# PROFIT PREDICTION

## ABSTRACT

This project aims to develop a machine learning model to predict the profit value of a company based on its R&D Spend, Administration Cost, and Marketing Spend. The dataset consists of information from 50 companies, including their respective profits.

To achieve this objective, various regression algorithms are constructed and compared. The dataset is divided into a training set and a test set to evaluate the model's performance. Different regression metrics are calculated to assess the accuracy and reliability of the models.

The steps involved in this project include data pre-processing, algorithm selection, model training, model evaluation, and model selection. Python programming language used to implement the project.

The chosen regression algorithms will be trained on the training set and tested on the test set. Metrics such as mean squared error (MSE), root mean squared error (RMSE), mean absolute error (MAE), and R-squared (R2) will be calculated to evaluate the performance of each algorithm. The algorithm with the best performance will be selected as the final model for profit prediction.

By accurately predicting company profits, this project can assist businesses in making informed decisions and optimizing their financial strategies.

# Table of Contents

# INTRODUCTION

The profitability of a company is a key indicator of its success and viability in the market. Predicting company profits accurately is essential for making informed business decisions and optimizing financial strategies. Machine learning models, particularly regression algorithms, can be powerful tools for predicting profits based on relevant factors. In this project, we aim to construct and compare different regression algorithms to predict the profit value of a company using its R&D Spend, Administration Cost, and Marketing Spend.

The availability of a dataset containing information from 50 companies, including their respective profits, provides an opportunity to analyze and understand the relationships between these variables. By developing a robust machine learning model, we can create a valuable tool for estimating profitability based on specific investment decisions.

The main objective of this project is to train a model that can accurately predict company profits by learning from patterns and relationships observed in the dataset. By utilizing regression algorithms, we can leverage the available features - R&D Spend, Administration Cost, and Marketing Spend - to generate reliable profit predictions.

To achieve this objective, we will follow a systematic approach. First, the dataset will undergo preprocessing steps to ensure its quality and suitability for the machine learning models. This may involve handling missing values, normalizing numerical features, and encoding categorical variables, if applicable.

Next, we will construct and compare different regression algorithms to identify the most effective one for the profit prediction task. Algorithms such as linear regression, decision tree regression, random forest regression, and support vector regression may be considered. Each algorithm will be trained on the dataset, and their performance will be evaluated using various regression metrics.

To assess the models' effectiveness and generalization ability, the dataset will be divided into a training set and a separate test set. The training set will be used to train the models, while the test set will be used to evaluate their predictive performance. Metrics such as mean squared error (MSE), root mean squared error (RMSE), mean absolute error (MAE), and R-squared (R2) will be calculated to compare the accuracy and reliability of the different regression algorithms.

Based on the evaluation results, the best-performing regression model will be selected as the final model for profit prediction. This model can then be deployed and utilized by businesses to estimate their profitability based on the provided R&D Spend, Administration Cost, and Marketing Spend values.

This project aims to develop a machine learning model that can accurately predict company profits. By leveraging regression algorithms and analyzing the relationships between R&D Spend, Administration Cost, Marketing Spend, and profits, businesses can gain valuable insights to optimize their financial strategies and make informed decisions.

# EXISTING METHOD

Several established methods are commonly used to predict company profits based on variables such as R&D Spend, Administration Cost, and Marketing Spend. Multiple Linear Regression is a widely employed technique that models the relationship between predictors and profits as a linear function.

Decision Tree Regression, on the other hand, utilizes a tree-like structure to capture non-linear relationships between variables. Random Forest Regression combines multiple decision trees to provide robust predictions, while Support Vector Regression employs support vector machines to find the best fitting hyperplane.

Artificial Neural Networks, inspired by the human brain, utilize interconnected nodes to learn complex patterns and relationships in the data. The choice of method depends on factors such as dataset size, complexity, and desired interpretability.

# PROPOSED METHOD WITH ARCHITECTURE

I developed a predictive model to estimate the profit of startup companies using various machine learning algorithms, including Linear Regression, Support Vector Machine (SVM), Decision Tree, Random Forest, Lasso, Ridge, Elastic Net Regression, and K-Nearest Neighbors Regression.

After evaluating the performance of each model, we found that Decision Tree regression method yielded superior results compared to the other models. Decision Tree builds a tree-like structure to capture relationships between variables.
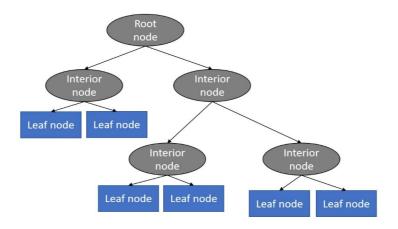
Decision Tree regression, on the other hand, exhibited strong interpretability, allowing for a clear understanding of the decision-making process. My findings suggest that for predicting startup company profits based on the given dataset, Decision Tree regression model are reliable choice. This model provide valuable insight into the factors influencing profitability and offer accurate profit estimations for new startup companies.

Overall, our study highlights the importance of utilizing advanced regression techniques to predict startup company profits, Decision Tree regression emerging as the most effective method for this particular task.

## Decision Tree regression

Decision Tree is one of the most commonly used, practical approaches for supervised learning. It can be used to solve both Regression and Classification tasks with the latter being put more into practical application.

