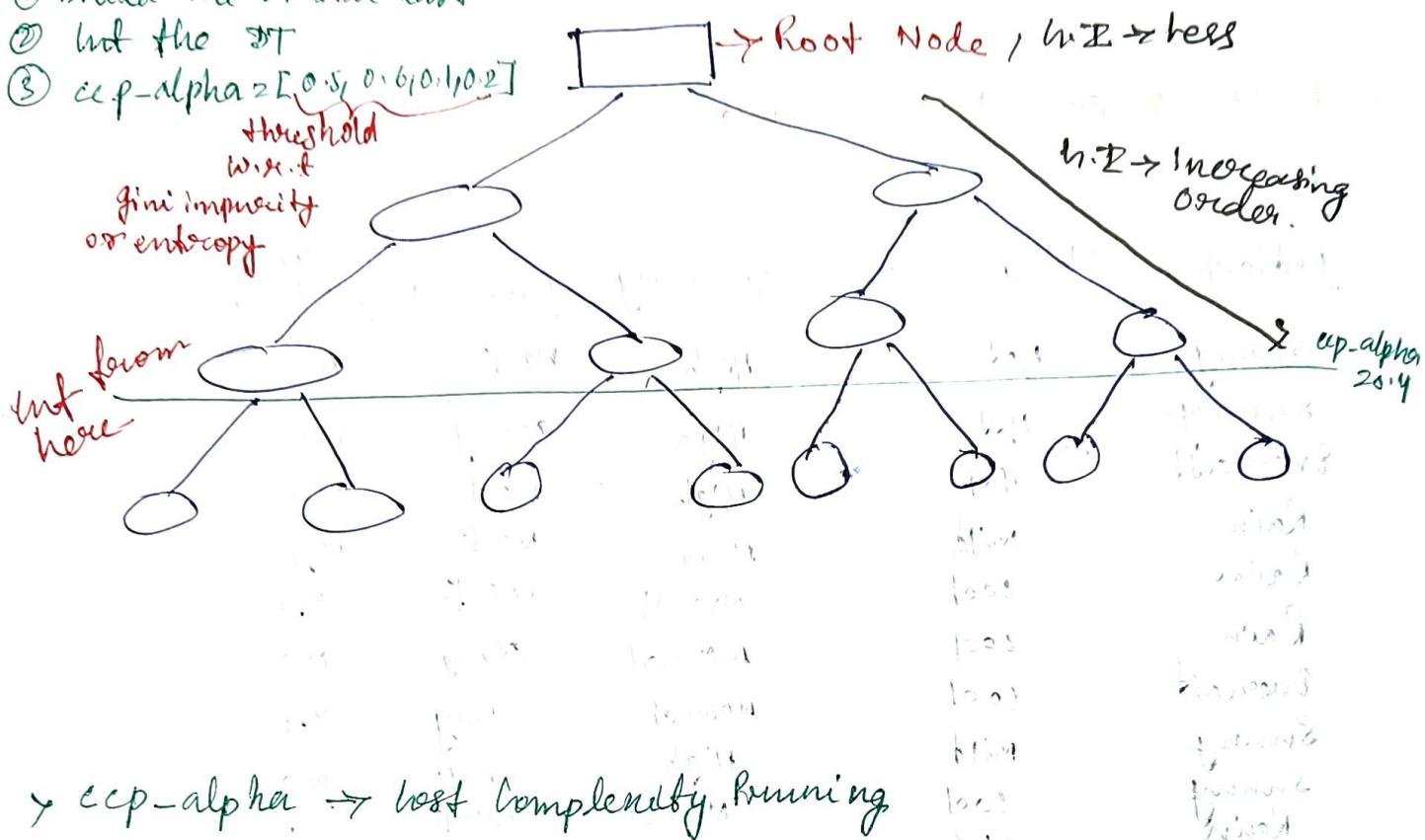
When ratio will be high we can pruning ~~at the full that point~~ perform

Post Penning

→ First build the decision and then act the process.

Step 1

- ① Build the DT first
 - ② fit the DT
 - ③ $\alpha, \rho, \text{alpha} = [0.5, 0.6, 0.1, 0.2]$



↳ CCP-alpha → best complementarityunning

→ If α value doesn't match with α^* , then we can't go further.

> Up-alpha value & height of the PT

\propto
Directly
Proportional