

TITANIC DATASET ANALYSIS: KEY FINDINGS AND RECOMMENDATIONS

INTRODUCTION

In this project, I explored the famous Titanic dataset available on Kaggle using Python. My goal was to perform an exploratory data analysis (EDA) using simple tools—pandas for data manipulation and Plotly along with Dash for interactive visualizations. I generated 12 distinct visual outputs that provided different insights into the data. I did not use advanced machine learning techniques but focused on understanding the dataset through data cleaning, visualization, and descriptive analysis. In this document, I summarize the key findings of my analysis and offer recommendations for future work. The work ran successfully in both a Python script and a Jupyter Notebook environment.

DATA PREPARATION AND CLEANING

The dataset contains information about the passengers aboard the Titanic, including details such as age, passenger class, gender, fare, and survival status. Since real-world data rarely comes pre-cleaned, the first step was to address missing values and convert categorical variables to numerical values. For example:

- I filled missing values in the Age column with the median age of the dataset.
- The Cabin column, which had too many missing entries, was dropped entirely.
- Categorical values for Sex and Embarked were converted into numerical codes to simplify analysis.

This careful data cleaning was crucial as it ensured that the subsequent visualizations provided clear,

trustworthy insights into the patterns hidden in the dataset.

EXPLORATORY DATA ANALYSIS (EDA)

Once the data was cleaned, I used Plotly to generate a series of interactive graphs. Each of the 12 visualizations focuses on a different aspect of the dataset:

1. AGE DISTRIBUTION (HISTOGRAM):

This visualization shows how the passengers' ages are distributed on the Titanic. It reveals that most passengers were within certain age ranges (often concentrated between 20 and 40 years). Such insights help in understanding if age played any role in survival chances.

2. SURVIVAL RATE BY PASSENGER CLASS (BAR CHART):

One of the first patterns discovered was that survival rates varied significantly between classes. The graph clearly indicates that passengers in first class had a higher chance of survival compared to those in second and third classes. This finding points to possible socio-economic factors influencing survival.

3. SURVIVAL RATE BY GENDER (PIE CHART):

The pie chart focusing on gender demonstrates that women had a significantly higher survival rate compared to men. This observation supports the common narrative that “women and children first” influenced evacuation priorities during the disaster.

4. FARE VS. AGE SURVIVAL COMPARISON (SCATTER PLOT):

In this visualization, fare and age were plotted with survival as a color-coded factor. This graph not only shows variations in fare prices across

different age groups but also suggests that older passengers often paid higher fares, and these factors might correlate with their survival outcomes.

5. AGE DISTRIBUTION BY CLASS (BOX PLOT):

By dividing passengers by class, the box plot effectively highlights differences in age distribution. Outliers in each class were easily spotted, offering insights into the range of ages in each cabin class and hinting at the distinct demographics within each group.

6. SURVIVAL RATE BY EMBARKATION POINT (BAR CHART):

The port from which passengers embarked (using the Embarked variable) might also play a role in their survival. This visualization shows survival variations based on embarkation points (e.g., Southampton, Cherbourg, Queenstown). Although the differences are not as stark as those

seen for gender or class, they add another layer of understanding to the data.

7. FARE DISTRIBUTION BY CLASS (VIOLIN PLOT):

This visualization shows the distribution of fares paid by passengers, segmented by their class. The violin plot reveals that higher-class passengers not only paid more on average but also had a wider range of fare values. This can help explain differences in service quality or amenities provided during the voyage.

8. FARE DISTRIBUTION AMONG PASSENGERS (HISTOGRAM):

An overall histogram of fares shows the skewness in fare values. There is a clustering of fares at the lower end, which reflects the large number of third-class passengers, while a few outliers indicate the high fares paid by first-class travelers.

9. SURVIVAL RATE BASED ON SIBLINGS/SPOUSES ABOARD (BAR CHART):

Family relationships sometimes affect survival chances. This chart examines survival rates based on the number of siblings or spouses a passenger had aboard. Patterns here can indicate whether traveling with family increased a passenger's likelihood of surviving.

10. AGE VS. SURVIVAL COLORED BY CLASS (SCATTER PLOT):

Another angle was to look at survival in relation to both age and class. In this plot, data points are colored by passenger class. The visualization provides a more nuanced view, showing that younger passengers in higher classes had better survival rates, while variations in older age groups provide additional context.

11. SURVIVAL RATE BASED ON PARENTS/CHILDREN (BAR CHART):

Similar to the chart for siblings/spouses, this graph explores whether having parents or children on board influenced survival. Family size and structure can affect the dynamics during evacuation, and these trends are critical for understanding the role of group behavior in survival.

12. FARE VS. AGE RELATIONSHIP BY CLASS (SCATTER PLOT) AND A SUNBURST CHART:

The final visualizations combine multiple dimensions. One scatter plot examines how fare and age relate across different classes, while a sunburst chart illustrates relations among passenger class, gender, and survival. This multi-layered approach offers a holistic view of the dynamics among different variables.

Each of these visualizations offers unique insights. Together, they paint a picture of how factors like

socioeconomic status (represented by class and fare), age, gender, and family presence played roles in the survival of Titanic passengers.

KEY FINDINGS AND INSIGHTS

After generating and analyzing the visualizations, several clear patterns emerged:

- INFLUENCE OF SOCIOECONOMIC STATUS:

First-class passengers not only enjoyed higher fares but also had much higher survival rates. This suggests that having more resources or a higher social standing likely provided better access to safety measures during the disaster.

- GENDER DIFFERENCES:

Women consistently showed higher survival rates, which aligns with historical records suggesting that

women were given priority during the rescue. This result was clearly depicted in the pie and bar charts.

- AGE-RELATED TRENDS:

While the age distribution reveals that most passengers were young to middle-aged adults, the survival advantage for younger passengers (particularly those traveling in first class) stood out in the scatter plots and box plots. This trend may be attributed to general physical resilience and preferential treatment during evacuation protocols.

- FAMILY INFLUENCE:

The analyses based on the number of siblings/spouses and parents/children on board indicate that family groups might have influenced survival chances. In some cases, passengers traveling with a larger family group had better survival outcomes, possibly due to mutual assistance during the crisis.

- **BOARDING LOCATION:**

Although not as significant as class or gender, the embarkation point did show slight variations in survival rates. This could be due to differences in passenger demographics or even in the distribution of cabins among passengers from different ports.

- **FARE DISTRIBUTION AND ITS IMPACT:**

The fare-related visualizations highlight a clear disparity between different classes, suggesting that financial means played an important role in shaping the travel experience. High fares among first-class travelers correlated with higher survival rates, whereas the clustering of lower fares among third-class passengers reminds us of the financial hardships many of these travelers faced.

Overall, the analysis provides compelling evidence that, even with basic EDA techniques, important

relationships and trends in the Titanic dataset can be uncovered. The visualizations not only serve to tell the story of what happened aboard the Titanic but also provide insights into broader themes like inequality, social structure, and the importance of family and support networks in crisis situations.

RECOMMENDATIONS FOR FUTURE WORK

Based on the findings from this project, I have several recommendations:

1. DEEPER ANALYSIS USING ADVANCED TECHNIQUES:

Although this project used simple EDA techniques, future work could explore predictive modeling to estimate survival chances. Algorithms like logistic regression, decision trees, or random forests could be applied with proper feature

engineering to see if survival can be predicted with higher accuracy.

2. FEATURE ENGINEERING:

Future work might benefit from creating new features that capture relationships between existing variables—for example, calculating family size from the combined number of siblings/spouses and parents/children, or even exploring interaction terms between age and fare. This could offer a deeper understanding of the factors influencing survival.

3. GEOGRAPHICAL ANALYSIS:

While I looked at the embarkation points, further analysis could include a geographical perspective. Mapping the origins or even some socio-economic metrics of the ports of embarkation might provide additional layers of insight into the passenger demographics.

4. INTERACTIVE DASHBOARDS FOR BROADER ENGAGEMENT:

The use of Dash and Plotly has already transformed static analysis into interactive dashboards. Enhancing these dashboards with additional filters (such as gender, family size, or class) could allow non-technical users to explore the data themselves, making the insights more accessible to a wider audience.

5. COMPARATIVE ANALYSIS OVER TIME:

Although the Titanic dataset is historical, comparing it with more recent datasets on travel safety could uncover trends in how crises are handled today. This could lead to a broader study on how technology, policy, and social factors have evolved.

6. VISUALIZATION ENHANCEMENTS:

Future projects could integrate animated visualizations or even

geospatial mapping to visualize the disaster's impact more dynamically. Such enhancements can serve to both educate and engage viewers in a more immersive manner.

CONCLUSION

This Titanic dataset analysis has been a valuable learning experience. By leveraging Python libraries such as pandas for data manipulation, along with Plotly and Dash for visualization and interactivity, I was able to uncover several meaningful insights without resorting to complex machine learning algorithms. The project reinforces key historical narratives—such as the influence of social class and gender on survival—and highlights the importance of good data cleaning and visualization practices.

I believe that even basic EDA techniques can tell a powerful story about historical events. The insights gathered from the Titanic dataset remind us of the many factors that

come into play during a crisis and underscore how data-when carefully analyzed-can help us understand and learn from the past. The recommendations provided here point towards exciting opportunities for further exploration and analysis.

I look forward to applying these lessons in future projects and using more advanced techniques to deepen my understanding. This project not only demonstrates my technical proficiency with tools like Pandas, Plotly, and Dash but also highlights my ability to derive actionable insights from historical data. I am excited to share this work on LinkedIn and engage with professionals who are equally passionate about data-driven insights and storytelling.

This summary was prepared as part of my personal project to explore the Titanic dataset. The findings and recommendations presented here provide

a solid foundation for future data analysis projects and serve as an example of clear, human-friendly communication of complex data insights.