

Unit 2

11
IA

Part A :-

Simple Linear Regression & Correlation

1. Introduction to Linear Regression
2. The Simple Linear Regression Model
3. Least Squares
4. Fitted Model
5. Properties of Least Squares estimators
6. Inferences Concerning the Regression Coefficients
7. Prediction
8. Simple Linear Regression Case Study.

Unit-2 Part-B

Correlation and Regression:-

Introduction:- We studied in part-A, the concept of a single R.V.

i.e. → We got an idea about the characteristics of R.V by means of various statistical averages or moments like

Mean

Variance

S.D

Skewness

Peakedness - - -

But in many real world problems @ many practical problems:- Several Random Variables interact with each other.

Ex:- Doctors may measure many parameters like height, weight, blood pressure, sugar level etc.

Properties :- 1. $0 \leq P(x_i, y_j) \leq 1$

2. $\sum_{i,j} P(x_i, y_j) = 1$

3. Joint Cumulative distribution

f_{xy}

$$F_{xy}(x, y) = P(X \leq x, Y \leq y)$$

$$\text{Ex:- } F_{xy}(2, 3) = P(X \leq 2, Y \leq 3)$$

X
Y
(1, 1) pairs
(1, 2)
(1, 3)

$$= P(X=1, Y=1) + P(X=1, Y=2) + P(X=1, Y=3)$$

+

$$P(X=2, Y=1) + P(X=2, Y=2) + P(X=2, Y=3)$$



What is Multiple R.V. -

Say X & Y are R.V \Rightarrow 2D R.V
 $\Rightarrow (X, Y)$

The outcome of trials $X = x, Y = y$

i.e. $(X, Y) = (x, y)$.

n-dimensional R.V \Rightarrow n-dimensional
Random Vector
(or)

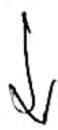
Vector of Random Variables

i.e. n-dimensional Euclidean
Space.

$$\Rightarrow \mathbb{R}^n$$

Say $\mathbb{R}^2 \Rightarrow$ Random Vector (X, Y) i.e. B.Variate

(X, Y) has discrete \Rightarrow It has only finite number



of pairs

PMF

Then Joint probability $f_{xy} \Rightarrow f_{(x,y)} = P(X=x, Y=y)$

Joint probability distributions \Rightarrow

$P(X,Y) @ X/Y$	1/2	1/2
Y\X	1/2	1/2
0	1/2	1/2

Properties

1. $0 \leq f_{xy}(x, y) \leq 1$ $\Leftrightarrow 0 \leq f_{xy} \leq 1$

2. $\int_{-\infty}^{\infty} f_{xy} = 1$ $\Leftrightarrow f_{xy}(0, 0) = 1$

3. f_{xy} is non decreasing

Marginal probability distribution functions
Components of distribution functions are
called P.P. P.d.

$$\boxed{f_{xy}(x, y)}$$
$$f_{xy}(x, y)$$

Joint Probability density function $f_{XY}(x, y)$

$$\int_{-\infty}^{+\infty} \int_{-\infty}^{+\infty} f_{XY}(x, y) dx dy = 1$$

Properties 1, $f_{XY} \geq 0$

$$2 - \int_{-\infty}^{+\infty} \int_{-\infty}^{+\infty} f_{XY}(x, y) dx dy = 1$$

$$3 - \int_{-\infty}^{+\infty} \int_{-\infty}^{+\infty} f_{XY}(x, y) dx dy = 1$$

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

Empirical
Model

where ϵ = Random error term.

$$\textcircled{2} \quad E(\epsilon) = 0$$

$$V(\epsilon) = \sigma^2$$

$$\textcircled{3} \quad \text{prove } E(Y) = 0$$

in Regression?

w.l.o.g. Regress (SLR)

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

$$E(Y) = E(\beta_0 + \beta_1 x + \epsilon)$$

$$= E(\beta_0) + E(\beta_1 x) + E(\epsilon)$$

$$= \beta_0 + \beta_1 x + 0$$

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

$$\textcircled{4} \quad \text{Prove } V(Y) = \sigma^2$$

Regression?

$$\text{w.l.o.g. SLR} \Rightarrow Y = \beta_0 + \beta_1 x + \epsilon$$

Apply Expectation

$$V(Y) = V(\beta_0 + \beta_1 x + \epsilon)$$

$$= V(\beta_0) + V(\beta_1 x) + V(\epsilon)$$

$$= \beta_0 + \beta_1 x + \sigma^2$$

$$= 0 + \sigma^2$$

$V(Y) = \sigma^2$

①

Regression Analysis ② Estimation)

Introduction:-

→ Many engineering and scientific problems are concerned with determining a relationship between a set of variables. (especially in Regression) OR

✓ Many problems in engineering and science involve exploring the relationships between two or more variables. //

Regression Analysis Regression Analysis is a statistical

technique that is very useful for these types of problems

→ The Concept of Regression Analysis was 1st

introduced by Sir Francis Galton in 19th

Century, according his words

"STEPPING BACK TOWARDS AVERAGE"

"Regression means"

Now a day's Regression mea.

"It is a mathematical measure of a given relationship between of 2 or more Variables Under Study."

Meaning In Regression Analysis, we can estimate the Value of One Variable with the Value of other Variable, which is known.

Definition The Statistical Method which helps us to estimate the Unknown Value of One Variable from the Known Value of other Variable is called Regression.

Regression

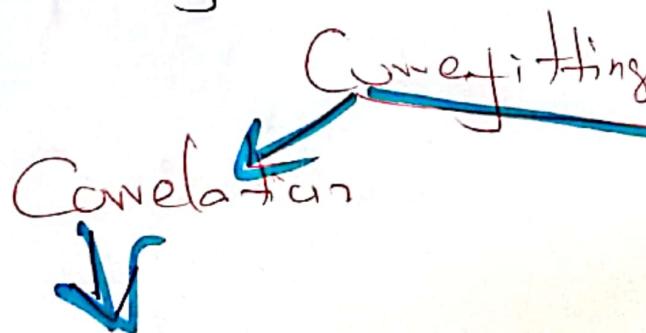
Unit-2
part-A

History & Introduction

→ Fitting of Curves to a set of numerical data is of considerable importance ~~to science~~
 theoretical as well as practical.

Theoretically:-

It is useful in study of



To measure the strength
 the linear relationship b/w
 two R.V say X & Y.

To describe the
 relationships b/w
 the mean of \bar{Y}
 random Variable Y
 and One or more other variables.

Houses in the same part of the country that are the same square footage of living space will not all be sold for the same price.

Dependent Variables
④ Y
Responses

Tar Content
Gas mileage (mpg)
price of houses

↓
Single Value (Y)

Expected Response.

Independent Variables

④ X
Regressors

Inlet temperature

engine volume

Square feet of living space

Natural Independent Variables.

Generally, these variables

are Input Variables

④ Input Values.

Set of Input Values (x_1, x_2, \dots)

Example — In Industry, it may be known that "Tar Content" ^{in the} outlet stream in a chemical process is related ⁱⁿ to the inlet temperature.

It may be of interest to develop a method of prediction i.e. A procedure for estimating the Tar Content for various levels of the inlet temperature from experimental information.

Now, of course, it is highly likely that for many example runs in which the inlet temperature is the same, say 130°C . The outlet Tar Content will not be the same engine volume.

They will not all have the same gas mileage.

If the relationship is exact, then it is a deterministic relation between two scientific variables.

i.e. here there is no random or probabilistic component to it.

But in general, the relationship is not deterministic (i.e. \square does not always give the same value for \square).

But, in real world problems are probabilistic in nature. So, we cannot expect exact solution.

So, the concept of Regression Analysis deals with finding the best relationship b/w y & x .

But sometimes multiple regressions also

will be there. Then

$$Y = a + b_1 x_1 + b_2 x_2 + b_3 x_3 + \dots$$

(a)

$$Y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 + \dots$$

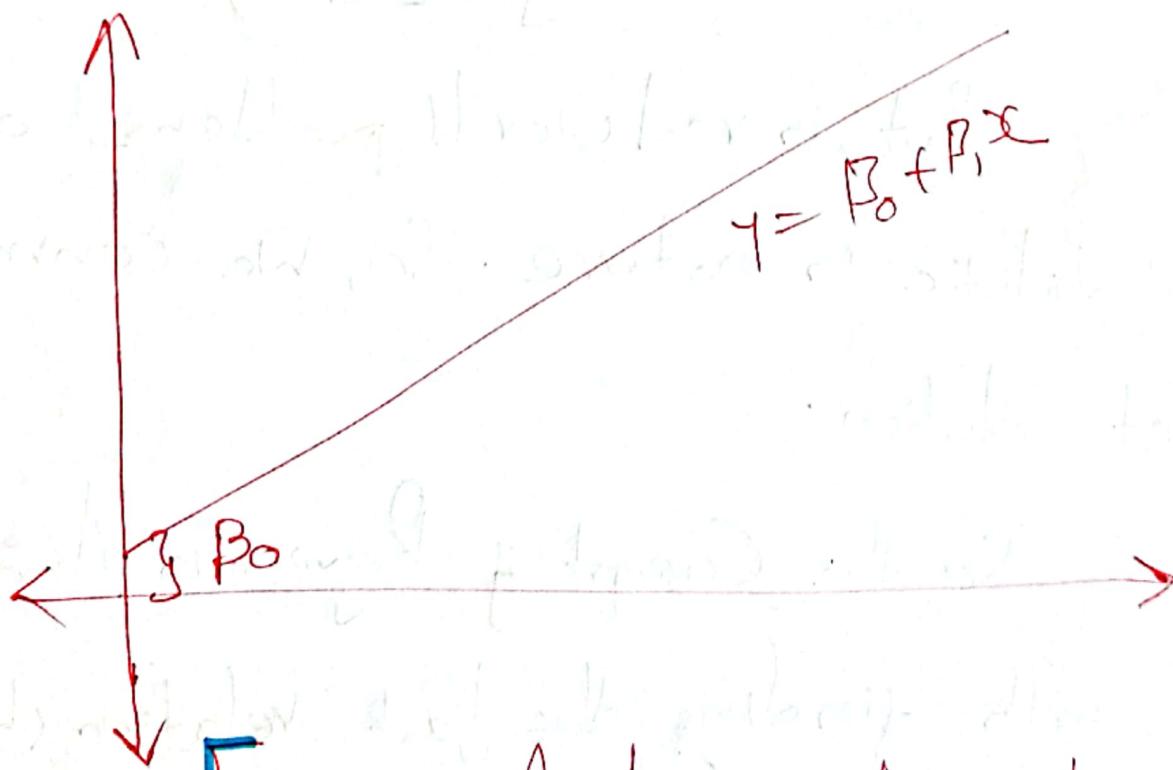


Figure: A Linear relationship β_0 ,
intercept, β_1 : slope.

