

21BCE7151

Import the Libraries

```
In [2]: import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns

Import the Dataset

In [3]: df=pd.read_csv("Titanic-Dataset.csv")

In [4]: df.head()

Out[4]:
   PassengerId  Survived  Pclass    Name  Sex  SibSp  Parch    Ticket   Fare Cabin Embarked
0            0         0       3  Braund, Mr. Owen Harris   male    0       1       0  A/5 21171   7.2500   NaN      S
1            1         0       3  Cumings, Mrs. John Bradley (Florence Briggs Th... female   38.0    1       0  PC 17599  71.2833   C85      C
2            2         1       3  Heikkinen, Miss. Laina   female   26.0    0       0  STON/O2 3101282  7.9250   NaN      S
3            3         1       1  Futrelle, Mrs. Jacques Heath (Lily May Peel) female   35.0    1       0       0  113803  53.1000  C123      S
4            4         0       3        Allen, Mr. William Henry   male   35.0    0       0       0  373450   8.0500   NaN      S

In [5]: 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              284 non-null    object
11  Embarked            891 non-null    object
dtypes: float64(2), int64(5), object(5)
memory usage: 85.7+ KB

In [6]: df.describe()

Out[6]:
   PassengerId  Survived  Pclass    Age    SibSp  Parch    Fare
count  891.000000  891.000000  891.000000  714.000000  891.000000  891.000000  891.000000
mean    446.000000  0.383838  2.308642  29.699118  0.523008  0.381594  32.204208
std     257.353842  0.486592  0.636071  13.002015  1.102743  0.806057  49.693429
min       1.000000  0.000000  1.000000  0.420000  0.000000  0.000000  7.910400
25%    223.500000  0.000000  2.000000  20.125000  0.000000  0.000000  14.454200
50%    446.000000  0.000000  3.000000  28.000000  1.000000  0.000000  31.000000
75%    668.500000  1.000000  3.000000  38.000000  1.000000  0.000000  51.225000
max    891.000000  1.000000  3.000000  80.000000  9.000000  6.000000  512.329000

In [7]: df.corr()
<ipython-input-7-2f6f688a2c>:1: FutureWarning: The default value of numeric_only in DataFrame.corr is deprecated. In a future version, it will default to False. Select only valid columns or specify the value of numeric_only to silence this warning.
df.corr()

Out[7]:
   PassengerId  Survived  Pclass    Age    SibSp  Parch    Fare
Survived  -0.005007  1.000000  -0.338481  0.036847  -0.07527  -0.01182  0.012658
Pclass     -0.035144  -0.338481  1.000000  -0.369226  0.030322  -0.018493  -0.249500
Age         0.036847  -0.077221  -0.369226  1.000000  -0.308247  -0.189119  0.096067
SibSp       -0.075227  -0.03322  -0.030321  -0.308247  1.000000  -0.414939  0.199601
Parch       -0.001652  0.081629  0.018443  -0.189119  0.414938  1.000000  0.216225
Fare        0.012658  0.012629  -0.249500  0.096067  0.199601  0.216225  1.000000

In [8]: df.corr().Survived.sort_values(ascending = False)
<ipython-input-8-936bcb2ae2>:1: FutureWarning: The default value of numeric_only in DataFrame.corr is deprecated. In a future version, it will default to False. Select only valid columns or specify the value of numeric_only to silence this warning.
df.corr().Survived.sort_values(ascending = False)

Out[8]:
Survived    1.000000
Fare        0.257387
Parch       0.081629
PassengerId -0.005007
SibSp       -0.035322
Age         -0.077221
Pclass      -0.338481
Name: Survived, dtype: float64

Handling Missing/Null Values

In [9]: df.isnull().any()

Out[9]:
PassengerId    False
Survived        False
Pclass          False
Name            False
Sex             False
Age             True
SibSp           False
Parch           False
Ticket          False
Fare            False
Cabin           True
Embarked        False
dtype: bool

In [10]: sum(df.cabin.isnull())

Out[10]: 687

In [11]: sum(df.Age.isnull())

Out[11]: 177

In [12]: df["Age"].fillna(df["Age"].mean(),inplace=True)

In [13]: sum(df.Embarked.isnull())

Out[13]: 2

In [14]: df["Embarked"].fillna(df["Embarked"].mode()[0],inplace=True)

In [15]: df.describe()

Out[15]:
   PassengerId  Survived  Pclass    Age    SibSp  Parch    Fare
count  891.000000  891.000000  891.000000  891.000000  891.000000  891.000000  891.000000
mean    446.000000  0.383838  2.308642  29.699118  0.523008  0.381594  32.204208
std     257.353842  0.486592  0.636071  13.002015  1.102743  0.806057  49.693429
min       1.000000  0.000000  1.000000  0.420000  0.000000  0.000000  7.910400
25%    223.500000  0.000000  2.000000  20.125000  0.000000  0.000000  14.454200
50%    446.000000  0.000000  3.000000  28.000000  1.000000  0.000000  31.000000
75%    668.500000  1.000000  3.000000  38.000000  1.000000  0.000000  51.225000
max    891.000000  1.000000  3.000000  80.000000  9.000000  6.000000  512.329000

Data Visualization

In [16]: plt.scatter(df["Fare"],df["Survived"])

Out[16]: <matplotlib.collections.PathCollection at 0x795caf2e14e0>

