

Part :1 (filter method used to get the unique features)

```
In [3]: import numpy as np
import matplotlib.pyplot as plt
import pandas as pd
import seaborn as sns
matplotlib.rcParams['font.family'] = 'serif'

In [4]: from sklearn.feature_selection import VarianceThreshold
from sklearn.ensemble import RandomForestClassifier
from sklearn.model_selection import train_test_split
from sklearn.metrics import accuracy_score

In [5]: data = pd.read_csv('sanctander_train.csv',nrows=2000)
data.head()
```

ID	var3	var15	imp_ent	var16	ult1	imp_op	var39	comer	ult1	imp_op	var39	comer	ult1	imp_op	var40	comer	ult1	imp_op	var40	comer	ult1
0	1	2	23		0.0			0.0				0.0				0.0				0.0	
1	3	2	34		0.0			0.0				0.0				0.0				0.0	
2	4	2	23		0.0			0.0				0.0				0.0				0.0	
3	8	2	37		0.0		1950				1950					0.0				0.0	
4	10	2	39		0.0			0.0				0.0				0.0				0.0	

5 rows × 21 columns

```
In [6]: x = data.drop( labels="TARGET",axis=1)
y = data["TARGET"]
x_train,x_test,y_train,y_test = train_test_split(x,y,train_size =.80,random_state=0,stratify=y)
x_train.shape,x_test.shape

Out[6]: (12000, 370), (2000, 1)
```

```
In [7]: x_train,x_test,y_train,y_test = train_test_split(x,y,train_size =.80,random_state=0,stratify=y)
x_train.shape,x_test.shape

Out[7]: (12000, 370), (400, 370)
```

Remove constant,quasi and duplicate

```
In [8]: #remove constant and quasi constant

In [9]: constant_filter = VarianceThreshold(threshold=0.01)
constant_filter.fit(x_train)
x_train_filter = constant_filter.transform(x_train)
x_test_filter = constant_filter.transform(x_test)
x_train_filter.shape,x_test_filter.shape

Out[9]: (11600, 223), (400, 223)
```

#Duplicate Feature removal

```
In [10]:

In [11]: x_train_T = x_train_filter.T
x_test_T = x_test_filter.T
x_train_T = pd.DataFrame(x_train_T)
x_test_T = pd.DataFrame(x_test_T)
duplicate_features = x_train_T.duplicated()
keep_them = [not index for index in duplicate_features]
x_train_unique = x_train_T[keep_them].T
x_test_unique = x_test_T[keep_them].T
x_train_unique.shape,x_test_unique.shape

Out[11]: (11600, 200), (400, 200)
```

```
In [12]: %time
def randomforest(x_train,x_test,y_train,y_test):
    clf = RandomForestClassifier(n_estimators=1000,n_jobs=-1,random_state=0)
    y_pred = clf.predict(x_test)
    p_fitted=accuracy_score(y_test,y_pred)
    randomforest(x_train,x_test,y_train,y_test)

Accuracy : 0.955
Wall time: 11.2 s

In [13]: %time
randomforest(x_train,x_test,y_train,y_test)

Accuracy : 0.955
Wall time: 12.5 s
```

However we got necessary model with the accuracy of 95 %

Part :2(by using the correlation coefficient remove the correlated features)

Feature Selection with Filtering Method - Correlated Feature Removal

In that we are going to remove the feature those who are correlated by certain coefficient or correlated.

If the correlation coefficient is more than 0.85 or less than -0.85 in that case features are highly correlated and these features might not give us the extra information to make the final prediction of classification so we remove those features and by doing this ultimately improve our model and it will take less time to train our model.

A dataset can be correlated features.Two or more than two features are correlated if they are close to each other in the linear space.

Correlation between the output observations and input features is very important and such features should be retained.

Our aim is to reduce the correlation function and reduce dimensionality of the overall feature space.

Because our model gives us the better performance for the particular position at time performance of classifier high. But after the some time as soon as dimensionality of the feature space is increase that time mode start to decreasing its performance.

So that is always recommended or desirable that we should have select optimal number of features.

In the feature space there are lots of the features are available that can predicted themselves we use one of them at one time.So this kind of feature not giving any information for prediction of the model so that is why its very necessary to remove this kind features.Those feature also called as co-linear feature.

