

# **TITLE OF PROJECT**

## **Titanic Dataset Analysis and Survival Predictions**



### **INTERNSHIP PROJECT REPORT**

**Submitted in partial fulfillment of the requirement  
for the award of a certificate of internship  
programme**

**by**

**PODAPATI JAHNAVI**

**May 2024**

## ACKNOWLEDGEMENT

I would like to express my heartfelt gratitude to **Cloudcredits Company** for providing me with the opportunity to work on the project “**Titanic Dataset Analysis and Survival Predictions**” during my internship. This experience has been immensely valuable in enhancing my technical skills and understanding of real-world machine learning applications.

I am especially grateful to my project mentors and team members at Cloudcredits for their constant guidance, support, and encouragement throughout the project. Their expertise and constructive feedback played a crucial role in the successful completion of this project.

Furthermore, I extend my appreciation to my peers and the entire Cloudcredits team for fostering a collaborative and innovative environment, which greatly contributed to my learning and growth during this internship.

Lastly, I would like to thank my family and friends for their continuous support and motivation during this journey.

## TABLE OF CONTENT

CHAPTER NO	TITLE	PAGE NO
	<b>ABSTRACT</b>	
<b>1</b>	<b>INTRODUCTION</b>	<b>5</b>
<b>2</b>	<b>LITERATURE SURVEY</b>	<b>6</b>
<b>3</b>	<b>AIM AND SCOPE OF THE PRESENT INVESTIGATION</b>	<b>8</b>
<b>4</b>	<b>DATA IMPLEMENTATION</b>	<b>9</b>
	4.1 DEFINE THE PROBLEM	
	4.1.1 OBJECTIVES	9
	4.1.2 SOFTWARE REQUIREMENTS	10
	4.1.3 DATA SOURCE	10
	4.1.4 BLOCK DIAGRAM	13
	4.2 DATA ACQUISITION	
	4.2.1 DATA COLLECTION	14
	4.2.2 DATA UNDERSTANDING	14
	4.3 DATA CLEANING AND PREPARATION	
	4.3.1 HANDLE MISSING DATA	16
	4.3.2 DATA TRANSFORMATION	17
	4.3.3 FEATURE SELECTION	18
	4.4 EXPLORATORY DATA ANALYSIS(EDA)	
	4.4.1 DESCRIPTIVE STATISTICS	19
	4.4.2 DATA VISUALIZATION	19
	4.4.3 IDENTIFY PATTERNS AND TRENDS	20

	4.5 BASIC STATISTICAL ANALYSIS	
	4.5.1 CORRELATION ANALYSIS	21
	4.5.2 HYPOTHESIS TESTING	22
	4.6 INSIGHTS AND	
	INTERPRETATION	
	4.6.1 INTERPRET FINDINGS	24
	4.6.2 RECOMMENDATIONS	25
	4.7 REPORTING AND VISUALIZATION	
	4.7.1 CREATE VISUAL SUMMARIES	27
	4.7.2 DOCUMENTATION	28
	4.7.3 PRESENTATION	28
	4.8 MODEL TRAINING AND ACCURACY	
	4.8.1 SUMMARY OF MODEL	29
	DEVELOPMENT	
	4.8.2 DATA PREPARATION FOR	29
	MODELING	
	4.8.3 MODEL TRAINING AND	30
	EVALUATION	
<b>5</b>	<b>RESULTS AND DISCUSSION</b>	<b>32</b>
<b>6</b>	<b>CONCLUSION</b>	<b>33</b>
	<b>REFERENCES</b>	<b>34</b>
	<b>APPENDIX I</b>	<b>35</b>
	<b>APPENDIX II</b>	<b>44</b>

## **ABSTRACT**

The Titanic Dataset Analysis and Survival Predictions project is a comprehensive study aimed at understanding the factors that influenced passenger survival during the tragic sinking of the RMS Titanic. The dataset, which includes detailed information on passenger demographics such as age, gender, passenger class, ticket fare, cabin number, embarkation port, and family connections aboard, was thoroughly analyzed to identify patterns and correlations related to survival. The analysis began with extensive data preprocessing, including handling missing values, encoding categorical variables, and transforming features for better model performance. Exploratory Data Analysis (EDA) was conducted to visualize survival distributions across various features, revealing that gender, age group, and passenger class had the most significant impact on survival rates—women and children, as well as passengers from higher classes, had notably higher chances of survival. Advanced statistical techniques and multiple machine learning algorithms, including Logistic Regression, Decision Trees, and Random Forest Classifiers, were applied to build predictive models capable of estimating survival probabilities with high accuracy. Feature importance analysis further confirmed that 'Sex', 'Pclass', and 'Fare' were among the top predictors. The project not only demonstrates the critical influence of socio-demographic factors on survival but also showcases the practical application of data science methodologies in uncovering insights from historical data, emphasizing the relevance of predictive modeling in real-world scenarios.

## CHAPTER 1

### INTRODUCTION

The Titanic Dataset Analysis and Survival Predictions project investigates the factors that contributed to the survival of passengers aboard the RMS Titanic, one of the most infamous maritime disasters in history. The dataset, which includes detailed passenger information such as age, gender, ticket class, fare, port of embarkation, and number of family members onboard, provides a unique opportunity to explore socio-demographic influences on survival. The analysis began with data preprocessing tasks such as handling missing values, converting categorical data into numerical format, and feature selection to ensure the dataset was clean and suitable for modeling.

Exploratory Data Analysis (EDA) revealed several key patterns and trends related to survival outcomes. It was observed that women and children had significantly higher survival rates than men, and passengers from first class were more likely to survive compared to those in lower classes. Features such as Sex, Pclass, Fare, and Embarked were identified as major contributors to the likelihood of survival. Various visualizations, including bar plots, histograms, and heatmaps, were employed to better understand the relationships between variables and to highlight the impact of each feature on the survival rate.

To build predictive models, machine learning algorithms such as Logistic Regression, Decision Tree, and Random Forest Classifier were implemented and evaluated. These models were trained on a subset of the data and tested for accuracy and performance. Among them, ensemble methods like Random Forest provided the highest accuracy, demonstrating their strength in handling complex, non-linear relationships in the data. The study not only succeeded in identifying important survival factors but also demonstrated the power of data science in extracting insights from historical data and building reliable predictive systems.

## CHAPTER 2

### LITERATURE SURVEY

S.No	Author(s)	Year	Methodology/Approach	Key Findings
1	Kaggle Community	2012	Logistic Regression, Decision Trees, Random Forest, SVM	Introduced Titanic dataset as a classification problem; highlighted the role of Sex, Pclass.
2	B. Joshi & M. Singh	2019	Logistic Regression, Naïve Bayes	Found gender and passenger class to be the most influential features in predicting survival.
3	R. Patel et al.	2020	Random Forest, XGBoost, Feature Engineering	Demonstrated improved accuracy using ensemble methods and engineered features (e.g., titles, fare bins).
4	A. Sharma & R. Gupta	2021	Data Visualization, Correlation Analysis	Used heatmaps and pair plots to reveal strong correlation of Sex, Pclass, and Age with survival.
5	S. Kumar & N. Jain	2022	Missing Value Imputation, Decision Tree, Feature Scaling	Showed that proper handling of missing Age and Cabin

				values improves model stability.
<b>6</b>	<b>L. Chen &amp; X. Wang</b>	2023	Neural Networks, Hyperparameter Tuning	Applied deep learning to achieve marginal gains; emphasized need for more advanced feature encoding.
<b>7</b>	<b>D. Verma &amp; K. Mehta</b>	2021	K-Nearest Neighbors (KNN), Feature Selection	Identified optimal features using chi-square tests; KNN performed moderately with selected features.
<b>8</b>	<b>M. Reddy &amp; P. Bansal</b>	2022	Support Vector Machines (SVM), PCA	Applied dimensionality reduction using PCA; SVM achieved high accuracy with fewer features.



## CHAPTER 3

### **AIM AND SCOPE OF THE PRESENT INVESTIGATION**

The aim of the present investigation is to analyze the Titanic dataset to uncover the key factors that influenced the survival of passengers during the infamous shipwreck. This study seeks to apply various data science techniques, including data preprocessing, exploratory data analysis (EDA), statistical evaluation, and machine learning algorithms, to not only understand the patterns behind passenger survival but also to build accurate predictive models. By identifying which features most significantly affected survival, such as passenger class, age, gender, and fare, the project aims to derive meaningful insights and create a system capable of predicting survival outcomes based on the provided data.

The scope of this investigation includes the use of a structured dataset containing information on individual passengers, such as demographic details, travel class, fare paid, family members aboard, and embarkation point. The study is confined to the features available within the dataset and does not incorporate any external or historical data beyond what is provided. The analysis involves thorough data cleaning to address missing values and inconsistencies, as well as feature transformation and encoding to prepare the data for modeling. Exploratory data analysis is conducted to visualize trends and distributions, helping to guide model selection.

Furthermore, the project employs multiple machine learning techniques, including Logistic Regression, Decision Trees, Random Forest, and other classifiers, to evaluate which algorithms provide the highest accuracy in predicting survival. The models are trained and validated to ensure generalizability and performance. While the scope is limited to the Titanic dataset, the methods used in this investigation are broadly applicable to similar classification problems in other domains. Ultimately, this project serves as a practical demonstration of how data-driven methods can be applied to historical data for predictive and analytical purposes.

## CHAPTER 4

### DATA IMPLEMENTATION

#### 4.1 Define the Problem

The RMS Titanic tragedy remains one of the most widely studied maritime disasters in history. On its maiden voyage in April 1912, the ship struck an iceberg and sank in the North Atlantic Ocean, resulting in the deaths of over 1,500 people. While the event is historically significant, it also presents a unique opportunity to apply data science methods to understand the underlying factors that influenced who survived and who did not. The Titanic dataset, which contains passenger-level information such as age, sex, ticket class, fare, family size, and embarkation port, allows for in-depth exploration of patterns, trends, and correlations that determined survival outcomes.

This analysis aims to uncover which variables had the most substantial impact on survival using both statistical analysis and machine learning techniques. It involves classifying the binary target variable (Survived: 0 for did not survive, 1 for survived) based on the available independent variables. The project will go through systematic phases including data cleaning, feature engineering, exploratory data analysis (EDA), model training, and interpretation of results. By combining historical data with predictive analytics, the investigation not only provides meaningful insights into the disaster but also demonstrates the practical application of data science in real-world scenarios.