```

sns.heatmap(df.corr(),annot=True)

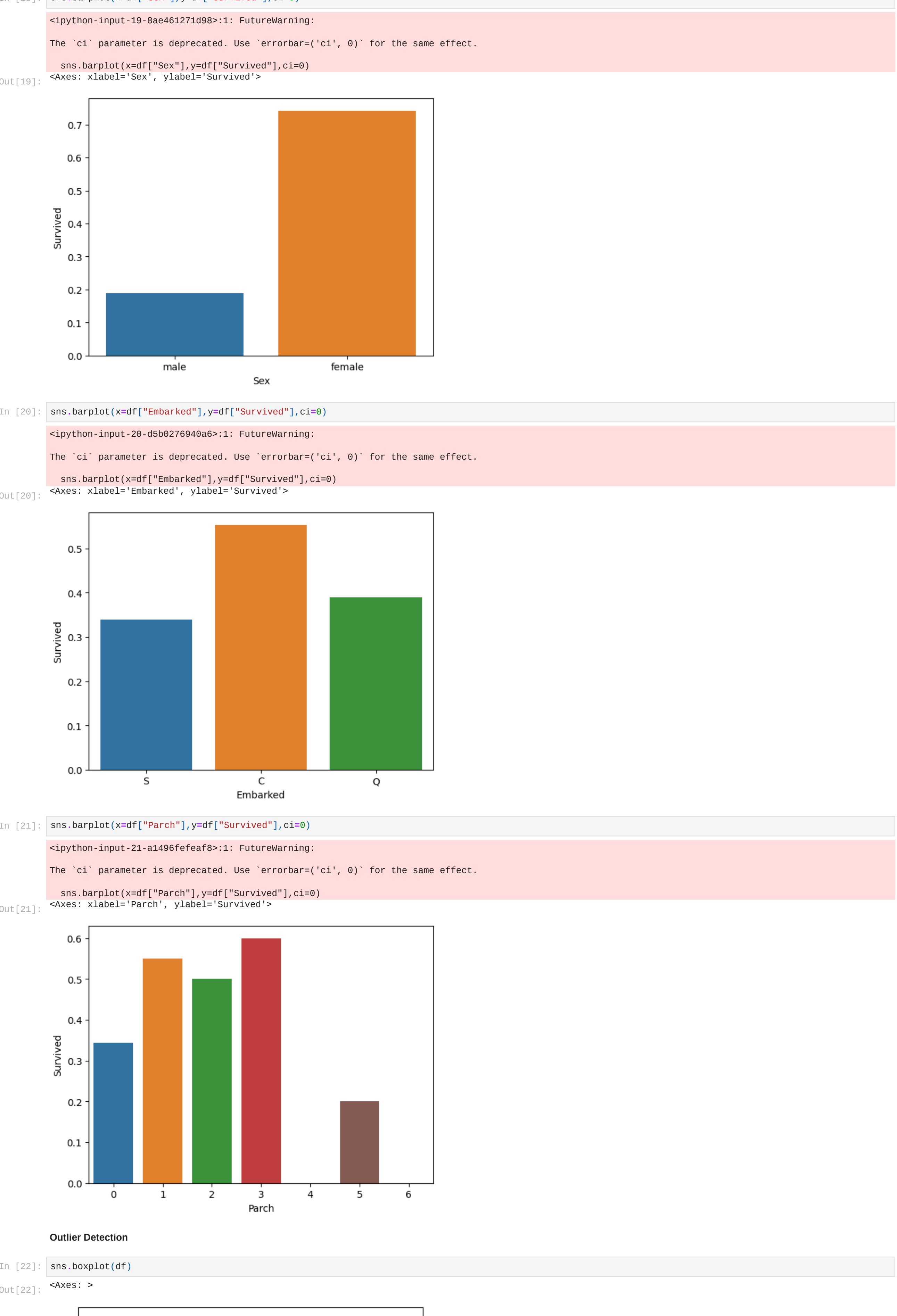
<ipython-input-17-8d7fbca528>:1: FutureWarning: The default value of numeric_only in DataFrame.corr is deprecated. In a future version, it will default to False. Select only valid columns or specify the value of numeric_only to silence this warning.

<Axes: >



sns.pairplot(df)

<seaborn.axisgrid.PairGrid at 0x795cf8707eb>

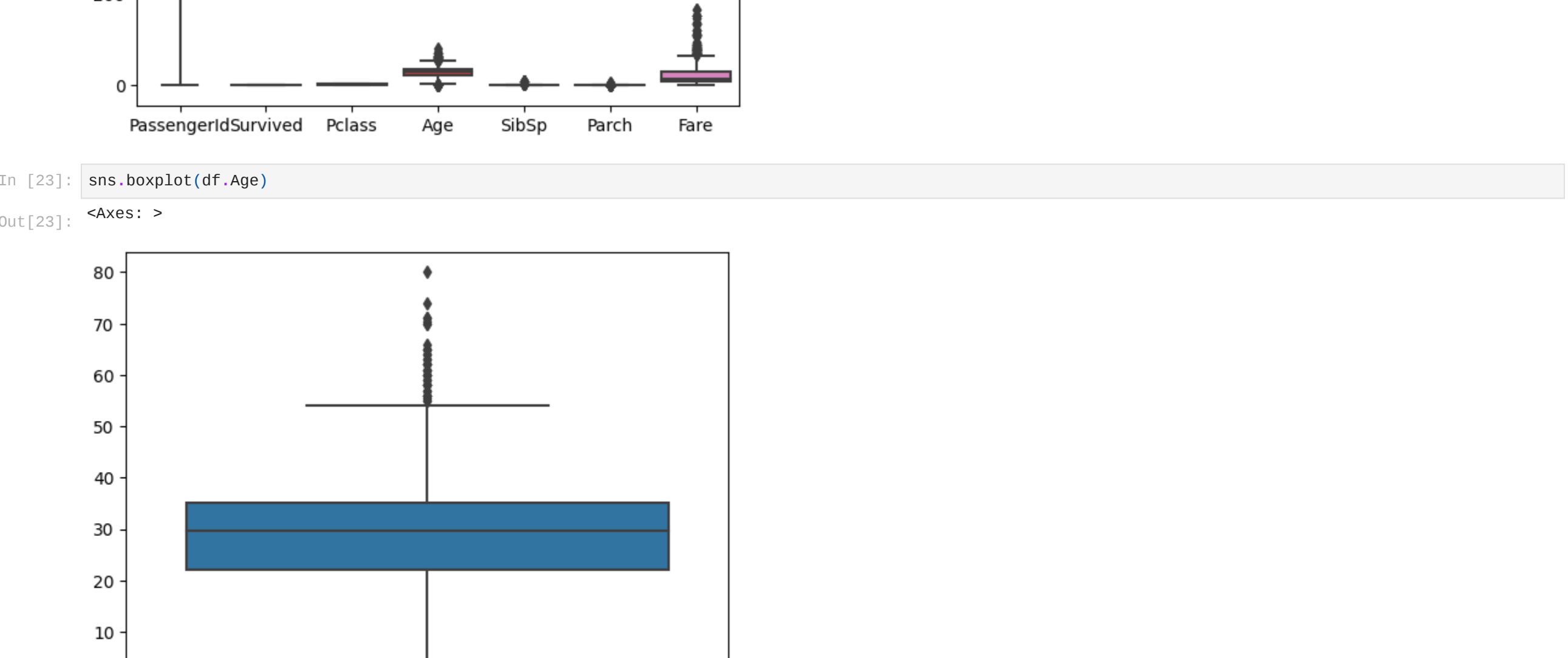


sns.barplot(x=df["Sex"],y=df["Survived"],ci=0)

<ipython-input-19-8ae461271d9b>:1: FutureWarning: The 'ci' parameter is deprecated. Use 'errorbar=('ci', 0)' for the same effect.

sns.barplot(x=df["Sex"],y=df["Survived"],ci=0)

