

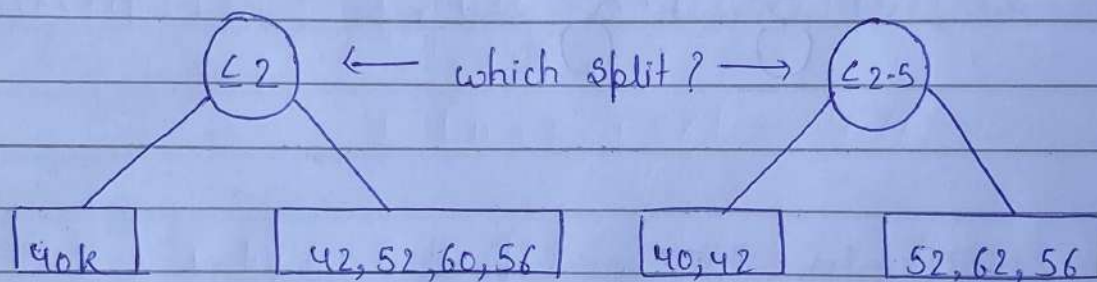
# Decision Tree Regressor

→ Dataset :

I/P		O/P
Exp	Gap	Salary
→ 2	Yes	40k
2-5	Yes	42
3	No	52
4	No	60
4-5	Yes	56

↓  
Sorted

$\hat{y} = 50k$



Note: In case of continuous feature, variance reduction is used to select split to choose.

\* Variance Reduction :

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

↓  
Average of o/p

$$VR = \text{Var}(\text{root}) - \sum w_i \text{Var}(\text{child})$$

where, Var → Variance

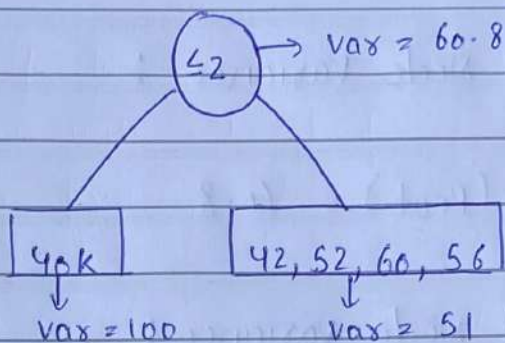
w → weightage →  $\frac{\text{child Node count}}{\text{Root Node count}}$

→ Selecting the split :

# split 1 :

• Root Node Variance :

$$\text{Var}(\text{Root}) = \frac{1}{n} \sum (y - \hat{y})^2$$



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

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

$$= \frac{304}{5} = 60.8$$

• Child Node Variance :

$$\text{Var}(\text{child 1}) = \frac{1}{1} \times [(40-50)^2] = 100$$

$$\text{Var}(\text{child 2}) = \frac{1}{4} \times [(42-50)^2 + (52-50)^2 + (60-50)^2 + (56-50)^2]$$

$$= \frac{1}{4} \times [64 + 4 + 100 + 36] = 51$$

• Variance Reduction of split 1 :

$$\text{VR}_1 = \text{Var}(\text{Root}) - \sum w_i \text{Var}(\text{child})$$

$$= 60.8 - \left[ \frac{1}{5} \times 100 + \frac{4}{5} \times 51 \right]$$

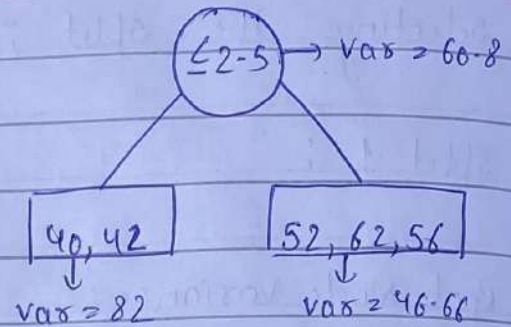
$$\text{VR}_1 = 0$$



# Split 2 :

- Root Node Variance :

$$\text{Var}(\text{Root}) = 60.8$$



- Child Node Variance :

$$\text{Var}(\text{child 1}) = \frac{1}{2} \times [(40-50)^2 + (42-50)^2]$$

$$= \frac{164}{2} = 82$$

$$\text{Var}(\text{child 2}) = \frac{1}{3} \times [(52-50)^2 + (62-50)^2 + (56-50)^2]$$

$$= \frac{1}{3} \times [4 + 100 + 36] = 46.66$$

- Variance Reduction of Split 2 :

$$\text{VR}_2 = 60.8 - \left[ \frac{2}{5} \times 82 + \frac{3}{5} \times 46.66 \right]$$

$$\text{VR}_2 = 0.304$$

$$\text{VR}_2 > \text{VR}_1$$

$\therefore$  Split 2 will be selected for splitting.

Note :

