# REGRESSION ANALYSIS AND TIME SERIES

## UNIT-I

Simple Linear Regression Model :

the functional relation of a dependent variable with one or more independent variable is called regression equation. If the relationship between the independent variable x and dependent variable y. it is known as simple linear regression.

The Straight line equation is

$$Y = \beta_0 + \beta_1 x$$

Where,

$\beta_0$ is intercept.

$\beta_1$ is slope.

$x$ is independent variable

$Y$ is dependent variable.

Let the difference between the observed value of y and the straight line $\beta_0 + \beta_1 x$ is an error, it is known and random error and is denoted by E,

Thus the simple linear regression model is written as

$$Y = \beta_0 + \beta_1 x + E$$

the errors are assumed to have mean Zero and Variance $\sigma^2$. the parameters $\beta_0$ and $\beta_2$ are regression Co-efficient.

Estimation of $\hat{\beta_0}$ and $\hat{\beta_1}$ using Linear Square Method:

Suppose that we have n pairs data say $(x_1, y_1)$, $(x_2, y_2) \ldots (x_n, y_n)$. the method of least square is used to estimate $\beta_0$ and $\beta_1$., so that the sum of the squares of the difference between |t| observation $y_i$ and the straight line minimum thus the least square estimation is

$$S(\beta_0, \beta_1) = \sum_{i=1}^{n} (y_i - \hat{\beta_0} - \hat{\beta_1} x_i)^2$$

the LS estimates $\beta_0$ and $\beta_1$ say $\hat{\beta_0}$ and $\hat{\beta_1}$ must Satisfy differentiating,

(ii) $\quad \dfrac{\partial S}{\partial \beta_0} \Big|_{\hat{\beta_0}, \hat{\beta_1}} = -2 \sum_{i=1}^{n} (y_i - \hat{\beta_0} - \hat{\beta_1} x_i) = 0 \longrightarrow (1)$

and

$\quad \dfrac{\partial S}{\partial \beta_1} \Big|_{\hat{\beta_0}, \hat{\beta_1}} = -2 \sum_{i=1}^{n} (y_i - \hat{\beta_0} - \hat{\beta_1} x_i) = 0 \longrightarrow (2)$

from (1)

$$\sum_{i=1}^{n} (y_i - \hat{\beta_0} - \hat{\beta_1} x_i) = 0$$

$$\Rightarrow \Sigma y_i = n\hat{\beta_0} + \hat{\beta_1} \Sigma x_i \longrightarrow (3)$$

From (2) multiply by $x_i$

$$\sum_{i=1}^{n} (y_i - \hat{\beta_0} - \hat{\beta_1} x_i) x_i = 0$$

$$\Rightarrow \Sigma x_i y_i = \hat{\beta_0} \Sigma x_i + \hat{\beta_1} \cdot \sum_{i=1}^{n} x_i^2 \rightarrow (4)$$

Equ (3) & (4) are called least square normal equations.

To estimate the value of $\hat{\beta_0}$ and $\hat{\beta_1}$

Consider,

$$\Sigma y_i = n\hat{\beta_0} + \hat{\beta_1} \Sigma x_i \rightarrow (4)$$

Multiply $\Sigma x_i$ on both sides,

$$\Sigma x_i \Sigma y_i = n\hat{\beta_0} \Sigma x_i + \hat{\beta_1} (\Sigma x_i)^2 \rightarrow (5)$$

Multiply Equ (5) by $n$ we get.

$$n \Sigma x_i y_i = n\hat{\beta_0} \Sigma x_i + n\hat{\beta_1} (\Sigma x_i)^2 \rightarrow (6)$$

$(6) - (5) \Rightarrow$

$$n \Sigma x_i y_i - \Sigma x_i \Sigma y_i = n\hat{\beta_1} \left[ \Sigma x_i^2 - y_n (\Sigma x_i)^2 \right]$$

$$\cancel{n} \left[ \Sigma x_i y_i - y_n \Sigma x_i y_i \right] = \cancel{n} \hat{\beta_1} \left[ \Sigma x_i^2 - y_n (\Sigma x_i)^2 \right]$$

$$\therefore \hat{\beta_1} = \frac{\Sigma x_i y_i - y_n \Sigma x_i y_i}{\Sigma x_i^2 - y_n (\Sigma x_i)^2}$$

$\div$ equation (3) by $n$, we get,

$$\frac{\Sigma y_i}{n} = \hat{\beta_0} + \hat{\beta_1} \frac{\Sigma x_i}{n} \qquad \therefore \bar{y} = \frac{\Sigma y_i}{n}$$

$$\hat{\beta_0} = \bar{y} - \hat{\beta_1} \bar{x} \qquad \therefore \bar{x} = \frac{\Sigma x_i}{n}$$

$$S_{xx} = \Sigma x_i^2 - \frac{(\Sigma x_i)^2}{n}$$

$$= \Sigma (x_i - \bar{x})^2$$

$$S_{xy} = \Sigma x_i y_i - \frac{\Sigma x_i \Sigma y_i}{n}$$

$$= \Sigma y_i (x_i - \bar{x})$$

$$\therefore \hat{\beta_1} = \frac{S_{xy}}{S_{xx}} - \sum_{i=1}^{n} c_i y_i$$

where,

$$c_i = \left( \frac{x_i - \bar{x}}{S_{xx}} \right) \qquad (i = 1, 2 \cdots n).$$

Statistical Assumptions in linear model:

Let $E_i = M_i - \bar{M}$ where $M_i$ is the disturbance term.

i) Normality:

The disturbance term are normally distributed
$$M_i \sim N(\sigma, \sigma^2)$$

ii) Zero mean:

The mean of $M_i$ is zero (ie) $E(M_i) = 0$
$$\forall i = 1, 2 \cdots n$$

iii) Homo cedasiticity:

Every disturbance variance has the same variance, whose value is unknown.

(ie) $E(u_i^2) = \sigma_u^2$, $\forall i = 1, 2 \cdots n$.

iv) Non Auto Regression:

The various disturbance terms are uncorrelated.

(ie) $E(M_i M_j) = 0$ for $i \neq j$ $i, j = 1, 2, \cdots n$.

v) Non - Stochastic:

If $x - x_i$ is a non-Stochastic variable Such that for each sample of size (n).

Properties of least Square estimators for Simple linear regression model.

i) the least Square estimators $\hat{\beta}_0$ and $\hat{\beta}_1$ are unbiased estimators of the model parameters $\beta_0$ and $\beta_1$.

Proof :

$$E(\hat{\beta}_1) = E\left[\sum_{i=1}^{n} c_i y_i\right]$$

$$= \sum_{i=1}^{n} c_i \, E(y_i)$$

$$= \sum_{i=1}^{n} c_i \, (\beta_0 + \beta_1 x_i)$$

$$= \beta_0 \sum_{i=1}^{n} c_i + \beta_1 \sum_{i=1}^{n} c_i x_i$$

But We know that,

$$\frac{\sum_{i=1}^{n} (x_i - \bar{x})}{S_{xx}} = \sum_{i=1}^{n} c_i = 0 \quad \text{and}$$

$$\frac{\sum_{i=1}^{n} (x_i - \bar{x}) x_i}{S_{xx}} = \frac{\sum x_i^2 - \bar{x} \sum x_i}{S_{xx}}$$

$$= \frac{\sum x_i^2 - \frac{\sum x_i \, \sum x_i}{n}}{}$$

$$= 1$$

$$\therefore E(\hat{\beta_1}) = 0 + \beta_1 v_1$$

$$E(\hat{\beta_1}) = \beta_1$$

$$E(\hat{\beta_0}) = E[\bar{Y} - \hat{\beta_1}\bar{x}]$$

$$= E(\bar{Y}) - \hat{\beta_1} E(\bar{x})$$

$$= \bar{y} - \beta_1 \bar{x}$$

$$= \beta_0$$

ii) the least square estimators are linear estimators.

We know that,

$$\hat{\beta_0} = \bar{Y} - \hat{\beta_1}.\bar{x} \text{ and}$$

$$\hat{\beta_1} = \dfrac{\Sigma y_i x_i - \dfrac{(\Sigma x_i)(\Sigma y_i)}{n}}{\Sigma x_i^2 - \dfrac{(\Sigma x_i)^2}{n}}$$

are linear combinations of the observations $y_i$, thus $\hat{\beta_0}$ and $\hat{\beta_1}$ are linear estimators.

iii) Least Square estimators are Best estimators.

$$V(\hat{\beta_1}) = V\left(\sum_{i=1}^{n} c_i y_i\right)$$

$$= \sum_{i=1}^{n} c_i^2 V(y_i)$$

because the observations $y_i$ are uncorrelated and so the variance of the sum is just the sum of the variance of each term in the sum is

$$c_i^2 V(y_i) \text{ and}$$

We have assumed $V(y_i) = \sigma^2$

$$V(\hat{\beta_1}) = \sigma^2 \, \Sigma \, c_i^2$$

$$= \sigma^2 \, \Sigma \, \frac{(x_i - \bar{x})^2}{S_{xx}^2}$$

$$= \frac{\sigma^2}{S_{xx}}, \quad S_{xx} = \Sigma (x_i - \bar{x})^2$$

The variance of $\hat{\beta_0}$ is,

$$V(\hat{\beta_0}) = V(\bar{y} - \hat{\beta_1} \, \bar{x})$$

$$= V(\bar{y}) + \bar{x}^2 \, V(\hat{\beta_1}) - 2\bar{x} \, \text{Cov}(\bar{y}, \hat{\beta_1})$$

We know that,

$$V(\bar{y}) = \sigma^2/n \quad \text{and} \quad \text{Cov}(\bar{y}, \hat{\beta_1}) = 0$$

$$\therefore V(\hat{\beta_0}) = V(\bar{y}) + \bar{x}^2 \, V(\hat{\beta_1})$$

$$= \sigma^2 \left( \frac{1}{n} + \frac{\bar{x}^2}{S_{xx}} \right)$$

Thus the least square estimators are unbiased linear and have minimum variance when compared with all other unbiased estimators.

Thus the least Square estimators are best linear unbiased estimators.

Other properties of the least Square fit.

1, the Sum of the residues is any regression model that Contains an intercept $\beta_0$ is always zero.

(ie) $\sum(y_i - \hat{y}_i) = \sum\limits_{i=1}^{n} e_i = 0$

2) the Sum of the observed values $y_i$ equals the Sum of the fields values $\hat{y}$.

(ie) $\sum\limits_{i=1}^{n} y_i = \sum\limits_{i=1}^{n} \hat{y}_i$

3) the Least Square regression line always passes through the Centeroid of the data.

4) the Sum of the residuals weighted by the corresponding value of the regression variate always equal.

(ie) $\sum\limits_{i-1}^{n} x_i e_i = 0$

5) the Sum of the residuals weighted by the corresponding fitted value always zero, that is

$\sum\limits_{i=1}^{n} \hat{y} e_i = 0.$

## Estimation of $\sigma^2$.

the estimate of $\sigma^2$ is obtained from the residual (or) error sum of squares.

$$SS_{Res} = \Sigma e_i^2 = \Sigma (y_i - \hat{y_i})^2$$

W.K.T

$$\hat{y_i} = \hat{\beta_i}, \hat{\beta_i} x_i$$

$$SS_{Res} = \sum_{i=1}^{n} (y_i - \hat{\beta_0} - \hat{\beta_i} x_i)^2$$

It can be written as,

$$= \Sigma y_i^2 - n\bar{y}^2 - \hat{\beta_i} S_{xy}$$

But,

$$\Sigma y_i^2 - n\bar{y}^2 = \Sigma (y_i - \bar{y})^2 = SS_T$$

is the corrected sum of squares of the response observations s

$$SS_{Res} = SS_T - \hat{\beta_i} S_{xy}$$

The residuals sum of squares has $n-2$ d.f because two degrees of freedom are associated with the estimate $\hat{\beta_0}$ and $\hat{\beta_i}$ involved in obtaining $\hat{y_i}$

So an unbiased estimators of $\sigma^2$ is $\hat{\sigma}^2 = \dfrac{SS_{Res}}{n-2} = MS_{Res}$

Where $M_{Res}$ is called the standard error of regression. The square root $r\hat{\sigma}^2$ is called the standard error of regression.

# Hypothesis Testing on the Slope.

Null hypothesis :

$H_0$ : The Slope equals a Constant Say $\beta_0$

(ie) $H_0 : \beta_1 = \beta_{10}$

Alternative hypothesis :

$H_1$ : The Slope does not equals a Constant Say $\beta_1$

(ie) $H_1 : \beta_1 \neq \beta_{10}$

Test Statistic,

under $H_0$,

$$t = \frac{\hat{\beta_1} - \beta_{10}}{\sqrt{MS_{Res}/S_{xx}}}$$

Where,

$$MS_{Res} = \frac{\Sigma y_i^2 - n\bar{y}^2 - \beta_1 S_{xy}}{n-2}$$

$$S_{xy} = \Sigma y_i x_i - \frac{\Sigma y_i \, \Sigma x_i}{n}$$

$$= \Sigma y_i (x_i - \bar{x})^2$$

$$S_{xx} = \Sigma (x_i - \bar{x})^2$$

Table Value:

Refer Student t-table for $(n-2)$ of f at $\alpha$% Level.

Inference :

If Calculated Value $<$ table Value do not reject $H_0$
(ie) $H_0 : \beta_1 = \beta_{10}$ , otherwise reject $H_0$

Hypothesis Testing for intercept $\beta_0$:

Null hypothesis:
$$H_0 : \beta_0 = \beta_{00}$$

Alternative hypothesis:
$$H_1 : \beta_0 \neq \beta_{00}$$

Test Statistic:

under $H_0$,

$$t = \frac{\hat{\beta_0} - \beta_{00}}{\sqrt{M_{Res} \left( \frac{1}{n} + \frac{\bar{x}^2}{S_{xx}} \right)}}$$

Table Value

Refer Studen's $t$-table for $(n-2)$ d.f at $\alpha\%$ level

Inference:

If the Calculated Value $<$ table Value do not reject $H_0$.

(Le) $H_0 : \beta_0 = \beta_{00}$

otherwise reject $H_0$.

# Estimation By Maximum Likelihood Method:

The method of least square can be used to estimate the parameters in a linear regression model regard sets of the form of the distribution of the errors E.

If the form of the distribution of the errors is known as alternative method of Parameter estimation is method of maximum likelihood.

Consider the data $(y_i, x_i)$ $(i=1,2\ldots n)$

if we assume that the errors in the regression models are Normally identically distributed $(0, \sigma^2)$. then the observations $r_i$ in the Sample are normally and identically distributed random Variables with means $\beta_0 + \beta_1 x$ and variance $\sigma^2$ For simple linear regression model with normal errors, the Likelihood function is,

$$L(y_i, x_i, \beta_0, \beta_1, \sigma^2) = \prod_{i=1}^{n} (2\pi\sigma^2)^{-1/2}$$

$$\exp\left\{ -\frac{1}{2\sigma^2} (y_i - \beta_0 - \beta_1 x_i)^2 \right\}$$

$$= (2\pi\sigma^2)^{-n/2} \exp\left\{ \frac{-1}{2\sigma^2} \Sigma(y_i - \beta_0 - \beta_1 x_i^2) \right\}$$

the MLE are the parameters values say,

$\hat{\beta_0}, \hat{\beta_1}$ and $\hat{\sigma_1}^2$ that maximise $L$ or $\ln L$

thus,

$$\ln L(y_i, x_i, \beta_0, \beta_1, \sigma^2) - (n/2) \ln \cdot 2\pi$$

$$- n/2 \ln \sigma^2 - \left( \frac{1}{2\sigma_2} \sum (y_i - \beta_0 - \beta_1 x_i)^2 \right\}$$

$$\frac{\partial \ln L}{\partial \beta_0} \bigg|_{\hat{\beta_0}, \hat{\beta_1}, \hat{\sigma}^2} = \frac{1}{\hat{\sigma}^2} \sum_{i=1}^{n} (\bar{y_i} - \hat{\beta_0} - \hat{\beta_1} x_i)^2 = 0$$

$$\frac{\partial \ln L}{\partial \beta_1} \bigg|_{\hat{\beta_0}, \hat{\beta_1}, \hat{\sigma}^2} = \frac{1}{\hat{\sigma}^2} \sum_{i=1}^{n} (y_i - \hat{\beta_0} - \hat{\beta_1} x_i) = 0$$

and,

$$\frac{\partial \ln L}{\partial \sigma^2} \bigg|_{\hat{\beta_0}, \hat{\beta_1}, \hat{\sigma}^2} = \frac{-n}{2\hat{\sigma}^2} = \frac{-n}{2\hat{\sigma}^2} + \frac{1}{2\hat{\sigma}} \sum_{i=1}^{n} (y_i - \hat{\beta_0} - \hat{\beta_1} x_i)$$

# UNIT - II

## Multiple Linear Regression Model.

A Regression model that involves more than one regression Variable is called a multiple regression model

the multiple regression model is,

$$Y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + E$$

this is a multiple linear regression model with two regression Variables.

