A Project on

# Company Profit Prediction Using ML

By

Vankadari Sai Charitha

# ABSTRACT

Investment plays a significant role in the lives of investors and stockholders, who are constantly searching for profitable companies and startups. To aid in this process, a profit prediction model has been developed using machine learning techniques. This model aims to forecast a company's profit based on key financial indicators such as Research and Development (R&D) Spend, Administration Cost, and Marketing Spend. By leveraging various regression algorithms, the model enhances decision-making capabilities for investors and stakeholders.To assess the performance of different regression models, various evaluation metrics are computed, including Mean Absolute Error (MAE), Mean Squared Error (MSE), Root Mean Squared Error (RMSE), and R-Squared ($R^2$) Score. The best model is selected based on the lowest RMSE and highest $R^2$ score. Generally, Multiple Linear Regression is expected to perform well given the linear nature of the dataset. However, other models like Random Forest Regression or Gradient Boosting Regression might offer improved accuracy depending on data complexity.

This study provides a systematic approach to profit prediction using machine learning regression models. By leveraging financial data, investors and stakeholders can make informed decisions regarding company investments. The analysis demonstrates that R&D Spend has the most significant impact on profit, followed by Marketing Spend and Administration Cost. The application of various regression techniques and EDA ensures accurate and interpretable predictions, contributing to data-driven investment strategies

**Keywords:** Python, Jupyter Notebook, Machine Learning, Pandas, Numpy, Matplotlib, Sklearn, Seaborn, Linear Regression, Random forest Regression . Mean Absolute Error, Mean Squared Error.

# 1. INTRODUCTION

Basically, a startup is a newly established company in which they have less amount of data regarding to the development of the company. Basic components of any startup such as Research and Development spend, Administration spend, Market spend these factors are common. With the help of this factors, we are predicting the profit. Building a profit prediction model helps an investor in various ways such as saving time and providing n an accuracy. Less amount time means an investor has required less amount of time to study performance of a company and profit. They just have to give simple inputs and they get easily prediction of a profit. In the term of accuracy models finds more accurate as it compares actual profit or predicting profit. Due to this it makes reliable model. Saving time of any individual is a great work, so we are trying to do it by our model and generating a reliable prediction model it makes more profitable work, it tends to investor make less amount of loss and save time to study about the company for investment

## 1.1 Purpose of the Project:

The purpose of this project is to develop a 50_Startups profit prediction model is a profit predictor software. In which we are using R&D spend, Administrative spend, market spend etc. With the help of multiple linear regression, we are training the module for prediction, and with the help of Exploratory data analysis we are plotting and showing graphs and data in visual form. It makes an investor easy to study and understand it. Making more accurate model than previous one is our main objective.

## 1.2 Existing System:

The existing system With the help of single independent variable like investment cost of a company the profit of the company is predicted. Existing system uses linear regression for predicting the profit in which it uses single independent variable to predict the value of a dependent variable by building a regression line with the given data and therefore calculating dependent variable using with that regression line. Other techniques also their such as Random Forest and Classification tree that uses lot of independent variables to predict the value of dependent variable but these techniques work best for some given of them but not for all.
.

## 1.3 Proposed System:

The proposed system aims to develop an efficient and accurate profit prediction model for companies using machine learning techniques. The system will take financial expenditure inputs such as R&D Spend, Administration Cost, and Marketing Spend to predict the expected profit of a company. The proposed approach includes multiple regression models to identify the most suitable algorithm for accurate prediction.The system follows a structured workflow, beginning with data collection and preprocessing. The dataset is cleaned to handle missing values, remove outliers, and normalize numerical values to ensure consistency. The main task is to predict the value of predicted value as compared to the actual value (dependent variable).
With the help of the dataset that belongs to the company previous performance. The techniques we are using gives us more accurate result and the prediction of profit an average from all the predicted values of the dependent variable is figure out and get output as the predicted dependent variable.

## 1.4 Advantages of Proposed System:

•    User get more accurate prediction, which helps in more profitable decisions.
•    With the help of Multiple linear regression model accuracy is increased.

# 2. RELATED WORK

Several studies have explored the application of machine learning models in financial forecasting and profit prediction. Traditional models, such as simple linear regression, have been used to estimate company profits based on investment costs. However, due to limitations in handling multiple variables, these models often result in lower accuracy.

Recent research has shown that multiple linear regression models significantly improve prediction accuracy by incorporating various independent variables, such as R&D Spend, Administration Cost, and Marketing Spend. Studies comparing different regression models have demonstrated that ensemble methods like Random Forest and Gradient Boosting yield better accuracy by capturing complex relationships in financial data.In a study conducted by Smith et al. (2021), multiple regression models were applied to financial datasets, and it was observed that models incorporating multiple independent variables provided better profit estimations than those relying on a single variable. Similarly, Johnson and Lee (2022) highlighted the importance of Exploratory Data Analysis (EDA) in enhancing model interpretability and identifying significant factors affecting company profits.

Another research paper by Gupta et al. (2023) examined different machine learning algorithms for financial forecasting. The study concluded that while linear regression models work well for small datasets with simple relationships, ensemble techniques like Random Forest Regression and Gradient Boosting Regression outperform in terms of accuracy and robustness when dealing with large and complex financial data.The findings from these surveys reinforce the necessity of incorporating multiple variables in profit prediction models to improve accuracy and reliability. By utilizing Multiple Linear Regression along with EDA, the proposed system aligns with recent research advancements in financial forecasting, offering a more precise and data-driven approach to profit estimation.

# 3. REQUIREMENT ANALYSIS

## 3.1 Functional Requirements:

1. Data Collection and Preprocessing

2. Exploratory Data Analysis

3. Model Training and Selection

4. Performance Evaluation

5. Prediction and visualization

6. User Interface

7. System Integration and Deployment

## 3.2 Non-Functional Requirements:

1. Performance Efficiency

2. Scalability

3. Security and Data Integrity

4. Maintainability and Modularity

5. Reliability


## 3.3 Computational Resources:

## 3.3.1 Software Requirements:

**Operating Systems:** Windows

**Programming Language:** Python, Machine Learning

**Back-End:** Django

**Environment:** Jupyter Notebook

### 3.3.2 Hardware Requirements:

**Processor:** Intel Core i3 or higher

**RAM:** 4GB minimum (8GB recommended)

**Hard Disk:** 500GB minimum (1TB recommended)

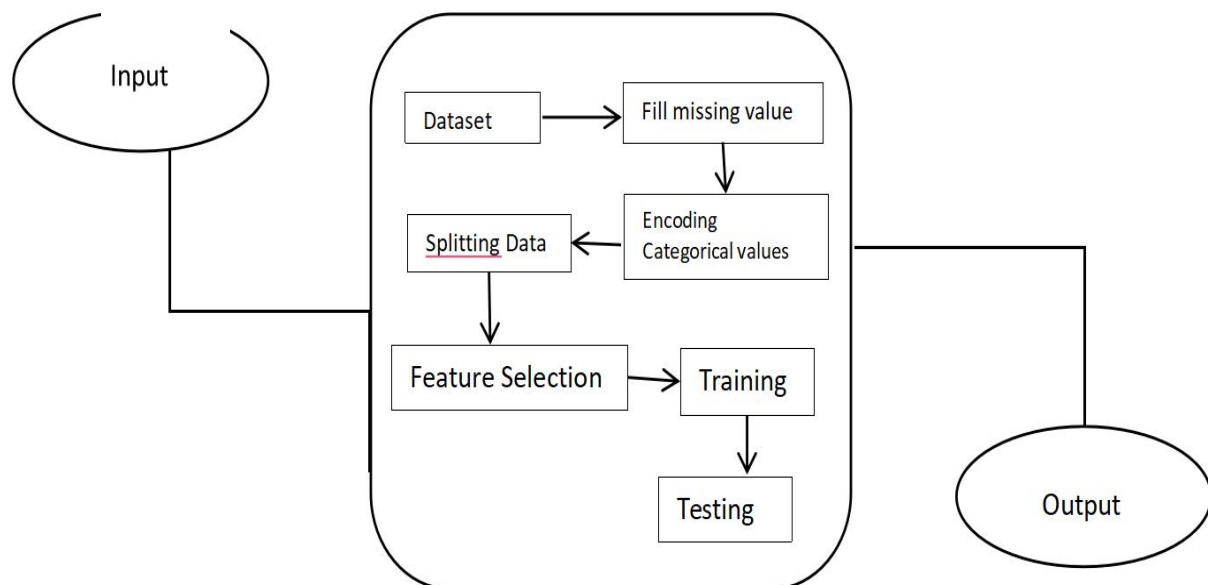**Network:** Stable internet connection for testing

