

SreeLekha
21bce7030
Vitap

```
import numpy as np
import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt
```

```
data=pd.read_csv("Titanic-Dataset.csv")
```

```
data.head()
```

	PassengerId	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare
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 T. B.)	female	38.0	1	0	PC 17599	71.2834

```
data.tail()
```

	PassengerId	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare
886	887	0	2	Montvila, Rev. Juozas	male	27.0	0	0	211536	13.00
887	888	1	1	Graham, Miss. Margaret Edith	female	19.0	0	0	112053	30.00

```
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
```

```
data.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	

```
data.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
```

```
data.isnull().sum()
```

```
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
```

```
mean=data["Age"].mean()
```

```
data["Age"]=data["Age"].fillna(mean)
```

```
data["Age"].tail()
```

```
886    27.000000
887    19.000000
888    29.699118
889    26.000000
890    32.000000
Name: Age, dtype: float64
```

```
data["Age"].isnull().sum()
```

```
0
```

```
data["Cabin"]
```

```
0      NaN
1      C85
2      NaN
3      C123
4      NaN
...
886    NaN
887     B42
888    NaN
889    C148
890    NaN
Name: Cabin, Length: 891, dtype: object
```

```
data.isnull().sum()
```

```

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

```

```
cor=data.corr()
```

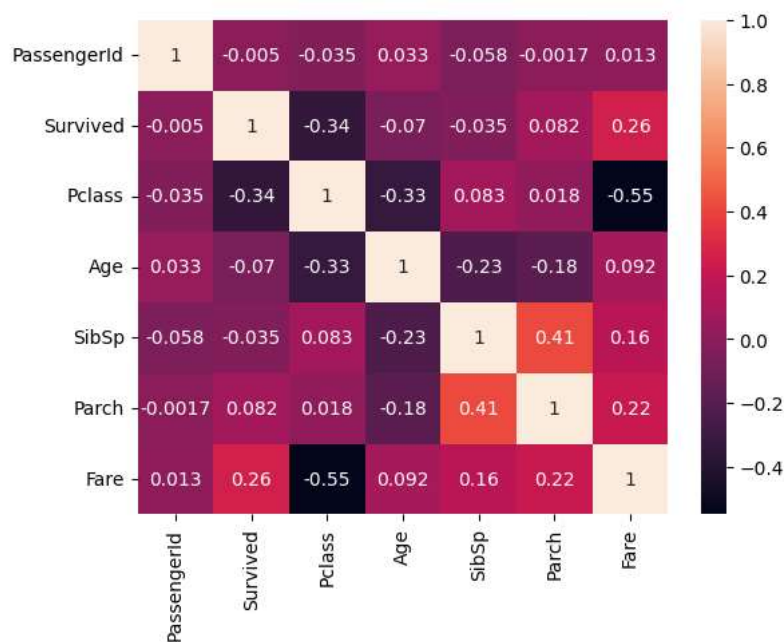
```

<ipython-input-15-410fe4458127>:1: FutureWarning: The default value of numeric_only in DataFrame.corr is deprecated. In a future version
cor=data.corr()

```

```
sns.heatmap(cor,annot=True)
```

<Axes: >



```
sns.boxplot(data["Age"])
```

```

<Axes: >
80
Age_q1 = data.Age.quantile(0.25)
Age_q3 = data.Age.quantile(0.75)
print(Age_q1)
print(Age_q3)

