

## Appendix C: TimeSeriesKMean on single sensor data

```
In [1]: 1 import pandas as pd
        2 import numpy as np
        3 from sklearn.cluster import KMeans
```

```
In [2]: 1 import matplotlib
        2 import matplotlib.pyplot as plt
```

```
In [4]: 1 from tslearn.clustering import TimeSeriesKMeans
        2 #tslearn requires numpy 1.22
        3 #pip install --upgrade threadpoolctl --user
```

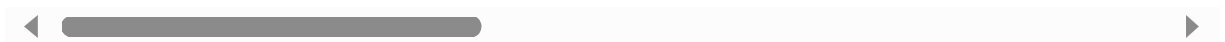
```
In [5]: 1 # Read the CSV data into a Pandas DataFrame, please update the path accora
        2 df = pd.read_csv('D:\\Download\\Tool_Sensor_Data.csv')
        3 print(df.shape)
        4 df.head()
```

(14844, 90)

Out[5]:

	TimeStamp	ToolName	TOOL_ID	Run	RunStartTime	DATA_QUALITY	EQPType	HasComme
0	21/3/2023 19:14	A	A	62301	12:14:19 AM	63.49	A	
1	21/3/2023 19:14	A	A	62301	12:14:19 AM	63.49	A	
2	21/3/2023 19:14	A	A	62301	12:14:19 AM	63.49	A	
3	21/3/2023 19:14	A	A	62301	12:14:19 AM	63.49	A	
4	21/3/2023 19:14	A	A	62301	12:14:19 AM	63.49	A	

5 rows × 90 columns



In [6]:

```
1 # ALL data cleaning in one Cell
2 # Remove Duplicated Rows, the timeline helps to ensure duplicated rows wou
3 print('Before = ', df.shape)
4 df = df.drop_duplicates()
5 print('After = ', df.shape)
6 # The TimeStamp indicates it is time series data, but since it is not a sa
7 # The more important info is the Run number for identifying which Wafer Ru
8 if 'TimeStamp' in df.columns:
9     df = df.drop('TimeStamp', axis=1)
10 if 'RunStartTime' in df.columns:
11     df = df.drop('RunStartTime', axis=1)
12 # Find/Check for any other non-integer columns
13 for column in df.columns:
14     if df[column].dtypes == 'object':
15         print('column name =', column, 'type =', df[column].dtypes, 'conter
16 # Verify manually if those columns are not useful (no value for ML), drop
17 for column in df.columns:
18     if df[column].dtypes == 'object':
19         print(column, df[column].dtypes, '= DROP')
20         df = df.drop(column, axis=1)
21 # Calculate how many cells inside a row is filled with data
22 non_empty_columns_per_row = df.count(axis=1)
23 non_empty_columns_per_row
24 # Pick a cut off threshold, in this case, based on the chart above, we can
25 print('Total length = ', len(df.columns))
26 threshold = len(df.columns)/2
27 print('Threshold =', threshold)
28 # Remove empty rows below the threshold
29 print('Before remove=', df.shape)
30 df = df.dropna(thresh=threshold)
31 print('After remove=', df.shape)
32 # Fill in empty cells with Zero? Number zero allows calculation of standar
33 df = df.fillna(0)
34 df.head()
35 # Remove Useless Columns by using standard deviation check, likely those c
36 column_stddev = {}
37 for column in df.columns:
38     std_deviation = df[column].std()
39     if std_deviation == 0:
40         print(column, std_deviation, "= DROP")
41         df = df.drop(column, axis=1)
42     else:
43         column_stddev[column]=std_deviation
44 # Save a copy of the clean dataframe
45 df_clean=df.copy()
46 df.to_csv('d:\\download\\1.3_Cleaning.csv', index=False)
```

Before = (14844, 90)  
After = (14113, 90)  
column name = ToolName type = object content = A  
column name = TOOL\_ID type = object content = A  
column name = EQPType type = object content = A  
column name = LOT\_ID type = object content = A  
column name = LogicalRecipeID type = object content = A  
column name = LotPurposeType type = object content = Process Lot  
column name = LotType type = object content = Production  
column name = MachineRecipeID type = object content = A  
column name = PhysicalRecipeID type = object content = A  
column name = PortID type = object content = A  
column name = ProcessOpNum type = object content = A  
column name = ProductGrpID type = object content = A  
column name = ProductID type = object content = A  
column name = RECIPE\_ID type = object content = A  
column name = RouteID type = object content = A  
column name = Technology type = object content = A  
column name = EventType type = object content = StartOfRun  
column name = EventName type = object content = WaferStart  
column name = EventId type = object content = WaferStart  
ToolName object = DROP  
TOOL\_ID object = DROP  
EQPType object = DROP  
LOT\_ID object = DROP  
LogicalRecipeID object = DROP  
LotPurposeType object = DROP  
LotType object = DROP  
MachineRecipeID object = DROP  
PhysicalRecipeID object = DROP  
PortID object = DROP  
ProcessOpNum object = DROP  
ProductGrpID object = DROP  
ProductID object = DROP  
RECIPE\_ID object = DROP  
RouteID object = DROP  
Technology object = DROP  
EventType object = DROP  
EventName object = DROP  
EventId object = DROP  
Total length = 69  
Threshold = 34.5  
Before remove= (14113, 69)  
After remove= (13657, 69)  
HasComments 0.0 = DROP  
ReticleID 0.0 = DROP  
GPTmghByqMSY 0.0 = DROP  
pZXcGFNpzPf 0.0 = DROP  
EHVtYhnRGb 0.0 = DROP  
XSOeMfJAB 0.0 = DROP  
SYklrMAXe 0.0 = DROP  
jQVGDTFl 0.0 = DROP  
YDlkDLfFEEi 0.0 = DROP  
CYiycrAoYbg 0.0 = DROP  
oUWQRhjudAd 0.0 = DROP  
WTcLPnkDtRwBuCou 0.0 = DROP  
EventSource 0.0 = DROP