<Axes: xlabel='Sex', ylabel='Survived'>

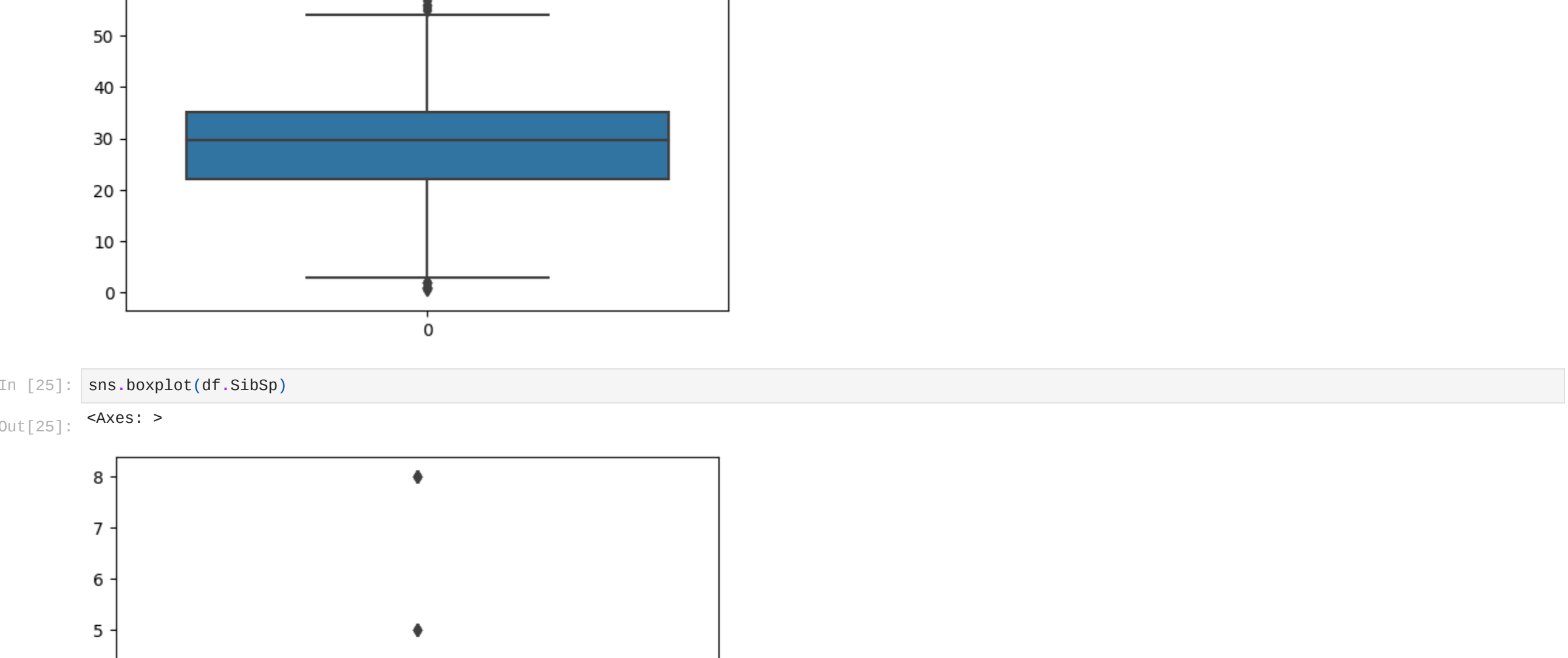


sns.barplot(x=df["Embarked"],y=df["Survived"],ci=0)

<ipython-input-20-d5b6276848a>:1: FutureWarning: The 'ci' parameter is deprecated. Use 'errorbar=('ci', 0)' for the same effect.

sns.barplot(x=df["Embarked"],y=df["Survived"],ci=0)

