

```
# Importing Libraries
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns

# Load the Dataset
df = pd.read_csv(r"C:\Users\Rohit\Documents\Rohit\Internship\Elevate Labs\Task 5\tr

# BASIC OVERVIEW

print("Shape:", df.shape)
print("\nColumn Types:\n", df.dtypes)
print("\nMissing Values:\n", df.isnull().sum())
print("\nUnique Values:\n", df.nunique())
print("\nDescriptive Stats:\n", df.describe(include='all'))
```

Shape: (891, 12)

Column Types:

```
PassengerId    int64
Survived        int64
Pclass          int64
Name            object
Sex             object
Age            float64
SibSp           int64
Parch           int64
Ticket          object
Fare            float64
Cabin           object
Embarked        object
dtype: object
```

Missing Values:

```
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
```

Unique Values:

```
PassengerId    891
Survived        2
Pclass          3
Name            891
Sex             2
Age            88
SibSp           7
Parch           7
Ticket          681
Fare            248
Cabin          147
Embarked        3
dtype: int64
```

Descriptive Stats:

	PassengerId	Survived	Pclass	Name	Sex	\
count	891.000000	891.000000	891.000000	891	891	
unique	NaN	NaN	NaN	891	2	
top	NaN	NaN	NaN	Braund, Mr. Owen Harris	male	
freq	NaN	NaN	NaN	1	577	
mean	446.000000	0.383838	2.308642	NaN	NaN	
std	257.353842	0.486592	0.836071	NaN	NaN	
min	1.000000	0.000000	1.000000	NaN	NaN	
25%	223.500000	0.000000	2.000000	NaN	NaN	
50%	446.000000	0.000000	3.000000	NaN	NaN	
75%	668.500000	1.000000	3.000000	NaN	NaN	
max	891.000000	1.000000	3.000000	NaN	NaN	

UNIVARIATE ANALYSIS

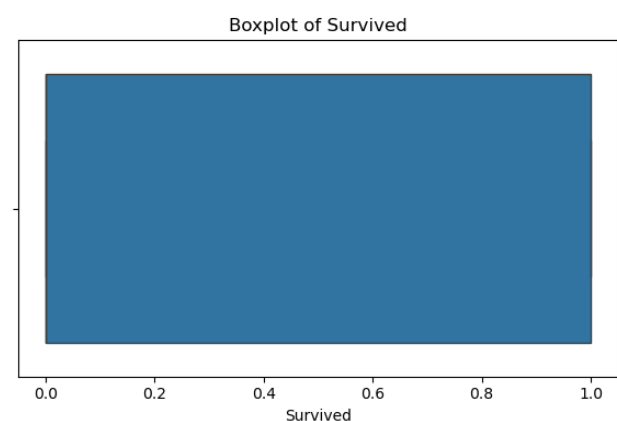
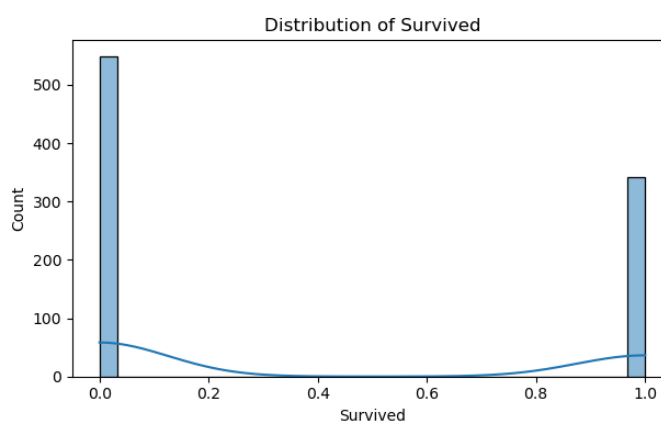
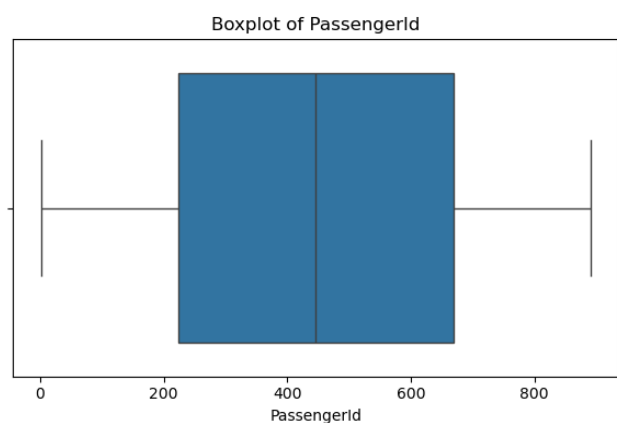
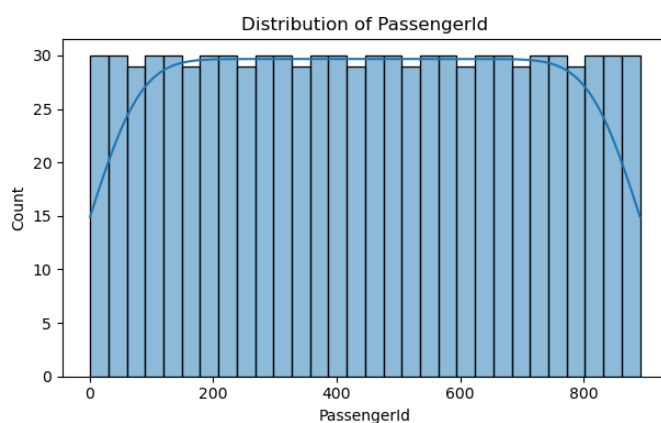
```
numeric_cols = df.select_dtypes(include=['int64', 'float64']).columns  
categorical_cols = df.select_dtypes(include='object').columns
```

