

Report of Anomaly Detection HW2

Name: 高煒軒

Student ID: 112062646

Implementation & explanation

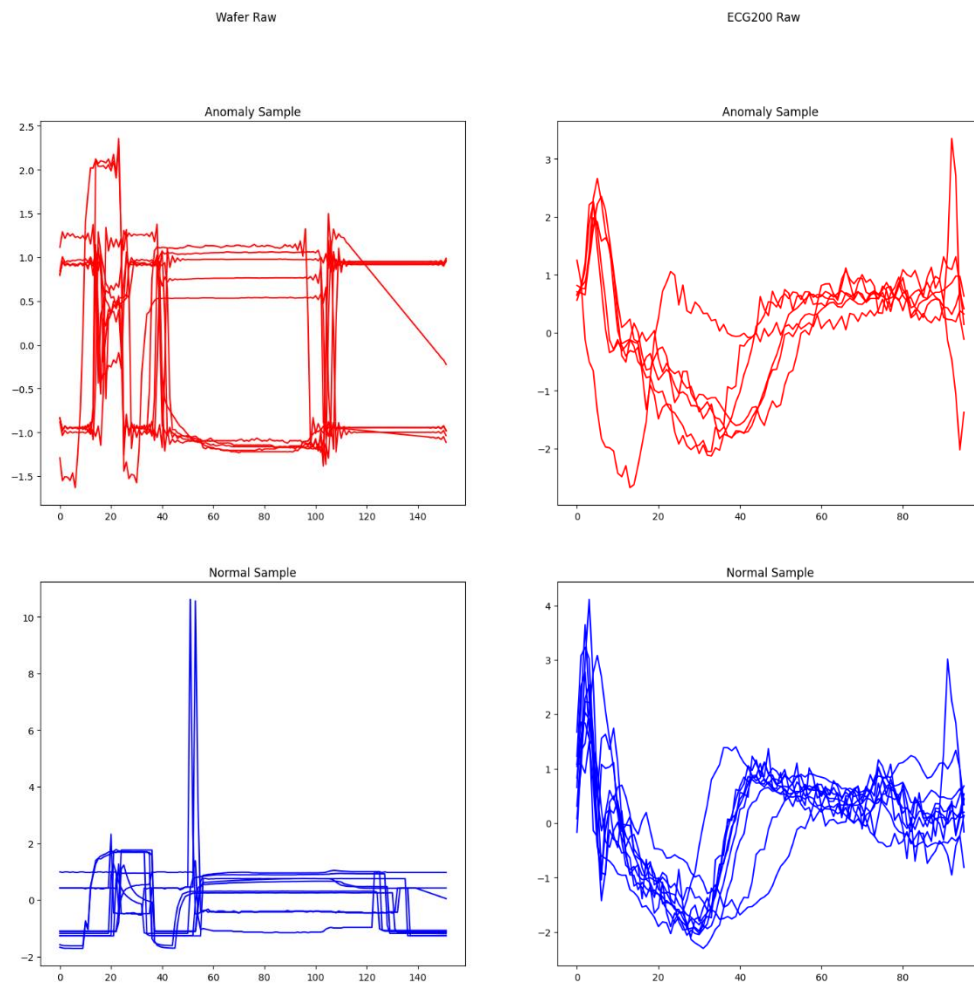
本次作業中的 5 個 Problem implementation code & explanation，會依據 5 個不同的 Problem 各自解釋。

(1) Visualization

Function 呼叫為:

```
data_plot(test_data, test_label, category+" Raw")
```

結果如下:



■ Implementation code

```
21 def data_plot(data, label, title):
22     normal_data = data[label==0]
23     anomaly_data = data[label==1]
24
25     normal_idx = np.random.choice(normal_data.shape[0], size=10, replace=False)
26
27     # Selecting "ALL" anomaly data, if the number of anomaly data <= 10
28     if anomaly_data.shape[0] > 10:
29         anomaly_idx = np.random.choice(anomaly_data.shape[0], size=10, replace=False)
30     else:
31         anomaly_idx = np.arange(anomaly_data.shape[0])
32
33
34     fig, axs = plt.subplots(2,1, figsize=(8,16))
35     x = np.arange(0, anomaly_data.shape[1])
36     y = anomaly_data[anomaly_idx]
37     for i in range(anomaly_idx.shape[0]):
38         axs[0].plot(x, y[i], 'r')
39     axs[0].set_title("Anomaly Sample")
40
41     y = normal_data[normal_idx]
42     for i in range(10):
43         axs[1].plot(x, y[i], 'b')
44     axs[1].set_title("Normal Sample")
45     fig.suptitle(title)
46     plt.show()
```

(2) Raw Data

Function 呼叫順序為:

```
Raw_AD = AD()  
Raw_AD(train_data, train_label, test_data, test_label
```

KNN_AD(): based on K-neighbors classification 的 Anomaly Detection algorithm。回傳 roc-auc 值。

AD():

