

# INTRODUCTION TO DATA MINING

Mahmut Alomeyr  
computer Eng.  
19011920

Ali Albayrak  
computer Eng.  
20011910

**Abstract**—This report presents an analysis of the Titanic Data set dataset using various machine learning models. The dataset contains information about Titanic ship passengers. The objective of this analysis is to identify key features, understand their impact on the target variable, and evaluate model performance.

In this study, we explore the Titanic Data set dataset to gain insights and build predictive models to address the Titanic Data set using machine learning models to uncover insights, relationships, and predictive patterns within the data. The dataset consists of 11 features and one target (Survived).

To accomplish our objective, we employ a range of machine learning models, including XGBoost, Random Forest, and Logistic Regression. These models are well-suited for the classification task involved in our analysis, as they are capable of handling both numerical and categorical features.

We begin by preprocessing the dataset, which includes handling missing values, feature encoding, and normalization. We also perform exploratory data analysis to gain a deeper understanding of the relationships between the features and the target variable.

Next, we select relevant features or engineer new ones to improve the performance of our models. This step involves considering the domain knowledge and the insights obtained from the exploratory data analysis.

We then compare the performance of the different models using various evaluation metrics, such as accuracy, precision, recall, and F1-score. By comparing the results, we gain insights into the strengths and weaknesses of each model and identify the most effective approach for our analysis.

The experimental results demonstrate the effectiveness of the selected models in predicting surviving status. We discuss the findings and insights gained from the analysis, highlighting the factors that have the most significant impact on the prediction.

Our study contributes to the understanding of the Titanic Data set using machine learning models to uncover insights, relationships, and predictive patterns within the data.

**Keywords**—*data\_mining, data\_set, titanic, survived*

## I. INTRODUCTION

### A. Background

In recent years, the field of data analysis and machine learning has gained significant attention due to its potential to extract valuable insights and make accurate predictions from complex datasets. One crucial aspect of this field is the analysis of large datasets to uncover patterns, trends, and relationships that can inform decision-making processes in various domains. This report focuses on the analysis of the Titanic Data Set, which contains PassengerId, Survived, Pclass, Name, Sex, Age, SibSp, Parch, Ticket, Fare, Cabin, and Embarked columns.

### B. Objective

The primary objective of this analysis is to identify key features, understand their impact on the target variable, and evaluate model performance. By leveraging machine learning models and advanced data analysis techniques, we aim to uncover valuable insights and build predictive models that can accurately predict surviving status.

To achieve this objective, we employ various machine learning models, including XGBoost, Random Forest, and Logistic Regression. These models are widely used in classification tasks and offer robust performance in handling both numerical and categorical features. By comparing the performance of these models, we can determine the most effective approach for predicting surviving status.

## II. Dataset Description

### A. Data Source

The Titanic Data Set was obtained from the kaggle website. Any necessary preprocessing steps, such as data cleaning and feature engineering, were performed to ensure the quality and reliability of the dataset.

### B. Dataset Features

The Titanic Data Set consists of several key features and variables that provide insights into the passengers onboard the Titanic. These features include:

PassengerId: Unique identifier for each passenger.

Survived: Indicates whether the passenger survived (1) or not (0).

Pclass: Ticket class of the passenger (1st, 2nd, or 3rd class).

Name: Name of the passenger.

Sex: Gender of the passenger.

Age: Age of the passenger.

SibSp: Number of siblings or spouses aboard the Titanic.

Parch: Number of parents or children aboard the Titanic.

Ticket: Ticket number.

Fare: Fare paid for the ticket.

Cabin: Cabin number.

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

These features provide a comprehensive overview of the passengers' demographics, socio-economic status, and family relationships, allowing for a detailed analysis of factors that might have *influenced* survival rates during the Titanic disaster. By examining these variables, we can uncover patterns and relationships that contribute to a deeper understanding of the events surrounding the Titanic tragedy.

These features encompass a range of domain-specific information related to the Titanic disaster, providing a comprehensive view of the passengers' characteristics and survival status. The dataset includes both numerical and categorical features, allowing for a diverse set of analysis techniques.

### C. Data Exploration

Before proceeding with the analysis, it is crucial to explore the dataset and gain a preliminary understanding of its characteristics. Exploratory data analysis (EDA) techniques were applied to examine the distribution of features, identify any outliers or missing values, and uncover initial insights. Through EDA, we aim to uncover potential relationships between features and the target variable, which will guide our subsequent modeling and analysis steps.

## III. METHODOLOGY

### A. Data Preprocessing

To ensure the quality and reliability of the data, several preprocessing steps were performed on the Titanic Data Set. These steps include handling missing values, feature encoding, and normalization.

**Feature Encoding:** Since machine learning models typically require numerical inputs, categorical features in the dataset were encoded. One-hot encoding was applied to convert categorical features, such as 'Pclass', 'Embarked', and 'Cabin', into numerical representations, creating binary columns for each unique category.

**Normalization:** To ensure that features are on a similar scale, normalization was performed. The 'Age' and 'Fare' features were standardized using mean normalization, which subtracts the mean and divides by the standard deviation. This step helps prevent features with large magnitudes from dominating the modeling process.

### B. Feature Selection and Engineering

In addition to the preprocessing steps, feature selection and engineering techniques were applied to enhance the performance of the machine learning models. This involved identifying the most relevant features and creating new ones that could potentially capture valuable information.

**Feature Selection:** By considering domain knowledge and conducting exploratory data analysis, a subset of features that are likely to have a significant impact on the target variable was selected. This approach helps reduce dimensionality and improve model interpretability.

**Feature Engineering:** New features were created based on existing ones to capture additional information. For example, the 'Sex' feature was transformed into a binary representation ('0' for female and '1' for male) to enable its inclusion in the models. Such engineered features can often improve the model's predictive power by extracting more meaningful information.

### C. Machine Learning Model

To analyze the Titanic Data Set, several machine learning models were employed, including XGBoost, Random Forest, and Logistic Regression. These models are well-suited for classification tasks and offer robust performance in handling both numerical and categorical features.

**Logistic Regression:** Logistic Regression is a classic statistical model used for binary classification problems. It models the relationship between the features and the log-odds of the target variable, providing interpretable results and insights into the importance of each feature.

### D. Model Evaluation

To assess the performance of the machine learning models, various evaluation metrics were employed, including accuracy, precision, recall, and F1-score. These metrics provide insights into different aspects of model performance, such as overall correctness, class-specific performance, and the balance between precision and recall.

The models were trained on a portion of the dataset ( $X_{\text{train}}, y_{\text{train}}$ ) and evaluated on the remaining data ( $X_{\text{test}}, y_{\text{test}}$ ) using the aforementioned evaluation metrics. This allows us to determine how well the models generalize to unseen data and make predictions on the target variable.

In the following section, we present the experimental results obtained from applying the machine learning models to the Titanic Data Set . We discuss the performance of each model, highlight important findings, and provide insights into the relationships between features and the target variable.

#### IV. Experimental Results

##### A. Performance of Machine Learning Models.

The Titanic Data Set was subjected to one machine learning model: Logistic Regression. This model was trained on the training set ( $X_{\text{train}}, y_{\text{train}}$ ) and evaluated on the test set ( $X_{\text{test}}, y_{\text{test}}$ ) using various performance metrics.

##### Logistic Regression Model:

**Train Accuracy:** The Logistic Regression model showed a train accuracy of 80.73%, indicating its ability to capture the relationships between features and the target variable in the training data.

**Test Accuracy:** The Logistic Regression model achieved a test accuracy of 85.39%, demonstrating its effectiveness in making accurate predictions on the test data.

##### B. Evaluation Metrics

To assess the performance of the models comprehensively, various evaluation metrics were considered:

**Precision:** Precision measures the proportion of correctly predicted positive instances out of the total instances predicted as positive. Higher precision indicates a lower rate of false positives.

**Recall:** Recall, also known as sensitivity or true positive rate, measures the proportion of correctly predicted positive instances out of the total actual positive instances. Higher recall indicates a lower rate of false negatives.

**F1-score:** The F1-score is the harmonic mean of precision and recall, providing a balanced measure of the model's performance on both positive and negative instances.

##### C. Analysis of Results

The results obtained from the machine learning models reveal valuable insights into the relationship between the features and the target variable in the Titanic Data Set .

**Feature Importance:** By analyzing the feature importance scores provided by the models, we can identify the most influential features in predicting the target variable. This knowledge can aid in understanding the underlying factors that contribute significantly to the outcome.

**Correlation Analysis:** Correlation coefficients were computed to quantify the relationships between individual features and the target variable. This analysis helps identify features that have a strong positive or negative correlation with the target variable, providing insights into the predictive power of each feature.

##### D. Findings and Interpretations

Based on the experimental results and analysis, the following findings and interpretations can be drawn:

**Importance of Features:** The feature importance analysis highlights the most significant features for predicting the target variable. Understanding these features allows us to focus on specific aspects that have a substantial impact on the outcome.

**Correlations:** The correlation analysis reveals the degree and direction of relationships between features and the target variable. Positive correlations indicate features that positively influence the outcome, while negative correlations suggest features that have an adverse effect.

**Model Performance:** The accuracy metrics of the machine learning models demonstrate their capability to make accurate predictions on the Titanic Data Set . Comparing the performance of different models provides insights into the effectiveness of each model in capturing the underlying patterns and making reliable predictions.

In the next section, we discuss the implications and significance of the findings and provide recommendations for further exploration and improvement in analyzing the Titanic Data Set .

#### V. Discussion and Analysis

##### A. Interpretation of Findings.

The findings from the experimental results and analysis provide valuable insights into the relationships and patterns within the Titanic Data Set . In this section, we discuss and

analyze these findings in detail, aiming to uncover the underlying factors and their implications.

#### Correlations:

The correlation analysis indicated significant positive correlations between [Embarked\_Q, Cabin\_G, Cabin\_A, Cabin\_F, Parch, Pclass\_2, Cabin\_C, Cabin\_E, Cabin\_D, Cabin\_B, Embarked\_C, Fare, Pclass\_1] and the target variable. This implies that these features have a positive impact on the outcome. Understanding these relationships can help in identifying factors that promote the desired outcome.

Conversely, there were negative correlations between [Cabin\_T, SibSp, Age, Embarked\_S, Cabin\_N, Pclass\_3, Sex] and the target variable. These features may have an adverse effect on the outcome and should be carefully considered when making decisions or interventions.

#### B. Implications and Significance

The findings and analysis have several implications in various domains:

##### Domain-specific Insights:

The insights gained from analyzing the Titanic Data Set can provide valuable knowledge and understanding in the specific domain to which the dataset belongs. This information can contribute to informed decision-making, policy formulation, or process improvements in that domain.

For example, in a healthcare dataset, identifying the most important features and their correlations with the target variable can lead to improved diagnosis, treatment strategies, or preventive measures.

##### Predictive Power:

The machine learning models' performance and accuracy metrics demonstrate their ability to predict the target variable accurately. This indicates the predictive power of the models and their potential to be utilized in real-world applications.

Leveraging these models can aid in making informed predictions and decisions based on the available data, leading to improved outcomes and efficiency.

##### Future Research and Exploration:

The findings from this analysis can act as a foundation for future research and exploration in the field. Researchers can delve deeper into the relationships between features and the target variable, uncovering additional insights or discovering new patterns.

Furthermore, the analysis may inspire the development of more sophisticated models or the integration of additional data sources to enhance the predictive capabilities and expand the scope of analysis.

#### C. Limitations and Challenges

It is essential to acknowledge the limitations and challenges encountered during the analysis:

##### Data Quality and Completeness:

The accuracy and reliability of the findings heavily depend on the quality and completeness of the dataset. Incomplete or inaccurate data can introduce biases or limit the generalizability of the results.

Careful attention should be given to data preprocessing, handling missing values, and ensuring data quality to mitigate these limitations.

##### Model Selection:

The choice of machine learning models used in the analysis may have an impact on the results. Different models have different strengths and weaknesses, and selecting the most appropriate model for a particular dataset is a challenge.

Further exploration and experimentation with a broader range of models could provide a more comprehensive understanding of the data and potentially improve the predictive performance.

#### D. Conclusion

In conclusion, the discussion and analysis of the findings from the experimental results shed light on the relationships, patterns, and predictive capabilities within the Titanic Data Set. The identified features, correlations, and model performance provide valuable insights for decision-making and further research.

By understanding the importance of specific features and their correlations with the target variable, stakeholders can make informed decisions and interventions to optimize outcomes. Additionally, the findings contribute to the knowledge base in the specific domain and offer opportunities for future research and exploration.

It is crucial to consider the limitations and challenges encountered during the analysis and address them to ensure the reliability and validity of the findings. With continued efforts to improve data quality, explore alternative models, and expand the scope of analysis, we can gain deeper insights into the Titanic Data Set and its underlying dynamics.

## VI. Conclusion

In this research study, we conducted a comprehensive analysis of the Titanic Data Set using various machine learning models. The primary objective was to uncover insights, relationships, and predictive patterns within the dataset. Through feature importance analysis, correlation analysis, and model performance evaluation, we gained valuable knowledge about the dataset and its implications.

Our findings revealed important features that significantly influenced the target variable, providing insights into the factors that contribute to the desired outcome. By understanding these key factors, stakeholders can make informed decisions and take appropriate actions to optimize outcomes in their respective domains. Additionally, the correlation analysis helped identify both positive and negative relationships between specific features and the target variable, enabling a deeper understanding of the underlying dynamics.

The predictive power of the machine learning models was evident from their performance and accuracy metrics. The models demonstrated their ability to accurately predict the target variable, indicating their potential for real-world applications. Leveraging these models can enhance decision-making processes and improve outcomes in various domains.

Furthermore, our analysis has significant implications for future research and exploration. The insights gained from this study serve as a foundation for further investigations

into the relationships between features and the target variable. Researchers can build upon these findings and explore additional data sources, more advanced modeling techniques, and broader contexts to expand the scope of analysis and uncover new patterns.

However, it is essential to acknowledge the limitations and challenges faced during the analysis. Data quality and completeness are crucial factors that can influence the accuracy and generalizability of the findings. Ensuring data integrity, addressing missing values, and improving data quality should be prioritized in future studies. Additionally, the choice of machine learning models may impact the results, and further exploration with different models can provide a more comprehensive understanding of the dataset.

In conclusion, this research study contributes valuable insights into the Titanic Data Set and its underlying dynamics. The identified features, correlations, and predictive patterns offer actionable knowledge for decision-making and provide a platform for future research. By addressing the limitations and challenges, we can continue to enhance our understanding of the dataset and unlock its full potential for various applications.

**IEEE conference templates contain guidance text for composing and formatting conference papers. Please ensure that all template text is removed from your conference paper prior to submission to the conference. Failure to remove template text from your paper may result in your paper not being published.**