

Titanic Dataset

15_Titanic Dataset1_22june23

```
In [124]: import pandas as pd
```

```
In [125]: data=pd.read_csv("/home/placement/Desktop/saimohan data/csv files/Titanic Dataset.csv")
```

```
In [126]: import warnings  
warnings.filterwarnings("ignore")
```

```
In [127]: data
#we find the rows X columns and total data
```

```
Out[127]:
```

	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 [128]: #finding the count,mean,std,min,max etc  
data.describe()
```

Out[128]:

	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 [129]: #finding the datatype in the titanic data
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
```

```
In [130]: #finding the only sum of the columns in the data
data.isna().sum()
```

```
Out[130]: 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 [131]: data.head(5)
#we can find the top of the data
```

```
Out[131]:
```

	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 [132]: data.tail(5)
#we can find the ending of the data
```

```
Out[132]:
```

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

```
In [133]: data['Pclass'].unique()
```

```
Out[133]: array([3, 1, 2])
```

```
In [134]: data['Age'].unique()
```

```
Out[134]: array([22. , 38. , 26. , 35. , nan, 54. , 2. , 27. , 14. ,  
 4. , 58. , 20. , 39. , 55. , 31. , 34. , 15. , 28. ,  
 8. , 19. , 40. , 66. , 42. , 21. , 18. , 3. , 7. ,  
 49. , 29. , 65. , 28.5 , 5. , 11. , 45. , 17. , 32. ,  
 16. , 25. , 0.83, 30. , 33. , 23. , 24. , 46. , 59. ,  
 71. , 37. , 47. , 14.5 , 70.5 , 32.5 , 12. , 9. , 36.5 ,  
 51. , 55.5 , 40.5 , 44. , 1. , 61. , 56. , 50. , 36. ,  
 45.5 , 20.5 , 62. , 41. , 52. , 63. , 23.5 , 0.92, 43. ,  
 60. , 10. , 64. , 13. , 48. , 0.75, 53. , 57. , 80. ,  
 70. , 24.5 , 6. , 0.67, 30.5 , 0.42, 34.5 , 74. ])
```

```
In [135]: data['SibSp'].unique()
```

```
Out[135]: array([1, 0, 3, 4, 2, 5, 8])
```

```
In [136]: #now i want to analyze the
```

```
In [137]: data1=data.drop(['PassengerId', 'Cabin', 'Name', 'Ticket', 'SibSp', 'Parch'],axis=1)
```

```
In [138]: data1
```

```
Out[138]:
```

	Survived	Pclass	Sex	Age	Fare	Embarked
0	0	3	male	22.0	7.2500	S
1	1	1	female	38.0	71.2833	C
2	1	3	female	26.0	7.9250	S
3	1	1	female	35.0	53.1000	S
4	0	3	male	35.0	8.0500	S
...
886	0	2	male	27.0	13.0000	S
887	1	1	female	19.0	30.0000	S
888	0	3	female	NaN	23.4500	S
889	1	1	male	26.0	30.0000	C
890	0	3	male	32.0	7.7500	Q

891 rows × 6 columns

```
In [139]: list(data1)
```

```
Out[139]: ['Survived', 'Pclass', 'Sex', 'Age', 'Fare', 'Embarked']
```

```
In [140]: data1['Sex']=data1['Sex'].map({'male':1,'female':0})  
data1['Pclass'].unique()
```

```
Out[140]: array([3, 1, 2])
```

```
In [141]: data1
```

```
Out[141]:
```

	Survived	Pclass	Sex	Age	Fare	Embarked
0	0	3	1	22.0	7.2500	S
1	1	1	0	38.0	71.2833	C
2	1	3	0	26.0	7.9250	S
3	1	1	0	35.0	53.1000	S
4	0	3	1	35.0	8.0500	S
...
886	0	2	1	27.0	13.0000	S
887	1	1	0	19.0	30.0000	S
888	0	3	0	NaN	23.4500	S
889	1	1	1	26.0	30.0000	C
890	0	3	1	32.0	7.7500	Q

891 rows × 6 columns

