



Assessment: Mercedes-Benz Greener Manufacturing

The comments/sections provided are your cues to perform the assignment. You don't need to limit yourself to the number of rows/cells provided. You can add additional rows in each section to add more lines of code.

If at any point in time you need help on solving this assignment, view our demo video to understand the different steps of the code.

Happy coding!

Project 1: Mercedes-Benz Greener Manufacturing

DESCRIPTION

Reduce the time a Mercedes-Benz spends on the test bench.

Problem Statement Scenario:

Since the first automobile, the Benz Patent Motor Car in 1886, Mercedes-Benz has stood for important automotive innovations. These include the passenger safety cell with a crumple zone, the airbag, and intelligent assistance systems. Mercedes-Benz applies for nearly 2000 patents per year, making the brand the European leader among premium carmakers. Mercedes-Benz is the leader in the premium car industry. With a huge selection of features and options, customers can choose the customized Mercedes-Benz of their dreams.

To ensure the safety and reliability of every unique car configuration before they hit the road, the company's engineers have developed a robust testing system. As one of the world's biggest manufacturers of premium cars, safety and efficiency are paramount on Mercedes-Benz's production lines. However, optimizing the speed of their testing system for many possible feature combinations is complex and time-consuming without a powerful algorithmic approach.

You are required to reduce the time that cars spend on the test bench. Others will work with a dataset representing different permutations of features in a Mercedes-Benz car to predict the time it takes to pass testing. Optimal algorithms will contribute to faster testing, resulting in lower carbon dioxide emissions without reducing Mercedes-Benz's standards.

Following actions should be performed:

- If for any column(s), the variance is equal to zero, then you need to remove those variable(s).
- Check for null and unique values for test and train sets.
- Apply label encoder.
- Perform dimensionality reduction.
- Predict your test_df values using XGBBoost.

Import the required libraries

```
In [1]: # importing the required libraries
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
from matplotlib import style
import seaborn as sns
import matplotlib
%matplotlib inline

# importing the library to ignore warnings
import warnings
warnings.filterwarnings('ignore')
```

Read the train dataset

```
In [2]: # importing the train dataset
train_original = pd.read_csv('train.csv')

# let us understand the train data
print('Size of training set: {} rows and {} columns'.format(*train_original.shape))

# print first few rows of training set
train_original.head()
```

ID	y	X0	X1	X2	X3	X4	X5	X6	X8	...	X375	X376	X377	X378	X379	X380	X382	X383	X384	
0	0	130.81	k	v	a	t	a	d	y	j	o	...	0	0	1	0	0	0	0	0
1	6	88.53	k	t	a	v	e	d	y	i	o	...	1	0	0	0	0	0	0	0
2	7	76.26	az	w	n	c	d	x	j	x	...	0	0	0	0	0	0	0	1	0
3	9	80.62	az	t	n	f	d	x	i	e	...	0	0	0	0	0	0	0	0	0
4	13	78.02	az	v	n	f	d	h	d	n	...	0	0	0	0	0	0	0	0	0

5 rows × 378 columns

- There are 4209 data points and 378 features in the dataset.
- It is observed that there are more number of features with binary values. Hence sparsity exists in the train data.
- Target feature 'y' from the training data is 'the time that cars spend on the test bench.'

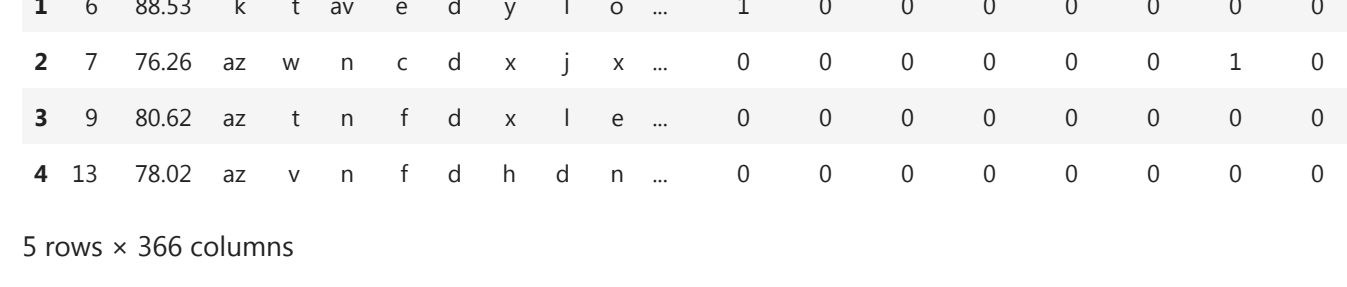
```
In [3]: # check the info of the train data
train_original.info()
```

```
<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: 12.0+ MB
```

