

\*WADI code file name\* (Note: variables adjusted during algorithm assessment and code may not run without changing variable names. Best to run cell by cell to adjust variables as needed. Required packages and libraries listed in cells)

- KMEANS\_WADI, (KMeans.ipynb, KMeans2.ipynb, KMeans3.ipynb)– novelty detection, subsystem grouping on WADI data set
  - Ensure that “WADI\_14days\_new.csv” and “WADI\_attackdataLABEL.csv “ are available.
- SWAT\_2015\_CUSUM (SWAT\_2015\_CUSUM.ipynb) – CUSUM anomaly detection on SWAT 2015
  - Ensure that “SWaT\_Dataset\_Attack\_v0.csv” is available. Initial code for KMeans not used for project.
- LSTM\_KNN (EDA\_LSTM.ipynb, LSTM.ipynb, LSTM2.ipynb, PyOD.ipynb, PyOD2.ipynb, PyOD3.ipynb) - LSTM and KNN anomaly detection WADI (EDA included)
  - Ensure that “WADI\_14days\_new.csv” and “WADI\_attackdataLABEL.csv “ are available.

### **main\_Clustering.py**

Usage: main\_Clustering.py

The main\_Clustering.py will run on SWaT 2015 and SWaT 2019 datasets with the below inputs and variable changes. The script is designed to run single or numerous test cases and produce visualization and metrics to assess the performance of the PCA and K-Means clustering outlier detection implementation.

Input:

- Modify the “test\_cases\_str” variable to a .csv with desired test cases to execute. The Columns and rows of the .csv as formatted as follows:

Case	clusters	PCA_Components	Threshold
1	5	5	2.5
2	6	5	2.5
3	7	5	2.5

- Modify “sub\_folder” variable to directory where data file exists.
  - Example: “/SWaT2019/”
- Modify “file\_name” variable to specific data file in directory above
  - Example: "SWaT\_Dataset\_2019.csv"
- Modify “trial” variable to specify if “test case” functions will run vs. All functions will run
  - If trial = True, an variable labeled “options” is set to equal an array as follows: [False, False, False, False, False, True, False]
    - This will execute only the confusion matrix scripts to grade the accuracy of all test cases
  - If Trial = False, the “options” is set to equal [True, True, True, True, True, True, True, True], indicating all functions and below plots will be produced.

Output:

- For each True/False value in the options array, the below outputs will be produced:

Index	Option array value	Output
0	True	Variance vs. PCA plot, PCA 1 vs. PCA 2 plot
1	True	Cluster vs. Inertia Plot (Elbow Plot)
2	True	PCA 1 vs. PCA 2 plot, PCA 2 vs. PCA 3 plot seaborn plots
3	True	3D Scatter Plot of clusters out of K-Means Model
4	True	3D Scatter Plot of original Normal vs. Attack data points/labels
5	True	Execute grade_Outlier function that will create confusion matrix with TP,TN,FP,FN and accuracy metrics
6	True	Produce Results.csv with confusion matrix and accuracy metrics

#### **"SWaTJUL2019 EDA CUSUM K-means match.ipynb"**

- Ensure that you have "SWaT\_dataset\_Jul 19 v3 (edited).csv" and "SWaT\_clusters\_2019.csv" loaded in the same working directory as this file
- Press "run all" to view dataframe clean up, normalization/standardization, early data analysis, CUSUM for SWaT2019, and CUSUM + K-means time matching for SWaT2019

#### **"SWaTDEC2015 CUSUM K-Mean match.ipynb"**

- Ensure that you have "SWaT\_Dataset\_Attack\_v0.csv", "SWaT\_Dataset\_Normal\_v0.csv", and "Case\_0\_SWaT2015\_P1\_clusters.csv" loaded in the same working directory as this file
- Press "run all" to view CUSUM for SWaT2015, and CUSUM + K-means time matching for SWaT2015