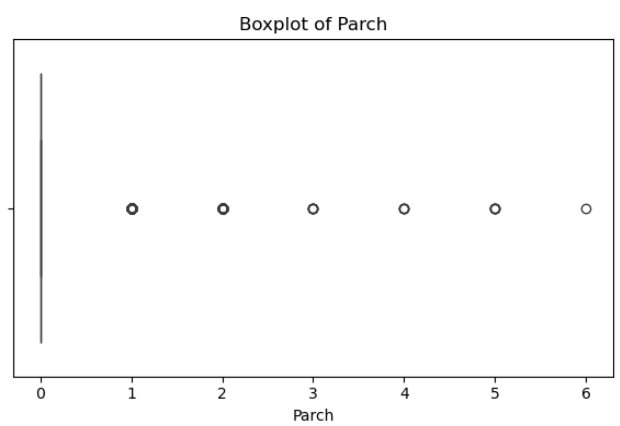
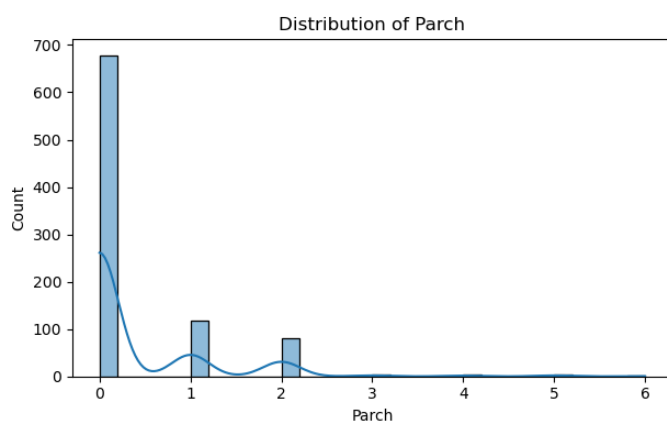
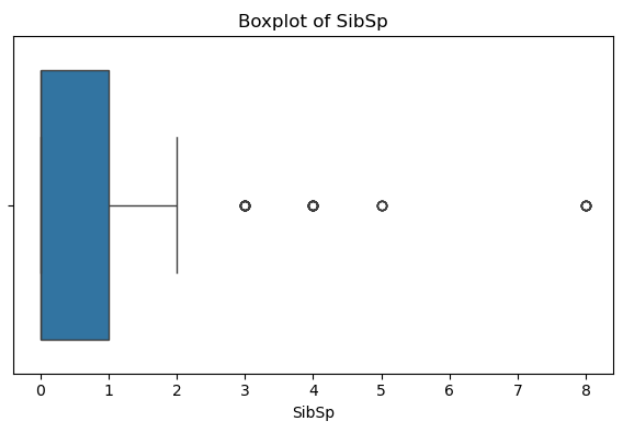
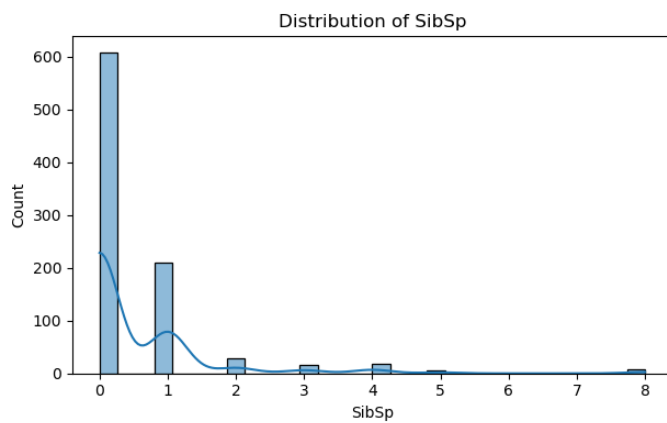
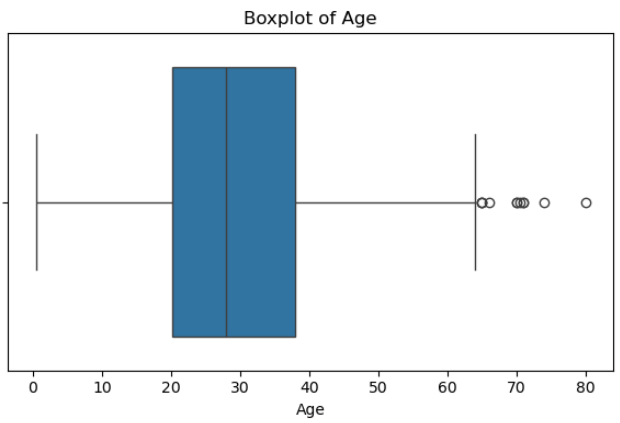
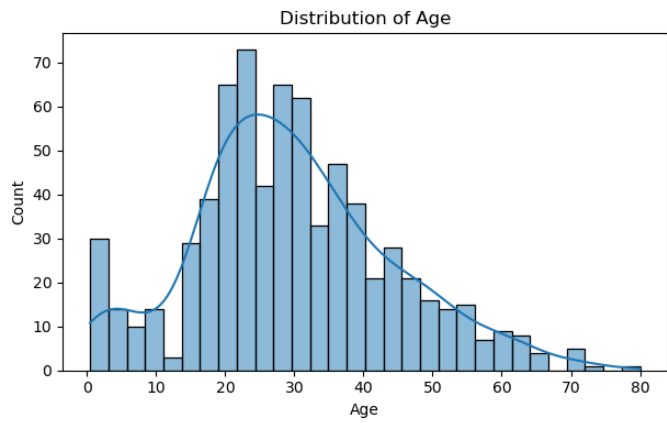
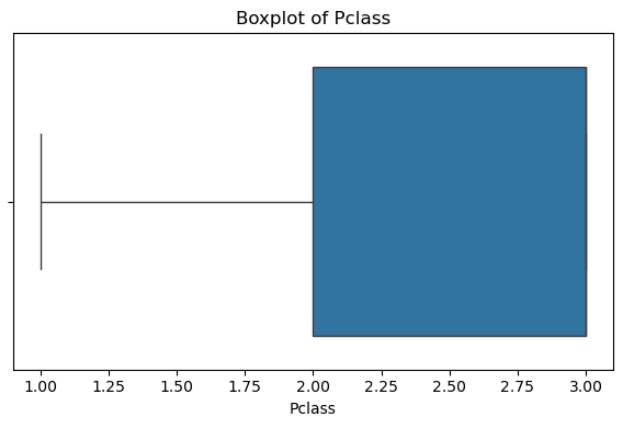
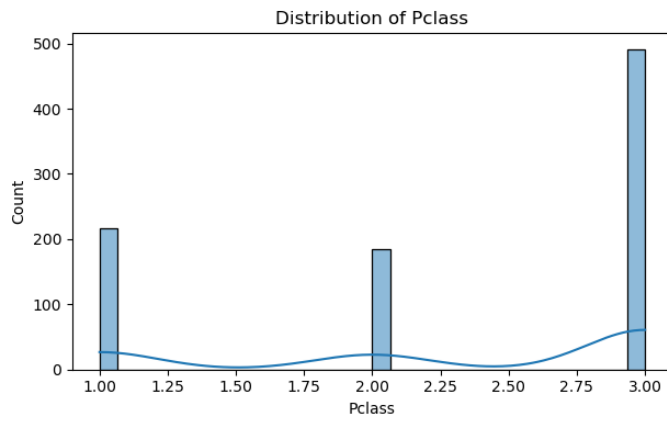
```
# Plot numeric columns
```

