# Linear Regression

1) What is Linear Regression.

2) What are the Types of Linear Regression.

3) Linear Regression Line.

4) Finding the Best Fit line.

5) Cost function

6) Gradient Decent

7) Model Performance

8) Assumptions of Linear Regression.

$\Rightarrow$ Linear Regression is one of the Easiest and most popular Machine Learning Algorithms.

$\Rightarrow$ It is a statistical method that is used for predictive analysis.

$\Rightarrow$ Linear Regression makes predictions for continuous / Real or numeric variables Such as Sales, Salary, age, product price, etc.

$\Rightarrow$ Linear Regression algorithm, shows a linear Relationship between a dependent (y) and one or more independent (y) variables, hence called as linear regression.

## Types of Linear Regression.

### i) Simple Linear Regression:

If a single independent variable is used to predict the value of a numerical dependant variable, then such a linear Regression algorithm is called simple Linear Regression.

### ii) Multiple Linear Regression:

If more than one independent variable is used to predict the value of a numerical dependent variable, then such a lineal Regression algorithm is called Multiple Linear Regression.
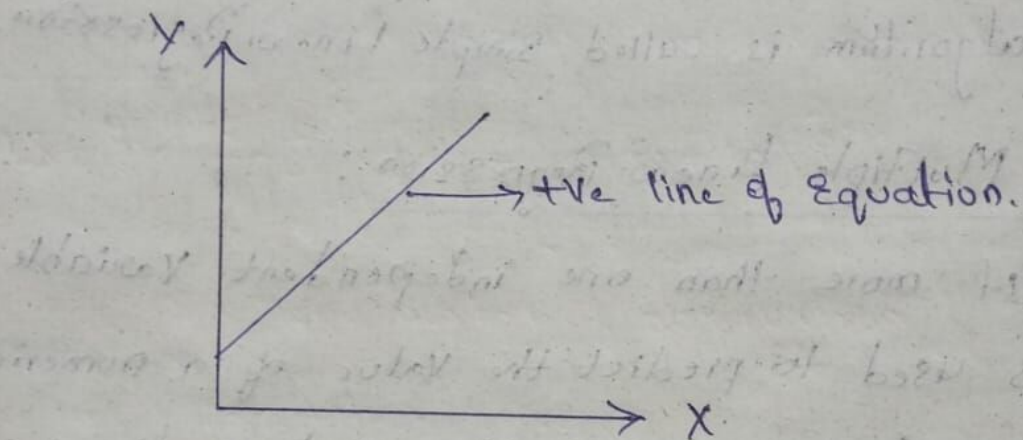
## Linear Regression Line

A linear line showing the Relationship between the dependent and independent variables is called a Regression line.

### Two Types of Regression line:-

(i) Positive Linear Relationship

(ii) Negative Lineal Relationship.
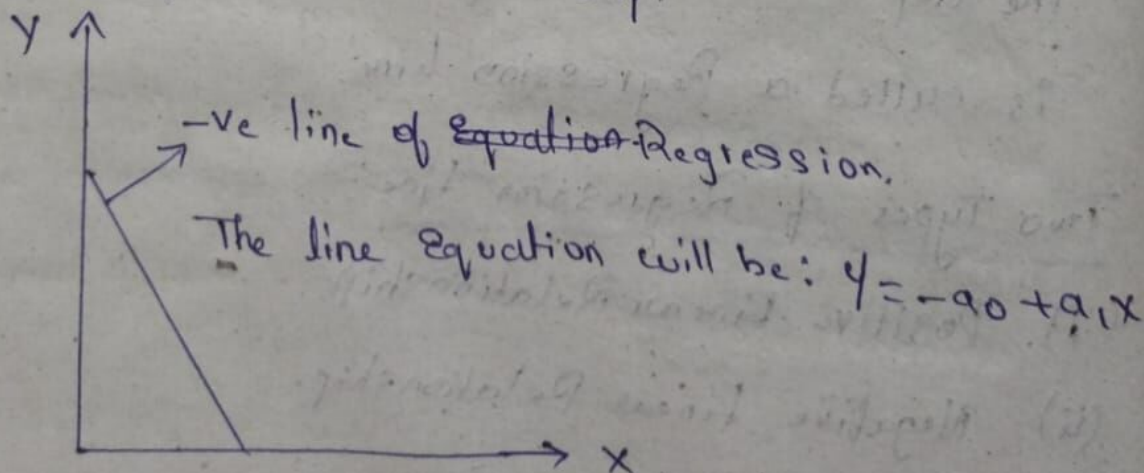
(i) Positive Linear Relationship.

if the dependent Variable increases, on the
Y-axis and independent Variable increases
on X-axis, then such a Relationship is
termed as a positive linear Relationship.

Y

→ +ve line of Equation.

→ X

∴ The line Equation will be : $Y = a_0 + a_1 x$

ii) Negative Linear Relationship:

if the dependent Variable decreases on the
Y-axis and independent Variable increases on
the X-axis, then such a relationship is called
a Negative linear Relationship.

Y

-ve line of ~~Equation~~ Regression.

The line Equation will be: $Y = -a_0 + a_1 x$

→ X

# Mean Squared Error (MSE)

For Linear Regression, we use the (MSE)
Mean Squared Error Cost function, which is
the average of Squared Error occurred b/w
the predicted values and actual values.

$$MSE = \frac{1}{N} \sum_{i=1}^{n} \left( y_i - (a_1 x_i + a_0) \right)^2$$

$N$ = Total number of observation

$y_i$ = Actual Value

$(a_1 x_i + a_0)$ = Predicted Value.

## Residuals:

The distance between the actual value and predicted
value is called Residual. If the observed points
are far from the Regression line, then the residual will be high, and so Cost function
will high. if the Scatter point are close to
the Regression line, then the residual will be
so small and hence the cost function.

## Gradient Decent

⟹ Gradient decent is used to minimize the MSE by calculating the gradient of the cost function.

⟹ A regression model uses gradient descent to upgra update the coefficients of the line by reducing the cost function.

⟹ It is done by random selection of values of coefficients and then iteratively update the values to reach the minimum cost function.

## Model Performance

⟹ The Goodness of fit determines how the line of regression fits the set of observations. The process of finding the best model out of various models is called optimization.

# Assumptions of Linear Regression :-

⟹ **Linear relationship between the features and target:**

- Linear regression assumes the linear relationship between the dependend & independent variables.

⟹ **Small or no Multicollinearity between the features:**
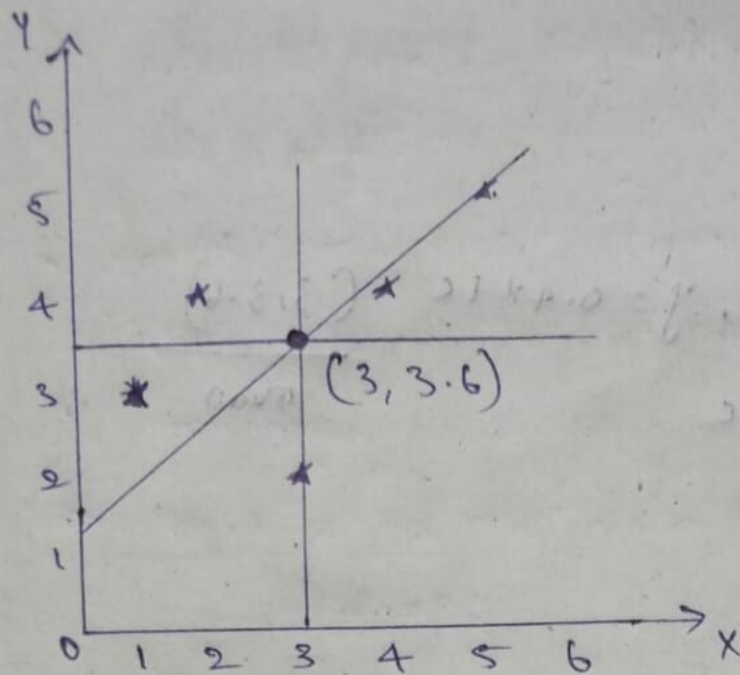
- Multicollinearity means high correlation b/w the independent variables. Due to multicollinearity, it may difficult to find the true relationship b/w the predictors and target variables. Or we can say, it is difficult to determine which predictor variable is affecting the target variables and, which is not. So, the model assumes either little or no multicollinearity b/w the features or independent variables.

⟹ **Homoscedasticity Assumption:**

Homoscedasticity is a situation when the error term is the same for all the values of independent variables. with homoscedasticity, there should be no clear pattern distibution. of data in the scatter plot.

## Linear Regression Example:-



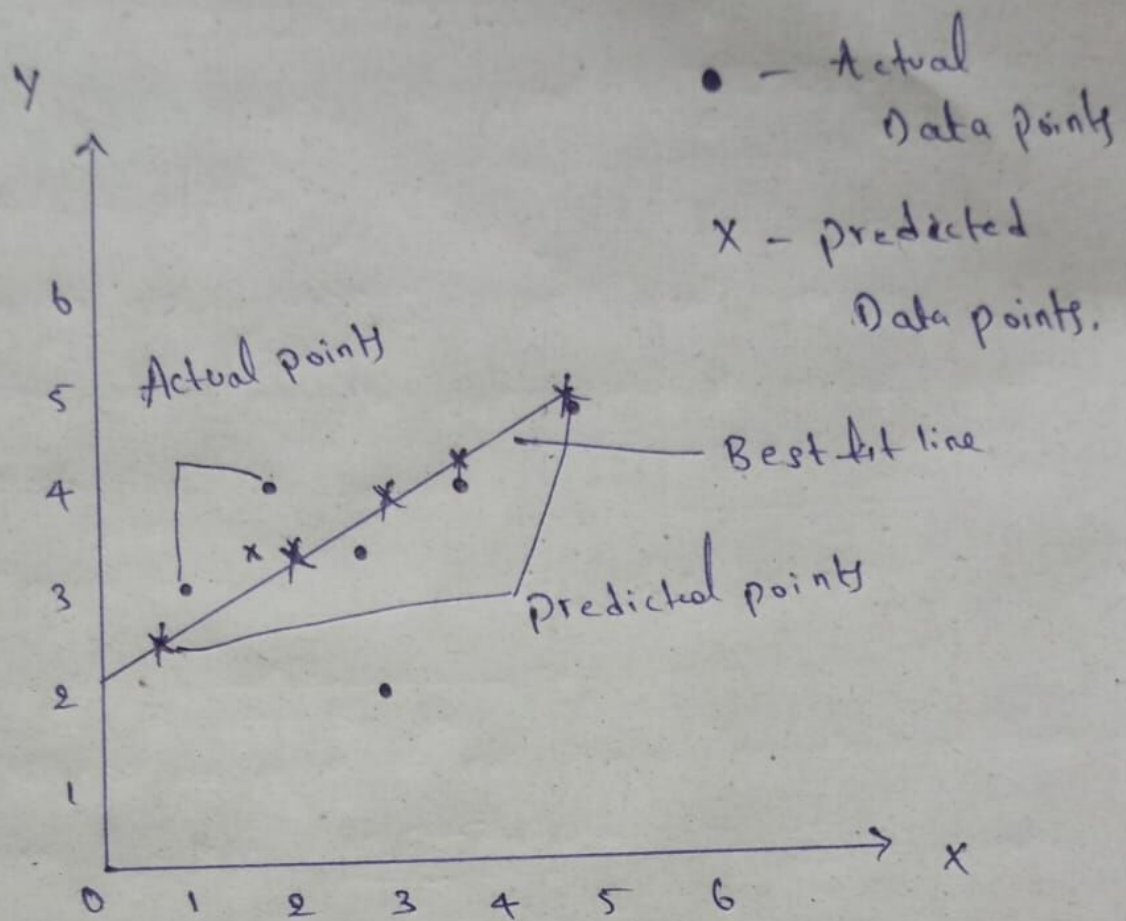| X | Y |
|---|---|
| 1 | 3 |
| 2 | 4 |
| 3 | 2 |
| 4 | 4 |
| 5 | 5 |

Mean of $x = \dfrac{1+2+3+4+5}{5} = \dfrac{15}{5} = 3$

Mean of $Y = \dfrac{3+4+2+4+5}{5} = \dfrac{18}{5} = 3.6$

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

$$\boxed{m = \dfrac{\Sigma\,(x-\bar{x})(y-\bar{y})}{\Sigma\,(x-\bar{x})^2}}$$

| X | y | $x-\bar{x}$ | $y-\bar{y}$ | $(x-\bar{x})^2$ | $(x-\bar{x})(y-\bar{y})$ |
|---|---|---|---|---|---|
| 1 | 3 | -2 | -0.6 | 4 | 1.2 |
| 2 | 4 | -1 | 0.4 | 1 | -0.4 |
| 3 | 2 | 0 | -1.6 | 0 | 0 |
| 4 | 4 | 1 | 0.4 | 1 | 0.4 |
| 5 | 5 | 2 | 1.4 | 4  $\Sigma=10$ | 2.8  $\Sigma=4$ |

Therefore we got Best fit line, The distance

Between Actual points and predicted points is

error.

$$m = \frac{4}{10} = 0.4$$

$$m = 0.4$$

$$y = mx + c$$

$$y = \textcircled{m} \rightarrow 0.4, \quad y = 0.4x + c \qquad (3, 3.6)$$
$$\underline{\qquad\qquad} \text{mean}$$

$$3.6 = 0.4(3) + c$$

$$3.6 = 1.2 + c$$

$$3.6 - 1.2 = c$$

$$\underline{\underline{c = 2.4}}$$

$$\underline{y = 0.4x + 2.4} \qquad \text{Regression line}$$

(i) $x = 1$

$$y = 0.4(1) + 2.4 = 2.8$$

(ii) $x = 2$

$$y = 0.4(2) + 2.4 = 3.2$$

(iii) $x = 3$

$$y = 0.4(3) + 2.4 = 3.6$$

(iv) $x = 4$

$$y = 0.4(4) + 2.4 = 4.0$$

(v) $x = 5$

$$y = 0.4(5) + 2.4 = 4.4$$

⟹ Normal distribution of error terms:

• Linear regression assumes that the error term should follow the normal distribution pattern.

• if error terms are not normally distributed, then confidence intervals will become either too wide or too narrow, which may cause difficulties in finding coefficients.

• It can be checked using the q-q plot.

• If the plot shows a straight line without any deviation, which means the error is normally distributed.

⟹ No Auto correlations:

• The linear regression model assumes no autocorrelation in error terms. if there will be any correlation in the error term, then it will drastically reduce the accuracy of the # model. Autocorrelation usually occurs if there is a dependency between residual errors.

# R-Squared method

→ R-Squared is a statistical method that determines the goodness of fit.

→ It measures the strength of the relationship between the dependent and independent variables on a scale of 0 - 100%.

→ The high value of R-square determines the less difference between the predicted values and actual values and hence represents a good model.

→ It is also called a coefficient of determination or coefficient of multiple determination for multiple regression.

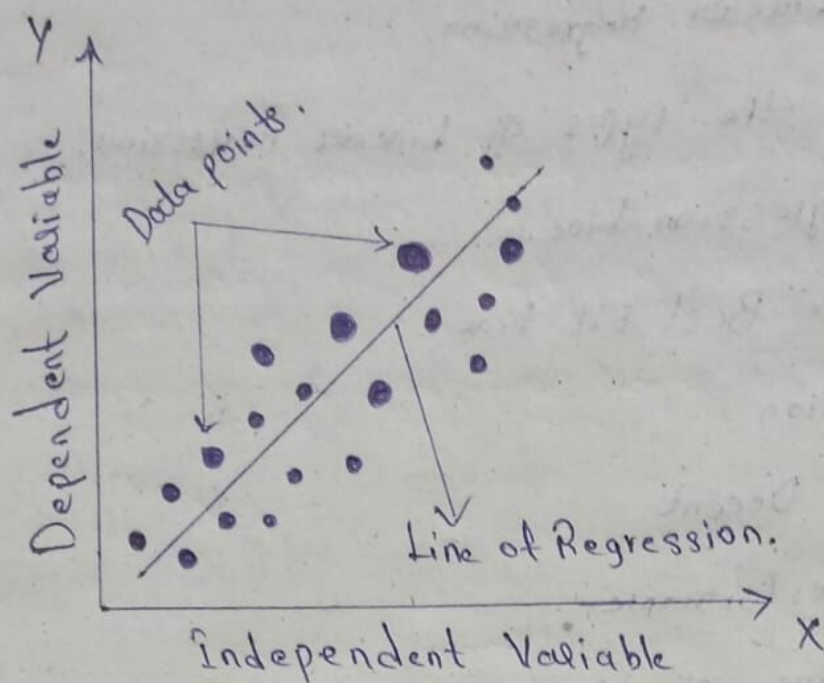$$R\text{-Squared} = \frac{\text{Explained Variation}}{\text{Total Variation}}$$

# Finding the Best Fit line

* our main goal is to find the Best fit line it is the Error between predicted Values and actual Values should be minimized.

* The Best fit line will have the least Error.

* The different Value for weights or the coefficient of lines, $(a_0, a_1)$ gives a different line of regression, So we need to calculate the Best Values for $a_0$ & $a_1$ to find the Best fit line, So to calculate this we use cost function.

## Cost function.

⇒ The cost function is used to find the accuracy of mapping function, which maps the input variable to the output variable.

This mapping function is also known as Hypothesis function.

Mathematical Representation.

$$y = a_0 + a_1 x + \varepsilon$$

- $y$ = Dependent Variable (Target Variable)

- $x$ = Independent Variable (Predictor Variable)

- $a_0$ = intercept of the line
  (Gives an additional degree of freedom)

- $a_1$ = Linear regression Coefficient
  (Scale factor to Each input value)

- $\varepsilon$ = Random Error.
  ($\varepsilon$ = epsuloin)