

NAME :U PAVAN KALYAN

REG NO :21BCE9090

ASSIGNMENT-3

1.Import the Libraries

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

2.Importing the Dataset

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

```
In [3]: df
```

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	
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	
...	
886	887	0	2	Montvila, Rev. Juozas	male	27.0	0	0	211536	13.0000	NaN	
887	888	1	1	Graham, Miss. Margaret Edith	female	19.0	0	0	112053	30.0000	B42	
888	889	0	3	Johnston, Miss. Catherine Helen "Carrie"	female	NaN	1	2	W./C. 6607	23.4500	NaN	
889	890	1	1	Behr, Mr. Karl Howell	male	26.0	0	0	111369	30.0000	C148	
890	891	0	3	Dooley, Mr. Patrick	male	32.0	0	0	370376	7.7500	NaN	

891 rows × 12 columns

In [4]:

df.head()

Out[4]:

	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 [5]:

df.shape

Out[5]: (891, 12)

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]:

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 [8]: df.corr()
```

```
Out[8]:
```

	PassengerId	Survived	Pclass	Age	SibSp	Parch	Fare
PassengerId	1.000000	-0.005007	-0.035144	0.036847	-0.057527	-0.001652	0.012658
Survived	-0.005007	1.000000	-0.338481	-0.077221	-0.035322	0.081629	0.257307
Pclass	-0.035144	-0.338481	1.000000	-0.369226	0.083081	0.018443	-0.549500
Age	0.036847	-0.077221	-0.369226	1.000000	-0.308247	-0.189119	0.096067
SibSp	-0.057527	-0.035322	0.083081	-0.308247	1.000000	0.414838	0.159651
Parch	-0.001652	0.081629	0.018443	-0.189119	0.414838	1.000000	0.216225
Fare	0.012658	0.257307	-0.549500	0.096067	0.159651	0.216225	1.000000

```
In [9]: df.corr().Age.sort_values(ascending=False)
```

```
Out[9]: Age          1.000000
Fare          0.096067
PassengerId   0.036847
Survived     -0.077221
Parch        -0.189119
SibSp        -0.308247
Pclass       -0.369226
Name: Age, dtype: float64
```

```
In [10]: df.corr().Fare.sort_values(ascending=False)
```

```
Out[10]: Fare          1.000000
Survived   0.257307
Parch      0.216225
SibSp      0.159651
Age        0.096067
PassengerId 0.012658
Pclass     -0.549500
Name: Fare, dtype: float64
```

3. Checking for Null Values

```
In [11]: df.isnull().any()
```

```
Out[11]: 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 [12]: df.isnull().sum()
```

```
Out[12]: 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 [13]: df.corr()
```

```
Out[13]:
```

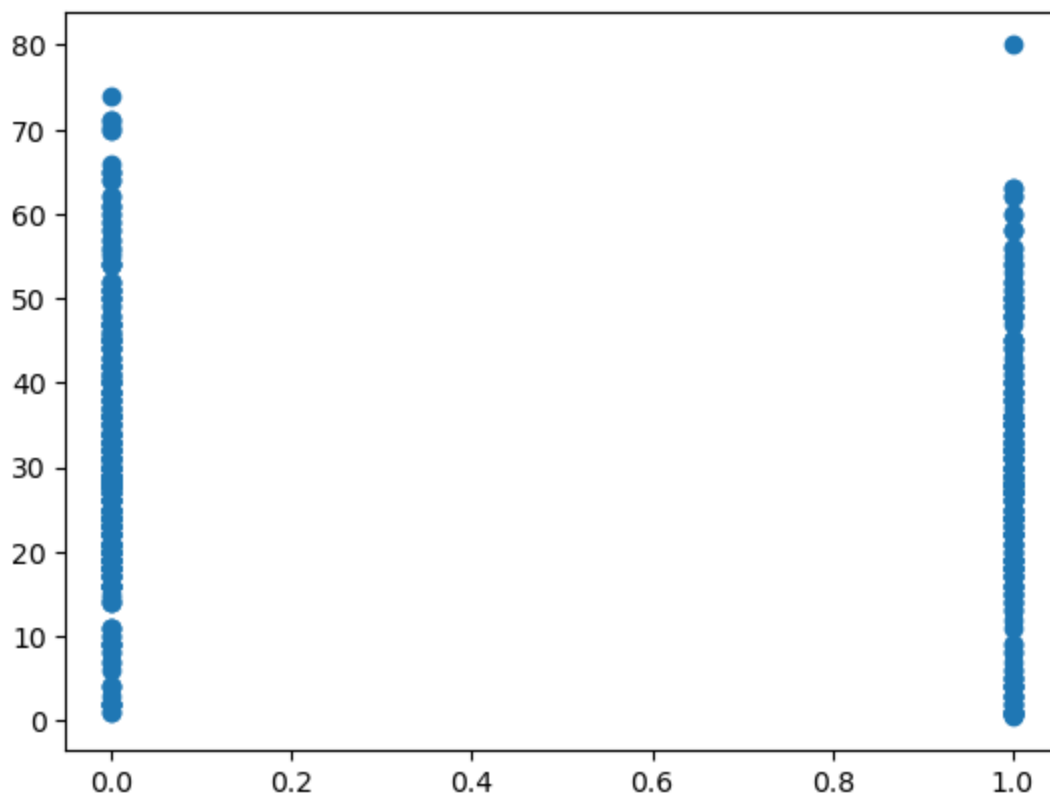
	PassengerId	Survived	Pclass	Age	SibSp	Parch	Fare
PassengerId	1.000000	-0.005007	-0.035144	0.036847	-0.057527	-0.001652	0.012658
Survived	-0.005007	1.000000	-0.338481	-0.077221	-0.035322	0.081629	0.257307
Pclass	-0.035144	-0.338481	1.000000	-0.369226	0.083081	0.018443	-0.549500
Age	0.036847	-0.077221	-0.369226	1.000000	-0.308247	-0.189119	0.096067
SibSp	-0.057527	-0.035322	0.083081	-0.308247	1.000000	0.414838	0.159651
Parch	-0.001652	0.081629	0.018443	-0.189119	0.414838	1.000000	0.216225
Fare	0.012658	0.257307	-0.549500	0.096067	0.159651	0.216225	1.000000

