Data Analytics Internship

# Day5  Exploratory Data Analysis (EDA)

## Dataset Used: Titanic_Datatset.csv

## 1. Basic Exploration using NumPy and Pandas
- Used functions: `.head()`, `.info()`, `.describe()`, `.isnull().sum()`
- Outcome: Helped understand the basic structure of the dataset including data types, number of rows/columns, summary statistics, and missing values.
- Inference: Detected missing values in **'Age'**, **'Fare'**, and **'Embarked'** columns. Noted that data had a mix of categorical and numerical variables.
- Additional Computations:
  - Calculated percentage of missing values to prioritize cleaning tasks.
  - Performed survival count and survival rate to understand class imbalance.
  - Distribution of Categorical Variables using `.value_counts()` provided insights into the passenger distribution by gender, class, embarkation point, etc.
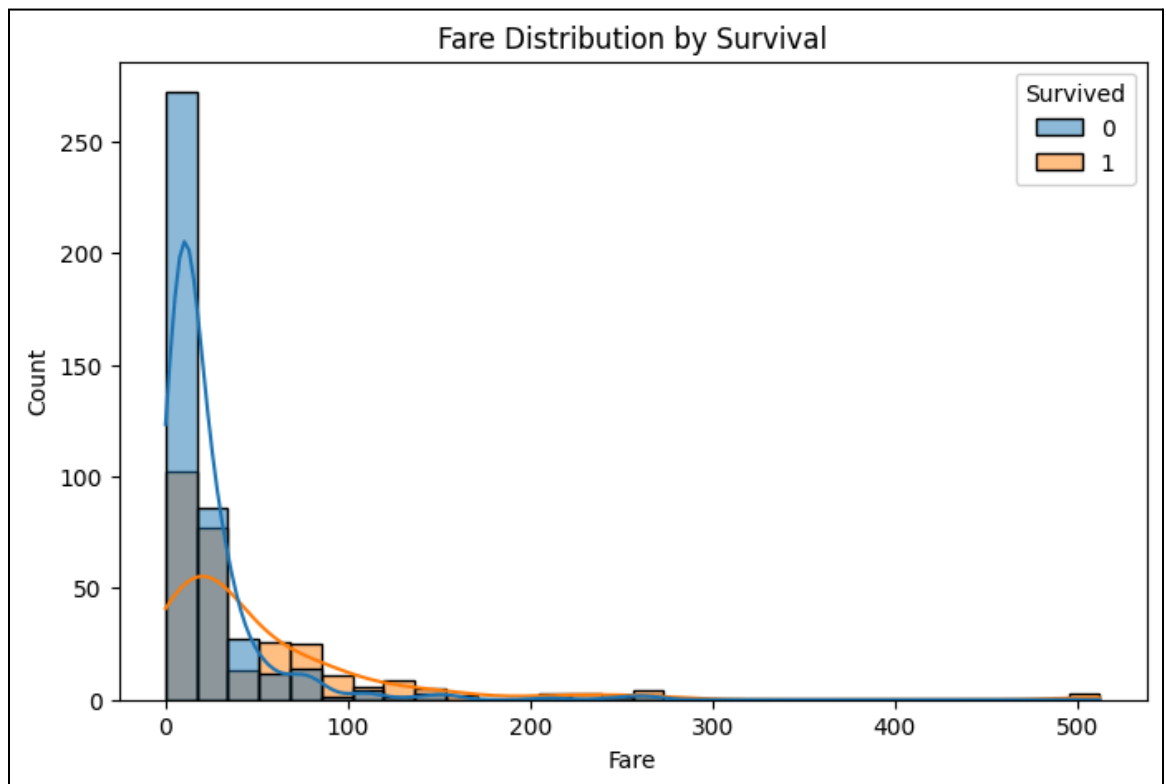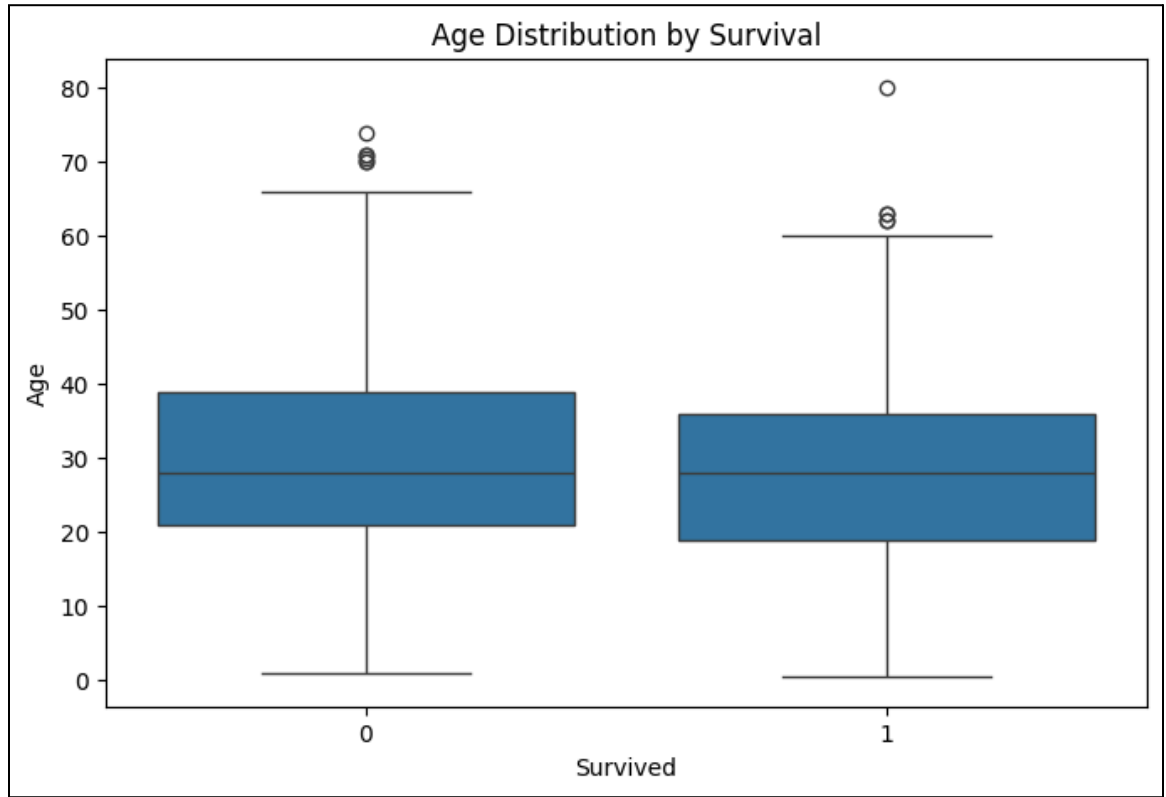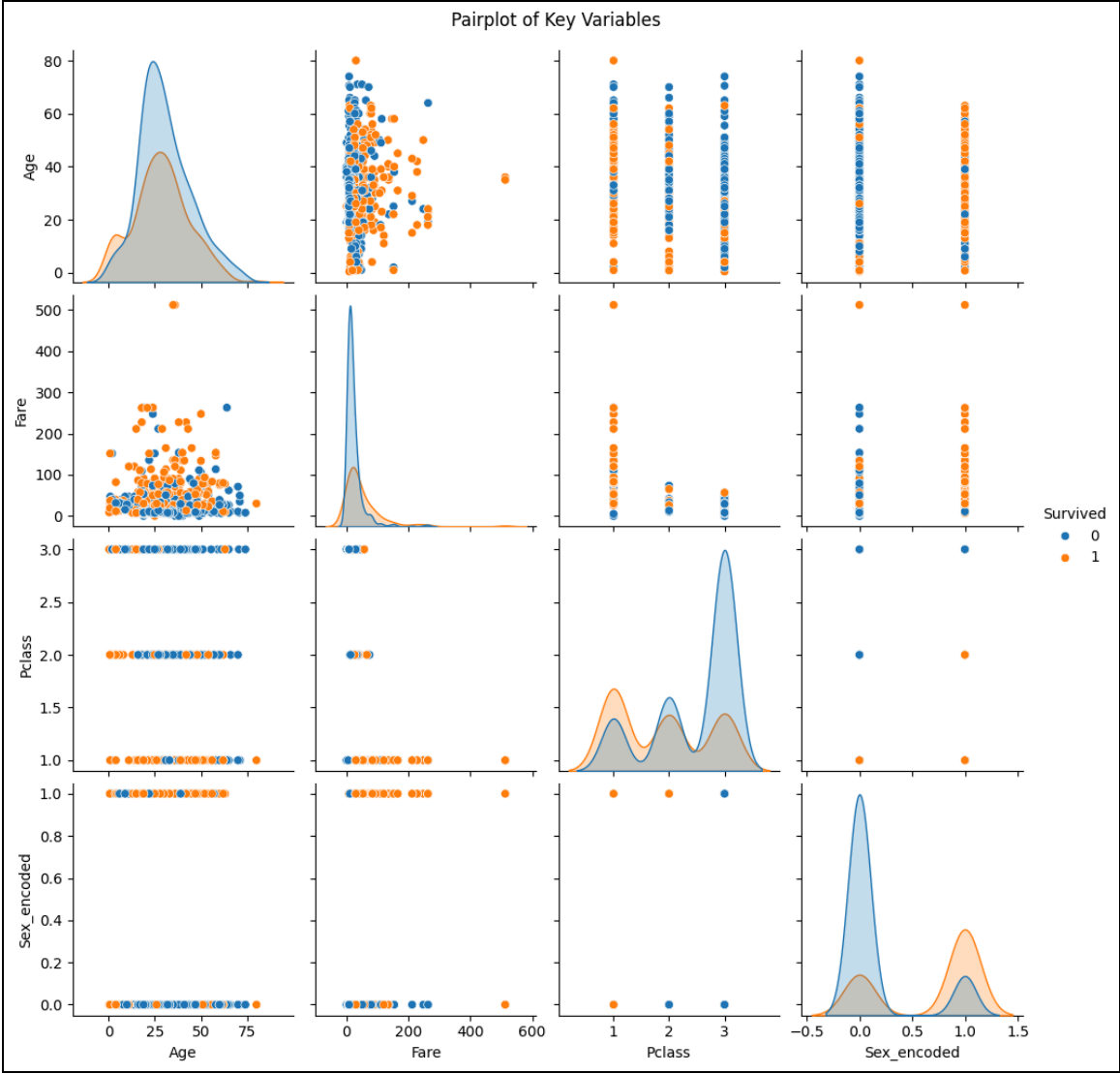
## 2. Data Cleaning
- Filled missing **'Embarked'** values using mode.
- Encoded categorical variables such as **'Sex'** and **'Embarked'** using `.map()`.
- Dropped rows with missing **'Age'** and **'Fare'** using `.dropna()`.
- Outcome: Cleaned dataset ready for analysis and visualization.
- Inference: Encoding allowed easier plotting and modeling; cleaning ensured no skew from missing data.
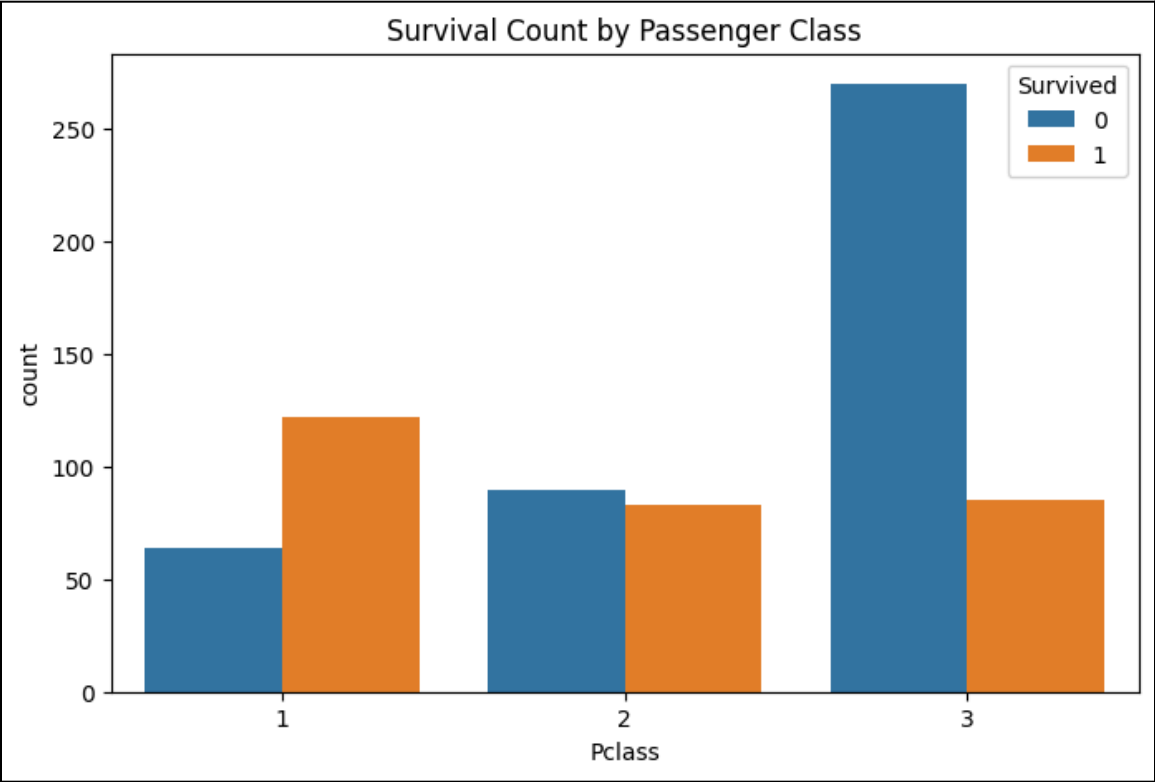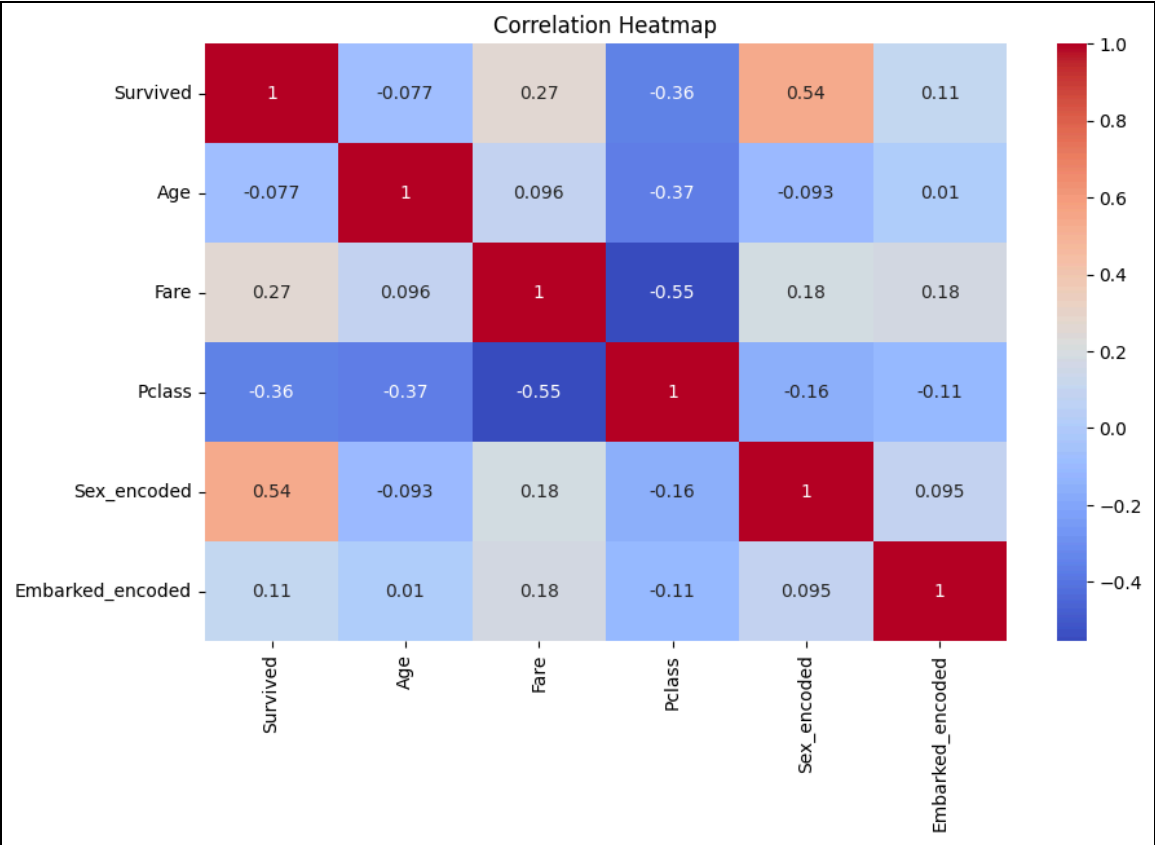
## 3. Graphical Exploration
- **Boxplot (Age vs Survival)**, **Histogram (Fare with Survival hue)**, **Pairplot,** and **Heatmap using Seaborn**:
  - Outcome: Visualized survival relationships with other features.
  - Inference: Females, higher class passengers, and younger individuals had higher survival chances.
- **Countplot of Passenger Class vs Survival**:
  - Outcome: Clear comparison between survival rate across passenger classes.

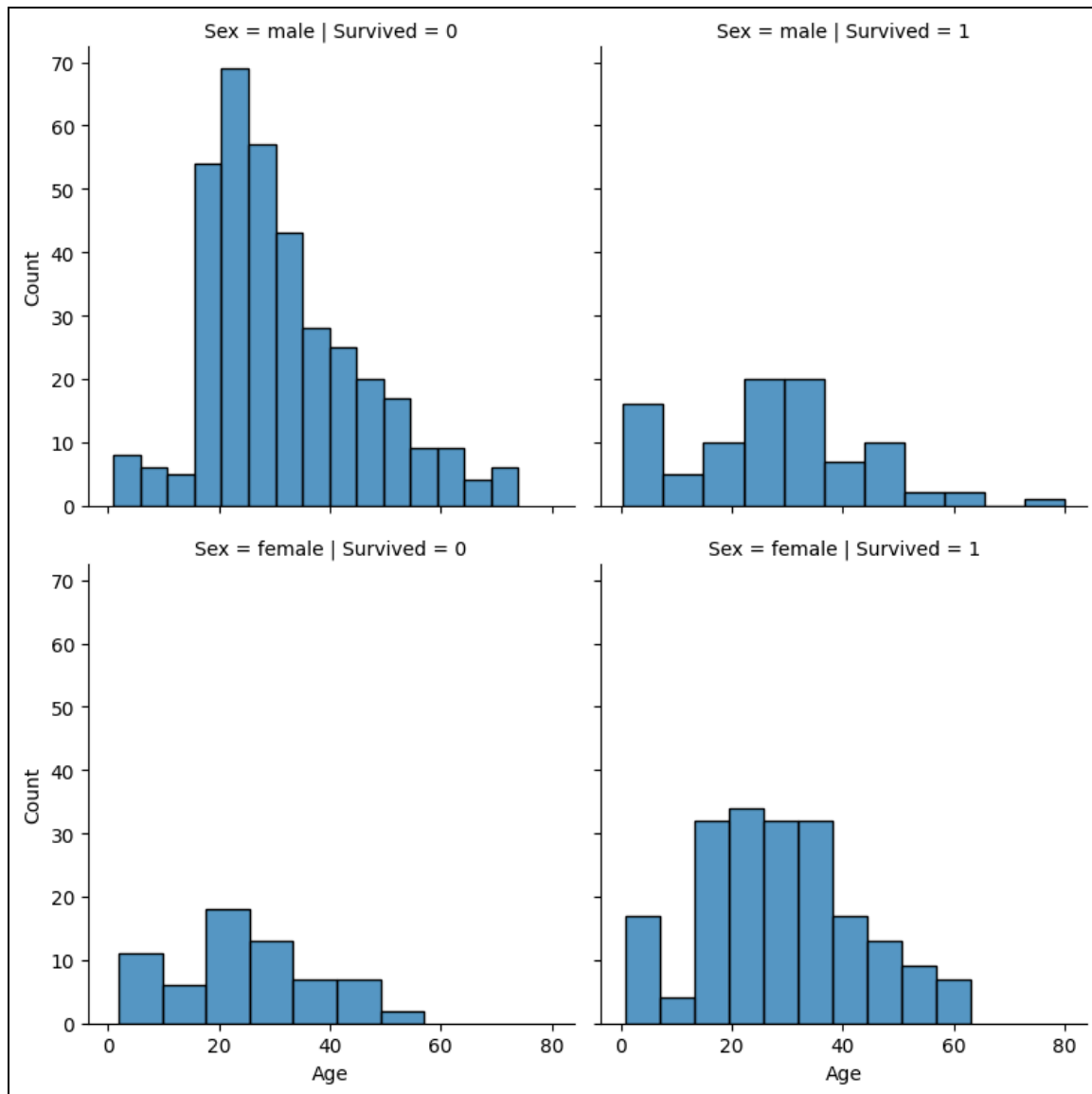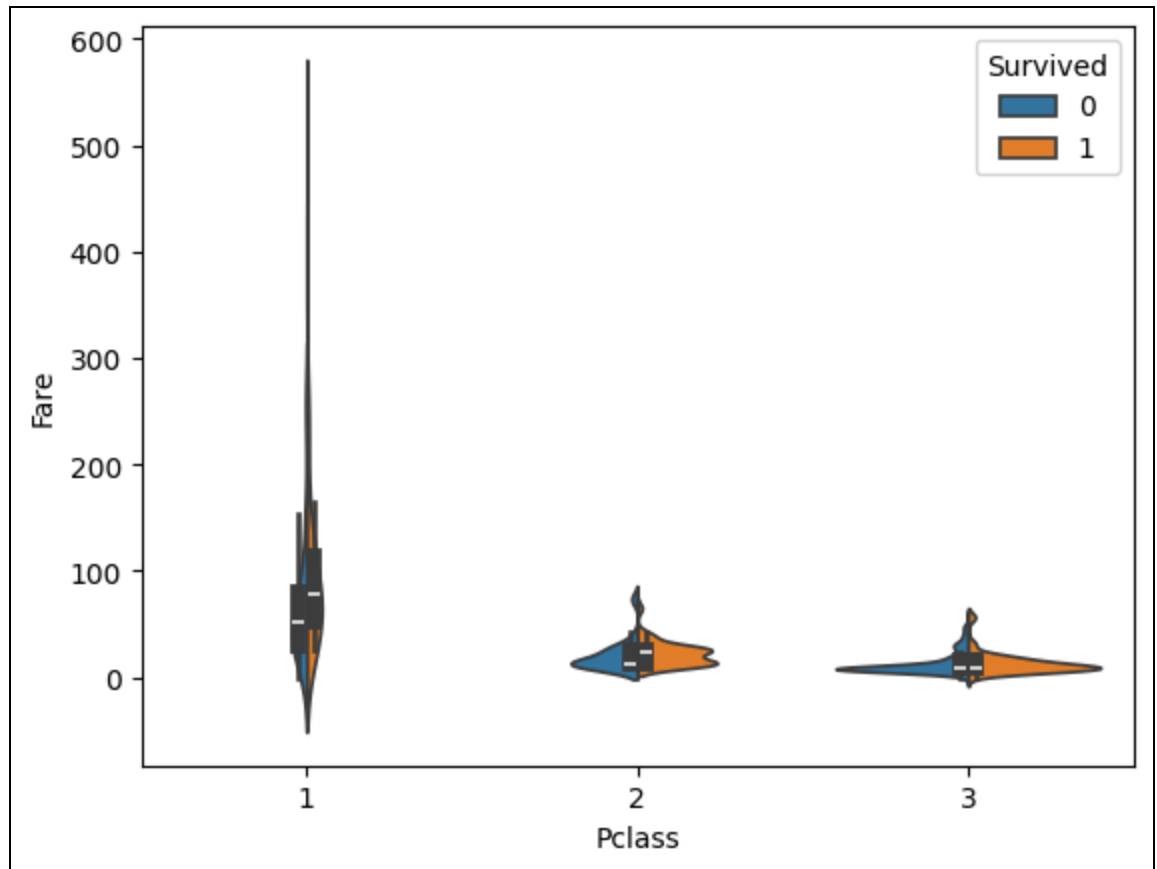○ Inference: Passengers from 1st class had significantly higher survival rates.

**Age Distribution by Survival**



**Fare Distribution by Survival**

Pairplot of Key Variables

## Correlation Heatmap



## Survival Count by Passenger Class

## 4. Extra Insights & Feature Engineering

- Created new features:
  - **Family Size** (SibSp + Parch + 1)
  - **Is Alone** (Binary feature if FamilySize == 1)
  - **Title Extraction from Name**
  - **Age Group Categorization**
- **ChiSquare Tests:** Assessed relationship between categorical features (e.g., Gender, Pclass) and Survival.
- **Violin Plot (Fare vs Class vs Survival)** and **FacetGrid (Age & Sex vs Survival)** for deeper insights.
  - Outcome: Derived new columns and revealed more meaningful survival patterns.
  - Inference: Certain titles (like 'Mrs', 'Miss') and family configurations were more likely to survive.

## 5. Final Summary & Observations

• Data cleaning and preprocessing were critical to prepare the dataset.

• Visualizations provided both univariate and multivariate insights.

• Engineered features enhanced understanding of survival patterns.

• Gender, class, family structure, and embark location strongly influenced survival.

• Navigation techniques, clear layout, and structured visuals were used for better storytelling.