

Project to Submit – Project 1

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 XGBoost.

Write Up:

Analysis Tasks to be performed:

1. Data was already divided into two files for training and test dataset.
2. Both had equal number of rows

```
In [5]: train_df.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
```

```
In [6]: test_df.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
```

```
In [7]: train_df.describe()
```

```
Out[7]:
```

	ID	y	X10	X11	X12	X13	X14	X15	X16	X17	...	X375	X385
count	4209.000000	4209.000000	4209.000000	4209.0	4209.000000	4209.000000	4209.000000	4209.000000	4209.000000	4209.000000	...	4209.000000	4209.000000
mean	4205.960798	100.669318	0.013305	0.0	0.075077	0.057971	0.428130	0.000475	0.002613	0.007603	...	0.318841	0.057
std	2437.608688	12.679381	0.114590	0.0	0.263547	0.233716	0.494867	0.021796	0.051061	0.086872	...	0.466082	0.232
min	0.000000	72.110000	0.000000	0.0	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	...	0.000000	0.000
25%	2095.000000	90.820000	0.000000	0.0	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	...	0.000000	0.000
50%	4220.000000	99.150000	0.000000	0.0	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	...	0.000000	0.000
75%	6314.000000	109.010000	0.000000	0.0	0.000000	0.000000	1.000000	0.000000	0.000000	0.000000	...	1.000000	0.000
max	8417.000000	265.320000	1.000000	0.0	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000	...	1.000000	1.000

8 rows x 370 columns

```
In [8]: test_df.describe()
```

```
Out[8]:
```

	ID	X10	X11	X12	X13	X14	X15	X16	X17	X18	...	X375
count	4209.000000	4209.000000	4209.000000	4209.000000	4209.000000	4209.000000	4209.000000	4209.000000	4209.000000	4209.000000	...	4209.000000
mean	4211.039202	0.019007	0.000238	0.074364	0.061060	0.427893	0.000713	0.002613	0.008791	0.010216	...	0.325968
std	2423.078926	0.136565	0.015414	0.262394	0.239468	0.494832	0.026691	0.051061	0.093357	0.100570	...	0.468791
min	1.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	...	0.000000
25%	2115.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	...	0.000000

-
-
- 3.
4. Training dataset as usual had y column which is the target column which was missing in test data set

```
In [9]: print("Number of datapoints: ", train_df.shape[0])
        print("Number of features: ", train_df.shape[1])
        train_df.head()
```

```
Number of datapoints: 4209
Number of features: 378
```

Out[9]:

	ID	y	X0	X1	X2	X3	X4	X5	X6	X8	...	X375	X376	X377	X378	X379	X380	X382	X383	X384	X385
0	0	130.81	k	v	at	a	d	u	j	o	...	0	0	1	0	0	0	0	0	0	0
1	6	88.53	k	t	av	e	d	y	l	o	...	1	0	0	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	1	0	0	0
3	9	80.62	az	t	n	f	d	x	l	e	...	0	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	0

5 rows x 378 columns

```
In [10]: print("Number of datapoints: ", test_df.shape[0])
          print("Number of features: ", test_df.shape[1])
          test_df.head()
```

```
Number of datapoints: 4209
Number of features: 377
```

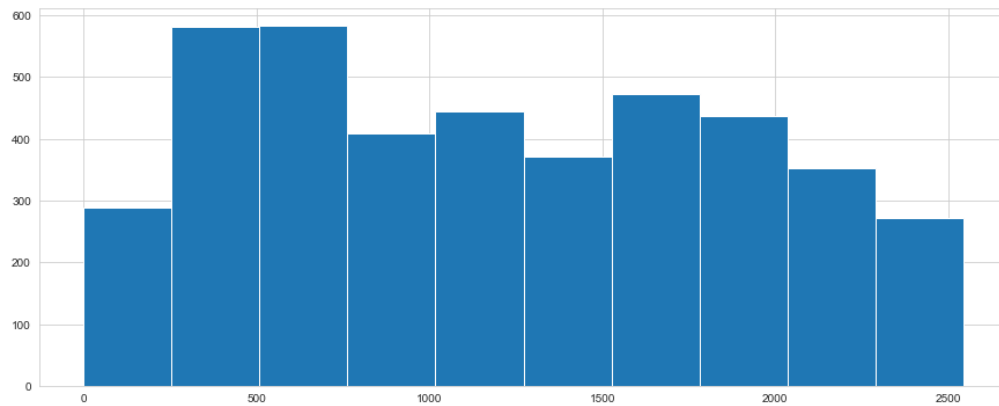
Out[10]:

	ID	X0	X1	X2	X3	X4	X5	X6	X8	X10	...	X375	X376	X377	X378	X379	X380	X382	X383	X384	X385
0	1	az	v	n	f	d	t	a	w	0	...	0	0	0	1	0	0	0	0	0	0
1	2	t	b	ai	a	d	b	g	y	0	...	0	0	1	0	0	0	0	0	0	0
2	3	az	v	as	f	d	a	j	j	0	...	0	0	0	1	0	0	0	0	0	0
3	4	az	l	n	f	d	z	l	n	0	...	0	0	0	1	0	0	0	0	0	0
4	5	w	s	as	c	d	y	i	m	0	...	1	0	0	0	0	0	0	0	0	0

5 rows x 377 columns

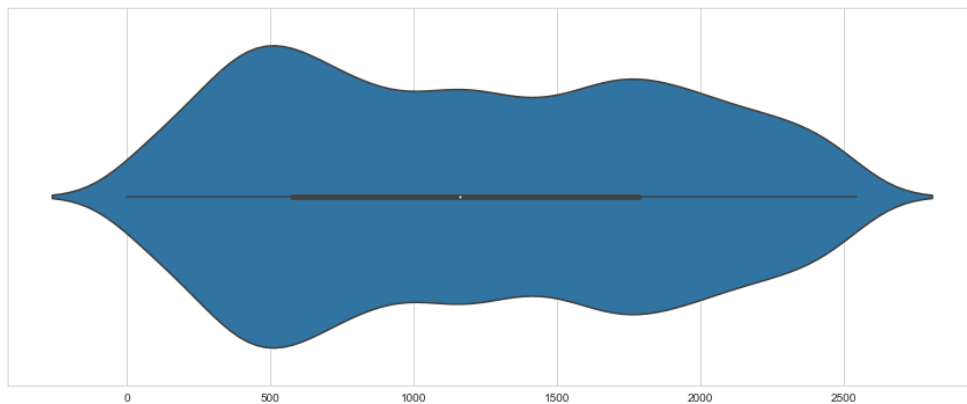
5. Training data set had 369 binary features, 8 features which have datatype = 'object' is most probably categorical features and 1 remaining feature is our target variable i.e. 'y'.
6. Performing univariate analysis on categorical features, to get the insight out of it. Any feature that has very low variance as compared to other categorical features, will be removed
7. Found 12 columns which had no variance. We removed them.
 1. Name = X11
 2. Name = X93
 3. Name = X107
 4. Name = X233
 5. Name = X235
 6. Name = X268
 7. Name = X289
 8. Name = X290
 9. Name = X293
 10. Name = X297
 11. Name = X330
 12. Name = X347
 13. No of columns which has zero variance = 12
8. Also checked for null, duplicate rows and unique values
9. Got the details of categorical columns and integers columns separatel
10. Then applied encoder
11. Then Summarize outcome (testing time) in training dataset (created box-plot to see how the data is spread)
12. On y, we plotted histogram & violin plot to see if there are any outliers.

Out[34]: <AxesSubplot:>



In [35]: `sns.violinplot(train_df['y'].values)`

Out[35]: <AxesSubplot:>



13. Perform dimensionality reduction.
14. The methods at our disposal using linear algebra are:
15. Principal Components Analysis Singular Value Decomposition Non-Negative Matrix Factorization. Identified 6 main components.
16. Created the dataset to include only these components for further analysis
17. Before using XGBoost, checked the model with the following methods
 - Logistic Regression
 - KNN
 - SVM
 - Random Forest
18. None of the model was good. We tried XGBoost.
19. For the first time we got good >95% accuracy.
20. We further tested with different booster='dart', 'gbliner', 'gbtree'. Gbtree gave the best result, and
21. We enhanced it further with K-fold.
22. We got predicted y with accuracy of 99.9%

Code:

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 XGBoost.