Uses:-

1. It is used to estimate the relationship b/w economic Variables income & expenditure
2. It is widely used for prediction purpose
3. It is highly Valuable in economics & business
4. It is useful in statistical estimation of production function, Cost function, Constant function, etc

Simple Linear Probabilistic Model (SLP)

① Linear Regression Equation =
② SLR
Linear Regression Model

Consider the problem of trying to predict the value of a response variable y based on the value of an independent variable x . Then, the best fitting line

$$y = a + bx$$

where a = intercept
 b = slope.

$$\text{Example: } Y = \alpha + \beta_1 x_1 + \beta_2 x_2 \quad (2)$$

$$\text{where } Y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 \quad (3)$$

$$Y = \beta_0 + \beta_1 x_1 + \epsilon \quad (4) \quad Y = a + b x + \epsilon$$

where β_0, β_1 = Unknown intercept parameter

β_1 = Slope parameter.

ϵ = Random Variable

i.e. we need to find β_0, β_1 and

estimate (4) finding from the data. $E(\epsilon) = 0$

assumed
 $E(\epsilon) = 0$

$V(\epsilon) = \sigma^2$

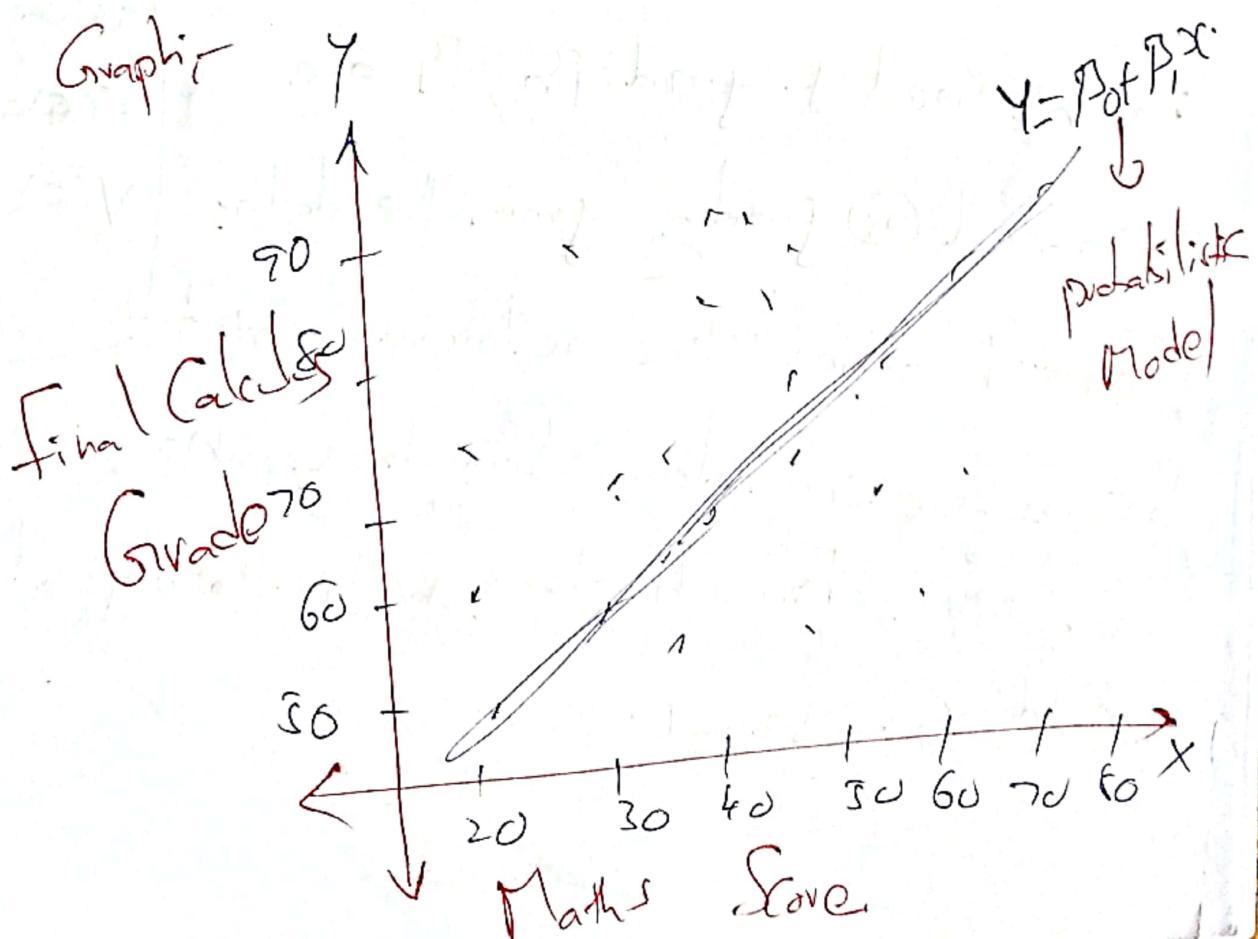
Ex: Mathematics achievement test

Scores for a random Sample of $n=5$ in
CSE AIML branch along with their final marks

In 2nd Year Sem-1

Now we need to see bivariate.

Student Roll No	Maths achieved	Final calculus Grade
1	39	65
2	43	78
3	21	52
4	64	82
5	57	92



What is Regression and briefly explain about Linear Regression or Line of Regression & their types?

→ Usually, the primary purpose of a regression study is prediction.

(estimation, finding the values, approximating - - -)

→ We want to develop an equation by which the value of \boxed{Y} can be predicted randomly well,

Based on knowledge of the values assumed by other variables in equation.

→ Suppose, if the variables in a bivariate distribution are related. We can find the points in Scatter diagram will

cluster around some curve called

"Curve of Regression"

II

x_1, x_2, x_3, \dots

only

Let $X =$ Mathematical Variable

Capital Y $y =$ Random Variable.

The Curve of regression of y on X is the graph of mean Value of y for Various Values of x .

$$Y = a + bx$$

Linear Regression

Line of Regression

II If the Curve is a straight line it is called Linear Regression. (or) Line of Regression

Equation

$$Y = a + bx$$

(or)

$$Y = a + b$$

→ First test is "attempting to derive up a regression equation to decide whether a linear regression curve appears to be appropriate or not."

Ensure that it should be appropriate!

Types of Linear Regression

(i) Linear Regression

There are two approaches to

above problem of determining, whether linear regression is appropriate or not.

① Graphic Method

① Scatter diagram

② Algebraic Method

② Method of Least Squares

① Graphic Method

Here points representing the points values of the variables are plotted on graph.

Graphic Method ① Scatter diagram @ Scattergram

Dependent Variable

Independent Variable

→ This is a picture of data in which the pairs (x, y) are plotted as points in a coordinate plane.

→ If points tend to form a linear trend, then we use a simple linear regression method to obtain a prediction equation.

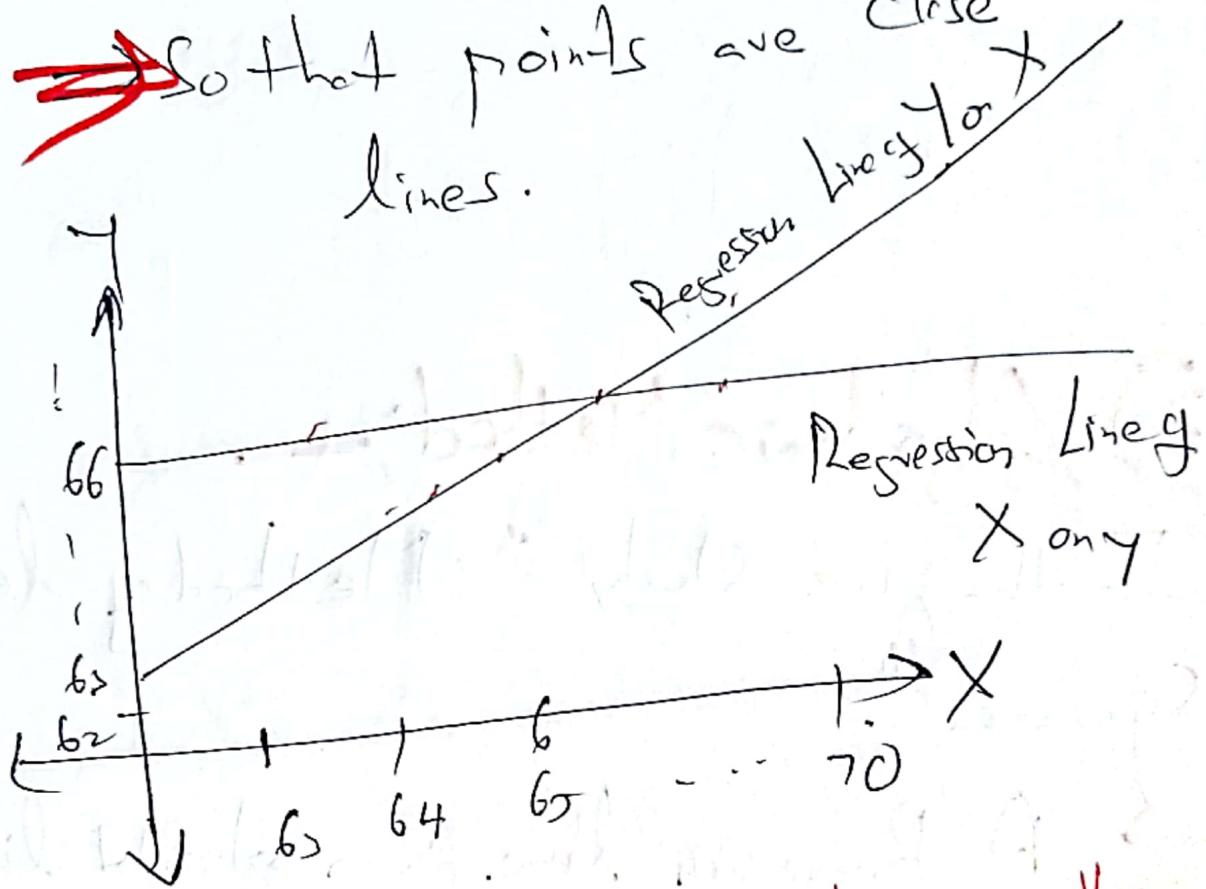
→ Then, our problem is to estimate the equation for straight line by means of observed data. ② Theoretical data

finally we will set Expected data ③ Desired data

i.e. Several straight lines can be drawn to include some of data points.