<Axes: xlabel='Embarked', ylabel='Survived'>

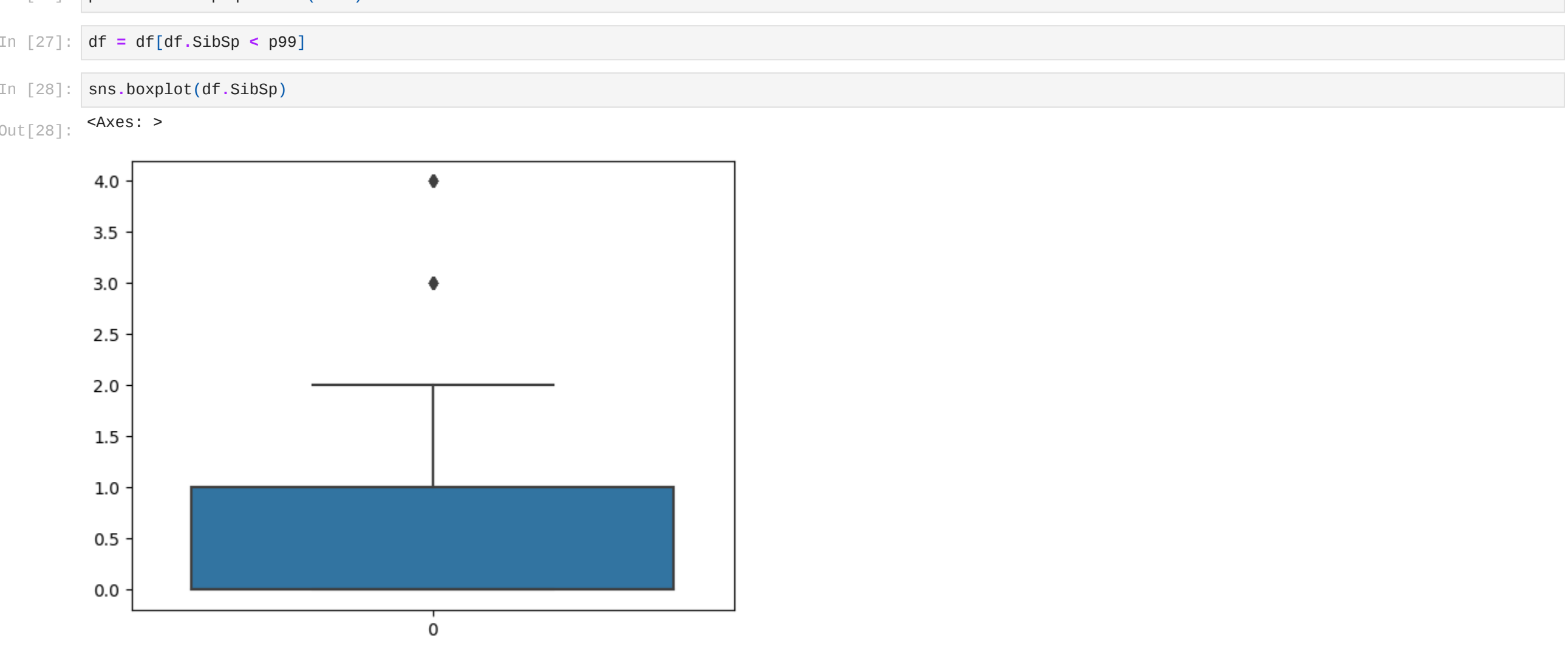


sns.barplot(x=df["Parch"],y=df["Survived"],ci=0)

<ipython-input-21-a1496fefeaf>:1: FutureWarning: The 'ci' parameter is deprecated. Use 'errorbar=('ci', 0)' for the same effect.

sns.barplot(x=df["Parch"],y=df["Survived"],ci=0)

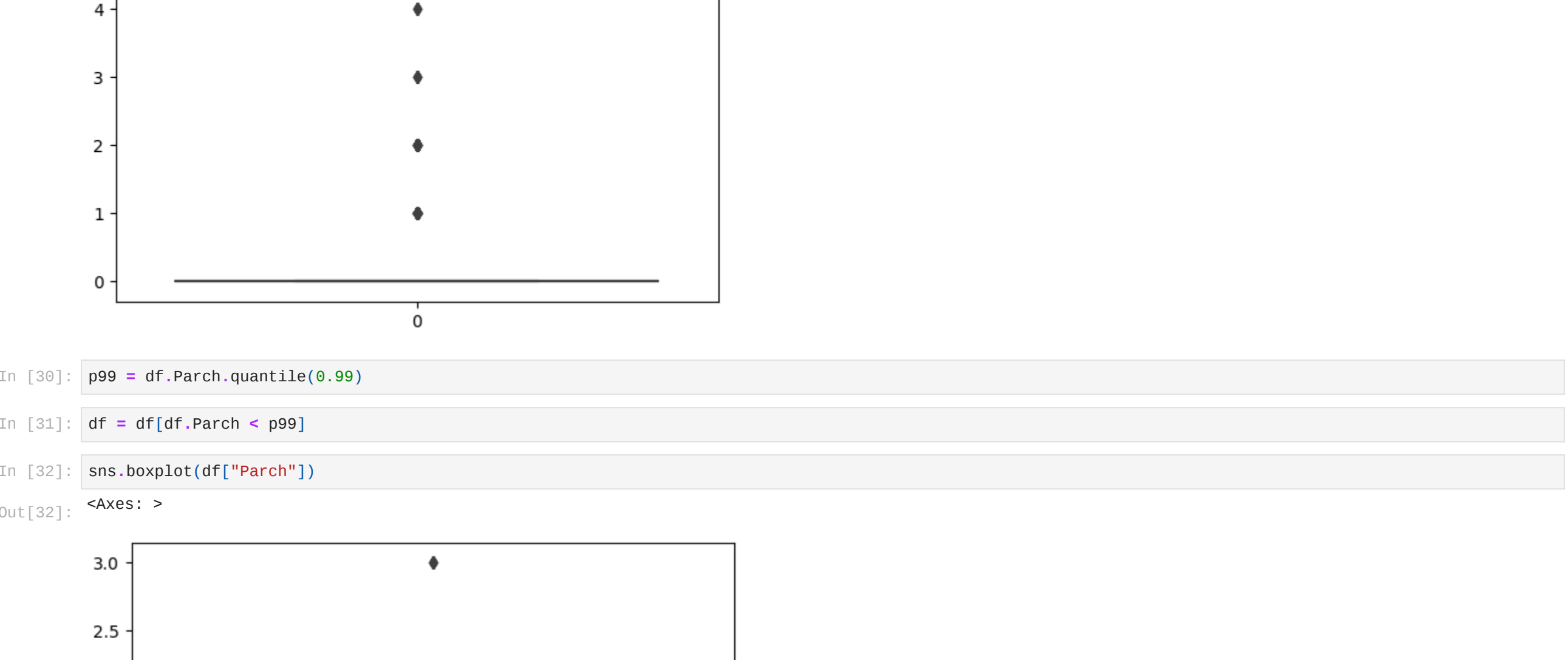
<Axes: xlabel='Parch', ylabel='Survived'>



Outlier Detection

sns.boxplot(df)

<Axes: >



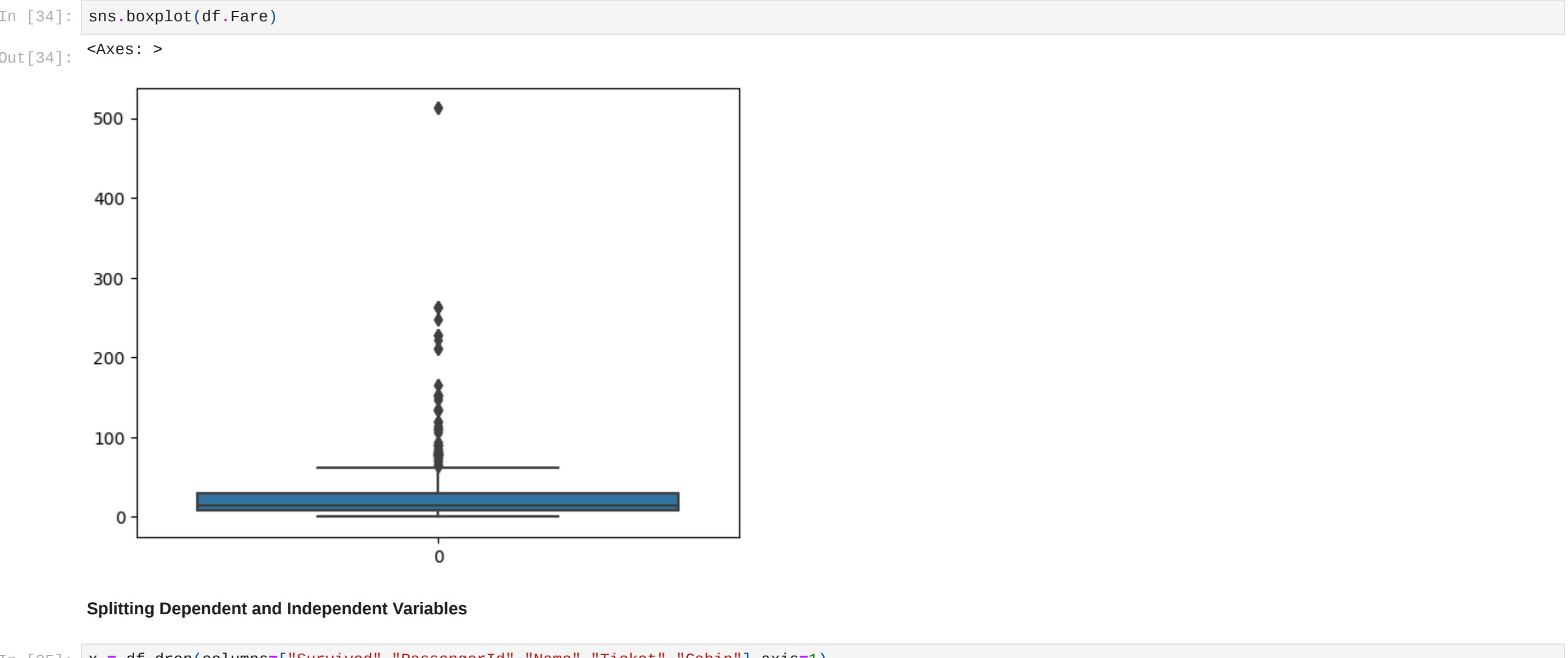
sns.boxplot(df.Age)