```
In [14]: df['Age'].fillna(df['Age'].median(), inplace=True)
         df['Embarked'].fillna(df['Embarked'].mode()[0], inplace=True)
         df['Cabin'].fillna(df['Cabin'].mode()[0], inplace=True)
```

4.Data Visualisation

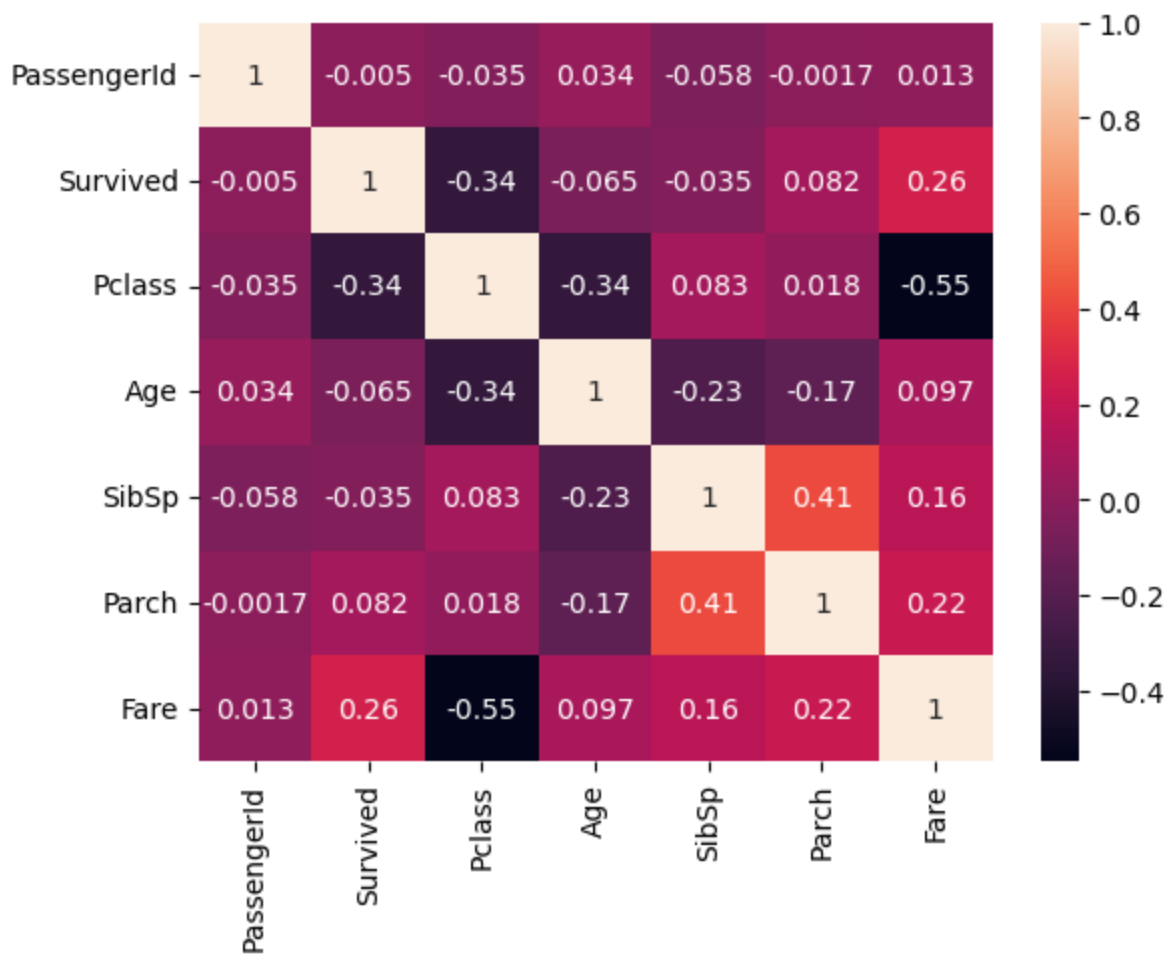
```
In [15]: plt.scatter(df["Survived"],df["Age"])
```

