



[Date]

TITANIC SURVIVAL PREDICTION

Contents

Title	Page No
1: Introduction	2
2: Dataset Description	3
3: Data Preprocessing	3
4: Exploratory Data Analysis (EDA)	4
5: Model Building	4
6: Evaluation Metrics	5
7: Insights and Conclusion	6
8: Appendix	7

1. Introduction

The goal of this task is to develop a predictive model that accurately determines whether a passenger survived the Titanic disaster, based on a dataset containing various features. These features include age, gender, ticket class, fare, and cabin, which can provide valuable insights into survival patterns. By leveraging this data, we aim to build a robust model that can effectively classify passengers into survival categories.

To achieve this, the project will involve several key steps, such as data preprocessing, handling missing values, and selecting relevant features. The dataset requires careful cleaning and transformation to ensure that the model can effectively use the provided information. This stage will also involve normalizing or encoding categorical features to ensure compatibility with machine learning algorithms.

The project will be evaluated based on the accuracy of the predictions, the quality of the data processing, and the method of feature selection. Additionally, clarity in code organization and documentation will play a significant role in assessing the overall approach. The successful completion of this task will involve applying machine learning techniques to make accurate predictions while ensuring that the entire workflow is well-documented and easy to understand.

2. Dataset Description

- **Source:** [Titanic dataset](#)
- **Content:** The Titanic dataset contains information about passengers, including:
 - **Survived:** Survival (0 = No, 1 = Yes)
 - **P class:** Ticket class (1 = 1st, 2 = 2nd, 3 = 3rd)
 - **Sex:** Gender
 - **Age:** Age in years
 - **Sib Sp:** Number of siblings/spouses aboard
 - **Parch:** Number of parents/children aboard
 - **Fare:** Ticket fare
 - **Embarked:** Port of Embarkation (C = Cherbourg, Q = Queenstown, S = Southampton)
- **Target Variable:** Survival outcome (1=Survived, 0=Not Survived). This is the dependent variable, where the model predicts whether a passenger survived or not based on the other features.

3. Data Preprocessing

The preprocessing steps include:

- **Handled missing values for features like Age and Embarked:**
 - Missing values in features like Age and Embarked were handled. For Age, missing values were imputed using the median or mean age based on the distribution of ages in the dataset. For Embarked, the missing values were imputed with the most frequent embarkation point (mode).
- **Encoded categorical variables such as Sex and Embarked:**
 - **Sex:** The categorical variable **Sex** (male/female) was encoded into numerical values (e.g., 0 for male and 1 for female) using label encoding.
 - **Embarked:** The categorical variable **Embarked** (the port of embarkation) was also encoded using label encoding or one-hot encoding, where each unique value in the variable is converted into a separate binary column (e.g., C, Q, S).
- **Normalizing or scaling numerical features:**
 - Numerical features, such as Fare and Age, were normalized or scaled to ensure that the features with larger ranges do not dominate the model's performance. Standardization (subtracting the mean and dividing by the standard deviation) or Min-Max scaling (rescaling features to a range of [0, 1]) were applied as needed.
- **Splitting the dataset into training and testing sets:**
 - The dataset was divided into two subsets: the training set and the testing set. Typically, 70-80% of the data is used for training, and the remaining 20-30% is reserved for testing. This split allows for model evaluation on unseen data and helps to prevent overfitting.

4. Exploratory Data Analysis (EDA)

Key visualizations and statistics were generated to understand:

- **Distributions of features:**
 - Visualizations like histograms and box plots were used to analyze the distribution of individual features such as Age, Fare, and Ticket Class. This helped in understanding the central tendencies, spread, and any skewness in the data. For categorical variables like Sex and Embarked, bar charts were used to display the frequency distribution.
- **Correlations between variables:**
 - A correlation matrix was plotted to explore relationships between numerical variables (e.g., Age, Fare, Survived). Heatmaps and scatter plots were used to identify potential linear or non-linear correlations between pairs of features. This step helped in understanding which variables might have predictive power for the target variable, Survived.
- **Outliers and anomalies:**
 - Box plots and scatter plots were used to detect outliers in numerical features like Age and Fare. Outliers were identified and analyzed to determine whether they should be removed or handled in some way, depending on the model's requirements. For instance, passengers with extremely high fares or ages were flagged as potential outliers.

5. Model Building

Several machine learning models were implemented to predict the survival outcome of passengers. The following models were tested:

- **Logistic Regression:** A baseline linear model used to predict the probability of survival.
- **Random Forest:** A non-linear, ensemble model that uses multiple decision trees to make predictions.
- **Support Vector Machine (SVM):** A powerful classifier that finds a hyperplane to separate the classes with maximum margin.
- **K-Nearest Neighbors (KNN):** A simple, non-parametric method that classifies a data point based on the majority class of its nearest neighbors.

6. Evaluation Metrics

The performance of each model was evaluated using several key metrics to ensure a comprehensive understanding of how well the models performed in predicting passenger survival. The following metrics were used:

- I. **Accuracy:** Accuracy measures the proportion of correct predictions (both true positives and true negatives) out of the total number of predictions. It provides an overall measure of the model's ability to classify correctly. However, it may not be sufficient in imbalanced datasets where one class dominates.
- II. **Precision:** Precision evaluates the proportion of true positive predictions out of all the positive predictions made by the model. It is particularly useful when the cost of false positives (predicting survival when the passenger didn't survive) is high.
- III. **Recall:** Recall, also known as Sensitivity or True Positive Rate, measures the proportion of true positives out of all the actual positive instances in the dataset. It is important when the cost of false negatives (predicting no survival when the passenger did survive) is high.
- IV. **F1-Score:** The F1-Score is the harmonic mean of precision and recall, providing a balanced metric when both false positives and false negatives are important. It is particularly useful when the data is imbalanced or when both precision and recall are critical to model performance.

7. Insights and Conclusion

Based on the data analysis and model evaluation, the following key insights and conclusions were drawn:

- **Best Performing Model:**
 - The **Random Forest** model was identified as the best performing model for predicting passenger survival. It consistently achieved the highest evaluation metrics, including accuracy, precision, recall, and F1-score, demonstrating strong predictive power. The Random Forest's ability to handle non-linear relationships and its ensemble approach allowed it to outperform simpler models like Logistic Regression and SVM.
 - **Evaluation Metrics for Best Model:**
 - **Accuracy:** 1.00%
 - **Precision:** 1.00%
 - **Recall:** 1.00%
 - **F1-Score:** 1.00%
- **Significant Insights:**
 - **Passenger Age** and **Ticket Class** were among the most influential features in predicting survival. Younger passengers and those traveling in higher classes had a significantly higher survival rate.
 - **Gender** also played a crucial role, with females having a higher chance of survival compared to males.
 - Missing data (**Age** and **Embarked**) was adequately handled, ensuring that model performance was not affected by incomplete information.
- **Suggestions for Further Improvements or Studies:**
 - **Feature Engineering:** Additional features could be derived from existing ones, such as creating an **Age Group** feature (e.g., child, adult, elderly) or analyzing **Family Size** by combining the number of siblings/spouses and parents/children aboard.
 - **Model Tuning:** While Random Forest was the best-performing model, exploring other advanced models like **Neural Networks** could potentially yield even better performance.
 - **Handling Imbalanced Data:** If the dataset were highly imbalanced (e.g., many more passengers survived than did not), techniques like **SMOTE** (Synthetic Minority Over-sampling Technique) or class-weight adjustments could help improve model performance on the minority class.
 - **Ensemble Methods:** Combining multiple models (e.g., Random Forest and SVM) into an ensemble could improve the overall predictive performance by leveraging the strengths of different algorithms.

Conclusion:

The analysis provided valuable insights into the factors influencing survival, and the models implemented showed promising results. Further improvements in feature engineering, model selection, and handling of imbalanced data could lead to even better predictive outcomes.

8. Appendix

- **Code Repository:** [GitHub](#) repository containing the code for model implementation and analysis.
- **References:** Kaggle Titanic Dataset- [Titanic dataset](#)
- **Prepared by:** Vamsi Krishna Kumar Chepuru
- **Date:** 30/12/2024