

```
In [1]: import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns

sns.set(style="whitegrid")
```

```
In [2]: df = pd.read_csv("train.csv")
```

```
In [3]: df.head()
df.tail()
df.shape
df.info()
df.describe(include="all")
df.isnull().sum()
```

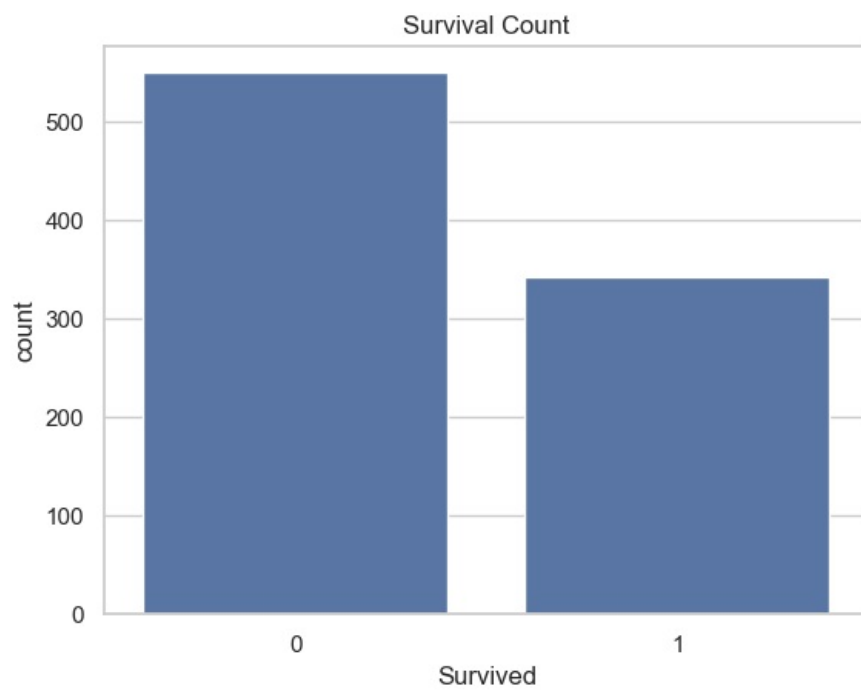
```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 891 entries, 0 to 890
Data columns (total 12 columns):
#   Column          Non-Null Count  Dtype
---  -
0   PassengerId      891 non-null   int64
1   Survived         891 non-null   int64
2   Pclass           891 non-null   int64
3   Name             891 non-null   object
4   Sex              891 non-null   object
5   Age              714 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(2), int64(5), object(5)
memory usage: 83.7+ KB
```

```
Out[3]: 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
```

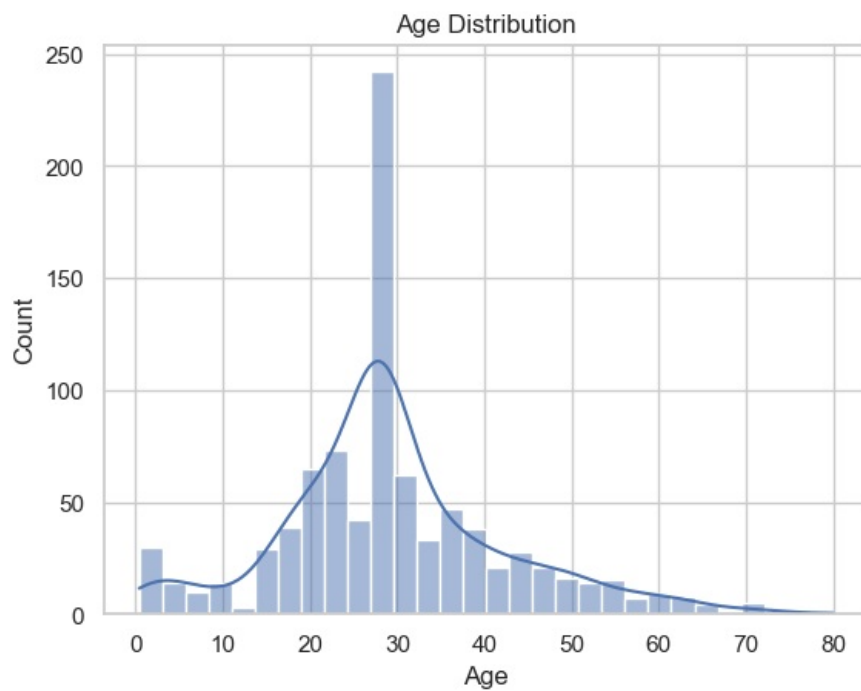
```
In [4]: df['Age'] = df['Age'].fillna(df['Age'].median())
df['Embarked'] = df['Embarked'].fillna(df['Embarked'].mode()[0])
df.drop(columns=['Cabin'], inplace=True)
df.isnull().sum()
```

```
Out[4]: PassengerId      0
Survived                0
Pclass                  0
Name                    0
Sex                     0
Age                     0
SibSp                   0
Parch                   0
Ticket                  0
Fare                     0
Embarked                0
dtype: int64
```

