# lab-08

April 13, 2024

Data Visualization - I

Problem statement :

Use the inbuilt dataset 'titanic'. The dataset contains 891 rows and contains information about the passengers who boarded the unfortunate Titanic ship. 1. Use the Seaborn library to see if we can find any patterns in the data. 2. Write a code to check how the price of the ticket (column name: 'fare') for each passenger is distributed by plotting a histogram.

```python
[1]: #imports
     import pandas as pd
     import numpy as np
     import seaborn as sns
     import matplotlib.pyplot as plt
```

```python
[2]: data = pd.read_csv('train.csv')
     data
```

```
[2]:      PassengerId  Survived  Pclass  \
     0              1         0       3
     1              2         1       1
     2              3         1       3
     3              4         1       1
     4              5         0       3
     ..           ...       ...     ...
     886          887         0       2
     887          888         1       1
     888          889         0       3
     889          890         1       1
     890          891         0       3

                                                        Name     Sex   Age  SibSp  \
     0                              Braund, Mr. Owen Harris    male  22.0      1
     1    Cumings, Mrs. John Bradley (Florence Briggs Th…  female  38.0      1
     2                               Heikkinen, Miss. Laina  female  26.0      0
     3         Futrelle, Mrs. Jacques Heath (Lily May Peel)  female  35.0      1
     4                             Allen, Mr. William Henry    male  35.0      0
     ..                                                 …      …     …      …
     886                              Montvila, Rev. Juozas    male  27.0      0
```

```
887                    Graham, Miss. Margaret Edith  female  19.0      0
888         Johnston, Miss. Catherine Helen "Carrie"  female   NaN      1
889                          Behr, Mr. Karl Howell    male  26.0      0
890                              Dooley, Mr. Patrick    male  32.0      0

     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
..     ...              ...      ...   ...      ...
886      0           211536  13.0000   NaN        S
887      0           112053  30.0000   B42        S
888      2       W./C. 6607  23.4500   NaN        S
889      0           111369  30.0000  C148        C
890      0           370376   7.7500   NaN        Q

[891 rows x 12 columns]
```

[4]: `data.head(5)`

```
[4]:    PassengerId  Survived  Pclass  \
0            1         0       3
1            2         1       1
2            3         1       3
3            4         1       1
4            5         0       3

                                                Name     Sex   Age  SibSp  \
0                            Braund, Mr. Owen Harris    male  22.0      1
1  Cumings, Mrs. John Bradley (Florence Briggs Th…  female  38.0      1
2                             Heikkinen, Miss. Laina  female  26.0      0
3       Futrelle, Mrs. Jacques Heath (Lily May Peel)  female  35.0      1
4                           Allen, Mr. William Henry    male  35.0      0

   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
```

[5]: `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
```

[20]: `data.isna().sum().sum()`

[20]: 866

[21]: `data.isnull().sum()`

[21]:
```
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
```
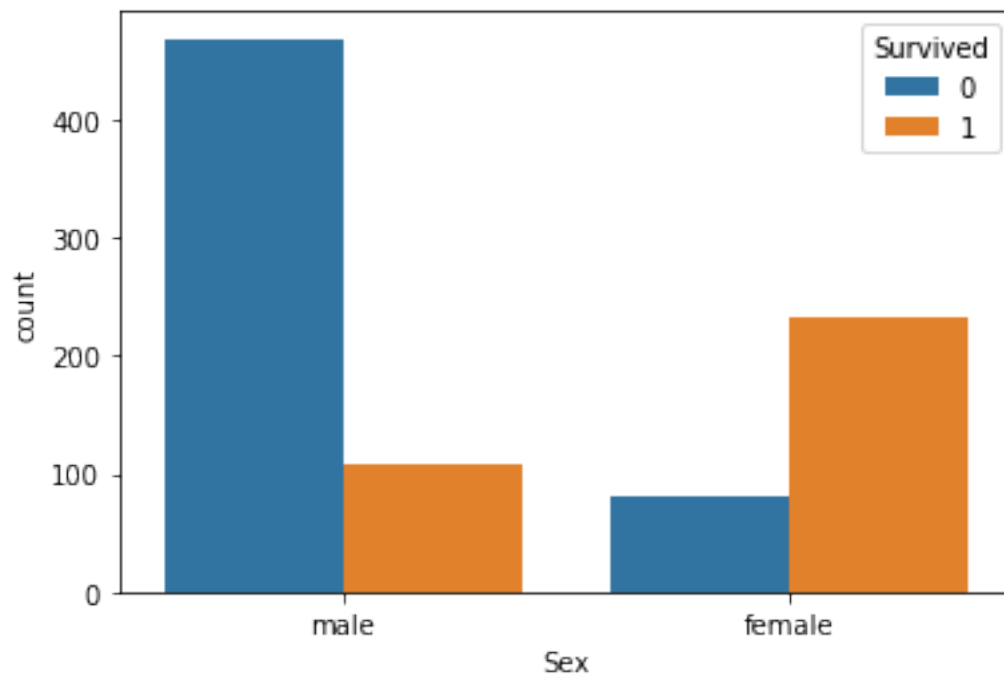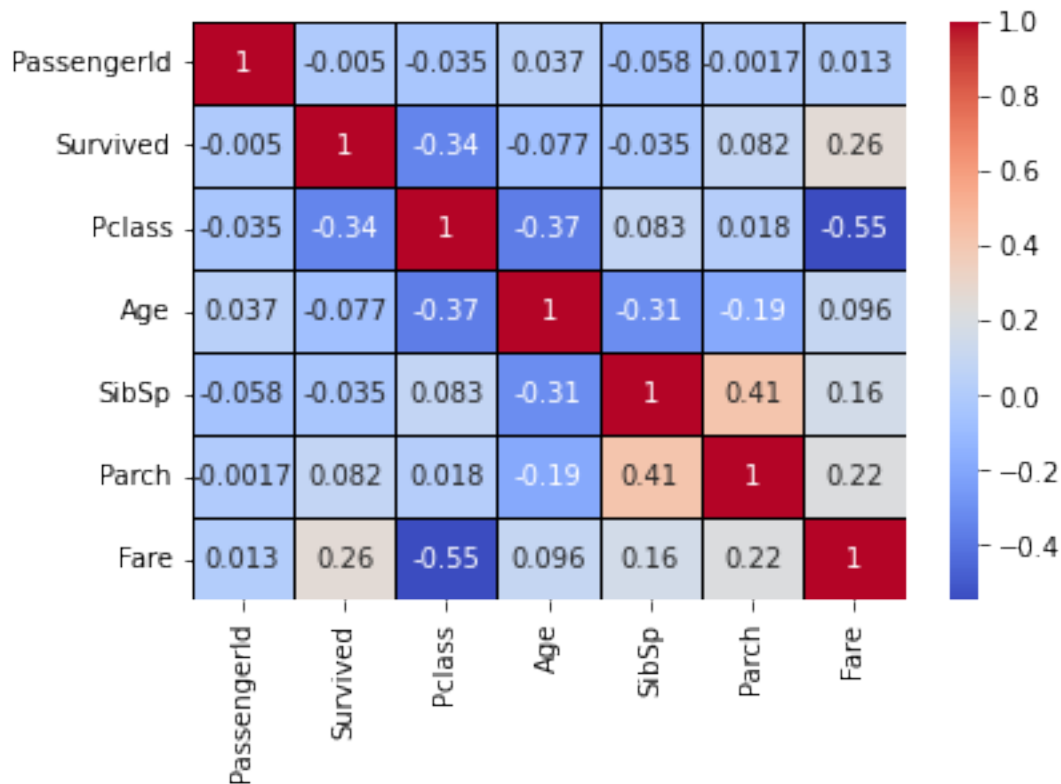
[19]: `sns.countplot(x ='Survived', data = data)`

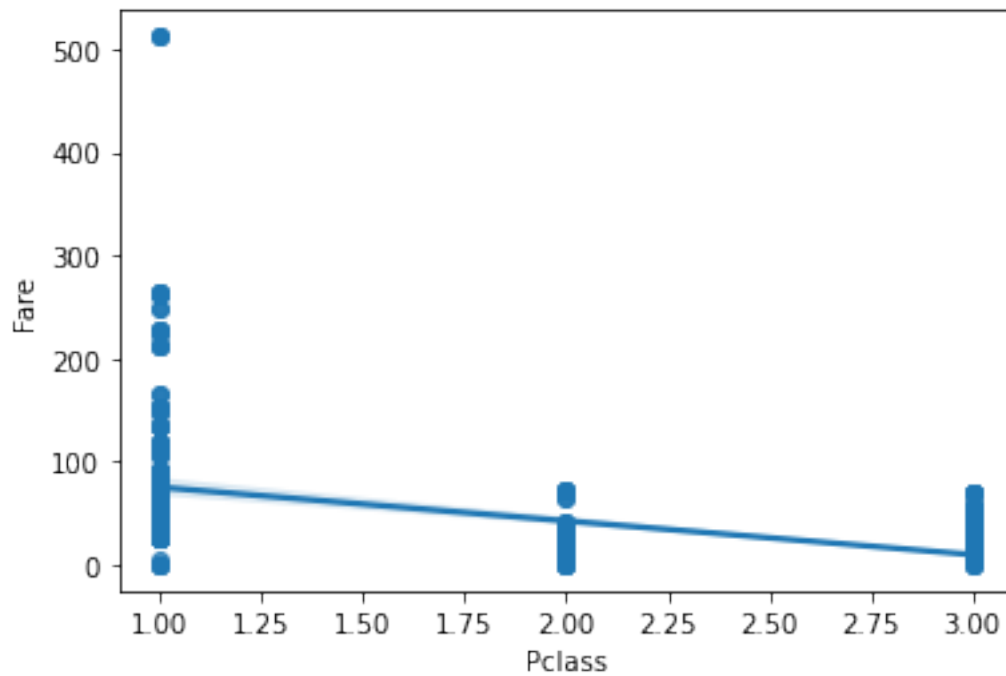[19]: <AxesSubplot:xlabel='Survived', ylabel='count'>

```
[8]: sns.heatmap(data.corr(), annot= True, cmap= 'coolwarm', linewidths = 1,␣
     ↪linecolor = 'black');
```



From the above corelation matrix, it is clear that 'Fare' and 'Survived' have a positive corelation.
Meaning higher the cost of the ticket, higher is the chance of survival.

```
[22]: sns.regplot(data=data,x='Pclass',y='Fare')
```
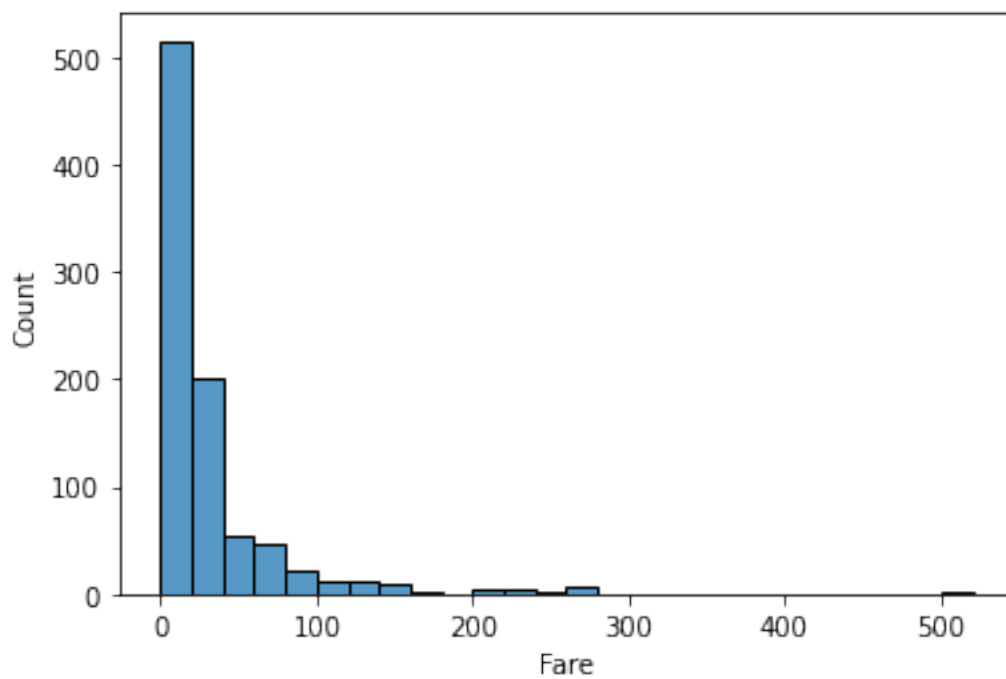
```
[22]: <AxesSubplot:xlabel='Pclass', ylabel='Fare'>
```

```
[18]: sns.histplot(data,x="Fare",bins=15,binwidth=20)
```

```
[18]: <AxesSubplot:xlabel='Fare', ylabel='Count'>
```

```
[9]: sns.histplot(data = data, x = 'Fare', hue = 'Survived',kde = True);
```