

```
In [ ]: import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
import warnings
warnings.filterwarnings("ignore")
%matplotlib inline
```

Load data and basic stats

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

```
In [ ]: df.shape
```

```
Out[ ]: (891, 12)
```

```
In [ ]: df.head()
```

```
Out[ ]:
```

	PassengerId	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Tic
0	1	0	3	Braund, Mr. Owen Harris	male	22.0	1	0	21
1	2	1	1	Cumings, Mrs. John Bradley (Florence Briggs Th...	female	38.0	1	0	PC 17
2	3	1	3	Heikkinen, Miss. Laina	female	26.0	0	0	STON/ 3101
3	4	1	1	Futrelle, Mrs. Jacques Heath (Lily May Peel)	female	35.0	1	0	113
4	5	0	3	Allen, Mr. William Henry	male	35.0	0	0	373

```
In [ ]: df.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

```

```
In [ ]: df.describe()
```

```
Out[ ]:
```

	PassengerId	Survived	Pclass	Age	SibSp	Parc
count	891.000000	891.000000	891.000000	714.000000	891.000000	891.000000
mean	446.000000	0.383838	2.308642	29.699118	0.523008	0.381590
std	257.353842	0.486592	0.836071	14.526497	1.102743	0.806050
min	1.000000	0.000000	1.000000	0.420000	0.000000	0.000000
25%	223.500000	0.000000	2.000000	20.125000	0.000000	0.000000
50%	446.000000	0.000000	3.000000	28.000000	0.000000	0.000000
75%	668.500000	1.000000	3.000000	38.000000	1.000000	0.000000
max	891.000000	1.000000	3.000000	80.000000	8.000000	6.000000

```
In [ ]: df.isna().sum()
```

```
Out[ ]: 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 [ ]: df["Age"] = df["Age"].fillna(df["Age"].mean())
```

```
In [ ]: df.isna().sum()
```

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

Visualization

```
In [ ]: df["Name"]
```

```
Out[ ]: 0      Braund, Mr. Owen Harris
1      Cumings, Mrs. John Bradley (Florence Briggs Th...
2      Heikkinen, Miss. Laina
3      Futrelle, Mrs. Jacques Heath (Lily May Peel)
4      Allen, Mr. William Henry
...
886     Montvila, Rev. Juozas
887     Graham, Miss. Margaret Edith
888     Johnston, Miss. Catherine Helen "Carrie"
889     Behr, Mr. Karl Howell
890     Dooley, Mr. Patrick
Name: Name, Length: 891, dtype: object
```

```
In [ ]: df["Sex"].value_counts()
```

```
Out[ ]: male      577
female    314
Name: Sex, dtype: int64
```

```
In [ ]: df["Ticket"].value_counts()
```

```
Out[ ]: 347082      7
CA. 2343          7
1601              7
3101295           6
CA 2144           6
..
9234              1
19988             1
2693              1
PC 17612          1
370376            1
Name: Ticket, Length: 681, dtype: int64
```

```
In [ ]: df["Cabin"].value_counts()
```

```
Out[ ]: B96 B98      4
        G6         4
        C23 C25 C27  4
        C22 C26     3
        F33         3
        ..
        E34         1
        C7          1
        C54         1
        E36         1
        C148        1
        Name: Cabin, Length: 147, dtype: int64
```

```
In [ ]: df["Embarked"].value_counts()
```

```
Out[ ]: S      644
        C      168
        Q       77
        Name: Embarked, dtype: int64
```

```
In [ ]: def fun1(value):
        if (value == "male"):
            return 1
        else:
            return 0
```

```
In [ ]: def fun2(value):
        if (value == 'S'):
            return 0
        elif (value == 'C'):
            return 1
        elif (value == 'Q'):
            return 2
        else:
            return 0
```

```
In [ ]: df["Sex"] = df["Sex"].apply(fun1)
```

```
In [ ]: df["Embarked"] = df["Embarked"].apply(fun2)
```

```
In [ ]: df.isna().sum()
```

