# Project_Mercedes-Benz Greener Manufacturing

May 25, 2022

## 1 Mercedes-Benz Greener Manufacturing

```
[1]: # Import Require librarary

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

```
[2]: import warnings
warnings.filterwarnings('ignore')
```

### 1.0.1 Load the Train and Test dataset

```
[3]: # Load the "Train.csv" file

data_train = pd.read_csv('train.csv')
data_train.head()
```

```
[3]:    ID       y X0 X1  X2 X3 X4 X5 X6 X8  …  X375  X376  X377  X378  X379  \
    0   0  130.81  k  v  at  a  d  u  j  o  …     0     0     1     0     0
    1   6   88.53  k  t  av  e  d  y  l  o  …     1     0     0     0     0
    2   7   76.26 az  w   n  c  d  x  j  x  …     0     0     0     0     0
    3   9   80.62 az  t   n  f  d  x  l  e  …     0     0     0     0     0
    4  13   78.02 az  v   n  f  d  h  d  n  …     0     0     0     0     0

       X380  X382  X383  X384  X385
    0     0     0     0     0     0
    1     0     0     0     0     0
    2     0     1     0     0     0
    3     0     0     0     0     0
    4     0     0     0     0     0

    [5 rows x 378 columns]
```

```
[4]: # Load the "Test.csv" file
```

```python
data_test = pd.read_csv('test.csv')
data_test.head()
```

[4]:
```
    ID  X0 X1  X2 X3 X4 X5 X6 X8  X10 …  X375  X376  X377  X378  X379  X380  \
0    1  az  v   n  f  d  t  a  w    0 …     0     0     0     1     0     0
1    2   t  b  ai  a  d  b  g  y    0 …     0     0     1     0     0     0
2    3  az  v  as  f  d  a  j  j    0 …     0     0     0     1     0     0
3    4  az  l   n  f  d  z  l  n    0 …     0     0     0     1     0     0
4    5   w  s  as  c  d  y  i  m    0 …     1     0     0     0     0     0

   X382  X383  X384  X385
0     0     0     0     0
1     0     0     0     0
2     0     0     0     0
3     0     0     0     0
4     0     0     0     0

[5 rows x 377 columns]
```

[5]:
```python
print('Train Data Size:',data_train.shape)
print('Test Data Size:',data_test.shape)
```

```
Train Data Size: (4209, 378)
Test Data Size: (4209, 377)
```

### 1.0.2 Exploratory Data Analysis:

[6]:
```python
# Check the columns in Train dataset

data_train.columns
```

[6]:
```
Index(['ID', 'y', 'X0', 'X1', 'X2', 'X3', 'X4', 'X5', 'X6', 'X8',
       …
       'X375', 'X376', 'X377', 'X378', 'X379', 'X380', 'X382', 'X383', 'X384',
       'X385'],
      dtype='object', length=378)
```

[7]:
```python
# Check the columns in Test dataset

data_test.columns
```

[7]:
```
Index(['ID', 'X0', 'X1', 'X2', 'X3', 'X4', 'X5', 'X6', 'X8', 'X10',
       …
       'X375', 'X376', 'X377', 'X378', 'X379', 'X380', 'X382', 'X383', 'X384',
       'X385'],
      dtype='object', length=377)
```

```
[8]:  # Check the information of train dataset

      data_train.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 4209 entries, 0 to 4208
Columns: 378 entries, ID to X385
dtypes: float64(1), int64(369), object(8)
memory usage: 12.1+ MB
```

```
[9]:  # Check the information of test dataset

      data_test.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 4209 entries, 0 to 4208
Columns: 377 entries, ID to X385
dtypes: int64(369), object(8)
memory usage: 12.1+ MB
```

```
[10]:  data_train.describe()
```

[10]:
|       | ID          | y           | X10         | X11    | X12         | \ |
|-------|-------------|-------------|-------------|--------|-------------|---|
| count | 4209.000000 | 4209.000000 | 4209.000000 | 4209.0 | 4209.000000 |   |
| mean  | 4205.960798 | 100.669318  | 0.013305    | 0.0    | 0.075077    |   |
| std   | 2437.608688 | 12.679381   | 0.114590    | 0.0    | 0.263547    |   |
| min   | 0.000000    | 72.110000   | 0.000000    | 0.0    | 0.000000    |   |
| 25%   | 2095.000000 | 90.820000   | 0.000000    | 0.0    | 0.000000    |   |
| 50%   | 4220.000000 | 99.150000   | 0.000000    | 0.0    | 0.000000    |   |
| 75%   | 6314.000000 | 109.010000  | 0.000000    | 0.0    | 0.000000    |   |
| max   | 8417.000000 | 265.320000  | 1.000000    | 0.0    | 1.000000    |   |

|       | X13         | X14         | X15         | X16         | X17         | … | \ |
|-------|-------------|-------------|-------------|-------------|-------------|---|---|
| count | 4209.000000 | 4209.000000 | 4209.000000 | 4209.000000 | 4209.000000 | … |   |
| mean  | 0.057971    | 0.428130    | 0.000475    | 0.002613    | 0.007603    | … |   |
| std   | 0.233716    | 0.494867    | 0.021796    | 0.051061    | 0.086872    | … |   |
| min   | 0.000000    | 0.000000    | 0.000000    | 0.000000    | 0.000000    | … |   |
| 25%   | 0.000000    | 0.000000    | 0.000000    | 0.000000    | 0.000000    | … |   |
| 50%   | 0.000000    | 0.000000    | 0.000000    | 0.000000    | 0.000000    | … |   |
| 75%   | 0.000000    | 1.000000    | 0.000000    | 0.000000    | 0.000000    | … |   |
| max   | 1.000000    | 1.000000    | 1.000000    | 1.000000    | 1.000000    | … |   |

|       | X375        | X376        | X377        | X378        | X379        | \ |
|-------|-------------|-------------|-------------|-------------|-------------|---|
| count | 4209.000000 | 4209.000000 | 4209.000000 | 4209.000000 | 4209.000000 |   |
| mean  | 0.318841    | 0.057258    | 0.314802    | 0.020670    | 0.009503    |   |
| std   | 0.466082    | 0.232363    | 0.464492    | 0.142294    | 0.097033    |   |
| min   | 0.000000    | 0.000000    | 0.000000    | 0.000000    | 0.000000    |   |

|      | 25%      |          |          |          |          |
|------|----------|----------|----------|----------|----------|
| 25%  | 0.000000 | 0.000000 | 0.000000 | 0.000000 | 0.000000 |
| 50%  | 0.000000 | 0.000000 | 0.000000 | 0.000000 | 0.000000 |
| 75%  | 1.000000 | 0.000000 | 1.000000 | 0.000000 | 0.000000 |
| max  | 1.000000 | 1.000000 | 1.000000 | 1.000000 | 1.000000 |

|       | X380        | X382        | X383        | X384        | X385        |
|-------|-------------|-------------|-------------|-------------|-------------|
| count | 4209.000000 | 4209.000000 | 4209.000000 | 4209.000000 | 4209.000000 |
| mean  | 0.008078    | 0.007603    | 0.001663    | 0.000475    | 0.001426    |
| std   | 0.089524    | 0.086872    | 0.040752    | 0.021796    | 0.037734    |
| min   | 0.000000    | 0.000000    | 0.000000    | 0.000000    | 0.000000    |
| 25%   | 0.000000    | 0.000000    | 0.000000    | 0.000000    | 0.000000    |
| 50%   | 0.000000    | 0.000000    | 0.000000    | 0.000000    | 0.000000    |
| 75%   | 0.000000    | 0.000000    | 0.000000    | 0.000000    | 0.000000    |
| max   | 1.000000    | 1.000000    | 1.000000    | 1.000000    | 1.000000    |

[8 rows x 370 columns]

