

Data Preprocessing

Import libraries

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

Import the dataset

```
In [2]: df = pd.read_csv("Titanic-Dataset.csv")  
df
```

Out[2]:	PassengerId	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare	C
0	1	0	3	Braund, Mr. Owen Harris	male	22.0	1	0	A/5 21171	7.2500	N
1	2	1	1	Cumings, Mrs. John Bradley (Florence Briggs Th...)	female	38.0	1	0	PC 17599	71.2833	C
2	3	1	3	Heikkinen, Miss. Laina	female	26.0	0	0	STON/O2. 3101282	7.9250	N
3	4	1	1	Futrelle, Mrs. Jacques Heath (Lily May Peel)	female	35.0	1	0	113803	53.1000	C
4	5	0	3	Allen, Mr. William Henry	male	35.0	0	0	373450	8.0500	N
...
886	887	0	2	Montvila, Rev. Juozas	male	27.0	0	0	211536	13.0000	N
887	888	1	1	Graham, Miss. Margaret Edith	female	19.0	0	0	112053	30.0000	I
888	889	0	3	Johnston, Miss. Catherine Helen "Carrie"	female	NaN	1	2	W./C. 6607	23.4500	N
889	890	1	1	Behr, Mr. Karl Howell	male	26.0	0	0	111369	30.0000	C
890	891	0	3	Dooley, Mr. Patrick	male	32.0	0	0	370376	7.7500	N

891 rows × 12 columns

In [3]: df.head()

Out[3]:	PassengerId	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare	Cabin
0	1	0	3	Braund, Mr. Owen Harris	male	22.0	1	0	A/5 21171	7.2500	NaN
1	2	1	1	Cumings, Mrs. John Bradley (Florence Briggs Th...)	female	38.0	1	0	PC 17599	71.2833	C85
2	3	1	3	Heikkinen, Miss. Laina	female	26.0	0	0	STON/O2. 3101282	7.9250	NaN
3	4	1	1	Futrelle, Mrs. Jacques Heath (Lily May Peel)	female	35.0	1	0	113803	53.1000	C123
4	5	0	3	Allen, Mr. William Henry	male	35.0	0	0	373450	8.0500	NaN

In [4]: `df.tail()`

Out[4]:	PassengerId	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare	Cabin
886	887	0	2	Montvila, Rev. Juozas	male	27.0	0	0	211536	13.00	NaN
887	888	1	1	Graham, Miss. Margaret Edith	female	19.0	0	0	112053	30.00	B42
888	889	0	3	Johnston, Miss. Catherine Helen "Carrie"	female	NaN	1	2	W./C. 6607	23.45	NaN
889	890	1	1	Behr, Mr. Karl Howell	male	26.0	0	0	111369	30.00	C148
890	891	0	3	Dooley, Mr. Patrick	male	32.0	0	0	370376	7.75	NaN

In [5]: `df.shape`

Out[5]: (891, 12)

In [6]: `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 [7]: `df.describe()`

	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	0.836071	14.526497	1.102743	0.806057	49.693429
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	7.910400
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

Checking the values

In [8]: `df.isnull().any()`

```
PassengerId    False
Survived      False
Pclass        False
Name          False
Sex           False
Age           True
SibSp         False
Parch         False
Ticket        False
Fare          False
Cabin         True
Embarked      True
dtype: bool
```

In [9]: `df.isnull().sum()`

```
Out[9]: 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 [10]: df["Age"].mean()
```

```
Out[10]: 29.69911764705882
```

```
In [11]: df['Age'].fillna(df['Age'].mean(), inplace=True)
```

```
In [12]: df.isnull().any()
```

```
Out[12]: PassengerId    False
          Survived      False
          Pclass        False
          Name          False
          Sex           False
          Age           False
          SibSp         False
          Parch         False
          Ticket        False
          Fare          False
          Cabin         True
          Embarked      True
          dtype: bool
```

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

```
Out[13]: 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
```

```
In [14]: df["Embarked"].mode()
```

```
Out[14]: 0    S
          Name: Embarked, dtype: object
```

```
In [15]: df['Embarked'].fillna(df['Embarked'].mode()[0], inplace=True)
```

```
In [16]: df.isnull().any()
```

```
Out[16]: PassengerId    False
          Survived      False
          Pclass        False
          Name         False
          Sex           False
          Age           False
          SibSp         False
          Parch         False
          Ticket        False
          Fare          False
          Cabin         True
          Embarked      False
          dtype: bool
```

```
In [17]: df.isnull().sum()
```

```
Out[17]: 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 [18]: df.drop(['Cabin'],axis=1,inplace= True)
```

```
In [19]: df.head(12)
```

Out[19]:	PassengerId	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare
0	1	0	3	Braund, Mr. Owen Harris	male	22.000000	1	0	A/5 21171	7.2500
1	2	1	1	Cumings, Mrs. John Bradley (Florence Briggs Th...)	female	38.000000	1	0	PC 17599	71.2833
2	3	1	3	Heikkinen, Miss. Laina	female	26.000000	0	0	STON/O2. 3101282	7.9250
3	4	1	1	Futrelle, Mrs. Jacques Heath (Lily May Peel)	female	35.000000	1	0	113803	53.1000
4	5	0	3	Allen, Mr. William Henry	male	35.000000	0	0	373450	8.0500
5	6	0	3	Moran, Mr. James	male	29.699118	0	0	330877	8.4583
6	7	0	1	McCarthy, Mr. Timothy J	male	54.000000	0	0	17463	51.8625
7	8	0	3	Palsson, Master. Gosta Leonard	male	2.000000	3	1	349909	21.0750
8	9	1	3	Johnson, Mrs. Oscar W (Elisabeth Vilhelmina Berg)	female	27.000000	0	2	347742	11.1333
9	10	1	2	Nasser, Mrs. Nicholas (Adele Achem)	female	14.000000	1	0	237736	30.0708
10	11	1	3	Sandstrom, Miss. Marguerite Rut	female	4.000000	1	1	PP 9549	16.7000
11	12	1	1	Bonnell, Miss. Elizabeth	female	58.000000	0	0	113783	26.5500



Data Visualization

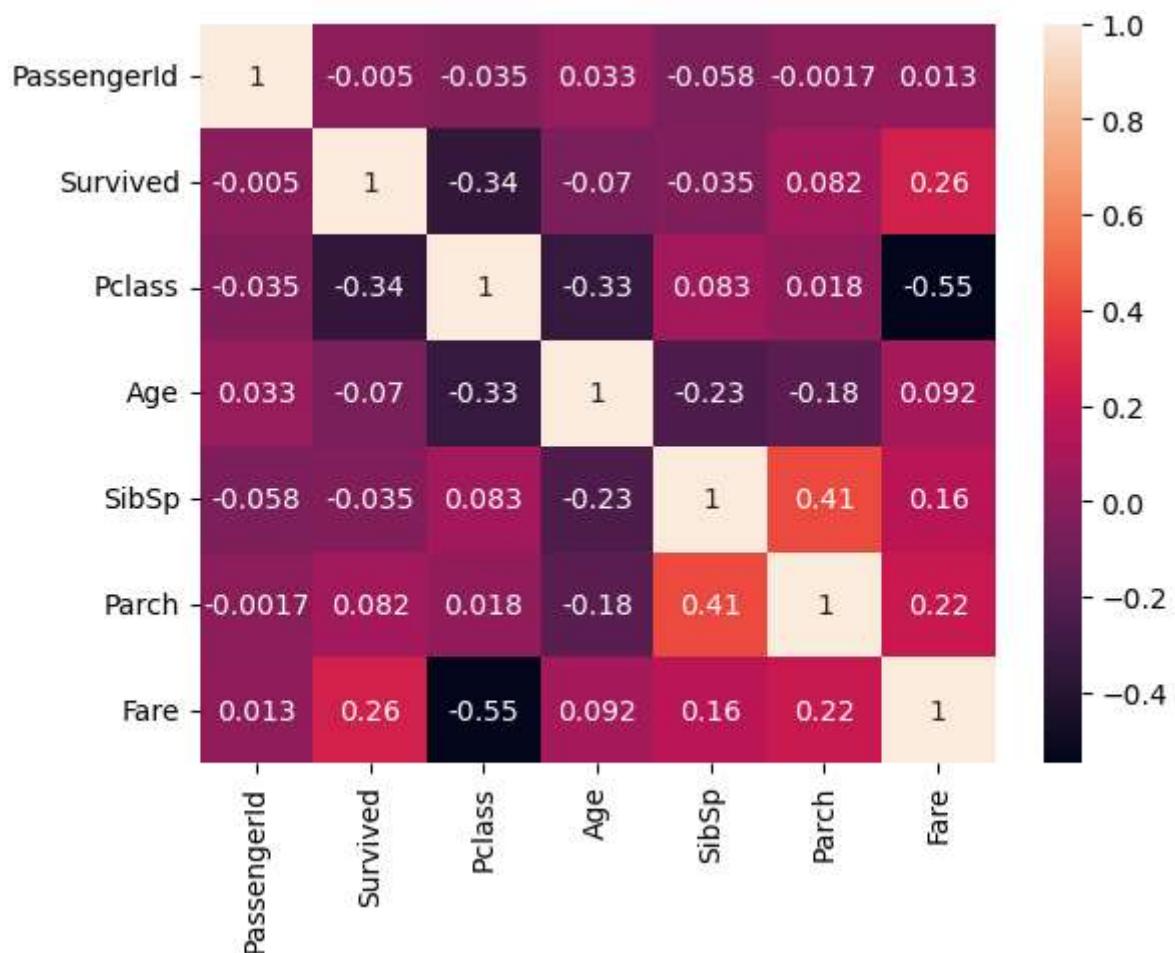
In [20]: corr = df.corr()
corr

Out[20]:

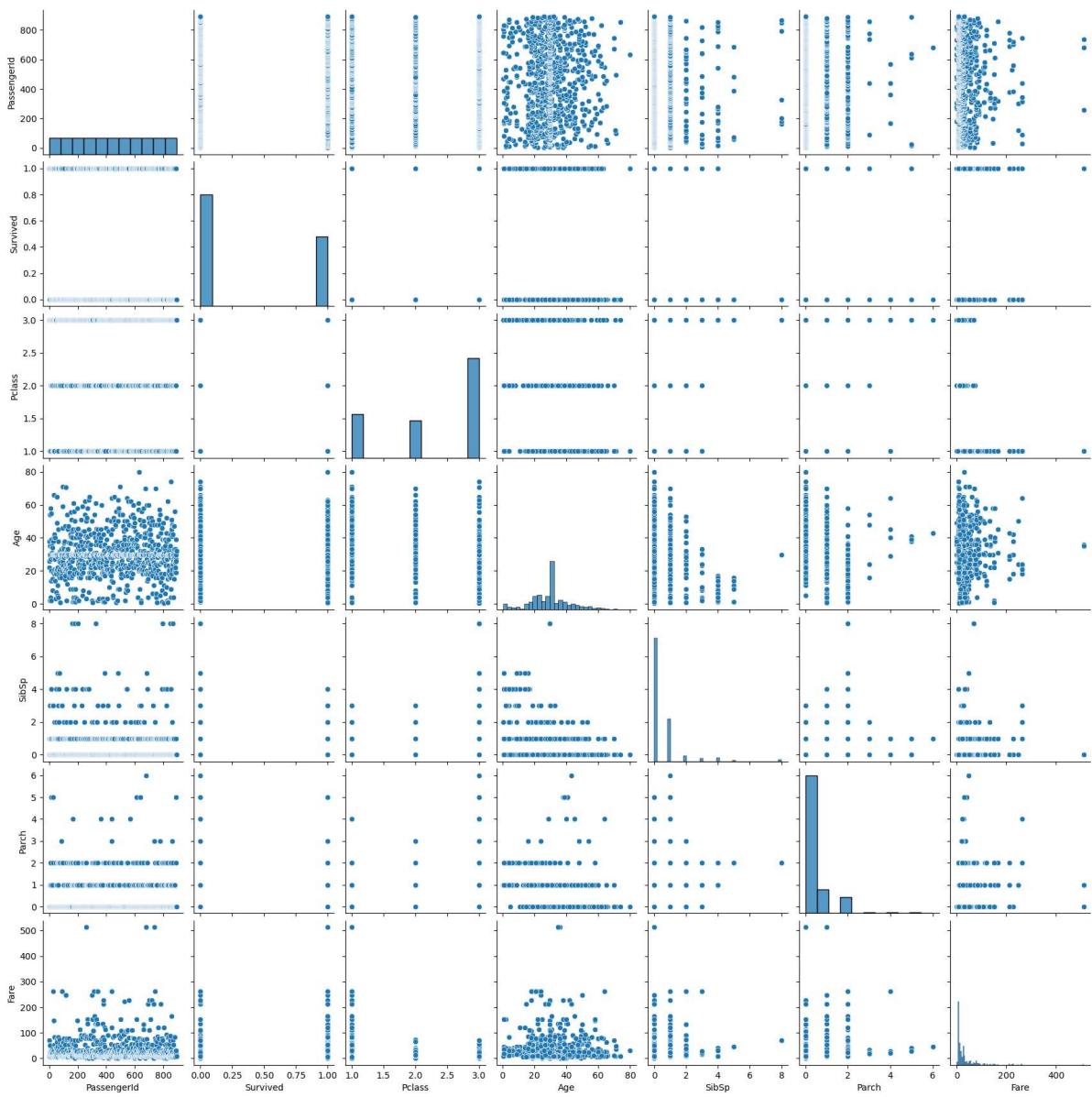
	PassengerId	Survived	Pclass	Age	SibSp	Parch	Fare
PassengerId	1.000000	-0.005007	-0.035144	0.033207	-0.057527	-0.001652	0.012658
Survived	-0.005007	1.000000	-0.338481	-0.069809	-0.035322	0.081629	0.257307
Pclass	-0.035144	-0.338481	1.000000	-0.331339	0.083081	0.018443	-0.549500
Age	0.033207	-0.069809	-0.331339	1.000000	-0.232625	-0.179191	0.091566
SibSp	-0.057527	-0.035322	0.083081	-0.232625	1.000000	0.414838	0.159651
Parch	-0.001652	0.081629	0.018443	-0.179191	0.414838	1.000000	0.216225
Fare	0.012658	0.257307	-0.549500	0.091566	0.159651	0.216225	1.000000

In [21]: `sns.heatmap(corr, annot = True)`

Out[21]: <AxesSubplot:>

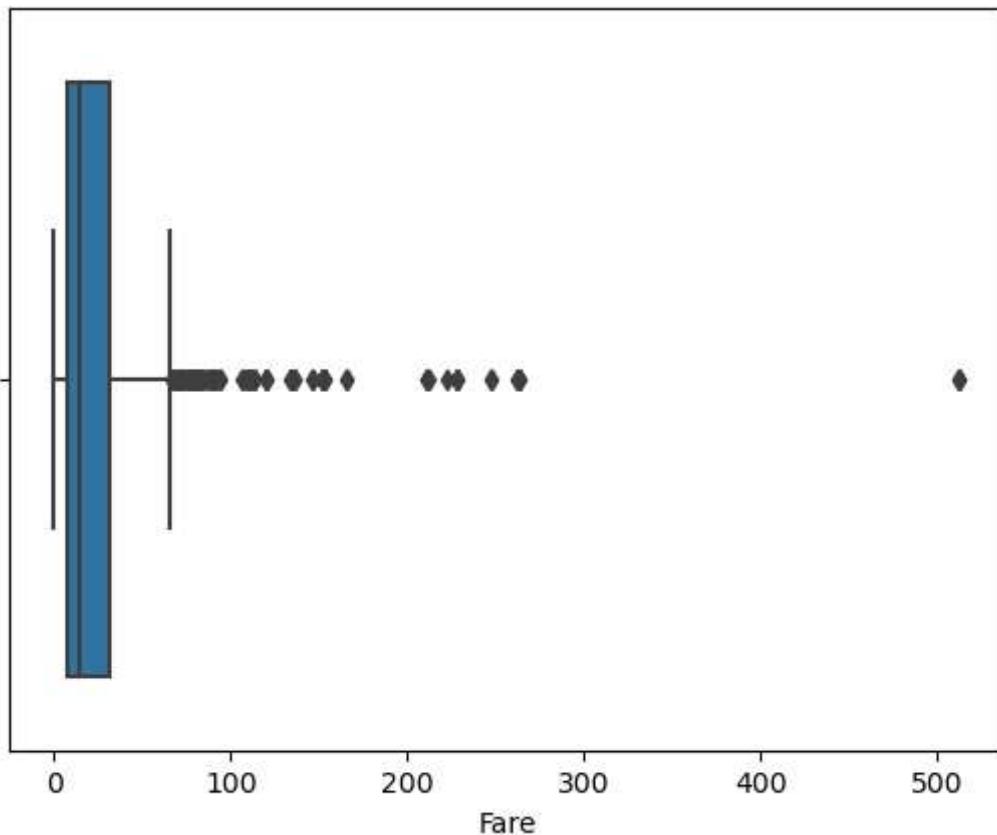
In [22]: `sns.pairplot(df)`

Out[22]: <seaborn.axisgrid.PairGrid at 0x22ae0f31f0>



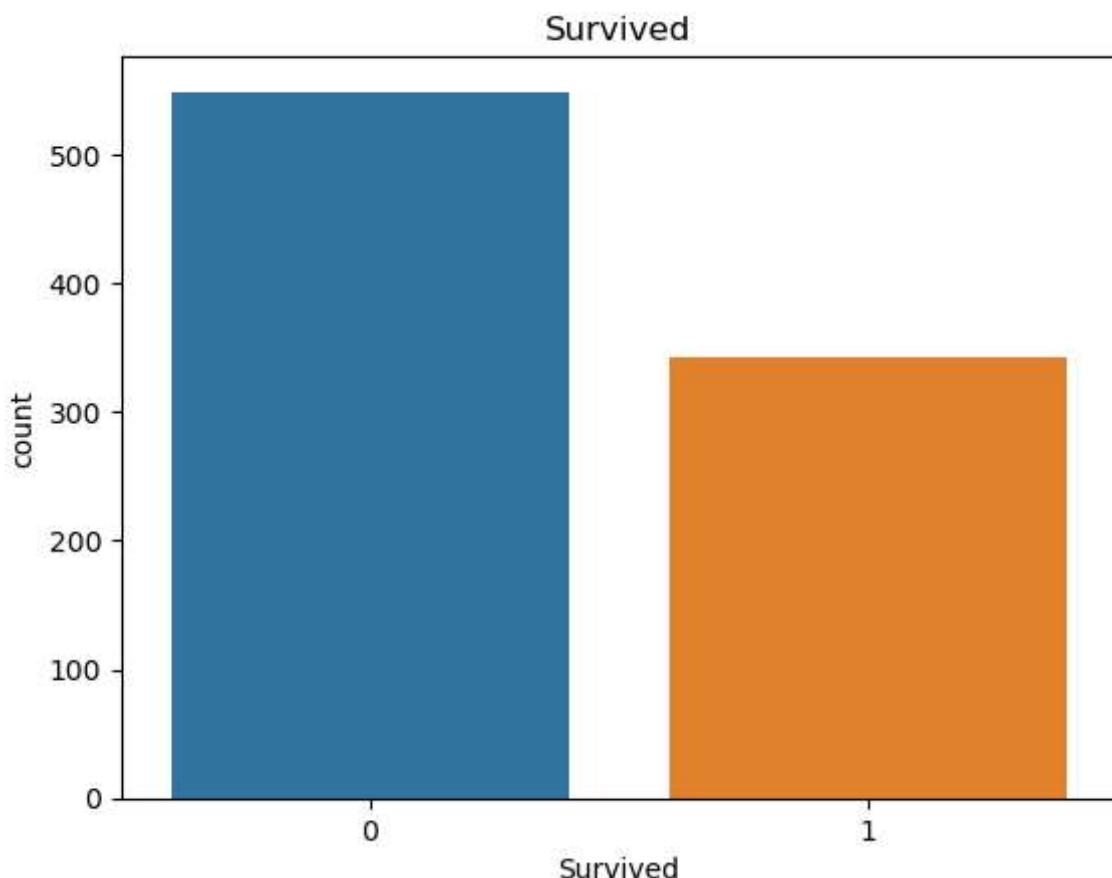
```
In [23]: sns.boxplot(x='Fare', data = df)
```

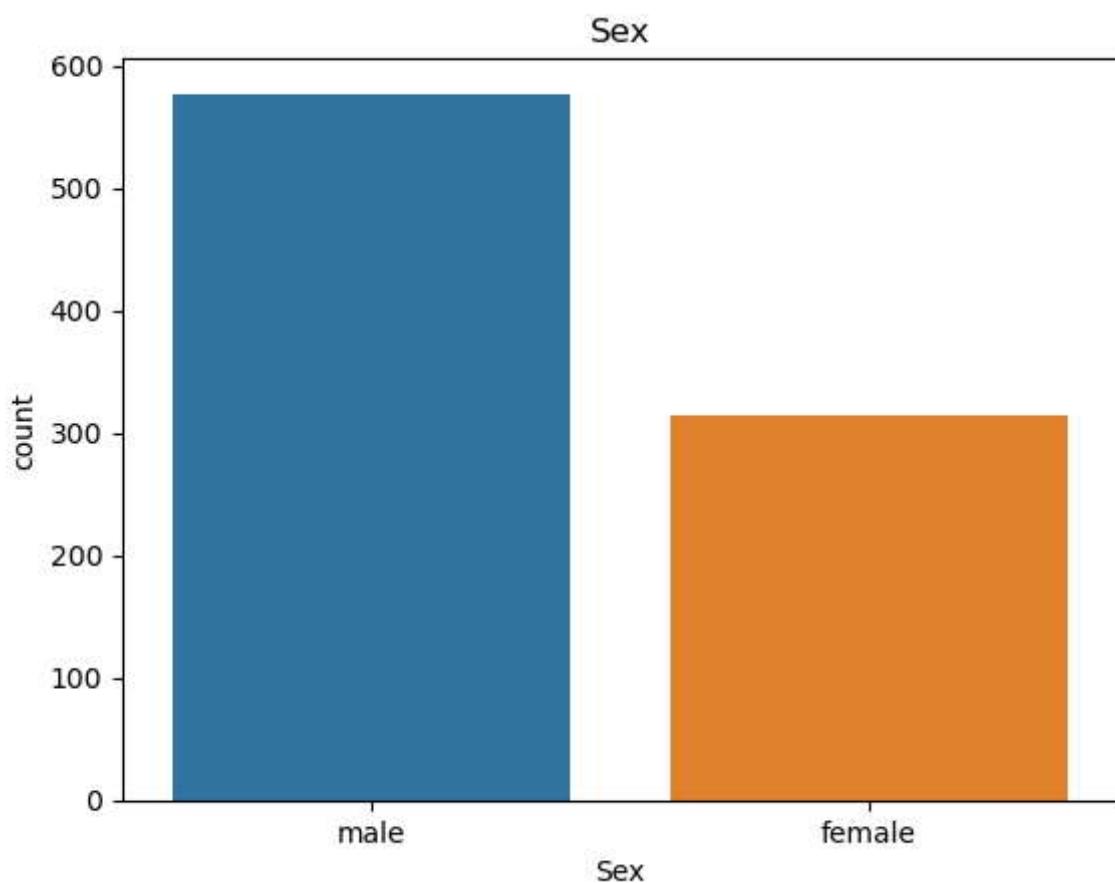
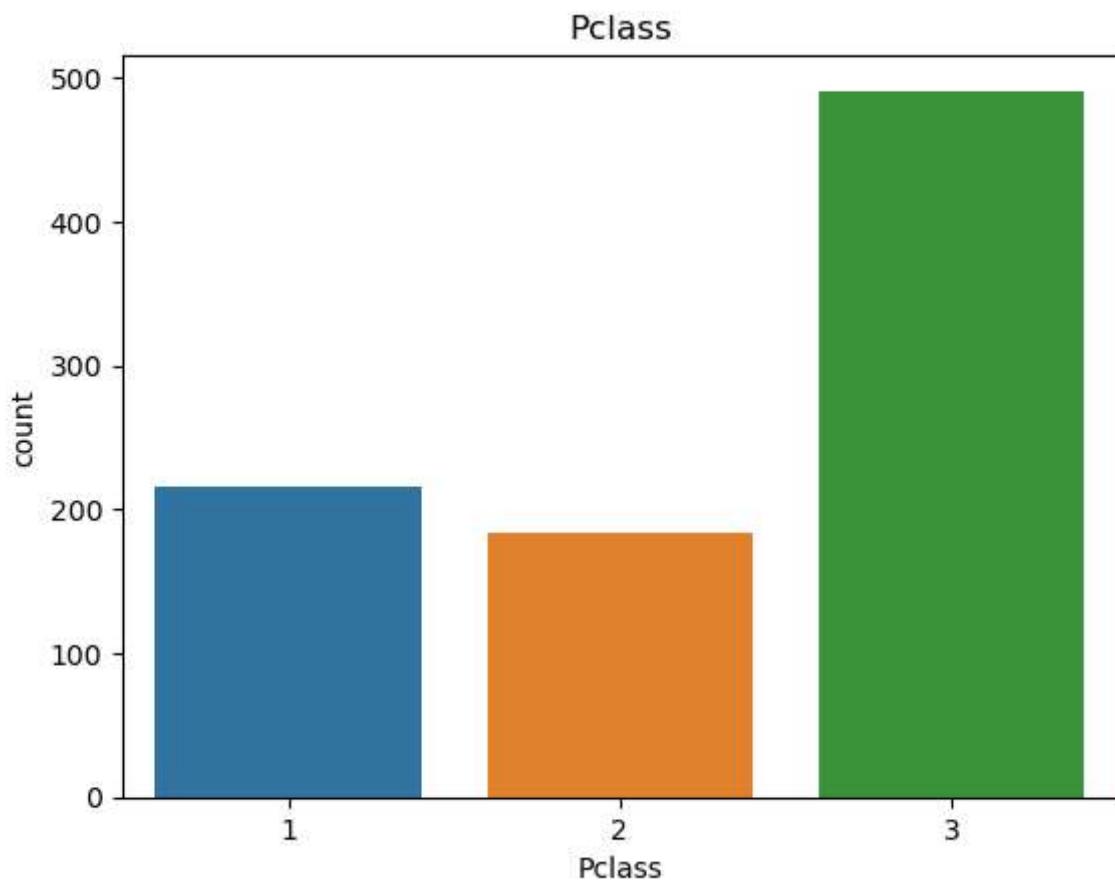
```
Out[23]: <AxesSubplot:xlabel='Fare'>
```

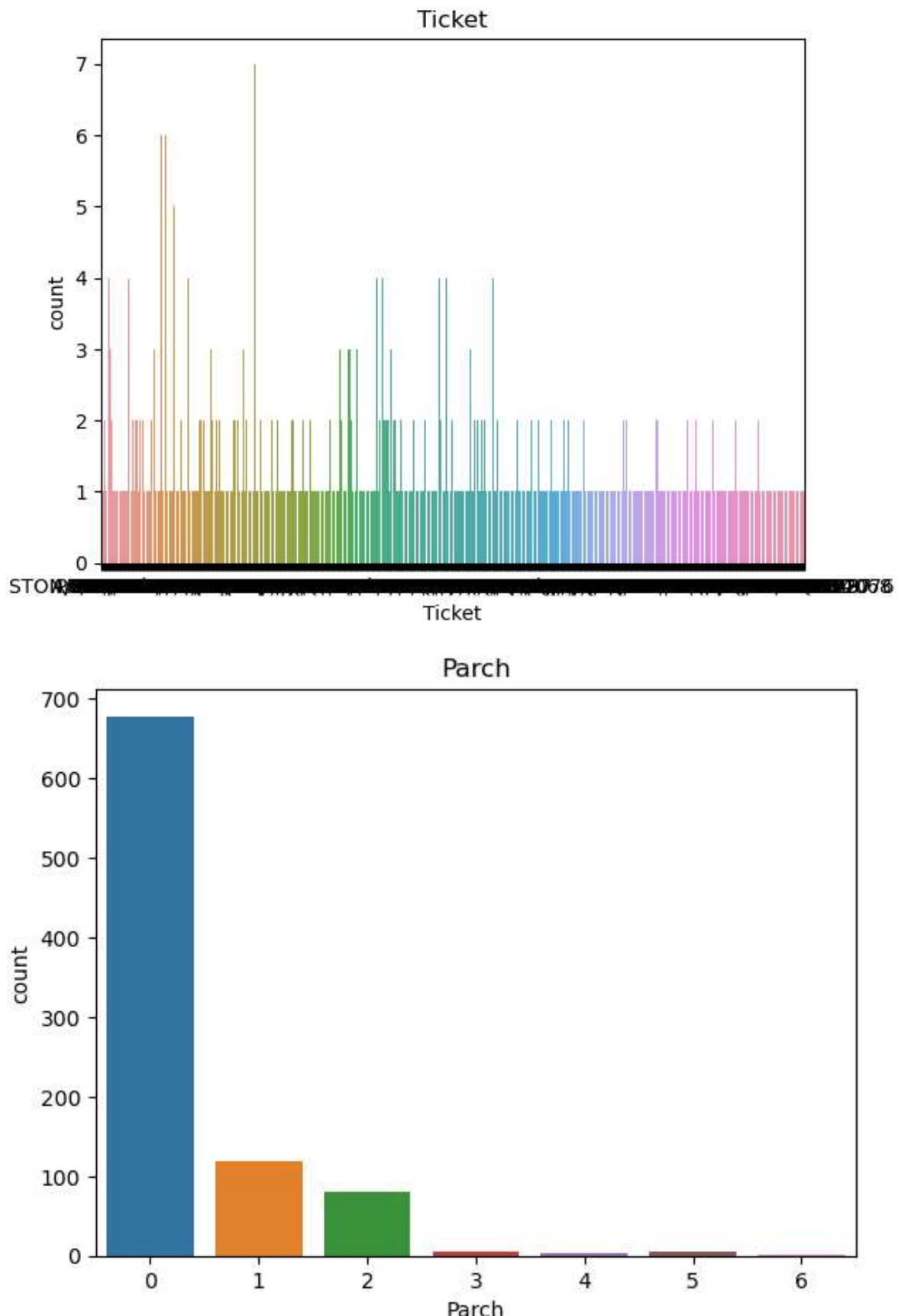


```
In [24]: df1 =df[['Survived','Pclass','Sex','Ticket','Parch']]
```

```
In [25]: for i in df1.columns:  
    sns.countplot(x=i,data=df1)  
    plt.title(i)  
    plt.show()
```







Outliners Detection

In [26]: `df.head()`

Out[26]:	PassengerId	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare	Emba
0	1	0	3	Braund, Mr. Owen Harris	male	22.0	1	0	A/5 21171	7.2500	
1	2	1	1	Cumings, Mrs. John Bradley (Florence Briggs Th... Heikkinen, Miss. Laina	female	38.0	1	0	PC 17599	71.2833	
2	3	1	3	Futrelle, Mrs. Jacques Heath (Lily May Peel)	female	26.0	0	0	STON/O2. 3101282	7.9250	
3	4	1	1	Allen, Mr. William Henry	male	35.0	1	0	113803	53.1000	
4	5	0	3				0	0	373450	8.0500	

In [27]: `from scipy import stats
z_scores = np.abs(stats.zscore(df["Age"]))`

In [28]: `outliers = df["Age"][z_scores>3]
outliers`

Out[28]:

96	71.0
116	70.5
493	71.0
630	80.0
672	70.0
745	70.0
851	74.0

Name: Age, dtype: float64

In [29]: `z_score = np.abs(stats.zscore(df["Fare"]))
outlier = df["Fare"][z_score>3]
outlier`

```
Out[29]: 27    263.0000
          88    263.0000
         118    247.5208
         258    512.3292
         299    247.5208
         311    262.3750
         341    263.0000
         377    211.5000
         380    227.5250
         438    263.0000
         527    221.7792
         557    227.5250
         679    512.3292
         689    211.3375
         700    227.5250
         716    227.5250
         730    211.3375
         737    512.3292
         742    262.3750
         779    211.3375
Name: Fare, dtype: float64
```

```
In [30]: q1 = df["Fare"].quantile(0.25)
q3 = df["Fare"].quantile(0.75)
IQR = q3-q1
lower_bound = q1-1.5*IQR
upper_bound = q3+1.5*IQR
df_cleaned = df[(df["Fare"] > lower_bound)&(df["Fare"]<upper_bound)]
print(f"Original Dataframe size:{df.shape}")
print(f"Cleaned Dataframe size:{df_cleaned.shape}")
df_cleaned
```

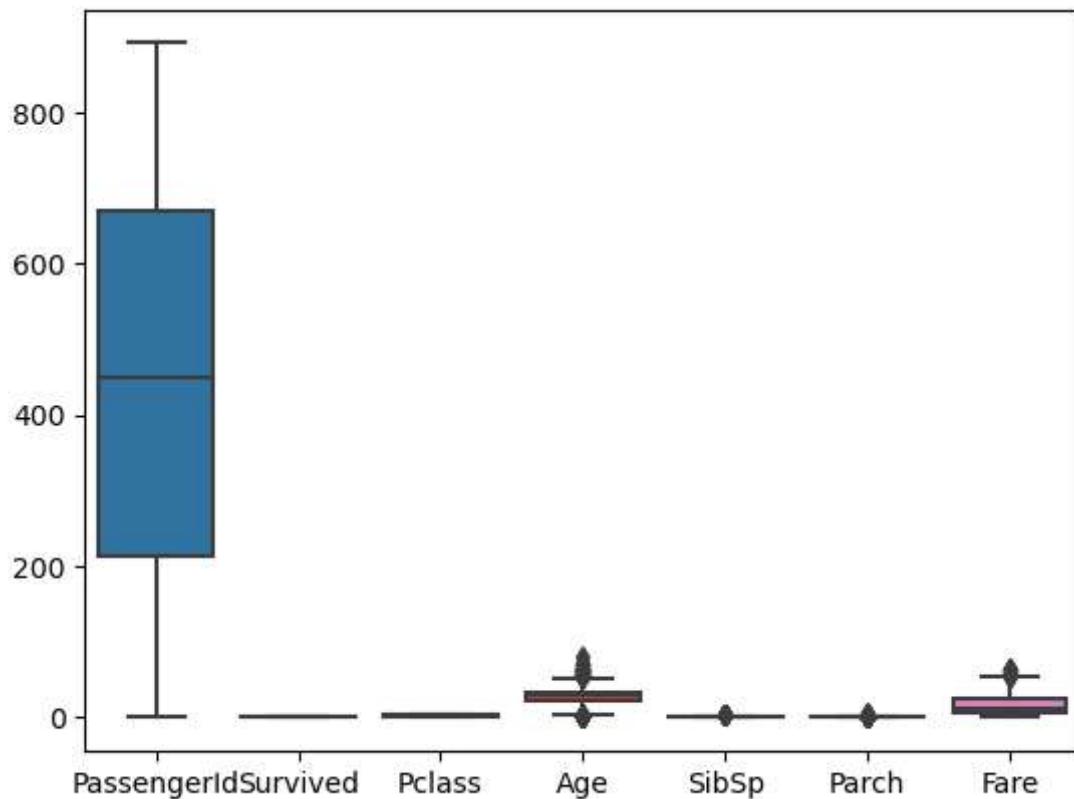
Original Dataframe size:(891, 11)
Cleaned Dataframe size:(775, 11)

Out[30]:

	PassengerId	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare
0	1	0	3	Braund, Mr. Owen Harris	male	22.000000	1	0	A/5 21171	7.2500
2	3	1	3	Heikkinen, Miss. Laina	female	26.000000	0	0	STON/O2. 3101282	7.9250
3	4	1	1	Futrelle, Mrs. Jacques Heath (Lily May Peel)	female	35.000000	1	0	113803	53.1000
4	5	0	3	Allen, Mr. William Henry	male	35.000000	0	0	373450	8.0500
5	6	0	3	Moran, Mr. James	male	29.699118	0	0	330877	8.4583
...
886	887	0	2	Montvila, Rev. Juozas	male	27.000000	0	0	211536	13.0000
887	888	1	1	Graham, Miss. Margaret Edith	female	19.000000	0	0	112053	30.0000
888	889	0	3	Johnston, Miss. Catherine Helen "Carrie"	female	29.699118	1	2	W./C. 6607	23.4500
889	890	1	1	Behr, Mr. Karl Howell	male	26.000000	0	0	111369	30.0000
890	891	0	3	Dooley, Mr. Patrick	male	32.000000	0	0	370376	7.7500

775 rows × 11 columns

In [31]: `sns.boxplot(data=df_cleaned)`Out[31]: `<AxesSubplot:>`



Splitting Dependent and Independent Variables

In [32]: `df.head()`

	PassengerId	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare	Emba
0	1	0	3	Braund, Mr. Owen Harris	male	22.0	1	0	A/5 21171	7.2500	
1	2	1	1	Cumings, Mrs. John Bradley (Florence Briggs Th... Heikkinen, Miss. Laina	female	38.0	1	0	PC 17599 STON/O2. 3101282	71.2833 7.9250	
2	3	1	3	Futrelle, Mrs. Jacques Heath (Lily May Peel)	female	26.0	0	0	113803	53.1000	
4	5	0	3	Allen, Mr. William Henry	male	35.0	0	0	373450	8.0500	

In [33]: `df.tail()`

Out[33]:

	PassengerId	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare	Embarked
886	887	0	2	Montvila, Rev. Juozas	male	27.000000	0	0	211536	13.00	
887	888	1	1	Graham, Miss. Margaret Edith	female	19.000000	0	0	112053	30.00	
888	889	0	3	Johnston, Miss. Catherine Helen "Carrie"	female	29.699118	1	2	W./C. 6607	23.45	
889	890	1	1	Behr, Mr. Karl Howell	male	26.000000	0	0	111369	30.00	
890	891	0	3	Dooley, Mr. Patrick	male	32.000000	0	0	370376	7.75	

In [34]: `x = df.drop(columns=["Survived"], axis = 1)`

In [35]: `x.head()`

Out[35]:

	PassengerId	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare	Embarked
0	1	3	Braund, Mr. Owen Harris	male	22.0	1	0	A/5 21171	7.2500	S
1	2	1	Cumings, Mrs. John Bradley (Florence Briggs Th...)	female	38.0	1	0	PC 17599	71.2833	C
2	3	3	Heikkinen, Miss. Laina	female	26.0	0	0	STON/O2. 3101282	7.9250	S
3	4	1	Futrelle, Mrs. Jacques Heath (Lily May Peel)	female	35.0	1	0	113803	53.1000	S
4	5	3	Allen, Mr. William Henry	male	35.0	0	0	373450	8.0500	S

In [36]: `type(x)`

Out[36]: `pandas.core.frame.DataFrame`

In [37]: `x.shape`

Out[37]: `(891, 10)`

```
In [38]: y = df["Survived"]
```

```
In [39]: y.head()
```

```
Out[39]: 0    0
1    1
2    1
3    1
4    0
Name: Survived, dtype: int64
```

```
In [40]: type(y)
```

```
Out[40]: pandas.core.series.Series
```

```
In [41]: y.shape
```

```
Out[41]: (891,)
```

Encoding

```
In [42]: x.head()
```

	PassengerId	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare	Embarked
0	1	3	Braund, Mr. Owen Harris	male	22.0	1	0	A/5 21171	7.2500	S
1	2	1	Cumings, Mrs. John Bradley (Florence Briggs Th... Heikkinen, Miss. Laina	female	38.0	1	0	PC 17599 STON/O2. 3101282	71.2833 7.9250	C S
2	3	3	Futrelle, Mrs. Jacques Heath (Lily May Peel)	female	26.0	0	0	113803	53.1000	S
3	4	1	Allen, Mr. William Henry	male	35.0	0	0	373450	8.0500	S
4	5	3								

```
In [43]: from sklearn.preprocessing import LabelEncoder
```

```
In [44]: le = LabelEncoder()
```

```
In [45]: x["Embarked"] = le.fit_transform(x["Embarked"])
```

```
In [46]: x.head()
```

Out[46]:

	PassengerId	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare	Embarked
0	1	3	Braund, Mr. Owen Harris	male	22.0	1	0	A/5 21171	7.2500	2
1	2	1	Cumings, Mrs. John Bradley (Florence Briggs Th... Heikkinen, Miss. Laina	female	38.0	1	0	PC 17599 3101282	71.2833 7.9250	0 2
2	3	3	Futrelle, Mrs. Jacques Heath (Lily May Peel) Allen, Mr. William Henry	female	26.0	0	0	113803 373450	53.1000 8.0500	2 2
3	4	1	Allen, Mr. William Henry	male	35.0	1	0	373450	8.0500	2
4	5	3								

In [47]: `print(le.classes_)`

```
['C' 'Q' 'S']
```

In [48]: `x["Sex"] = le.fit_transform(x["Sex"])`

In [49]: `x.head()`

Out[49]:

	PassengerId	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare	Embarked
0	1	3	Braund, Mr. Owen Harris	1	22.0	1	0	A/5 21171	7.2500	2
1	2	1	Cumings, Mrs. John Bradley (Florence Briggs Th... Heikkinen, Miss. Laina	0	38.0	1	0	PC 17599 3101282	71.2833 7.9250	0 2
2	3	3	Futrelle, Mrs. Jacques Heath (Lily May Peel) Allen, Mr. William Henry	0	26.0	0	0	113803 373450	53.1000 8.0500	2 2
3	4	1	Allen, Mr. William Henry	1	35.0	0	0			
4	5	3								

Splitting data into train and test dataset

In [50]: `from sklearn.model_selection import train_test_split`

In [51]: `x = df.drop('Survived', axis=1)
y = df['Survived']
x_train, x_test, y_train, y_test = train_test_split(x, y, test_size=0.2, random_state=42)`

```
In [52]: x_train.shape,x_test.shape,y_train.shape,y_test.shape
```

```
Out[52]: ((712, 10), (179, 10), (712,), (179,))
```

Feature Scaling

```
In [53]: from sklearn.preprocessing import StandardScaler  
sc = StandardScaler()
```

```
In [54]: x[['Age', 'Fare']] = sc.fit_transform(x[['Age', 'Fare']])
```

```
In [55]: x.head()
```

	PassengerId	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare	Embarke
0	1	3	Braund, Mr. Owen Harris	male	-0.592481	1	0	A/5 21171	-0.502445	
1	2	1	Cumings, Mrs. John Bradley (Florence Briggs Th...	female	0.638789	1	0	PC 17599	0.786845	
2	3	3	Heikkinen, Miss. Laina	female	-0.284663	0	0	STON/O2. 3101282	-0.488854	
3	4	1	Futrelle, Mrs. Jacques Heath (Lily May Peel)	female	0.407926	1	0	113803	0.420730	
4	5	3	Allen, Mr. William Henry	male	0.407926	0	0	373450	-0.486337	

```
In [ ]:
```

```
In [ ]:
```