

Titanic Dataset Analysis Report

1. Introduction

The Titanic tragedy is still one of the best-researched maritime disasters in history, and the dataset here provides significant insights on what predicts survival among passengers. The project entailed an analysis of the Titanic dataset to extract survival patterns based on passenger data as well as ticket data. The main aims were to clean the data, manage missing values efficiently, and conduct exploratory data analysis (EDA) to discern influential variables in survival rates.

2. Tools Utilized

- **Python:** Primary programming language utilized for data analysis.
- **Pandas:** Utilized for data cleaning and manipulation.
- **NumPy:** Utilized for numerical computation and missing data handling.
- **Matplotlib & Seaborn:** Utilized for plotting distributions and data visualization.
- **Jupyter Notebook:** Utilized for interactive coding and step-by-step analysis.

3. Steps in Creating the Project

Data Cleaning

- Recognized missing values in columns like *Age*, *Cabin*, and *Embarked*.
- Imputed missing *Age* values with the median age to preserve distribution integrity.
- Removed the *Cabin* column due to high levels of missing data.
- Imputed missing *Embarked* values with the mode (most frequent embarkation point).
- Transformed categorical variables such as *Sex* and *Embarked* into suitable formats for analysis.

Charts used:

- **Bar Plots / Count Plots:** To show counts of categorical variables such as survival counts by gender, passenger class (Pclass), and embarkation port. These plots helped reveal differences in survival rates across groups.

- **Histograms with KDE (Kernel Density Estimate):** To visualize the age distribution of passengers, comparing survivors and non-survivors.
- **Pie Charts:** To depict proportions of passengers by embarkation port.
- **Box Plots:** To explore the distribution of fares paid across different passenger classes and survival status.
- **Scatter Plots:** To examine relationships between continuous variables, such as Age vs. Fare.
- **Heatmaps:** To visualize correlation matrices among numerical features, identifying potential relationships.

These visualizations, created using libraries like Matplotlib and Seaborn, were instrumental in uncovering patterns such as higher survival rates among females, first-class passengers, and children, as well as fare-related trends.

Exploratory Data Analysis (EDA)

- Computed overall survival rate (~38%).
- Examined survival rates by gender, and found females had much higher survival probabilities than males.
- Analysed passenger class (Pclass), reporting greater survival in first- and second-class passengers.
- Researched age profile of survivors and non-survivors:
 - Survivors were on average younger than non-survivors.
- Researched fare paid, reporting passengers who paid more were more likely to survive.

4. Conclusion

The Titanic dataset analysis, following extensive data cleaning and exploratory analysis, identified a number of key factors determining survival. The most important predictors were gender and passenger class, with the highest survival rates for females and first-class passengers. Age was also an important factor, with children having significantly higher survival probabilities, consistent with evacuation priorities at the time of the disaster. The project illustrates the capacity of data cleansing and EDA to reveal rich insights from old datasets, leading to a potential platform for deeper statistical exploration or further predictive modelling.