```
Out[15]: <matplotlib.collections.PathCollection at 0x226c8ab5f70>
```



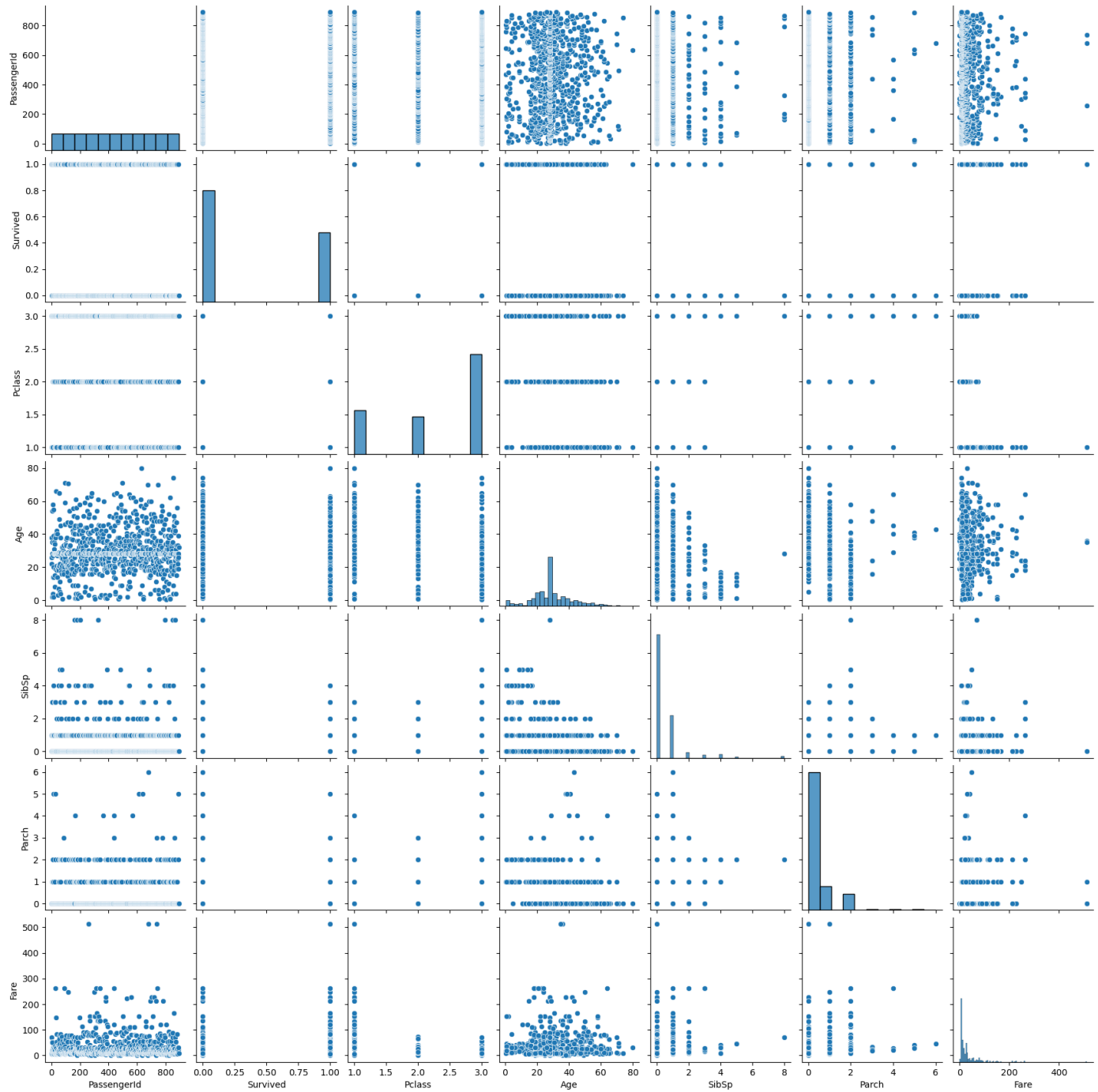
```
In [16]: sns.heatmap(df.corr(),annot=True)
```

```
Out[16]: <AxesSubplot:>
```



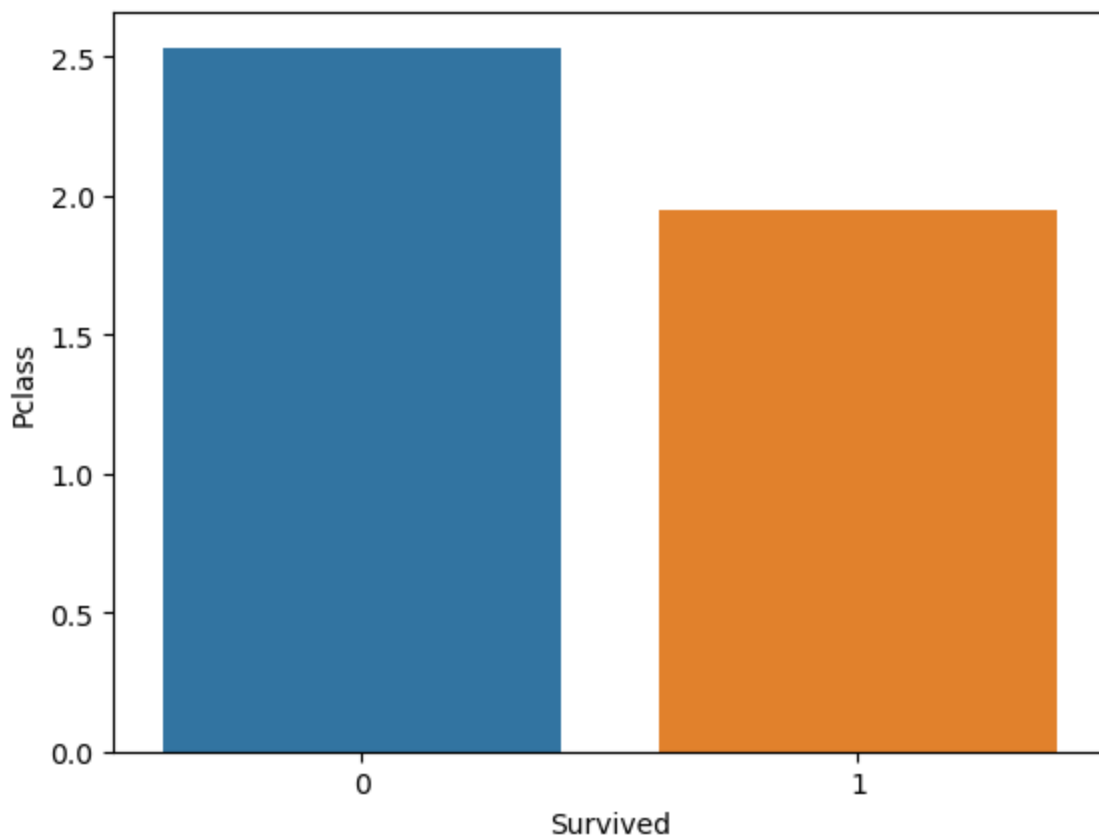
```
In [17]: sns.pairplot(df)
```

```
<seaborn.axisgrid.PairGrid at 0x226c8bf7340>
```



```
In [18]: sns.barplot(x=df["Survived"],y=df["Pclass"],ci=0)
```

```
Out[18]: <AxesSubplot:xlabel='Survived', ylabel='Pclass'>
```



5.Outlier Detection

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

```
Out[19]:
```

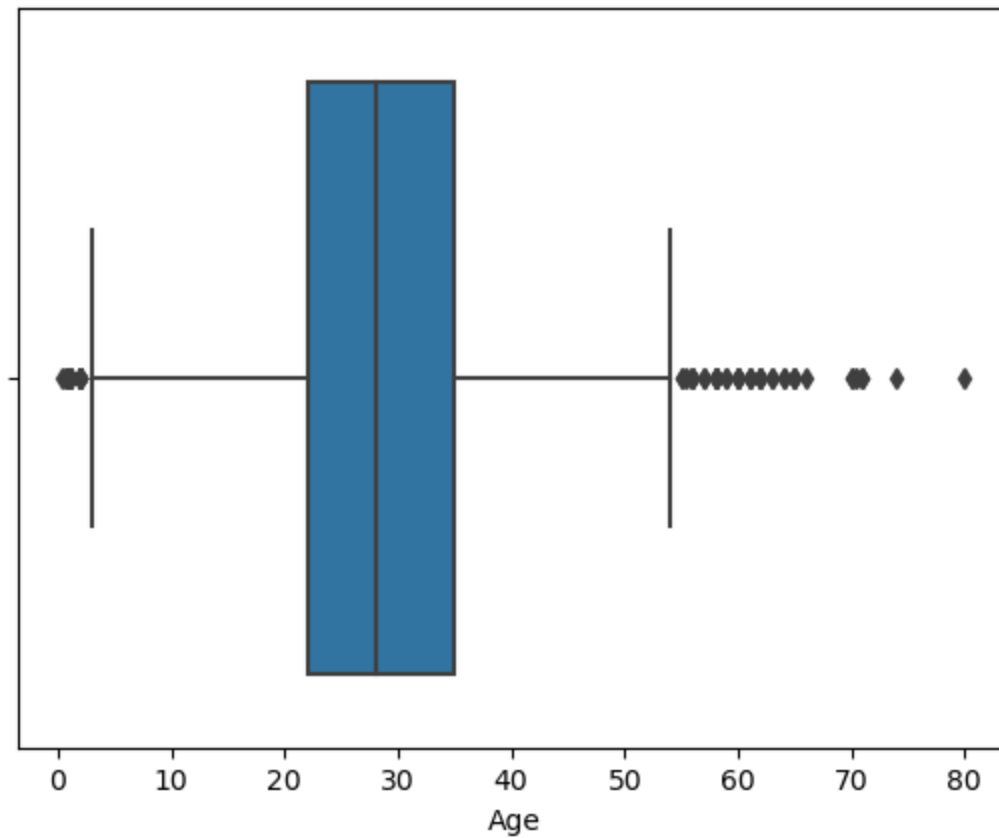
	PassengerId	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare	Cabin	Embarked
0	1	0	3	Braund, Mr. Owen Harris	male	22.0	1	0	A/5 21171	7.2500	B96 B98	S
1	2	1	1	Cumings, Mrs. John Bradley (Florence Briggs Th...	female	38.0	1	0	PC 17599	71.2833	C85	C
2	3	1	3	Heikkinen, Miss. Laina	female	26.0	0	0	STON/O2. 3101282	7.9250	B96 B98	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	B96 B98	S

```
In [20]: sns.boxplot(df["Age"])
```


C:\ProgramData\Anaconda3\lib\site-packages\seaborn_decorators.py:36: FutureWarning: Pass the following variable as a keyword arg: x. From version 0.12, the only valid positional argument will be `data`, and passing other arguments without an explicit keyword will result in an error or misinterpretation.

warnings.warn(

Out[20]: <AxesSubplot:xlabel='Age'>

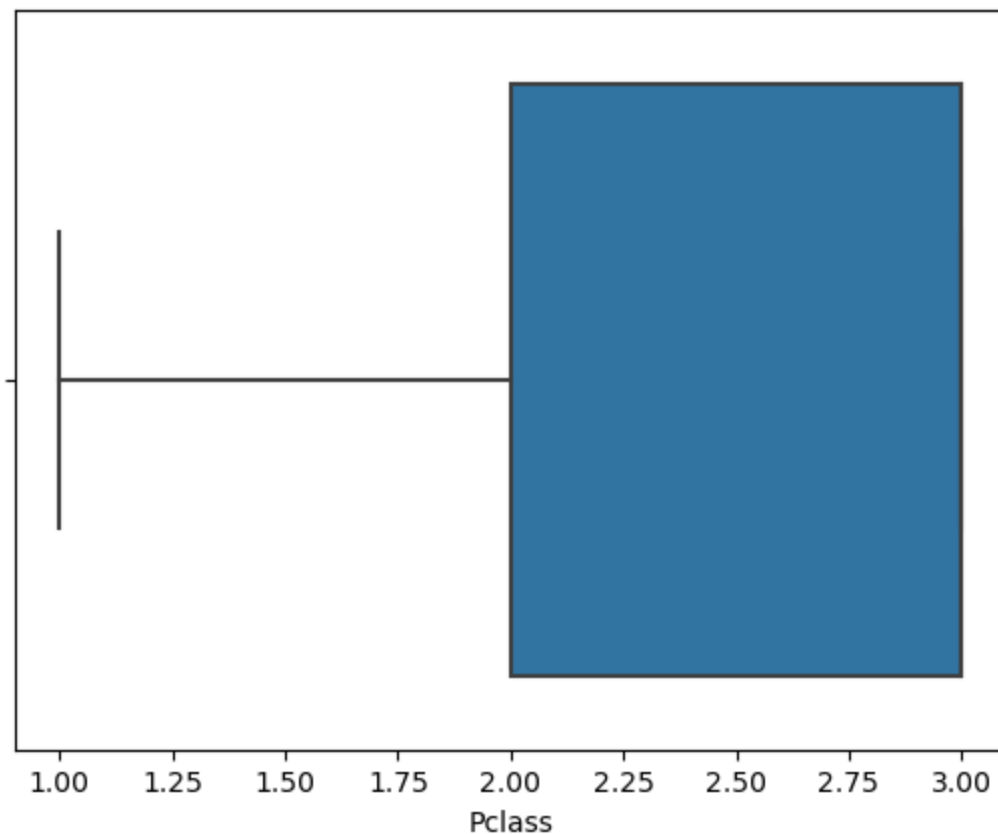


In [21]: sns.boxplot(df["Pclass"])

C:\ProgramData\Anaconda3\lib\site-packages\seaborn_decorators.py:36: FutureWarning: Pass the following variable as a keyword arg: x. From version 0.12, the only valid positional argument will be `data`, and passing other arguments without an explicit keyword will result in an error or misinterpretation.

warnings.warn(

Out[21]: <AxesSubplot:xlabel='Pclass'>



6.Splitting Dependent and Independent Variables

In [22]: `df.head()`

Out[22]:

	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	B96 B98	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	B96 B98	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	B96 B98	S

In [23]: `#independent variable should be 2d array or dataframe
x=df.drop(columns=["Survived", "PassengerId", "Name", "Ticket", "Cabin"], axis=1)
x.head()`

```
Out[23]:
```

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

```
In [24]: x.shape
```

```
Out[24]: (891, 7)
```

```
In [25]: type(x)
```

```
Out[25]: pandas.core.frame.DataFrame
```

```
In [26]: y=df["Survived"]
y.head()
```

```
Out[26]:
```

0	0
1	1
2	1
3	1
4	0

Name: Survived, dtype: int64

7.Encoding

```
In [27]: from sklearn.preprocessing import LabelEncoder
le=LabelEncoder()
```

```
In [28]: x["Sex"]=le.fit_transform(x["Sex"])
```

```
In [29]: x.head()
```

```
Out[29]:
```

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

```
In [30]: print(le.classes_)

['female' 'male']
```

```
In [31]: mapping=dict(zip(le.classes_,range(len(le.classes_))))
mapping
```

```
Out[31]: {'female': 0, 'male': 1}
```

```
In [32]: x["Embarked"]=le.fit_transform(x["Embarked"])
```

```
Out[32]:
```

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

```
In [33]: print(le.classes_)  
['C' 'Q' 'S']
```

```
In [34]: mapping=dict(zip(le.classes_, range(len(le.classes_))))  
mapping
```

```
Out[34]: {'C': 0, 'Q': 1, 'S': 2}
```

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

```
Out[35]:
```

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

8.Feature Scaling

```
In [37]: from sklearn.preprocessing import MinMaxScaler  
ms=MinMaxScaler()
```

```
In [38]: x_scaled=pd.DataFrame(ms.fit_transform(x),columns=x.columns)
```

```
In [39]: x_scaled
```

Out[39]:

	Pclass	Sex	Age	SibSp	Parch	Fare	Embarked
0	1.0	1.0	0.271174	0.125	0.000000	0.014151	1.0
1	0.0	0.0	0.472229	0.125	0.000000	0.139136	0.0
2	1.0	0.0	0.321438	0.000	0.000000	0.015469	1.0
3	0.0	0.0	0.434531	0.125	0.000000	0.103644	1.0
4	1.0	1.0	0.434531	0.000	0.000000	0.015713	1.0
...
886	0.5	1.0	0.334004	0.000	0.000000	0.025374	1.0
887	0.0	0.0	0.233476	0.000	0.000000	0.058556	1.0
888	1.0	0.0	0.346569	0.125	0.333333	0.045771	1.0
889	0.0	1.0	0.321438	0.000	0.000000	0.058556	0.0
890	1.0	1.0	0.396833	0.000	0.000000	0.015127	0.5

891 rows × 7 columns

9.Splitting Data into Train and Test

```
In [40]: from sklearn.model_selection import train_test_split
x_train,x_test,y_train,y_test=train_test_split(x_scaled,y,test_size=0.2,random_state=0)
```

```
In [41]: print(x_train.shape,x_test.shape,y_train.shape,y_test.shape)
(712, 7) (179, 7) (712,) (179,)
```

```
In [42]: x_train,x_test,y_train,y_test=train_test_split(x,y,test_size=0.2,random_state=0)
```

```
In [43]: x_train
```

Out[43]:

	Pclass	Sex	Age	SibSp	Parch	Fare	Embarked
140	3	0	28.0	0	2	15.2458	0
439	2	1	31.0	0	0	10.5000	2
817	2	1	31.0	1	1	37.0042	0
378	3	1	20.0	0	0	4.0125	0
491	3	1	21.0	0	0	7.2500	2
...
835	1	0	39.0	1	1	83.1583	0
192	3	0	19.0	1	0	7.8542	2
629	3	1	28.0	0	0	7.7333	1
559	3	0	36.0	1	0	17.4000	2
684	2	1	60.0	1	1	39.0000	2

712 rows × 7 columns