Co-linearity in the feature space and output is always desirable but the co-linearity amongs the features space space is not desirable.

Summary:-

Feature Space to target coefficient is desired.

Feature to feature is not desirable.

If two features are highly correlated then either features is redundant (either we have to remove feature 1 or feature 2)

Correlation is feature space increases the model complexity.

Removing Correlated features improve the model performance.

Different model shows the different performance over the correlated features.

Before the work on the correlation feature we need to think about the previous part-1 there we have got the x_train_unique,x_test_unique.

Afterwards we have to think about the correlation feature.First we will calculate pearson correlation coefficient for training and testing datasets and then we will select only those features which have correlation less than 0.85.

```
In [14]: #first we are going to use the correlation matrix by using cormat

In [15]: cormat = x_train_unique.corr()

In [16]: plt.figure(figsize = (12,8))
sns.heatmap(cormat)

Out[16]: <matplotlib.axes._subplots.AxesSubplot at 0xc823b3b>
```

Here if see the black and white plot then which the features are intersecting those are the features are very highly correlated.

White is positive correlation and black represent the negative correlation.

If we observed that the black is below -0.85% and white is above the +85%.

```
In [17]: #we are using the python method to get the correlation features
#here we are defining function that is correlation with attributes are data,threshold
#data = x_train_unique to reduce the chances of overfitting
#threshold is the limit upto which we count as no correlative feature but above which is counted as correlative
#threshold means gives the limit after which we discard the features.
def correlation(data,threshold):
    #there could be the possibility of the same features occurrence on multiple time.
    #here we are using set with the corr_col the set method is avoid to use the same feature multipletime.
    corr_col = set()
    #then we use the cormat which is presenting the correlation between the features of the data which we will provide
    #data = data.cormat
    #after we use the for loop for the range(len(cormat.columns))
    #that means it is iterating features from features range it will start from 0 to final number from 0 to final.
    #range(0,213)
    for i in range(len(cormat.columns)):
        #by using above its possible to get the situation range(1,1),range(1,2) to avoid that kind of situation
        #if we use the range in that we containing the same features that time the correlation will become 1 because we
        #that kind of correlation we not to consider.
        for j in range(1):
            #after this condition it will run whenever i and j values in cormat like when the i is 1 that time j will be 0 as
            #time j is will be the 1 and also it take care if we discard the 0 and 1 then it will not repeat the same iteration.
            #we mention the condition when the abs(absolute values) of the cormat.iloc[ilocation][i,j] will be greater than
            if abs(cormat.iloc[i,j])>threshold:
                #that time,what feature will above the threshold they will add in corr_col and it will add those features which
                colname = cormat.columns[i]
                #if we have list is will come in the cormat which are available in correlation format and it will come add to the corr_col
                corr_col.add(colname)
    return corr_col

In [18]: def correlation(data,threshold):
    corr_col = set()
    cormat = data.corr()
    for i in range(len(cormat.columns)):
        for j in range(i):
            if abs(cormat.iloc[i,j]) > threshold:
                colname = cormat.columns[i]
                corr_col.add(colname)
    return corr_col
length: 2000

In [19]: correlated_features = correlation(x_train_unique,0.85)
correlated_features

Out[19]: (5,
1,
11,
12,
14,
16,
17,
18,
23,
26,
27,
32,
34,
35,
37,
41,
46,
49,
50,
51,
52,
54,
55,
56,
58,
60,
63,
65,
66,
67,
68,
70,
73,
74,
76,
78,
82,
84,
86,
87,
89,
91,
96,
99,
102,
103,
104,
105,
109,
110,
112,
113,
115,
116,
119,
120,
122,
123,
127,
128,
129,
130,
131,
132,
133,
134,
135,
136,
137,
138,
139,
140,
141,
142,
143,
144,
145,
146,
147,
148,
149,
150,
151,
154,
156,
157,
158,
159,
160,
161,
162,
163,
164,
165,
166,
167,
168,
169,
170,
171,
172,
173,
174,
175,
176,
177,
178,
179,
180,
181,
182,
183,
184,
185,
186,
187,
188,
189,
190,
191,
192,
193,
194,
195,
196,
197,
198,
199,
200,
201,
202,
203,
204,
205,
206,
207,
208,
209,
210,
211,
212,
213,
214,
215,
216,
218,
219,
220,
221,
222)
```

```
In [20]: len(correlated_features)

Out[20]: 121

We have got the correlated feature then our next task is to remove these features and then strat to train model.

In [21]: x_train_uncorr = x_train_unique.drop(labels=correlated_features,axis=1)
x_test_uncorr = x_test_unique.drop(labels=correlated_features,axis=1)

In [22]: x_train_uncorr.shape,x_test_uncorr.shape

Out[22]: (11600, 79), (400, 79)
```

```
In [26]: %time
randomforest(x_train_uncorr,x_test_uncorr,y_train,y_test)

Accuracy : 0.9755
Wall time: 9.75 s

In [24]: %time
randomforest(x_train,x_test,y_train,y_test)

Accuracy : 0.955
Wall time: 10.9 s

In [28]: (10.9-9.75)*100/9.75

Out[28]: 11.794871794871797
```

From the above we can only say that the it does not making any on the accuracy but it impacting on the training time which is equired less than the original datasets.

Here we use the broot force technique to reduce the dimensionality of the datasets by removing the correlated features.

Till now we can make any statement that our accuracy doesn't improve pretty much are very high but the training time has reduced the effect source around 11% of the training time has been reduced to increase if you are doing a very big data set in that place where you can remove uncorrelated features and uncorrelated features and you can reduce dimensions of your data and you can handle curse of dimension what we did we actually we did kind of the brute force that we had removed all the correlated features but you know this is not always recommended there could be a places where some group of features are correlated with each others but that mean they only need to remove a certain features we can keep some features by testing the features importance for example it's a year in out of 121 features there would be kind of the 10 groups where they are correlated that we can avoid these all the correlated features in kind of the groups and then finally we can select a single feature in each of those groups after performing feature importance so here now we are going to do what feature selection and the grouping of these correlated features.

```
In [40]: cormat

Out[40]:
```

	0	1	2	3	4	5	6	7	8	9	...	212	213	
0	1.000000	-0.027387	-0.005389	0.035624	0.031620	0.014071	0.022963	0.023790	0.031051	-0.002832	...	-0.007409	-0.017482	-0.015
1	-0.027387	1.000000	0.129000	-0.003721	0.162310	0.183237	0.026635	0.027339	-0.002888	-0.002831	...	0.093521	-0.002286	-0.026
2	-0.005389	0.129000	1.000000	0.039272	0.094338	0.098071	0.041147	0.033639	0.020170	0.018505	...	0.122776	-0.021590	-0.002
3	0.035624	-0.003721	0.039272	1.000000	0.084736	0.044504	-0.004283	-0.003979	-0.003148	-0.003086	...	-0.006308	-0.002492	-0.003
4	0.031620	0.162310	0.094338	0.084736	1.000000	0.871544	0.388872	0.368375	0.242827	0.216412	...	-0.014652	-0.005787	-0.007
5	0.014071	0.183237	0.098071	0.044504	0.871544	1.000000	-0.001592	-0.001479	-0.001770	-0.001147	...	-0.002345	-0.000926	-0.001
6	0.022963	0.026635	0.027339	-0.004283	-0.003979	-0.001592	1.000000	-0.001196	-0.001111	-0.000879	...	-0.001762	-0.000696	-0.000
7	0.023790	0.033639	0.020170	0.003148	0.003086	0.001147	-0.001196	1.000000	-0.001233	-0.001215	...	-0.002483	-0.000976	-0.001
8	0.031051	-0.002832	0.018505	-0.003086	-0.002492	-0.001147	-0.001233	-0.001215	1.000000	-0.001208	...	-0.002470	-0.000981	-0.001
9	-0.002832	0.002831	0.018505	-0.003086	-0.002492	-0.001147	-0.001208	-0.001215	-0.002470	1.000000	...	0.001086	-0.001878	-0.003
...
212	-0.007409	0.093521	-0.006308	-0.014652	-0.005787	-0.002345	-0.001762	-0.000976	-0.002483	0.001086	...	1.000000	-0.001878	-0.003
213	-0.017482	-0.002286	-0.002492	-0.002470	-0.000926	-0.000976	-0.000696	-0.000981	-0.000981	-0.001878	...	-0.001878	1.000000	-0.003

```
In [41]: #lets do the the vertically stacking (make the 0th row into the form 0th column)
cormat = cormat.abc().stack()
cormat

Out[41]:
```