```
In [142]: data2=data1.fillna(data1.median)
```



```
In [143]: data2
```

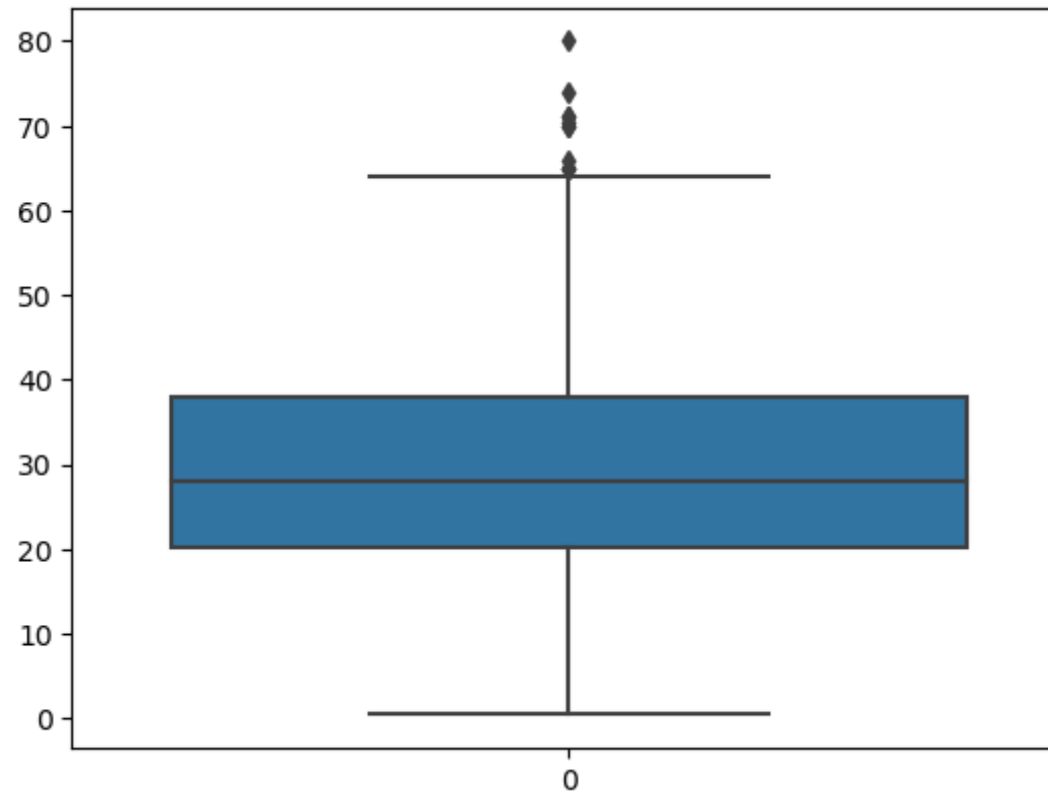
```
Out[143]:
```

	Survived	Pclass	Sex		Age	Fare	Embarked
0	0	3	1		22.0	7.2500	S
1	1	1	0		38.0	71.2833	C
2	1	3	0		26.0	7.9250	S
3	1	1	0		35.0	53.1000	S
4	0	3	1		35.0	8.0500	S
...
886	0	2	1		27.0	13.0000	S
887	1	1	0		19.0	30.0000	S
888	0	3	0	<bound method NDFrame._add_numeric_operations....		23.4500	S
889	1	1	1		26.0	30.0000	C
890	0	3	1		32.0	7.7500	Q

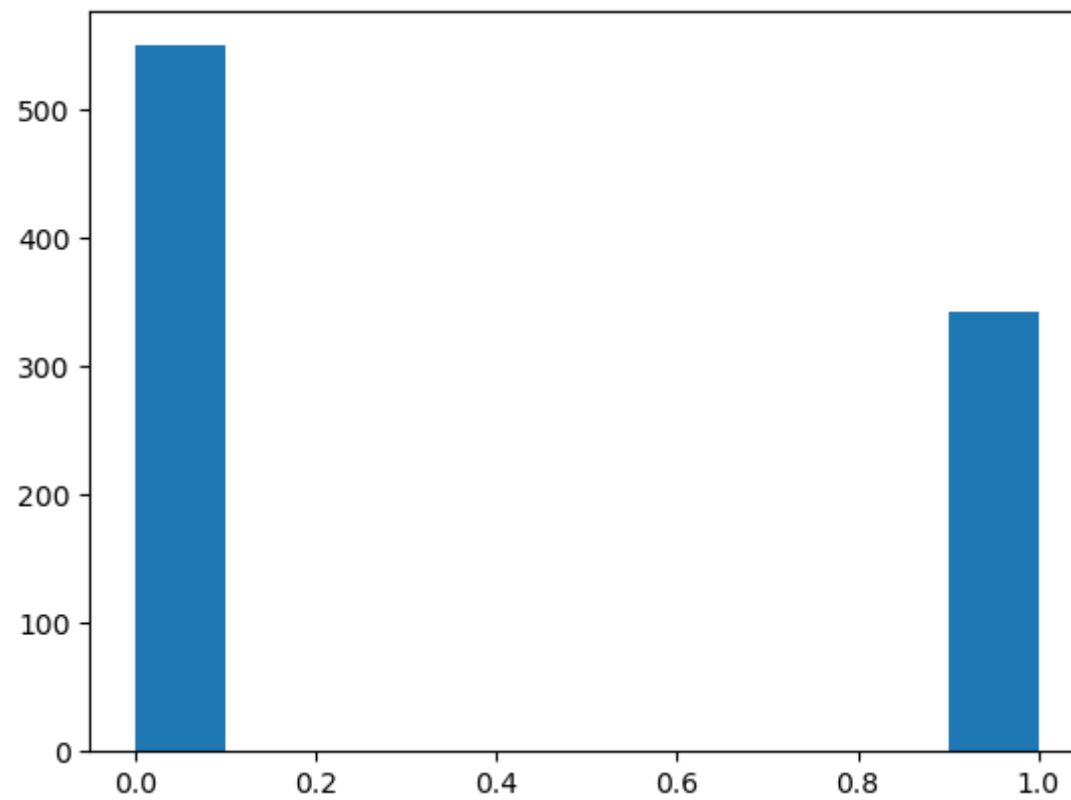
891 rows × 6 columns

```
In [144]: import seaborn as sns  
import matplotlib.pyplot as plt  
sns.boxplot(data.Age)
```

Out[144]: <Axes: >

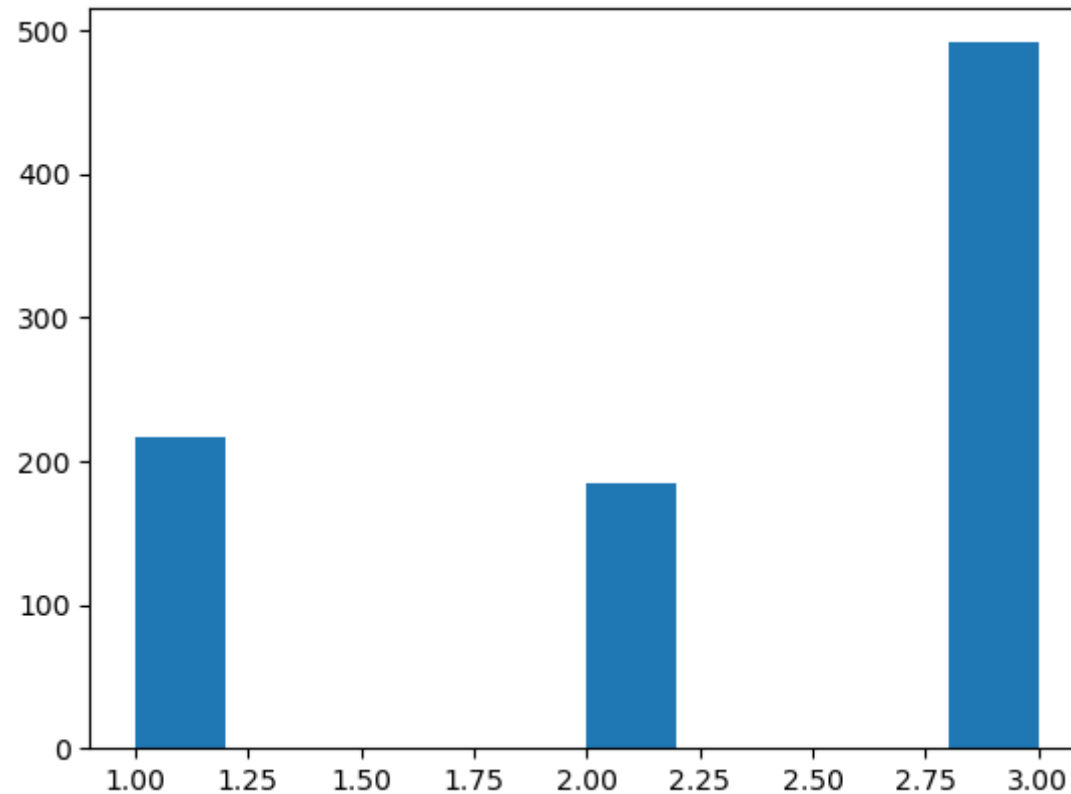


