

Automated EDA - ClfAutoEDA

November 14, 2022

```
[1]: from ClfAutoEDA import *
df=pd.read_csv('titanic_train.csv')
#Dropping Id related columns
df.drop(['PassengerId','Ticket'],axis=1,inplace=True)
#Setting parameter values
labels=["not survived","survived"]
target_variable_name='Survived'
df_processed,num_features,cat_features=EDA(df,labels,
                                           target_variable_name,
                                           data_summary_figsize=(6,6),
                                           corr_matrix_figsize=(6,6),
                                           corr_matrix_annot=True,
                                           pairplt=True)
```

The data looks like this:

	Survived	Pclass	Name \
0	0	3	Braund, Mr. Owen Harris
1	1	1	Cumings, Mrs. John Bradley (Florence Briggs Th...
2	1	3	Heikkinen, Miss. Laina
3	1	1	Futrelle, Mrs. Jacques Heath (Lily May Peel)
4	0	3	Allen, Mr. William Henry

	Sex	Age	SibSp	Parch	Fare	Cabin	Embarked
0	male	22.0	1	0	7.2500	NaN	S
1	female	38.0	1	0	71.2833	C85	C
2	female	26.0	0	0	7.9250	NaN	S
3	female	35.0	1	0	53.1000	C123	S
4	male	35.0	0	0	8.0500	NaN	S

The shape of data is: (891, 10)

The missing values in data are:

Cabin	687
Age	177
Embarked	2
Survived	0
Pclass	0
Name	0

```
Sex          0
SibSp        0
Parch        0
Fare         0
dtype: int64
```

The summary of data is:

	Survived	Pclass	Age	SibSp	Parch	Fare
count	891.000000	891.000000	714.000000	891.000000	891.000000	891.000000
mean	0.383838	2.308642	29.699118	0.523008	0.381594	32.204208
std	0.486592	0.836071	14.526497	1.102743	0.806057	49.693429
min	0.000000	1.000000	0.420000	0.000000	0.000000	0.000000
25%	0.000000	2.000000	20.125000	0.000000	0.000000	7.910400
50%	0.000000	3.000000	28.000000	0.000000	0.000000	14.454200
75%	1.000000	3.000000	38.000000	1.000000	0.000000	31.000000
max	1.000000	3.000000	80.000000	8.000000	6.000000	512.329200

Some useful data information:

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 891 entries, 0 to 890
Data columns (total 10 columns):
#   Column      Non-Null Count  Dtype
---  -
0   Survived    891 non-null    int64
1   Pclass      891 non-null    int64
2   Name        891 non-null    object
3   Sex         891 non-null    object
4   Age         714 non-null    float64
5   SibSp       891 non-null    int64
6   Parch       891 non-null    int64
7   Fare        891 non-null    float64
8   Cabin       204 non-null    object
9   Embarked    889 non-null    object
dtypes: float64(2), int64(4), object(4)
memory usage: 69.7+ KB
None
```

The columns in data are:

```
['Survived' 'Pclass' 'Name' 'Sex' 'Age' 'SibSp' 'Parch' 'Fare' 'Cabin'
 'Embarked']
```

The target variable is divided into:

```
0    424
1    288
Name: Survived, dtype: int64
```

The numerical features are:

```
['Survived', 'Pclass', 'Age', 'SibSp', 'Parch', 'Fare']
```

The categorical features are:

```
['Name', 'Sex', 'Embarked']
```

The categorical variable is divided into:

```
Braund, Mr. Owen Harris      1
West, Miss. Constance Mirium 1
Palsson, Miss. Torborg Danira 1
Bonnell, Miss. Elizabeth     1
Heikkinen, Miss. Laina       1
..
Sutehall, Mr. Henry Jr       1
Rice, Mrs. William (Margaret Norton) 1
Montvila, Rev. Juozas        1
Graham, Miss. Margaret Edith 1
Dooley, Mr. Patrick          1
```