The data set is already divided into train and test

```
import numpy as np
```

```
import pandas as pd
```

```
from datetime import datetime as dt

import seaborn as sns

import matplotlib.pyplot as plt

%matplotlib inline

from matplotlib.pylab import rcParams

rcParams['figure.figsize'] = 15, 6


import warnings

warnings.filterwarnings('ignore')


# importing csv module

import csv

# csv file name

train_df = pd.read_csv(r'D:\OneDrive\Studies\AI - ML\Python\Examples\ML Pracs\train.csv')


# importing csv module

import csv

# csv file name

test_df = pd.read_csv(r'D:\OneDrive\Studies\AI - ML\Python\Examples\ML Pracs\test.csv')


train_df.info()

test_df.info()

train_df.describe()

test_df.describe()

print("Number of datapoints: ", train_df.shape[0])

print("Number of features: ", train_df.shape[1])

train_df.head()
```

```
print("Number of datapoints: ", test_df.shape[0])
print("Number of features: ", test_df.shape[1])
test_df.head()
```

```
dtype_df = train_df.dtypes.reset_index()
dtype_df.columns = ["feature name", "dtypes"]
dtype_df.groupby("dtypes").agg("count").reset_index()
```

there are 369 binary features,
8 features which have datatype = 'object' is most probably categorical features and
1 remaining feature is our target variable i.e. 'y'.

Performing univariate analysis on categorical features, to get the insight out of it.
Any feature that has very low variance as compared to other categorical features, will be removed

```
dtype_df = test_df.dtypes.reset_index()
dtype_df.columns = ["feature name", "dtypes"]
dtype_df.groupby("dtypes").agg("count").reset_index()
```

Question 1:

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

Starting with train and then with test data

```
variance = pow(train_df.drop(columns={'ID', 'y'}).std(), 2).to_dict()
```

```
null_cnt = 0
```

```
for key, value in variance.items():
```

```
    if(value==0):
```

```
        print('Name = ', key)
```

```
        null_cnt = null_cnt+1
```

```
print('No of columns which has zero variance = ',null_cnt)
```

```
train_df = train_df.drop(columns={'X11','X93','X107','X233','X235','X268','X289','X290','X293','X297','X330','X347'})
```

```
train_df.shape
```

```
variance = pow(test_df.drop(columns={'ID'}).std(),2).to_dict()
```

```
null_cnt = 0
```

```
for key, value in variance.items():
```

```
    if(value==0):
```

```
        print('Name = ',key)
```

```
        null_cnt = null_cnt+1
```

```
print('No of columns which has zero variance = ',null_cnt)
```

```
train_df = train_df.drop(columns={'X257','X258','X295','X296','X369'})
```

```
train_df.shape
```

Question 2:

Check for null and unique values for test and train sets.

```
print(train_df.nunique())
```

```
print(test_df.nunique())
```

```
#Check for null value
```

```
print(train_df.isnull().sum().any())
```

```
print(test_df.isnull().sum().any())
```

```
train_df.describe(include='object')
```

```
test_df.describe(include='object')
```



```
dup_ID = train_df['ID'].duplicated().sum()
print(f"Here we have {dup_ID} duplicate IDs")
```

No null data, all unique values across the file listed.

Henceforth working with Train data only as it is the data that we would use for our model.

Question 3:

Apply label encoder.

No null variable. All the variables are categorical applying encoder

```
from sklearn.preprocessing import LabelEncoder
le=LabelEncoder()
for i in train_df.columns:
    train_df[i]=le.fit_transform(train_df[i])
```

```
train_df.head()
```

```
train_df.corr()
```

Summarize outcome (testing time) in training dataset

Draw a vertical boxplot grouped

by a categorical variable: X0

```
sns.set_style("whitegrid")
```

```
object_columns = test_df.describe(include='object').columns
```

```
print('\nobject columns:\n',object_columns)
```

```
cols = len(object_columns)
```

```
sns.boxplot(x = 'X0', y = 'y', data = train_df)
```

```
sns.boxplot(x = 'X1', y = 'y', data = train_df)
```

```
sns.boxplot(x = 'X2', y = 'y', data = train_df)
```

```
sns.boxplot(x = 'X3', y = 'y', data = train_df)
```

```
sns.boxplot(x = 'X4', y = 'y', data = train_df)
```

```
sns.boxplot(x = 'X5', y = 'y', data = train_df)
```

```
sns.boxplot(x = 'X6', y = 'y', data = train_df)
```

```
sns.boxplot(x = 'X8', y = 'y', data = train_df)
```

```
#Now the target y
```

```
train_df['y'].hist()
```

```
sns.violinplot(train_df['y'].values)
```

The data seems optimized. The removal of few data points which had no variance had optimized the y. No reason to test for any more outliers.

Ideally, this test is done first, but if variance 0 is removed, it increases the chances of y being optimized, with no outliers.

Dimensionality reduction refers to techniques for reducing the number of input variables in training data.

Fewer input dimensions often means correspondingly fewer parameters or a simpler structure in the machine learning model, referred to as degrees of freedom. A model with too many degrees of freedom is likely to overfit the training dataset and may not perform well on new data.

It is desirable to have simple models that generalize well, and in turn, input data with few input variables. This is particularly true for linear models where the number of inputs and the degrees of freedom of the model are often closely related.

Dimensionality reduction is a data preparation technique performed on data prior to modeling. It might be performed after data cleaning and data scaling and before training a predictive model.

Question 4: Perform dimensionality reduction.

The methods at our disposal using linear algebra are:

Principal Components Analysis

Singular Value Decomposition

Non-Negative Matrix Factorization

```
# Draw a vertical boxplot grouped
```

```
# by a categorical variable: X0
```

```
train_df.describe(include='int64')
```

```
bin_columns = train_df['ID']
```

```
print('\nobject columns:\n',bin_columns)
```

```
from sklearn.preprocessing import StandardScaler
```

```
scaler = StandardScaler()
```

```
scaler.fit(train_df)
```

```
sdata = scaler.transform(train_df)
```

```
from sklearn.decomposition import PCA
```

```
sdata.shape
```

```
# lets take top 6 pca components
```

```
pca = PCA(n_components=6)
```

```
pca.fit(sdata)
```

```
x_pca = pca.transform(sdata)
```

```
x_pca.shape
```

```
# number of components
```

```
n_pcs= pca.components_.shape[0]
```

```
n_pcs
```

```
# get the index of the most important feature on EACH component i.e. largest  $\lambda$ ,  $\rightarrow$  absolute value
```

```
# using LIST COMPREHENSION HERE
```

```
most_important = [np.abs(pca.components_[i]).argmax() for i in range(n_pcs)]
```

```
initial_feature_names = bin_columns
```

```
most_important
```

```
# using LIST COMPREHENSION HERE AGAIN
```

```
dic = {'PC{}'.format(i): most_important[i] for i in range(n_pcs)}
```

```
dic
```

```
pca.components_
```

```
explained_variance = pca.explained_variance_ratio_
```

```
explained_variance
```

```
#it is a measure of the variance of the data when projected onto that axis. The projection of each data point onto the
```

#principal axes are the “principal components” of the data. .4 is the var of PCA

#and .179 is the var of PCA2

#Creating training and test data with only these columns

```
selected_columns = train_df[['ID', 'X325', 'X27', 'X180', 'X161', 'X309', 'X85', 'y']]
```

```
X_train = selected_columns.copy()
```

```
X_train.shape
```

```
X_train
```

#Now creating a df with only these 5 components

```
selected_columns = test_df[['ID', 'X325', 'X27', 'X180', 'X161', 'X309', 'X85']]
```

```
X_test = selected_columns.copy()
```

```
X_test.shape
```

```
X_test
```

#now will perform XGBoost

Predict your test_df values using XGBoost.