```python
[11]: data_test.describe()
```

```
[11]:
```

|       | ID          | X10         | X11         | X12         | X13         | \ |
|-------|-------------|-------------|-------------|-------------|-------------|---|
| count | 4209.000000 | 4209.000000 | 4209.000000 | 4209.000000 | 4209.000000 |   |
| mean  | 4211.039202 | 0.019007    | 0.000238    | 0.074364    | 0.061060    |   |
| std   | 2423.078926 | 0.136565    | 0.015414    | 0.262394    | 0.239468    |   |
| min   | 1.000000    | 0.000000    | 0.000000    | 0.000000    | 0.000000    |   |
| 25%   | 2115.000000 | 0.000000    | 0.000000    | 0.000000    | 0.000000    |   |
| 50%   | 4202.000000 | 0.000000    | 0.000000    | 0.000000    | 0.000000    |   |
| 75%   | 6310.000000 | 0.000000    | 0.000000    | 0.000000    | 0.000000    |   |
| max   | 8416.000000 | 1.000000    | 1.000000    | 1.000000    | 1.000000    |   |

|       | X14         | X15         | X16         | X17         | X18         | … | \ |
|-------|-------------|-------------|-------------|-------------|-------------|---|---|
| count | 4209.000000 | 4209.000000 | 4209.000000 | 4209.000000 | 4209.000000 | … |   |
| mean  | 0.427893    | 0.000713    | 0.002613    | 0.008791    | 0.010216    | … |   |
| std   | 0.494832    | 0.026691    | 0.051061    | 0.093357    | 0.100570    | … |   |
| min   | 0.000000    | 0.000000    | 0.000000    | 0.000000    | 0.000000    | … |   |
| 25%   | 0.000000    | 0.000000    | 0.000000    | 0.000000    | 0.000000    | … |   |
| 50%   | 0.000000    | 0.000000    | 0.000000    | 0.000000    | 0.000000    | … |   |
| 75%   | 1.000000    | 0.000000    | 0.000000    | 0.000000    | 0.000000    | … |   |
| max   | 1.000000    | 1.000000    | 1.000000    | 1.000000    | 1.000000    | … |   |

|       | X375        | X376        | X377        | X378        | X379        | \ |
|-------|-------------|-------------|-------------|-------------|-------------|---|
| count | 4209.000000 | 4209.000000 | 4209.000000 | 4209.000000 | 4209.000000 |   |
| mean  | 0.325968    | 0.049656    | 0.311951    | 0.019244    | 0.011879    |   |
| std   | 0.468791    | 0.217258    | 0.463345    | 0.137399    | 0.108356    |   |
| min   | 0.000000    | 0.000000    | 0.000000    | 0.000000    | 0.000000    |   |
| 25%   | 0.000000    | 0.000000    | 0.000000    | 0.000000    | 0.000000    |   |
| 50%   | 0.000000    | 0.000000    | 0.000000    | 0.000000    | 0.000000    |   |
| 75%   | 1.000000    | 0.000000    | 1.000000    | 0.000000    | 0.000000    |   |

```
max         1.000000        1.000000        1.000000        1.000000        1.000000

                  X380            X382            X383            X384            X385
count    4209.000000     4209.000000     4209.000000     4209.000000     4209.000000
mean        0.008078        0.008791        0.000475        0.000713        0.001663
std         0.089524        0.093357        0.021796        0.026691        0.040752
min         0.000000        0.000000        0.000000        0.000000        0.000000
25%         0.000000        0.000000        0.000000        0.000000        0.000000
50%         0.000000        0.000000        0.000000        0.000000        0.000000
75%         0.000000        0.000000        0.000000        0.000000        0.000000
max         1.000000        1.000000        1.000000        1.000000        1.000000

[8 rows x 369 columns]
```

# 2   1. If for any column(s), the variance is equal to zero, then you need to remove those variable(s).

```
[12]: # Check the variance of train dataset

      data_train.var()
```

```
[12]: ID       5.941936e+06
      y        1.607667e+02
      X10      1.313092e-02
      X11      0.000000e+00
      X12      6.945713e-02

                   ...
      X380     8.014579e-03
      X382     7.546747e-03
      X383     1.660732e-03
      X384     4.750593e-04
      X385     1.423823e-03
      Length: 370, dtype: float64
```

```
[13]: # Check the variance of test dataset

      data_test.var()
```

```
[13]: ID       5.871311e+06
      X10      1.865006e-02
      X11      2.375861e-04
      X12      6.885074e-02
      X13      5.734498e-02

                   ...
      X380     8.014579e-03
      X382     8.715481e-03
```

```
X383    4.750593e-04
X384    7.124196e-04
X385    1.660732e-03
Length: 369, dtype: float64
```

[14]: ```python
# Check variance is equal to zero variables from train dataset

data_train.var()[data_train.var() == 0].index.values
```

[14]: ```
array(['X11', 'X93', 'X107', 'X233', 'X235', 'X268', 'X289', 'X290',
       'X293', 'X297', 'X330', 'X347'], dtype=object)
```

[15]: ```python
# Check variance is equal to zero variables from test dataset

data_test.var()[data_test.var() == 0].index.values
```

[15]: ```
array(['X257', 'X258', 'X295', 'X296', 'X369'], dtype=object)
```

[16]: ```python
# Removed the variance is equal to zero variables from train dataset

data_train = data_train.drop(['X11', 'X93', 'X107', 'X233', 'X235', 'X268',
 ↪'X289', 'X290',
                              'X293', 'X297', 'X330', 'X347'],axis=1)
data_train.head()
```

[16]:
```
     ID       y  X0 X1  X2 X3 X4 X5 X6 X8  …   X375  X376  X377  X378  X379  \
0    0  130.81   k  v  at  a  d  u  j  o  …      0     0     1     0     0
1    6   88.53   k  t  av  e  d  y  l  o  …      1     0     0     0     0
2    7   76.26  az  w   n  c  d  x  j  x  …      0     0     0     0     0
3    9   80.62  az  t   n  f  d  x  l  e  …      0     0     0     0     0
4   13   78.02  az  v   n  f  d  h  d  n  …      0     0     0     0     0

   X380  X382  X383  X384  X385
0     0     0     0     0     0
1     0     0     0     0     0
2     0     1     0     0     0
3     0     0     0     0     0
4     0     0     0     0     0

[5 rows x 366 columns]
```

[17]: ```python
# Removed the variance is equal to zero variables from test datase

data_test = data_test.drop(['X11', 'X93', 'X107', 'X233', 'X235', 'X268',
 ↪'X289', 'X290',
                            'X293', 'X297', 'X330', 'X347'],axis=1)
data_test.head()
```

```
[17]:     ID  X0 X1  X2 X3 X4 X5 X6 X8  X10  …   X375  X376  X377  X378  X379  X380  \
      0    1  az  v   n  f  d  t  a  w    0  …      0     0     0     1     0     0
      1    2   t  b  ai  a  d  b  g  y    0  …      0     0     1     0     0     0
      2    3  az  v  as  f  d  a  j  j    0  …      0     0     0     1     0     0
      3    4  az  l   n  f  d  z  l  n    0  …      0     0     0     1     0     0
      4    5   w  s  as  c  d  y  i  m    0  …      1     0     0     0     0     0

         X382  X383  X384  X385
      0     0     0     0     0
      1     0     0     0     0
      2     0     0     0     0
      3     0     0     0     0
      4     0     0     0     0

      [5 rows x 365 columns]
```

```
[18]:  # After removal variance is equal to zero check the shape of the train dataset

       data_train.shape
```

```
[18]:  (4209, 366)
```

```
[19]:  # After removal variance is equal to zero check the shape of the test dataset

       data_test.shape
```

```
[19]:  (4209, 365)
```

```
[20]:  data_train.info()
```

```
       <class 'pandas.core.frame.DataFrame'>
       RangeIndex: 4209 entries, 0 to 4208
       Columns: 366 entries, ID to X385
       dtypes: float64(1), int64(357), object(8)
       memory usage: 11.8+ MB
```

# 3  2. Check for null and unique values for test and train sets.

```
[21]:  # Check the unique value of train dataset.

       for columns in data_train:
           print('Train data Unique Values:',columns,data_train[columns].unique(),
                 'Train data Unique Value Shape:',data_train[columns].unique().shape)
```

```
       Train data Unique Values: ID [   0    6    7 … 8412 8415 8417] Train data
       Unique Value Shape: (4209,)
       Train data Unique Values: y [130.81  88.53  76.26 …  85.71 108.77  87.48]
```

```
Train data Unique Value Shape: (2545,)
Train data Unique Values: X0 ['k' 'az' 't' 'al' 'o' 'w' 'j' 'h' 's' 'n' 'ay' 'f'
'x' 'y' 'aj' 'ak' 'am'
 'z' 'q' 'at' 'ap' 'v' 'af' 'a' 'e' 'ai' 'd' 'aq' 'c' 'aa' 'ba' 'as' 'i'
 'r' 'b' 'ax' 'bc' 'u' 'ad' 'au' 'm' 'l' 'aw' 'ao' 'ac' 'g' 'ab'] Train data
Unique Value Shape: (47,)
Train data Unique Values: X1 ['v' 't' 'w' 'b' 'r' 'l' 's' 'aa' 'c' 'a' 'e' 'h'
'z' 'j' 'o' 'u' 'p' 'n'
 'i' 'y' 'd' 'f' 'm' 'k' 'g' 'q' 'ab'] Train data Unique Value Shape: (27,)
Train data Unique Values: X2 ['at' 'av' 'n' 'e' 'as' 'aq' 'r' 'ai' 'ak' 'm' 'a'
'k' 'ae' 's' 'f' 'd'
 'ag' 'ay' 'ac' 'ap' 'g' 'i' 'aw' 'y' 'b' 'ao' 'al' 'h' 'x' 'au' 't' 'an'
 'z' 'ah' 'p' 'am' 'j' 'q' 'af' 'l' 'aa' 'c' 'o' 'ar'] Train data Unique Value
Shape: (44,)
Train data Unique Values: X3 ['a' 'e' 'c' 'f' 'd' 'b' 'g'] Train data Unique
Value Shape: (7,)
Train data Unique Values: X4 ['d' 'b' 'c' 'a'] Train data Unique Value Shape:
(4,)
Train data Unique Values: X5 ['u' 'y' 'x' 'h' 'g' 'f' 'j' 'i' 'd' 'c' 'af' 'ag'
'ab' 'ac' 'ad' 'ae'
 'ah' 'l' 'k' 'n' 'm' 'p' 'q' 's' 'r' 'v' 'w' 'o' 'aa'] Train data Unique Value
Shape: (29,)
Train data Unique Values: X6 ['j' 'l' 'd' 'h' 'i' 'a' 'g' 'c' 'k' 'e' 'f' 'b']
Train data Unique Value Shape: (12,)
Train data Unique Values: X8 ['o' 'x' 'e' 'n' 's' 'a' 'h' 'p' 'm' 'k' 'd' 'i'
'v' 'j' 'b' 'q' 'w' 'g'
 'y' 'l' 'f' 'u' 'r' 't' 'c'] Train data Unique Value Shape: (25,)
Train data Unique Values: X10 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X12 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X13 [1 0] Train data Unique Value Shape: (2,)
Train data Unique Values: X14 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X15 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X16 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X17 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X18 [1 0] Train data Unique Value Shape: (2,)
Train data Unique Values: X19 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X20 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X21 [1 0] Train data Unique Value Shape: (2,)
Train data Unique Values: X22 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X23 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X24 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X26 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X27 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X28 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X29 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X30 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X31 [1 0] Train data Unique Value Shape: (2,)
Train data Unique Values: X32 [0 1] Train data Unique Value Shape: (2,)
```

```
Train data Unique Values: X33 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X34 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X35 [1 0] Train data Unique Value Shape: (2,)
Train data Unique Values: X36 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X37 [1 0] Train data Unique Value Shape: (2,)
Train data Unique Values: X38 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X39 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X40 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X41 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X42 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X43 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X44 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X45 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X46 [1 0] Train data Unique Value Shape: (2,)
Train data Unique Values: X47 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X48 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X49 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X50 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X51 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X52 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X53 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X54 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X55 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X56 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X57 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X58 [1 0] Train data Unique Value Shape: (2,)
Train data Unique Values: X59 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X60 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X61 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X62 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X63 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X64 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X65 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X66 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X67 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X68 [1 0] Train data Unique Value Shape: (2,)
Train data Unique Values: X69 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X70 [1 0] Train data Unique Value Shape: (2,)
Train data Unique Values: X71 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X73 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X74 [1 0] Train data Unique Value Shape: (2,)
Train data Unique Values: X75 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X76 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X77 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X78 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X79 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X80 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X81 [0 1] Train data Unique Value Shape: (2,)
```

```
Train data Unique Values: X82 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X83 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X84 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X85 [1 0] Train data Unique Value Shape: (2,)
Train data Unique Values: X86 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X87 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X88 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X89 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X90 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X91 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X92 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X94 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X95 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X96 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X97 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X98 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X99 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X100 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X101 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X102 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X103 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X104 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X105 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X106 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X108 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X109 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X110 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X111 [1 0] Train data Unique Value Shape: (2,)
Train data Unique Values: X112 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X113 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X114 [1 0] Train data Unique Value Shape: (2,)
Train data Unique Values: X115 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X116 [1 0] Train data Unique Value Shape: (2,)
Train data Unique Values: X117 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X118 [1 0] Train data Unique Value Shape: (2,)
Train data Unique Values: X119 [1 0] Train data Unique Value Shape: (2,)
Train data Unique Values: X120 [1 0] Train data Unique Value Shape: (2,)
Train data Unique Values: X122 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X123 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X124 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X125 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X126 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X127 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X128 [1 0] Train data Unique Value Shape: (2,)
Train data Unique Values: X129 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X130 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X131 [1 0] Train data Unique Value Shape: (2,)
Train data Unique Values: X132 [0 1] Train data Unique Value Shape: (2,)
```

```
Train data Unique Values: X133 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X134 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X135 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X136 [1 0] Train data Unique Value Shape: (2,)
Train data Unique Values: X137 [1 0] Train data Unique Value Shape: (2,)
Train data Unique Values: X138 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X139 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X140 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X141 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X142 [1 0] Train data Unique Value Shape: (2,)
Train data Unique Values: X143 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X144 [1 0] Train data Unique Value Shape: (2,)
Train data Unique Values: X145 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X146 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X147 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X148 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X150 [1 0] Train data Unique Value Shape: (2,)
Train data Unique Values: X151 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X152 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X153 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X154 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X155 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X156 [1 0] Train data Unique Value Shape: (2,)
Train data Unique Values: X157 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X158 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X159 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X160 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X161 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X162 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X163 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X164 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X165 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X166 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X167 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X168 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X169 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X170 [1 0] Train data Unique Value Shape: (2,)
Train data Unique Values: X171 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X172 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X173 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X174 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X175 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X176 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X177 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X178 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X179 [1 0] Train data Unique Value Shape: (2,)
Train data Unique Values: X180 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X181 [0 1] Train data Unique Value Shape: (2,)
```

```
Train data Unique Values: X182 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X183 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X184 [1 0] Train data Unique Value Shape: (2,)
Train data Unique Values: X185 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X186 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X187 [1 0] Train data Unique Value Shape: (2,)
Train data Unique Values: X189 [1 0] Train data Unique Value Shape: (2,)
Train data Unique Values: X190 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X191 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X192 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X194 [1 0] Train data Unique Value Shape: (2,)
Train data Unique Values: X195 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X196 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X197 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X198 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X199 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X200 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X201 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X202 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X203 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X204 [1 0] Train data Unique Value Shape: (2,)
Train data Unique Values: X205 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X206 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X207 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X208 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X209 [1 0] Train data Unique Value Shape: (2,)
Train data Unique Values: X210 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X211 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X212 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X213 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X214 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X215 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X216 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X217 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X218 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X219 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X220 [1 0] Train data Unique Value Shape: (2,)
Train data Unique Values: X221 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X222 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X223 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X224 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X225 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X226 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X227 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X228 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X229 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X230 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X231 [0 1] Train data Unique Value Shape: (2,)
```

```
Train data Unique Values: X232 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X234 [1 0] Train data Unique Value Shape: (2,)
Train data Unique Values: X236 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X237 [1 0] Train data Unique Value Shape: (2,)
Train data Unique Values: X238 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X239 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X240 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X241 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X242 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X243 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X244 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X245 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X246 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X247 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X248 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X249 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X250 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X251 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X252 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X253 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X254 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X255 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X256 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X257 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X258 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X259 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X260 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X261 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X262 [1 0] Train data Unique Value Shape: (2,)
Train data Unique Values: X263 [1 0] Train data Unique Value Shape: (2,)
Train data Unique Values: X264 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X265 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X266 [1 0] Train data Unique Value Shape: (2,)
Train data Unique Values: X267 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X269 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X270 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X271 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X272 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X273 [1 0] Train data Unique Value Shape: (2,)
Train data Unique Values: X274 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X275 [1 0] Train data Unique Value Shape: (2,)
Train data Unique Values: X276 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X277 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X278 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X279 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X280 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X281 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X282 [0 1] Train data Unique Value Shape: (2,)
```

```
Train data Unique Values: X283 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X284 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X285 [1 0] Train data Unique Value Shape: (2,)
Train data Unique Values: X286 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X287 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X288 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X291 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X292 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X294 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X295 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X296 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X298 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X299 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X300 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X301 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X302 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X304 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X305 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X306 [1 0] Train data Unique Value Shape: (2,)
Train data Unique Values: X307 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X308 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X309 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X310 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X311 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X312 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X313 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X314 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X315 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X316 [1 0] Train data Unique Value Shape: (2,)
Train data Unique Values: X317 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X318 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X319 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X320 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X321 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X322 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X323 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X324 [1 0] Train data Unique Value Shape: (2,)
Train data Unique Values: X325 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X326 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X327 [1 0] Train data Unique Value Shape: (2,)
Train data Unique Values: X328 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X329 [1 0] Train data Unique Value Shape: (2,)
Train data Unique Values: X331 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X332 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X333 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X334 [1 0] Train data Unique Value Shape: (2,)
Train data Unique Values: X335 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X336 [0 1] Train data Unique Value Shape: (2,)
```

```
Train data Unique Values: X337 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X338 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X339 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X340 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X341 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X342 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X343 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X344 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X345 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X346 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X348 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X349 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X350 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X351 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X352 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X353 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X354 [1 0] Train data Unique Value Shape: (2,)
Train data Unique Values: X355 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X356 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X357 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X358 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X359 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X360 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X361 [1 0] Train data Unique Value Shape: (2,)
Train data Unique Values: X362 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X363 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X364 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X365 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X366 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X367 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X368 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X369 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X370 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X371 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X372 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X373 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X374 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X375 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X376 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X377 [1 0] Train data Unique Value Shape: (2,)
Train data Unique Values: X378 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X379 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X380 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X382 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X383 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X384 [0 1] Train data Unique Value Shape: (2,)
Train data Unique Values: X385 [0 1] Train data Unique Value Shape: (2,)
```

```
[22]:  # Check the unique value of test dataset.

       for columns in data_test:
           print('Test Data Unique Values:',columns,data_test[columns].unique(),
                 'Test Data Unique Value Shape:',data_test[columns].unique().shape)
```

Test Data Unique Values: ID [   1    2    3 … 8413 8414 8416] Test Data Unique
Value Shape: (4209,)
Test Data Unique Values: X0 ['az' 't' 'w' 'y' 'x' 'f' 'ap' 'o' 'ay' 'al' 'h' 'z'
'aj' 'd' 'v' 'ak'
 'ba' 'n' 'j' 's' 'af' 'ax' 'at' 'aq' 'av' 'm' 'k' 'a' 'e' 'ai' 'i' 'ag'
 'b' 'am' 'aw' 'as' 'r' 'ao' 'u' 'l' 'c' 'ad' 'au' 'bc' 'g' 'an' 'ae' 'p'
 'bb'] Test Data Unique Value Shape: (49,)
Test Data Unique Values: X1 ['v' 'b' 'l' 's' 'aa' 'r' 'a' 'i' 'p' 'c' 'o' 'm'
'z' 'e' 'h' 'w' 'g' 'k'
 'y' 't' 'u' 'd' 'j' 'q' 'n' 'f' 'ab'] Test Data Unique Value Shape: (27,)
Test Data Unique Values: X2 ['n' 'ai' 'as' 'ae' 's' 'b' 'e' 'ak' 'm' 'a' 'aq'
'ag' 'r' 'k' 'aj' 'ay'
 'ao' 'an' 'ac' 'af' 'ax' 'h' 'i' 'f' 'ap' 'p' 'au' 't' 'z' 'y' 'aw' 'd'
 'at' 'g' 'am' 'j' 'x' 'ab' 'w' 'q' 'ah' 'ad' 'al' 'av' 'u'] Test Data Unique
Value Shape: (45,)
Test Data Unique Values: X3 ['f' 'a' 'c' 'e' 'd' 'g' 'b'] Test Data Unique Value
Shape: (7,)
Test Data Unique Values: X4 ['d' 'b' 'a' 'c'] Test Data Unique Value Shape: (4,)
Test Data Unique Values: X5 ['t' 'b' 'a' 'z' 'y' 'x' 'h' 'g' 'f' 'j' 'i' 'd' 'c'
'af' 'ag' 'ab' 'ac'
 'ad' 'ae' 'ah' 'l' 'k' 'n' 'm' 'p' 'q' 's' 'r' 'v' 'w' 'o' 'aa'] Test Data
Unique Value Shape: (32,)
Test Data Unique Values: X6 ['a' 'g' 'j' 'l' 'i' 'd' 'f' 'h' 'c' 'k' 'e' 'b']
Test Data Unique Value Shape: (12,)
Test Data Unique Values: X8 ['w' 'y' 'j' 'n' 'm' 's' 'a' 'v' 'r' 'o' 't' 'h' 'c'
'k' 'p' 'u' 'd' 'g'
 'b' 'q' 'e' 'l' 'f' 'i' 'x'] Test Data Unique Value Shape: (25,)
Test Data Unique Values: X10 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X12 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X13 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X14 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X15 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X16 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X17 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X18 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X19 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X20 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X21 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X22 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X23 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X24 [0 1] Test Data Unique Value Shape: (2,)

```
Test Data Unique Values: X26 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X27 [1 0] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X28 [1 0] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X29 [1 0] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X30 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X31 [1 0] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X32 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X33 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X34 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X35 [1 0] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X36 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X37 [1 0] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X38 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X39 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X40 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X41 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X42 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X43 [1 0] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X44 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X45 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X46 [1 0] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X47 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X48 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X49 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X50 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X51 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X52 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X53 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X54 [1 0] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X55 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X56 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X57 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X58 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X59 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X60 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X61 [1 0] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X62 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X63 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X64 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X65 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X66 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X67 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X68 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X69 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X70 [1 0] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X71 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X73 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X74 [1 0] Test Data Unique Value Shape: (2,)
```

