

Decision Tree Regression

Data set

O/P

Exp	Gap	Salary
2	Yes	40 K
2.5	Yes	42 K
3	No	52 K
4	No	60 K
4.5	Yes	56 K

here the O/P feature is continuous, hence we will have to use regression

In decision tree we cannot use

- Entropy
- Gini Impurity
- Information gain



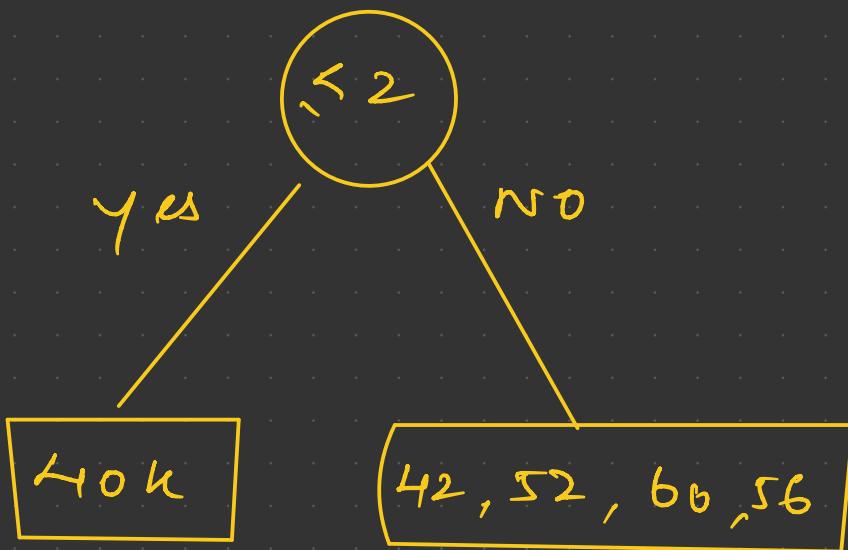
In order to use this output has to be fixed number of categories, but in case of regression output will be continuous.

Data Set

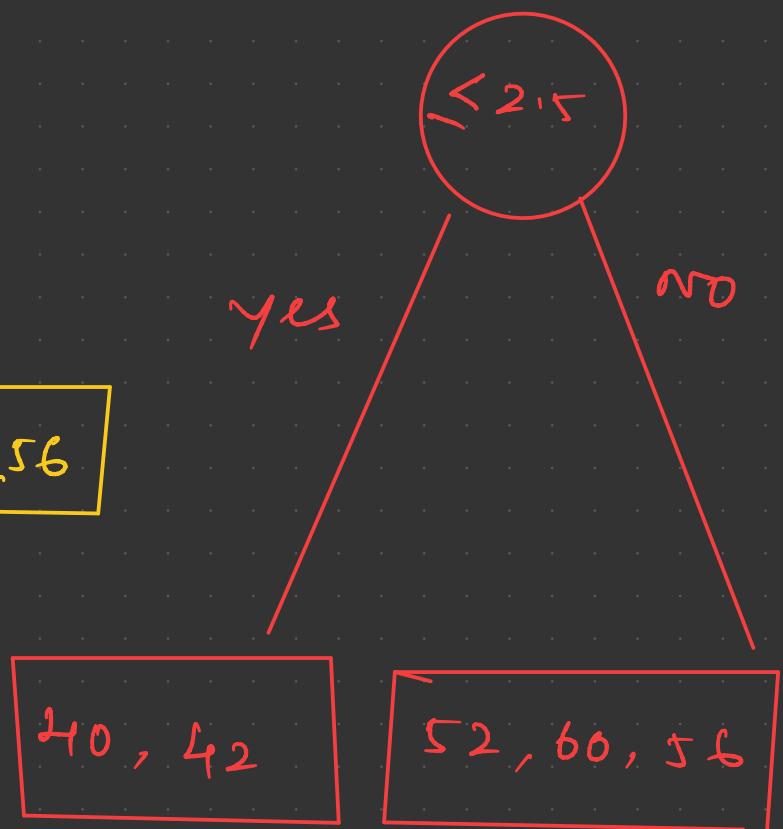
O/P

Exp	Hab	Salary
2	yes	40 K
2.5	yes	42 K
3	no	52 K
4	no	60 K
4.5	yes	56 K

Threshold : 2



Threshold : 2.5



how we select, which split needs to be deleted?

→ we will be using Variance Reduction to select the best split for the regression Problem.

In Variance Reduction, we need to

① compute the variance,

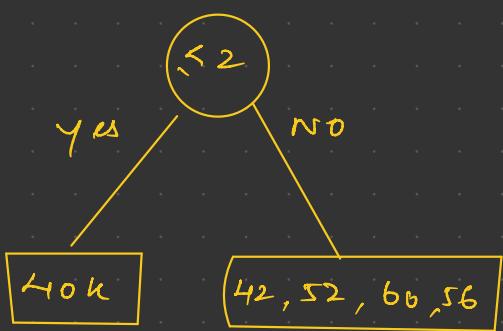
$$\text{Variance} = \frac{1}{n} \sum_{i=1}^n (y - \bar{y})^2$$

\hookrightarrow average.

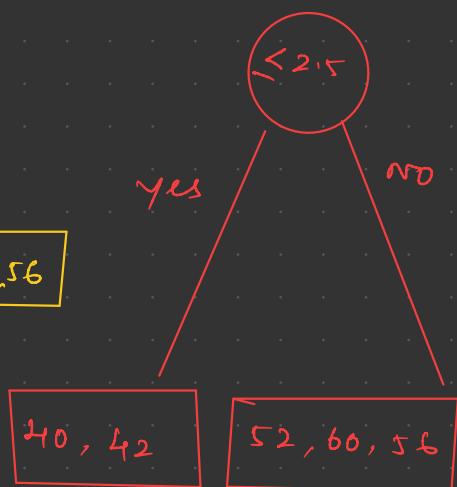
② we will be finding the mean square error for each and every node and combine them to find out the variance reduction

③ Select the split which has maximum variance reduction.

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



$[40, 42, 32, 60, 56]$



$$\bar{Y} = 50 \text{ k}$$



mean of o/p

Variance of root:

$$\text{Variance}(\text{root}) = \frac{1}{5} \left[(40 - 50)^2 + \right.$$

$$(42 - 50)^2 + (52 - 50)^2 + (60 - 50)^2 \\ + (56 - 50)^2 \left. \right]$$

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

$$= 60.8$$

Variance of child 1

$$\text{Variance } (c_1) = \frac{1}{n} \sum_{i=1}^n (y - \bar{y})^2$$

$$= \frac{1}{4} (40 - 50)^2$$

$$\text{Variance } (c_1) = 100$$

Variance of child 2

$$\text{Variance } (c_2) = \frac{1}{n} \sum_{i=1}^n (y - \bar{y})^2$$

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

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

$$\text{Variance } (c_2) = 51$$

Variance Reduction

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

$w_i \rightarrow$ Ratio of the split

$$\begin{aligned} &= 60.8 - \left[\frac{1}{5} \times 100 + \frac{4}{5} \times 51 \right] \\ &= 60.8 - [20 + 40.8] \end{aligned}$$

$\Rightarrow 0 \rightarrow$ Variance reduction of first split

Similarly, after computing the variance reduction of the second split.

Variance Reduction of split 2 is

0.004.

hence we will be choosing the split 2 for Decision Tree Regression.

For example

