



# Data Science and Artificial Intelligence

## Machine Learning

Regression

Lecture No. 06



GATE WALLAH

By- SIDDHARTH SABHARWAL SIR



# Recap of Previous Lecture



Topic

Questions

Topic

$R^2 \Rightarrow$  Coeff of determination

Goodness of fit

Topic

MSE

RSS

TSS

Topic

$R_{MSE}$

\*Residual plot

Topic

$\tilde{f}_w$

# Topics to be Covered



- Topic
- Topic
- Topic
- Topic
- Topic

effect of outlier

Assumption in L.R

→ Multicollinearity

VIF

→ Homo/Heteroscedasticity





Positivity

**Optimism is the one  
quality more  
associated with  
success and happiness  
than any other.**

BRIAN TRACY

BRIAN TRACY  
INTERNATIONAL



### What is Coeff of Determination

$$R^2 = 1 - \frac{RSS}{TSS} \Rightarrow (0 \text{ to } 1)$$

1 → overfitting  
0 → the model = be KAR

- goodness of fit
- It determine how much our model detect/learn the variation/pattern of y.
-



## Basics of Machine Learning

### What is MSE and RMSE

$$\left( \frac{1}{N} \text{RSS} \right), \sqrt{\text{MSE}}$$



# REVISION



### What is Hat Matrix

$$\hat{Y} = X\hat{\beta}$$
$$\hat{Y} = X \underbrace{(X^T X)^{-1} X^T Y}_{\text{Hat matrix}}$$

# REVISION

If one dimensional data

$$\cdot \bar{y} = \bar{\omega}x \Rightarrow \bar{\omega} = \frac{\sum_{i=1}^N x_i y_i}{\sum_{i=1}^N x_i^2}$$

$$\cdot \bar{y} = \bar{\omega}_0 \Rightarrow \bar{\omega}_0 = \frac{\sum_{i=1}^N y_i^0}{N}$$



## Linear Regression



Question 15: What is the purpose of the coefficient of determination (R-squared) in simple linear regression?

- A. To determine the slope of the regression line
- B. To measure the strength of the linear relationship
- C. To calculate the p-value of the regression
- D. To identify outliers in the dataset

How much the  
The linear model  
Understand pattern of data



## Linear Regression



Question 19: If the R-squared value in simple linear regression is 0.75, what does it indicate?



- A. A strong linear relationship between the variables
- B. A weak linear relationship between the variables
- C. No linear relationship between the variables
- D ~~X~~ The model is overfitting  $\leftarrow R^2 = 1$

Simple LR

↳ Single dimension  
data

$$y = \beta_0 + \beta_1 x$$

Multiple LR

↳ data is of more  
than one dimension

$$y = \beta_0 + \beta_1 x^1 + \beta_2 x^2 - \dots$$



## Linear Regression



Question 2: What does the coefficient of determination (R-squared) measure in multiple linear regression?

- A. The correlation between predictor variables
- B. The percentage of variance in the dependent variable explained by the model
- C. The significance of the intercept term
- D. The number of predictor variables in the model



## Linear Regression



In multiple linear regression, what is the key difference between simple linear regression and multiple linear regression?

- A) Simple linear regression has one independent variable, while multiple linear regression has two or more.
- B) Simple linear regression uses categorical variables, while multiple linear regression uses continuous variables.
- C) Simple linear regression is used for classification, while multiple linear regression is used for prediction.
- D) There is no difference between simple and multiple linear regression.



## Linear Regression



Which statistic is used to assess the strength and direction of the relationship between the dependent variable and each independent variable in multiple linear regression?

$y$

$x^1, x^2, x^3 \dots$

- A) Mean absolute error (MAE)
- B) R-squared ( $R^2$ ) ✓
- C) Standard error
- D) Confidence interval



## Linear Regression



What is the purpose of the residual plot in multiple linear regression analysis?

- A) To visualize the relationship between independent variables.
- B) To check for homoscedasticity and the presence of outliers.
- C) To calculate the correlation coefficient ( $r$ ).
- D) To assess multicollinearity.



## Linear Regression



What is the main purpose of the intercept term in a multiple linear regression model?

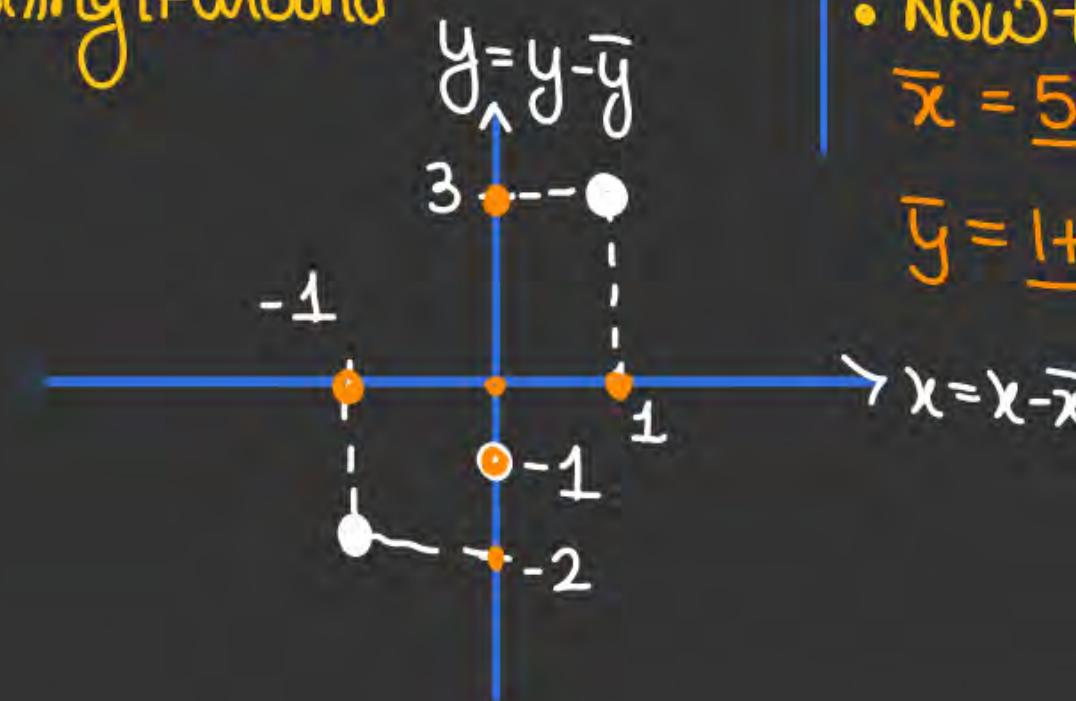
C

→  $\beta_0 \Rightarrow$  Biastem.

- A) It represents the slope of the regression line.
- B) It is used to control for multicollinearity.
- C) It represents the expected value of the dependent variable when all independent variables are zero.  
$$y = \beta_0 + \beta_1 x^1 + \beta_2 x^2 - \dots$$
- D) It is not used in multiple linear regression.

What is  
the centred  
data

- So One Dimension data  $\Rightarrow$
- The data is plotted  $\Rightarrow$
- The Centre of the 2D Plane  $\Rightarrow$  origin
- Now Centering of data means, we pull the data and bring it around origin



data  
 $(5,1), (6,2), (7,3)$

- Now to Centre the data

$$\bar{x} = \frac{5+6+7}{3} = 6$$

$$\bar{y} = \frac{1+2+3}{3} \Rightarrow 3$$

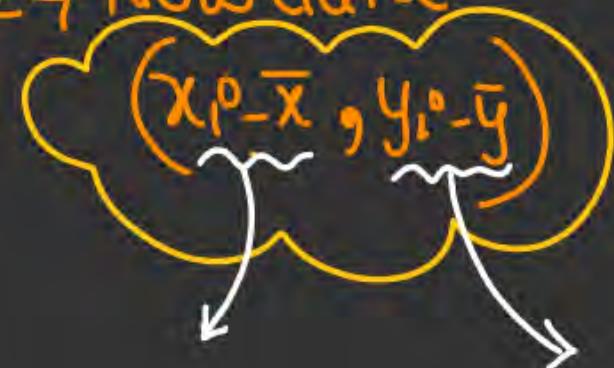
So Centreing of any data  $\Rightarrow$

$(x_i^o, y_i)$  are datapoints

Step 1  $\Rightarrow \bar{x} \Rightarrow \frac{\sum x_i^o}{N}$

$$\bar{y} = \frac{\sum y_i^o}{N}$$

Step 2  $\Rightarrow$  Now data



This new Centred data  $\Rightarrow$

has mean of x and y = 0

- If we have a centred data  $\Rightarrow$  then in LR the Bias or intercept term  $B_0 \Rightarrow$  should be  $0$ .

