

```
In [2]: import pandas as pd
import numpy as np
train=pd.read_csv("F:/train.csv")
train
```

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
	890	891	0	3	Dooley, Mr. Patrick	male	32.0	0	0		370376	7.7500	NaN	Q

891 rows × 12 columns

```
In [3]: train.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]: train.shape
```

Out[4]: (891, 12)

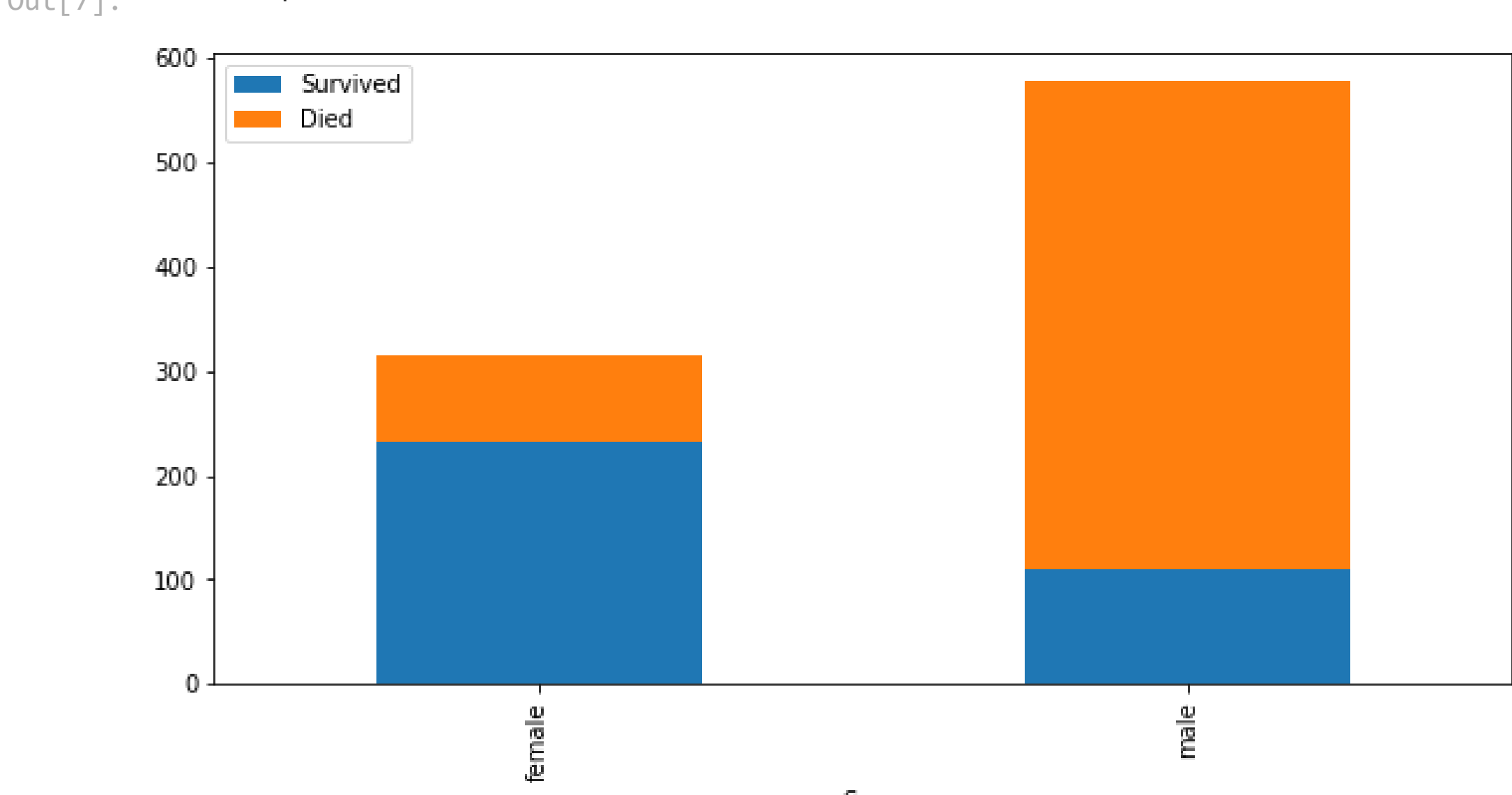
```
In [5]: train.columns
```

Out[5]: Index(['PassengerId', 'Survived', 'Pclass', 'Name', 'Sex', 'Age', 'SibSp', 'Parch', 'Ticket', 'Fare', 'Cabin', 'Embarked'], dtype='object')

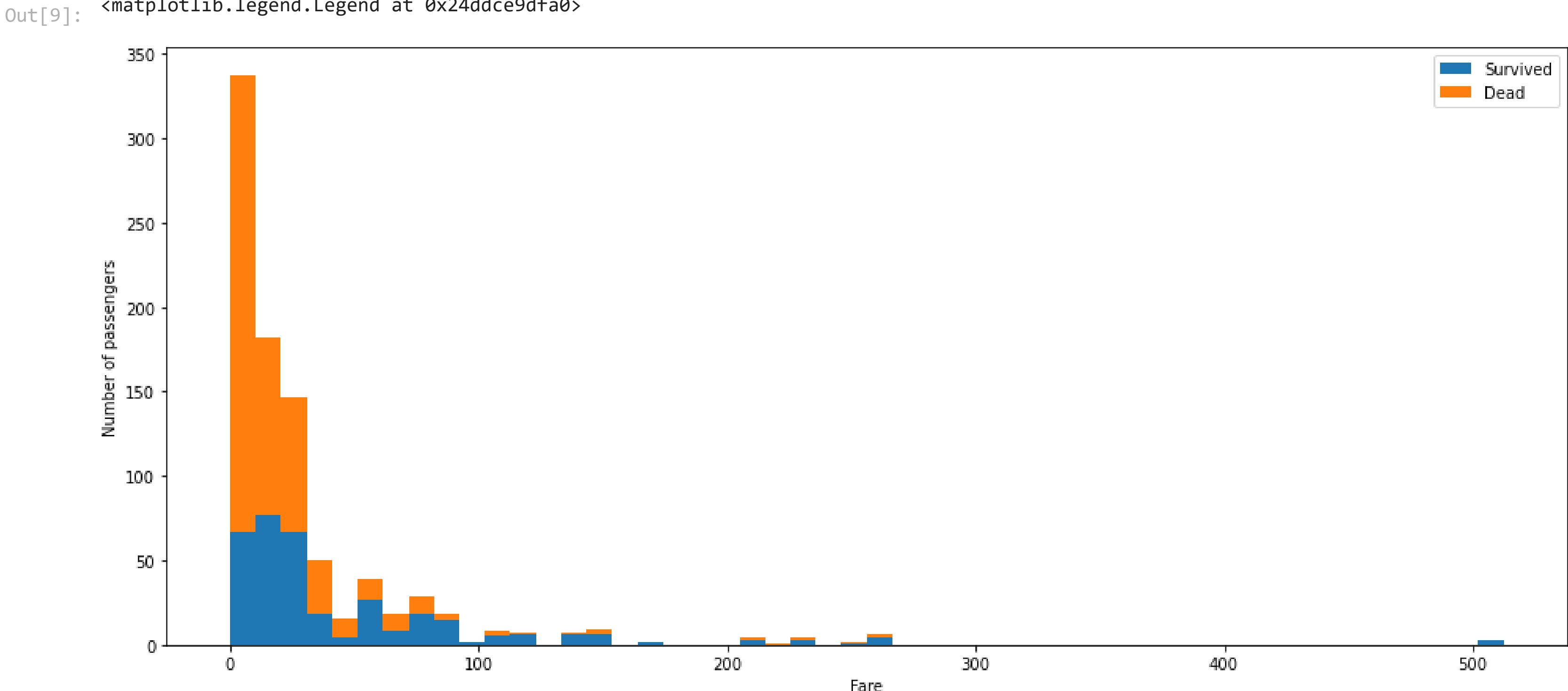
```
In [6]: train['Sex'].value_counts()
```

Out[6]: male 577  
female 314  
Name: Sex, dtype: int64

```
In [7]: #Visualizing survivals based on gender
train['Died'] = 1 - train['Survived']
train.groupby('Sex').agg('sum')['Survived', 'Died'].plot(kind='bar',
                                                         figsize=(10, 5),
                                                         stacked=True)
```



```
In [9]: ##Visualizing survivals based on fare
import matplotlib.pyplot as plt
figure = plt.figure(figsize=(16, 7))
plt.hist([train[train['Survived'] == 1]['Fare'], train[train['Survived'] == 0]['Fare']],
         stacked=True, bins = 50, label = ['Survived','Dead'])
plt.xlabel('Fare')
plt.ylabel('Number of passengers')
plt.legend()
```



```
In [10]: #Cleaning the data by removing irrelevant columns
df1=train.drop(['Name','Ticket','Cabin','PassengerId','Died'], axis=1)
df1.head(10)
```

Survived	Pclass	Sex	Age	SibSp	Parch	Fare	Embarked	
0	0	3	male	22.0	1	0	7.2500	S
1	1	1	female	38.0	1	0	71.2833	C
2	1	3	female	26.0	0	0	7.9250	S
3	1	1	female	35.0	1	0	53.1000	S
4	0	3	male	35.0	0	0	8.0500	S
5	0	3	male	NaN	0	0	8.4583	Q
6	0	1	male	54.0	0	0	51.8625	S
7	0	3	male	2.0	3	1	21.0750	S
8	1	3	female	27.0	0	2	11.1333	S
9	1	2	female	14.0	1	0	30.0708	C

```
In [11]: df1.isnull().sum()
```

Out[11]: Survived 0  
Pclass 0  
Sex 0  
Age 177  
SibSp 0  
Parch 0  
Fare 0  
Embarked 2  
dtype: int64

```
In [12]: #Converting the categorical features 'Sex' and 'Embarked' into numerical values 0 & 1
df1.Sex=df1.Sex.map({'female':0, 'male':1})
df1.Embarked=df1.Embarked.map({'S':0, 'C':1, 'Q':2, 'nan':NaN})
df1.head()
```

```
Parch      0
Fare       0
Embarked   2
dtype: int64
```

```
In [12]: #Converting the categorical features 'Sex' and 'Embarked'
df1.Sex=df1.Sex.map({'female':0, 'male':1})
df1.Embarked=df1.Embarked.map({'S':0, 'C':1, 'Q':2, 'nan':3})
df1.head()
```

```
Out[12]:
```

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

```
In [13]: #Mean age of each sex
mean_age_men=df1[df1['Sex']==1]['Age'].mean()
mean_age_women=df1[df1['Sex']==0]['Age'].mean()
```

```
In [14]: #Filling all the null values in 'Age' with respective mean age
df1.loc[(df1.Age.isnull()) & (df1['Sex']==0),'Age']=mean_age_women
df1.loc[(df1.Age.isnull()) & (df1['Sex']==1),'Age']=mean_age_men
```

```
In [15]: df1.isnull().sum()
```

Out[15]: Survived 0  
Pclass 0  
Sex 0  
Age 0  
SibSp 0  
Parch 0  
Fare 0  
Embarked 2  
dtype: int64

```
In [16]: #Since there exist 2 null values in the Embarked column, Let's drop those rows containing null values
df1.dropna(inplace=True)
```

```
In [17]: df1.isnull().sum()
```

Out[17]: Survived 0  
Pclass 0  
Sex 0  
Age 0  
SibSp 0  
Parch 0  
Fare 0  
Embarked 0  
dtype: int64

```
In [18]: #Doing Feature Scaling to standardize the independent features present in the data in a fixed range
df1.Age = (df1.Age-min(df1.Age))/(max(df1.Age)-min(df1.Age))
df1.Fare = (df1.Fare-min(df1.Fare))/(max(df1.Fare)-min(df1.Fare))
df1.describe()
```

Out[18]:		Survived	Pclass	Sex	Age	SibSp	Parch	Fare
	count	889.000000	889.000000	889.000000	889.000000	889.000000	889.000000	889.000000
	mean	0.382452	2.311586	0.649044	0.367812	0.524184	0.382452	0.062649
	std	0.486260	0.834700	0.477538	0.163124	1.103705	0.806761	0.097003
	min	0.000000	1.000000	0.000000	0.000000	0.000000	0.000000	0.000000
	25%	0.000000	2.000000	0.000000	0.271174	0.000000	0.000000	0.015412
	50%	0.000000	3.000000	1.000000	0.371701	0.000000	0.000000	0.028213
	75%	1.000000	3.000000	1.000000	0.434531	1.000000	0.000000	0.060508
	max	1.000000	3.000000	1.000000	1.000000	8.000000	6.000000	1.000000

```
In [19]: #Splitting the data for training and testing
from sklearn.model_selection import train_test_split
X_train, X_test, y_train, y_test = train_test_split(
    df1.drop(['Survived'], axis=1),
    df1.Survived,
    test_size= 0.2,
    random_state=0,
    stratify=df1.Survived)
```

```
In [20]: from sklearn.linear_model import LogisticRegression
lrmmod = LogisticRegression()
lrmmod.fit(X_train, y_train)
from sklearn.metrics import accuracy_score
y_predict = lrmmod.predict(X_test)
accuracy_score(y_test, y_predict)
```

Out[20]: 0.8426966292134831

```
In [22]: #Confusion Matrix
import seaborn as sns
from sklearn.metrics import confusion_matrix
cma=confusion_matrix(y_test, y_predict)
sns.heatmap(cma,annot=True)
```



```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```