```
1    import pandas as pd
2    import numpy as np
3    import time
```

```
1    from google.colab import drive
2    drive.mount('/content/drive/')
```

Drive already mounted at /content/drive/; to attempt to forcibly remount, call drive.mou

## Injesting The Data

```
1    train_data = pd.read_csv('/content/drive/MyDrive/Colab Notebooks/data/benz/train.csv')
2    test_data = pd.read_csv('/content/drive/MyDrive/Colab Notebooks/data/benz/test.csv')
```

## Data Analysis

```
1    train_data.shape
```
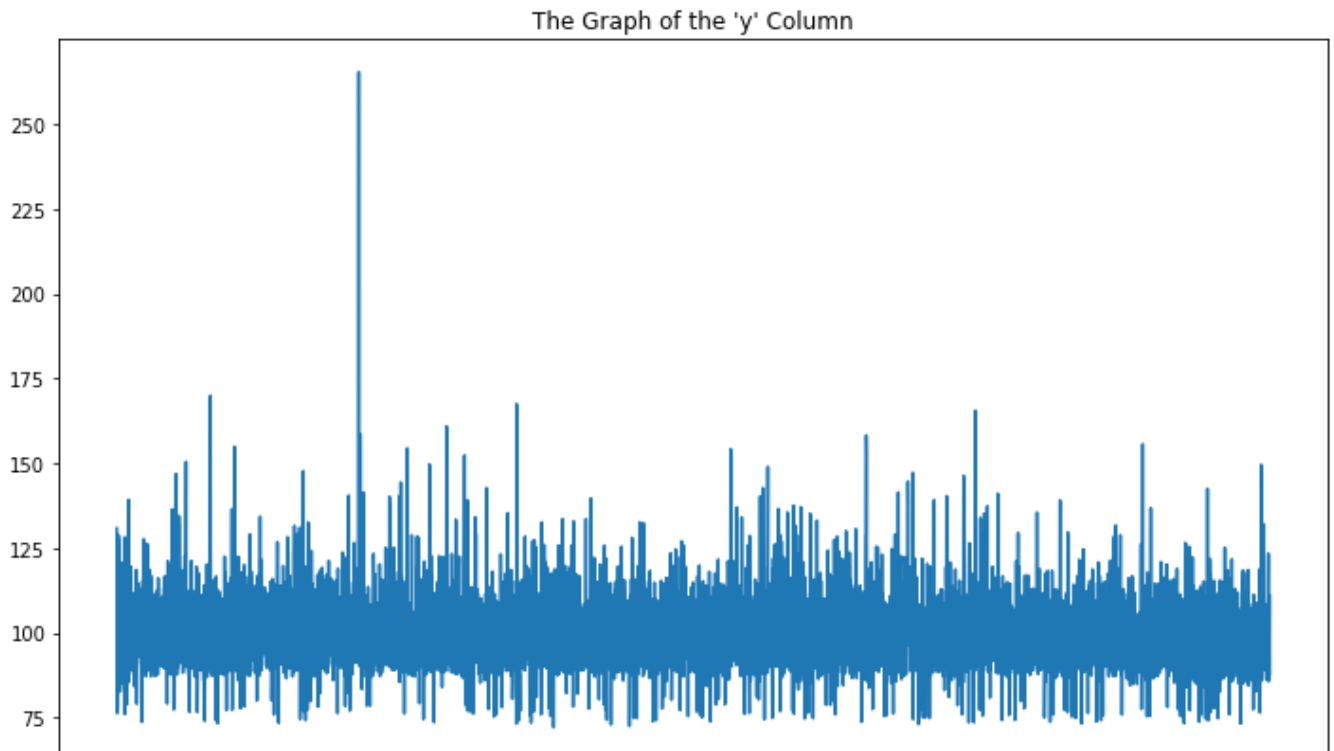
(4209, 378)

```
1    test_data.shape
```

(4209, 377)

```
1    train_data.head()
```

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

5 rows × 378 columns

```
1 train_data['y'].plot(figsize=(12,7), title="The Graph of the 'y' Column");
```