```
Test Data Unique Values: X75 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X76 [1 0] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X77 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X78 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X79 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X80 [1 0] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X81 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X82 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X83 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X84 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X85 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X86 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X87 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X88 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X89 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X90 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X91 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X92 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X94 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X95 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X96 [1 0] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X97 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X98 [1 0] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X99 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X100 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X101 [1 0] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X102 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X103 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X104 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X105 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X106 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X108 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X109 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X110 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X111 [1 0] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X112 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X113 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X114 [1 0] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X115 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X116 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X117 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X118 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X119 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X120 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X122 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X123 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X124 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X125 [0 1] Test Data Unique Value Shape: (2,)
```

```
Test Data Unique Values: X126 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X127 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X128 [1 0] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X129 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X130 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X131 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X132 [1 0] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X133 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X134 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X135 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X136 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X137 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X138 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X139 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X140 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X141 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X142 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X143 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X144 [1 0] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X145 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X146 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X147 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X148 [1 0] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X150 [1 0] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X151 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X152 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X153 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X154 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X155 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X156 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X157 [1 0] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X158 [1 0] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X159 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X160 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X161 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X162 [1 0] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X163 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X164 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X165 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X166 [1 0] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X167 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X168 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X169 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X170 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X171 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X172 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X173 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X174 [0 1] Test Data Unique Value Shape: (2,)
```

```
Test Data Unique Values: X175 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X176 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X177 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X178 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X179 [1 0] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X180 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X181 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X182 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X183 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X184 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X185 [1 0] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X186 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X187 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X189 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X190 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X191 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X192 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X194 [1 0] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X195 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X196 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X197 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X198 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X199 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X200 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X201 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X202 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X203 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X204 [1 0] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X205 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X206 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X207 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X208 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X209 [1 0] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X210 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X211 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X212 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X213 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X214 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X215 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X216 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X217 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X218 [1 0] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X219 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X220 [1 0] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X221 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X222 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X223 [1 0] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X224 [0 1] Test Data Unique Value Shape: (2,)
```

```
Test Data Unique Values: X225 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X226 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X227 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X228 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X229 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X230 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X231 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X232 [1 0] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X234 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X236 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X237 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X238 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X239 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X240 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X241 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X242 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X243 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X244 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X245 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X246 [1 0] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X247 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X248 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X249 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X250 [1 0] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X251 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X252 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X253 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X254 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X255 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X256 [1 0] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X257 [0] Test Data Unique Value Shape: (1,)
Test Data Unique Values: X258 [0] Test Data Unique Value Shape: (1,)
Test Data Unique Values: X259 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X260 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X261 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X262 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X263 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X264 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X265 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X266 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X267 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X269 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X270 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X271 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X272 [1 0] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X273 [1 0] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X274 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X275 [0 1] Test Data Unique Value Shape: (2,)
```

```
Test Data Unique Values: X276 [1 0] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X277 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X278 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X279 [1 0] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X280 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X281 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X282 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X283 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X284 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X285 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X286 [1 0] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X287 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X288 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X291 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X292 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X294 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X295 [0] Test Data Unique Value Shape: (1,)
Test Data Unique Values: X296 [0] Test Data Unique Value Shape: (1,)
Test Data Unique Values: X298 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X299 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X300 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X301 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X302 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X304 [1 0] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X305 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X306 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X307 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X308 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X309 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X310 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X311 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X312 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X313 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X314 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X315 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X316 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X317 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X318 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X319 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X320 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X321 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X322 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X323 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X324 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X325 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X326 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X327 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X328 [1 0] Test Data Unique Value Shape: (2,)
```

```
Test Data Unique Values: X329 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X331 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X332 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X333 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X334 [1 0] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X335 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X336 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X337 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X338 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X339 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X340 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X341 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X342 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X343 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X344 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X345 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X346 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X348 [1 0] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X349 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X350 [1 0] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X351 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X352 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X353 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X354 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X355 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X356 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X357 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X358 [1 0] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X359 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X360 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X361 [1 0] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X362 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X363 [1 0] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X364 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X365 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X366 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X367 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X368 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X369 [0] Test Data Unique Value Shape: (1,)
Test Data Unique Values: X370 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X371 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X372 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X373 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X374 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X375 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X376 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X377 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X378 [1 0] Test Data Unique Value Shape: (2,)
```

```
Test Data Unique Values: X379 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X380 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X382 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X383 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X384 [0 1] Test Data Unique Value Shape: (2,)
Test Data Unique Values: X385 [0 1] Test Data Unique Value Shape: (2,)
```

[23]: 
```python
# Check null (NaN) or Missing value in train dataset.

data_train.isna().any()
```

[23]: 
```
ID      False
y       False
X0      False
X1      False
X2      False
        …
X380    False
X382    False
X383    False
X384    False
X385    False
Length: 366, dtype: bool
```

[24]: 
```python
data_train.isna().sum().any()
```

[24]: False

[25]: 
```python
# Check null (NaN) of Missing value in test dataset.

data_test.isna().any()
```

[25]: 
```
ID      False
X0      False
X1      False
X2      False
X3      False
        …
X380    False
X382    False
X383    False
X384    False
X385    False
Length: 365, dtype: bool
```

[26]: 
```python
data_test.isna().sum().any()
```

[26]: False

- As we can say that there is no Null (NaN or Missing values) values available in Train as well as Test dataset.

```
[27]: data_train.isna().any()[data_train.isna().any()].index.values
```

```
[27]: array([], dtype=object)
```

```
[28]: data_test.isna().any()[data_test.isna().any()].index.values
```

```
[28]: array([], dtype=object)
```

```
[29]: # Find out the outliers in dataset

      sns.boxplot(data_train['y'])
```

```
[29]: <AxesSubplot:xlabel='y'>
```



```
[30]: # Treatment on outliers

      #Import required library

      from scipy import stats
```

```
iqr = stats.iqr(data_train['y'])

q1 = data_train['y'].quantile(0.25)
q3 = data_train['y'].quantile(0.75)

upperbound = q3 + 1.5*(iqr)
lowerbound = q1 - 1.5*(iqr)

print('Q1 is:',q1)
print('Q3 is:',q3)
print('Upper Bound is:',upperbound)
print('Lower Bound is:',lowerbound)
```

```
Q1 is: 90.82
Q3 is: 109.01
Upper Bound is: 136.29500000000002
Lower Bound is: 63.534999999999975
```

[31]: `data_train = data_train[(data_train['y']<136.295)]`

[32]: `sns.boxplot(data_train['y'])`

[32]: `<AxesSubplot:xlabel='y'>`



- Outliers found in higher side in y column hence removed the outliers from the train dataset

26

```
[33]: # Check the correlation between the train dataset.

      data_train.corr()
```

```
[33]:            ID         y        X10       X12       X13       X14       X15  \
      ID    1.000000 -0.050630  0.001390  0.059408 -0.035614 -0.026758  0.002207
      y    -0.050630  1.000000 -0.024244  0.089837  0.052033  0.216901  0.026882
      X10   0.001390 -0.024244  1.000000 -0.033155 -0.028975 -0.101005 -0.002563
      X12   0.059408  0.089837 -0.033155  1.000000  0.215648 -0.245361 -0.006225
      X13  -0.035614  0.052033 -0.028975  0.215648  1.000000 -0.085463 -0.005440
      ...        ...       ...       ...       ...       ...       ...       ...
      X380 -0.013824  0.050207 -0.010606 -0.005432  0.023198  0.007861 -0.001991
      X382 -0.038549 -0.173845 -0.010287 -0.024990 -0.021839  0.012860 -0.001931
      X383 -0.006331 -0.002986 -0.004053 -0.009846 -0.008605  0.026104 -0.000761
      X384 -0.015481 -0.004021 -0.002563 -0.006225  0.041500  0.025370 -0.000481
      X385  0.029146 -0.022945 -0.004441 -0.010787 -0.009427  0.043964 -0.000834

