

Report On
Predict Profit Value of the Company

By,
Samarth Pandey

EXPOSYS DATA LABS
Data Science Internship

Table of Contents

1.) Abstract

2.) Introduction

3.) Comparison between Linear and Multiple Regression

4.) Proposed method with architecture

5.) Methodology

6.) Implementation

7.) Conclusion

1.) Abstract: -

A company should always set a goal that should be achievable, otherwise, employees will not be able to work to their best potential if they find that the goal set by the company is unachievable. The task of profit prediction for a particular period is the same as setting goals. If you know how much profit you can make with the amount of R&D and marketing you do, then a business can make more than the predicted profit provided the predicted value is achievable.

2.) Introduction: -

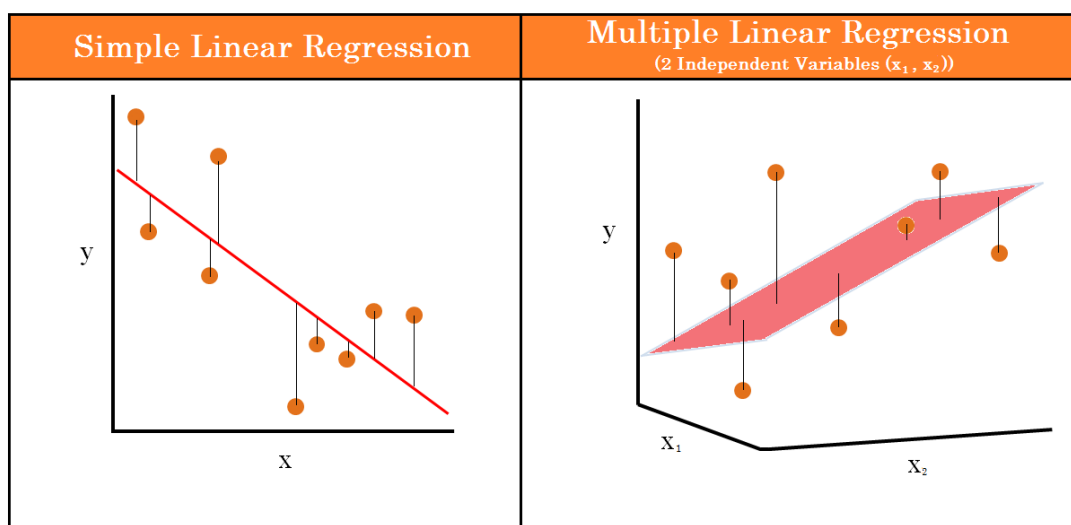
The profit earned by a company for a particular period depends on several factors like how much time and money a company spends on R&D, marketing and many more. So, for predicting the profit of a company for a particular period we need to train a machine learning model with a dataset that contains historical data about the profit generated by the company. The task of predicting profit is an important task for every business to set an achievable goal. For example, if the business spends \$500 on marketing, it can't expect a profit of \$20,000. Likewise, there are many other factors on which the profit of a business depends.

3.) Comparison between Linear and Multiple Regression: -

The major difference between linear regression and multiple linear regression is that in linear regression there is only one independent variable while when we check out, Multiple linear regression there is more than one independent variable. Let us take an example of both the scenarios

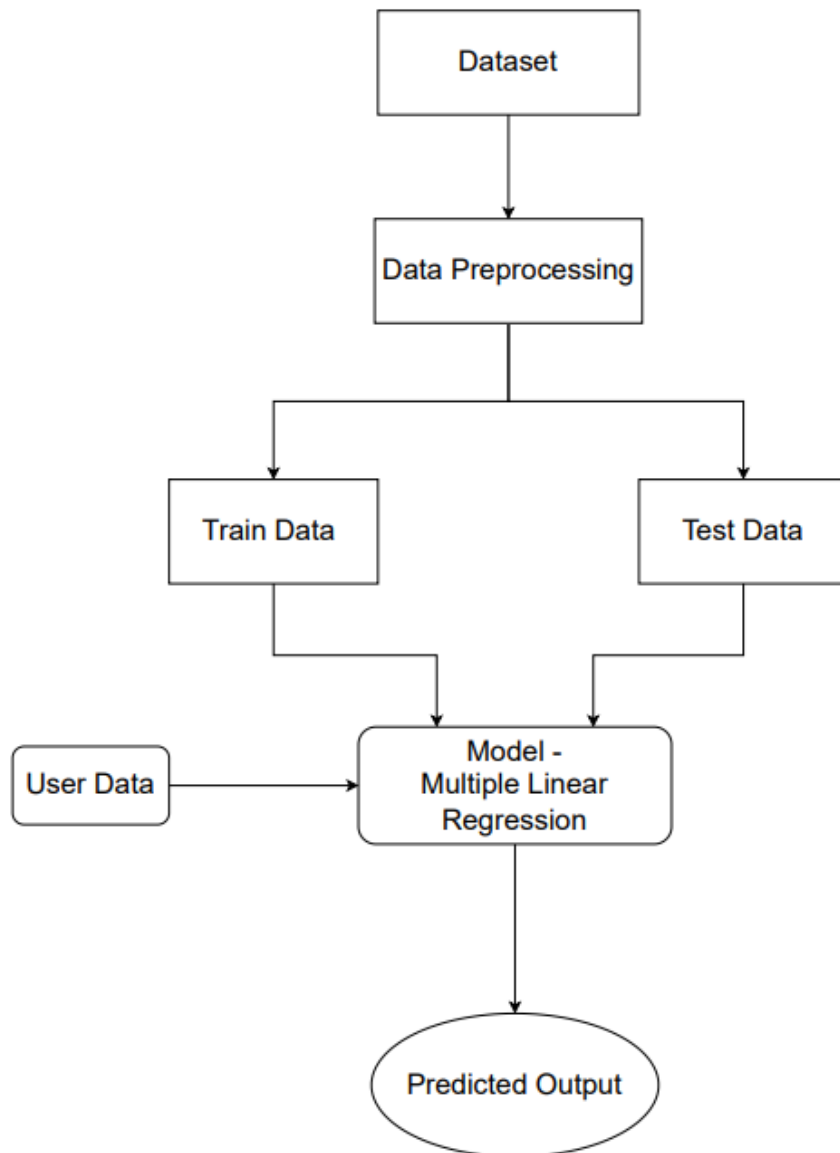
→Linear regression: When we want to predict the height of one person just from the weight of that person.

→Multiple Linear regression: If we alter the above problem statement just a little bit like, if we have the features like height, age, and gender of the person and we have to predict the weight of the person then we have to use the concept of multiple linear regression.

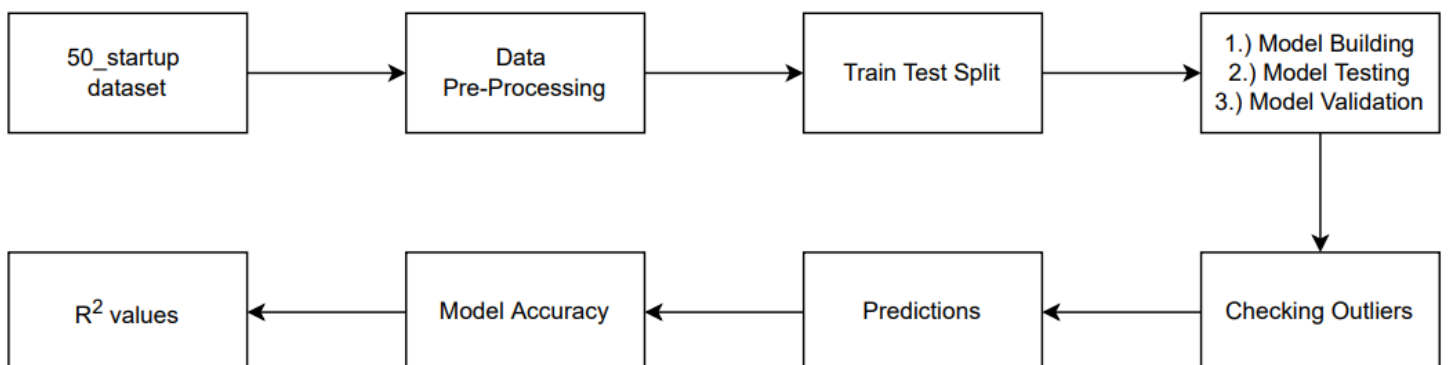


4.) Proposed Method with Architecture: -

The competition goal is to predict the profit of start-up profit on the bases of data provided which are on the bases of Research and Development Spend (R&D Spend), Administration Spend, Marketing Spend and State. We use multiple regression in this model because we have to predict profit (dependent variable) on bases of multiple field(independent variables) rather than one field just like we done in Simple Linear Regression. This model can help those people who want to invest in start-up company by analysing profit of the company.

