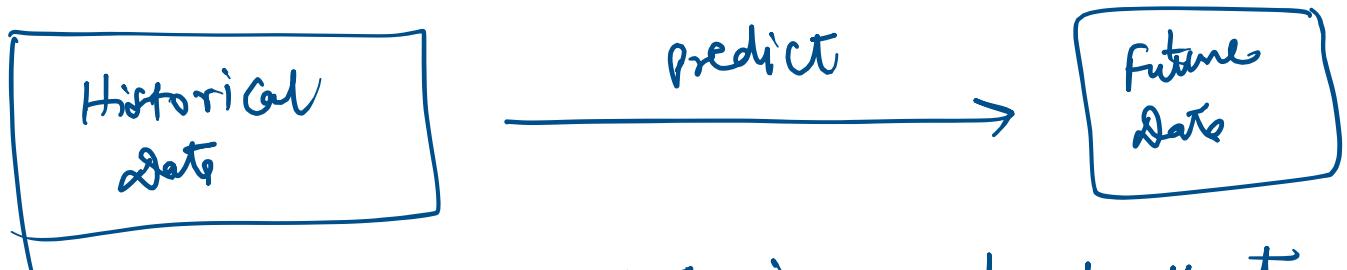
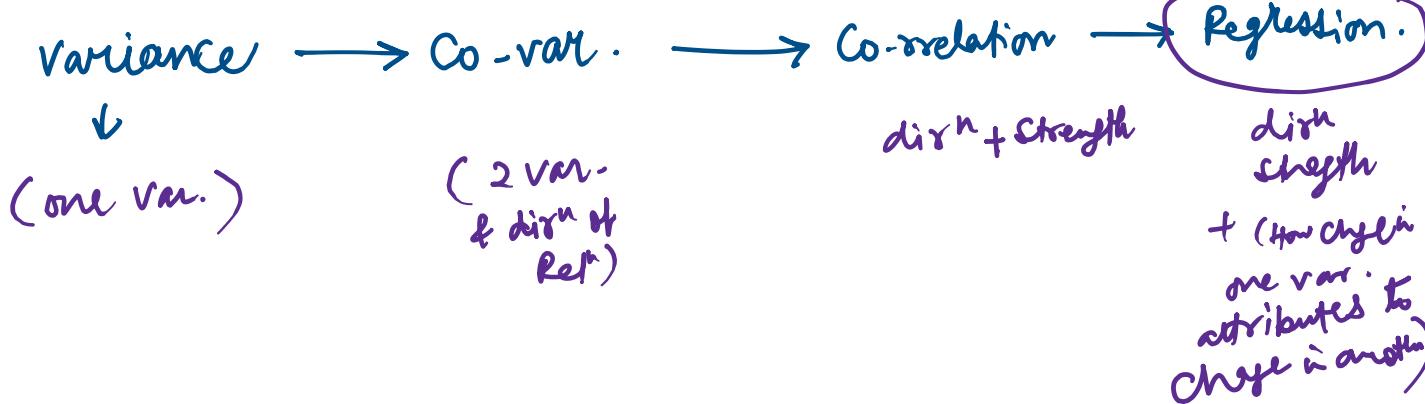


## Linear Regression



e.g. → from historical data I've found out that Price of house is dependent upon its size.

$$\begin{aligned}
 P &\propto S \\
 P = kS & \\
 80,000,000 & \xrightarrow{8000 \text{ s.f.}} 8000
 \end{aligned}$$

If I give you a new size  $\Rightarrow$  Price?

1500 Sq. ft.

$1500 \times 8000$

120,00,000

1.2 crores.