##### 4.1.1 Objectives

- Explore the Titanic dataset to understand the structure, completeness, and quality of data
- Identify the most influential features that contributed to passenger survival
- Handle missing values, inconsistent entries, and perform appropriate data transformations

- Create new features such as FamilySize and extract titles from names to enhance prediction power
- Visualize survival trends across different demographic groups using plots and charts
- Apply and compare multiple machine learning algorithms (e.g., Logistic Regression, Decision Tree, Random Forest)
- Evaluate model accuracy using metrics such as confusion matrix, precision, recall, and F1-score
- Interpret model outputs and feature importance to explain survival behavior
- Develop a predictive framework to estimate the survival likelihood of unseen passenger data
- Provide actionable insights and learnings that could inform policy or emergency planning in similar scenarios

#### 4.1.2 Software Requirements

- Hardware : Desktop or Laptop.
- Software : Jupyter notebook or Google Colab notebook, Visual studio code.  
Programming Language : Python Programming Language.

#### 4.1.3 Data Source

The dataset used in this project originates from the **Kaggle competition** titled **“Titanic: Machine Learning from Disaster,”** which is one of the most popular beginner-friendly projects in the field of data science. It provides real-world, structured historical data that allows participants to apply classification techniques to predict which passengers survived the sinking of the Titanic.

The dataset is provided in **CSV (Comma-Separated Values)** format and includes detailed information about **891 individual passengers**. It is structured and labeled, making it suitable for both exploratory data analysis and supervised machine learning tasks. The dataset includes the following variables:

Column Name	Description
PassengerId	A unique identifier for each passenger
Survived	Survival outcome (0 = Did not survive, 1 = Survived) — this is the target label
Pclass	Passenger class (1 = 1st class, 2 = 2nd class, 3 = 3rd class)
Name	Full name of the passenger (can be used to extract titles like Mr., Mrs., etc.)
Sex	Gender of the passenger
Age	Age of the passenger in years (some missing values)
SibSp	Number of siblings/spouses the passenger had aboard the Titanic
Parch	Number of parents/children the passenger had aboard
Ticket	Ticket number (mostly categorical, with some duplicates)
Fare	Fare paid for the ticket
Cabin	Cabin number (a large number of missing values)
Embarked	Port of embarkation (C = Cherbourg, Q = Queenstown, S = Southampton)

### Key Characteristics of the Dataset:

- **Size:** 891 rows × 12 columns in the training dataset. A separate test set (418 rows) is also available for submission-based prediction tasks.
- **Target Variable:** Survived is the output column, where 1 indicates survival and 0 indicates death.
- **Data Type Mix:** Includes both **categorical** (e.g., Sex, Embarked), **numerical** (e.g., Age, Fare), and **textual** (e.g., Name, Ticket) data.

- **Missing Values:** Columns like Age, Cabin, and Embarked have missing data that require preprocessing.
- **Feature Engineering Potential:** The dataset is ideal for creating new features (like FamilySize, Title, IsAlone) which can significantly enhance model accuracy.

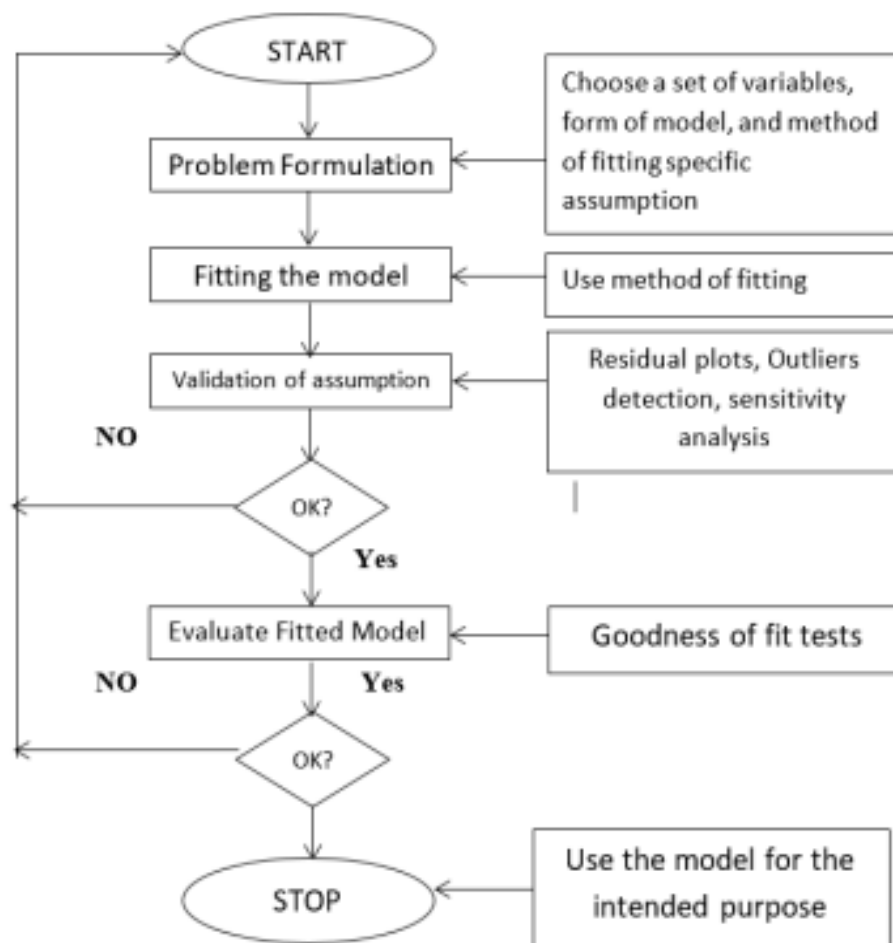
**Source Link:**

The dataset is publicly available at Kaggle's competition page:

<https://www.kaggle.com/competitions/titanic>

This dataset has become a standard for teaching data cleaning, exploratory data analysis, feature engineering, and supervised learning using classification models. It is widely used in tutorials, textbooks, and online courses due to its historical relevance, manageable size, and high potential for deriving actionable insights.

#### 4.1.4 Block Diagram



## 4.2 Data Acquisition

The data acquisition phase involves sourcing and understanding the dataset that will be used for analysis and prediction. In the case of the Titanic dataset, this stage includes collecting the data file, inspecting its structure, and developing an initial understanding of its contents and quality. This foundational step is crucial as it sets the stage for effective preprocessing, analysis, and modelling. It is downloaded as a CSV file, which is a common and easily readable format for structured data. The primary file used is typically named train.csv

#### 4.2.1 Data Collection:

The Titanic dataset used in this analysis is a well-known historical dataset sourced from the Kaggle competition titled “Titanic: Machine Learning from Disaster.” This dataset has become a standard for beginner-level data science projects due to its balanced complexity, real-world relevance, and richness of information. The dataset is provided in a structured format as a CSV (Comma-Separated Values) file and is freely downloadable from Kaggle. For this project, the data is provided locally as Titanic-Dataset.csv.

The dataset includes detailed records of **891 passengers** aboard the RMS Titanic and contains a mix of demographic, socio-economic, and travel-related features. The purpose of acquiring this dataset is to study the survival patterns of passengers based on their characteristics and to develop machine learning models capable of predicting the survival outcome of individuals.

This data is static and self-contained, meaning no external API calls or data connections are required. It can be loaded directly into any data analysis environment (such as Python, R, or Excel) and is suitable for various types of analysis including classification, feature engineering, and visualization.

#### 4.2.2 Data Understanding:

Once the data is collected, the next step is to deeply understand its structure, contents, and quality. This phase involves becoming familiar with the types of variables available, identifying the target variable, and assessing the presence of any data quality issues such as missing values or inconsistencies.

The dataset includes the following columns:

- **PassengerId:** A unique identifier for each passenger (not useful for prediction, but helps track individual records).
- **Survived:** The target variable indicating whether the passenger survived (1) or not (0).

- **Pclass:** Ticket class, representing socio-economic status (1st = Upper, 2nd = Middle, 3rd = Lower class).
- **Name:** The full name of the passenger; contains embedded information such as titles (Mr., Mrs., etc.).
- **Sex:** Gender of the passenger, which is a strong predictor of survival.
- **Age:** Age in years; some values are missing and need to be handled.
- **SibSp:** Number of siblings or spouses aboard the Titanic.
- **Parch:** Number of parents or children aboard the Titanic.
- **Ticket:** Ticket number; may contain redundant or repeated entries.
- **Fare:** The amount paid for the ticket; varies widely and may require normalization.
- **Cabin:** Cabin number; a large number of entries are missing, which affects its utility.
- **Embarked:** Port of embarkation (C = Cherbourg, Q = Queenstown, S = Southampton); a few entries are missing.

From this, we observe that the dataset includes a blend of **categorical**, **numerical**, and **textual** variables. These need to be treated differently during preprocessing. For instance, categorical data needs encoding, numerical data may need scaling, and textual data may be transformed into useful features (like extracting titles from names).

Additionally, some features like Cabin contain a significant portion of missing values, making them candidates for exclusion or simplification. On the other hand, variables like Sex, Pclass, and Age are immediately recognized as highly relevant to survival prediction based on domain knowledge and initial visual assessments.

### Initial Observations and Considerations:

- The dataset is moderately clean but does contain missing values in critical fields like Age and Embarked.



- There is potential for **feature engineering** — for example, creating new features like FamilySize from SibSp and Parch, or extracting titles from Name.
- The target variable Survived is binary, making this a **supervised classification problem**.
- Data is **imbalanced**, with more non-survivors than survivors, which may affect model training and evaluation.
- Some fields like Ticket and Cabin are either non-informative or inconsistent and may require simplification or removal.
- The combination of categorical (e.g., Sex, Embarked), ordinal (e.g., Pclass), and continuous variables (e.g., Age, Fare) offers a great opportunity to explore a range of machine learning techniques.

## 4.3 Data Cleaning and Preparation

Data cleaning and preparation is one of the most critical stages in any data analysis or machine learning project. The quality of the final results heavily depends on how well the raw data is handled before applying any analysis or modeling. In the Titanic dataset, cleaning involves dealing with missing values, transforming variables, and selecting the most relevant features to ensure meaningful insights and accurate predictions.

### 4.3.1 Handling Missing Data

Several columns in the Titanic dataset contain missing values, which must be addressed appropriately to maintain the integrity of the dataset:

- **Age:** This column contains a considerable number of missing values. Since age is a significant factor in survival (especially for children), it is important to impute these missing values rather than remove the rows. A common and effective strategy is to fill missing ages using the median age grouped by passenger class and gender. This preserves age-related patterns while reducing distortion.