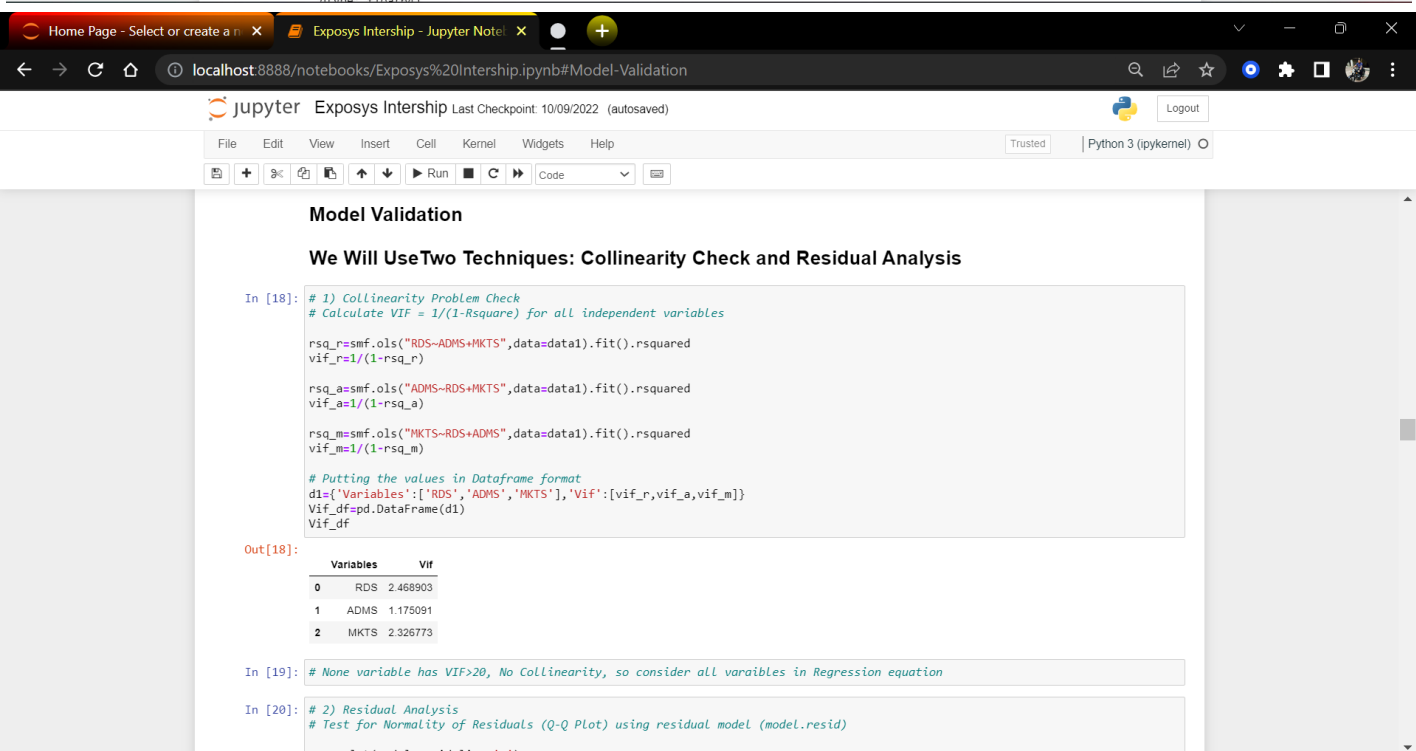
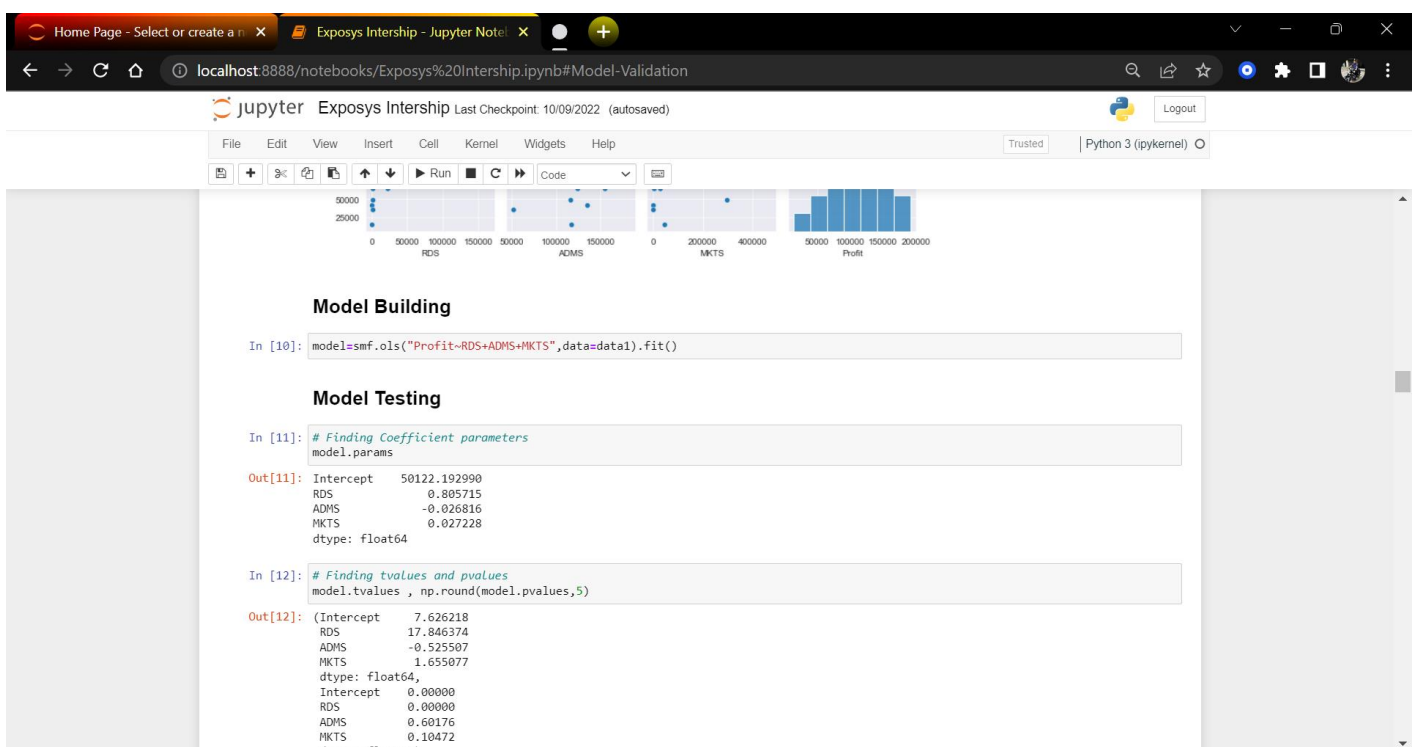