Model Selection

Logistic Regression

KNN

SVM

Random Forest

#Now splitting the data into train & test. Before that, identifying all input parameters as X, and output parameter as y

```
y_train=train_df['y']
```

```
y_train
```

```
y_train.shape
```

```
from sklearn.model_selection import learning_curve  
from sklearn.linear_model import LogisticRegression  
from sklearn.neighbors import KNeighborsClassifier  
from sklearn.model_selection import cross_val_score  
from sklearn.metrics import classification_report, confusion_matrix
```

```
import xgboost as xgb  
from sklearn.metrics import r2_score  
from sklearn.model_selection import train_test_split
```

```
X_train.shape
```

```
X_test.shape
```

```
y_train.shape
```

```
#We do not have X & y. Creating X
```

```
X = X_train  
#X = pd.concat([X_train,X_test])  
print(X)
```

```
X.shape
```

```
#y = pd.concat([y_train,y_train])  
y = y_train
```

```
y.shape
```

```
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_state=72)
```

```
# Logistic Regression
```

```
logreg=LogisticRegression(solver='liblinear',multi_class='ovr')
```

```
logreg.fit(X_train,y_train)
```

```
y_pred=logreg.predict(X_test)
```

```
y_pred
```

```
#Accuracy Score
```

```
#print(metrics.accuracy_score(y_pred,y_train))
```

```
accuracy = (logreg.score(X_train,y_train))
```

```
print(accuracy)
```

```
# Logistic Regression
```

```
logreg=LogisticRegression(solver='lbfgs',multi_class='auto')
```

```
logreg.fit(X_train,y_train)
```

```
y_pred=logreg.predict(X_test)
```

```
y_pred
```

```
#Accuracy Score
```

```
#print(metrics.accuracy_score(y_pred,y_train))
```

```
accuracy = (logreg.score(X_train,y_train))
```

```
print(accuracy)
```

```
#SVM "Support Vector Classifier"
```

```
from sklearn.svm import SVC
```

```
svm = SVC(kernel='linear')
```

```
# fitting x samples and y classes
```

```
svm.fit(X_train,y_train)
```

```
y_pred = svm.predict(X_test)
```

```
from sklearn import metrics

accuracy = metrics.accuracy_score(y_test, y_pred)

print(accuracy)
```

```
#KNN with 5 neighbours

knn = KNeighborsClassifier(n_neighbors=5)

knn.fit(X_train, y_train)

y_pred = knn.predict(X_test)

print(metrics.accuracy_score(y_test, y_pred))
```

```
print(classification_report(y_test,pred))
```

```
knn = KNeighborsClassifier(n_neighbors=1)

knn.fit(X_train, y_train)

y_pred = knn.predict(X_test)

print(metrics.accuracy_score(y_test, y_pred))
```

```
print(classification_report(y_test,pred))
```

```
avg_score=[]

for k in range(2,30):

    knn=KNeighborsClassifier(n_jobs=-1,n_neighbors=k)

    score=cross_val_score(knn,X_train,y_train,cv=5,n_jobs=-1,scoring='accuracy')

    avg_score.append(score.mean())
```

```
plt.figure(figsize=(12,8))

plt.plot(range(2,30),avg_score)

plt.xlabel("n_neighbours")

plt.ylabel("accuracy")

#plt.xticks(range(2,30,2))
```



```
#Random Forests Classifier

from sklearn.ensemble import RandomForestClassifier

from sklearn.metrics import f1_score


rfc=RandomForestClassifier(n_jobs=-1,random_state=51)

rfc.fit(X_train,y_train)

print(rfc.score(X_test,y_test))

print(f1_score(y_test,rfc.predict(X_test),average='macro'))


#Till now SVM followed by Random forest is the leading model
```

```
from xgboost import XGBClassifier

from sklearn.metrics import mean_squared_error

from sklearn import svm

from xgboost import XGBClassifier

import xgboost as xgb
```

```
#Here, we are using XGBRegressor as a Machine Learning model to fit the data.

model = xgb.XGBRegressor(booster='dart', objective='reg:squarederror', num_class = 1, eval_metric = 'merror',
n_estimators = 10, seed = 123)

model.fit(X_train, y_train)

print(); print(model)
```

```
# Predict the model

pred = model.predict(X_test)
```

```
# RMSE Computation

rmse = np.sqrt(mean_squared_error(y_test, pred))

print("RMSE : % f" %(rmse))
```

```
expected_y = y_test
predicted_y = model.predict(X_test)
print(metrics.r2_score(y_test, predicted_y))
```

```
predicted_y
```

```
#Here, we are using XGBRegressor as a Machine Learning model to fit the data.
```

```
model = xgb.XGBRegressor(booster='gblinear', objective='reg:squarederror', num_class = 1, eval_metric = 'merror',
n_estimators = 10, seed = 123)
```

```
model.fit(X_train, y_train)
```

```
print(); print(model)
```

```
# Predict the model
```

```
pred = model.predict(X_test)
```

```
# RMSE Computation
```

```
rmse = np.sqrt(mean_squared_error(y_test, pred))
```

```
print("RMSE : % f" %(rmse))
```

```
expected_y = y_test
```

```
predicted_y = model.predict(X_test)
```

```
print(metrics.r2_score(y_test, predicted_y))
```

```
#As we see gblinear is not the right model.
```

```
#Here, we are using XGBRegressor as a Machine Learning model to fit the data.
```

```
model = xgb.XGBRegressor(booster='gbtree', objective='reg:squarederror', num_class = 1, eval_metric = 'merror',
n_estimators = 10, seed = 123)
```

```
model.fit(X_train, y_train)
```

```
print(); print(model)
```

```
# Predict the model
```

```

pred = model.predict(X_test)

# RMSE Computation
rmse = np.sqrt(mean_squared_error(y_test, pred))
print("RMSE : % f" %(rmse))

expected_y = y_test
predicted_y = model.predict(X_test)
print(metrics.r2_score(y_test, predicted_y))

model = xgb.XGBRegressor()
model.fit(X_train, y_train)
print(); print(model)

plt.figure(figsize=(10,10))
sns.regplot(expected_y, predicted_y, fit_reg=True, scatter_kws={"s": 100})

#e'll check the training accuracy with cross-validation and k-fold methods.
# Applying k-Fold Cross Validation
from sklearn.model_selection import cross_val_score, KFold

kfold = KFold(n_splits=10, shuffle=True)

kf_cv_scores = cross_val_score(model, X_train, y_train, cv=kfold )
print("K-fold CV average score: %.2f" % kf_cv_scores.mean())

#Now we have predicted the output by passing X_test and also stored real target in expected_y.
expected_y = y_test
predicted_y = model.predict(X_test)
print(metrics.r2_score(expected_y, predicted_y))

#This is the best model so far

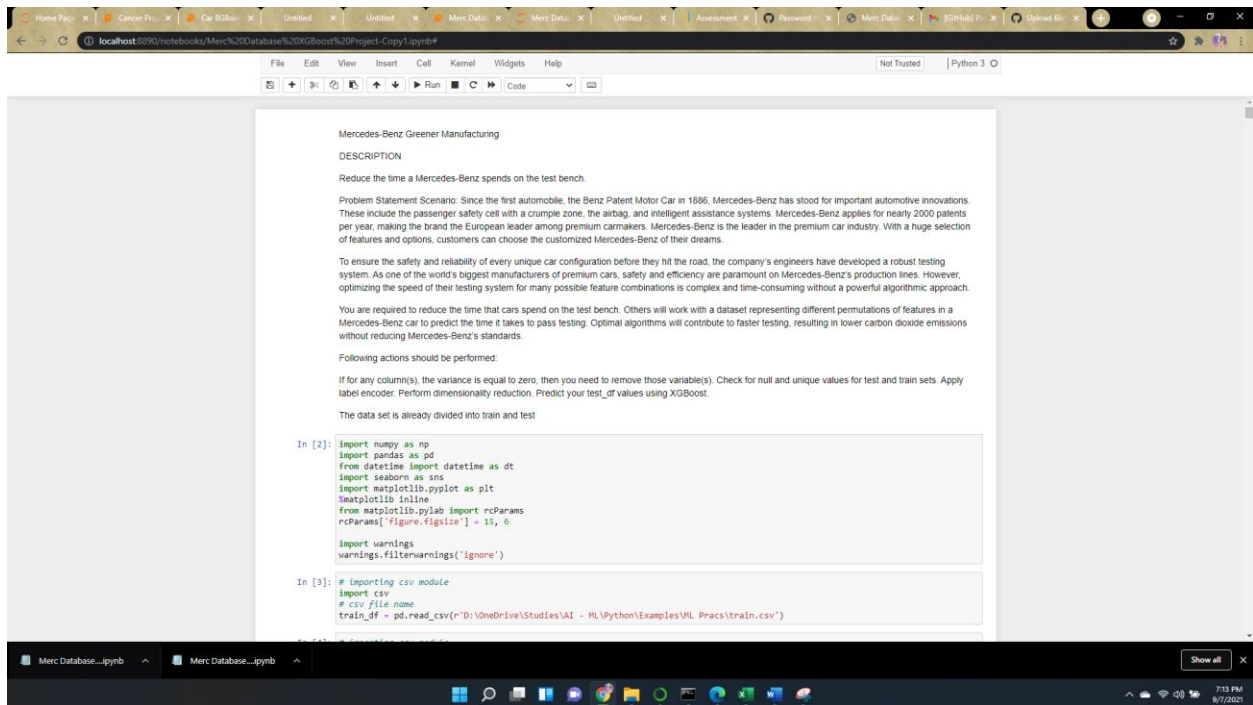
```

#Q 5. Predicting y with XGBoost

#Displaying predicted values

predicted_y

ScreenShots:



