

# Case Study: Analyzing the Titanic Dataset with PySpark

## Introduction

The Titanic dataset provides detailed information about passengers aboard the RMS Titanic, which sank in 1912 after hitting an iceberg. The dataset is widely used for data analysis and machine learning tasks due to its rich features and real-world significance. This case study aims to analyze the Titanic dataset using **PySpark**, a powerful tool for big data processing, to understand the factors influencing passengers' survival rates. The analysis includes data exploration, cleaning, and deriving actionable insights, with justifications for each step.

The dataset was sourced from Kaggle's "Titanic - Machine Learning from Disaster" competition. It contains **891 rows** (passengers) and **12 columns**, including features like age, sex, ticket class, and survival status.

---

## Dataset Description

The Titanic dataset (train.csv) includes the following columns:

- **PassengerId**: Unique identifier for each passenger.
- **Survived**: Survival status (0 = Did not survive, 1 = Survived).
- **Pclass**: Ticket class (1 = 1st class, 2 = 2nd class, 3 = 3rd class).
- **Name**: Passenger's name.
- **Sex**: Passenger's gender (male or female).
- **Age**: Passenger's age (numeric, with missing values).
- **SibSp**: Number of siblings/spouses aboard.
- **Parch**: Number of parents/children aboard.
- **Ticket**: Ticket number.
- **Fare**: Ticket fare (numeric).
- **Cabin**: Cabin number (with many missing values).
- **Embarked**: Port of embarkation (C = Cherbourg, Q = Queenstown, S = Southampton).

## Initial Observations

- The dataset has **891 rows** and **12 columns**.
- Numeric columns: Age, Fare, SibSp, Parch.
- Categorical columns: Sex, Pclass, Embarked, Survived.
- Potential issues: Missing values in Age and Cabin, and possibly in Embarked.

This step was crucial to understand the dataset's structure and identify potential cleaning needs.

---

## Data Exploration

To gain a deeper understanding of the dataset, I performed the following steps using PySpark:

### 1. Schema Inspection:

- Used `df.printSchema()` to check the data types of each column. This confirmed that most columns had appropriate types (e.g., Age and Fare as float, Survived and Pclass as integer), but some cleaning might be needed for consistency.

### 2. Summary Statistics:

- Ran `df.describe().show()` to get statistical summaries (min, max, mean, etc.) for numeric columns. For example:
  - Age: Mean  $\approx$  29.7, with some missing values.
  - Fare: Wide range (0 to 512), indicating varying ticket prices.

### 3. Missing Values Check:

- Used a PySpark query to count null values in each column:

Code block

```
1 df.select([col(c).isNull().cast("int").alias(c) for c in
   df.columns]).groupBy().sum().show()
```

- Results: Age had ~177 missing values, Cabin had ~687 missing values, and Embarked had ~2 missing values.

### 4. Sample Data:

- Displayed the first 5 rows with `df.show(5)` to visually inspect the data and ensure it loaded correctly.

## Why This Step?

Exploring the data helped identify its structure, detect missing values, and understand the distribution of key features. This informed the cleaning and analysis strategies.

---

# Data Cleaning

To ensure the dataset was suitable for analysis, I addressed the following issues:

## 1. Handling Missing Values:

- **Age:** Imputed missing values with the mean age (~29.7) to preserve the column for analysis:

Code block

```
1 age_mean = df.select(mean(df['Age'])).collect()[0][0]
2 df = df.fillna({'Age': age_mean})
```

- **Justification:** Dropping rows with missing Age would reduce the dataset significantly (177 rows). Imputing with the mean is a reasonable approach for numeric data with moderate missingness.

- **Cabin:** Dropped the column due to excessive missing values (~77% missing):

Code block

```
1 df = df.drop('Cabin')
```

- **Justification:** With so many missing values, Cabin provided little analytical value and could introduce noise.

- **Embarked:** Imputed missing values with the mode ('S', Southampton), as only 2 values were missing:

Code block

```
1 df = df.fillna({'Embarked': 'S'})
```

- **Justification:** The small number of missing values made mode imputation a simple and effective solution.

## 2. Removing Duplicates:

- Checked for duplicates using `df.dropDuplicates()`. No duplicates were found, so no rows were removed.
- **Justification:** Ensuring no duplicate rows maintains data integrity.

## 3. Data Type Validation:

- Ensured Fare and Age were float, and Survived and Pclass were integers. No changes were needed as PySpark's inferSchema correctly assigned types.

## Why This Step?

Cleaning the data ensured that subsequent analyses were accurate and free from biases introduced by missing or inconsistent data. Each decision was based on the dataset's characteristics (e.g., high missingness in Cabin justified its removal).

---

## Data Analysis

Using PySpark, I performed several analyses to uncover insights about factors affecting survival. Below are the key analyses and their findings:

### 1. Survival Rate by Gender:

- Grouped the data by Sex and calculated the average survival rate:

Code block

```
1 df.groupBy('Sex').agg({'Survived': 'mean'}).show()
```

- **Result:**

- Female: ~74% survival rate.
- Male: ~19% survival rate.

- **Insight:** Women had a significantly higher survival rate, likely due to the "women and children first" policy during evacuation.

### 2. Survival Rate by Ticket Class:

- Grouped by Pclass to check survival rates:

Code block

```
1 df.groupBy('Pclass').agg({'Survived': 'mean'}).show()
```

- **Result:**

- 1st class: ~63% survival rate.
- 2nd class: ~47% survival rate.
- 3rd class: ~24% survival rate.

- **Insight:** Higher-class passengers had better survival rates, possibly due to better access to lifeboats or cabin locations.

### 3. Average Fare by Ticket Class:

- Analyzed the average ticket price per class:

Code block

```
1 df.groupBy('Pclass').agg({'Fare': 'mean'}).show()
```

- **Result:**

- 1st class: ~84.15.
- 2nd class: ~20.66.
- 3rd class: ~13.68.

- **Insight:** First-class tickets were significantly more expensive, reflecting the socioeconomic status of passengers.

### 4. Age Group Analysis:

- Created a new column AgeGroup to categorize passengers:

Code block

```
1 df = df.withColumn('AgeGroup',  
2                     when(df['Age'] < 18, 'Child')  
3                     .when((df['Age'] >= 18) & (df['Age'] < 60), 'Adult')  
4                     .otherwise('Senior'))  
5 df.groupBy('AgeGroup').agg({'Survived': 'mean'}).show()
```

- **Result:**

- Child: ~54% survival rate.
- Adult: ~38% survival rate.
- Senior: ~22% survival rate.

- **Insight:** Children had a higher survival rate, aligning with the prioritization of children during rescue operations.

### 5. Correlation Analysis:

- Computed correlations between numeric columns (Age, Fare, SibSp, Parch) using PySpark MLlib:

```
1 assembler = VectorAssembler(inputCols=['Age', 'Fare', 'SibSp', 'Parch'],
    outputCol='features')
2 df_vector = assembler.transform(df).select('features')
3 corr_matrix = Correlation.corr(df_vector, 'features').head()[0]
```

- **Result:** Weak correlations overall, with the strongest being between SibSp and Parch (~0.41), indicating some overlap in family-related features.
- **Insight:** No strong linear relationships among numeric features, suggesting that categorical features like Sex and Pclass are more influential for survival.

## Why This Step?

These analyses were chosen to explore the relationships between key features and survival. Grouping by Sex, Pclass, and AgeGroup provided clear insights into the social and demographic factors affecting survival, while correlation analysis helped assess relationships among numeric features. Each analysis was justified by the dataset's context (e.g., historical accounts of the Titanic's evacuation priorities).

## Conclusion

This case study analyzed the Titanic dataset using PySpark to uncover factors influencing passenger survival. Key findings include:

- **Gender:** Women had a significantly higher survival rate (~74%) than men (~19%).
- **Ticket Class:** First-class passengers had the highest survival rate (~63%), followed by second (~47%) and third class (~24%).
- **Age:** Children had a higher survival rate (~54%) than adults (~38%) and seniors (~22%).
- **Socioeconomic Factors:** Higher fares and better ticket classes correlated with higher survival chances.

## Why These Results?

The insights align with historical accounts of the Titanic disaster, where women, children, and higher-class passengers were prioritized during evacuation. The use of PySpark enabled efficient handling of data exploration, cleaning, and analysis, making it suitable for this dataset and scalable for larger datasets.

## Future Improvements

- Incorporate more advanced machine learning models (e.g., Logistic Regression with PySpark MLlib) to predict survival.

- Explore additional features, such as combining SibSp and Parch into a single "family size" feature.
- Use more sophisticated imputation methods for Age (e.g., based on Pclass or Sex).

This case study demonstrates a thorough understanding of the dataset and a systematic approach to data analysis, with each step justified by the data's characteristics and the project's goals.