#### **EX-05-Feature-Generation**

### 'AIM

To read the given data and perform Feature Generation process and save the data to a file.

# **Explanation**

Feature Generation (also known as feature construction, feature extraction or feature engineering) is the process of transforming features into new features that better relate to the target.

#### 'FEATURE ENCODING:

- 1. Ordinal Encoding An ordinal encoding involves mapping each unique label to an integer value. This type of encoding is really only appropriate if there is a known relationship between the categories. This relationship does exist for some of the variables in our dataset, and ideally, this should be harnessed when preparing the data.
- 2. Label Encoding Label encoding is a simple and straight forward approach. This converts each value in a categorical column into a numerical value. Each value in a categorical column is called Label.
- 3. Binary Encoding Binary encoding converts a category into binary digits. Each binary digit creates one feature column. If there are n unique categories, then binary encoding results in the only log(base 2)<sup>n</sup> features.
- 4. One Hot Encoding We use this categorical data encoding technique when the features are nominal(do not have any order). In one hot encoding, for each level of a categorical feature, we create a new variable. Each category is mapped with a binary variable containing either 0 or 1. Here, 0 represents the absence, and 1 represents the presence of that category.

#### 'FEATURE SCALING:

- 1. Standard Scaler It is also called Z-score normalization. It calculates the z-score of each value and replaces the value with the calculated Z-score. The features are then rescaled with  $\bar{x}=0$  and  $\sigma=1$
- 2. MinMaxScaler It is also referred to as Normalization. The features are scaled between 0 and 1. Here, the mean value remains same as in Standardization, that is, 0.
- 3. Maximum absolute scaling Maximum absolute scaling scales the data to its maximum value; that is, it divides every observation by the maximum value of the variable. The result of the preceding transformation is a distribution in which the values vary approximately within the range of -1 to 1.

4. RobustScaler RobustScaler transforms the feature vector by subtracting the median and then dividing by the interquartile range (75% value — 25% value).

#### **ALGORITHM**

#### 'STEP 1

Read the given Data

#### STEP 2

Clean the Data Set using Data Cleaning Process

#### STEP 3

Apply Feature Generation techniques to all the feature of the data set

#### STEP 4

Save the data to the file

### CODE:

#### <sup>1</sup>1.FEATURE GENERATION FOR Data.csv

## CODE FOR FEATURE ENCODING AND FEATURE SCALING:

```
import pandas as pd
df=pd.read csv("data.csv")
from sklearn.preprocessing import LabelEncoder,OrdinalEncoder
temp=["Cold","Warm","Hot","Very Hot"]
enc=OrdinalEncoder(categories=[temp])
df["Ord_1"]=enc.fit_transform(df[["Ord_1"]])
df
education=["High School","Diploma","Bachelors","Masters","PhD"]
ed=OrdinalEncoder(categories=[education])
df["Ord_2"]=ed.fit_transform(df[["Ord_2"]])
df
pip install category_encoders
from category_encoders import BinaryEncoder
be=BinaryEncoder()
be.fit_transform(df["bin_1"])
df["bin_1"]=be.fit_transform(df[["bin_1"]])
df
```