	0	1	2	3	4	5	6	7	8	9	...	212	213	
0	1.000000	-0.027387	-0.005389	0.035624	0.031620	0.014071	0.022963	0.023790	0.031051	-0.002832	...	-0.007409	-0.017482	-0.015
1	-0.027387	1.000000	0.129000	-0.003721	0.162310	0.183237	0.026635	0.027339	-0.002888	-0.002831	...	0.093521	-0.002286	-0.026
2	-0.005389	0.129000	1.000000	0.039272	0.094338	0.098071	0.041147	0.033639	0.020170	0.018505	...	0.122776	-0.021590	-0.002
3	0.035624	-0.003721	0.039272	1.000000	0.084736	0.044504	-0.004283	-0.003979	-0.003148	-0.003086	...	-0.006308	-0.002492	-0.003
4	0.031620	0.162310	0.094338	0.084736	1.000000	0.871544	0.388872	0.368375	0.242827	0.216412	...	-0.014652	-0.005787	-0.007
5	0.014071	0.183237	0.098071	0.044504	0.871544	1.000000	-0.001592	-0.001479	-0.001770	-0.001147	...	-0.002345	-0.000926	-0.001
6	0.022963	0.026635	0.027339	-0.004283	-0.003979	-0.001592	1.000000	-0.001196	-0.001111	-0.000879	...	-0.001762	-0.000696	-0.000
7	0.023790	0.033639	0.020170	0.003148	0.003086	0.001147	-0.001196	1.000000	-0.001233	-0.001215	...	-0.002483	-0.000976	-0.001
8	0.031051	-0.002832	0.018505	-0.003086	-0.002492	-0.001147	-0.001233	-0.001215	1.000000	-0.001208	...	-0.002470	-0.000981	-0.001
9	-0.002832	0.002831	0.018505	-0.003086	-0.002492	-0.001147	-0.001208	-0.001215	-0.002470	1.000000	...	0.001086	-0.001878	-0.003
...
212	-0.007409	0.093521	-0.006308	-0.014652	-0.005787	-0.002345	-0.001762	-0.000976	-0.002483	0.001086	...	1.000000	-0.001878	-0.003
213	-0.017482	-0.002286	-0.002492	-0.002470	-0.000926	-0.000976	-0.000696	-0.000981	-0.000981	-0.001878	...	-0.001878	1.000000	-0.003

```
In [42]: #lets do the the vertically stacking (make the 0th row into the form 0th column)
cormat = cormat.abc().stack()
cormat

Out[42]:
```