& also many lines can be drawn

~~So that points are close to the lines.~~



Example: Fit a Regression Line on the

Scatter diagram for the following data.

These above points form a scatter diagram. A Regression Lines are shown in above graph by free hand.

② Algebraic Method:-

In this, we study "Method of least squares".

→ A Regression line is a straight line fitted to the data by the method of "least squares".

→ It indicates the best possible mean value of one variable.

~~(2) Method of Least Squares~~

to the mean value of the other.

➤ There are always two regression lines constructed for the relationship

in two variables X & Y.

➤ Thus,

Thus (SLR) Regression Equations

or
Simple Linear Regression Equations

Y on X

X on Y

① Straight Line Equations

$$Y = a + bx \quad \text{or} \quad Y = ax + b$$

$$X = a + by \quad \text{or} \quad X = ay + b$$

② Second degree Parabola:-

$$Y = a + bx + cx^2 \quad \text{or} \quad Y = ax^2 + bx + c$$

$$x = a + by + cy^2 \quad \text{or} \quad x = ay^2 + by + c$$

③ Regression Coefficients $r = \sqrt{b_{yx} - b_{xy}}$

Y on X

X on Y

$$(Y - \bar{Y}) = b_{yx} (X - \bar{x})$$

$$(X - \bar{x}) = b_{xy} (Y - \bar{Y})$$

$$(Y - \bar{Y}) = \frac{r \sigma_y}{\sigma_x} (X - \bar{x})$$

$$(X - \bar{x}) = \frac{r \sigma_x}{\sigma_y} (Y - \bar{Y})$$

$$(Y - \bar{Y}) = \frac{\sum xy}{\sum x^2} (X - \bar{x})$$

$$(X - \bar{x}) = \frac{\sum xy}{\sum y^2} (Y - \bar{Y})$$

Principle of Least Squares

(2)

Why we are using principle of Least Squares?

(a)

Disadvantage of Curve of fit or Curve fitting?

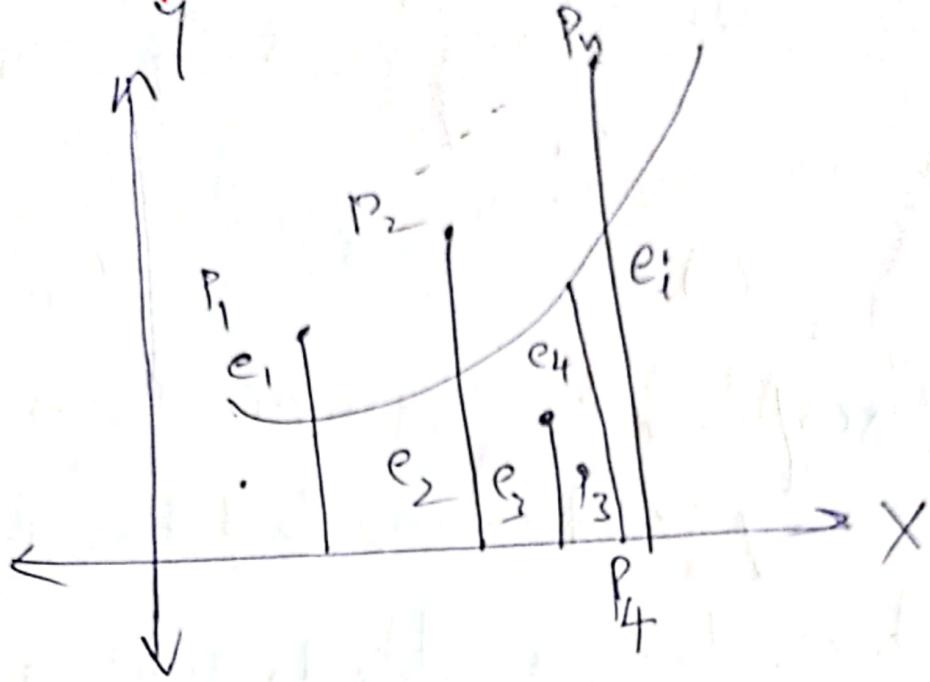
Legendre in 1806

The graphical method has the obvious
drawback of being unable to give a unique

Curve fit "So, it is necessary to get Unique
Curve of fit for any problem. So, we need one
platform to get Unique One.

" The principle of least Squares , it provides
an elegant procedure for fitting a Unique
Curve to a given data. "

Explanation:-



Let the Curve $y = a + bx + cx^2 + dx^3 + \dots$
n-data points $(x_1, y_1), (x_2, y_2), (x_3, y_3), (x_4, y_4), \dots$

Now, our Aim is to determine the Constants

a, b, c, d, \dots

i.e It represents the Curves of Best fit
if we find Values a, b, c, d .

→ If $n=m$ in the above equation.

We get Unique Set of Constants.

→ If $n > m \Rightarrow$ we cannot solve for \downarrow
 \downarrow
more than
equation
constants these constants.

So, we try to determine these values

of a, b, c, \dots, k , which satisfy all the equations as nearly as possible. So, it gives the best fit.

In these cases, we will apply the principle of Least Squares.

$e_1, e_2, e_3, \dots =$ errors (are positive).

We are Considering here, but some.

times it may -ve also. But, we're only +ve values. So, Thus, to give eqn from weightage to each error, we square each of these & form their sum.

$$E = e_1^2 + e_2^2 + \dots + e_n^2$$

Conclusion:-

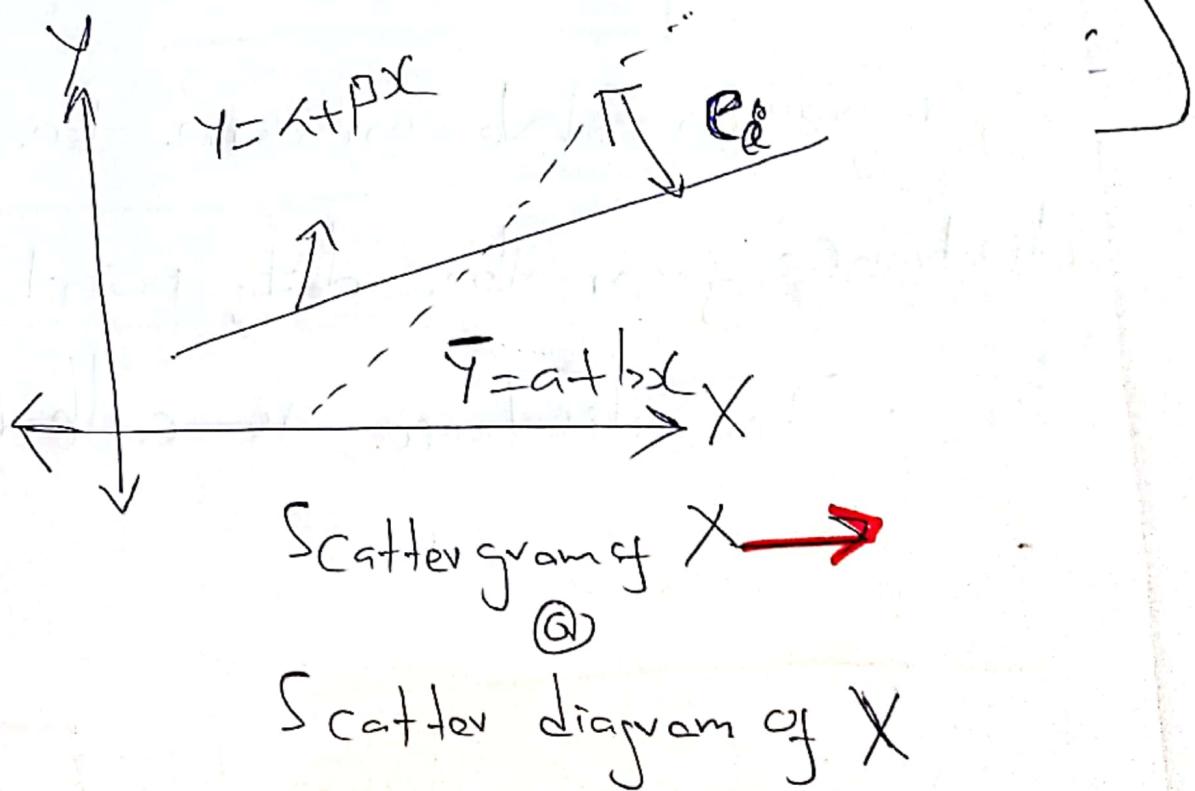
- The curve of best fit is that for which e 's are as small as possible.
- E , the sum of the squares of the errors is a minimum.

This is called the principle of least squares.
Suggested by Legendre in 1806.

be reasoning behind the method of least squares is quite simple.

47/57
50/115
wave veg

(ii) From among the many straight line drawn through the Scattergram, we pick out the one that is closest to (◎) best fits the observations."



Where $y = \alpha + \beta x$ is theoretical Regression line & ideal prediction line

$\bar{y} = a + b x$ is estimated regression line & equation used to make predictions.

where e_i = Difference b/w estimate & actual value

Value of γ at x_i

To determine how close the i^{th} data point is to the estimated line of regression, we measure the vertical distance from the data point to this line. This distance is called i^{th} residual.