                 X16       X17       X18  ...      X375      X376      X377  \
      ID   -0.036778 -0.038549 -0.030057  ...  0.046756 -0.083687 -0.023784
      y     0.057175 -0.173845 -0.008593  ...  0.031648  0.125295  0.065368
      X10  -0.006016 -0.010287 -0.010287  ...  0.166474 -0.028719 -0.074593
      X12  -0.014614 -0.024990 -0.024990  ... -0.111403 -0.069763  0.030725
      X13  -0.012772 -0.021839 -0.010061  ... -0.169490 -0.060967  0.357497
      ...        ...       ...       ... ...       ...       ...       ...
      X380 -0.004675 -0.007994 -0.007994  ... -0.062043 -0.022318 -0.061460
      X382 -0.004535  1.000000  0.086723  ... -0.060176 -0.021646 -0.059610
      X383 -0.001787 -0.003055 -0.003055  ... -0.008814 -0.008528  0.021355
      X384 -0.001130 -0.001931 -0.001931  ... -0.014990 -0.005392  0.008776
      X385 -0.001957 -0.003347 -0.003347  ...  0.055620 -0.009344 -0.025731

                 X378      X379      X380      X382      X383      X384      X385
      ID    0.030264  0.022330 -0.013824 -0.038549 -0.006331 -0.015481  0.029146
      y    -0.281230  0.069033  0.050207 -0.173845 -0.002986 -0.004021 -0.022945
      X10  -0.017077 -0.011367 -0.010606 -0.010287 -0.004053 -0.002563 -0.004441
      X12  -0.015895 -0.027612 -0.005432 -0.024990 -0.009846 -0.006225 -0.010787
      X13  -0.036252 -0.024130  0.023198 -0.021839 -0.008605  0.041500 -0.009427
      ...        ...       ...       ...       ...       ...       ...       ...
      X380 -0.013270 -0.008833  1.000000 -0.007994 -0.003150 -0.001991 -0.003451
      X382 -0.012871 -0.008567 -0.007994  1.000000 -0.003055 -0.001931 -0.003347
      X383 -0.005071 -0.003375 -0.003150 -0.003055  1.000000 -0.000761 -0.001319
      X384 -0.003206 -0.002134 -0.001991 -0.001931 -0.000761  1.000000 -0.000834
      X385 -0.005556 -0.003698 -0.003451 -0.003347 -0.001319 -0.000834  1.000000

      [358 rows x 358 columns]
```

```
[34]: # Check the correlation between the test dataset.
```

```
data_test.corr()
```

[34]:
```
            ID        X10        X12        X13        X14        X15        X16  \
ID    1.000000  -0.016166   0.043162   0.017910  -0.036099   0.005100  -0.024482
X10  -0.016166   1.000000  -0.039453  -0.035496  -0.120379  -0.003717  -0.007125
X12   0.043162  -0.039453   1.000000   0.283228  -0.245127  -0.007570  -0.014509
X13   0.017910  -0.035496   0.283228   1.000000  -0.076145  -0.006811  -0.013054
X14  -0.036099  -0.120379  -0.245127  -0.076145   1.000000  -0.023097  -0.044269
...        ...        ...        ...        ...        ...        ...        ...
X380  0.012520  -0.012561  -0.025578   0.054582   0.007787  -0.002410  -0.004619
X382 -0.021581  -0.013108  -0.016991  -0.024015   0.000864  -0.002515  -0.004821
X383 -0.001625  -0.003035  -0.006180  -0.005560   0.025212  -0.000582  -0.001116
X384  0.013948  -0.003717  -0.007570  -0.006811   0.030881  -0.000713  -0.001367
X385  0.026357  -0.005681  -0.011569  -0.010408   0.047195  -0.001090  -0.002089

            X17        X18        X19    ...        X375       X376       X377  \
ID    -0.021581  -0.010920  -0.020800   ...    0.024007  -0.087891   0.004271
X10   -0.013108  -0.014142  -0.049351   ...    0.189023  -0.031817  -0.086214
X12   -0.016991  -0.028796  -0.100493   ...   -0.148812  -0.064790   0.080843
X13   -0.024015  -0.025908  -0.090413   ...   -0.177340  -0.058291   0.359450
X14    0.000864  -0.087862  -0.306620   ...    0.107496   0.043260  -0.139742
...        ...        ...        ...    ...        ...        ...        ...
X380  -0.008498   0.043621  -0.023568   ...   -0.062756  -0.020628  -0.060764
X382   1.000000   0.066366  -0.033389   ...   -0.065490  -0.021526  -0.063411
X383  -0.002053  -0.002215  -0.007730   ...   -0.015163  -0.004984   0.008850
X384  -0.002515  -0.002713  -0.009469   ...    0.000420  -0.006105   0.020448
X385  -0.003844  -0.004147  -0.014471   ...    0.058691  -0.009330  -0.027482

            X378       X379       X380       X382       X383       X384       X385
ID    -0.002939   0.031762   0.012520  -0.021581  -0.001625   0.013948   0.026357
X10   -0.019498  -0.015262  -0.012561  -0.013108  -0.003035  -0.003717  -0.005681
X12   -0.006747  -0.022720  -0.025578  -0.016991  -0.006180  -0.007570  -0.011569
X13   -0.035722  -0.027961   0.054582  -0.024015  -0.005560  -0.006811  -0.010408
X14   -0.051238   0.113487   0.007787   0.000864   0.025212   0.030881   0.047195
...        ...        ...        ...        ...        ...        ...        ...
X380  -0.012641  -0.009895   1.000000  -0.008498  -0.001968  -0.002410  -0.003683
X382  -0.013192  -0.010326  -0.008498   1.000000  -0.002053  -0.002515  -0.003844
X383  -0.003054  -0.002391  -0.001968  -0.002053   1.000000  -0.000582  -0.000890
X384  -0.003741  -0.002928  -0.002410  -0.002515  -0.000582   1.000000  -0.001090
X385  -0.005717  -0.004475  -0.003683  -0.003844  -0.000890  -0.001090   1.000000

[357 rows x 357 columns]
```

[35]:
```python
# ID column not required for modeling hence removed from both dataset.

data_train.drop('ID',axis=1,inplace=True)
data_test.drop('ID',axis=1,inplace=True)
```

```
[36]: # Shuffle or Split the train dataset in input and output feature.

      features = data_train.drop('y',axis=1)
      target = data_train[['y']]
```

# 4  3. Apply label encoder.

### 4.0.1  Train Dataset

```
[37]: # Find out the categorical values in train datatset.

      dictionary={}
      dictionary['categorical']=features.dtypes[features.dtypes=='object'].index
      dictionary
```

```
[37]: {'categorical': Index(['X0', 'X1', 'X2', 'X3', 'X4', 'X5', 'X6', 'X8'],
      dtype='object')}
```

- In train data found some columns in categorical form so convert in numerical form by using Label Encoder

```
[38]: # Categorical value convert in numeric form for modeling by using Label Encoder␣
      ↪(Ordinal)

      from sklearn.preprocessing import LabelEncoder
```

```
[39]: le_X0 = LabelEncoder()
      le_X1 = LabelEncoder()
      le_X2 = LabelEncoder()
      le_X3 = LabelEncoder()
      le_X4 = LabelEncoder()
      le_X5 = LabelEncoder()
      le_X6 = LabelEncoder()
      le_X8 = LabelEncoder()
```

```
[40]: features['X0'] = le_X0.fit_transform(features['X0'])
      features['X1'] = le_X1.fit_transform(features['X1'])
      features['X2'] = le_X2.fit_transform(features['X2'])
      features['X3'] = le_X3.fit_transform(features['X3'])
      features['X4'] = le_X4.fit_transform(features['X4'])
      features['X5'] = le_X5.fit_transform(features['X5'])
      features['X6'] = le_X6.fit_transform(features['X6'])
      features['X8'] = le_X8.fit_transform(features['X8'])
```