22.0
35.0
IQR_Age=Age_q3-Age_q1
IQR_Age

13.0
upperlimit_Age=Age_q3+1.5*IQR_Age
upperlimit_Age

54.5

lower_limit_Age = Age_q1-1.5*IQR_Age
lower_limit_Age

2.5

median_Age=data["Age"].median()
median_Age

29.69911764705882

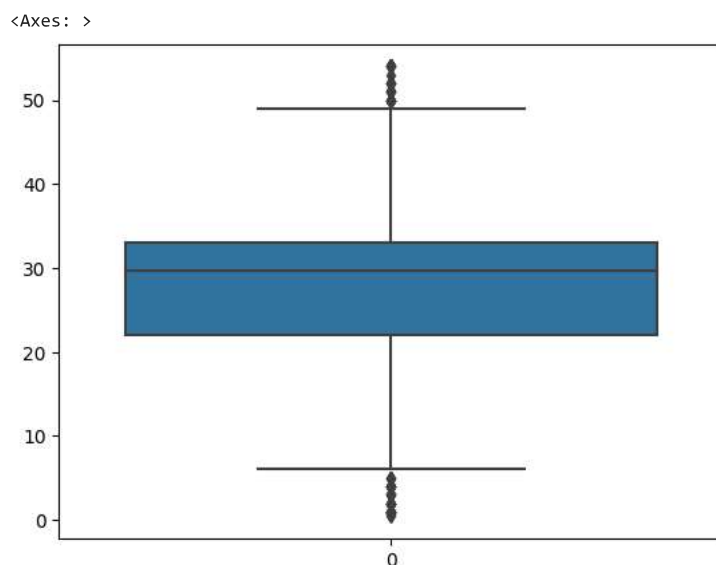
data["Age"]=np.where(data["Age"]>upperlimit_Age,median_Age,data["Age"])

(data["Age"]>54.5).sum()

0

```

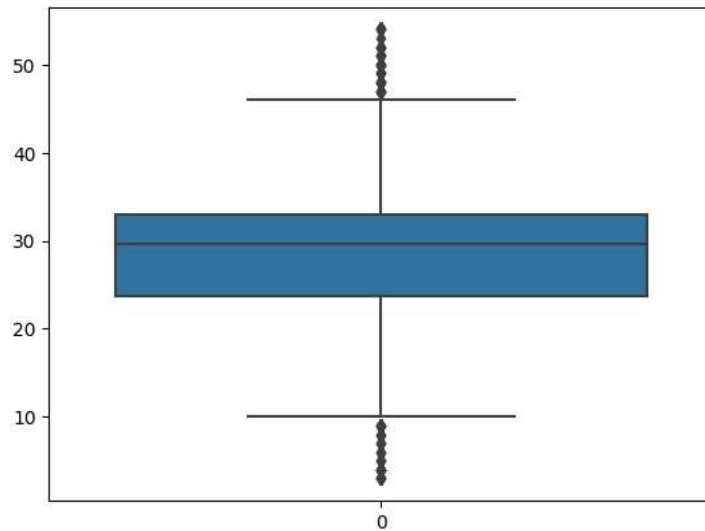
```
sns.boxplot(data["Age"])
```



```
data["Age"]=np.where(data["Age"]<lower_limit_Age,median_Age,data["Age"])
```

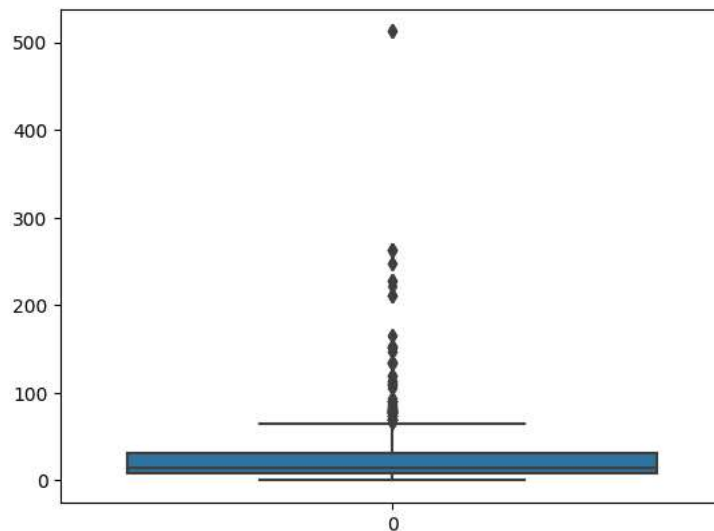
```
sns.boxplot(data["Age"])
```

<Axes: >