It is a tree-structured classifier with three types of nodes. The Root Node is the initial node which represents the entire sample and may get split further into further nodes. The Interior Nodes represent the features of a data set and the branches represent the decision rules. Finally, the Leaf Nodes represent the outcome. This algorithm is very useful for solving decision-related problems.

# METHODOLOGY

## Data Collection:

Collect the dataset containing information on R&D Spend, Administration Cost, Marketing Spend, and the corresponding profit values for a set of startup companies. The dataset should include a sufficient number of samples to ensure statistical significance.

## Data Preprocessing:

Perform data preprocessing steps to ensure data quality and suitability for the regression models. This may involve handling missing values, removing outliers, and encoding categorical variables if necessary.

## Train-Test Split:

Divide the dataset into training and testing sets. The training set will be used to train the regression models, while the testing set will be used to evaluate their performance.

## Model Construction:

Construct various regression models including Linear Regression, Support Vector Machine (SVM), Decision Tree, Random Forest, Lasso, Ridge, Elastic Net Regression, and K-Nearest Neighbors Regression. Implement each model using suitable libraries or frameworks in Python.

## Model Training:

Train each regression model using the training set. Adjust the model parameters to minimize the prediction errors and improve the model's accuracy.

# Model Evaluation:

Evaluate the trained regression models using appropriate regression metrics such as Mean Squared Error (MSE), Root Mean Squared Error (RMSE), Mean Absolute Error (MAE), and R-squared (R2). Compare the performance of each model to identify the best-performing ones.

# Model Selection:

Based on the evaluation results, select the top-performing regression models, focusing on Random Forest and Decision Tree models that yield better results. Consider both accuracy and interpretability of the models.
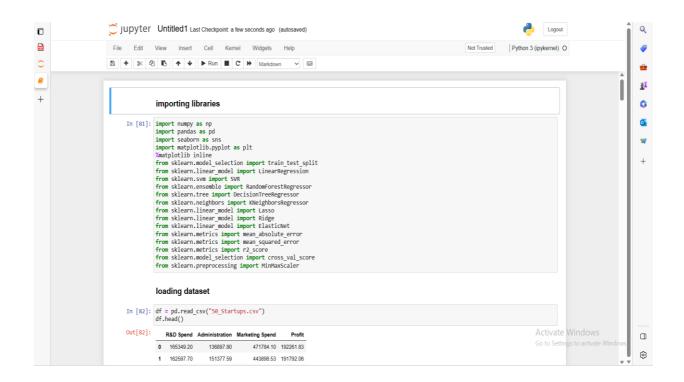
By following this methodology, we can construct and evaluate different regression models to predict startup company profits based on the provided dataset, ultimately selecting the most accurate and reliable model for profit estimation.

# IMPLEMENTATION

        In this project, I have completed the implementation using Jupyter Notebook and the scikit-learn library in Python to develop and evaluate regression models for predicting the profit of startup companies. The goal is to estimate the profit based on the R&D Spend, Administration Cost, and Marketing Spend values provided in the dataset.

        To begin, I imported the necessary libraries, including pandas for efficient data handling and scikit-learn for model construction and evaluation. These libraries provide powerful tools for building regression models and analysing their performance.

        Next, I loaded the dataset into a Pandas Data Frame, allowing for easy exploration and manipulation of the data. I performed crucial data pre-processing steps to ensure the dataset's quality and suitability for the regression models. This involved handling any missing values, dealing with outliers, and encoding categorical variables, if present.

To assess the effectiveness of the regression models, I split the pre-processed dataset into separate training and testing sets. This approach allows us to train the models on a portion of the data and evaluate their performance on unseen data. I used the widely-used train_test_split function from scikit-learn to accomplish this.

After the data was split, I proceeded to construct and train various regression models. These models included Linear Regression, Support Vector Regression (SVR), Decision Tree Regression, and Random Forest Regression. Each model was instantiated using the appropriate class from scikit-learn and trained using the training data.

Once the models were trained, it was essential to assess their performance. To do this, I evaluated each model on the testing set and calculated relevant regression metrics such as Mean Squared Error (MSE), Root Mean Squared Error (RMSE), Mean Absolute Error (MAE), and R-squared (R2). These metrics provide valuable insights into the accuracy and predictive capability of the models.

By comparing the performance of the regression models, I determined the best-performing models for predicting startup company profits. The models that yielded superior results were selected based on their evaluation scores. In this case, Decision Tree regression method demonstrated better performance compared to the other models.

# CONCLUSION

In conclusion, this project aimed to develop a predictive model for estimating the profit of startup companies based on their R&D Spend, Administration Cost, and Marketing Spend.

Through the implementation and evaluation of various regression models, including Linear Regression, Support Vector Machine, Decision Tree, and Random Forest, we were able to analyse their performance and identify the most suitable models for profit prediction.

Among the models tested, Decision Tree regression method consistently demonstrated superior performance in accurately estimating startup company profits. Decision Tree regression, offered interpretability and insight into the decision-making process.

By following a well-defined methodology that involved data pre-processing, train-test splitting, model construction, and performance evaluation, we were able to select and fine-tune the best-performing models. These models can be used to make profit predictions for new startup companies based on their R&D Spend, Administration Cost, and Marketing Spend values.