```
Name: Name, Length: 712, dtype: int64
```

The categorical variable Name has too many divisions to plot

The categorical variable is divided into:

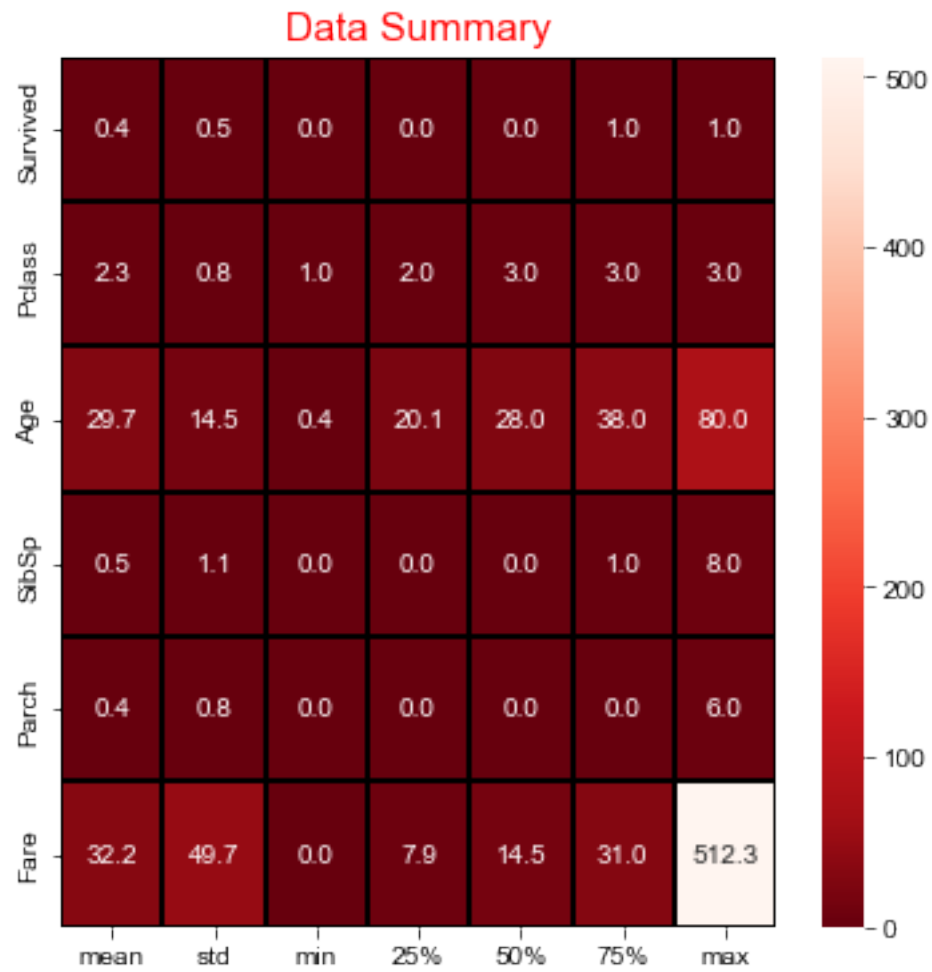
```