```
sns.boxplot(data["Fare"])
```

<Axes: >



```
Fare_q1 = data.Fare.quantile(0.25)
Fare_q3 = data.Fare.quantile(0.75)
print(Fare_q1)
print(Fare_q3)
```

```
7.9104
31.0
```

```
IQR_Fare=Fare_q3-Fare_q1
IQR_Fare
```

```
23.0896
```

```
upperlimit_Fare=Fare_q3+1.5*IQR_Fare
upperlimit_Fare
```

```
65.6344
```

```
lower_limit_Fare = Fare_q1-1.5*IQR_Fare
lower_limit_Fare
```

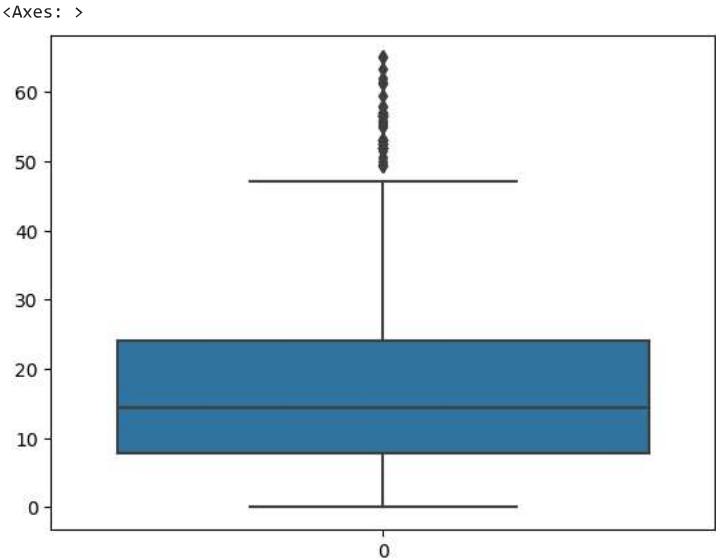
```
-26.724
```

```
median_Fare=data["Fare"].median()
median_Fare
```

14.4542

```
data['Fare'] = np.where(
    (data['Fare'] > upperlimit_Fare),
    median_Fare,
    data['Fare']
)
```

```
sns.boxplot(data["Fare"])
```



```
(data["Fare"]>65).sum()
```

0

```
data.drop(['Name'],axis=1,inplace=True)
```

data

	PassengerId	Survived	Pclass	Sex	Age	SibSp	Parch	Ticket	Fare	Cabin
0	1	0	3	male	22.000000	1	0	A/5 21171	7.2500	NaN
1	2	1	1	female	38.000000	1	0	PC 17599	14.4542	C85
2	3	1	3	female	26.000000	0	0	STON/O2. 3101282	7.9250	NaN
3	4	1	1	female	35.000000	1	0	113803	53.1000	C123
4	5	0	3	male	35.000000	0	0	373450	8.0500	NaN
...
886	887	0	2	male	27.000000	0	0	211536	13.0000	NaN
887	888	1	1	female	19.000000	0	0	112053	30.0000	B42
888	889	0	3	female	29.699118	1	2	W./C. 6607	23.4500	NaN
889	890	1	1	male	26.000000	0	0	111369	30.0000	C148
890	891	0	3	male	32.000000	0	0	370376	7.7500	NaN

```
data.drop(['Ticket'],axis=1,inplace=True)
```



data

	PassengerId	Survived	Pclass	Sex	Age	SibSp	Parch	Fare	Cabin	Embarked
0	1	0	3	male	22.000000	1	0	7.2500	NaN	S
1	2	1	1	female	38.000000	1	0	14.4542	C85	C
2	3	1	3	female	26.000000	0	0	7.9250	NaN	S
3	4	1	1	female	35.000000	1	0	53.1000	C123	S
4	5	0	3	male	35.000000	0	0	8.0500	NaN	S
...
886	887	0	2	male	27.000000	0	0	13.0000	NaN	S
887	888	1	1	female	19.000000	0	0	30.0000	B42	S
888	889	0	3	female	29.699118	1	2	23.4500	NaN	S
889	890	1	1	male	26.000000	0	0	30.0000	C148	C
890	891	0	3	male	32.000000	0	0	7.7500	NaN	Q

891 rows × 10 columns

```
data.drop(["PassengerId"],axis=1,inplace=True)
```

data

	Survived	Pclass	Sex	Age	SibSp	Parch	Fare	Cabin	Embarked	
0	0	3	male	22.000000	1	0	7.2500	NaN	S	
1	1	1	female	38.000000	1	0	14.4542	C85	C	
2	1	3	female	26.000000	0	0	7.9250	NaN	S	
3	1	1	female	35.000000	1	0	53.1000	C123	S	
4	0	3	male	35.000000	0	0	8.0500	NaN	S	
...	
886	0	2	male	27.000000	0	0	13.0000	NaN	S	
887	1	1	female	19.000000	0	0	30.0000	B42	S	
888	0	3	female	29.699118	1	2	23.4500	NaN	S	
889	1	1	male	26.000000	0	0	30.0000	C148	C	
890	0	3	male	32.000000	0	0	7.7500	NaN	Q	

891 rows × 9 columns

data

```

Survived  Pclass  Sex      Age  SibSp  Parch    Fare  Cabin  Embarked
0         0       3   male  22.000000  1      0   7.2500   NaN    S
y=data["Survived"]

y.head()

0      0
1      1
2      1
3      1
4      0
Name: Survived, dtype: int64

```

data

	Survived	Pclass	Sex	Age	SibSp	Parch	Fare	Cabin	Embarked
0	0	3	male	22.000000	1	0	7.2500	NaN	S
1	1	1	female	38.000000	1	0	14.4542	C85	C
2	1	3	female	26.000000	0	0	7.9250	NaN	S
3	1	1	female	35.000000	1	0	53.1000	C123	S
4	0	3	male	35.000000	0	0	8.0500	NaN	S
...
886	0	2	male	27.000000	0	0	13.0000	NaN	S
887	1	1	female	19.000000	0	0	30.0000	B42	S
888	0	3	female	29.699118	1	2	23.4500	NaN	S
889	1	1	male	26.000000	0	0	30.0000	C148	C
890	0	3	male	32.000000	0	0	7.7500	NaN	Q

891 rows × 9 columns

```
from sklearn.preprocessing import LabelEncoder
```

```
le=LabelEncoder()
```

```
data["Sex"]=le.fit_transform(data["Sex"])
```

```
data["Sex"]
```

```

0      1
1      0
2      0
3      0
4      1
..
886    1
887    0
888    0
889    1
890    1
Name: Sex, Length: 891, dtype: int64


```

```
data.head()
```


	Survived	Pclass	Sex	Age	SibSp	Parch	Fare	Cabin	Embarked	
0	0	3	1	22.0	1	0	7.2500	NaN	S	

```
data["Embarked"]=le.fit_transform(data["Embarked"])
```

```
data.head()
```

	Survived	Pclass	Sex	Age	SibSp	Parch	Fare	Cabin	Embarked	
0	0	3	1	22.0	1	0	7.2500	NaN	2	
1	1	1	0	38.0	1	0	14.4542	C85	0	
2	1	3	0	26.0	0	0	7.9250	NaN	2	
3	1	1	0	35.0	1	0	53.1000	C123	2	
4	0	3	1	35.0	0	0	8.0500	NaN	2	

```
data["Pclass"].nunique()
```

```
3
```

```
data["Pclass"].unique()
```

```
array([3, 1, 2])
```

```
data["Sex"].unique()
```

```
array([1, 0])
```

```
data["Embarked"].unique()
```

```
array([2, 0, 1, 3])
```

```
from sklearn.model_selection import train_test_split
x_train,x_test,y_train,y_test=train_test_split(data,y,test_size=0.3,random_state=0)
```

```
x_train.shape,x_test.shape,y_train.shape,y_test.shape
```

```
((623, 9), (268, 9), (623,), (268,))
```

```
from sklearn.preprocessing import StandardScaler
```

```
sc=StandardScaler()
```

