

29-day29-simple-basic-eda

October 26, 2023

****Simple Basic EDA****

By: Loga Aswin

```
[19]: #Load the required libraries
import pandas as pd
import numpy as np
import seaborn as sns

#Load the data
df = pd.read_csv('/content/titanic.csv')

#View the data
df.head()
```

```
[19]: PassengerId  Survived  Pclass  \
0             1         0         3
1             2         1         1
2             3         1         3
3             4         1         1
4             5         0         3

                                     Name    Sex  Age  SibSp  \
0                        Braund, Mr. Owen Harris    male  22.0      1
1  Cumings, Mrs. John Bradley (Florence Briggs Th...  female  38.0      1
2                        Heikkinen, Miss. Laina  female  26.0      0
3  Futrelle, Mrs. Jacques Heath (Lily May Peel)  female  35.0      1
4                        Allen, Mr. William Henry    male  35.0      0

   Parch    Ticket   Fare Cabin Embarked
0      0   A/5 21171   7.2500   NaN        S
1      0    PC 17599  71.2833   C85        C
2      0  STON/O2. 3101282   7.9250   NaN        S
3      0    113803  53.1000  C123        S
4      0    373450   8.0500   NaN        S
```

```
[20]: #Basic information
```

```
df.info()
```

```
#Describe the data
```

```
df.describe()
```

```
<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
```

```
[20]:
```

	PassengerId	Survived	Pclass	Age	SibSp \
count	891.000000	891.000000	891.000000	714.000000	891.000000
mean	446.000000	0.383838	2.308642	29.699118	0.523008
std	257.353842	0.486592	0.836071	14.526497	1.102743
min	1.000000	0.000000	1.000000	0.420000	0.000000
25%	223.500000	0.000000	2.000000	20.125000	0.000000
50%	446.000000	0.000000	3.000000	28.000000	0.000000
75%	668.500000	1.000000	3.000000	38.000000	1.000000
max	891.000000	1.000000	3.000000	80.000000	8.000000

	Parch	Fare
count	891.000000	891.000000
mean	0.381594	32.204208
std	0.806057	49.693429
min	0.000000	0.000000
25%	0.000000	7.910400
50%	0.000000	14.454200
75%	0.000000	31.000000
max	6.000000	512.329200

```
[ ]: #Find the duplicates
```

```
df.duplicated().sum()
```

```
[ ]: 0
```

```
[ ]: #unique values

df['Pclass'].unique()

df['Survived'].unique()

df['Sex'].unique()
```

```
[ ]: array(['male', 'female'], dtype=object)
```

```
[ ]: #Find null values

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
```

```
[13]: #Replace null values

df.replace(np.nan, '0', inplace = True)

#Check the changes now
df.isnull().sum()
```

```
[13]: 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
```

```
[14]: #Datatypes
```

```
df.dtypes
```

```
[14]: PassengerId      int64
Survived             int64
Pclass              int64
Name                object
Sex                 object
Age                 object
SibSp               int64
Parch              int64
Ticket             object
Fare               float64
Cabin              object
Embarked            object
dtype: object
```

```
[15]: #Filter data
```

```
df[df['Pclass']==1].head()
```

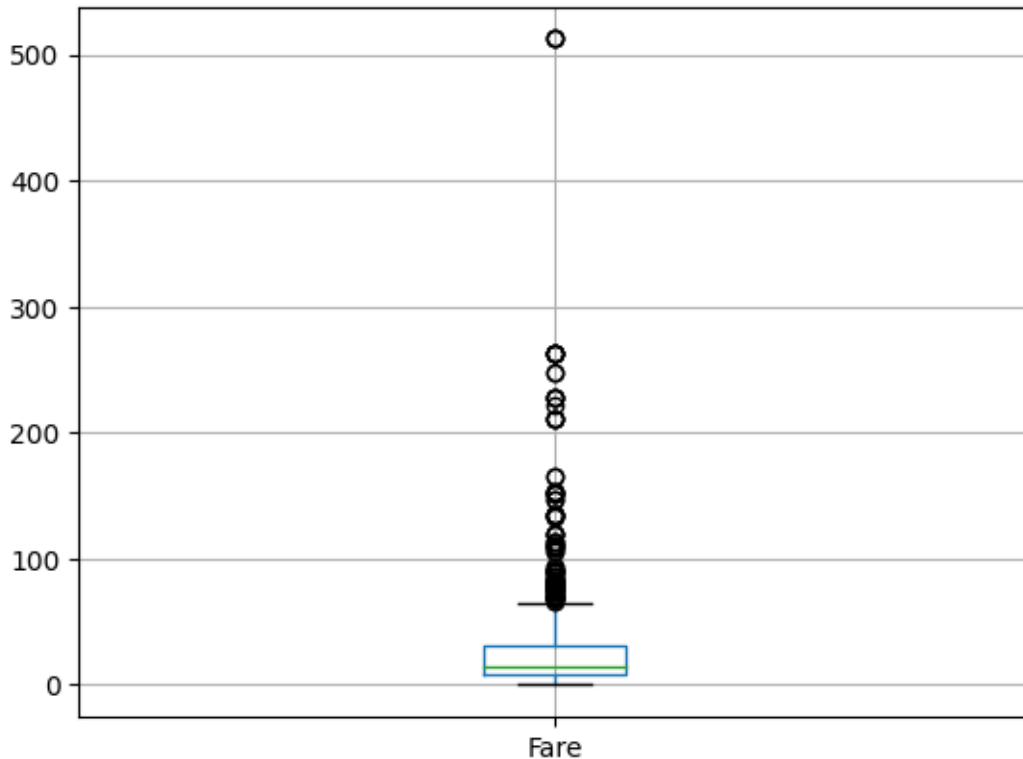
```
[15]:   PassengerId  Survived  Pclass  \
1             2         1        1
3             4         1        1
6             7         0        1
11            12         1        1
23            24         1        1
```

```
      Name      Sex  Age  SibSp  \
1  Cumings, Mrs. John Bradley (Florence Briggs Th...  female  38.0      1
3    Futrelle, Mrs. Jacques Heath (Lily May Peel)  female  35.0      1
6                McCarthy, Mr. Timothy J      male  54.0      0
11              Bonnell, Miss. Elizabeth  female  58.0      0
23            Sloper, Mr. William Thompson      male  28.0      0
```

```
      Parch  Ticket   Fare Cabin Embarked
1         0  PC 17599  71.2833   C85        C
3         0  113803  53.1000  C123        S
6         0   17463  51.8625   E46        S
11        0  113783  26.5500  C103        S
23        0  113788  35.5000   A6        S
```

```
[16]: #Boxplot
df[['Fare']].boxplot()
```

```
[16]: <Axes: >
```



```
[17]: #Correlation
df.corr()
```

<ipython-input-17-4431e0efce13>:3: 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()
```

```
[17]:
```

	PassengerId	Survived	Pclass	SibSp	Parch	Fare
PassengerId	1.000000	-0.005007	-0.035144	-0.057527	-0.001652	0.012658
Survived	-0.005007	1.000000	-0.338481	-0.035322	0.081629	0.257307
Pclass	-0.035144	-0.338481	1.000000	0.083081	0.018443	-0.549500
SibSp	-0.057527	-0.035322	0.083081	1.000000	0.414838	0.159651
Parch	-0.001652	0.081629	0.018443	0.414838	1.000000	0.216225

```
Fare          0.012658  0.257307 -0.549500  0.159651  0.216225  1.000000
```

```
[18]: #Correlation plot
```

```
sns.heatmap(df.corr())
```

<ipython-input-18-5bb0f5d05dad>:3: 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())
```

```
[18]: <Axes: >
```

