

## Assignment 3 II

Date: / / 20

Title - Correlation and Linear Regression in R

Problem Statement - Use of R for correlation & regression analysis.

Prelab - A basic understanding of correlation & regression concepts

Theory :-

Linear Regression -

In data analytics we come across the term "Regression" very frequently. Regression is a statistical way to establish a relationship between a dependent variable & a set of independent variable(s) e.g. if we say that,  
 $\text{Age} = 5 + \text{Height} * 10 + \text{Weight} * 13.$

Simple Linear Regression -

"Linear Regression" is a statistical method to regress having continuous values whereas independent variables can have their continuous or categorical values. In other words "Linear regression" is a method to predict dependent variable (Y) based on values of independent variable.

e.g. Predicting traffic in retail store, predicting users dwell time or number of pages visited

Prerequisites -

To start with Linear Regression, few basic concepts are required.

Pooja



- Correlation ( $r$ ) - Explain the relationship between two variables, possible 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 -

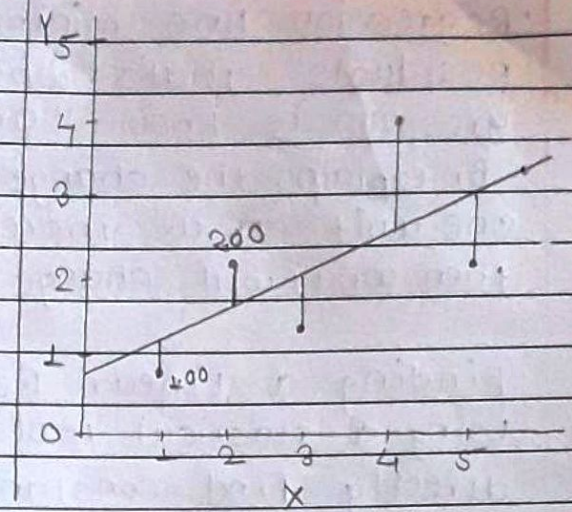
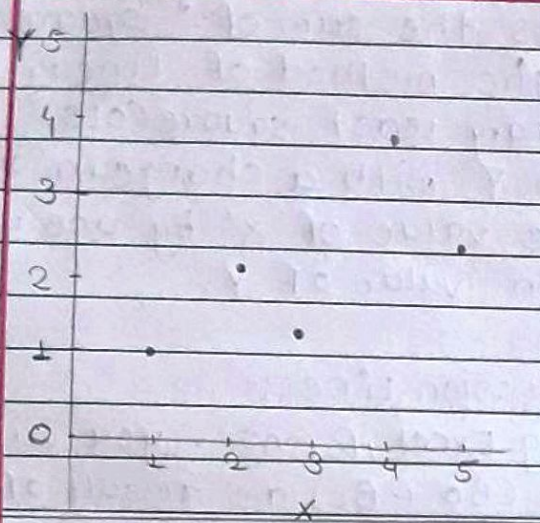
Not a single size fits for all, the same is true for Linear Regression as well as misleading

- i) Linearity & Additive: there should be linear relationship between dependent & independent variable & the impact of change in independent variable should have additive impact on dependent variable.
- ii) Normality of error distribution. Distribution of differences between actual & predicted values.
- iii) Homoscedasticity: variance of errors should be constant versus
  - a) Time
  - b) The predictions
  - c) Independent variable values.
- iv) Statistical independence of errors: the error terms 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 fitted line





For examples in above figure (left) dots represent various data points & line (right) represent an approximate line which can explain relationship between 'x' & 'y' axes. Through, linear regression we try to find out search a line. For e.g if we have dependent variable 'y' and one independent variable 'x' - relation between 'x' & 'y' can be represented in a form of foll<sup>n</sup>.

$$Y = B_0 + B_1 X$$

- where
- Y = Dependent variable
  - X = Independent variable.
  - $B_0$  = Constant term / Intercept
  - $B_1$  = Coefficient of relationship between '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". That's why the method of Linear Regression is known "Ordinary least square (OLS)".
- $B_1$  explains the change in  $Y$  with a change in  $X$  by one unit. If we increase value of ' $x$ ' by one unit then what will change in value of  $Y$ .