The Graph of the 'y' Column



```
1 test_data.head()
```

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

5 rows × 377 columns
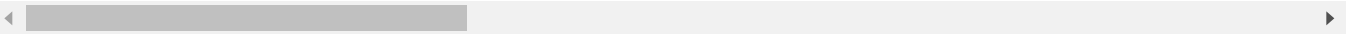
```
1 train_data.dtypes
```

```
ID          int64
y         float64
X0         object
X1         object
X2         object
           ...
X380        int64
X382        int64
X383        int64
X384        int64
X385        int64
Length: 378, dtype: object
```

```
1 train_data.describe()
```

| | ID | y | X10 | X11 | X12 | X13 | X1 |
|---|---|---|---|---|---|---|---|
| **count** | 4209.000000 | 4209.000000 | 4209.000000 | 4209.0 | 4209.000000 | 4209.000000 | 4209.00000 |
| **mean** | 4205.960798 | 100.669318 | 0.013305 | 0.0 | 0.075077 | 0.057971 | 0.42813 |
| **std** | 2437.608688 | 12.679381 | 0.114590 | 0.0 | 0.263547 | 0.233716 | 0.49486 |
| **min** | 0.000000 | 72.110000 | 0.000000 | 0.0 | 0.000000 | 0.000000 | 0.00000 |
| **25%** | 2095.000000 | 90.820000 | 0.000000 | 0.0 | 0.000000 | 0.000000 | 0.00000 |
| **50%** | 4220.000000 | 99.150000 | 0.000000 | 0.0 | 0.000000 | 0.000000 | 0.00000 |
| **75%** | 6314.000000 | 109.010000 | 0.000000 | 0.0 | 0.000000 | 0.000000 | 1.00000 |
| **max** | 8417.000000 | 265.320000 | 1.000000 | 0.0 | 1.000000 | 1.000000 | 1.00000 |

8 rows × 370 columns

```
1 test_data.dtypes
```

```
ID       int64
X0       object
X1       object
X2       object
X3       object
         ...
X380     int64
X382     int64
X383     int64
X384     int64
X385     int64
Length: 377, dtype: object
```

```
1   test_data.describe()
```

|       | ID          | X10         | X11         | X12         | X13         | X14         |        |
|-------|-------------|-------------|-------------|-------------|-------------|-------------|--------|
| count | 4209.000000 | 4209.000000 | 4209.000000 | 4209.000000 | 4209.000000 | 4209.000000 | 4209.0 |

## Checking Null Values

| std | 2425.078920 | 0.136565 | 0.015414 | 0.202594 | 0.239468 | 0.494852 | 0. |

```
1  train_data.isnull().sum()
```

```
ID      0
y       0
X0      0
X1      0
X2      0
       ..
X380    0
X382    0
X383    0
X384    0
X385    0
Length: 378, dtype: int64
```

```
1 test_data.isnull().sum()
```

```
ID      0
X0      0
X1      0
X2      0
X3      0
       ..
X380    0
X382    0
X383    0
X384    0
X385    0
Length: 377, dtype: int64
```

```
1 train_data.describe(include='object')
```

|        | 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  |

## Separating categorical and numeric feature

```
1 dictionary={}
2 dictionary['num'] = train_data.dtypes[train_data.dtypes=='int64'].index
3 dictionary['cat'] = train_data.dtypes[train_data.dtypes=='object'].index
4 dictionary
```

```
{'cat': Index(['X0', 'X1', 'X2', 'X3', 'X4', 'X5', 'X6', 'X8'], dtype='object'),
 'num': Index(['ID', 'X10', 'X11', 'X12', 'X13', 'X14', 'X15', 'X16', 'X17', 'X18',
        ...
        'X375', 'X376', 'X377', 'X378', 'X379', 'X380', 'X382', 'X383', 'X384',
        'X385'],
      dtype='object', length=369)}
```

```
1 train_data['X376'].value_counts()
```

```
0    3968
1     241
Name: X376, dtype: int64
```

