

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

```
df = pd.read_csv("/content/sample_data/Titanic-Dataset.csv")
print(df.head()) # first 5 rows
print(df.info()) # data types + null values
print(df.describe()) # statistical summary
print(df.isnull().sum()) # missing values count
```

```

Parch      Ticket      Fare Cabin Embarked
0      0      A/5 21171   7.2500   NaN        S
1      0      PC 17599  71.2833   C85        C
2      0  STON/O2. 3101282   7.9250   NaN        S
3      0      113803  53.1000  C123        S
4      0      373450   8.0500   NaN        S
<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
None

      PassengerId  Survived  Pclass    Age  SibSp  \
count  891.000000  891.000000  891.000000  714.000000  891.000000
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

      Parch      Fare
count  891.000000  891.000000
mean     0.381594   32.204208
std     0.806057   49.693429
min     0.000000    0.000000
25%     0.000000    7.910400
50%     0.000000   14.454200
75%     0.000000   31.000000
max     6.000000  512.329200
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
```

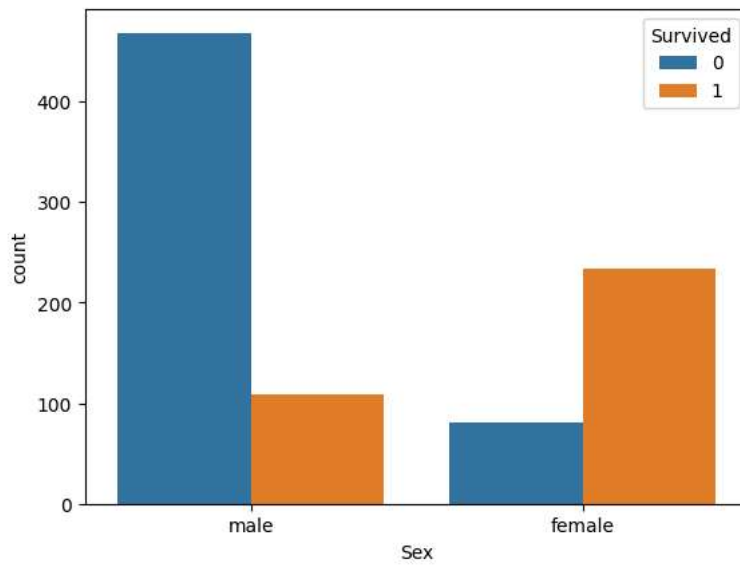
```
sns.countplot(x="Sex", data=df)
plt.show()

sns.countplot(x="Survived", data=df)
plt.show()
sns.histplot(df["Age"].dropna(), bins=20, kde=True)
plt.show()

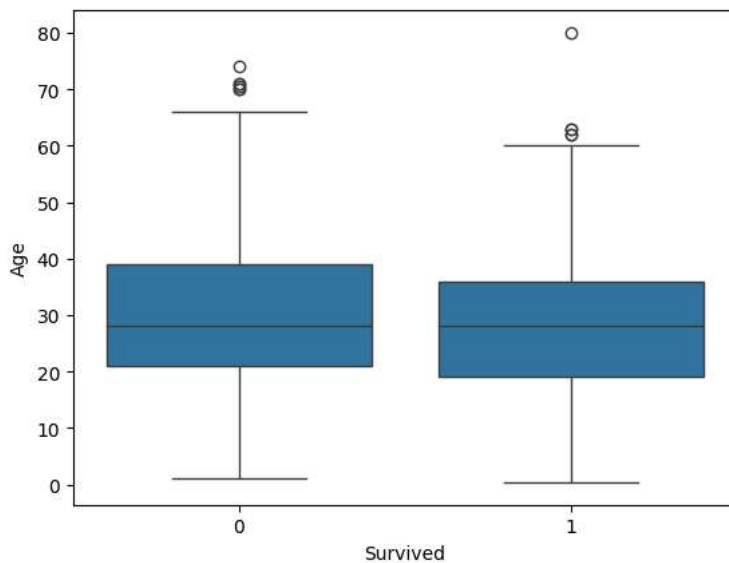
sns.boxplot(x=df["Fare"])
```

```
plt.show()
```

```
sns.countplot(x="Sex", hue="Survived", data=df)  
plt.show()
```

