

12/11/22

Decision Tree Regress

1. Build Decision Tree with numeric values

2. DT regress

3. Visualize DT

4. Pre Pruning and Post Pruning

5. DTR Practical Implementation.

Black Box model :- The model in which we can't visualize the mathematical analysis.

White Box model :- we can clearly visualize the mathematical analysis.

like, max_depth, min_sample_split, criterion, mse, etc

Steps

1. Split the values eg DT.

2. Avg. of the adjacent value

3. w.r. to every avg. value, need to find out "gini impurity" / entropy

I.G.

Ex

Weight	heart disease
220	Y
180	Y
225	Y
190	N
155	N

After sorting values

weight	heart disease
155	N
180	Y
190	N
220	Y
225	Y

$0.16 \times 0.85 = 0.136$

$0.16 \times 0.85 = 0.136$

$0.16 \times 0.85 = 0.136$

$0.16 \times 0.85 = 0.136$

$0.16 \times 0.85 = 0.136$

$0.16 \times 0.85 = 0.136$

$0.16 \times 0.85 = 0.136$

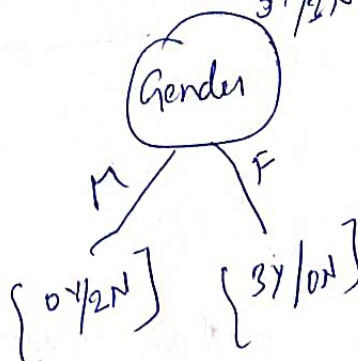
$0.16 \times 0.85 = 0.136$

$0.16 \times 0.85 = 0.136$

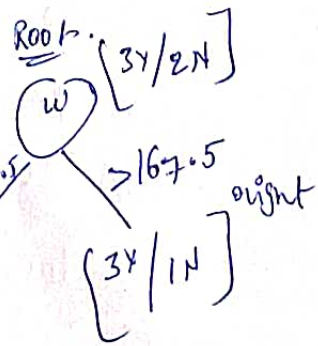
$0.16 \times 0.85 = 0.136$

Gender heart disease

M N
F Y
M N
F Y
F Y
F Y



correct threshold with low impurity



left leaf node

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

Root node :- weight : $1 - \left[\left(\frac{3}{5} \right)^2 + \left(\frac{2}{5} \right)^2 \right]$

$$= 1 - \frac{9+4}{25} = \frac{25-13}{25} = \frac{12}{25} = 0.48.$$

Left node

0Y/1N : $1 - \left[\left(\frac{0}{1} \right)^2 + \left(\frac{1}{1} \right)^2 \right]$
 $= 1 - 1 = 0$

Right node :-

3Y/1N : $1 - \left[\left(\frac{3}{4} \right)^2 + \left(\frac{1}{4} \right)^2 \right]$

$$= 1 - \left[\frac{9+1}{16} \right] = 1 - \frac{10}{16} = \frac{6}{16} = 0.375$$

Information Gain = $H(S) - \sum \frac{H(S_v)}{|S_v|} H(S_v)$ → weighted impurity

Information Gain = $G_2(\text{Root}) - \sum_{v=\text{value}(S)} \frac{|S_v|}{|S|} G_2(\text{Child})$

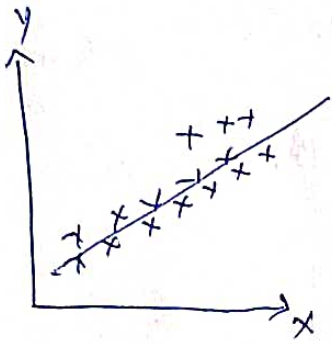
$$= 0.48 - 0.48 - \left[\frac{1}{5} \times 0 + \frac{4}{5} \times 0.375 \right]$$

$$IG(167.5) = 0.48 - \left[\frac{4 \times 0.375}{5} \right] = 0.18$$

G.I	weight	height
	155	N
	180	Y
0.18	167.5	N
	190	Y
0.20	185	Y
	220	Y
0.30	205	Y
	225	

0.30 → low impurity
Conclusion :- Gain should be Highest ; impurity is low (less high).

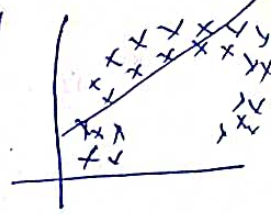
Decision tree - Regressor :-



when R^2

Svm

little good



if $R^2 \ll 1$
worse

R^2 will near to 1
good



height	weight
165	65
160	50
180	90
170	80
175	70

1. Sorting values

	height	weight
avg.		
162.5	160	50
	165	65
167.5	170	80
	175	70
172.5	180	90
177.5		