Finding and removing features with zero variance

```
1 da=[]
2 count=0
3 for i in dictionary['num']:
4   if(np.var(train_data[i])==0):
5     da.append(i)
6 print(da)
```

```
['X11', 'X93', 'X107', 'X233', 'X235', 'X268', 'X289', 'X290', 'X293', 'X297', 'X330',
```

```
1 z=['ID','y']
2 da.extend(z)
3 print(da)
```

```
['X11', 'X93', 'X107', 'X233', 'X235', 'X268', 'X289', 'X290', 'X293', 'X297', 'X330',
```

```
1 y_train=train_data['y'].values
2 ID_train=train_data['ID'].values
3 ID_test=test_data['ID'].values
```

```
1 da_test=['X11', 'X93', 'X107', 'X233', 'X235', 'X268', 'X289', 'X290', 'X293', 'X297', '>
```

```
1 X_train=train_data.drop(da,axis=1)
2 X_test=test_data.drop(da_test,axis=1)
```

```
1 X_train.shape
```

```
    (4209, 364)
```

```
1 X_test.shape
```

```
    (4209, 364)
```

Applying Label encoding on the categorical variables

```
1 from sklearn.preprocessing import LabelEncoder
```

```
1 #Train dataset
2 le0=LabelEncoder()
3 le1=LabelEncoder()
4 le2=LabelEncoder()
5 le3=LabelEncoder()
6 le4=LabelEncoder()
7 le5=LabelEncoder()
8 le6=LabelEncoder()
9 le8=LabelEncoder()
```

```
1 le0.fit(X_train['X0'])
2 le1.fit(X_train['X1'])
3 le2.fit(X_train['X2'])
4 le3.fit(X_train['X3'])
5 le4.fit(X_train['X4'])
6 le5.fit(X_train['X5'])
7 le6.fit(X_train['X6'])
8 le8.fit(X_train['X8'])
```

```
    LabelEncoder()
```

```
1 #Test Dataset
2 let0=LabelEncoder()
3 let1=LabelEncoder()
4 let2=LabelEncoder()
5 let3=LabelEncoder()
6 let4=LabelEncoder()
7 let5=LabelEncoder()
8 let6=LabelEncoder()
9 let8=LabelEncoder()
```

```
1 #Test Dataset
2 let0.fit(X_test['X0'])
3 let1.fit(X_test['X1'])
4 let2.fit(X_test['X2'])
5 let3.fit(X_test['X3'])
```

```
6 let4.fit(X_test['X4'])
7 let5.fit(X_test['X5'])
8 let6.fit(X_test['X6'])
9 let8.fit(X_test['X8'])
```

```
LabelEncoder()
```

Transforming and replacing the categorical variables into 0s and 1s

```
1 #Train dataset
2 X_train['X0'] = le0.transform(X_train['X0'])
3 X_train['X1'] = le1.transform(X_train['X1'])
4 X_train['X2'] = le2.transform(X_train['X2'])
5 X_train['X3'] = le3.transform(X_train['X3'])
6 X_train['X4'] = le4.transform(X_train['X4'])
7 X_train['X5'] = le5.transform(X_train['X5'])
8 X_train['X6'] = le6.transform(X_train['X6'])
9 X_train['X8'] = le8.transform(X_train['X8'])
```

```
1 #Test dataset
2 X_test['X0'] = let0.transform(X_test['X0'])
3 X_test['X1'] = let1.transform(X_test['X1'])
4 X_test['X2'] = let2.transform(X_test['X2'])
5 X_test['X3'] = let3.transform(X_test['X3'])
6 X_test['X4'] = let4.transform(X_test['X4'])
7 X_test['X5'] = let5.transform(X_test['X5'])
8 X_test['X6'] = let6.transform(X_test['X6'])
9 X_test['X8'] = let8.transform(X_test['X8'])
```

```
1 X_train.head()
```

| | X0 | X1 | X2 | X3 | X4 | X5 | X6 | X8 | X10 | X12 | ... | X375 | X376 | X377 | X378 | X379 | X380 | X38 |
|---|----|----|----|----|----|----|----|----|-----|-----|-----|------|------|------|------|------|------|-----|
| 0 | 32 | 23 | 17 | 0 | 3 | 24 | 9 | 14 | 0 | 0 | ... | 0 | 0 | 1 | 0 | 0 | 0 | |
| 1 | 32 | 21 | 19 | 4 | 3 | 28 | 11 | 14 | 0 | 0 | ... | 1 | 0 | 0 | 0 | 0 | 0 | |
| 2 | 20 | 24 | 34 | 2 | 3 | 27 | 9 | 23 | 0 | 0 | ... | 0 | 0 | 0 | 0 | 0 | 0 | |
| 3 | 20 | 21 | 34 | 5 | 3 | 27 | 11 | 4 | 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 | |