In general the model,

$$Y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \cdots + \beta_k x_k + E$$

is Called a multiple regression model with k regression.

the parameters $\beta_j$ $(j=1,2 \dots k)$ are called regression Co-efficient.

Data for multiple Regression Model.

| Observation | frequency | Regressions |
|---|---|---|
| 1 | $y_1$ | $x_1, x_2 \dots x_k$ |
| 2 | $y_2$ | $x_{11}, x_{12} \dots x_{1k}$ |
| $\vdots$ | $\vdots$ | $\vdots$ |
| $n$ | $y_n$ | $x_{n1} \quad x_{n2} \quad x_{nk}$ |

The regression model Can be written as,

$$y_1 = \beta_0 + \beta_1 x_i + \beta_2 x_{12} + \dots + \beta_k x_{ik} + \epsilon_i .,$$

$$= \beta_0 + \sum_{j=1}^{k} \beta_j x_{ij} + \epsilon_i \quad (i=1,2\dots n).$$

Estimation of the model Parameters least Square estimation of the regression Co-efficient.

the multiple linear regression model is

$$y_i = \beta_0 + \sum_{j=1}^{k} \beta_j x_j + \epsilon_j \quad (i=1,2\dots n)$$

The linear square function is,

$$S(\beta_0, \beta_1, \dots \beta_k) = \sum_{i=1}^{n} \epsilon_i^2$$

$$= \sum_{i=1}^{n} (y_i - \beta_0 - \sum_{j=1}^{n} \beta_j x_{ij})$$

$$\frac{\partial S}{\partial \beta_0}\bigg|\hat{\beta_0},\hat{\beta_1}, \ldots \hat{\beta_k} = -2 \sum_{i=1}^{n} (y_i - \hat{\beta_0} - \sum_{j=1}^{k} \hat{\beta_j} x_{ij}) = 0 \quad \rightarrow (1)$$

and

$$\frac{\partial S}{\partial \beta_j}\bigg|\hat{\beta_0},\hat{\beta_1}, \ldots \hat{\beta_k} = -2 \sum_{i=1}^{n} (y_i - \hat{\beta_0} = \sum_{j=1}^{k} \hat{\beta_j} x_{ij}) x_{ij} = 0 \rightarrow (2)$$

when.

$$Y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_k \end{bmatrix}; \quad X = \begin{bmatrix} 1 & x_{11} & x_{12} & \cdots & x_{1k} \\ 1 & x_{21} & x_{22} & \cdots & x_{2k} \\ & & \vdots & & \\ 1 & x_{n1} & x_{n2} & \cdots & x_{nk} \end{bmatrix}$$

$$\beta = \begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_k \end{bmatrix} \qquad \epsilon = \begin{bmatrix} \epsilon_1 \\ \epsilon_2 \\ \vdots \\ \epsilon_n \end{bmatrix}$$

$$(\beta) = \sum_{i=1}^{n} \epsilon_i^2 = \epsilon' \epsilon (y - x\beta)' (y - x\beta)$$

$$= y'y - 2\beta' x' y + \beta' y' x_\beta$$

the linear Square estimators must satisfy,

$$\frac{\partial S}{\partial \beta}\bigg|\hat{\beta} = -2x'y + 2x' x \hat{\beta} = 0,$$

$$\Rightarrow -2(x'y - x'x\beta) = 0$$

$$\Rightarrow x'y \, x'x = x'x\beta'.$$

$$\hat{\beta} = (x'x)^{-1} (x'y)$$

Properties of least Square estimators for multiple linear model.

i) The least Square estimators are unbiased estimators.

We know That,

$$\hat{\beta} = (x'x)^{-1} x'y$$

$$E(\hat{\beta}) = E((x'y)^{-1} xy)$$

$$= E\left[(x'x)^{-1} x' x\beta' (x'x)^{-1} x'\epsilon\right]$$

$$= \beta$$

$$\therefore E(\epsilon) = 0.$$

ii) The least Square estimators are linear estimator.

Since the least square estimators are,

$$\hat{\beta} = (x'x)^{-1} x'y$$

$$\Rightarrow (x'x)^{-1} x' \left[(x\beta + \epsilon)\right]$$

$$\Rightarrow (x'x)^{-1} x'x\beta + (xx^{-1})^{-1} x'\epsilon$$

$$= \beta + (x'x)^{-1} x'\epsilon$$

iii) Least square estimates are the last estimators.

