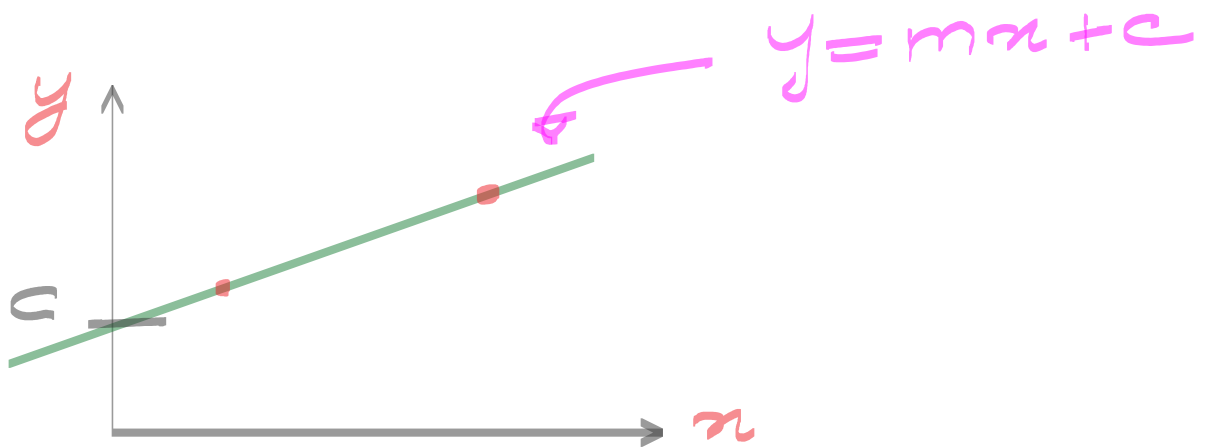


* Equation of Line .

$$y = mx + c$$

$$y = a + bx$$

$$y = \theta_0 + \beta \theta_1$$



$m \rightarrow$ slope of line
 $c \rightarrow$ intercept

$$y = 5x + 3$$

If $x = 2$, $y = 13$

$x = 5$, $y = 28$

* Types of LR

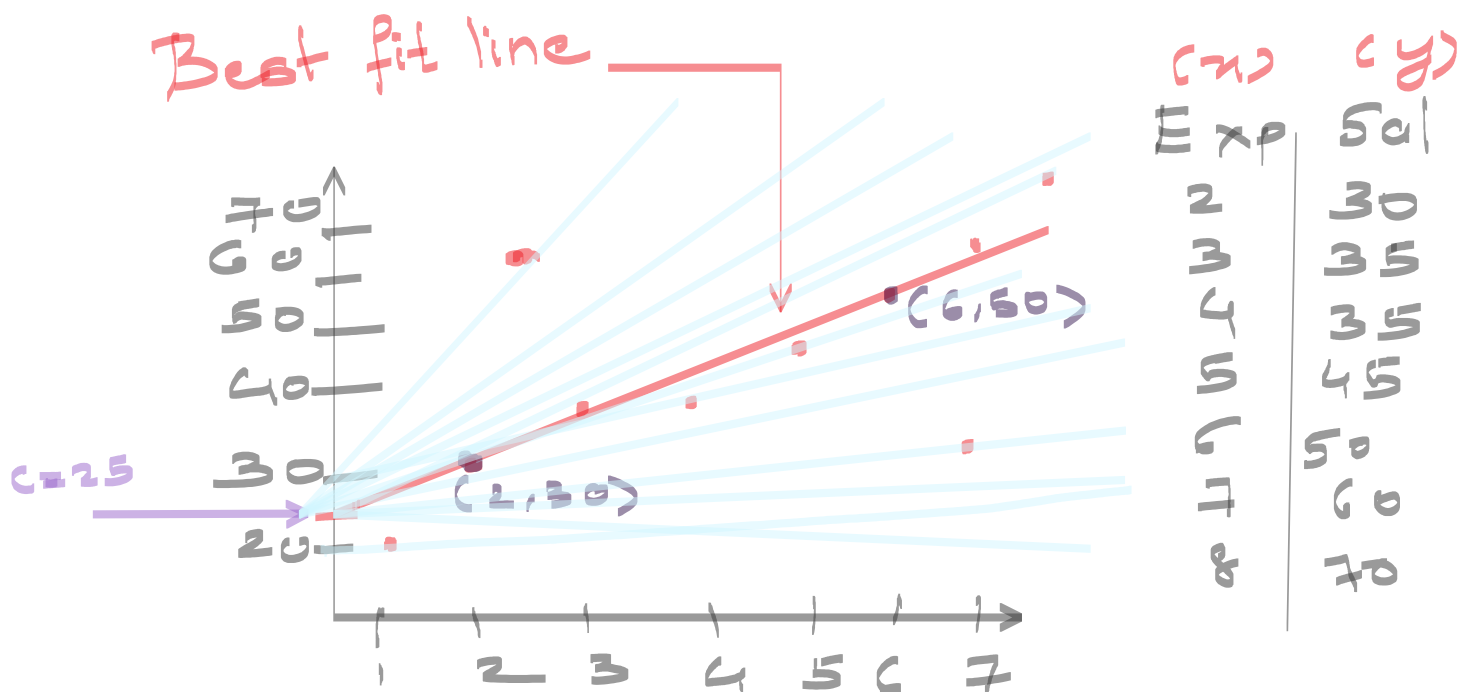
1. Simple LR

1. Simple LR

$$y = mx + c$$

2. Multiple LR

$$y = m_1x_1 + m_2x_2 + \dots + m_nx_n + c$$



$$m = \frac{y_2 - y_1}{x_2 - x_1} = \frac{50 - 30}{6 - 2} = \frac{20}{4} = 5$$

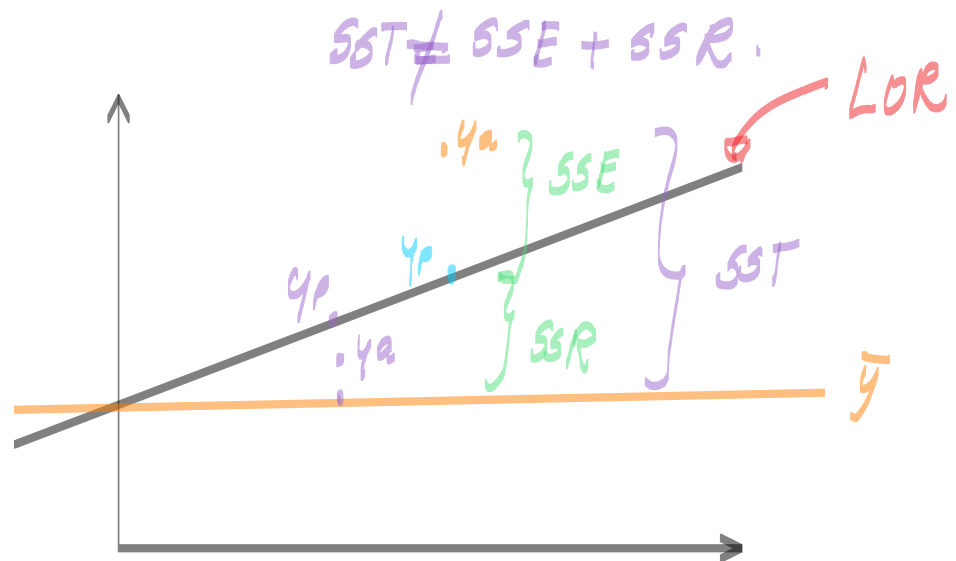
$$y = 5x + 25$$

Linear Model

* Blue lines → infinite number of possibilities

Linear Regression

08 December 2022 07:31



MSE \rightarrow Scale Variant

Ht.	Temp.
5.1	150
5.2	160
5.4	170
6	183
6.1	185

MSE 25-36

MSE 28900

R-squared $\rightarrow R^2 \rightarrow$ Scale Invariant

0 to 1, -ve

1. $R^2 = 1$