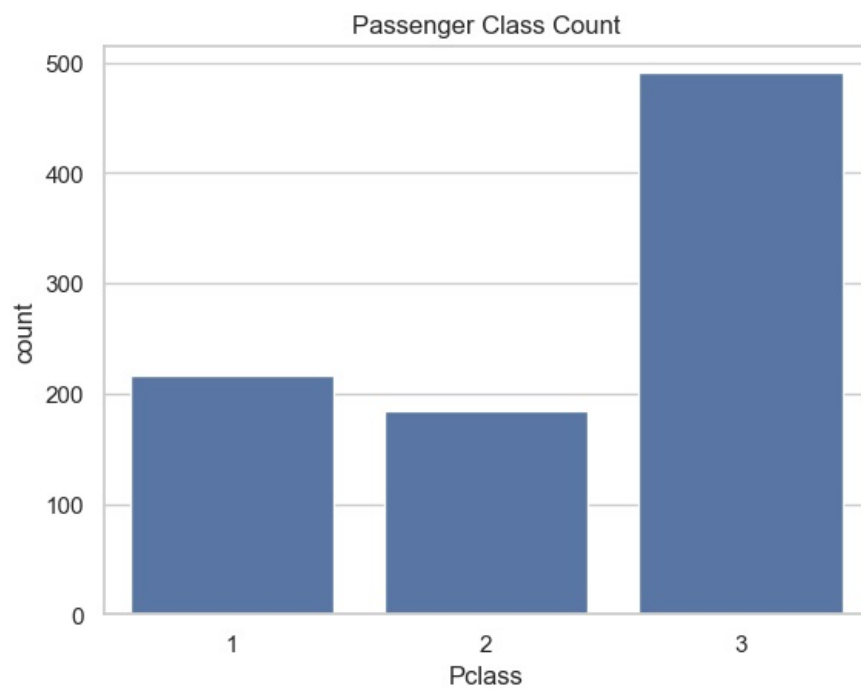
```
In [5]: sns.countplot(x='Survived', data=df)
plt.title("Survival Count")
plt.show()
```



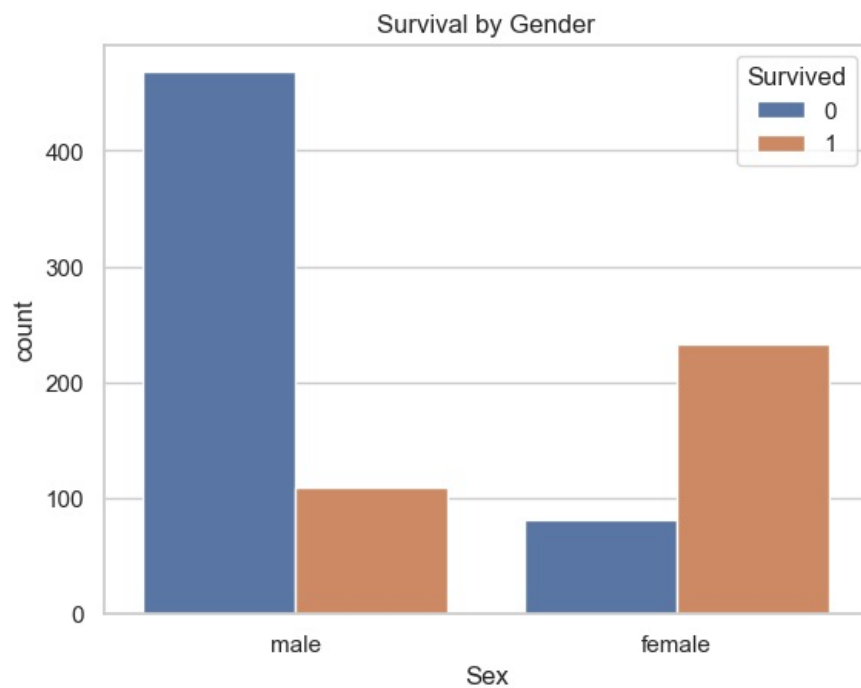
```
In [6]: sns.histplot(df['Age'], kde=True)
plt.title("Age Distribution")
plt.show()
```



