

```
In [168]: import pandas as pd
import numpy as np
import matplotlib.pyplot as plt

In [169]: train=pd.read_csv("ML/train.csv")

In [170]: test=pd.read_csv("ML/test.csv")

In [171]: train.head()

Out[171]:
   ID   y  X0  X1  X2  X3  X4  X5  X6  X8  ...  X375  X376  X377  X378  X379  X380  X382  X383  X384  X386
0  0  130.01  k  v  w  e  a  d  u  j  o  ...  0  0  1  0  0  0  0  0  0  0  0
1  13  78.02  k  t  w  e  d  y  i  o  ...  1  0  0  0  0  0  0  0  0  0  0
2  2  78.26  az  w  n  e  c  d  x  j  x  ...  0  0  0  0  0  0  0  1  0  0  0
3  9  80.02  az  t  n  f  d  x  i  e  ...  0  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  0
5 rows x 378 columns

In [172]: train.info()
<class 'pandas.core.frame.DataFrame'>
Int64Index: 4209 entries, 0 to 4208
Columns: 378 entries, ID to X386
dtypes: float64(1), int64(80), object(8)
memory usage: 12.1+ MB

In [173]: train.shape

Out[173]: (4209, 378)

In [174]: test.shape

Out[174]: (4209, 377)

In [175]: train_data=train.drop(['ID','y'],axis=1)
train_target=train.y

In [177]: test_data=test.drop(['ID'],axis=1)

REMOVING COLUMN WHERE VARIANCE IS ZERO

In [126]: train_data.var().sort_values(ascending=True)

C:\Users\B1940\AppData\Local\Temp\ipykernel_8688\3852086680.py:1: FutureWarning: Dropping of nuisance columns in DataFrame reductions (with 'numeric_only=None') is deprecated; in a future version t
His will raise TypeError. Select only valid columns before calling the reduction.
labelled_data1.var().sort_values(ascending=True)

Out[126]:
X297    0.000000
X296    0.000000
X290    0.000000
X295    0.000000
...
X384    0.248764
X395    0.249222
X382    0.248647
X337    0.248787
X217    0.250836
Length: 366, dtype: float64

In [129]: keep_columns = train_data.columns[train_data.nunique()[>1]
labelled_data1= train_data.loc[:,keep_columns].copy()

In [130]: keep_columns = test_data.columns[test_data.nunique()[>1]
labelled_data1= test_data.loc[:,keep_columns].copy()

In [131]: labelled_data1.var().sort_values(ascending=True)

C:\Users\B1940\AppData\Local\Temp\ipykernel_8688\588475679.py:1: FutureWarning: Dropping of nuisance columns in DataFrame reductions (with 'numeric_only=None') is deprecated; in a future version t
His will raise TypeError. Select only valid columns before calling the reduction.
labelled_data1.var().sort_values(ascending=True)

Out[131]:
X297    0.000238
X296    0.000238
X290    0.000238
X295    0.000238
...
X386    0.248764
X395    0.249222
X382    0.248647
X337    0.248787
X217    0.250836
Length: 366, dtype: float64

In [132]: labelled_data1.var().sort_values(ascending=True)

C:\Users\B1940\AppData\Local\Temp\ipykernel_8688\728021611.py:1: FutureWarning: Dropping of nuisance columns in DataFrame reductions (with 'numeric_only=None') is deprecated; in a future version t
His will raise TypeError. Select only valid columns before calling the reduction.
labelled_data1.var().sort_values(ascending=True)

Out[132]:
X297    0.000238
X296    0.000238
X290    0.000238
X295    0.000238
...
X394    0.248953
X337    0.249323
X392    0.248489
X391    0.248738
X317    0.248866
Length: 369, dtype: float64

In [133]: labelled_data1.shape

Out[133]: (4209, 373)

In [134]: labelled_data1.shape

Out[134]: (4209, 364)

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

In [135]: labelled_data1.isna().sum().sort_values(ascending=False).head(20)

Out[135]:
X0      0
X263    0
X261    0
X260    0
X269    0
X259    0
X257    0
X255    0
X254    0
X253    0
X251    0
X250    0
X249    0
X248    0
X247    0
X246    0
X244    0
X244    0
dtype: int64

In [136]: labelled_data1.isna().sum().sort_values(ascending=False).head(20)

Out[136]:
X0      0
X262    0
X261    0
X260    0
X269    0
X259    0
X258    0
X254    0
X253    0
X251    0
X250    0
X249    0
X248    0
X247    0
X246    0
X244    0
X244    0
dtype: int64

In [137]: labelled_data1.nunique().sort_values(ascending=False)

Out[137]:
X0      47
X2      44
X5      28
X1      27
X8      25
...
X133    1
X130    2
X136    2
X135    2
X128    2
X385    2
Length: 364, dtype: int64

In [138]: labelled_data1.nunique().sort_values(ascending=False)

Out[138]:
X0      49
X2      45
X5      32
X1      27
X8      25
...
X130    1
X136    2
X135    2
X127    2
X385    2
Length: 373, dtype: int64

Applying label encoder

In [139]: labelled_data.select_dtypes(include='object')

Out[139]:
   X0  X1  X2  X3  X4  X5  X6  X8
0  k  v  w  e  a  d  u  j  o
1  k  t  w  e  d  y  i  o
2  az  w  n  e  c  d  x  j  x
3  az  t  n  f  d  x  i  e
...
4204 ak  s  as  c  d  aa  d  h
4205 j  o  t  d  d  aa  h  h
4206 ak  v  r  a  d  aa  g  e
4207 al  t  e  f  d  aa  i  u
4208 z  i  w  e  d  aa  g  w
4209 rows x 8 columns

In [140]: from sklearn import preprocessing

In [141]: for col in labelled_data.columns:
if labelled_data.dtypes[col] == "object":
    le = preprocessing.LabelEncoder()
    le.fit(labelled_data[col])
    labelled_data[col] = le.transform(labelled_data[col])

In [142]: labelled_data[["X0","X1","X2","X3","X4","X5","X6","X8","X1"]]

Out[142]:
   X0  X1  X2  X3  X4  X5  X6  X8
0  32  37  0  3  21  9  14
1  32  21  19  4  3  28  11  14
2  20  24  34  2  3  27  9  23
3  20  21  34  5  3  27  11  4
4  20  23  34  5  3  12  3  13
...
4204  8  20  18  2  3  0  3  16
4205  31  15  40  3  3  0  7  7
4206  8  23  38  0  3  0  6  4
4207  9  19  25  5  3  0  11  20
4208  46  19  3  2  3  0  6  22
4209 rows x 8 columns

In [143]: final_train_data=labelled_data

In [144]: for col in labelled_data.columns:
if labelled_data.dtypes[col] == "object":
    le = preprocessing.LabelEncoder()
    le.fit(labelled_data[col])
    labelled_data1[col] = le.transform(labelled_data1[col])

In [145]: final_test_data=labelled_data1

In [146]: final_train_data.nunique().sort_values(ascending=True)

Out[146]:
X193    2
X260    2
X296    2
X258    2
X297    2
...
X8      25
X1      27
X5      29
X2      44
X0      47
Length: 364, dtype: int64

In [147]: final_test_data.nunique().sort_values(ascending=True)

Out[147]:
X192    2
X262    2
X295    2
X260    2
X299    2
...