$$B_0 = \bar{y} - m\bar{x}$$

$$B_0 = 0 - 0$$

already done in  
Maths

$\Rightarrow Y, X$  are R✓

$$Y = ax + b$$

$$\Rightarrow \bar{Y} = a\bar{x} + b$$

If  $X_{\text{new}} = X_{\text{old}} - \bar{X}_{\text{old}}$

$$\frac{\bar{X}_{\text{new}}}{X_{\text{new}}} = \frac{\bar{X}_{\text{old}}}{X_{\text{old}}} - \frac{\bar{X}_{\text{old}}}{X_{\text{old}}}$$
$$= 0$$

2D data

$$(x_i^1, x_i^2, y_i)$$

→ Centring →

$$\text{Step 1} \cdot \bar{y} = \frac{\sum y_i}{N}, \quad \bar{x}^2 = \frac{\sum x_i^2}{N}, \quad \bar{x}^1 = \frac{\sum x_i^1}{N}$$

$$\text{Step 2} \Rightarrow y_{i\text{new}} = y_i - \bar{y}$$

$$x_{i\text{new}}^2 = x_i^2 - \bar{x}^2$$

$$x_{i\text{new}}^1 = x_i^1 - \bar{x}^1$$



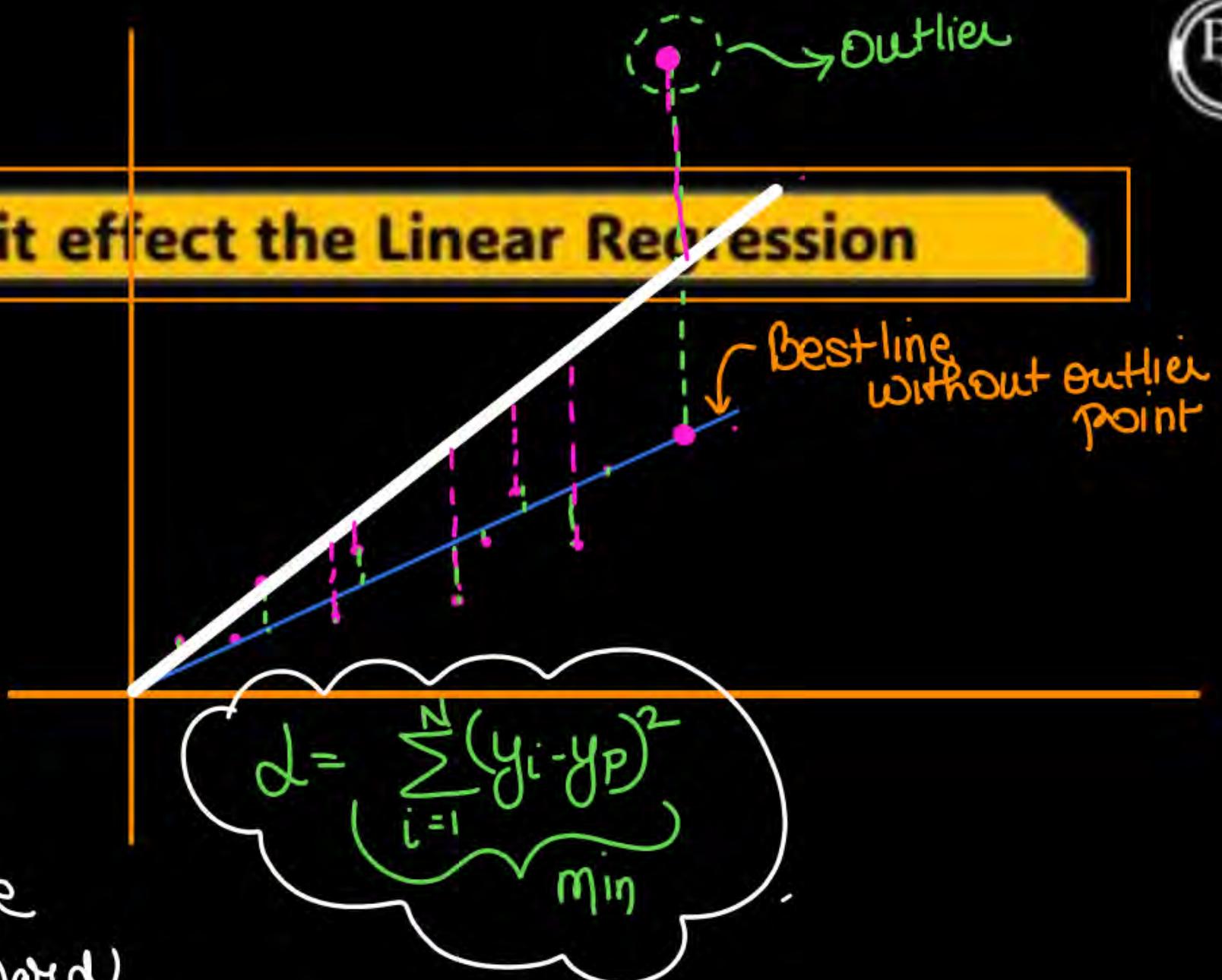
# Linear Regression



## What is an outlier and how it effect the Linear Regression

- data point which has huge noise and unfollow the data pattern

\* So the LR try to minimise the RSS, if the data has outlier then the LR will be effected by this and the Resultant line will be shifted toward the outlier



So LR algorithm donot neglect the outliers.

Assumptions in LR  $\Rightarrow$

①

we assume that actual pattern of  
data is linear



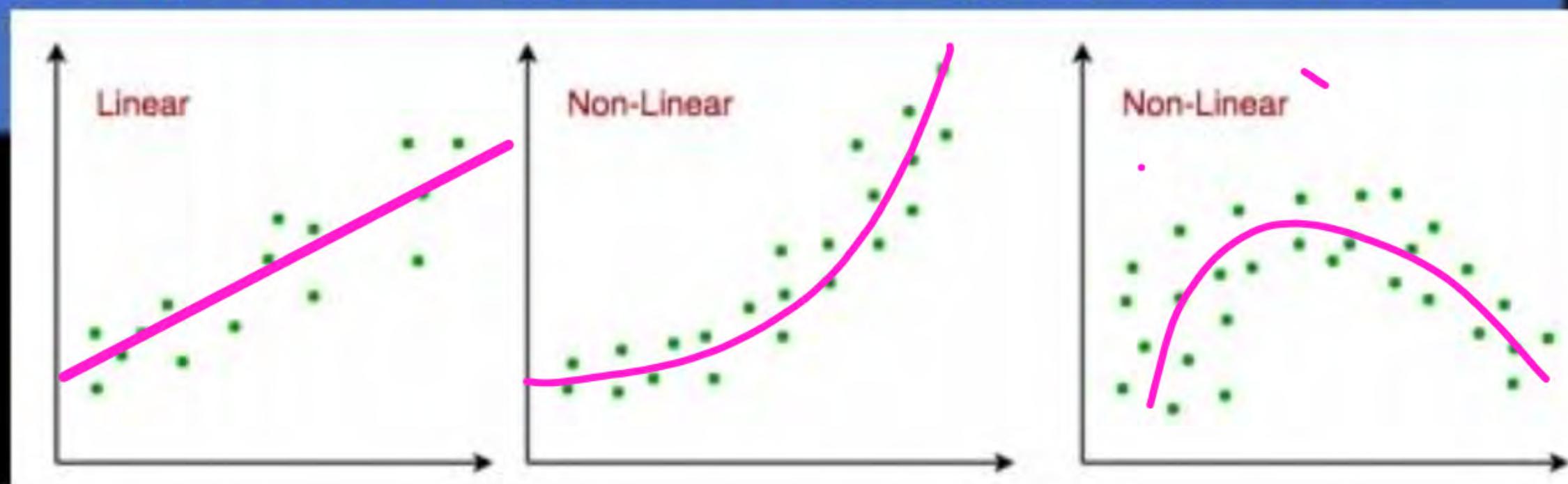
# Linear Regression



## Assumptions in Linear Regression

Linear regression needs to meet a few conditions in order to be accurate and dependable solutions.

**1. Linearity:** The independent and dependent variables have a linear relationship with one another. This implies that changes in the dependent variable follow those in the independent variable(s) in a linear fashion. This means that there should be a straight line that can be drawn through the data points. If the relationship is not linear, then linear regression will not be an accurate model.



2) In L-R the model  $\Rightarrow$

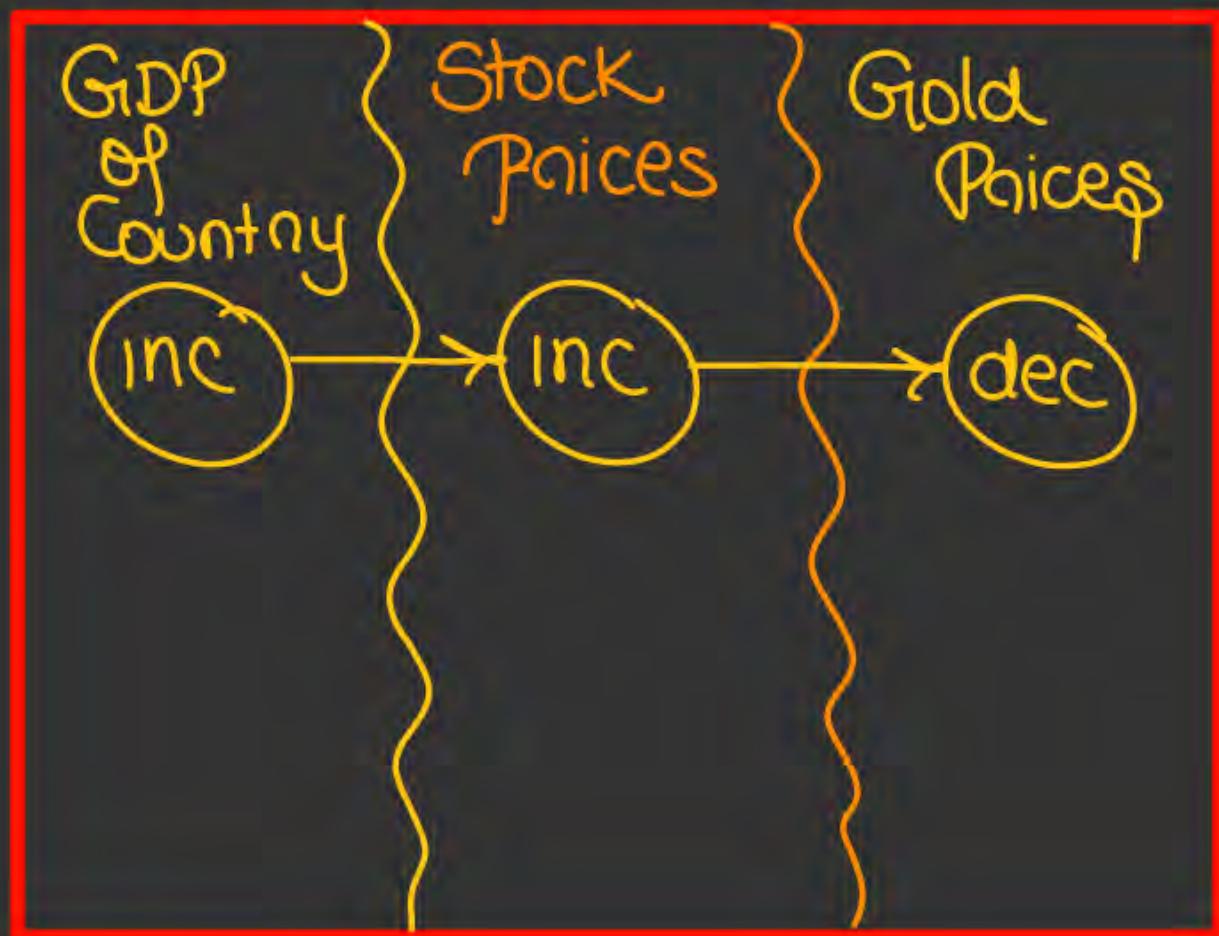
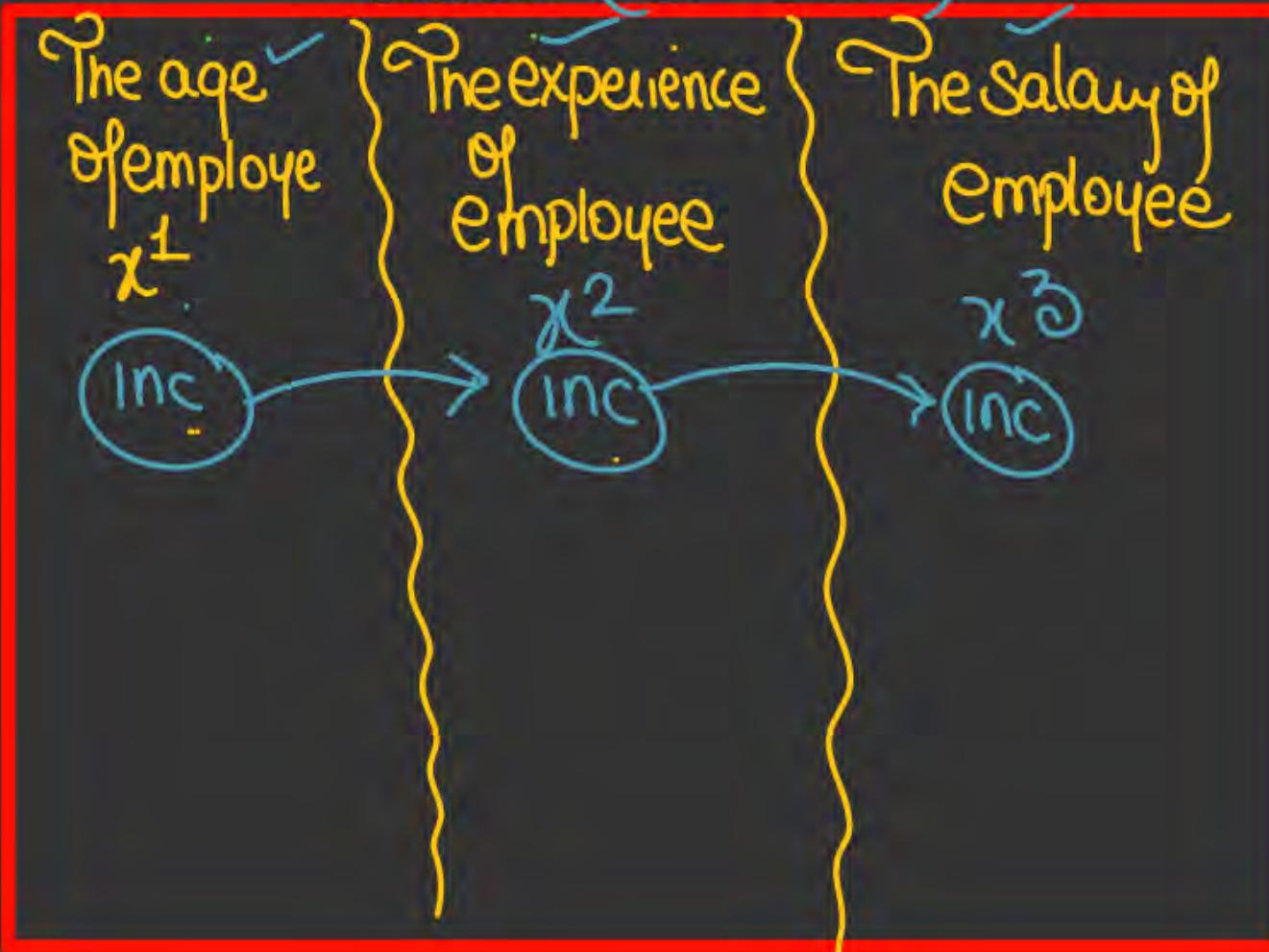
$$(y = \beta_0 + \beta_1 x^1 + \beta_2 x^2 + \beta_3 x^3 \dots \dots)$$

$\Rightarrow$  So from this eq, if  $x^1$  changes by 1 unit then  $y$  changes by  $\beta_1$

So we can say that  $\beta_1, \beta_2, \beta_3, \dots$  shows contribution of 1st, 2nd, 3rd dimension when these dimension changes by 1 unit.

ex

data (GovJob)

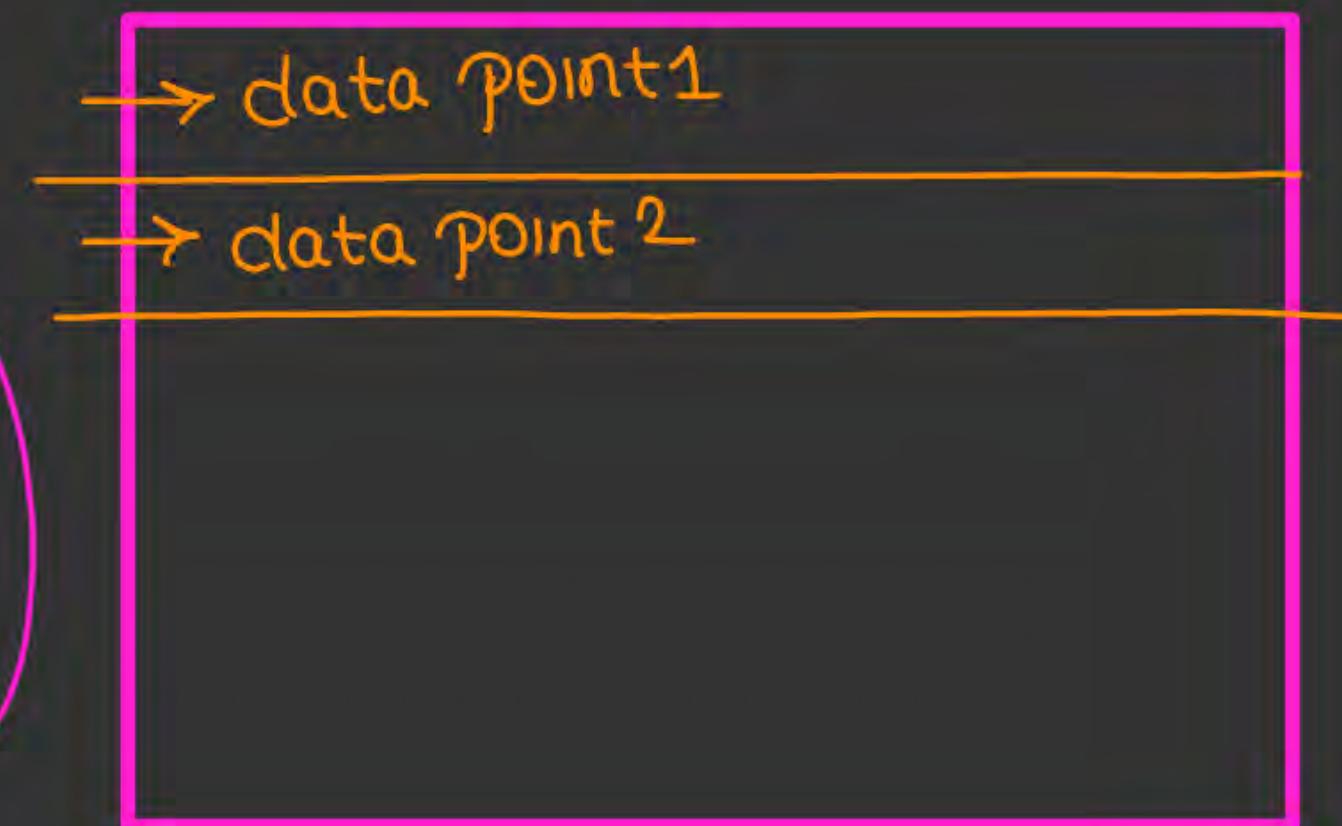


- So in LR we need that the dimensions in the data should be independent of each other.
- If data dimensions are dependent on each other then we say that data has multi collinearity
- So LR cannot work as desired in case of multi collinearity.

- also we assume that the data points are also independent

dependent ??

There shd be no Relationship  
btw the value of one data point  
with any other data point.





## Linear Regression



### Assumptions in Linear Regression

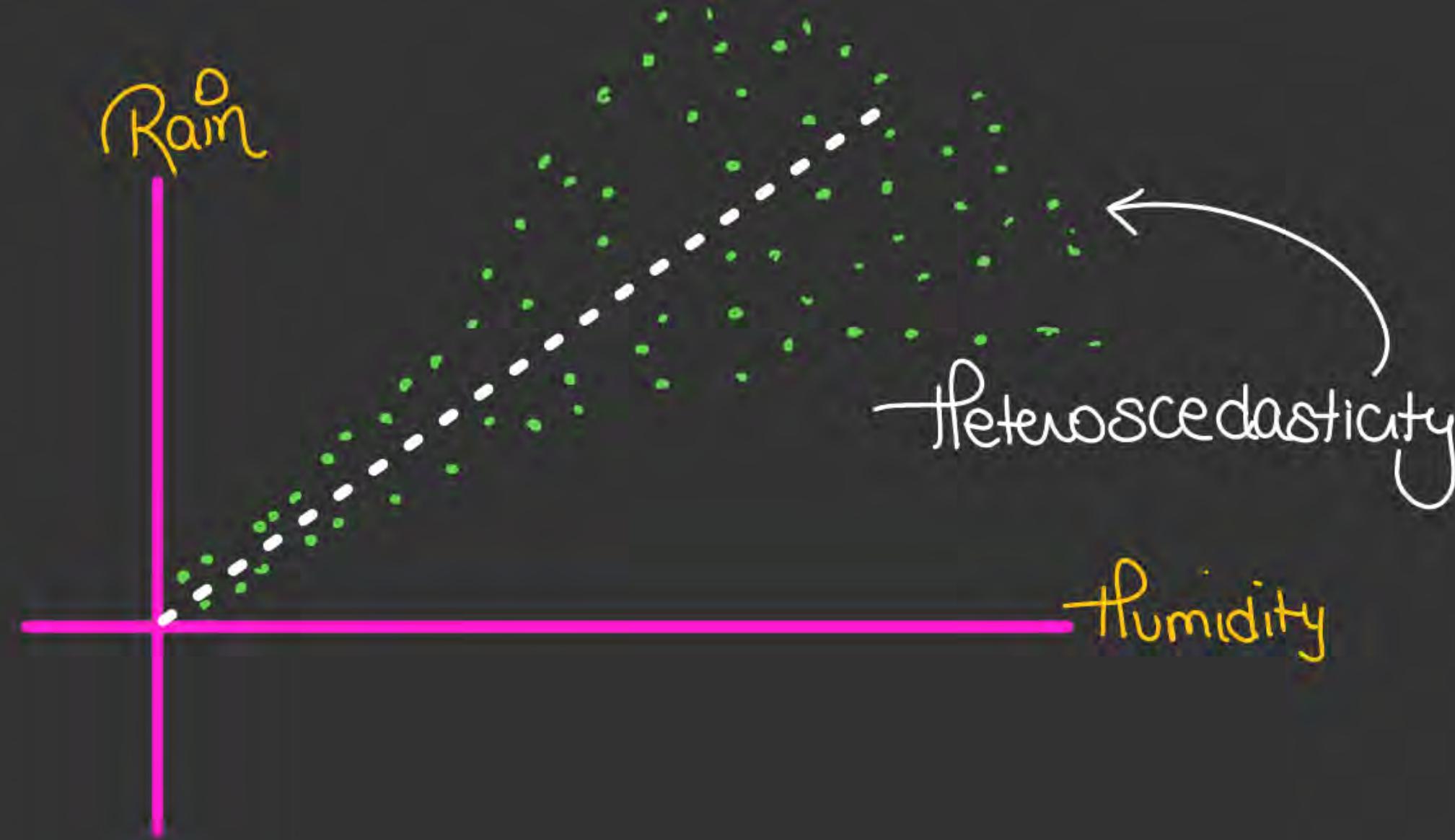
→ datapoints shd be  
independent of each other.

2. Independence: The observations in the dataset are independent of each other. This means that the value of the dependent variable for one observation does not depend on the value of the dependent variable for another observation. If the observations are not independent, then linear regression will not be an accurate model.

### 3. Homoscedasticity

- we assume that the noise or the error is symmetric and the amount of noise/error is independent of value of  $x$





So L-R will work when data has

homoscedasticity.

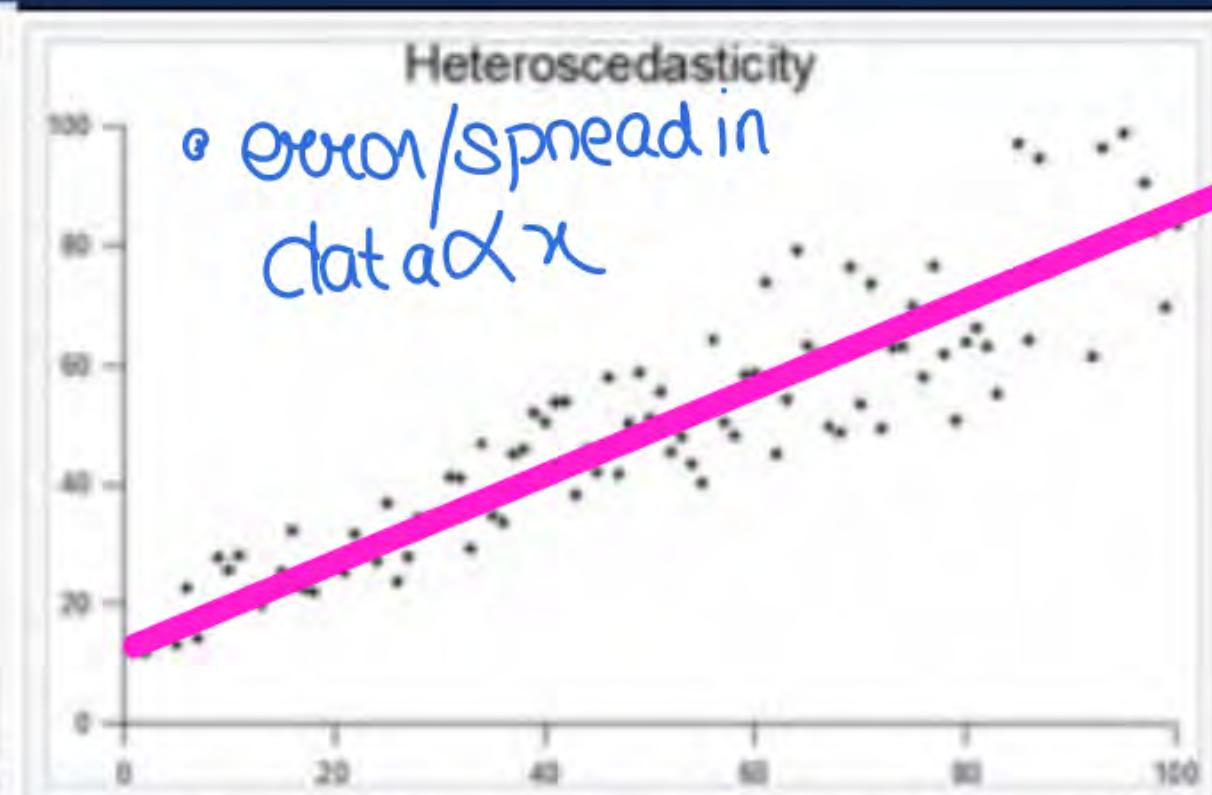
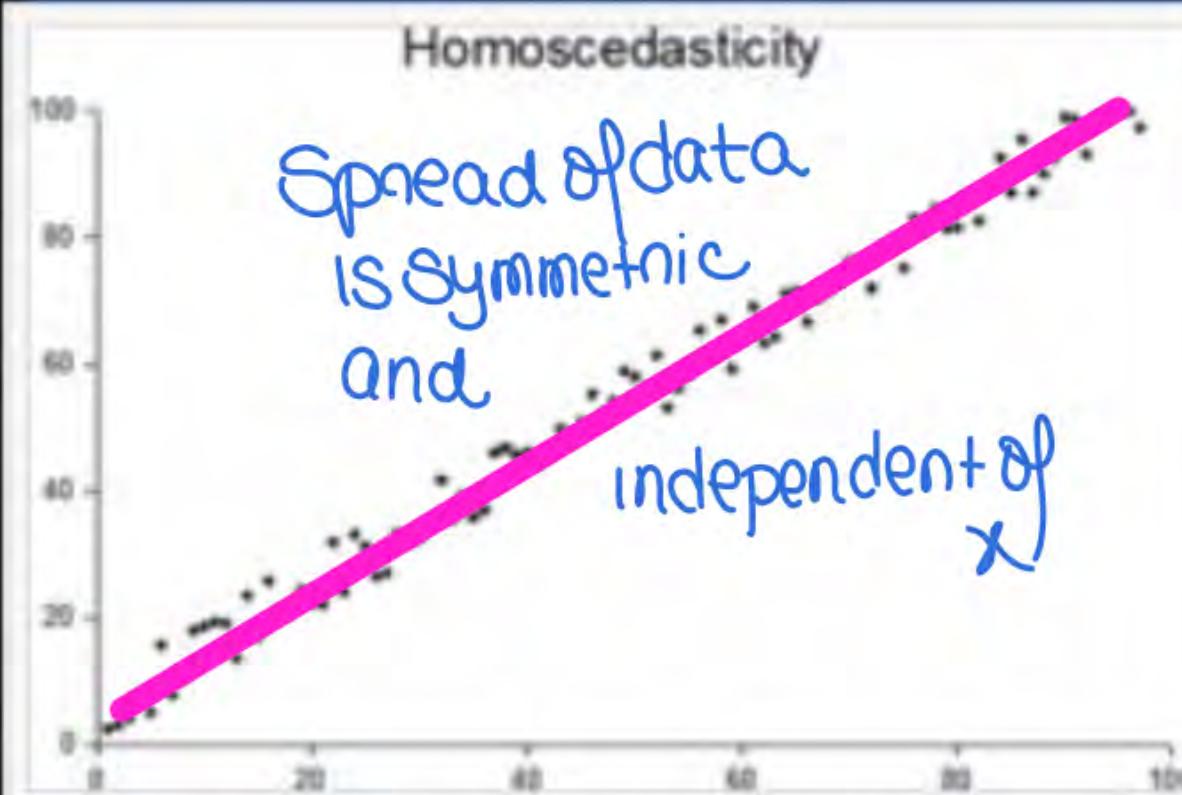


# Linear Regression



## Assumptions in Linear Regression

### 3. Homoscedasticity:





## Linear Regression

### Assumptions in Linear Regression

3. we assume that data in L-R has Homoscedasticity.

In Heteroscedasticity  $\Rightarrow$  due to huge noise and outlier the L-R will fail.



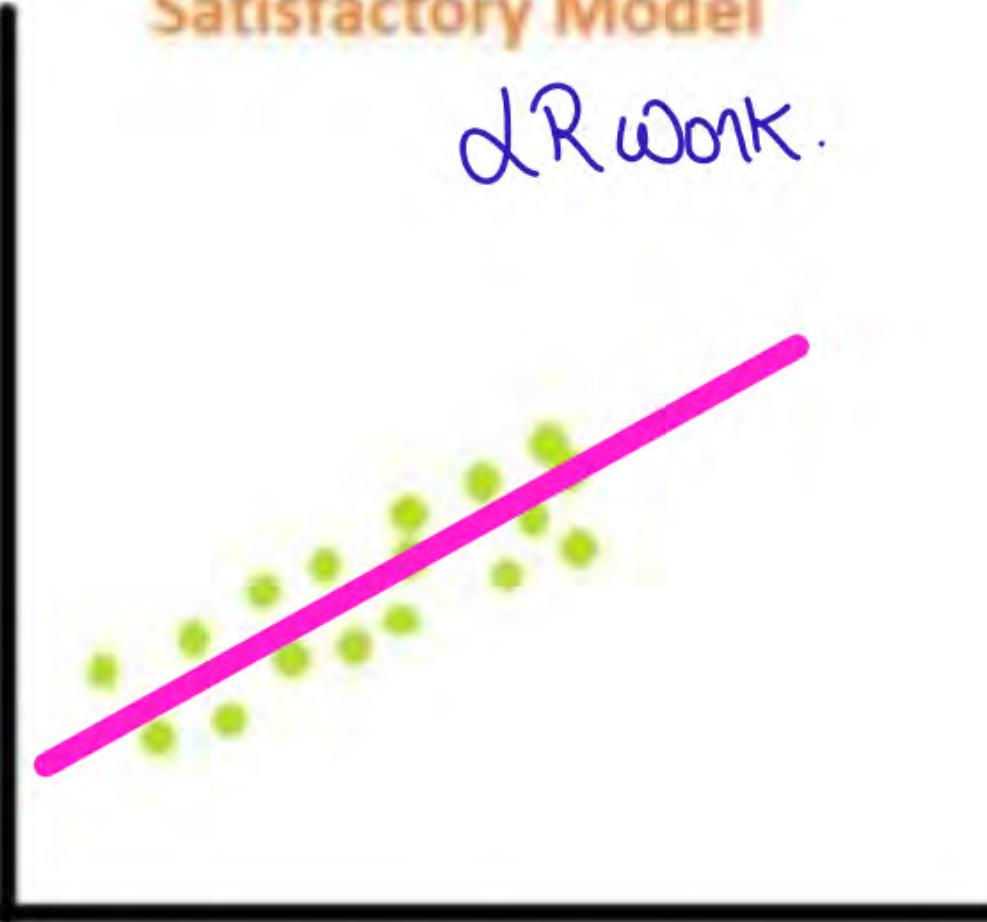
# Linear Regression



## Assumptions in Linear Regression

Satisfactory Model

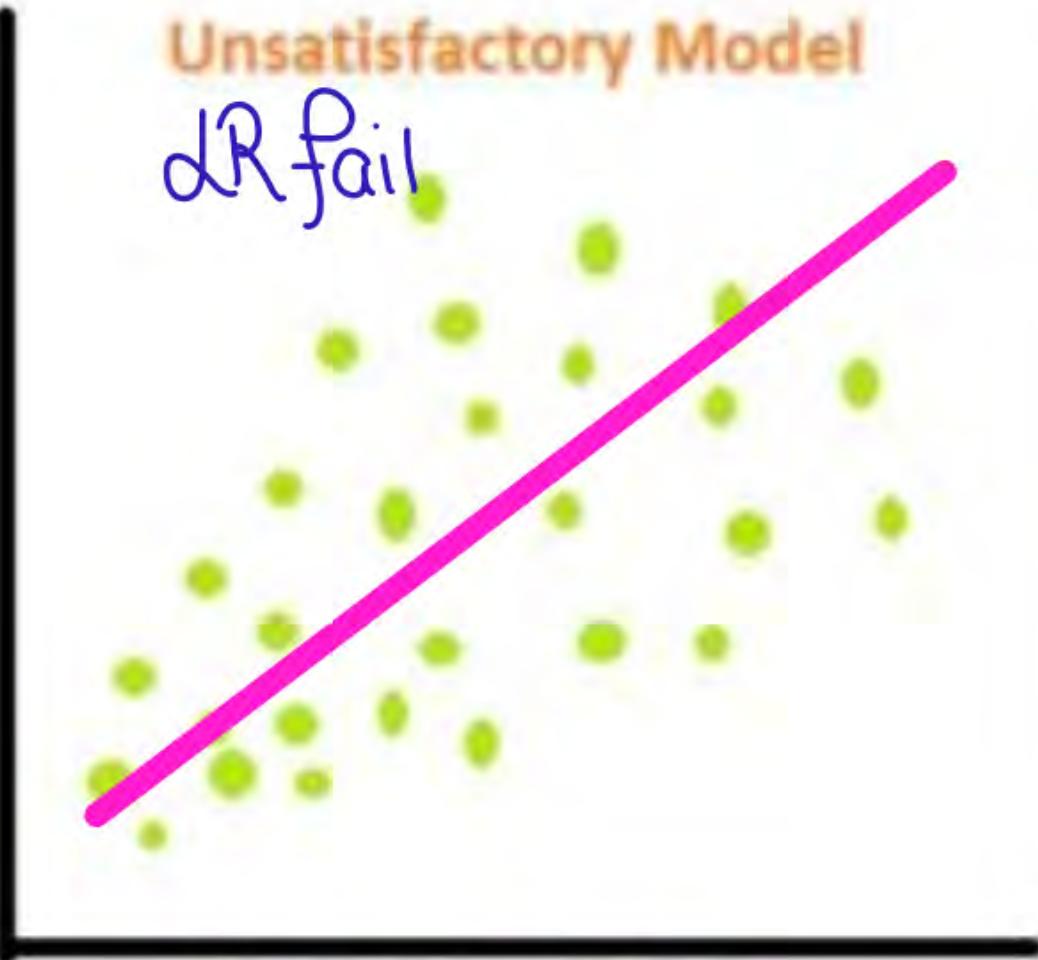
LR work.



Homoscedasticity

Unsatisfactory Model

LR fail



Heteroscedasticity



## Linear Regression



### Assumptions in Linear Regression

already told in 27) independence of dimension

4. No multicollinearity: There is no high correlation between the independent variables. This indicates that there is little or no correlation between the independent variables. Multicollinearity occurs when two or more independent variables are highly correlated with each other, which can make it difficult to determine the individual effect of each variable on the dependent variable. If there is multicollinearity, then multiple linear regression will not be an accurate model.



## Linear Regression



### Assumptions in Linear Regression

Detecting Multicollinearity includes two techniques:

- **Correlation Matrix:** Examining the correlation matrix among the independent variables is a common way to detect multicollinearity. High correlations (close to 1 or -1) indicate potential multicollinearity.
- **VIF (Variance Inflation Factor):** VIF is a measure that quantifies how much the variance of an estimated regression coefficient increases if your predictors are correlated. A high VIF (typically above 10) suggests multicollinearity.

Correlation



Correlation b/w any two

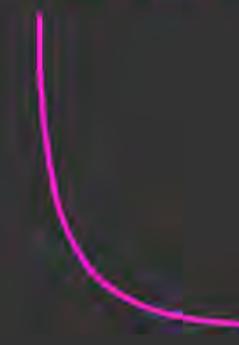
RV's  $\Rightarrow$  i.e dependence b/w

two RV's

Correlation Coef  $\Rightarrow \rho_{x^1 x^2} = \frac{\text{Cov}(x^1, x^2)}{\sigma_{x^1} \sigma_{x^2}}$

$$\text{Cov}(x^1, x^2) = \frac{1}{N-1} \sum_{i=1}^N (x_i^1 - \bar{x}^1)(x_i^2 - \bar{x}^2)$$

$$\rho_{x^1 x^1} \Rightarrow \frac{\text{Cov}(x^1, x^1)}{\sigma_{x^1} \cdot \sigma_{x^1}} \Rightarrow \frac{\text{Var}(x^1)}{\text{Var}(x^1)} = 1$$

- Value of Correlation Coef  $\Rightarrow$  -1 to 1 , if this is zero  $\Rightarrow$  the RV's are independent of each other  
  
→ If this is 1, -1  $\Rightarrow$  then RV's are highly dependent on each other.

How to find multicollinearity in data  $\Rightarrow$  Correlation Coeff  $\Rightarrow$  VIF

- \* So if we have 3 dimensions in data  
then to check Correlation b/w dimensions

we create a matrix

$\Rightarrow$  So if data has no multicollinearity then values in the Correlation Matrix shd be close to zero.

$$\rho_{x^2 x^1} = \rho_{x^1 x^2}$$

$\rho_{xx}$	$x^1$	$x^2$	$x^3$
$x^1$	$\rho_{x^1 x^1} = 1$	$\rho_{x^1 x^2}$	$\rho_{x^1 x^3}$
$x^2$	$\rho_{x^2 x^1}$	$\rho_{x^2 x^2} = 1$	$\rho_{x^2 x^3}$
$x^3$	$\rho_{x^3 x^1}$	$\rho_{x^3 x^2}$	$\rho_{x^3 x^3} = 1$



## Linear Regression

### Assumptions in Linear Regression

Correlation between two variables :

$$\text{Correlation} = \rho = \frac{\text{cov}(X,Y)}{\sigma_X \sigma_Y}$$

done



# Linear Regression



## Considering data of P Dimensions

### Lets Practice

Based on the data provided below, answer questions from (7-10). We consider a function we wish to minimize.

$J(w) = \frac{1}{10} \sum_{i=1}^5 (y^{(i)} - w_1 x^{(i)} - w_0)^2$  where the constants  $x^{(i)}$ ,  $y^{(i)}$  are provided in the table below

$$= \frac{1}{10} \sum_{i=1}^5 (y_i - w_1 x_i - w_0)^2$$

$$\underbrace{w_0, w_1 \geq 0.0}$$

$i$	$x^{(i)}$	$y^{(i)}$
1	0	1.4822
2	0.25	1.8165
3	0.50	1.9171
4	0.75	2.3930
5	1.00	2.5826

Dataset

- 7) The dimension of  $w$  is \_\_\_\_\_



## Linear Regression



### Considering data of P Dimensions

#### Lets Practice

- 8) Start with the initial guess of  $[w_0, w_1] = [0, 0]$ . Take the value of learning rate = 1. The value of  $w_0$  after 2 iterations of gradient descent will be \_\_\_\_\_.

Find  $w_0$  after 2 iteration



# Linear Regression

## What is Multicollinearity

- One crucial assumption in regression models is that independent variables should not correlate among themselves. This is essential for isolating the individual impact of each variable on the target variable, as indicated by regression coefficients.
- Multicollinearity arises when variables are correlated, making it challenging to discern their separate effects on the target variable.





# Linear Regression

What is  
Multicollinearity

- Example of multicollinearity :

Example



# Linear Regression



## What is Multicollinearity

- Why this is a problem ?
- Because in regression we are looking at how the independent variables are individually effecting the output label.

done

So Why we need VIF (Variance Inflation factor)

- \* in Correlation matrix we can see dependence of one dimension on other individually
- \* To see whole data together we use VIF

$x^1$	$x^2$	$x^3$	$y$
			No need

• Here we want to check linear dependence of one dimension on others

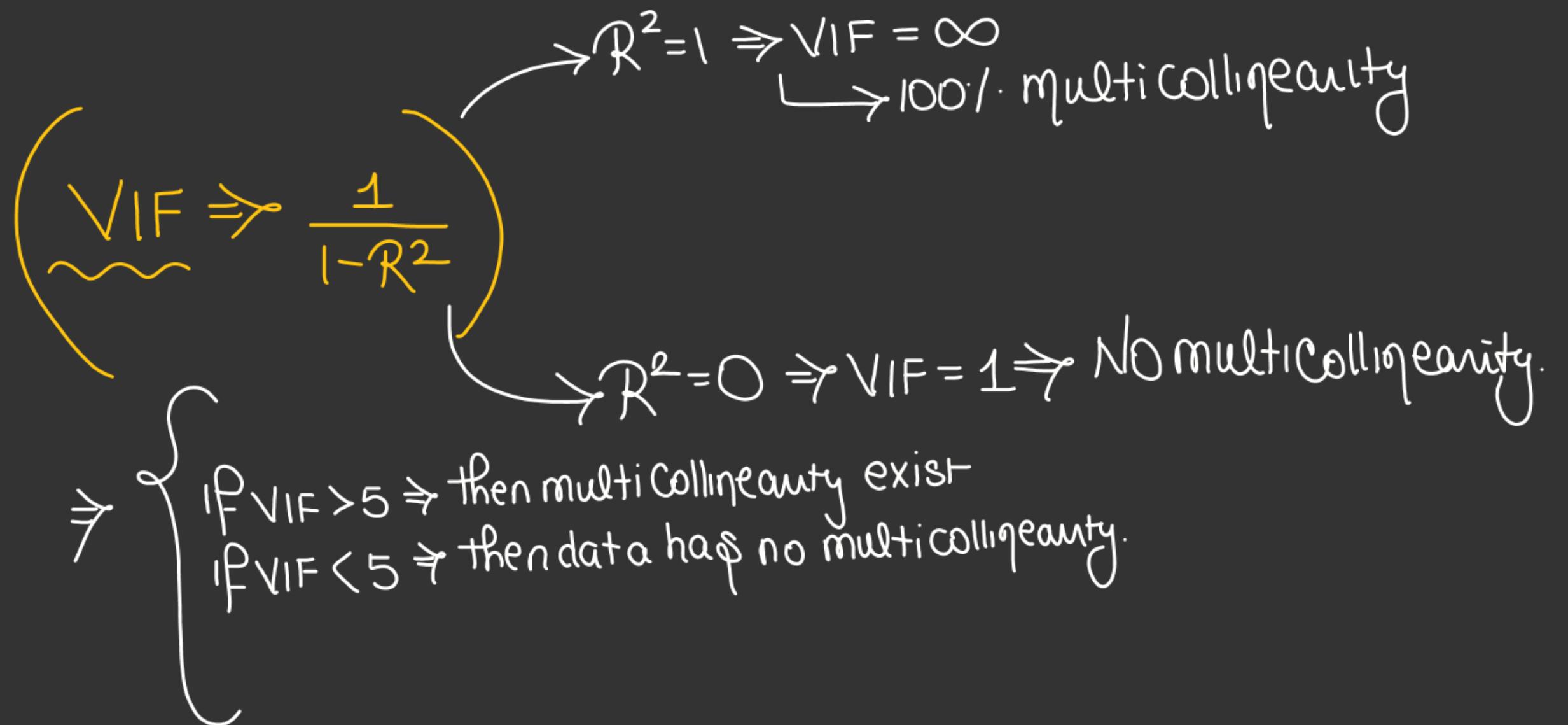
$$\text{If } x^1 = a + b x^2 + c x^3$$

So taking  $x^1$  as label now perform L-R  
 $(\underline{x}^1 = a + bx^2 + cx^3)$

↓  
model

Now  $R^2$  of this model  $(R^2 = 1) \Rightarrow$  100% dependent data  
If  $R^2 \neq 0$  the model is not able to understand pattern of  $x^1$ .

So independent data





# Linear Regression



## What is Multicollinearity

- How to solve the problem of multicollinearity ?

Dimension  $\geq$  features

\* By feature selection method we remove the highly dependent dimensions



## Linear Regression

### What is Multicollinearity

done

- Multicollinearity creates a problem in the multiple regression model because the inputs are all influencing each other. Therefore, they are not actually independent, and it is difficult to test how much the combination of the independent variables affects the dependent variable, or outcome, within the regression model.



## Linear Regression

How do we  
measure  
Multicollinearity?

done

- A very simple test known as the VIF test is used to assess multicollinearity in our regression model. The variance inflation factor (VIF) identifies the strength of correlation among the predictors.
- VIF help in predicting that which variable in the data is more correlated with other variables

## Linear Regression

How do we  
measure  
Multicollinearity?

### Formula and Calculation of VIF

The formula for VIF is:

$$\text{VIF}_i = \frac{1}{1 - R_i^2}$$

where:

$R_i^2 = \cancel{\text{coefficient of determination}}$  coefficient of determination for  
regressing the ith independent variable on the  
remaining ones

$\infty = M \cdot C \text{ exist}$

$1 \neq M \cdot C \text{ do not exist.}$



# Linear Regression



## Advantage & Disadvantage of Linear Regression

### Advantages

Linear Regression is simple to implement and easier to interpret the output coefficients.

When you know the relationship between the independent and dependent variable have a linear relationship, this algorithm is the best to use because of its less complexity compared to other algorithms.

Linear Regression is susceptible to over-fitting but it can be avoided using some dimensionality reduction techniques, regularization (L1 and L2) techniques and cross-validation.

### Disadvantages

On the other hand in linear regression technique outliers can have huge effects on the regression and boundaries are linear in this technique.

Diversely, linear regression assumes a linear relationship between dependent and independent variables. That means it assumes that there is a straight-line relationship between them. It assumes independence between attributes.

But then linear regression also looks at a relationship between the mean of the dependent variables and the independent variables. Just as the mean is not a complete description of a single variable, linear regression is not a complete description of relationships among variables.



## Linear Regression

**Why Linear Regression is Important**

**The interpretability of linear regression is a notable strength.**

**The model's equation provides clear coefficients that elucidate the impact of each independent variable on the dependent variable, facilitating a deeper understanding of the underlying dynamics.**

**Its simplicity is a virtue, as linear regression is transparent, easy to implement, and serves as a foundational concept for more complex algorithms.**



## Linear Regression



### Space and Time Complexity of Linear Regression

Assumptions:

**n** = number of training examples, **m** = number of features, **n'** = number of support vectors,  
**k** = number of neighbors, **k'** = number of trees

- **Linear Regression**

- Train Time Complexity= $O(n*m^2 + m^3)$
- Test Time Complexity= $O(m)$
- Space Complexity =  $O(m)$

Question 12: In simple linear regression, which variable is considered the independent variable? H.W.

- A. The variable being predicted
- B. The response variable
- C. The predictor variable
- D. There is no independent variable in simple linear regression

Question 19: If the R-squared value in simple linear regression is 0.75, what does it indicate?

- A. A strong linear relationship between the variables
- B. A weak linear relationship between the variables
- C. No linear relationship between the variables
- D. The model is overfitting

Question 20: Which of the following statements is true regarding the residual plot in simple linear regression?

- A. Residuals should exhibit a clear linear pattern.
- B. Residuals should be randomly scattered around the horizontal line.
- C. Residuals should be negatively correlated with the predictor variable.
- D. Residuals should have a positive correlation with the dependent variable.

5. For a given N independent input variables ( $X_1, X_2, \dots, X_n$ ) and dependent (target) variable Y a linear regression is fitted for the best fit line using least square error on this data. The correlation coefficient for one of its variables (Say  $X_1$ ) with Y is -0.97. Which of the following is true for  $X_1$ ?

- A) Relation between the  $X_1$  and Y is weak
- B) Relation between the  $X_1$  and Y is strong
- C) Relation between the  $X_1$  and Y is neutral
- D) Correlation does not imply relationship

6. Given below characteristics which of the following option is the correct for Pearson correlation between V1 and V2? If you are given the two variables V1 and V2 and they are following below two characteristics. 1. If V1 increases then V2 also increases 2. If V1 decreases then V2 behavior is unknown ?

- A) Pearson correlation will be close to 1
- B) Pearson correlation will be close to -1
- C) Pearson correlation will be close to 0
- D) None of these

- 1) A regression analysis is inappropriate when;
  - a) you have two variables that are measured on an interval or ratio scale.
  - b) you want to make predictions for one variable based on information about another variable.
  - c) the pattern of data points forms a reasonably straight line.
  - d) **there is heteroscedasticity in the scatter plot.**

- 2) In regression analysis, the variable that is being predicted is;
- a) the independent variable
  - b) the dependent variable**
  - c) usually denoted by x
  - d) usually denoted by r

- 3) In the regression equation  $y = b_0 + b_1x$ ,  $b_0$  is the;
- a) slope of the line
  - b) independent variable
  - c) y **intercept**
  - d) coefficient of determination

- 6) Least square method calculates the best-fitting line for the observed data by minimizing the sum of the squares of the \_\_\_\_\_ deviations.
- a) **Vertical**
  - b) Horizontal
  - c) Both of these
  - d) None of these

- 7) Which one is the least square method formula;
- a)  $\min \sum(y_i - \hat{y}_i)^2$
  - b)  $\min \sum(\hat{y}_i - y_i)$
  - c)  $\min \sum(y_i - \hat{y}_i)^2$
  - d)  $\min \sum(y_i - \hat{y}_i)$

13) Below you are given a summary of the output from a simple linear regression analysis from a sample of size 15,  $\text{SSR}=100$ ,  $\text{SST} = 152$ . The coefficient of determination is;

- a) 0.5200
- b) **0.6579**
- c) 0.8111
- d) 1.52

10) A residual is defined as

- a) The difference between the actual Y values and the mean of Y.
- b) The difference between the actual Y values and the predicted Y values.
- c) The predicted value of Y for the average X value.
- d) The square root of the slope.

11) If the regression equation is equal to  $y=23.6-54.2x$ , then 23.6 is the \_\_\_\_\_ while -54.2 is the \_\_\_\_\_ of the regression line.

- a) Slope, intercept
- b) Slope, regression coefficient
- c) Intercept, slope
- d) Radius, intercept

Q8. Suppose we have N independent variables ( $X_1, X_2 \dots X_n$ ) and Y's dependent variable.

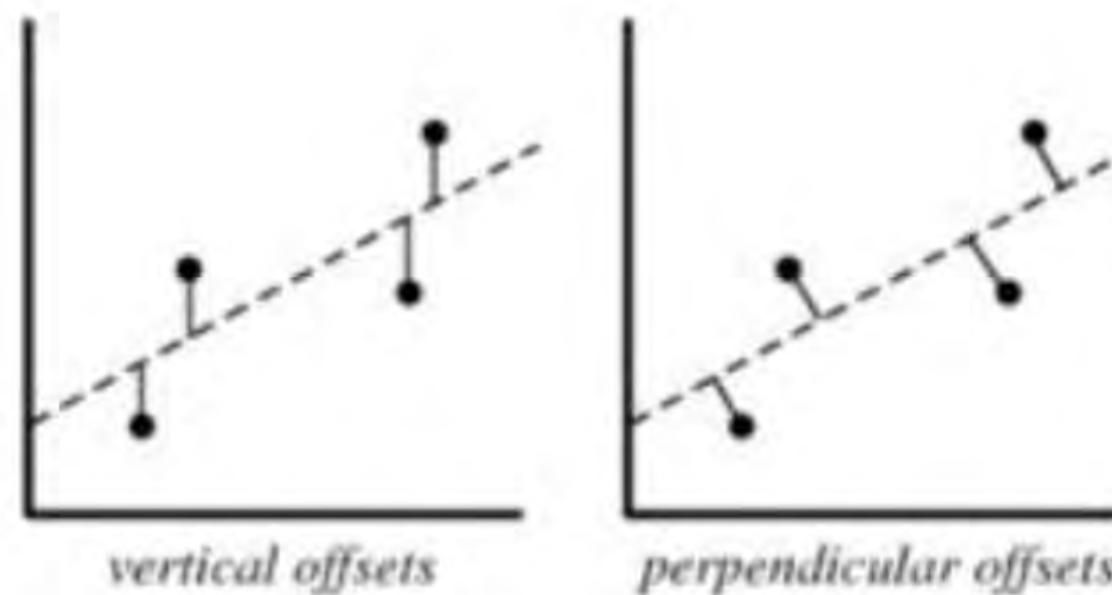
Now Imagine that you are applying linear [regression](#) by fitting the best-fit line using the least square error on this data. You found that the correlation coefficient for one of its variables (Say  $X_1$ ) with Y is -0.95.

**Which of the following is true for  $X_1$ ?**

- A) Relation between the  $X_1$  and Y is weak
- B) Relation between the  $X_1$  and Y is strong
- C) Relation between the  $X_1$  and Y is neutral
- D) Correlation can't judge the relationship

**Solution: (B)**

Q11. Suppose the horizontal axis is an independent variable and the vertical axis is a dependent variable. Which of the following offsets do we use in linear regression's least square line fit?



- B) Perpendicular offset
- C) Both, depending on the situation
- D) None of above

Q12. True- False: Overfitting is more likely when you have a huge amount of data to train.

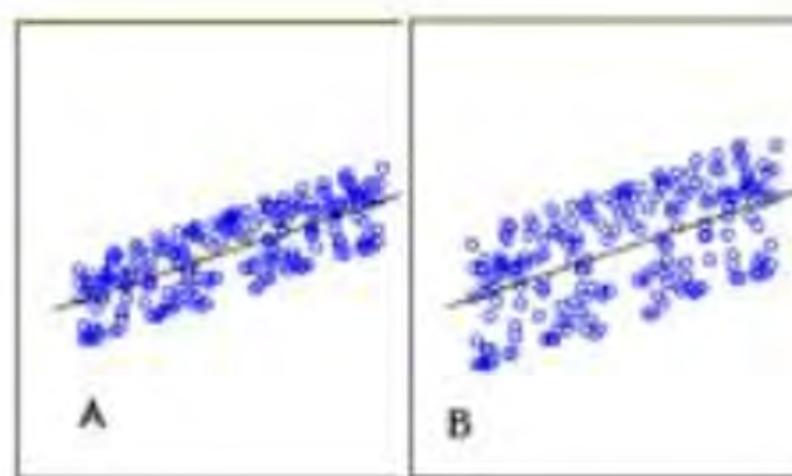
- A) TRUE
- B) FALSE

**Solution: (B)**

---

Q14. Which of the following statement is true about the sum of residuals of A and B?

Below graphs show two fitted regression lines (A & B) on randomly generated data. Now, I want to find the sum of residuals in both cases, A and B.



- A) A has a higher sum of residuals than B
- B) A has a lower sum of residual than B
- C) Both have the same sum of residuals
- D) None of these

Q18. Which of the following statement is true about outliers in Linear regression?

- A) Linear regression is sensitive to outliers
- B) Linear regression is not sensitive to outliers
- C) Can't say
- D) None of these

Q19. Suppose you plotted a scatter plot between the residuals and predicted values in linear regression and found a relationship between them. Which of the following conclusion do you make about this situation?

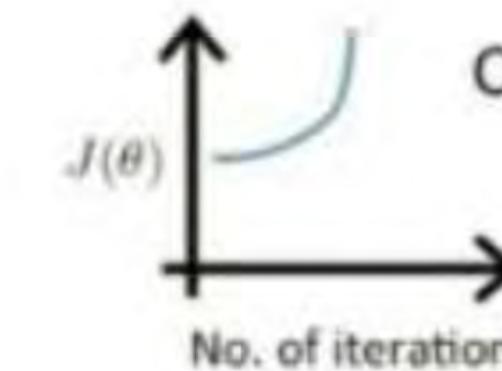
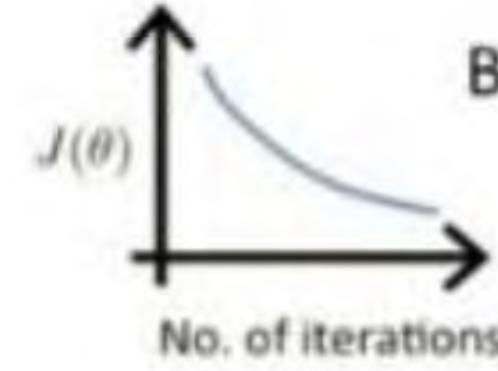
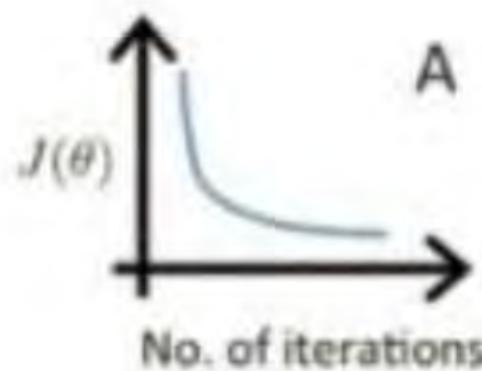
- A) Since there is a relationship means our model is not good
- B) Since there is a relationship means our model is good
- C) Can't say
- D) None of these

Suppose that you have a dataset D1 and you design a linear model of degree 3 polynomial and find that the training and testing error is "0" or, in other words, it perfectly fits the data.

**Q20. What will happen when you fit a degree 4 polynomial in linear regression?**

- A) There is a high chance that degree 4 polynomial will overfit the data
- B) There is a high chance that degree 4 polynomial will underfit the data
- C) Can't say
- D) None of these

Below are three graphs, A, B, and C, between the cost function and the number of iterations, I<sub>1</sub>, I<sub>2</sub>, and I<sub>3</sub>, respectively.



Q23. Suppose I<sub>1</sub>, I<sub>2</sub>, and I<sub>3</sub> are the three learning rates for A, B, and C, respectively. Which of the following is true about I<sub>1</sub>, I<sub>2</sub>, and I<sub>3</sub>?

- A) I<sub>2</sub> < I<sub>1</sub> < I<sub>3</sub>
- B) I<sub>1</sub> > I<sub>2</sub> > I<sub>3</sub>
- C) I<sub>1</sub> = I<sub>2</sub> = I<sub>3</sub>
- D) None of these

Intercept

$$SE(b_0) = \sqrt{\frac{\sum_{i=1}^n e_i^2}{n-2}} * \sqrt{\frac{1}{n} + \frac{\bar{x}^2}{\sum_{i=1}^n (x_i - \bar{x})^2}}$$

Slope

$$SE(b_1) = \sqrt{\frac{\sum_{i=1}^n e_i^2}{n-2} / \sum_{i=1}^n (x_i - \bar{x})^2}$$



THANK - YOU