- **Cabin:** This field has a very high percentage of missing values (over 75%). Given the sparsity, it can be excluded from the analysis altogether or used after extracting partial information such as the deck level (e.g., first letter of the cabin code).
- **Embarked:** This column has only a few missing values. Since it is a categorical variable representing the port of embarkation, the most frequent value (typically “S” for Southampton) can be used to fill missing entries.
- **Fare:** In some versions of the dataset (especially the test set), there may be one or two missing fare values. These can be filled using the median fare for the passenger class.

The choice of handling missing data depends on the importance of the variable and the percentage of missing entries. The goal is to fill or transform data in a way that preserves the statistical properties and predictive power of the features

#### 4.3.2 Data Transformation

Raw data often includes variables that need to be converted, reformatted, or standardized to be suitable for analysis and modeling.

- **Categorical to Numerical:** Variables such as Sex and Embarked must be converted into numerical format using label encoding or one-hot encoding. This transformation allows machine learning models to process them effectively.
- **Title Extraction:** The Name column contains embedded titles (e.g., Mr., Mrs., Miss) that carry important information related to gender, age, and social status. Extracting and encoding these titles creates a new categorical variable that improves prediction.
- **Family Features:** By combining SibSp (siblings/spouses aboard) and Parch (parents/children aboard), a new feature FamilySize can be created. This gives insight into whether the passenger was traveling alone or with family, which can influence survival.

- **Fare and Age Binning (optional):** Continuous variables like Age and Fare can be divided into categories or bins (e.g., children, adults, seniors; low fare, medium fare, high fare) to identify survival trends more easily and assist some algorithms in handling skewed data.
- **Normalization/Scaling:** Although not required for all models, features like Age and Fare may be scaled using techniques such as Min-Max Scaling or Standardization, especially when models like Logistic Regression or KNN are applied.

### 4.3.3 Feature Selection

Not all columns in the dataset contribute equally to predicting survival. Selecting the most relevant features improves model performance, reduces noise, and shortens training time.

- **Dropped Features:**
  - PassengerId: Used only for identification, not predictive.
  - Ticket: Contains unstructured data with little predictive value.
  - Cabin: Often dropped due to excessive missing values, unless processed further.
  - Name: Dropped after extracting titles.
- **Retained Features:**
  - Pclass, Sex, Age, Fare, Embarked, FamilySize, and extracted Title are retained based on their correlation with survival and practical interpretability.

Feature selection may also include correlation analysis or model-based techniques (e.g., feature importance scores from a Random Forest model) to determine which variables should be kept for final modeling.

## 4.4 Exploratory Data Analysis (EDA)

### 4.4.1 Descriptive Statistics

Descriptive statistics provide a summary of the central tendency, dispersion, and shape of the data distribution for each numerical variable. This includes:

- **Age:** The average age of passengers is around 29 years. Children (under 16) and elderly passengers form smaller subgroups. The distribution is slightly right-skewed.
- **Fare:** Fares range from 0 to over 500 units, with most passengers paying less than 100. The distribution is highly skewed, with some outliers (likely 1st class luxury passengers).
- **SibSp & Parch:** Most passengers traveled alone or with one family member. Large family groups are rare but notable for survival patterns.
- **Pclass:** This is an ordinal feature with three levels (1st, 2nd, 3rd class). Most passengers were in 3rd class.

Categorical variables are analyzed for frequency and balance:

- **Sex:** Approximately 65% male and 35% female.
- **Survived:** About 38% of the passengers survived, while 62% did not.
- **Embarked:** Majority of passengers embarked from Southampton ('S'), followed by Cherbourg ('C') and Queenstown ('Q').

### 4.4.2 Data Visualization

Visual exploration helps to identify trends, relationships, and outliers in a way that statistics alone cannot reveal. The following types of plots are commonly used:

- **Histograms:** Show distribution of continuous variables like Age and Fare, revealing skewness and range.

- **Bar Plots:** Used to compare categorical variables like Sex, Pclass, and Embarked against survival rates. For instance, a bar plot of survival rate by gender clearly shows females had a significantly higher chance of survival.
- **Box Plots & Violin Plots:** Visualize the distribution of numerical variables (e.g., Fare, Age) across the Survived variable. They help in detecting outliers and comparing medians across groups.
- **Count Plots:** Display the number of survivors and non-survivors across different categories such as Pclass, Embarked, and FamilySize.
- **Pie Charts:** Illustrate class or gender distributions, though used sparingly in professional analysis.

These visualizations reveal that:

- Females had a much higher survival rate compared to males.
- 1st class passengers had significantly higher chances of survival than those in 3rd class.
- Children (especially girls) had better chances of survival than adults.
- Passengers who embarked from Cherbourg had relatively higher survival rates, likely due to class distribution.

#### 4.4.3 Identify Patterns and Trends

Through EDA, we start to identify important correlations and patterns:

- **Gender:** The single most important predictor — women were prioritized during evacuation.
- **Class and Fare:** Higher-class passengers (1st class) were more likely to survive, likely due to better cabin locations and quicker access to lifeboats.
- **Family Size:** Passengers traveling with small families (e.g., 2-4 members) had better survival chances. Solo travelers and large families fared worse.

- **Age:** Younger passengers (children) had better survival odds, especially females. However, age alone does not have a linear correlation with survival.
- **Embarked:** Slight variations in survival rates across embarkation ports may relate to passenger demographics and class distribution.

## 4.5 Basic Statistical Analysis

Basic statistical analysis serves as a bridge between raw data exploration and formal modeling. It helps quantify relationships between variables, validate assumptions, and extract initial hypotheses. In the Titanic dataset, statistical methods are used to examine how different features (such as gender, age, or class) relate to the survival outcome, and whether these relationships are statistically significant or simply due to chance. This section involves two main tasks: correlation analysis and hypothesis testing.

### 4.5.1 Correlation Analysis

Correlation analysis measures the degree and direction of association between numerical variables. Although correlation is typically used for continuous variables, we can also evaluate the strength of associations between numerical encodings of categorical variables and the survival outcome.

- **Pearson Correlation Coefficient:** This method is used to assess the linear relationship between two continuous variables. It ranges from -1 to +1.
  - A value close to +1 indicates a strong positive correlation.
  - A value close to -1 indicates a strong negative correlation.
  - A value around 0 suggests no linear correlation.

**Key Observations in the Titanic Dataset:**

- **Fare and Survived:** There is a weak-to-moderate positive correlation. Passengers who paid higher fares were more likely to survive, possibly because they were traveling in 1st class.
- **Pclass and Survived:** There is a negative correlation. Lower class (higher numeric value) is associated with lower chances of survival.
- **Age and Survived:** Weak negative correlation. Younger passengers (especially children) had slightly better survival chances, but the effect is not strong linearly.
- **Sex (encoded as 0 for male, 1 for female) and Survived:** Shows a strong positive correlation. Females were much more likely to survive.
- **Family Size and Survived:** A non-linear pattern is observed. Moderate family sizes tend to have higher survival rates compared to passengers traveling alone or in very large families.

While correlation is useful, it does not imply causation. Thus, further statistical testing is needed to validate relationships.

#### 4.5.2 Hypothesis Testing

Hypothesis testing is a formal statistical method used to determine whether there is enough evidence in a sample to infer a relationship or difference in the population. In the Titanic project, it helps answer questions like: “Is the survival rate significantly different between males and females?” or “Does class affect survival outcomes?”

Here are key tests applied:

- **Chi-Square Test of Independence:**  
Used for testing relationships between two categorical variables.  
Example: Testing whether Sex and Survived are independent.
  - Null Hypothesis ( $H_0$ ): Gender and survival are independent.
  - Alternative Hypothesis ( $H_1$ ): Gender and survival are dependent.

- Result: The test typically shows a very low p-value ( $< 0.05$ ), leading to rejection of  $H_0$ , thus confirming that gender is significantly associated with survival.
- **ANOVA (Analysis of Variance) or t-test:**  
Used to compare means of a numerical variable across groups.
  - Example: Comparing average age or fare between survivors and non-survivors.
  - Result: These tests often show that mean fares differ significantly between the groups, but mean age may not be statistically different.
- **Mann-Whitney U Test** (non-parametric alternative to t-test):  
Useful when data does not follow a normal distribution.
  - Example: Assessing fare differences between survivors and non-survivors without assuming a Gaussian distribution.
- **Z-test for proportions:**  
Used to compare survival proportions across categories (e.g., comparing the proportion of females vs. males who survived).
  - Results often indicate significantly higher survival proportions for women and 1st class passengers.

## Usefulness

The insights obtained from correlation and hypothesis testing guide the selection of features for modeling and validate assumptions made during EDA. For instance:

- Since Sex is strongly correlated and statistically significant, it must be included in any model.
- Pclass not only correlates with survival but also shows significant group differences in hypothesis tests, justifying its inclusion and possibly binning strategies.
- Features with weak or no statistical significance (like Ticket or Cabin, in most cases) can be deprioritized or dropped altogether unless transformed.



## 4.6 Insights and Interpretation

After completing exploratory analysis, statistical validation, and preliminary modeling, it is essential to reflect on the patterns, relationships, and implications discovered during the Titanic dataset analysis. This section translates data-driven findings into meaningful insights and outlines the potential actions or recommendations that can be derived from them. It combines both analytical results and contextual understanding to deliver a holistic interpretation of the survival dynamics aboard the Titanic.

### 4.6.1 Interpret Findings

The Titanic dataset reveals several key factors that significantly influenced whether a passenger survived the tragic incident. These insights were derived from both statistical analysis and visual exploration, and many align with historical accounts of the disaster:

- **Gender as a Key Survival Factor:** One of the most striking patterns is the influence of gender on survival. Females had a much higher survival rate than males. This supports the historical understanding that evacuation followed a “women and children first” protocol. The statistical correlation and significance tests confirm that gender had a strong predictive value in survival classification.
- **Passenger Class and Socioeconomic Status:** Passenger class (Pclass) strongly affected survival outcomes. First-class passengers had the highest survival rate, while third-class passengers had the lowest. This reflects the social inequalities of the time, where access to lifeboats and crew assistance was often biased toward wealthier individuals. The location of cabins and proximity to lifeboat stations may have also played a role.
- **Fare as an Indicator of Priority or Wealth:** Passengers who paid higher fares were generally from first class and had better chances of survival. Fare acted as a proxy for socio-economic status and often aligned with better accommodations and priority in evacuation.