male      453
female    259
Name: Sex, dtype: int64
```

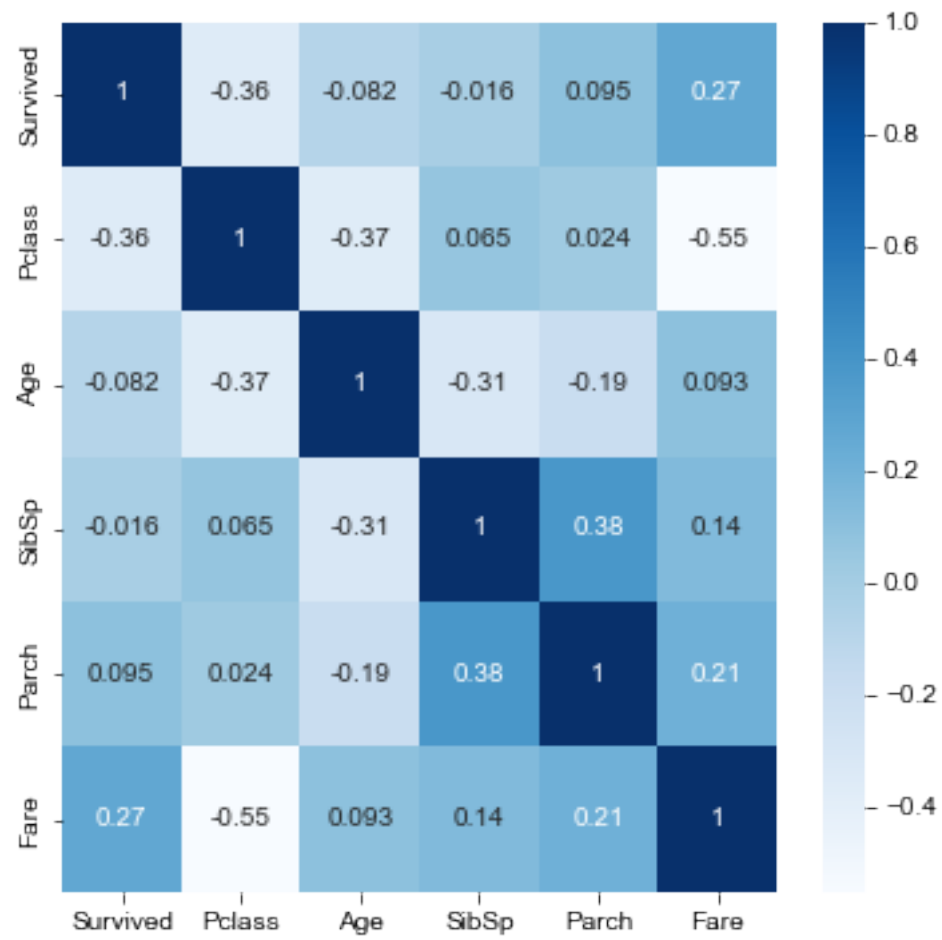
The categorical variable is divided into:

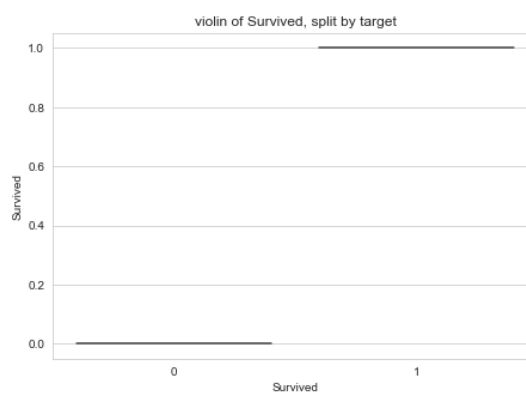
