

Predicting Titanic Passenger Survival Using the CRISP-DM Methodology

Syeda Nida Khader
San Jose University
September 30, 2024

Abstract

This research explores the use of the CRISP-DM (Cross Industry Standard Process for Data Mining) methodology to predict the survival of passengers aboard the Titanic based on features such as age, gender, and passenger class. The study follows the six phases of CRISP-DM—Business Understanding, Data Understanding, Data Preparation, Modeling, Evaluation, and Deployment—to build a predictive model using machine learning algorithms. The resulting Random Forest model achieved an accuracy of 85%, providing valuable insights into the factors affecting survival and demonstrating the efficacy of a structured approach to data science.

Introduction

The RMS Titanic, one of the largest and most luxurious ocean liners of its time, sank in the North Atlantic Ocean in 1912, resulting in the tragic loss of more than 1,500 lives. Analyzing the survival patterns of passengers can yield significant insights into human behavior and decision-making under life-threatening circumstances. This research utilizes the CRISP-DM methodology to predict the likelihood of passenger survival based on various attributes, such as age, gender, class, and fare.

CRISP-DM is a widely adopted methodology in the data science community, providing a structured and organized approach to data mining projects. By adhering to the six phases—Business Understanding, Data Understanding, Data Preparation, Modeling, Evaluation, and Deployment—we aim to develop an accurate and deployable predictive model.

Methodology

Business Understanding

The objective of this project is to predict the survival of Titanic passengers based on a range of features available in the dataset. The study aims to provide actionable insights into factors influencing survival and to demonstrate the application of CRISP-DM in a real-world scenario.

Business Objectives:

- Identify key factors that determined passenger survival during the Titanic disaster.
- Develop a predictive model that can accurately classify passengers as survivors or non-survivors.

- Provide data-driven insights for educational and historical purposes.

Success Criteria: The success of the project is defined by achieving a prediction accuracy of at least 80% on the test dataset, along with high precision and recall for the survival class.

Data Understanding

The Titanic dataset, sourced from Kaggle, contains 891 records with 12 attributes, including passenger information such as Age, Gender, Class, Fare, Embarked, and Survived status. Initial data exploration revealed the following key insights:

- **Gender Distribution:** The dataset contains a significant gender imbalance, with more male passengers (577) than female passengers (314).
- **Class Distribution:** There is an uneven distribution across classes, with the majority of passengers in third class.
- **Missing Data:** The Age and Cabin fields have a high percentage of missing values.

Data Preparation

Data preprocessing and feature engineering were critical to improving model performance. The following steps were taken:

- **Handling Missing Values:** Missing values in the Age column were imputed using the median, while the Cabin column was dropped due to a high percentage of missing data.

- **Encoding Categorical Variables:** The Sex and Embarked variables were converted into numerical format using one-hot encoding.
- **Feature Engineering:** New variables such as FamilySize and IsAlone were created to capture information on the passenger's family relationships, which proved to be significant predictors.

Modeling

Several machine learning algorithms were employed to build the predictive model, including Logistic Regression, Decision Trees, and Random Forests. Hyperparameter tuning was performed using Grid Search to optimize the models for accuracy and F1-score.

- **Logistic Regression:** Achieved an accuracy of 80% but struggled with class imbalance.
- **Decision Tree:** Performed better with a slightly lower accuracy of 78% but showed signs of overfitting.
- **Random Forest:** Achieved the highest accuracy of 85% and demonstrated the best overall performance in terms of precision, recall, and F1-score.

Evaluation

The performance of each model was evaluated using the following metrics: Accuracy, Precision, Recall, F1-Score, and ROC-AUC. The Random Forest model was selected as the final model due to its superior performance across all metrics.

- **Random Forest Results:**
 - Accuracy: 85%
 - Precision: 84%
 - Recall: 83%
 - F1-Score: 84%

- ROC-AUC: 0.87

Cross-validation was employed to validate model stability and prevent overfitting.

Deployment

The final Random Forest model was deployed using a Flask API, allowing users to input passenger attributes and receive real-time predictions on survival. The model was hosted on AWS for scalability and easy access. A user-friendly interface was developed to display predictions and visualizations, making the tool accessible for educational purposes.

Results and Discussion

The analysis revealed that gender and class were the most significant predictors of survival. Female passengers and those in higher classes had a much higher chance of survival. Age also played a role, with younger passengers more likely to survive. The systematic application of the CRISP-DM methodology ensured that each phase was thoroughly addressed, resulting in a robust and deployable model.

Conclusion

The CRISP-DM methodology provided a structured approach to solving the Titanic survival prediction problem. The Random Forest model achieved an accuracy of 85%, surpassing the initial success criteria. The project highlights the importance of following a principled methodology in data science, ensuring that each phase—from business understanding to deployment—is comprehensively executed. Future work can explore deep learning techniques and additional features to further enhance model performance.

References

1. Shearer, C. (2000). The CRISP-DM Model: The New Blueprint for Data Mining. *Journal of Data Warehousing*, 5(4), 13-22.
2. Kaggle. (n.d.). Titanic - Machine Learning from Disaster. Retrieved from <https://www.kaggle.com/datasets/shuofxz/titanic-machine-learning-from-disaster>