- The train dataset has one feature with float64 datatype which is target (Y) features.
- The train dataset has 369 features with int64 datatype which are features having binary values (0 & 1).
- The train dataset has 8 features with object datatype which are features having categorical data.

```
In [4]: colors = ['f00157', '#50248f', '#a6a6a6', '#38d1ff', '#ca7beb', '#c8ae8f', '#9154f8', '#ce6b28']
plt.figure(figsize=(15,5))
plt.subplot(121)
sns.distplot(train_original.y.values, bins=20, color=colors[4])
plt.title('Target Value Distribution', fontstyle='italic', fontweight='bold', fontsize=15)
plt.xlabel('Target Value in Seconds'); plt.ylabel('Occurrences');
```

```
plt.subplot(122)
sns.boxplot(train_original.y.values, color=colors[0])
plt.title('Target Value Distribution', fontstyle='italic', fontweight='bold', fontsize=15)
plt.xlabel('Target Value in Seconds');
```



If for any column(s), the variance is equal to zero, then we need to remove those variable(s)

```
In [5]: #checking the variance of all the features in the train dataset and storing it to new
train_original_var=pd.DataFrame(train_original.var(axis=0),columns=['Variance'])
train_original_var
```

ID	Variance
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

370 rows × 1 columns

- Define a function to remove the features of train dataset with zero variance.
- Add the object train_original to the function. i.e., features_zero_var(train_original).
- The variance with respect to each features in the train_original dataset will be stored in the object df_original_var.

```
In [6]: # Define a function to remove the features of train dataset with zero variance
def features_zero_var(df):
    df_original_var=pd.DataFrame(df.var(axis=0),columns=['Variance'])
    return(df_original_var[df_original_var.Variance==0])
```

```
In [7]: # Call the function to return the train dataset features having zero variance.
features_zero_var(train_original)
```

Variance	
X11	0.0
X93	0.0
X107	0.0
X233	0.0
X235	0.0
X268	0.0
X289	0.0
X290	0.0
X293	0.0
X297	0.0
X330	0.0
X347	0.0

Above columns are the list of features having zero variance in training dataset.

```
In [8]: # Remove the features with zero variance from train dataset and store the data in the
train_original_modified=train_original.drop(columns=train_original_var[train_original_
```

```
In [9]: # Display the modified train dataset after removing the features having zero variance
train_original_modified.head()
```

ID	y	X0	X1	X2	X3	X4	X5	X6	X8	...	X375	X376	X377	X378	X379	X380	X382	X383	X384	
0	0	130.81	k	v	a	t	a	d	y	j	o	...	0	0	1	0	0	0	0	0
1	6	88.53	k	t	a	v	e	d	y	i	o	...	1	0	0	0	0	0	0	0
2	7	76.26	az	w	n	c	d	x	j	x	...	0	0	0	0	0	0	0	1	0
3	9	80.62	az	t	n	f	d	x	i	e	...	0	0	0	0	0	0	0	0	0
4	13	78.02	az	v	n	f	d	h	d	n	...	0	0	0	0	0	0	0	0	0

5 rows × 366 columns

```
In [10]: # Print the original and modified shape of the train dataset
print('Shape of original train dataset is:', train_original.shape)
print('Shape of modified train dataset after removing features having zero variance
```

```
Shape of original train dataset is: (4209, 378)
Shape of modified train dataset after removing features having zero variance is: (4209, 366)

The modified train dataset contains 366 features, which means the 12 features with zero variance from the original dataset is removed.
```

```
In [11]: # Save the modified train dataset
train_original_modified.to_csv('train_original_modified.csv',index=False)
```

Similarly, Read the Test Dataset and remove the zero variance features from the dataset

```
In [12]: #import the test dataset
test_original = pd.read_csv('test.csv')

# let us understand the test data
print('Size of testing set: {} rows and {} columns'.format(*test_original.shape))

# print few rows and see how the data looks like
test_original.head()
```

ID	X0	X1	X2	X3	X4	X5	X6	X8	X10	...	X375	X376	X377	X378	X379	X380	X382	X383	X384	
0	1	az	v	n	f	d	t	a	w	o	...	0	0	0	1	0	0	0	0	0
1	2	t	b	ai	a	d	b	g	y	o	...	0	0	1	0	0	0	0	0	0
2	3	az	v	as	f	d	a	j	j	o	...	0	0	0	1	0	0	0	0	0
3	4	az	i	n	f	d	z	i	n	o	...	0	0	0	1	0	0	0	0	0
4	5	w	s	as	c	d	y	i	m	o	...	1	0	0	0	0	0	0	0	0

5 rows × 377 columns

- There are 4209 data points and 377 features in the test dataset.
- Number of features in the test dataset is 377 because, the target feature 'y' is intentionally removed from the test dataset compared to train dataset.

```
In [13]: # Check the original variance of all the features in the test dataset and store it to
test_original_var=pd.DataFrame(test_original.var(axis=0),columns=['Variance'])
test_original_var
```

Variance	
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

360 rows × 1 columns

```
In [14]: # Call the function to return the test dataset features having zero variance.
features_zero_var(test_original)
```

Variance	
X257	0.0
X258	0.0
X295	0.0
X296	0.0
X369	0.0

- Above listed features have zero variance.
- However, since test dataset is not considered for training and only used for testing, we can remove the same features with zero variance in train dataset.
- This will ensure the same size and shape of the train and test dataset.

```
In [15]: # In test dataset, remove the same features of train dataset having zero variance.
test_original_modified=test_original.drop(columns=['X11', 'X93', 'X107', 'X233', 'X235', 'X268', 'X289', 'X290', 'X293', 'X297', 'X330', 'X347'])
```

```
In [16]: # Display the modified test dataset after removing the features having zero variance
test_original_modified.head()
```

ID	X0	X1	X2	X3	X4	X5	X6	X8	X10	...	X375	X376	X377	X378	X379	X380	X382	X383	X384	
0	1	az	v	n	f	d	t	a	w	o	...	0	0	0	1	0	0	0	0	0
1	2	t	b	ai	a	d	b	g	y	o	...	0	0	1	0	0	0	0	0	0
2	3	az	v	as	f	d	a	j	j	o	...	0	0	0	1	0	0	0	0	0
3	4	az	i	n	f	d	z	i	n	o	...	0	0	0	1	0	0	0	0	0
4	5	w	s	as	c	d	y	i	m	o	...	1	0	0	0	0	0	0	0	0

5 rows × 365 columns

```
In [17]: # Print the original and modified shape of the test dataset
print('Shape of original test dataset is:', test_original.shape)
print('Shape of modified test dataset after removing features having zero variance
```

```
Shape of original test dataset is: (4209, 377)
Shape of modified test dataset after removing features having zero variance is: (4209, 365)

The modified test dataset contains 365 features. This is because, we had removed the same 12 features having zero variance in the train dataset.
```

- The modified test dataset contains 365 features. This is because, we had removed the same 12 features having zero variance in the train dataset.
- The modified test dataset has 365 features, whereas modified train dataset has 366 features. This is because target feature 'y' is missing in the train dataset.

```
In [18]: # Save the modified test dataset
test_original_modified.to_csv('test_original_modified.csv',index=False)
```

Check for the null and unique values in the train dataset

```
In [19]: # Check the null values in the train dataset
print('The sum of null values in the train dataset is:', train_original_modified.isnull().sum())
```

The sum of null values in the train dataset is: 0

There are no null values in the train dataset.

The unique() function includes the missing value. But the nunique() function excludes the missing value as the default parameter is dropna=True. Since, there are no missing values in the train dataset, we can use nunique().

```
In [20]: # Check the unique values in the train dataset
train_original_modified_UV=pd.DataFrame(train_original_modified.nunique(),columns=['Unique_Value', 'train_original_modified_UV'])
train_original_modified_UV
```

Unique_Value	ID
Y	2545
X0	47
X1	27
X2	44
...	...
X380	2
X382	2
X383	2
X384	2
X385	2

366 rows × 1 columns

- It is observed that all the values in ID column are unique. In the provided dataset, the ID represents the unique car configuration. So, this feature must be ignored for training as it will not make any sense to the prediction.
- 'y' feature is the target feature.
- Features X0,X1,X2,X3,X4,X5,X6,X8 are the categorical features which must be converted to numerical values/one hot encoded values.
- All the features after X8 are having binary values.

```
In [21]: # Print the train dataset features unique values where the values = 2 and values >2.
print('Train Features with unique values greater than 2 are as follows:\n',train_original_modified_UV[train_original_modified_UV.Unique_Value>2])
```

```
Train Features with unique values greater than 2 are as follows:
Unique_Value  ID
Y            2545
X0            47
X1            27
X2            44
X3            7
X4            4
X5            29
X6            12
X8            25
dtype: int64
```

```
Test Features with unique values equal to 2 are as follows:
Unique_Value  ID
X10           2
X12           2
X13           2
X14           2
X15           2
X380          2
X382          2
X383          2
X384          2
X385          2
Length: 356, dtype: int64
```

- Out of 366 train dataset features, 10 features are having greater than 2 unique values and remaining features are having only 2 unique values (0 and 1).

Similarly, check the null and unique values in the test dataset

```
In [22]: # Check the null values in the test dataset
print('The sum of null values in the test dataset is:', test_original_modified.isnull().sum())
```

The sum of null values in the test dataset is: 0

- There are no null values in the test dataset.

```
In [23]: # Check the unique values in the test dataset
test_original_modified_UV=pd.DataFrame(test_original_modified.nunique(),columns=['Unique_Value', 'test_original_modified_UV'])
test_original_modified_UV
```

Unique_Value	ID
X0	49
X1	27
X2	45
X3	7
...	...
X380	2
X382	2
X383	2
X384	2
X385	2

365 rows × 1 columns

- It is observed that all the value in ID column are unique. In the provided dataset, the ID represents the unique car configuration. So, this feature must be ignored for testing as it will not make any sense to the prediction.
- Features X0,X1,X2,X3,X4,X5,X6,X8 are the categorical features which must be converted to numerical values/one hot encoded values.
- All the features after X8 are having binary values.

```
In [24]: # Print the test dataset features unique values where the values = 2 and values >2.
print('Test Features with unique values greater than 2 are as follows:\n',test_original_modified_UV[test_original_modified_UV.Unique_Value>2])
```

```
Test Features with unique values greater than 2 are as follows:
Unique_Value  ID
X0            49
X1            27
X2            45
X3            7
X4            4
X5            32
X6            12
X8            25
dtype: int64
```

```
Test Features with unique values equal to 2 are as follows:
Unique_Value  ID
X10           2
X12           2
X14           2
X15           2
X380          2
X382          2
X383          2
X384          2
X385          2
Length: 351, dtype: int64
```

- Out of 365 test dataset features, 9 features are having greater than 2 unique values and remaining features are having only 2 unique values (0 and 1).

Apply the label encoder/one hot encoding for the train dataset

Before applying label encoder, separate the ID and 'y' features from the train dataset.

```
In [25]: #Drop the columns 'ID' and 'y' and store the data into new object train_X_check and y
train_X_check = train_original_modified.drop(columns=['ID','y'])
# train_X_check to verify how the one hot encoding works for the features with multi
train_X_check.shape
```

Out[25]: (4209, 364)

```
In [26]: # Perform label encoder/one hot encoding for the categories features of train_X_check
# import the required library
from sklearn.preprocessing import LabelEncoder
```

```
# Define the function to apply the label encoder for the categories features of train_X_check
def label_encoder(df,x):
    # select only the features with datatype Object
    features_cat=df.select_dtypes(include='object').columns
    # instantiate the label encoder
    le=LabelEncoder()
    for i in features_cat:
        # Fit,transform
        # And replace with label-encoded data for the existing data in the object data
        x[i]=le.fit_transform(x[i])
```

```
In [27]: # Call the function to apply label encoder for the train_X_check data
label_encoder(train_original_modified_X_train_X_check)
```

```
In [28]: # verify the train_X_check dataset
train_X_check.head()
```

X0	X1	X2	X3	X4	X5	X6	X8	X10	X12	...	X375	X376	X377	X378	X379	X380	X382	X383	X384
0	32	17	0	4	3	28	1	14	0	...	0	0	0	1	0	0	0	0	0
1	32	21	19	4	3	24	1	9	0	...	1	0	0	0	0	0	0	0	0
2	20	24	34	2	3	27	9	23	0	...	0	0	0	0	0	0	0	0	1
3	20	21	34	5	3	27	11	4	0	...	0	0	0	0	0	0	0	0	0
4	20	23	34	5	3	12	3	13	0	...	0	0	0	0	0	0	0	0	0

