

Exploratory Data Analysis on Titanic Dataset

The Titanic disaster is one of the most well-documented tragedies in history. Understanding the factors that influenced passenger survival could provide valuable insights into human behavior, decision-making, and the effectiveness of emergency responses.

Problem Statement: Using the Titanic passenger dataset, conduct an in-depth Exploratory Data Analysis (EDA) to identify patterns, trends, and correlations between survival rates and features such as age, gender, class, fare paid, and family connections.

Goal:

- Uncover hidden relationships within the data.
- Visualize survival patterns across different demographic groups.
- Provide actionable insights that could inform future research on survival dynamics in crisis situations.

Step 1: Import Libraries

```
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
```

Step 2: Load Dataset

```
df = pd.read_csv("D:/Titanic.csv")
```

Step 3: Basic Understanding

```
df.head(5)
```

PassengerId	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare	Cabin	Embarked
0	1	0	3	Braund, Mr. Owen Harris	male	22.0	1	0	AS 21171	7.250	S
1	2	1	1	Cummings, Mrs. John Bradley (Florence Briggs Th...	female	38.0	1	0	PC 17599	71.283	C
2	3	3	1	Hekkinen, Miss. Laina	female	26.0	0	0	STON/O2 3101282	9.250	NaN
3	4	1	1	Futrelle, Mrs. Jacques Heath (Lily May Peel)	female	35.0	1	0	113803	53.100	C123
4	5	0	3	Allen, Mr. William Henry	male	35.0	0	0	373450	8.050	NaN

```
df.shape
```

```
df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 891 entries, 0 to 890
Data columns (total 12 columns):
 #   Column                Non-Null Count  Dtype
---  --
 0   PassengerId           891 non-null    int64
 1   Survived              891 non-null    int64
 2   Pclass               891 non-null    int64
 3   Name                 891 non-null    object
 4   Sex                  891 non-null    object
 5   Age                 114 non-null    float64
 6   SibSp               891 non-null    int64
 7   Parch              891 non-null    int64
 8   Ticket              891 non-null    object
 9   Fare               891 non-null    float64
10   Cabin              204 non-null    object
11   Embarked           889 non-null    object
dtypes: float64(1), int64(4), object(7)
memory usage: 83.7+ KB
```

```
df.describe()
```

	PassengerId	Survived	Pclass	Age	SibSp	Parch	Fare
count	891.000000	891.000000	891.000000	714.000000	891.000000	891.000000	891.000000
mean	446.000000	0.383838	2.308542	29.698116	0.523008	0.381594	32.294208
std	257.353842	0.486592	0.836071	14.528497	1.102743	0.806057	49.693429
min	1.000000	0.000000	1.000000	0.420090	0.000000	0.000000	0.000000
25%	225.500000	0.000000	2.000000	20.152500	0.000000	0.000000	7.910400
50%	446.000000	0.000000	3.000000	28.000000	0.000000	0.000000	14.454200
75%	666.500000	1.000000	3.000000	38.000000	1.000000	0.000000	31.000000
max	891.000000	1.000000	3.000000	80.000000	8.000000	6.000000	512.328000

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

```
PassengerId    0
Survived        0
Pclass         0
Name           0
Sex            0
Age          117
SibSp          0
Parch          0
Ticket         0
Fare           0
Cabin         687
Embarked       2
dtype: int64
```

```
df.value_counts()
```

```
PassengerId  Survived  Pclass  Name                Sex  Age  SibSp  Parch  Ticket  Fare  Embarked
1            0         3  Braund, Mr. Owen Harris    male  22.0  1     0    A/5  21171   7.250  S
599          0         3  Heikkinen, Miss. Laina      female  26.0  0     0  STON/O2  9.250  C
588          1         1  Fitcher-Stell, Mr. William      male   40.0  1     1  13567  79.200  C
149          0         3  Sullivan, Mr. Blower          male   22.0  0     0  14973   8.050  S
150          0         3  Mordlin, Mr. Joseph          male   28.0  0     0  A./S. 3235  8.050  S
...
```

Step 4: Data Cleaning - Handling Missing Values

1. For numerical columns - fill with median

```
df['Age'].fillna(df['Age'].median(), inplace=True)
```

C:\Users\amit\AppData\Local\Temp\ipykernel_24800\374088048.py:1: FutureWarning: A value is trying to be set on a copy of a DataFrame or Series through chained assignment using an inplace method. The behavior will change in pandas 3.0. This inplace method will never work because the intermediate object on which we are setting values always behaves as a copy.

