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

path='/content/drive/MyDrive/prodigy ds/ Titanic-Dataset.csv'
data=pd.read_csv(path)
```

```
data.head()
```

| | PassengerId | Survived | Pclass | Name | Sex | Age | SibSp | Parc |
|---|---|---|---|---|---|---|---|---|
| **0** | 1 | 0 | 3 | Braund, Mr. Owen Harris | male | 22.0 | 1 | |
| **1** | 2 | 1 | 1 | Cumings, Mrs. John Bradley (Florence Briggs Th... | female | 38.0 | 1 | |
| **2** | 3 | 1 | 3 | Heikkinen, Miss. Laina | female | 26.0 | 0 | |
| | | | | Futrelle, | | | | |

Next steps: 🔘 **View recommended plots**

```
data.describe()
```

| | PassengerId | Survived | Pclass | Age | SibSp | |
|---|---|---|---|---|---|---|
| **count** | 891.000000 | 891.000000 | 891.000000 | 714.000000 | 891.000000 | 89 |
| **mean** | 446.000000 | 0.383838 | 2.308642 | 29.699118 | 0.523008 | |
| **std** | 257.353842 | 0.486592 | 0.836071 | 14.526497 | 1.102743 | |
| **min** | 1.000000 | 0.000000 | 1.000000 | 0.420000 | 0.000000 | |
| **25%** | 223.500000 | 0.000000 | 2.000000 | 20.125000 | 0.000000 | |
| **50%** | 446.000000 | 0.000000 | 3.000000 | 28.000000 | 0.000000 | |
| **75%** | 668.500000 | 1.000000 | 3.000000 | 38.000000 | 1.000000 | |
| **max** | 891.000000 | 1.000000 | 3.000000 | 80.000000 | 8.000000 | |

```
data.info()
```

```
<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
```

data.isnull().sum()

```
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
```

data.dropna(subset=["Embarked"], inplace=True)
data["Cabin"].fillna("Unknown", inplace=True)
data["Age"].fillna(data["Age"].mean(), inplace=True)

data.isnull().sum()

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

data.duplicated().sum()

```
0
```

print(data.dtypes)

```
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_data = data.isnull().sum()

```
missing_data
```

```
        PassengerId     0
        Survived        0
        Pclass          0
        Name            0
        Sex             0
        Age             0
        SibSp           0
        Parch           0
        Ticket          0
        Fare            0
        Cabin           0
        Embarked        0
        dtype: int64
```
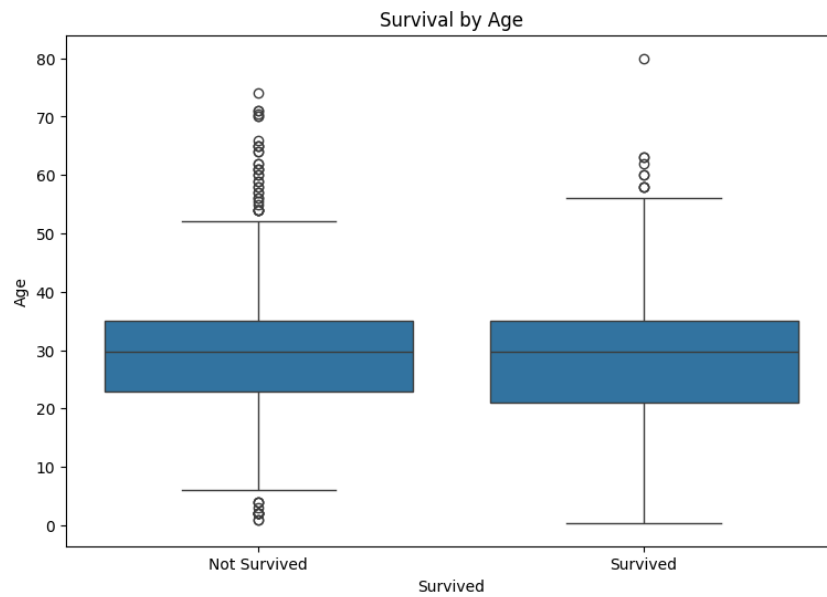
```
data = data.drop_duplicates()
data.head(8)
```

