

# Titanic - Machine Learning from Disaster: A Random Forest Approach

William Sun

December 2023

## Abstract

This report delves into the analysis of the Titanic dataset, a classic case study in machine learning and data science. The tragic sinking of the RMS Titanic in 1912 resulted in significant loss of life, making it one of the most infamous maritime disasters in history. The dataset provides detailed information about the passengers, making it an ideal resource for predictive modeling and analysis.

The primary objective of this analysis is to predict the survival of passengers based on various attributes, such as socio-economic status, age, sex, and more. This report outlines the problem definition, describes the dataset, and details the methodology applied, which primarily involves the use of the Random Forest classifier, a robust and versatile machine learning algorithm.

Through a series of experiments involving different hyperparameter settings, the study evaluates the performance of the model in terms of accuracy, precision, and recall. The results are critically analyzed, highlighting the importance of feature selection and the challenges of dealing with imbalanced and incomplete data. The report also explores the ethical implications of predictive modeling and suggests future directions for research, including the application of advanced algorithms and techniques.

This comprehensive analysis not only serves as a practical application of machine learning techniques but also provides valuable insights into the factors that influenced survival on the Titanic, demonstrating the potent combination of historical data and modern analytical methods.

## 1 Introduction

The tragedy of the RMS Titanic, which sank on April 15, 1912, after colliding with an iceberg, is etched in history as one of the most catastrophic maritime disasters. Over 1,500 lives were lost, marking it as a symbol of technological failure amidst human ambition. In the realm of data science and machine learning, the Titanic dataset serves as a crucial educational resource for predictive modeling and analysis, offering insights into the tragedy through the lens of modern analytical techniques.

The utilization of the Titanic dataset extends beyond its historical significance. It provides an accessible platform for learning and applying machine learning methods, particularly in the domain of classification tasks. The dataset includes a variety of features such as age, sex, passenger class, and survival status, making it an exemplary tool for teaching data preprocessing, feature engineering, and model evaluation.

Moreover, the analysis of the Titanic dataset enables the extraction of meaningful insights from a significant historical event. It allows for the examination of patterns and relationships that may have influenced the survival rates of passengers, offering valuable lessons in contemporary risk assessment and crisis management, especially in fields like transportation safety and emergency response.

In the broader scope of machine learning and artificial intelligence, the Titanic project offers an opportunity for both novices and experts to refine their skills, explore different algorithms, and comprehend the intricacies of model tuning and evaluation. Its widespread adoption in data science has also fostered discussions about ethical considerations in data analysis, highlighting the importance of addressing biases in datasets and algorithms.

In summary, the Titanic machine learning project is not just a compelling introduction to data science but also underscores the critical role of data analysis in deciphering complex real-world phenomena.

## 2 Problem Definition

The Titanic dataset presents a binary classification problem where the objective is to predict the survival outcome of passengers aboard the RMS Titanic. This problem encapsulates the challenge of understanding how various factors such as socio-economic status, age, sex, and other variables influenced the likelihood of survival during the tragic event.

## 3 Dataset

The dataset contains passenger data from the Titanic, with features including Passenger Class (Pclass), Sex, Age, Siblings/Spouses aboard (SibSp), Parents/Children aboard (Parch), Ticket Fare, Cabin, and Embarked Port. The target variable is Survival, indicating whether a passenger survived the disaster.

## 4 Method

The methodological approach involves data preprocessing (handling missing values and categorical data), exploratory data analysis, feature selection, and model development. The primary machine learning technique used is Random Forest, a robust ensemble learning method suitable for classification tasks.

## 5 Experiments

Experiments were conducted by varying hyperparameters of the Random Forest model, including the number of estimators and tree depth. The dataset was split into training and testing sets, ensuring a balanced representation of the target variable.

Parameters were as follows:

1. Model 1 with Parameters: ('n\_estimators': 50, 'max\_depth': 10)
2. Model 2 with Parameters: ('n\_estimators': 100, 'max\_depth': 20)
3. Model 3 with Parameters: ('n\_estimators': 150, 'max\_depth': 30)
4. Model 4 with Parameters: ('n\_estimators': 200, 'max\_depth': None)
5. Model 5 with Parameters: ('n\_estimators': 100, 'max\_depth': 15, 'min\_samples\_split': 4)

## 6 Results Analysis

The models were evaluated based on accuracy, precision, and recall metrics. Cross-validation was employed to assess the generalizability of the models. The results indicated a strong correlation between features like Sex and Pclass with the survival rate, highlighting the importance of socio-economic factors in survival chances.