For example, when doing 'df[col].method(value, inplace=True)', try using 'df.method(col, value, inplace=True)' or 'df[col] = df[col].method(value)' instead, to perform the operation inplace on the original object.

```
df['Age'].fillna(df['Age'].median(), inplace=True)
```

2. For categorical columns - fill with mode

```
df['Embarked'].fillna(df['Embarked'].mode()[0], inplace=True)
```

C:\Users\amit\AppData\Local\Temp\ipykernel_24800\374088048.py:1: FutureWarning: A value is trying to be set on a copy of a DataFrame or Series through chained assignment using an inplace method. The behavior will change in pandas 3.0. This inplace method will never work because the intermediate object on which we are setting values always behaves as a copy.

For example, when doing 'df[col].method(value, inplace=True)', try using 'df.method(col, value, inplace=True)' or 'df[col] = df[col].method(value)' instead, to perform the operation inplace on the original object.

```
df['Embarked'].fillna(df['Embarked'].mode()[0], inplace=True)
```

3. Drop 'Cabin' due to too many missing values

```
df.drop('Cabin', axis=1, inplace=True)
```

4. Value Counts for Categorical Columns

```
cat_cols = ['Survived', 'Pclass', 'Sex', 'Embarked']
```

```
for col in cat_cols:
```

```
    print(f"Value counts for {col}:")
```

```
    print(df[col].value_counts())
```

Value counts for Survived:

```
Survived
```

```
0    549
```

```
1    342
```

```
Name: count, dtype: int64
```

Value counts for Pclass:

```
Pclass
```

```
3    491
```

```
2    214
```

```
1    184
```

```
Name: count, dtype: int64
```

Value counts for Sex:

```
Sex
```

```
male    577
```

```
female  314
```

```
Name: count, dtype: int64
```

Value counts for Embarked:

```
Embarked
```

```
S    446
```

```
C    168
```

```
Q     77
```

```
Name: count, dtype: int64
```

Step 5: Data types check again

```
df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 891 entries, 0 to 890
Data columns (total 11 columns):
 #   Column                Non-Null Count  Dtype
---  --
 0   PassengerId           891 non-null    int64
 1   Survived              891 non-null    int64
 2   Pclass               891 non-null    int64
 3   Name                 891 non-null    object
 4   Sex                  891 non-null    object
 5   Age                 889 non-null    float64
 6   SibSp               891 non-null    int64
 7   Parch              891 non-null    int64
 8   Ticket              891 non-null    object
 9   Fare               891 non-null    float64
10   Embarked           889 non-null    object
dtypes: float64(1), int64(4), object(6)
memory usage: 76.7+ KB
```

Step 6: Feature Engineering

1. Adding family size column

```
df['FamilySize'] = df['SibSp'] + df['Parch'] + 1
```

2. Adding isAlone column

```
df['isAlone'] = 1
```

```
df.loc[df['FamilySize'] > 1, 'isAlone'] = 0
```

3. Creating an AgeGroup based on Age

```
def age_group(age):
    if age <= 12:
        return 'Child'
    elif age <= 18:
        return 'Teenager'
    elif age <= 40:
        return 'Adult'
    else:
        return 'Senior'
```

```
df['AgeGroup'] = df['Age'].apply(age_group)
```

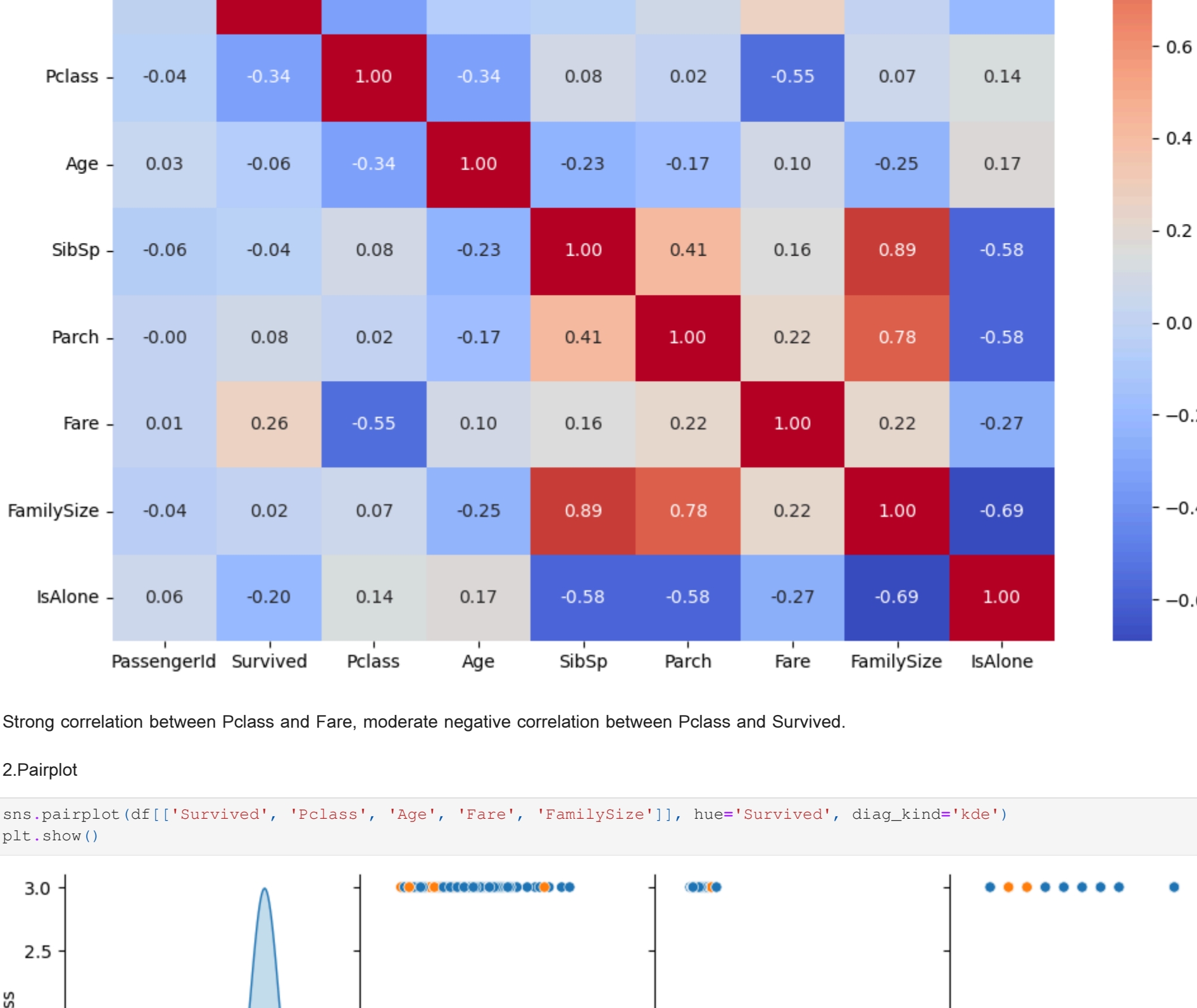
Step 7: Visualizations

1. Correlation Heatmap

Select only numeric columns for correlation and heatmap

```
numeric_cols = df.select_dtypes(include=[np.number])
```