③ i^{th} error

which is denoted by e_i .

$$Y = a + bx + e$$

$$Y_i = a + b x_i + e_i$$

$$e_i = Y_i - (a + b x_i)$$

63	47	55
149	150	
Sum of Squares		

This residual error is +ve

But in general, this is not possible

(i.e. always e is +ve values). We may still get +ve or -ve. So, the method of least squares utilizes not the sum of residuals but

Sum of their Squares.

Basic Aim of Least Squares Principle-

The technique selects the line through the data that is best in the sense that it minimizes the sum of squares of the residuals.

$$\sum e_i^2 = \sum [(y_i - (a + bx_i))^2]$$

minimum

We can use Normal Equations.

$$Y = ax + b$$

$$\sum Y = a \sum x + nb$$

Multiply $\sum x$ on both sides

$$\sum xy = a \sum x^2 + b \sum x$$

① & ② give two set of Linear Equations with
a & b are two Unknowns.

Solutions of a & b can also find out
without solving ① ② equations.

$$a = \frac{\sum y - b \sum x}{n} \quad \text{&} \quad b = \frac{n \sum xy - \sum x \cdot \sum y}{n \sum x^2 - (\sum x)^2}$$

Disadvantage:-

- The principle of least Squares does not help us to determine, the form of the appropriate Curve, which can fit a given data.
- It only determines, the best possible Values of the Constants in the equation when the form of the Curve is known before hand.
- The Selection of the Curve is a matter of experience & practical Considerations.

Inferences based on Least Squares

Estimator:-

In order to extend point estimates
(which are found by least squares),
initially to estimate the Comman Variate

$$\boxed{\hat{y}}$$

σ^2 is usually estimated in terms of
the vertical deviations of sample points
from least squares line.

Sample points = data points

$$S_e^2 = \frac{1}{n-2} \sum_{i=1}^n (y_i - (\hat{a} + \hat{b}x_i))^2$$

$$S_e^2 = \frac{1}{n-2} \sum (y - (\hat{a} + \hat{b}x))^2$$

where s_e = Standard error

plain

$(n-2) s_e^2$ = Residual sum of squares.

EVRCV \oplus Sum of Squares.

Explain Non-Linear Regression Equations

(Q)

Explain briefly Curvilinear Regression?

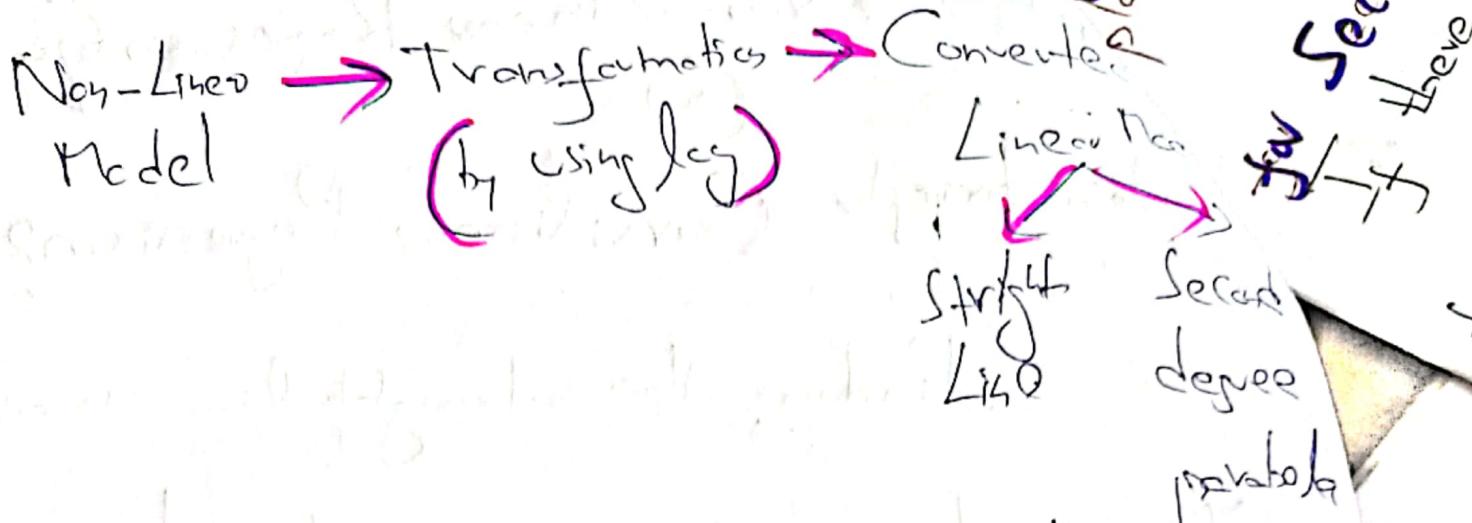
Y) finding the straight line model

$y = a + bx$ is inappropriate, because
the true regression curve is Non-Linear
and also called Curvilinear Regression)

Sometimes non-linearity is visibly
determined from Scatter diagram. These

Non-Linear Regression Equation can be ~~obtained~~
Converted to Linear Regression Equation (SLR).

Then, we can find easily for getting
the solution.



We can solve, easily

Now.

Non-Linear Modeling

$$1. Y = a \cdot b^x$$

$$2. Y = a \cdot x^b$$

$$3. Y = \frac{1}{ax + bx}$$

$$4. Y = a e^{bx}$$

Explain Polynomial Regression?

47/57 | 69
150/145 | 115
Java, Regre

Briefly & write Normal Equations

for Second degree polynomial Regression (Ans)

If there is no clear indication about

functional form of regression of Y on X .
we often assume that the underlying relationship is atleast, we behaved to the extent that it has a Taylor Series expansion and that the first terms of this expansion will yield a fairly good approximation.

Thus, we fit our data as a polynomial i.e predicting equation of form.

$$Y = a + bx + cx^2 + dx^3 + \dots$$

W.L.Fit Second degree parabola (y) *Method*

$$Y = a + bX + cX^2$$

$$\Sigma Y = na + b\Sigma X + c\Sigma X^2 \quad \text{--- (1)}$$

multiply ΣX on L.S.

$$\Sigma X Y = a\Sigma X + b\Sigma X^2 + c\Sigma X^3 \quad \text{--- (2)}$$

multiply ΣX^2 on L.S.

$$\Sigma X^2 Y = a\Sigma X^2 + b\Sigma X^3 + c\Sigma X^4 \quad \text{--- (3)}$$

(1)(2)(3) are called Normal Equations

for Second degree Polynomial Regression

C Parabola)

Procedure

Method of Least Squares

Procedure:-

To fit the straight
Line $y = a + bx$

Procedure:-

To fit the
parabola

$$y = a + bx + cx^2$$

Say

→ Substitute the observed
set of x values in this
eqn.

→ Write Normal Equations.

$$y = a + bx$$

$$y = a + bx + cx^2$$

→ Solve these Nf as simultaneous
equations for a & b

→ Substitute the values of a & b

$$\text{in } y = a + bx$$

→ which is required Line & best fit

Ans.

① Fit a Second degree parabola to the following data

x	0	1	2	3	4
y	1	1.8	1.2	2.5	6.3

$$a = 1.18$$

$$b = 1.13$$

$$c = 0.35$$

② Fit a 2nd degree parabola to the following data

x = 1.0	1.5	2.0	2.5	3.0	3.5	4.0
y = 1.1	1.3	1.6	2.0	2.7	3.4	4.1

$$a = 2.07$$

$$c = 0.061$$

$$b = 0.511$$

~~x₀~~ ~~3~~ fits 2nd degree parabola

x	1989	1990	1991	1992	1993	1994	1995	1996	1997
y	352	356	357	358	360	361	361	360	359
y	352	356	357	358	360	361	361	360	359

$$X = x - \bar{x}$$

$$Y = y - \bar{y}$$

$$Y = -1000106.41 + 1034.29X -$$

$$\bar{a} = \frac{694}{231}$$

$$0.267 \cancel{+^2}$$

$$\bar{b} = \frac{17}{20}$$

$$C = -\frac{247}{924}$$

H.W problems @ straight

line

- ① The following are data on the drying time of a certain point and the amount of an additive that is intended to reduce the drying time.

Amount of point additive x - grams	0	1	2	3	4	5	6	7	8
Drying time (G+S)Y	12	10.5	10	8.7	8.75	8.5	9.		

By using Method of least Squares to

predict the drying time of the point when 6.5 grams of additive is being used

① Fit a Second degree polynomial
(Parabola) $a = 12.2$
 $b = -14.5$
 $c = 0.183$
final $y = 7.9 \text{ hrs} @ x = 6.5$

② Fit a straight Line

② The following are data on the number
of twists required to break a certain
kind of forged alloy bar and percentage
of two alloying elements present in
metal

Number of twists (Y)	41	49	69	65	40	50	51	57	31
Percent of element A (x ₁)	1	2	3	4	1	2	3	4	1
Percent of element B (x ₂)	5	5	5	5	10	10	10	10	15

fit a least Squares Regression plane
 and use its equation to estimate the
 number of twists required to break one
 of the bars when $x_1 = 2.5$ &
 $x_2 = 12$?

~~Hint~~

$$Y = a + bx_1 + cx_2$$

$$\sum Y = na + b \sum x_1 + c \sum x_2 \quad (1)$$

multiply $\sum x_1$ on B.S

$$\sum Y x_1 = a \sum x_1 + b \sum x_1^2 + c \sum x_1 x_2 \quad (2)$$

multiply $\sum x_2$ on B.S

$$\sum Y x_2 = a \sum x_1 x_2 + b \sum x_1^2 \cdot \sum x_2 + c \sum x_1 x_2^2$$

drop box } & solve it ~~why~~

$$x_1 = \sqrt{0.5}$$

X

$$x_2 = 12$$

3 Fit a straight Line for the following

data Using (a) $y = a + bx$

(b) $y = ax + b$

(c) $x = a + by$

(d) $x = ay + b$

(e) $y = a + bx + cx^2$

(f) ~~$y = a + bx + c$~~

x	53	98	95	81	75	61	59	55
y	47	25	32	37	30	40	39	45

④ In the accompanying table [x] is the tensile force applied to a steel specimen in thousand of pounds and [y] is the resulting elongation of an inch.

x	1	2	3	4	5	6
y	14	33	40	63	76	85

Find the equation of least squares line and use it to predict the elongation when the tensile force is [3.5] of thousands pounds?