Home Page x Cancer Pro... x Car B... x Untitled x Untitled x Merc Data... x Merc Data... x Assessment x Password x Merc Data... x GitHub P... x Upload B... x

localhost:8090/notebooks/Merc%20Database%20XGBoost%20Project-Copy1.ipynb#

File Edit View Insert Cell Kernel Widgets Help Not Trusted Python 3

```
In [3]: # importing csv module
import csv
# csv file name
train_df = pd.read_csv(r'D:\OneDrive\Studies\AI - ML\Python\Examples\VL Pracs\train.csv')
```

```
In [4]: # importing csv module
import csv
# csv file name
test_df = pd.read_csv(r'D:\OneDrive\Studies\AI - ML\Python\Examples\VL Pracs\test.csv')
```

```
In [5]: train_df.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
```

```
In [6]: test_df.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
```

```
In [7]: train_df.describe()

Out[7]:
```

	ID	y	X10	X11	X12	X13	X14	X15	X16	X17	...	X375	...
count	4209.000000	4209.000000	4209.000000	4209.0	4209.000000	4209.000000	4209.000000	4209.000000	4209.000000	4209.000000	...	4209.000000	4209.000000
mean	4205.960798	100.680318	0.013305	0.0	0.075077	0.057971	0.428130	0.000475	0.002613	0.007603	...	0.318841	0.057
std	2437.008888	12.679381	0.114990	0.0	0.263547	0.233716	0.494867	0.021796	0.051061	0.086872	...	0.486082	0.232
min	0.000000	72.110000	0.000000	0.0	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	...	0.000000	0.000
25%	2095.000000	90.820000	0.000000	0.0	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	...	0.000000	0.000
50%	4220.000000	99.150000	0.000000	0.0	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	...	0.000000	0.000
75%	6314.000000	108.010000	0.000000	0.0	0.000000	0.000000	1.000000	0.000000	0.000000	0.000000	...	1.000000	0.000
max	6417.000000	265.320000	1.000000	0.0	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000	...	1.000000	1.000

8 rows x 370 columns

```
In [8]: test_df.describe()
```

Home Page x Cancer Pro... x Car B... x Untitled x Untitled x Merc Data... x Merc Data... x Assessment x Password x Merc Data... x GitHub P... x Upload B... x

localhost:8090/notebooks/Merc%20Database%20XGBoost%20Project-Copy1.ipynb#

File Edit View Insert Cell Kernel Widgets Help Not Trusted Python 3

8 rows x 370 columns

```
In [8]: test_df.describe()

Out[8]:
```

	ID	X10	X11	X12	X13	X14	X15	X16	X17	X18	...	X375	...
count	4209.000000	4209.000000	4209.000000	4209.000000	4209.000000	4209.000000	4209.000000	4209.000000	4209.000000	4209.000000	...	4209.000000	4209.000000
mean	4211.030202	0.019007	0.000238	0.074364	0.061065	0.427893	0.000713	0.002613	0.006701	0.010216	...	0.325968	...
std	2423.078926	0.136565	0.015414	0.262384	0.239466	0.494832	0.026991	0.051061	0.093357	0.109570	...	0.486791	...
min	1.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	...	0.000000	...
25%	2115.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	...	0.000000	...
50%	4202.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	...	0.000000	...
75%	6310.000000	0.000000	0.000000	0.000000	0.000000	1.000000	0.000000	0.000000	0.000000	0.000000	...	1.000000	...
max	6415.000000	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000	...	1.000000	...

8 rows x 369 columns

```
In [9]: print("Number of datapoints: ", train_df.shape[0])
print("Number of features: ", train_df.shape[1])
train_df.head()

Number of datapoints: 4209
Number of features: 378
```

```
Out[9]:
```

	ID	y	X0	X1	X2	X3	X4	X5	X6	X8	...	X375	X376	X377	X378	X379	X380	X382	X383	X384	X385
0	0	130.81	k	v	a	t	a	d	u	j	o	...	0	0	1	0	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	0
2	7	78.26	a	w	n	c	d	x	j	x	...	0	0	0	0	0	0	1	0	0	0
3	9	80.62	a	t	n	f	d	x	i	e	...	0	0	0	0	0	0	0	0	0	0
4	13	78.02	a	v	n	f	d	h	d	n	...	0	0	0	0	0	0	0	0	0	0

5 rows x 378 columns

```
In [10]: print("Number of datapoints: ", test_df.shape[0])
print("Number of features: ", test_df.shape[1])
test_df.head()

Number of datapoints: 4209
Number of features: 377
```



```
File Edit View Insert Cell Kernel Widgets Help
Not Trusted Python 3

print('No of columns which has zero variance = ',null_cnt)

Name = X257
Name = X258
Name = X295
Name = X296
Name = X369
No of columns which has zero variance = 5

In [16]: train_df = train_df.drop(columns=['X257','X258','X295','X296','X369'])
train_df.shape
Out[16]: (4289, 361)

Question 2: Check for null and unique values for test and train sets.

In [17]: print(train_df.nunique())
print(test_df.nunique())

ID
4289
y
2545
X0
47
X1
27
X2
44
...
X380
2
X382
2
X383
2
X384
2
X385
2
Length: 361, dtype: int64
ID
4289
y
49
X0
27
X2
45
X3
7
...
X380
2
X382
2
X383
2
X384
2
X385
2
Length: 377, dtype: int64

In [18]: #Check for null value
print(train_df.isnull().sum().any())
print(test_df.isnull().sum().any())
```

```
File Edit View Insert Cell Kernel Widgets Help
Not Trusted Python 3

Length: 377, dtype: int64

In [18]: #Check for null value
print(train_df.isnull().sum().any())
print(test_df.isnull().sum().any())

False
False

In [19]: train_df.describe(include='object')
Out[19]:
   X0  X1  X2  X3  X4  X5  X6  X8
count 4209 4209 4209 4209 4209 4209 4209 4209
unique  47  27  44   7   4  29  12  25
top     z  aa  as  c  d  w  g  j
freq  360  833 1659 1942 4205 231 1042 277

In [20]: test_df.describe(include='object')
Out[20]:
   X0  X1  X2  X3  X4  X5  X6  X8
count 4209 4209 4209 4209 4209 4209 4209 4209
unique  49  27  45   7   4  32  12  25
top    ak  aa  as  c  d  v  g  *
freq  432  826 1658 1900 4203 246 1073 274

In [21]: dup_ID = train_df['ID'].duplicated().sum()
print(f"Here we have {dup_ID} duplicate IDs")

Here we have 0 duplicate IDs

No null data, all unique values across the file listed. Henceforth working with Train data only as it is the data that we would use for our model.

Question 3: Apply label encoder

No null variable. All the variables are categorical applying encoder

In [22]: from sklearn.preprocessing import LabelEncoder
le=LabelEncoder()
for i in train_df.columns:
    train_df[i]=le.fit_transform(train_df[i])
```

Home Page | Cancer Pro... | Car B... | Untitled | Merc Data... | Merc Data... | Assessment | Password | Merc Data... | EDRHub | Upload file...

localhost:8790/notebooks/Merc%20Database%20XGBoost%20Project-Copy1.ipynb#

File Edit View Insert Cell Kernel Widgets Help Not Trusted Python 3

```

In [22]: from sklearn.preprocessing import LabelEncoder
         le=LabelEncoder()
         for i in train_df.columns:
             train_df[i]=le.fit_transform(train_df[i])

In [23]: train_df.head()
Out[23]:
```

	ID	y	X0	X1	X2	X3	X4	X5	X6	X8	X375	X376	X377	X378	X379	X380	X382	X383	X384	X385
0	0	2466	32	23	17	0	3	24	9	14	...	0	0	1	0	0	0	0	0	0
1	1	366	32	21	19	4	3	26	11	14	...	1	0	0	0	0	0	0	0	0
2	2	69	20	24	34	2	3	27	9	23	...	0	0	0	0	0	0	1	0	0
3	3	133	20	21	34	5	3	27	11	4	...	0	0	0	0	0	0	0	0	0
4	4	106	20	23	34	5	3	12	3	13	...	0	0	0	0	0	0	0	0	0

5 rows x 361 columns

```

In [24]: train_df.corr()
Out[24]:
```

	ID	y	X0	X1	X2	X3	X4	X5	X6	X8	X375	X376	X377	X378	
ID	1.000000	-0.053835	-0.012938	0.085511	-0.022195	-0.032842	0.018940	0.649727	-0.017728	0.006444	...	0.045307	-0.079988	-0.022990	0.030387
y	-0.053835	1.000000	-0.235347	0.012417	0.111803	-0.153171	-0.015482	-0.035435	0.001477	-0.003473	...	0.021055	0.131811	0.055796	-0.226880
X0	-0.012938	-0.235347	1.000000	-0.271123	-0.139804	-0.079645	0.017968	0.012293	0.037549	0.047735	...	0.113272	-0.070546	0.045173	-0.192136
X1	0.085511	0.012417	-0.271123	1.000000	0.088266	0.295957	-0.020724	0.048417	-0.079119	-0.000306	...	0.056874	-0.162424	-0.248791	0.145382
X2	-0.022195	0.111803	-0.139804	0.088266	1.000000	-0.083546	0.002289	-0.017722	0.085778	-0.069932	...	-0.174308	0.033697	0.125603	0.131974
X3	-0.032842	0.017968	0.012293	-0.020724	-0.083546	1.000000	0.004195	0.002911	0.010434	-0.014059	...	-0.061741	-0.022240	-0.061188	-0.013110
X4	0.649727	-0.035435	0.037549	0.048417	-0.017722	0.004195	1.000000	0.005533	-0.031128	0.054548	...	-0.000995	-0.059883	-0.021571	-0.059327
X5	-0.017728	0.001477	0.047735	-0.079119	-0.017722	0.002911	-0.031128	1.000000	0.007030	0.023867	...	0.008950	-0.014917	-0.005373	0.008884
X6	0.006444	-0.003473	0.047735	-0.079119	-0.017722	0.002911	0.007030	0.023867	1.000000	0.008950	...	0.008950	-0.014917	-0.005373	0.008884
X8	0.045307	0.021055	0.113272	0.056874	-0.174308	-0.061741	-0.000995	-0.059883	-0.021571	-0.059327	1.000000	...	0.025810	-0.025810	-0.005488
X375	-0.022990	0.055796	0.045173	-0.162424	0.033697	0.125603	0.061188	0.021571	0.005373	0.008884	0.025810	1.000000	...	0.025810	-0.005488
X376	-0.022990	0.055796	0.045173	-0.162424	0.033697	0.125603	0.061188	0.021571	0.005373	0.008884	0.025810	0.025810	1.000000	...	0.005488
X377	-0.022990	0.055796	0.045173	-0.162424	0.033697	0.125603	0.061188	0.021571	0.005373	0.008884	0.025810	0.025810	0.025810	1.000000	...
X378	-0.022990	0.055796	0.045173	-0.162424	0.033697	0.125603	0.061188	0.021571	0.005373	0.008884	0.025810	0.025810	0.025810	0.025810	1.000000

361 rows x 361 columns

Summarize outcome (testing time) in training dataset

```

In [25]: # Draw a vertical boxplot grouped
         # by a categorical variable: X0
         sns.set_style("whitegrid")

         object_columns = test_df.describe(include='object').columns
         print("\nobject columns:\n",object_columns)
         cols = len(object_columns)

         object_columns:
         Index(['X0', 'X1', 'X2', 'X3', 'X4', 'X5', 'X6', 'X8'], dtype='object')

In [26]: sns.boxplot(x = 'X0', y = 'y', data = train_df)
Out[26]: <AxesSubplot: xlabel='X0', ylabel='y'>
```

```

In [27]: sns.boxplot(x = 'X1', y = 'y', data = train_df)
Out[27]: <AxesSubplot: xlabel='X1', ylabel='y'>
```

Home Page | Cancer Pro... | Car B... | Untitled | Merc Data... | Merc Data... | Assessment | Password | Merc Data... | EDRHub | Upload file...