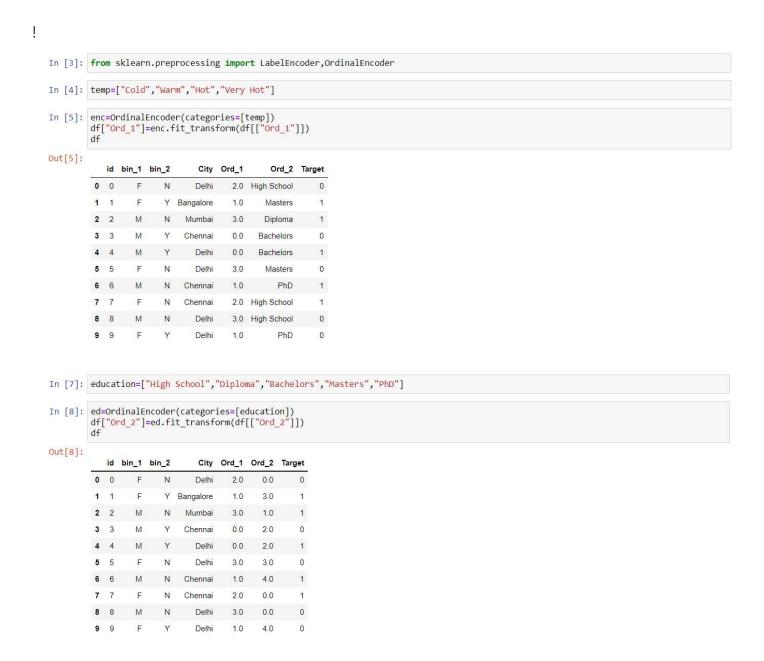
```
be1=BinaryEncoder()
be1.fit_transform(df["bin_2"])
df["bin_2"]=be1.fit_transform(df[["bin_2"]])
df
from sklearn.preprocessing import OneHotEncoder
ohe=OneHotEncoder(sparse=False)
ohe.fit_transform(df[["City"]])
from category_encoders import TargetEncoder
te=TargetEncoder()
te.fit_transform(X=df['Ord_1'],y=df["Target"])
le=LabelEncoder()
df["City"]=le.fit transform(df[["City"]])
from sklearn.preprocessing import StandardScaler
sc=StandardScaler()
df1=pd.DataFrame(sc.fit_transform(df),columns=
['id','bin_1','bin_2','City','Ord_1','Ord_2','Target'])
from sklearn.preprocessing import MaxAbsScaler
mas=MaxAbsScaler()
df2=pd.DataFrame(mas.fit_transform(df),columns=
['id','bin_1','bin_2','City','Ord_1','Ord_2','Target'])
df2
from sklearn.preprocessing import MinMaxScaler
mms=MinMaxScaler()
df3=pd.DataFrame(mms.fit_transform(df),columns=
['id','bin 1','bin 2','City','Ord 1','Ord 2','Target'])
df3
from sklearn.preprocessing import RobustScaler
rs=RobustScaler()
df4=pd.DataFrame(rs.fit_transform(df),columns=
['id','bin 1','bin 2','City','Ord 1','Ord 2','Target'])
df4
```

# **OUPUT** for Feauture Encoding and scaling:

#### 'Given data:

```
In [1]: import pandas as pd
        df=pd.read_csv("data.csv")
In [2]: df
Out[2]:
                                              Ord_2 Target
         0 0
                       Ν
                             Delhi
                                      Hot High School
                       Y Bangalore
                       N Mumbai Very Hot
                                             Diploma
                 M
                             Delhi
                                     Cold
                                           Bachelors
                             Delhi Very Hot
                                             Masters
           6
                M
                      N Chennai
                                               PhD
                                    Warm
                 M N
                             Delhi Very Hot High School
                             Delhi
                                    Warm
                                               PhD
                                                        0
```

## <sup>2</sup> Feature encoding using Ordinal Encoder:



'Feature encoding using Binary Encoder:

```
In [10]: from category_encoders import BinaryEncoder
be=BinaryEncoder()
In [11]: be.fit_transform(df["bin_1"])
Out[11]:
           bin_1_0 bin_1_1
         1
                0
         2
                      0
         3
         4
                1
                      0
         5
         6
                      0
         7
                0
                1
                      0
         9
                0
In [13]: df["bin_1"]=be.fit_transform(df[["bin_1"]])
Out[13]:
           id bin_1 bin_2
                           City Ord_1 Ord_2 Target
         0 0 0 N
                           Delhi
                                 2.0
                                       0.0
                                              0
         1 1
                0
                     Y Bangalore
                                 1.0
                                       3.0
         2 2 1 N Mumbai
                                3.0 1.0 1
         3 3
                     Y
                                              0
                1
                         Chennai
                                 0.0
                                       2.0
                                 0.0
                                           1
                1
                     Y
                           Delhi
                                      2.0
         5 5
                           Delhi
                                       3.0
                0
                     N
                                 3.0
                                              0
                     N Chennai 1.0 4.0 1
         7 7
                                       0.0
                                              1
                0
                     N Chennai
                                 2.0
                                       0.0 0
         8 8
                1
                     N
                           Delhi
                                 3.0
         9 9
                0
                           Delhi
                                 1.0
                                      4.0
                                             0
 In [14]: be1=BinaryEncoder()
 In [16]: be1.fit_transform(df["bin_2"])
Out[16]:
            bin_2_0 bin_2_1
         0
                   1
              0
         1
                1
         2
                0
                      1
                      0
         3
                1
         5
         6
                0
         7
                       1
         8
             0
                   1
In [17]: df["bin_2"]=be1.fit_transform(df[["bin_2"]])
        df
Out[17]:
           id bin_1 bin_2
                            City Ord_1 Ord_2 Target
         0 0
                0 0
                           Delhi
                                  2.0
                                       0.0
                 0
         1 1
                                  1.0
                                       3.0
                                              1
                     1 Bangalore
         2 2 1 0 Mumbai
                                 3.0
                                      1.0 1
           3
                         Chennai
                                  0.0
                                       2.0
                 1 1
         4 4
                                       2.0 1
                           Delhi
                                 0.0
         5
           5
                      0
                           Delhi
                                  3.0
                                       3.0
                                              0
         6
           6
                1 0 Chennai
                                 1.0
                                       4.0 1
         7 7
                 0
                      0 Chennai
                                 2.0
                                       0.0
                                              1
         8 8
                 1 0
                           Delhi
                                 3.0
                                       0.0
```

9 9

0 1

Delhi

1.0 4.0

0

#### Feature encoding using One Hot Encoder:

# 'Feature encoding using Target Encoder:

#### Feature encoding using Label Encoder:

## <sup>2</sup> Feature scaling using Standard Scaler:

```
In [27]: from sklearn.preprocessing import StandardScaler
        sc=StandardScaler()
        df1=pd.DataFrame(sc.fit_transform(df),columns=['id','bin_1','bin_2','City','Ord_1','Ord_2','Target'])
Out[27]:
            id bin_1 bin_2 City Ord_1 Ord_2 Target
        0 -1.566699 -1.0 -0.816497 0.50 0.359211 -1.255555 -1.0
        1 -1.218544 -1.0 1.224745 -2.00 -0.538816 0.726900
        4 -0.174078
                   1.0 1.224745 0.50 -1.436842 0.066082
        5 0 174078 -1 0 -0 816497 0.50 1.257237 0.726900 -1 0
        6 0.522233 1.0 -0.816497 -0.75 -0.538816 1.387719 1.0
        7 0.870388 -1.0 -0.816497 -0.75 0.359211 -1.255555
        8 1.218544
                   1.0 -0.816497 0.50 1.257237 -1.255555
                                                   -1.0
         9 1.566699 -1.0 1.224745 0.50 -0.538816 1.387719
```

# 'Feature scaling using MaxAbs Scaler:

```
In [29]: from sklearn.preprocessing import MaxAbsScaler
         mas=MaxAbsScaler()
         df2=pd.DataFrame(mas.fit_transform(df),columns=['id','bin_1','bin_2','City','Ord_1','Ord_2','Target'])
Out[29]:
                 id bin_1 bin_2
                                 City Ord_1 Ord_2 Target
         0 0.000000 0.0 0.0 0.666667 0.666667 0.00 0.0
         1 0.111111 0.0 1.0 0.000000 0.333333 0.75
         2 0.222222 1.0 0.0 1.000000 1.000000 0.25 1.0
          3 0.333333 1.0 1.0 0.333333 0.000000 0.50 0.0
         4 0.444444 1.0 1.0 0.666667 0.000000 0.50 1.0
                    0.0 0.0 0.666667 1.000000 0.75 0.0
         6 0.666667 1.0 0.0 0.333333 0.333333 1.00 1.0
          7 0.777778 0.0 0.0 0.333333 0.666667 0.00 1.0
         8 0.888889 1.0 0.0 0.666667 1.000000 0.00 0.0
          9 1.000000 0.0 1.0 0.666667 0.333333 1.00 0.0
```

# 'Feature scaling using MinMax Scaler:

'Feature scaling using Robust Scaler:

```
In [44]: from sklearn.preprocessing import RobustScaler
       df4=pd.DataFrame(rs.fit_transform(df),columns=['id','bin_1','bin_2','City','Ord_1','Ord_2','Target'])
Out[44]:
              id bin_1 bin_2 City
                             Ord_1
                                    Ord_2 Target
       0 -1.000000 -0.5 0.0 0.0 0.285714 -0.727273
       1 -0.777778 -0.5 1.0 -2.0 -0.285714 0.363636
       4 -0.111111 0.5
                      1.0 0.0 -0.857143 0.000000
        5 0.111111 -0.5 0.0 0.0 0.857143 0.363636
                                            -0.5
       7 0.555556 -0.5 0.0 -1.0 0.285714 -0.727273
                                            0.5
        8 0.777778  0.5  0.0  0.0  0.857143  -0.727273
                                            -0.5
```