$$R^2 = 1 - \frac{SSE}{SST}$$

$$= 1 - \frac{0}{SST} = 1$$

2. $R^2 = 0$

R^2

1

-

$SSE = SST$

$$2. R^2 = 0$$

$$R^2 = 1 - 0 = \underline{\underline{1}}$$

$$SSE = SST$$

$$3. R^2 = +ve, \quad SSE < SST$$

$$R^2 = 1 - 0.1 = \underline{\underline{0.9}}$$

$$4. R^2 = -ve, \quad SSE > SST$$

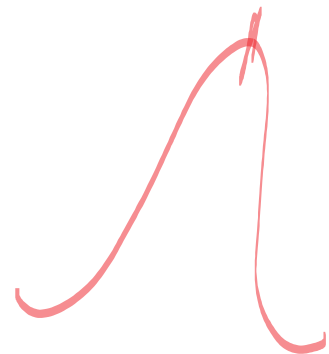
$$R^2 = 1 - 1.4 = \underline{\underline{-0.4}}$$

Non-linear.

$$R^2 = 1 - \frac{SSE}{SST} \rightarrow \frac{(y_a - y_p)^2}{(y_a - \bar{y})^2}$$

Exp	Sal
2	30
3	40
4	50
5	60
6	70

$y_a, \bar{y} = 50$



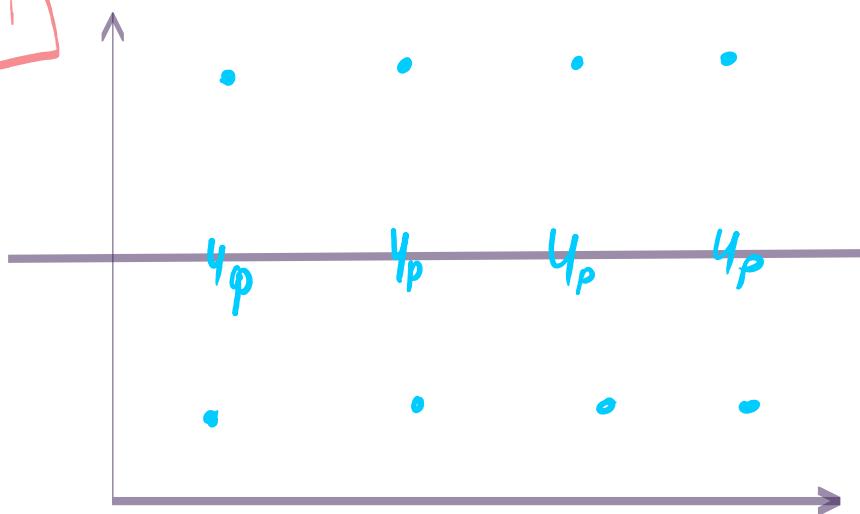
$$R^2 = 1 - \frac{SSE}{SST}$$

$$R^2 = 1 - \frac{\text{Unexplained Variation}}{\text{Total Variation}}$$

$$R^2 = \frac{SST - SSE}{SST} \rightarrow \frac{\text{Explained Var}}{\text{Total Var}}$$

$$SSE = 0$$

$$SSE = SST$$

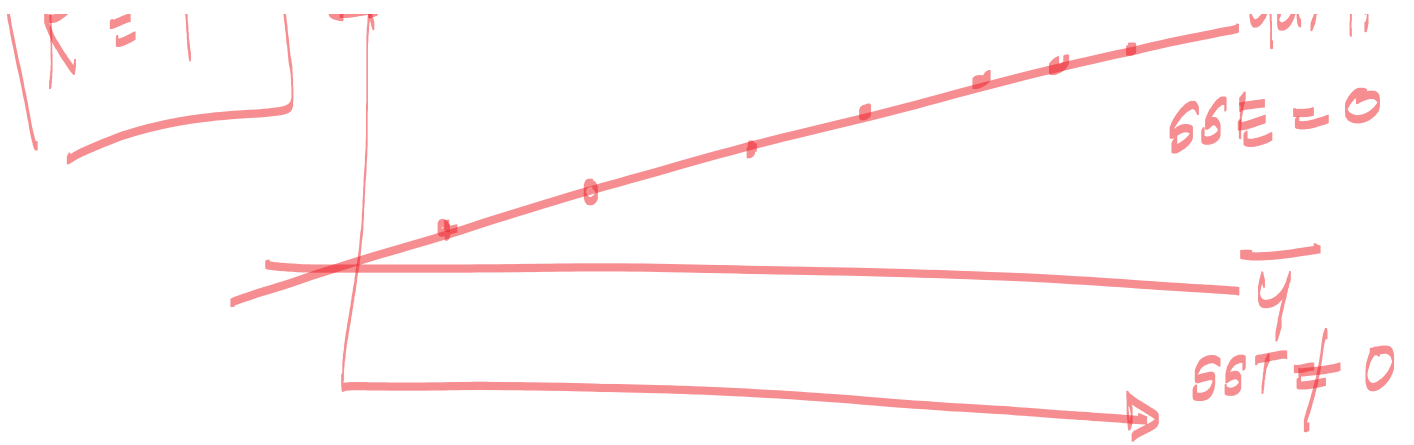


$$SSE = y_a - y_p$$

$$SST = y_a - \bar{y}$$

$$R^2 = 1$$





R-squared

Coeff. of Cor.

$R^2 = R \times R \rightarrow SLR$

$R^2 \neq R \times R \rightarrow MLR$

$R^2 \rightarrow$ (r^2 -score) Goodness of Best fit line

x_1	x_2	x_3	x_4	x_5	y
0.9	0.4	0.8	0.7	-0.5	

3 features $\rightarrow R^2 = 0.85$ } 0.852
 (x_1, x_3, x_4)

4 features $\rightarrow R^2 = 0.86$

$$(x_1, \underline{x_2}, x_3, x_4)$$

Adjusted R-squared $\rightarrow \bar{R}^2$

$$\bar{R}^2 = 1 - \frac{(1-R^2)(n-1)}{n-p-1}$$

n = No of samples (rows)

p = No of predictors

$$R^2 \geq \bar{R}^2$$

R^2	\bar{R}^2
0.85	0.85
0.86	0.84
0.87	0.86

① $R = 0.3$

② $R = 0.9$

$$R^2 = 1 - \frac{SSE}{SST}$$

$$\overline{R^2} = 1 - \frac{(1-R^2)(n-1)}{n-p-1} \leftarrow 999$$

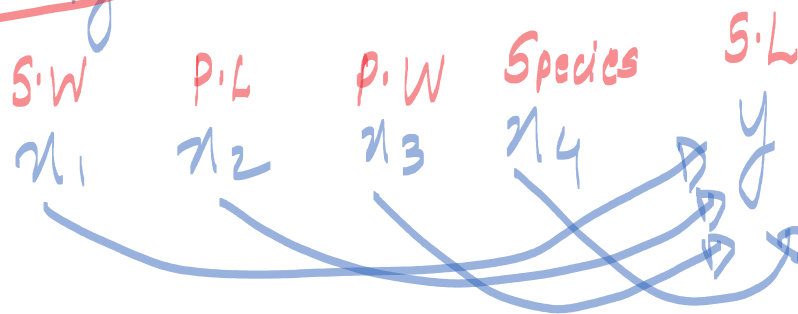
$n-p-1 \leftarrow 994$

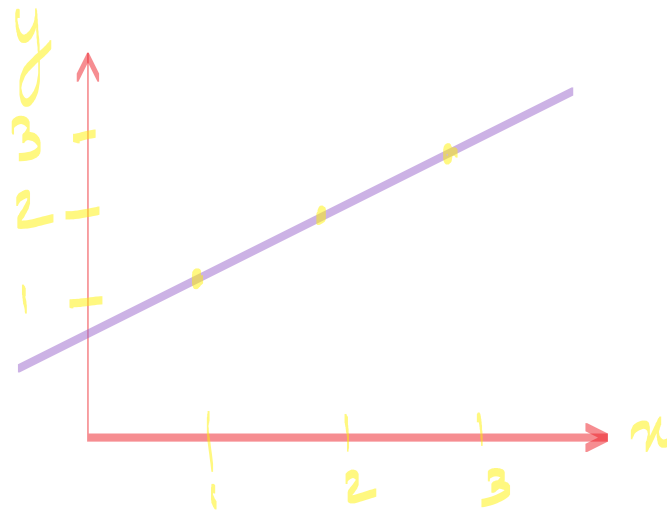
$$n = 1000, p = 5$$

$$\boxed{n > p}$$

✓ Rows
6 columns X

* Linearity .





R-value = 1

150 rows

df \rightarrow x_1 x_2 x_3 x_4 y

Tr \swarrow
Test \swarrow $x \rightarrow$ Df \rightarrow x_1 x_2 x_3 x_4

Tr \swarrow
Test \swarrow $y \rightarrow$ series \rightarrow y

Train
 \swarrow \searrow
 $x_{\text{-train}}$ $y_{\text{-train}}$
 120 rows 120 rows

Test
 \swarrow \searrow
 $x_{\text{-test}}$ $y_{\text{-test}}$
 30 rows 30 rows

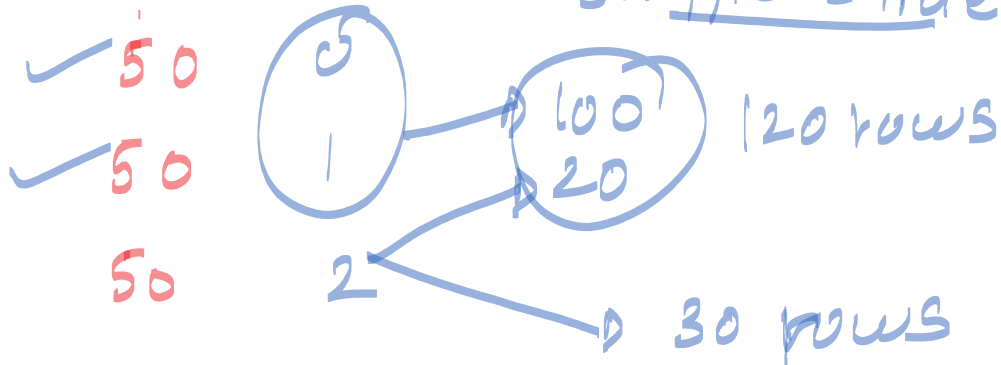
x
 $(x_1 \quad x_2)$
 y
 (y)

Tr \checkmark 1
 Tr \checkmark 2
 Tr \checkmark 3
 Tr \checkmark 4
 Tr \checkmark 5

Tr = 5

✓ Train → 3 row → x-train, 2, 4, 5
y-train, 2, 4, 5
✓ Test → 2 row → x-test 1, 3
y-test 1, 3

shuffle = True



100 row 80 → Train ✓

20 → Test

* What is Data Science?

Train — 4 lines → 3 correct 75%
evaluate 1 incorrect 25%

* What is KNN?

Test — 4 line → 2 correct 50%
Evaluate unseen

* What is RMSE
Test — 4 line \rightarrow 2 correct 50%
Evaluate
unseen

	Area	Bed	Bath	Location	Price
Train	1000	3	2	Baner	2cr
	700	2	1	Wakad	80
	800	1	1	Bale	75
Test Unseen	650	2	2	Baner	90L

$$y = m \times x + c$$