$P \propto$  Size ✓

$P \propto$  No. of Rooms ✓

$P \propto$  No. of Balcony ✓

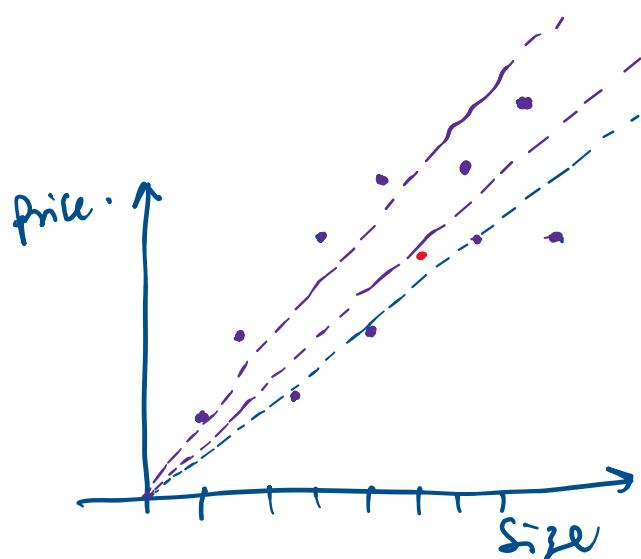
$\propto$  No. of Washrooms ✓

Mathematics -

multiple lines



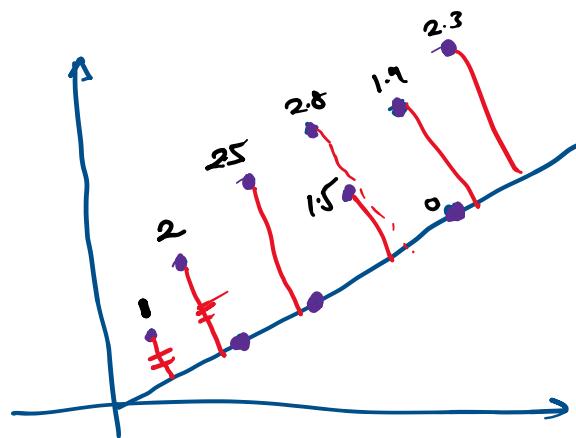
Best fit line



Case I

Error = dist<sup>n</sup> betw pt. & line.

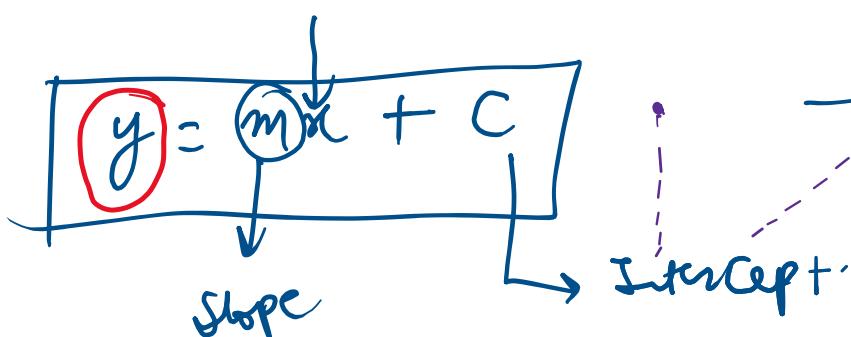
$$1+2+2.5+1.5+2.8+0+1.7+2.3$$



Aim = Least fit line

$$\boxed{\text{Error} = \text{Min}^m}$$

$$\text{Error} = \frac{1+2.5+2.8+2.9+1.7}{(0.3+0.9+2.3)}$$



Sales  $\propto$  Adv.