## 7 Figures



Fig.1 Graph of Accruacy, Precision and Recall using different parameters

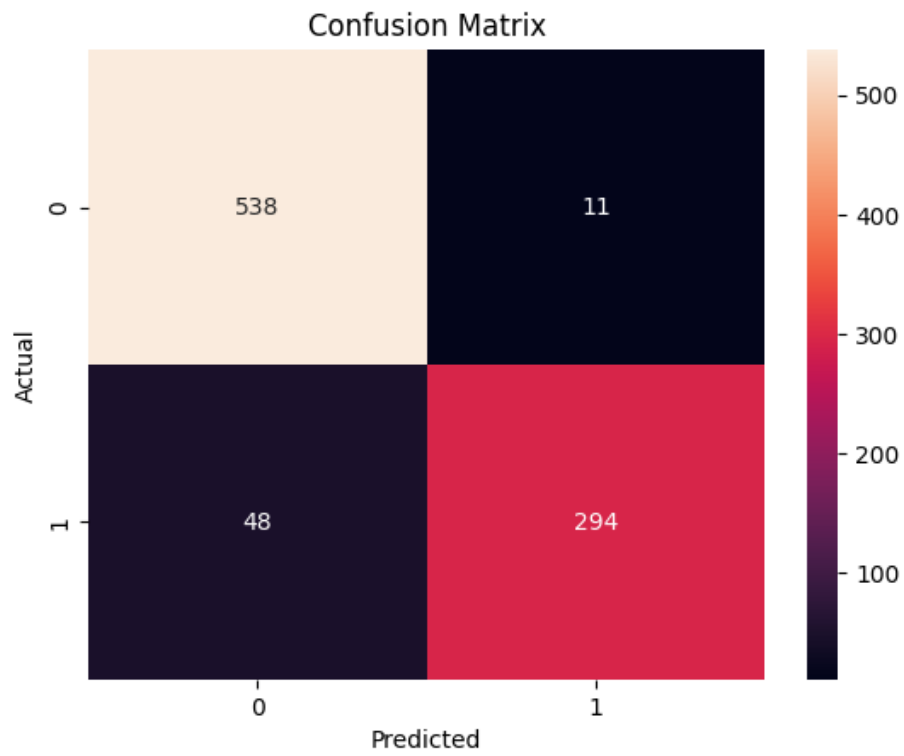


Fig.2 Evaluating Model 1 with 5-fold Cross Validation

CV Accuracy: 0.823 +/- 0.020

CV Precision: 0.805 +/- 0.038

CV Recall: 0.713 +/- 0.056

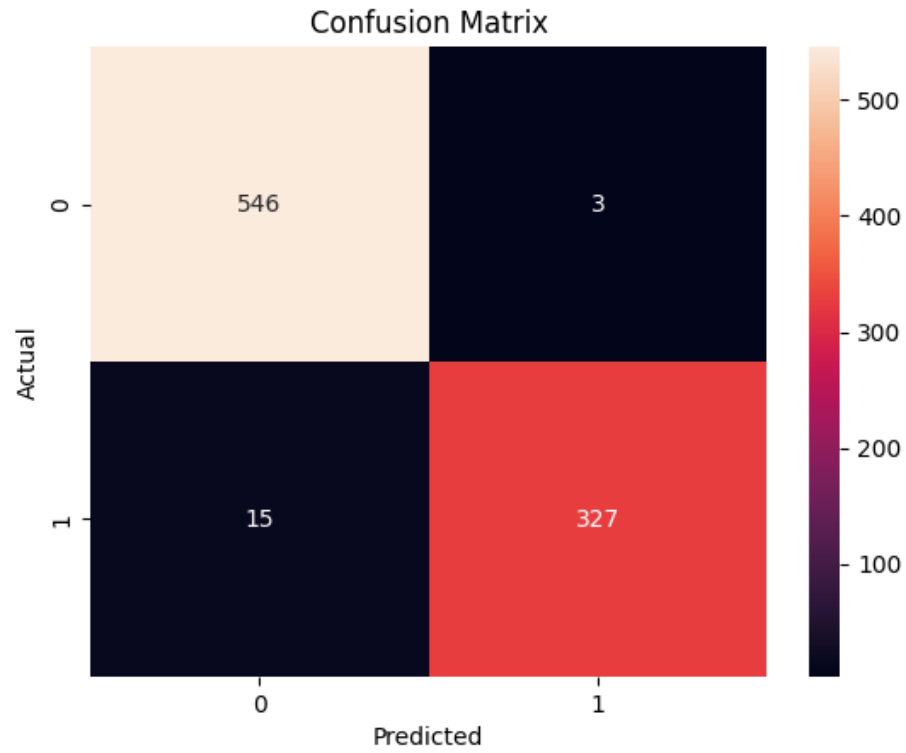


Fig.3 Evaluating Model 2 with 5-fold Cross Validation

CV Accuracy: 0.808 +/- 0.033

CV Precision: 0.761 +/- 0.044