localhost:8790/notebooks/Merc%20Database%20XGBoost%20Project-Copy1.ipynb#

File Edit View Insert Cell Kernel Widgets Help Not Trusted Python 3

```

Summarize outcome (testing time) in training dataset

In [25]: # Draw a vertical boxplot grouped
         # by a categorical variable: X0
         sns.set_style("whitegrid")

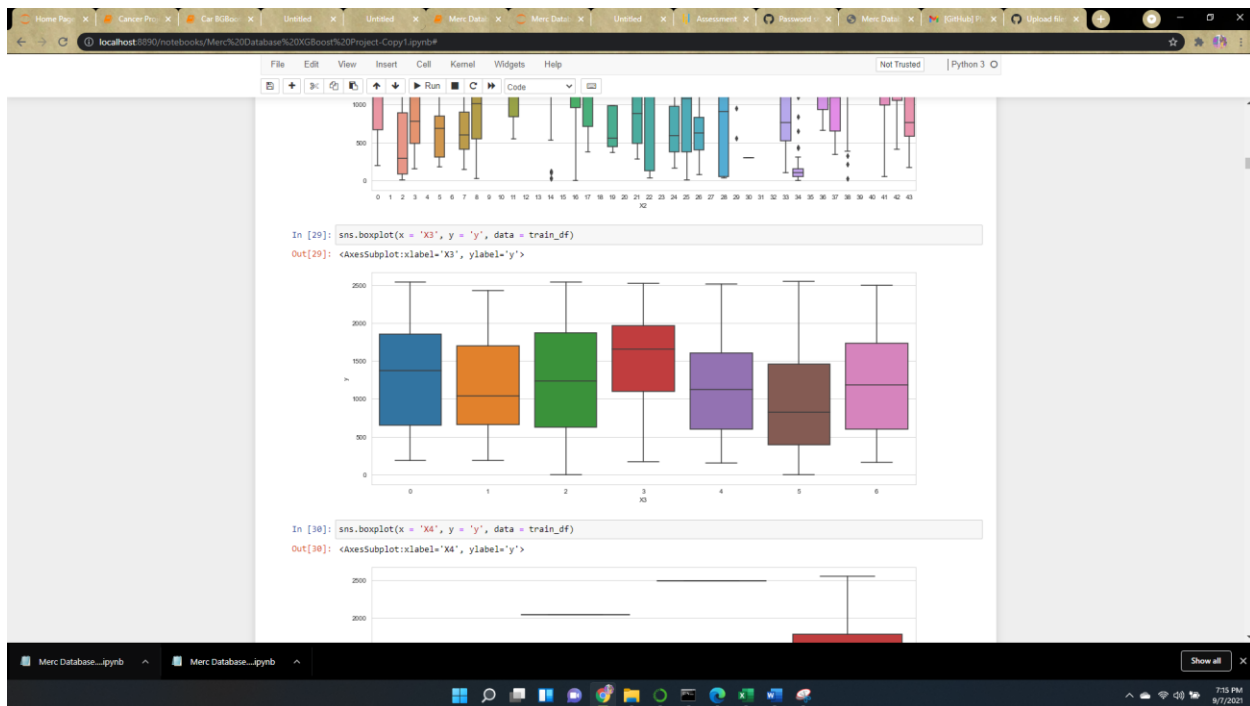
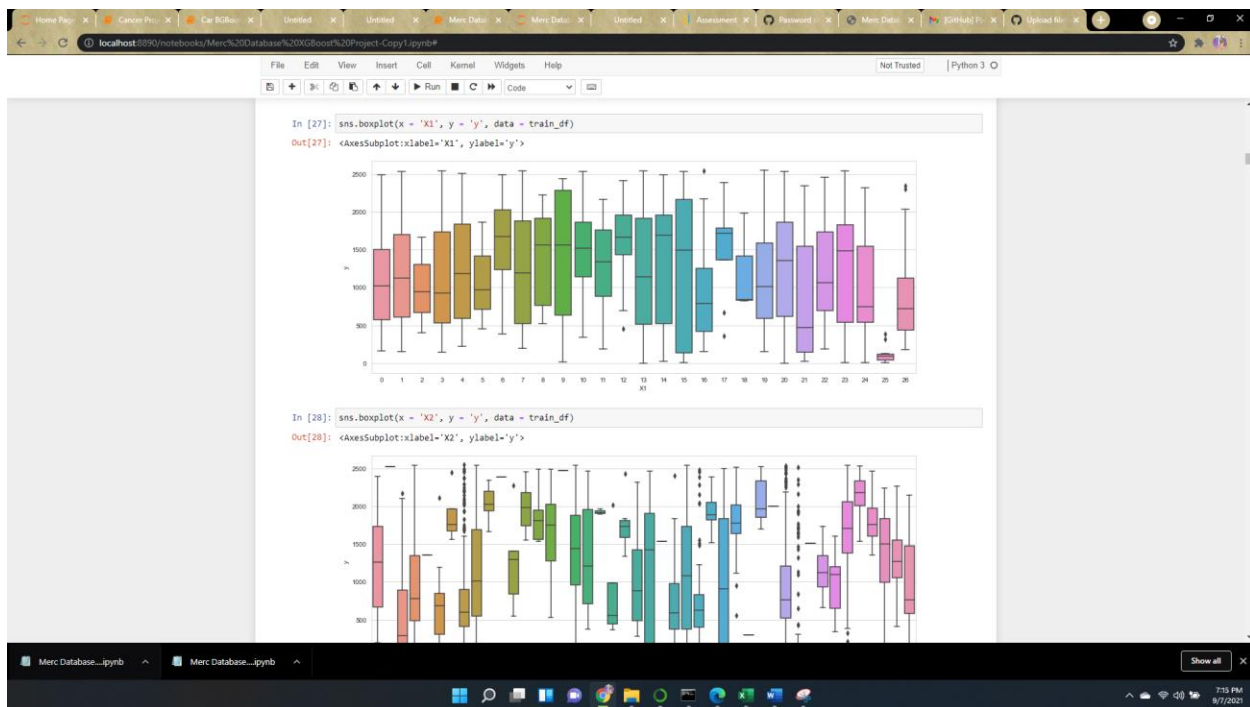
         object_columns = test_df.describe(include='object').columns
         print("\nobject columns:\n",object_columns)
         cols = len(object_columns)

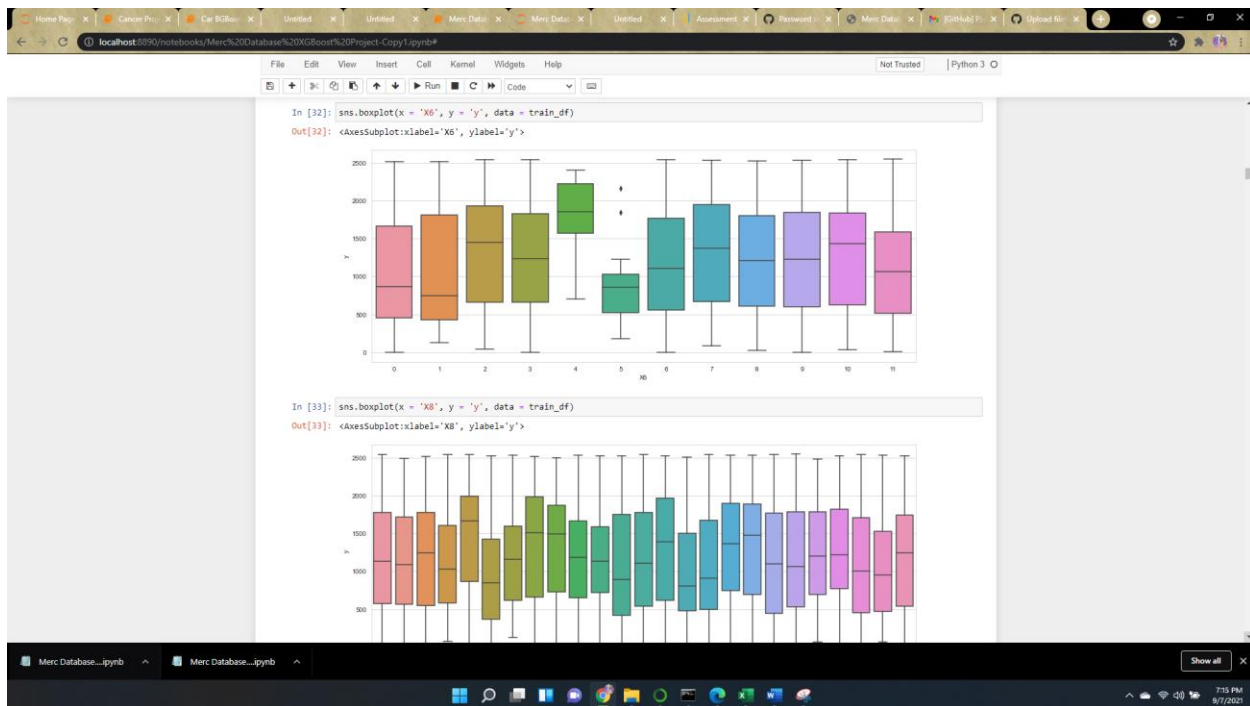
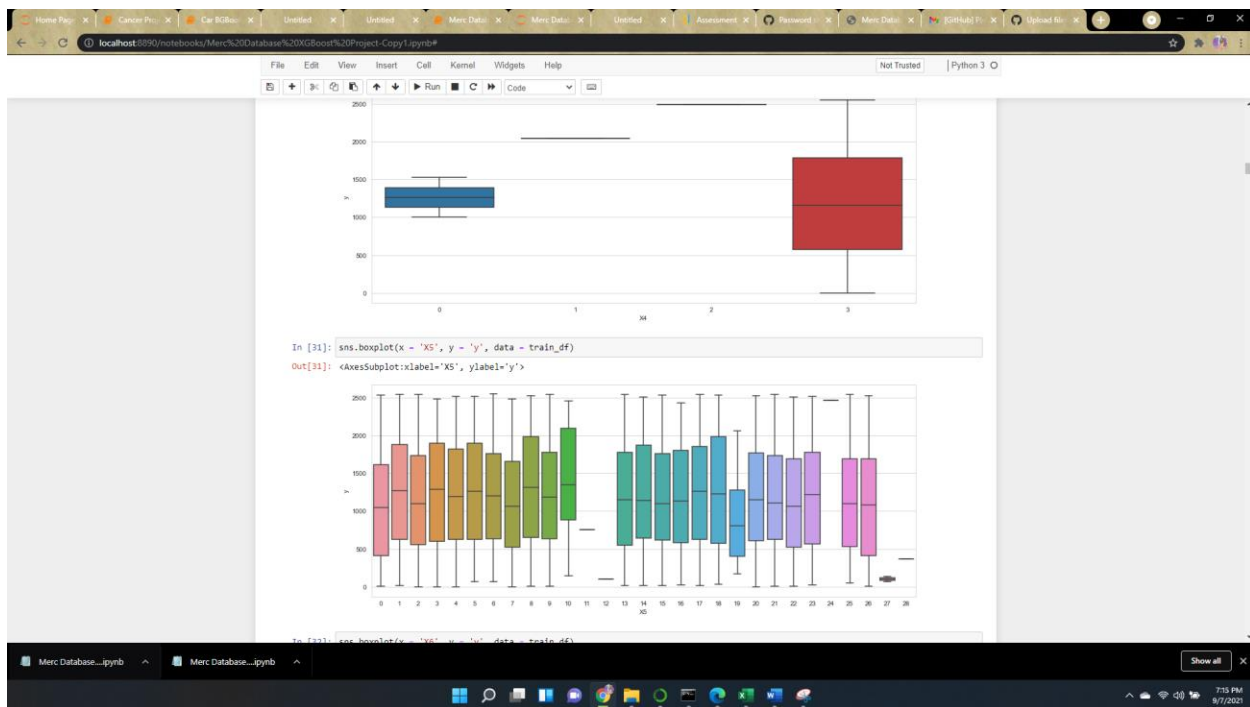
         object_columns:
         Index(['X0', 'X1', 'X2', 'X3', 'X4', 'X5', 'X6', 'X8'], dtype='object')

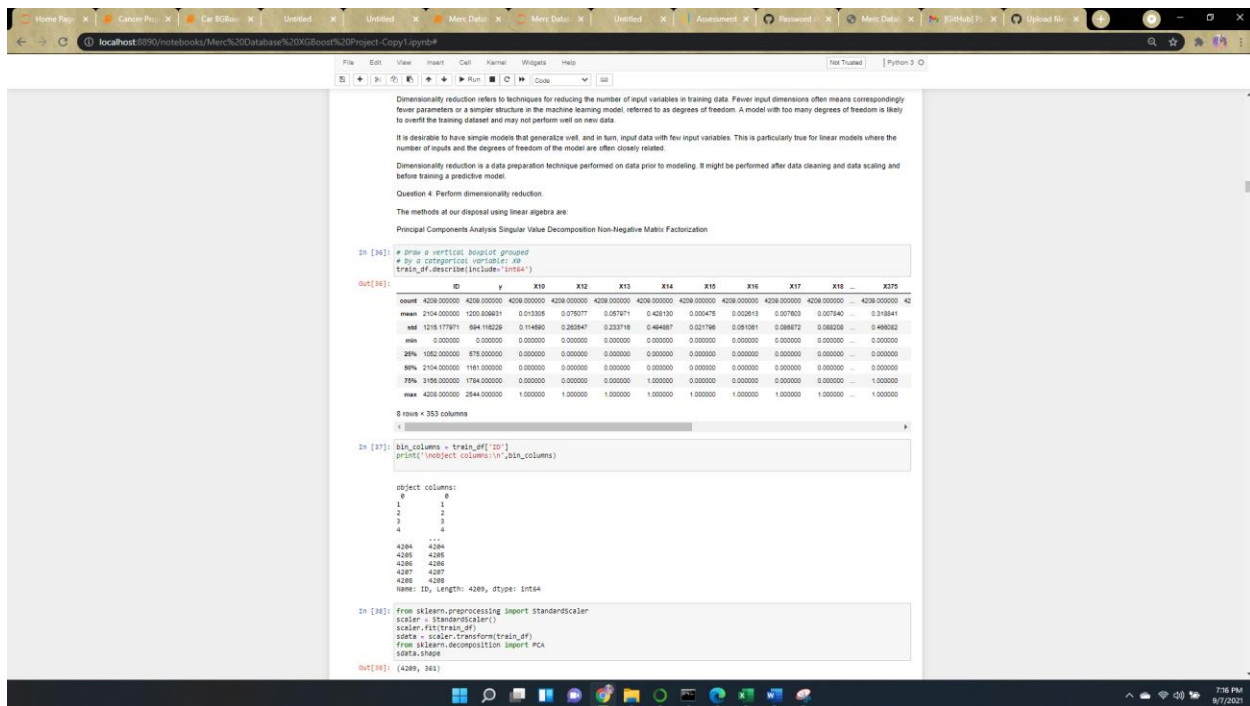
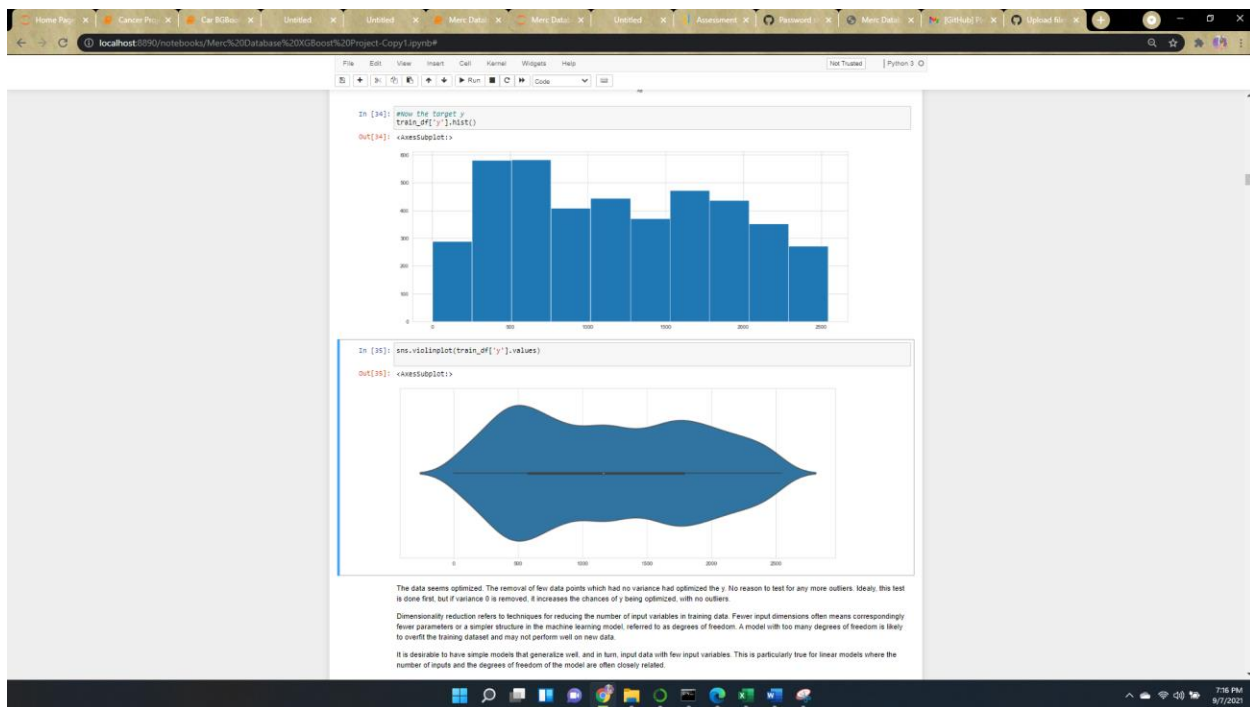
In [26]: sns.boxplot(x = 'X0', y = 'y', data = train_df)
Out[26]: <AxesSubplot: xlabel='X0', ylabel='y'>
```

```

In [27]: sns.boxplot(x = 'X1', y = 'y', data = train_df)
Out[27]: <AxesSubplot: xlabel='X1', ylabel='y'>
```





```
File Edit View Insert Cell Format Windows Help Not Trained Python 3
In [38]: from sklearn.preprocessing import StandardScaler
scaler = StandardScaler()
scaler.fit(train_df)
data = scaler.transform(train_df)
from sklearn.decomposition import PCA
pca = PCA(n_components=4)
pca.fit(data)
pca.components_
Out[38]: (4289, 262)

In [39]: # Lets take top 6 pca components
pca = PCA(n_components=6)
pca.fit(data)
X_pca = pca.transform(data)
X_pca.shape
# number of components
X_pca.shape[0]
X_pca.shape[1]
Out[39]: 6

In [40]: # get the index of the most important feature on each component i.e. largest absolute value
# using L27 decomposition here
most_important = [np.abs(pca.components_[i]).argmax() for i in range(n_pcs)]
initial_feature_names = data.columns
most_important
Out[40]: [325, 27, 188, 161, 389, 65]

In [41]: # using L27 decomposition here AGAIN
dic = {}
for i in range(n_pcs):
    dic['PC'+str(i+1)] = most_important[i]
Out[41]: {'PC0': 325, 'PC1': 27, 'PC2': 188, 'PC3': 161, 'PC4': 389, 'PC5': 65}

In [42]: pca.components_
Out[42]: array([[ 0.0020887,  0.0400885, -0.0335833, ...,  0.00244079,
                -0.0020312,  0.0003943],
               [ 0.0009532, -0.0774927, -0.0428254, ..., -0.00019554,
                -0.0003713,  0.0004221],
               [ 0.0004485,  0.0728521, -0.0777957, ..., -0.00023733,
                0.0024457,  0.0546821],
               [ 0.0020887,  0.0003334,  0.0703487, ..., -0.01077918,
                0.0008849,  0.0007179],
               [ 0.0022261,  0.0023862, -0.0212129, ...,  0.00476884,
                -0.0023327, -0.0007786],
               [-0.0000345,  0.0004625,  0.0226829, ...,  0.0004413,
                0.000352,  0.0003593]])

In [43]: explained_variance = pca.explained_variance_ratio_
explained_variance
# It is a measure of the variance of the data when projected onto that axis. The projection of each data point onto the
# principal axes are the "principal components" of the data. .4 is the var of PCA
# .129 is the var of PCA
Out[43]: array([0.00957844, 0.05767819, 0.04582234, 0.03464726, 0.03299807,
                0.03192876])

In [44]: #creating training and test data with only these columns
selected_columns = train_df[['ID', 'X325', 'X27', 'X188', 'X161', 'X389', 'X65', 'y']]
X_train = selected_columns.copy()
X_train.shape
Out[44]: (4289, 8)
```

```
File Edit View Insert Cell Format Windows Help Not Trained Python 3
Out[44]: array([0.00957844, 0.05767819, 0.04582234, 0.03464726, 0.03299807,
                0.03192876])

In [45]: #creating training and test data with only these columns
selected_columns = train_df[['ID', 'X325', 'X27', 'X188', 'X161', 'X389', 'X65', 'y']]
X_train = selected_columns.copy()
X_train.shape
X_train
Out[45]:
   ID  X325  X27  X188  X161  X389  X65  y
0    0    0    0    0    0    0    1  2485
1    1    0    1    0    0    0    1  305
2    2    0    1    0    0    0    1   80
3    3    0    1    0    0    0    0  133
4    4    0    1    0    0    0    0  106
...
4284 4204    0    1    0    0    0    1  1087
4285 4205    0    0    0    0    0    0  1786
4286 4206    0    1    0    0    0    1  1021
4287 4207    0    0    0    0    0    0   280
4288 4208    0    0    0    0    0    0  1021
4289 rows x 8 columns

In [46]: #now creating a df with only these 6 components
selected_columns = test_df[['ID', 'X325', 'X27', 'X188', 'X161', 'X389', 'X65']]
X_test = selected_columns.copy()
X_test.shape
X_test
Out[46]:
   ID  X325  X27  X188  X161  X389  X65
0    0    1    0    1    0    0    0
1    1    0    1    0    1    0    0
2    2    0    1    0    0    0    1
3    3    0    1    0    0    0    0
4    4    0    1    0    0    0    1
...
4284 8410    0    1    0    0    0    1
4285 8411    0    1    0    0    0    0
4286 8413    0    1    0    0    0    0
4287 8414    0    1    0    0    0    1
4288 8415    0    1    0    0    0    0
4289 rows x 7 columns

#now will perform XGBoost
Predict your test_df values using XGBoost Model Selection
Logistic Regression KNN SVM Random Forest
```

```
4209 rows x 7 columns

#now will perform XGBoost
Predict your test_of values using XGBoost Model Selection
Logistic Regression KNN SVM Random Forest

In [47]: #now splitting the data into train & test. Before that, identifying all input parameters as X, and output parameter as y
         y_train=y_off[y]
         y_train

Out[47]: 0    2466
         1     366
         2      69
         3     133
         4     186
         ...
         4204  1657
         4205  1756
         4206  1881
         4207   288
         4208  1921
         Name: y, length: 4209, dtype: int64

In [48]: y_train.shape
Out[48]: (4209,)

In [49]: from sklearn.model_selection import learning_curve
         from sklearn.linear_model import LogisticRegression
         from sklearn.neighbors import KNeighborsClassifier
         from sklearn.model_selection import cross_val_score
         from sklearn.metrics import classification_report, confusion_matrix
         import xgboost as xgb
         from sklearn.metrics import r2_score
         from sklearn.model_selection import train_test_split

In [50]: X_train.shape
Out[50]: (4209, 6)

In [51]: X_test.shape
Out[51]: (4209, 7)

In [52]: y_train.shape
Out[52]: (4209,)

In [53]: #we do not have X & y. Creating X

In [54]: X = X_train
         #X = pd.concat([X_train,X_test])
         print(X)

0    ID  X325  X27  X100  X361  X309  X365  y
1    1    0    0    0    0    0    0    1  2466
2    2    0    1    0    0    0    0    1    366
3    3    0    1    0    0    0    0    1    69
4    4    0    1    0    0    0    0    1    133
...
4204  4204  0    1    0    0    0    0    1  1657
4205  4205  0    0    0    0    0    0    0    1  1756
4206  4206  0    1    0    0    0    0    1  1881
4207  4207  0    0    0    0    0    0    0    1    288
4208  4208  0    0    0    0    0    0    0    1  1921

[4209 rows x 8 columns]
```

```
print(X)

0    ID  X325  X27  X100  X361  X309  X365  y
1    1    0    0    0    0    0    0    1  2466
2    2    0    1    0    0    0    0    1    366
3    3    0    1    0    0    0    0    1    69
4    4    0    1    0    0    0    0    1    133
...
4204  4204  0    1    0    0    0    0    1  1657
4205  4205  0    0    0    0    0    0    0    1  1756
4206  4206  0    1    0    0    0    0    1  1881
4207  4207  0    0    0    0    0    0    0    1    288
4208  4208  0    0    0    0    0    0    0    1  1921

[4209 rows x 8 columns]

In [55]: X.shape
Out[55]: (4209, 8)

In [56]: #y = pd.concat([y_train,y_test])
         y = y_train
         y.shape

Out[56]: (4209,)

In [57]: X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_state=72)

In [58]: # Logistic Regression
         logreg=LogisticRegression(solver='liblinear',multi_class='ovr')
         logreg.fit(X_train,y_train)
         y_pred=logreg.predict(X_test)
         #Accuracy score
         #print(metrics.accuracy_score(y_pred,y_train))
         accuracy = (logreg.score(X_test,y_train))
         print(accuracy)

0.61731160896138343

In [61]: # Logistic Regression
         logreg=LogisticRegression(solver='lbfgs',multi_class='auto')
         logreg.fit(X_train,y_train)
         y_pred=logreg.predict(X_test)
         #Accuracy score
         #print(metrics.accuracy_score(y_pred,y_train))
         accuracy = (logreg.score(X_train,y_train))
         print(accuracy)

0.6074677285526816

In [62]: #Now "Support Vector Classifier"
         from sklearn.svm import SVC
         svm = SVC(kernel='linear')

         # fitting a samples and y classes
         svm.fit(X_train,y_train)
         y_pred = svm.predict(X_test)

