

Decision Trees

16 February 2022 20:10

CART - Classification and Regression Trees.

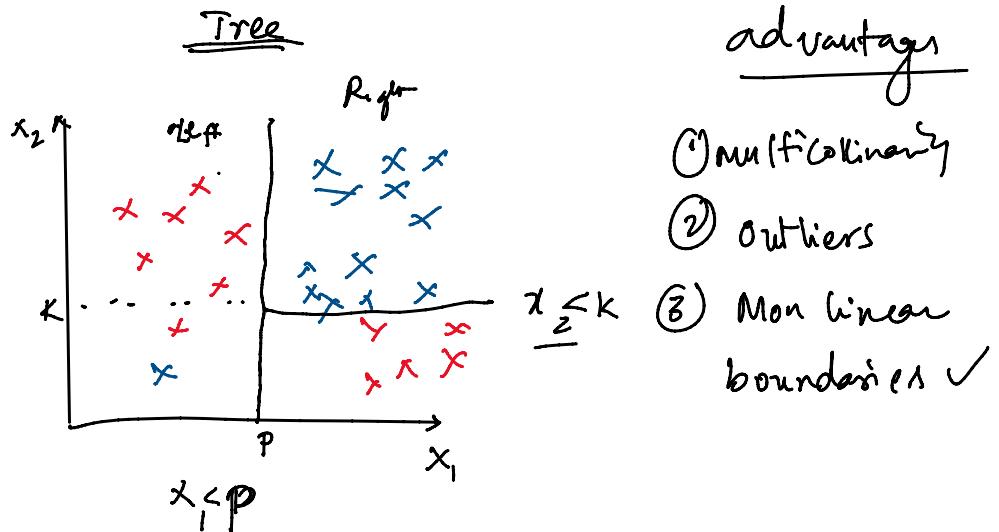
Classification

Linear

① Multicollinearity

② Outliers

③ Non linear relationships



advantages

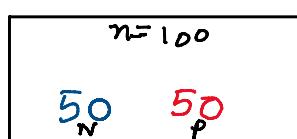
① Multicollinearity

② Outliers

③ Non linear boundaries ✓

Decision Tree - Root Node

$$\{x_1, x_2, \dots, x_m\}$$



← ?

y

100 < 50 & 10
50 Yes

Measures of Impurity

Purity - When all obs belong to Same class

Impurity - When mix of both the classes

Metrics :- Entropy, Gini Index.

$$\text{Entropy} = - \sum_{i=1}^k p_i \log_2 p_i \quad k - \text{no. of classes}$$

Down Cut

$$T = - \frac{50}{100} \log_2 \frac{50}{100} - \frac{50}{100} \log_2 \frac{50}{100}$$

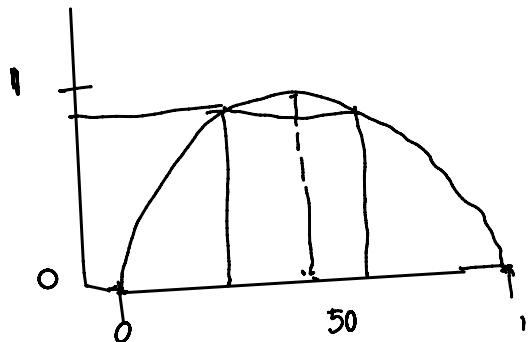
Perfect

Impure $E = 1$

$$\begin{aligned}
 & \text{Class} \\
 & \boxed{50 \quad 50} = -\frac{50}{100} \log_2 \frac{50}{100} - \frac{50}{100} \log_2 \frac{50}{100} \\
 & = -\frac{1}{2} \log_2 \frac{1}{2} - \frac{1}{2} \log_2 \frac{1}{2} \\
 & = \frac{1}{2} \log_2^2 + \frac{1}{2} \log_2^2 \\
 & = \frac{1}{2} + \frac{1}{2} = 1
 \end{aligned}$$