```
[41]: features.head(10)
```

```
[41]:      X0   X1   X2   X3   X4   X5   X6   X8   X10   X12   …   X375   X376   X377   X378   \
      0   32   23   16    0    3   24    9   14     0     0   …      0      0      1      0
      1   32   21   18    4    3   28   11   14     0     0   …      1      0      0      0
      2   20   24   33    2    3   27    9   23     0     0   …      0      0      0      0
      3   20   21   33    5    3   27   11    4     0     0   …      0      0      0      0
      4   20   23   33    5    3   12    3   13     0     0   …      0      0      0      0
      5   40    3   24    2    3   11    7   18     0     0   …      0      0      1      0
      6    9   19   24    5    3   10    7   18     0     0   …      0      0      0      0
      7   36   13   15    5    3   10    9    0     0     0   …      0      0      0      0
      8   43   20   15    4    3   10    8    7     0     0   …      1      0      0      0
      9   31    3   13    2    3   10    0    4     0     0   …      0      0      1      0

         X379   X380   X382   X383   X384   X385
      0     0      0      0      0      0      0
      1     0      0      0      0      0      0
      2     0      0      1      0      0      0
      3     0      0      0      0      0      0
      4     0      0      0      0      0      0
      5     0      0      0      0      0      0
      6     0      0      0      0      0      0
      7     0      0      0      0      0      0
      8     0      0      0      0      0      0
      9     0      0      0      0      0      0

      [10 rows x 364 columns]
```

### 4.0.2 Test Dataset

```
[42]: dictionary={}
      dictionary['cat']=data_test.dtypes[data_test.dtypes=='object'].index
      dictionary
```

```
[42]: {'cat': Index(['X0', 'X1', 'X2', 'X3', 'X4', 'X5', 'X6', 'X8'], dtype='object')}
```

- In test data found some columns in categorical form so convert in numerical form by using Label Encoder

```
[43]: le_X0 = LabelEncoder()
      le_X1 = LabelEncoder()
      le_X2 = LabelEncoder()
      le_X3 = LabelEncoder()
      le_X4 = LabelEncoder()
      le_X5 = LabelEncoder()
      le_X6 = LabelEncoder()
      le_X8 = LabelEncoder()
```

```
[44]: data_test['X0'] = le_X0.fit_transform(data_test['X0'])
      data_test['X1'] = le_X1.fit_transform(data_test['X1'])
      data_test['X2'] = le_X2.fit_transform(data_test['X2'])
      data_test['X3'] = le_X3.fit_transform(data_test['X3'])
      data_test['X4'] = le_X4.fit_transform(data_test['X4'])
      data_test['X5'] = le_X5.fit_transform(data_test['X5'])
      data_test['X6'] = le_X6.fit_transform(data_test['X6'])
      data_test['X8'] = le_X8.fit_transform(data_test['X8'])
```

```
[45]: data_test.head(10)
```

```
[45]:     X0  X1  X2  X3  X4  X5  X6  X8  X10  X12  …  X375  X376  X377  X378  \
      0   21  23  34   5   3  26   0  22    0    0  …     0     0     0     1
      1   42   3   8   0   3   9   6  24    0    0  …     0     0     1     0
      2   21  23  17   5   3   0   9   9    0    0  …     0     0     0     1
      3   21  13  34   5   3  31  11  13    0    0  …     0     0     0     1
      4   45  20  17   2   3  30   8  12    0    0  …     1     0     0     0
      5   47   1   8   4   3  29   6  18    0    0  …     1     0     0     0
      6   46   3   4   3   3  29   3  24    0    0  …     0     0     0     0
      7   29  20   4   2   3  14   3   0    0    0  …     0     0     1     0
      8   12  13  38   2   3  14   9  13    0    0  …     0     0     0     0
      9   38  23  17   5   3  13   5  21    0    0  …     0     0     0     0

         X379  X380  X382  X383  X384  X385
      0     0     0     0     0     0     0
      1     0     0     0     0     0     0
      2     0     0     0     0     0     0
      3     0     0     0     0     0     0
      4     0     0     0     0     0     0
      5     0     0     0     0     0     0
      6     0     1     0     0     0     0
      7     0     0     0     0     0     0
      8     0     0     0     0     0     0
      9     0     0     0     0     0     0

      [10 rows x 364 columns]
```

```
[46]: # More variation in dataset hence rescaling the values by using standard scaler␣
      ↪(range 0 to 1).

      from sklearn.preprocessing import StandardScaler

      scaler = StandardScaler()

      features = scaler.fit_transform(features)
```

```
[47]: # Import required library

      from sklearn.model_selection import train_test_split

      # Split the train dataset in traning and testing from.

      xtrain,xtest,ytrain,ytest = train_test_split(features,target,test_size=0.
       →25,random_state=10)
```

```
[48]: # Check shape of the data.

      print(xtrain.shape)
      print(xtest.shape)
      print(ytrain.shape)
      print(ytest.shape)
```

```
(3119, 364)
(1040, 364)
(3119, 1)
(1040, 1)
```

# 5  4. Perform dimensionality reduction (PCA).

```
[49]: # Import required library

      from sklearn.decomposition import PCA
```

```
[50]: pca = PCA(n_components=0.95)
```

```
[51]: xtrain_transform = pca.fit_transform(xtrain)
      xtest_transform = pca.transform(xtest)
```

```
[52]: pca.explained_variance_ratio_
```

```
[52]: array([0.07019274, 0.05729392, 0.04609237, 0.03515298, 0.03375958,
             0.03292353, 0.02905596, 0.02185832, 0.02104466, 0.0177154 ,
             0.01656399, 0.01553658, 0.01492109, 0.01405254, 0.01393377,
             0.01305241, 0.01184265, 0.01118177, 0.01083652, 0.01068139,
             0.01012464, 0.00943591, 0.0088153 , 0.00873972, 0.00813892,
             0.00795704, 0.00786851, 0.00752673, 0.0074265 , 0.00727647,
             0.00693622, 0.00674369, 0.00669274, 0.00652471, 0.00644761,
             0.00623725, 0.00609012, 0.00592114, 0.00580969, 0.00577976,
             0.00556441, 0.00549373, 0.00542699, 0.00520633, 0.00516627,
             0.0050948 , 0.00497358, 0.00495062, 0.00489009, 0.00471821,
             0.00466741, 0.00463163, 0.00453932, 0.00444688, 0.00435878,
             0.00430327, 0.00423754, 0.00416412, 0.00412734, 0.00409017,
```

```
          0.00406457, 0.00400174, 0.00396979, 0.0039042 , 0.00385469,
          0.00381901, 0.0037737 , 0.00373333, 0.00366122, 0.00358136,
          0.0035554 , 0.00349281, 0.00348558, 0.00340535, 0.00337646,
          0.00334585, 0.00333005, 0.00328155, 0.00320043, 0.00317289,
          0.00314952, 0.00311033, 0.00309291, 0.0030509 , 0.00302206,
          0.00297985, 0.00292358, 0.00287127, 0.00282779, 0.00281652,
          0.00278533, 0.00275123, 0.00272304, 0.00268221, 0.00266552,
          0.00263688, 0.00260685, 0.00256342, 0.00255046, 0.00252926,
          0.00251445, 0.0024949 , 0.00245533, 0.00244842, 0.00242278,
          0.00239232, 0.00236028, 0.00234539, 0.00230195, 0.00228297,
          0.00226582, 0.00221066, 0.00220342, 0.00215847, 0.00211229,
          0.00211003, 0.00207315, 0.00204617, 0.00199182, 0.00195448,
          0.00194303, 0.00191564, 0.00187583, 0.0018657 , 0.0018458 ,
          0.00181042, 0.00177961, 0.0017544 , 0.0017379 , 0.00171343,
          0.00168156, 0.00167204, 0.00164478, 0.00161469, 0.00159322,
          0.00155612, 0.0015177 , 0.00149996, 0.00146936, 0.00143128,
          0.00140798, 0.00139064, 0.00137898])
```

```python
[53]: print(xtrain_transform.shape)
      print(xtest_transform.shape)
```

```
      (3119, 143)
      (1040, 143)
```

```python
[54]: xtrain_transform
```

```
[54]: array([[10.61283975, -2.47702013,  0.82357088, …,  0.53083803,
               0.29588805, -1.14707981],
             [-2.17584734, -1.08495605,  0.59481737, …, -0.12867998,
               0.01131092,  0.4324378 ],
             [-2.73604686,  0.15089878,  2.7461672 , …,  0.40397308,
              -0.07231659, -0.28202867],
             …,
             [-2.12381216, -3.91014979, -6.81306636, …, -0.53600069,
               0.27311935,  0.23946768],
             [-1.09391938,  0.32794205,  2.96897123, …, -0.56746158,
              -0.6645795 , -1.05995913],
             [-0.29338825, -1.43396707, -3.31466989, …,  0.85925952,
               0.24446146, -0.60080401]])
```

```python
[55]: # Linear Regression applied
      #Import required library

      from sklearn.linear_model import LinearRegression
```

```python
[56]: model = LinearRegression()
```

```
model.fit(xtrain_transform,ytrain)
ypred = model.predict(xtest_transform)
```

[57]:
```
from sklearn.metrics import␣
 ↪mean_absolute_error,mean_squared_error,r2_score,accuracy_score
```

[58]:
```
print('Mean Absolute Error = ',mean_absolute_error(ytest,ypred))
print('Mean Squared Error = ',mean_squared_error(ytest,ypred))
print('Root Mean Squared Error = ',np.sqrt(mean_squared_error(ytest,ypred)))
print('R2 Score = ',r2_score(ytest,ypred))
```

```
Mean Absolute Error =  5.159089114357101
Mean Squared Error =  51.137287254097004
Root Mean Squared Error =  7.151033998947075
R2 Score =  0.6069512538830951
```

- Got R2 score is 0.6069 by applying Linear Regression.

# 6 5. Predict your test_df values using XGBoost.

[59]:
```
#Import required library

from xgboost import XGBRegressor
```

[60]:
```
xgb_rg = XGBRegressor(booster = 'gblinear')
```

[61]:
```
xgb_rg.fit(xtrain_transform, ytrain)
```

[61]:
```
XGBRegressor(base_score=0.5, booster='gblinear', colsample_bylevel=None,
             colsample_bynode=None, colsample_bytree=None,
             enable_categorical=False, gamma=None, gpu_id=-1,
             importance_type=None, interaction_constraints=None,
             learning_rate=0.5, max_delta_step=None, max_depth=None,
             min_child_weight=None, missing=nan, monotone_constraints=None,
             n_estimators=100, n_jobs=4, num_parallel_tree=None, predictor=None,
             random_state=0, reg_alpha=0, reg_lambda=0, scale_pos_weight=1,
             subsample=None, tree_method=None, validate_parameters=1,
             verbosity=None)
```

[62]:
```
xgb_preds = xgb_rg.predict(xtest_transform)
```

[63]:
```
print('Mean Absolute Error = ',mean_absolute_error(xgb_preds,ytest))
print('Mean Squared Error = ',mean_squared_error(xgb_preds, ytest))
print('Root Mean Squared Error = ',np.sqrt(mean_squared_error(xgb_preds,␣
 ↪ytest)))
print('R2 Score = ',r2_score(ytest,xgb_preds))
```

```
Mean Absolute Error =  5.159087237724891
Mean Squared Error =  51.13726496280257
Root Mean Squared Error =  7.151032440340526
R2 Score =  0.6069514252172781
```

- Got R2 score is 0.6069 by applying XGBoost (XGBRegressor).

[64]: 
```python
# Applied Ridge Regressor (for checking R2 score)

from sklearn.linear_model import Ridge
```

[65]: 
```python
ridge_model = Ridge()
```

[66]: 
```python
ridge_model.fit(xtrain_transform,ytrain)
```

[66]: 
```
Ridge()
```

[67]: 
```python
ridge_preds = ridge_model.predict(xtest_transform)
```

[68]: 
```python
print('Mean Absolute Error = ',mean_absolute_error(ridge_preds,ytest))
print('Mean Squared Error = ',mean_squared_error(ridge_preds, ytest))
print('Root Mean Squared Error = ',np.sqrt(mean_squared_error(ridge_preds,
 →ytest)))
print('R2 Score = ',r2_score(ytest,ridge_preds))
```

```
Mean Absolute Error =  5.159036828040949
Mean Squared Error =  51.135617930582214
Root Mean Squared Error =  7.1509172789637425
R2 Score =  0.6069640845502967
```

- Got R2 score is 0.6069 by applying Ridge Regressor.
- R2 score values found same in Linear Regressor, XGBoost and Ridge Regressor.

### 6.0.1 Find out the predicted values by using XGBoost on test data

[69]: 
```python
pca.fit(data_test)
```

[69]: 
```
PCA(n_components=0.95)
```

[70]: 
```python
test_transform = pca.transform(data_test)
```

[71]: 
```python
pred_xgbr = xgb_rg.predict(test_transform)
```

[72]: 
```python
pred_xgbr
```

[72]: 
```
array([112.495605, 103.38444 ,  85.74017 , …,  81.08352 , 104.7797  ,
        91.61079 ], dtype=float32)
```

```
[73]: pred_xgbr.shape
```

```
[73]: (4209,)
```

```
[74]: preds_xgbreg = pd.DataFrame(pred_xgbr, index = data_test.index)
       preds_xgbreg = preds_xgbreg.rename(columns = {0:'Predict_values'} )
```

```
[75]: # Predicted values by using XGBoost

       preds_xgbreg
```

```
[75]:       Predict_values
       0          112.495605
       1          103.384438
       2           85.740173
       3          107.230255
       4          104.222031
       ...               ...
       4204        97.012054
       4205        93.501297
       4206        81.083519
       4207       104.779701
       4208        91.610786

       [4209 rows x 1 columns]
```

```
[76]: test_predict = pd.concat([data_test,preds_xgbreg],axis = 1)
       test_predict
```

```
[76]:        X0  X1  X2  X3  X4  X5  X6  X8  X10  X12  …  X376  X377  X378  X379  \
       0      21  23  34   5   3  26   0  22    0    0  …     0     0     1     0
       1      42   3   8   0   3   9   6  24    0    0  …     0     1     0     0
       2      21  23  17   5   3   0   9   9    0    0  …     0     0     1     0
       3      21  13  34   5   3  31  11  13    0    0  …     0     0     1     0
       4      45  20  17   2   3  30   8  12    0    0  …     0     0     0     0

       ...    ..  ..  ..  ..  ..  ..  ..  ..  ...  ...  …   ...   ...   ...   ...
       4204    6   9  17   5   3   1   9   4    0    0  …     0     0     0     0
       4205   42   1   8   3   3   1   9  24    0    0  …     1     0     0     0
       4206   47  23  17   5   3   1   3  22    0    0  …     0     0     0     0
       4207    7  23  17   0   3   1   2  16    0    0  …     0     1     0     0
       4208   42   1   8   2   3   1   6  17    0    0  …     0     0     0     0

              X380  X382  X383  X384  X385  Predict_values
       0          0     0     0     0     0      112.495605
       1          0     0     0     0     0      103.384438
       2          0     0     0     0     0       85.740173
       3          0     0     0     0     0      107.230255
```

```
4          0    0    0    0    0       104.222031
...        ...  ...  ...  ...          ...
4204       0    0    0    0    0        97.012054
4205       0    0    0    0    0        93.501297
4206       0    0    0    0    0        81.083519
4207       0    0    0    0    0       104.779701
4208       0    0    0    0    0        91.610786
```

[4209 rows x 365 columns]

[ ]: