

Project Report: Titanic Survival Analysis

1. Introduction

The Titanic Survival Analysis project aims to examine the factors that influenced the survival of passengers on the RMS Titanic. By analyzing data from the Titanic disaster, we can identify key variables that impacted the likelihood of survival. This report summarizes the objectives, methodology, analysis, and findings of the project.

2. Objectives

To understand the demographic distribution of passengers on the Titanic.

To identify significant factors that influenced survival rates.

To build predictive models to estimate the probability of survival for different passenger profiles.

3. Data Description

The dataset used for this analysis is the Titanic dataset, which includes the following variables:

PassengerId: Unique ID for each passenger.

Survived: Survival indicator (0 = No, 1 = Yes).

Pclass: Ticket class (1 = First, 2 = Second, 3 = Third).

Name: Passenger's name.

Sex: Gender of the passenger.

Age: Age of the passenger.

SibSp: Number of siblings/spouses aboard.

Parch: Number of parents/children aboard.

Ticket: Ticket number.

Fare: Passenger fare.

Cabin: Cabin number.

Embarked: Port of embarkation (C = Cherbourg, Q = Queenstown, S = Southampton).

4. Methodology

Data Cleaning:

Handling missing values.

Converting categorical variables into numerical formats.

Feature engineering (e.g., creating new variables like 'FamilySize' by combining 'SibSp' and 'Parch').

Exploratory Data Analysis (EDA):

Analyzing survival rates based on different features such as gender, age, passenger class, etc.

Visualizing data using graphs and charts.

Model Building:

Splitting the dataset into training and testing sets.

Building and evaluating models such as Logistic Regression, Decision Trees, and Random Forests.

Using accuracy, precision, recall, and F1-score to assess model performance.

5. Analysis and Findings

Gender: Females had a significantly higher survival rate than males.

Passenger Class: Passengers in first class had the highest survival rates, followed by second class, with third class having the lowest.

Age: Younger passengers, especially children, had higher survival rates.

Family Size: Passengers with smaller family sizes (including those traveling alone) had better chances of survival compared to those with larger families.

Embarkation Point: Passengers who embarked from Cherbourg had higher survival rates than those who

boarded at Queenstown or Southampton.

6. Predictive Modeling

Logistic Regression: Provided good interpretability and identified significant predictors like gender, age, and passenger class.

Decision Trees: Offered a clear visualization of decision rules but were prone to overfitting.

Random Forests: Achieved the best performance metrics by reducing overfitting and capturing interactions between variables.

7. Conclusion

The Titanic Survival Analysis project reveals that gender, passenger class, age, and family size significantly influenced survival chances. Predictive models built using these features can accurately estimate the probability of survival. This analysis highlights the importance of demographic factors in survival scenarios and can be extended to other similar datasets for broader insights.

8. Recommendations

Further research could explore the interaction effects between variables in more detail.

Incorporating more advanced machine learning techniques like Gradient Boosting or Neural Networks could improve prediction accuracy.

Historical context and qualitative data can enhance the understanding of survival dynamics beyond quantitative analysis.

Appendices

Appendix A: Detailed Data Cleaning Steps

Appendix B: Code Snippets for Data Analysis

Appendix C: Model Evaluation Metrics and Confusion Matrices

References

Titanic Dataset: Kaggle Titanic Dataset

Data Analysis Techniques: Data Science Handbook by Jake VanderPlas

Machine Learning Algorithms: Introduction to Statistical Learning by Gareth James, Daniela Witten, Trevor Hastie, and Robert Tibshirani