



INTERN PROJECT REPORT – DATA SCIENCE

26.06.23 – 26.07.23

EXPOSYS DATA LABS

DOMAIN:DATASCIENCE

LANGUAGE: PYTHON

LIBRARIES USED:NUMPY,PANDAS,SEABORN,MATPLOTLIN.PYPLOTT,SKLEARN.

- SOWMITHRAN S

(B.TECH ARTIFICIAL INTELLIGENCE AND DATA SCIENCE)

ABSTRACT

From the given dataset ("50-Startups.csv"). I have done some models which can give a best score of the given dataset .This project aims to develop a robust machine learning (ML) model capable of predicting the profit value of companies based on their R&D Spend, Administration Cost, and Marketing Spend. The dataset contains information from 50 companies, including their respective profits. By employing advanced ML techniques, such as regression analysis, feature engineering, and model optimization, we seek to create an accurate and reliable predictive model. The model's performance will be evaluated using various metrics to ensure its effectiveness and generalization to unseen data. Ultimately, the successful implementation of this ML model will provide valuable insights for businesses to make informed decisions and optimize their financial strategies, contributing to enhanced profitability and competitiveness in the market.

TABLE OF CONTENT

CHAPTER	TITLE	PAGE NO
01.	Introduction	05
02.	Existing Method	07
03.	Proposed method with Architecture	09
04.	Methodology	11
05.	Implementation	13
06.	Conclusion	15

CHAPTER - 1

INTRODUCTION

In today's competitive business landscape, companies are constantly seeking ways to enhance their profitability and stay ahead of the competition. Making informed decisions based on data-driven insights has become a crucial aspect of achieving sustained success. In this context, the use of Machine Learning (ML) techniques has emerged as a powerful tool for predicting and understanding various business outcomes. The goal of this project is to develop a robust ML model that can accurately predict the profit value of companies based on their R&D Spend, Administration Cost, and Marketing Spend. These three key factors have a significant impact on a company's financial performance and understanding their relationships is vital for effective resource allocation and business planning. The dataset at hand contains essential information from 50 companies, providing insights into their R&D investments, administrative expenses, marketing strategies, and corresponding profits. Leveraging this data, we aim to construct a predictive model that can generalize well to unseen data, enabling businesses to forecast their expected profits with greater accuracy. By employing state-of-the-art ML techniques, including regression analysis, feature engineering, and model optimization, we aim to extract meaningful patterns and correlations from the dataset. The project's success will be measured not only by the model's ability to predict profits accurately but also by its ability to

uncover valuable insights into the factors that drive profitability in different industries and business contexts. The implications of this project are significant. A reliable profit prediction model will empower companies to make informed decisions regarding their financial strategies, identify areas for cost optimization, and allocate resources effectively. Additionally, entrepreneurs and investors can utilize the model to assess the potential profitability of new ventures, aiding them in making sound investment choices. In conclusion, this project endeavors to leverage the power of machine learning to construct a predictive model that can aid businesses in understanding and predicting their profitability. By unlocking valuable insights from the data, companies can chart their course for sustainable growth and prosperity in an increasingly competitive business landscape

CHAPTER – 2

EXISTING METHOD

In the realm of business and finance, predicting company profits has long been a topic of interest, and various traditional statistical methods have been employed for this purpose. Some of the conventional approaches include:

1. Linear Regression: One of the simplest methods used to predict profit is linear regression. It establishes a linear relationship between the input features (R&D Spend, Administration Cost, Marketing Spend) and the target variable (Profit). While it serves as a good starting point, it may not capture complex non-linear relationships present in the data.

2. Multiple Regression: This method extends linear regression to incorporate multiple input features simultaneously. It attempts to find the best-fit plane in the multi-dimensional feature space to predict the target variable. However, like linear regression, it may lack the ability to handle non-linear relationships effectively

3. Decision Trees: Decision trees are a popular method for both regression and classification tasks. They recursively split the data based on the input features, creating a tree-like structure to make predictions. While decision trees can capture non-linear relationships, they are prone to overfitting and may not generalize well to new data

4. Ensemble Methods: Techniques like Random Forest and Gradient Boosting are popular ensemble methods that combine multiple models to make predictions. These methods aim to mitigate the weaknesses of individual models and improve prediction accuracy. However, these traditional methods may not fully exploit the complexities and non-linearities present in the data, which can limit their predictive power.

Moreover, feature engineering and model hyperparameter tuning require substantial manual effort and expertise. In recent years, with the advent of deep learning and more advanced ML techniques, researchers have explored the use of neural networks and other sophisticated algorithms for profit prediction tasks. These models can learn intricate patterns from the data and potentially offer improved prediction accuracy compared to conventional methods. In this project, we will leverage the power of advanced ML techniques, possibly including deep learning algorithms, to develop a state-of-the-art profit prediction model that can outperform traditional methods and deliver more accurate and reliable forecasts for companies based on their R&D Spend, Administration Cost, and Marketing Spend

CHAPTER – 3

PROPOSED METHOD WITH ARCHITECTURE

Let's propose a method using traditional Machine Learning (ML) algorithms for profit prediction: Proposed Method using Traditional ML Algorithms:

1. Data Preprocessing: The first step is to preprocess the dataset by handling missing values, performing feature scaling (if required), and encoding categorical variables. Proper data preprocessing ensures that the ML algorithms can effectively learn from the data.

2. Feature Selection: Since we are using traditional ML algorithms, feature selection is crucial to identify the most relevant features that contribute significantly to profit prediction. Techniques like correlation analysis, recursive feature elimination, or feature importance from tree-based models can be employed to select the most important features.

3. Model Selection: - Given that the dataset is relatively small (with only 50 companies), we can experiment with various ML algorithms suitable for regression tasks. Some potential algorithms to consider are -
Linear Regression: As a baseline model for its simplicity and interpretability. - Random Forest: A powerful ensemble method that can handle non-linear relationships and provide feature importances. - Gradient Boosting: Another ensemble method that can capture complex interactions between features and often delivers strong predictive performance.

4. Model Training and Evaluation: - The dataset will be split into training and testing sets (e.g., 80-20 or 70-30 split). - Each ML algorithm will be trained on the training data and evaluated on the testing data using appropriate metrics such as Mean Absolute Error (MAE), Mean Squared Error (MSE), and R-squared (R^2) to assess the model's accuracy and performance.

5. Hyperparameter Tuning: - To improve the model's performance, hyperparameter tuning will be conducted using techniques like grid search or random search. This process involves exploring different combinations of hyperparameters for each algorithm to find the optimal configuration that yields the best results.

6. Model Comparison: - Once all the models are trained, their performance metrics will be compared to identify the most accurate and robust model for profit prediction. The chosen model will be the one with the lowest error and the highest R-squared value.

7. Profit Prediction: - After selecting the best model, it will be used to predict the profit values for new, unseen data

CHAPTER – 4

METHODOLOGY

The methodology used in the code for the project report involves the following steps:

- 1. Data Collection:** The dataset "50_Startups.csv" is read using pandas to gather information about R&D Spend, Administration Cost, Marketing Spend, and Profit for 50 companies.
- 2. Data Manipulation:** Various operations like checking for missing values, displaying the head and tail of the dataset, and using the describe() function to get basic statistical information about the data are performed.
- 3. Data Visualization:** Different data visualization techniques are used to explore the relationships and distributions between variables. These visualizations include line plots, histograms, area plots, scatter plots, stem plots, polar plots, distplots, violin plots, and countplots.
- 4. Data Preprocessing:** The data is divided into independent variables (R&D Spend, Administration, Marketing Spend) and the dependent variable (Profit).
- 5. Training and Testing Data Split:** The dataset is split into training and testing sets using the train_test_split function from sklearn.model_selection.
- 6. Model Training and Testing:** Various regression algorithms are implemented and trained on the training data, and predictions are made on the testing data.