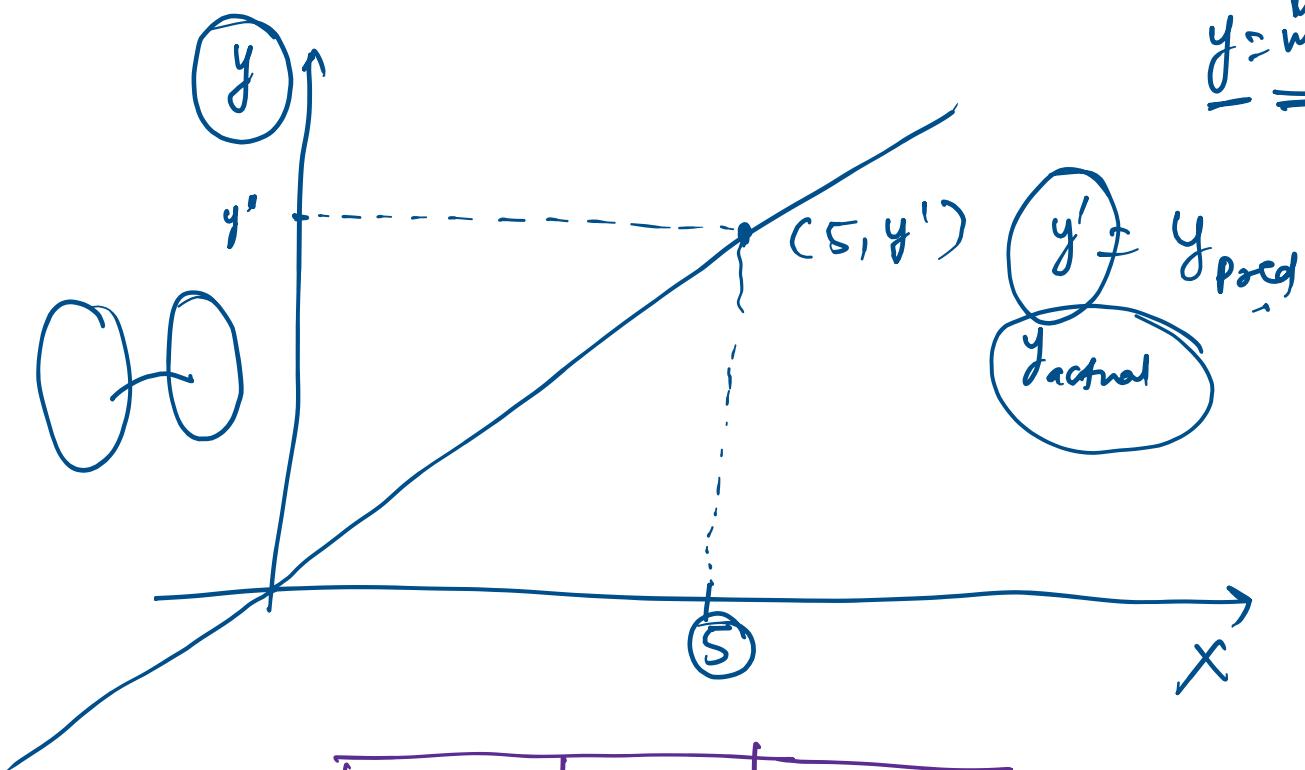
$$\text{Sales} = K(\text{Adv}) + C$$

If Adv = 0, Sales = 0.

our goal is to find statistically significant values of  $\underset{m}{\uparrow}$  slope &  $\underset{C}{\uparrow}$  Intercept that minimises

diff bet<sup>n</sup>  $y_{pred}$  &  $y$ .

$$\underline{y = mx + c}$$



	$y_{pred}$	$y_{actual}$	Error
10	65	68	(3)
60	60	0	0

↓  
Total.

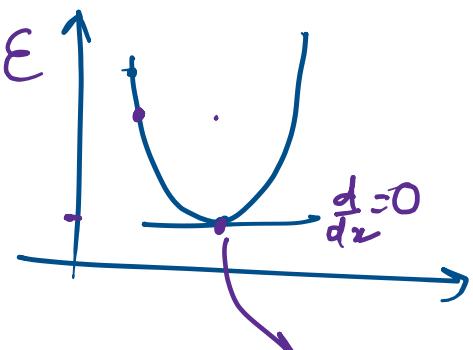
Ways to think abt the error.

- ① Absolute
- ② Square

$\times \textcircled{3}$  4<sup>th</sup> Power  $\times$  (more complex)  
 $\times \textcircled{4}$  6<sup>th</sup> Power  $\times$  (more complex)

Error needs to be minimised  $E$

$$\frac{d}{dx}(\text{error}) = 0$$



derivative = slope

Sq of error:

$$\frac{d}{dx} x^2 = 2x = 0$$

$x = 0$

$$\frac{d}{dx} (\underline{x}) = \underline{1}$$

$$\frac{d}{dx} (\underline{x^2}) = \underline{2x}$$

Error

$$SSR = \sum (y - \hat{y})^2 \longrightarrow \text{Reduce this error.}$$

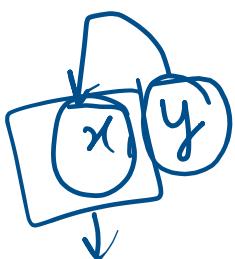
Sum of Sq. residuals

Algorithms

- ① ordinary Least Square  $\rightarrow$  Best fit Line.
- ② Gradient descent.  $\rightarrow$  Step wise-

① OLS:

$$\text{Least Sq.} = \frac{\text{Cov.}(x, y)}{\text{var}(x)} \quad \begin{array}{l} \xrightarrow{\text{var. of } y \text{ w.r.t } x} \\ \xrightarrow{\text{Co-var. of } x \text{ w.r.t. } y} \\ \downarrow \\ \text{var}(x) \end{array}$$



$$= \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

$$= \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2}$$

How would you become great / best

↓  
become better.  
↓  
best

2, 3

(OLS)  
best fit line

↓  
minimise your fault.

↓  
(Gradient descent)

minimising the error  
(step wise)

## Gradient Descent:

Aim is to minimise the error.

↓  
Cost

↓  
Cost function

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

(wx + b)

$$h(\theta) = (y - (wx + b))^2$$

$$= (y - wx - b)^2$$

$$\frac{d h(\theta)}{dw} = 2(y - wx - b) \cdot (-x)$$

$$\Downarrow$$

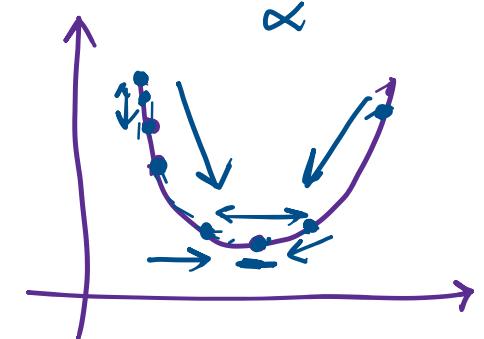
$$\frac{d h(\theta)}{dw} = -2(y - wx - b)(x)$$

$$\frac{d h(\theta)}{db} = 2(y - wx - b) \cdot (-1)$$

$$db \leftarrow \frac{d h(\theta)}{db} = -2(y - wx - b)$$


How?

$$w_{next} = w \ominus \alpha (dw)$$



$$b_{next} = b \ominus \alpha (db)$$

e.g.  $\rightarrow$

$$w_{next} = 5 - 0.1 \cancel{(-0.25)}$$

$$= 5 + \boxed{\phantom{0.00}}$$

first diff.

$$W_{next} = 5 - 0.1(+0.25)$$

$$= 5 - \boxed{ }$$

↓  
backward diff  $\neq 0$  = -ve  
 $\leftarrow 0$  = +ve  
 $\rightarrow -$  = 0

ACCenture

"ACCY"

MLR

y

$x_1$  = interest rates

$x_2$  = GDP

$x_3$  = value of NIFTY 50 Index

$x_4$  = Spend. on Ad.

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 + \beta_4 x_4$$

Project

1

## Evaluation Metrics

### \* Error Metrics

#### ① Mean Absolute Error : (MAE)

Mean of the absolute value of errors.

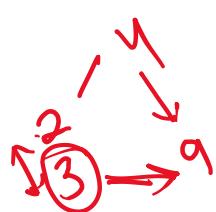
$$\frac{1}{n} \sum_{i=1}^n |y - \hat{y}|$$

②

#### Mean Squared Error (MSE)

Mean of squared Errors.

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

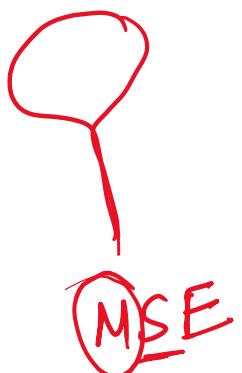
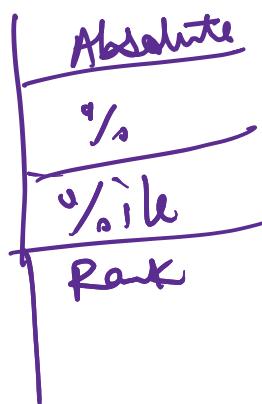


③

#### Root Mean Square Error (RMSE)

Sq. Root of Mean of Squared Error.

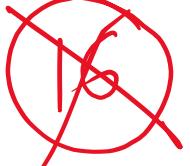
$$\sqrt{\frac{1}{n} \sum_{i=1}^n (y - \hat{y})^2}$$



small errors



large errors



# MAE V/S RMSE

E	E	$E^2$
2	2	4
2	2	4
2	2	4
2	2	4
2	2	4
2	2	4
2	2	4
2	2	4
2	2	4
2	2	4

T-1

E	E	$E^2$
1	1	1
1	1	1
1	1	1
1	1	1
3	3	9
3	3	9
3	3	9
3	3	9
3	3	9
3	3	9

T-2

E	E	$E^2$
0	0	0
0	0	0
0	0	0
0	0	0
0	0	0
0	0	0
0	0	0
0	0	0
0	0	0
20	20	400

T-3

$$\begin{aligned} \text{MAE} &= 2 \\ \text{RMSE} &= 2 \end{aligned}$$

$$\begin{aligned} \text{MAE} &= 2 \\ \text{RMSE} &\approx 2.25 \end{aligned}$$

$$\begin{aligned} \text{MAE} &= 2 \checkmark \\ \text{RMSE} &= 6.32 \end{aligned}$$

Evenly distributed error.

Since RMSE gives relatively high weights to large errors hence should be used when ~~large~~ large errors are undesirable.

↳ Interpretability =  $\text{MAE} > \text{RMSE} > \text{MSE}$

↳ sensitive to outliers =  $\text{MSE} > \text{RMSE} > \text{MAE}$

## # $R^2$ & Adjusted $R^2$

$R^2$  = Coeff. of determination

$$-1 \leq R^2 \leq 1$$

$$0 \leq R^2 \leq 1$$

$$R^2 = 1 - \frac{SS_{res}}{SS_{total}}$$

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

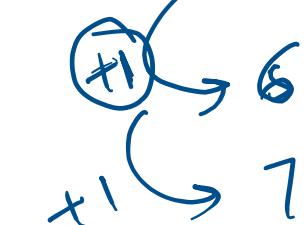
1 or 100%  
0 or 0%

~~68%~~

$$R^2 = 1 - \frac{\text{unexplained variance}}{\text{Total variance}}$$

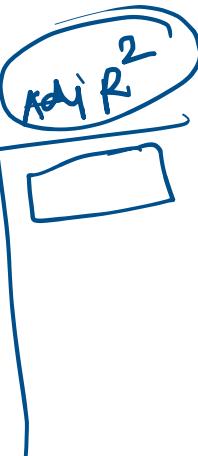
## Adjusted $R^2$

Iteration 1: 5 var.



$R^2$   
68 %

68.9 %  
69.2 %



y

$R^2$

32

Adj  $R^2$  can decrease also.

EDA

select variables



77.8

77.0%

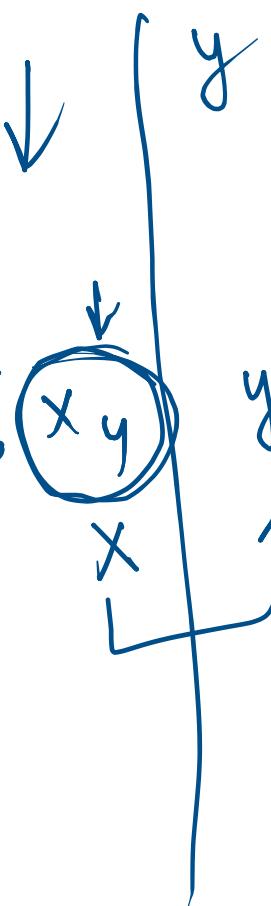
$x_1, x_2, x_3$

$R^2$

78%

$x_1, x_2, x_3, x_4$

$$R^2 = 79.5\%$$

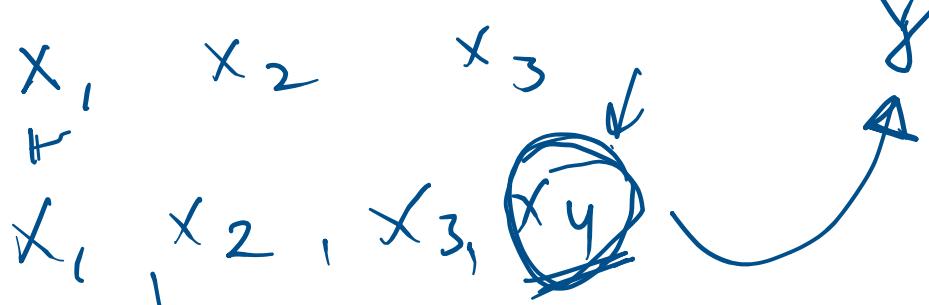


$$\bar{R}^2 = 1 - \frac{SS_{\text{res}} / df_e}{SS_{\text{total}} / df_t}$$

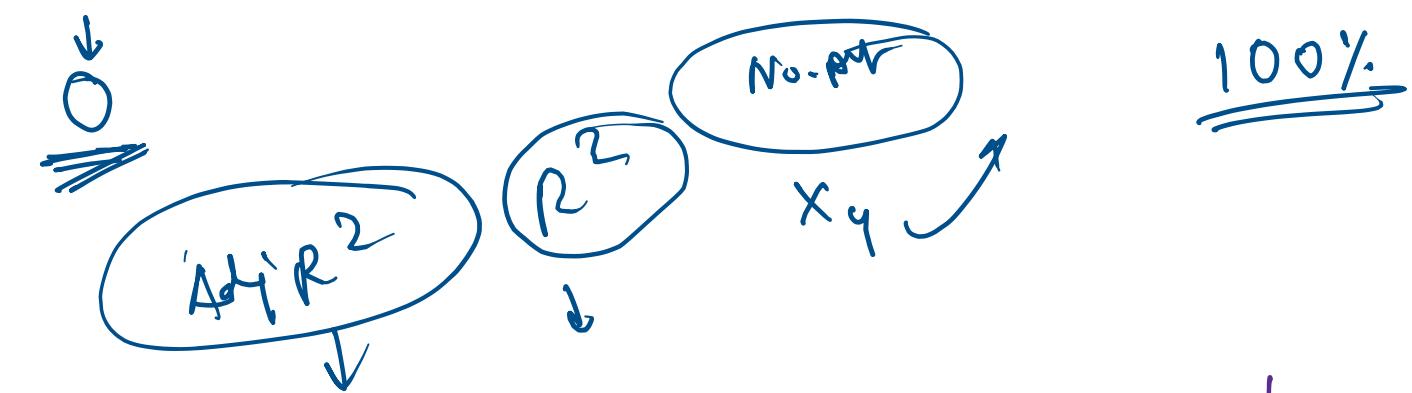
1.5

Adj.  $R^2$

Penalises for every increase in  
variable -  
spur



$R^2 \downarrow$   
-



Q. what are the evaluation metrics you have used for your Regression Model?

- 
- 
- ①  $R^2$
  - ②  $\text{Adj-}R^2$
  - ③ MSE
  - ④ MAE
  - ⑤ RMSE

# Multicollinearity & How to Avoid it?

Independent Variable

Dependent Variable  
□

What happens if one independent var. is highly correlated with one or more other independent variables?

Independent var. 1

Depend. var.

Model (generalization)  
Learning ("")

2 1 2  
2 2 9  
2 5 9  
9

Detection :

① Correlation Matrix

	$x_1$	$x_2$	$x_3$	$x_4$	$x_5$	$y$
$x_1$	1					
$x_2$		1				
$x_3$			1			
$x_4$				1		
$x_5$					1	
$y$						1

Correlation matrix showing values for  $x_1, x_2, x_3, x_4, x_5, y$ . The diagonal elements are 1. The off-diagonal elements are highlighted in pink:  $x_1x_2 = 0.8$ ,  $x_1x_3 = 0.9$ ,  $x_2x_3 = 0.9$ ,  $x_2x_4 = 1$ ,  $x_3x_4 = 1$ ,  $x_3x_5 = 1$ ,  $x_4x_5 = 1$ .

## ② variance Inflation factor (VIF) :-

VIF tells how much the variance of an estimated regression co-eff increases if your predictors are correlated.

$$VIF = \frac{1}{1 - R^2} \quad VIF \in (1, \infty)$$

VIF (now about 5 or more).

## Prevent

- ① drop one of it .
  - ② Feature Engineering . ( $x_3$  &  $x_4$ )

