# Mercedes-Benz

February 7, 2023

```python
[31]: import pandas as pd
      import numpy as np
      import matplotlib.pyplot as plt
      %matplotlib inline
```

```python
[32]: #importing test and train dataset
      df_test=pd.read_csv('test[1].csv')
      df_train=pd.read_csv('train[1].csv')
```

```python
[33]: #droping ID columns
      df_train.drop('ID',inplace=True,axis=1)
      df_test.drop('ID',inplace=True,axis=1)
```

```python
[34]: print(df_train.columns)
      print(df_test.columns)
```

```
Index(['y', 'X0', 'X1', 'X2', 'X3', 'X4', 'X5', 'X6', 'X8', 'X10',
       ...
       'X375', 'X376', 'X377', 'X378', 'X379', 'X380', 'X382', 'X383', 'X384',
       'X385'],
      dtype='object', length=377)
Index(['X0', 'X1', 'X2', 'X3', 'X4', 'X5', 'X6', 'X8', 'X10', 'X11',
       ...
       'X375', 'X376', 'X377', 'X378', 'X379', 'X380', 'X382', 'X383', 'X384',
       'X385'],
      dtype='object', length=376)
```

```python
[35]: print(df_train.shape)
      print(df_test.shape)
```

```
(4209, 377)
(4209, 376)
```

```python
[36]: #Task1:If for any columns,the varience is equal to zero,then you need to remove␣
      ↪the variables
      zero_var_col=df_train.var()[df_train.var()==0].index.values
```

```python
[37]: zero_var_col
```

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

```
[38]: #Droping the var==0 columns
      df_train.drop(['X11', 'X93', 'X107', 'X233', 'X235', 'X268', 'X289', 'X290',
       ↪'X293', 'X297', 'X330', 'X347'],axis=1,inplace=True)
      df_test.drop(['X11', 'X93', 'X107', 'X233', 'X235', 'X268', 'X289', 'X290',
       ↪'X293', 'X297', 'X330', 'X347'],axis=1,inplace=True)
```

```
[39]: df_train.columns
```

```
[39]: Index(['y', 'X0', 'X1', 'X2', 'X3', 'X4', 'X5', 'X6', 'X8', 'X10',
             ...
             'X375', 'X376', 'X377', 'X378', 'X379', 'X380', 'X382', 'X383', 'X384',
             'X385'],
            dtype='object', length=365)
```

```
[40]: #Task2:check for null and unique values in test and train set
```

```
[41]: #checking Null values
      df_train.isna().sum()
```

```
[41]: y        0
      X0       0
      X1       0
      X2       0
      X3       0
               ..
      X380     0
      X382     0
      X383     0
      X384     0
      X385     0
      Length: 365, dtype: int64
```

```
[42]: df_test.isna().sum()
```

```
[42]: X0       0
      X1       0
      X2       0
      X3       0
      X4       0
               ..
      X380     0
      X382     0
      X383     0
      X384     0
```

```
X385        0
Length: 364, dtype: int64
```

[43]:
```python
#Checking for unique columns
for i in df_train.columns:
    print(df_train[i].unique())
```

```
[130.81  88.53  76.26 …  85.71 108.77  87.48]
['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']
['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']
['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']
['a' 'e' 'c' 'f' 'd' 'b' 'g']
['d' 'b' 'c' 'a']
['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']
['j' 'l' 'd' 'h' 'i' 'a' 'g' 'c' 'k' 'e' 'f' 'b']
['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']
[0 1]
[0 1]
[1 0]
[0 1]
[0 1]
[0 1]
[0 1]
[1 0]
[0 1]
[0 1]
[1 0]
[0 1]
[0 1]
[0 1]
[0 1]
[0 1]
[0 1]
[0 1]
[0 1]
[1 0]
[0 1]
[0 1]
[0 1]
[1 0]
```

```
[0 1]
[1 0]
[0 1]
[0 1]
[0 1]
[0 1]
[0 1]
[0 1]
[0 1]
[0 1]
[1 0]
[0 1]
[0 1]
[0 1]
[0 1]
[0 1]
[0 1]
[0 1]
[0 1]
[0 1]
[0 1]
[0 1]
[1 0]
[0 1]
[0 1]
[0 1]
[0 1]
[0 1]
[0 1]
[0 1]
[0 1]
[0 1]
[1 0]
[0 1]
[1 0]
[0 1]
[0 1]
[1 0]
[0 1]
[0 1]
[0 1]
[0 1]
[0 1]
[0 1]
[0 1]
[0 1]
[0 1]
```

```
[1 0]
[0 1]
[0 1]
[0 1]
[0 1]
[0 1]
[0 1]
[0 1]
[0 1]
[0 1]
[0 1]
[0 1]
[0 1]
[0 1]
[0 1]
[0 1]
[0 1]
[0 1]
[0 1]
[0 1]
[0 1]
[0 1]
[0 1]
[1 0]
[0 1]
[0 1]
[1 0]
[0 1]
[1 0]
[0 1]
[1 0]
[1 0]
[1 0]
[0 1]
[0 1]
[0 1]
[0 1]
[0 1]
[0 1]
[1 0]
[0 1]
[0 1]
[1 0]
[0 1]
[0 1]
[0 1]
[0 1]
```

```
[1 0]
[1 0]
[0 1]
[0 1]
[0 1]
[0 1]
[1 0]
[0 1]
[1 0]
[0 1]
[0 1]
[0 1]
[0 1]
[1 0]
[0 1]
[0 1]
[0 1]
[0 1]
[0 1]
[1 0]
[0 1]
[0 1]
[0 1]
[0 1]
[0 1]
[0 1]
[0 1]
[0 1]
[0 1]
[0 1]
[0 1]
[0 1]
[1 0]
[0 1]
[0 1]
[0 1]
[0 1]
[0 1]
[0 1]
[0 1]
[0 1]
[1 0]
[0 1]
[0 1]
[0 1]
[0 1]
[1 0]
```

```
[0 1]
[0 1]
[1 0]
[1 0]
[0 1]
[0 1]
[0 1]
[1 0]
[0 1]
[0 1]
[0 1]
[0 1]
[0 1]
[0 1]
[0 1]
[0 1]
[0 1]
[1 0]
[0 1]
[0 1]
[0 1]
[0 1]
[1 0]
[0 1]
[0 1]
[0 1]
[0 1]
[0 1]
[0 1]
[0 1]
[0 1]
[0 1]
[0 1]
[1 0]
[0 1]
[0 1]
[0 1]
[0 1]
[0 1]
[0 1]
[0 1]
[0 1]
[0 1]
[0 1]
[0 1]
[0 1]
[1 0]
[0 1]
```

```
[1 0]
[0 1]
[0 1]
[0 1]
[0 1]
[0 1]
[0 1]
[0 1]
[0 1]
[0 1]
[0 1]
[0 1]
[0 1]
[0 1]
[0 1]
[0 1]
[0 1]
[0 1]
[0 1]
[0 1]
[0 1]
[0 1]
[0 1]
[0 1]
[1 0]
[1 0]
[0 1]
[0 1]
[1 0]
[0 1]
[0 1]
[0 1]
[0 1]
[0 1]
[1 0]
[0 1]
[1 0]
[0 1]
[0 1]
[0 1]
[0 1]
[0 1]
[0 1]
[0 1]
[0 1]
[0 1]
[1 0]
```

```
[0 1]
[0 1]
[0 1]
[0 1]
[0 1]
[0 1]
[0 1]
[0 1]
[0 1]
[0 1]
[0 1]
[0 1]
[0 1]
[0 1]
[0 1]
[1 0]
[0 1]
[0 1]
[0 1]
[0 1]
[0 1]
[0 1]
[0 1]
[0 1]
[0 1]
[1 0]
[0 1]
[0 1]
[0 1]
[0 1]
[0 1]
[0 1]
[0 1]
[1 0]
[0 1]
[0 1]
[1 0]
[0 1]
[1 0]
[0 1]
[0 1]
[0 1]
[1 0]
[0 1]
[0 1]
[0 1]
[0 1]
[0 1]
```

```
[0 1]
[0 1]
[0 1]
[0 1]
[0 1]
[0 1]
[0 1]
[0 1]
[0 1]
[0 1]
[0 1]
[0 1]
[0 1]
[1 0]
[0 1]
[0 1]
[0 1]
[0 1]
[0 1]
[0 1]
[1 0]
[0 1]
[0 1]
[0 1]
[0 1]
[0 1]
[0 1]
[0 1]
[0 1]
[0 1]
[0 1]
[0 1]
[0 1]
[0 1]
[0 1]
[0 1]
[1 0]
[0 1]
[0 1]
[0 1]
[0 1]
[0 1]
[0 1]
[0 1]
```

[44]: #Task3:Applying labelencoder

```python
[45]: from sklearn.preprocessing import LabelEncoder
      label_col=df_train.describe(include=['object']).columns.values
      label_col
```

```
[45]: array(['X0', 'X1', 'X2', 'X3', 'X4', 'X5', 'X6', 'X8'], dtype=object)
```

```python
[46]: le=LabelEncoder()
      for col in label_col:
          le.fit(df_train[col].append(df_test[col]).values)
          df_train[col]=le.transform(df_train[col])
          df_test[col]=le.transform(df_test[col])
```

```python
[47]: df_train.columns
```

```
[47]: Index(['y', 'X0', 'X1', 'X2', 'X3', 'X4', 'X5', 'X6', 'X8', 'X10',
             ...
             'X375', 'X376', 'X377', 'X378', 'X379', 'X380', 'X382', 'X383', 'X384',
             'X385'],
            dtype='object', length=365)
```

```python
[48]: df_train['X0'].unique()
```

```
[48]: array([37, 24, 46, 11, 41, 49, 36, 34, 45, 40, 23, 32, 50, 51,  9, 10, 12,
             52, 43, 18, 15, 48,  6,  0, 31,  8, 30, 16, 29,  1, 26, 17, 35, 44,
             25, 22, 28, 47,  4, 19, 39, 38, 21, 14,  3, 33,  2])
```

```python
[49]: #Task4:perform demensionlity reduction
```

```python
[50]: from sklearn.model_selection import train_test_split
      from sklearn.decomposition import PCA
      pca=PCA(n_components=0.98,random_state=1)
      X=df_train.drop('y',axis=1)
      y=df_train['y']
```

```python
[51]: X_train , X_test , y_train , y_test = train_test_split(X,y,test_size=0.
      →2,random_state=1)
```

```python
[52]: print(X_train.shape)
      print(X_test.shape)
      print(y_train.shape)
      print(y_test.shape)
```

```
(3367, 364)
(842, 364)
(3367,)
(842,)
```

```
[53]: pca.fit(X)
```

```
[53]: PCA(n_components=0.98, random_state=1)
```

```
[54]: pca.n_components_
```

```
[54]: 12
```

```
[55]: pca.explained_variance_ratio_
```

```
[55]: array([0.40868988, 0.21758508, 0.13120081, 0.10783522, 0.08165248,
              0.0140934 , 0.00660951, 0.00384659, 0.00260289, 0.00214378,
              0.00209857, 0.00180388])
```

```
[56]: #Task5:predict your test df values using XGBOOST
```

```
[57]: import xgboost as xgb
      from sklearn.ensemble import RandomForestRegressor
      from sklearn.metrics import mean_squared_error
```

```
[58]: pca_X_train = pd.DataFrame(pca.transform(X_train))
      pca_X_test = pd.DataFrame(pca.transform(X_test))
      pca_test = pd.DataFrame(pca.transform(df_test))
```

```
[59]:  model = xgb.XGBRegressor(objective='reg:linear',learning_rate=0.1)
```

```
[60]: model.fit(pca_X_train,y_train)
```

    [10:44:38] WARNING: /workspace/src/objective/regression_obj.cu:167: reg:linear
    is now deprecated in favor of reg:squarederror.

```
[60]: XGBRegressor(base_score=0.5, booster=None, colsample_bylevel=1,
                   colsample_bynode=1, colsample_bytree=1, gamma=0, gpu_id=-1,
                   importance_type='gain', interaction_constraints=None,
                   learning_rate=0.1, max_delta_step=0, max_depth=6,
                   min_child_weight=1, missing=nan, monotone_constraints=None,
                   n_estimators=100, n_jobs=0, num_parallel_tree=1,
                   objective='reg:linear', random_state=0, reg_alpha=0, reg_lambda=1,
                   scale_pos_weight=1, subsample=1, tree_method=None,
                   validate_parameters=False, verbosity=None)
```

```
[61]: pred_y_test = model.predict(pca_X_test)
```

```
[62]: mse_score = mean_squared_error(y_test,pred_y_test)
```

```
[69]: print( mse_score)
```

    83.93976003748757
```

```
[ ]:
```