Regression Analysis:                    chapter-11

① Regression Analysis is a statistical tool for investing the relationship b/w a dependent variable & one or more independent variables.

Note: Regression Analysis is widely used for prediction and forecasting.

Applications: Economies,

Ex: Suppose you are marketing Analyst for Toys we gather the following data.

($) Advertising sels          Sales.

| (X) | | Sales (Y) |
|---|---|---|
| 1 | | 1 |
| 2 | | 1 |
| 3 | | 2 |
| 34 | | 2 |
| 5 | | 4 |

regressor variable ↙          ← Repense variable.

Now we want to relationship b/w sales & adverting cost. Since we can decide How much money we want to spend (ie. sort of control variable) for advertising. but sales amount is not a control variable ie. we can not control the sales amount. So sales amount is dependent variable it depends on advertisement. So sales amount is dependent variable and Advertising cost is independent variable.
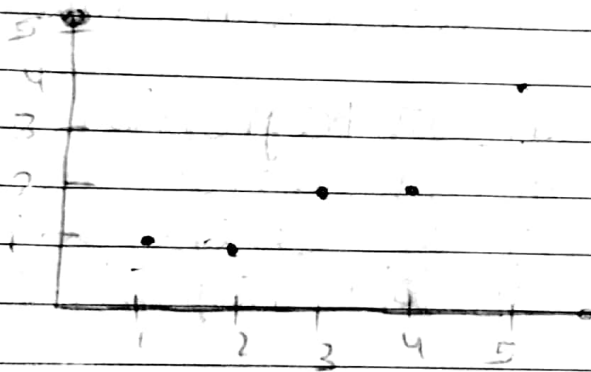
## Scatter Plot

② Let $X_i$ stands for regressor variable,
  $Y_i$ " " Response variable.

| $X_i$ | $Y_i$ |
|-------|-------|
| 1 | 1 |
| 2 | 1 |
| 3 | 2 |
| 4 | 2 |
| 5 | 4 |

Let here we have 5 observation
ie. $(x_1, y_1) = (1, 1)$
$(x_2, y_2) = (2, 1)$
$(x_3, y_3) = (3, 2)$
$(x_4, y_4) = (4, 2)$
$(x_5, y_5) = (5, 4)$

{ A scatter plot is a mathematical diagram to display values or two vate variables for a set of data }

{ A scatter plot are used to investigate the possible relationship b/w the variables. }



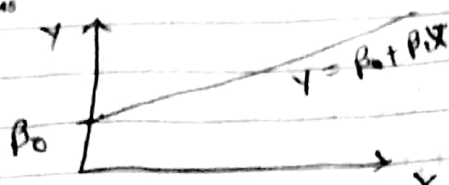This scatter plot indicate linear relationship b/w regressor variable or response variable

If Scatter plot indicates linear relationship b/w the variables then we go for linear model.
If Scatter plot plot indicates non-linear relationship b/w x and y then we go for quadratic or cubic or higher order polynomial model.

③



linear relationship b/w x & y

$\beta_1$ - slop

$\beta_0$ - Y-intercept.

Note!:- for a given x does not always give the same value for Y.

* a random sample of size n by the set $\{(x_i, y_i) | i=1,...n\}$ If additional samples were taken using exactly the same values of $x$, we should expect the y values to vary. Hence, the value $y_i$ in the ordered pair $(x_i, y_i)$ is a value of some random variable $Y_i$.

The Simple Linear Regression Model:- (SLR):-
Simple linear regression model is a model with a single regressor x that has a linear relationship with response variable y.

The Simple linear regression model is

$$Y = \beta_0 + \beta_1 X + \varepsilon$$

where
$\varepsilon \rightarrow$ is random error component
$\beta_0 \rightarrow$ Y-intercept
$\beta_1 \rightarrow$ slop.
$X \rightarrow$ Regressor variable.
$Y \rightarrow$ Response variable.