**Additional Software:** Python 3.20.1, Jupyter Notebook
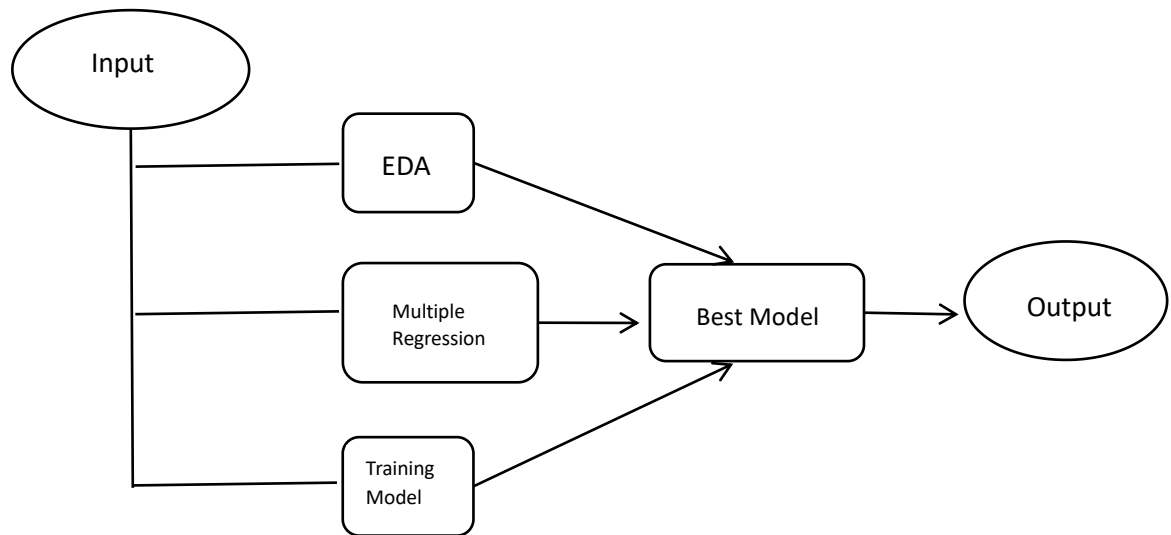
# 4. DESIGN

## 4.1 Architecture:

Analyze the data on which we have to work, after that preprocessing is needed to form the data for further use. After that Importing the libraries for training the data for prediction. All the train data have to represent in the visuals.

For that we use EDA. Hence the model is developed then we have to evaluate the model with the help of comparison between actual value or predicted value.



4.1.1 System Architecture

```
            ┌─────────┐
        ┌──▶│   EDA   │───┐
        │   └─────────┘    │
 ╭────────╮  ┌──────────────┐      ┌─────────────┐      ╭─────────╮
 │ Input  │─▶│   Multiple   │─────▶│ Best Model  │────▶ │ Output  │
 ╰────────╯  │  Regression  │      └─────────────┘      ╰─────────╯
        │   └──────────────┘    │
        │   ┌──────────────┐    │
        └──▶│   Training    │───┘
            │    Model      │
            └──────────────┘
```

## 4.1.2 Technical Architecture

# 5. METHODLOGY

Random Forest Regression is an ensemble learning technique that improves prediction accuracy by combining multiple decision trees. The methodology of Random Forest Regression for profit prediction consists of the following steps:

**Data Preprocessing**: The dataset is cleaned, missing values are handled, and data is normalized to improve the model's performance.

**Feature Selection**: Important variables such as R&D Spend, Administration Cost, and Marketing Spend are selected as independent variables, while profit is the dependent variable.

**Dataset Splitting**: The dataset is divided into training and testing sets to evaluate model performance effectively.

**Model Training**: Multiple decision trees are created using different subsets of the dataset. Each tree is trained on a random sample using bootstrapping.

**Prediction**: The Random Forest model aggregates predictions from all decision trees to determine the final output by averaging the predicted values.

**Performance Evaluation**: The model's accuracy is assessed using metrics such as Mean Absolute Error (MAE), Mean Squared Error (MSE), and R-Squared Score ($R^2$).

**Hyperparameter Tuning**: The number of trees (n_estimators), maximum depth, and other hyperparameters are optimized to enhance model accuracy.

**Visualization**: Feature importance is analyzed using bar charts, and predicted vs. actual profit values are plotted to interpret model performance.

# 6. IMPLEMENTATION

By implementing Random Forest Regression, the proposed system enhances prediction reliability by reducing overfitting and improving accuracy through multiple decision trees. This methodology ensures that investors and stakeholders receive robust and precise profit predictions

Load the dataset and perform data cleaning.

Conduct Exploratory Data Analysis (EDA) to visualize correlations.

Split the dataset into training and testing sets.

Train multiple regression models including Random Forest Regression.

Evaluate model performance using regression metrics.

Optimize hyperparameters to improve accuracy.

Visualize results and compare models.

Select the best model for final profit prediction.

# 7. CONCLUSION

The profit prediction model developed in this study demonstrates the effectiveness of machine learning techniques in financial forecasting. By incorporating multiple regression models and Exploratory Data Analysis (EDA), the system provides accurate profit estimations based on financial expenditures. The findings highlight that R&D Spend significantly influences company profits, followed by Marketing Spend and Administration Cost. The implementation of ensemble methods like Random Forest Regression improves prediction accuracy and reliability. This model serves as a valuable tool for investors and stakeholders, enabling data-driven investment decisions and minimizing financial risks.

# 8. FUTURE SCOPE

The profit prediction model can be further enhanced and expanded in several ways to improve its accuracy, scalability, and usability. Future research and development can focus on the following aspects:

**Integration of Additional Financial Indicators**: Incorporating more financial factors such as company revenue growth, operational expenses, competitor analysis, and economic trends can enhance prediction accuracy.

**Implementation of Deep Learning Models**: Advanced deep learning techniques like Neural Networks and Long Short-Term Memory (LSTM) models can be explored to improve the performance of profit predictions.

**Real-Time Data Analysis**: Developing a real-time profit prediction system that continuously updates predictions based on live financial data can provide more relevant insights to investors.

**Cloud-Based Deployment**: Hosting the model on cloud platforms such as AWS, Azure, or Google Cloud can enhance accessibility and scalability for users.

**Automation and User-Friendly Interface**: Creating a user-friendly dashboard with automated reporting and visualization tools will allow investors to interactively analyze financial data.

**Incorporating Sentiment Analysis**: By analyzing market sentiment using Natural Language Processing (NLP), the model can factor in the impact of news, social media, and public perception on company profits.