

In [1]:

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

In [2]:

```
df=pd.read_csv("https://raw.githubusercontent.com/dsrscientist/dataset1/master/titanic_train.csv")
df
```

Out[2]:

	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	Cumings, Mrs. John Bradley (Florence Briggs Th...	female	38.0	1	0	PC 17599	71.2833	C85	C
2	3	1	3	Heikkinen, Miss. Laina	female	26.0	0	0	STON/O2. 3101282	7.9250	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
...
886	887	0	2	Montvila, Rev. Juozas	male	27.0	0	0	211536	13.0000	NaN	S
887	888	1	1	Graham, Miss. Margaret Edith	female	19.0	0	0	112053	30.0000	B42	S
888	889	0	3	Johnston, Miss. Catherine Helen "Carrie"	female	NaN	1	2	W./C. 6607	23.4500	NaN	S
889	890	1	1	Behr, Mr. Karl Howell	male	26.0	0	0	111369	30.0000	C148	C

	PassengerId	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare	Cabin	Embarked
890	891	0	3	Dooley, Mr. Patrick	male	32.0	0	0	370376	7.7500	NaN	Q

891 rows × 12 columns

In [3]:

```
df.head()
```

Out[3]:

	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	Cumings, Mrs. John Bradley (Florence Briggs Th...	female	38.0	1	0	PC 17599	71.2833	C85	C
2	3	1	3	Heikkinen, Miss. Laina	female	26.0	0	0	STON/O2. 3101282	7.9250	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

In [4]:

```
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 [5]:

```
df.isnull().sum()
```

Out[5]:

```
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 [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	0.836071	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

In [7]:

```
#dropping Cabin as more than 50% Null values are available
#replacing nul value of age with mean
```

In [8]:

```
df=df.drop(columns='Cabin',axis=1)
df['Age'].fillna(df['Age'].mean(), inplace=True)
df['Embarked'].fillna(df['Embarked'].mode()[0], inplace=True)
```

In [9]:

```
df.head()
```

Out[9]:

	PassengerId	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare	Embarked
0	1	0	3	Braund, Mr. Owen Harris	male	22.0	1	0	A/5 21171	7.2500	S
1	2	1	1	Cumings, Mrs. John Bradley (Florence	female	38.0	1	0	PC 17599	71.2833	C

	PassengerId	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare	Embarked
				Briggs Th...							
2	3	1	3	Heikkinen, Miss. Laina	female	26.0	0	0	STON/O2. 3101282	7.9250	S
3	4	1	1	Futrelle, Mrs. Jacques Heath (Lily May Peel)	female	35.0	1	0	113803	53.1000	S
4	5	0	3	Allen, Mr. William Henry	male	35.0	0	0	373450	8.0500	S

In [10]:

```
categorical=[]
for i in df.dtypes.index:
    if df.dtypes[i]=="object":
        categorical.append(i)
print("Categorical columns:",categorical)
print("\n")

numerical=[]
for i in df.dtypes.index:
    if df.dtypes[i!="object":
        numerical.append(i)
print("Numerical columns:",numerical)
print("\n")
Categorical columns: ['Name', 'Sex', 'Ticket', 'Embarked']
```

```
Numerical columns: ['PassengerId', 'Survived', 'Pclass', 'Age', 'SibSp',
'Parch', 'Fare']
```

In [11]:

```
df["Survived"].unique()
```

Out[11]:

```
array([0, 1], dtype=int64)
```

In [12]:

```
df["Ticket"].nunique()
```

Out[12]:

```
681
```

In [13]:

```
df["Embarked"].nunique()
```

Out[13]:

```
3
```

In [14]:

```
#Now we can replace Survived and Embarked to numerical column
```

```
df.replace({'Sex':{'male':0,'female':1}, 'Embarked':{'S':0,'C':1,'Q':2}},
inplace=True)
```

```
df.head()
```

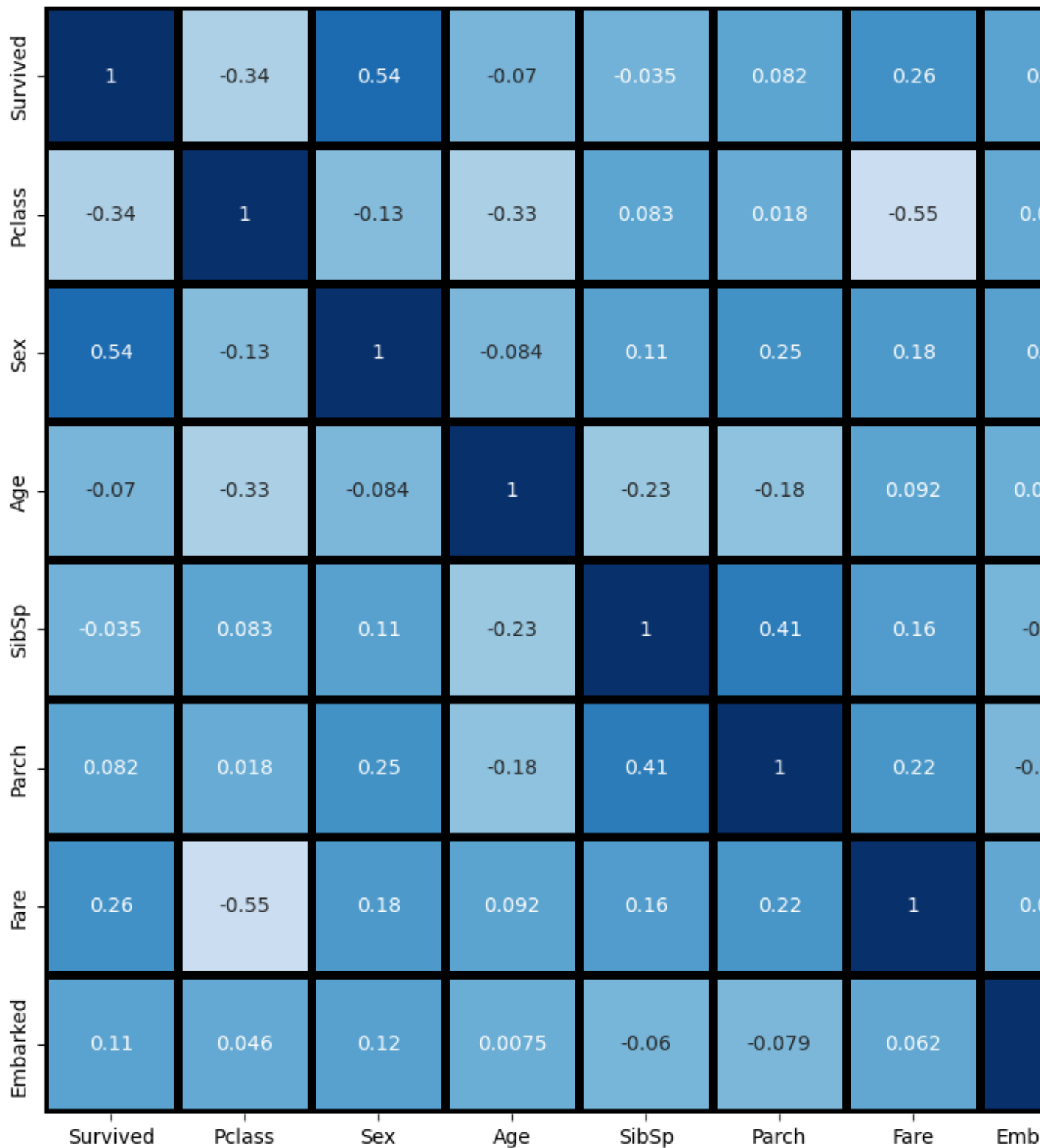
	PassengerId	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare	Embarked
0	1	0	3	Braund, Mr. Owen Harris	0	22.0	1	0	A/5 21171	7.2500	0
1	2	1	1	Cumings, Mrs. John Bradley (Florence Briggs Th...	1	38.0	1	0	PC 17599	71.2833	1
2	3	1	3	Heikkinen, Miss. Laina	1	26.0	0	0	STON/O2. 3101282	7.9250	0
3	4	1	1	Futrelle,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

```
dfnew=df.drop(['PassengerId','Name','Ticket'], axis=1)
dfnew.head()
```

	Survived	Pclass	Sex	Age	SibSp	Parch	Fare	Embarked
0	0	3	0	22.0	1	0	7.2500	0
1	1	1	1	38.0	1	0	71.2833	1
2	1	3	1	26.0	0	0	7.9250	0
3	1	1	1	35.0	1	0	53.1000	0
4	0	3	0	35.0	0	0	8.0500	0

```
plt.figure(figsize=(15,10))
sns.heatmap(dfnew.corr())
sns.heatmap(dfnew.corr(), annot = True, vmin=-1, vmax=1, center= 0, cmap=
'Blues', linewidths=3, linecolor='black')

plt.show()
```



In [19]:

```
# Splitting of data into training(80%) and testing(20%) sets
```

```
from sklearn.model_selection import train_test_split
```

In [20]:

```
#taking high correlation value
```

```
# X = features, y = target variable
```

```
X = dfnew[['Pclass', 'Sex', 'Age', 'SibSp', 'Parch', 'Fare', 'Embarked']]
y = dfnew['Survived']
```

```
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size = 0.2,
random_state = 0)
```

In [21]:

```
from sklearn.preprocessing import StandardScaler
scale=StandardScaler()
dfnew=scale.fit_transform(dfnew)
```

In [22]:

```
from sklearn.linear_model import LogisticRegression
```

In [23]:

```
lr = LogisticRegression()
lr.fit(X_train,y_train)
```

Out[23]:

```
LogisticRegression()
```

In [24]:

```
y_pred=lr.predict(X_test)
actual_vs_pred=pd.DataFrame({'Actual': y_test, 'Predicted': y_pred})
actual_vs_pred
```

Out[24]:

	Actual	Predicted
495	0	0
648	0	0
278	0	0
31	1	1
255	1	1
...
780	1	1
837	0	0
215	1	1
833	0	0
372	0	0

179 rows × 2 columns

In []:

In []:

In [36]:

```
from sklearn.metrics import accuracy_score
```

In [37]:

```
X_test_prediction=lr.predict(X_test)
test_data_accuracy=accuracy_score(y_test, X_test_prediction)
print('Accuracy score of test data : ', test_data_accuracy)
Accuracy score of test data : 0.8044692737430168
```

In []:

In []:

In []:

In []:

In []:

In []:

In []:

In []:

In []:

In []:

In []:

In []: