

Data Preprocessing on Titanic Dataset

Data Preprocessing involves the following steps: 1. Importing the Libraries, 2. Importing the dataset, 3. Data Visualization, 4. Outlier Detection 5. Splitting Dependent and Independent variables 6. Perform Encoding 7. Feature Scaling, 8. Splitting Data into Train and Test

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
```

	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		Cummings, 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		Morreim, Rev. Jozsef	male	27.0	0	0		211536	13.0000	NaN	S
887	888	1	1		Graham, Miss. Margaret Edith	female	19.0	0	0		110503	30.0000	B42	S
888	889	0	3		Johnson, Mrs. Catherine Helen "Carnie"	female	NA	1	2		W/C 6007	23.4500	NaN	S
889	890	1	1		Behr, Mr. Karl Howell	male	26.0	0	0		111369	53.1000	C148	C
890	891	0	3		Dooley, Mr. Patrick	male	32.0	0	0		370376	7.7500	NaN	Q

891 rows x 12 columns

```
In [4]: df.head()
```

	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		Cummings, 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
```

```
(891, 12)
```

```
In [6]: 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 Fare           891 non-null    float64
 9 Cabin         204 non-null    object
10 Embarked      889 non-null    object
dtypes: float64(2), int64(5), object(5)
memory usage: 63.7+ KB
```

```
In [7]: df.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	0.520000	0.381594	32.204208
std	257.353442	0.486592	0.830771	14.526497	1.102743	0.800957	49.693429
min	1.000000	0.000000	1.000000	0.420000	0.000000	0.000000	0.000000
25%	223.000000	0.000000	2.000000	20.125000	0.000000	0.000000	7.912400
50%	446.000000	0.000000	3.000000	29.000000	0.000000	0.000000	14.454200
75%	668.500000	1.000000	3.000000	35.000000	1.000000	0.000000	31.000000
max	891.000000	1.000000	3.000000	80.000000	8.000000	6.000000	512.329200

3.Handling Null Values.

```
In [8]: df.isnull().any()
```

```
PassengerId    False
Survived        False
Pclass         False
Name            False
Sex            False
Age            False
SibSp          False
Parch          False
Ticket         False
Fare           False
Cabin          True
Embarked       True
dtype: bool
```

```
In [9]: df.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
```

```
In [10]: #There are null values in Age,Cabin and Embarked columns
#Handling null values in Age column- numerical type
df["Age"].fillna(df["Age"].mean(),inplace=True)
```

```
In [11]: #Handling null values in Cabin column- categorical type
df["Cabin"].fillna(df["Cabin"].mode()[0],inplace=True)
```

```
In [12]: #Handling null values in Embarked column- categorical type
df["Embarked"].fillna(df["Embarked"].mode()[0],inplace=True)
```

```
In [13]: df.isnull().any()
```

```
PassengerId    False
Survived        False
Pclass         False
Name            False
Sex            False
Age            False
SibSp          False
Parch          False
Ticket         False
Fare           False
Cabin          False
Embarked       False
dtype: bool
```