         from sklearn import metrics
         accuracy = metrics.accuracy_score(y_test, y_pred)
         print(accuracy)

0.6958936421129319
```



```
File Edit View Insert Cell Format Windows Help Not Trained Python 3

XGBRegressor(base_score=0.5, booster='dart', colsample_bylevel=1,
              colsample_bytree=1, eval_metric='error',
              gamma=0, gpu_id=-1, importance_type='test',
              interaction_constraints='', learning_rate=0.30000001,
              max_delta_step=0, max_depth=6, min_child_weight=1, missing=nan,
              monotone_constraints=(), n_estimators=10, n_jobs=4, num_class=1,
              num_parallel_tree=1, random_state=11, reg_alpha=0, reg_lambda=1,
              scale_pos_weight=1, seed=11, subsample=1, tree_method='exact',
              validate_parameters=1, verbosity=None)

In [87]: # RMSE Computation
rmse = np.sqrt(mean_squared_error(y_test, pred))
print("RMSE : % f" % (rmse))
RMSE : 41.334797

In [88]: expected_y = y_test
predicted_y = model.predict(x_test)
print(metrics.r2_score(y_test, predicted_y))
0.99645593972794

In [89]: predicted_y

Out[89]: array([1310.2015, 2051.473 , 856.6799 , ..., 793.9099 , 1325.2378 ,
                211.49994], dtype=float32)

In [90]: # Here, we are using XGBRegressor as a machine learning model to fit the data.
model = xgb.XGBRegressor(booster='gblinear', objective='reg:squarederror', num_class = 1, eval_metric = 'error', n_estimators =
model.fit(x_train, y_train)
print() print(model)

# Predict the model
pred = model.predict(x_test)

XGBRegressor(base_score=0.5, booster='gblinear', colsample_bylevel=None,
              colsample_bytree=None, colsample_bynode=None, eval_metric='error',
              gamma=None, gpu_id=-1, importance_type='gain',
              interaction_constraints=None, learning_rate=0.3,
              max_delta_step=None, max_depth=None, min_child_weight=None,
              missing=None, monotone_constraints=None, n_estimators=10, n_jobs=4,
              num_class=1, num_parallel_tree=None, random_state=11, reg_alpha=0,
              reg_lambda=1, scale_pos_weight=1, seed=11, subsample=None,
              tree_method=None, validate_parameters=1, verbosity=None)

In [91]: # RMSE Computation
rmse = np.sqrt(mean_squared_error(y_test, pred))
print("RMSE : % f" % (rmse))
RMSE : 196.394913

In [92]: expected_y = y_test
predicted_y = model.predict(x_test)
print(metrics.r2_score(y_test, predicted_y))
0.919993195162286

In [93]: # As we see gbLinear is not the right model.
# Here, we are using XGBRegressor as a machine learning model to fit the data.
model = xgb.XGBRegressor(booster='gbtree', objective='reg:squarederror', num_class = 1, eval_metric = 'error', n_estimators =
model.fit(x_train, y_train)
print() print(model)

# Predict the model
```

```
File Edit View Insert Cell Format Windows Help Not Trained Python 3

rmse_cv_scores = cross_validation.cross_val_score(x,
                                                  y,
                                                  XGBRegressor(booster='gbtree', objective='reg:squarederror', num_class = 1, eval_metric = 'error', n_estimators =
                                                  reg_alpha=0, reg_lambda=1, scale_pos_weight=1, subsample=1,
                                                  tree_method='exact', validate_parameters=1, verbosity=None))

In [97]: plt.figure(figsize=(10,10))
sns.regplot(expected_y, predicted_y, fit_reg=True, scatter_kws={'s': 100})
Out[97]: <axes.SubplotZmpl>

In [98]: # We'll check the training accuracy with cross-validation and k-fold methods.
# Applying k-fold cross validation
from sklearn.model_selection import cross_val_score, Kfold

kfold = Kfold(n_splits=10, shuffle=True)
kf_cv_scores = cross_val_score(model, x_train, y_train, cv=kfold)
print("k-fold cv average score: %.2f" % kf_cv_scores.mean())

k-fold cv average score: 1.00

In [99]: # Now we have predicted the output by passing x_test and also stored real target in expected_y.
expected_y = y_test
predicted_y = model.predict(x_test)
print(metrics.r2_score(expected_y, predicted_y))
0.999987747167113

# This is the best model so far RMSE Predicting y with XGBoost

In [107]: # Displaying predicted values
```

Embedded File:

Python File:



Merc Database XGBoost Project.ipynb

Code hosted on GitHub :

https://github.com/ks-alokranjan/Simplilearn_ML