5 rows × 364 columns

```
In [29]: train_X_check.X5.nunique()
```

Out[29]: 29

- It is observed from the above table, that the label encoder is applied to categorical features of train_X_check data, the encoded labels are not binary (0 and 1) since the features has 4 or more different categories.
- Example: X5 feature has 29(0 to 28) unique categories. So, when label encoder is applied, the categories will be replaced with 28 unique values starting from 0 to 28.
- This may impact the accuracy level.
- In order to fix this issue, we have to perform one hot encoding for multi-categorical variables, the procedure is follows as
 - Identify the top 10 most frequent categories from each feature.
 - Perform one hot encoding only for the top 10 most frequent categories.
 - All Top 10 most frequent categories will be considered as '1' and all the remaining categories will be considered as '0' in each feature.
 - By performing above 3 steps ensures only binary values (0 and 1) in all the features

```
In [30]: # Identify the top 10 most frequent categories of features X0,X1,X2,X3,X4,X5,X6,X8
#X0
train_original_modified.X0.value_counts().sort_values(ascending=False).head(10)
```

```
Out[30]: z            360
ak            349
y             324
ay            313
t             306
x             300
o             269
f             227
n             195
w             182
Name: X0, dtype: int64
```

```
In [31]: # create a list for the top 10 most frequent categories of feature X0

```


verify the train dataset after applying one hot encoding for the features X0,X1,X2,...

```
In [46]: top_10_test_ID = [x for x in test_original_modified.X0.value_counts().sort_values(ascending=False).head(10)]\ntest_original_modified.head(2)
```

```
Out[46]: ID    X0   X1   X2   X3   X4   X5   X6   X8   X10 ... Xak     Xy   XQ_z   XQ_x   X0_y   X0_t   X0_e   X0_f   X0_w\n1      0  az v n f d t a w o - 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0\n1      1  t b ai a d b g y o - 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0\n\n2 rows × 375 columns
```

In [47]: # X1\ntop_10_test_X1 = [x for x in test_original_modified.X1.value_counts().sort_values(ascending=False).head(10)]\ntest_original_modified.head(2)

```
Out[47]: ID    X0   X1   X2   X3   X4   X5   X6   X8   X10 ... X1_ae   X1_s   X1_b   X1_v   X1_r   X1_l   X1_a   X1_e   X1_c\n0      0  az v n f d t a w o - 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0\n1      1  t b ai a d b g y o - 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0\n\n2 rows × 385 columns
```

In [48]: # X2\ntop_10_test_X2 = [x for x in test_original_modified.X2.value_counts().sort_values(ascending=False).head(10)]\ntest_original_modified.head(2)

```
Out[48]: ID    X0   X1   X2   X3   X4   X5   X6   X8   X10 ... X2_as   X2_ae   X2.ai   X2.m   X2.ak   X2.r   X2.n   X2_s   X2_c\n0      0  az v n f d t a w o - 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0\n1      1  t b ai a d b g y o - 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0\n\n2 rows × 395 columns
```

In [49]: # X3\ntop_10_test_X3 = [x for x in test_original_modified.X3.value_counts().sort_values(ascending=False).head(10)]\ntest_original_modified.head(2)

```
Out[49]: ID    X0   X1   X2   X3   X4   X5   X6   X8   X10 ... X3_s   X3.f   X3.e   X3.c   X3.f   X3.a   X3.d   X3.g   X3.e   X3_g\n0      0  az v n f d t a w o - 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0\n1      1  t b ai a d b g y o - 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 1\n\n2 rows × 402 columns
```

- In case X3, only 7 columns were added as the maximum categorical variable for this feature is only 7.

In [50]: # X4\ntop_10_test_X4 = [x for x in test_original_modified.X4.value_counts().sort_values(ascending=False).head(10)]\ntest_original_modified.head(2)

```
Out[50]: ID    X0   X1   X2   X3   X4   X5   X6   X8   X10 ... X3_f   X3_a   X3.d   X3.g   X3.e   X3.b   X4.d   X4.b   X4.a : \n0      0  az v n f d t a w o - 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0\n1      1  t b ai a d b g y o - 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0\n\n2 rows × 406 columns
```

- In case X4, only 4 columns were added as the maximum categorical variable for this feature is only 4.

In [51]: # X5\ntop_10_test_X5 = [x for x in test_original_modified.X5.value_counts().sort_values(ascending=False).head(10)]\ntest_original_modified.head(2)

```
Out[51]: ID    X0   X1   X2   X3   X4   X5   X6   X8   X10 ... X5_v   X5_f   X5.p   X5.w   X5.af   X5.ad   X5.ac   X5_n   X5.l :\n0      0  az v n f d t a w o - 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0\n1      1  t b ai a d b g y o - 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0\n\n2 rows × 416 columns
```

In [52]: # X6\ntop_10_test_X6 = [x for x in test_original_modified.X6.value_counts().sort_values(ascending=False).head(10)]\ntest_original_modified.head(2)

```
Out[52]: ID    X0   X1   X2   X3   X4   X5   X6   X8   X10 ... X6.g   X6.i   X6.d   X6.j   X6.J   X6.h   X6.a   X6.k   X6.c   X6_e\n0      0  az v n f d t a w o - 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0\n1      1  t b ai a d b g y o - 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0\n\n2 rows × 426 columns
```

In [53]: # X8\ntop_10_test_X8 = [x for x in test_original_modified.X8.value_counts().sort_values(ascending=False).head(10)]\ntest_original_modified.head(2)

```
Out[53]: ID    X0   X1   X2   X3   X4   X5   X6   X8   X10 ... X8_e   X8_j   X8.s   X8.f   X8.n   X8_i   X8.r   X8.a   X8.w   X8_l :\n0      0  az v n f d t a w o - 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0\n1      1  t b ai a d b g y o - 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0\n\n2 rows × 436 columns
```

- The label encoding/one hot encoding is successfully applied to the test dataset

In [54]: # Store the data set in test_original_modified to new object test_original_modified_OH
drop the columns which are not required after performing one hot encoding\ntest_original_modified_OHE = test_original_modified.drop(columns=['X0','X1','X2','X3','X4','One hot encoded_OHE_data'])\ntest_original_modified_OHE

```
Out[54]: ID    X0   X1   X2   X3   X4   X5   X6   X8   X10 ... X8_e   X8_j   X8.s   X8.f   X8.n   X8_i   X8.r   X8.a   X8.w   X8_l :\n0      0  1  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0\n1      1  2  0  0  0  0  0  0  0  0  0  0  0  0  1  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0\n2      2  3  0  0  0  0  0  0  0  0  0  0  0  0  1  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0\n3      3  4  0  0  0  0  0  0  0  0  0  0  0  0  0  1  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0\n4      4  5  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0\n...\n...          .....\n4204   8410  0  0  0  0  0  0  1  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0\n4205   8411  0  0  0  0  0  0  0  0  0  0  0  0  0  0  1  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0\n4206   8413  0  0  0  0  0  0  1  0  0  0  0  0  0  0  0  1  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0\n4207   8414  0  0  0  0  0  1  1  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0\n4208   8416  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  1  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0\n\n4209 rows × 428 columns
```

In [55]: # Save the modified one hot encoded train dataset\ntest_original_modified_OHE.to_csv('test_original_modified_OHE.csv',index=False)

In [56]: # Check the shape of the modified one hot encoded train dataset\ntest_original_modified_OHE.shape

```
Out[56]: (4209, 428)
```

- The shape of the train dataset is reduced to 428 from 436 since 8 features were dropped after performing one hot encoding

Perform dimensionality reduction.

Before performing dimensional reduction (aka) principal component analysis (PCA), separate the following:

- Separate the features 'ID' and 'Y' from train_original_modified_OHE dataset and store it in the new object.
- Separate the feature 'ID' from test_original_modified_OHE.shape and store it in the new object.

In [57]: # store the 'ID' values into the new object train_ID and test_ID and verify the shape

```
# Train dataset\ntrain_ID = train_original_modified_OHE.ID\nprint('The shape of the ID feature in train dataset is:',train_ID.shape)\n\n# Test dataset\ntest_ID = test_original_modified_OHE.ID\nprint('The shape of the ID feature in test dataset is:',test_ID.shape)
```

The shape of the ID feature in train dataset is: (4209,)
The shape of the ID feature in test dataset is: (4209,)

In [58]: # store the remaining values into the new object train_X and test_X and verify the shape

```
# Train dataset\ntrain_X=train_original_modified_OHE.drop(columns=['ID'],'y')\nprint('The shape of the final train dataset is:',train_X.shape)\n\n# Test dataset\ntest_X=test_original_modified_OHE.drop(columns=['ID'])\nprint('The shape of the final test dataset is:',test_X.shape)
```

The shape of the final train dataset is: (4209, 427)
The shape of the final test dataset is: (4209, 427)

In [59]: # store the 'Y' values into the new object train_Y and verify the shape

```
# store the Y values into the
```