

Import Libraries

```
import pandas as pd
import numpy as np
import seaborn as sns
import matplotlib.pyplot as plt
import warnings
warnings.filterwarnings('ignore')
%matplotlib inline
```

Loading dataset

```
In [3]: train = pd.read_csv('kaggle/input/titanic/train.csv')
test = pd.read_csv('kaggle/input/titanic/test.csv')
```

| PassengerId Survived Pclass | | | | | | | | | | Name | Sex | Age | SibSp | Parch | Ticket | Fare | Cabin | Embarked |
|-----------------------------|---|---|---|--|--|--|--|--|--|---|--------|------|-------|-------|-----------------|---------|-------|----------|
| 0 | 1 | 0 | 3 | | | | | | | Braund, Mr. Owen Harris | male | 22.0 | 1 | 0 | A/5 21171 | 7.2500 | NaN | S |
| 1 | 2 | 1 | 1 | | | | | | | Cummings, Mrs. John Bradley (Florence Briggs Th...) | female | 38.0 | 1 | 0 | PC 17599 | 71.2833 | C85 | C |
| 2 | 3 | 1 | 3 | | | | | | | Wheeler, Mrs. Joseph Heath (Lily May Peel) | female | 26.0 | 0 | 0 | STON/OJ 3101282 | 7.2500 | NaN | S |
| 3 | 4 | 1 | 1 | | | | | | | Futrelle, Mrs. Jacques Heath (Lily May Peel) | female | 35.0 | 1 | 0 | 113803 | 53.1000 | C123 | S |
| 4 | 5 | 0 | 3 | | | | | | | Allen, Mr. William Henry | male | 35.0 | 0 | 0 | 373450 | 8.0500 | NaN | S |

| # Statistical Info | | | | | | | | | | train.describe() | | | | | | | | |
|--------------------|--|--|--|--|--|--|--|--|--|------------------|-------------|------------|------------|------------|------------|------------|------------|--|
| | | | | | | | | | | count | 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 | 0.383839 | 2.308642 | 29.699118 | 0.523008 | 0.546854 | 32.254208 | |
| | | | | | | | | | | std | 257.308046 | 0.486692 | 0.806091 | 14.526697 | 1.103743 | 0.806091 | 49.693479 | |
| | | | | | | | | | | min | 1.000000 | 0.000000 | 1.000000 | 0.420000 | 0.000000 | 0.000000 | 0.000000 | |
| | | | | | | | | | | 25% | 223.500000 | 0.000000 | 2.000000 | 20.125000 | 0.000000 | 0.000000 | 7.915000 | |
| | | | | | | | | | | 50% | 446.000000 | 0.000000 | 3.000000 | 28.000000 | 0.000000 | 0.000000 | 14.454200 | |
| | | | | | | | | | | 75% | 668.500000 | 1.000000 | 3.000000 | 38.000000 | 1.000000 | 0.000000 | 31.000000 | |
| | | | | | | | | | | max | 891.000000 | 1.000000 | 3.000000 | 80.000000 | 8.000000 | 6.000000 | 512.329200 | |

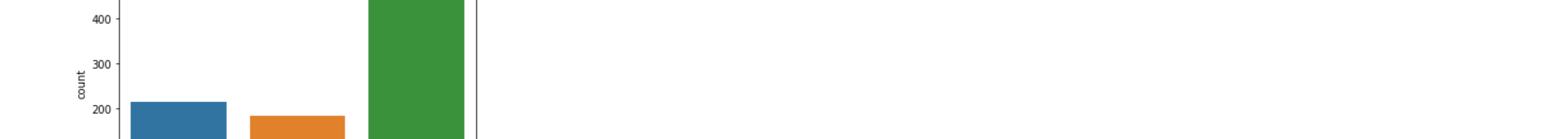
| # datatye info | | | | | | | | | | train.info() | | | | | | | | |
|----------------|--|--|--|--|--|--|--|--|--|---|----------------|--------------|---------|--|--|--|--|--|
| | | | | | | | | | | <class 'pandas.core.frame.DataFrame'> | | | | | | | | |
| | | | | | | | | | | RangeIndex: 891 entries, 0 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: 53.7+ KB | | | | | | | | |