Pure $E = 0$

$$\boxed{0 \quad 50} = \frac{0}{50} \log_2 \frac{0}{50} - \frac{50}{50} \log_2 \frac{50}{50} = 0$$



Gini Index

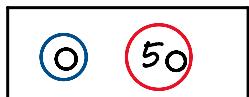


perfect
Impure $G = 0.5$

$$G = 1 - \sum_{i=1}^k (P_i)^2$$

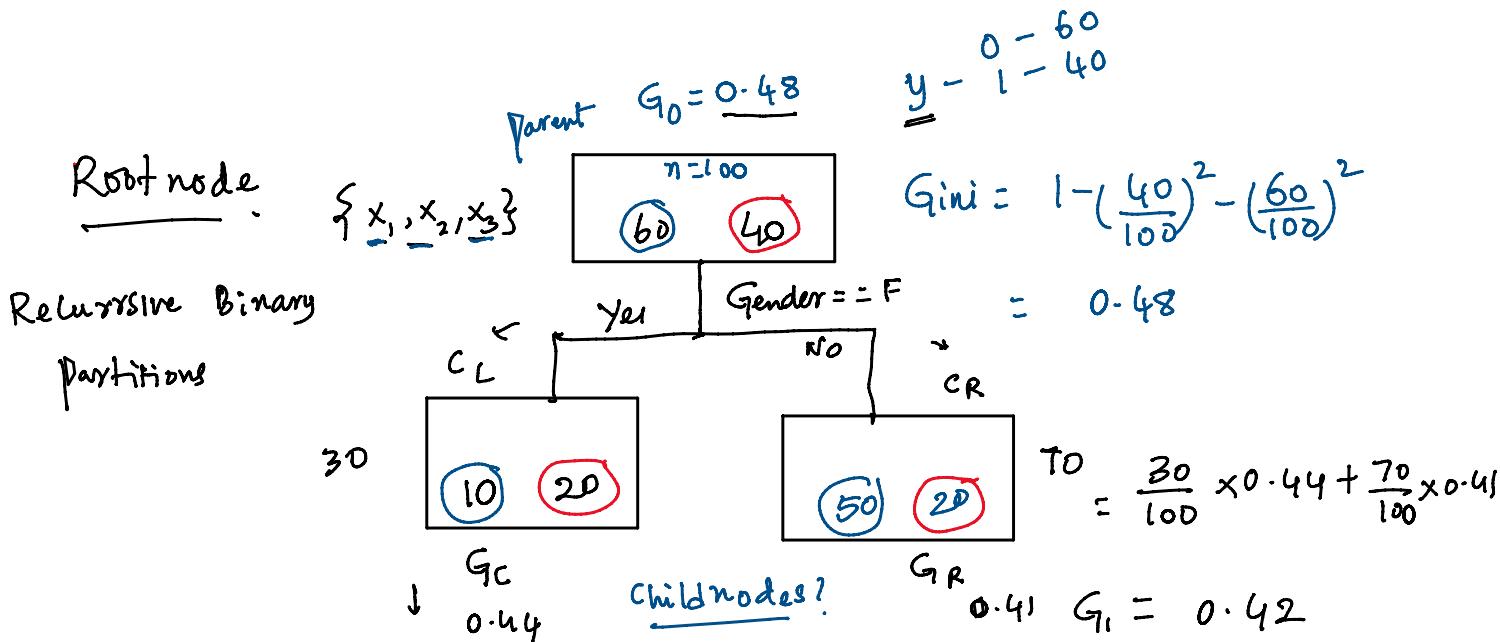
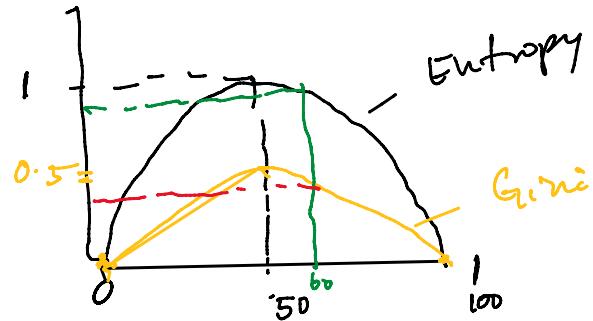
$$= 1 - \left(\frac{50}{100}\right)^2 - \left(\frac{50}{100}\right)^2$$

$$= 1 - \frac{1}{4} - \frac{1}{4} = 0.5$$



$$= 1 - \left(\frac{50}{50}\right)^2 - \left(\frac{0}{50}\right)^2 \quad Gini = (0 - 0.5)$$

$$\therefore (-) = 0$$



Target Variable	Gender		Total
	Male	Female	
0 (Stays)	50	10	
1 (Churns)	20	20	
Total	70	30	100

Target Variable	Age		Total
	<35	>35	
0 (Stays)	50	10	
1 (Churns)	10	30	
Total	60	40	100

Information Gain, $Gini$ decrease

$$G_C = 1 - \left(\frac{10}{30}\right)^2 - \left(\frac{20}{30}\right)^2$$

$$G_R = 1 - \left(\frac{50}{70}\right)^2 - \left(\frac{20}{70}\right)^2$$

$$G_C = 1 - \left(\frac{1}{30} \right) - \left(\frac{2}{30} \right)$$

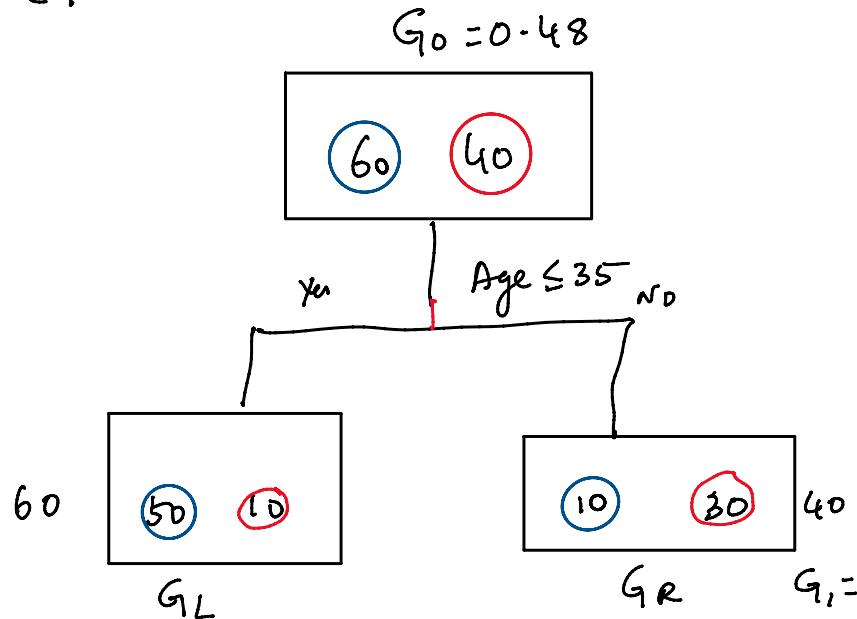
$$= 0.44$$

$$= 0.41$$

-IR- 5/70/ 4/70/

$$\text{Gini decrease} = G_0 - G_1 = 0.48 - 0.42 = 0.06$$

(Gender) = F



$$= 1 - \left(\frac{50}{60} \right)^2 - \left(\frac{10}{60} \right)^2$$

$$= 0.27$$

$$1 - \left(\frac{10}{40} \right)^2 - \left(\frac{30}{40} \right)^2$$

$$= 0.37$$

$$G_1 = G_L \times \frac{\#L}{T} + G_R \frac{\#R}{T}$$

$$= 0.27 \times \frac{60}{100} + \underline{0.37} \times \frac{40}{100} = 0.18$$

$$\text{Gini decrease} = G_0 - G_1 = 0.48 - 0.18 = 0.3$$

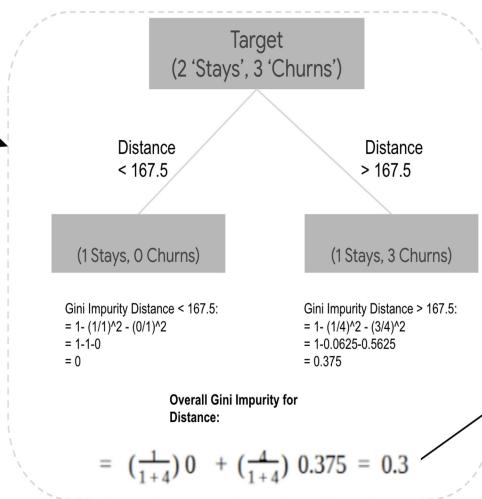
(Age < 35)

Variable	Gini decrease	(Max decrease in Gini)
Gender	0.166	
Age < 35	0.3 ✓	
distance < 167.5	0.21	

① which variable to split

② Split Criteria

Distance	Target
155	0
180	1
190	0
220	1
225	1



Distance	Target	Gini
155	0	0.48
180	1	0.47
190	0	0.48
220	1	0.27
225	1	0.4

distance ≤ 167.5

Gini down

$$0.48 - 0.3$$

$$0.48 - 0.47$$

$$0.48 - 0.27 \downarrow$$

$$0.48 - 0.4$$

distance < 167.5 Gini = 0.3

distance < 180 Gini = 0.4

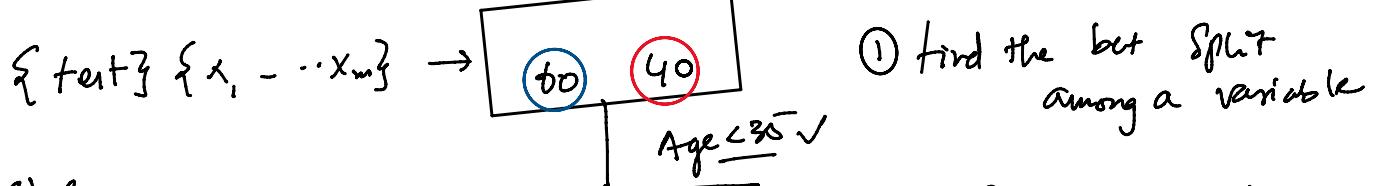
dist < 205 Gini = 0.27 ✓

dist < 225 Gini = 0.4

Decision Tree

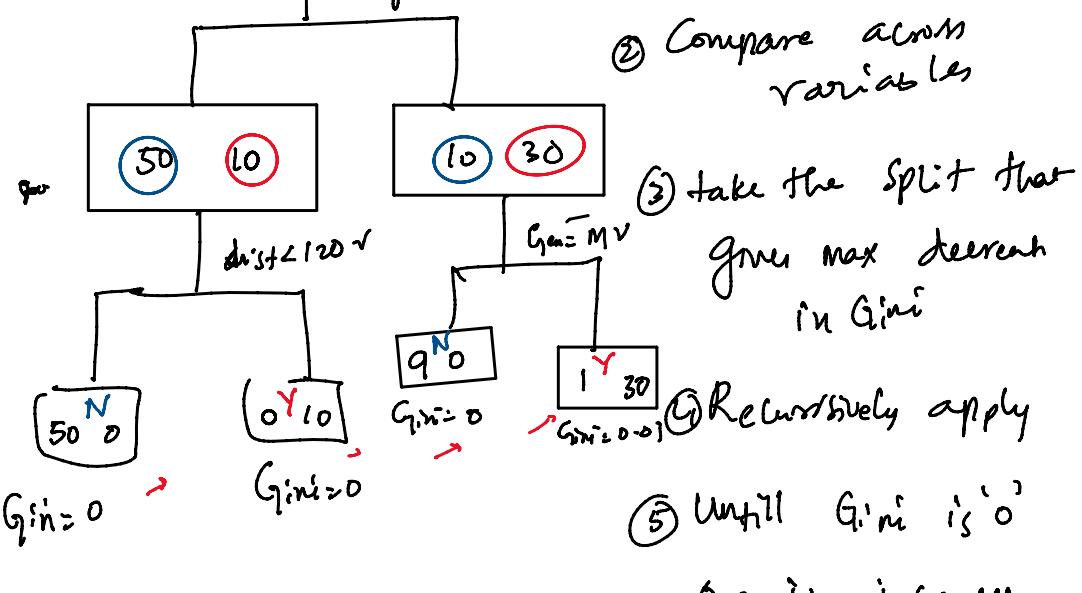
{test} { x_1, \dots, x_m } \rightarrow 60 40

① find the best split variable



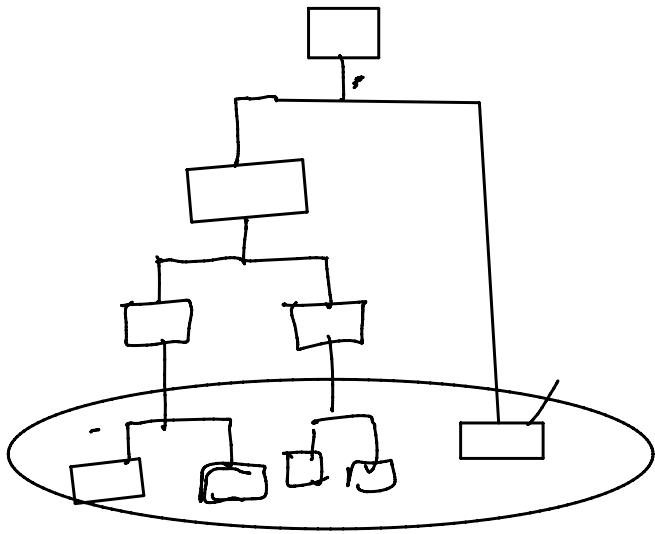
Recursive
[Binary partition]

Terminal
No dev
leaf nodes
decision nodes

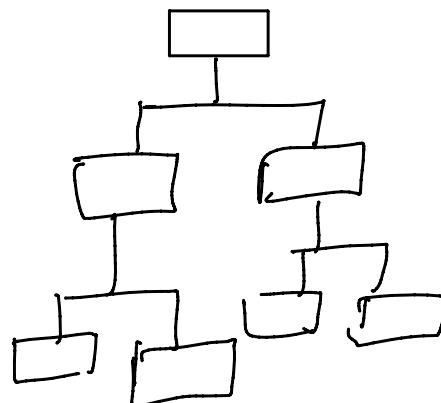


Entropy

Tree 1



Tree 2



Entropy - (Computationally complex) X

Gini - very fast ✓

{0 1 2 3}

Split Criteria

Cont

X sort y

—
—
—
—
—
—
—

Binary

$$x_1 = 0$$

Categorical

$$(A, B, C)$$

$$x_1 = A$$

$$x_1 \notin A, B$$

$$x_1 \in A, C$$

$$A \quad | \quad B \quad C$$

$$\{0, 1, 2, 3\}$$

Ordinal

$$x_1 \leq 1$$

$$x_1 \leq 2$$

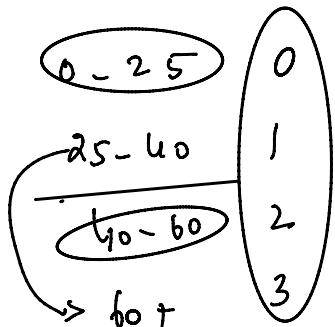
$$x_1 \leq 0$$

$$(0, 2 | 1, 3)$$

label encoder is

$$\begin{matrix} A & B & C & D \\ (0 & 1 & 2 & 3) \end{matrix}$$

Age



$$0, 1 \quad | \quad 2, 3$$

$$0, 2 \quad | \quad 1, 3$$