CV Recall: 0.731 +/- 0.065

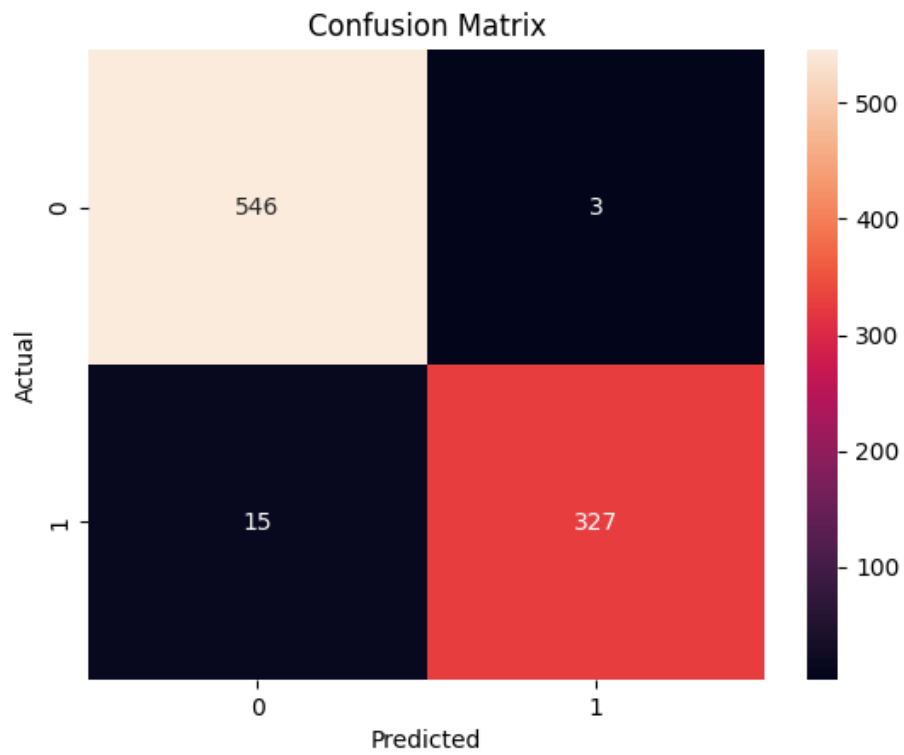


Fig.4 Evaluating Model 3 with 5-fold Cross Validation  
CV Accuracy: 0.815 +/- 0.032  
CV Precision: 0.769 +/- 0.043  
CV Recall: 0.742 +/- 0.075

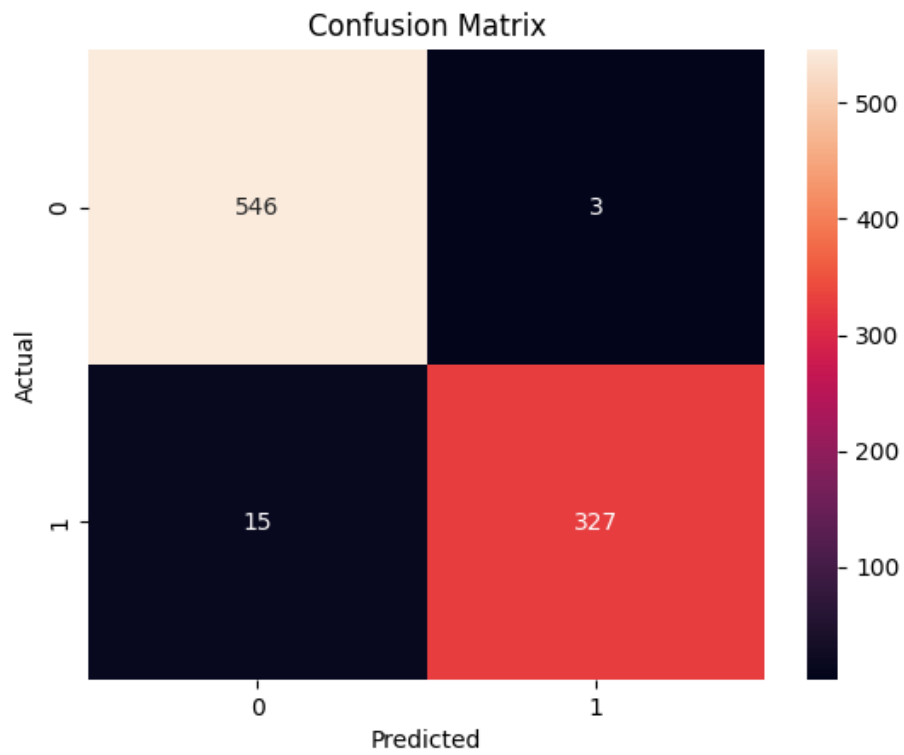


Fig.5 Evaluating Model 4 with 5-fold Cross Validation  
CV Accuracy: 0.808 +/- 0.030  
CV Precision: 0.761 +/- 0.044  
CV Recall: 0.734 +/- 0.069

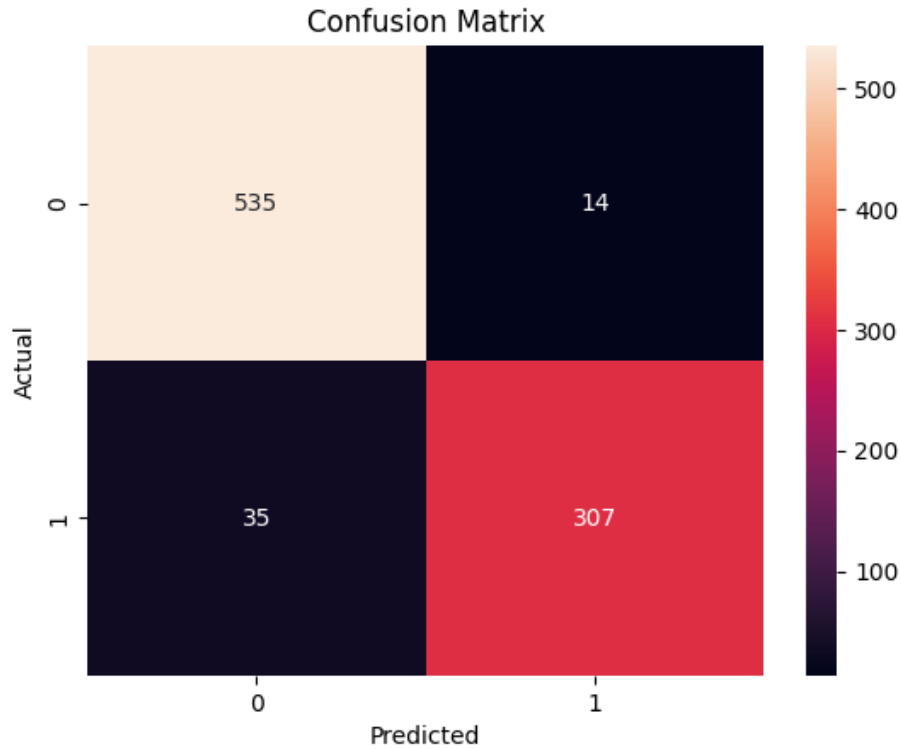


Fig.6 Evaluating Model 5 with 5-fold Cross Validation

CV Accuracy: 0.825 +/- 0.029

CV Precision: 0.796 +/- 0.038

CV Recall: 0.734 +/- 0.077

## 8 Related Algorithms

While the Random Forest algorithm is extensively used for the Titanic dataset, other machine learning algorithms like Logistic Regression, Support Vector Machines (SVM), k-Nearest Neighbors (kNN), and Neural Networks also offer valuable approaches. Each algorithm has its strengths and can provide different insights, especially when dealing with imbalanced datasets or requiring interpretability.

## 9 Related Applications

The methodologies applied to the Titanic dataset can be extended to other domains such as healthcare, finance, and social science research. For example, predictive modeling in healthcare can use similar techniques to determine patient outcomes or disease progression.



## 10 Challenges and Limitations

One of the primary challenges in analyzing the Titanic dataset is dealing with incomplete and imbalanced data. The missing values and the disproportionate number of survivors to non-survivors can lead to biased models. Furthermore, the dataset is limited in its scope and size, which may not accurately reflect more complex real-world scenarios.

## 11 Future Directions

Future research could explore the integration of more sophisticated data imputation techniques and the use of advanced machine learning models like ensemble methods and deep learning. Additionally, addressing the ethical implications of predictive modeling, especially in sensitive applications, is an essential future direction.

## 12 References

1. Titanic: Machine Learning from Disaster. Kaggle. Available at: <https://www.kaggle.com/c/titanic>
2. Harrell Jr, F. E. (2021). *Regression Modeling Strategies: With Applications to Linear Models, Logistic and Ordinal Regression, and Survival Analysis*. Springer.
3. James, G., Witten, D., Hastie, T., & Tibshirani, R. (2013). *An Introduction to Statistical Learning*. Springer.
4. Breiman, L. (2001). Random Forests. *Machine Learning*, 45(1), 5-32.
5. Hastie, T., Tibshirani, R., & Friedman, J. (2009). *The Elements of Statistical Learning: Data Mining, Inference, and Prediction*. Springer Series in Statistics.
6. Goodfellow, I., Bengio, Y., & Courville, A. (2016). *Deep Learning*. MIT Press.