





```
[23]: def get_redundant_cols(df):
    """Get diagonal and lower triangular pairs of correlation matrix"""
    pairs_to_drop = set()
    cols = df.columns
    for i in range(0, df.shape[1]):
        for j in range(0, i+1):
            pairs_to_drop.add((cols[i], cols[j]))
    return pairs_to_drop

def get_top_abs_correlations(df, n=5):
    au_corr = df.corr().abs().unstack()
    labels_to_drop = get_redundant_pairs(df)
    au_corr = au_corr.drop(labels=labels_to_drop).sort_values(ascending=False)
    return au_corr[0:n]

print("Top Absolute Correlations")
print(get_top_abs_correlations(training_data, 30))

Top Absolute Correlations
E  R    0.99715
K  N    0.99218
A  C    0.99203
H  K    0.99081
A  F    0.990704
F  J    0.989863
E  K    0.989220
N  O    0.988922
A  H    0.988808
C  G    0.988346
K  O    0.982992
H  N    0.982406
E  N    0.979926
A  G    0.972212
C  E    0.971813
I  O    0.970961
F  G    0.969058
A  K    0.968928
C  R    0.968355
E  M    0.964798
H  M    0.964671
O  O    0.962886
A  M    0.958953
E  O    0.958899
K  W    0.956621
A  N    0.953094
D  L    0.949514
G  J    0.948416
M  N    0.947426
C  F    0.943501
dtype: float64
```

After observing the pair of features and correlation heatmap, I decided to follow the approach where if one feature is correlated with other features by 95% or above then we retain only one feature among them, and the rest of the features will be dropped.

With above mentioned approach I divided the features into 3 groups

- A,C,E,F,G,H,I,K,M,N,O
- D,L
- B

I am selecting E from first group and D from second group.

I decided to keep (B,D,E) features (from each group I am taking one feature) and finding the accuracy of the model.

```
In [24]: to_keep = ['B','D','E']

pipeline_keep=PMMLPipeline(((('mapper',DataFrameMapper([(X_train[to_keep].columns.values,
                                                             [StandardScaler()])))],
                              ('pca',PCA(n_components=2)),
                              ('classifier',RandomForestClassifier(max_depth=3)))

pipeline_keep.fit(X_train,y_train)
results=pipeline_keep.predict(X_test)
actual=np.concatenate(y_test.values)
print("Accuracy:", metrics.accuracy_score(actual,results))

Accuracy: 0.5007416666666666

I decided to drop B because B have a outliers (Figure - 1)
```

```
In [25]: to_keep = ['D','E']

pipeline_keep=PMMLPipeline(((('mapper',DataFrameMapper([(X_train[to_keep].columns.values,
                                                             [StandardScaler()])))],
                              ('pca',PCA(n_components=2)),
                              ('classifier',RandomForestClassifier(max_depth=3)))

pipeline_keep.fit(X_train,y_train)
results=pipeline_keep.predict(X_test)
actual=np.concatenate(y_test.values)
print("Accuracy:", metrics.accuracy_score(actual,results))

Accuracy: 0.5007458333333333

Still we are getting the accuracy around 50%. So, manual feature selection approach is not working.
```

Now I am going to apply some feature selection techniques which are explain in the below link and try to improve my accuracy. [https://scikit-learn.org/stable/modules/feature\\_selection.html](https://scikit-learn.org/stable/modules/feature_selection.html)

## L1-based feature selection

```
In [26]: from sklearn.svm import LinearSVC

pipeline_l1=PMMLPipeline(((('mapper',DataFrameMapper([(X_train.columns.values,
                                                         [StandardScaler()])))],
                              ('pca',PCA(n_components=15)),
                              ('feature_selection',SelectFromModel(LinearSVC(C=0.01,penalty="l1",dual=False)
                              ))),
                              ('classifier',RandomForestClassifier(max_depth=3)))

pipeline_l1.fit(X_train,y_train)
results=pipeline_l1.predict(X_test)
actual=np.concatenate(y_test.values)
print("Accuracy:", metrics.accuracy_score(actual,results))

Accuracy: 0.5007416666666666

Univariate feature selection - SelectKBest
```

```
In [27]: from sklearn.feature_selection import SelectKBest

pipeline_kBest=PMMLPipeline(((('mapper',DataFrameMapper([(X_train.columns.values,
                                                             [StandardScaler()])))],
                              ('pca',PCA(n_components=15)),
                              ('feature_selection',SelectKBest(k=3)),
                              ('classifier',RandomForestClassifier(max_depth=3)))

pipeline_kBest.fit(X_train,y_train)
results=pipeline_kBest.predict(X_test)
actual=np.concatenate(y_test.values)
print("Accuracy:", metrics.accuracy_score(actual,results))

Accuracy: 0.5007416666666666
```

## Suitable classifier

```
In [28]: classifier = [
    LogisticRegression(),
    DecisionTreeClassifier(max_depth = 3),
    RandomForestClassifier(max_depth = 3),
    KNeighborsClassifier(),
    GaussianNB(),
    SGDClassifier(alpha=0.001, max_iter=100)]

for cls in classifier:
    print(cls)
    pipeline = PMMLPipeline([
        ('mapper',
         DataFrameMapper([
             (X_train.columns.values,StandardScaler())])),
        ('pca',
         PCA(n_components=15)),
        ('selector',
         SelectKBest(k=3)),
        ('classifier',cls)
    ])

    pipeline.fit(X_train,y_train)
    results=pipeline.predict(X_test)
    actual=np.concatenate(y_test.values)
    print("Accuracy:", metrics.accuracy_score(actual,results))
    print('\n')

LogisticRegression()
Accuracy: 0.5007416666666666

DecisionTreeClassifier(max_depth=3)
Accuracy: 0.5007458333333333

RandomForestClassifier(max_depth=3)
Accuracy: 0.5007416666666666

KNeighborsClassifier()
Accuracy: 0.4139958333333334

GaussianNB()
Accuracy: 0.5007416666666666

SGDClassifier(alpha=0.001, max_iter=100)
Accuracy: 0.5007416666666666
```

All the classifier gives me almost same accuracy except KNearest Neighbour (Which is lowest). So, I am using DecisionTreeClassifier with max\_depth = 3 in my final pipeline.

## Making model on undersampling of the data

Now I am using totally different approach. I am splitting the data into 3 different dataframe.

- First dataframe contains the rows having 1 and 2 class labels that is called df12.
- Second dataframe contains the rows having 1 and 3 class labels that is called df13.
- Third dataframe contains the rows having 2 and 3 class labels that is called df23.

Then I am making three separate models and find the accuracy for each one. Which will give me insight that which label has more weightage on the data.

```
In [29]: df1 = df[df['Class'] == 1]
df2 = df[df['Class'] == 2]
df3 = df[df['Class'] == 3]
df12=pd.concat([df1,df2])
df13=pd.concat([df1,df3])
df23=pd.concat([df2,df3])

• Now I am implementing some functions to split the data, scale the data.
• Implement the function which gives me the optimal component of the dataframe using PCA.
• Implement the function to find accuracy using pipeline
```

```
In [30]: def SplitTheData(df):
    X = df.iloc[:, 0:15]
    y = df.iloc[:, 15:16]
    X_train,X_test,y_train,y_test = train_test_split(X,y_test,test_size=0.2)
    return X_train,X_test,y_train,y_test

def ScaleTheData(df):
    features = df.columns
    scaler = StandardScaler()
    df_scaled = pd.DataFrame(scaler.fit_transform(df.drop('Class',axis=1)))
    return df_scaled

def PCA_RequiredComponents(df_scaled):
    pca = decomposition.PCA(n_components=0.95)
    pca.fit(df_scaled)
    return pca.n_components_

def FindAccuracy(X_train,y_train,X_test,y_test):
    pipeline=PMMLPipeline(((('mapper',DataFrameMapper([(X_train.columns.values,
                                                             [StandardScaler()])))],
                              ('pca',PCA(n_components=2)),
                              ('classifier',DecisionTreeClassifier(max_depth=3)))

    pipeline.fit(X_train,y_train)
    results=pipeline.predict(X_test)
    actual=np.concatenate(y_test.values)
    print("Confusion Matrix")
    print(metrics.confusion_matrix(actual,results))
    print("\n Classification Report")
    print(metrics.classification_report(actual,results))
    return metrics.accuracy_score(actual,results)
```

```
In [31]: X12_train,X12_test,y12_train,y12_test = SplitTheData(df12)
X13_train,X13_test,y13_train,y13_test = SplitTheData(df13)
X23_train,X23_test,y23_train,y23_test = SplitTheData(df23)

print("The optimal number of components to maintain 95% variance for df12: {}".format(PCA_RequiredComp
nent(ScaleTheData(df12))))
print("The optimal number of components to maintain 95% variance for df13: {}".format(PCA_RequiredComp
nent(ScaleTheData(df13))))
print("The optimal number of components to maintain 95% variance for df23: {}".format(PCA_RequiredComp
nent(ScaleTheData(df23))))

print("\n Accuracy of the model with for df12 dataframe where we consider the rows having 1 and 2 class
labels :",FindAccuracy(X12_train,y12_train,X12_test,y12_test))
print("\n Accuracy of the model with for df13 dataframe where we consider the rows having 1 and 3 class
labels :",FindAccuracy(X13_train,y13_train,X13_test,y13_test))
print("\n Accuracy of the model with for df23 dataframe where we consider the rows having 2 and 3 class
labels :",FindAccuracy(X23_train,y23_train,X23_test,y23_test))

The optimal number of components to maintain 95% variance for df12: 2
The optimal number of components to maintain 95% variance for df13: 2
The optimal number of components to maintain 95% variance for df23: 2

Confusion Matrix
[[ 0 39995]
 [ 4 119845]]

classification Report
              precision    recall  f1-score   support

     1         0.00         0.00         0.00      39995
     2         0.75         1.00         0.86     119849

 accuracy         0.37         0.50         0.75     159844
 macro avg        0.37         0.50         0.43     159844
weighted avg        0.56         0.75         0.64     159844

Accuracy of the model with for df12 dataframe where we consider the rows having 1 and 2 class labels
: 0.7497622682115063

Confusion Matrix
[[ 2 39845]
 [ 4 80304]]

classification Report
              precision    recall  f1-score   support

     1         0.33         0.00         0.00      39847
     3         0.67         1.00         0.80      80308

 accuracy         0.50         0.67         0.67     120155
 macro avg        0.50         0.50         0.40     120155
weighted avg        0.56         0.67         0.54     120155

Accuracy of the model with for df13 dataframe where we consider the rows having 1 and 3 class labels
: 0.6683533768881861

Confusion Matrix
[[119764    1]
 [ 80236    1]]

classification Report
              precision    recall  f1-score   support

     2         0.60         1.00         0.75     119765
     3         0.50         0.00         0.00      80237

 accuracy         0.60         0.50         0.55     200002
 macro avg        0.55         0.50         0.37     200002
weighted avg        0.56         0.60         0.45     200002

Accuracy of the model with for df23 dataframe where we consider the rows having 2 and 3 class labels
: 0.558913911809882
```

I got the higher accuracy for the model which we have made from the df12 dataframe.

```
In [32]: print(df1.shape)
print(df2.shape)
print(df3.shape)
print(df12.shape)
print(df13.shape)
print(df23.shape)

(199992, 16)
(59998, 16)
(400780, 16)
(799220, 16)
(600772, 16)
(1000095, 16)
```

As we can see the highest number of the rows have class label 2.

So, Undersampling of the data is random.

```
In [33]: # Undersampling of the data by random selection.
class1=df1
class3=df3.sample(n=199992) # select 199992 rows randomly
class2=df2.sample(n=199992) # select 199992 rows randomly
df_with_sampling=pd.concat([class1,class2,class3])
df_with_sampling
```

```
Out [33]:
```

	A	B	C	D	E	F	G	H	I	J	
11	-23.413125	-11.119531	16.910592	18.915184	-25.170026	-26.504337	-2.371616	-26.557941	-4.756554	20.160979	3.4
12	-29.249384	-13.496606	9.719229	18.708490	-21.304944	-34.440508	-1.884057	-26.843994	-9.007811	22.341808	0.9
15	-39.238896	-16.804404	15.525642	17.059754	-23.125142	-34.300873	6.034805	-25.488902	-1.199718	28.432906	0.9
24	214.501120	-12.386835	220.170600	-11.482754	136.266760	84.174828	113.271746	202.299392	79.514846	142.781926	206.9
28	244.410014	-11.742173	224.945971	-11.031026	131.058794	83.866057	131.314049	192.609849	88.551907	139.349595	198.4
...	...	...	...	...	...	...	...	...	...	...	...
41498	220.252792	-12.445251	228.549679	-15.188049	136.968363	86.931645	127.123562	217.430116	82.055721	146.100794	19.0
79870	-34.976800	-11.454732	15.273580	20.770186	-29.434631	-30.651293	-2.727978	-22.530576	-13.057073	28.944970	-14.0
166548	221.387757	-15.864191	168.715416	-10.862085	137.443604	83.521395	127.981545	213.514370	86.458663	123.687976	213.9
158074	-56.821837	-47.891036	-53.794490	-122.438180	-13.757492	-131.461159	-54.954805	-13.652408	51.232848	-95.968326	37.0
383866	-32.976390	-14.504423	2.872470	18.671767	-23.295589	-30.323475	-4.277570	-25.204396	-6.444951	22.500665	4.5

599976 rows x 16 columns

- For sampling of the data for all the labels, we should use

from imblearn.under\_sampling import NearMiss

smk=NearMiss(X\_res,y\_res = smk\_fit\_resample(df.drop('Class,axis=1),df['Class'])

- I have never implemented above method but according to the below link this is the correct way to do undersampling specially when your data is imbalanced Classification data.

<https://machinelearningmastery.com/undersampling-algorithms-for-imbalanced-classification/>

```
In [34]: print("The optimal number of components to maintain 95% variance in newly created sampled data frame:
({})".format(PCA_RequiredComponents(ScaleTheData(df_with_sampling))))

The optimal number of components to maintain 95% variance in newly created sampled data frame: 2
```

```
In [35]: XN_train,XN_test,yN_train,yN_test = SplitTheData(df_with_sampling)
print("Accuracy of the model which we have made from the dataframe with sampling data :",FindAccuracy(X
N_train,yN_train,XN_test,yN_test))

Confusion Matrix
[[ 361 3216 36383]
 [ 307 3358 36448]
 [ 309 3317 36277]]

classification Report
              precision    recall  f1-score   support

     1         0.37         0.01         0.02      39980
     2         0.34         0.08         0.13     40113
     3         0.33         0.91         0.49      39903

 accuracy         0.33         0.33         0.21     119996
 macro avg        0.35         0.33         0.21     119996
weighted avg        0.35         0.33         0.21     119996

Accuracy of the model which we have made from the dataframe with sampling data : 0.3333111037034566
```

- We are not getting higher accuracy because we are doing undersampling of the data using the random selection.
- According to my point of view and understanding, if we are doing undersampling of the data using NearMiss() and use those train and test data to make a model then we will get more decent model and will get higher accuracy.

## Final Model: ONNX Implementation

```
In [36]: transformer = Pipeline(steps=[
    ('scaler',
     StandardScaler())
])

preprocessor = ColumnTransformer(transformers=[
    ('feature',
     transformer,
     df.columns[0:15])
])

classifier = DecisionTreeClassifier(max_depth = 3)
```

```
In [37]: pipeline = Pipeline([
    ('preprocessor',
     preprocessor),
    ('pca',
     PCA(n_components=3)),
    ('selector',
     SelectKBest(k=2)),
    ('classifier',
     classifier)
])

pipeline.fit(X_train,y_train)
```

```
In [38]: print(metrics.accuracy_score(y_test.values.ravel(),pipeline.predict(X_test)))

0.5007206333333333

Covert pipeline to ONNX file.
```

```
In [39]: input_types = dict([(X, FloatTensorType([None, 1])) for x in X_train.columns.values])

try:
    model_onnx = convert_sklearn(pipeline,
                                'pipeline_onnx',
                                initial_types=input_types.items())
except Exception as e:
    print(e)

with open("f:\pipeline\pipeline.onnx", "wb") as f:
    f.write(model_onnx.SerializeToString())
```

```
In [40]: inputs_onnx = {k: np.array(v).astype(np.float32)[:, np.newaxis] for k, v in X_test.to_dict(orient='list
t').items()}

session_onnx = rt.InferenceSession("f:\pipeline\pipeline.onnx")
predict_onnx = session_onnx.run([None], inputs_onnx)
print("Predict", predict_onnx[0])

predict [[ 2  2  2 ...  2  2  2]
```

```
In [41]: np.unique(np.array(predict_onnx[0]), return_counts=True)

Out [41]: (array([1, 2, 3], dtype=int64), array([ 3, 239988, 9], dtype=int64))
```

```
In [42]: np.unique(predict_onnx[0] == testing_data.iloc[:,15], return_counts=True)

Out [42]: (array([False,  True]), array([119827, 120173], dtype=int64))
```

## Conclusion

- Data after cleaning and EDA : Data is multimodal with few outliers and doesn't follow normal distribution
- Optimal Dimensionality is 2
- Scaling and transformation improve accuracy
- This is a multiclass classification problem and models which are better in multiclassification improve accuracy.
- Results of feature selection at a certain extent are inconclusive.
- Use of Pipeline can save time and streamline process.
- High correlation found between the features.
- We can achieve the higher accuracy with sampling of the data.
- Segregating data on basis of class and using binary classifiers improve accuracy and provide better estimate

## Reference :

- <https://www.dualsight.com/guides/evaluating-a-data-mining-model>
- <https://analyticsvidhya.com/blog/2020/01/build-your-first-machine-learning-pipeline-using-scikit-learn/>
- <https://hazlecast.com/glossary/data-pipeline/>
- <https://www.analyticsvidhya.com/blog/2020/10/feature-selection-techniques-in-machine-learning/>
- [https://scikit-learn.org/stable/modules/feature\\_selection.html](https://scikit-learn.org/stable/modules/feature_selection.html)
- <https://stackoverflow.com/questions/1778394/list-highestcorrelation-pairs-from-a-large-correlation-matrix-in-pandas>

Notebook: Individual Project - Pipeline - ONNX - Example.ipynb

Notebook: Individual Project - Pipeline - PMML - Example.ipynb

Example: Example Project - Spring 2019 - Victoria Belotti.pdf

```
In [ ]:
```