	0	1	2	3	4	5	6	7	8	9	...	212	213	
0	1.000000	-0.027387	-0.005389	0.035624	0.031620	0.014071	0.022963	0.023790	0.031051	-0.002832	...	-0.007409	-0.017482	-0.015
1	-0.027387	1.000000	0.129000	-0.003721	0.162310	0.183237	0.026635	0.027339	-0.002888	-0.002831	...	0.093521	-0.002286	-0.026
2	-0.005389	0.129000	1.000000	0.039272	0.094338	0.098071	0.041147	0.033639	0.020170	0.018505	...	0.122776	-0.021590	-0.002
3	0.035624	-0.003721	0.039272	1.000000	0.084736	0.044504	-0.004283	-0.003979	-0.003148	-0.003086	...	-0.006308	-0.002492	-0.003
4	0.031620	0.162310	0.094338	0.084736	1.000000	0.871544	0.388872	0.368375	0.242827	0.216412	...	-0.014652	-0.005787	-0.007
5	0.014071	0.183237	0.098071	0.044504	0.871544	1.000000	-0.001592	-0.001479	-0.001770	-0.001147	...	-0.002345	-0.000926	-0.001
6	0.022963	0.026635	0.027339	-0.004283	-0.003979	-0.001592	1.000000	-0.001196	-0.001111	-0.000879	...	-0.001762	-0.000696	-0.000
7	0.023790	0.033639	0.020170	0.003148	0.003086	0.001147	-0.001196	1.000000	-0.001233	-0.001215	...	-0.002483	-0.000976	-0.001
8	0.031051	-0.002832	0.018505	-0.003086	-0.002492	-0.001147	-0.001233	-0.001215	1.000000	-0.001208	...	-0.002470	-0.000981	-0.001
9	-0.002832	0.002831	0.018505	-0.003086	-0.002492	-0.001147	-0.001208	-0.001215	-0.002470	1.000000	...	0.001086	-0.001878	-0.003
...
212	-0.007409	0.093521	-0.006308	-0.014652	-0.005787	-0.002345	-0.001762	-0.000976	-0.002483	0.001086	...	1.000000	-0.001878	-0.003
213	-0.017482	-0.002286	-0.002492	-0.002470	-0.000926	-0.000976	-0.000696	-0.000981	-0.000981	-0.001878	...	-0.001878	1.000000	-0.003


```
[features      23.00000
importance    0.290988      features      25.00000
Name: 2, dtype: float64, features
importance    0.322323
Name: 4, dtype: float64, features      104.000000
importance    0.522979
Name: 1, dtype: float64, features      168.000000
importance    0.384278
Name: 4, dtype: float64, features      80.000000
importance    0.506065
Name: 0, dtype: float64, features      180.000000
importance    0.437209
Name: 5, dtype: float64, features      147.000000
importance    0.160047
Name: 1, dtype: float64, features      14.000000
importance    0.360863
Name: 0, dtype: float64, features      82.000000
importance    0.186798
Name: 5, dtype: float64, features      192.000000
importance    0.359525
Name: 0, dtype: float64, features      114.000000
importance    0.523164
Name: 1, dtype: float64, features      109.000000
importance    0.270626
Name: 3, dtype: float64, features      165.000000
importance    0.476715
Name: 2, dtype: float64, features      9.000000
importance    0.30866
Name: 1, dtype: float64, features      19.000000
importance    0.443676
Name: 1, dtype: float64, features      19.000000
importance    0.443676
Name: 1, dtype: float64, features      18.000000
importance    0.621775
Name: 4, dtype: float64, features      122.000000
importance    0.995193
Name: 3, dtype: float64, features      10.000000
importance    0.345847
Name: 1, dtype: float64, features      15.000000
importance    0.505214
Name: 1, dtype: float64, features      200.000000
importance    0.356489
Name: 2, dtype: float64, features      85.000000
importance    0.484196
Name: 0, dtype: float64, features      212.0000
importance    0.356
Name: 1, dtype: float64, features      6.000000
importance    0.266495
Name: 1, dtype: float64, features      40.000000
importance    0.332024
Name: 4, dtype: float64, features      37.000000
importance    0.351427
Name: 2, dtype: float64, features      47.000000
importance    0.541774
Name: 0, dtype: float64, features      178.000000
importance    0.196509
Name: 5, dtype: float64, features      96.000000
importance    0.476923
Name: 2, dtype: float64, features      70.000000
importance    0.484083
Name: 3, dtype: float64, features      108.000
importance    0.354
Name: 1, dtype: float64, features      25.0000
importance    0.356
Name: 1, dtype: float64, features      48.000000
importance    0.683515
Name: 1, dtype: float64, features      42.000000
importance    0.288148
Name: 1, dtype: float64, features      162.00000
importance    0.81676
Name: 1, dtype: float64, features      91.000000
importance    0.581767
Name: 1, dtype: float64, features      143.000000
importance    0.510613
Name: 1, dtype: float64, features      143.000000
importance    0.510613
Name: 1, dtype: float64, features      147.000000
importance    0.445206
Name: 1, dtype: float64, features      147.000000
importance    0.389848
Name: 1, dtype: float64, features      11.000000
importance    0.785356
Name: 0, dtype: float64, features      166.000000
importance    0.357875
Name: 0, dtype: float64, features      84.000000
importance    0.374779
Name: 0, dtype: float64, features      218.000000
importance    0.763193
Name: 1, dtype: float64, features      218.000000
importance    0.763193
Name: 1, dtype: float64, features      125.000000
importance    0.484199
Name: 2, dtype: float64, features      196.000000
importance    0.57198
Name: 1, dtype: float64, features      64.00000
importance    0.5807
Name: 0, dtype: float64, features      154.000000
importance    0.805808
Name: 0, dtype: float64, features      137.000000
importance    0.607398
Name: 1, dtype: float64, features      4.000000
importance    0.332152
Name: 2, dtype: float64, features      140.000000
importance    0.608853
Name: 0, dtype: float64, features      185.000000
importance    0.788932
Name: 1, dtype: float64, features      145.000000
importance    0.629176
Name: 0, dtype: float64]
```

```
In [64]: #now make the DataFrame of important features
important_features = pd.DataFrame(important_features)
important_features.reset_index(inplace = True,drop = True)
important_features
```