2011
| 5 | 5 | WK |
| 5 | 6 | 45 |
| 12 | 13 | 46 |
| 19 | 20 | 47 |
| 26 | 27 | 48 |
| | | 49 |

NOTES: Where $\beta_0$ & $\beta_1$ are the parameter and they are typically unknown. ~~$\beta_0$ be cool~~ are called regression cofficient.
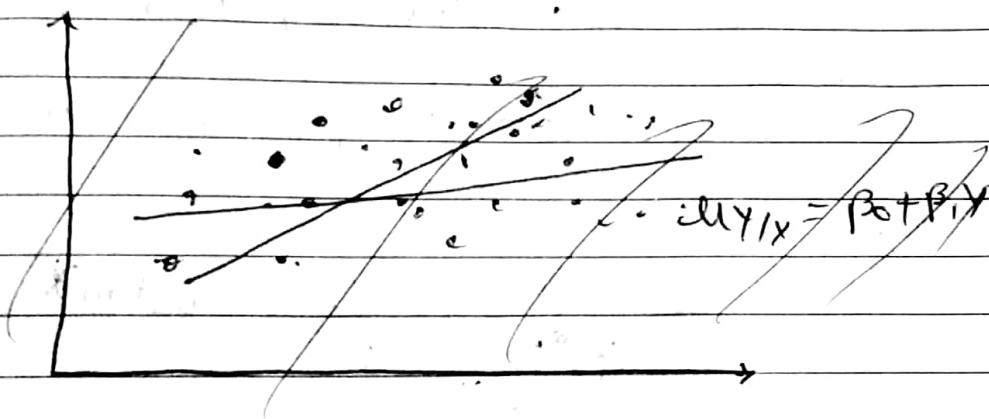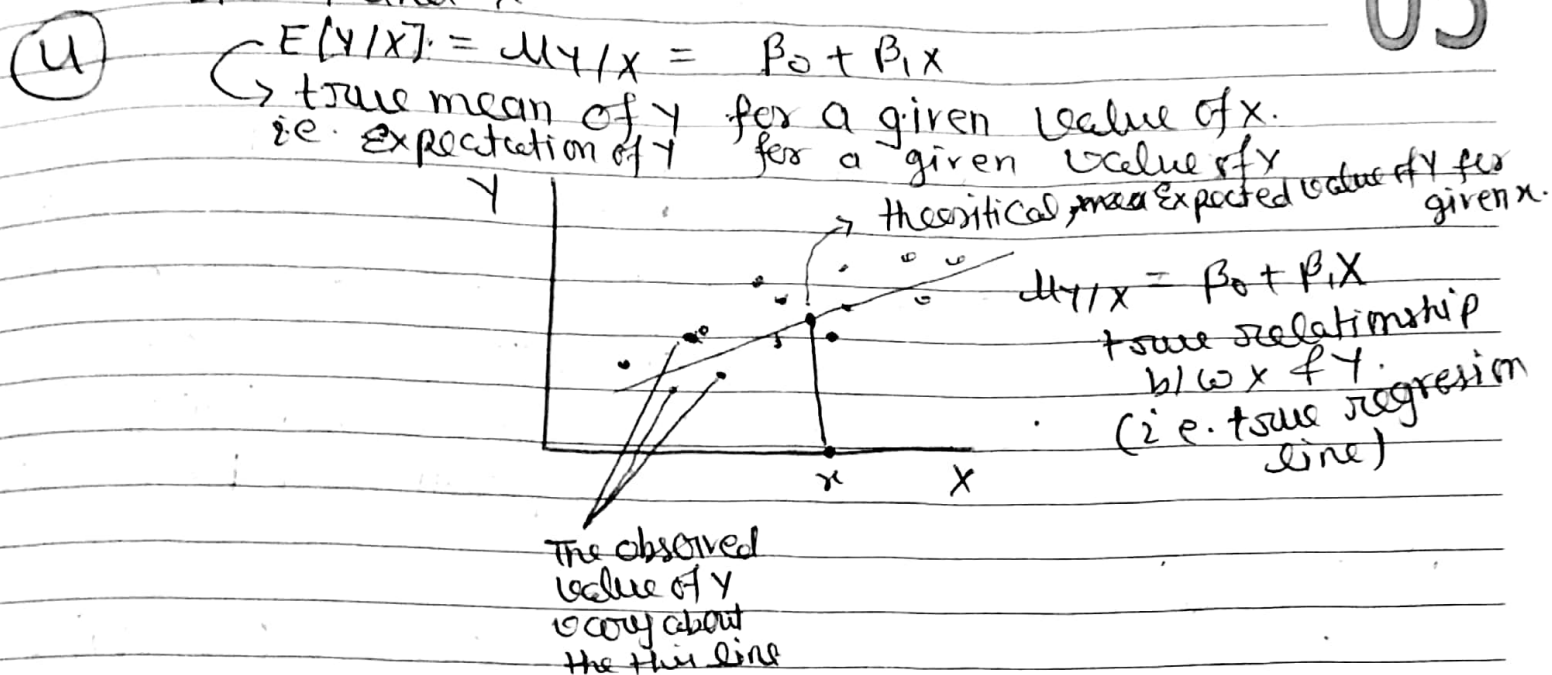
* We will assume a linear relationship b/w Y and X

(u) $\quad\hookrightarrow E(Y/X) = \mu_{Y/X} = \beta_0 + \beta_1 X$
$\quad\hookrightarrow$ true mean of Y for a given value of X.
$\quad$ i.e. expectation of Y for a given value of Y

$\hookrightarrow$ theoritical mean Expected value of Y for given x.

$\mu_{Y/X} = \beta_0 + \beta_1 X$
true relationship
b/w x + Y
(i.e. true regression line)

The observed value of Y vary about the this line

$\mu_{Y/X} = \beta_0 + \beta_1 X$

$(x_i, y_i) \to$ observed data point.

$\mu_{Y/X} = \beta_0 + \beta_1 X.$

$\hat{y}_i$

$6$

$\to \varepsilon_i$

here $\varepsilon_i$ represents that Y values will vary about the line and value of Y do not fall precisely on the line $\mu_{Y/X} = \beta_0 + \beta_1 X$

* We will use sample data to obtain the estimated regression line ( fitted Regression line)
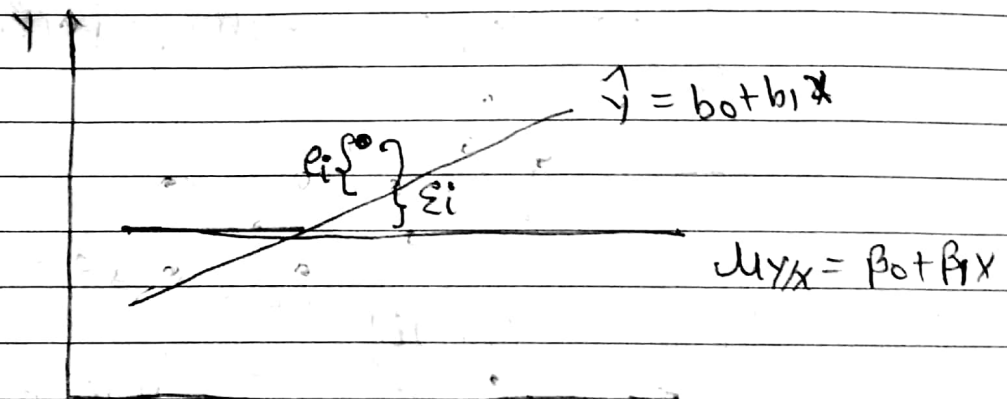
$$\hat{y} = b_0 + b_1 x$$

⑦  & Where $b_0$ is estimation of $\beta_0$ (i.e. $\beta_0$ estimates $b_0$)

$b_1$ is estimation of $\beta_1$ (i.e. $\beta_1$ estimates $b_1$).

{ where $\hat{y}$ represents predicted value of $y$ for a given value of $x$.