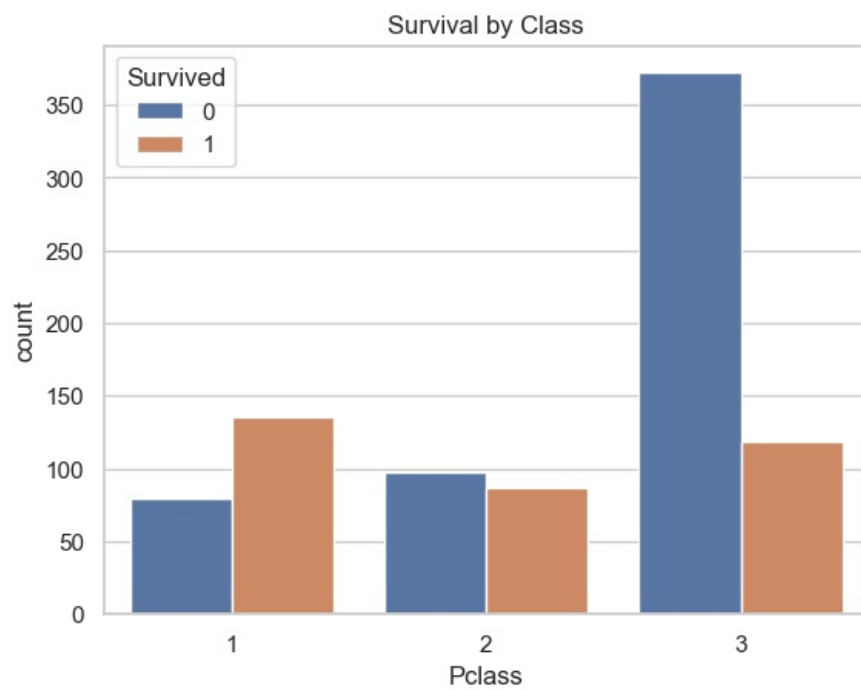
```
In [7]: sns.countplot(x='Pclass', data=df)
plt.title("Passenger Class Count")
plt.show()
```



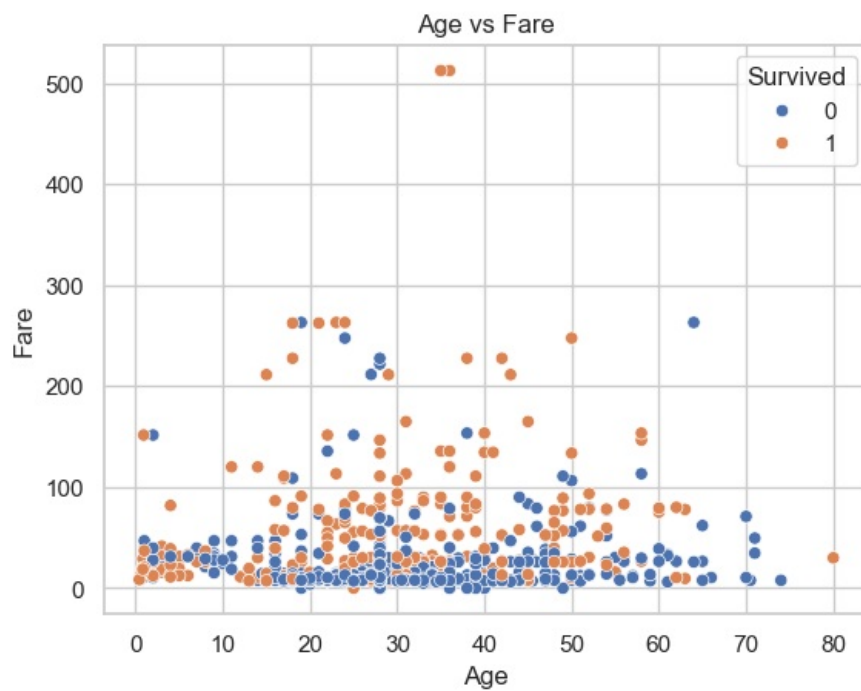
```
In [8]: sns.countplot(x='Sex', hue='Survived', data=df)
plt.title("Survival by Gender")
plt.show()
```