5 rows × 364 columns

```
1 X_test.head()
```

|   | X0 | X1 | X2 | X3 | X4 | X5 | X6 | X8 | X10 | X12 | ... | X375 | X376 | X377 | X378 | X379 | X380 | X38 |
|---|----|----|----|----|----|----|----|----|-----|-----|-----|------|------|------|------|------|------|-----|
| 0 | 21 | 23 | 34 | 5 | 3 | 26 | 0 | 22 | 0 | 0 | ... | 0 | 0 | 0 | 1 | 0 | 0 | |
| 1 | 42 | 3 | 8 | 0 | 3 | 9 | 6 | 24 | 0 | 0 | ... | 0 | 0 | 1 | 0 | 0 | 0 | |
| 2 | 21 | 23 | 17 | 5 | 3 | 0 | 9 | 9 | 0 | 0 | ... | 0 | 0 | 0 | 1 | 0 | 0 | |
| 3 | 21 | 13 | 34 | 5 | 3 | 31 | 11 | 13 | 0 | 0 | ... | 0 | 0 | 0 | 1 | 0 | 0 | |
| 4 | 45 | 20 | 17 | 2 | 3 | 30 | 8 | 12 | 0 | 0 | ... | 1 | 0 | 0 | 0 | 0 | 0 | |

## Perform dimensionality reduction on Train dataset

```
1 from sklearn.decomposition import PCA
2 pca = PCA(n_components=0.7)
3 pca.fit(X_train)
```

```
PCA(n_components=0.7)
```

```
1 pca.explained_variance_ratio_
```

```
array([0.38334782, 0.21388033, 0.13261866])
```

```
1 X_train_transformed = pca.transform(X_train)
2 X_test_transformed = pca.transform(X_test)
```

```
1 X_train_transformed.shape, X_test_transformed.shape
```

```
((4209, 3), (4209, 3))
```

```
1 pca.inverse_transform(X_train_transformed).shape
```

```
(4209, 364)
```

## Split the datasets

```
1 from sklearn.model_selection import train_test_split
2 X_train, X_test, y_train, y_test = train_test_split(X_train_transformed,y_train , random_
```

```
1 print(X_train.shape)
2 print(X_test.shape)
3 print(y_train.shape)
4 print(y_test.shape)
```

```
(3156, 3)
(1053, 3)
```

```
(3156,)
(1053,)
```

## Evaluate with ensemble learning algorithms

```
1 import xgboost
2 from sklearn.ensemble import GradientBoostingRegressor
3 from sklearn.metrics import mean_squared_error
```

```
1 #Using XGBoost - The most powerful ensemble learning algorithm
2 start = time.time()
3 xgbr = xgboost.XGBRegressor()
4 xgbr.fit(X_train, y_train)
5
6 y_pred= xgbr.predict(X_test)
7 score = mean_squared_error(y_pred, y_test)
8 print('Score : ', score)
9 end = time.time()
10
11 print('\nRunTime : ', end - start)
```

```
[02:53:42] WARNING: /workspace/src/objective/regression_obj.cu:152: reg:linear is now de
Score :   110.00614368130348

RunTime :   0.1572103500366211
```

```
1 #Using GradientBoost
2 start = time.time()
3 gbr = GradientBoostingRegressor()
4 gbr.fit(X_train, y_train)
5
6 y_pred= gbr.predict(X_test)
7 score = mean_squared_error(y_pred, y_test)
8 print('Score : ', score)
9 end = time.time()
10
11 print('\nRunTime : ', end - start)
```

```
Score :   111.54442824296098

RunTime :   0.3748013973236084
```

As seen above, XGBoost performs more efficiently than GradientBoost

✓  0s    completed at 10:53 PM    ● ✕