Classification Tree :-

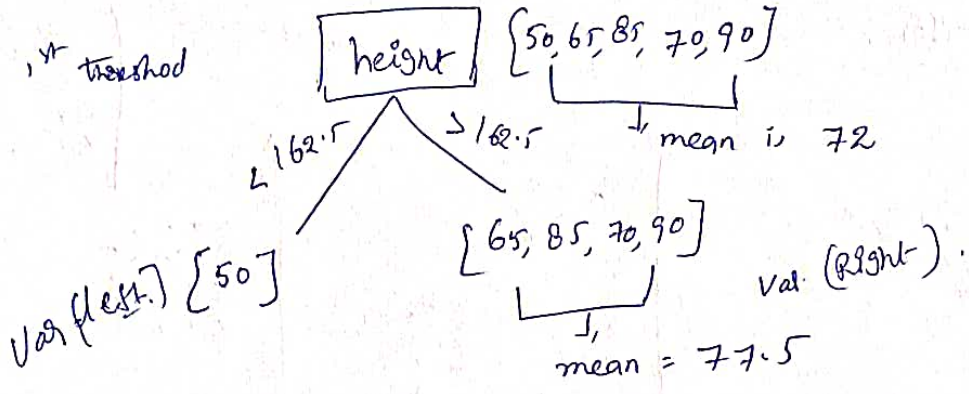
- Entropy
- Gini Impurity
- Information gain

Regression :- Target variable in numeric value

1. Mean
2. MSE/MAE/RMSE
3. Reduction variance

$$MSE = \text{Variance}$$

best threshold



error

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



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

$$= 206$$

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

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

$$= 106.25$$

$$\text{Reduction in variance} = \text{Var (root)} - \sum_{i=1}^n w_i \cdot \text{Var (child)}$$

used '0' instead of 50

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

$$RV = 121$$

Conclusion :- Here we need to choose MSE ↓ RV ↑
low in err (MSE) ⇒

Best threshold → Min MSE

→ cutting
Prepruning and Post Pruning

Overfitting

↓
Trimming the data

Pre-pruning :- before building a decision tree.

while I am going to
create a decision tree.

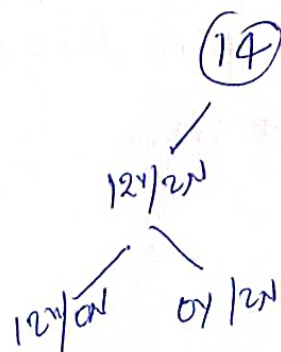
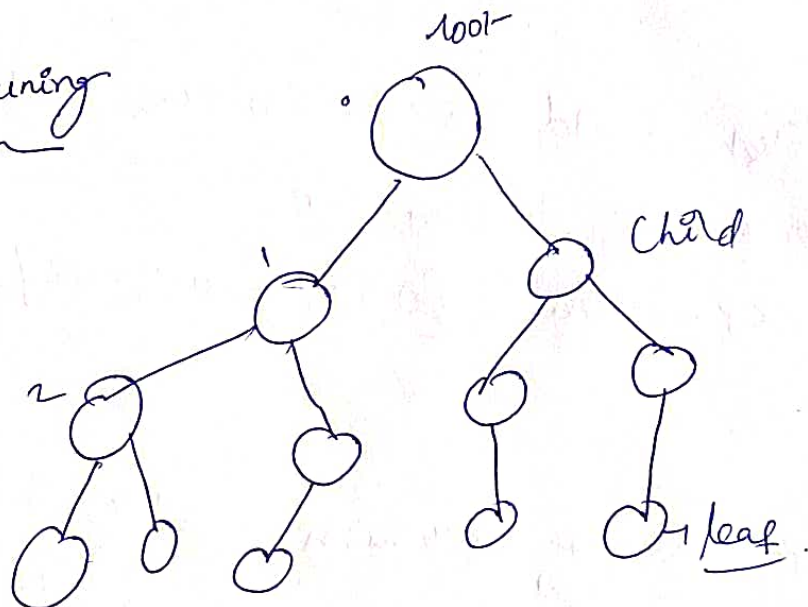
max-depth = 2

min-sample-leaf =

min-sample-split, max-feature

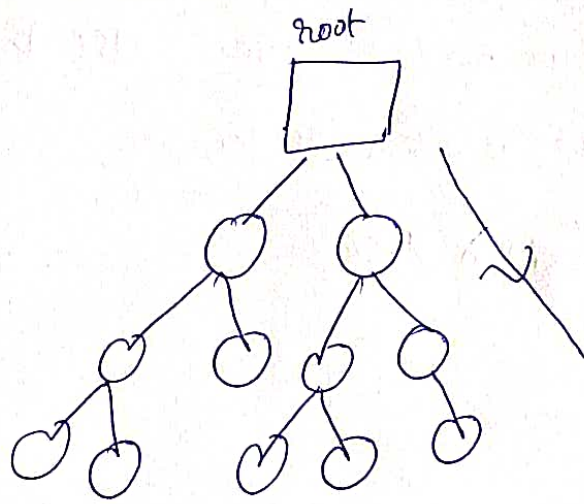
} Hyper Parameters

Responsible for Pre Pruning =



→ overfitting → Pattern well with training data
low bias

Post Pruning :-



1. First build the decision tree till last
2. Cut the decision tree
3. $ccp - \alpha = [0.5, 0.6, 0.1, 0.2] \Rightarrow$ threshold w.r.t. gini / entropy.

if $ccp - \alpha = 0.4$ \Rightarrow very high - \Rightarrow DT - height \rightarrow high (depth)

my - DT height (low)
0.2

Assignment :- Graduation Admission

Individual household Prediction

Census income dataset

Decision tree regression ;