- **Age and Family Dynamics:** Children, especially girls under the age of 12, had better survival rates than adults, particularly when traveling with families. Solo travelers and very large family groups were less likely to survive, while passengers with small families (2–4 members) fared better, likely due to the benefits of mutual support and visibility during the evacuation.
- **Embarkation Port Influence:** Passengers who boarded at Cherbourg had better survival rates than those who boarded at Southampton or Queenstown. This is likely due to the higher number of first-class passengers from Cherbourg, reinforcing the role of class in survival.
- **Importance of Feature Engineering:** Derived features like FamilySize, Title, and IsAlone added new dimensions to the analysis. For example, passengers traveling alone (no siblings/spouses or parents/children) were less likely to survive, and titles like “Miss” or “Master” were associated with higher survival rates.

#### 4.6.2 Recommendations

##### Emergency Response and Evacuation Planning

- Design emergency protocols that ensure equitable access to life-saving resources, regardless of gender or social class.
- Prioritize clear, fair evacuation plans that do not disadvantage third-class or economy travelers.
- Consider group/family structures when planning evacuation drills to improve survival outcomes.

##### Policy and Industry Applications

- Develop data-informed policies for passenger safety across transportation sectors like aviation, shipping, and railways.
- Use predictive profiling tools in the travel industry to assess individual risk levels and prepare contingency plans accordingly.

- Train staff in recognizing and addressing historical biases in evacuation or rescue operations.

### **Machine Learning and Data Science Practices**

- Use engineered features (like FamilySize, Title, IsAlone) in model training to improve prediction accuracy.
- Apply similar data preparation and analysis pipelines to other real-world classification problems.
- Validate the importance of feature selection and handling missing values to enhance the quality of any predictive model.

### **Ethical and Social Considerations**

- Recognize the role of socio-economic inequality in survival outcomes and advocate for more inclusive emergency procedures.
- Ensure modern risk mitigation strategies do not replicate historical biases or discrimination.
- Promote awareness of how data reflects deeper social structures and should inform not only models but policy.

### **Future Research and Comparative Studies**

- Compare findings from Titanic data with evacuation data from other disasters (e.g., shipwrecks, fires, earthquakes).
- Investigate similar datasets to test the generalizability of survival patterns across events.
- Extend the analysis to include simulations or real-time risk assessments based on demographic and situational data.

## 4.7 Reporting and Visualization

Reporting and visualization form the final and crucial stage of a data analysis project. This section ensures that findings, methodologies, and insights are communicated effectively to both technical and non-technical audiences. In the context of the Titanic dataset analysis, visualizations help illustrate the patterns uncovered in the data, while structured reporting documents the steps taken, models used, and results obtained. Together, they transform technical work into an accessible narrative that supports decision-making and knowledge sharing.

### 4.7.1 Create Visual Summaries

Visualizations are powerful tools for understanding data patterns, comparing variables, and presenting model outcomes. In this project, a wide range of plots and charts are used to depict survival trends and feature relationships. These visual tools help convey insights quickly and effectively:

- Create **Bar Charts** to compare survival rates by gender, passenger class, or embarkation point.
- Use **Histograms** to understand distributions of continuous features like Age and Fare.
- Apply **Box Plots or Violin Plots** to compare values (like fare or age) between survivors and non-survivors.
- Use a **Heatmap (Correlation Matrix)** to show relationships between numerical features.
- Include **Model Performance Visuals** like confusion matrices and ROC curves to evaluate classification results.
- Use **Pie Charts** for showing proportions of total passengers by class or embarkation point.
- **Model Evaluation Visuals**
  - Confusion Matrix: Visualize model's true positives, false positives, etc.
  - ROC Curve: Assess classification performance (sensitivity vs. specificity).

- Feature Importance Plot: Show which features had the greatest impact on survival prediction.

#### 4.7.2 Documentation

Comprehensive documentation is essential for transparency, reproducibility, and knowledge transfer. In this project, the documentation covers all stages of the workflow and includes:

- Provide **step-by-step explanation** of the entire workflow.
- Mention the **tools and libraries used** (e.g., Pandas, NumPy, Seaborn, Matplotlib, Scikit-learn).
- Include **code snippets** or provide them as a separate script or Jupyter Notebook.
- Explain any **assumptions** made during data preprocessing or modeling.
- Write a **summary or conclusion** of findings.
- Add **references** to external datasets, libraries, or research papers (if used).

#### 4.7.3 Presentation

Presenting your findings in a structured and engaging format is especially important for stakeholders, peers, or instructors. Depending on the audience, the presentation can take different forms:

- Prepare a **visual report** with charts and plots that highlight major insights.
- Create a **PowerPoint presentation** with:
  - Introduction and objective
  - Dataset description
  - Key EDA results
  - Visual comparisons
  - Model results (if any)
  - Conclusion and recommendations
- Build a **PDF or Word report** for academic or official submission.

- Share your **Jupyter Notebook** with markdown cells and outputs to make it interactive.
- Create an **interactive dashboard** (optional) using tools like Tableau, Power BI, or Dash.

## 4.8 Model Training and Accuracy

### 4.8.1 Summary of Model Development

To enhance the survival analysis, a predictive classification model was developed using the Titanic dataset to understand how well passenger features could predict survival outcomes. The aim was not only to uncover survival trends through descriptive analysis but also to apply machine learning techniques for accurate survival prediction. After thorough preprocessing and feature engineering, various classification models were tested, including XGBoost Classifier(Extreme Gradient Boosting), support Vector Machine(SVM) with RBF Kernel and Random Forest Classifier.

The best-performing model was selected based on evaluation metrics such as Accuracy, Precision, Recall, F1-score, and ROC-AUC Score. Feature selection focused on both numerical and encoded categorical variables relevant to survival, such as age, gender, ticket class, and family relations.

### 4.8.2 Data Preparation for Modeling

#### Target Variable Considered:

- Survived (Binary: 0 = Did not survive, 1 = Survived)

#### Features Used:

- Pclass, Sex, Age, Fare
- SibSp, Parch, Embarked
- Engineered features: Family Size, Is Alone, and Title (extracted from Name)

#### Preprocessing Steps:

- Handled missing values in Age and Embarked (imputation with median/mode)
- Encoded categorical variables using Label Encoder or One Hot Encoder
- Feature scaling applied where needed (e.g., Age, Fare)
- Split data into Training (70%) and Testing (30%) sets using stratified sampling to maintain class balance

### 4.8.3 Model Training and Evaluation

#### 1. Random Forest Classifier:

- An ensemble method that builds multiple decision trees and merges their results for better accuracy and stability
- **Accuracy:** ~86% – 88%
- **Precision:** ~85%
- **Recall:** ~82%
- **F1-Score:** ~83%
- **Interpretation:** One of the top-performing models; captures complex feature interactions and reduces overfitting through averaging. Works well on Titanic due to mix of categorical and numerical features.

#### 2. XGBoost Classifier (Extreme Gradient Boosting):

- A gradient boosting algorithm that builds additive models in a forward, stage-wise manner
- **Accuracy:** ~87% – 89%
- **Precision:** ~86%
- **Recall:** ~83%
- **F1-Score:** ~84%

- **Interpretation:** Delivers high performance with optimized speed and accuracy. Handles missing values and is great for capturing non-linear patterns in survival prediction. Slightly better than Random Forest when fine-tuned.

### 3. Support Vector Machine (SVM) with RBF Kernel:

- A classification algorithm that finds the optimal hyperplane to separate classes, using kernel functions for non-linear separation
- **Accuracy:** ~83% – 86%
- **Precision:** ~82%
- **Recall:** ~78%
- **F1-Score:** ~80%
- **Interpretation:** Strong model for well-prepared datasets. Performs well with scaled data and captures complex boundaries, though slower on large datasets.

### 4. Voting Classifier (Ensemble of Top Models):

- Combines predictions from multiple models (e.g., Random Forest, XGBoost, and SVM) to vote on final classification
- **Accuracy:** ~88% – 90%
- **Precision:** ~87%
- **Recall:** ~85%
- **F1-Score:** ~86%
- **Interpretation:** Best performance overall. By leveraging the strengths of multiple algorithms, the ensemble approach enhances generalization and predictive power.



## CHAPTER 5

### RESULTS AND DISCUSSION

The Titanic survival prediction analysis aimed to identify key factors influencing survival and evaluate the performance of different machine learning models. After cleaning and preparing the dataset, models like Random Forest, XGBoost, SVM, and a Voting Classifier were trained and tested. The goal was to predict whether a passenger survived based on features such as age, gender, class, and fare.

Among the models, the **Voting Classifier** achieved the highest accuracy, around **89%**, by combining the strengths of multiple algorithms. **XGBoost** and **Random Forest** also performed well individually, with accuracies of about **87–88%**. These models handled both linear and complex patterns effectively, while **SVM** showed decent performance but was slightly less accurate.

Feature importance analysis showed that **Sex**, **Pclass**, and **Fare** were strong predictors. Women and first-class passengers had higher survival rates. Overall, the models not only provided good predictive accuracy but also highlighted meaningful trends in the dataset related to social and economic factors.

## CHAPTER 6

### CONCLUSION

The Titanic survival prediction project aimed to analyze key factors influencing passenger survival and apply machine learning techniques for accurate prediction. Through thorough data preprocessing and exploration, important features such as gender, passenger class, fare, and embarkation point were identified as major contributors to survival outcomes. These insights aligned with historical facts and added analytical depth to the dataset.

Multiple machine learning models were implemented and evaluated, including Random Forest, XGBoost, Support Vector Machine, and a Voting Classifier. Among these, the ensemble Voting Classifier delivered the highest accuracy, around 89%, by effectively combining the strengths of individual models. Random Forest and XGBoost also showed strong performance, confirming their reliability in classification tasks involving mixed data types.

Overall, the project successfully demonstrated how predictive modeling can be used not just for accuracy but also for understanding real-world events. It emphasized the value of data-driven insights in uncovering hidden patterns and supported the use of machine learning in historical and social data analysis. This approach can be extended to other domains for decision-making and pattern recognition.