```
In [14]: df.isnull().sum()
```

```
PassengerId    0
Survived        0
Pclass         0
Name            0
Sex            0
Age            0
SibSp          0
Parch          0
Ticket         0
Fare           0
Cabin          0
Embarked       0
dtype: int64
```

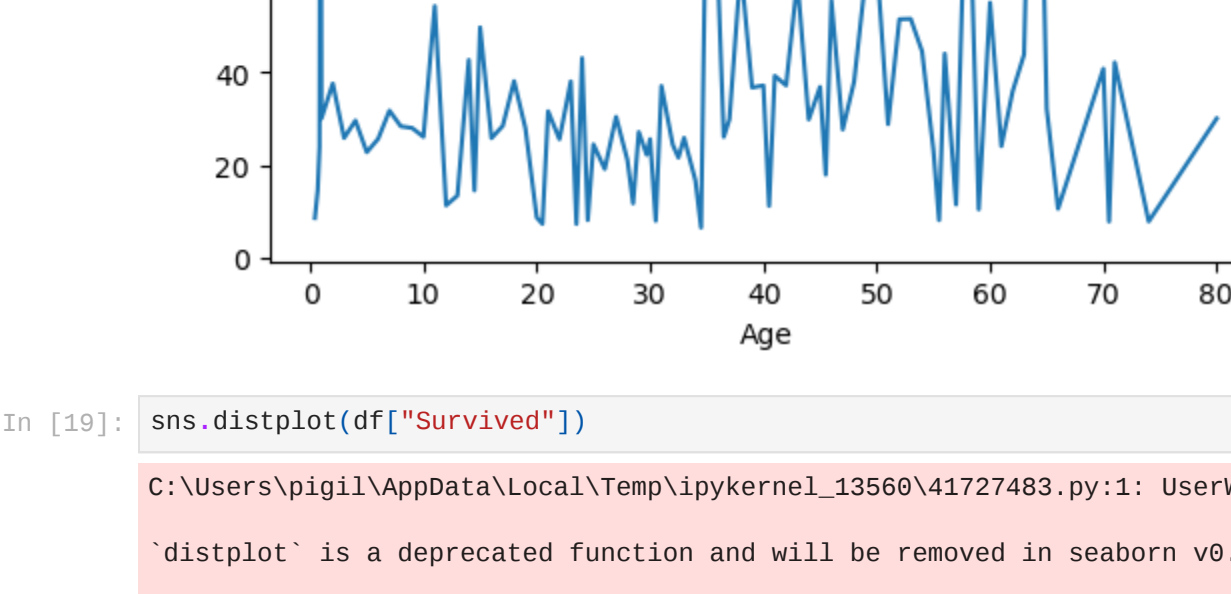
```
In [15]: df.head()
```

	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	896 B98	S
1	2	1	1		Cummings, 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	896 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	896 B98	S

4.Data Visualization

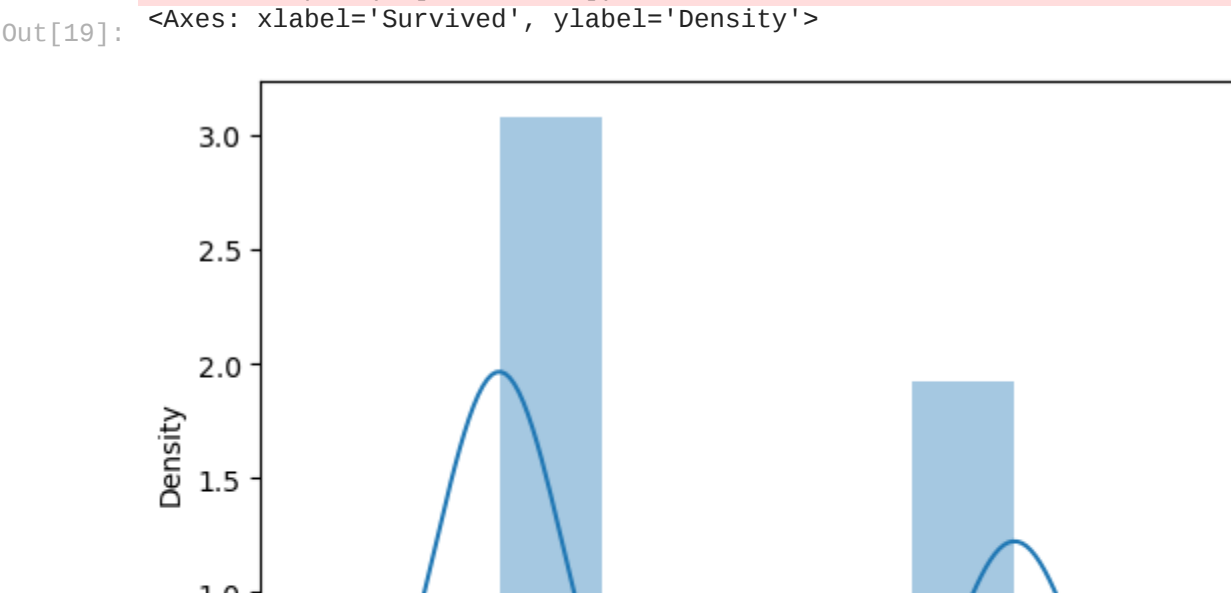
```
In [16]: sns.scatterplot(x="Age",y="Fare",data=df)
```

```
<Axes: xlabel='Age', ylabel='Fare'>
```



```
In [17]: plt.scatter(df["Survived"],df["Fare"])
```

```
<matplotlib.collections.PathCollection at 0x28e387591ab>
```

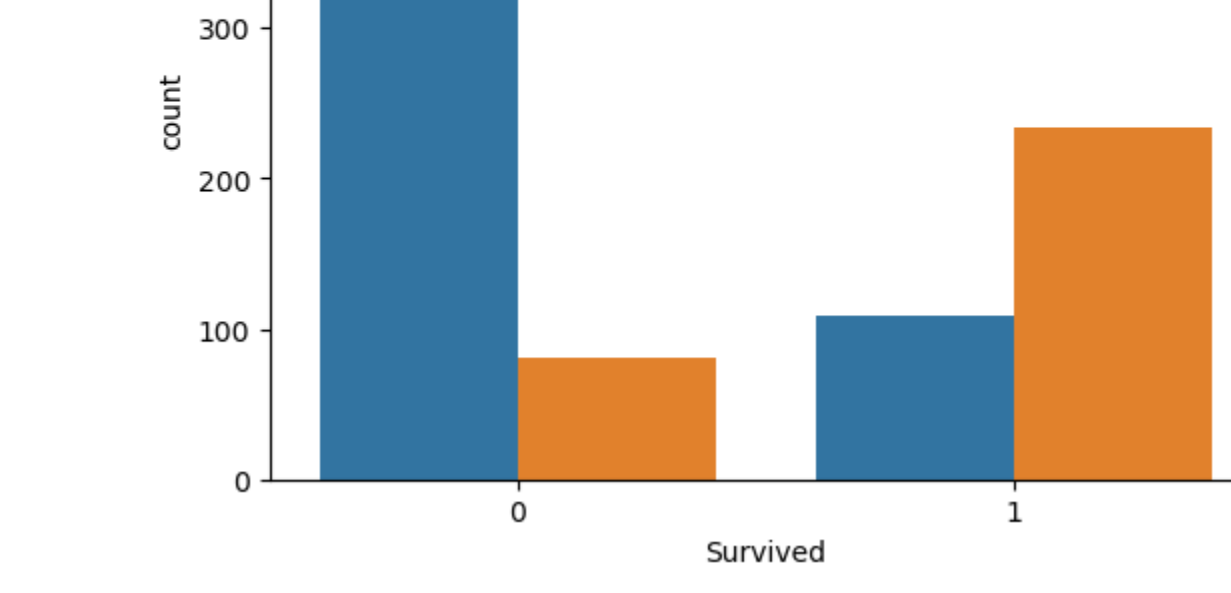


```
In [18]: sns.lmplot(x="Age",y="Fare",data=df,ci=None)
```

```
C:\Users\vgill\AppData\Local\Temp\ipykernel_13560\2712183228.py:1: FutureWarning:
The 'ci' parameter is deprecated. Use 'errorbar=None' for the same effect.

sns.lmplot(x="Age",y="Fare",data=df,ci=None)
```

```
<Axes: xlabel='Age', ylabel='Fare'>
```



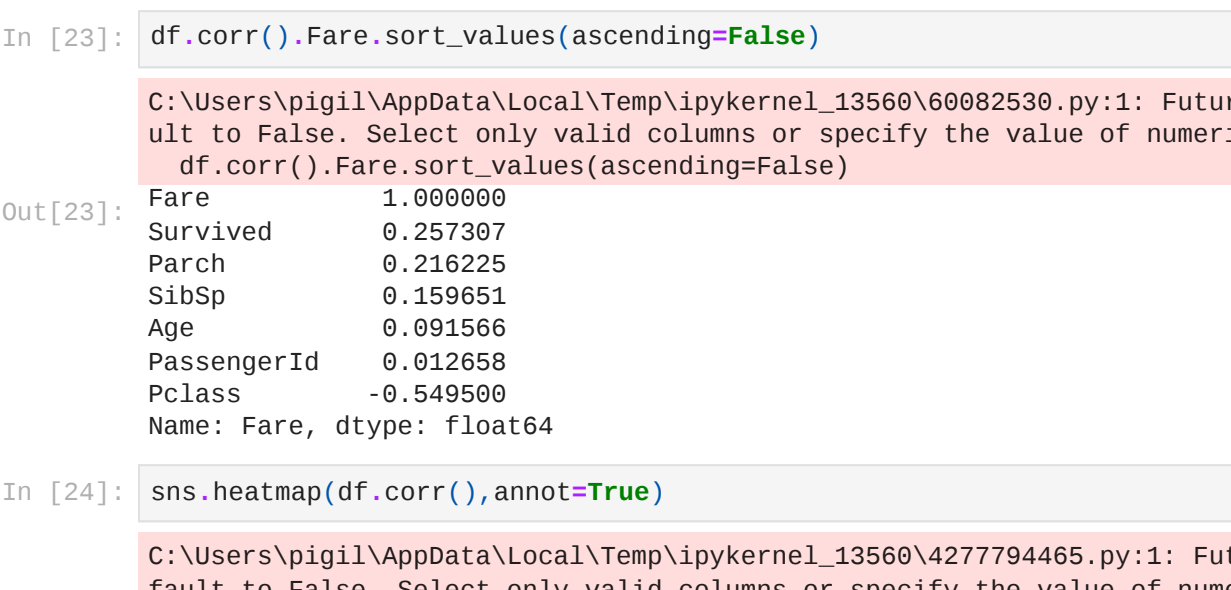
```
In [19]: sns.distplot(df["Survived"])
```

```
C:\Users\vgill\AppData\Local\Temp\ipykernel_13560\41727483.py:1: UserWarning:
'distplot' is a deprecated function and will be removed in seaborn v0.14.0.

Please adapt your code to use either 'displot' (a figure-level function with similar flexibility) or 'histplot' (an axes-level function for histograms).

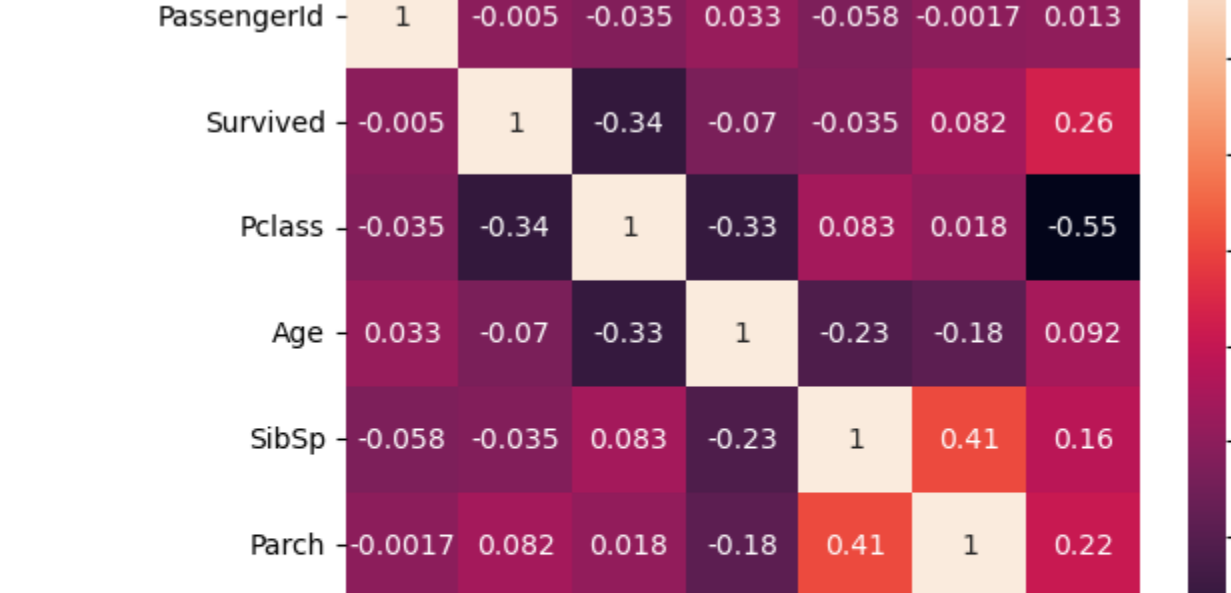
For a guide to updating your code to use the new functions, please see
https://github.com/mwaskom/dea4147ed2974457ad03727560b0e9751

<Axes: xlabel='Survived', ylabel='Density'>
```



```
In [20]: sns.countplot(x="Survived",data=df,hue="Sex")
```

```
<Axes: xlabel='Survived', ylabel='count'>
```



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

