

```
In [8]: import numpy as np
import pandas as pd

In [11]: train = pd.read_csv("C:\\Users\\VHP\\Desktop\\original projects\\impl\\machine learning\\project mercedes benz\\train.csv")
train.head()
```

Out[11]:

	ID	X0	X1	X2	X3	X4	X5	X6	X8	...	X375	X376	X377	X378	X379	X380	X382	X383	X384	X385	
0	0	130.81	k	v	at	a	d	u	j	o	...	0	0	1	0	0	0	0	0	0	0
1	6	88.53	k	t	av	e	d	y	i	e	...	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	78.92	az	v	n	f	d	h	d	n	...	0	0	0	0	0	0	0	0	0	0

5 rows × 378 columns

```
In [12]: test = pd.read_csv("C:\\Users\\VHP\\Desktop\\original projects\\impl\\machine learning\\project mercedes benz\\test.csv")
test.head()
```

Out[12]:

	ID	X0	X1	X2	X3	X4	X5	X6	X8	X10	...	X375	X376	X377	X378	X379	X380	X382	X383	X384	X385
0	1	az	v	n	f	d	t	a	w	o	...	0	0	0	1	0	0	0	0	0	0
1	2	t	b	ai	a	d	b	g	y	o	...	0	0	1	0	0	0	0	0	0	0
2	3	az	v	as	f	d	a	j	i	o	...	0	0	0	1	0	0	0	0	0	0
3	4	az	i	n	f	d	z	i	n	o	...	0	0	0	1	0	0	0	0	0	0
4	5	w	s	as	c	d	y	i	m	o	...	1	0	0	0	0	0	0	0	0	0

5 rows × 377 columns

```
In [14]: train.isnull().sum()
```

Out[14]:

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

```
In [15]: test.isnull().sum()
```

Out[15]:

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

```
In [17]: train.describe()
```

Out[17]:

	ID	y	X10	X11	X12	X13	X14	X15	X16	X17	...	X375	X376	X377	X378	X379	X380	X382	X383	X384	X385
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	4209.000000	4209.000000	4209.000000	4209.000000	4209.000000
mean	4205.960798	100.669318	0.013305	0.0	0.075077	0.057971	0.428130	0.000475	0.002613	0.007603	...	0.318841	0.057258	0.314802	0.020670	0.009503	0.008078	0.009678	0.009678	0.009678	0.009678
std	2437.630888	12.679381	0.114590	0.0	0.263547	0.237316	0.494867	0.021796	0.051061	0.086872	...	0.466082	0.232363	0.464492	0.142294	0.097033	0.089524	0.089524	0.089524	0.089524	0.089524
min	0.000000	72.110000	0.000000	0.0	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
25%	2095.000000	90.820000	0.000000	0.0	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%	4230.000000	99.150000	0.000000	0.0	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
75%	6314.000000	109.010000	0.000000	0.0	0.000000	0.000000	1.000000	0.000000	0.000000	0.000000	...	1.000000	0.000000	1.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000
max	8417.000000	265.320000	1.000000	0.0	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000	...	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000

8 rows × 370 columns

```
In [18]: test.describe()
```

Out[18]:

	ID	X10	X11	X12	X13	X14	X15	X16	X17	...	X375	X376	X377	X378	X379	X380	X382	X383	X384	X385
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	4209.000000	4209.000000	4209.000000	4209.000000
mean	4211.039202	0.019007	0.000238	0.074364	0.061060	0.427893	0.000713	0.002613	0.008791	0.010216	...	0.325968	0.049556	0.311951	0.019244	0.011679	0.009678	0.009678	0.009678	0.009678
std	2423.079826	0.136565	0.015414	0.262394	0.239466	0.494832	0.026991	0.051061	0.093257	0.100570	...	0.468791	0.217258	0.462345	0.127399	0.106256	0.089524	0.089524	0.089524	0.089524
min	1.000000	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	0.000000	0.000000	0.000000
25%	2115.000000	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	0.000000	0.000000	0.000000
50%	4230.000000	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	0.000000	0.000000	0.000000
75%	6314.000000	0.000000	0.000000	0.000000	0.000000	1.000000	0.000000	0.000000	0.000000	0.000000	...	1.000000	0.000000	1.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000
max	8416.000000	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000	...	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000

8 rows × 369 columns

```
In [31]: train_target = train["y"]
df_train = train.drop(["y","ID"],axis=1)
df_train.head(5)
```

Out[31]:

	X0	X1	X2	X3	X4	X5	X6	X8	X10	X11	...	X375	X376	X377	X378	X379	X380	X382	X383	X384	X385
0	k	v	at	a	d	u	j	o	0	0	...	0	0	1	0	0	0	0	0	0	0
1	k	t	av	e	d	y	i	o	0	0	...	1	0	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	1	0	0	0
3	az	t	n	f	d	x	i	e	o	0	...	0	0	0	0	0	0	0	0	0	0
4	az	v	n	f	d	h	d	n	o	0	...	0	0	0	0	0	0	0	0	0	0

5 rows × 376 columns

```
In [33]: df_train.var().sort_values().head(5)
```

Out[33]:

```
X330    0.0
X297    0.0
X268    0.0
X259    0.0
X235    0.0
dtype: float64
```

```
In [68]: import seaborn as sns
from sklearn.feature_selection import VarianceThreshold
variance = VarianceThreshold(threshold=0)
from sklearn.preprocessing import LabelEncoder
label = LabelEncoder()
import matplotlib.pyplot as plt
%matplotlib inline
```

```
In [34]: train_without_zero_var = variance.fit_transform(df_train.iloc[:,9:])
train_without_zero_var
```

Out[34]:

```
array([[0, 1, 0, ..., 0, 0, 0],
       [0, 0, 0, ..., 0, 0, 0],
       [0, 0, 0, ..., 0, 0, 0],
       ...,
       [1, 1, 0, ..., 0, 0, 0],
       [0, 0, 1, ..., 0, 0, 0],
       [0, 0, ..., 0, 0, 0]], dtype=int64)
```

```
In [35]: labeled = df_train.iloc[:,9:8]
labeled.head()
```

Out[35]:

	X0	X1	X2	X3	X4	X5	X6	X8
0	k	v	at	a	d	u	j	o
1	k	t	av	e	d	y	i	o
2	az	w	n	c	d	x	j	x
3	az	t	n	f	d	x	i	e
4	az	v	n	f	d	h	d	n

```
In [36]: labeled.nunique()
```

Out[36]:

```
X0      47
X1      27
X2      44
X3       7
X4       4
X5      29
X6      12
X8      25
dtype: int64
```

```
In [38]: labeled_data = labeled.apply(label().fit_transform)
labeled_data.head()
```

Out[38]:

	X0	X1	X2	X3	X4	X5	X6	X8
0	32	23	17	0	3	24	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

```
In [39]: labeled_data.var()
```

Out[39]:

```
X0      188.741938
X1      72.777974
X2     118.808135
X3      3.627295
X4      0.805461
X5      68.676236
X6      8.580720
X8     49.531868
dtype: float64
```

```
In [48]: train_zero_var = pd.DataFrame(train_without_zero_var)
train_zero_var.head()
```

Out[48]:

	0	1	2	3	4	5	6	7	8	9	...	345	346	347	348	349	350	351	352	353	354
0	0	1	0	0	0	0	1	0	0	1	...	0	0	1	0	0	0	0	0	0	0
1	0	0	0	0	0	1	0	0	0	1	0	0	0	0	0	0	0	0	0
2	0	0	0	0	0	1	0	0	0	0	...	0	0	0	0	0	0	1	0	0	0
3	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0	0	0	0	0
4	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0	0	0	0	0

5 rows × 355 columns

```
In [42]: train_data = pd.concat([labeled_data,train_zero_var],axis=1)
train_data.head()
```

Out[42]:

	X0	X1	X2	X3	X4	X5	X6	X8	0	1	...	345	346	347	348	349	350	351	352	353	354
0	32	23	17	0	3	24	9	14	0	1	...	0	0	1	0	0	0	0	0	0	0
1	32	21	19	4	3	28	11	14	0	1	0	0	0	0	0	0	0	0	0
2	20	24	34	2	3	27	9	23	0	0	0	0	0	0	0	1	0	0	0
3	20	21	34	5	3	27	11	4	0	0	0	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	0	0	0

5 rows × 363 columns

```
In [44]: train_data.isnull().any()
```

Out[44]:

```
X0      False
X1      False
X2      False
X3      False
X4      False
...
X350      False
X351      False
X352      False
X353      False
X354      False
Length: 363, dtype: bool
```

```
In [45]: test = test.drop(["ID"],axis=1)
test.head()
```

Out[45]:

	X0	X1	X2	X3	X4	X5	X6	X8	X10	X11	...	X375	X376	X377	X378	X379	X380	X382	X383	X384	X385
0	az	v	n	f	d	t	a	w	o	0	...	0	0	0	1	0	0	0	0	0	0
1	t	b	ai	a	d	b	g	y	o	0	...	0	0	1	0	0	0	0	0	0	0
2	az	v	as	f	d	a	j	i	o	0	...	0	0	0	1	0	0	0	0	0	0
3	az	i	n	f	d	z	i	n	o	0	...	0	0	0	1	0	0	0	0	0	0
4	w	s	as	c	d	y	i	m	o	0	...	1	0	0	0	0	0	0	0	0	0

5 rows × 376 columns

```
In [46]: test.nunique()
```

Out[46]:

```
X0      49
X1      27
X2      45
X3       7
X4       4
...
X380     2
X382     2
X383     2
X384     2
X385     2
Length: 376, dtype: int64
```

```
In [47]: test.isnull().any()
```

Out[47]:

```
X0      False
X1      False
X2      False
X3      False
X4      False
...
X380      False
X382      False
X383      False
X384      False
X385      False
Length: 376, dtype: bool
```

```
In [48]: test.var().sort_values().head()
```

Out[48]:

```
X255    0.0
X369    0.0
X296    0.0
X257    0.0
X258    0.0
dtype: float64
```

```
In [49]: test_without_zero_var = variance.fit_transform(test.iloc[:,9:])
test_without_zero_var
```

Out[49]:

```
array([[0, 0, 0, ..., 0, 0, 0],
       [0, 0, 0, ..., 0, 0, 0],
       [0, 0, 0, ..., 0, 0, 0],
       ...,
       [0, 0, 0, ..., 0, 0, 0],
       [0, 0, 1, ..., 0, 0, 0],
       [0, 0, ..., 0, 0, 0]], dtype=int64)
```

```
In [58]: test_data = pd.DataFrame(test_without_zero_var)
test_data.head()
```

Out[58]:

	0	1	2	3	4	5	6	7	8	9	...</
--	---	---	---	---	---	---	---	---	---	---	-------