

$$\frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2$$

## Multiple Linear Regression

House pricing independent feature  
No. of Rooms      Size      Location      Z

$$h_\theta(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3$$

→ multiple linear regression

$\theta_1, \theta_2, \theta_3$  are Coefficient

$\theta_0$  is intercept



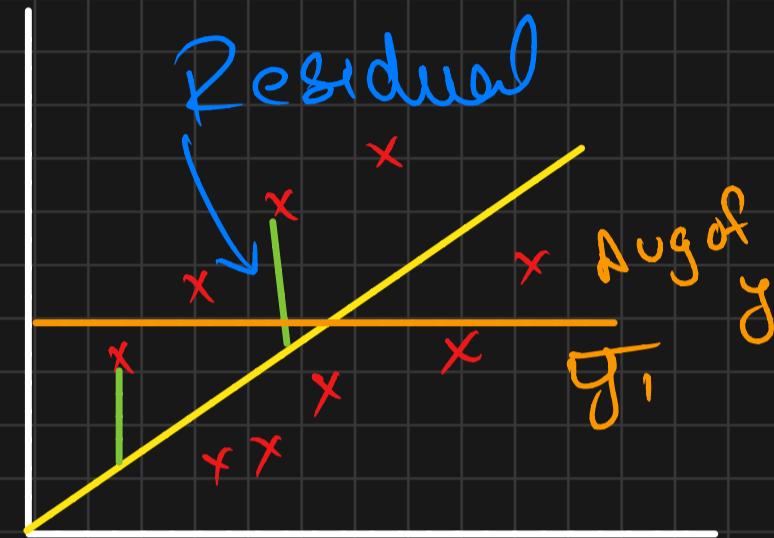
# Performance Matrix

① R squared

② Adjusted R squared

$$R_{sq} = \frac{1 - SS_{\text{residual}}}{SS_{\text{total}}}$$

$$R_{sq} = \frac{1 - \sum (y_i - \hat{y}_i)^2}{\sum (y_i - \bar{y}_i)^2}$$



$$1 - \frac{\text{Smaller no.}}{\text{Bigger no.}} \approx 1$$

- Value more near to 1 the model performs better well

② Adjusted R squared

Size

Price  $(R)^2 \propto \text{Size} \propto \text{Price}$   
• true correlation

$$R_{sq} = 0.75$$

↓  
Size

no. of Room Price  $\Rightarrow \text{Size} \propto \text{no. of Room}$   
 $\propto \text{Price}$

$$R_{sq} = .80 \text{ (80%)} \downarrow$$

Size      no. of Room      Location      Price

$$R_{sq} = .85 \text{ (85%)} \downarrow$$

Size      no. of Room      Location      gender      Price

$$R_{sq} = .87 \text{ (87%)} \downarrow$$

$\Rightarrow R_{sq}$  value increased but gender & Price does not relate with each other

$\Rightarrow$  adjusted  $R_{sq}$

$$\text{Adj. } R_{sq} = 1 - \frac{(1 - R^2)(N-1)}{N-p-1}$$

N = No. of data point

P = No. of independent feature

$\Rightarrow$  MSE, MAE, RMSE

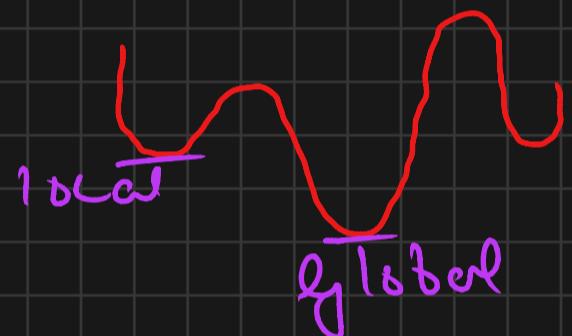
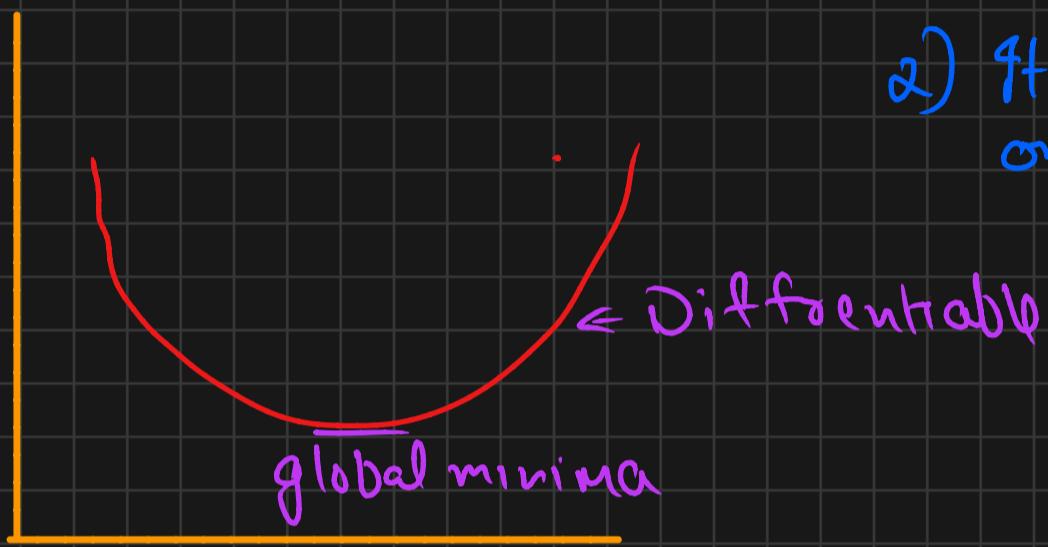
Exp	Salary



$$\text{Mean Square Error} = \frac{1}{n} \sum_{i=1}^n (y - \hat{y})^2 = \text{cost fn} \downarrow$$

advantage

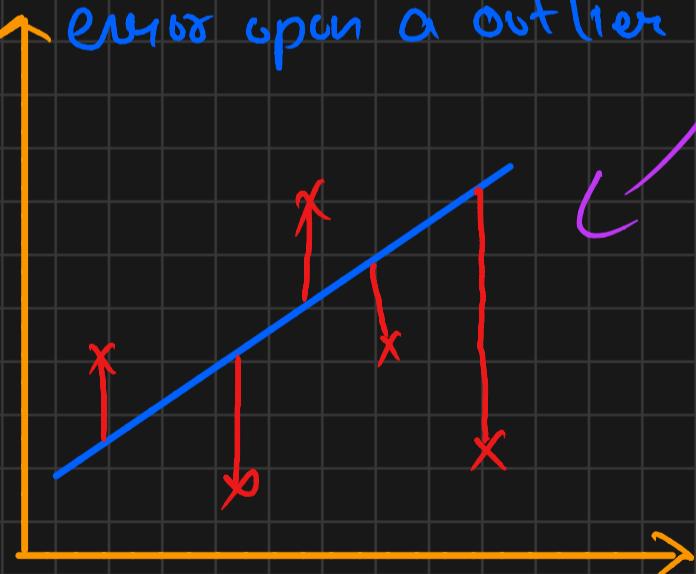
- 1) it is differentiable
- 2) it has one local & one global minima



Disadvantage

- 1) it is not Robust

it will increase error upon an outlier



if no global it will converge at local



$$2e \quad y_i - \hat{y}_i \quad (\text{lakhs})^2 \quad (\text{Units changed})$$

## Mean Absolute Error

$$\frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i|$$

Disadvantage

- Not differentiable  
(we need sub gradient  
⇒ convergence take  
more time)
- Time Consuming

## Advantage

- Robust to outliers
- $|y_i - \hat{y}_i|$  will be in  
same units

## Root Mean Squared Error

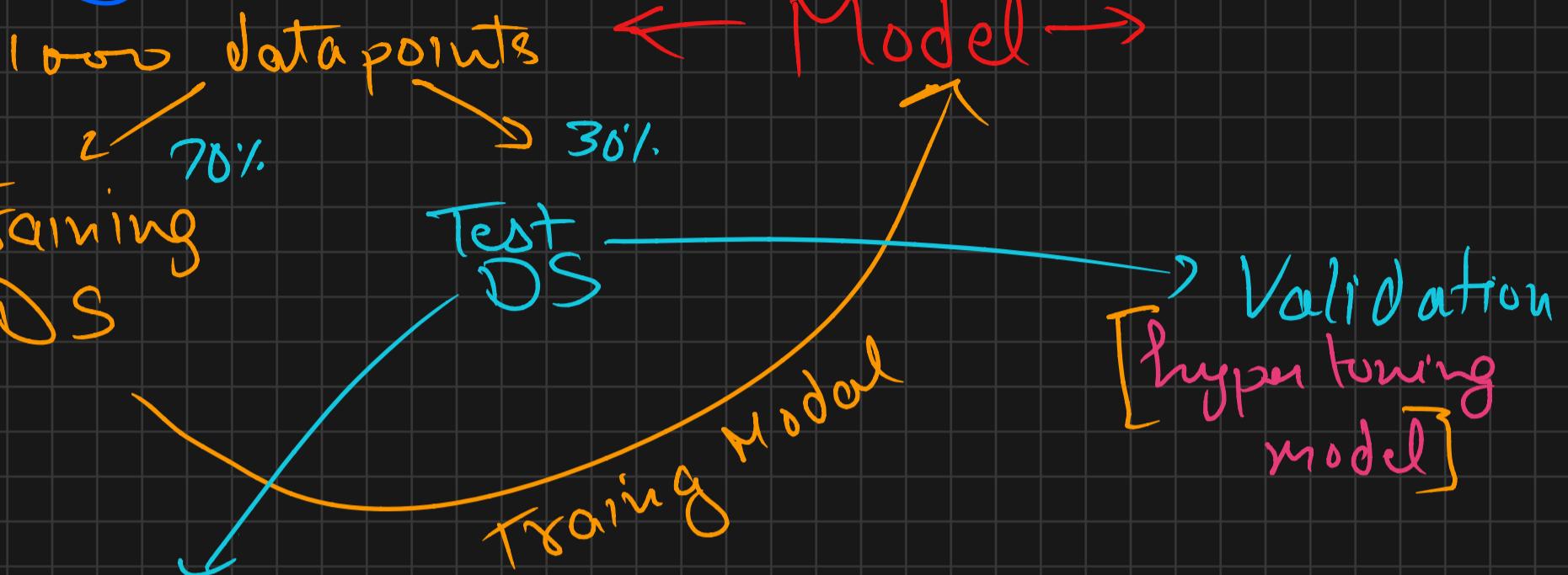
Advantage

- ① Same unit
- ② Differentiable

⇒ We should consider all these will  
Creating Regression Model

# Overfitting & Underfitting

DataSet



Train

low bias & low variance

Train DS

good accuracy

good accuracy

Test DS good accuracy

Bad accuracy

$\Rightarrow$  low bias high variance

model is overfitting

Train accuracy is low

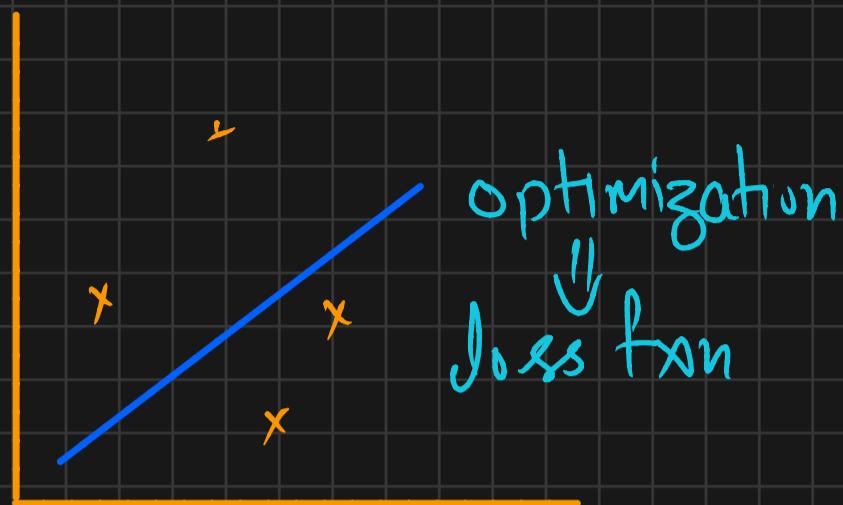
Bias Variance Tradeoff

Test accuracy is low

$\Rightarrow$  high bias high variance

model is underfitting

# Linear Regression Using OLS



OLS = formula to calculate  
Reduce error  $h_{\theta}(x) = \beta_0 + \beta_1 x$ ,

$$S(\beta_0, \beta_1) =$$

$$\frac{1}{n} \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i)^2$$

$$\frac{\partial S(\beta_0, \beta_1)}{\partial \beta_0} = \frac{2}{n} \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i)$$

$$-2 \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i) \leftarrow$$

$$\frac{\partial S(\beta_0, \beta_1)}{\partial \beta_1} = \frac{2}{n} \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i) (-x_i) = 0$$

$$-2 \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i) (x_i) = 0$$

$$\frac{1}{n} \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i) = 0$$

$$-\sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i) = 0$$

$$\boxed{\beta_0 = \bar{y} - \beta_1 \bar{x}}$$

Eq(2)

$$\frac{-2}{n} \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i)(x_i) = 0$$

$$\sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i)(x_i) = 0$$

$$\sum_{i=1}^n x_i y_i - \beta_0 \sum_{i=1}^n x_i - \beta_1 \sum_{i=1}^n (x_i)^2 = 0$$

$$\beta_1 = \frac{\sum_{i=1}^n (y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})}$$