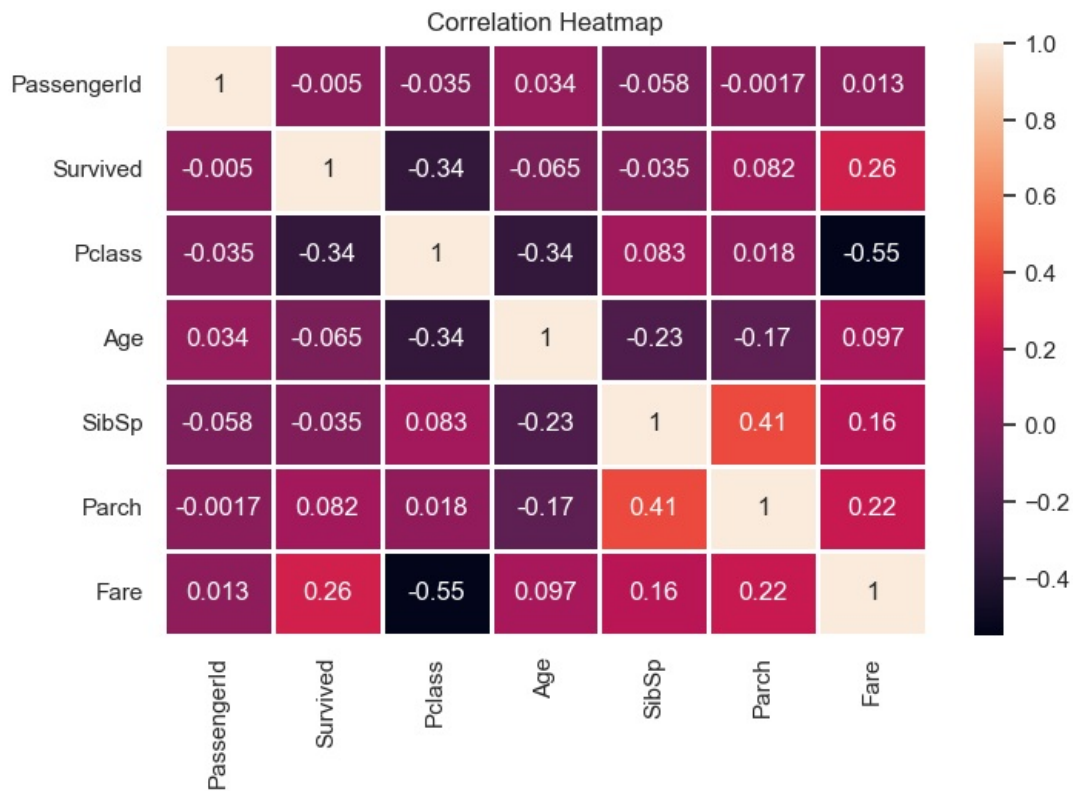
```
In [9]: sns.countplot(x='Pclass', hue='Survived', data=df)
plt.title("Survival by Class")
plt.show()
```