Out[64]:

	features	importance
0	23.0	0.290988
1	25.0	0.322323
2	104.0	0.522979
3	168.0	0.384278
4	80.0	0.506065
5	180.0	0.437209
6	147.0	0.160047
7	14.0	0.360863
8	82.0	0.186798
9	192.0	0.359525
10	114.0	0.523164
11	109.0	0.270626
12	165.0	0.476715
13	9.0	0.308660
14	19.0	0.443676
15	19.0	0.443676
16	118.0	0.621775
17	122.0	0.995193
18	10.0	0.345847
19	15.0	0.505214
20	200.0	0.356489
21	85.0	0.484196
22	212.0	0.356000
23	6.0	0.266495
24	40.0	0.333204
25	37.0	0.351427
26	47.0	0.541774
27	178.0	0.196509
28	96.0	0.476923
29	70.0	0.484083
30	108.0	0.354000
31	25.0	0.356000
32	48.0	0.683515
33	40.0	0.288148
34	162.0	0.816760
35	91.0	0.591767
36	143.0	0.510613
37	143.0	0.510613
38	147.0	0.445206
39	147.0	0.389848
40	11.0	0.357875
41	166.0	0.785356
42	84.0	0.374779
43	218.0	0.761991
44	218.0	0.761991
45	125.0	0.484199
46	196.0	0.571980
47	64.0	0.586700
48	154.0	0.805808
49	137.0	0.607398
50	4.0	0.332152
51	140.0	0.608851
52	185.0	0.788932
53	145.0	0.629176

We had already to know that we had 54 group then we can see here also 0 to 53 that means 54.

```
In [69]: #feature are to be consider means important features : features selected from the important_feature
features_to_consider = set(important_features['features'])

#feature which are to be discarded it means the correlated features and feature that was accepted due to high im
Features_discarded = set(correlated_features) - set(features_to_consider)

#make list of features which are discarded.
features_discarded = list(features_discarded)
```

Now we have to take a new datasets because we had discared correlated features and keeping the atleast the one features from one group.

We had taken the new datasets in that we pass the x_train_unique,x_test_unique data with drop of the feature which are discarded.

```
In [71]: x_train_grouped_uncorr = x_train_unique.drop(labels= features_discarded,axis =1)
x_train_grouped_uncorr.shape

Out[71]: (1600, 110)
```

```
In [75]: x_test_grouped_uncorr = x_test_unique.drop(labels= features_discarded,axis=1)
x_test_grouped_uncorr.shape

Out[75]: (400, 110)
```

We had taken the imporant features from each group then we use these features to get the accuracy.

```
In [76]: %time
randomforest(x_train_grouped_uncorr,x_test_grouped_uncorr,y_train,y_test)

Accuracy : 0.955
Wall time: 9.79 s

Lets check the same accuracy and computing time with the original dataset for the comparison.
```

```
In [76]: %time
randomforest(x_train,x_test,y_train,y_test)

Accuracy : 0.955
Wall time: 10.7 s
```

```
In [78]: %time
randomforest(x_train_uncorr,x_test_uncorr,y_train,y_test)

Accuracy : 0.9575
Wall time: 9.68 s

Time GAIN is 9% Decrease.
```

```
In [79]: (10.7-9.79)*100/9.79

Out[79]: 9.295199182839635
```

```
In [ ]:

In [ ]:

In [ ]:
```