EventDescription 0.0 = DROP  
AlarmCode 0.0 = DROP  
AlarmStatus 0.0 = DROP  
VcnDxeRWfK 0.0 = DROP  
HxXxxvS 0.0 = DROP  
LqnSjZtcJs 0.0 = DROP  
vmZiUljnYP 0.0 = DROP  
yqScxEFPLde 0.0 = DROP  
jrsnDLYHnMHD 0.0 = DROP  
RunTag 0.0 = DROP  
Unnamed: 85 0.0 = DROP  
Unnamed: 86 0.0 = DROP  
Unnamed: 87 0.0 = DROP  
Unnamed: 88 0.0 = DROP  
Unnamed: 89 0.0 = DROP

```

In [7]: 1 # Get unique run numbers
2 runs = df['Run'].unique()
3 print('Total Unique Runs=',len(runs))
4 #runs
5
6 df_run = pd.DataFrame()
7 i=0
8 # Calculate each Run's Length and record inside df_run
9 for run in runs:
10     count = (df['Run'] == run).sum()
11     df_run.loc[i, 'Run'] = run
12     df_run.loc[i, 'Count'] = count
13     i=i+1
14 df_run.head()
15 # SwpYipezsdueC - Remove all columns except SwpYipezsdueC and Run column
16 df_check=df
17 check_column = 'SwpYipezsdueC'
18 for column in df.columns:
19     if column != check_column and column != 'Run':
20         df_check=df_check.drop(column, axis=1)
21 print(df.shape)
22 # Split data into dicts by run number
23 data_by_run = {run: df[df['Run'] == run] for run in runs}
24 df.head()
25 # SwpYipezsdueC - Get the max length for resampling
26 i=0
27 mydict={}
28 for run, df1 in data_by_run.items():
29     i = i + 1
30     mylist = df1[check_column].tolist()
31     mydict[run]=len(mylist)
32 max_key = max(mydict, key=lambda key: mydict[key])
33 max_value = mydict[max_key]
34 print('Maximum Run record length =', max_value)
35 # SwpYipezsdueC - Convert column to rows and Resample the time series and
36 new_df = pd.DataFrame()
37 i=0
38 for run, df1 in data_by_run.items():
39     i = i + 1
40     mylist = df1[check_column].tolist()
41     # Resample mylist using linear interpolation to match the length of th
42     if len(mylist) < max_value:
43         x1 = np.arange(len(mylist))
44         x2 = np.linspace(0, len(mylist) - 1, max_value)
45         mylist = np.interp(x2, x1, mylist)
46         mylist = np.insert(mylist, 0, run)
47         new_row = pd.DataFrame([mylist])
48         new_df = pd.concat([new_df, new_row], ignore_index=True)
49 new_df.set_index(0, inplace=True)
50 print(new_df.shape)
51 new_df.head()

```

```

Total Unique Runs= 228
(13657, 41)
Maximum Run record length = 87
(228, 87)

```

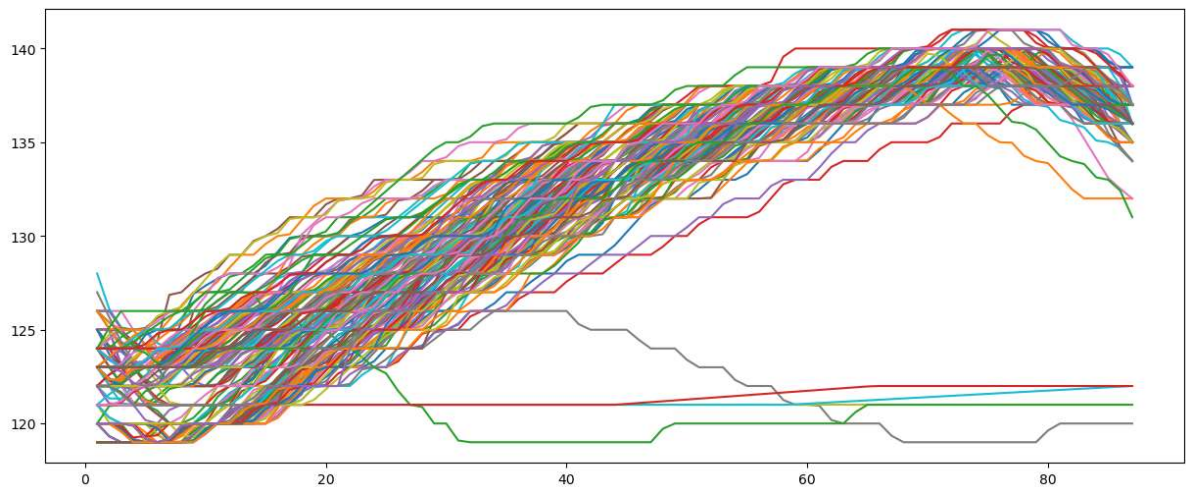
Out[7]:

	1	2	3	4	5	6	7	8	
0									
62301.0	125.0	123.697674	123.0	123.00000	123.000000	123.000000	123.000000	123.558140	12
62302.0	120.0	120.000000	120.0	120.00000	120.000000	119.395349	119.000000	119.046512	11
62303.0	121.0	121.000000	121.0	121.00000	121.000000	121.197674	121.837209	122.000000	12
62304.0	123.0	123.000000	123.0	123.00000	123.000000	123.000000	123.488372	124.000000	12
62305.0	120.0	120.000000	120.0	119.94186	119.255814	119.000000	119.116279	119.802326	12

5 rows × 87 columns



```
In [8]: 1 plt.figure(figsize=(15, 6))
2 for index, row in new_df.iterrows():
3     plt.plot(row, label=f'Line {index}') # Plot each row with a label
4 plt.show()
```