Proof :-

We know that

$$\hat{\beta} = (x'x)^{-1} x'y$$

$$= (x'x)^{-1} x' (x\beta + \epsilon)$$

$$\cdot [(x'x)^{-1} x'x\beta + (x'x)^{-1} x] \epsilon$$

$$= \beta + (xx^{-1}) x' \epsilon$$

$$\hat{\beta} - \beta = (x'x)^{-1} x' \epsilon$$

$$r(\hat{\beta}) = E[(\hat{\beta} - \beta)(\hat{\beta} - \beta')]$$

$$= E[(x'x)^{-1} x' \epsilon] \quad (x'x)^{-1} x' \epsilon$$

$$\therefore E(\epsilon\epsilon)' = \sigma_\epsilon^2$$

$$V(\hat{\beta} = \sigma_\epsilon^2 (x'x)^{-1} (x'x) (x'x)^{-1}$$

$$\Rightarrow r(\hat{\beta}) < V(\epsilon)'$$

Estimation of $\sigma^2$

$$SS_{res} = \sum_{i=1}^{n} (y_i - \hat{y}) \Rightarrow \sum_{i=1}^{n} e_i^2 \Rightarrow e'e$$

Sub $e = y - x\hat{\beta}$

$$SS_{res} = (y - x\hat{\beta})' (y - x\hat{\beta})$$

$$= y'y = \hat{\beta}' y' y - y' x\hat{\beta} + \hat{\beta} x' x\beta$$

$$= y'y - 2\hat{\beta} x' y + \hat{\beta}' x' x\hat{\beta}$$

We know that,

$$((x'x)\hat{\beta}) = x'y$$

$$SS_{res} = y'y - 2\hat{\beta}' y + \hat{\beta}' x' y$$

$$= y'y - \hat{\beta}' x' y$$

the residual mean square is,

$$MS_{res} = \frac{SS_{res}}{n-p}$$

$$\Rightarrow \sigma^2 = MS_{res}.$$

# Assumption in multiple linear regression model.

Some assumption are needed in the model

$$Y = X\beta, \epsilon$$ for drawing the satisfied inference.

the following assumption are made

1) $E(\epsilon) = 0$

ii) $E(\epsilon\epsilon') = \sigma^2 I_n$

iii) $Rank(x) = K$

iv) $x$ is a non-Stochastic matrix

v) $\epsilon \sim N(0, \sigma^2 I_n)$

vi) $\lim_{n \to \infty} \left( \frac{x'x}{n} \right) = \Delta$

exists and it is a non-stochastic and non-singular matrix (with finite elements).

# Estimation of Parameters.

A general procedure for the estimation of regression co-efficient is to minimize.

$$\sum_{i=1}^{n} M(\varepsilon i) = \sum_{i=1}^{n} H(yi - xi_1 \beta_1 - xi_2 \beta_2 \cdots - xin \beta_{ik})$$

for a suitably chosen function $H$.

Some examples of choice of $M$ are

$$M(x) = |x|$$

$$M(x) = x^2$$

$$M(x) = |x|^P \text{ is general,}$$

We consider the principal of Least Square which is related to $M(x) = x^2$ and the method of maximum likelihood estimation for the estimation of Parameters.

Using Regression Models for forecasting :-

• Forcasting and estimation of Casual effects are quite different objectives.

For forecasting,

• $R^2$ matters a lot.

• Omitted Variable bias isn't a problem.

* We will not worry about interpreting Co-fficients in forecasting models.

* External Validity is Paramount : the model estimated using historical data must hold into the near future.

* External validity is paramount : the model estimated using historical data must be hold into the (near) future.

## Generalized Least Squares :-

Until now we assumed that $\Sigma = \sigma^2 I$ but it can happen that the errors have non-constant variance or are correlated. Suppose instead that $\Sigma = \sigma^2 \underset{=}{\Sigma}$ where $\sigma^2$ is known but $\Sigma$ is known - in other words we know that the correlation and relative variance between the errors but we don't know the absolute scale

$$(Y - X\beta)^T \Sigma^{-1} (Y - X\beta)$$

which is solved by

$$\hat{\beta} = (X^T \Sigma^{-1} X)^{-1} X^T \Sigma^{-1} Y$$

Since we can write $\Sigma = SS^T$, where S is triangular matrix using the choleski Decomposition, we have

$$(Y - X\beta)^T S^{-T} S^{-1} (Y - X\beta) = (S^{-1} Y - S^{-1} X\beta)^T$$

So GLS is like regressing $S^{-1} x$ on $s^{-1} y$.

Furthermore

$$y = X\beta + \epsilon$$
$$S^{-1} y = S^{-1} X\beta + S^{-1} \epsilon$$
$$y' = x'\beta + \epsilon'$$

. 

So we have a new regression equation $y' = x'\beta + \epsilon'$ where if we examine the variance of the new errors $\epsilon'$ we find

$$\text{Var } \varepsilon' = \text{var}(S^{-1}\varepsilon) = S^{-1}(\text{Var } \varepsilon)S^{-T}$$
$$= S^{-1}\sigma^2 S S^T S^{-T} = \sigma^2 I$$

So the new variable $y'$ and $x'$ are related by a regression equation which has uncorrelated errors with equal variance. of course, the practical problem is that $\varepsilon$ may not be known.

We find that
$$\text{Var } \hat{\beta} = (x^T \Sigma^{-1} x)^{-1} \sigma^2$$

Weighted Least Squares:—

Sometimes the errors are uncorrelated, but have unequal variance where the from the inequality is known. Weighted least squares (WLS) can be used in this situation. When $\Sigma$ is diagonal, the errors are uncorrelated but do not neccessarily we have equal variance. we can write $\Sigma = \text{diag}$
$(1/w_1, \cdots, 1/w_n)$, where the $w_i$ are the weights
So $S = \text{diag}(\sqrt{1/w_1}, \cdots, \sqrt{1/w_n})$. So we can regress
$\sqrt{w_i} x_i$ on $\sqrt{w_i} y_i$ although the column of ones in the $x$-matrix needs to be replaced with $\sqrt{w_i}$. cases with low variability should get a high weight, high variability a low weight. Some example:

1. Errors proportional to a predictor. $\text{var}(\varepsilon_i) \propto x_i$
   Suggests $w_i = \bar{x}_i^{-1}$

2. $y_i$ are the average of $n_i$ observations then $y_i = \text{Var } \varepsilon_i$
   $= \sigma^2/n_i$ Suggests $w_i = n_i$.

The experiment was designed to test certain theories about the nature of a the strong interaction. The cross-section variable is believed to be linearly related to the inverse of the energy. At each level of the momentum, a very large number of observations were taken to that it was possible to accurately estimate the standard of the response (s.d).

## Robust Regression:-

In robust statistics robust regression is a form of regression analysis designed to overcome some limitations of traditional parametric and nonparametric methods. Regression analysis seeks to find the relationship between one or more independent variables and a dependent variable. Certain widely used methods of regression, such as ordinary least squares, have favourable properties in their underlying assumptions are true, but can give misleading results if those assumptions are not true, thus ordinary least squares is said to assumptions by the underlying data-generating process.

In particular least square estimates for regression models are highly sensitive to outliers. while there is no precise definition of an outlier, outliers are observations that do not follow the pattern of other observations. This is not normally a pattern if the outlier is simply an extreme observation drawn from the tail of a normal distribution, but if the outlier results from non-normal measurement error or some other violation standard

ordinary least squares assumptions, then it compromises (4) the validity of the regression result if a non-robust regual regression technique is used.

## Relationship between Analysis of Variance and regression:-

### (or)

### Difference between Regression and anova :-

Both the regression and anova are the that statistical methods which are used in order to predict the continuous outcome but in case of the regression, continuous outcome is predicted on basis of the one or more than one continuous Predictor Variable whearreas in case of anova continuous outcome is predicted on basis of the one or more than one categorical predictor Variables.

Regression is a statistical method to establish the relationship between Sets of Variables in order to make predictions of the dependent Variable with the help of independent variables. Anova On the other hand, is a statistical tool applied to uncorrelated group to find out whether they have a Common mean.

What is Regression Analysis?

Regression is a very effective statistical method to establish the relationship between set of variables. The variables for which the regression analysis is done are the dependent variable and one or more independent variable. It is a method to understand the effect on a dependent variable of one or more than one independent variable.

(i) Suppose, the example; a paint company uses one of the derivatives of crude solvent & monomers as its raw meterial. We can run a regression analysis between the price of that raw meterial and the price of Brent Crude Prices.

(ii) In this example, the price of the raw metrial is the dependent variable and the price of Brent Prices is the independent variable.

(iii) As the price of solvents and monomers increases and decreases in price with the rise and fall of Brent Prices, the price of the raw meterial is the dependent variable.

(iv) Similarly for any business decision in order to validate a hypothesis that a particular action will lead to the increases in the profitability of a division can be validated based on the result of the regression between the dependent and independent variable.

## What is ANOVA?

ANOVA is the short from of analysis of variance. ANOVA is a statistical tool that is generally used on random variables. It involves group not directly related to each other in order to find out whether there exist any common means.

(i) A simple example to understand this point to run ANOVA for the series of marks of student from different college in order in to try to find out whether one student from one school is better than the other.

(ii) Another example can be if two separate research team is researching different product not related to each other. ANOVA will help to find which one is providing better results. the three popular techniques of ANOVA are a random effect, fixed effect, and mixed effect.

# Generalized Linear Model :-

The generalized linear model (GLm) is a flexible generalization of ordinary linear regression that allows for response variables that have error distribution models other a normal distribution.

The GLm generalizes linear regression by allowing the linear model to be related to the response variable via a link function and by allowing the magnitude of the variance of each measurement to be a function of its predicted value.

Generalized Linear models were formulated by John Nelder and Robert Weddertum as a way of unifying various other statistical models, including Linear regression, logistic regression and poisson regression. The proposed an alternatively reweighted least squares method for maximum likelihood estimation esti of the model parameters. Maximum-likelihood estimation remains popular and is the default method of many statistical computing Package, other approaches including Bayesian approaches and least square fits to variance stabilized responses have been developed.

# Comparison to generalized Linear model:-

The general linear mode (GLm)[2,3] and generalized Linear Model[4,5] are two commonly used families of statistical methods to relate some number of continuous and/or categorical predictor to single outcome variable.

The main difference between the two approaches is that the GLm strictly assumes that the residuals will follow a conditionally normal distribution[3] with the GLm loosens this assumption and allows for a variety of the other distributions from the exponential family for the residuals. Of note, the GLm is a special case of the GLm in which the distribution of the residuals follow a conditionally normal distribution.

# Non-Linear Regression Model :-

NonLinear regression is a form of regression analysis in which observations data are modeled by a function which is a non-Linear common combination of the model parameters and depends on one or more independent variables. The data are fitted by a method of successive approximation.

# MULTICOLINEARITY - Definition :-

Multicolinearity generally occurs when there are high correlations between two or more predictor variables. In other words, one predictor variable can be used to predictor the other. This creates redundant information skewing the results in a regression model.

For example :-

A person's hight and weight, age and sales price of a car, or year of education and annual income.