Hint $y = a + bx$

$$b = 14.69 \quad a = 1.12$$

$$x = 3.5 @ \boxed{y = 51.83}$$

⑥ Given the following data

x	10	30	40	60	80	90	110	140
y	10	20	40	40	50	70	80	90

Compute least square Regression Line?

~~hint~~ here not specific Geotrig

on y on x or x on y .

→ All these case use y on x

$$y = a + bx \quad \text{or} \quad \text{scratched}$$

$$n=8, \quad a=38.66 \quad b=0.162$$

$$\boxed{y = 38.66 + (0.162)x}$$

x	10	30	40	60	80	90	110	140
---	----	----	----	----	----	----	-----	-----

~~Ques.~~ The following table shows the ages X and systolic blood pressures Y of 12 women;

Age (X)	56	42	72	36	63	47	55	49	38	42	68
Blood Pressure (Y)	147	125	160	118	149	150	145	115	141	150	16
Determine the least square regression eqn of Y on X .											

~~Hint:~~

$$Y = a + bX$$

$$a = 150.06$$

$$b = 5.546$$

$$Y = 150.062 + 5.546X$$

Problems

- ① A panel of two judges P & Q graded Seven check performances by independently according mark as follows:

Performance	1	2	3	4	5	6	7
Marks by P	46	42	44	40	43	41	45
Marks by Q	40	38	36	35	39	37	41

The Eight performance, which judge Q could not attend, was awarded 37 Marks by judge P. If judge Q had also been present. How many marks would be expected to have been awarded by him to the eight performances?

- These type problems has been

selected by using Regression Co-efficients.

Performance	Marks by P X	$X = x - \bar{x}$	Marks by Q Y	$y = y - \bar{y}$	Xy
1	46	+1	9	-1	46
2	42	-1	1	1	38
3	44	+1	1	1	36
4	40	-3	0	0	35
5	43	0	4	-1	39
6	41	-2	4	-1	37
7	45	+2	25	2	41
	30				266

$$\bar{x} = \frac{\sum x}{n} = \frac{30}{7} = 4.3, \bar{y} = \frac{\sum y}{n} = \frac{266}{7} = 38.$$

$\therefore \bar{x}, \bar{y}$ can not fractions.

Now we have 2 Q-F's so, If we see the data, we have to calculate $Q_{1e}(Y)$.

Regression Equation y on X :

$$y - \bar{y} = b_{yx} (x - \bar{x})$$

$$\text{where } b_{yx} = \frac{\sum xy}{\sum x^2} = \frac{21}{28} = 0.75$$

$$y - 38 = 0.75 (x - 43)$$

$$y = 0.89 x + 37.36$$

But by data $x = 37$ by P

$$\therefore y = 0.89(37) + 37.36$$

$$\boxed{y = 33.5}$$

Conclusion: If Q would have been present, he would have awarded 33.5 marks to the eight persons.

(2)

In the following Table S is weight of Potassium bromide, which will dissolve in 100grams of water at $V^{\circ}\text{C}$. Fit an equation of the form $S = mT + b$. By the method of Least Squares.

$$\boxed{S = mT + b}$$

Use this relation to estimate S

$T = 50^\circ C$?

T	0	20	40	60	80
S	54	65	75	85	96

$$y = ax + b \quad S = mT + b \Rightarrow S \text{ is a straight line.}$$

$$N.E. \rightarrow \sum S = m \sum T + nb$$

Multiply $\sum T$

$$\sum TS = m \sum T^2 + b \sum T$$

by $T = 50^\circ C$ seeing this we
m, b are constant.

@ by data T, S

T	S	T^2	$+2$	TS
0	54			
20	65			
40	75			
60	85			
80	96			
200	375			
			12000	17080

$$\bar{T} = \frac{\sum T}{n}$$

$$\bar{S} = \frac{\sum S}{n}$$

Substitute Their Values in ①

$$\text{We will get. } m = 0.52, b = 54.2$$

$$\text{But } T = 50^\circ C \text{ given by data. } S = 0.52 \times 50 + 54.2$$

$$S = 80.2$$

(3)

Find the most likely production ~~Convs per day~~
 due to Rainfall 40 from the following data:-

	Rain fall (X)	Production (Y)
Average	$\bar{x} = 30$	$\bar{y} = 500 \text{ kgs}$
Standard Deviation	$\sigma_x = 5$	$\sigma_y = 100 \text{ kgs}$
Coefficient of Correlation	$r = 0.8$	-

i.e. $X = \text{Rainfall} = 40$ i.e we have to

find y .
 i.e. $R.E. Y \text{ on } X = (Y - \bar{y}) = b_{yx} (x - \bar{x})$

where $b_{yx} = \frac{\sum xy}{\sum x^2}$ or $r \cdot \frac{\bar{y}}{\sigma_x}$

data related to 2nd one formula.

Data: $\bar{x} = 30$, $\bar{y} = 500$
 $\sigma_x = 5$, $\sigma_y = 100$

$$b_{yx} = 16$$

$$y = 16x + 20$$

$$(Y - 500) = 16 \cdot 5 (x - 30)$$

When $x = 40$ i.e.

$$y = 660$$

$$(Y - 500) = 16 \cdot \frac{4}{100} (40 - 30)$$

$$(Y - 500) = \frac{4}{100} \cdot 10$$

$$\therefore Y = 500 + 4$$

(4) From a sample of 200 pairs of observations
following quantities were calculated.

$$\sum x = 11.34, \sum y = 20.78, \sum x^2 = 12.16,$$

$$\sum xy = 84.96, \sum x y = 22.13$$

From the above data. Show how to calculate
the coefficients of the equation $y = a + bx$.

$$y = a + bx$$

$$N \cdot \bar{c} + \sum y = n a + b \sum x$$

$$\underline{\sum xy = a \sum x + b \sum x^2}$$

Substituting above data,

$$20.78 = 200 a + 11.34$$

$$\underline{22.13 = a 11.34 + 12.16} \quad \text{Solving}$$

$$a = 0.0005$$

$$b = 1.82$$

$$\therefore y = 0.0005 + 1.82x$$

Topic:

Lecture NO:
 Link to Session
 Planner (SP): S.No....of SP
 Date Conducted:
 Page No:

(5) Criticise the following

Regression Coefficient of Y on X is +0.7 & X on Y is 3.2

$$r = \sqrt{b_{YX} \cdot b_{XY}} = \sqrt{0.7 \times 0.32} = \sqrt{2.24}$$

But Correlation Coefficient can not exceed 1

there is some inconsistency in the information given.

(6)

price	10	12	13	12	16	15
Demand	40	38	43	45	37	43

Calculate Regression Equations of Y on X from the data

Given below ~~assuming~~ from actual Means of

X & Y. Estimate the Likely demand when price=20?

$$x - \bar{x} = 20 - 13 = 7 \quad \text{we have to find } y \\ i.e. y - \bar{y} = b_{YX} (x - \bar{x})$$

$$b_{YX} = \frac{\sum xy}{n} = \frac{40 \times 10 + 38 \times 12 + 43 \times 13 + 45 \times 12 + 37 \times 16 + 43 \times 15}{6} = 1.25$$

Find out we will substitute
 $x = 20$

①

Price indices of Cotton and Wool are given
 for the 12 Months of a year. Calculate the
 Two-Stage Trend Indices by
 Method of Regression.

X (Price Index of Cotton)	78	77	82	85	88	83	87	82	81	77	76	83	97	93
Y (Price Index of Wool)	81	82	82	85	89	90	90	88	82	85	85	85	95	95

$$\begin{aligned}
 & \text{Cotton Price Index} \\
 & X = 82.5 \\
 & Y = 85.5 \\
 & \text{Wool Price Index} \\
 & X = 85.5 \\
 & Y = 88.5
 \end{aligned}$$

$$\begin{aligned}
 & \text{Cotton Price Index} \\
 & X = 82.5 \\
 & Y = 85.5 \\
 & \text{Wool Price Index} \\
 & X = 85.5 \\
 & Y = 88.5
 \end{aligned}$$

①

Price indices of Cotton and wool are given

for the 12 Months of a year.

Obtain the eqn of the Two lines of Regression of the Indices.

of Lines of Regression

X (Price Index of Cotton)	78	77	85	88	87	82	81	77	76	83	97	93
X (Price Index of Wool)	81	82	82	85	89	90	88	92	83	89	98	95

$\sum X$	948	$\sum Y$	1024	$\sum XY$	1024	$\sum X^2$	774	$\sum Y^2$	952	$\sum XY^2$	1024
\bar{X}	79	\bar{Y}	84	\bar{XY}	84	\bar{X}^2	64.25	\bar{Y}^2	70.5	\bar{XY}^2	70.5

$$\bar{Y} = 83.9$$

$$\bar{X} = 81.4$$

Subject: Pr
Faculty: Dr.K.K
Topic:

Angle b/w Two Regression Lines

$$\tan \theta = \frac{1-r^2}{\sqrt{r^2+2}} \cdot \frac{\sigma_x \sigma_y}{\sigma_x^2 + \sigma_y^2}$$

$\rightarrow \theta = \text{Acute Angle} \Rightarrow \tan \theta = \frac{1-r^2}{r} \cdot \frac{\sigma_y}{\sigma_x}$

$$r = -1 \text{ to } +1$$

$\rightarrow \theta = \text{Obtuse Angle} \Rightarrow \tan \theta = \frac{r^2-1}{\sqrt{r^2+2}} \cdot \frac{\sigma_y}{\sigma_x}$

$\Leftrightarrow r=0; \theta = \pi/2 \cdot \text{No relation b/w } x \& y \cdot$
ie They are Independent.

$\Rightarrow r=\pm 1 \Rightarrow \theta=0 \text{ or } \pi \cdot \text{Dependent}$

Two Regression Lines are parallel or coincide

ie Correlation b/w X & Y Variables are perfect.



Problems

① If α is the angle b/w two Regression Lines & S.I.
 y is twice the S.I. of x , $r = 0.25$. Find $\tan \alpha$?

$$\text{Sol} : \frac{\sigma_y}{\sigma_x} = 2 \quad (\text{given})$$

$$r = 0.25$$

$$\tan \alpha = \frac{1-r^2}{r} \cdot \frac{\sigma_x \sigma_y}{\sigma_x^2 + \sigma_y^2}$$

$$= \frac{1-r^2}{r} \cdot \frac{\sigma_x \cdot 2\sigma_y}{\sigma_x^2 + 4\sigma_x^2}$$

$$\text{Taking } \frac{\sigma_x^2}{\sigma_x^2 + 4\sigma_x^2} = \frac{1}{5}$$

$$= \frac{1-r^2}{r} \cdot \frac{2\sigma_y}{4(1+4)} = \frac{1-r^2}{r} \cdot \frac{2}{5}$$

$$= \frac{1-(0.25)^2}{0.25} \cdot \frac{2}{5}$$

$$= 1.5$$

② If $\sigma_x = \sigma_y = \sigma$. The angle b/w two Regression Lines is tan $(\frac{\pi}{4})$?

$$\tan \alpha = \frac{1-r^2}{r} \cdot \frac{\sigma_x \sigma_y}{\sigma_x^2 + \sigma_y^2}$$

$$\sigma_x = \sigma_y = \sigma$$

$$\tan \theta = \left[\frac{r^2}{r} \cdot \frac{\frac{4}{3}}{2r} \right]$$

$$\tan \theta = \frac{r^2}{2r}$$

$$\theta = \tan^{-1} \left(\frac{r^2}{2r} \right) \rightarrow \textcircled{1}$$

$$\text{But } \tan \theta = \tan^{-1} \left(\frac{4}{3} \right) \rightarrow \textcircled{2}$$

$$\textcircled{1} = \textcircled{2}$$

$$\tan \left(\frac{r^2}{2r} \right) = \tan \left(\frac{4}{3} \right)$$

$$\frac{r^2}{2r} < \frac{4}{3}$$

$$3 - 3r^2 = 8V$$

$$3r^2 + 8V - 3 = 0$$

$$r = \frac{1}{2} \quad \textcircled{3} \quad 1 = \pm 3$$

But $r = -3$ \times not possible

$$r^2 + 8V - 3 = 0$$
$$r = \frac{1}{2} \quad \textcircled{3} \quad V = -3$$

$$\therefore r = \frac{1}{2}$$

(2) The tangent of angle between regression lines is
or if $\alpha = \frac{1}{2} \pi - \sqrt{1 - \text{Correlation coefficient}}$
What is it?

$$\sigma_x = \frac{5}{2}$$

$$r = \frac{1}{2}$$

$$r = \frac{-3}{2} \times \frac{1}{2} = -\frac{3}{4}$$

Finding Mean Values & r_{xy} by using
Two Regression functions

Q Find the Mean Values of Variable X & Y & Correlation Coefficient from the following Regression Functions

$$2Y - X - 50 = 0 \quad \left\{ \begin{array}{l} \text{Eqn 1} \\ \text{Eqn 2} \end{array} \right. \quad (1)$$

$$3Y - 2X - 10 = 0$$

$$\bar{x} = 130, \bar{y} = 90$$

~~Mean~~

From (1) i.e. $2Y - X - 50 = 0$

$$2Y = X + 50$$

$$Y = \frac{X}{2} + 25$$

$$3Y - 2X - 10 = 0$$

(2) This is Your XRF

$$2X = 3Y - 10$$

$$X = \frac{3Y}{2} - 5$$

(3) This is XYRF

$$\therefore b_{yx} = \frac{1}{2}, \quad b_{xy} = \frac{3}{2}$$

$$V = \sqrt{b_{xy} \cdot b_{yx}} = \sqrt{\frac{3}{4} \cdot \frac{1}{2}} = \sqrt{\frac{3}{4}}$$

$$\therefore r = 0.866$$

Ques Test whether the equations $2x + 3y = 4$ & $x - y = 5$ represent valid Regression Lines?

$r^2 = -\frac{1}{2}$ $r^2 = -ve$ so, it gives equation
does not represent valid Regression Lines.

- Q) For 20 army personal Regression of weight of kidneys (y) on weight of heart X is $y = 359.8 + 6.394x$
& the Regression of weight of heart on weight of kidneys is $X = 1.212Y + 2.461$. find the correlation coefficient between the two variables also their means?

$$\bar{x} = 197720.0$$

$$\bar{y} = 78890.28$$

$$V = 0.695$$

Theorem 10.2.1. If Ω is a bounded domain in \mathbb{R}^n with smooth boundary, then the classification of eigenvalues and the mapping of every λ to its corresponding eigenfunction u_λ is continuous.

The continuity of $\lambda \mapsto u_\lambda$ follows from the fact that

$$\|u_\lambda - u_{\lambda'}\|_{H^1(\Omega)} \leq C \sqrt{\lambda - \lambda'}$$

for some constant $C > 0$.

Continuity of the mapping $\lambda \mapsto \lambda^{-1/2} u_\lambda$ follows from the fact that

the mapping $\lambda \mapsto \lambda^{-1/2}$ is continuous and the mapping $\lambda \mapsto u_\lambda$ is continuous.

Continuity of the mapping $\lambda \mapsto \lambda^{1/2} u_\lambda$ follows from the fact that

the mapping $\lambda \mapsto \lambda^{1/2}$ is continuous and the mapping $\lambda \mapsto u_\lambda$ is continuous.

$$\begin{aligned} (\partial_{x_i})_{\Omega_D}^2 u_\lambda &= \lambda u_\lambda \\ \partial_{x_i} \partial_{x_j} u_\lambda &= \lambda^{1/2} \delta_{ij} u_\lambda \\ \partial_{x_i} \partial_{x_j} u_\lambda &= \lambda^{1/2} \delta_{ij} u_\lambda \end{aligned}$$

Continuity of the mapping $\lambda \mapsto \lambda^{1/2} u_\lambda$ follows from the fact that

the mapping $\lambda \mapsto \lambda^{1/2}$ is continuous and the mapping $\lambda \mapsto u_\lambda$ is continuous.

Continuity of the mapping $\lambda \mapsto \lambda^{-1/2} u_\lambda$ follows from the fact that

the mapping $\lambda \mapsto \lambda^{-1/2}$ is continuous and the mapping $\lambda \mapsto u_\lambda$ is continuous.

(3) The equation of two R. Lines are $7x - 16y + 15 = 0$

$S_y - 4 > 1 - 3 = 0$ - find the Coefficient of Correlation

& the Means of \bar{x} & \bar{y} ?

$$\bar{x} = 0.1034, \bar{y} = 0.5172 \quad r = \frac{5}{\sqrt{35}}$$

Ans: $\bar{x} = 0.1034, \bar{y} = 0.5172$

and removal bilab. tongue & molar

(31-a)

370

34-6

32-5

= part C

$$0.05501 = R$$

$$= 0.055 = R$$

Properties of Regression Coefficients

① The Correlation Coefficient r is the geometric mean b/w the two regression Coefficients.

$$r^2 = b_{yx} \cdot b_{xy}$$

(i)

$$r^2 = r \frac{\sigma_y}{\sigma_x} \cdot r \frac{\sigma_x}{\sigma_y}$$

(ii)

$$r = \sqrt{b_{yx} \cdot b_{xy}}$$

$$(iii) r = \sqrt{r \frac{\sigma_y}{\sigma_x} \cdot r \frac{\sigma_x}{\sigma_y}}$$

Where b_{yx} or $r \frac{\sigma_y}{\sigma_x}$ = The regression Coefficient of Y on X

b_{xy} or $r \frac{\sigma_x}{\sigma_y}$ = The regression Coefficient of X on Y .

(2) $r = \frac{\sum xy}{\sqrt{\sum x^2 \sum y^2}}$ for y on x

$r = \frac{\sum xy}{\sqrt{\sum x^2 \sum y^2}}$ for x on y .

(3) If θ is the angle b/w two regression lines. Then $\tan \theta = \frac{1-r^2}{\sqrt{\frac{\sigma_x^2 \sigma_y^2}{\sigma_x^2 + \sigma_y^2}}}$

$$\tan \theta = \frac{1-r^2}{\sqrt{\frac{\sigma_x^2 \sigma_y^2}{\sigma_x^2 + \sigma_y^2}}}$$

$\Rightarrow r=0, \tan \theta \rightarrow \infty (\theta > \pi/2)$. Then

two lines of regression are perpendicular to each other

$\Rightarrow r=\pm 1, \tan \theta = 0 (\theta = 0 \text{ or } \pi)$. Then, two lines of regression coincide.

i.e. Here, there is a perfect Correlation b/w the two Variables.

$$④ r = \frac{\sigma_x^2 + \sigma_y^2 - \frac{2}{n} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{2\sigma_x \sigma_y}$$

$$⑤ r = \frac{\sum xy}{n \bar{x} \bar{y}}$$

8

$$r = \frac{\sum (x - \bar{x})(y - \bar{y})}{n \bar{x} \bar{y}}$$

↓

$\bar{x}\bar{y}$ ave
not fractions

$\bar{x}\bar{y}$ ave
fractions.

$$⑥ \sigma_z^2 = \sigma_x^2 + \sigma_y^2 - 2r\sigma_x\sigma_y$$

⑦ If one of the regression coefficients is greater than Unity, the other must be less than

Unity.

$$b_{yx} \cdot b_{xy} = r^2 < 1$$

$$b_{xy} \leq \frac{1}{b_{yx}} < 1$$

$$1 - e^{b_{xy}} > 1 \quad b_{xy} < 1$$

(8) $r^2 \leq 1$ & σ_x, σ_y as +ve.

(9) +ve sign gives the arc clockwise.
Hence lines.

$$S_{xx} = \sum (x - \bar{x}) \cdot \sum (x - \bar{x})$$

Least Squares

Notations

in Regression

$$S_{xy} = \sum (x - \bar{x}) \cdot \sum (y - \bar{y})$$

$$S_{yy} = \sum (y - \bar{y}) \cdot \sum (y - \bar{y})$$

$$\Rightarrow S_{SR} = \frac{S_{xx} S_{yy} - S_{xy}^2}{S_{xx}}$$

where S_{SR} = Sum of Squares of the residuals.

$$\Rightarrow \text{Total SS} = SSR + SSE$$

where SSR = Sum of Squares for Regressions

SSE = Sum of Squares for error.

r^2 Values if Coefficient Correlation

Prop

Correlation Coefficient Values

① $0 < r < 0.5 \rightarrow$ low degree +vely Correlated

$0 < r < 0.5 \rightarrow$ " " -Vely "

② $0.5 < r < 0.75 \rightarrow$ Moderate +vely "

$-0.5 < r < -0.75 \rightarrow$ " " -Vely "

③ $0.75 < r < 1 \rightarrow$ high degree +vely "

$-0.75 < r < -1 \rightarrow$ high degree -Vely "

④ $0.9 < r < 1 \rightarrow$ Very high degree +vely "

$-0.9 < r < -1 \rightarrow$ " " -Vely "

$r = +1 \rightarrow$ perfectly +ve

$r = -1 \rightarrow$ perfectly -ve

$r = 0 \rightarrow$ Un Correlated

Properties of Least Squares Estimators

1. The least squares estimators for a straight line $y = a + bx$ are $\boxed{a \& b}$.
2. $y = a + bx + e$ is a RV with $E(x) = 0$ & $V(x) = \sigma^2$.
3. The values of a & b depend on the observed y 's.
4. The least squares estimators of the regression coefficients may be viewed as Random Variables.

Inferences Concerning the Regression

Coefficients;

→ We know that, for estimating the linear relationship b/w x & y for purposes of certain applications of prediction.

→ The experiments may also be interested in drawing certain inferences about the slope and intercept.

→ In order to allow for the testing of hypothesis and the construction of confidence interval on a & b .

→ Now, here we need to assume that e_i is normally distributed

$$y = a + bx + e_i$$

e_i = residual error $i = 1, 2, \dots, n$

→ In this, Inferential Concerning

Regression Coefficients, we use

χ^2 -test & also t-test

~~t-test~~
(Test statistic)

$$T = \frac{b - B}{\sqrt{\sum(x - \bar{x})^2}} \quad \text{S} \quad b - B$$

$$\text{S} \quad \frac{b - B}{\sqrt{\sum(x - \bar{x})^2}} \quad \text{S}_{xx}$$

&
Confidence Interval

$$\left(B - t_{\frac{\alpha}{2}} \frac{S}{\sqrt{S_{xx}}} \right) < b < \left(B + t_{\frac{\alpha}{2}} \frac{S}{\sqrt{S_{xx}}} \right)$$

where $t_{\frac{\alpha}{2}}$ = t-distribution with $(n-2)$ degrees of freedom.

Simple Linear Regression 5M

Case Study:-

Write briefly Simple Linear Regression Case Study? (SLR - Case Study);-

- In the manufacture of Commercial wood products, it is important to estimate relationship b/w the density of a wood product and its stiffness.
- A relation near type of particleboard is being considered that can be formed with considerations, more ease than the accepted Commercial product.
- It is necessary to know at what density the stiffness is comparable to that of the well-known, well documented product (Commercial product)

→ A study was done by Tervore \rightarrow
E. Corneve, they investigate of certain a
Mechanical properties of a wood foam \rightarrow
Composite.

→ These particle boards were produced at
densities ranging from roughly 8 to 26 pounds
& the stiffness was measured in pounds per
square inch.

→ Now, it is necessary for the data
analysis to focus on an appropriate fit to
the data & use inferential methods discussed

in this.
→ Hypothesis testing on the slope of the
regression, as well as confidence predictions
interval estimation may well be appropriate

→ Now, we will begin by demonstrating
a simple scatter plot of the data with
a simple linear regression superimposed.

For example:-

$$y = -25431.739 + 3884.9x$$

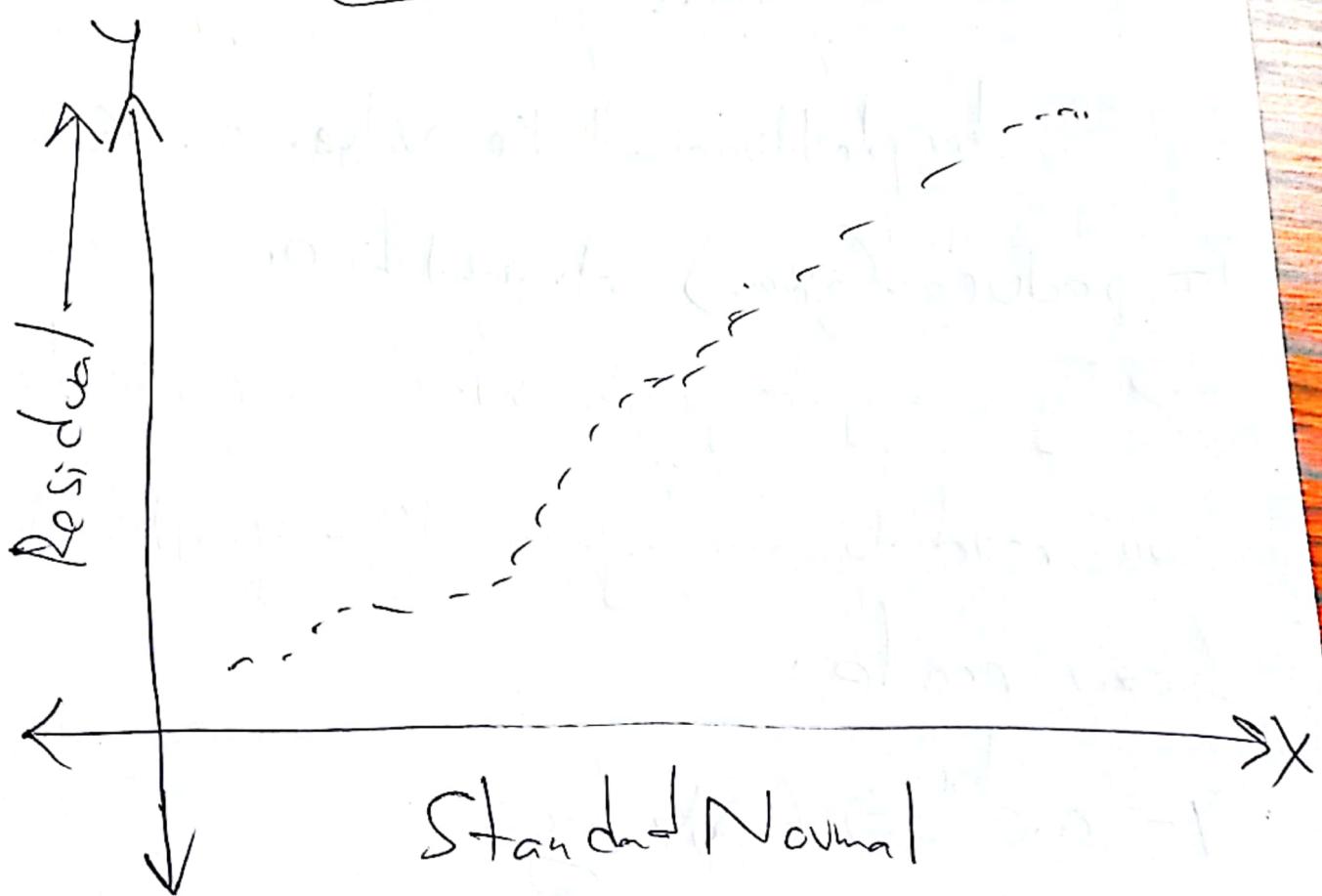


Figure:- Normal probability plot of residuals
for wood density data.

X-axis, which is horizontal axis
represents the empirical normal distribution

Y-axis, which is vertical axis represents the Residuals.

by plotting all the values on X & Y axis
it produces (gives) straight line.
→ If we get graph which is non-linear
we need to transform the graphs into
linear mode.

$$y = a \cdot e^{bx} \rightarrow \text{Apply } \log$$

$$\Rightarrow \log y = \log a + bx \log e$$

$$\Rightarrow y = A + Bx$$
 Then we have

to apply normal procedure of Regression Equations (SLR).

Case - Study - 2:-

Case - Study - 3:-

We can write
from Real time
applications.

This is left for H.W.

Fitted Model:-

Explain fitted Model in Regression?

→ write about SLR (which is Simple Linear Regression Equation &)
+

and also write Least Squares Methods
stating that $(Y = a + bx + \epsilon)$ ^{no residual error also.}

→ Write about few lines of

Curvilinear Regression .

~~Exponential~~ ^{bx}, $y = a \cdot b^x$ all these

$\rightarrow y = a \cdot e^{bx}$, ^{do Linear Regression}

Non-Linear Equations ^{to do} (ie Apply \log) for
Equations by transforming

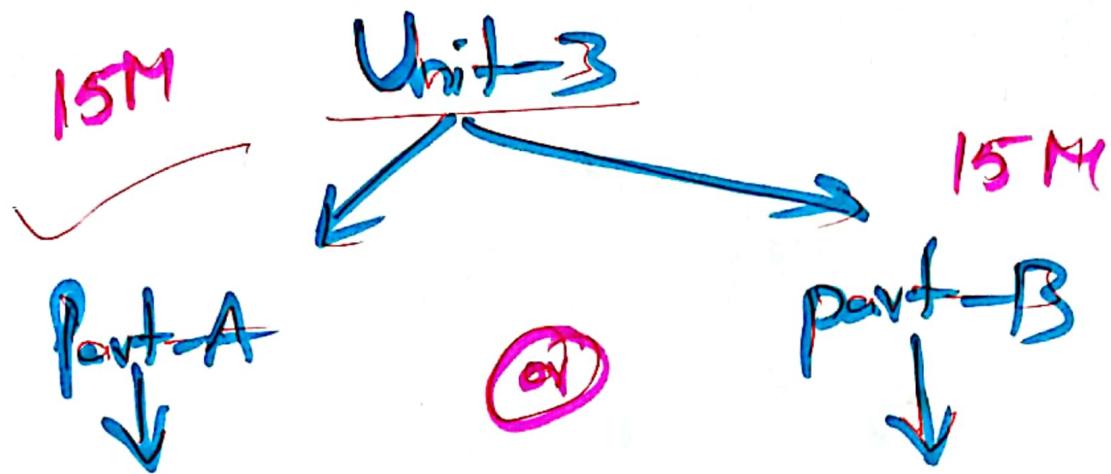
Convincing Questions

1. How do you think the world would change if there were no more oil?
2. What would happen if we stopped using oil?
3. What are the pros and cons of using oil?
4. What are the benefits of using oil?
5. What are the disadvantages of using oil?
6. What are the alternatives to oil?
7. What are the environmental impacts of oil?
8. What are the social impacts of oil?
9. What are the economic impacts of oil?
10. What are the political impacts of oil?

Prediction:

There are several reasons for building a linear regression. One, of course, is predict response values at one or more values of the independent variable. In this section, the focus is on errors associated with prediction.

The eqn $y = \hat{y} + e$, may be used to predict estimates the mean response. It may be used to predict a single value. We would expect the error of prediction to be higher in the case of a single predict value than in the case where mean is predicted. Then, we can use here confidence interval also by applying t-distribution.



Continuous Probability
Distributions

↓
Normal distributions

2. Areas under the Normal
Curve

3. Applications of N.P

4. Normal Approximation to
B.D

Fundamental Sampling
distributions

- ↓
1. Random Sampling
 2. Sampling distributions
 3. Sampling distributions
of Means
 4. Central Limit theorem
 5. Sampling distributions
of S^2

6. t - Distributions

7. F - distribution

of X on Y & Y on X .

Regression Equation → Equation of a straight line.

Regression Equation is an Algebraic expression of the Regression Line.

It can be classified into Regression Equations, Regression Coefficients.

Individual Observations & Group Discussion

Standard form
of Regression

Equation

$$Y = a + bX$$

where a, b = Constants

a_i - Value of y , when $x=0$ - It is called y -Intercept.

b_i - Value of slope of Regression Line.
called Regression Coefficient y on x .

To find a & b - with help of Normal Equations

$$\begin{array}{l} \textcircled{1} \text{ Regression Equation of } Y \text{ on } X_i : \\ y = a + b x \\ \sum y = n a + b \sum x \\ \sum xy = a \sum x + b \sum x^2 \end{array} \quad \begin{array}{l} \textcircled{2} \text{ R.E of } X \text{ on } Y_i : \\ x = a + b y \\ \sum x = n a + b \sum y \\ \sum x y = a \sum y + b \sum y^2 \end{array}$$

~~Two types of problems~~

① Determine the Equation of a straight line which best fits the data. (i.e Regression Equation.)

X	10	12	13	16	17	20	25
y	10	22	24	27	29	33	37

$$\text{St. Line } y = a + b x$$

$$\text{Normal Equations} \quad \sum y = b \sum x + n a$$

$$\sum xy = b \sum x^2 + a \sum x c$$

Topic:

What is Meant by Curve fitting? Explain the principle of Least Squares on Legendre principle of Least Squares?

Let (x_i, y_i) where $i = 1 \text{ to } n$
 $j = 1 \text{ to } n$.

Given set of n -pairs of values.

x_i = Independent Variable

y = Dependent Variable

Analytic expression of the form $y = f(x)$. It is useful in study of Correlation & Regression. In practical statistics enables us to represent the relation ship b/w two Variables by Simple Algebraic expressions

Ex:- polynomial & exponential functions

etc.

It is used to estimate the value of first

Variable which would correspond to specific values of other Variable.

Values of other Variable.

The Legendre principle of Least Square

Consists of minimising, the sum of square

deviation of the actual values from the estimated values given by the degree of fitting curve of like.

Non Linear Regression \leftrightarrow Curve Fitting

In this Method, we have to use Method of Least

Squares to fit straight line

$$① y = a \cdot b^x \quad \text{— Power Curve}$$

$$y = a \cdot e^{bx} = \text{use log base e}$$

$$② y = a \cdot e^{bx} \quad \text{— Exponential Curve}$$

$$y = a \cdot b^x = \text{use log base e}$$

$$③ y = a \cdot x^b \quad \text{— power law}$$

$$④ y = a + bx + cx^2 \quad \text{— Second degree parabola}$$

$$⑤ y = a \cdot b^x; \quad \text{apply log}$$

$$\log y = \log a + b \log x$$

$$y = A + Bx$$

$$y = \log y$$

$$⑥ y = a e^{bx}$$

$$\log y = \log a + b \log x$$

$$y = A + Bx$$

$$\Sigma y = nA + b \Sigma x$$

$$\Sigma xy = A \Sigma x + B \Sigma x^2$$

$$\Sigma x^2 = a \Sigma x^2 + b \Sigma x^3$$

$$⑦ y = a \cdot x^b$$

$$\log y = \log a + b \log x$$

$$y = A + Bx$$

$$\Sigma y = nA + b \Sigma x$$

$$\Sigma xy = A \Sigma x + B \Sigma x^2$$

$$\Sigma x^2 = a \Sigma x^2 + b \Sigma x^3$$

$$⑧ y = a + bx + cx^2$$

$$\Sigma y = nA + b \Sigma x + c \Sigma x^2$$

$$\Sigma xy = a \Sigma x + b \Sigma x^2 + c \Sigma x^3$$

$$\Sigma x^2 = a \Sigma x^2 + b \Sigma x^3 + c \Sigma x^4$$

$$y = a + bx \quad \therefore y = a + bx + c \\ = 0.82 + 1.56x$$

③ The following data pertains to the number of jobs per day and the total processes unit time required.

Fit a straight Line - Estimate the Mean CPU Time at TS

Fit a straight Line by Method of Least Squares.

$x = 3.5$	Fit a straight Line by Method of Least Squares.
No. of Jobs	1 2 3 4 5
CPU Time	2 5 4 9 10

$$y = a + bx$$

$$3a + 15b$$

$$a = 0, b = 2$$

$$10 = 15a + 55b$$

$$\therefore y = 0 + 2x \Rightarrow y = 2x$$

$$x = 3.5 \Rightarrow y = 2(3.5)$$

$$y = 7$$

④ Find Least Square Regression Equation of X_3 on X_1 & X_2 ?

$$\Sigma X_3 = a + bX_1 + cX_2 \Rightarrow$$

$$\Sigma X_1 X_3 =$$

Mutl. X_1, X_2

$$\Sigma X_2 X_3 =$$

X_1	X_2	X_3	X_1^2	X_2^2	$X_1 X_2$	$X_1 X_3$	$X_2 X_3$
1	2	2	1	4	2	2	4
2	5	5	4	25	10	10	25
3	4	4	9	16	12	12	16
4	9	9	16	81	36	36	81
5	10	10	25	100	50	50	100
15	30	30	55	145	90	90	145

$$\text{Substitution & solve} \\ a = 61.40, b = -3.65 \\ c = 2.54$$

$$\therefore X_3 = 61.40 - 3.65X_1 + 2.54X_2$$

Quantile Deviations : Mean : Standard Deviation = 10 : 12 : 15

Quantile Deviations = $\frac{Q_3 - Q_1}{2} = \frac{2}{5} \sigma$

Mean Deviations = $\frac{4}{5} \sigma$

2. Define probable error (σ)

"It is such that the probability of an error falling within the limits ($M-\sigma$ & $M+\sigma$) is exactly equal to the chance of an error falling outside these limits."

i.e. chance of an error lying within $M-\sigma$ & $M+\sigma$ = $\frac{1}{2}$

Example:- In a Industry, lot of articles manufactured to certain specifications is subject to small errors. In fact, measurement of any physical quantity should be light error.

1. Skewness :- $\beta_1 = \frac{\mu_3^2}{\mu_2^3} = 0$

2. Kurtosis:-

$$\beta_2 = \frac{\mu_4}{\mu_2^2} = \frac{3\sigma^4}{(\sigma^2)^2} = 3$$

3. Recurrence relation for Moments of N.D.:-

$$\mu_{2n} = \sigma^{(2n-1)} \mu_{2n-2}$$

all odd moments = Variances

all even moments = Excess.

Moments about Origin	Moments about Mean
-------------------------	-----------------------

1. 1st Moment = Mean = $\mu_1 = b$

1. Mean = $\mu_1 = b$
2) Variance = σ^2

Table of functions

3 types R.V. = X
a function.

Now we have to
consider R.V. function.

P.D.F

distribution
functions

Let X be a R.V.
Then function F
is defined for all x

$$F(x) = P(X \leq x) = \int_{-\infty}^x f(x) dx$$



is called P.D.F.

P.D.F Properties:-

1. If F is d.f. of
1D R.V. is X ,

Then

$$\text{① } 0 \leq F(x) \leq 1$$

$$\text{② } F(x_1) \leq F(x_2)$$

* **Ques**

All distribution functions
are monotonically
non-decreasing
& lie b/w 0 & 1.

P.m.f

Probability mass function

Let X = discrete R.V
with distinct values
 x_1, x_2, \dots, x_n . Then
function $p(x)$ is

$$P(x) = \begin{cases} p(X=x_i), & \text{if } x=x_i \\ 0, & \text{if } x \neq x_i \\ i = 1, 2, \dots \end{cases}$$

is called P.m.f.

P.m.f. properties:-

$$\text{1. } p(x_i) \geq 0$$

$\forall i$

$$\text{2. } \sum p(x) = 1$$

Ques

3. The set of
values, which
takes is

P.d.f

probability density function

Let X = Continuous R.V
with p.d.f $f(x)$.

Then the function

The pdf is the derivative
of the probability distribu-
tion function.

P.d.f Properties:-

$$\text{1. } f(x) \geq 0$$

$\forall x$

$$\text{2. } \int_{-\infty}^{+\infty} f(x) dx = 1$$

$-\infty$

Explain Mean, Median & Mode?

Mean: If X is a Continuous RV with probability $f(x)$. Then

$$M = \int_{-\infty}^{+\infty} x \cdot f(x) dx = \int_{-\infty}^b x \cdot f(x) dx$$

Median: Median is the point, which divides the entire distribution into two equal parts. Then

$$\int_a^M f(x) dx = \int_M^b f(x) dx = \frac{1}{2}$$

Mode: Mode is the value of x for which $f(x)$ is Maximum. Mode is given by

$$\begin{aligned} f'(x) &= 0 \\ f''(x) &< 0 \quad \text{for } a < x < b \end{aligned}$$