**Based on the pattern above, there are 4 lines which is out of normal pattern.**

## KMean Clustering

```
In [9]: 1 import numpy
2 import matplotlib.pyplot as plt
3
4 from tslearn.clustering import KShape
5 from tslearn.datasets import CachedDatasets
6 from tslearn.preprocessing import TimeSeriesScalerMeanVariance
7
8 import pandas as pd
9 from sklearn.model_selection import train_test_split
```

```
In [10]: 1 # Load the dataset
2 #X_train, y_train, X_test, y_test = CachedDatasets().load_dataset("Trace")
3
4 # Split the data into features (X) and Labels (y)
5 #X = df.drop('y_column_name', axis=1) # Replace 'y_column_name' with the
6 #y = df['y_column_name']
7 #y = new_df.iloc[0]
8 #y = y.to_numpy().flatten()
9 #print(y)
10 #X = list(range(1, len(y) + 1))
11 #X = np.arange(1, len(y) + 1)
12 #X
13 # above wrong.
14 y = new_df.index
15 y = y.astype(int)
16 y
17
18 #X = np.arange(1, len(y) + 1)
19 index_variable = new_df.index
20 index_variable = index_variable.astype(int)
21 X_df = new_df.reset_index(drop=True)
22 X = X_df.values
23 y = np.arange(1, len(X) + 1)
24 #print(y)
25 #X
26 print(X.shape)
27 print(y.shape)
```

(228, 87)

(228,)

```
In [14]: 1 # Split the data into training and testing sets
2 #X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,
3 #print(X_train.shape)
4 #print(X_train)
5 #print(y_train.shape)
6 #print(y_train)
```

```
In [15]: 1 #X_train = X_train[y_train < 4]
2 #print(X_train.shape)
3 #X_train
```

```
In [16]: 1 #X_train = X_train[:50]
         2 #X_train
```

```
In [18]: 1 #numpy.random.shuffle(X_train)
         2 #print(X_train.shape)
         3 #print(y_train.shape)
```

```
In [19]: 1 # For this method to operate properly, prior scaling is required
         2 X_train = X
         3 X_train = TimeSeriesScalerMeanVariance().fit_transform(X_train)
         4 sz = X_train.shape[1]
```



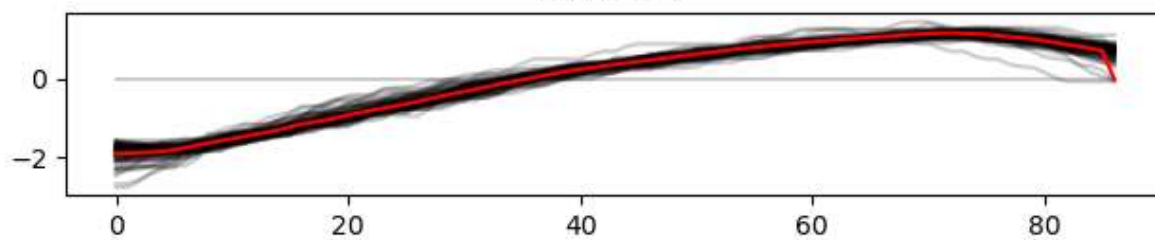
## KShape Clustering

```
In [21]: 1 #X_train = X_train[y_train < 4]
2 #X_train = X_train[:50]
3 #numpy.random.shuffle(X_train)
4 # For this method to operate properly, prior scaling is required
5 #X_train = TimeSeriesScalerMeanVariance().fit_transform(X_train)
6 #sz = X_train.shape[1]
7
8 # kShape clustering
9 seed = 0
10 numpy.random.seed(seed)
11 ks = KShape(n_clusters=3, verbose=True, random_state=seed)
12 y_pred = ks.fit_predict(X_train)
13 print(y_pred)
14
15 plt.figure()
16 for yi in range(3):
17     plt.subplot(3, 1, 1 + yi)
18     for xx in X_train[y_pred == yi]:
19         plt.plot(xx.ravel(), "k-", alpha=.2)
20     plt.plot(ks.cluster_centers_[yi].ravel(), "r-")
21     #plt.xlim(0, sz)
22     #plt.ylim(-4, 4)
23     plt.title("Cluster %d" % (yi + 1))
24
25 plt.tight_layout()
26 plt.show()
```

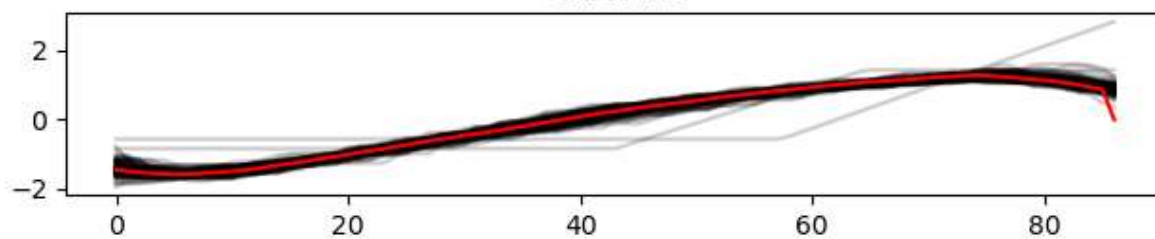
0.005 --> 0.005 -->

```
[1 0 1 1 1 1 1 1 1 0 0 1 1 0 1 0 1 1 1 0 1 1 0 1 0 0 1 1 1 1 1 1 1 1 1
 1 1 0 1 1 0 1 1 0 1 1 0 1 1 0 1 1 0 1 1 0 1 1 0 1 1 1 1 1 0 1 1 0 1
 1 1 1 1 1 1 1 1 1 1 1 1 1 0 1 1 0 1 1 1 0 1 0 1 1 1 0 0 0 0 1 1 1 1 1 0
 1 1 1 1 0 0 2 0 1 0 1 0 0 0 0 0 0 0 0 0 0 0 2 1 0 0 0 0 0 1 1 1 0 1 1 0 0 1 0
 0 1 0 0 1 1 1 1 1 1 0 1 1 1 1 1 0 1 0 1 1 0 0 1 0 1 1 0 0 1 0 0 1 1 0 1 0
 0 1 0 0 1 1 0 0 1 0 0 1 1 1 1 1 1 0 1 0 1 1 1 0 1 1 0 0 1 1 0 0 0 1 1 0 1
 0 0 1 1 1 1]
```

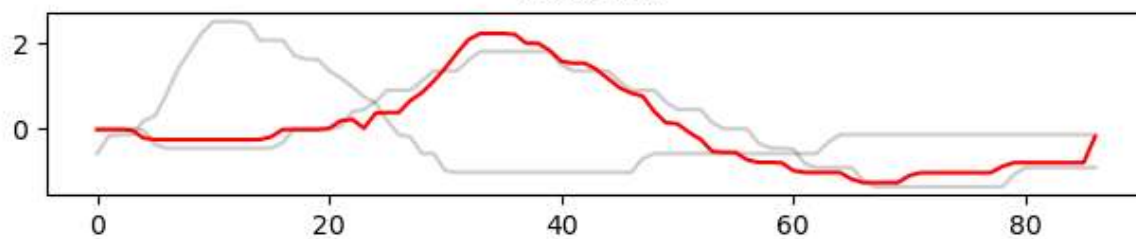
Cluster 1



Cluster 2



Cluster 3



In [22]:

```
1 print('X_train')
2 print(X_train)
3 print('y_pred')
4 print(y_pred)
```

```

X_train
[[-1.34464998]
 [-1.58249252]
 [-1.70990816]
 ...
 [ 0.90211251]
 [ 0.84689907]
 [ 0.84689907]]

[[-1.61992522]
 [-1.61992522]
 [-1.61992522]
 ...
 [ 0.89738317]
 [ 0.8193271 ]
 [ 0.61768224]]

[[-1.79545174]
 [-1.79545174]
 [-1.79545174]
 ...
 [ 0.89863921]
 [ 0.80782715]
 [ 0.80782715]]

...

[[-1.55746254]
 [-1.55746254]
 [-1.55746254]
 ...
 [ 0.98059225]
 [ 0.91102763]
 [ 0.79930262]]

[[-1.59519967]
 [-1.59519967]
 [-1.59519967]
 ...
 [ 1.02032427]
 [ 0.90306542]
 [ 0.78580656]]

[[-1.38152918]
 [-1.38152918]
 [-1.452847  ]
 ...
 [ 1.20092642]
 [ 1.15588359]
 [ 1.03952294]]]

y_pred
[1 0 1 1 1 1 1 1 1 0 0 1 1 0 1 0 1 1 1 0 1 1 0 1 0 0 1 1 1 1 1 1 1 1 1 1 1
 1 1 0 1 1 0 1 1 0 1 1 0 1 1 0 1 1 0 1 1 0 1 1 0 1 1 1 1 1 1 0 1 1 0 1
 1 1 1 1 1 1 1 1 1 1 1 1 1 0 1 1 0 1 1 1 0 1 0 1 1 1 0 0 0 0 1 1 1 1 1 0
 1 1 1 1 0 0 2 0 1 0 1 0 0 0 0 0 0 0 0 0 0 2 1 0 0 0 0 1 1 1 0 1 1 0 0 1 0
 0 1 0 0 1 1 1 1 1 0 1 1 1 1 1 0 1 0 1 1 0 0 1 0 1 1 0 0 1 0 0 1 0 0 1 1 0 1 0

```

```
0 1 0 0 1 1 0 0 1 0 0 1 1 1 1 1 0 1 0 1 1 1 0 1 1 0 0 1 1 0 0 0 1 1 0 1
0 0 1 1 1 1]
```

```
In [23]: 1 X_train_df = pd.DataFrame(X_train[:, 0, 0]) # Assuming X_train is a 3D ar
2 y_pred_df = pd.DataFrame(y_pred, columns=['y_pred'])
3 #merged_df = pd.concat([X_train_df, y_pred_df], axis=1)
4 merged_df = pd.concat([X_df, y_pred_df], axis=1)
5 merged_df = merged_df.set_index(index_variable)
6 merged_df
```

```
Out[23]:
```

	1	2	3	4	5	6	7	8
0								
62301	125.0	123.697674	123.00000	123.000000	123.000000	123.000000	123.000000	123.558140
62302	120.0	120.000000	120.00000	120.000000	120.000000	119.395349	119.000000	119.046512
62303	121.0	121.000000	121.00000	121.000000	121.000000	121.197674	121.837209	122.000000
62304	123.0	123.000000	123.00000	123.000000	123.000000	123.000000	123.488372	124.000000
62305	120.0	120.000000	120.00000	119.941860	119.255814	119.000000	119.116279	119.802326
...	...	...	...	...	...	...	...	...
62596	124.0	124.000000	124.00000	124.000000	124.232558	124.790698	125.000000	125.000000
62597	120.0	120.000000	120.00000	119.906977	119.209302	119.000000	119.186047	119.883721
62598	123.0	123.000000	123.00000	123.000000	123.000000	123.000000	123.000000	123.000000
62599	121.0	121.000000	121.00000	121.000000	121.000000	121.000000	121.000000	121.000000
62600	122.0	122.000000	121.55814	121.000000	121.000000	121.000000	121.000000	121.000000

228 rows × 88 columns



```
In [24]: 1 df=merged_df
2 print('0',(df['y_pred'] == 0).sum())
3 print('1',(df['y_pred'] == 1).sum())
4 print('2',(df['y_pred'] == 2).sum())
5 print('3',(df['y_pred'] == 3).sum())
```

```
0 84
1 142
2 2
3 0
```

```
In [25]: 1 df.to_csv('d:\\download\\timeserieskmean.csv', index=False)
```

**Conclusion: KMean KShape clustering does provide some accurate clustering but missed out a few abnormal patterns. Increasing the number of cluster up to 6 can help improve the accuracy.**

---

---

## **TimeSeriesKMeans - DTW**

```
In [26]: 1 index_variable = new_df.index
          2 index_variable = index_variable.astype(int)
          3 X_df = new_df.reset_index(drop=True)
          4 X = X_df.values
          5 y = np.arange(1, len(X) + 1)
          6 print(X.shape)
```

(228, 87)

```
In [27]: 1 # For this method to operate properly, prior scaling is required
2 X_train = X
3 print(X_train.shape)
4
5 from tslearn.preprocessing import TimeSeriesResampler
6 # Don't scale it because the size of the pattern is important
7 #X_train = TimeSeriesScalerMeanVariance().fit_transform(X_train)
8 #X_train = TimeSeriesResampler().fit_transform(X_train)
9 X_train_scale = X_train.copy()
10 sz = X_train.shape[1]
11 print(sz)
12 X_train = TimeSeriesResampler(sz=sz).fit_transform(X_train)
13 print(X_train.shape)
14 X_train
```

(228, 87)

87

(228, 87, 1)

```
Out[27]: array([[125.        ],
                [123.69767442],
                [123.        ],
                ...,
                [137.30232558],
                [137.        ],
                [137.        ]],

               [[120.        ],
                [120.        ],
                [120.        ],
                ...,
                [138.        ],
                [137.44186047],
                [136.        ]],

               [[121.        ],
                [121.        ],
                [121.        ],
                ...,
                [137.55813953],
                [137.        ],
                [137.        ]],

               ...,

               [[123.        ],
                [123.        ],
                [123.        ],
                ...,
                [137.        ],
                [136.61627907],
                [136.        ]],

               [[121.        ],
                [121.        ],
                [121.        ],
                ...,
                [139.6744186 ],
                [138.8372093 ],
                [138.        ]],

               [[122.        ],
                [122.        ],
                [121.55813953],
                ...,
                [138.        ],
                [137.72093023],
                [137.        ]]])
```



```
In [28]: 1 n_cluster=6
2 model = TimeSeriesKMeans(n_clusters=n_cluster, metric="dtw", max_iter=10,
3 model.fit(X_train)
4 y_pred = model.predict(X_train)
5 print(y_pred)
```

```
[4 2 0 4 2 0 0 2 0 0 4 2 0 0 2 0 4 2 0 4 0 0 0 0 0 0 0 3 3 0 3 3 0 3 0 0 3
4 0 2 3 3 2 3 3 2 0 3 2 0 4 0 0 3 2 0 4 2 0 3 2 3 4 2 3 0 0 3 3 2 3 3 2 0
3 2 0 0 2 0 4 2 0 0 2 0 3 2 0 0 2 0 0 2 0 4 2 0 0 2 0 0 2 0 4 2 0 4 2 0 4
0 4 4 0 3 0 5 1 1 0 0 4 0 3 0 0 0 4 0 0 4 5 1 2 0 3 2 0 4 2 3 4 2 0 4 2 0
4 2 0 4 2 0 4 2 0 4 2 2 2 4 2 0 4 2 3 3 2 0 4 2 0 2 2 0 4 2 4 4 2 4 4 2 0
0 2 0 0 2 0 4 2 0 3 2 0 4 2 0 0 2 4 4 2 0 4 2 0 0 2 0 4 2 0 2 0 0 2 0 0 2
3 4 2 0 2 0]
```

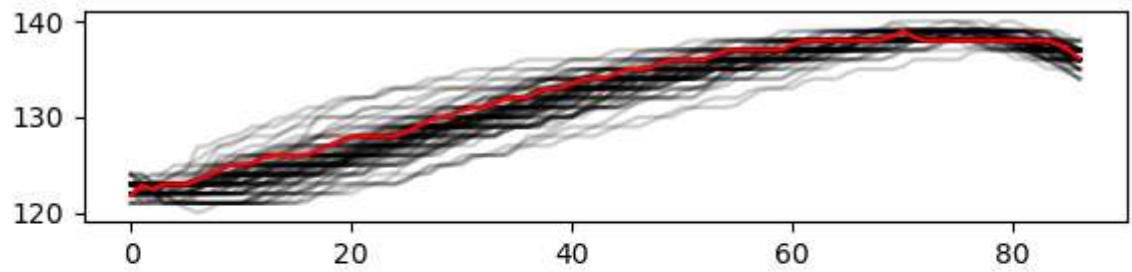
```
In [29]: 1 #model.cluster_centers_[0]
```

In [30]:

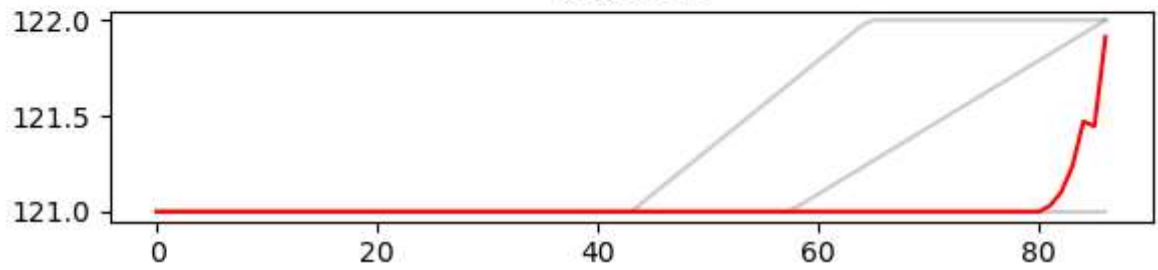
```
1 #plt.figure()
2 plt.figure(figsize=(6, 10))
3 for yi in range(n_cluster):
4     plt.subplot(n_cluster, 1, 1 + yi)
5     for xx in X_train[y_pred == yi]:
6         plt.plot(xx.ravel(), "k-", alpha=.2)
7     plt.plot(model.cluster_centers_[yi].ravel(), "r-")
8     #plt.xlim(0, sz)
9     #plt.ylim(-4, 4)
10    plt.title("Cluster %d" % (yi + 1))
11
12 plt.tight_layout()
13 plt.show()
```