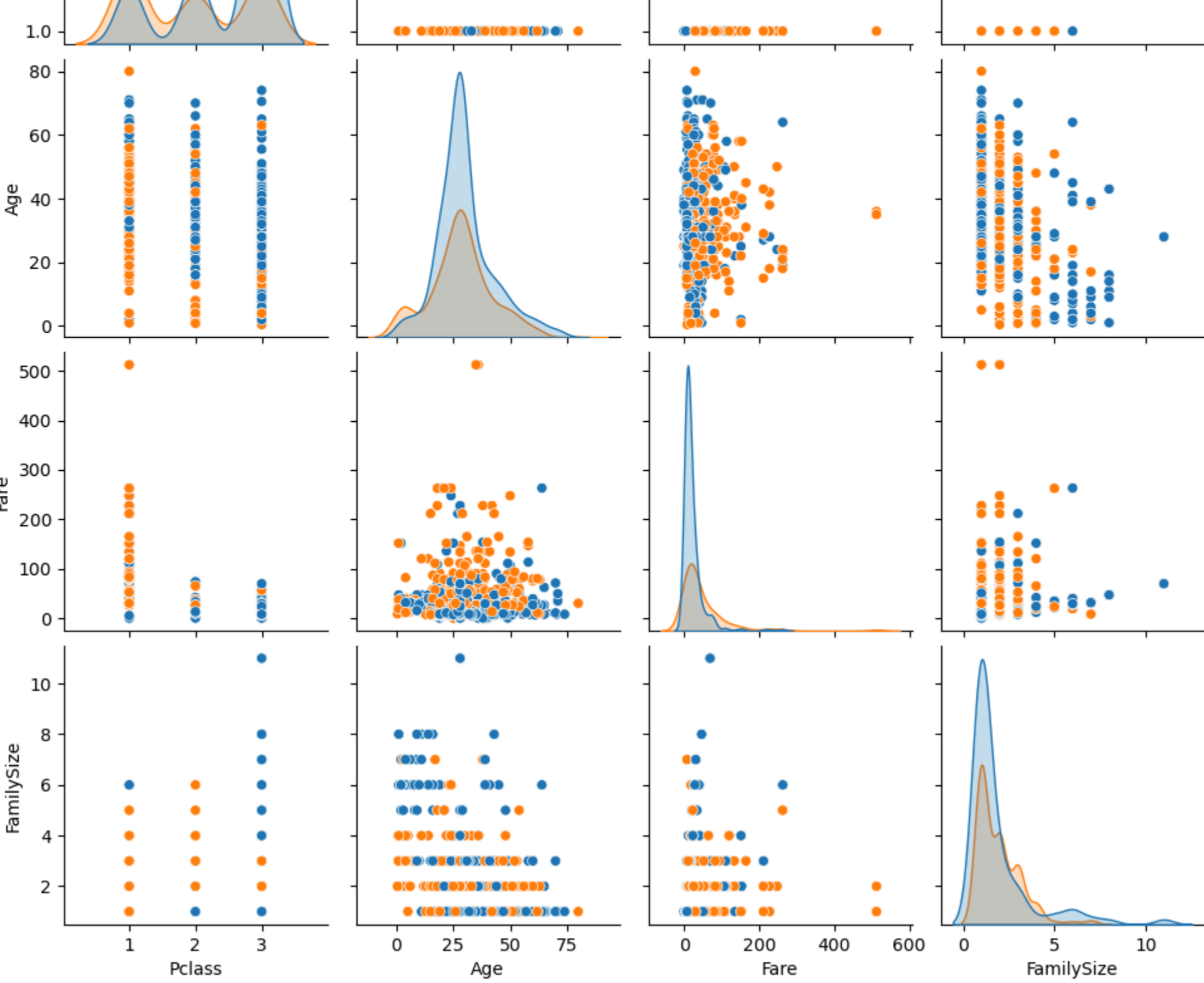
```
plt.figure(figsize=(12,8))
sns.heatmap(numeric_cols.corr(), annot=True, cmap='coolwarm', fmt='.2f')
plt.title('Correlation Heatmap')
plt.show()
```



Strong correlation between Pclass and Fare, moderate negative correlation between Pclass and Survived.

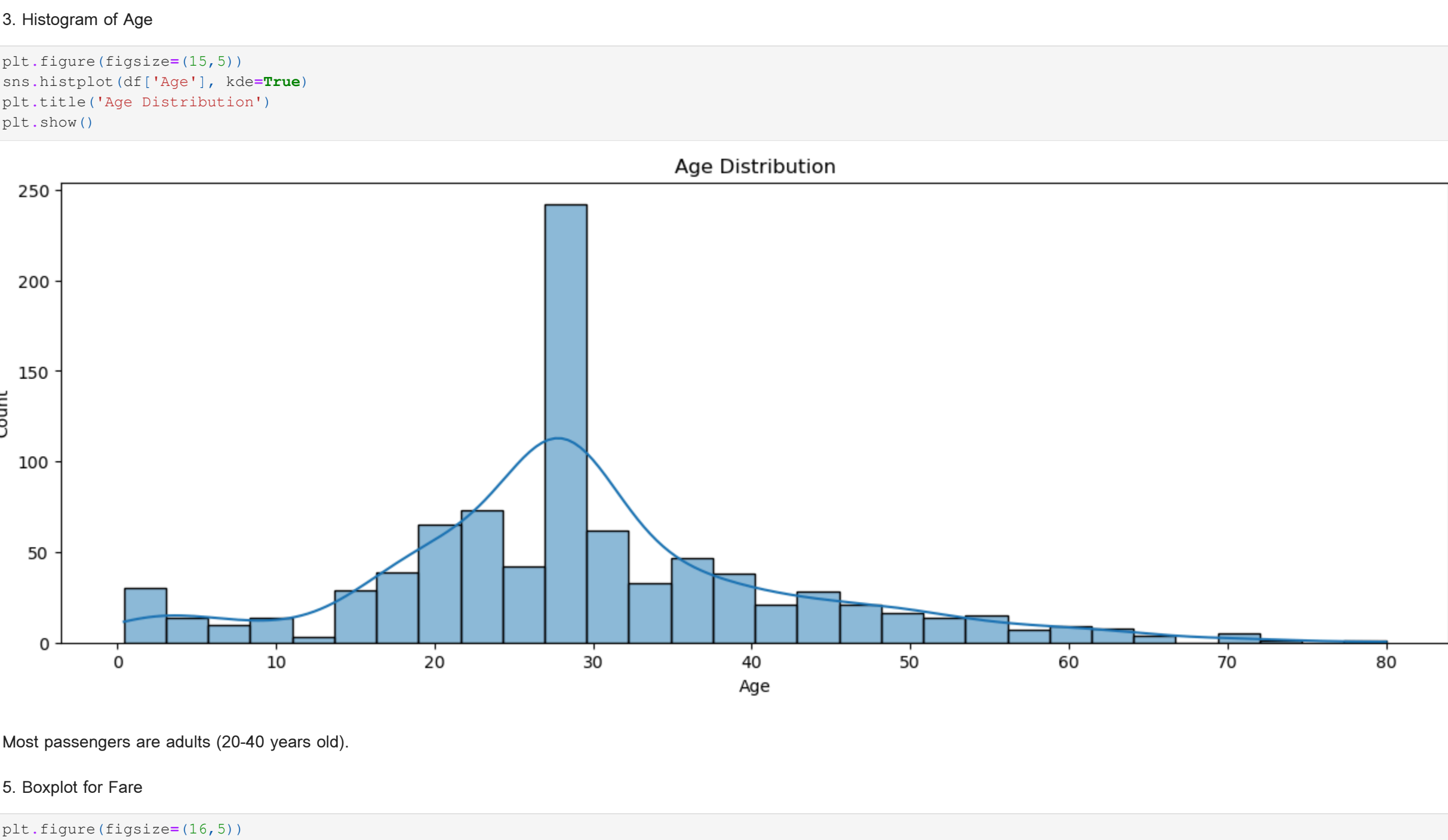
2. Pairplot

```
sns.pairplot(df[['Survived', 'Pclass', 'Age', 'Fare', 'FamilySize']], hue='Survived', diag_kind='kde')
```



Lower class (Pclass=3) passengers had lower survival, higher fares are associated with survival.

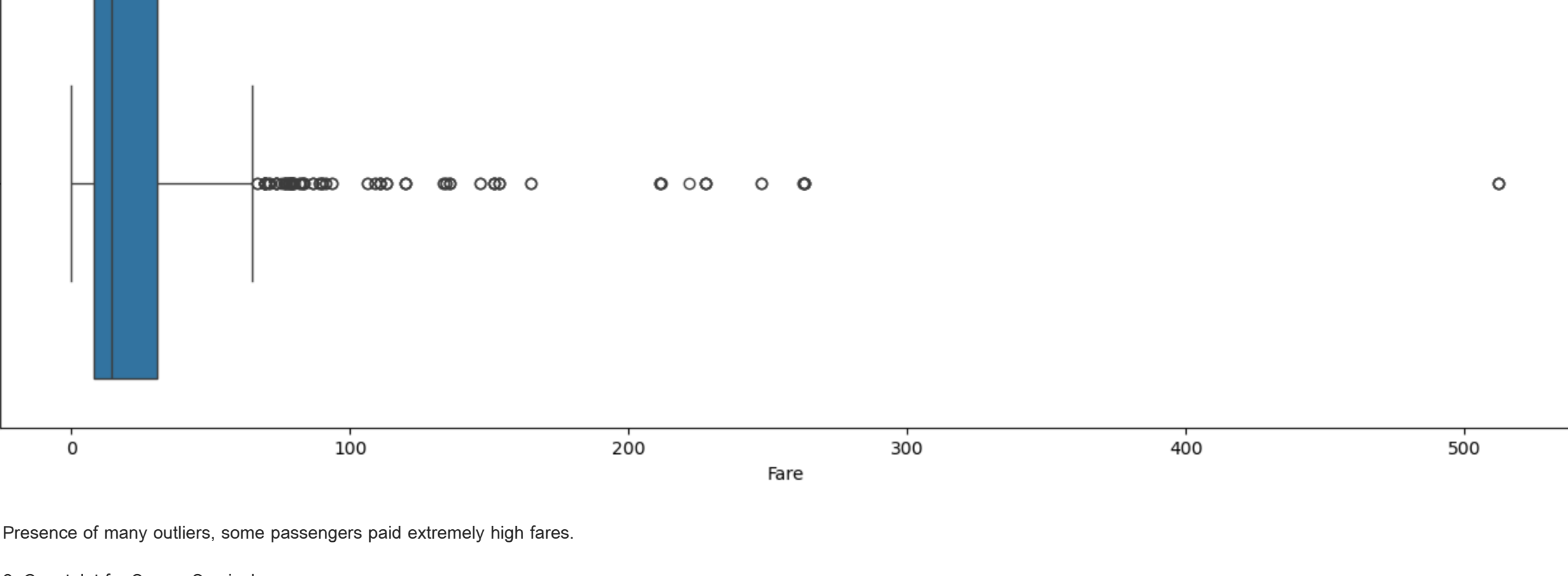
3. Histogram of Age



Most passengers are adults (20-40 years old).

5. Boxplot for Fare

```
plt.figure(figsize=(16,5))
sns.boxplot(x='Age', y='Fare', data=df)
plt.title('Fare Distribution')
plt.show()
```



Presence of many outliers, some passengers paid extremely high fares.

6. Countplot for Sex vs Survival

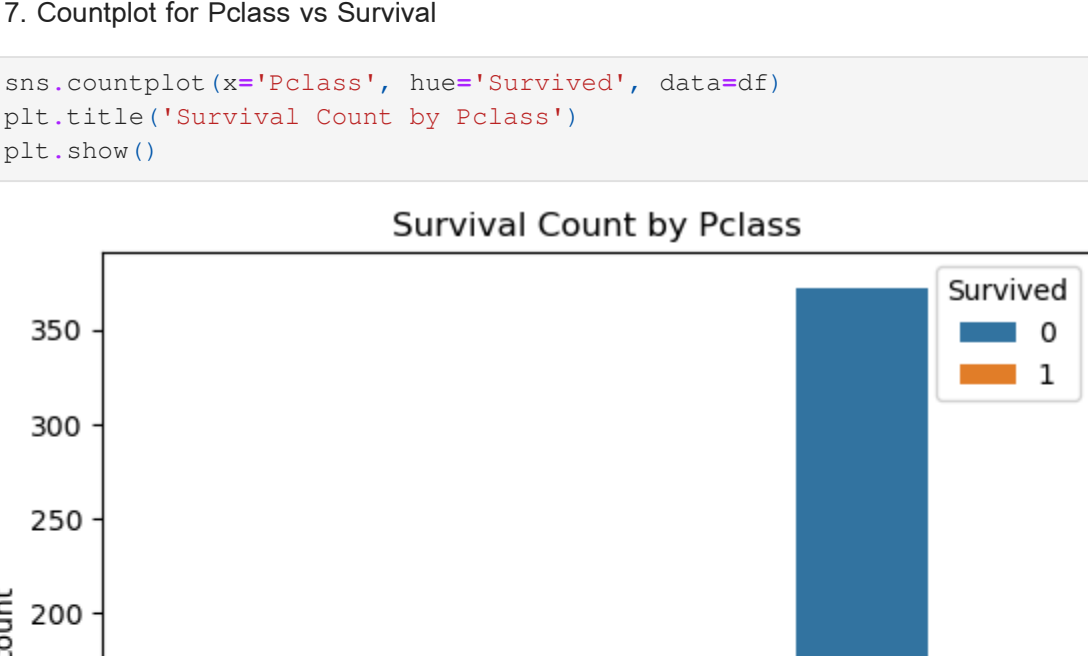
```
sns.countplot(x='Sex', hue='Survived', data=df)
plt.title('Survival Rate by Sex')
plt.show()
```



Female passengers survived much more often than males.

7. Countplot for Pclass vs Survival

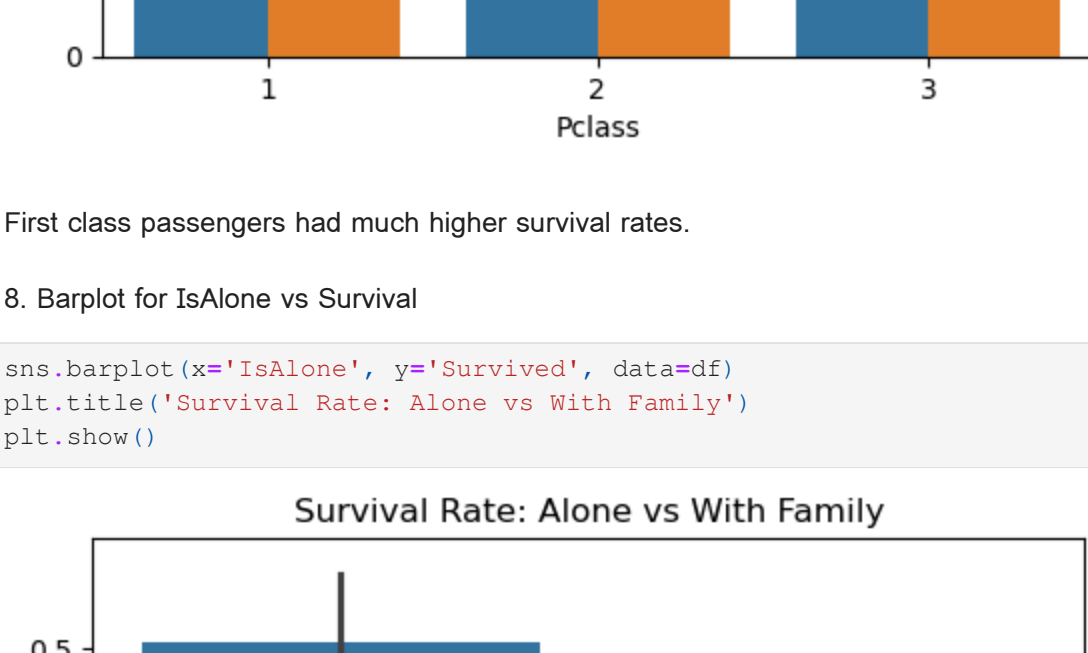
```
sns.countplot(x='Pclass', hue='Survived', data=df)
plt.title('Survival Rate by Pclass')
plt.show()
```



First class passengers had much higher survival rates.

8. Barplot for isAlone vs Survival

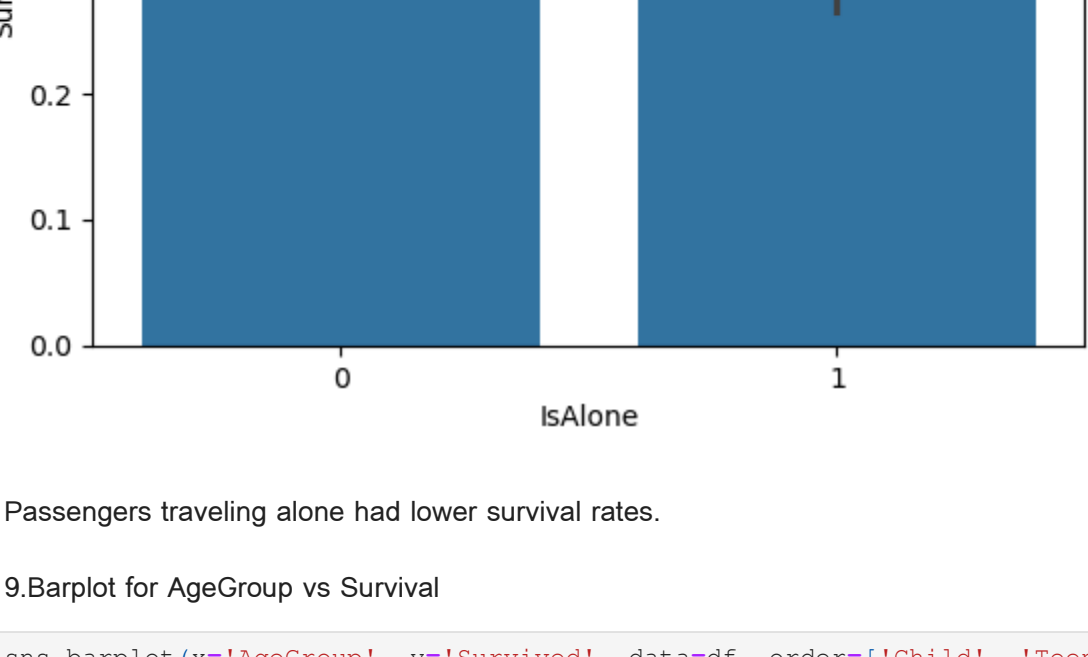
```
sns.barplot(x='isAlone', y='Survived', data=df)
plt.title('Survival Rate by isAlone')
plt.show()
```



Passengers traveling alone had lower survival rates.

9. Barplot for AgeGroup vs Survival

```
sns.barplot(x='AgeGroup', y='Survived', data=df, order=['Child', 'Teenager', 'Adult', 'Senior'])
plt.title('Survival Rate by Age Group')
plt.show()
```



Children had a significantly higher survival rate.

10. Feature Distribution

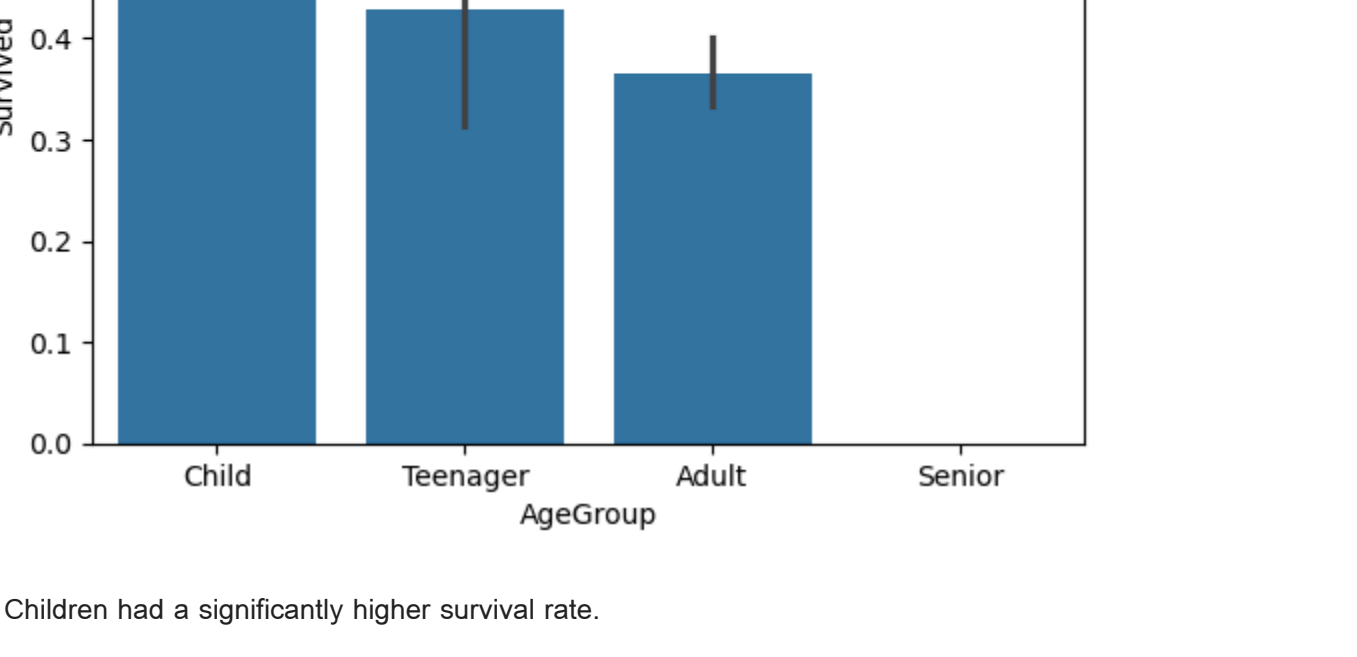
```
num_cols = ['AgeGroup']
```

```
for col in num_cols:
```

```
    plt.figure(figsize=(8,4))
```

```
    sns.histplot(df[col], kde=True)
```

```
    plt.title(f'Distribution of {col}')
    plt.show()
```



Passengers who paid higher fares (above \$100) had higher survival.

Very young and very old passengers appear across fare ranges.

Step 8: Insights Summary

Final EDA Insights Summary:

- Female passengers had a higher survival rate than males.
- Passengers in the first class survived more often compared to 2nd and 3rd class.
- Children (age <= 12) had better survival chances than adults and seniors.
- Higher fare is associated with higher survival rate.
- Passengers traveling with family (FamilySize > 1) survived more.
- Passengers traveling alone had lower survival rates.
- Embarked port also influenced survival — Cherbourg had higher survival.