X8      25
X1      27
X5      32
X2      45
X0      49
Length: 373, dtype: int64

In [148]: common_cols = list(set(final_test_data.columns).intersection(final_train_data.columns))
common_cols

Out[148]: ['X162',
'X211',
'X388',
'X229',
'X240',
'X251',
'X340',
'X373',
'X84',
'X331',
'X174',
'X125',
'X219',
'X110',
'X271',
'X373',
'X264',
'X263',
'X248',
'X249',
'X246',
'X247',
'X244',
'X245',
'X243',
'X242',
'X241',
'X238',
'X237',
'X236',
'X235',
'X234',
'X233',
'X232',
'X231',
'X230',
'X229',
'X228',
'X227',
'X226',
'X225',
'X224',
'X223',
'X222',
'X221',
'X220',
'X219',
'X218',
'X217',
'X216',
'X215',
'X214',
'X213',
'X212',
'X211',
'X210',
'X209',
'X208',
'X207',
'X206',
'X205',
'X204',
'X203',
'X202',
'X201',
'X200',
'X199',
'X198',
'X197',
'X196',
'X195',
'X194',
'X193',
'X192',
'X191',
'X190',
'X189',
'X188',
'X187',
'X186',
'X185',
'X184',
'X183',
'X182',
'X181',
'X180',
'X179',
'X178',
'X177',
'X176',
'X175',
'X174',
'X173',
'X172',
'X171',
'X170',
'X169',
'X168',
'X167',
'X166',
'X165',
'X164',
'X163',
'X162',
'X161',
'X160',
'X159',
'X158',
'X157',
'X156',
'X155',
'X154',
'X153',
'X152',
'X151',
'X150',
'X149',
'X148',
'X147',
'X146',
'X145',
'X144',
'X143',
'X142',
'X141',
'X140',
'X139',
'X138',
'X137',
'X136',
'X135',
'X134',
'X133',
'X132',
'X131',
'X130',
'X129',
'X128',
'X127',
'X126',
'X125',
'X124',
'X123',
'X122',
'X121',
'X120',
'X119',
'X118',
'X117',
'X116',
'X115',
'X114',
'X113',
'X112',
'X111',
'X110',
'X109',
'X108',
'X107',
'X106',
'X105',
'X104',
'X103',
'X102',
'X101',
'X100',
'X99',
'X98',
'X97',
'X96',
'X95',
'X94',
'X93',
'X92',
'X91',
'X90',
'X89',
'X88',
'X87',
'X86',
'X85',
'X84',
'X83',
'X82',
'X81',
'X80',
'X79',
'X78',
'X77',
'X76',
'X75',
'X74',
'X73',
'X72',
'X71',
'X70',
'X69',
'X68',
'X67',
'X66',
'X65',
'X64',
'X63',
'X62',
'X61',
'X60',
'X59',
'X58',
'X57',
'X56',
'X55',
'X54',
'X53',
'X52',
'X51',
'X50',
'X49',
'X48',
'X47',
'X46',
'X45',
'X44',
'X43',
'X42',
'X41',
'X40',
'X39',
'X38',
'X37',
'X36',
'X35',
'X34',
'X33',
'X32',
'X31',
'X30',
'X29',
'X28',
'X27',
'X26',
'X25',
'X24',
'X23',
'X22',
'X21',
'X20',
'X19',
'X18',
'X17',
'X16',
'X15',
'X14',
'X13',
'X12',
'X11',
'X10',
'X9',
'X8']

In [149]: final_test_data[final_test_data[common_cols]]
final_train_data[final_train_data[common_cols]]

In [150]: final_train_data.head()

Out[150]:
   X162  X21  X30  X35  X39  X28  X83  X240  X340  X373  ...  X100  X111  X45  X267  X273  X100  X3  X183  X333  X371
0  0  1  0  0  0  1  0  0  0  0  0  0  ...  0  1  0  0  0  1  0  34  0  0  0
1  1  0  0  0  0  1  0  0  0  0  0  0  ...  1  0  0  0  0  1  0  14  0  0  0
2  1  0  0  0  0  1  0  0  0  0  0  0  ...  0  0  0  0  0  1  0  23  0  0  0
3  1  0  0  0  0  0  1  0  0  0  0  0  ...  0  0  0  0  0  1  0  4  0  0  0
4  1  0  0  0  0  1  0  0  0  0  0  0  ...  0  0  0  0  0  1  0  13  0  0  0
5 rows x 359 columns

In [151]: final_test_data.head()

Out[151]:
   X162  X21  X30  X35  X39  X28  X83  X240  X340  X373  ...  X100  X111  X45  X267  X273  X100  X3  X183  X333  X371
0  1  0  0  0  0  1  0  0  0  0  0  0  ...  0  0  0  0  0  1  0  34  0  0  0
1  0  0  0  0  0  0  0  0  0  0  0  0  ...  0  0  0  0  0  1  0  14  0  0  0
2  1  0  0  0  0  1  1  0  0  0  0  0  ...  1  0  0  0  0  0  9  0  0  0  0
3  1  0  0  0  0  1  0  0  0  0  0  0  ...  0  0  0  0  0  1  13  0  0  0  1
4  0  0  0  0  0  1  0  0  0  0  0  0  ...  1  0  0  0  0  1  0  12  0  0  0
5 rows x 359 columns

Performing dimensionality reduction.

In [152]: final_train_data.describe()

Out[152]:
   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  4209.000000  4209.000000  4209.000000
std      0.060965  0.000013  0.002722  0.000702  0.002628  0.000162  0.000186  0.000161  0.002338  0.000044  ...  0.000188  0.000018  0.000028  0.000028  0.000028  0.000028
min      0.000000  0.000000  0.000000  0.000000  0.000000  0.000000  0.000000  0.000000  0.000000  0.000000  ...  0.000000  0.000000  0.000000  0.000000  0.000000  0.000000
50%      0.000000  0.000000  0.000000  0.000000  0.000000  0.000000  0.000000  0.000000  0.000000  0.000000  ...  1.000000  0.000000  0.000000  0.000000  0.000000  0.000000
75%      0.000000  0.000000  0.000000  0.000000  0.000000  0.000000  0.000000  0.000000  0.000000  0.000000  ...  1.000000  0.000000  0.000000  0.000000  0.000000  0.000000
max      1.000000  1.000000  1.000000  1.000000  1.000000  1.000000  1.000000  1.000000  1.000000  1.000000  ...  1.000000  0.000000  1.000000  1.000000  1.000000  1.000000
5 rows x 359 columns

In [153]: from sklearn.preprocessing import StandardScaler # to standardize the features
from sklearn.decomposition import PCA # to apply PCA
import seaborn as sns # to plot the heat map

In [154]: scaler = StandardScaler()
final_train_data = pd.DataFrame(scaler.fit_transform(final_train_data)) #scaling the data
final_test_data = pd.DataFrame(scaler.fit_transform(final_test_data))

In [155]: final_train_data.shape

Out[155]: (4209, 359)

In [156]: final_train_data.shape

Out[156]: (4209, 359)

In [157]: from sklearn.model_selection import train_test_split
x_train, x_test, y_train, y_test = train_test_split(final_train_data,train_target,test_size=0.25, random_state=0)

In [158]: pca = PCA(n_components = 3)
x_train_pca = pca.transform(x_train)
x_test_pca = pca.transform(x_test)
final_train_data_pca = pca.transform(final_train_data)
final_test_data_pca = pca.transform(final_test_data)

Predicting test_df values using XGBBoost.

In [159]: pip install xgboost

Requirement already satisfied: xgboost in c:\users\b1940\anaconda11\site-packages (1.7.2)
Requirement already satisfied: scipy in c:\users\b1940\anaconda11\site-packages (from xgboost) (1.7.3)
Note: you may need to restart the kernel to use updated packages.

In [160]: from sklearn.metrics import r2_score,mean_squared_error
from xgboost import XGBRegressor
xgb=XGBRegressor(random_state=0)

In [161]: model=xgb.fit(x_train,y_train)

In [162]: y_pred_train=model.predict(x_train)
y_pred_train

Out[162]: array([104.73666, 96.26712, 92.62421, ..., 99.69573, 110.29155,
-134.31129], dtype=float32)

In [163]: y_train

Out[163]:
98071    108.62
951      93.41
187      91.21
87      84.80
97      108.37
...
1023    140.43
1044    104.73
103      102.77
2907    113.94
2732    116.19
Name: y, Length: 4209, dtype: float64

In [164]: r2_score(y_test,y_train,y_pred_train)

Out[164]: 0.860219917490613

In [165]: y_pred_test=model.predict(x_test)
y_pred_test

Out[165]: array([107.25132, 95.88886, 108.21682, ..., 90.55357, 98.342894,
-105.28956], dtype=float32)

In [166]: y_test

Out[166]:
3433    96.49
2333    96.93
2089    114.22
185      88.18
88      92.63
...
2358    111.21
71      133.40
1023    97.17
2850    110.80
189      110.41
Name: y, Length: 4209, dtype: float64

In [167]: from sklearn.metrics import mean_absolute_error, mean_squared_error

In [168]: print("MAE:", mean_absolute_error(y_train,y_pred_train))

MAE: 3.54529242296785
```

