

Lesson 3:- Linear RegressionD The Business Problem

Ques:- What decision we need to be made?

What information is need to inform that decision?

What type of analysis is need to get the information need to make that decision?

How do we supply enough resources we need?

How much ticket we need per week?

We need ~~predic~~ predictive analysis to help us obtain the data to inform the decision we need to make.

D Approach the Business Problem :-

What decision need to be make?

Do we have enough capacity on the support team to handle the support tickets from the new customer?

And if not, how many people do we need to add to the support team to reach the desired capacity.

What information do we need to inform this decision?

We need to calculate the average number of tickets per customer per week. We can then aggregate the average no. of tickets for each customer to get

Data
 mod.
 Vaex

A ~~#~~ total average no. of support tickets that we predict will be submitted per week. Once we have this information, we need to compare the predicted average no. of tickets with the current capacity of the support staff, specifically, the average no. of tickets team member can handle.

P Data Understanding Ques :-

Linear Example Data :-

Client ID	Average No. of tickets	No. of employees	Value contract	Industry
TK669	90	561	200000	Retail

Q. Which methodology would best assist in solving the business problem?

Continuous Numeric Model.

Q. Describe the steps you took to determine the appropriate model.

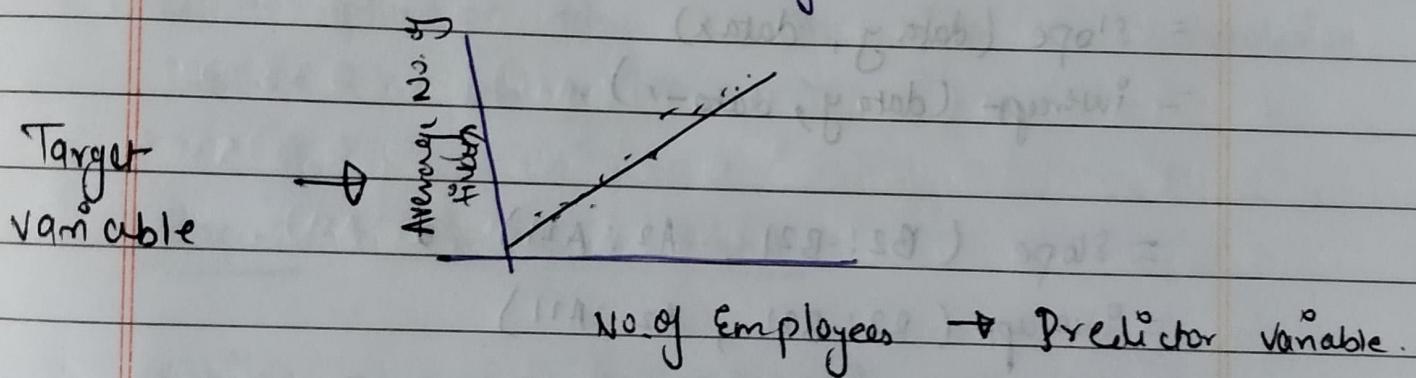
We need to calculate the Average no. of tickets per customer per week (which was given) → Data Rich →

We need numeric value (i.e. Numeric Analysis → Continuous

based Numeric model. (Because we are not trying to predict what will happen over time).

▷ Introduction to Linear Regression:-

Data displayed in the scatter plot, It appears that we have a linear relationship between the number of employees and the number of tickets.



$$y = mx + b.$$

Target Variable:- If is the variable we are trying to understand and predict. It is also referred to as dependent variable. for ex:- we are trying to predict Y (or the average no. of tickets).

Predictor Variable:- These are used to try to predict the target variable and are also known as independent Variable.

D Linear equation in Google Sheets :-

No. of Employees Average No. of tickets.

s_1

s_2

68

5

= slope (data_y, data_x)

= intercept (data_y, data_x).

= slope (B2:B21, A2:A21)

= intercept (B2:B21, A2:A21)

$$\hat{y} = 0.1833x - 11.055$$

D Linear Regression Validation :-

- Validation:-

Now that we've performed the analysis and run the linear regression model, we need to validate the results of the model. In other words, is there a way to measure how good the model is? or in this case, is the linear expression we calculated a

good fit of our data ?