## <sup>2</sup>2.FEATURE GENERATION FOR Encoding.csv

**9** 1.000000 -0.5 1.0 0.0 -0.285714 0.727273

#### CODE FOR FEATURE ENCODING AND FEATURE SCALING:

```
import pandas as pd
df=pd.read_csv("Encoding Data.csv")
df
from sklearn.preprocessing import LabelEncoder,OrdinalEncoder
temp=["Cold","Warm","Hot"]
enc=OrdinalEncoder(categories=[temp])
df["ord_2"]=enc.fit_transform(df[["ord_2"]])
from category encoders import BinaryEncoder
be=BinaryEncoder()
df['bin_1']=be.fit_transform(df[['bin_1']])
df
be1=BinaryEncoder()
df['bin_2']=be1.fit_transform(df[['bin_2']])
df
from sklearn.preprocessing import OneHotEncoder
ohe=OneHotEncoder(sparse=False)
ohe.fit transform(df[["nom 0"]])
le=LabelEncoder()
df["nom_0"]=le.fit_transform(df[["nom_0"]])
df
from sklearn.preprocessing import StandardScaler
sc=StandardScaler()
df1=pd.DataFrame(sc.fit_transform(df),columns=['id','bin_1','bin_2','nom_0','Ord_2'])
df1
from sklearn.preprocessing import MinMaxScaler
mms=MinMaxScaler()
df2=pd.DataFrame(mms.fit_transform(df),columns=['id','bin_1','bin_2','nom_0','Ord_2'])
df2
```

```
from sklearn.preprocessing import MaxAbsScaler mas=MaxAbsScaler()
df3=pd.DataFrame(mas.fit_transform(df),columns=['id','bin_1','bin_2','nom_0','Ord_2'])
df3

from sklearn.preprocessing import RobustScaler
rs=RobustScaler()
df4=pd.DataFrame(rs.fit_transform(df),columns=['id','bin_1','bin_2','nom_0','Ord_2'])
df4
```

#### Output for FEATURE GENERATION FOR Encoding.csv:

### 'Given DataFrame:

```
In [1]: import pandas as pd df=pd.read_csv("Encoding Data.csv")

In [2]: 

id bin_1 bin_2 nom_0 ord_2

0 0 F N Red Hot
1 1 F Y Blue Warm
2 2 F N Blue Cold
3 3 F N Green Warm
4 4 T N Red Cold
5 5 T N Green Hot
6 6 F N Red Cold
7 7 T N Red Cold
8 8 F N Blue Warm
9 9 F Y Red Hot
```

# 'Feature encoding using Ordinal Encoder:

```
In [3]: from sklearn.preprocessing import LabelEncoder,OrdinalEncoder
In [4]: temp=["Cold","Warm","Hot"]
In [5]: enc=OrdinalEncoder(categories=[temp])
In [6]: df["ord_2"]=enc.fit_transform(df[["ord_2"]])
Out[6]:
          id bin_1 bin_2 nom_0 ord_2
        0 0 F N Red
        2 2 F N
                         Blue
                             0.0
                    N Green
                               1.0
        4 4 T
                    N
                         Red
                              0.0
                         Red
                               0.0
                               0.0
        8 8 F
                         Blue
                              1.0
```

#### 'Feature encoding using Binary Encoder:

```
In [7]: from category_encoders import BinaryEncoder
In [9]: be=BinaryEncoder()
      df['bin_1']=be.fit_transform(df[['bin_1']])
Out[9]:
        id bin_1 bin_2 nom_0 ord_2
       0 0 0 N Red
             0
       2 2 0 N
                          0.0
                    Blue
                N Green
                          1.0
       4 4 1 N Red 0.0
       5 5
                          2.0
            0 N
                     Red
                          0.0
                     Red
                          0.0
                           1.0
                     Red
                          20
In [10]: be1=BinaryEncoder()
       df['bin_2']=be1.fit_transform(df[['bin_2']])
Out[10]:
         id bin_1 bin_2 nom_0 ord_2
       0 0 0 0 Red 2.0
        1 1 0 1 Blue 1.0
        2 2 0 0 Blue
        3 3
             0 0 Green
                           1.0
        4 4 1 0 Red 0.0
                           2.0
                  0 Green
        6 6 0 0 Red
                           0.0
             1 0 Red
                          0.0
        8 8 0 0 Blue 1.0
```

# 'Feature encoding using One Hot Encoder:

'Feature encoding using Label Encoder:

## 'Feature scaling using Standard Scaler:

```
In [17]: from sklearn.preprocessing import StandardScaler sc=StandardScaler() df1=pd.DataFrame(sc.fit_transform(df),columns=['id','bin_1','bin_2','nom_0','Ord_2']) df1

Out[17]: id bin_1 bin_2 nom_0 Ord_2

0 -1.566699 -0.654654 -0.5 0.917663 1.324244

1 -1.218544 -0.654654 -0.5 -1.376494 0.120386

2 -0.870388 -0.654654 -0.5 -0.229416 0.120386

4 -0.174078 1.527525 -0.5 0.917663 -1.083473

5 0.174078 1.527525 -0.5 0.917663 -1.083473

7 0.870388 1.527525 -0.5 0.917663 -1.083473

8 1.218544 -0.654654 -0.5 -1.376494 0.120386

9 1.566699 -0.654654 -0.5 -1.376494 0.120386
```

#### Feature scaling using MinMax Scaler:

#### 'Feature scaling using MaxAbs Scaler:

# 'Feature scaling using Robust Scaler:

```
In [20]: from sklearn.preprocessing import RobustScaler
rs=RobustScaler()
df4=pd.DataFrame(rs.fit_transform(df),columns=['id','bin_1','bin_2','nom_0','Ord_2'])
df4

Out[20]:

id bin_1 bin_2 nom_0 Ord_2

0 -1.000000 0.000000 0.0 0.285714 0.571429

1 -0.777778 0.000000 1.0 -0.857143 0.000000

2 -0.555556 0.000000 0.0 -0.857143 -0.571429

3 -0.333333 0.000000 0.0 -0.285714 0.000000

4 -0.111111 1.333333 0.0 0.285714 -0.571429

5 0.111111 1.333333 0.0 0.285714 -0.571429

6 0.333333 0.00000 0.0 0.285714 -0.571429

7 0.555556 1.333333 0.00 0.0 0.285714 -0.571429

8 0.777778 0.00000 0.0 0.0857143 0.000000

9 1.000000 0.000000 1.0 0.285714 0.571429
```

## 3.FEATURE GENERATION FOR titanic\_dataset.csv

#### CODE FOR FEATURE ENCODING AND FEATURE SCALING:

```
import pandas as pd

df=pd.read_csv("titanic_dataset.csv")

df

df.isnull().sum()

df["Age"]=df["Age"].fillna(df["Age"].median())

df["Embarked"]=df["Embarked"].fillna(df["Embarked"].mode()[0])

df

df.drop("Cabin",axis=1,inplace=True)

df.drop("Ticket",axis=1,inplace=True)

df.drop("Name",axis=1,inplace=True)

df.isnull().sum()

df

from sklearn.preprocessing import LabelEncoder,OrdinalEncoder

embark=["C","S","Q"]

emb=OrdinalEncoder(categories=[embark])
```

```
df["Embarked"]=emb.fit_transform(df[["Embarked"]])
df
from category_encoders import BinaryEncoder
be=BinaryEncoder()
df["Sex"]=be.fit_transform(df[["Sex"]])
df
from sklearn.preprocessing import StandardScaler
ss=StandardScaler()
df1=pd.DataFrame(ss.fit_transform(df),columns=
['Passenger', 'Survived', 'Pclass', 'Sex', 'Age', 'SibSp', 'Parch', 'Fare', 'Embarked', 'PClass'])
df1
from sklearn.preprocessing import MinMaxScaler
mms=MinMaxScaler()
df2=pd.DataFrame(mms.fit_transform(df),columns=
['Passenger','Survived','Pclass','Sex','Age','SibSp','Parch','Fare','Embarked','PClass'])
df2
from sklearn.preprocessing import MaxAbsScaler
mas=MaxAbsScaler()
df3=pd.DataFrame(mas.fit transform(df),columns=
['Passenger','Survived','Pclass','Sex','Age','SibSp','Parch','Fare','Embarked','Pclass'])
df3
from sklearn.preprocessing import RobustScaler
rs = RobustScaler()
df4=pd.DataFrame(rs.fit_transform(df),columns=
['Passenger', 'Survived', 'Pclass', 'Sex', 'Age', 'SibSp', 'Parch', 'Fare', 'Embarked', 'Pclass'])
df4
```

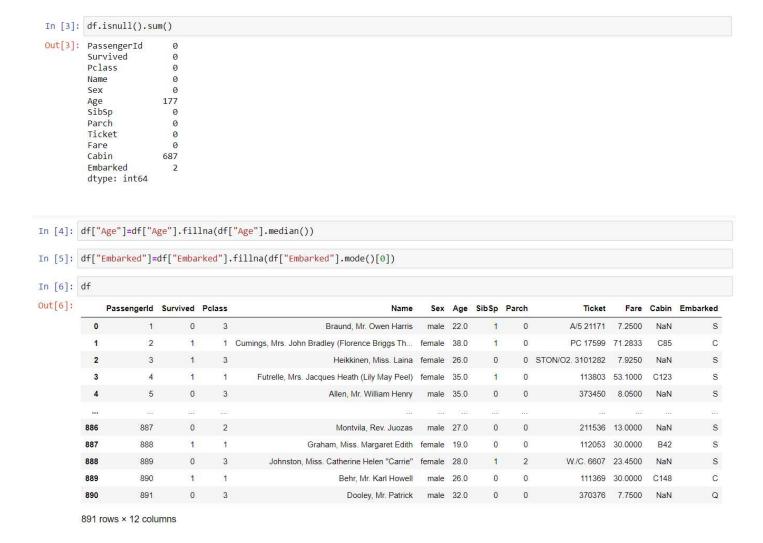
# Output for FEATURE GENERATION FOR titanic\_dataset.csv:

### 'Given DataFrame:

```
In [1]: import pandas as pd
          df=pd.read_csv("titanic_dataset.csv")
In [2]: df
Out[2]:
                Passengerld Survived Pclass
                                                                                              Sex Age SibSp Parch
                                                                                                                                   Ticket
                                                                                                                                             Fare Cabin Embarked
             0
                                                                                                                                A/5 21171
                                                                     Braund, Mr. Owen Harris
                                                                                                   22.0
                                                                                                                                           7.2500
                                                                                                                                                    NaN
                                                                                             male
                                               Cumings, Mrs. John Bradley (Florence Briggs Th...
                                                                                                                                PC 17599
                                                                                                                                          71.2833
                                                                                                                                                     C85
                                                                                                                                                                  C
             2
                                                                       Heikkinen, Miss. Laina
                                                                                                   26.0
                                                                                                                    0
                                                                                                                       STON/O2. 3101282
                                                                                                                                           7.9250
                                                                                                                                                    NaN
                                                                                                                                                                   S
                                                                                            female
             3
                                                    Futrelle, Mrs. Jacques Heath (Lily May Peel)
                                                                                            female 35.0
                                                                                                                                  113803 53.1000
                                                                                                                                                    C123
                                                                                                                                                                   S
                                                                      Allen, Mr. William Henry
                                                                                                                                  373450
                                                                                                                                           8.0500
                                                                                                                                                    NaN
           886
                        887
                                                                        Montvila, Rev. Juozas
                                                                                             male 27.0
                                                                                                             0
                                                                                                                                  211536 13.0000
                                                                                                                                                    NaN
           887
                        888
                                                                 Graham, Miss. Margaret Edith female
                                                                                                                                  112053 30 0000
                                                                                                                                                     B42
                                                                                                                                                                   S
           888
                        889
                                                        Johnston, Miss. Catherine Helen "Carrie" female NaN
                                                                                                                               W./C. 6607 23.4500
                                                                                                                                                    NaN
                                                                                                                                                                  C
           889
                        890
                                                                        Behr, Mr. Karl Howell
                                                                                             male 26.0
                                                                                                                                  111369 30.0000
                                                                                                                                                    C148
           890
                        891
                                                                          Dooley, Mr. Patrick
                                                                                                                                  370376
                                                                                                                                           7.7500
                                                                                                                                                    NaN
```

891 rows × 12 columns

#### Resolving null value data:



<sup>2</sup> Dropping unnecessary columns:

```
In [7]: df.drop("Cabin",axis=1,inplace=True)

In [8]: df.drop("Ticket",axis=1,inplace=True)

In [11]: df.drop("Name",axis=1,inplace=True)

In [12]: df.isnull().sum()

Out[12]: PassengerId 0
Survived 0
Pclass 0
Sex 0
Age 0
Sibsp 0
Parch 0
Parch 0
Fare 0
Embarked 0
dtype: int64
```

# 'Feature encoding using Ordinal Encoder:

		k=["C", <mark>"S","Q</mark> "] rdinalEncoder(categories=[embark])												
		barked"]=emb.fit_transform(df[["Embarked"]])												
df														
	Passengerld	Survived	Pclass	Sex	Age	SibSp	Parch	Fare	Embarked	PClass				
0	1	0	3	male	22.0	1	0	7.2500	1.0	2.0				
1	2	1	1	female	38.0	1	0	71.2833	0.0	0.0				
2	3	1	3	female	26.0	0	0	7.9250	1.0	2.0				
3	4	1	1	female	35.0	1	0	53.1000	1.0	0.0				
4	5	0	3	male	35.0	0	0	8.0500	1.0	2.0				
			W.		1755	1555	15%	8757		1757				
886	887	0	2	male	27.0	0	0	13.0000	1.0	1.0				
887	888	1	1	female	19.0	0	0	30.0000	1.0	0.0				
888	889	0	3	female	28.0	1	2	23.4500	1.0	2.0				
889	890	1	1	male	26.0	0	0	30.0000	0.0	0.0				
890	891	0	3	male	32.0	0	0	7.7500	2.0	2.0				

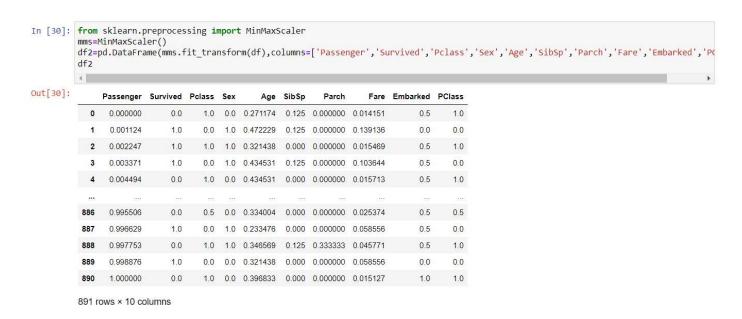
# 'Feature encoding using Binary Encoder:

L D		- //								
pe=B	inaryEncode	r()								
df["	Sex"]=be.fi	t_transf	orm(df	[["Se	x"]])	)				
df										
	Passengerld	Survived	Pclass	Sex	Age	SibSp	Parch	Fare	Embarked	PClass
0	1	0	3	0	22.0	1	0	7.2500	1.0	2.0
1	2	1	1	1	38.0	1	0	71.2833	0.0	0.0
2	3	1	3	1	26.0	0	0	7.9250	1.0	2.0
3	4	1	1	1	35.0	1	0	53.1000	1.0	0.0
4	5	0	3	0	35.0	0	0	8.0500	1.0	2.0
	411	300	100	122		227	222			
886	887	0	2	0	27.0	0	0	13.0000	1.0	1.0
887	888	1	1	1	19.0	0	0	30.0000	1.0	0.0
888	889	0	3	1	28.0	1	2	23.4500	1.0	2.0
889	890	1	1	0	26.0	0	0	30.0000	0.0	0.0
890	891	0	3	0	32.0	0	0	7.7500	2.0	2.0

#### Feature scaling using Standard Scaler:



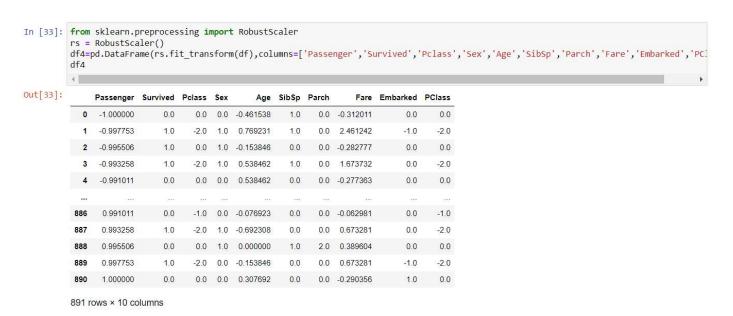
#### 'Feature scaling using MinMax Scaler



'Feature scaling using MaxAbs Scaler:

```
In [32]: from sklearn.preprocessing import MaxAbsScaler
                           mas=MaxAbsScaler()
                          df3=pd.DataFrame(mas.fit transform(df),columns=['Passenger','Survived','Pclass','Sex','Age','SibSp','Parch','Fare','Embarked','Pclass','Sex','Age','SibSp','Parch','Fare','Embarked','Pclass','Sex','Age','SibSp','Parch','Fare','Embarked','Pclass','Sex','Age','SibSp','Parch','Fare','Embarked','Pclass','Sex','Age','SibSp','Parch','Fare','Embarked','Pclass','Sex','Age','SibSp','Parch','Fare','Embarked','Pclass','Sex','Age','SibSp','Parch','Fare','Embarked','Pclass','Sex','Age','SibSp','Parch','Fare','Embarked','Pclass','Sex','Age','SibSp','Parch','Fare','Embarked','Pclass','Sex','Age','SibSp','Parch','Fare','Embarked','Pclass','Sex','Age','SibSp','Parch','Fare','Embarked','Pclass','Sex','Age','SibSp','Parch','Fare','Embarked','Pclass','Barked','Pclass','Barked','Pclass','Barked','Pclass','Barked','Pclass','Barked','Pclass','Barked','Pclass','Barked','Pclass','Barked','Pclass','Barked','Pclass','Barked','Barked','Barked','Barked','Barked','Barked','Barked','Barked','Barked','Barked','Barked','Barked','Barked','Barked','Barked','Barked','Barked','Barked','Barked','Barked','Barked','Barked','Barked','Barked','Barked','Barked','Barked','Barked','Barked','Barked','Barked','Barked','Barked','Barked','Barked','Barked','Barked','Barked','Barked','Barked','Barked','Barked','Barked','Barked','Barked','Barked','Barked','Barked','Barked','Barked','Barked','Barked','Barked','Barked','Barked','Barked','Barked','Barked','Barked','Barked','Barked','Barked','Barked','Barked','Barked','Barked','Barked','Barked','Barked','Barked','Barked','Barked','Barked','Barked','Barked','Barked','Barked','Barked','Barked','Barked','Barked','Barked','Barked','Barked','Barked','Barked','Barked','Barked','Barked','Barked','Barked','Barked','Barked','Barked','Barked','Barked','Barked','Barked','Barked','Barked','Barked','Barked','Barked','Barked','Barked','Barked','Barked','Barked','Barked','Barked','Barked','Barked','Barked','Barked','Barked','Barked','Barked','Barked','Barked','Barked','Barked','Barked','Barked','Barked','Barked','Barked','Barked','
                         4
Out[32]:
                                        Passenger Survived Pclass Sex
                                                                                                                             Age SibSp Parch
                                                                                                                                                                                                Fare Embarked PClass
                                          0.001122 0.0 1.000000 0.0 0.2750 0.125 0.000000 0.014151 0.5
                                                                    1.0 0.333333 1.0 0.4750 0.125 0.000000 0.139136
                                         0.003367 1.0 1.000000 1.0 0.3250 0.000 0.000000 0.015469
                                                                                                                                                                                                                         0.5 1.0
                                           0.004489 1.0 0.333333 1.0 0.4375 0.125 0.000000 0.103644
                                                                                                                                                                                                                                                  0.0
                                            0.0 0.666667 0.0 0.3375 0.000 0.000000 0.025374 0.5 0.5
                                                                                                                                                                                                                             0.5 0.0
                             887
                                           0.996633
                                                                              1.0 0.333333 1.0 0.2375 0.000 0.000000 0.058556
                                                                     0.0 1.000000 1.0 0.3500 0.125 0.333333 0.045771
                                                                                                                                                                                                                            0.5
                                            0.997755
                                                                                                                                                                                                                                               1.0
                                            0.998878
                                                                              1.0 0.333333 0.0 0.3250 0.000 0.000000 0.058556
                                                                                                                                                                                                                              0.0
                                                                                                                                                                                                                                                 0.0
                             889
                                           891 rows × 10 columns
```

#### Feature scaling using Robust Scaler



### RESULT:

Feature Encoding process and Feature Scaling process is applied to the given data frame sucessfully.