## Assignment NO:- 07

**Title :—**

Correlation & linear Regression in R.

**Problem statement:—** Use of R for correlation & regression analysis.

**Pre-Lab:—** A basic understanding of the correlation & regression concept is required.

## Theory:—

**Linear Regression: —**

In data analysis we come across the term "Regression" very frequently e.g. if we say that.
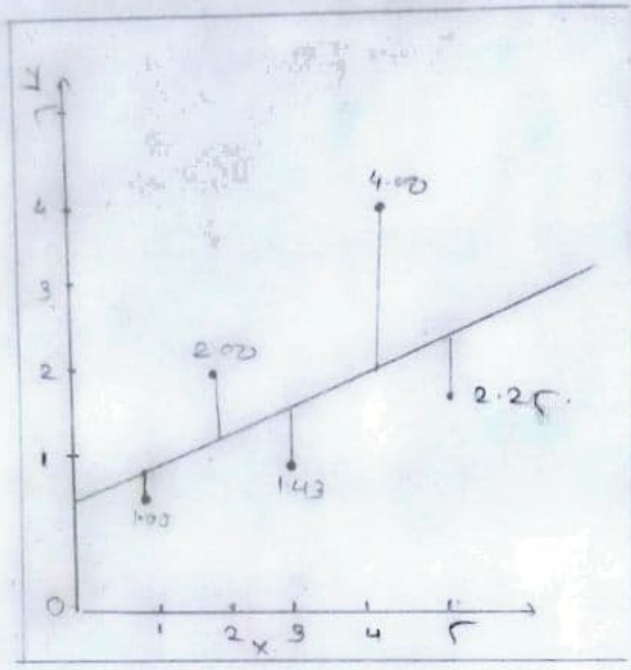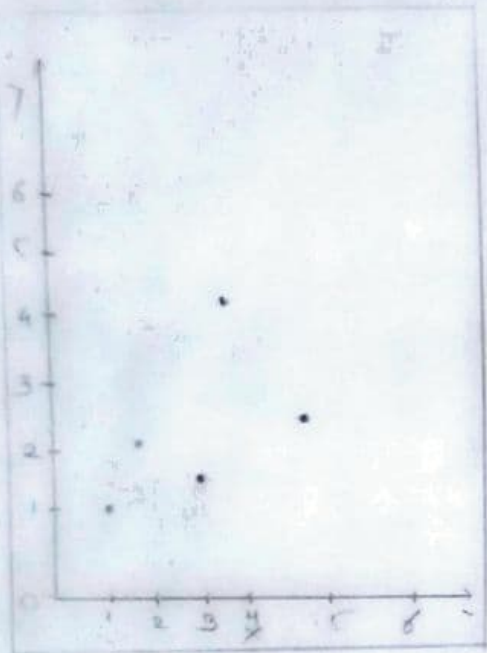
$$Age = 5 + Height * 10 + weight * 13$$

**simple Linear Regression :—**

"Linear Regression" is a statistical method to regress the data with depend variable having contineous values. E.g. Predicting traffic in a retail store, predicting a user's dwell time or number of pages visited on Dezyre.com etc.

**Prerequisites:—**

- correlation (r)- Explains the relationship betn two variables. posible values -1 to +1.
- variance $(\sigma^2)$ - Measure of spread in your data.
- standard deviation $(\sigma)$ - Measure of spread in your data (square root of variance).
- Normal distribution.
- Residual (error term) - { Actual value - Predicted value}

## Assumption of Linear Regression :—

i) Linearity &
Additive :- There should be a linear relationship bet$^n$ dependent & independent variables & the impact of change in independent variable values should have additive impact on dependent variable.

ii) Normality of error distribution :- Distribution of differences bet$^n$ actual & predicted values (Residuals) should be normally distribution.

iii) Homoscedastisity :-

      a) Time

      b) The predictions.

      c) Independent variable values.

iv) Statistical independance of errors :- The error terms (residuals) should not have any correlation among themselves.

## Linear Regression Line :—

While doing linear regression our objective is to fit a line through the distribution which is nearest to most of the points. Hence reducing the distance (error term) of data points from the fitted line.

following. equations.

$$Y = B_0 + B_1 X .$$

where,

    $Y$ = Dependent variable.

    $X$ = Independent variable.

    $B_0$ = Constant term / Intercept.

    $B_1$ = coefficient of relationship bet$^n$ 'x' & 'y'

Few properties of linear regression line :-

- Regression line always passes through mean of independent variable ($x$) as well as mean of dependent variable ($y$).
- Regression line, minimizes the sum of "square of Residuals".
- $B_1$ explains the change in $y$ with a change in $x$ by one unit.