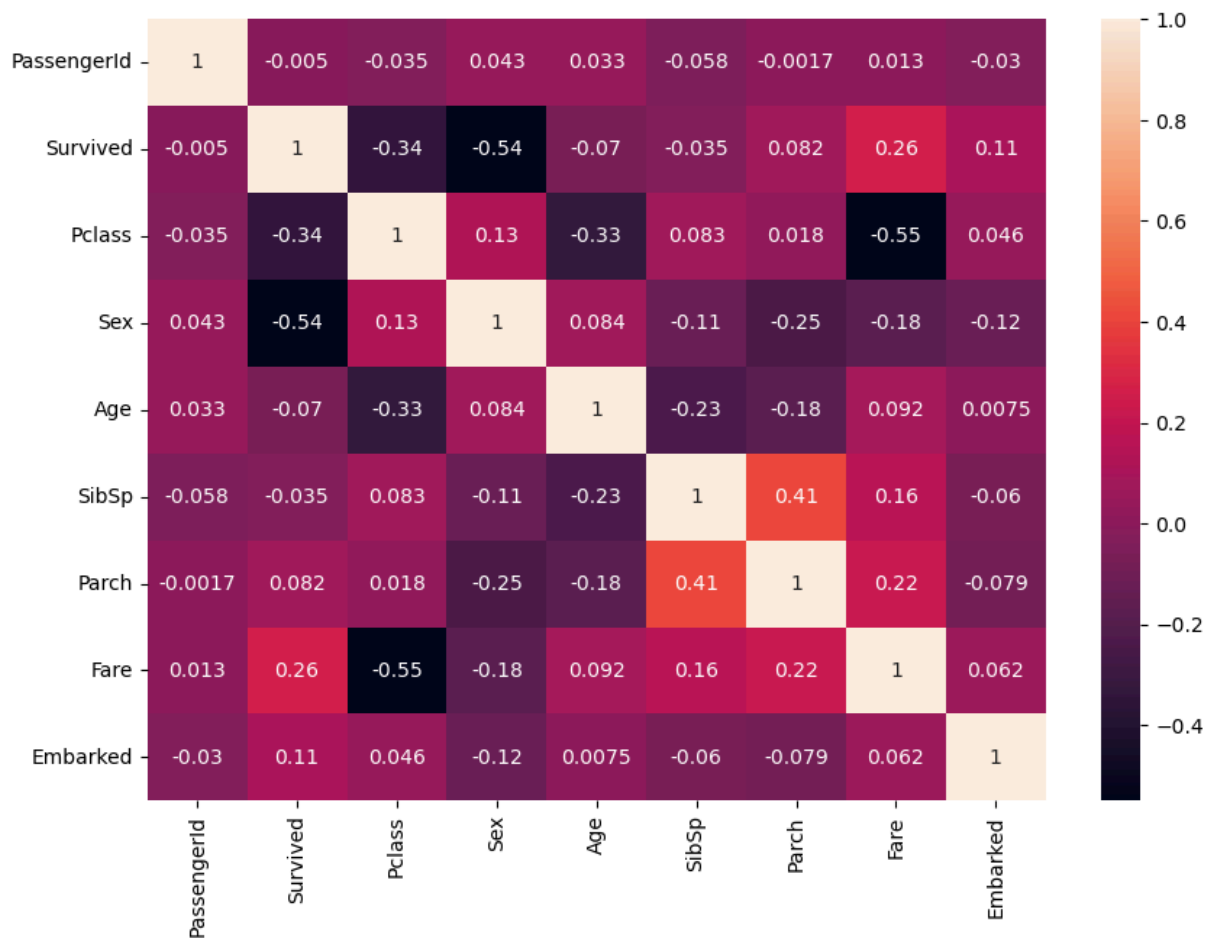
```
Out[ ]: PassengerId      0
Survived      0
Pclass        0
Name          0
Sex           0
Age           0
SibSp         0
Parch         0
Ticket        0
Fare          0
Cabin         687
Embarked      0
dtype: int64
```

```
In [ ]: df = df.drop("Cabin", axis=1)
```

```
In [ ]: df.shape
```

```
Out[ ]: (891, 11)
```

```
In [ ]: plt.figure(figsize=(10,7))
sns.heatmap(df.corr(), annot=True)
plt.show()
```

