

Decision Tree Regression

O/P is continuous feature.

→ Exp

→ 2

→ 2.5

→ 3

→ 4

4.5

Dataset

Career

Gap

Op

Salary

40K

42K

52K

60K

56K

[40K, 42K, 52K, 60K, 56K]

58K

60.8

≤ 2

Yes

No

40K

42, 52, 60, 56

$V(c1) = 100$

$V(c2) = 51$

Variance Reduction (Regression Problem)
It determines which value should select to split.

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

Average

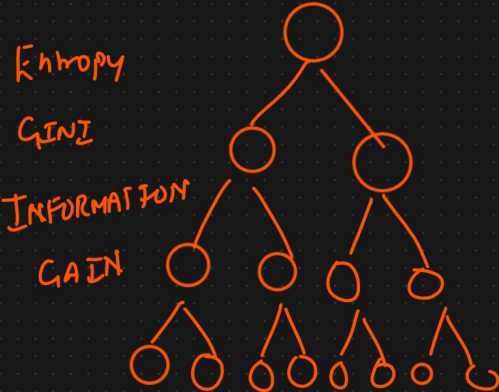
Variance & Mean Squared Error is same in regression.

$$\text{Variance (Root)} = \frac{1}{5} \left[(40-50)^2 + (42-50)^2 + (52-50)^2 + (60-50)^2 + (56-50)^2 \right]$$

$$= \frac{1}{5} [100 + 64 + 4 + 100 + 36]$$

$$= 60.8$$

$$\text{Variance}(c1) = \frac{1}{n} \sum_{i=1}^n (y_i - \bar{y})^2$$



CART

[40K, 42K, 52K, 60K, 56K]

≤ 2.5

60.8

Decision Tree Regression

$\text{Var}(c1) =$

$$\frac{1}{2} \left[(40-50)^2 + (42-50)^2 \right]$$

$$= \frac{1}{2} [100 + 64] = \frac{164}{2} = 82$$

$$\text{Var}(child2) = \frac{1}{3} [4 + 100 + 36] = \frac{140}{3} = 46.66$$

Variance Reduction

$$= 60.8 - \left[\frac{2}{5} \times 82 + \frac{3}{5} \times 46.66 \right]$$

$$= 60.8 - [32.8 + 27.996]$$

$$= +0.004$$

$$= \frac{1}{1} (40-50)^2$$

$$= 100$$

$$\text{Variance (c2)} = \frac{1}{n} \sum_{i=1}^n (y_i - \bar{y})^2$$

$$= \frac{1}{4} \left[(42-50)^2 + (52-50)^2 + (60-50)^2 + (56-50)^2 \right]$$

$$= \frac{1}{4} [64 + 4 + 100 + 36]$$

$$= 51$$

Number of samples in child node.
Total number of samples in parent node.

Variance Reduction

$$= \text{Var}(\text{Root}) - \sum_{i=1}^k w_i \text{Var}(\text{child})$$

$$= 60.8 - \left[\frac{1}{5} \times 20 + \frac{4}{5} \times 51 \right]$$

$$= 60.8 - 20 - 40.8$$

$$\text{Variance Reduction} = 0$$

$$\text{Variance Reduction (Left Split)} < \text{VR (Right Split)}$$

Left Tree



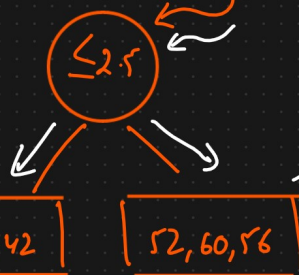
Leaf Node

$$\text{O/P} = \left[\frac{40+42}{2} \right]$$

$$= 41$$

0.004

Test data



Decision Tree
O/P Reg.

$$\frac{52+60+56}{3} = 56$$