7. Regression Metrics Calculation: Different regression metrics like Mean Absolute Error (MAE), Mean Squared Error (MSE), Mean Absolute Percentage Error (MAPE), and R-squared (R^2) are calculated to evaluate the performance of each regression model.

8. Choosing the Best Model: The R-squared values for each regression model are compared, and the model with the highest R-squared value is considered the best model for predicting the profit value of the company. 9. Visualization of Model Comparison: A bar plot is used to visualize and compare the R-squared values of different regression models to identify the best performing model. 10. Conclusion: Based on the comparison of R-squared values, the best regression model is chosen for profit prediction

CHAPTER – 5

IMPLEMENTATION

The code I have done here is an implementation of different regression algorithms for predicting the profit value of a company based on its R&D Spend, Administration Cost, and Marketing Spend. It also calculates various regression metrics and then compares the performance of the models to choose the best one. Here's a summary of the steps taken in the code:

1. **Data Collection:** - The code reads the dataset "50_Startups.csv" using pandas and stores it in the DataFrame `df`.
2. **Data Manipulating:** - Various operations are performed to understand the dataset better, such as checking for missing values (`isnull()`), describing the dataset (`describe()`), etc.
3. **Data Visualization:** - Different visualization techniques are used to explore the relationships between the input features and the target variable (Profit).
4. **Regression Algorithms:** - The code implements different regression algorithms, including Linear Regression, Random Forest Regression, Support Vector Regression (SVR), Ridge Regression, Lasso Regression, Decision Tree Regression and Gradient Boosting.
5. **Model Training:** - Each regression model is trained using the corresponding training set.

6. Model Prediction: - The trained models are used to make predictions on the test set.

7. Regression Metrics: - Various regression metrics, such as Mean Absolute Error (MAE), Mean Squared Error (MSE), Mean Absolute Percentage Error (MAPE), and R-squared (R^2), are calculated for each model.

8. Choosing the Best Model: - The R-squared values of the different models are plotted to compare their performance, and the model with the highest R^2 value is considered the best

CHAPTER – 6

CONCLUSION

Based on the implementation and analysis of different regression models, we can draw the following conclusions:

Model Performance: The code implemented and evaluated various regression models, including Linear Regression, Random Forest Regression, Support Vector Regression (SVR), Ridge Regression, Lasso Regression, and Decision Tree Regression. The performance of each model was assessed using common regression metrics such as Mean Absolute Error (MAE), Mean Squared Error (MSE), and R-squared (R^2).

Best Model Selection: To choose the best model, the code compared the R-squared values of all the models. The model with the highest R-squared value is generally considered the best as it explains the highest proportion of the variance in the target variable (Profit).

Model Comparison: The R-squared values for each model were plotted to visualize and c The plot helped in understanding which model performed better in predicting the profit based on the given input features (R&D Spend, Administration Cost, and Marketing Spend).

Predictive Ability: The selected best model can be used to predict the profit value of a company when its R&D Spend, Administration Cost, and Marketing Spend are provided as input. Insights: The scatter plots and other visualizations provided insights into the relationships between the input features and the target variable (Profit). The conditional analysis provided a deeper understanding of how the profit varies concerning different ranges of R&D Spend, Administration Cost, and Marketing Spend.

Recommendations: Based on the analysis, it is recommended to use the best-performing regression model for predicting profits in similar scenarios with R&D Spend, Administration Cost, and Marketing Spend as input features.

Limitations: It's important to note that the quality of predictions depends on the quality of the data and the assumptions made by each regression mode Other factors not included in the dataset could also influence the profit of companies, and these should be considered for a more comprehensive analysis.

Future Work: Further analysis and feature engineering could be performed to improve the predictive performance of the models.

Additional data or external factors that might affect profit could be incorporated into the analysis to enhance the model's accuracy.

Overall, the code provides a comprehensive approach to constructing, training, and evaluating regression models for profit prediction based on given input features. The chosen best model, as indicated by its R-squared value, can be used for future profit predictions in similar scenarios

MY RESULT

