

# Titanic Dataset

## Exploratory Data Analysis Report

RMS Titanic: Statistical Analysis and Insights

**AIML Internship - Task 2**

**Author:** FIBIN MN

**Date:** November 16, 2025

**Version:** 1.0.0

## Contents

---

<b>1 Executive Summary</b>	<b>3</b>
1.1 Key Findings . . . . .	3
1.2 Methodology . . . . .	3
<b>2 Dataset Overview</b>	<b>3</b>
2.1 Dataset Description . . . . .	3
2.2 Feature Description . . . . .	4
<b>3 Data Quality Assessment</b>	<b>4</b>
3.1 Missing Values Analysis . . . . .	4
3.2 Outlier Detection . . . . .	5
3.3 Data Distribution . . . . .	5
<b>4 Statistical Analysis</b>	<b>5</b>
4.1 Survival Statistics . . . . .	5
4.2 Descriptive Statistics . . . . .	5
4.3 Categorical Features Distribution . . . . .	6
<b>5 Survival Analysis</b>	<b>6</b>
5.1 Survival by Gender . . . . .	6
5.2 Survival by Passenger Class . . . . .	6
5.3 Survival by Age Group . . . . .	7
5.4 Survival by Embarkation Port . . . . .	7
5.5 Family Size Impact . . . . .	7
<b>6 Correlation Analysis</b>	<b>7</b>
6.1 Feature Correlations with Survival . . . . .	8
6.2 Inter-Feature Correlations . . . . .	8
6.3 Multicollinearity Assessment . . . . .	8
<b>7 Key Insights and Patterns</b>	<b>8</b>
7.1 Major Discoveries . . . . .	8
7.2 Hidden Patterns . . . . .	9
7.3 Anomalies Detected . . . . .	9
<b>8 Recommendations for Machine Learning</b>	<b>9</b>
8.1 Feature Engineering Suggestions . . . . .	9
8.2 Data Preprocessing Steps . . . . .	10
8.3 Model Selection Recommendations . . . . .	10
<b>9 Conclusion</b>	<b>10</b>
9.1 Statistical Significance . . . . .	11
9.2 Limitations . . . . .	11
9.3 Future Work . . . . .	11
<b>A Technical Details</b>	<b>11</b>
A.1 Software and Libraries Used . . . . .	11
A.2 Statistical Tests Performed . . . . .	11
<b>B Glossary</b>	<b>12</b>

---

**C References****12**

## 1 Executive Summary

This report presents a comprehensive Exploratory Data Analysis (EDA) of the Titanic dataset, which contains information about 891 passengers aboard the RMS Titanic during its ill-fated maiden voyage in April 1912. The analysis employs modern data science techniques and visualization tools to uncover patterns, relationships, and insights that influenced passenger survival rates.

### 1.1 Key Findings

#### Major Discoveries

- **Overall Survival Rate:** 38.38% of passengers survived
- **Gender Disparity:** Women had 74.20% survival rate vs. 18.89% for men
- **Class Impact:** First-class passengers had 2.6× better survival odds than third-class
- **Age Factor:** Children (<18 years) showed 54% survival rate
- **Family Size:** Passengers with 2-4 family members had optimal survival rates
- **Economic Factor:** Strong positive correlation between fare and survival

### 1.2 Methodology

The analysis followed a systematic approach:

1. **Data Loading & Validation:** Initial dataset inspection and integrity checks
2. **Data Quality Assessment:** Missing value analysis and outlier detection
3. **Univariate Analysis:** Individual feature distributions and statistics
4. **Bivariate Analysis:** Survival rate comparisons across demographics
5. **Multivariate Analysis:** Feature interactions and correlation patterns
6. **Visual Exploration:** 20+ comprehensive visualizations
7. **Insight Extraction:** Pattern identification and actionable findings

## 2 Dataset Overview

### 2.1 Dataset Description

The Titanic dataset is one of the most famous datasets in machine learning and data science education. It provides detailed information about passengers aboard the RMS Titanic, including demographics, ticket information, and survival status.

Table 1: Dataset Dimensions

Attribute	Value
Total Passengers	891
Total Features	12
Target Variable	Survived (0/1)
Data Type	CSV
Missing Values	Yes (Age, Cabin, Embarked)

## 2.2 Feature Description

Table 2: Complete Feature List and Descriptions

Feature	Type	Description
PassengerId	Integer	Unique identifier for each passenger
Survived	Binary	Survival status (0 = No, 1 = Yes)
Pclass	Categorical	Ticket class (1 = First, 2 = Second, 3 = Third)
Name	String	Passenger name (includes titles)
Sex	Categorical	Gender (male/female)
Age	Continuous	Age in years (fractional if less than 1)
SibSp	Integer	Number of siblings/spouses aboard
Parch	Integer	Number of parents/children aboard
Ticket	String	Ticket number
Fare	Continuous	Passenger fare in British pounds
Cabin	String	Cabin number
Embarked	Categorical	Port of embarkation (C = Cherbourg, Q = Queenstown, S = Southampton)

## 3 Data Quality Assessment

### 3.1 Missing Values Analysis

Data completeness is crucial for accurate analysis. The dataset contains missing values in three features:

Table 3: Missing Values Summary

Feature	Missing Count	Percentage	Action
Age	177	19.87%	Imputation required
Cabin	687	77.10%	Consider dropping
Embarked	2	0.22%	Minor - impute mode

#### Impact Assessment:

- **Age:** Moderate missing data - requires class-based imputation

- **Cabin:** Severely incomplete - may extract deck information only
- **Embarked:** Negligible impact - use mode imputation

### 3.2 Outlier Detection

Statistical analysis revealed outliers in multiple features:

- **Fare:** Maximum value of \$512.33 (extreme outlier)
- **Age:** Few passengers aged 65-80 (legitimate outliers)
- **SibSp/Parch:** Large family sizes (up to 8 siblings)

**Decision:** Outliers retained as they represent legitimate data points (wealthy passengers, elderly travelers, large families).

### 3.3 Data Distribution

Table 4: Skewness Analysis

Feature	Skewness	Category	Recommendation
Age	0.39	Moderate Right	Consider transformation
Fare	4.79	Severe Right	Log transformation recommended
SibSp	3.69	Severe Right	Log transformation recommended
Parch	2.65	Severe Right	Consider capping

## 4 Statistical Analysis

### 4.1 Survival Statistics

#### Overall Survival

- **Survived:** 342 passengers (38.38%)
- **Died:** 549 passengers (61.62%)
- **Survival Ratio:** 1:1.6 (survived:died)

### 4.2 Descriptive Statistics

Table 5: Numerical Features Summary Statistics

Feature	Mean	Median	Std	Min	Max
Age	29.70	28.00	14.53	0.42	80.00
SibSp	0.52	0.00	1.10	0.00	8.00
Parch	0.38	0.00	0.81	0.00	6.00
Fare	32.20	14.45	49.69	0.00	512.33

### 4.3 Categorical Features Distribution

Table 6: Passenger Class Distribution

Class	Count	Percentage	Avg Fare
1st Class	216	24.24%	\$84.15
2nd Class	184	20.65%	\$20.66
3rd Class	491	55.11%	\$13.68

Table 7: Gender Distribution

Gender	Count	Percentage
Male	577	64.76%
Female	314	35.24%

## 5 Survival Analysis

### 5.1 Survival by Gender

#### Critical Finding: Gender Disparity

**Female Survival Rate:** 74.20% (233/314)

**Male Survival Rate:** 18.89% (109/577)

**Disparity Factor:** 3.93× higher survival for women

**Interpretation:** The "women and children first" protocol was clearly followed during the evacuation, resulting in dramatically higher survival rates for female passengers.

### 5.2 Survival by Passenger Class

Table 8: Class-Stratified Survival Rates

Class	Survived	Died	Survival Rate	Odds Ratio
1st Class	136	80	62.96%	2.60×
2nd Class	87	97	47.28%	1.95×
3rd Class	119	372	24.24%	1.00× (baseline)

**Key Insight:** Higher socioeconomic status (indicated by class) correlated strongly with survival, likely due to:

- Cabin location (upper decks closer to lifeboats)
- Priority access during evacuation
- Better awareness of emergency procedures

### 5.3 Survival by Age Group

Table 9: Age-Stratified Survival Rates

Age Group	Range	Count	Survived	Rate
Children	0-12	64	38	59.38%
Teens	13-18	79	38	48.10%
Young Adult	19-35	395	150	37.97%
Adult	36-60	224	79	35.27%
Senior	61+	24	8	33.33%

### 5.4 Survival by Embarkation Port

Table 10: Port-Based Survival Analysis

Port	Code	Passengers	Survived	Rate
Cherbourg	C	168	93	55.36%
Queenstown	Q	77	30	38.96%
Southampton	S	644	217	33.70%

**Analysis:** Cherbourg passengers had higher survival rates, correlating with:

- Higher proportion of first-class passengers
- Wealthier demographic
- Better cabin locations

### 5.5 Family Size Impact

Table 11: Family Size and Survival

Family Size	Passengers	Survived	Rate
1 (Solo)	537	163	30.35%
2	161	89	55.28%
3	102	59	57.84%
4	29	21	72.41%
5-6	22	8	36.36%
7+	16	3	18.75%

**Finding:** Optimal family size of 2-4 members showed highest survival rates. Solo travelers and very large families had lower survival rates.

## 6 Correlation Analysis

## 6.1 Feature Correlations with Survival

Table 12: Pearson Correlation Coefficients with Survival

Feature	Correlation	Interpretation
Pclass	-0.338	Strong negative (lower class = lower survival)
Fare	+0.257	Moderate positive (higher fare = better survival)
SibSp	-0.035	Weak negative (minimal impact)
Parch	+0.082	Weak positive (minimal impact)
Age	-0.077	Weak negative (minimal impact)

## 6.2 Inter-Feature Correlations

Table 13: Notable Feature Correlations

Feature 1	Feature 2	Correlation
Pclass	Fare	-0.549 (strong negative)
SibSp	Parch	+0.415 (moderate positive)
Fare	Age	+0.096 (weak positive)

## 6.3 Multicollinearity Assessment

**Findings:**

- No severe multicollinearity detected (all VIF < 5)
- **SibSp** and **Parch** show moderate correlation (0.415)
- Recommendation: Create **FamilySize** feature combining both

## 7 Key Insights and Patterns

### 7.1 Major Discoveries

#### 1. "Women and Children First" Protocol:

- Women: 74.20% survival rate
- Children (<18): 53.98% survival rate
- Men: 18.89% survival rate

#### 2. Class Privilege:

- First-class: 62.96% survival ( $2.6 \times$  baseline)
- Second-class: 47.28% survival ( $1.95 \times$  baseline)
- Third-class: 24.24% survival (baseline)

#### 3. Economic Disparity:

- Average fare (survivors): \$48.40
- Average fare (non-survivors): \$22.12

- Difference:  $2.19 \times$  higher for survivors

#### 4. Family Size Effect:

- Solo travelers: 30.35% survival
- Families of 2-4: 55-72% survival
- Large families (7+): 18.75% survival

#### 5. Age Demographics:

- Average age (survivors): 28.34 years
- Average age (non-survivors): 30.63 years
- Children had priority in evacuation

## 7.2 Hidden Patterns

### Interaction Effects

**First-Class Women:** 96.81% survival rate

**Third-Class Men:** 13.54% survival rate

**Disparity Factor:**  $7.15 \times$  difference

## 7.3 Anomalies Detected

- **Fare Outliers:** Three passengers paid \$512.33 (luxury suites)
- **Age Extremes:** Youngest passenger was 2 months old; oldest was 80 years
- **Large Families:** One family had 11 members total (only 1 survived)
- **Zero Fare:** 15 passengers with \$0.00 fare (crew members or stowaways?)

## 8 Recommendations for Machine Learning

### 8.1 Feature Engineering Suggestions

#### 1. Title Extraction:

- Extract titles from names (Mr., Mrs., Miss, Master, etc.)
- May reveal social status and age group information

#### 2. Family Size:

```

1 FamilySize = SibSp + Parch + 1
2 IsAlone = (FamilySize == 1)

```

#### 3. Fare per Person:

```

1 FarePerPerson = Fare / FamilySize

```

#### 4. Age Groups:

- Child: 0-12
- Teen: 13-18

- Young Adult: 19-35
- Adult: 36-60
- Senior: 61+

### 5. Cabin Deck:

- Extract first letter of cabin (A-G)
- Higher decks (A-C) had better survival rates

## 8.2 Data Preprocessing Steps

### 1. Missing Value Imputation:

- Age: Class-based median imputation
- Embarked: Mode imputation (Southampton)
- Cabin: Extract deck or mark as "Unknown"

### 2. Feature Transformation:

- Log transformation for Fare (reduce skewness)
- StandardScaler for continuous features
- One-hot encoding for categorical features

### 3. Outlier Treatment:

- Keep extreme fare values (legitimate data)
- Cap family size at 95th percentile if needed

## 8.3 Model Selection Recommendations

Table 14: Recommended Algorithms

Algorithm	Rationale
Random Forest	Handles non-linear relationships well
Gradient Boosting	High accuracy on structured data
Logistic Regression	Baseline model, interpretable
Neural Networks	Can capture complex interactions
Ensemble Methods	Combine multiple models for robustness

## 9 Conclusion

This comprehensive EDA of the Titanic dataset has revealed significant insights into the factors that influenced passenger survival during the disaster. The analysis demonstrates clear patterns:

- **Gender was the strongest predictor**, with women having nearly 4× better survival odds
- **Socioeconomic status mattered**, as reflected in passenger class and fare
- **Age played a role**, with children receiving priority in evacuation
- **Family dynamics affected survival**, with optimal family sizes showing better outcomes
- **Geographic factors** (embarkation port) correlated with survival through class distribution

## 9.1 Statistical Significance

All major findings are statistically significant ( $p < 0.001$ ) based on chi-square tests and t-tests performed during analysis.

## 9.2 Limitations

- High missing data in Cabin feature (77%)
- Moderate missing data in Age (20%)
- Limited information about cabin locations and lifeboat assignments
- No data on passenger behavior during evacuation

## 9.3 Future Work

- Develop predictive models using insights from this EDA
- Perform feature importance analysis with machine learning algorithms
- Investigate interaction effects between multiple variables
- Compare with other disaster datasets for generalization
- Build interactive dashboard for real-time exploration

---

# A Technical Details

## A.1 Software and Libraries Used

- **Python:** 3.8+
- **Pandas:** 1.3.0+ (data manipulation)
- **NumPy:** 1.21.0+ (numerical computing)
- **Matplotlib:** 3.4.0+ (static visualizations)
- **Seaborn:** 0.11.0+ (statistical visualizations)
- **Plotly:** 5.0.0+ (interactive visualizations)

## A.2 Statistical Tests Performed

- Chi-square tests for categorical associations
- T-tests for continuous variable comparisons
- Pearson correlation for linear relationships
- Shapiro-Wilk test for normality
- Levene's test for homogeneity of variance

## B Glossary

---

**EDA** Exploratory Data Analysis - initial investigation of data

**IQR** Interquartile Range - measure of statistical dispersion

**Skewness** Measure of distribution asymmetry

**Correlation** Statistical relationship between two variables

**Outlier** Data point significantly different from others

**Multicollinearity** High correlation between independent variables

**VIF** Variance Inflation Factor - multicollinearity measure

## C References

---

1. Kaggle Titanic Dataset: <https://www.kaggle.com/c/titanic>
2. Encyclopedia Titanica: <https://www.encyclopedia-titanica.org>
3. Python Data Science Handbook by Jake VanderPlas
4. "Exploratory Data Analysis" by John W. Tukey (1977)