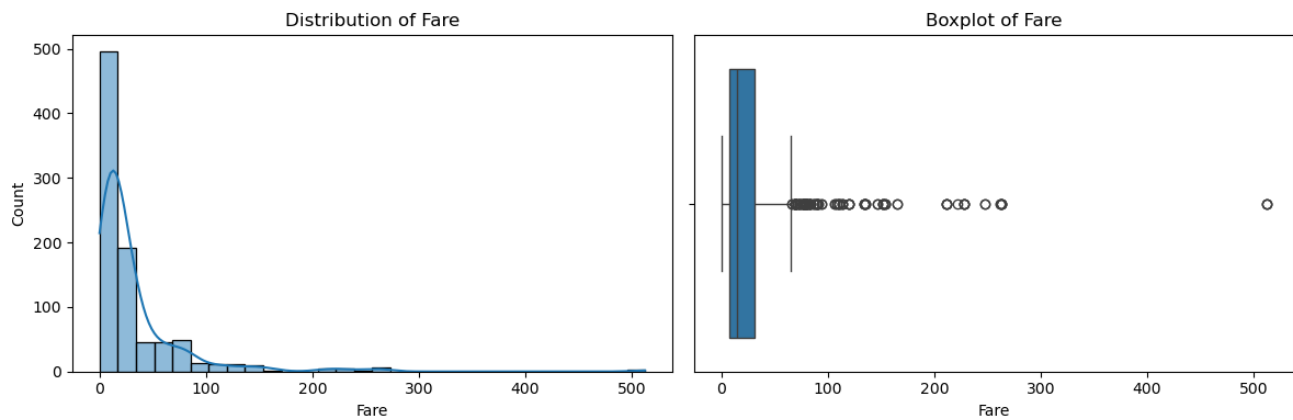
```
for col in numeric_cols:  
    plt.figure(figsize=(12, 4))  
    plt.subplot(1, 2, 1)  
    sns.histplot(df[col].dropna(), kde=True, bins=30)  
    plt.title(f"Distribution of {col}")  
  
    plt.subplot(1, 2, 2)  
    sns.boxplot(x=df[col])  
    plt.title(f"Boxplot of {col}")  
    plt.tight_layout()  
    plt.show()
```

```
# Plot categorical columns
```

```
for col in categorical_cols:  
    plt.figure(figsize=(8, 4))  
    df[col].value_counts().plot(kind='bar')  
    plt.title(f"Bar Plot of {col}")  
    plt.xticks(rotation=45)  
    plt.tight_layout()  
    plt.show()
```

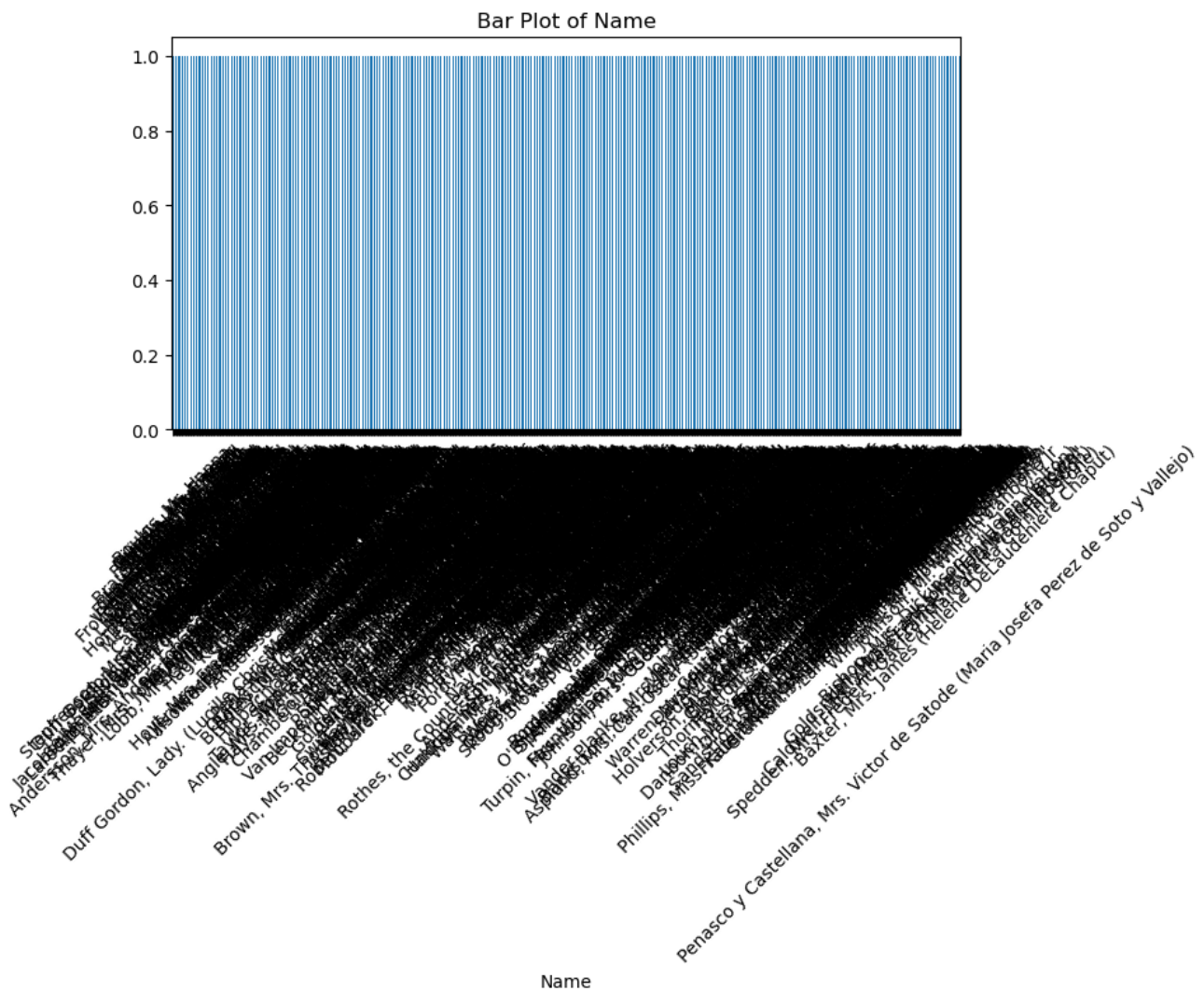







C:\Users\Rohit\AppData\Local\Temp\ipykernel_6244\2807690441.py:25: UserWarning:
Tight layout not applied. The bottom and top margins cannot be made large enough to
accommodate all Axes decorations.

```
plt.tight_layout()
```



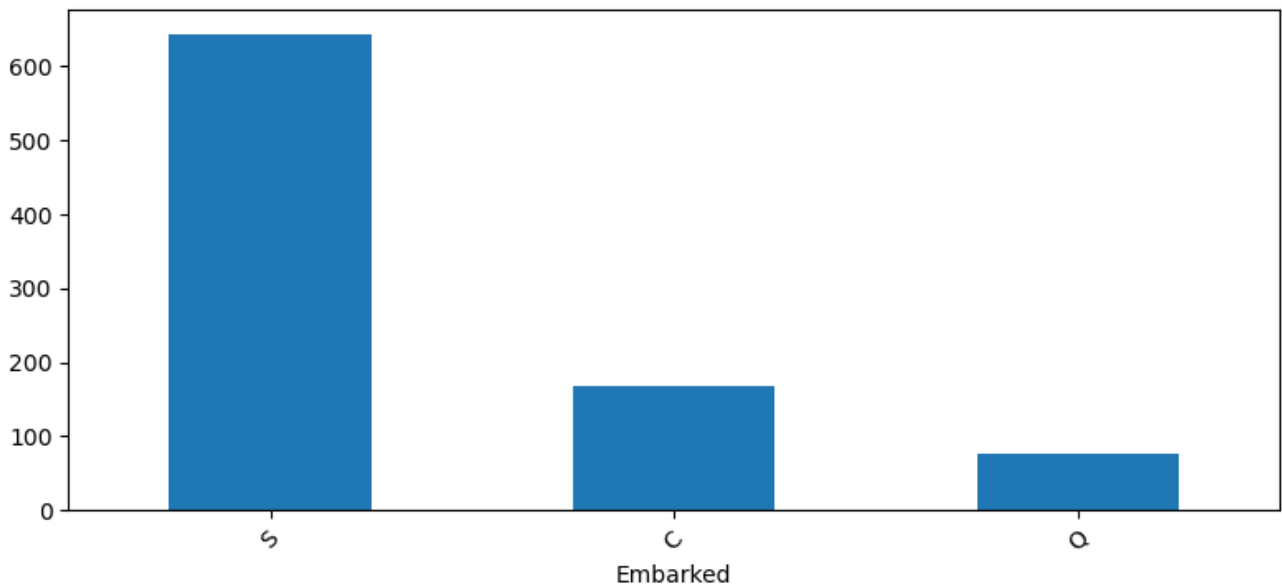


Sex	Frequency
male	580
female	310

Ticket	Frequency
S-00	4
S-01	3
S-02	3
S-03	2
S-04	2
S-05	2
S-06	2
S-07	2
S-08	1
S-09	1
S-10	1
S-11	1
S-12	1
S-13	1
S-14	1
S-15	1
S-16	1
S-17	1
S-18	1
S-19	1
S-20	1
S-21	1
S-22	1
S-23	1
S-24	1
S-25	1
S-26	1
S-27	1
S-28	1
S-29	1
S-30	1
S-31	1
S-32	1
S-33	1
S-34	1
S-35	1
S-36	1
S-37	1
S-38	1
S-39	1
S-40	1
S-41	1
S-42	1
S-43	1
S-44	1
S-45	1
S-46	1
S-47	1
S-48	1
S-49	1
S-50	1
S-51	1
S-52	1
S-53	1
S-54	1
S-55	1
S-56	1
S-57	1
S-58	1
S-59	1
S-60	1
S-61	1
S-62	1
S-63	1
S-64	1
S-65	1
S-66	1
S-67	1
S-68	1
S-69	1
S-70	1
S-71	1
S-72	1
S-73	1
S-74	1
S-75	1
S-76	1
S-77	1
S-78	1
S-79	1
S-80	1
S-81	1
S-82	1
S-83	1
S-84	1
S-85	1
S-86	1
S-87	1
S-88	1
S-89	1
S-90	1
S-91	1
S-92	1
S-93	1
S-94	1
S-95	1
S-96	1
S-97	1
S-98	1
S-99	1

A bar chart showing the frequency of cabin numbers. The y-axis represents frequency (0 to 4). The x-axis represents cabin numbers. The highest frequency is 4 for cabin B96. Other cabins with frequency 3 are C23, C22, and C11. Most other cabins have a frequency of 2 or 1.