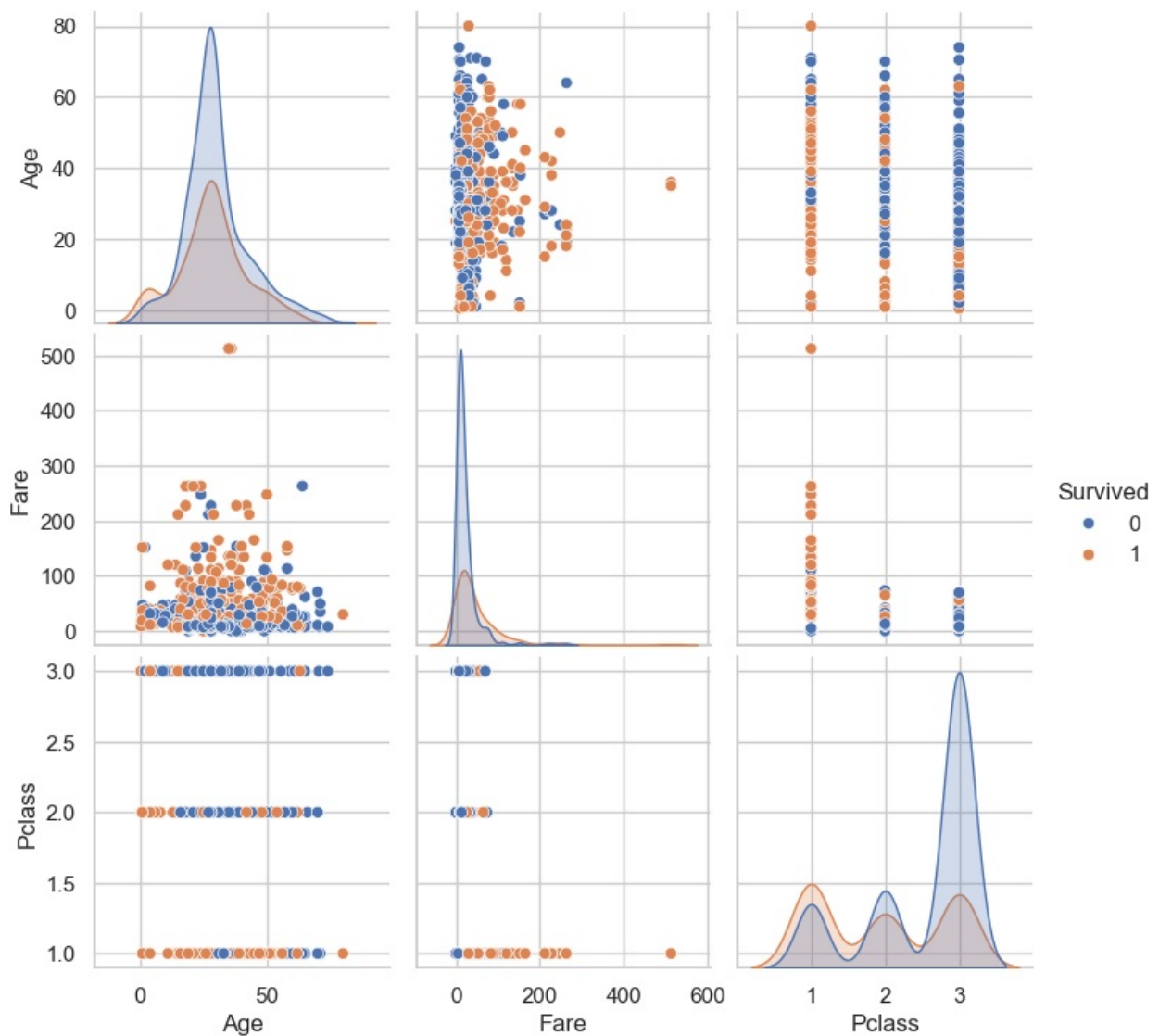
```
In [10]: sns.scatterplot(x='Age', y='Fare', hue='Survived', data=df)
plt.title("Age vs Fare")
plt.show()
```



```
In [11]: plt.figure(figsize=(8,5))
sns.heatmap(df.corr(numeric_only=True), annot=True, linewidths=1)
plt.title("Correlation Heatmap")
plt.show()
```



```
In [12]: sns.pairplot(df[['Survived', 'Age', 'Fare', 'Pclass']], hue='Survived')
plt.show()
```



```
In [13]: df.to_csv('train_analysis.csv', index=False)
```

```
In [ ]: ##Final Summary of EDA
```

1. Missing values Fixed

Missing Age filled using Median.

Missing Embarked filled using Mode.

2. Age Distribution

Most passengers belonged to the 3rd class. Age is mostly between 20-40.

3. Gender Distribution

There are more male passengers than female.

4. Survival Based on Gender

Women and first-class passengers had higher survival rates.

5. Survival Based on Class

1st Class has the Highest Survival rate.

3rd Class has the Lowest Survival rate.

6. Correlation Insights

Fare and Pclass negatively correlat

Fare has a positive Correlation with survival.

Pclass has a negative correlation with survival.

Overall Conclusion

Survival is strongly influenced by demographic and socio-economic factors.

Women, children, and wealthier passengers had significantly higher survival chances.

Therefore, the key determinants of survival were Gender, Passenger Class, and Fare.