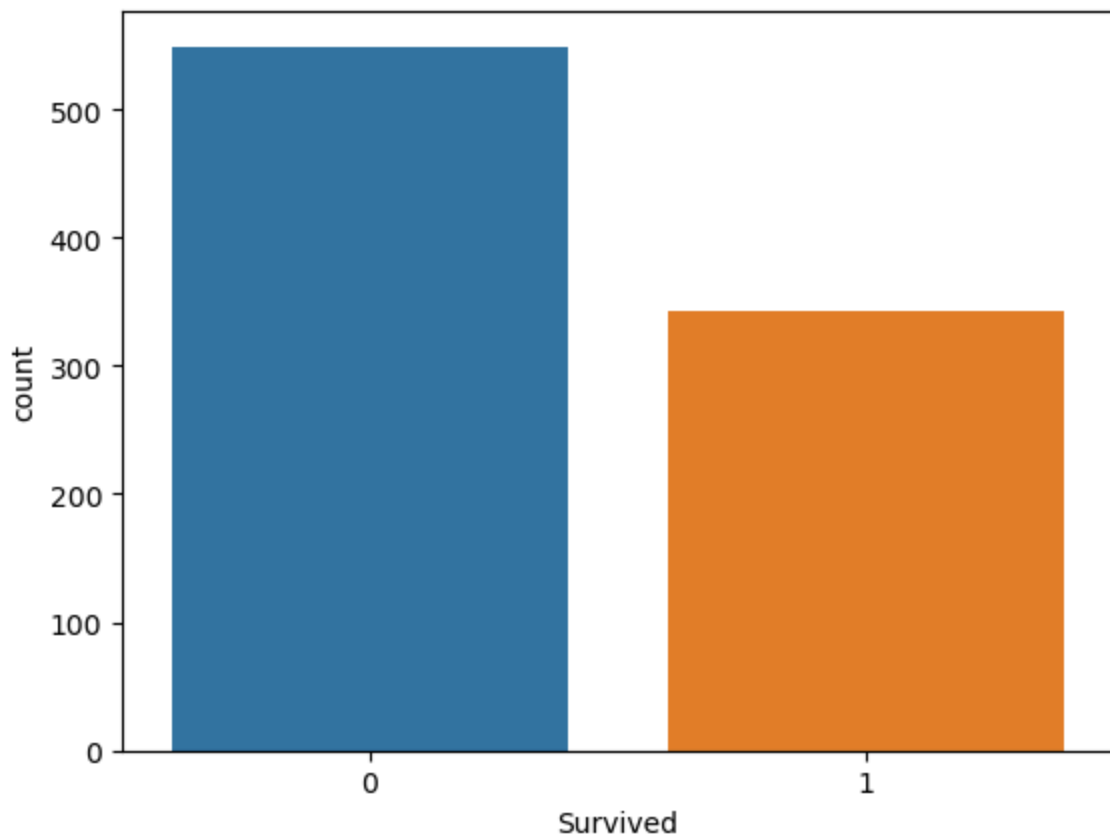


```
In [ ]: df.info()
```

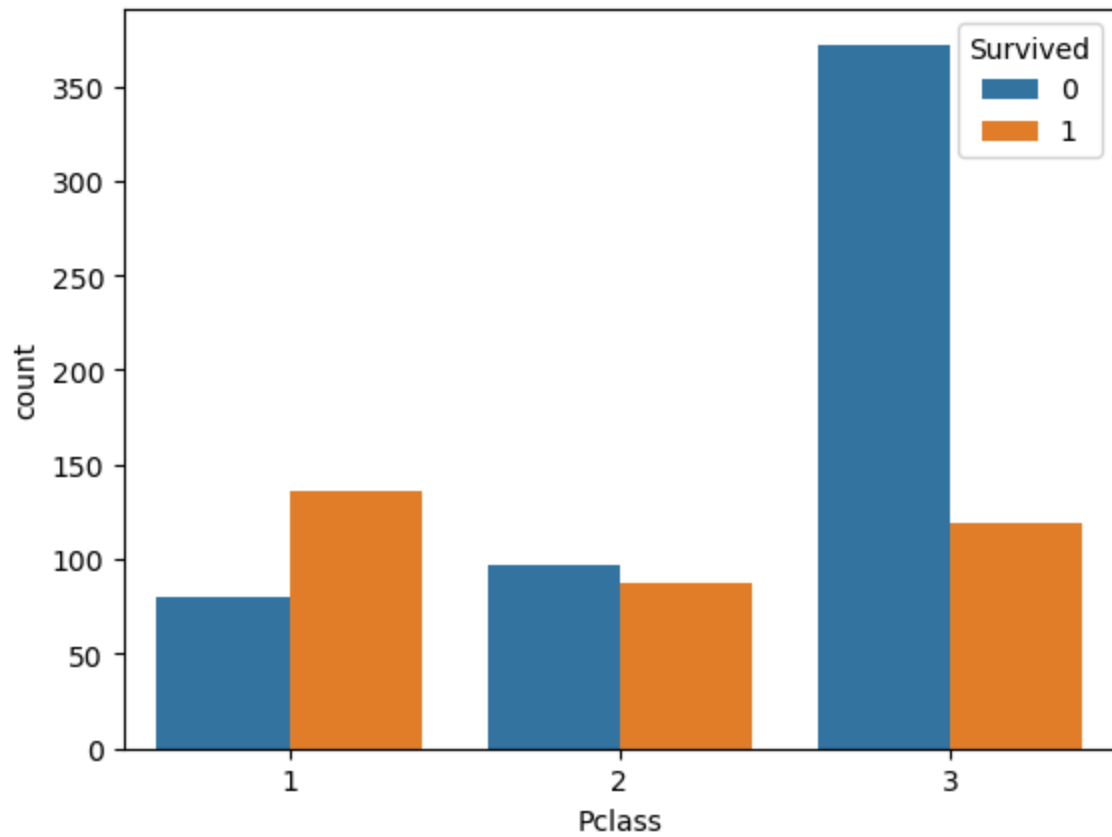
```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 891 entries, 0 to 890
Data columns (total 11 columns):
#   Column      Non-Null Count  Dtype
---  -
0   PassengerId  891 non-null    int64
1   Survived     891 non-null    int64
2   Pclass       891 non-null    int64
3   Name         891 non-null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

"Survived" is the label

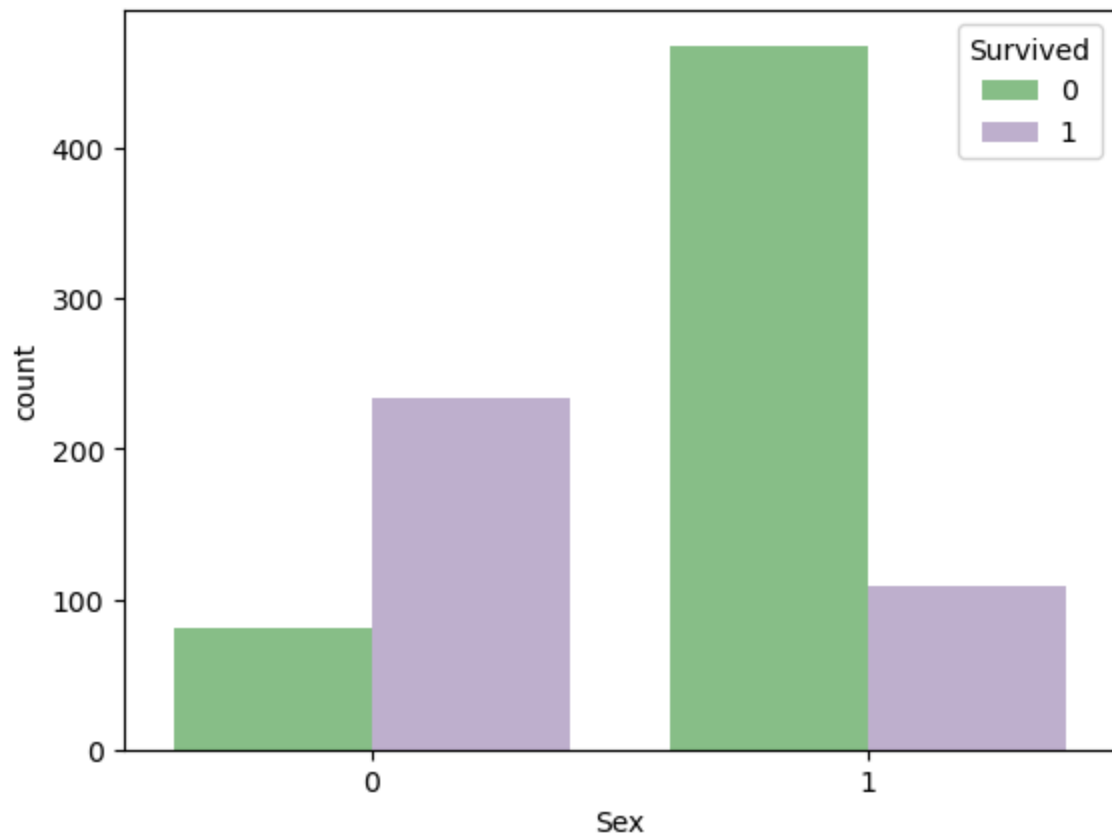
```
In [ ]: sns.countplot(df["Survived"])
plt.show()
```



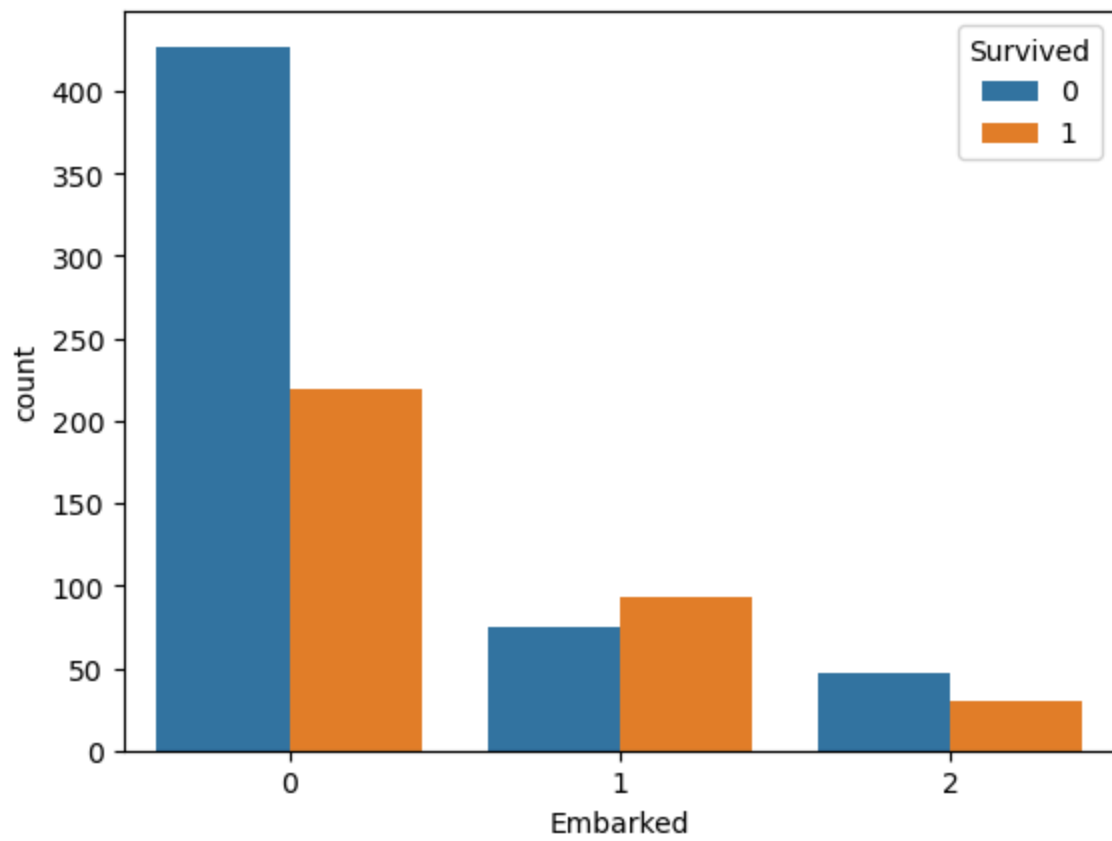
```
In [ ]: sns.countplot(df["Pclass"], hue=df["Survived"])
plt.show()
```



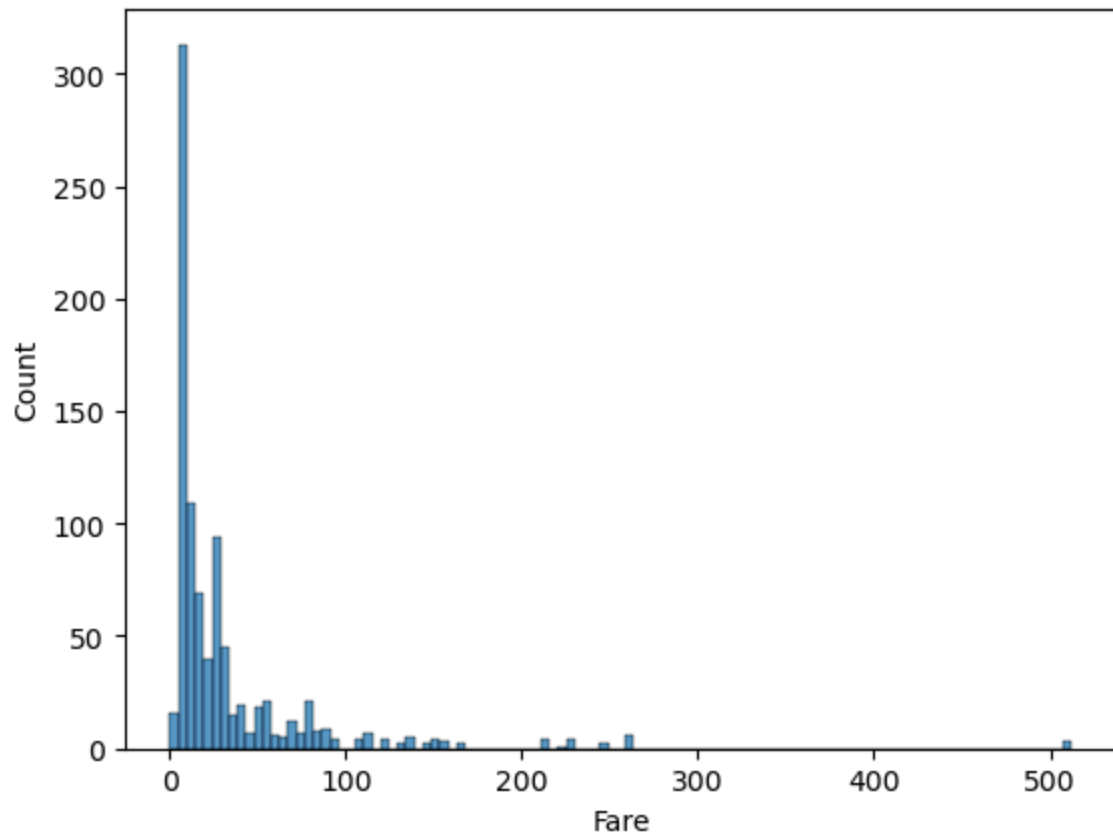
```
In [ ]: sns.countplot(df["Sex"], hue=df["Survived"], palette="Accent")  
plt.show()
```



```
In [ ]: sns.countplot(df["Embarked"], hue=df["Survived"])
plt.show()
```



```
In [ ]: sns.histplot(df["Fare"])
plt.show()
```

This notebook was converted with convert.ploomber.io