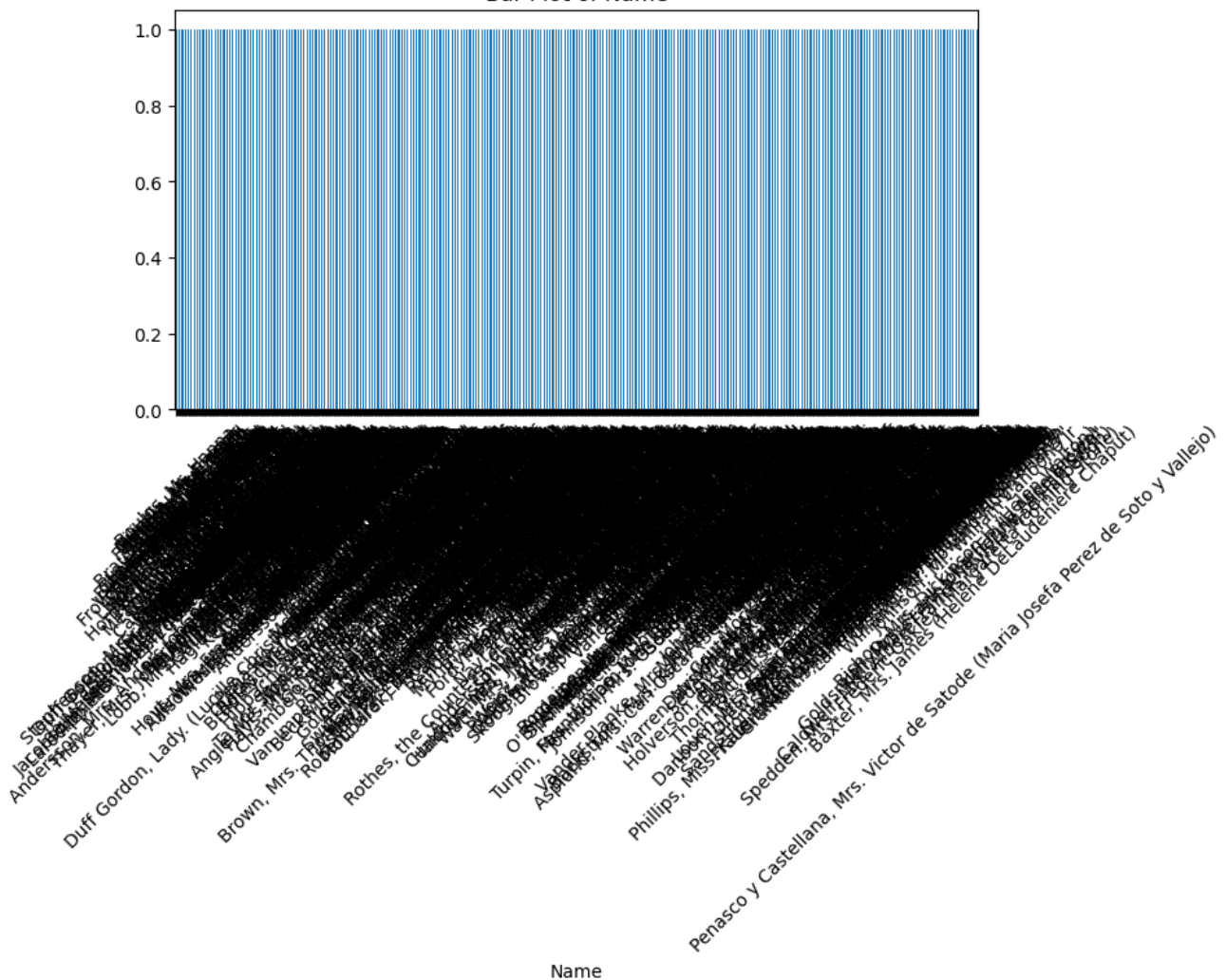
Bar Plot of Embarked




```
C:\Users\Rohit\AppData\Local\Temp\ipykernel_6244\1676552089.py:7: UserWarning:
Tight layout not applied. The bottom and top margins cannot be made large enough to
accommodate all Axes decorations.
```

```
plt.tight_layout()
```

Bar Plot of Name

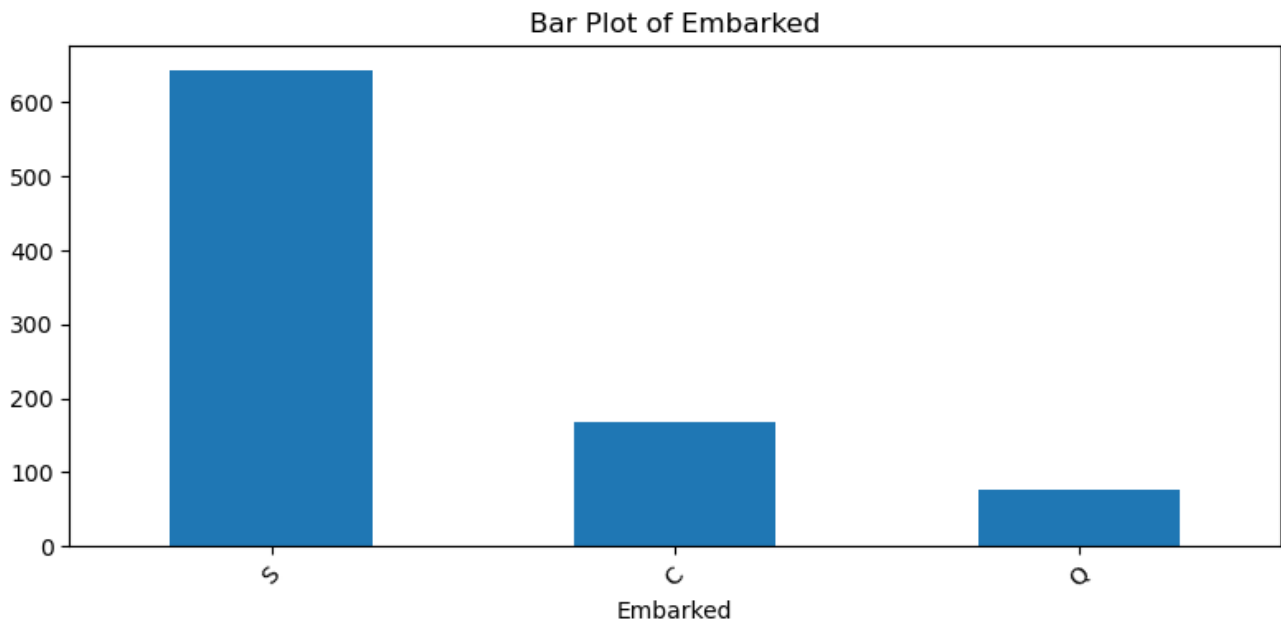




Sex	Count
male	580
female	310

Ticket	Frequency
S-00	4
S-01	3
S-02	3
S-03	2
S-04	2
S-05	2
S-06	2
S-07	2
S-08	1
S-09	1
S-10	1
S-11	1
S-12	1
S-13	1
S-14	1
S-15	1
S-16	1
S-17	1
S-18	1
S-19	1
S-20	1
S-21	1
S-22	1
S-23	1
S-24	1
S-25	1
S-26	1
S-27	1
S-28	1
S-29	1
S-30	1
S-31	1
S-32	1
S-33	1
S-34	1
S-35	1
S-36	1
S-37	1
S-38	1
S-39	1
S-40	1
S-41	1
S-42	1
S-43	1
S-44	1
S-45	1
S-46	1
S-47	1
S-48	1
S-49	1
S-50	1
S-51	1
S-52	1
S-53	1
S-54	1
S-55	1
S-56	1
S-57	1
S-58	1
S-59	1
S-60	1
S-61	1
S-62	1
S-63	1
S-64	1
S-65	1
S-66	1
S-67	1
S-68	1
S-69	1
S-70	1
S-71	1
S-72	1
S-73	1
S-74	1
S-75	1
S-76	1
S-77	1
S-78	1
S-79	1
S-80	1
S-81	1
S-82	1
S-83	1
S-84	1
S-85	1
S-86	1
S-87	1
S-88	1
S-89	1
S-90	1
S-91	1
S-92	1
S-93	1
S-94	1
S-95	1
S-96	1
S-97	1
S-98	1
S-99	1

A bar chart showing the frequency of cabin numbers. The y-axis represents frequency (0 to 4). The x-axis represents cabin numbers. The highest frequency is 4 for cabin B96. Other cabins with frequency 3 are C23, C22, and C11. Most other cabins have a frequency of 2 or 1.



```
# BIVARIATE / MULTIVARIATE ANALYSIS
# Correlation heatmap
plt.figure(figsize=(10, 6))
sns.heatmap(df[numeric_cols].corr(), annot=True, cmap='coolwarm')
plt.title("Correlation Matrix")
plt.show()

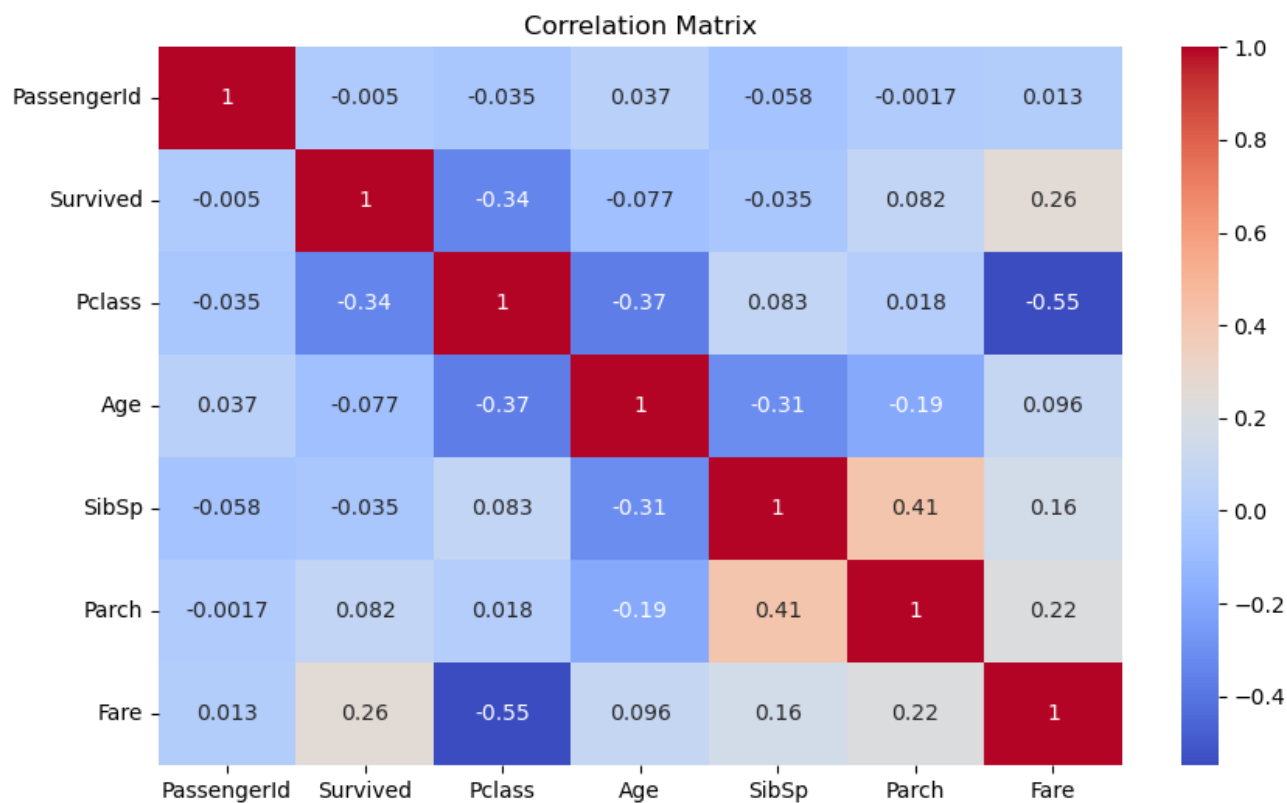
# Clean and prepare pairplot data
pairplot_data = df[['Age', 'Fare', 'Pclass', 'Survived']].copy()
for col in pairplot_data.columns:
    pairplot_data[col] = pd.to_numeric(pairplot_data[col], errors='coerce')
pairplot_data.dropna(inplace=True)

# Pairplot
sns.pairplot(pairplot_data, hue='Survived')
plt.suptitle("Pairplot of Age, Fare, Pclass, Survived", y=1.02)
plt.show()

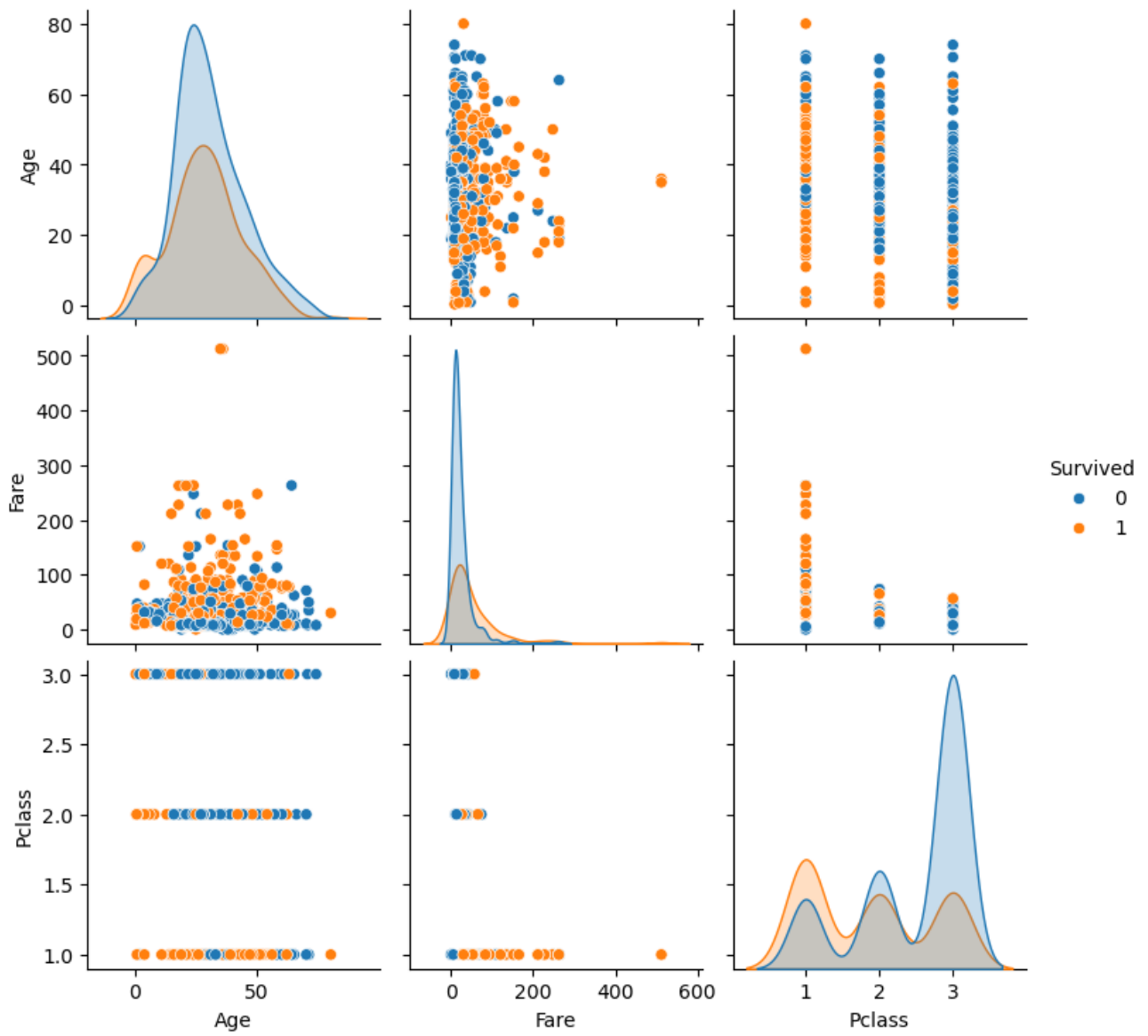
# Boxplot: Age by Survival
sns.boxplot(x='Survived', y='Age', data=df)
plt.title("Age Distribution by Survival")
plt.show()

# Boxplot: Fare by Pclass
sns.boxplot(x='Pclass', y='Fare', data=df)
plt.title("Fare Distribution by Pclass")
plt.show()

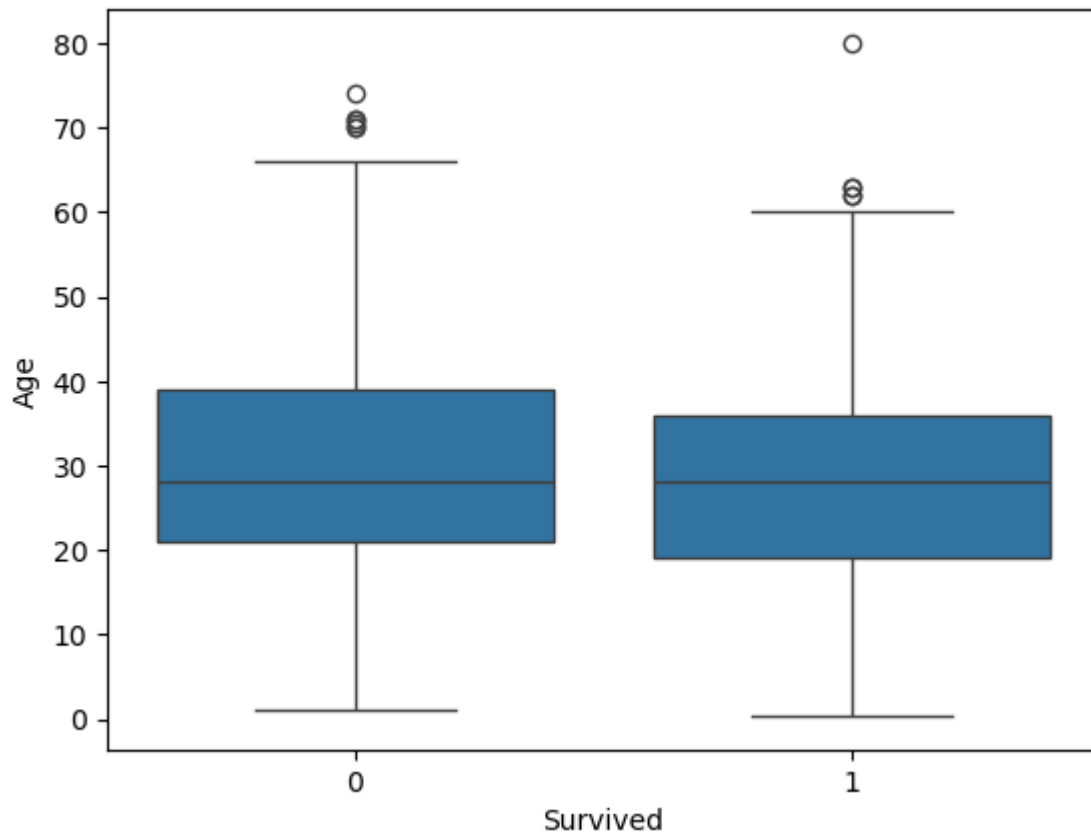
# Countplot: Survival by Sex
sns.countplot(x='Sex', hue='Survived', data=df)
plt.title("Survival Count by Sex")
plt.show()
```