```
In [145]: plt.hist(data1['Survived'])  
plt.show()
```



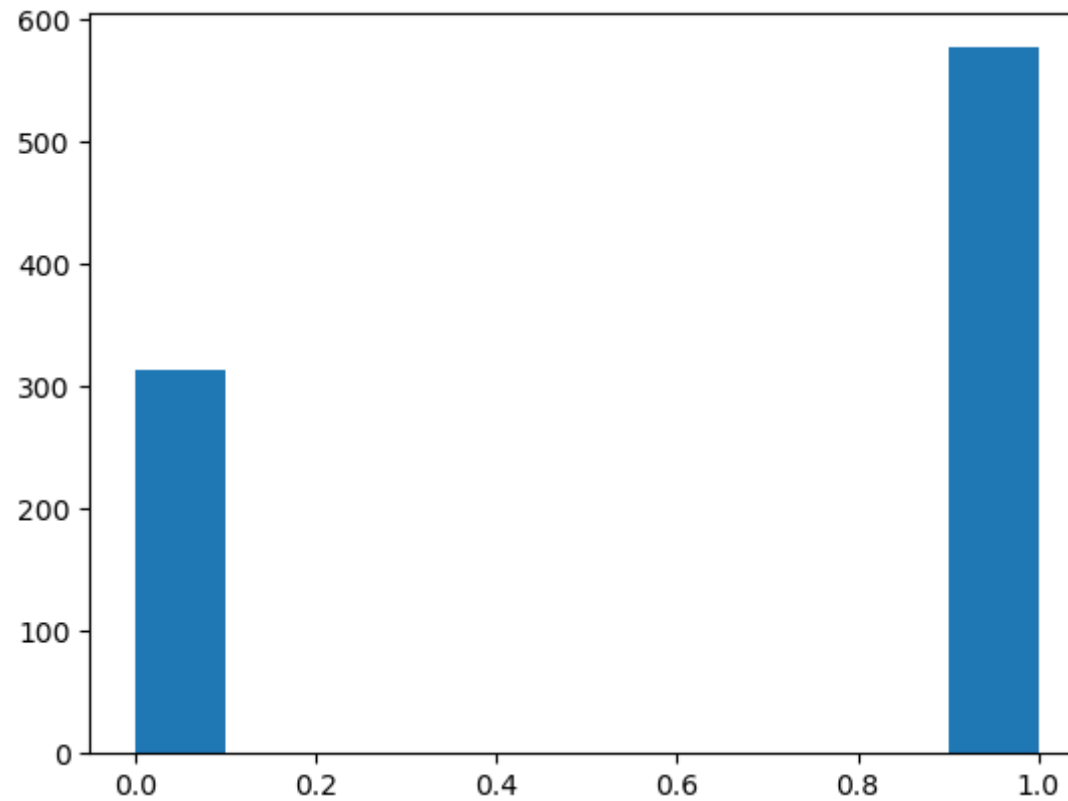
```
In [146]: plt.hist(data1['Pclass'])
```

```
Out[146]: (array([216.,  0.,  0.,  0.,  0., 184.,  0.,  0.,  0., 491.]),  
          array([1. , 1.2, 1.4, 1.6, 1.8, 2. , 2.2, 2.4, 2.6, 2.8, 3. ]),  
          <BarContainer object of 10 artists>)
```



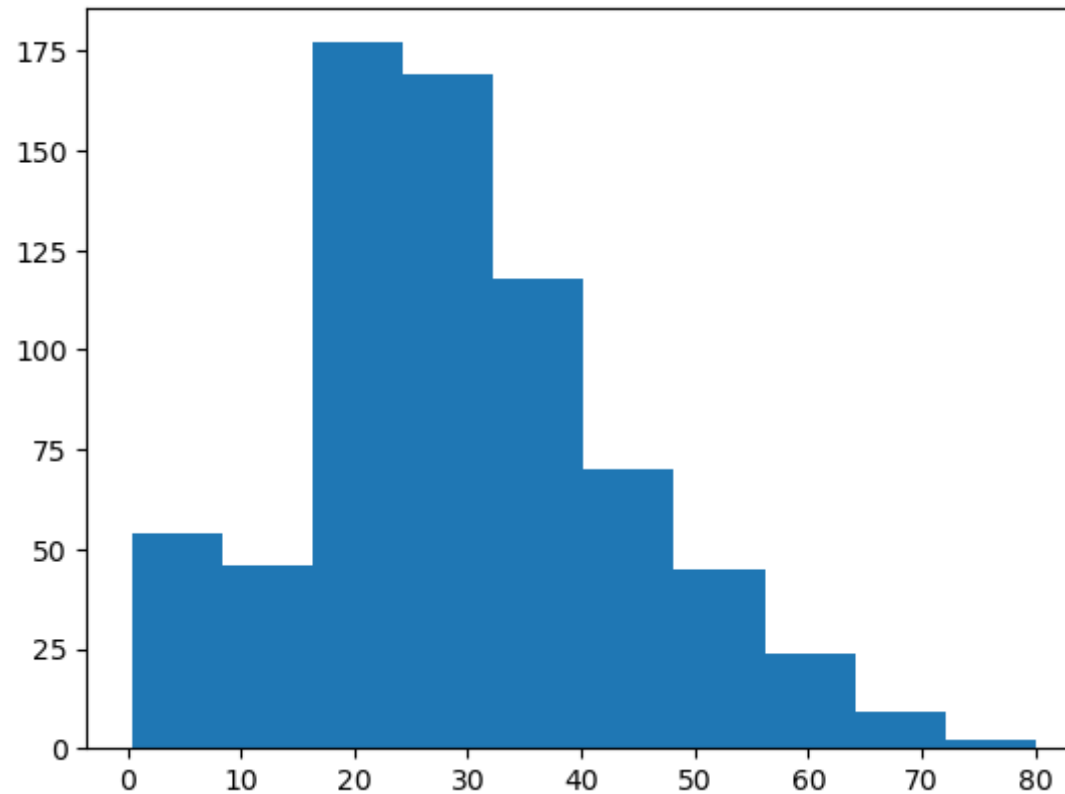
```
In [147]: plt.hist(data1['Sex'])
```

```
Out[147]: (array([314.,  0.,  0.,  0.,  0.,  0.,  0.,  0.,  0., 577.]),  
          array([0. , 0.1, 0.2, 0.3, 0.4, 0.5, 0.6, 0.7, 0.8, 0.9, 1. ]),  
          <BarContainer object of 10 artists>)
```



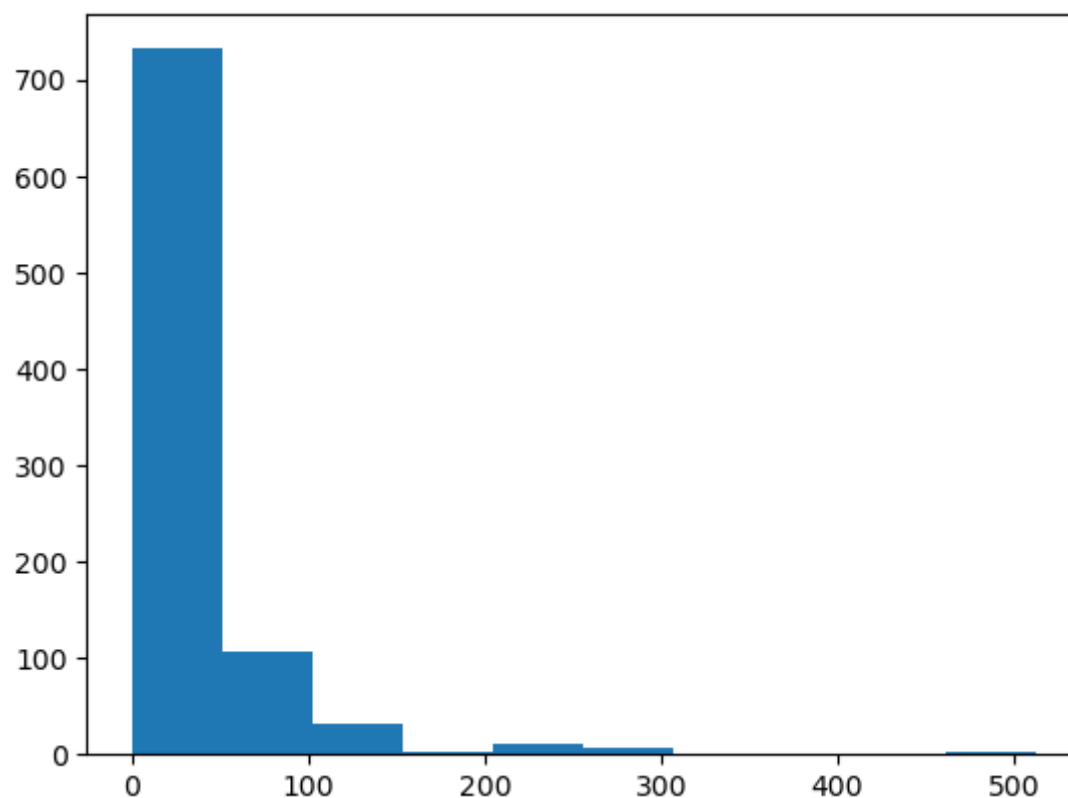
```
In [148]: plt.hist(data1['Age'])
```

```
Out[148]: (array([ 54.,  46., 177., 169., 118.,  70.,  45.,  24.,   9.,   2.]),  
array([ 0.42 ,  8.378, 16.336, 24.294, 32.252, 40.21 , 48.168, 56.126,  
        64.084, 72.042, 80.   ]),  
<BarContainer object of 10 artists>)
```



```
In [149]: plt.hist(data1['Fare'])
```

```
Out[149]: (array([732., 106., 31., 2., 11., 6., 0., 0., 0., 3.]),  
array([ 0., 51.23292, 102.46584, 153.69876, 204.93168, 256.1646 ,  
307.39752, 358.63044, 409.86336, 461.09628, 512.3292 ]),  
<BarContainer object of 10 artists>)
```



```
In [150]: # some time string and number will be same then we have to replace all the null values into 35  
#data1.fillna(35,inplace=True)
```

```
In [151]: data1.isna().sum()
```

```
Out[151]: Survived      0
          Pclass       0
          Sex          0
          Age        177
          Fare         0
          Embarked     2
          dtype: int64
```

```
In [152]: data1.describe()
```

```
Out[152]:
```

	Survived	Pclass	Sex	Age	Fare
count	891.000000	891.000000	891.000000	714.000000	891.000000
mean	0.383838	2.308642	0.647587	29.699118	32.204208
std	0.486592	0.836071	0.477990	14.526497	49.693429
min	0.000000	1.000000	0.000000	0.420000	0.000000
25%	0.000000	2.000000	0.000000	20.125000	7.910400
50%	0.000000	3.000000	1.000000	28.000000	14.454200
75%	1.000000	3.000000	1.000000	38.000000	31.000000
max	1.000000	3.000000	1.000000	80.000000	512.329200

```
In [153]: data1['Age'].unique()
          #the null values can be replaced by 35
          data1.fillna(35,inplace=True)
```


In [154]:

```
data1['Age'].unique()
```

```
Out[154]: array([22. , 38. , 26. , 35. , 54. ,  2. , 27. , 14. ,  4. ,
        58. , 20. , 39. , 55. , 31. , 34. , 15. , 28. ,  8. ,
        19. , 40. , 66. , 42. , 21. , 18. ,  3. ,  7. , 49. ,
        29. , 65. , 28.5,  5. , 11. , 45. , 17. , 32. , 16. ,
        25. ,  0.83, 30. , 33. , 23. , 24. , 46. , 59. , 71. ,
        37. , 47. , 14.5, 70.5, 32.5, 12. ,  9. , 36.5, 51. ,
        55.5, 40.5, 44. ,  1. , 61. , 56. , 50. , 36. , 45.5 ,
        20.5, 62. , 41. , 52. , 63. , 23.5,  0.92, 43. , 60. ,
        10. , 64. , 13. , 48. ,  0.75, 53. , 57. , 80. , 70. ,
        24.5 ,  6. ,  0.67, 30.5 ,  0.42, 34.5 , 74.  ])
```

In [155]: *#passenger class mapped into the numbers*

```
In [156]: data1["Pclass"]=data1["Pclass"].map({1:'F',2:'S',3:'Third'})
```

```
In [157]: data1.isna().sum()
```

```
Out[157]: Survived    0
Pclass      0
Sex         0
Age         0
Fare        0
Embarked    0
dtype: int64
```

```
In [158]: data1
```

```
Out[158]:
```

	Survived	Pclass	Sex	Age	Fare	Embarked
0	0	Third	1	22.0	7.2500	S
1	1	F	0	38.0	71.2833	C
2	1	Third	0	26.0	7.9250	S
3	1	F	0	35.0	53.1000	S
4	0	Third	1	35.0	8.0500	S
...
886	0	S	1	27.0	13.0000	S
887	1	F	0	19.0	30.0000	S
888	0	Third	0	35.0	23.4500	S
889	1	F	1	26.0	30.0000	C
890	0	Third	1	32.0	7.7500	Q

891 rows × 6 columns

```
In [159]: data1=pd.get_dummies(data1)
```

```
In [160]: data1.shape
```

```
Out[160]: (891, 11)
```

```
In [161]: data1.head(500)
```

```
Out[161]:
```

	Survived	Sex	Age	Fare	Pclass_F	Pclass_S	Pclass_Third	Embarked_35	Embarked_C	Embarked_Q	Embarked_S
0	0	1	22.0	7.2500	0	0	1	0	0	0	1
1	1	0	38.0	71.2833	1	0	0	0	1	0	0
2	1	0	26.0	7.9250	0	0	1	0	0	0	1
3	1	0	35.0	53.1000	1	0	0	0	0	0	1
4	0	1	35.0	8.0500	0	0	1	0	0	0	1
...
495	0	1	35.0	14.4583	0	0	1	0	1	0	0
496	1	0	54.0	78.2667	1	0	0	0	1	0	0
497	0	1	35.0	15.1000	0	0	1	0	0	0	1
498	0	0	25.0	151.5500	1	0	0	0	0	0	1
499	0	1	24.0	7.7958	0	0	1	0	0	0	1

500 rows × 11 columns

```
In [162]: # finding correletion
cor=data1.corr()
cor
```

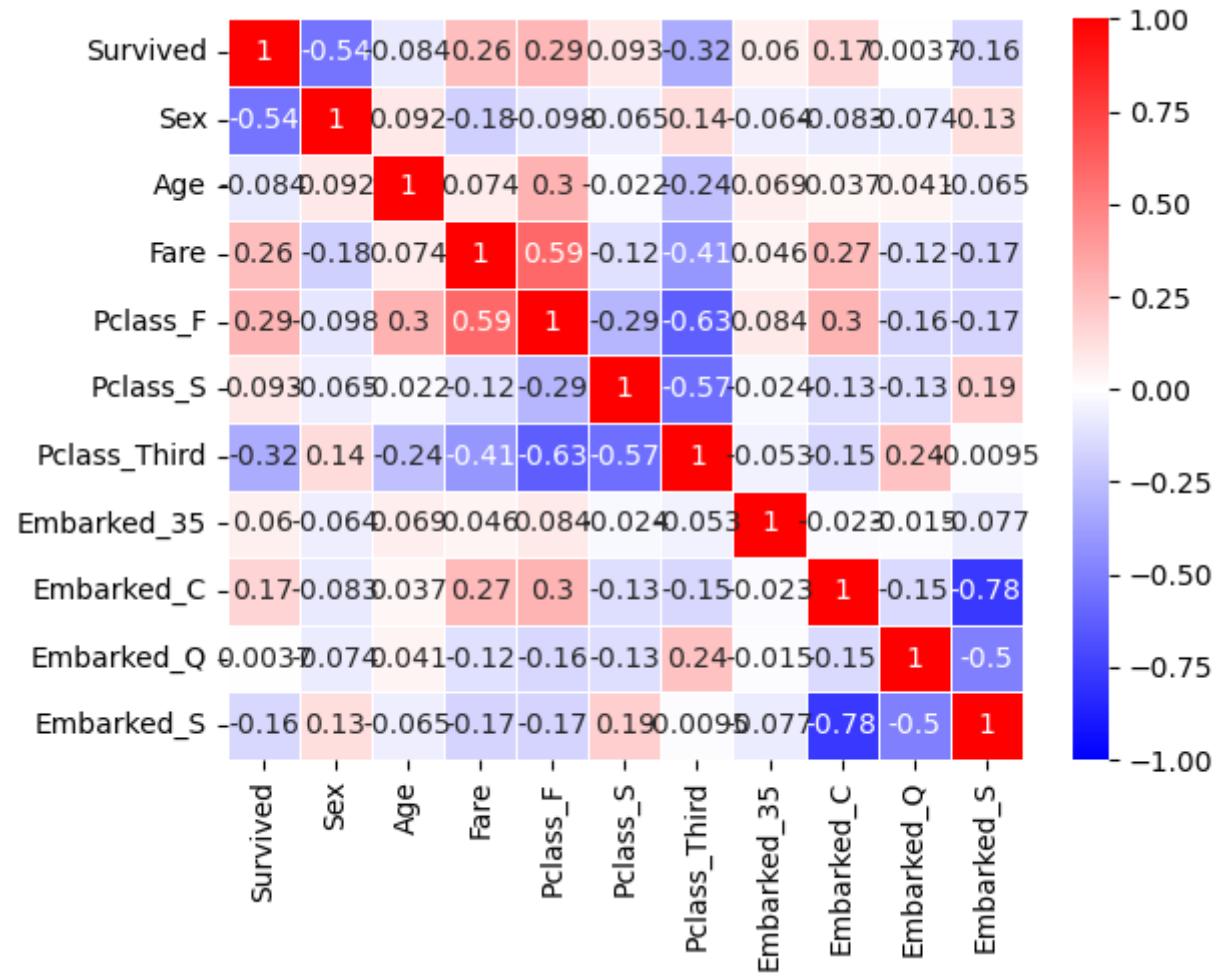
Out[162]:

	Survived	Sex	Age	Fare	Pclass_F	Pclass_S	Pclass_Third	Embarked_35	Embarked_C	Embarked_Q	Embarked_S
Survived	1.000000	-0.543351	-0.083713	0.257307	0.285904	0.093349	-0.322308	0.060095	0.168240	0.003650	-0.1556
Sex	-0.543351	1.000000	0.091930	-0.182333	-0.098013	-0.064746	0.137143	-0.064296	-0.082853	-0.074115	0.1257
Age	-0.083713	0.091930	1.000000	0.074199	0.302149	-0.022021	-0.242412	0.069343	0.036953	0.040528	-0.0650
Fare	0.257307	-0.182333	0.074199	1.000000	0.591711	-0.118557	-0.413333	0.045646	0.269335	-0.117216	-0.1666
Pclass_F	0.285904	-0.098013	0.302149	0.591711	1.000000	-0.288585	-0.626738	0.083847	0.296423	-0.155342	-0.1703
Pclass_S	0.093349	-0.064746	-0.022021	-0.118557	-0.288585	1.000000	-0.565210	-0.024197	-0.125416	-0.127301	0.1920
Pclass_Third	-0.322308	0.137143	-0.242412	-0.413333	-0.626738	-0.565210	1.000000	-0.052550	-0.153329	0.237449	-0.0095
Embarked_35	0.060095	-0.064296	0.069343	0.045646	0.083847	-0.024197	-0.052550	1.000000	-0.022864	-0.014588	-0.0765
Embarked_C	0.168240	-0.082853	0.036953	0.269335	0.296423	-0.125416	-0.153329	-0.022864	1.000000	-0.148258	-0.7783
Embarked_Q	0.003650	-0.074115	0.040528	-0.117216	-0.155342	-0.127301	0.237449	-0.014588	-0.148258	1.000000	-0.4966
Embarked_S	-0.155660	0.125722	-0.065062	-0.166603	-0.170379	0.192061	-0.009511	-0.076588	-0.778359	-0.496624	1.0000

```
In [163]: #finding heat map
import seaborn as sns

sns.heatmap(cor,vmax=1,vmin=-1,annot=True,linewidths=.5,cmap='bwr')
```

Out[163]: <Axes: >



```
In [164]: #####  
data.groupby('Survived').count()
```

```
Out[164]:
```

	PassengerId	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare	Cabin	Embarked
Survived											
0	549	549	549	549	424	549	549	549	549	68	549
1	342	342	342	342	290	342	342	342	342	136	340

```
In [165]: #which the parameter is predeected values can be removed from the data file  
#1) we copied the data into another data("y")  
#2) later we can removed those file from main data set  
y=data1['Survived']  
x=data1.drop('Survived',axis=1)
```

```
In [166]: #i am calling function to split  
#split enter data into ->67% traning , ->33% testing  
  
from sklearn.model_selection import train_test_split  
x_train,x_test,y_train,y_test=train_test_split(x,y,test_size=0.33,random_state=42)
```

```
In [167]: #in logistic we predicted 1 0r 2 only survival or not
from sklearn.linear_model import LogisticRegression
classifier=LogisticRegression()
classifier.fit(x_train,y_train)
```

Out[167]: LogisticRegression()

In a Jupyter environment, please rerun this cell to show the HTML representation or trust the notebook.
On GitHub, the HTML representation is unable to render, please try loading this page with nbviewer.org.

```
In [168]: y_pred=classifier.predict(x_test)
```

```
In [169]: y_pred
```

```
Out[169]: array([0, 0, 0, 1, 1, 1, 1, 0, 1, 1, 0, 0, 0, 0, 0, 1, 0, 1, 0, 0, 0, 0,
                1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 1, 1, 0, 0, 0, 0, 0,
                1, 0, 0, 0, 0, 0, 1, 1, 0, 0, 0, 1, 0, 1, 1, 1, 0, 1, 1, 0, 0, 1,
                0, 0, 0, 1, 1, 1, 1, 1, 0, 0, 1, 1, 1, 0, 0, 1, 1, 0, 0, 0, 1, 1,
                0, 0, 0, 1, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 1, 0, 0, 0, 1, 0, 0, 0,
                1, 0, 1, 0, 0, 0, 0, 0, 1, 0, 0, 1, 1, 0, 0, 0, 1, 1, 1, 0, 1, 0,
                0, 1, 0, 1, 1, 0, 0, 1, 0, 1, 0, 0, 1, 1, 0, 0, 1, 0, 0, 0, 0, 1,
                0, 0, 0, 1, 1, 1, 0, 0, 0, 1, 0, 0, 0, 1, 0, 0, 1, 1, 0, 1, 0, 0,
                0, 1, 1, 0, 0, 0, 0, 1, 1, 0, 0, 0, 0, 1, 0, 0, 0, 0, 1, 1, 1, 0,
                1, 1, 0, 0, 1, 0, 0, 1, 0, 0, 0, 0, 1, 0, 1, 0, 0, 0, 1, 0, 1, 0,
                0, 1, 0, 0, 0, 1, 0, 1, 1, 0, 0, 1, 0, 1, 0, 1, 1, 1, 1, 0, 0, 1,
                0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 1, 1, 0, 1, 0,
                0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 1, 0, 0, 0, 1, 0, 0, 0, 1, 0, 0, 0,
                1, 0, 0, 0, 0, 0, 1, 1, 0])
```

```
In [170]: from sklearn.metrics import confusion_matrix
confusion_matrix(y_test,y_pred)
```

```
Out[170]: array([[155, 20],
                [ 37, 83]])
```

```
In [171]: #effecency
          from sklearn.metrics import accuracy_score
          accuracy_score(y_test,y_pred)
```

```
Out[171]: 0.8067796610169492
```

```
In [ ]:
```

```
In [ ]:
```