$b_0, b_1$ are the sample statistics that estimates the parameter $\beta_0$ & $\beta_1$ respectively

## Residuals :- ( error in a fit) :-

Given a set of regression data $\{ (x_i, y_i); i=1, \cdots n\}$ and fitted model, $\hat{y}_i = b_0 + b_1 x_i$, then the $i^{th}$ residual $e_i$ is given by

$$e_i = y_i - \hat{y}_i , \quad i = 1, 2, \cdots, n$$



$\hat{y} = b_0 + b_1 x$

$e_i \{ \} \} \varepsilon_i$

$\mu_{Y/x} = \beta_0 + \beta_1 x$
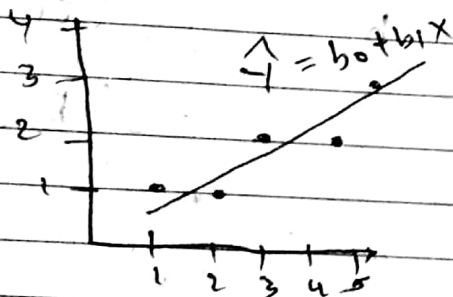
Comparing $\varepsilon_i$ with the residual $e_i$ :-

⑥ 1. **Gauss Markov assumption:**
$\varepsilon_i$ is a random variable with zero mean & variance $\sigma^2$ (which is unknown)
i.e. $E[\varepsilon_i] = 0$ & $Var[\varepsilon_i] = \sigma^2$

2. $\varepsilon_i$ & $\varepsilon_j$ are uncorrelated $i \neq j$, so covariance b/w $\varepsilon_i$ & $\varepsilon_j$ is $0$. i.e. $Cov(\varepsilon_i, \varepsilon_j) = 0$

3. $\varepsilon_i$ is a normally distributed random variable with mean zero and variance $\sigma^2$.
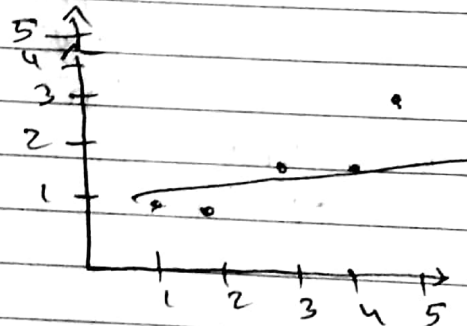i.e. $\varepsilon_i \sim N(0, \sigma^2)$.

* Since $\beta_0$ & $\beta_1$ are the unknown parameters so we want to estimate $\beta_0, \beta_1$. for this

**Least squares estimations of the parameters:-**
The parameter $\beta_0$ & $\beta_1$ are unknown and must be estimated using the data $(x_1, y_1), \ldots, (x_n, y_n)$



models 1.                    Model.

L.S.E it fits the model s.t $\sum_{i=1}^{5} e_i^2$ is minimum.

i.e. Sum of squares of the residuals is a minimum.
Note:- i.e. residuals sum of squares is often called the sum of squares of errors about the regression line and is denoted by S.S.E.

i.e. $SSE = \sum_{i=1}^{n} e_i^2$

DECEMBER 2011

| Wk | M | T | W | T | F | S | S |
|---|---|---|---|---|---|---|---|
| 49 | | | | 1 | 2 | 3 | 4 |
| 50 | 5 | 6 | 7 | 8 | 9 | 10 | 11 |
| 51 | 12 | 13 | 14 | 15 | 16 | 17 | 18 |
| 52 | 19 | 20 | 21 | 22 | 23 | 24 | 25 |
| 53 | 26 | 27 | 28 | 29 | 30 | 31 | |

NOTES :

Ⓧ

These two equation are called normal curve

$$SSE = \sum_{i=1}^{n} (e_i^2) = \sum_{i=1}^{n} (y_i - \hat{y}_i) = \sum_{i=1}^{n} (y_i - b_0 - b_1 x_i)^2$$

Differentiate SSE w.r.t. $b_0$ & $b_1$

we have:

$$\frac{\partial(SSE)}{\partial b_0} = -2 \sum_{i=1}^{n} (y_i - b_0 - b_1 x_i) = 0$$

$$\frac{\partial(SSE)}{\partial b_1} = -2 \sum_{i=1}^{n} (y_i - b_0 - b_1 x_i) \cdot x_i = 0$$

i.e. $b_0$ & $b_1$ are the sol. of the normal curve

$$\sum_{i=1}^{n} (y_i - b_0 - b_1 x_i) = 0 \quad \text{①}$$

$$\sum_{i=1}^{n} (y_i - b_0 - b_1 x_i) x_i = 0 \quad \text{②}$$

$$\frac{\sum_{i=1}^{n} y_i}{n} - b_0 \frac{\sum_{i=1}^{n} 1}{n} - b_1 \frac{\sum_{i=1}^{n} x_i}{n} = 0$$

$$\bar{y} = b_0 \frac{n}{n} - b_1 \bar{x} = 0$$

where $\bar{x} = \dfrac{\sum_{i=1}^{n} x_i}{n}$

$$\bar{y} = \frac{\sum_{i=1}^{n} y_i}{n}$$

$$\bar{y} - b_0 - b_1 \bar{x} = 0$$

$$\boxed{b_0 = \bar{y} - b_1 \bar{x}}$$

put the value of $b_0$ in ②

$$\sum_{i=1}^{n} (y_i - \bar{y} + b_1 \bar{x} - b_1 x_i) x_i = 0$$

$$\sum_{i=1}^{n} (y_i - \bar{y}) x_i = \sum_{i=1}^{n} b_1 (x_i - \bar{x}) x_i$$

$$b_1 = \left. \frac{\sum_{i=1}^{n} (y_i - \bar{y}) x_i}{\sum_{i=1}^{n} (x_i - \bar{x}) x_i} \right\}$$

$$\boxed{b_1 = \frac{\sum_{i=1}^{n} (y_i - \bar{y})(x_i - \bar{x})}{\sum_{i=1}^{n} (x_i - \bar{x})^2}}$$

$$\left( \because \sum_{i=1}^{n} (y_i - \bar{y}) \bar{x} = \bar{x} \sum_{i=1}^{n} y_i - \bar{y} \bar{x} \sum_{i=1}^{n} 1 \right.$$

$$= \bar{x} \cdot n \bar{y} - \bar{x} \bar{y} \cdot n$$

$$= 0$$

VEMBER 2011

| T | W | T | F | S | S | Wk |
|---|---|---|---|---|---|----|
| 1 | 2 | 3 | 4 | 5 | 6 | 45 |
| 8 | 9 | 10 | 11 | 12 | 13 | 46 |
| 5 | 16 | 17 | 18 | 19 | 20 | 47 |
| 2 | 23 | 24 | 25 | 26 | 27 | 48 |

NOTES :

(8) Estimating the Regression Coefficients

Given the sample $\{(x_i, y_i); i = 1, 2, \ldots, n\}$ the least squares estimates $b_0$ and $b_1$ of the regression coefficient $b_0$, $b_1$ are computed from the formulas:

$$b_1 = \frac{n \sum_{i=1}^{n} x_i y_i - \left(\sum_{i=1}^{n} x_i\right)\left(\sum_{i=1}^{n} y_i\right)}{n \sum_{i=1}^{n} x_i - \sum_{i=1}^{n} x_i}$$

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

$$b_0 = \frac{\sum_{i=1}^{n} y_i - b_1 \sum_{i=1}^{n} x_i}{n} = \bar{y} - b_1 \bar{x}$$

$$\Rightarrow b_1 = \frac{\sum_{i=1}^{n} x_i y_i - \bar{y} \sum_{i=1}^{n} x_i}{\sum_{i=1}^{n} x_i^2 - \bar{x} \sum_{i=1}^{n} x_i}$$

$$b_1 = \frac{\sum_{i=1}^{n} x_i y_i - n \bar{y} \bar{x}}{\sum_{i=1}^{n} x_i^2 - n \bar{x}^2}$$

Start Lec-2:

Parameter estimation solution Table

| $x_i$ | $y_i$ | $x_i^2$ | $y_i^2$ | $x_i y_i$ |
|-------|-------|---------|---------|-----------|
| 1 | 1 | 1 | 1 | 1 |
| 2 | 1 | 4 | 1 | 2 |
| 3 | 2 | 9 | 4 | 6 |
| 4 | 2 | 16 | 4 | 8 |
| 5 | 4 | 25 | 16 | 20 |
| Sum 15 | 10 | 55 | 26 | 37 |

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x} = 2 - (0.70)(3)$$
$$= -0.10$$

$$\bar{Y} = \sum_{i=1}^{5} Y_i = \frac{10}{5} = 2$$

$$\bar{X} = \sum_{i=1}^{5} X_i = \frac{15}{5} = 3$$

$$\hat{\beta}_1 = \frac{\sum_{i=1}^{n} x_i y_i - n\bar{x}\bar{y}}{\sum_{i=1}^{n} x_i^2 - n\bar{x}^2} = \frac{37 - 5 \times 3 \times 2}{55 - 5 \times 3^2} = 0.70$$

So fitted eqn $\boxed{\hat{Y} = -0.10 + 0.7X}$

Correlation officient:- The measure $\rho$ of linear association b/w two variable $x$ & $y$ is estimated by sample correlation cofficient $\gamma$

where

$$\gamma = b_1 \sqrt{\frac{SS_{xx}}{S_{yy}}} = \frac{S_{xy}}{\sqrt{S_{xx} S_{yy}}}$$

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

$$S_{xy} = \sum_{i=1}^{n} (x_i - \bar{x})(y_i - \bar{y}) = cov(x,y)$$

$$S_{yy} = \sum_{i=1}^{n} (y_i - \bar{x})^2 = var(y)$$

**Q. What is correlation ?**

**A.** Correlation is the statistical technique which used for knowing How strong pairs of variables are related.

Ex:- Height and weight are related. i.e. taller people is heavier than the shorter people. this relation is not perfect. because people of same height may very. So shorter one may have heavier than the taller.

**Correlation Cofficient :-** The main result of the correlation is called the correlation cofficient (or $r$) It's range is from $-1.0$ to $1.0$. & closer $r$ is to $+1$ or $-1$. In statistics, the correlation cofficient '$r$' measures the strength and direction of a linear relationship b/w two variables on a scatterplot.

If $r = -1$ means A perfect downhill (negative) linear relationship

$r = -0.7$ " A strong " " " "

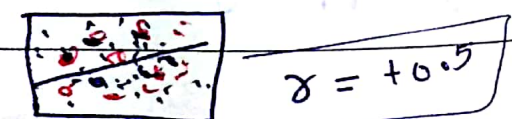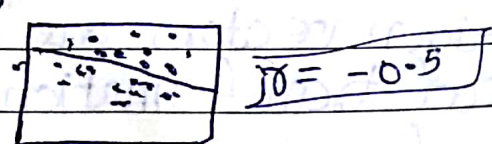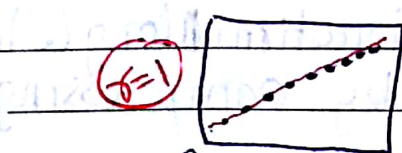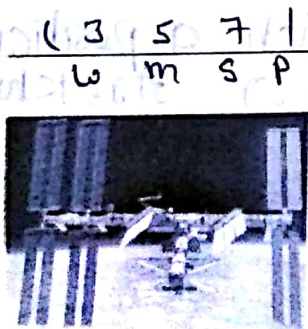$r = -0.5$ " " moderate " " " "

$r = -0.3$ " " weak " " " "

$r = 0$ " " No relationship

$r = +0.3$ " A weak uphill (positive) linear relationship

$r = +0.5$ " A moderate " " " "

$r = +0.7$ " " strong " " " "

$r = $ Exactly $+1$ " " perfect " " " "

1.3   5   7  1
w    m   s   p



$r = 1$

$r = -0.5$

$r = +0.75$

$r = +0.5$