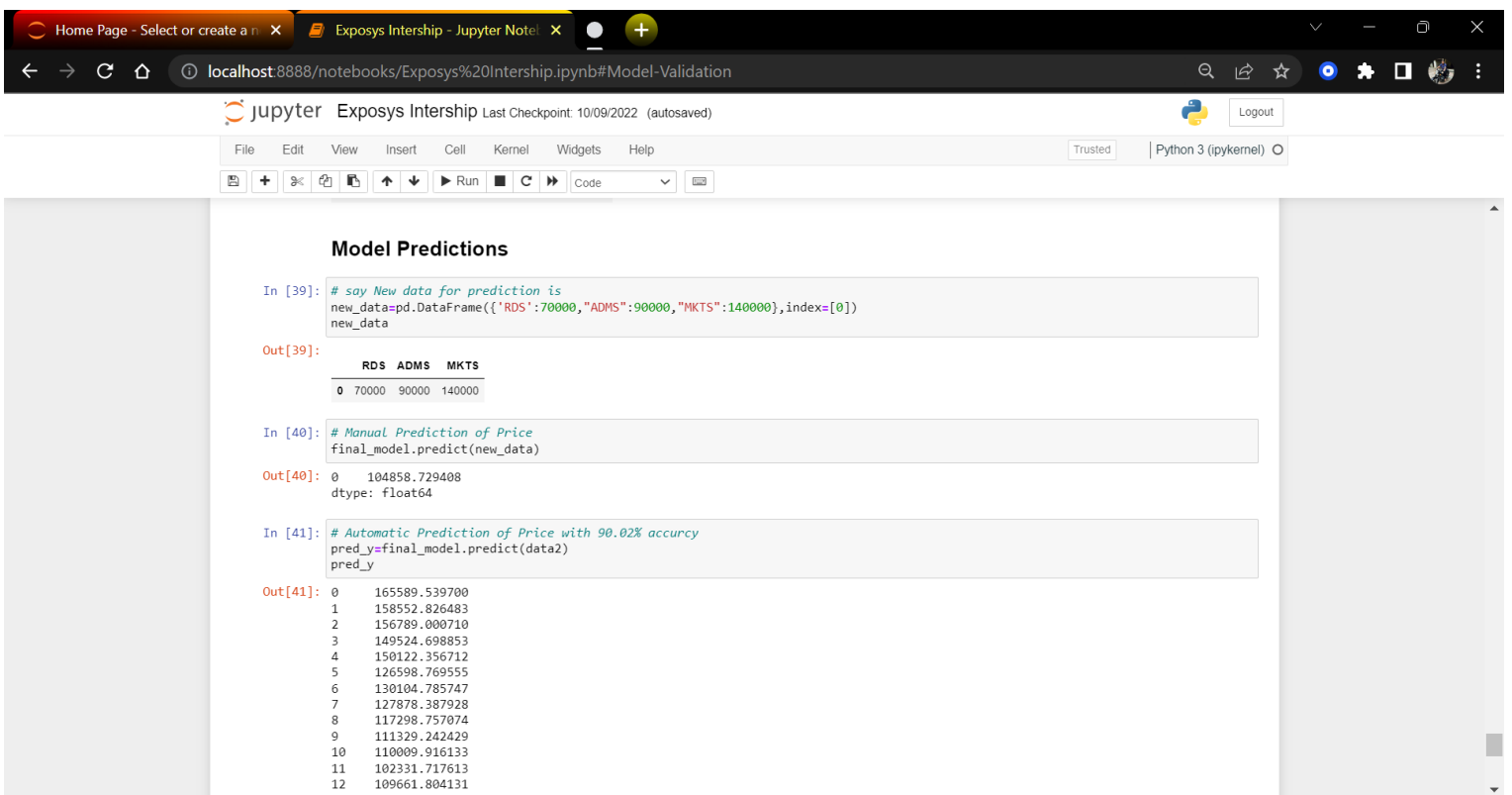
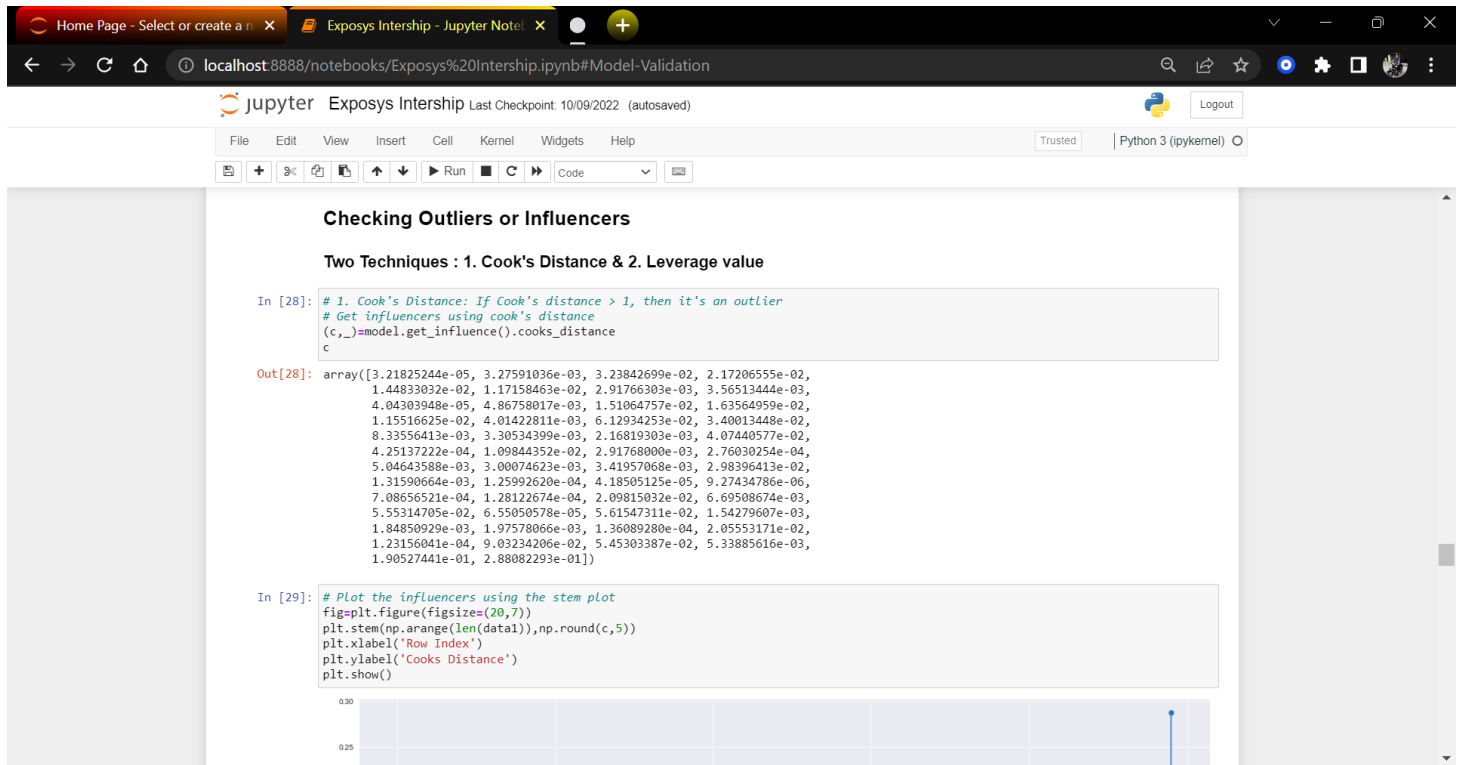
5.) Methodology



We first read the dataset and perform data exploration and data cleaning techniques, and make sure there are no null values or missing values in the dataset. And, then we split the dataset into train and test and then build the model, test it and then validate it. Checking of outliers is done and then the predictions are made, after that model accuracy is checked and the R^2 values are generated.

7.) Implementation: -





The screenshot shows a Jupyter Notebook titled "Exposys Internship" running on a local host. The notebook contains a data table and a code cell for model validation.

	Profit	RDS	ADMS	MKTS
46	1313.40	113016.21	237114.40	43430.73
47	0.00	135426.92	0.00	42559.73
48	542.05	51743.15	0.00	35673.41

Diagnostics and Final Model

```
In [35]: model2=smf.ols("Profit~RDS+ADMS+MKTS",data=data2).fit()

In [36]: while model2.rsquared < 0.99:
    for c in [np.max(c)>1]:
        model2=smf.ols("Profit~RDS+ADMS+MKTS",data=data2).fit()
        (c,_)=model2.get_influence().cooks_distance
        c
        np.argmax(c) , np.max(c)
        data2=data2.drop(data2.index[[np.argmax(c)]],axis=0).reset_index(drop=True)
        data2
    else:
        final_model=smf.ols("Profit~RDS+ADMS+MKTS",data=data2).fit()
        final_model.rsquared , final_model.aic
        print("Thus model accuracy is improved to",final_model.rsquared)

Thus model accuracy is improved to 0.9626766170294073
Thus model accuracy is improved to 0.9614129113440602
Thus model accuracy is improved to 0.962593650298269
Thus model accuracy is improved to 0.9638487279209413
Thus model accuracy is improved to 0.9663901957918793
Thus model accuracy is improved to 0.9706076169779905
Thus model accuracy is improved to 0.9727840588916423
Thus model accuracy is improved to 0.9734292907181952
Thus model accuracy is improved to 0.9785801571833451
Thus model accuracy is improved to 0.9777383743090915
Thus model accuracy is improved to 0.9790510088977512
Thus model accuracy is improved to 0.9790004461890552
Thus model accuracy is improved to 0.9807878666153609
Thus model accuracy is improved to 0.9838299343609735
```

7.) Conclusions: -

The project's primary goal was to prepare an ML model which can predict the profit value of a company if the value of its R&D Spend, Administration Cost and Marketing Spend are given.

Multiple Linear Regression Analysis consists of more than just fitting a linear line through a cloud of data points. It consists of three stages: 1) analysing the correlation and directionality of the data, 2) estimating the model, i.e., fitting the line, and 3) evaluating the validity and usefulness of the model.