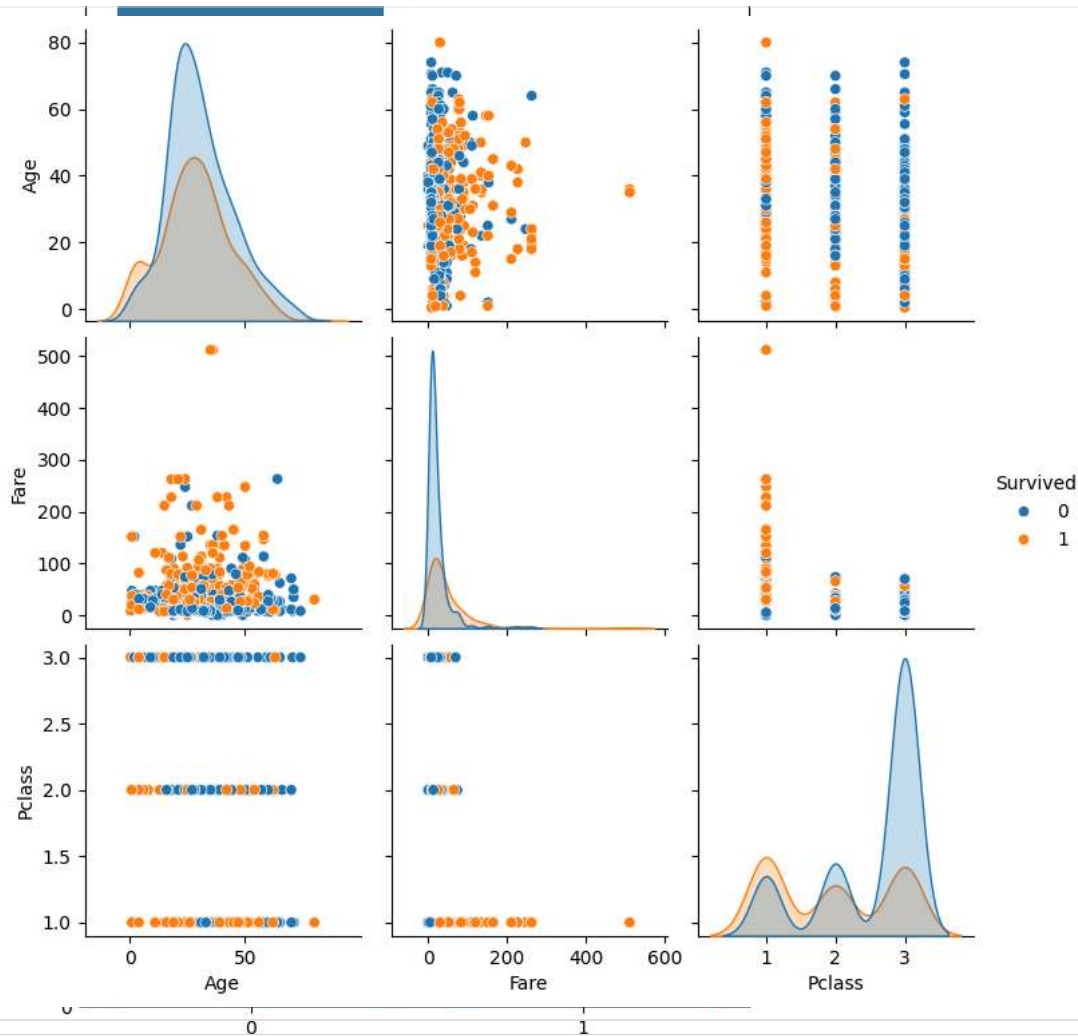


```
sns.boxplot(x="Survived", y="Age", data=df)  
plt.show()
```



Start coding or [generate](#) with AI.

```
sns.pairplot(df[["Survived", "Age", "Fare", "Pclass"]], hue="Survived")
plt.show()
```



```
df["Age"].fillna(df["Age"].median(), inplace=True)
df.drop("Cabin", axis=1, inplace=True) # too many missing
```

/tmp/ipython-input-2910835745.py:1: FutureWarning: A value is trying to be set on a copy of a DataFrame or Series through chained assignment. The behavior will change in pandas 3.0. This inplace method will never work because the intermediate object on which we are setting is 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, inplace=True)

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

