# Exploratory Data Analysis Report: Titanic Dataset

**Data Overview:** The Titanic dataset contains information about passengers aboard the Titanic, including demographic details, travel information, and survival status. The key variables include:

- **Survived:** Survival status (0 = No, 1 = Yes)

- **Pclass:** Passenger class (1st, 2nd, 3rd)

- **Name:** Passenger name

- **Sex:** Passenger sex (male, female)

- **Age:** Passenger age

- **SibSp:** Number of siblings/spouses aboard

- **Parch:** Number of parents/children aboard

- **Ticket:** Ticket number

- **Fare:** Passenger fare

- **Cabin:** Cabin number

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

**Key Findings:**

1. **Survival Rate:**

   - The overall survival rate is approximately 38%.

   - Females had a significantly higher survival rate than males.

   - Passengers in higher classes (1st and 2nd) had a greater chance of survival compared to those in 3rd class.

2. **Age and Fare:**

   - The age distribution is slightly right-skewed, with most passengers being between 20 and 40 years old.

- o Passenger fares are heavily right-skewed, indicating that most passengers paid lower fares, with a few paying significantly higher fares.

- o There is no strong linear correlation between age and fare.

3. **Passenger Class:**

- o Most passengers were in the 3rd class.

- o Survival rates decrease as passenger class decreases.

4. **Embarked:**

- o Most passengers embarked from Southampton (S).

- o The port of embarkation appears to have some correlation with survival, though less pronounced than sex or passenger class.
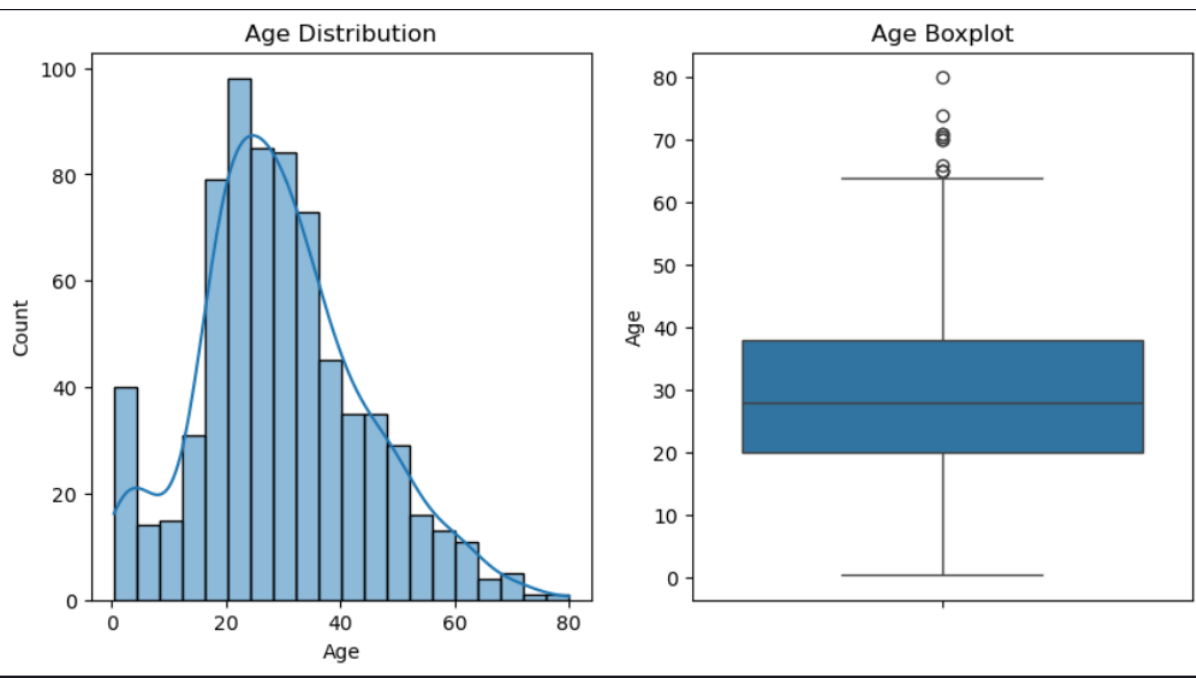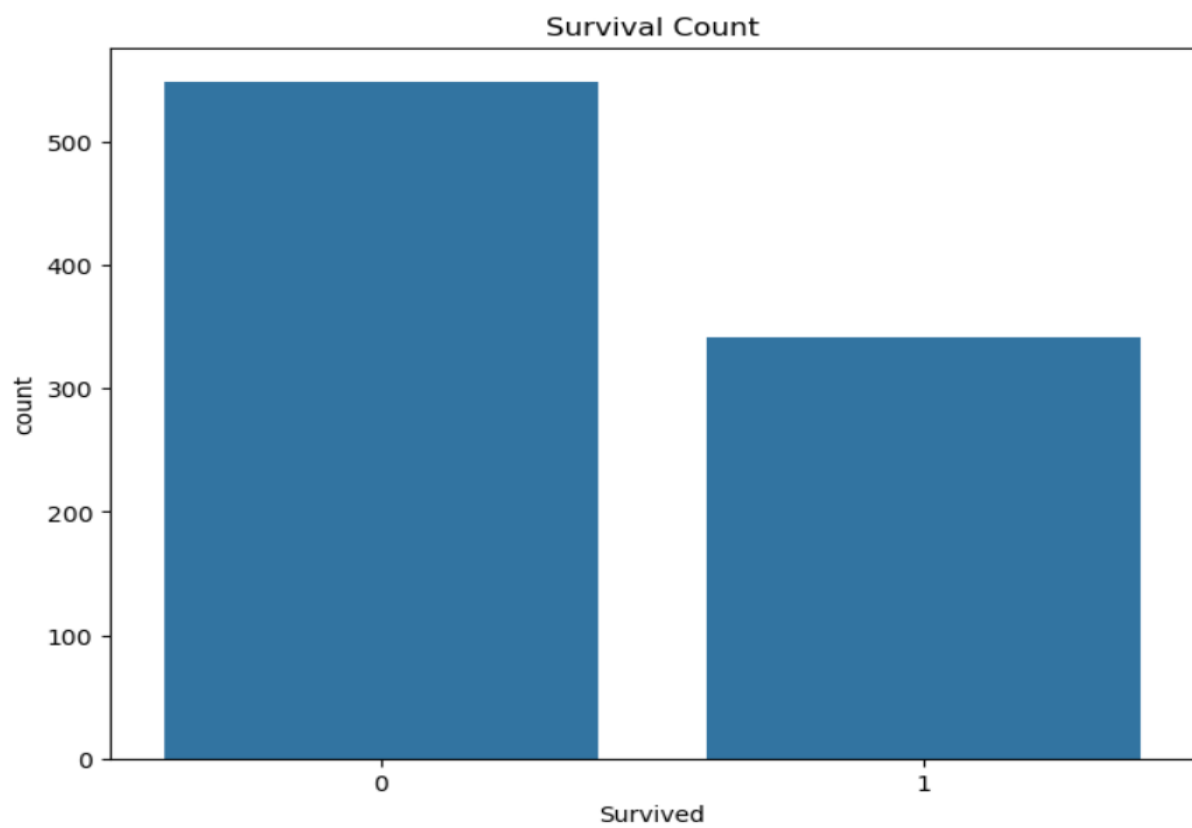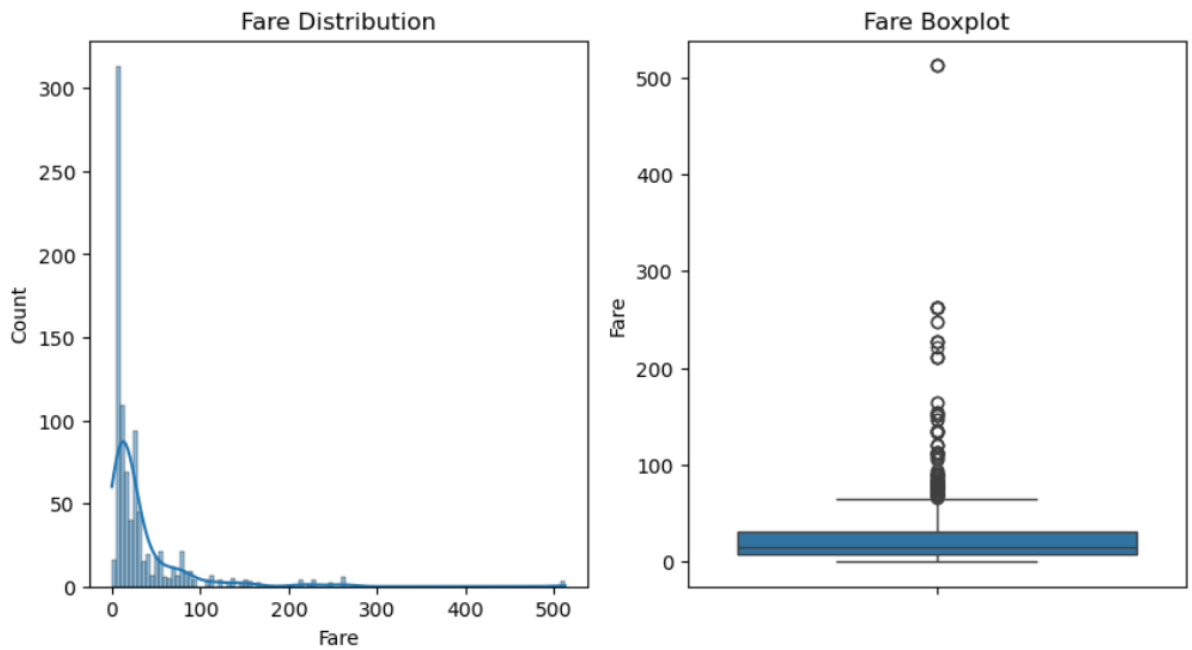
5. **Family Size:**

- o Most passengers traveled without siblings/spouses or parents/children.

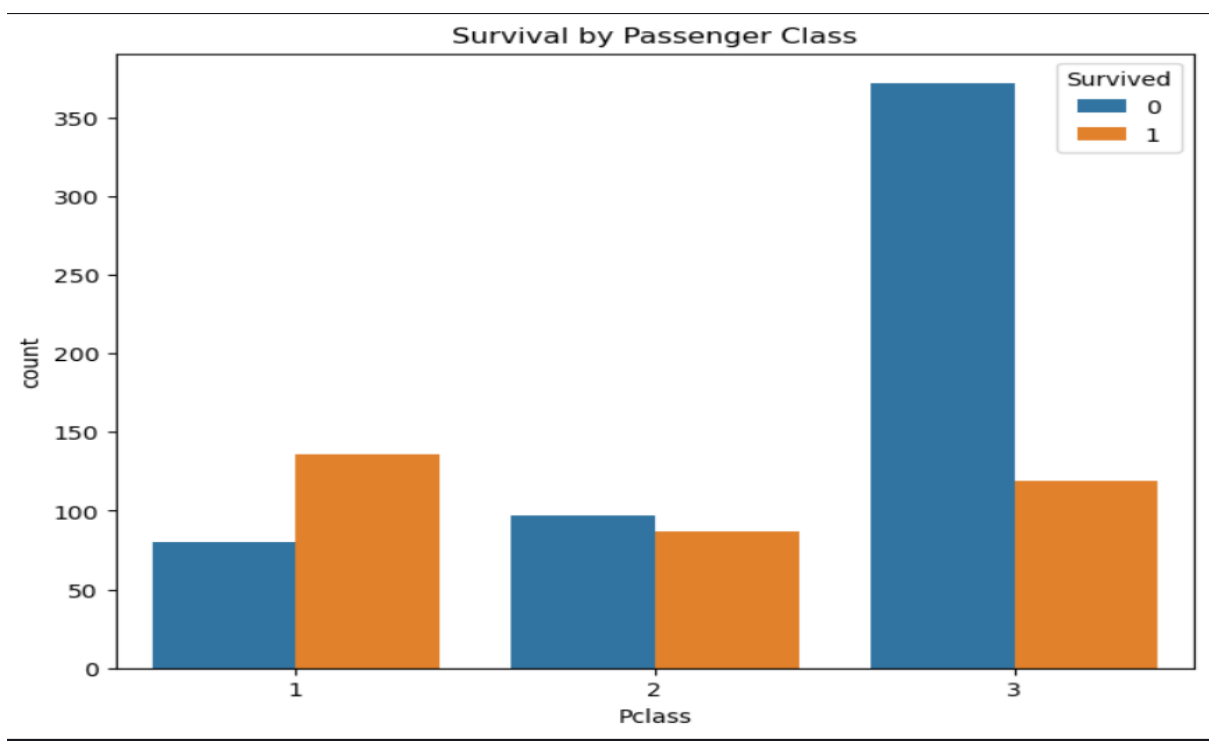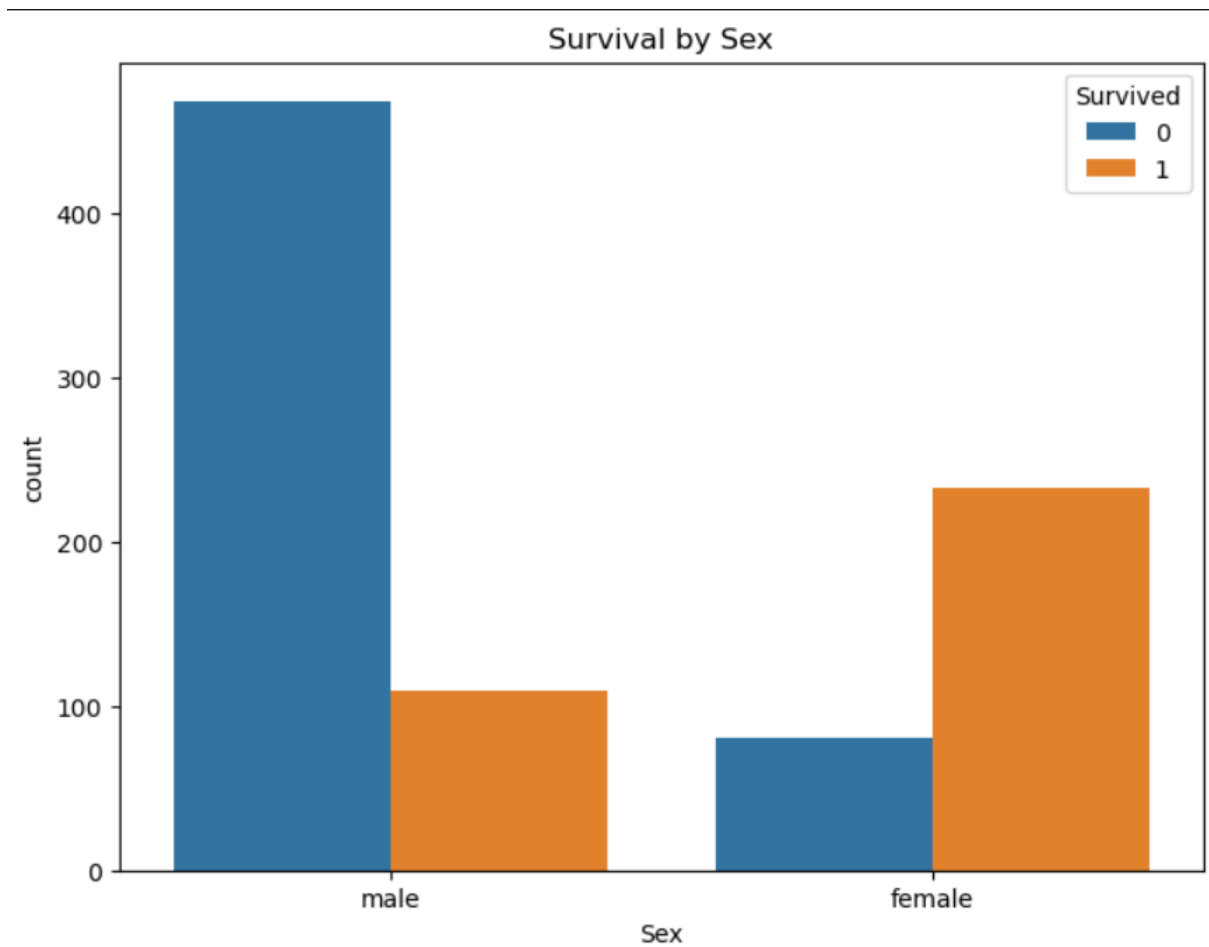- o Passengers with a small family size had a higher survival rate.
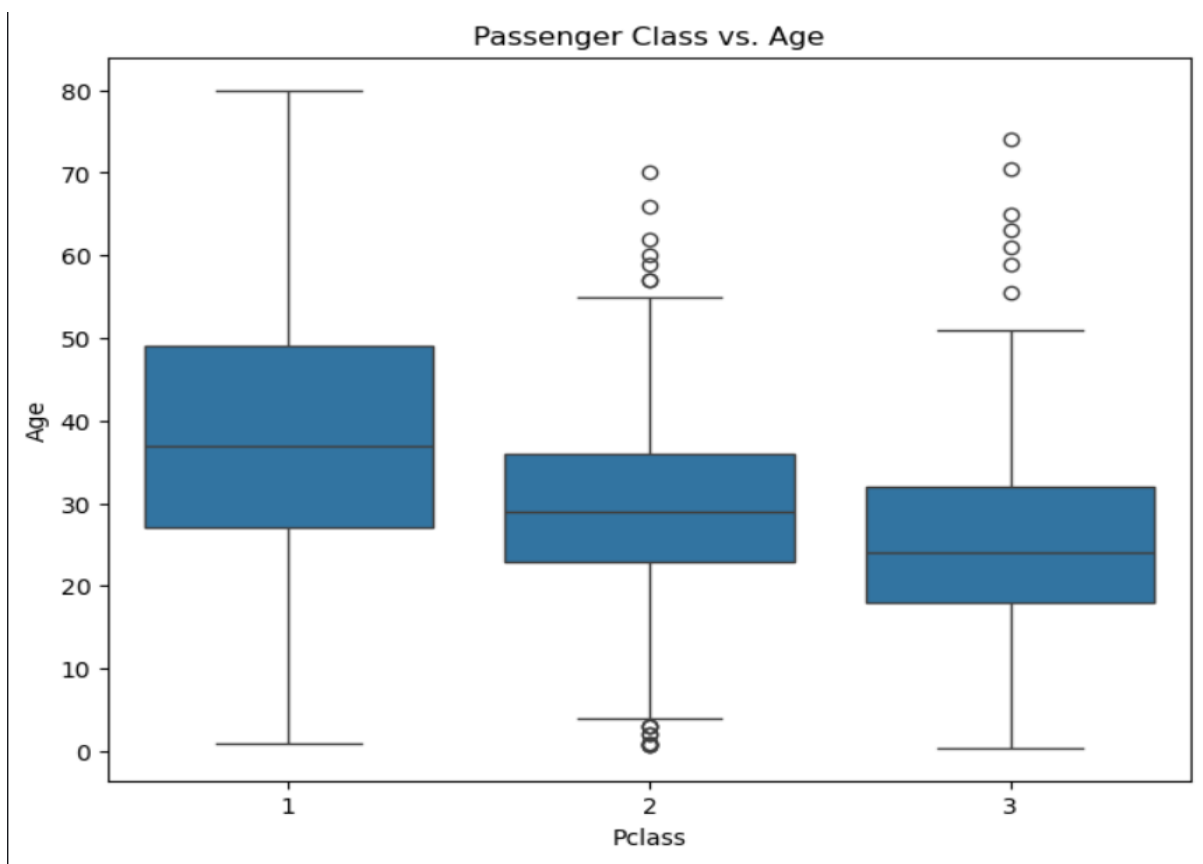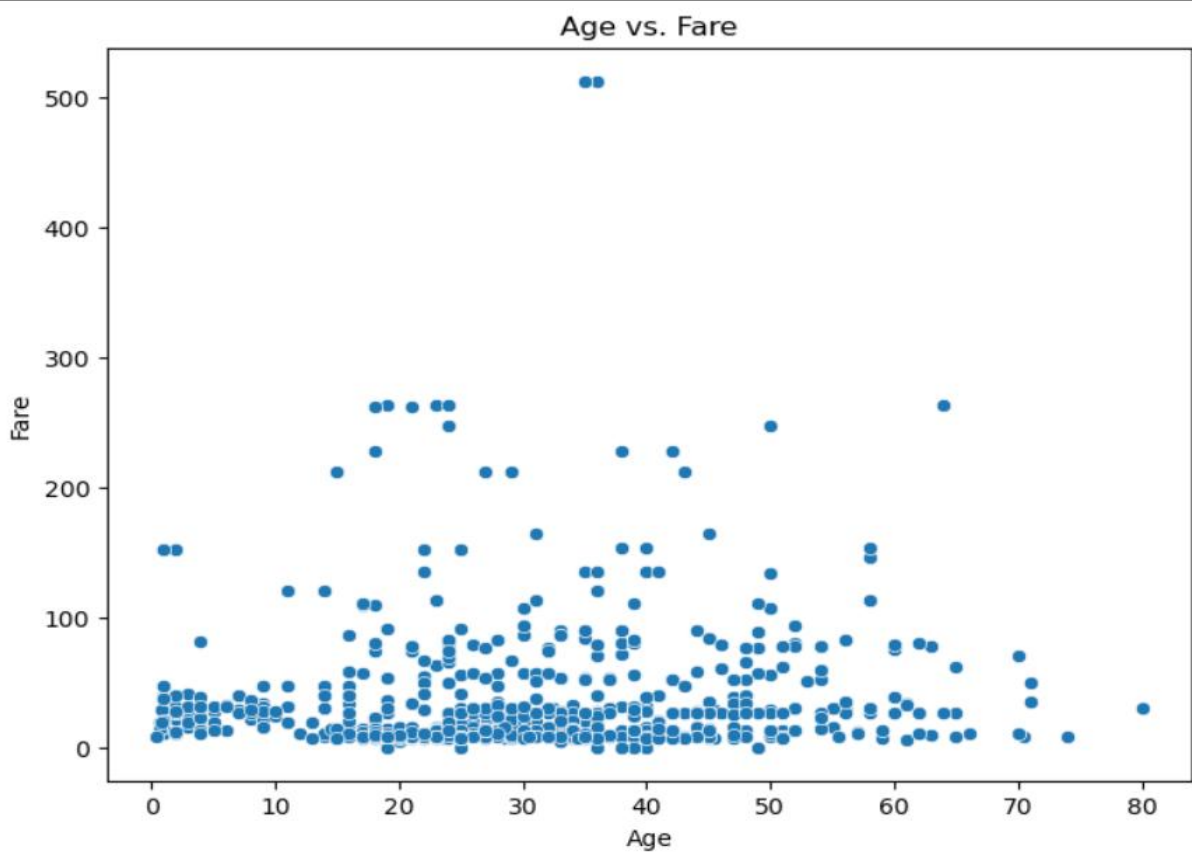
**Visualizations:**

- • **Histograms and Boxplots:** Used to visualize the distribution of numerical variables (Age, Fare) and identify potential outliers.

- • **Count Plots:** Used to visualize the frequency of categories in categorical variables (Survived, Sex, Pclass, Embarked) and their relationship with survival.

- • **Scatterplots:** Used to explore the relationship between two numerical variables (Age vs. Fare).
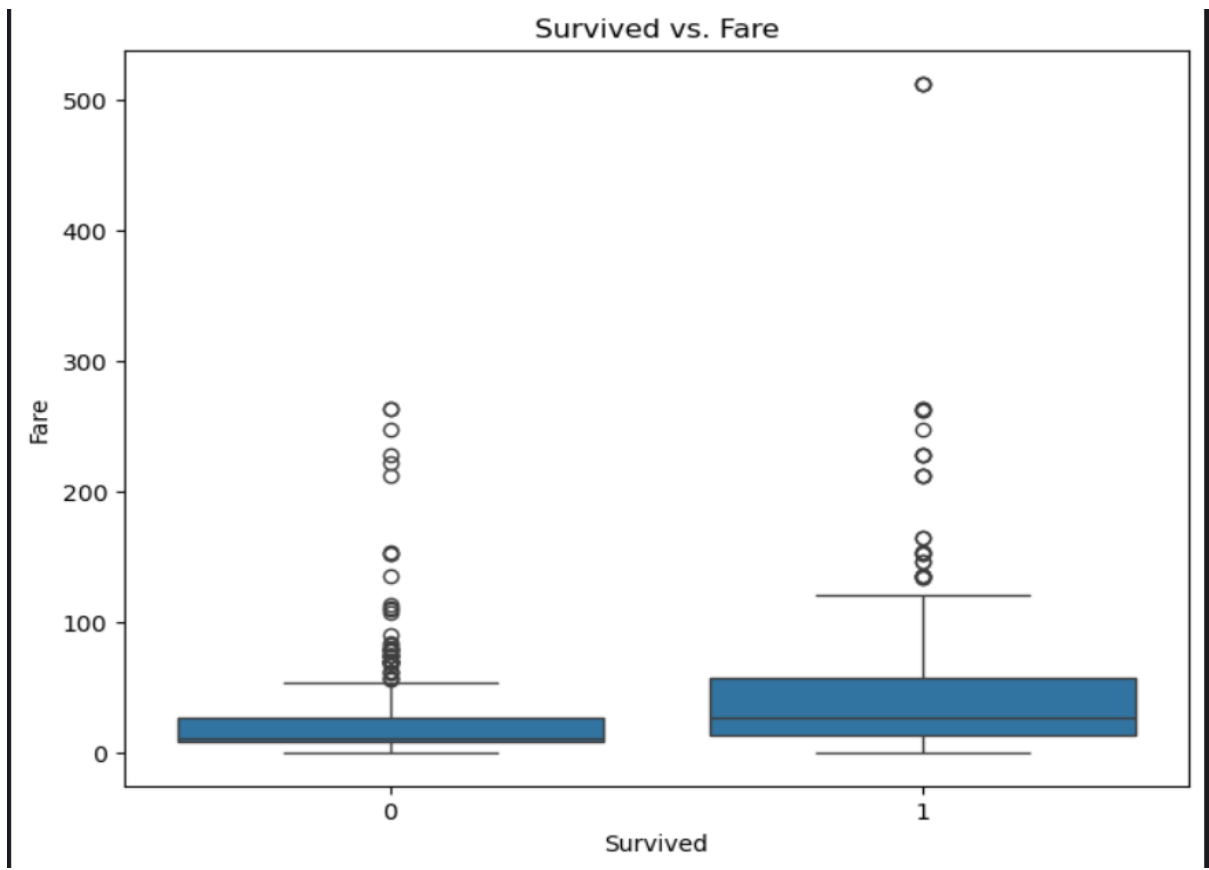
- **Crosstabs and Bar Plots:** Used to visualize the relationship between categorical variables (Sex vs. Survived, Pclass vs. Survived).

- **Pairplot:** Used to visualize pairwise relationships between multiple numerical variables.

- **Heatmap:** Used to display the correlation matrix and visualize correlations between numerical variables.
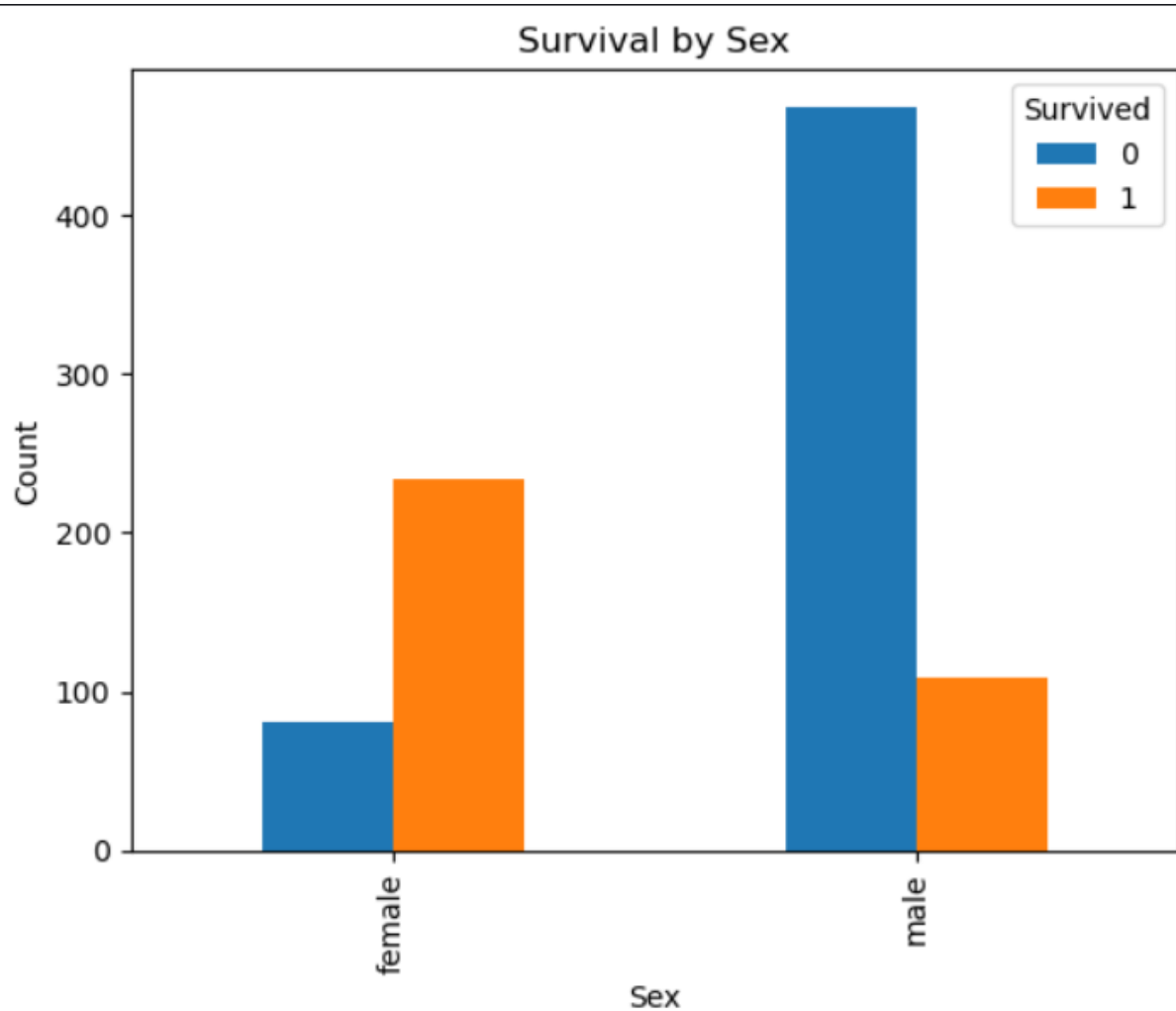
## Fare Distribution

## Fare Boxplot

## Survival Count

Survival by Sex



Survival by Passenger Class

Age vs. Fare



Passenger Class vs. Age

Survived vs. Fare

Survival by Sex

**Summary: The EDA reveals that survival on the Titanic was influenced by several factors, including sex, passenger class, and potentially family size and port of embarkation. Females and passengers in higher classes had a higher likelihood of survival. Age and fare did not appear to be strong predictors of survival.**