

# Decision Tree Regressor

DT Classifier

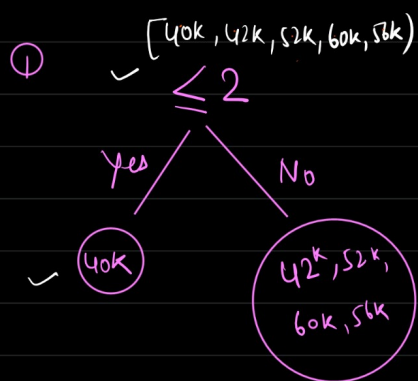
→ Entropy / Gini  
→ Information gain

DT Regressor

→ Variance.  
→ Variance reduction

Experience	Career Gap	Salary
2	Yes	40k
2.5	Yes	42k
3	No	52k
4	No	60k
4.5	Yes	56k

$$\bar{y} = 50k$$



\* final aim is variance reduction.

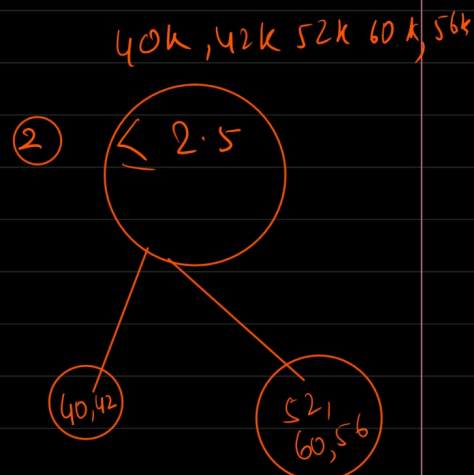
\* you will split whichever feature / feature value gives the highest variance reduction.

$$\text{Variance} = \frac{1}{n} \sum_{i=1}^n (y_i - \bar{y})^2 \quad [\text{Mean Squared Error}]$$

$$\begin{aligned} \text{Variance}(P(\text{root})) &= \frac{1}{5} \left( (40-50)^2 + (42-50)^2 + (52-50)^2 + (60-50)^2 + (56-50)^2 \right) \\ &= \frac{1}{5} (100 + 64 + 4 + 100 + 36) \\ &= 60.8 \end{aligned}$$

$$\begin{aligned} \text{Variance}(L.C) &= \frac{1}{1} (40-50)^2 \\ &= 100 \end{aligned}$$

$$\begin{aligned} \text{Variance}(R.C) &= \frac{1}{4} (8^2 + 2^2 + 10^2 + 6^2) \\ &= \frac{1}{4} (64 + 4 + 100 + 36) \\ &= 51 \end{aligned}$$



$$\text{Var}(P(\text{root})) = 60.8$$

$$\text{Var}(\text{left child}) = \frac{1}{2} (100 + 64) = 82$$

$$\text{Var}(\text{right child}) = \frac{1}{3} (4 + 36 + 100) = 46.66$$

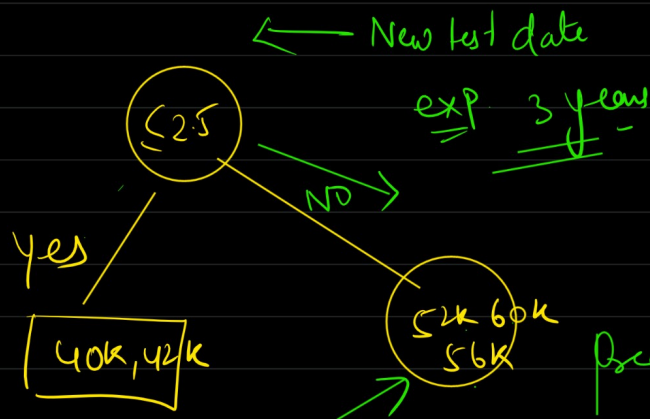
$$\begin{aligned}\text{Variance Reduction} &= \text{Varp} - \text{Vchild combined} \\ &= 60.8 - \left( \frac{1}{5} \times 100 + \frac{4}{5} \times 51 \right) \\ &= 0\end{aligned}$$

$$\begin{aligned}60.8 - \left( \frac{2}{5} \times 82 + \frac{3}{5} \times 46.6 \right) \\ = 0.004\end{aligned}$$

Variance reduction  
for 2

Variance Reduction  
for 2.5

Select this  
split of the  
feature.



$$\begin{aligned}\text{Pred for 3 years} &= \frac{52 + 60 + 56}{3} \\ &= 56k.\end{aligned}$$