<Axes: >



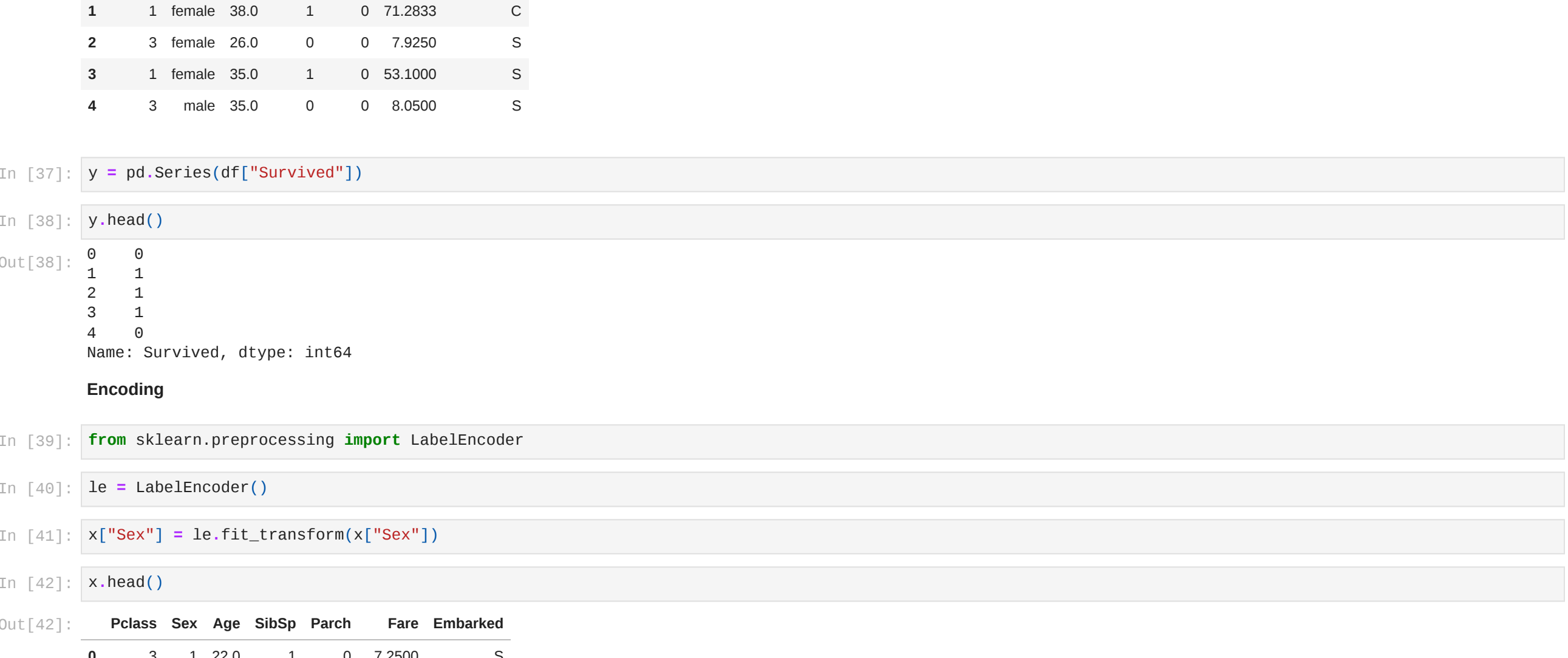
sns.boxplot(df.Age)