● Finding a linear Regression Line :—

Using a statistical tool e.g. Excel, R, SAS etc. For example, let say we want to predict '$y$' from '$x$' given in following table & let's assume that our regression eqⁿ will look like
"$y = B_0 + B_1 * x$".

| $x$ | $y$ | Predicted '$y$'. |
|-----|-----|------------------|
| 1 | 2 | $B_0 + B_1 * 1$ |
| 2 | 1 | $B_0 + B_1 * 2$ |
| 3 | 3 | $B_0 + B_1 * 3$ |
| 4 | 6 | $B_0 + B_1 * 4$ |
| 5 | 9 | $B_0 + B_1 * 5$ |
| 6 | 11 | $B_0 + B_1 * 6$ |
| 7 | 13 | $B_0 + B_1 * 7$ |
| 8 | 15 | $B_0 + B_1 * 8$ |
| 9 | 17 | $B_0 + B_1 * 9$ |
| 10 | 20 | $B_0 + B_1 * 10$ |

where,
Table 1,

| std. Dev of x | 3.02765 |
|---|---|
| std. Dev. of y | 6.617317 |
| Mean of x | 5.5 |
| Mean of y | 9.7 |
| correlation between x & y | .989938 |

$B_1$ = correlation * (std. Dev. of y / std. Dev. of x )

$B_0$ = Mean (y) - $B_1$ * Mean(x).

Putting values from table 1 into the above eqns.

$B_1 = 2.64$

$B_2 = -2.2$.

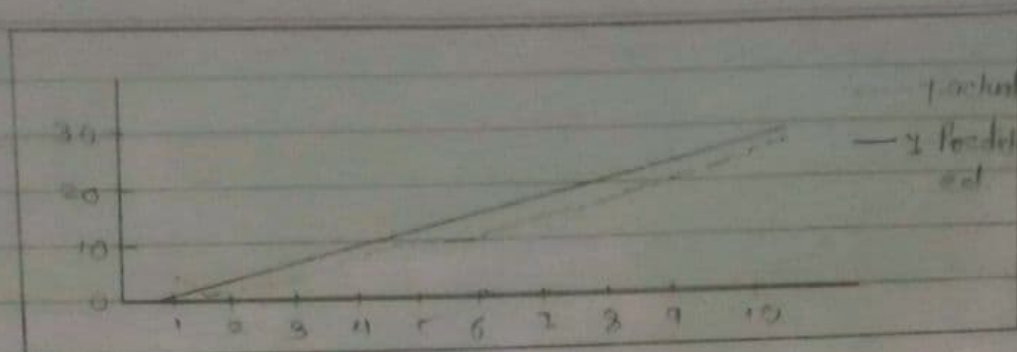Hence, the least regression eqn will become -

y = -2.2 + 2.64 * x.

let's see,

| x | y - Actual | y - Predicted. |
|---|---|---|
| 1 | 2 | 0.44 |
| 2 | 1 | 3.08 |
| 3 | 3 | 5.72 |
| 4 | 6 | 8.36 |
| 5 | 9 | 11 |
| 6 | 11 | 13.64 |
| 7 | 13 | 16.28 |
| 8 | 15 | 18.92 |
| 9 | 17 | 21.56 |
| 10 | 20 | 24.2 |

Linear Regression in R using lm() function :—
                                        It is easiest
way to find regression using lm() fun.
The syntax is:
        lm (formula, data).
- formula is a symbol presenting the relon betn $x$ & $y$.
- data is the vector on which the formula will be applied.

predict () function:
        The basic syntax for predict () in
in linear regression is —
predict(object, newdata).
- object is formula which is already created using
  lm() fun.
- newdata is the vector containing the new data value
  for predictor variable.

Multiple Regression :—
        The general mathematical eqn
for multiple regression is —
$y = a + b_1 \times x_1 + b_2 \times x_2 + \dots b_n x_n$.
following is discription of parameters used —
- $y$ is the response variable.
- $a, b_1, b_2, \dots$ the the coefficients.
- $x_1, x_2, \dots x_n$ are predictor variables.

we create the regression model using the lm() function in R.

The lm() fun creates the relationship model bet^n the predictors & the response variable.

The basic syntax for lm() fun in multiple regression is-

$lm(y \sim x_1 + x_2 + x_3 \ldots, data)$.

- formula is a symbol presenting the relan bet^n the response variable & predictors variables.
- data is the vector on which the formula will be applied

Create Equation for Regression Model :-

Based on the above intercept & coefficient values. we create the mathematical eqn.

Apply Equation for predicting New values :-

we can use the regression eqn created above to predict the new value of dependent variable for the given set of independent variables.

Logistic Regression :-

The general mathematical eqn for logistic regression is

$y = 1/(1 + e^{-(a+b_1 x_1 + b_2 x_2 + b_3 x_3 + \cdots)})$

- y is a responce variable -
- x is a predictor variable.
- a & b are the coefficient which are numeric constant

The fun used to create the regression model is a glm fun.

The basic syntax for glm() fun in logistic regression is -

glm(formula, data, family)

Following is the description of the parameters used -
- formula is the symbol presenting the relationship betn the variables.
- data is the data set giving the values of those variables.
- family is R object to specify the details of the model. It's value is binomial for logistic regression.

Post-Lab :-
            students will be able to find relation betn dependent & independent variables using training dataset & can predict values from the new dataset given.

Conclusion :-
            Thus exercised various commands relate to linear Regression in R.