

(Mid term) Sta 404 (Regression & Correlation Analysis) (Important Notes)

Lecture #01.

* Things to remember:-

* Regression Analysis is a statistical tool for investigation then relationship between a dependent variable & one or more independent variable.

* MCQs; Regression analysis is widely used for prediction & forecasting.

* Applications:-

- * Economics
- * Management
- * Engineering
- * Social Sciences
- * Medical fields
- * Sports & physical Sciences
- * Computer Sciences.

Relationship between variables.

* Exact relationship

* Fahrenheit ~ Celsius

$$F = 32 + \frac{9}{5}C$$

$$a \propto a^2$$

* Deterministic Model

$$Y = a + bX$$

Non-exact relationship

Weight ~ Height.

Probabilistic Model.

$$Y = a + bX + e \rightarrow \text{error}$$

Dependent vs Independent variables.

The dependent variable is the variable in regression that can not be controlled or Manipulated.

The Independent variables is the variable in regression that can be controlled or Manipulated.

Scatter plots and Relationship.

① Scatter plot is used to show the relationship between two variables.

② The Independent and dependent variables can plotted on a graph called a scatter plot.

③ The independent variables 'x' is plotted on the horizontal axis, and dependent variables on the vertical axis.

Lecture #02.

Regression line.

$$\hat{y} = a + bx \Rightarrow \text{equation of regression line.}$$

for finding 'a'

$$a = y - b(x) \quad ; \quad x = ? \quad , \quad y = ?$$

for finding 'b'

$$b = \frac{n(\sum xy) - (\sum x)(\sum y)}{n(\sum x^2) - (\sum x)^2}$$

$$n = ? \quad , \quad \sum xy = ? \quad , \quad \sum x = ? \quad \sum y = ? \quad , \quad \sum x^2 = ?$$

Lecture #03

Important points:-

- ① Simple linear regression is really a comparison of two methods;
 - * Where the independent variable does not even exist.
 - * The other uses the best fit regression line.
- ② If there is only one variable the best prediction for other values is the mean of the "dependent" variable.

③ The difference between the best fit-line and the observed value is called the residual. (or error)

④ The residuals are squared and then added together to generate sum of square residuals/error. SSE.

⑤ Simple linear regression is designed to find the best fitting line through the data that minimize the SSE.

⑥ By Determining the coefficient of determination we can find the estimated regression equation fit our data.

* Formula for coefficient of Determination = $R^2 = \frac{SSR}{SST}$ (COD)

↑ Residual.

* $SST = SSR + SSE$

Sum of \downarrow Sq. of Total → Sum of Sq. Error

"The Simple linear Regression Model"

Steps involves;

- ① The Mean of the responses is a linear function.
- ② The errors are independent.
- ③ The errors are normally distributed.
- ④ The errors have equal variances for all X values.

Properties of the least square Regression line.

- ① The least square regression line always passes through the centre of the data.
- ② The sum of deviations of the observed values of y from the LS regression line is always equal to zero. $\sum e_i = (\sum (y_i - \hat{y}_i)) = 0$
- ③ The sum of the squares deviations of the observed values of y from the LS regression line is Minimum. $\sum e_i^2 = \sum (y_i - \hat{y}_i)^2 \Rightarrow \text{Minimum.}$

(4) The line fitted by LS Method is "best" fit line.

Lecture # 04.

* Correlation:-

The word Correlation is made of Co- (together) and Relation.
"Two variables are said to be correlated if they tend to simultaneously vary together in some direction."

* Types of Correlation:-

* Linear

↓
When graph is plotted like straight line then it is linear correlation.

* Non-linear.

↓
When graph is not straight line.

Correlation might be

* Simple

Involve two and ~~one~~ one dependent variable.

* Multiple:-

one dependent and more

than one independent variables.

* Partial:-

one dependent variable and more than one independent variable but only one independent variable is considered and other independent variables are considered constant.

Correlation might be:-

* **Positive:-** A positive correlation is a relationship between two variables where if one variable increases, the other one also increases. $\nearrow \nearrow$

* **Negative:-** A Negative correlation means that there is an inverse relationship between two variables. When one variable decreases, the other increases. $\searrow \searrow$

* No Correlation:-

Their movement is independent from each other.

zero $\uparrow \longrightarrow$

Methods of Studying Correlation.

- ★ Scatter plot Method.
- ★ Karl Pearson's Coefficient.
- ★ Spearman's Rank Correlation Coefficient.

Correlation Coefficient - Important points.

- ① The range of the Correlation Coefficient is from -1 to $+1$.
- ② If there is a strong positive linear relationship between the variables, the value of ' r ' will be close to 1 .
- ③ If there is a strong Negative linear relationship between the variables, the value of ' r ' will be close to -1 .
- ④ When there is no linear relationship between the variables, the value of r will be close to 0 .

* How to Calculate Correlation Coefficient - r .

* Formula :-

$$r = \frac{n(\sum xy) - (\sum x)(\sum y)}{\sqrt{[n(\sum x^2) - (\sum x)^2][n(\sum y^2) - (\sum y)^2]}}$$

* There are no units associated with r .

* The value of r will remain unchanged. If the x and y values are switched.

* "Correlation is Not Causation" — which says that a correlation does not mean that one thing causes the other. (There could be other reasons the data has a good correlation).

Lecture #05.

Rank Correlation.

Developed By:-

Rank Correlation Method

was developed by the British Psychologist Charles Edward Spearman in 1904.

Rank Correlation Coefficient formula.

$$r_s = 1 - \frac{6\sum d^2}{n(n^2-1)}$$

where,

r_s = Coefficient of rank correlation.

d = difference.

n = no. of data

It deals with measuring Correspondence between two rankings and assessing the significance of this Correspondence.

Features.

* The range of rank Correlation is -1 to $+1$.

The rank Correlation can be interpreted in the same way as Karl Pearson's Correlation Coefficient.

Rank Correlation Coefficient is a distribution-free measure since no strict assumption is made about the population from which it is drawn.

Spearman's formula is the only formula available to find the Correlation between rank.

Types of Rank Methods.

Ranked data :- where ranks are given.

* Not Ranked data :- where ranks are not given.

Ranks are tied :- where repeated ranks occurs.

Lecture # 06.

Multiple Linear Regression Analysis.

Difference Between Simple and Multiple Regression.

Simple Linear Regression

In simple linear regression, the regression equation contains one independent variable x and one dependent variable y and it is written as;

$$\hat{Y} = a + bX$$

Multiple linear Regression.

In multiple linear regression equation with two independent variables (x_1 and x_2) and one dependent variable has the form;

$$\hat{Y} = a + b_1 X_1 + b_2 X_2$$

The general form of the multiple regression equation with k -independent variable is:-

$$\hat{Y} = a + b_1 X_1 + b_2 X_2 + \dots + b_k X_k$$

- * Y is called Dependent variable.
- * The X's are the independent variables.
- * The b's are called partial regression coefficients.
- * 'k' is number of independent variables.

In Multiple Regression;
the variables X_1 & X_2 are
Not correlated with each other.

Multiple Regression Analysis.

Normal Equation:-

$$\hat{Y} = a + b_1 X_1 + b_2 X_2$$

- * $\sum Y = na + b_1 \sum X_1 + b_2 \sum X_2$ ----- (1)
- * $\sum X_1 Y = a \sum X_1 + b_1 \sum X_1^2 + b_2 \sum X_1 X_2$ ----- (2)
- * $\sum X_2 Y = a \sum X_2 + b_1 \sum X_1 X_2 + b_2 \sum X_2^2$ ----- (3)

Method to solve Problems By using Normal Equations:-

- ① Take Multiple Model
- ② Make the 3 Normal equations.
- ③ Calculate the required values in the table.
- ④ Reduce the 3 Normal equations.
- ⑤ Reduce the 3 normal equation to the 2 equations
- ⑥ solve the New two equations to get unknown constants, step by step.
- ⑦ Put the constants in the Model.

* Method to find coefficient values by using formula.

$$\hat{Y} = a + b_1 X_1 + b_2 X_2$$

Area ↙ ↘ No. of data

* for calculating value of 'a'

$$a = \bar{Y} - b_1 \bar{X}_1 - b_2 \bar{X}_2$$

* For Calculating value of b_1 :

$$b_1 = \frac{(\sum x_1 y)(\sum x_2^2) - (\sum x_2 y)(\sum x_1 x_2)}{(\sum x_1^2)(\sum x_2^2) - (\sum x_1 x_2)^2}$$

* For Calculating value of b_2 :

$$b_2 = \frac{(\sum x_2 y)(\sum x_1^2) - (\sum x_1 y)(\sum x_1 x_2)}{(\sum x_1^2)(\sum x_2^2) - (\sum x_1 x_2)^2}$$

For Formula Method
First you have to convert
data into deviation form.

original data			Data in Deviation form.		
Y	X_1	X_2	$y = Y - \bar{Y}$	$x_1 = X_1 - \bar{X}_1$	$x_2 = X_2 - \bar{X}_2$
30	10	15	12.2	4	4.6
22	5	8	4.2	-1	-2.4
16	10	12	-1.8	4	1.6
14	3	7	-10.8	-3	-3.4
7	2	20	-3.8	-4	-0.4
14					
T	30	52	0	0	0
M	17.8	10.4			