An easy way to detect multicolinearity to calculate correlation coefficient for all pairs of predictor variables. If the correlation coefficient, $r$ is exactly +1 or -1, this is called predictor perfect multicolinearity. If $r$ is close to or exactly -1 or +1. one of the variables should be removed from the model if all possible

It's more common for multicolinearity to rear its ugly head in observational studies, its less common with experimental data.

When the condition is present, it can result (10)
in unstable and unreliable regression estimates.
Several other problems can interfere with
analysis of results, including

(i) The t-statistic will generally be very
small and coefficient confidence interval will be
very wide. The means that it is harder
to reject the null hypothesis.

(ii) The Partial regression coefficient
may be an imprecise estimate, standard
error may be very wide. This means
standard errors may very large.

(iii) Partial regression coefficient may
have sigh and/or magnitude change as
they pass from sample to sample.

(iv) Multicolinearity makes it difficult
to gauge the effect of independent
variables on dependent variables.

Comparing Regression Models to See the Effect
of Multicolinearity:-

The predictors effectively removed
the multicolinearity, we could run the same
model twice, once with severe multicolinearity
and once with moderate multicolinearity.

This provides a great head-to-head comparison and it reveals the classic effects of multicolinearity.

The standard error of the coefficient indicate the precision of the coefficient estimates. Smaller values represent more reliable estimate. In the second model, you can see that the SE coefficient is smaller for both % Fat and weight.

Also, % Fat is significant this time, while it was insignificant in the model with severe multicolinearity. Also its sign has switched from + 0.005 to - 0.005! The % Fat estimate in both models is about the same absolute distance from zero, but it is only significant in the second model because the estimate is more precise

Compare the summary of model statistics between the two models and you'll notice that, S, R & R- squared, adjusted R-required, and the other are all identical. Multicolinearity dosen't effect how well the model fits. In fact, if you want to use the model to make predictions, both models produce indentical results for fitted values and prediction interval

# Dealing with multicolinearity :-

Multicolinearity occurs when independent Variables in a regression model are correlated. This correlation is a problem because independent Variables should be independent. If the degree of correlation between Variable is high enough, it Can cause problems when you fit the model and interpret the result.

In Some Causes, multicolinearity isn't necessarily a problem, and I'll Show you how to make this determination. I'll work through an example datest which contains multicolinearity to big bring it all to life.

# Multicolinearity based on the potential Problem : -

A key goal of regression analysis is to isolate the relationship between each independent Variable and the dependent Variables The interpretation of a regression coefficient is that it represents the mean change in the dependent Variable for each 1 unit change in an independent Variable When you hold all of the other independent Variable Constant. That last Portion is Crucial for our discussion about multicolinearity.

The idea is that you can change the variable of one independent and not the others. However, when independent variables are correlated. it indicates that changes in one variable are associated with shifts in another variable.

The stronger the correlation, the more difficult it is to change one variable without changing another. It becomes difficult for the model to estimate the relationship between each independent variable and the dependent variable independent because the independent variable tend to change in ~~usi~~ unison.

There are two basic kinds of multicolinearity:-

i) Structural multicolinearity:-
This type occurs when we create a model term using other terms. In other words, it's a byproduct of the model that we specify rather the being present in the data itself.

For example:-
If you square term x to model curvature, clearly there is a correlation between x and $x^2$.

# data multicolinearity :-

This type of multicolinearity is present in the data itself rather than being an artifact of our model. Observational experiments are more likely to exhibit this kind of multicolinearity.

Multicolinearity causes the following two basic types of problems :-

(1) The coefficient estimate can swing based on which other independent variables are in the model. The coefficients become very sensitive to small changes in the model.

(2) Multicolinearity reduces the precision of the estimate coefficients which weakens the statistical power of your regression model. You might not be able to m trust the p-values to identify independent variables that are statistically significant.

## fix multicolinearity :-

The good news is that it is not always mandatory to fix the multicolinearity. It all depends on the primary goal of the regression model.

The degree of multicolinearity greatly impact the p-values and coefficient but not predictions

and goodness-of-fit test. If your goal is to perform the predictions and not neccessary to understand the significance of the independent variable, it is not a mandata to fix the multicolinearity. issue.

How to test multicolinearity?

(1) correlation matrix / correlation plot.

(2) Variation Inflation factor (VIF)

A correlation plot can be used to identify the correlation (or bivariate relationship between two independent variables VIF is used to identify the correlation of one independent variable with a group of other variables.

Hence, it is preferred to use VIF for better understanding.

$$VIF = 1 \longrightarrow No\ correlation$$
$$VIF = 1\ to\ 5 \longrightarrow Moderate\ correlation$$
$$VIF > 10 \longrightarrow High\ correlation.$$

Testing for Multicolinearity with Variance Inflation Factors (VIF)

If you can identify which variables are affected by multicolinearity and the Strength of the correlation, you're well on your way to determining whether you need to fix it. Fortunately, there is a very simple test to assess multicolinearity in your regression model.

The variance Inflation factor (VIF) identifies correlation Independent Variables and the Strength of that correlation.

Statistical software calculates VIF for each Independent Variable. VIF start at 1 and have no upper limit. A value of 1 Indicates that there is no correlation between this Independent Variable and any other. VIFS between 1 and 5 suggest that there is a moderate correlation, but it is not severe enough to warrant corrective measures VIFS greater than 5 represent critical levels of mut multicolinearity where the coefficients are poorly estimated, and the P-Values are questionable

Use VIFs to identify correlation between Variables and determine the strength of the relationships. most statistical software can display VIFS for you. Assessing VIFS is Particularly Important for observational studies because these Studies are more prone to having multicolinearity

## Models of Time Series

Time Series is Statistical data that we arrange and present in a chronological order spreading over a period of time. Time series analysis is a statistical technique dealing with time series data. According to Spiegel, "A time series is a set of observations taken at specified times, usually at equal intervals". In Statistics, for time series analysis two main categories of models are popular. Let us discuss the models of time series Analysis in details.

In time series quantitative data are arranged in the order of their occurrence and resulting statistical series. The quantitative values are usually recorded over equal time interval such as daily, weekly, monthly, quarterly, half-yearly, yearly or any other measure of time.

Some example are statistics of industrial production in india on a mothly basis, bith-rate figures annually, the yield on ordinary shares, and weekly wholesale Price of rice, etc.

# Components of Time Series :-

There is a different kind of forces which influence the time series analysis. Some are continuously effective while others make themselves felt at recurring time and intervals. So our first task is to divide the data and elements into components.

A time series consists of the following four components or basic elements :

1. Basic or secular or long-time trend.
2. Seasonal Variations.
3. Business cycles or cyclical movements.
4. Erratic or Irregular fluctuations.

## Models of Time Series Analysis :-

The following are the two models which we generally use for the decomposition of time series into its four components. The objectives is to estimate and seperate the four types of variations and to bring out the relative effect of each on the overall behaviour of the time series.

(1) Additive Model (2) Multiplicative model.

## (1) Additive Model :—

In the additive Model, we represent a particular observation in a time series as the sum of these four components.

ie., $O = T + S + C + I$

Where O represents the original data. T represents the trend. S represents the seasonal variations. C represents the cyclical variations and I represents the irregular variations. We can write $y(t) = T(t) + S(t) + C(t) + I(t)$

## (2) Multiplicative Model :—

In this model, four components have a multiplicative relationship. So, we represent a particular observation in a time series as the product of these four components.

ie $O = T \times S \times C \times I$

Where O, T, S, C and I represent the terms as in additive model.

In another way, we can write

$$y(t) = T(t) \times S(t) \times C(t) \times I(t)$$

This model is the most used model in the decomposition of time series. To remove any doubt between the two models, it should be made clear that in multiplicative model

S, C, and I are indices expressed as decimal
percentages whereas, in additive model s, c and I
are quantitative deviations about a trend that can
be expressed as seasonal, cyclical and irregular
in nature.

Example :-

If a multiplicative Model.

$$T = 500, S = 1.4, C = 1.20 \text{ and } I = 0.7$$

then $O = T \times S \times C \times I$

By substituting the values we get.

$$O = 500 \times 1.4 \times 1.20 \times 0.7 = 588$$

$$O = 588$$

If in a additive Model :-

$$T = 500, S = 100, C = 25, I = -60$$

then $O = 500 + 100 + 25 - 60 = 565$

Growth Curve —Definition :-

Let X be a p×n random matrix
corresponding to the observation. A a p×q within
design matrix with $q \leq p$, B a q×k parameter
matrix, C a k×n between individual design
matrix with rank $(c) + p \leq n$ and Let $\varepsilon$ be a
positive definite p×p matrix, Then

$$X = ABC + \Sigma^{1/2} E$$

definite the growth Curve model. Where A and c

are known. B and $\Sigma$ are unknown and E is a
random matrix distributed as $N_{p,n}(0, I_{p,n})$.
This differs from standard MANOVA by the
addition of G, a postmatrix $x''$ [3].

## Gompertz Curve :-

The gompertz Curve or gompertz function
is a type of Mathematical model for a time series
named after Benjamin Gompertz. It is a sigmoid
function which describes growth as being slowest at the
start and end of a given time period. The right
hand or future Value asymptote of the function is
approached much more ~~great~~ generally by the curve than the
left-hand (or) lower values asymptote. This is in contrast
to the simple logistic function in which both asymptotes are
approached by the curve Symmentrically It is a special
Case of the generalised logistic function. The function
was originally designed to describe human mortality.
but since has been modified to be applied in biology.
with regard to detailing Populations.

## Growth Curve :-

A growth curve is a graphical representation
of how a particular quantity increases over time.
Growth Curve are used in statistics to determine the
type of growth Pattern of the quantity-be it linear,
exponential, or cubic. once the type of growth is
determined, a business can create a mathematical model
to predict future Sales.

## Exponential Curve :-

### Definition :-

An Exponential function or Curve is a function that grows exponentially, or grows at an increasingly larger rate as you pick larger values of x, and usually takes the form

$Y = a^x$, where a is any real number.

$2^x, 3^x, 10^x$ etc are all Examples of Exponential function.

## Logistic Curve :-

A logistic function (or) logistic curve is a common s-shaped curve with equation.

$$f(z) = \frac{L}{L+e^{-k(x-x_0)^s}}$$

where $x_0$ the x value of the sigmoid's midpoint.
L the curve's maximum value.
k the logistic growth rate or steepness

as the curve

For values of x in the domain of real numbers from $-\infty$ to $+\infty$. the s-curve shown on the right is obtained, with the graph of f approaching L as a x approaches $+\infty$ and approaching zero as x approaches $-\infty$.

The logistic function finds applications in a range of fields, including biology, biomathematics, chemistry, demography, economics

Economics, geoscience, mathematical Psychology, Sociology, Political Science, linguistic, statistics and artificial natural neuron network.

A generalization of the logistic function is the hyperbolastic function of type-I.

## AutoRegression :-

A regression model, such as Linear regression, models an output value based on linear combination of input values.

For example:-

$$yhat = b0 + b1 * x1$$

where yhat is the prediction. bo and b1 are coefficient

Reg