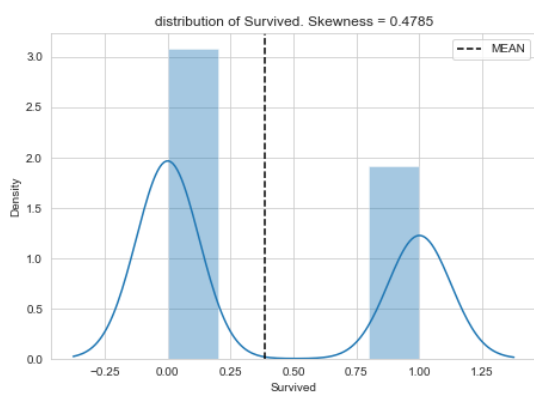
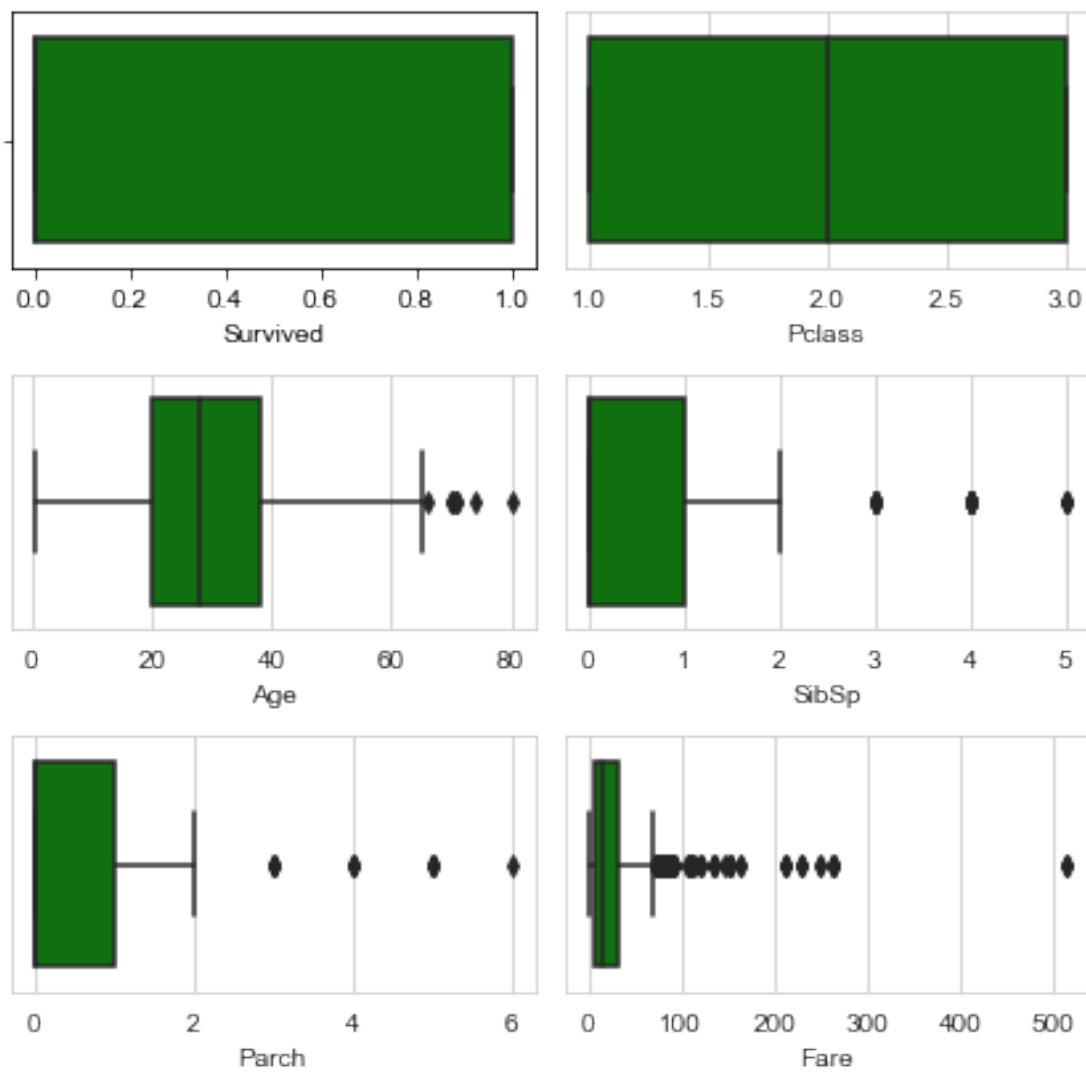
```
S      554
C      130
Q       28
Name: Embarked, dtype: int64
```

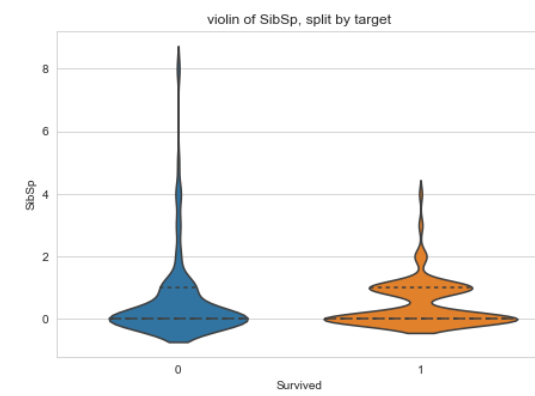
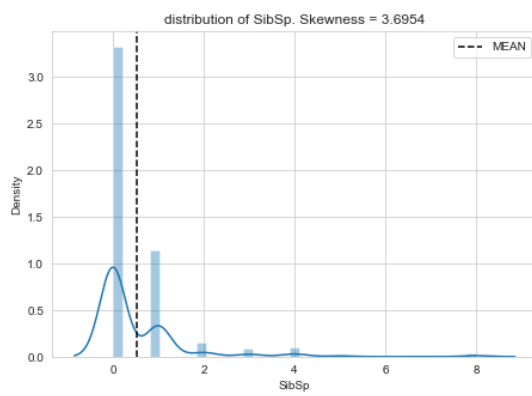