```
C:\Users\vgill\AppData\Local\Temp\ipykernel_13560\1134722465.py:1: FutureWarning: The default value of numeric_only in DataFrame.corr is deprecated. In a future version, it will default to False. Select only valid columns or specify the value of numeric_only to silence this warning.

df.corr()
```

	PassengerId	Survived	Pclass	Age	SibSp	Parch	Fare
PassengerId	1.000000	-0.090007	-0.025144	0.033207	-0.057307	-0.001662	0.023668
Survived	-0.090007	1.000000	-0.338481	-0.069809	-0.053032	0.001429	-0.257807
Pclass	-0.025144	-0.338481	1.000000	-0.331339	0.003031	0.018463	-0.549590
Age	0.033207	-0.069809	0.331339	1.000000	-0.232625	-0.179191	0.091566
SibSp	-0.057307	-0.053032	0.003031	-0.232625	1.000000	0.414838	0.159561
Parch	-0.001662	0.001429	0.018463	-0.179191	0.414838	1.000000	0.216225
Fare	0.023668	-0.257807	-0.549590	0.091566	0.159561	0.216225	1.000000

```
In [23]: df.corr().fare.sort_values(ascending=False)
```

```
C:\Users\vgill\AppData\Local\Temp\ipykernel_13560\68882530.py:1: FutureWarning: The default value of numeric_only in DataFrame.corr is deprecated. In a future version, it will default to False. Select only valid columns or specify the value of numeric_only to silence this warning.

df.corr().fare.sort_values(ascending=False)
```

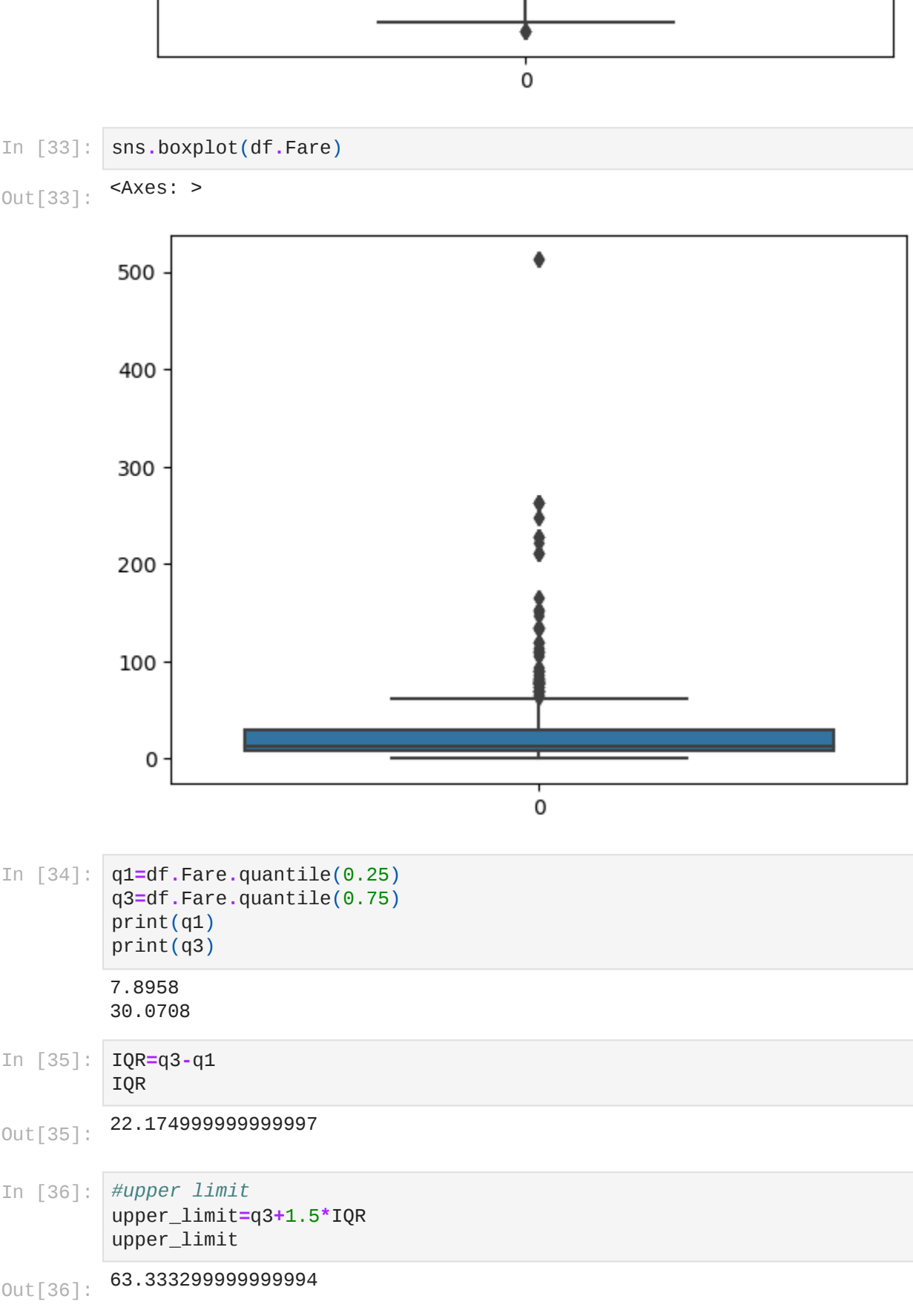
```
PassengerId    3.098880
Survived        0.257807
Parch          0.216225
SibSp          0.159561
Age            0.091566
PassengerId    0.023668
Pclass         -0.549590
Name: Fare, dtype: float64
```

```
In [24]: sns.heatmap(df.corr(),annot=True)
```

```
C:\Users\vgill\AppData\Local\Temp\ipykernel_13560\427794465.py:1: FutureWarning: The default value of numeric_only in DataFrame.corr is deprecated. In a future version, it will default to False. Select only valid columns or specify the value of numeric_only to silence this warning.

sns.heatmap(df.corr(),annot=True)
```

```
<Axes: >
```



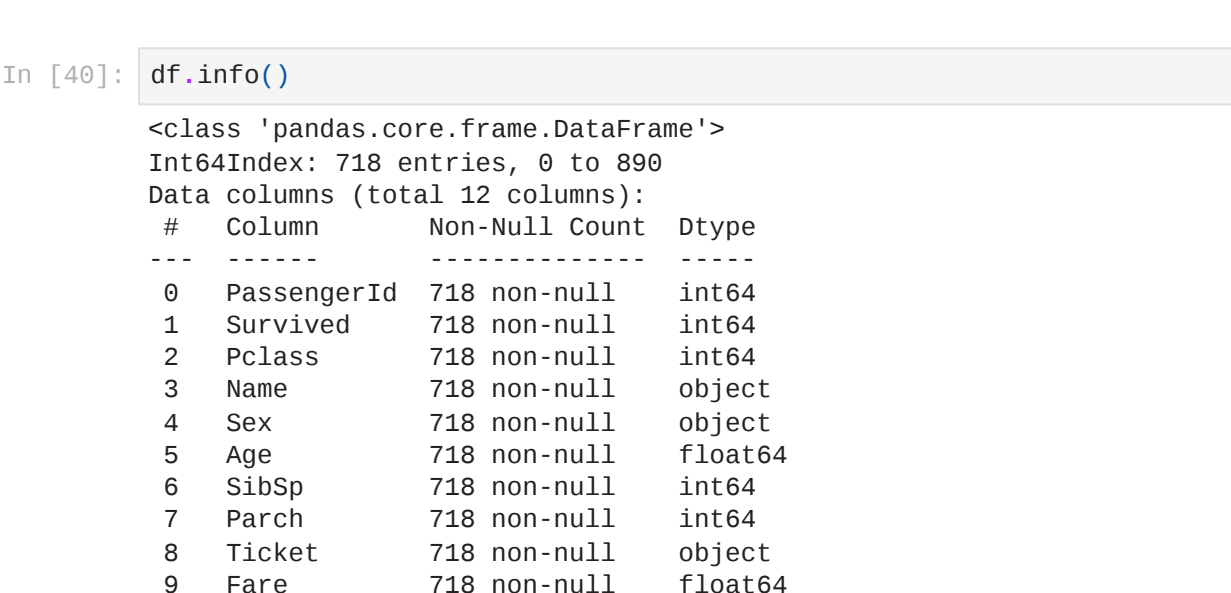
5.Outlier Detection

```
In [25]: df.head()
```

	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	896 B98	S
1	2	1	1		Cummings, 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	896 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	896 B98	S

```
In [26]: sns.boxplot(df.Age)
```

```
<Axes: >
```



```
In [27]: q3=df.Age.quantile(0.25)
q3=df.Age.quantile(0.75)
print(q1)
print(q2)
```

```
22.0
35.0

In [28]: IQR=q3-q1
100
```

```
Out[28]: 13.9
```

```
In [29]: upper_limit=upper_limit+1.5*IQR
upper_limit
```

```
54.5

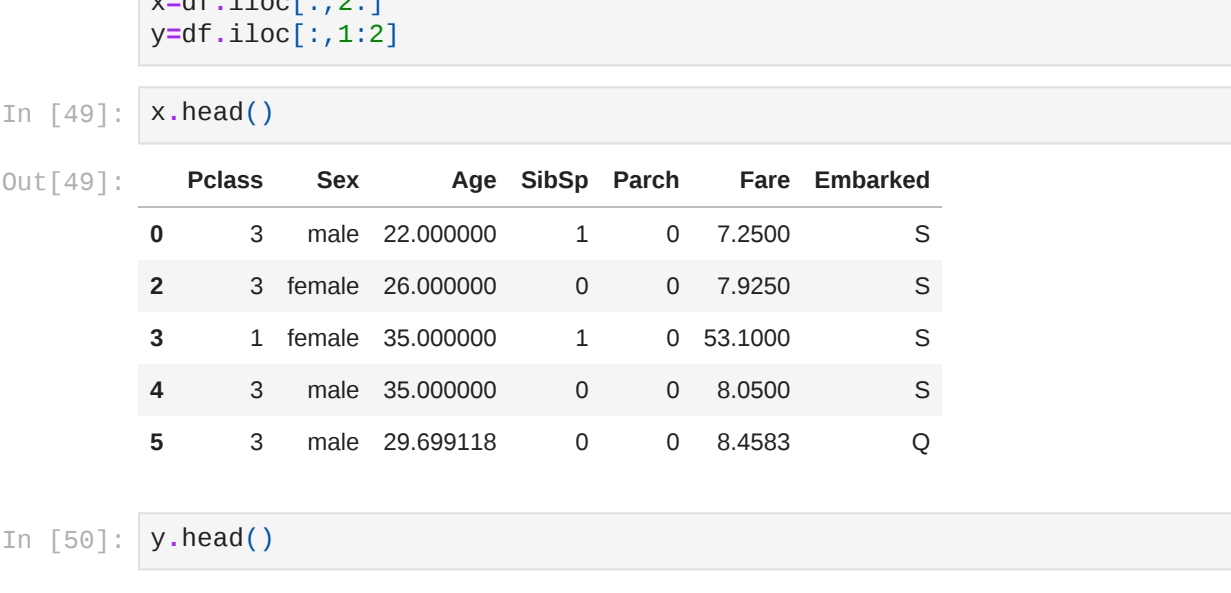
In [30]: #Lower limit
lower_limit=lower_limit-1.5*IQR
lower_limit
```

```
2.5
```

```
In [31]: df[df[(df.Age>upper_limit)&(df.Age<lower_limit)]]
```

```
In [32]: sns.boxplot(df.Age)
```

```
<Axes: >
```



```
In [33]: sns.boxplot(df.Fare)
```

```
<Axes: >
```



```
In [34]: q3=df.Fare.quantile(0.25)
q3=df.Fare.quantile(0.75)
print(q1)
print(q2)
```

```
7.2500
30.8788

In [35]: IQR=q3-q1
22.174999999999997
```

```
Out[35]: 22.174999999999997
```

```
In [36]: #Upper limit
upper_limit=upper_limit+1.5*IQR
upper_limit
```

```
63.332299999999994

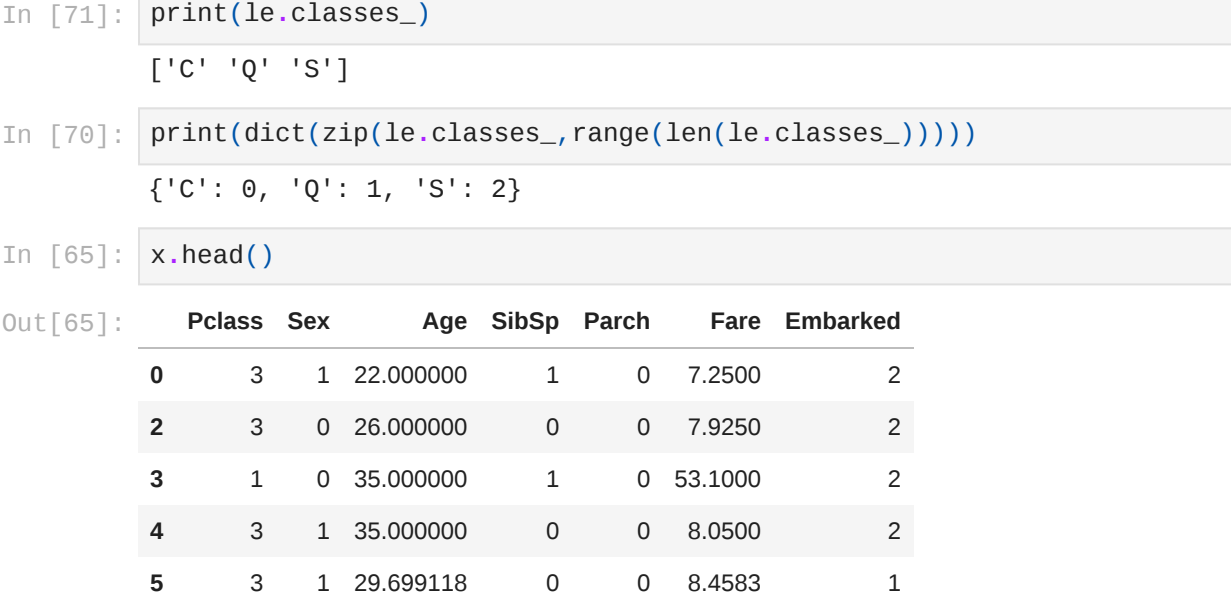
In [37]: #Lower limit
lower_limit=lower_limit-1.5*IQR
lower_limit
```

```
Out[37]: -25.366099999999994
```

```
In [38]: df[df[(df.Fare>upper_limit)&(df.Fare<lower_limit)]]
```

```
In [39]: sns.boxplot(df.Fare)
```

```
<Axes: >
```



```
In [40]: df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 718 entries, 0 to 890
Data columns (total 12 columns):
 # Column          Non-Null Count  Dtype
---  --
 0 PassengerId      718 non-null    int64
 1 Survived        718 non-null    int64
 2 Pclass          718 non-null    int64
 3 Name            718 non-null    object
 4 Sex             718 non-null    object
 5 Age            718 non-null    float64
 6 SibSp          718 non-null    int64
 7 Parch          718 non-null    int64
 8 Fare           718 non-null    float64
 9 Cabin         218 non-null    object
10 Embarked      718 non-null    object
dtypes: float64(2), int64(5), object(5)
memory usage: 89.1+ KB
```

6.Splitting Dependent and Independent Variables

```
In [41]: df.head()
```

```
In [78]: from sklearn.model_selection import train_test_split
x_train, x_test, y_train, y_test = train_test_split(x, y, test_size=0.2, random_state=0)

In [79]: x_train.shape, x_test.shape, y_train.shape, y_test.shape
Out[79]: ((574, 7), (244, 7), (574, 1), (244, 1))

In [80]: x_train.head()
Out[80]:
```

	Pclass	Sex	Age	SibSp	Parch	Fare	Embarked
593	3	0	29.699118	0	2	7.7500	1
289	3	0	22.000000	0	0	7.7500	1
12	3	0	15.000000	0	0	8.0500	1
338	3	1	45.000000	0	0	8.0500	2
451	3	1	29.699118	1	0	19.9667	2

```

In [81]: y_train.head()
Out[81]:
```

	Survived
593	0
289	1

```
In [42]: df.drop(["Name","Ticket"],axis=1,inplace=True)
```

```
In [43]: df.head()
```

	PassengerId	Survived	Pclass	Sex	Age	SibSp	Parch	Fare	Cabin	Embarked
0	1	0	3	male	22.000000	1	0	7.2500	896 B98	S
1	2	1	1	female	38.000000	1	0	7.9250	896 B98	S
2	3	1	3	female	26.000000	0	0	7.9250	896 B98	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	896 B98	S
5	6	0	3	male	29.699118	0	0	8.4583	896 B98	Q