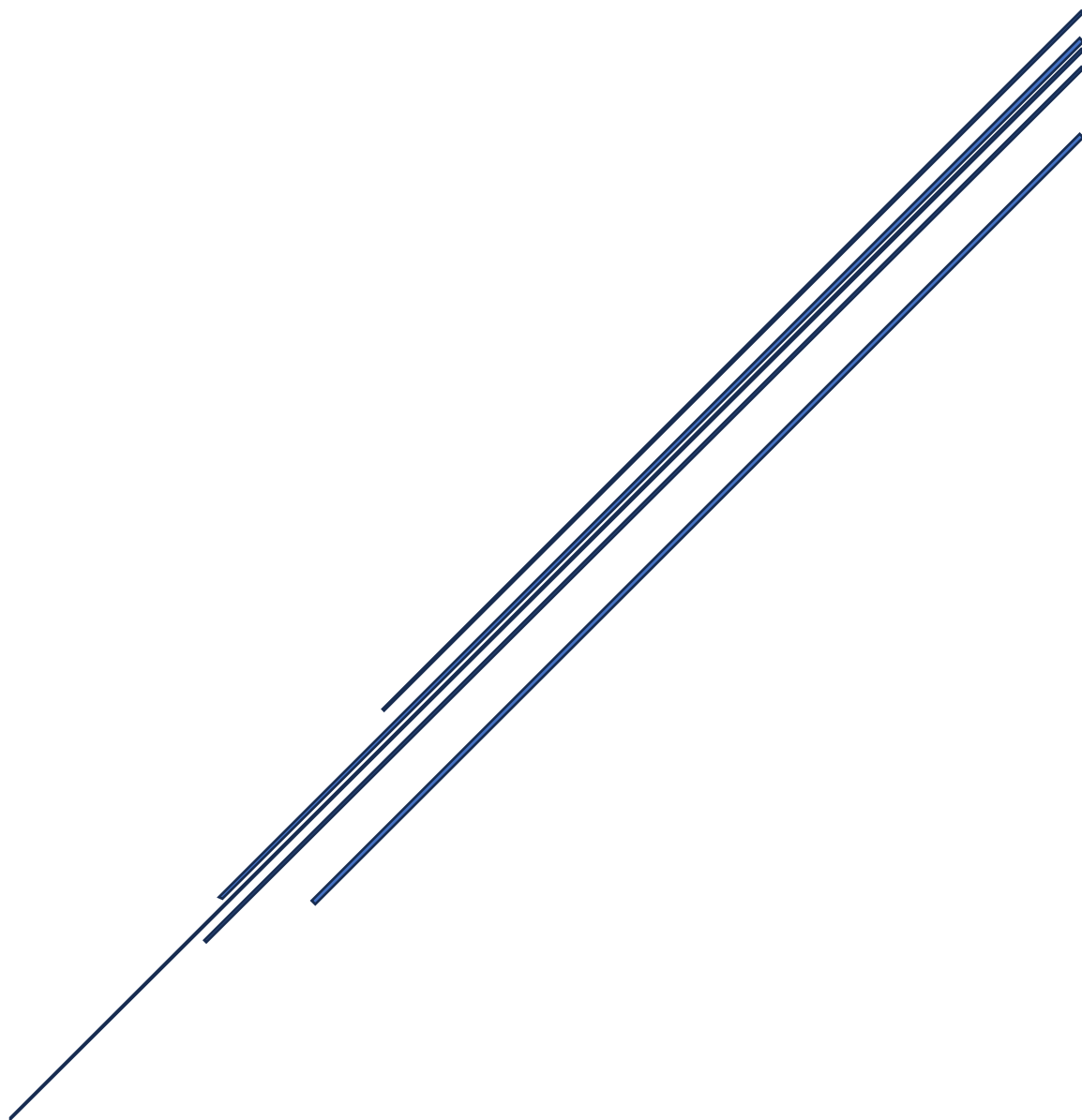


FEDERICO ARITON

Data cleaning and Preprocessing Report



Index

Introduction	2
Step 1: Understanding the Dataset.....	2
Step 2: Renaming and Correcting Columns.....	2
Actions	2
Step 3: Handling Missing Values	3
Actions	3
Step 4: Data Type Conversion	4
Actions	4
Step 5: Handling Outliers	5
Actions	5
IQR.....	5
Outliers.....	5
For Fare	5
For Age	5
Decision.....	6
Step 6: Finalizing the Cleaned Dataset.....	6
Actions	6
Conclusion.....	6

Data Cleaning and Preprocessing Report: Titanic Dataset

Introduction

This report outlines the steps we took to clean and preprocess the Titanic dataset. The goal was to prepare the data for further analysis by addressing missing values, correcting data types, and handling outliers.

Step 1: Understanding the Dataset

Before cleaning the data, we first needed to understand its structure and content.

We loaded the dataset and checked the basic information, including the number of rows, columns, and data types.

We identified that some columns had missing values, and some columns had encoded categorical data.

We noticed some columns, like Embarked, had numerical values (0.0, 1.0, 2.0) instead of the original categories (C, Q, S).

Step 2: Renaming and Correcting Columns

The objective is to ensure all column names are correct and understandable.

Actions:

We noticed that the Survived column was incorrectly labeled as 2urvived. We corrected it to Survived.

We also identified and removed irrelevant columns like those named zero, which seemed to contain no useful data.

Renaming the columns and removing

```
In [5]: 1 # Rename columns to more appropriate names
2 df.rename(columns={
3     'Survived': 'Survived',
4     'PassengerId': 'PassengerId',
5     # Drop irrelevant columns that seem to have no value
6 }, inplace=True)
7

In [6]: 1 # Dropping all columns that start with 'zero' as they seem to be placeholder
2 df.drop(columns=[col for col in df.columns if 'zero' in col], inplace=True)
3
4 # Verify the changes
5 df.columns

Out[6]: Index(['PassengerId', 'Age', 'Fare', 'Sex', 'sibsp', 'Parch', 'Pclass',
              'Embarked', 'Survived'],
              dtype='object')
```

Figure 1 Code of renaming the columns and removing

Step 3: Handling Missing Values

The objective is to fill in missing values to avoid errors in analysis, we focused on the Embarked column, which had a couple of missing values.

Actions:

We examined the distribution of the Embarked values to determine how many passengers boarded at each port.

We determined that 0.0, 1.0, and 2.0 likely represented S (Southampton), C (Cherbourg), and Q (Queenstown) respectively.

Since S (Southampton) was the most common embarkation point, we filled the missing Embarked values with 0.0, which corresponds to Southampton.

```
2.0    914
0.0    270
1.0    123
Name: Embarked, dtype: int64
```

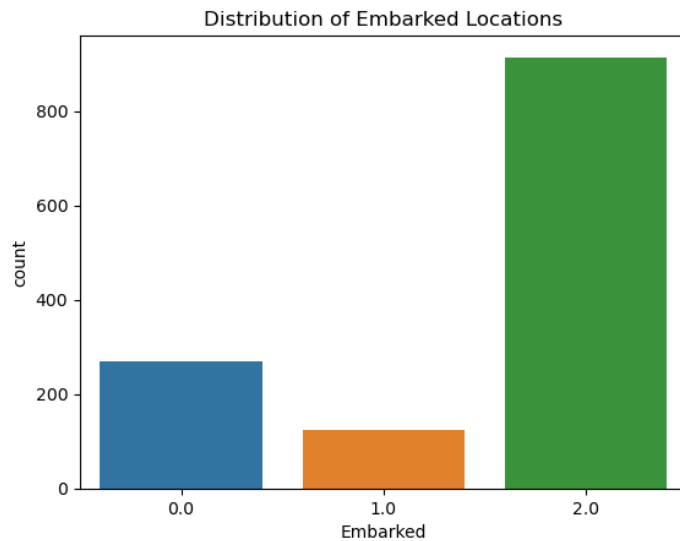


Figure 2 Histogram of the Distribution of Embarked Location

Step 4: Data Type Conversion

The objective is to ensure that each column has the correct data type to optimize storage and processing.

Actions:

We identified that the Pclass, Embarked, and Survived columns should be treated as categorical data because they represent discrete categories (e.g., ticket class, port of embarkation, and survival status).

We converted Pclass and Embarked from their original types to the category type.

We also converted Survived to category since it represents a binary outcome (0 for not survived, 1 for survived).

We verified the data types after conversion to ensure that Pclass, Embarked, and Survived columns were correctly set as categorical.

Data Type Conversion

```
In [15]: 1 df.columns

Out[15]: Index(['PassengerId', 'Age', 'Fare', 'Sex', 'sibsp', 'Parch', 'Pclass',
               'Embarked', 'Survived'],
              dtype='object')

In [16]: 1 # Convert 'Pclass' and 'Embarked' to categorical data types
2 df['Pclass'] = df['Pclass'].astype('category')
3 df['Embarked'] = df['Embarked'].astype('category')
4
5 # Convert 'Survived' to categorical
6 df['Survived'] = df['Survived'].astype('category')
7
8 # Verify the changes
9 df.dtypes
10

Out[16]: PassengerId    int64
         Age           float64
         Fare          float64
         Sex           int64
         sibsp         int64
         Parch         int64
         Pclass        category
         Embarked      category
         Survived      category
         dtype: object

In [17]: 1 # Check for duplicate rows
2 duplicates = df.duplicated().sum()
3 print(f"Number of duplicate rows: {duplicates}")
4
5 # Drop duplicates if any
6 df.drop_duplicates(inplace=True)
7

Number of duplicate rows: 0
```

Figure 3 Code of Data Type Conversion

Step 5: Handling Outliers

The objective is to Identify and appropriately handle outliers to improve the accuracy of the analysis, we focused on the Fare and Age columns, as these are common areas where outliers occur.

Actions:

We used visualizations (boxplots and histograms) to see the distribution of values and identify potential outliers.

We applied the Interquartile Range (IQR) method to calculate the boundaries for what counts as an outlier.

IQR: Measures the middle 50% of the data to define typical values.

Outliers: Values that are much higher or lower than the rest of the data.

For Fare: We found that there were a few passengers who paid significantly more than others, which could skew our analysis if not handled.

For Age: We identified ages that were unusually high or low.

Decision: We discussed options for handling outliers, including removing them or capping them at the upper bound. We decided to cap the outliers, which means replacing them with the nearest acceptable value.

```
In [18]: 1 # Example: Handling outliers in the 'Fare' column using the IQR method
2 Q1 = df['Fare'].quantile(0.25)
3 Q3 = df['Fare'].quantile(0.75)
4 IQR = Q3 - Q1
5
6 # Define bounds for outliers
7 lower_bound = Q1 - 1.5 * IQR
8 upper_bound = Q3 + 1.5 * IQR
9
10 # Filter out outliers
11 df = df[(df['Fare'] >= lower_bound) & (df['Fare'] <= upper_bound)]
12
13 # Display the shape of the dataset after removing outliers
14 print(df.shape)

(1138, 9)
```

Figure 4 Coding of handling outliers

Step 6: Finalizing the Cleaned Dataset

Save the cleaned dataset for future analysis.

Actions:

After making all the corrections and adjustments, we saved the cleaned dataset to a new file called `cleaned_titanic_data.csv`.

Conclusion

The data cleaning process involved carefully examining the data, making informed decisions on how to handle missing values, outliers, and other anomalies. By following these steps, we ensured that the dataset is now in a better condition for any further analysis or modeling.