

Page No.

Title : Correlation & Linear Regression R.

Problem Statement : Use of R for Correlation & Regression analysis.

Pre-Lab : A basic understanding of the Correlation & regression concepts is required.

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). eg. if we say that $\text{Age} = 5 + \text{Height} * 10 + \text{Weight} * 13$ Here we are establishing a relationship between Height & Weight of a person with his/her Age.

Simple Linear Regression

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

Prerequisites

- To start with Linear Regression, few basic concepts of statistics are required:
- Correlation (r) - Explains the relationship between two variables, possible values -1 to +1.
 - Variance (σ^2) - Measure of spread in your data
 - Standard Deviation (σ) - Measure of spread in your data (Square root of variance)
 - Normal distribution
 - Residual (error term) - (Actual value - Predicted value)

Assumptions of Linear Regression:

Not a single size fits or all, the same is true for Linear Regression as well. In order to fit a linear regression line data should satisfy few basic but important assumptions. If your data doesn't follow the assumptions, your results may be wrong as well as misleading.

- Linearity & Additive:** There should be a linear relationship between dependent & independent variables & the impact of change in independent variable value should have establish additive impact on dependent variable.
- Normality of error distribution:** Distribution of differences between Actual & predicted values (Residuals) should be normally distributed.

iii. Homoscedasticity: Variance of errors should be constant versus,

a. Time

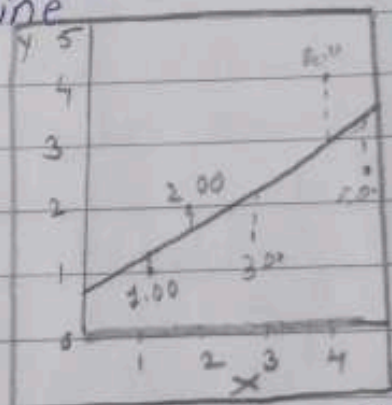
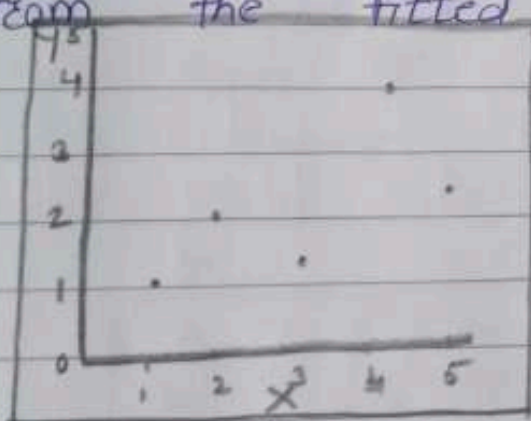
b. The predictions

c. Independent variable values.

iv. Statistical independent 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 of data points from the fitted line.



For example, in above figure dots represent various data points & line represents an approximate line which can explain the relationship between 'x' & 'y' axes. If we have one dependent variable 'y' & one independent variable 'x' - relationship between 'x' & 'y' can be represented in a form of following equation: $Y = B_0 + B_1 X$

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 as "Ordinary Least Square (OLS)".

Finding Linear Regression Line.

Using a Statistical tool e.g., Excel, R, SAS you will directly find constants as a result of linear Regression function.

For example, let us we want to predict 'y' from 'x' given in following table & let's assume that our Regression equation 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	$B_0 + B_1 * 10$

Table 1 :

Std. Dev. of x	
Std. Dev. of y	3.02765
Mean of x	6.617317
Mean of y	8.5
Correlation between x & y	9.7
	.989938

If we differentiate the Residual Sum of Square (RSS) w.r.t. B_0 & B_1 & equate the results to Zero, we get the following equations as a result:

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

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

Putting values from table 1 into the above equation

$$B_1 = 2.64$$

$$B_0 = -2.2$$

Hence, the least Regression equation will become — $Y = -2.2 + 2.64 * X$

Let see, how our prediction are looking like using this equation

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 the easiest way to find regression using `lm()` function.

The syntax is:

`lm(formula, data)`

Following is the description of the parameters used -

- `formula` is a symbol presenting the relationship between x & y

- `data` is vector on which the formula will be applied

`predict` function :

The syntax is:

`predict(object, newdata)`

Following is the description of the parameters used -

- `Object` is the formula which is already created using the `lm()` function

- `newdata` is the vector containing the new value for predictor variable.

Create Equation for Regression Model
Based on above intercept & coefficient values, we create the mathematical equation

Apply Equation for predicting New values:
We can use the regression equation created above to predict the new values of dependent variables for the given set.

Logistic Regression:

The Logistic Regression is Regression Model in which the response variable (dependent variable) has categorised value such as True/false or 0/1.

The general mathematical equation for logistic regression is -

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

`glm()` function.

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

`glm(formula, data, family)`.

Post-Lab: Students will be able to find relation between dependent & independent variables using training dataset and can predict values for the new dataset given.

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