<Axes: >



sns.boxplot(df.SibSp)

<Axes: >



p99 = df.SibSp.quantile(0.99)

df = df[df.SibSp <= p99]

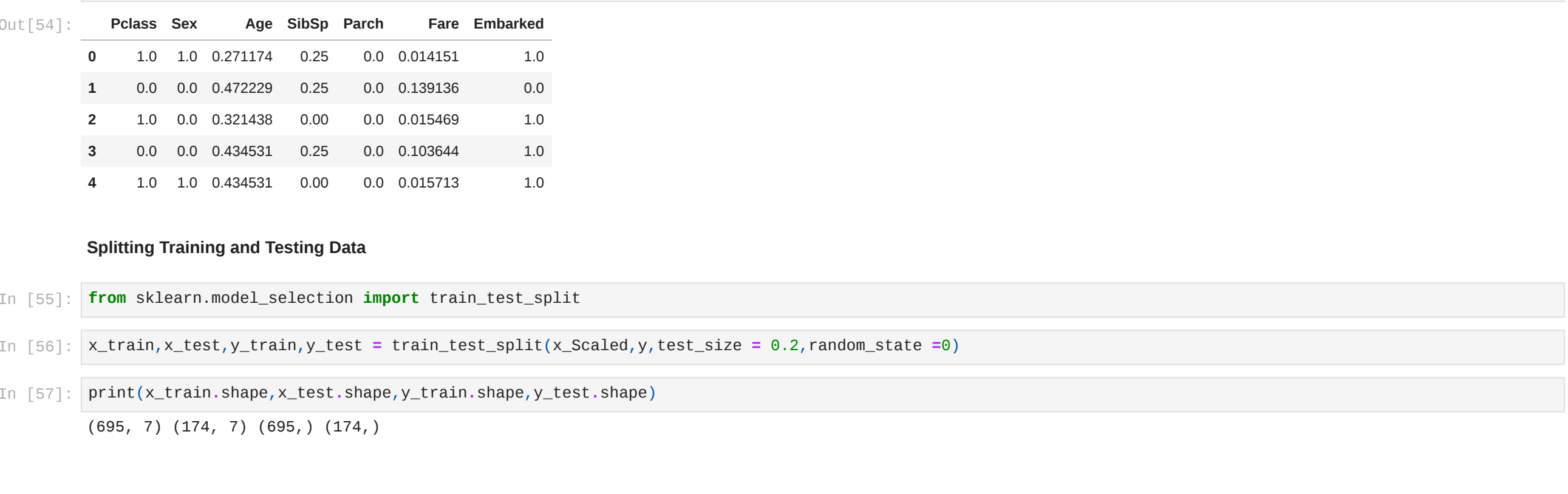
sns.boxplot(df.SibSp)

<Axes: >



sns.boxplot(df.Parch)

<Axes: >



p99 = df.Parch.quantile(0.99)

df = df[df.Parch <= p99]

sns.boxplot(df.Parch)

<Axes: >

sns.boxplot(df.Fare)

<Axes: >

sns.boxplot(df.Fare)

<Axes: >

Splitting Dependent and Independent Variables

x = df.drop(columns=["Survived","PassengerId","Name","Ticket","Cabin"],axis=1)

<Independent variables should be in df or df array

x.head()

<Out[36]:>

y.head()

<Out[38]:>

le = LabelEncoder()

x["Sex"] = le.fit_transform(x["Sex"])

x.head()

<Out[42]:>

print(le.classes_)

['female' 'male']

mapping=dict(zip(le.classes_,range(len(le.classes_))))

mapping

['female': 0, 'male': 1]

le = LabelEncoder()

x["Embarked"] = le.fit_transform(x["Embarked"])

x.head()

<Out[48]:>

print(le.classes_)

['C' 'Q' 'S']

mapping=dict(zip(le.classes_,range(len(le.classes_))))

mapping

['C': 0, 'Q': 1, 'S': 2]

Feature Scaling

from sklearn.preprocessing import MinMaxScaler

ms = MinMaxScaler()

X_Scaled = pd.DataFrame(ms.fit_transform(x),columns = x.columns)

X_Scaled.head()

<Out[54]:>

Splitting Training and Testing Data

from sklearn.model_selection import train_test_split

x_train,x_test,y_train,y_test = train_test_split(x_Scaled,y,test_size = 0.2,random_state =0)

print(x_train.shape,x_test.shape,y_train.shape,y_test.shape)

(695, 7) (114, 7) (695,) (114,)