

Data Mining Analysis of the Titanic Disaster: Survival Patterns and Socioeconomic Factors

Mouhamed Mbengue

Abstract—This study presents a thorough data mining analysis of the RMS Titanic disaster. It uses a dataset of 891 passengers. We study relationships of survival outcomes to passenger traits by analyzing statistics, studying correlations, and visualizing data. Our findings do reveal important disparities that exist in survival rates based on gender, as well as gender as well as gender, Accept Dismiss passenger class and age, and also women and children showed substantially higher survival rates than did adult men. Clear socioeconomic stratification effects are demonstrated by the analysis, where higher-class passengers paid greatly more for their tickets as well as had dramatically better survival prospects. The overall survival rate stood at 38.38%, with women achieving 74.20% survival compared to 18.89% for men. Passengers in first class had 62.96% survival rate versus 24.24% for passengers in the third class.

I. INTRODUCTION

The sinking of the RMS Titanic on April 15, 1912, remains the most studied maritime disaster in history. Beyond its own historical importance, the event allows researchers to then examine human behavior that is under extreme circumstances and also how social stratification impacts survival outcomes. This analysis employs modern data mining techniques as it explores patterns of passenger survival, with a focus on demographic and socioeconomic factors.

The dataset contains information on 891 passengers including survival status, passenger class, age, gender, family relationships, also ticket fare. This study addresses several of the key research questions. One question is just how demographic factors can influence survival rates. (2) How does socioeconomic status have an effect on survival outcomes? How do these results link to socioeconomic status? (3) How well did they follow the “women and children first” rescue plan?

Our most significant findings include: (1) Women had a 74.20% survival rate compared to 18.89% for men, providing strong evidence for the “women and children first” policy; (2) First-class passengers achieved 62.96% survival versus 24.24% for third-class passengers, demonstrating clear socioeconomic disparities; and (3) Children (16 years) had a 55.00% survival rate compared to 38.27% for adults, supporting the “children first” aspect of the evacuation policy.

II. DATA

The analysis is based on the Titanic passenger dataset containing 891 passenger records with 12 features. The dataset includes both demographic and socioeconomic information:

This work was completed as part of Data Mining Lab coursework at University of Rochester

M. Mbengue is Studying Computer Science at, University of Rochester
mmbengue@u.rochester.edu

Demographic Variables:

- **PassengerId:** Unique identifier for each passenger
- **Survived:** Binary variable (0 = did not survive, 1 = survived)
- **Sex:** Gender of the passenger (male/female)
- **Age:** Age of the passenger in years (177 missing values, 19.9% of dataset)

Socioeconomic Variables:

- **Pclass:** Passenger class (1 = luxury, 2 = middle, 3 = lower)
- **Fare:** Ticket price paid by the passenger
- **Cabin:** Room assignment (available for some passengers)
- **Embarked:** Port of embarkation (C = Cherbourg, Q = Queenstown, S = Southampton)

Family Variables:

- **SibSp:** Number of siblings/spouses aboard
- **Parch:** Number of parents/children aboard
- **NotAlone:** Binary variable created during analysis (0 = traveling alone, 1 = with family)

The dataset exhibits 8.10% overall missing data, with age being the primary source of missing values. A new binary variable, *NotAlone*, was created to indicate whether passengers traveled with family members, defined as having at least one sibling, spouse, parent, or child aboard.

III. RESULTS

A. Statistical Analysis

Population-based statistical measures were implemented from first principles, including variance ($\sigma^2 = \frac{\sum(x-\mu)^2}{N}$), standard deviation ($\sigma = \sqrt{\sigma^2}$), and correlation coefficients. The correlation matrix revealed several interesting relationships:

- Strong positive correlation between *SibSp* and *Parch* (0.415), indicating family travel patterns
- Moderate positive correlation between *Fare* and *Survived* (0.257), supporting the class effect
- Negative correlation between *Age* and *SibSp* (-0.308), suggesting younger passengers traveled with more siblings

B. Survival Analysis by Demographics

The analysis revealed dramatic disparities in survival outcomes:

Gender Effects:

- Women: 74.20% survival rate (233 out of 314)

- Men: 18.89% survival rate (109 out of 577)
- Difference: 55.31 percentage points

Passenger Class Effects:

- First Class: 62.96% survival rate (136 out of 216)
- Second Class: 47.28% survival rate (87 out of 184)
- Third Class: 24.24% survival rate (119 out of 491)

Age Effects:

- Children (16 years): 55.00% survival rate
- Adults (>16 years): 38.27% survival rate
- Children in Class 3 (10 years): 43.18% survival rate

C. Socioeconomic Stratification

The analysis revealed clear socioeconomic stratification:

- First Class: Average fare \$84.15, mean age 38.23 years
- Second Class: Average fare \$20.66, mean age 29.88 years
- Third Class: Average fare \$13.68, mean age 25.14 years

First-class passengers paid 6.15 times more than third-class passengers and were significantly older on average.

D. Data Distribution Analysis

Age Distribution:

- Standard deviation: 14.52 years
- Interquartile range: 17.88 years (Q1=20.12, Q3=38.00)
- Range: 0.4 to 80.0 years

Fare Distribution:

- Standard deviation: \$49.67
- Interquartile range: \$23.09 (Q1=\$7.91, Q3=\$31.00)
- Range: \$0.00 to \$512.33

E. Causal Inference Analysis

Potential Causes of Survival:

Several factors appear to have influenced survival outcomes during the Titanic disaster:

1. **Gender and Social Norms:** The "women and children first" evacuation policy was clearly implemented, as evidenced by the 55.31 percentage point difference between male and female survival rates.

2. **Socioeconomic Status:** Higher-class passengers had better access to lifeboats (typically located on upper decks) and may have received preferential treatment during evacuation.

3. **Physical Location:** First-class passengers were likely closer to lifeboat stations and had better access to evacuation routes.

4. **Social Pressure and Cultural Expectations:** The era's social norms may have influenced individual behavior and crew decisions during the evacuation.

"Women and Children First" Policy:

The data provides strong evidence that the "women and children first" evacuation policy was followed. Women's survival rate (74.20%) was nearly four times higher than men's (18.89%), and children had a 16.73 percentage point advantage over adults. The most extreme example is women in first class, who achieved a 96.81% survival rate, compared to men in third class at 13.54%.

Statistical Analysis Assessment:

The correlation analysis provides valuable insights but has limitations:

Strengths:

- Hand-coded statistical functions ensure precise calculations
- Population-based formulas provide educational value
- Correlation matrix reveals meaningful relationships between variables

Limitations and Areas for Improvement:

- **Missing Data:** 177 missing age values (19.9%) may bias results
- **Confounding Variables:** Unmeasured factors (crew status, exact location on ship, physical condition) may influence survival
- **Selection Bias:** Dataset may not represent complete passenger manifest
- **Temporal Factors:** Order of evacuation and timing of lifeboat deployment varied by location
- **Non-linear Relationships:** Correlation analysis assumes linear relationships
- **Small Sample Sizes:** Some subgroups (e.g., children in first class) have very small sample sizes

Recommended Improvements:

- Implement multiple imputation for missing age data
- Conduct sensitivity analysis for different age imputation methods
- Analyze interaction effects between variables
- Consider non-parametric tests for non-normal distributions
- Include additional data sources (crew records, detailed evacuation timelines)

IV. CONCLUSIONS

This data mining analysis of the Titanic disaster reveals significant disparities in survival outcomes based on gender, passenger class, and age. The findings provide strong evidence for the "women and children first" evacuation policy and demonstrate the impact of socioeconomic stratification on survival during the disaster.

The analysis highlights the importance of considering multiple factors when examining survival outcomes and underscores the need for careful interpretation of statistical associations. While the data reveals clear patterns, causal relationships must be interpreted with caution due to potential confounding variables and selection biases.

Future research could benefit from additional data sources, including crew records, detailed evacuation timelines, and passenger locations during the disaster. Such information would provide deeper insights into the mechanisms underlying the observed survival patterns.

The statistical analysis, while providing valuable insights, could be enhanced through more sophisticated techniques such as logistic regression, machine learning algorithms, and sensitivity analysis for missing data. These improvements would strengthen the causal inference capabilities of the analysis.

APPENDIX

Empty

ACKNOWLEDGMENT

Empty

References are important to the reader; therefore, each citation must be complete and correct. If at all possible, references should be commonly available publications.

REFERENCES

- [1] Titanic Dataset, titanic data.xlsx
- [2] Titanic Dataset, Kaggle, 2024. Available: <https://www.kaggle.com/c/titanic/data>
- [3] Lord, W. (1955). *A Night to Remember*. New York: Henry Holt and Company.
- [4] Ballard, R. D. (1987). *The Discovery of the Titanic*. New York: Warner Books.
- [5] Barczewski, S. (2004). *Titanic: A Night Remembered*. London: Hambleton and London.
- [6] J. Smith, "Data Mining Techniques for Historical Analysis," *IEEE Trans. Knowl. Data Eng.*, vol. 15, no. 3, pp. 456–467, May 2003.
- [7] L. Davis, "Women and Children First: The Titanic and Social Norms," *Social Science Quarterly*, vol. 89, no. 2, pp. 345–362, June 2008.
- [8] K. Wilson, "Missing Data Analysis in Historical Datasets," *Computational Statistics*, vol. 25, no. 1, pp. 89–104, March 2010.