Finding a Linear Regression Line :-

Using a statistical tool e.g. Excel, R, SAS. you will directly find constants ( $B_0$  &  $B_1$ ) as result of linear regression function. But conceptually as discussed it works on OLS concept & tries to reduce square root of errors, using very concept software package.

$x$	$y$	Predicted ' $y$ '
1	2	$B_0 + B_1 * 1$
2	4	$B_0 + B_1 * 2$
3	6	$B_0 + B_1 * 3$
4	8	$B_0 + B_1 * 4$
5	10	$B_0 + B_1 * 5$
6	12	$B_0 + B_1 * 6$
7	14	$B_0 + B_1 * 7$
8	16	$B_0 + B_1 * 8$
9	18	$B_0 + B_1 * 9$
10	20	$B_0 + B_1 * 10$



Table 1:

std. dev of x	3.02765
std. dev. of y	6.6137317
Mean of x	5.5
Mean of y	9.7
Correlation bet <sup>n</sup> x & y	.989938

If we differentiate the Residual sum of square (RSS) wrt  $B_0$  &  $B_1$  & equate results into zero

$$B_1 = \text{Correlation} * (\text{std. dev of } y / \text{std. dev. of } x)$$

$$B_0 = \text{Mean } (Y) - B_1 * \text{Mean } (X)$$

Putting values from table 1

$$B_1 = 2.64$$

$$B_0 = -2.2$$

$\therefore$  least regression equation will become

$$Y = -2.2 + 2.64 * x$$

how our prediction are looking like equation

X	Y-Actual	Y-Predicted
1	2	0.44
2	4	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

Give only 10 data points to fit a line our prediction are not pretty accurate but if see correlation bet<sup>n</sup> 'Y-Actual' & 'Y-Predicted' it will turn out very high

Pooja



## Linear Regression in R using `lm()` function.

It is easiest way to find regression using `lm()`

The syntax is

`lm(formula, data)`

- formula is a symbol presenting the relation between  $x$  &  $y$
- data is the vector on which formula will applied.

### \* `predict()` Function :-

The basic syntax for `predict()` in linear regression  
`predict(object, newdata)`

- object is formula which is already created using the `lm()` function
- newdata is the vector containing the new value for predictor variable

This function will be used to predict the new value of dependent variable using the new dataset & values found using `lm()` function.

## Multiple Regression.

Multiple regression is an extension of linear regression into relationship between more than two variables. In simple linear relation we have one predictor & one response variable, but in multiple regression we have more than one predictor variable & one response variable

The general mathematical eq<sup>n</sup> for multiple regression

$$y = a + b_1x_1 + b_2x_2$$

- $y$  is response variable.
- $a, b_1, b_2, \dots, b_n$  are coefficients.



We create the regression model using the `lm()` function in R. The model determines the value of coefficients using the data. We can predict the value of response for given set of predictor variables using these coefficient.

The `lm()` function create the relation model between the predictor & the response variable.

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

`lm (y ~ x1 + x2 + x3 ----, data)`

- Formula is a symbol presenting the relation between the response variable & predictor.
- 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 equation.

Apply Equation for predicting new values -

We can use regression equation created above to predict the new values of dependent variables for the given set of independent variables.

Logistic Regression:-

The logistic Regression is a regression model in which response variable has categorical values such as TRUE / FALSE or 0 or 1. It actually measure the probability of a binary response as value of

Pooja



based on mathematical equation.

The general mathematical eq<sup>n</sup> for logistic regression is -

$$y = 1 / (1 + e^{-(a + b_1x_1 + b_2x_2 + b_3x_3 + \dots)})$$

fol<sup>l</sup> is description of parameter used -

- $y$  is the response variable.
- $x$  is the predictor variable.
- $a$  &  $b$  are the coefficient which are numeric constants.

The basic syntax for `glm()` function is logistic regression  
`glm(formula, data, family)`

- formula is symbol presenting the relationship between the variables
- data is data set giving the values of these variables.
- family is R object to specify the details of model. It's value is binomial for logistic regression.

Post-Lab :- Students will be able to find relation bet<sup>n</sup> dependent & independent variables using training dataset & can predict value for new dataset given.

Conclusion- Thus exercised various commands related to linear regression in R.