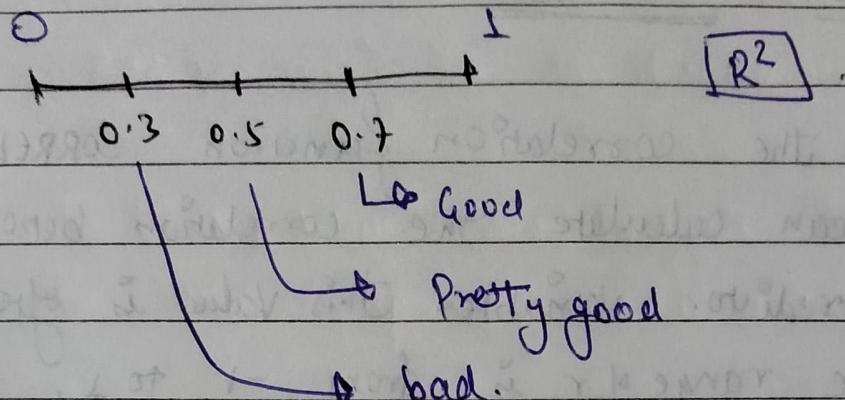
Step 1:- Correlation :-

Using the correlation function CORREL (datay, datax), we can calculate the correlation between the target and predictor variable. This value is often referred to as r . The range of r is from -1 to 1. The closer r is to plus or minus 1, the higher the correlation between x and y .

Step 2:- Calculate r-squared :-

While a strong correlation is good, we really want to know how well the data fits our line. Fortunately, we can get a sense of how good the formula is at approximating the data by calculating the coefficient of determination or r -squared: R-Squared is a coefficient between 0 and 1. R-Squared is interpreted as the percentage of variance in observation that is explained by the model, or the explanatory power of the model. An R-Squared value close to 1 would mean that nearly all variance in the target variable is explained by the model. An R-Squared value close to zero

would mean that nearly none of the variance in the target variable is explained by model.



$$\Rightarrow R^2 = \text{Correl}(B_2: B_{21}, A_2: A_{21}) = 0.987109.$$

$$\Rightarrow R^2 = \frac{C_5 * C_5}{\text{or} \quad = \text{RSG}(B_2: B_{21}, A_2: A_{21})} = 0.979384.$$

Simple linear Regression Ex:

$$\text{Slope} = 0.181736$$

$$\text{Intercept} = -7.54648$$

$$\text{Correlation} = 0.886683$$

$$R^2 = 0.786207$$

$$y = 0.181736 x - 7.54648$$

$$y(525) = 87.865 \approx 88$$

▷ Introduction to Multiple Linear Regression :-

Linear Regression:- $y = b + mx$.

$$y = b + mx \quad \begin{matrix} \downarrow \\ \text{Coefficient} \end{matrix}$$

Multiple linear Regression:- $y = b_0 + b_1 x_1 + b_2 x_2$

y = target variable.

b, b_0 = y -intercept

b, b_1, b_2, m = coefficient.

Initial Steps:- Preparing and Understanding your Data

Given any data set, we must make sure that the ~~given~~ data is not clean and not biased. Once we have clean data set, the next step to prepare for your multiple linear Regression is to understand the relationship between each of your predictor variable and your target variable.

Important!:- This is important because linear regression models assume that our numerical predictor variable should have a linear relationship with the target variable. It's good practice to analyze the individual variables first before you run your variable through the linear model.

Steps for Multiple linear Regression in Excel!-

Step 1:- Make sure you have the Analysis ToolPak Add-in active in Ms EXCEL.

Step 2:- In Excel, select Data Analysis. Select Regression in the pop-up window and select OK.

Step 3:- Input Y Range Should be the range of your target variable

Input X Range Should be the range of data of your predictor variable.

Step 4:- Click OK to see the Results.

D R Squared vs Adjusted R-squared:-

The adjusted r-squared value should be used with multiple linear regressions due to phenomenon that occurs when adding additional variables to the model. In a nutshell, the more variables that are included, the higher the r-squared value will be - even if there is no relationship between the additional variables and the target variable. Thus, we use the Adjusted R-squared value.