

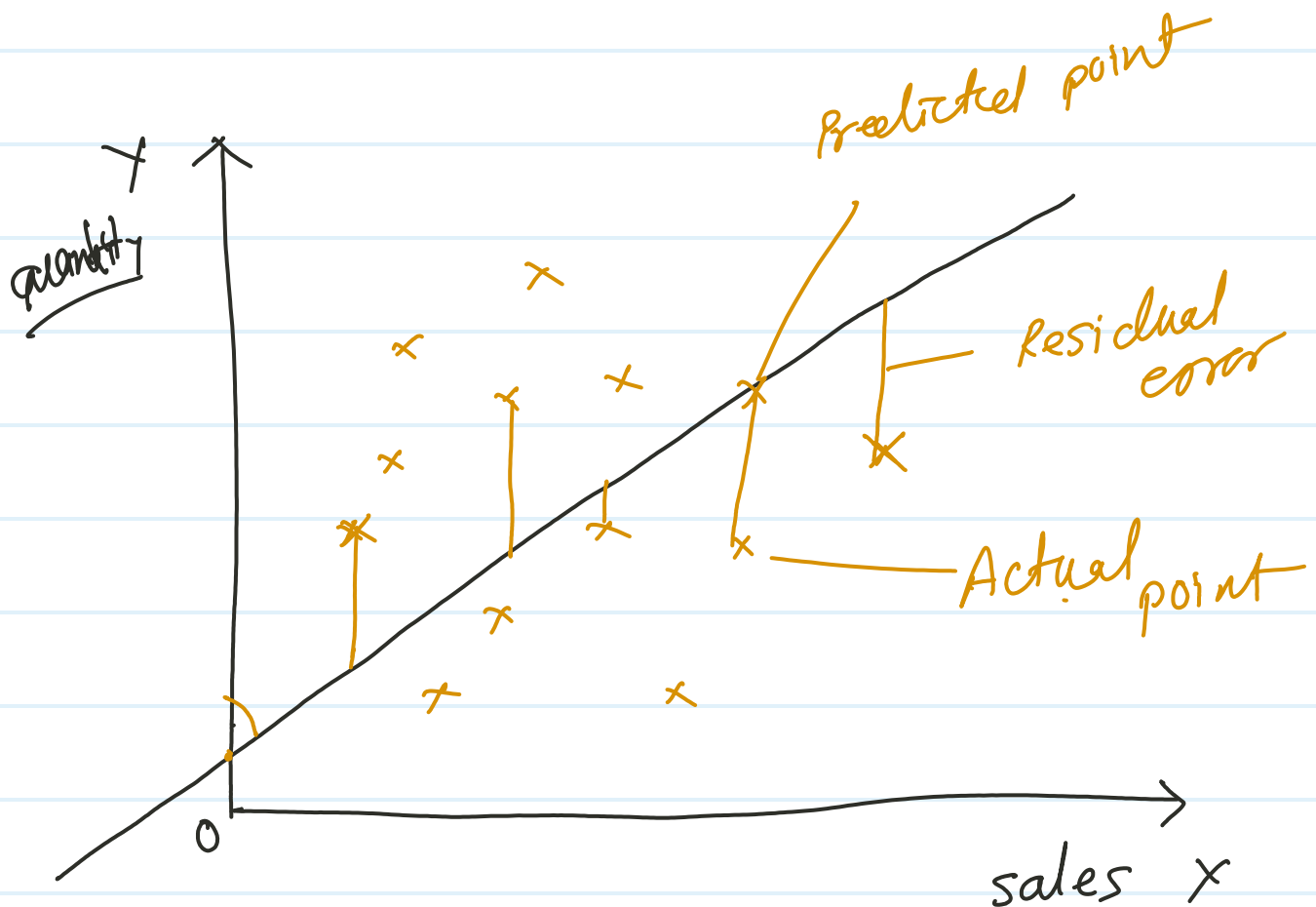
# \* Linear Regression

$$y = mx + c$$

$m$  = slope or coeff.

$x$  =

$c$  = intercept ( $0, x=0$ )



Residual error  $(y - \hat{y})$

To Find best fit line with minimal error.

$$y = mx + c$$

$$h_{\theta}(x) = \hat{y}$$

single linear Repr.

$$h_{\theta}(x) = \theta_0 + \theta_1 x$$

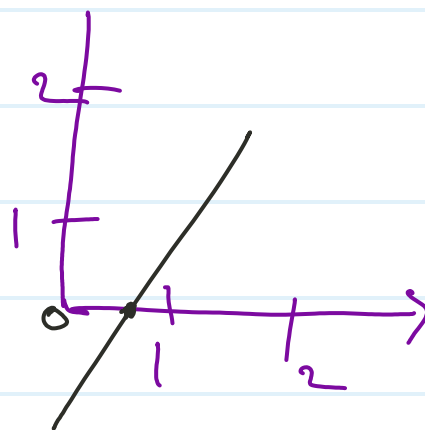
multipoint linear Regress

$$h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n$$



$$\theta_0 = 1.5$$

$$\theta_1 = 0$$



$$\theta_0 = 0$$

$$\theta_1 = 0.5$$

$$\hat{y} = 0 + (0.5) \times 1 \Rightarrow 0.5$$

$$\hat{y} = 0 + (0.5) \times 2 \Rightarrow 1$$

$$\hat{y} = 0 + (0.5) \times 3 \Rightarrow 1.5$$

\* Cost function =

$$J(x) \quad J(\theta_0, \theta_1)$$

$$J(\theta_0, \theta_1) = \frac{1}{n} \sum_{i=1}^m (h_{\theta}(x) - y)^2$$



Repeat convergen theorem

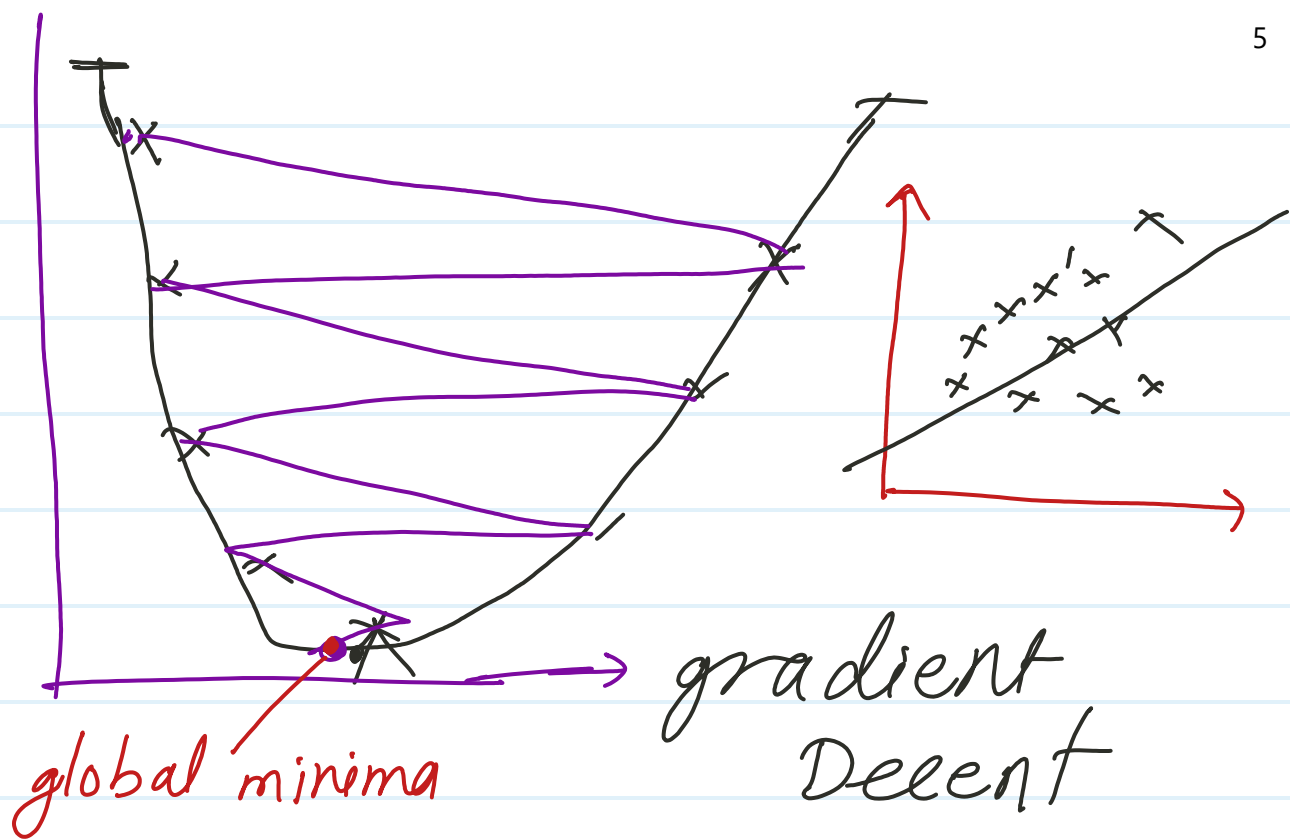
$$\theta_j = \theta_j - \alpha \frac{d}{d\theta_j} (J(\theta_j))$$

$\alpha$  = learning rate

0.05, 0.01

$$\theta_0 = \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^i) - y^i)$$

$$\theta_1 = \theta_1 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^i) - y^i) x^i$$



① MSE (mean square error)

② RMSE (Root mean square error)

③ MAE (mean Absolute error)

① MSE

$$MSE = \frac{\sum_{i=1}^n (Y - \hat{Y})^2}{n}$$

It create global minima

## ② RMSE

$$RMSE = \sqrt{\frac{1}{n} \sum_{i=1}^n [\{y - (\theta_0 + \theta_1 x)\}^2]}$$

It is not robust to outlier.  
create local minima

## ③ MAE

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

pros. Robust to outlier.

cons. - It take usually more time to optimization.

\* performance matrix

7

$R^2$  statistics  $\rightarrow$

$$R^2 = 1 - \frac{RSS}{TSS}$$

$RSS$  = Sum of squ. of residuals

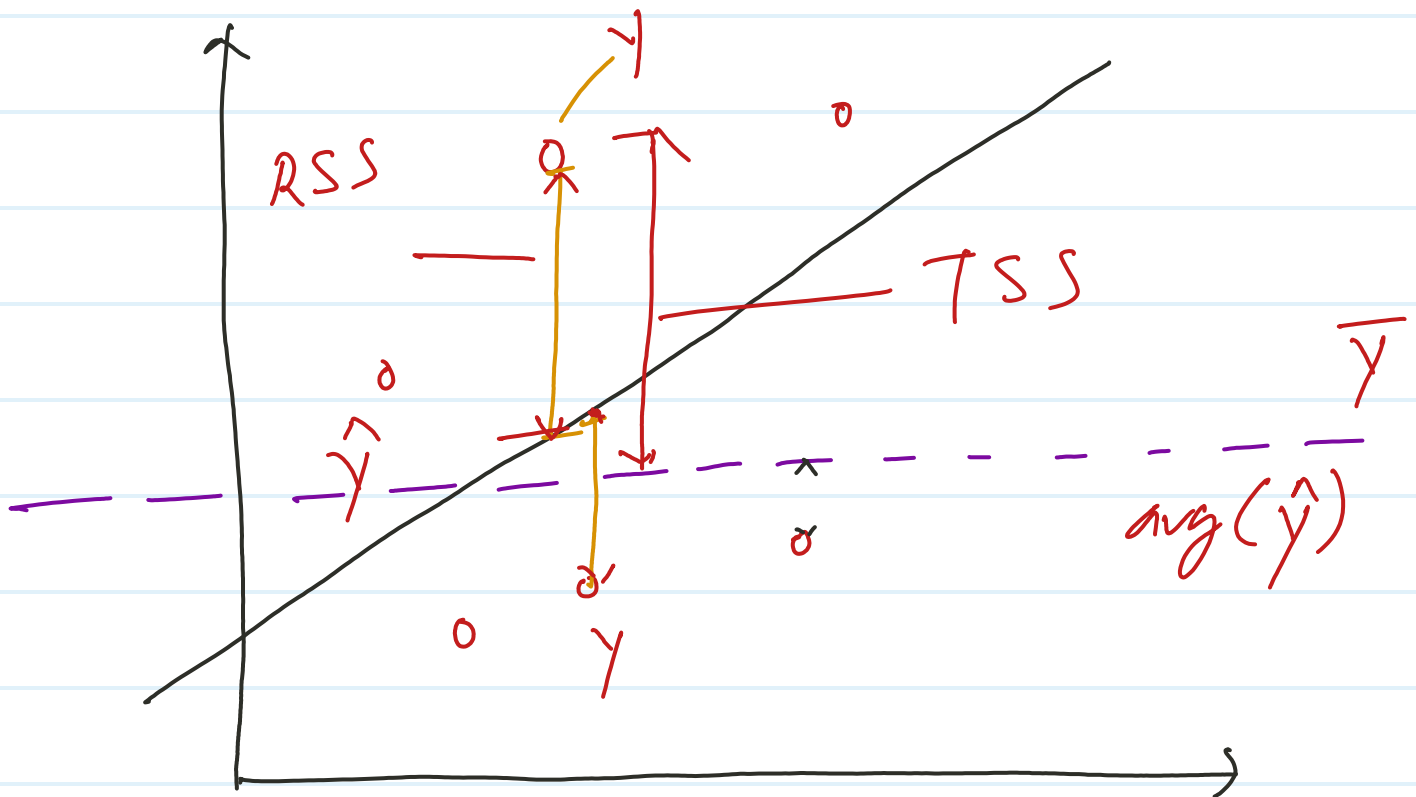
$TSS$  = Total sum of squ.