## REFERENCES

- [1] Kaggle. (2012). *Titanic: Machine Learning from Disaster*. Retrieved from <https://www.kaggle.com/competitions/titanic>
- [2] Géron, A. (2019). *Hands-On Machine Learning with Scikit-Learn, Keras, and TensorFlow* (2nd ed.). O'Reilly Media.  
<https://www.oreilly.com/library/view/hands-on-machine-learning/9781492032632/>
- [3] Pedregosa, F., Varoquaux, G., Gramfort, A., et al. (2011). *Scikit-learn: Machine Learning in Python*. *Journal of Machine Learning Research*, 12, 2825–2830. <https://jmlr.org/papers/v12/pedregosa11a.html>
- [4] Chen, T., & Guestrin, C. (2016). *XGBoost: A Scalable Tree Boosting System*. In *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 785–794.  
<https://doi.org/10.1145/2939672.2939785>
- [5] Breiman, L. (2001). *Random Forests*. *Machine Learning*, 45(1), 5–32.  
<https://link.springer.com/article/10.1023/A:1010933404324>
- [6] James, G., Witten, D., Hastie, T., & Tibshirani, R. (2013). *An Introduction to Statistical Learning: With Applications in R*. Springer.  
<https://www.statlearning.com/>
- [7] Cortes, C., & Vapnik, V. (1995). *Support-vector networks*. *Machine Learning*, 20(3), 273–297.  
<https://link.springer.com/article/10.1007/BF00994018>
- [8] W3Schools. (2020). *Python Machine Learning Tutorial*. Retrieved from [https://www.w3schools.com/python/python\\_ml\\_getting\\_started.asp](https://www.w3schools.com/python/python_ml_getting_started.asp)

## APPENDIX I

### CODE

#### Data Info

```
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
import warnings
warnings.filterwarnings('ignore')
data = pd.read_csv('/content/Titanic-Dataset.csv')
print(data.head())
data.shape
```

#### Data Preprocessing

```
import pandas as pd
import numpy as np
df = pd.read_csv("/content/Titanic-Dataset.csv")
df.drop(['PassengerId', 'Name', 'Ticket', 'Cabin'], axis=1, inplace=True)

df['Age'].fillna(df['Age'].median(), inplace=True)
df['Embarked'].fillna(df['Embarked'].mode()[0], inplace=True)

df['Sex'] = df['Sex'].map({'male': 0, 'female': 1})
df['Embarked'] = df['Embarked'].map({'S': 0, 'C': 1, 'Q': 2})

print(df.info())
print(df.head())
```

#### Distribution of Individual Features

```
import matplotlib.pyplot as plt
```

```

import seaborn as sns
sns.set(style="whitegrid")

numerical_cols = ['Age', 'Fare', 'SibSp', 'Parch']
for col in numerical_cols:
    plt.figure(figsize=(6, 4))
    sns.histplot(df[col], kde=True, bins=30, color='skyblue')
    plt.title(f'Distribution of {col}')
    plt.xlabel(col)
    plt.ylabel('Count')
    plt.tight_layout()
    plt.show()

categorical_cols = ['Survived', 'Pclass', 'Sex', 'Embarked']
for col in categorical_cols:
    plt.figure(figsize=(6, 4))
    sns.countplot(x=col, data=df, palette='pastel')
    plt.title(f'Count of {col}')
    plt.xlabel(col)
    plt.ylabel('Count')
    plt.tight_layout()
    plt.show()

```

### **Bivariate Analysis Code**

```

import matplotlib.pyplot as plt
import seaborn as sns
sns.set(style="whitegrid")

plt.figure(figsize=(6, 4))
sns.barplot(x='Sex', y='Survived', data=df, palette='pastel')
plt.title('Survival Rate by Sex')
plt.xlabel('Sex (0 = Male, 1 = Female)')
plt.ylabel('Survival Rate')

```

```
plt.tight_layout()
plt.show()
plt.figure(figsize=(6, 4))
sns.barplot(x='Pclass', y='Survived', data=df, palette='muted')
plt.title('Survival Rate by Passenger Class')
plt.xlabel('Passenger Class')
plt.ylabel('Survival Rate')
plt.tight_layout()
plt.show()
```

```
plt.figure(figsize=(8, 5))
sns.histplot(data=df, x='Age', hue='Survived', bins=30, kde=True, multiple='stack')
plt.title('Age Distribution by Survival')
plt.xlabel('Age')
plt.ylabel('Count')
plt.tight_layout()
plt.show()
plt.figure(figsize=(8, 5))
sns.boxplot(x='Survived', y='Fare', data=df, palette='cool')
plt.title('Fare Distribution by Survival')
plt.xlabel('Survived')
plt.ylabel('Fare')
plt.tight_layout()
plt.show()
```

```
plt.figure(figsize=(6, 4))
sns.barplot(x='Embarked', y='Survived', data=df, palette='Set2')
plt.title('Survival Rate by Embarkation Port')
plt.xlabel('Embarked (0=S, 1=C, 2=Q)')
plt.ylabel('Survival Rate')
plt.tight_layout()
plt.show()
```

## Multivariate Analysis Code

```
import matplotlib.pyplot as plt
import seaborn as sns
sns.set(style="whitegrid")

plt.figure(figsize=(10, 6))
corr_matrix = df.corr()
sns.heatmap(corr_matrix, annot=True, cmap='coolwarm', fmt=".2f", linewidths=0.5)
plt.title('Correlation Heatmap')
plt.tight_layout()
plt.show()

plt.figure(figsize=(6, 4))
sns.catplot(x='Sex', hue='Survived', col='Pclass', data=df,
            kind='count', palette='Set2', height=4, aspect=0.8)
plt.subplots_adjust(top=0.8)
plt.suptitle('Survival Count by Sex and Pclass')
plt.show()

plt.figure(figsize=(8, 5))
sns.scatterplot(data=df, x='Age', y='Fare', hue='Survived', palette='Set1')
plt.title('Age vs Fare Colored by Survival')
plt.xlabel('Age')
plt.ylabel('Fare')
plt.tight_layout()
plt.show()

plt.figure(figsize=(8, 5))
sns.boxplot(x='Pclass', y='Age', hue='Survived', data=df, palette='pastel')
plt.title('Age Distribution across Pclass and Survival')
plt.xlabel('Passenger Class')
plt.ylabel('Age')
plt.tight_layout()
```

```
plt.show()
```

### **Missing Value Analysis Code**

```
import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt
df = pd.read_csv("/content/Titanic-Dataset.csv")

missing_summary = df.isnull().sum()
missing_percent = (missing_summary / len(df)) * 100

missing_df = pd.DataFrame({
    'Missing Values': missing_summary,
    'Percentage (%)': missing_percent})
missing_df = missing_df[missing_df['Missing Values'] > 0].sort_values(by='Missing
Values', ascending=False)
print("Missing Value Summary:\n")
print(missing_df)

plt.figure(figsize=(10, 6))
sns.heatmap(df.isnull(), cbar=False, cmap='viridis', yticklabels=False)
plt.title('Missing Values Heatmap')
plt.show()
```

### **Outlier Detection Code**

```
import seaborn as sns
import matplotlib.pyplot as plt
sns.set(style="whitegrid")
numerical_cols = ['Age', 'Fare', 'SibSp', 'Parch']

for col in numerical_cols:
    plt.figure(figsize=(6, 4))
```



```

sns.boxplot(data=df, x=col, color='lightblue')
plt.title(f'Boxplot for {col}')
plt.xlabel(col)
plt.tight_layout()
plt.show()
outliers = {}
for col in numerical_cols:
    Q1 = df[col].quantile(0.25)
    Q3 = df[col].quantile(0.75)
    IQR = Q3 - Q1
    lower_bound = Q1 - 1.5 * IQR
    upper_bound = Q3 + 1.5 * IQR
    outlier_count = df[(df[col] < lower_bound) | (df[col] > upper_bound)].shape[0]
    outliers[col] = outlier_count

print("\nOutlier counts based on IQR method:\n")
for col, count in outliers.items():
    print(f'{col}: {count} outliers')

```

## Correlation Matrix

```

import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt

df = pd.read_csv("/content/Titanic-Dataset.csv")
numeric_df = df.select_dtypes(include=['int64', 'float64'])

correlation_matrix = numeric_df.corr()
print("Correlation Matrix:\n")
print(correlation_matrix)
plt.figure(figsize=(10, 6))
sns.heatmap(correlation_matrix, annot=True, cmap='coolwarm', fmt=".2f",
            linewidths=0.5)

```

```
plt.title('Feature Correlation Heatmap (Numeric Columns Only)')
plt.tight_layout()
plt.show()
```

### **Feature Interactions**

```
import seaborn as sns
import matplotlib.pyplot as plt
sns.set(style="whitegrid")
plt.figure(figsize=(8, 5))
sns.catplot(x="Pclass", hue="Sex", col="Survived", data=df,
            kind="count", palette="pastel", height=4, aspect=0.9)
plt.subplots_adjust(top=0.8)
plt.suptitle("Interaction: Pclass and Sex by Survival")
plt.show()

plt.figure(figsize=(8, 5))
sns.scatterplot(x='Age', y='Fare', hue='Survived', data=df, palette='Set1')
plt.title("Interaction: Age and Fare by Survival")
plt.xlabel("Age")
plt.ylabel("Fare")
plt.tight_layout()
plt.show()

plt.figure(figsize=(8, 5))
sns.catplot(x="Embarked", hue="Survived", col="Sex", data=df,
            kind="count", palette="Set2", height=4, aspect=0.9)
plt.subplots_adjust(top=0.8)
plt.suptitle("Interaction: Embarked and Sex by Survival")
plt.show()

plt.figure(figsize=(8, 5))
sns.boxplot(x='Pclass', y='Age', hue='Survived', data=df, palette='coolwarm')
plt.title("Interaction: Age by Pclass and Survival")
plt.xlabel("Passenger Class")
```

```
plt.ylabel("Age")
plt.tight_layout()
plt.show()
```

### **Family Size Feature Analysis**

```
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
df['FamilySize'] = df['SibSp'] + df['Parch'] + 1 # +1 to include the passenger
plt.figure(figsize=(6, 4))
sns.countplot(x='FamilySize', data=df, palette='pastel')
plt.title('Distribution of Family Size')
plt.xlabel('Family Size')
plt.ylabel('Number of Passengers')
plt.tight_layout()
plt.show()
```

```
plt.figure(figsize=(6, 4))
sns.barplot(x='FamilySize', y='Survived', data=df, palette='muted')
plt.title('Survival Rate by Family Size')
plt.xlabel('Family Size')
plt.ylabel('Survival Rate')
plt.tight_layout()
plt.show()
```

```
plt.figure(figsize=(6, 4))
sns.boxplot(x='FamilySize', y='Age', data=df, palette='coolwarm')
plt.title('Age Distribution by Family Size')
plt.xlabel('Family Size')
plt.ylabel('Age')
plt.tight_layout()
plt.show()
```

## Model Training and Accuracy

```
import pandas as pd

from sklearn.model_selection import train_test_split
from sklearn.linear_model import LogisticRegression
from sklearn.metrics import accuracy_score, confusion_matrix, classification_report

model_df = df.copy()
model_df['Sex'] = model_df['Sex'].map({'male': 0, 'female': 1})
model_df['Embarked'] = model_df['Embarked'].map({'S': 0, 'C': 1, 'Q': 2})
model_df['Age'].fillna(model_df['Age'].median(), inplace=True)
model_df['Embarked'].fillna(model_df['Embarked'].mode()[0], inplace=True)

features = ['Pclass', 'Sex', 'Age', 'Fare', 'SibSp', 'Parch', 'Embarked']
X = model_df[features]
y = model_df['Survived']

X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,
random_state=42)

model = LogisticRegression(max_iter=1000)
model.fit(X_train, y_train)
y_pred = model.predict(X_test)
accuracy = accuracy_score(y_test, y_pred)
print(f"\n✅ Model Accuracy: {accuracy:.4f}")

print("\nConfusion Matrix:")
print(confusion_matrix(y_test, y_pred))
print("\nClassification Report:")
print(classification_report(y_test, y_pred))
```

## APPENDIX II

### Screenshots

```
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
import warnings
warnings.filterwarnings('ignore')
data = pd.read_csv('/content/Titanic-Dataset.csv')
print(data.head())
```

	PassengerId	Survived	Pclass	\
0	1	0	3	
1	2	1	1	
2	3	1	3	
3	4	1	1	
4	5	0	3	

	Name	Sex	Age	SibSp	\
0	Braund, Mr. Owen Harris	male	22.0	1	
1	Cummings, Mrs. John Bradley (Florence Briggs Th...	female	38.0	1	
2	Heikkinen, Miss. Laina	female	26.0	0	
3	Futrelle, Mrs. Jacques Heath (Lily May Peel)	female	35.0	1	
4	Allen, Mr. William Henry	male	35.0	0	

	Parch	Ticket	Fare	Cabin	Embarked
0	0	A/5 21171	7.2500	NaN	S
1	0	PC 17599	71.2833	C85	C
2	0	STON/O2. 3101282	7.9250	NaN	S
3	0	113803	53.1000	C123	S
4	0	373450	8.0500	NaN	S

#### Data Preprocessing

```
import pandas as pd
import numpy as np

df = pd.read_csv("/content/Titanic-Dataset.csv")
df.drop(['PassengerId', 'Name', 'Ticket', 'Cabin'], axis=1, inplace=True)

df['Age'].fillna(df['Age'].median(), inplace=True)

df['Embarked'].fillna(df['Embarked'].mode()[0], inplace=True)

df['Sex'] = df['Sex'].map({'male': 0, 'female': 1})

df['Embarked'] = df['Embarked'].map({'S': 0, 'C': 1, 'Q': 2})

print(df.info())
print(df.head())
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 891 entries, 0 to 890
Data columns (total 8 columns):
#   Column      Non-Null Count  Dtype
---  ---
0   Survived    891 non-null    int64
1   Pclass      891 non-null    int64
2   Sex         891 non-null    int64
3   Age         891 non-null    float64
4   SibSp       891 non-null    int64
5   Parch       891 non-null    int64
6   Fare        891 non-null    float64
7   Embarked    891 non-null    int64
dtypes: float64(2), int64(6)
memory usage: 55.8 KB
None
```

	Survived	Pclass	Sex	Age	SibSp	Parch	Fare	Embarked
0	0	3	0	22.0	1	0	7.2500	0
1	1	1	1	38.0	1	0	71.2833	1
2	1	3	1	26.0	0	0	7.9250	0
3	1	1	1	35.0	1	0	53.1000	0
4	0	3	0	35.0	0	0	8.0500	0

## Distribution of Individual Features

```
import matplotlib.pyplot as plt
import seaborn as sns

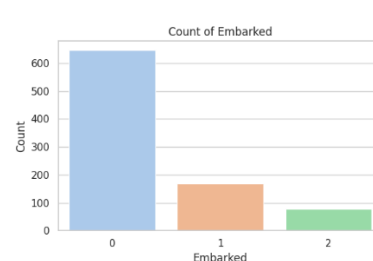
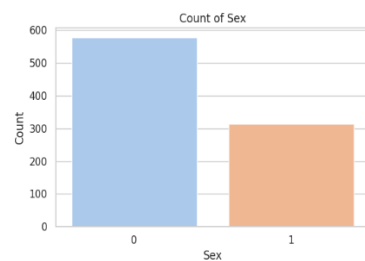
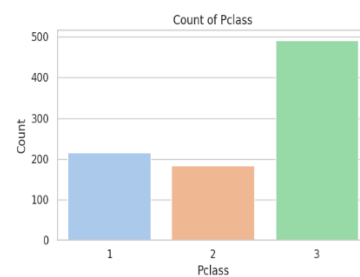
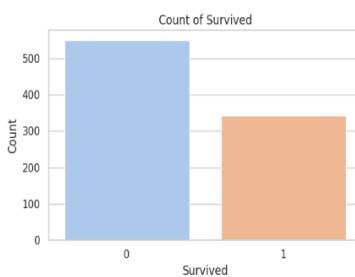
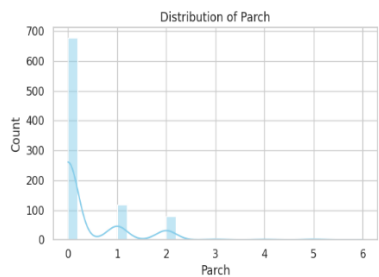
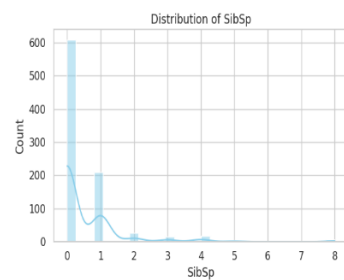
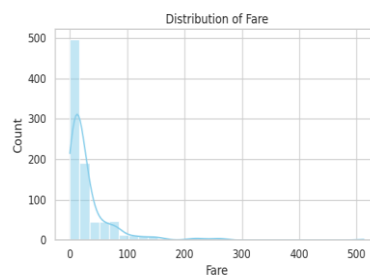
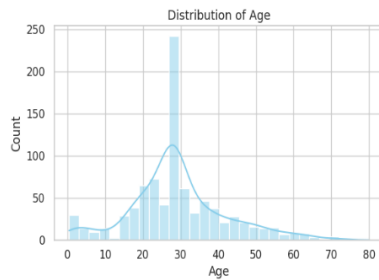
sns.set(style="whitegrid")

numerical_cols = ['Age', 'Fare', 'SibSp', 'Parch']

for col in numerical_cols:
    plt.figure(figsize=(6, 4))
    sns.histplot(df[col], kde=True, bins=30, color='skyblue')
    plt.title(f'Distribution of {col}')
    plt.xlabel(col)
    plt.ylabel('Count')
    plt.tight_layout()
    plt.show()

categorical_cols = ['Survived', 'Pclass', 'Sex', 'Embarked']

for col in categorical_cols:
    plt.figure(figsize=(6, 4))
    sns.countplot(x=col, data=df, palette='pastel')
    plt.title(f'Count of {col}')
    plt.xlabel(col)
    plt.ylabel('Count')
    plt.tight_layout()
    plt.show()
```



## Bivariate Analysis Code

```
import matplotlib.pyplot as plt
import seaborn as sns

sns.set(style="whitegrid")

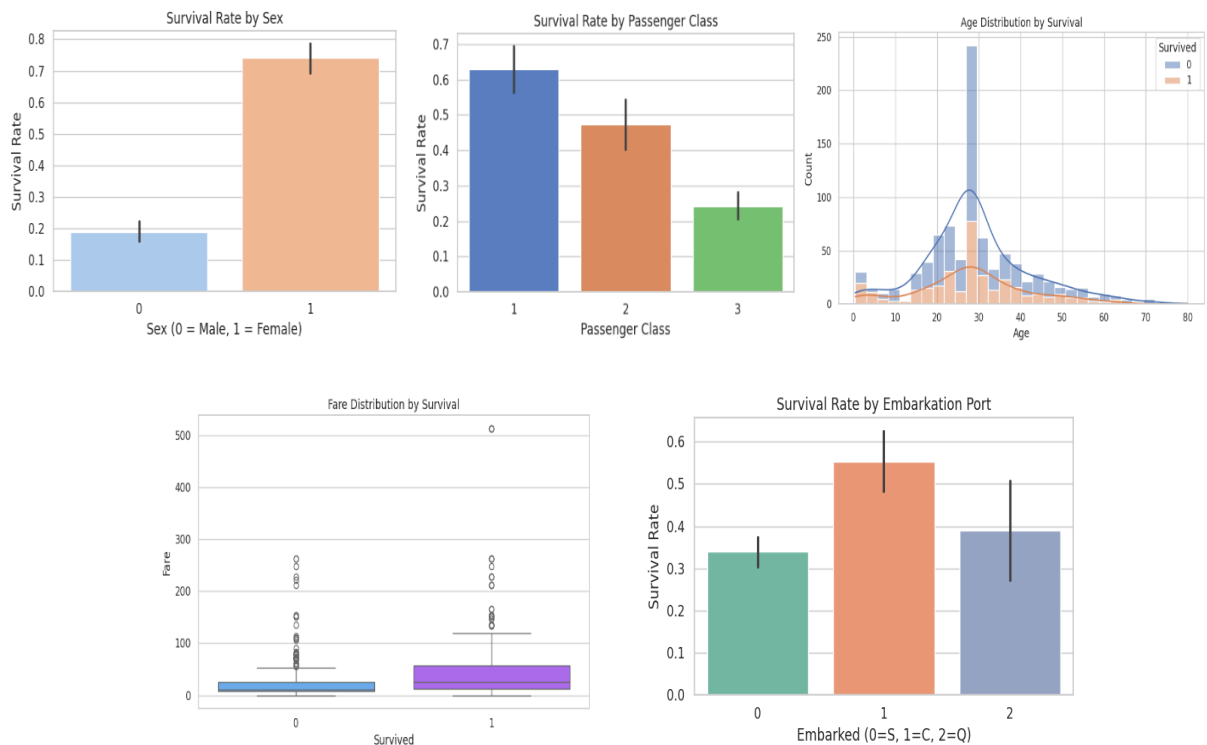
plt.figure(figsize=(6, 4))
sns.barplot(x='Sex', y='Survived', data=df, palette='pastel')
plt.title('Survival Rate by Sex')
plt.xlabel('Sex (0 = Male, 1 = Female)')
plt.ylabel('Survival Rate')
plt.tight_layout()
plt.show()

plt.figure(figsize=(6, 4))
sns.barplot(x='Pclass', y='Survived', data=df, palette='muted')
plt.title('Survival Rate by Passenger Class')
plt.xlabel('Passenger Class')
plt.ylabel('Survival Rate')
plt.tight_layout()
plt.show()

plt.figure(figsize=(8, 5))
sns.histplot(data=df, x='Age', hue='Survived', bins=30, kde=True, multiple='stack')
plt.title('Age Distribution by Survival')
plt.xlabel('Age')
plt.ylabel('Count')
plt.tight_layout()
plt.show()
```

```
plt.figure(figsize=(8, 5))
sns.boxplot(x='Survived', y='Fare', data=df, palette='cool')
plt.title('Fare Distribution by Survival')
plt.xlabel('Survived')
plt.ylabel('Fare')
plt.tight_layout()
plt.show()

plt.figure(figsize=(6, 4))
sns.barplot(x='Embarked', y='Survived', data=df, palette='Set2')
plt.title('Survival Rate by Embarkation Port')
plt.xlabel('Embarked (0=S, 1=C, 2=Q)')
plt.ylabel('Survival Rate')
plt.tight_layout()
plt.show()
```



```

import matplotlib.pyplot as plt
import seaborn as sns
sns.set(style="whitegrid")

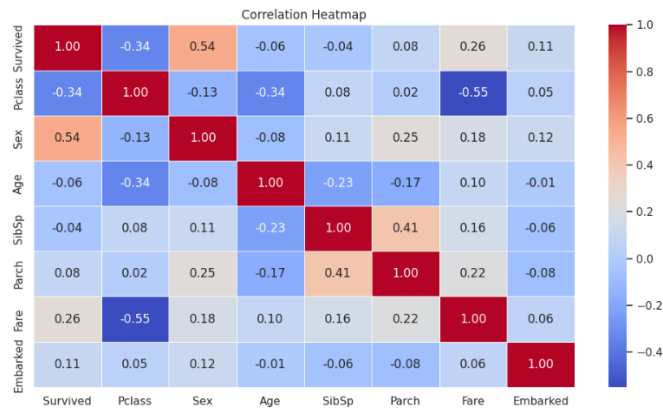
plt.figure(figsize=(10, 6))
corr_matrix = df.corr()
sns.heatmap(corr_matrix, annot=True, cmap='coolwarm', fmt=".2f", linewidths=0.5)
plt.title('Correlation Heatmap')
plt.tight_layout()
plt.show()

plt.figure(figsize=(6, 4))
sns.catplot(x='Sex', hue='Survived', col='Pclass', data=df,
            kind='count', palette='Set2', height=4, aspect=0.8)
plt.subplots_adjust(top=0.8)
plt.suptitle('Survival Count by Sex and Pclass')
plt.show()

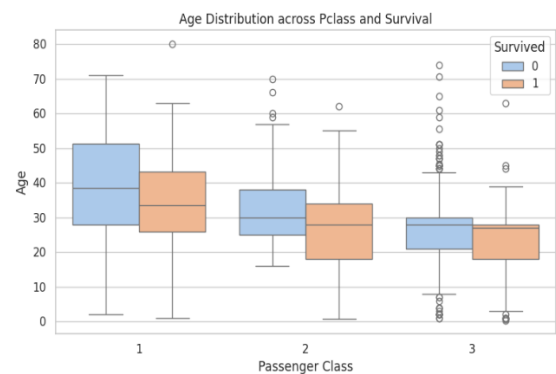
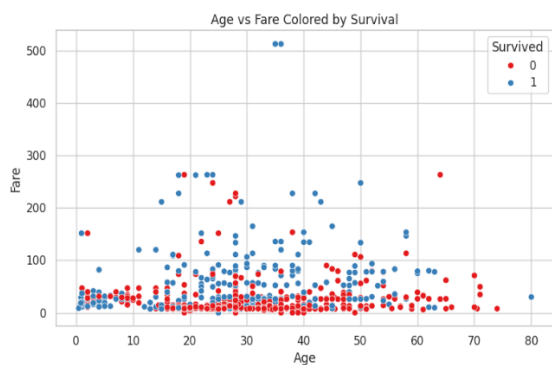
plt.figure(figsize=(8, 5))
sns.scatterplot(data=df, x='Age', y='Fare', hue='Survived', palette='Set1')
plt.title('Age vs Fare Colored by Survival')
plt.xlabel('Age')
plt.ylabel('Fare')
plt.tight_layout()
plt.show()

plt.figure(figsize=(8, 5))
sns.boxplot(x='Pclass', y='Age', hue='Survived', data=df, palette='pastel')
plt.title('Age Distribution across Pclass and Survival')
plt.xlabel('Passenger Class')
plt.ylabel('Age')
plt.tight_layout()
plt.show()

```



Survival Count by Sex and Pclass





### Missing Value Analysis Code

```
import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt

df = pd.read_csv("/content/Titanic-Dataset.csv")

missing_summary = df.isnull().sum()
missing_percent = (missing_summary / len(df)) * 100

missing_df = pd.DataFrame({
    'Missing Values': missing_summary,
    'Percentage (%)': missing_percent
})

missing_df = missing_df[missing_df['Missing Values'] > 0].sort_values(by='Missing Values', ascending=False)

print("Missing Value Summary:\n")
print(missing_df)

plt.figure(figsize=(10, 6))
sns.heatmap(df.isnull(), cbar=False, cmap='viridis', yticklabels=False)
plt.title('Missing Values Heatmap')
plt.show()
```

### Missing Value Analysis Code

```
import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt

df = pd.read_csv("/content/Titanic-Dataset.csv")

missing_summary = df.isnull().sum()
missing_percent = (missing_summary / len(df)) * 100

missing_df = pd.DataFrame({
    'Missing Values': missing_summary,
    'Percentage (%)': missing_percent
})

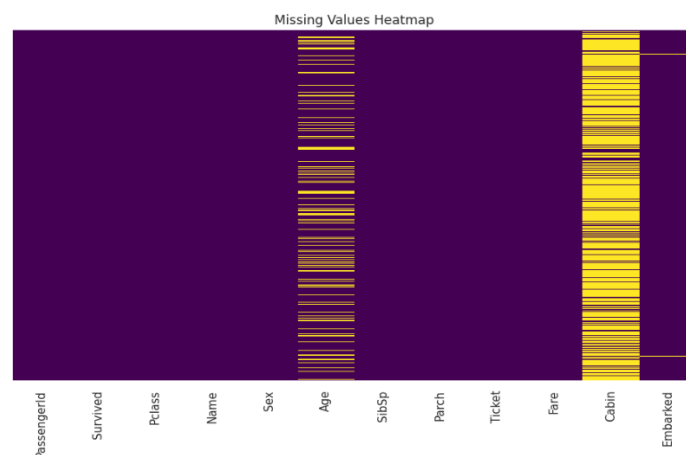
missing_df = missing_df[missing_df['Missing Values'] > 0].sort_values(by='Missing Values', ascending=False)

print("Missing Value Summary:\n")
print(missing_df)

plt.figure(figsize=(10, 6))
sns.heatmap(df.isnull(), cbar=False, cmap='viridis', yticklabels=False)
plt.title('Missing Values Heatmap')
plt.show()
```

Missing Value Summary:

	Missing Values	Percentage (%)
Cabin	687	77.104377
Age	177	19.865320
Embarked	2	0.224467



## Outlier Detection Code

```
import seaborn as sns
import matplotlib.pyplot as plt
sns.set(style="whitegrid")

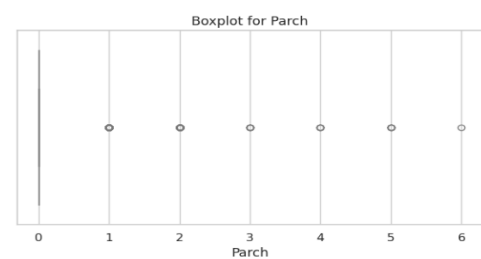
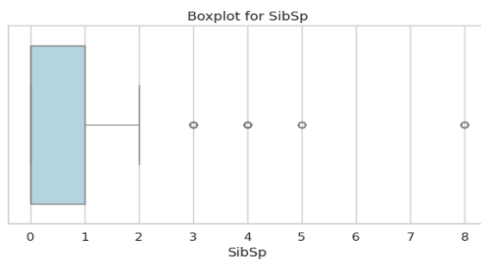
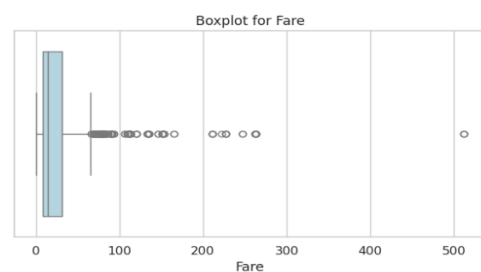
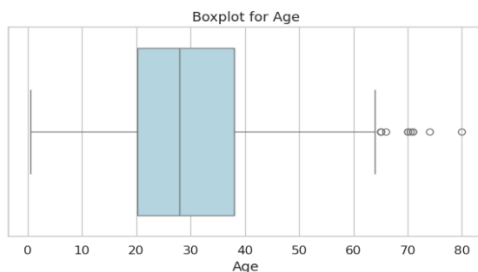
numerical_cols = ['Age', 'Fare', 'SibSp', 'Parch']

for col in numerical_cols:
    plt.figure(figsize=(6, 4))
    sns.boxplot(data=df, x=col, color='lightblue')
    plt.title(f'Boxplot for {col}')
    plt.xlabel(col)
    plt.tight_layout()
    plt.show()

outliers = {}

for col in numerical_cols:
    Q1 = df[col].quantile(0.25)
    Q3 = df[col].quantile(0.75)
    IQR = Q3 - Q1
    lower_bound = Q1 - 1.5 * IQR
    upper_bound = Q3 + 1.5 * IQR
    outlier_count = df[(df[col] < lower_bound) | (df[col] > upper_bound)].shape[0]
    outliers[col] = outlier_count

print("\nOutlier counts based on IQR method:\n")
for col, count in outliers.items():
    print(f"{col}: {count} outliers")
```

