

# **PROJECT REPORT**

## **STARTUP FUNDING ANALYSIS: Exploring Venture Capital Investments in Various Industries**

### **Project Overview:**

This data science project aims to predict the operating status of startups using machine learning techniques. The project involves analyzing a dataset containing information about startups, performing data preprocessing, conducting exploratory data analysis (EDA), feature engineering, model training, and deploying a predictive model using Streamlit.

### **Problem Statement:**

Despite the rapid growth of the startup scene and the increasing investments from venture capital firms, understanding the dynamics and characteristics of these newly established firms remains a challenge. The chosen dataset provides a valuable glimpse into this captivating realm, presenting an opportunity to delve deeper and unravel the underlying trends and patterns within the startup ecosystem.

### **Project Objective:**

The objective of this project is to analyze a startup dataset and derive valuable insights about the dynamics and characteristics of startups. By exploring funding trends, industry preferences, geographical distribution, and success metrics, we aim to provide actionable knowledge for entrepreneurs, investors, and policy-makers. The findings will contribute to a deeper understanding of the startup landscape and support informed decision-making to foster the growth of the startup ecosystem.

### **Dataset:**

The dataset used in this project consists of features related to startups, including funding details, investment amounts, funding rounds, and industry groups. The target variable is the operating status of the startups, which can be categorized as "closed," "operating," or "acquired." The dataset contains a combination of numerical and categorical features.

### **Data Preprocessing:**

During the data preprocessing phase, several preprocessing steps were performed to prepare the data for analysis and modeling. These steps included handling missing values, transforming categorical variables into numerical representations, and creating new categorical features

based on descriptive summaries. The dataset was then split into feature matrix  $X$  and target variable  $y$ .

#### Exploratory Data Analysis:

Exploratory data analysis was conducted to gain insights into the dataset and understand the relationships between different variables. Visualizations such as histograms, bar plots, and correlation matrices were used to explore the distribution of variables, identify patterns, and detect any potential correlations.

#### **Feature Engineering:**

Feature engineering techniques were applied to enhance the predictive power of the model. New categorical features were created based on the values of certain numerical features, allowing for a more nuanced representation of the data. This step aimed to capture relevant information and improve the model's ability to predict the operating status of startups accurately.

#### **Model Training and Evaluation:**

In this project, we evaluated several machine learning models for predicting the target variable. The models included Decision Tree, KNN, Random Forest, AdaBoost, Gradient Boosting, XGBoost, and SVM. We assessed the models based on their RMSE (Root Mean Squared Error) and R-squared metrics.

The results revealed that the XGBoost (Tuned) model outperformed the other models. It achieved the lowest RMSE on the test set (74,876,296) and the highest R-squared value (0.227206). This indicates its superior accuracy in predicting the target variable and its ability to explain a significant proportion of the variance.

However, it's worth noting that some models, such as AdaBoost and SVM, performed poorly in terms of RMSE and R-squared, suggesting limitations in capturing the target variable's variations.

Overall, the XGBoost (Tuned) model stands out as the recommended choice due to its strong predictive performance. Further enhancements and optimizations can be explored to improve the model's accuracy and generalization capabilities.

#### **Limitation of the Dataset:**

1. Limited Generalization: The model's poor performance on the test set suggests it struggles to apply learned patterns to new data.
2. Overfitting: The large train-test performance gap indicates the model has over-adapted to the training data, hindering its ability to generalize.

3. **Data Variability:** The dataset's lack of diversity may limit the model's exposure to various scenarios and real-world variations.
4. **Unaccounted Factors:** Missing relevant features or unconsidered factors in the dataset can limit the model's predictive capabilities and understanding.

### **Conclusion and Recommendations:**

1. Based on the analysis conducted in this project, several key findings and recommendations can be made to assist investors in making informed decisions
2. Investment-related variables such as total investment and funding rounds (e.g., round A, round B) have a significant impact on the operating status of startups. Investors should consider the overall investment levels and strategically plan and secure funding during critical funding rounds.
3. Equity crowdfunding presents an attractive opportunity for startups to raise funds and attract potential investors. Investors should explore equity crowdfunding platforms to diversify their investment portfolio.
4. Seed funding plays a crucial role in kick-starting startup operations and increasing their chances of success. Investors should pay attention to startups that have successfully secured seed funding.
5. Industry-specific factors and dynamics significantly influence the operating status of startups. Investors should consider the industry group of startups and analyze its potential for growth and success.
6. Venture funding and angel funding can provide valuable resources and expertise to fuel startup growth and expansion. Investors should actively seek opportunities to invest in startups that have secured venture capital or angel funding.
7. These recommendations serve as a guide for investors to navigate the startup landscape and optimize their investment strategies.

### **Future Enhancements:**

Future enhancements to this project could include:

- Incorporating real-time data feeds to keep the model up-to-date and improve its predictive capabilities. Expanding the feature set by incorporating additional external data sources to provide a more comprehensive analysis.
- Refining the model using advanced techniques such as ensemble learning or neural networks to enhance prediction accuracy.

- Conducting further analysis on the impact of specific features on startup operating status.

### **Deployment and Interactive Interface:**

The trained model was deployed using the Streamlit framework to create an interactive web application. The application allows users to input their investor preferences, such as desired investment amount, investment timeframe, and risk appetite. Based on these preferences and the selected features, the model generates predictions for the operating status of startups. Additionally, investment recommendations and additional information about recommended startups are provided to assist users in making informed investment decisions.

### **Dependencies:**

The project has the following dependencies:

- Python (version 3.7+)
- Pandas
- Numpy
- Matplotlib
- Seaborn
- Scikit-learn
- XGBoost
- Mlxtend
- Streamlit

### **Installation and Usage:**

1. Clone the repository to your local machine.
2. Install the necessary dependencies using pip or conda package manager.
3. Run the Streamlit application to launch the web interface.
4. Follow the instructions on the web interface to input your investor preferences and generate predictions for startup operating status.

### **How to Use**

1. Open the web application in your browser after running the Streamlit application.
2. On the sidebar, you will find the "Investor Preferences" section. Fill in the desired investment amount, investment timeframe, and risk appetite.

3. Click the "Generate Recommendations" button to generate predictions and recommendations based on your preferences.
4. The application will display the predicted operating status of startups based on the selected features and preferences. It will also provide investment recommendations and additional information about the recommended startups.
5. Explore the visualization and insights section to gain further understanding of the data and the factors influencing startup operating status.
6. You can modify your investor preferences and click the "Generate Recommendations" button again to obtain updated predictions and recommendations.

### **Acknowledgments**

The dataset used in this project is provided by Crunchbase. Special thanks to the open-source community for their valuable contributions to the libraries and tools used in this project.

### **Contributors**

This project was developed as part of [Data Science Bootcamp] at [Moringa School]. Contributions and suggestions are welcome. Feel free to submit a pull request or open an issue.

Special thanks to our Moringa School Data Science Technical Mentors for their guidance throughout the project. We would also like to acknowledge the contributions of the Elites team members:

1. Florence Nguuni
2. Joel Omondi
3. Isaac Muturi
4. Kennedy Juma
5. Diana Mwaura
6. Brian Chabari