# From Raw Data to Insights: A Complete Data Science Mini Project

## Objective

This project simulates a real-world data science workflow using the Titanic dataset from Kaggle (train.csv). The tasks include data preparation, transformation, exploratory data analysis (EDA), and visualization to derive meaningful insights about passenger survival.

# 1 Part 1: Data Preparation

## 1.1 Loading the Dataset

The dataset was loaded using pandas from the `train.csv` file.

```python
import pandas as pd
data = pd.read_csv('train.csv')
df = pd.DataFrame(data)
df.head()
```

**Output (First 5 Rows):**

| PassengerId | Survived | Pclass | Name | Sex | Age | SibSp | Parc |
|---|---|---|---|---|---|---|---|
| 1 | 0 | 3 | Braund, Mr. Owen Harris | male | 22.0 | 1 | |
| 2 | 1 | 1 | Cumings, Mrs. John Bradley | female | 38.0 | 1 | |
| 3 | 1 | 3 | Heikkinen, Miss. Laina | female | 26.0 | 0 | |
| 4 | 1 | 1 | Futrelle, Mrs. Jacques Heath | female | 35.0 | 1 | |
| 5 | 0 | 3 | Allen, Mr. William Henry | male | 35.0 | 0 | |

## 1.2 Inspecting Data Types and Identifying Features

The dataset contains 891 rows and 12 columns. Data types were inspected, and features were categorized.

```python
df = df.convert_dtypes()
df.dtypes
```

**Output:**

```
PassengerId      Int64
Survived         Int64
Pclass           Int64
Name             string
Sex              string
Age              Float64
SibSp            Int64
Parch            Int64
Ticket           string
Fare             Float64
Cabin            string
Embarked         string
dtype: object
```

**Categorical Features:**

- Name, Sex, Ticket, Cabin, Embarked

**Numerical Features:**

- PassengerId, Survived, Pclass, Age, SibSp, Parch, Fare

## 1.3  Handling Missing Values

Missing values were identified and handled appropriately.

```
df.isnull().sum()
```

**Output:**

```
PassengerId        0
Survived           0
Pclass             0
Name               0
Sex                0
Age              177
SibSp              0
Parch              0
Ticket             0
Fare               0
Cabin            687
Embarked           2
dtype: int64
```

**Missing Value Percentages:**

```
((df.isnull().sum() / df.shape[0]) * 100).round()
```

**Output:**

```
PassengerId      0.0
Survived         0.0
```

```
Pclass          0.0
Name            0.0
Sex             0.0
Age            20.0
SibSp           0.0
Parch           0.0
Ticket          0.0
Fare            0.0
Cabin          77.0
Embarked        0.0
dtype: float64
```

**Handling Missing Values:**

- **Cabin**: 77% missing, so the column was dropped.

```
1  df = df.drop(columns="Cabin", axis=1)
```

- **Age**: 20% missing. Recommended to fill with median age.

- **Embarked**: 0.2% missing. Recommended to fill with mode.

**Recommended Code:**

```
1  # Fill Age with median
2  df['Age'] = df['Age'].fillna(df['Age'].median())
3  # Fill Embarked with mode
4  df['Embarked'] = df['Embarked'].fillna(df['Embarked'].mode()[0])
```

## 1.4 Summary Statistics

Summary statistics were generated using `describe()`.

```
1  df.describe()
```

**Output:**

|  | PassengerId | Survived | Pclass | Age | SibSp | Parch | |
|---|---|---|---|---|---|---|---|
| count | 891.0 | 891.0 | 891.0 | 714.0 | 891.0 | 891.0 | |
| mean | 446.0 | 0.383838 | 2.308642 | 29.699118 | 0.523008 | 0.381594 | 32.20 |
| std | 257.353842 | 0.486592 | 0.836071 | 14.526497 | 1.102743 | 0.806057 | 49.69 |
| min | 1.0 | 0.0 | 1.0 | 0.42 | 0.0 | 0.0 | |
| 25% | 223.5 | 0.0 | 2.0 | 20.125 | 0.0 | 0.0 | 7 |
| 50% | 446.0 | 0.0 | 3.0 | 28.0 | 0.0 | 0.0 | 14 |
| 75% | 668.5 | 1.0 | 3.0 | 38.0 | 1.0 | 0.0 | |
| max | 891.0 | 1.0 | 3.0 | 80.0 | 8.0 | 6.0 | 512 |

**Insights:**

- Average survival rate: 38.4%.

- Average age: 29.7 years.

- Fares vary widely (0 to 512.33), indicating potential outliers.

# 2 Part 2: Data Transformation

## 2.1 Converting Categorical Columns

Recommended approach for encoding Sex and Embarked:

```python
from sklearn.preprocessing import LabelEncoder
# Label Encoding for Sex
le = LabelEncoder()
df['Sex'] = le.fit_transform(df['Sex'])  # male: 1, female: 0
# One-Hot Encoding for Embarked
df = pd.get_dummies(df, columns=['Embarked'], prefix='Embarked')
```

## 2.2 Normalizing/Standardizing Fare and Age

Recommended code for scaling Age and Fare:

```python
from sklearn.preprocessing import MinMaxScaler
scaler = MinMaxScaler()
df[['Age', 'Fare']] = scaler.fit_transform(df[['Age', 'Fare']])
```

## 2.3 Creating New Column (FamilySize)

Recommended code to create FamilySize:

```python
df['FamilySize'] = df['SibSp'] + df['Parch']
```

**Note:** The notebook lacks these transformations, which are critical for analysis.

# 3 Part 3: Exploratory Data Analysis (EDA)

## 3.1 Survival Rate by Gender

```python
survival_by_gender = df.groupby('Sex')['Survived'].mean()
print(survival_by_gender)
```

**Output:**

```
Sex
female    0.747573
male      0.188908
Name: Survived, dtype: float64
```

**Insight:**

- Females had a significantly higher survival rate (74.8%) than males (18.9%).

## 3.2  Survival Rate by Passenger Class

```
survival_by_class = df.groupby('Pclass')['Survived'].mean()
print(survival_by_class)
```

**Output:**

```
Pclass
1    0.629630
2    0.472826
3    0.242363
Name: Survived, dtype: float64
```

**Insight:**

- First-class passengers had the highest survival rate (62.9%), followed by second-class (47.3%) and third-class (24.2%).

## 3.3  Survival Rate by Age Groups

```
df['AgeGroup'] = pd.cut(df['Age'], bins=[0, 12, 18, 30, 50, 80],
    labels=['Child', 'Teen', 'Young Adult', 'Adult', 'Senior'])
survival_by_age = df.groupby('AgeGroup')['Survived'].mean()
print(survival_by_age)
```

**Output:**

```
AgeGroup
Child         0.590909
Teen          0.428571
Young Adult   0.365854
Adult         0.391304
Senior        0.200000
Name: Survived, dtype: float64
```

**Insight:**

- Children had the highest survival rate (59.1%), while seniors had the lowest (20%).

## 3.4  Patterns and Anomalies

- **Pattern**: Higher-class passengers and females had better survival odds, likely due to prioritization during evacuation.

- **Anomaly**: Wide range in fares (0 to 512.33) suggests potential outliers in first-class tickets.

- **Insight**: Family size (SibSp + Parch) may influence survival, as larger families might have faced challenges during evacuation.

# 4 Part 4: Data Visualization

## 4.1 Bar Chart: Survival by Passenger Class

```python
import seaborn as sns
import matplotlib.pyplot as plt
sns.barplot(x='Pclass', y='Survived', data=df, palette='pastel')
plt.title('Survival Rate by Passenger Class')
plt.xlabel('Passenger Class')
plt.ylabel('Survival Rate')
plt.show()
```

**Insight:**

- Confirms first-class passengers had the highest survival rate.

## 4.2 Histogram: Age Distribution

```python
plt.hist(df['Age'].dropna(), bins=20, color='skyblue', edgecolor='
    black')
plt.title('Age Distribution of Passengers')
plt.xlabel('Age')
plt.ylabel('Frequency')
plt.show()
```

**Insight:**

- Age distribution is slightly right-skewed, with most passengers being young adults (20–40 years).

## 4.3 Boxplot: Fare by Passenger Class

```python
plt.figure(figsize=(12, 6))
sns.boxplot(data=df, x="Pclass", y="Fare", palette="pastel")
plt.title("Fare Distribution by Passenger Class")
plt.xlabel("Passenger Class")
plt.ylabel("Fare")
plt.show()
```

**Insight:**

- First-class fares show significant variability and outliers, with some paying extremely high amounts (e.g., >500).

## 4.4 Heatmap: Correlation Matrix

```python
plt.figure(figsize=(10, 8))
sns.heatmap(df.corr(), annot=True, cmap='coolwarm', center=0)
plt.title('Correlation Matrix')
plt.show()
```

**Insight:**

- Strong negative correlation between Pclass and Fare (-0.55).
- Moderate correlation between Survived and Sex (after encoding).

# 5 Key Findings

1. **Survival Disparities:**
   - Females had a higher survival rate (74.8%) than males (18.9%).
   - First-class passengers had a higher survival rate (62.9%) than third-class (24.2%).
   - Children (0–12 years) had the highest survival rate (59.1%).

2. **Data Quality:**
   - Cabin dropped due to 77% missing values.
   - Age (20% missing) should be imputed with median.
   - Embarked (0.2% missing) can be filled with mode.

3. **Feature Relationships:**
   - Higher fares associated with first-class tickets, with outliers.
   - Family size may impact survival due to logistical challenges.

4. **Recommendations:**
   - Investigate FamilySize impact using statistical tests.
   - Explore Sex, Pclass, and Age interactions with machine learning.
   - Complete missing transformations (encoding, scaling).

# 6 Conclusion

This project cleaned and analyzed the Titanic dataset, revealing key factors influencing survival (gender, class, age). Visualizations highlighted disparities in survival rates and fare distributions. The notebook lacks complete transformations, which should be addressed. Future work could involve predictive modeling.

# Submission Notes

- Notebook (assignment_2.ipynb) included in a ZIP file with plots.
- Missing transformations (encoding, scaling, FamilySize) should be implemented.

- This report addresses all tasks with code, visualizations, and findings.