

```
In [1]: #Import required libraries
import numpy as np
import pandas as pd

In [3]: #Load train dataset
data=pd.read_csv('train.csv')
datatest=pd.read_csv('test.csv')
data.head()
```

```
Out[3]:
```

	ID	y	X0	X1	X2	X3	X4	X5	X6	X8	...	X375	X376	X377	X378	X379	X380	X382	X383	X384	X385	
0	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	0
2	7	76.26	az	w	n	c	d	x	j	x	...	0	0	0	0	0	0	0	1	0	0	
3	9	80.62	az	t	n	f	d	x	i	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 × 378 columns

```
In [4]: #Get details for data
data.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]: #Remove columns with variance 0 as those columns will not impact the prediction and mostly those will have constant data
#to reduce number of cols remove those var=0 cols

data.var()[data.var()==0].index.values
```

```
C:\Users\dell\AppData\Local\Temp\ipykernel_7712\312826853.py:2: FutureWarning: Dropping of nuisance columns in DataFrame reductions (with 'numeric_only=None')
is deprecated; in a future version this will raise TypeError. Select only valid columns before calling the reduction.
  data.var()[data.var()==0].index.values
array(['X11', 'X93', 'X107', 'X233', 'X235', 'X268', 'X289', 'X290',
       'X293', 'X297', 'X330', 'X347'], dtype=object)
```

```
Out[6]:
```

```
In [8]: #drop 0 var columns
datafinal=data.drop(data.var()[data.var()==0].index.values, axis=1)
datafinal.info()
```

```
C:\Users\dell\AppData\Local\Temp\ipykernel_7712\247266354.py:2: FutureWarning: Dropping of nuisance columns in DataFrame reductions (with 'numeric_only=None')
is deprecated; in a future version this will raise TypeError. Select only valid columns before calling the reduction.
  datafinal=data.drop(data.var()[data.var()==0].index.values, axis=1)
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 4209 entries, 0 to 4208
Columns: 366 entries, ID to X385
dtypes: float64(1), int64(357), object(8)
memory usage: 11.8+ MB
```

```
In [9]: datafinal.head()
```

```
Out[9]:
```

	ID	y	X0	X1	X2	X3	X4	X5	X6	X8	...	X375	X376	X377	X378	X379	X380	X382	X383	X384	X385	
0	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	0
2	7	76.26	az	w	n	c	d	x	j	x	...	0	0	0	0	0	0	0	1	0	0	
3	9	80.62	az	t	n	f	d	x	i	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 × 366 columns

```
In [10]: #get features and labels
features=datafinal.iloc[:,2:]
label=datafinal.iloc[:,1].values
```

```
In [11]: features.head()
```

```
Out[11]:
```

	X0	X1	X2	X3	X4	X5	X6	X8	X10	X12	...	X375	X376	X377	X378	X379	X380	X382	X383	X384	X385
0	k	v	a	t	a	d	u	j	o	0	0	...	0	0	1	0	0	0	0	0	0
1	k	t	a	v	e	d	y	i	o	0	0	...	1	0	0	0	0	0	0	0	0
2	az	w	n	c	d	x	j	x	0	0	...	0	0	0	0	0	0	0	1	0	0
3	az	t	n	f	d	x	i	e	0	0	...	0	0	0	0	0	0	0	0	0	0
4	az	v	n	f	d	h	d	n	0	0	...	0	0	0	0	0	0	0	0	0	0

5 rows × 364 columns

```
In [12]: label
```

```
Out[12]: array([130.81,  88.53,  76.26, ..., 109.22,  87.48, 110.85])
```

```
In [13]: #Get object columns from feature column for LE and OHE
features.describe(include="O")#Get object columns from feature column for LE and OHE
features.describe(include="O")
```

```
Out[13]:
```

	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 [14]: #Get only object column names
objcols=features.describe(include='object').columns.values
objcols
```

```
Out[14]: array(['X0', 'X1', 'X2', 'X3', 'X4', 'X5', 'X6', 'X8'], dtype=object)
```

```
In [15]: #Apply Label Encoder for object columns. OHE is not use das this will create PC's instead of using direct variables

from sklearn.preprocessing import LabelEncoder
#from sklearn.preprocessing import OneHotEncoder

le=LabelEncoder()

for i in objcols:
    features[i]= le.fit_transform(features[i])

#ohe=OneHotEncoder(categorical_features=[0,1,2,3,4,5,6,8])

fea=features.values
#featureohe=fea #ohe.fit_transform(fea).toarray()

#featureohe
fea
```

```
Out[15]: array([[32, 23, 17, ..., 0, 0, 0],
       [32, 21, 19, ..., 0, 0, 0],
       [20, 24, 34, ..., 0, 0, 0],
       ...,
       [ 8, 23, 38, ..., 0, 0, 0],
       [ 9, 19, 25, ..., 0, 0, 0],
       [46, 19,  3, ..., 0, 0, 0]], dtype=int64)
```

```
In [16]: #Standardize the data using StandardScaler

from sklearn.preprocessing import StandardScaler
sc = StandardScaler()
featuresstd = sc.fit
```