


```
[80]: to_keep = X_train.columns.drop(['E', 'H', 'M', 'O', 'K', 'A', 'N', 'I', 'C', 'F', 'D', 'G']).values
      pipeline5 = PMMLPipeline({
          'mapper',
          DataFrameMapper({
              (X_train.to_keep().columns.values,
               ContinuousDomain()),
              (ContinuousScaler())}),
          ('pca',
           PCA(n_components=3)),
          ('cluster', kMeans(n_clusters=3))
      })
      pipeline5.fit(X_train, y_train)
      results = pipeline5.predict(X_test)
      actual = np.concatenate(y_test.values)
      print("Homogeneity_score=", metrics.homogeneity_score(actual, results))

C:\Users\dokur\Anaconda3\lib\site-packages\sklearn\base.py:197: FutureWarning: From version 0.24, get_params will raise an AttributeError if a parameter cannot be retrieved as an instance attribute. Pre viously it would return None.
  FutureWarning)

Homogeneity_score: 1.0
```

```
In [136]: to_keep = X_train.columns.drop(['E', 'H', 'M', 'O', 'K', 'A', 'N', 'I', 'C', 'F', 'D', 'G']).values
      pipeline6 = PMMLPipeline({
          'mapper',
          DataFrameMapper({
              (X_train.to_keep().columns.values,
               ContinuousDomain()),
              (StandardScaler())}),
          ('cluster', kMeans(n_clusters=3))
      })
      pipeline6.fit(X_train, y_train)
      results = pipeline6.predict(X_test)
      actual = np.concatenate(y_test.values)
      print("Homogeneity_score=", metrics.homogeneity_score(actual, results))

Homogeneity_score: 1.0
```

At this point, I'm not really sure what else I can do to not overfit the model. I know there must be more to the data pre-processing stage that I can do but I'm not so sure I'm capable of completing and understanding it. Since there is nothing more than I can really do, I'll just call it quits for now. Maybe some day I'll be angry enough at myself that I'll proceed to torture myself even more to accomplish my original goal of this project.

Final Model Validation

```
In [129]: df.head()
```

	A	B	C	D	E	F	G	H	I	J	K	L
0	31.628860	-4.925617	-26.828238	-5.772780	20.129709	7.072375	19.335662	-27.283873	9.375416	15.336531	34.086952	-3.845316
1	-24.874820	-12.140354	-0.638120	6.712162	-18.044057	-12.486812	-5.060765	-23.392410	-9.336210	6.583095	-2.308664	-3.626419
2	-25.711282	-16.540789	7.450853	3.964190	-18.273375	-17.082651	-7.082651	-17.443779	-8.103270	8.845905	-1.319906	-1.234682
2	27.837458	-2.436980	29.778957	-3.360897	14.997228	10.254915	10.530714	-17.082219	9.254091	26.860616	29.592860	-2.926678
4	30.520663	-5.017300	24.798154	-4.868743	19.897230	12.570798	13.141632	31.598069	12.120745	16.478925	25.257392	-2.13821

```
In [104]: pipeline6.fit(X, y)

Out[104]: PMMLPipeline(steps=[('mapper', DataFrameMapper(default=False, df_out=False, features=(array(['B', 'J', 'L'], dtype=object), [ContinuousDomain(dtype=None, high_value=None, invalid_value_replacement=None, invalid_value_treatment='return_invalid', low_value=None, missing_value_replacement=None, missing_value_treatment='as_is', missing_values=None, outlier_treatment='as_is', with_data=True, with_statistics=True), StandardScaler(copy=True, with_mean=True, with_std=True)]), input_df=False, sparse=False)), ('cluster', kMeans(algorithm='auto', copy_x=True, init='k-means++', max_iter=300, n_clusters=3, n_init=10, n_jobs=None, precompute_distances='auto', random_state=None, tol=0.0001, verbose=0)))]

In [102]: from sklearn.common.data_types import FloatTensorType
from sklearn import convert_sklearn

input_types = dict(((x, FloatTensorType((None, 1))) for x in df.columns.values))

model_onnx = convert_sklearn(pipeline6, 'pipeline_project_onnx', initial_types=list(input_types.items()))

with open("projectpipeline.onnx", "wb") as f:
    f.write(model_onnx.SerializeToString())

-----
MissingShapeCalculator: Traceback (most recent call last)
<ipython-input-102-6a576fb5d0d39> in <module>
      5 input_types = dict(((x, FloatTensorType((None, 1))) for x in df.columns.values))
----> 7 model_onnx = convert_sklearn(pipeline6, 'pipeline_project_onnx', initial_types=list(input_types.items()))
      8
      9 with open("projectpipeline.onnx", "wb") as f:

~\Anaconda3\lib\site-packages\sklearn\convert.py in convert_sklearn(model, name, initial_types, doc_string, target_opset, custom_conversion_functions, custom_shape_calculators, custom_parsers, options, dtype, intermediate)
    132
    133     # Infer variable shapes
--> 134     topology.compile()
    135
    136     # Convert our Topology object into ONNX. The outcome is an ONNX model.

~\Anaconda3\lib\site-packages\sklearn\common\_topology.py in compile(self)
    842     self._resolve_duplicates()
    843     self._fix_shapes()
--> 844     self._infer_all_types()
    845     self._check_structure()
    846

~\Anaconda3\lib\site-packages\sklearn\common\_topology.py in _infer_all_types(self)
    695     shape_calc(operator)
    696     else:
--> 697         operator.infer_types()
    698
    699     def _resolve_duplicates(self):

~\Anaconda3\lib\site-packages\sklearn\common\_topology.py in infer_types(self)
    211     raise MissingShapeCalculator(
    212         "Unable to find a shape calculator for type '{}'. ".format(
--> 213             type(self.raw_operator)))
    214     try:
    215         shape_calc = _registration.get_shape_calculator(self.type)

MissingShapeCalculator: Unable to find a shape calculator for type '<class 'sklearn_pandas.dataframe_mapper.DataFrameMapper'>'.
```

It usually means the pipeline being converted contains a transformer or a predictor with no corresponding converter implemented in sklearn-onnx. If the converted is implemented in another library (i.e. onmmatcols), you need to register the converted so that it can be used by sklearn-onnx (function update_registered_converter). If the model is not yet covered by sklearn-onnx, you may raise an issue to <https://github.com/onnx/sklearn-onnx/issues> to get the converter implemented or even contribute to the project. If the model is a custom model, a new converter must be implemented. Examples can be found in the gallery.

```
In [105]: # Create final pipeline that doesn't include DataFrameMapper...
finalPipeline = PMMLPipeline({
    ('scaler', StandardScaler()),
    ('pca', PCA(n_components=3)),
    ('cluster', kMeans(n_clusters=3))
})
final_pipeline.fit(df, df['Class'])

C:\Users\dokur\Anaconda3\lib\site-packages\sklearn\base.py:197: FutureWarning: From version 0.24, get_params will raise an AttributeError if a parameter cannot be retrieved as an instance attribute. Pre viously it would return None.
  FutureWarning)

Out[105]: Pipeline(memory=None,
      steps=[('column_selector',
              ColumnSelector(cols=['B', 'J', 'L'], drop_axis=False)),
             ('pca',
              PCA(copy=True, iterated_power='auto', n_components=3,
                  random_state=None, svd_solver='auto', tol=0.0,
                  whiten=False)),
             ('cluster',
              kMeans(algorithm='auto', copy_x=True, init='k-means++',
                     max_iter=300, n_clusters=3, n_init=10, n_jobs=None,
                     precompute_distances='auto', random_state=None,
                     tol=0.0001, verbose=0))],
      verbose=False)
```

```
In [107]: from sklearn.common.data_types import FloatTensorType
from sklearn import convert_sklearn

input_types = dict(((x, FloatTensorType((None, 1))) for x in df.columns.values))

model_onnx = convert_sklearn(finalPipeline, 'pipeline_project_onnx', initial_types=list(input_types.items()))

with open("projectpipeline.onnx", "wb") as f:
    f.write(model_onnx.SerializeToString())

-----
AttributeError: Traceback (most recent call last)
<ipython-input-107-f65d03e78864> in <module>
      4 input_types = dict(((x, FloatTensorType((None, 1))) for x in df.columns.values))
----> 6 model_onnx = convert_sklearn(finalPipeline, 'pipeline_project_onnx', initial_types=list(input_types.items()))
      7
      8 with open("projectpipeline.onnx", "wb") as f:

~\Anaconda3\lib\site-packages\sklearn\convert.py in convert_sklearn(model, name, initial_types, doc_string, target_opset, custom_conversion_functions, custom_shape_calculators, custom_parsers, options, dtype, intermediate)
    132
    133     # Infer variable shapes
--> 134     topology.compile()
    135
    136     # Convert our Topology object into ONNX. The outcome is an ONNX model.

~\Anaconda3\lib\site-packages\sklearn\common\_topology.py in compile(self)
    842     self._resolve_duplicates()
    843     self._fix_shapes()
--> 844     self._infer_all_types()
    845     self._check_structure()
    846

~\Anaconda3\lib\site-packages\sklearn\common\_topology.py in _infer_all_types(self)
    695     shape_calc(operator)
    696     else:
--> 697         operator.infer_types()
    698
    699     def _resolve_duplicates(self):

~\Anaconda3\lib\site-packages\sklearn\common\_topology.py in infer_types(self)
    218         "Unable to find a shape calculator for alias '{}'. "
    219         "and type '{}'. ".format(self.type, type(self.raw_operator)))
--> 220     shape_calc(self)
    221
    222     @property

~\Anaconda3\lib\site-packages\sklearn\shape_calculators\svd.py in calculate_sklearn_truncated_svd_output_shapes(operator)
    33     K = operator.raw_operator.n_components
    34     if operator.type == 'SklearnTruncatedSVD':
--> 35         else operator.raw_operator.n_components'
    36
    37     operator.outputs[0].type.shape = [N, K]

AttributeError: 'PCA' object has no attribute 'n_components'
```

```
In [133]: # Another attempt, but I guess I'll do pca before hand
scaler = StandardScaler()
df_values = df.drop(['E', 'H', 'M', 'O', 'K', 'A', 'N', 'I', 'C', 'F', 'D', 'G', 'Class'], axis=1).val
ues
scaled = scaler.fit_transform(df_values)
df_scaled = pd.DataFrame(scaled, columns = ['B', 'J', 'L'])

pca_ds = pca.fit_transform(df_scaled)
pca_df = pd.DataFrame(data=pca_ds, columns=['PC1', 'PC2', 'PC3'])
pca_df.head()
```

```
Out[133]:
```

	PC1	PC2	PC3
0	-0.482524	0.882670	0.022636
1	-0.554739	0.768167	-0.063692
2	-0.600938	1.401382	0.855717
4	-0.224123	-0.986485	0.421984

```
In [140]: # Create final pipeline that doesn't include DataFrameMapper...
finalPipeline2 = PMMLPipeline({
    ('cluster', kMeans(n_clusters=3))
})
final_pipeline2.fit(pca_df, df['Class'])

C:\Users\dokur\Anaconda3\lib\site-packages\sklearn\base.py:197: FutureWarning: From version 0.24, get_params will raise an AttributeError if a parameter cannot be retrieved as an instance attribute. Pre viously it would return None.
  FutureWarning)
```

```
Out[140]: PMMLPipeline(steps=[('cluster', kMeans(algorithm='auto', copy_x=True, init='k-means++', max_iter=300, n_clusters=3, n_init=10, n_jobs=None, precompute_distances='auto', random_state=None, tol=0.0001, verbose=0)))]

In [142]: from sklearn.common.data_types import FloatTensorType
from sklearn import convert_sklearn

input_types = dict(((x, FloatTensorType((None, 1))) for x in pca_df.columns.values))

model_onnx = convert_sklearn(finalPipeline2, 'pipeline_project_onnx', initial_types=list(input_types.items()))

with open("projectpipeline.onnx", "wb") as f:
    f.write(model_onnx.SerializeToString())

-----
AttributeError: Traceback (most recent call last)
<ipython-input-142-31b4ea20ddfd> in <module>
      4 input_types = dict(((x, FloatTensorType((None, 1))) for x in pca_df.columns.values))
----> 6 model_onnx = convert_sklearn(finalPipeline2, 'pipeline_project_onnx', initial_types=list(input_types.items()))
      7
      8 with open("projectpipeline.onnx", "wb") as f:

~\Anaconda3\lib\site-packages\sklearn\convert.py in convert_sklearn(model, name, initial_types, doc_string, target_opset, custom_conversion_functions, custom_shape_calculators, custom_parsers, options, dtype, intermediate)
    132
    133     # Infer variable shapes
--> 134     topology.compile()
    135
    136     # Convert our Topology object into ONNX. The outcome is an ONNX model.

~\Anaconda3\lib\site-packages\sklearn\common\_topology.py in compile(self)
    842     self._resolve_duplicates()
    843     self._fix_shapes()
--> 844     self._infer_all_types()
    845     self._check_structure()
    846

~\Anaconda3\lib\site-packages\sklearn\common\_topology.py in _infer_all_types(self)
    695     shape_calc(operator)
    696     else:
--> 697         operator.infer_types()
    698
    699     def _resolve_duplicates(self):

~\Anaconda3\lib\site-packages\sklearn\common\_topology.py in infer_types(self)
    218         "Unable to find a shape calculator for alias '{}'. "
    219         "and type '{}'. ".format(self.type, type(self.raw_operator)))
--> 220     shape_calc(self)
    221
    222     @property

~\Anaconda3\lib\site-packages\sklearn\shape_calculators\k_means.py in calculate_sklearn_kmeans_output_shapes(operator)
    18     yood_output_types=[Int64TensorType, FloatTensorType, DoubleTensorType])
    19     if len(operator.inputs) != 1:
--> 20         raise RuntimeError("Only one input vector is allowed for KMeans.")
    21     if len(operator.outputs) != 2:
    22         raise RuntimeError("Two outputs are expected for KMeans.")

RuntimeError: Only one input vector is allowed for KMeans.
```

Conclusion

I tried to create an onnx file, but it kept popping error after error, so the best I can do is turn in this .ipynb. Anyway, for this project, I thought I had a solid plan for the data pre-processing section, but it seems like what I did was not enough. Even after removing over 12 features, the I kept overfitting the model. Maybe if I had more knowledge in feature selection techniques I would have been able to properly complete my objective. I thought for sure I was on the right track with clustering since the data looked as if it was heading in that direction, but perhaps I was wrong. If I ever turn to this data file, I'll for sure try using other methods, like decision tree. I guess the original sample with the decision tree I meant to guide us in that direction, but like I said, I really believed the data was pushing for clusters.

I have to be admit though, I had just a tiny bit of fun. At first, it was very confusing and scary since I didn't know where to start or what the examples files were even doing. But after doing some bit of research and looking at other examples, I really felt as if I was able to understand what was going on. Or at least to some degree. Writing out my thought process also helped.

References

<https://towardsdatascience.com/pca-using-python-sikit-learn-6b53f9983e60> <https://towardsdatascience.com/dimensionality-reduction-the-best-feature-selection-https://stackoverflow.com/questions/1778394/the-highest-correlation-pairs-from-a-large-correlation-matrix-in-pandas> <https://blog.cambridgespark.com/how-to-determine-the-optimal-number-of-clusters-for-k-means-clustering-1427070046/>

```
In [ ]:
```