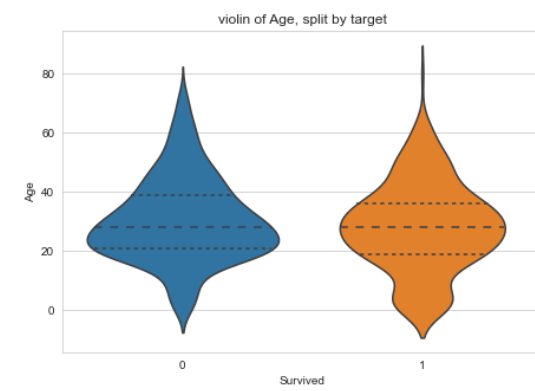
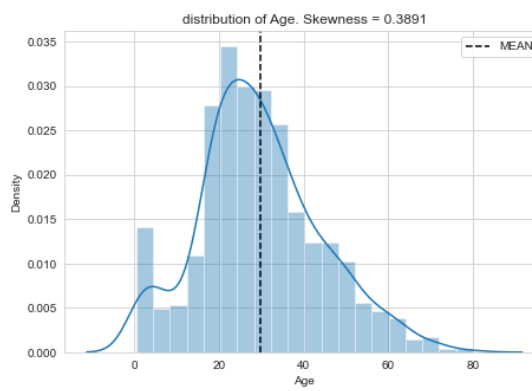
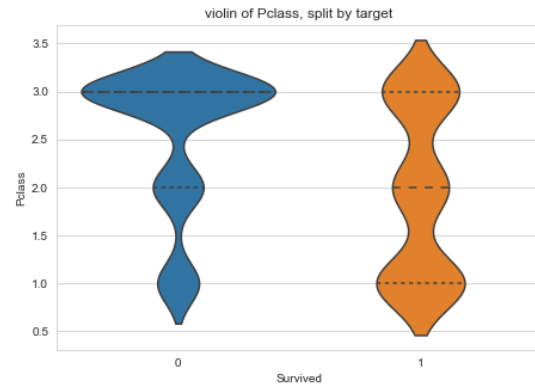
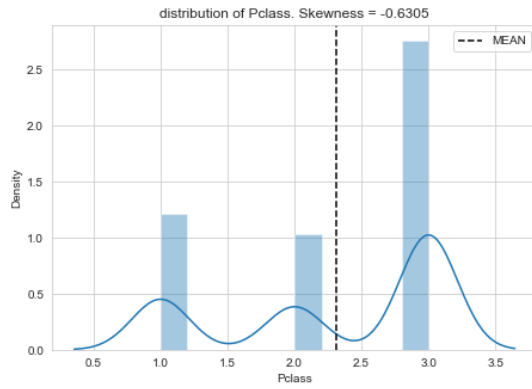
Execution Time for EDA: 0.10 minutes

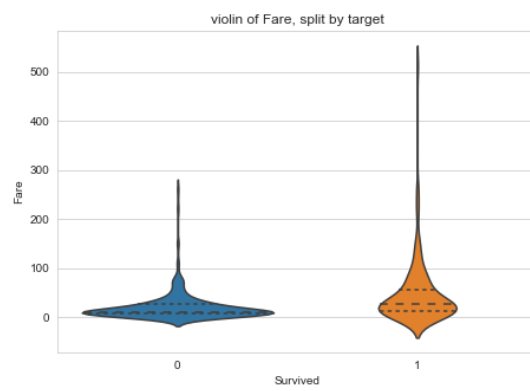
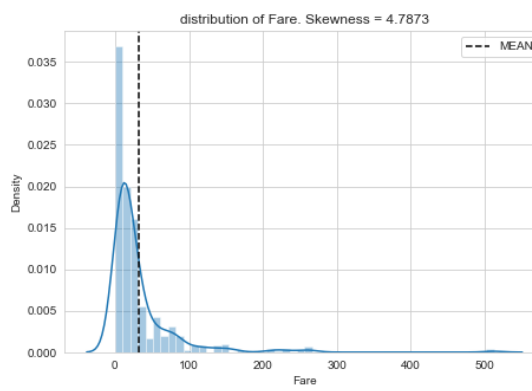
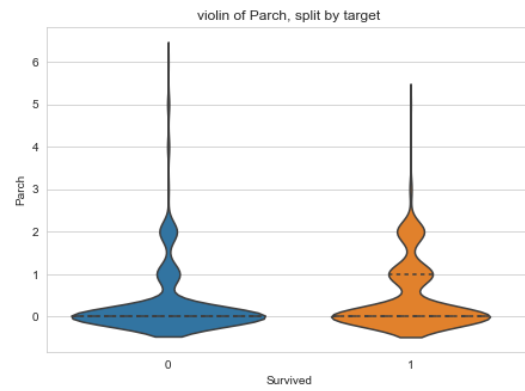
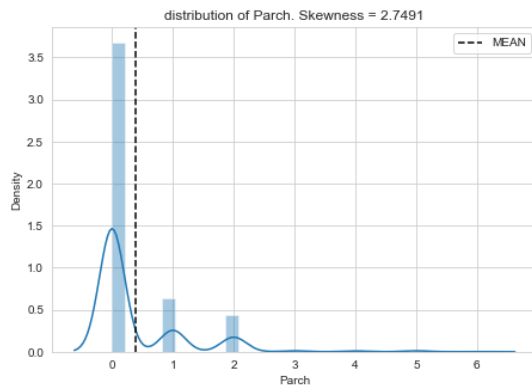