Exploratory Data Analysis

```
In [6]: sns.pairplot(train)
sns.countplot(train['Survived'])
```



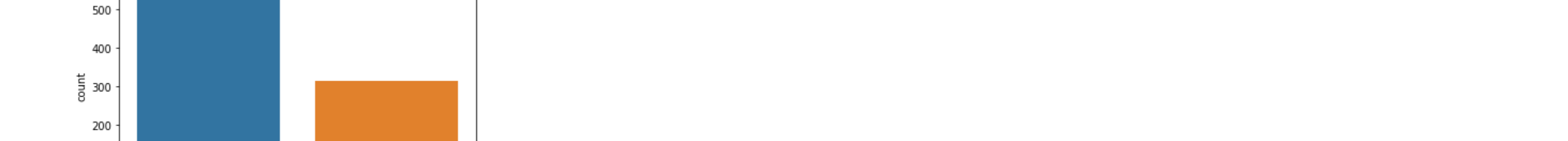
```
In [7]: sns.countplot(train['Pclass'])
```



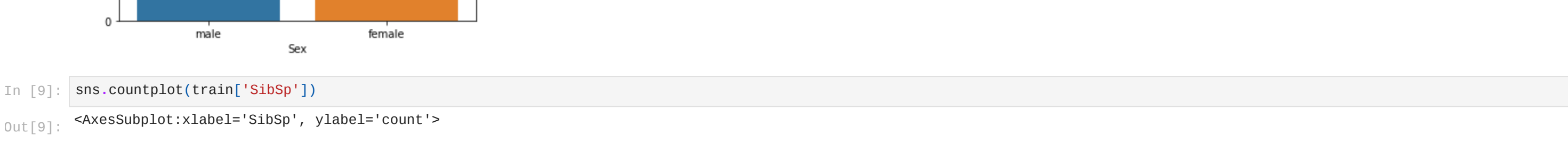
```
In [8]: sns.countplot(train['Sex'])
```



```
In [9]: sns.countplot(train['SibSp'])
```



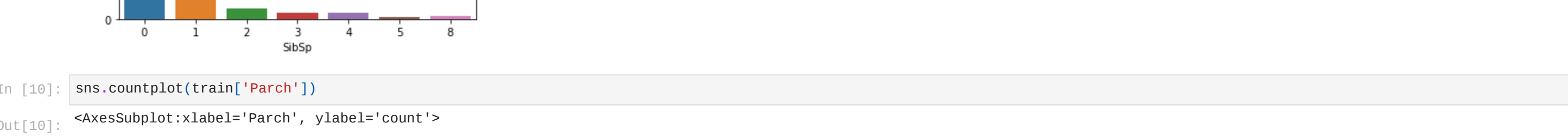
```
In [10]: sns.countplot(train['Parch'])
```



```
In [11]: sns.countplot(train['Embarked'])
```



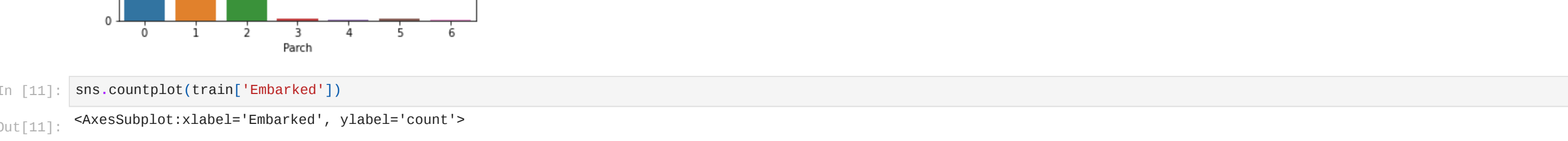
```
In [12]: sns.pairplot(train[['Age', 'Fare']])
```



```
In [13]: sns.distplot(train['Fare'])
```



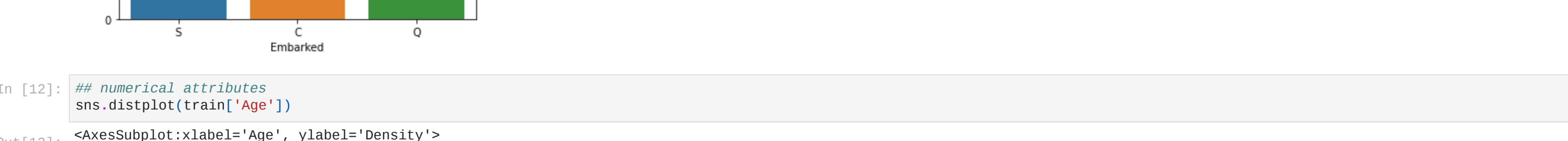
```
In [14]: class_fare = train.pivot_table(index='Pclass', values='Fare')
```



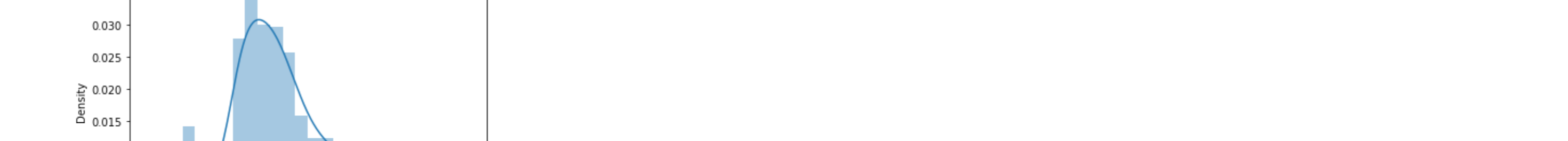
```
In [15]: class_fare = train.pivot_table(index='Pclass', values='Fare', aggfunc=np.sum)
```



```
In [16]: sns.barplot(data=train, x='Pclass', y='Fare', hue='Survived')
```



```
In [17]: sns.barplot(data=train, x='Survived', y='Fare', hue='Pclass')
```



Data Preprocessing

```
In [18]: train_len = len(train)
# combine two datasets
df = pd.concat([train, test], axis=0)
df = df.reset_index(drop=True)
```

| PassengerId Survived Pclass | | | | | | | | | | Name | Sex | Age | SibSp | Parch | Ticket | Fare | Cabin | Embarked |
|-----------------------------|---|-----|---|--|--|--|--|--|--|--|--------|------|-------|-------|-----------------|---------|-------|----------|
| 0 | 1 | 0.0 | 3 | | | | | | | Braund, Mr. Owen Harris | male | 22.0 | 1 | 0 | A/5 21171 | 7.2500 | NaN | S |
| 1 | 2 | 1.0 | 1 | | | | | | | Cummings, Mrs. John Bradley (Florence Briggs Th... | female | 38.0 | 1 | 0 | PC 17599 | 71.2833 | C85 | C |
| 2 | 3 | 1.0 | 3 | | | | | | | Hekkinen, Miss. Laina | female | 26.0 | 0 | 0 | STON/OJ 3101282 | 7.2500 | NaN | S |
| 3 | 4 | 1.0 | 1 | | | | | | | Futrelle, Mrs. Jacques Heath (Lily May Peel) | female | 35.0 | 1 | 0 | 113803 | 53.1000 | C123 | S |
| 4 | 5 | 0.0 | 3 | | | | | | | Allen, Mr. William Henry | male | 35.0 | 0 | 0 | 373450 | 8.0500 | NaN | S |

```
In [19]: df.tail()
```

| PassengerId Survived Pclass | | | | | | | | | | Name | Sex | Age | SibSp | Parch | Ticket | Fare | Cabin | Embarked |
|-----------------------------|------|-----|---|--|--|--|--|--|--|-----------------------------|--------|------|-------|-------|------------------|----------|-------|----------|
| 1384 | 1385 | NaN | 3 | | | | | | | Specter, Mr. Howell | male | NaN | 0 | 0 | A/5 3205 | 8.0500 | NaN | S |
| 1385 | 1386 | NaN | 1 | | | | | | | O'Heay, Orlan. Dona. Femina | female | 38.0 | 0 | 0 | PC 17758 | 108.9000 | C106 | C |
| 1386 | 1387 | NaN | 3 | | | | | | | Sawyer, Mr. Simon Swerten | male | 38.5 | 0 | 0 | 5070/OJQ 3101282 | 7.2500 | NaN | S |
| 1387 | 1388 | NaN | 3 | | | | | | | Waters, Mr. Frederick | male | NaN | 0 | 0 | 369309 | 8.0500 | NaN | S |
| 1388 | 1389 | NaN | 3 | | | | | | | Peter, Master Michael J | male | NaN | 1 | 1 | 2068 | 22.3583 | NaN | C |

```
In [20]: # find the null values
df.isnull().sum()
```

| | |
|-------------|-------|
| PassengerId | 0 |
| Survived | 418 |
| Pclass | 0 |
| Name | 0 |
| Sex | 0 |
| Age | 263 |
| SibSp | 0 |
| Parch | 0 |
| Ticket | 0 |
| Fare | 1 |
| Cabin | 1014 |
| Embarked | 2 |
| dtype: | int64 |

```
In [21]: # drop or delete the column
df = df.drop(columns='Cabin', axis=1)
```

```
In [22]: df['Age'].mean()
```

Out[22]: 29.8813767384014

```
In [23]: # fill missing values using mean of the numerical column
df['Age'] = df['Age'].fillna(df['Age'].mean())
```

```
In [24]: df['Fare'] = df['Fare'].fillna(df['Fare'].mean())
```

```
In [25]: df['Embarked'].mode()[0]
```

Out[25]: 'S'

```
In [26]: # fill missing values using mode of the categorical column
df['Embarked'] = df['Embarked'].fillna(df['Embarked'].mode()[0])
```

Log transformation for uniform data distribution

```
In [27]: sns.distplot(df['Fare'])
```



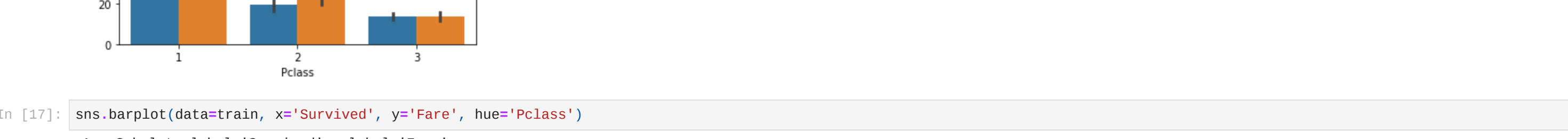
```
In [28]: df['Fare'] = np.log(df['Fare']+1)
```

Out[28]:

<AxesSubplot: xlabel='Fare', ylabel='Density'>

Correlation Matrix

```
In [29]: corr = df.corr()
plt.figure(figsize=(15, 9))
sns.heatmap(corr, annot=True, cmap='coolwarm')
```



```
In [30]: df.head()
```

| PassengerId Survived Pclass | | | | | | | | | | Name | Sex | Age | SibSp | Parch | Ticket | Fare | Cabin | Embarked |
|-----------------------------|---|-----|---|--|--|--|--|--|--|-------------------------|------|------|-------|-------|--------|------|-------|----------|
| 0 | 1 | 0.0 | 3 | | | | | | | Braund, Mr. Owen Harris | male | 22.0 | 1 | 0 | 15160 | 53.1 | C | S |