$RSS$  = Dist. b/w  $y$  and  $\hat{y}$

$TSS$  = Dist. b/w  $y$  and  $\text{avg}(\hat{y})$

$$RSS = \sum (y - \hat{y})^2$$

$$TSS = \sum (y - \bar{y})^2$$



\* Adjusted  $R^2$  - statistics

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

$N$  = number of datapoint in dataset  
 $p$  = number of independent variable



# overfitting and underfitting

## \* overfitting

$\Rightarrow$  Train = 90%

$\Rightarrow$  Test = 10%

high variance  
low biased

## underfitting

$\Rightarrow$  Train = 50%

$\Rightarrow$  Test = 90% / 40%

low variance  
high biased

1—  
2  
3—  
4—  
5—  
6  
7  
8  
9  
10

6 train

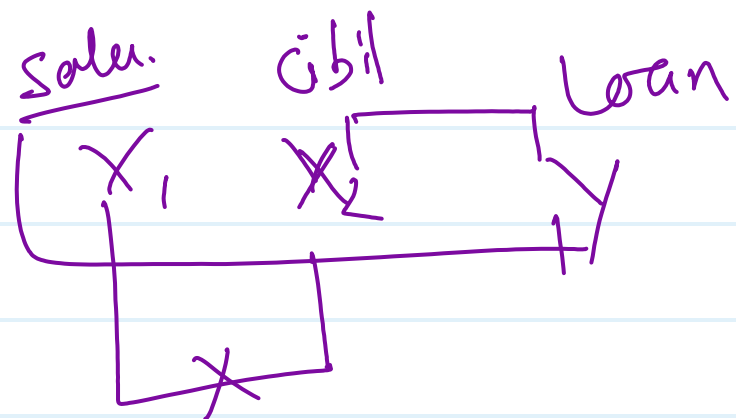
4 test

## Best fit model

low variance  
low biased

## Important Assumption of LR

- ① There should be linear relationship b/w dependent and independent variable.
- ② Error term are not suppose to co-related.
- ③ Ind. variable ( $x$ ) and residual error should be uncorrelated.
- ④ no - multicollinearity



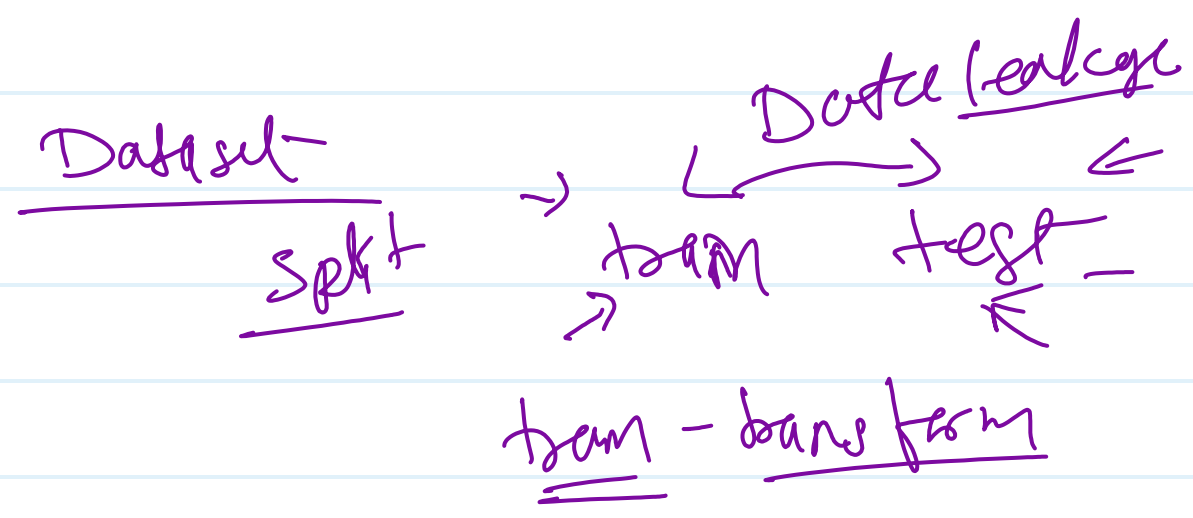
$X_1 \quad X_2 \quad X_3 \quad X_4$

line eqn  $\checkmark$   $h_0(x) = \theta_0 + \theta_1 x_1$

Cost function  $\leftarrow$   $\begin{matrix} \text{MSE} \\ R^2 \\ \text{MAE} \end{matrix}$

evaluation / Performance metrics

$R^2$  or Adj  $R^2$



$R^2$

