

```
In [1]: import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns

In [2]: train_data = pd.read_csv(r'C:\Users\abha mohan\Desktop\train.csv')
test_data = pd.read_csv(r'C:\Users\abha mohan\Desktop\test.csv')

In [3]: train_data.head()
```

	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	0
1	6	88.53	k	t	a	v	e	d	y	i	o	...	1	0	0	0	0	0	0	0	0	0
2	7	76.26	a	z	w	n	c	d	x	j	x	...	0	0	0	0	0	0	1	0	0	0
3	9	80.62	a	z	t	n	f	d	x	i	e	...	0	0	0	0	0	0	0	0	0	0
4	13	78.02	a	z	v	n	f	d	h	d	n	...	0	0	0	0	0	0	0	0	0	0

5 rows × 378 columns


```
In [4]: train_data.shape

Out[4]: (4299, 378)
```



```
In [5]: train_data.dtypes

Out[5]: ID          int64
y          float64
X0         object
X1         object
X2         object
...         .
X380        int64
X382        int64
X383        int64
X384        int64
X385        int64
Length: 378, dtype: object

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



```
In [6]: train_data.var()
```



```
Out[6]: ID          5.941936e+06
y           1.687657e+02
X10         1.313892e-02
X11         0.808080e+00
X12         6.945713e-02
...         .
X380         8.814570e-03
X382         7.546747e-03
X383         1.660732e-03
X384         4.759536e-04
X385         1.423823e-03
Length: 378, dtype: float64
```



```
In [7]: for i in train_data.columns:
data_type = train_data[i].dtype
if data_type == 'object':
    print(i)
```

X0
X1
X2
X3
X4
X5
X6
X8


```
In [8]: variance = pow(train_data.drop(columns=['ID','y']).std(),2).to_dict()

null_cnt = 0
for key, value in variance.items():
    if(value==0):
        print(key)
    null_cnt = null_cnt+1
print('No of columns which has zero variance = ',null_cnt)
```

X11
X93
X197
X233
X235
X268
X289
X290
X293
X297
X330
X347
No of columns which has zero variance = 12


```
In [9]: train_data = train_data.drop(columns=['X11','X93','X107','X233','X235','X268','X289','X290','X293','X297','X330','X347'])
train_data.shape

Out[9]: (4299, 366)
```

check for null values


```
In [10]: train_data.isnull().sum().sum()

Out[10]: 0
```



```
In [11]: test_data.isnull().sum().sum()

Out[11]: 0
```

unique values in train and test data


```
In [12]: train_data['X0'].unique()

Out[12]: array(['k', 'az', 't', 'ai', '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'],
      dtype=object)
```



```
In [13]: train_data['X1'].unique()

Out[13]: array(['v', 't', 'w', 'b', 'r', 'i', 's', 'aa', 'c', 'a', 'e', 'h', 'z',
      'j', 'u', 'p', 'm', 'i', 'y', 'd', 'f', 'm', 'k', 'g', 'q',
      'ab'], dtype=object)
```



```
In [14]: train_data['X2'].unique()

Out[14]: array(['at', 'av', 'n', 'e', 'as', 'aq', 'r', 'ai', 'ak', 'm', 'a', 'aq', 'ag',
      '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'], dtype=object)
```



```
In [15]: train_data['X3'].unique()

Out[15]: array(['a', 'e', 'c', 'f', 'd', 'b', 'g'], dtype=object)
```



```
In [16]: train_data['X4'].unique()

Out[16]: array(['d', 'b', 'c', 'a'], dtype=object)
```



```
In [17]: train_data['X5'].unique()

Out[17]: array(['v', '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'], dtype=object)
```



```
In [18]: train_data['X6'].unique()

Out[18]: array(['j', 'i', 'd', 'h', 'i', 'a', 'g', 'c', 'k', 'e', 'f', 'b'],
      dtype=object)
```



```
In [19]: train_data['X8'].unique()

Out[19]: array(['o', 'x', 'e', 'n', 's', 'a', 'h', 'p', 'm', 'k', 'd', 'i', 'v',
      'j', 'b', 'q', 'w', 'g', 'y', 'l', 'r', 'u', 'r', 't', 'c'],
      dtype=object)
```



```
In [20]: test_data['X0'].unique()

Out[20]: array(['az', 't', 'w', 'y', 'x', 'f', 'ap', 'o', 'ay', 'ai', 'h', 'z',
      'aj', 'd', 'v', 'ak', 'ba', 'n', 'j', 's', 'af', 'ax', 'at', 'aq',
      'aw', 'm', 'k', 'a', 'e', 'ai', 'l', 'ag', 'b', 'am', 'aw', 'as',
      'r', 'ao', 'u', 'i', 'c', 'ad', 'au', 'bc', 'g', 'an', 'ae', 'p',
      'bb'], dtype=object)
```



```
In [21]: test_data['X1'].unique()

Out[21]: array(['v', 'b', 'i', 's', 'aa', 'r', 'a', 'i', 'p', 'c', 'o', 'm', 'z',
      'e', 'h', 'w', 'g', 'k', 'y', 't', 'u', 'd', 'j', 'q', 'n', 'f',
      'ab'], dtype=object)
```



```
In [22]: test_data['X2'].unique()

Out[22]: array(['n', 'ai', 'as', 'ae', 's', 'b', 'e', 'ak', 'm', 'a', 'aq', 'ag',
      'r', 'k', 'aj', 'ay', 'ao', 'an', 'ac', 'af', 'ax', 'h', 's', 'f',
      'ap', 'p', 'au', 't', 'z', 'y', 'aw', 'd', 'at', 'g', 'am', 'j',
      'x', 'ab', 'w', 'q', 'ah', 'ad', 'al', 'av', 'u'], dtype=object)
```



```
In [23]: test_data['X3'].unique()

Out[23]: array(['f', 'a', 'c', 'e', 'd', 'g', 'b'], dtype=object)
```



```
In [24]: test_data['X4'].unique()

Out[24]: array(['d', 'b', 'a', 'c'], dtype=object)
```



```
In [25]: test_data['X5'].unique()

Out[25]: array(['t', 'b', 'a', 'z', 'y', 'x', 'h', 'g', 'f', 'j', 'i', 'd', 'c',
      'q', 's', 'r', 'v', 'w', 'o', 'aa'], dtype=object)
```



```
In [26]: test_data['X6'].unique()

Out[26]: array(['a', 'g', 'j', 'l', 'i', 'd', 'f', 'h', 'c', 'k', 'e', 'b'],
      dtype=object)
```



```
In [27]: test_data['X8'].unique()

Out[27]: array(['w', 'y', 'j', 'n', 'm', 's', 'a', 'g', 'v', 'r', 'g', 't', 'h', 'c',
      'k', 'p', 'u', 'd', 'g', 'b', 'q', 'e', 'l', 'r', 'i', 'x'],
      dtype=object)
```

Apply Label encoder


```
In [28]: from sklearn.preprocessing import LabelEncoder
le = LabelEncoder()
```



```
In [29]: train_data['X0'] = le.fit_transform(train_data.X0)
train_data['X1'] = le.fit_transform(train_data.X1)
train_data['X2'] = le.fit_transform(train_data.X2)
train_data['X3'] = le.fit_transform(train_data.X3)
train_data['X4'] = le.fit_transform(train_data.X4)
train_data['X5'] = le.fit_transform(train_data.X5)
train_data['X6'] = le.fit_transform(train_data.X6)
train_data['X8'] = le.fit_transform(train_data.X8)
```



```
In [30]: features = train_data.drop(['y','ID'],axis=1)
target = train_data[['y']]

features.shape,target.shape

Out[30]: ((4299, 364), (4299, 1))
```



```
In [31]: from sklearn.model_selection import train_test_split
```



```
In [32]: X_train, X_test, y_train, y_test = train_test_split(features,target,train_size=0.8)
```



```
In [33]: print(X_train.shape)
print(X_test.shape)
print(y_train.shape)
print(y_test.shape)

(3367, 364)
(842, 364)
(3367, 1)
(842, 1)
```



```
In [34]: train_data.head()
```

	ID	y	X0	X1	X2	X3	X4	X5	X6	X8	...	X375	X376	X377	X378	X379	X380	X382	X383	X384	X385
0	0	130.81	32	23	17	0	3	24	9	14	...	0	0	1	0	0	0	0	0	0	0
1	6	88.53	32	21	19	4	3	28	11	14	...	1	0	0	0	0	0	0	0	0	0
2	7	76.26	20	24	34	2	3	27	9	23	...	0	0	0	0	0	0	1	0	0	0
3	9	80.62	20	21	34	5	3	27	11	4	...	0	0	0	0	0	0	0	0	0	0
4	13	78.02	20	23	34	5	3	12	3	13	...	0	0	0	0	0	0	0	0	0	0

5 rows × 366 columns


```
In [35]: train_data.dtypes

Out[35]: ID          int64
y          float64
X0          int32
X1          int32
X2          int32
...         .
X380        int64
X382        int64
X383        int64
X384        int64
X385        int64
Length: 366, dtype: object

Perform Dimensionality reduction
```



```
In [36]: from sklearn.decomposition import PCA
pca = PCA(n_components=9)
```



```
In [37]: pca.fit(X_train)

Out[37]: PCA(n_components=9.0)
```



```
In [38]: pca.explained_variance_ratio_

Out[38]: array([0.38849891, 0.21558917, 0.13119024, 0.11911883, 0.09364883])
```



```
In [39]: X_train_pca = pca.transform(X_train)
X_test_pca = pca.transform(X_test)
```



```
In [40]: X_train_pca.shape, X_test_pca.shape

Out[40]: ((3367, 5), (842, 5))
```



```
In [41]: X_train_pca

Out[41]: array([[ -9.19887269, -13.84530882, 12.33491469, -1.65654752,
      8.52639149],
      [  8.10739277, -2.86338488, -14.23373394, -5.04823912,
     -10.40350464],
      [-15.75807568,  1.24447381,  0.21806228, -6.37402159,
       -11.81788529],
      [ 12.33694493, -13.02810897, -4.83217329, 14.28191569,
      11.40132384],
      [-16.1222119 ,  7.15897352, -1.06019843, -5.78130533,
       7.73132272],
      [-5.5375811 , -0.11408029,  5.11700565,  1.43491012,
       0.79881494]])
```

Predict your test_df values using XGBoost


```
In [42]: import xgboost as xgb
from sklearn.model_selection import train_test_split
from sklearn.metrics import r2_score, mean_squared_error
from math import sqrt
```



```
In [43]: X_train,X_test,Y_train,Y_test = train_test_split(features,target,test_size=0.2)
print(X_train.shape)
print(X_test.shape)
print(X_test.shape)
print(Y_test.shape)

(3367, 364)
(3367, 1)
(842, 364)
(842, 1)
```



```
In [44]: xgb_reg = xgb.XGBRegressor(objective='reg:linear', colsample_bytree = 0.3, learning_rate = 0.4, max_depth = 10, alpha = 6,
                                n_estimators = 20)
model = xgb_reg.fit(X_train,y_train)
print('RMSE = ',sqrt(mean_squared_error(model.predict(X_test),y_test)))

[11:24:34] WARNING: C:\xgboost-split_1619728435298\work\src\objective\regression_obj.cu:178: reg:linear is now deprecated in favor of reg:squarederror.
RMSE = 13.701496319881393
```



```
In [45]: pred_y_test= model.predict(X_test)

plt.figure(figsize=(10,5))

sns.distplot(y_test[y_test<160], color="skyblue", label="Actual value")
sns.distplot(pred_y_test[pred_y_test<160], color="red", label="Predicted value")
plt.legend()

plt.tight_layout()
```

C:\Users\abha mohan\anaconda3\lib\site-packages\seaborn\distributions.py:2557: FutureWarning: 'distplot' is a deprecated function and will be removed in a future version. Please use plt.figure and use either 'displot' (a figure-level function with similar flexibility) or 'histplot' (an axes-level function for histograms).

Warning: warn(msg, FutureWarning)

C:\Users\abha mohan\anaconda3\lib\site-packages\seaborn\distributions.py:2557: FutureWarning: 'distplot' is a deprecated function and will be removed in a future version. Please use plt.figure and use either 'displot' (a figure-level function with similar flexibility) or 'histplot' (an axes-level function for histograms).

Warning: warn(msg, FutureWarning)