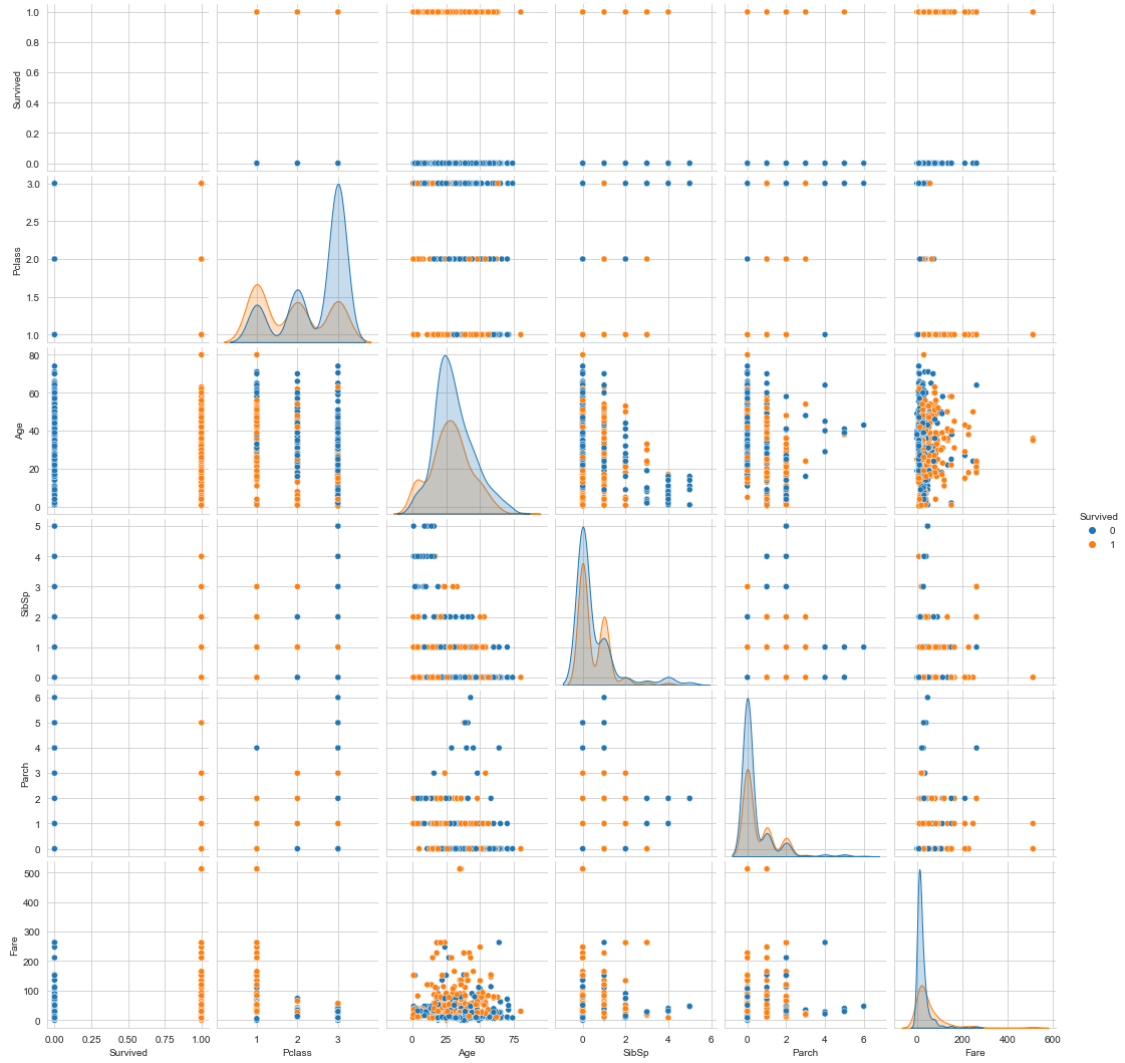


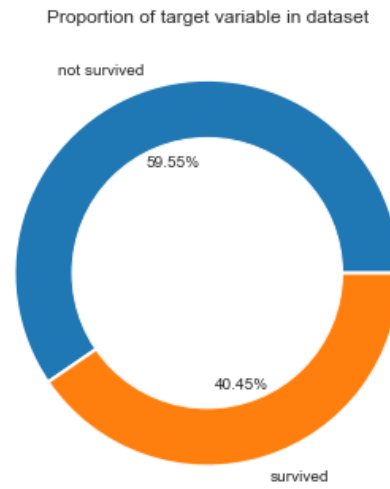
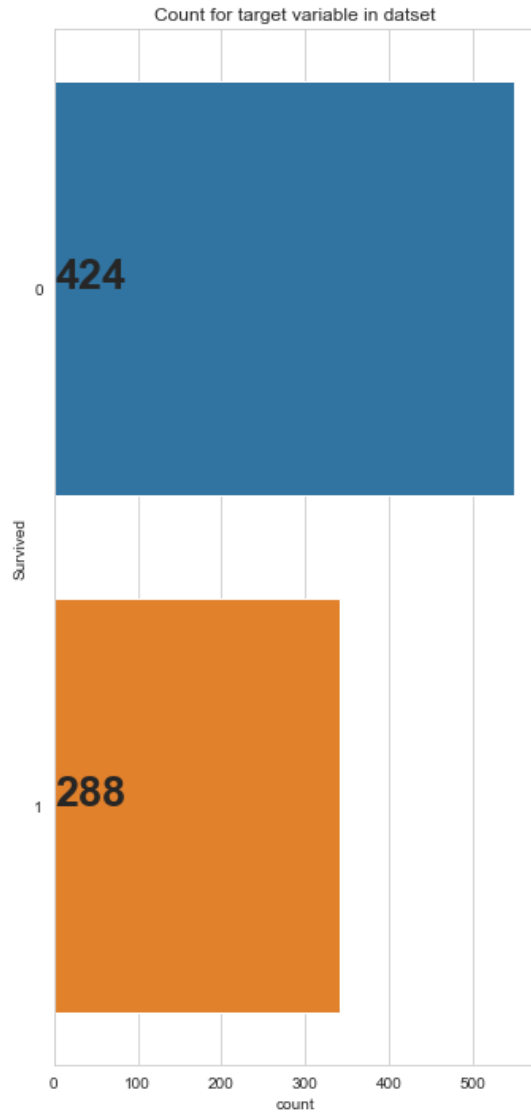


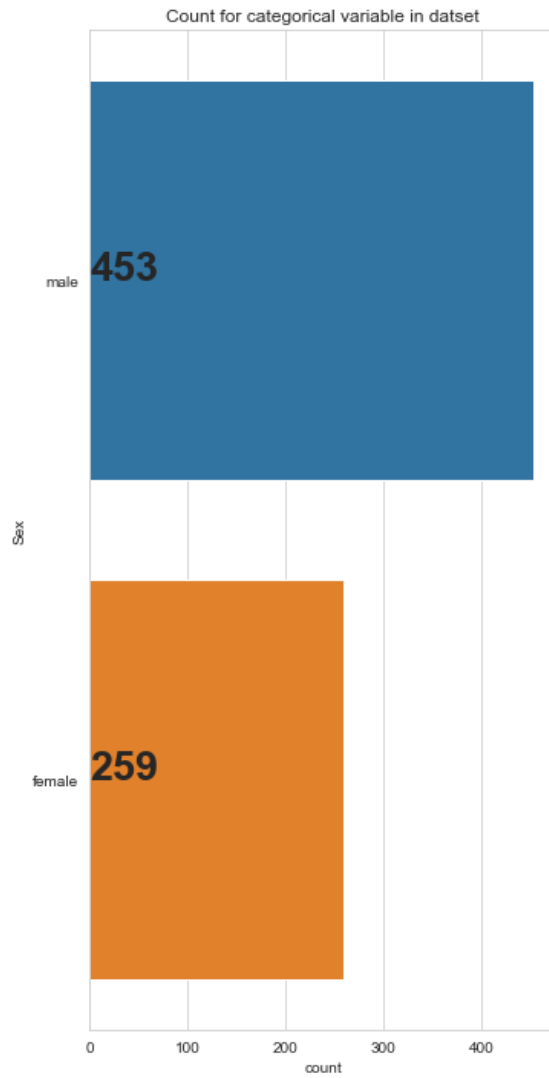




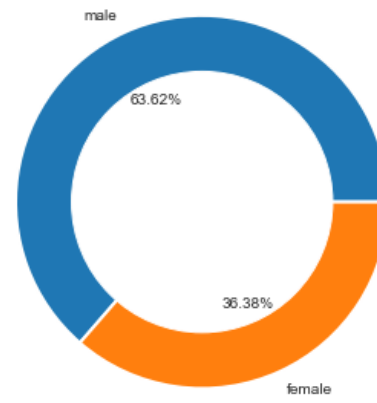


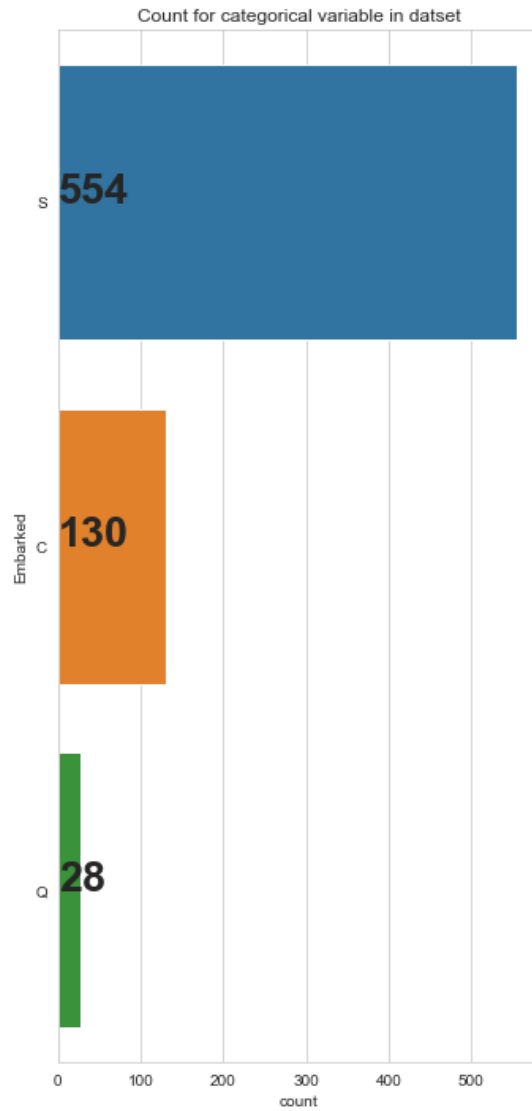




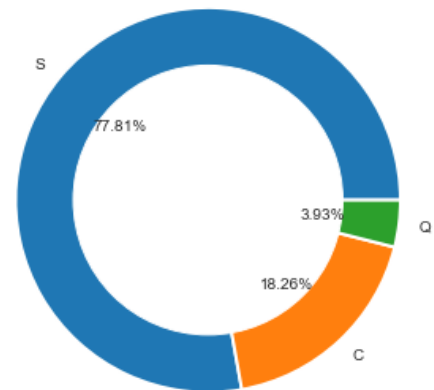


Proportion of categorical variable in dataset





Proportion of categorical variable in dataset



[]: