

⑦ Assignment 3 II

Date: / / 20

Title :- Correlation & Linear Regression in R

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

Theory:-

Linear regression -

In data analysis we come across the term "Regression" very frequently. Regression is a statistical way to establish a relationship b/w a dependent variable + a set of independent variable (s) e.g if we say Age = 5 + height * 10 + weight * L3.

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) base 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.

Give only 10 data points to fit a line our prediction are not pretty accurate but if we see correlation bet' 'y-actual' & 'y-Predicted' it will turn out very high.

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 relation bet' 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 formula which is containing new value for predictor variable.

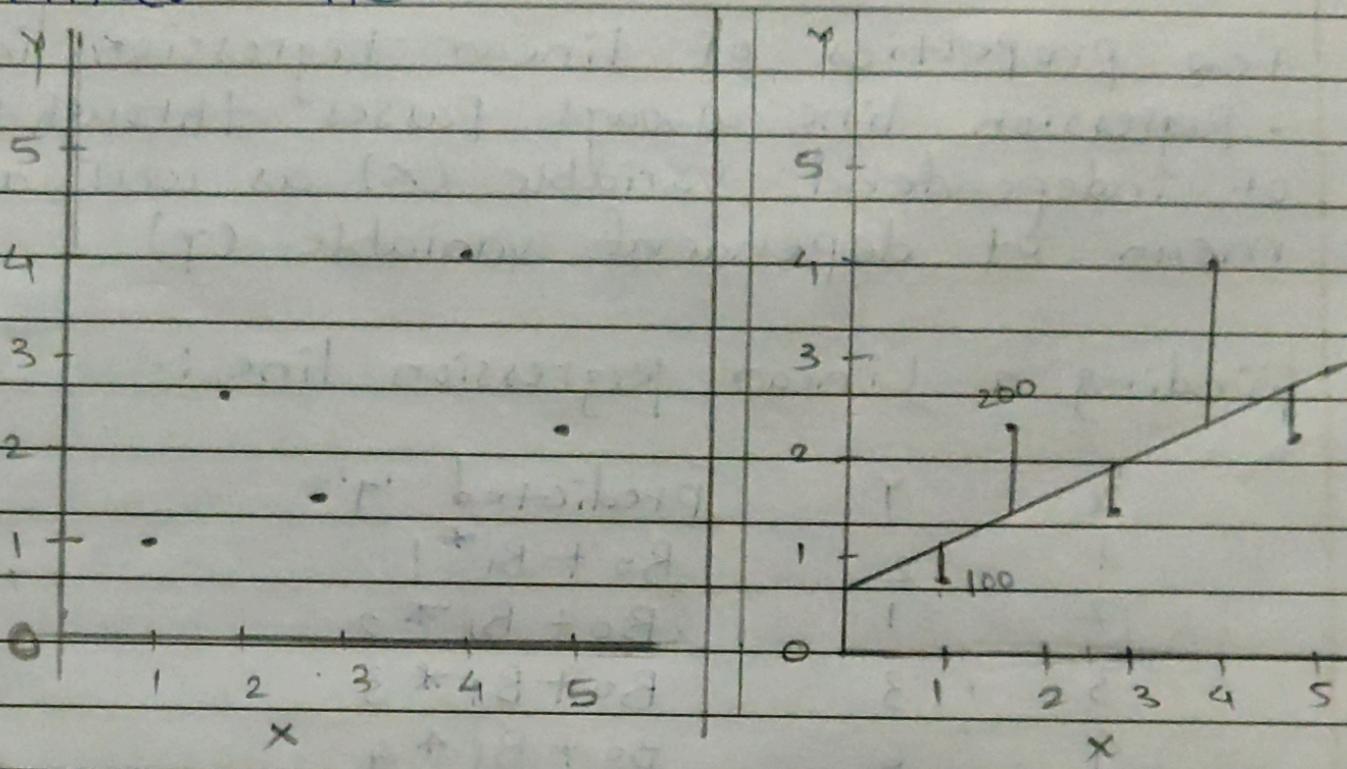
Value found using lm() function.

Multiple Regression -

multiple regression is an extension of linear regression into relationship bet' more than two variables. In Simple linear relation we have one predictor & one

Linear regression line -

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



For example in above figure (left) dots represent various data points & line (right) represent an approximate line which can explain relationship bet 'x' & 'y' arises. Through linear regression we try to find out search a line.

for. eg if we have re dependent variable 'y' & on independent variable 'x'. relation bet 'x' & 'y' can be represented in a form of following:

$$Y = B_0 + B_1 X$$

Where Y = Dependent variable

X = Independent variable.

B_0 = constant term / Intercept

B_1 = coefficient of relationship
bet '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)

Finding a Linear Regression line :-

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	$B_0 + B_1 * 10$

Table 1 :

Std. Dev of \bar{x} 3.02765

Std. Dev of \bar{y} 6.6137317

Mean of \bar{x} 5.5

Mean of \bar{y} 9.1

Correlation betw $x \& y$.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

$$B_1 = 2.64$$

$$B_0 = -2.2$$

∴ least regression eqn will be.

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

now our prediction are looking like equation

X	y -actual	y -predicted
1	2	0.44
2	3	3.08
3	6	5.72
4	9	8.36
5	11	
6	13	13.64
7	15	16.28
8	17	18.92
9	20	21.56
10		24.2

- Correlation (r) - Explain relationship between two variables, possible values $-1 \text{ to } +1$
- Variance (s^2) - Measure of spread in your data
- Standard deviation (s) - 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. Should have additive impact on dependent variable.
- ii) Normality of error distribution. Distribution or 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.

one response variable.

The general mathematical eqⁿ for multiple regression

$$y = a + b_1 x_1 + b_2 x_2$$

• y is response variable

• a, b₁, b₂... b_n are coefficients.

Create Equation for Regression model -

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

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 based on mathematical eqⁿ.

The general mathematical eqⁿ for logistic regression is -

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

formⁿ is description of parameter used.

• y is response variable

• x is 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 relationship betⁿ 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.

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