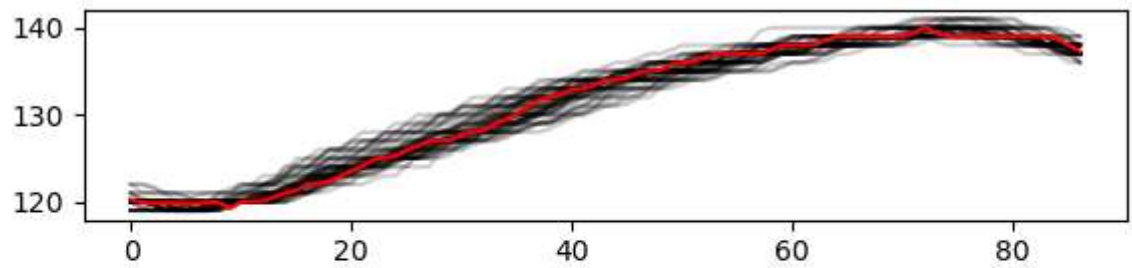
Cluster 1



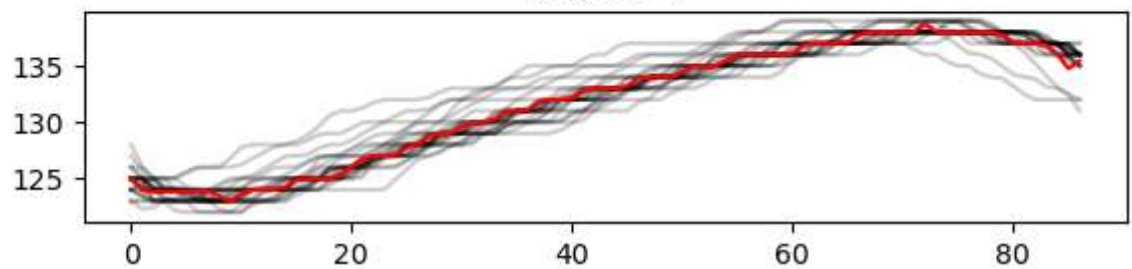
Cluster 2



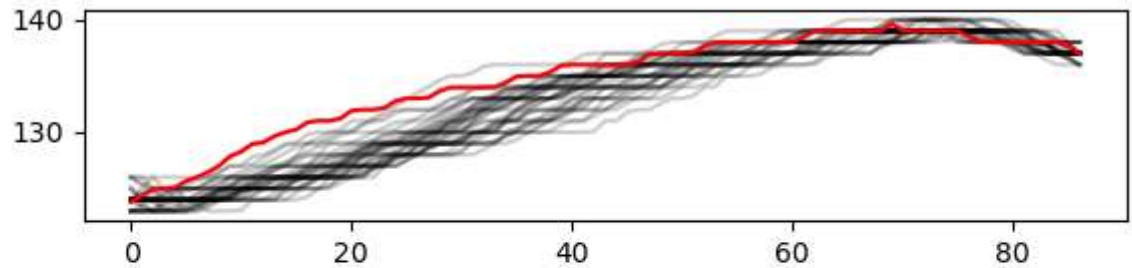
Cluster 3



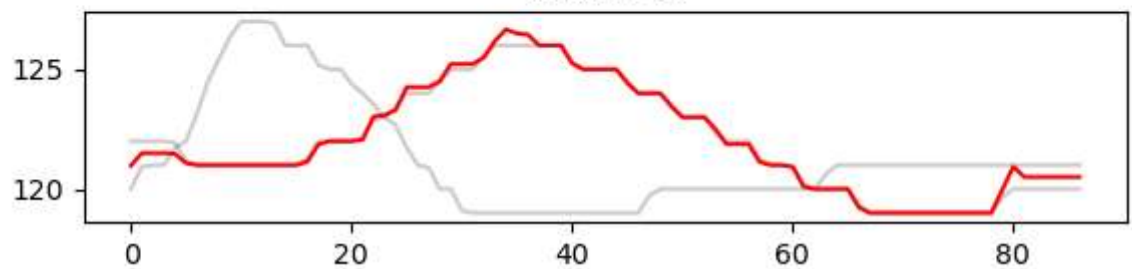
Cluster 4



Cluster 5



Cluster 6



## Cluster 2 and Cluster 6 seems to identified all the abnormal patterns.

```
In [31]: 1 #X_train_df = pd.DataFrame(X_train[:, 0, 0]) # Assuming X_train is a 3D a
2 y_pred_df = pd.DataFrame(y_pred, columns=['y_pred'])
3 #merged_df = pd.concat([X_train_df, y_pred_df], axis=1)
4 merged_df = pd.concat([X_df, y_pred_df], axis=1)
5 merged_df = merged_df.set_index(index_variable)
6 #merged_df
```

```
In [32]: 1 df=merged_df
2 print('0', (df['y_pred'] == 0).sum())
3 print('1', (df['y_pred'] == 1).sum())
4 print('2', (df['y_pred'] == 2).sum())
5 print('3', (df['y_pred'] == 3).sum())
6 print('4', (df['y_pred'] == 4).sum())
7 print('5', (df['y_pred'] == 5).sum())
```

```
0 90
1 3
2 63
3 29
4 41
5 2
```

```
In [33]: 1 print(X_train_scale.shape)
2 #X_train_scale
3 #X_train_scale[0]
4 X_train_scale
```

```
(228, 87)
```

```
Out[33]: array([[125.          , 123.69767442, 123.          , ..., 137.30232558,
137.          , 137.          ],
[120.          , 120.          , 120.          , ..., 138.          ,
137.44186047, 136.          ],
[121.          , 121.          , 121.          , ..., 137.55813953,
137.          , 137.          ],
...,
[123.          , 123.          , 123.          , ..., 137.          ,
136.61627907, 136.          ],
[121.          , 121.          , 121.          , ..., 139.6744186 ,
138.8372093 , 138.          ],
[122.          , 122.          , 121.55813953, ..., 138.          ,
137.72093023, 137.          ]])
```

```
In [34]: 1 X_train_scale2 = X_train_scale.reshape(X_train_scale.shape[0], X_train_scale.shape[1]*8)
2 X_train_scale2
```

```
Out[34]: array([[125.        , 123.69767442, 123.        , ..., 137.30232558,
        137.        , 137.        ],
       [120.        , 120.        , 120.        , ..., 138.        ,
        137.44186047, 136.        ],
       [121.        , 121.        , 121.        , ..., 137.55813953,
        137.        , 137.        ],
       ...,
       [123.        , 123.        , 123.        , ..., 137.        ,
        136.61627907, 136.        ],
       [121.        , 121.        , 121.        , ..., 139.6744186 ,
        138.8372093 , 138.        ],
       [122.        , 122.        , 121.55813953, ..., 138.        ,
        137.72093023, 137.        ]])
```

```
In [35]: 1 #X_train_chart = pd.DataFrame(X_train)
2 X_train_chart = pd.DataFrame(X_train_scale2)
3 X_train_chart
```

```
Out[35]:
```

	0	1	2	3	4	5	6	7
0	125.0	123.697674	123.000000	123.000000	123.000000	123.000000	123.000000	123.558140
1	120.0	120.000000	120.000000	120.000000	120.000000	119.395349	119.000000	119.046512
2	121.0	121.000000	121.000000	121.000000	121.000000	121.197674	121.837209	122.000000
3	123.0	123.000000	123.000000	123.000000	123.000000	123.000000	123.488372	124.000000
4	120.0	120.000000	120.000000	119.941860	119.255814	119.000000	119.116279	119.802326
...	...	...	...	...	...	...	...	...
223	124.0	124.000000	124.000000	124.000000	124.232558	124.790698	125.000000	125.000000
224	120.0	120.000000	120.000000	119.906977	119.209302	119.000000	119.186047	119.883721
225	123.0	123.000000	123.000000	123.000000	123.000000	123.000000	123.000000	123.000000
226	121.0	121.000000	121.000000	121.000000	121.000000	121.000000	121.000000	121.000000
227	122.0	122.000000	121.55814	121.000000	121.000000	121.000000	121.000000	121.000000

228 rows × 87 columns



In [36]:

```
1 merged_df_chart = X_train_chart.set_index(index_variable)
2 X_train_chart
```

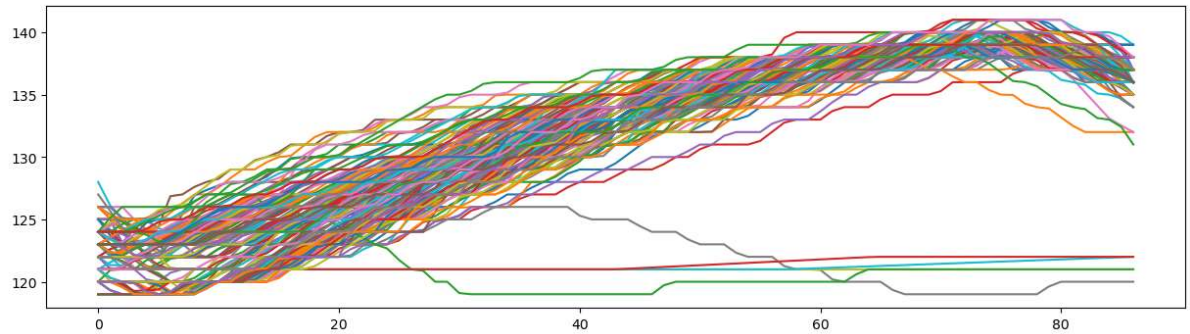
Out[36]:

	0	1	2	3	4	5	6	7
0	125.0	123.697674	123.000000	123.000000	123.000000	123.000000	123.000000	123.558140
1	120.0	120.000000	120.000000	120.000000	120.000000	119.395349	119.000000	119.046512
2	121.0	121.000000	121.000000	121.000000	121.000000	121.197674	121.837209	122.000000
3	123.0	123.000000	123.000000	123.000000	123.000000	123.000000	123.488372	124.000000
4	120.0	120.000000	120.000000	119.941860	119.255814	119.000000	119.116279	119.802326
...	...	...	...	...	...	...	...	...
223	124.0	124.000000	124.000000	124.000000	124.232558	124.790698	125.000000	125.000000
224	120.0	120.000000	120.000000	119.906977	119.209302	119.000000	119.186047	119.883721
225	123.0	123.000000	123.000000	123.000000	123.000000	123.000000	123.000000	123.000000
226	121.0	121.000000	121.000000	121.000000	121.000000	121.000000	121.000000	121.000000
227	122.0	122.000000	121.55814	121.000000	121.000000	121.000000	121.000000	121.000000

228 rows × 87 columns



```
In [38]: 1 plt.figure(figsize=(15, 4))
2 for index, row in merged_df_chart.iterrows():
3     plt.plot(row, label=f'Line {index}') # Plot each row with a label
4 plt.show()
5 merged_df_chart
```



Out[38]:

	0	1	2	3	4	5	6	7
0								
62301	125.0	123.697674	123.000000	123.000000	123.000000	123.000000	123.000000	123.558140
62302	120.0	120.000000	120.000000	120.000000	120.000000	119.395349	119.000000	119.046512
62303	121.0	121.000000	121.000000	121.000000	121.000000	121.197674	121.837209	122.000000
62304	123.0	123.000000	123.000000	123.000000	123.000000	123.000000	123.488372	124.000000
62305	120.0	120.000000	120.000000	119.941860	119.255814	119.000000	119.116279	119.802326
...	...	...	...	...	...	...	...	...
62596	124.0	124.000000	124.000000	124.000000	124.232558	124.790698	125.000000	125.000000
62597	120.0	120.000000	120.000000	119.906977	119.209302	119.000000	119.186047	119.883721
62598	123.0	123.000000	123.000000	123.000000	123.000000	123.000000	123.000000	123.000000
62599	121.0	121.000000	121.000000	121.000000	121.000000	121.000000	121.000000	121.000000
62600	122.0	122.000000	121.55814	121.000000	121.000000	121.000000	121.000000	121.000000

228 rows × 87 columns





```

1 # Optimization/Hyperparameter Tuning
2 import numpy as np
3 from tslearn.clustering import TimeSeriesKMeans
4 from sklearn.metrics import silhouette_score
5 X_train2 = X_train.reshape(X_train.shape[0], X_train.shape[1])
6 print(X_train2.shape)
7
8 max_clusters = 10 # Adjust the maximum number of clusters as needed
9 silhouette_scores = []
10
11 for n_clusters in range(2, max_clusters + 1):
12     # Create a TimeSeriesKMeans model with the current number of clusters
13     model = TimeSeriesKMeans(n_clusters=n_clusters, verbose=False, random_
14
15     # Fit the model to your time series data
16     model.fit(X_train2)
17
18     # Predict cluster labels for your data
19     cluster_labels = model.predict(X_train2)
20
21     # Calculate the silhouette score for the current number of clusters
22     silhouette = silhouette_score(X_train2, cluster_labels)
23
24     silhouette_scores.append(silhouette)
25 silhouette_scores

```



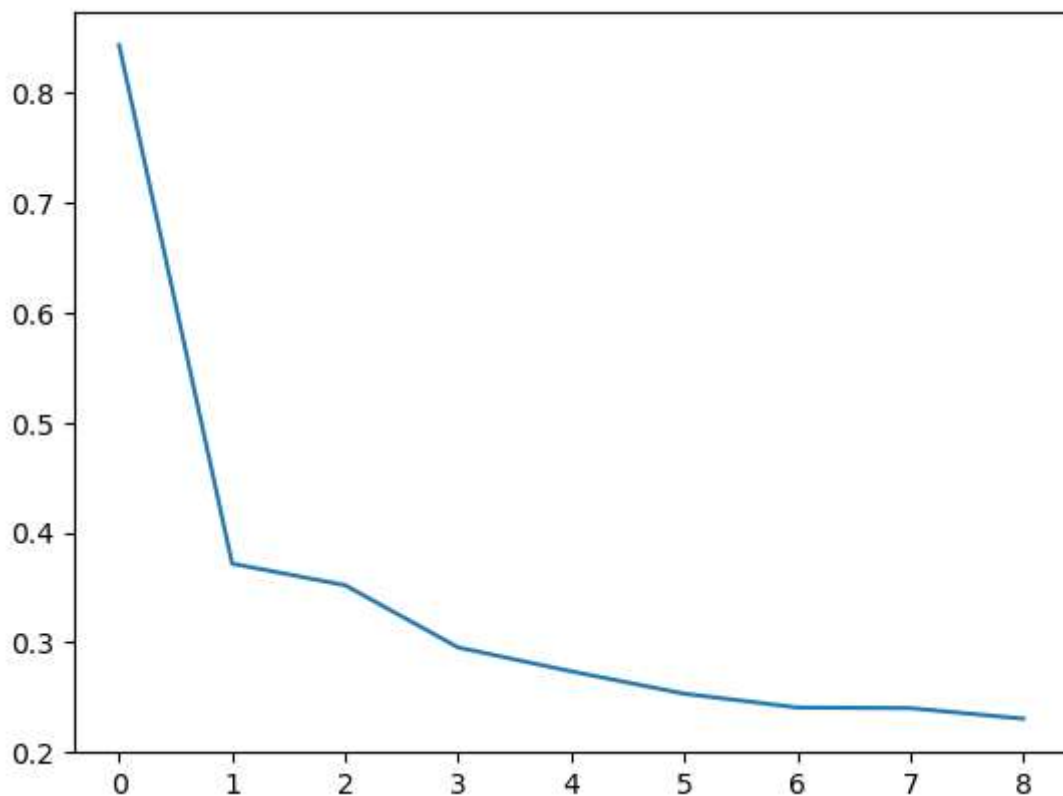
```
In [40]: 1 print(type(silhouette_scores))
         2 silhouette_scores
```

```
<class 'list'>
```

```
Out[40]: [0.8435585406819868,
          0.37147139080283603,
          0.3518400762253112,
          0.2951402633521269,
          0.27339617881284395,
          0.2530568447035266,
          0.24049922433566023,
          0.23991270989282823,
          0.23030860525031174]
```

```
In [41]: 1 plt.plot(silhouette_scores)
         2
         3     # plt.plot(row, label=f'Line {index}') # Plot each row with a label
         4 #plt.show()
```

```
Out[41]: [<matplotlib.lines.Line2D at 0x1f2b45e5e20>]
```



**Conclusion: The Elbow analysis shows the clustering is at 1 cluster but 3 and 6 is the more accurate clustering number.**

---

```
In [ ]: 1 # Reference
        2 #https://scikit-learn.org/stable/modules/generated/sklearn.metrics.silhouette_score.html
        3 #sklearn.metrics.silhouette_score(X, labels, *, metric='euclidean', sample_weight=None)
```

```
In [ ]: 1 # Example:
        2 from sklearn.datasets import make_blobs
        3 X, y = make_blobs(
        4     n_samples=500,
        5     n_features=2,
        6     centers=4,
        7     cluster_std=1,
        8     center_box=(-10.0, 10.0),
        9     shuffle=True,
       10     random_state=1,
       11 ) # For reproducibility
       12 print(X.shape)
       13 X
```

---

## Reference sample code

```
In [ ]: 1 import numpy as np
        2 from tslearn.clustering import TimeSeriesKMeans
        3 from tslearn.datasets import CachedDatasets
        4 from tslearn.preprocessing import TimeSeriesScalerMeanVariance
```

```
In [ ]: 1 mylist = CachedDatasets().list_datasets()
        2 mylist
```

```
In [ ]: 1 data_loader = CachedDatasets()
        2 X_train, y_train, X_test, y_test = data_loader.load_dataset("Trace")
        3 #https://tslearn.readthedocs.io/en/stable/gen_modules/datasets/tsLearn.datasets.html
```

```
In [ ]: 1 X_train
```

```
In [ ]: 1 # Author: Romain Tavenard
2 # License: BSD 3 clause
3 #https://tslearn.readthedocs.io/en/stable/auto_examples/clustering/plot_ks
4 import numpy
5 import matplotlib.pyplot as plt
6
7 from tslearn.clustering import KShape
8 from tslearn.datasets import CachedDatasets
9 from tslearn.preprocessing import TimeSeriesScalerMeanVariance
10
11 seed = 0
12 numpy.random.seed(seed)
13 X_train, y_train, X_test, y_test = CachedDatasets().load_dataset("Trace")
14 # Keep first 3 classes and 50 first time series
15 X_train = X_train[y_train < 4]
16 X_train = X_train[:50]
17 numpy.random.shuffle(X_train)
18 # For this method to operate properly, prior scaling is required
19 X_train = TimeSeriesScalerMeanVariance().fit_transform(X_train)
20 sz = X_train.shape[1]
21
22 # kShape clustering
23 ks = KShape(n_clusters=3, verbose=True, random_state=seed)
24 y_pred = ks.fit_predict(X_train)
25
26 plt.figure()
27 for yi in range(3):
28     plt.subplot(3, 1, 1 + yi)
29     for xx in X_train[y_pred == yi]:
30         plt.plot(xx.ravel(), "k-", alpha=.2)
31     plt.plot(ks.cluster_centers_[yi].ravel(), "r-")
32     plt.xlim(0, sz)
33     plt.ylim(-4, 4)
34     plt.title("Cluster %d" % (yi + 1))
35
36 plt.tight_layout()
37 plt.show()
```

```
In [ ]: 1
```