$$0, 1, \quad 2, 3$$

$$\{62, 13\}$$

Root node

33 → (X) inputs

① Greedy ✓

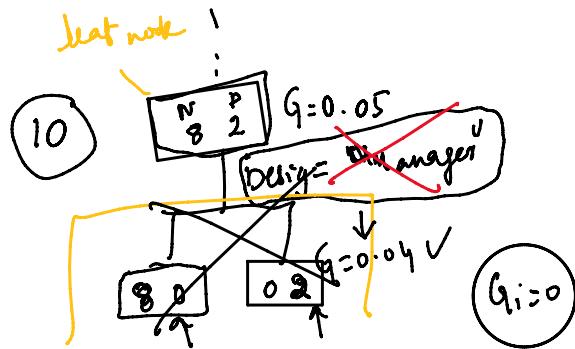
② Overfitting



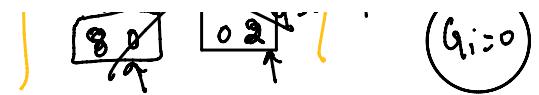
Mean decrease in Gini

① Pruning \leftarrow Post \rightarrow Pre

Stop tree growth if any of

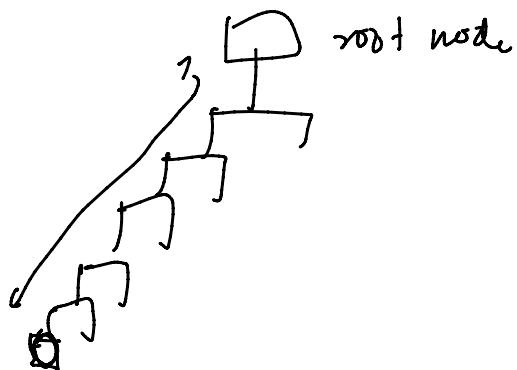


Stop tree growth if any of
these conditions are satisfied

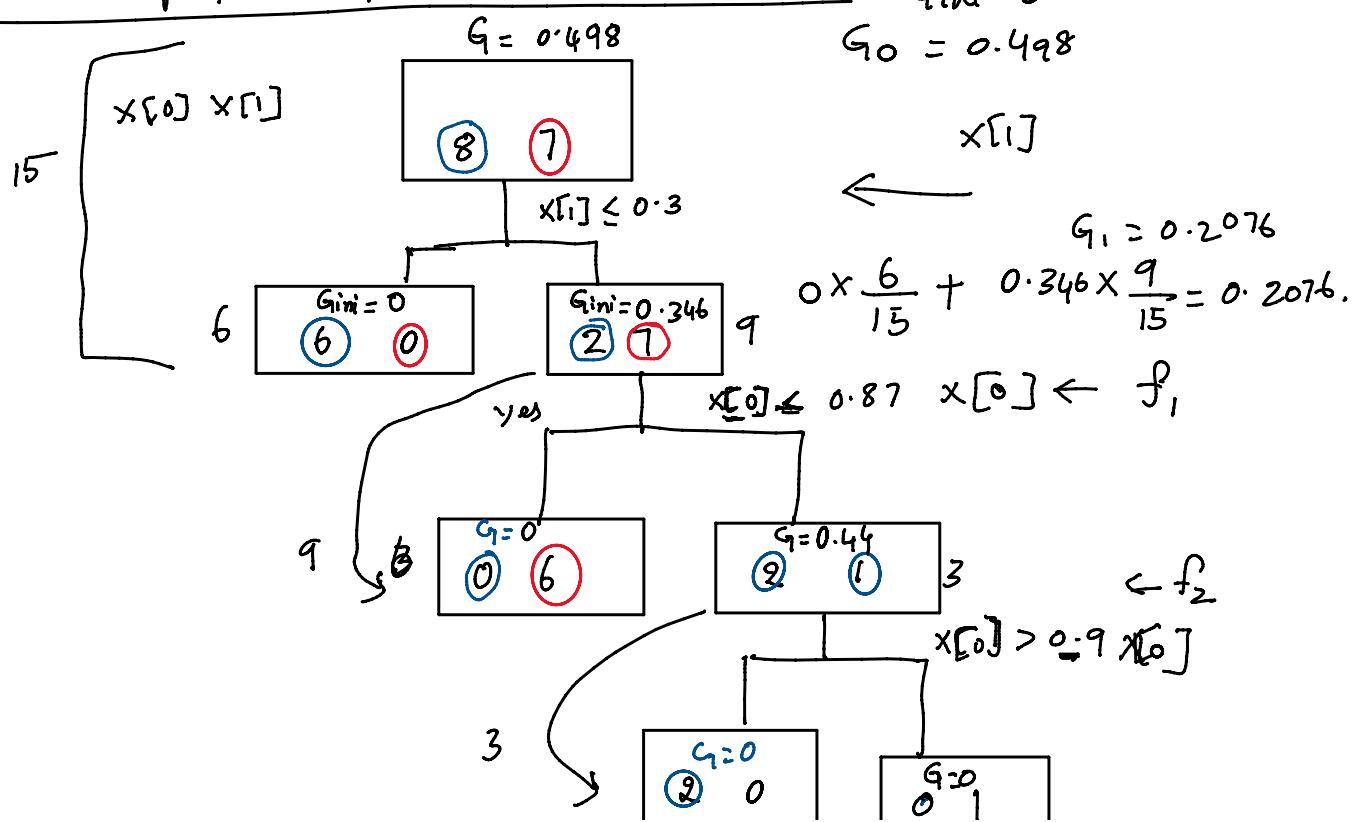


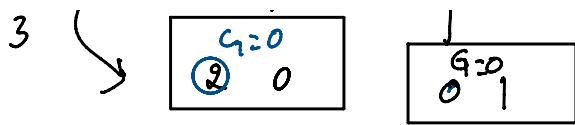
$$G_{i=0}$$

- ① if Max depth is reached Max depth = 10
 - ② if no. of samples in a node is less than 20
 - ③ Min no. of samples in a leaf node = 5
- will stop tree growth
reduce overfitting.



feature importance from a decision tree





Variable Importance of $X[1]$

$$= 0.498 - 0.2076 = 0.2904$$

$$= 0.2904 \times \frac{15}{15} = 0.2904$$

Variable Importance $X[0] f_1 G_0 = 0.346$

$$G_0 \quad (G_L + G_R)_{W..}$$

$$0.346 - 0.44 \times \frac{6}{9} + 0.44 \times \frac{3}{9} = 0.1981$$

$$f_2: \quad G_0 = 0.44$$

$$0.44 - 0 = 0.44$$

$$0.1981 \times \frac{9}{15} + \frac{3}{15} \times 0.44 = 0.2074$$

$$\boxed{X[1]} = 0.2904 = \frac{0.2904}{0.2904 + 0.2074} = \underline{\underline{0.5833}}$$

$$\underline{X[0]} = 0.2074 = \frac{0.2074}{0.2074 + 0.2904} = \underline{\underline{0.4166}}$$

3, 50L $\overbrace{(u=25)}$ $\overbrace{8 - \frac{25^2}{(4-u)^2} (5-1-u)^2 + (50-25)^2}$

$$3, 50L \quad (\mu = 25)$$

$$\frac{2}{3} - \frac{2}{25} \times (4-25)^2 + (5-1-25)^2 + (50-25)^2$$

Regression

X
Age, deg, Gen, dep + 3