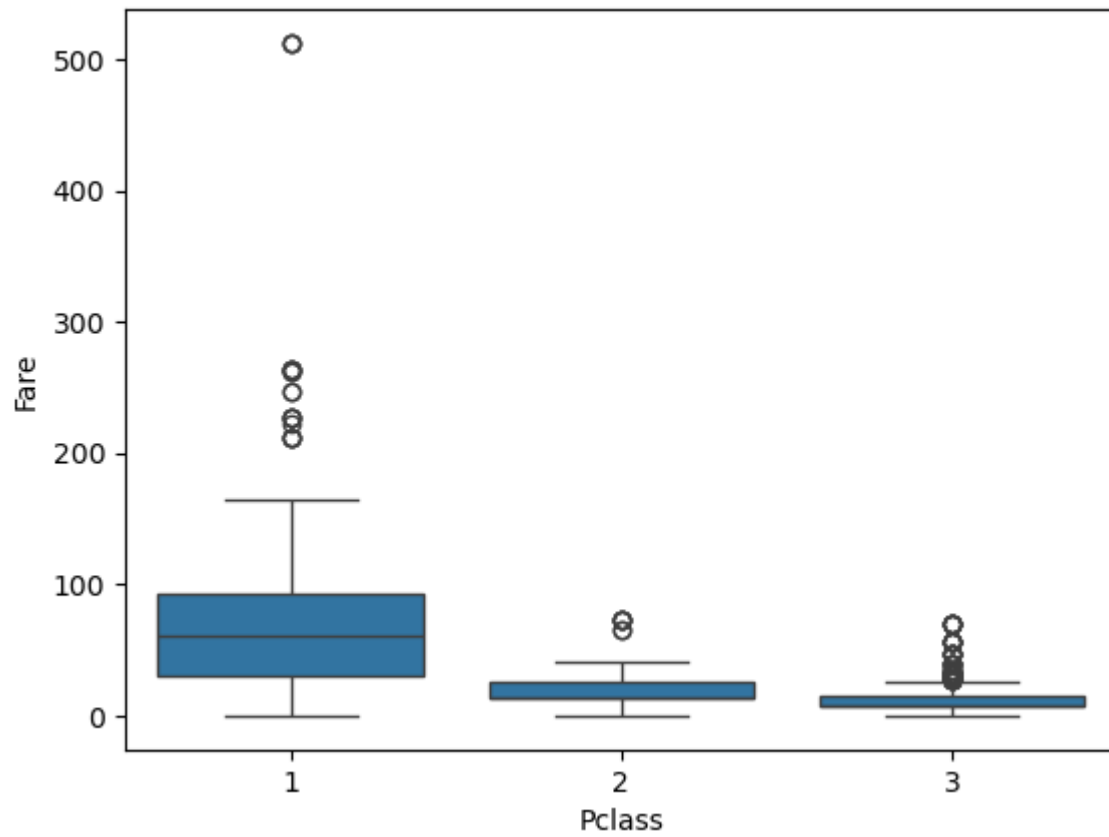
Pairplot of Age, Fare, Pclass, Survived

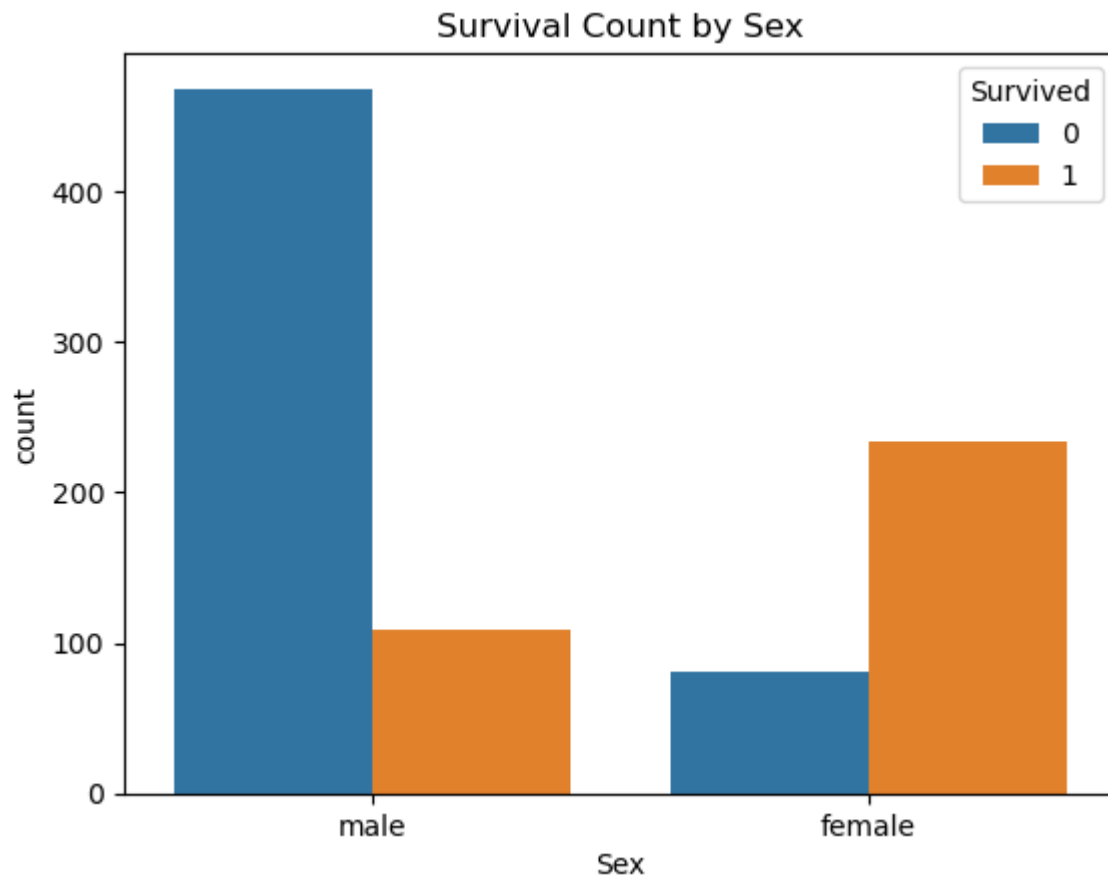


Age Distribution by Survival



Fare Distribution by Pclass





MISSING VALUES VISUALIZATION

```
plt.figure(figsize=(10, 5))
sns.heatmap(df.isnull(), cbar=False, cmap='viridis')
plt.title("Missing Values Heatmap")
plt.show()
```