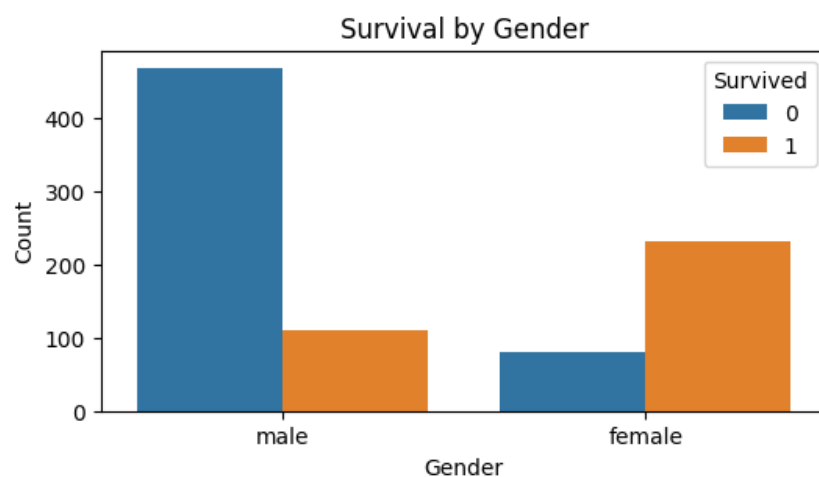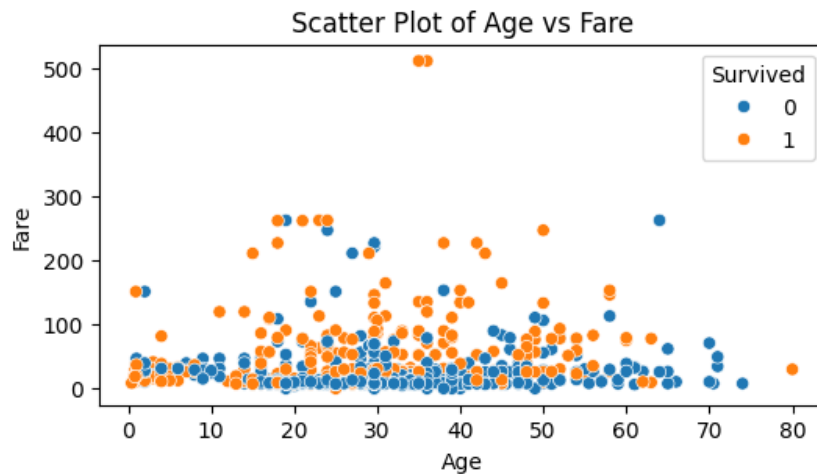| | PassengerId | Survived | Pclass | Name | Sex | Age | SibSp |
|---|---|---|---|---|---|---|---|
| **0** | 1 | 0 | 3 | Braund, Mr. Owen Harris | male | 22.000000 | 1 |
| **1** | 2 | 1 | 1 | Cumings, Mrs. John Bradley (Florence Briggs Th... | female | 38.000000 | 1 |
| **2** | 3 | 1 | 3 | Heikkinen, Miss. Laina | female | 26.000000 | 0 |
| **3** | 4 | 1 | 1 | Futrelle, Mrs. Jacques Heath (Lily May Peel) | female | 35.000000 | 1 |
| **4** | 5 | 0 | 3 | Allen, Mr. William Henry | male | 35.000000 | 0 |
| **5** | 6 | 0 | 3 | Moran, Mr. James | male | 29.642093 | 0 |

Next steps:　　🔘 **View recommended plots**

```
# Visualizing Survival by Age:
plt.figure(figsize=(9, 6))
sns.boxplot(x='Survived', y='Age', data=data)
plt.title('Survival by Age')
plt.xticks([0, 1], ['Not Survived', 'Survived'])
plt.show()
```
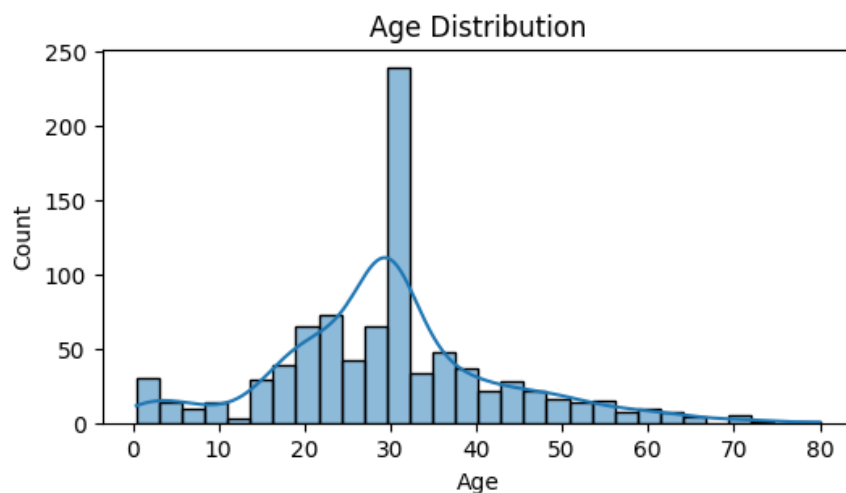
Survival by Age

```
plt.figure(figsize=(6, 3))
sns.countplot(data=data, x="Sex", hue="Survived")
plt.title("Survival by Gender")
plt.xlabel("Gender")
plt.ylabel("Count")
plt.legend(title="Survived", loc="upper right")
plt.show()
```



Survival by Gender

```python
plt.figure(figsize=(6, 3))
sns.scatterplot(data=data, x="Age", y="Fare", hue="Survived")
plt.title("Scatter Plot of Age vs Fare")
plt.xlabel("Age")
plt.ylabel("Fare")
plt.legend(title="Survived")
plt.show()
```

### Scatter Plot of Age vs Fare



```python
plt.figure(figsize=(6, 3))
sns.histplot(data["Age"], kde=True)
plt.title("Age Distribution")
plt.xlabel("Age")
plt.ylabel("Count")
plt.show()
```

### Age Distribution



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

sns.pairplot(data, vars=['Age', 'Fare', 'FamilySize'], hue='Survived')
plt.suptitle('Pair Plot of Age, Fare, and FamilySize')
plt.show()
```

Pair Plot of Age, Fare, and FamilySize