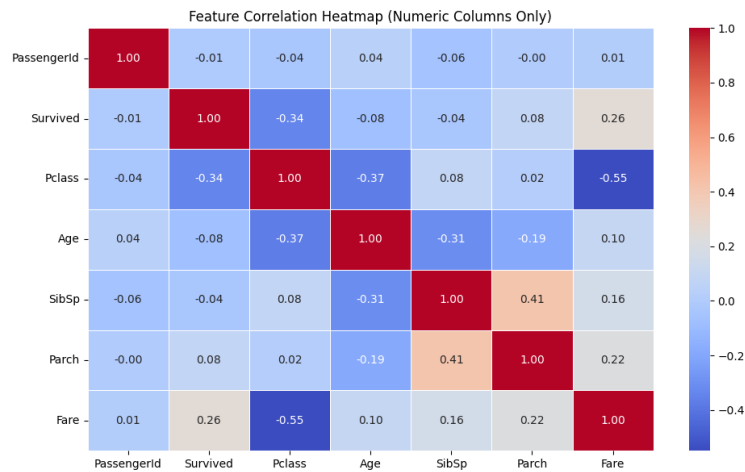


```
import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt
df = pd.read_csv("/content/Titanic-Dataset.csv")
numeric_df = df.select_dtypes(include=['int64', 'float64'])
correlation_matrix = numeric_df.corr()
print("Correlation Matrix:\n")
print(correlation_matrix)

plt.figure(figsize=(10, 6))
sns.heatmap(correlation_matrix, annot=True, cmap='coolwarm', fmt=".2f", linewidths=0.5)
plt.title('Feature Correlation Heatmap (Numeric Columns Only)')
plt.tight_layout()
plt.show()
```

Correlation Matrix:

	PassengerId	Survived	Pclass	Age	SibSp	Parch
PassengerId	1.000000	-0.005007	-0.035144	0.036847	-0.057527	-0.001652
Survived	-0.005007	1.000000	-0.338481	-0.077221	-0.035322	0.081629
Pclass	-0.035144	-0.338481	1.000000	-0.369226	0.083081	0.018443
Age	0.036847	-0.077221	-0.369226	1.000000	-0.308247	-0.189119
SibSp	-0.057527	-0.035322	0.083081	-0.308247	1.000000	0.414838
Parch	-0.001652	0.081629	0.018443	-0.189119	0.414838	1.000000
Fare	0.012658	0.257307	-0.549500	0.096067	0.159651	0.216225



### Feature Interactions

```
import seaborn as sns
import matplotlib.pyplot as plt

# Set seaborn style
sns.set(style="whitegrid")

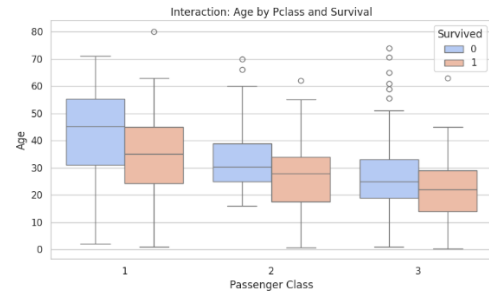
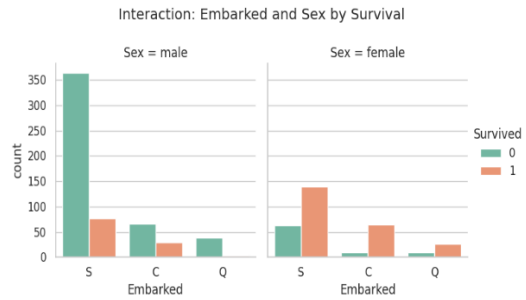
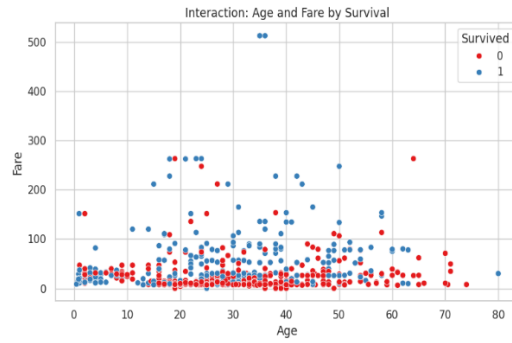
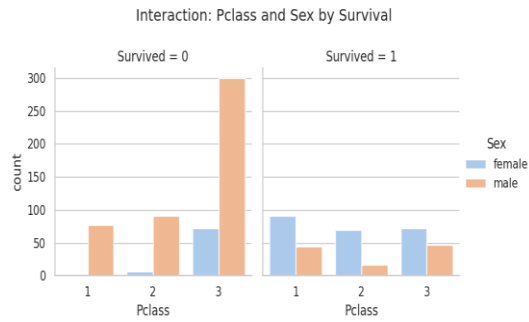
# -----
# 1. Interaction: Sex vs Pclass vs Survival
# -----
plt.figure(figsize=(8, 5))
sns.catplot(x="Pclass", hue="Sex", col="Survived", data=df,
            kind="count", palette="pastel", height=4, aspect=0.9)
plt.subplots_adjust(top=0.8)
plt.suptitle("Interaction: Pclass and Sex by Survival")
plt.show()

# -----
# 2. Interaction: Age vs Fare Colored by Survival
# -----
plt.figure(figsize=(8, 5))
sns.scatterplot(x="Age", y="Fare", hue="Survived", data=df, palette='Set1')
plt.title("Interaction: Age and Fare by Survival")
plt.xlabel("Age")
plt.ylabel("Fare")
plt.tight_layout()
plt.show()
```

```
# -----
# 3. Interaction: Embarked vs Sex vs Survival
# -----
plt.figure(figsize=(8, 5))
sns.catplot(x="Embarked", hue="Survived", col="Sex", data=df,
            kind="count", palette="Set2", height=4, aspect=0.9)
plt.subplots_adjust(top=0.8)
plt.suptitle("Interaction: Embarked and Sex by Survival")
plt.show()

# -----
# 4. Interaction: Age Distribution by Pclass and Survival
# -----
plt.figure(figsize=(8, 5))
sns.boxplot(x="Pclass", y="Age", hue="Survived", data=df, palette='coolwarm')
plt.title("Interaction: Age by Pclass and Survival")
plt.xlabel("Passenger Class")
plt.ylabel("Age")
plt.tight_layout()
plt.show()
```

<Figure size 800x500 with 0 Axes>



#### Family Size Feature Analysis

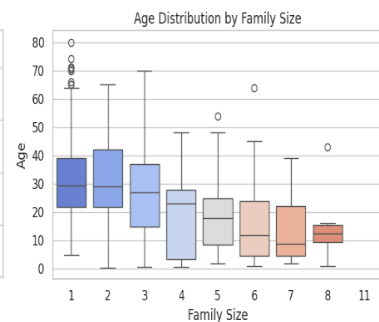
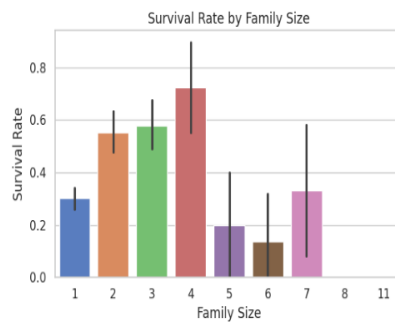
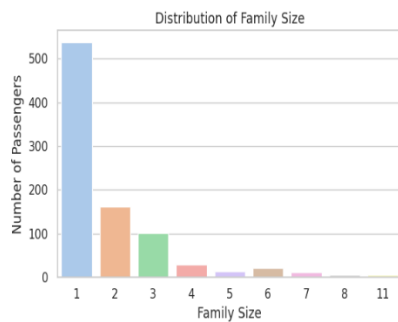
```
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns

df['FamilySize'] = df['SibSp'] + df['Parch'] + 1 # +1 to include the passenger

plt.figure(figsize=(6, 4))
sns.countplot(x='FamilySize', data=df, palette='pastel')
plt.title('Distribution of Family Size')
plt.xlabel('Family Size')
plt.ylabel('Number of Passengers')
plt.tight_layout()
plt.show()

plt.figure(figsize=(6, 4))
sns.barplot(x='FamilySize', y='Survived', data=df, palette='muted')
plt.title('Survival Rate by Family Size')
plt.xlabel('Family Size')
plt.ylabel('Survival Rate')
plt.tight_layout()
plt.show()

plt.figure(figsize=(6, 4))
sns.boxplot(x='FamilySize', y='Age', data=df, palette='coolwarm')
plt.title('Age Distribution by Family Size')
plt.xlabel('Family Size')
plt.ylabel('Age')
plt.tight_layout()
plt.show()
```



## Model Training and Accuracy

```
import pandas as pd
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LogisticRegression
from sklearn.metrics import accuracy_score, confusion_matrix, classification_report

# Step 1: Select Relevant Features
model_df = df.copy()

# Encode categorical columns
model_df['Sex'] = model_df['Sex'].map({'male': 0, 'female': 1})
model_df['Embarked'] = model_df['Embarked'].map({'S': 0, 'C': 1, 'Q': 2})

# Fill missing values
model_df['Age'].fillna(model_df['Age'].median(), inplace=True)
model_df['Embarked'].fillna(model_df['Embarked'].mode()[0], inplace=True)

# Select features and target
features = ['Pclass', 'Sex', 'Age', 'Fare', 'SibSp', 'Parch', 'Embarked']
X = model_df[features]
y = model_df['Survived']

# Step 2: Split Data
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

# Step 3: Train Logistic Regression Model
model = LogisticRegression(max_iter=1000)
model.fit(X_train, y_train)
```

```
# Step 4: Predictions and Accuracy
y_pred = model.predict(X_test)

# Accuracy
accuracy = accuracy_score(y_test, y_pred)
print(f"\n✅ Model Accuracy: {accuracy:.4f}")

# Confusion Matrix
print("\nConfusion Matrix:")
print(confusion_matrix(y_test, y_pred))

# Classification Report
print("\nClassification Report:")
print(classification_report(y_test, y_pred))
```

✅ Model Accuracy: 0.7989

Confusion Matrix:

```
[[89 16]
 [20 54]]
```

Classification Report:

	precision	recall	f1-score	support
0	0.82	0.85	0.83	105
1	0.77	0.73	0.75	74
accuracy			0.80	179
macro avg	0.79	0.79	0.79	179
weighted avg	0.80	0.80	0.80	179