__inti__: 定義特別的 feature processing function e.g. DFT(), DWT()。

__call__: 使用定義好的 function 處理資料、取得 features 後，將處理後的 data 丟入 KNN_AD。

■ Implementation code

```
48 def KNN_AD(train_X, train_y, test_X, test_y, k=5):  
49     KNN = KNeighborsClassifier(n_neighbors=k)  
50     KNN.fit(train_X, train_y)  
51     dist, _ = KNN.kneighbors(test_X)  
52     avg_dist = np.mean(dist, axis=1)  
53  
54     return roc_auc_score(y_true=test_y, y_score=avg_dist)
```

```
115 class AD:  
116     def __init__(self, func=None) -> None:  
117         self.process_func = func  
118  
119     def __call__(self, train_X, train_y, test_X, test_y, k=5, arg=None) -> float:  
120         if self.process_func == None:  
121             train_proc = train_X  
122             test_proc = test_X  
123         else:  
124             train_proc = self.process_func(train_X, arg)  
125             test_proc = self.process_func(test_X, arg)  
126  
127         return KNN_AD(train_proc, train_y, test_proc, test_y, k=k)
```

(3) PCA

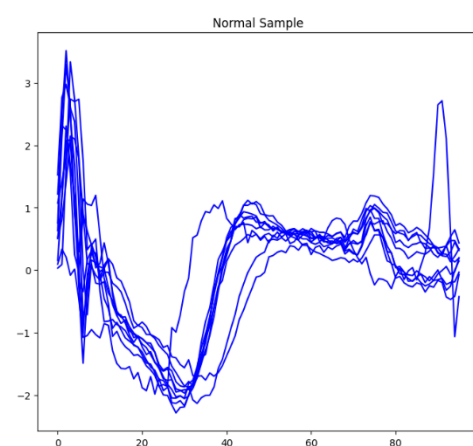
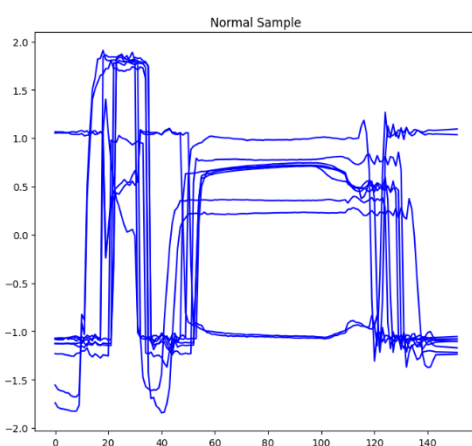
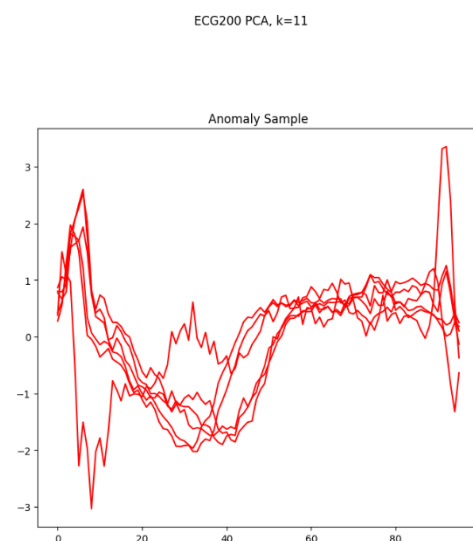
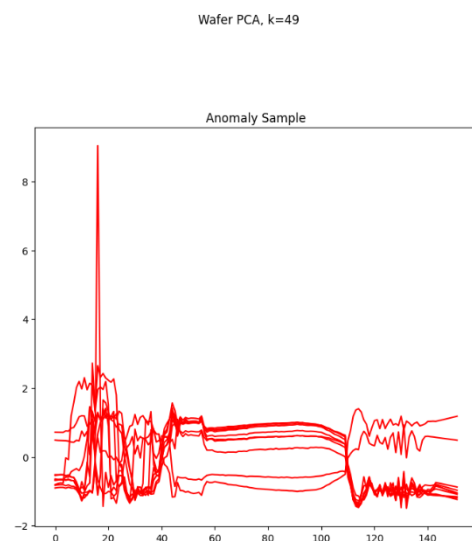
Function 呼叫順序為:

```
PCA_AD(train_data, test_data, test_label, k=i)
```

此處 **N** 值嘗試 1 到 **data** 的 **attribute** 數的 1/4 。

PCA_AD():

1. 使用 **training data** 找出 **PCA** 的 **N** 個 **components** 。
2. 使用這 **N** 個 **components** 將 **testing data** 轉換成 **feature**，再將 **feature** 重新 **reconstruct** 回 **testing data** 。
3. **Reconstruction** 與原 **testing data** 之間的距離即為 **anomaly score**



■ Implementation code

```
56 def PCA_AD(train_X, test_X, test_y, k, return_recons=False):
57     pca = PCA(n_components=k)
58     pca.fit(train_X)
59     test_recons = pca.inverse_transform(pca.transform(test_X))
60     score = np.linalg.norm(test_X-test_recons, axis=1)
61     if return_recons:
62         return test_recons, roc_auc_score(y_true=test_y, y_score=score)
63
64     return roc_auc_score(y_true=test_y, y_score=score)
```

(4) Discrete Fourier Transform

Function 呼叫順序為:

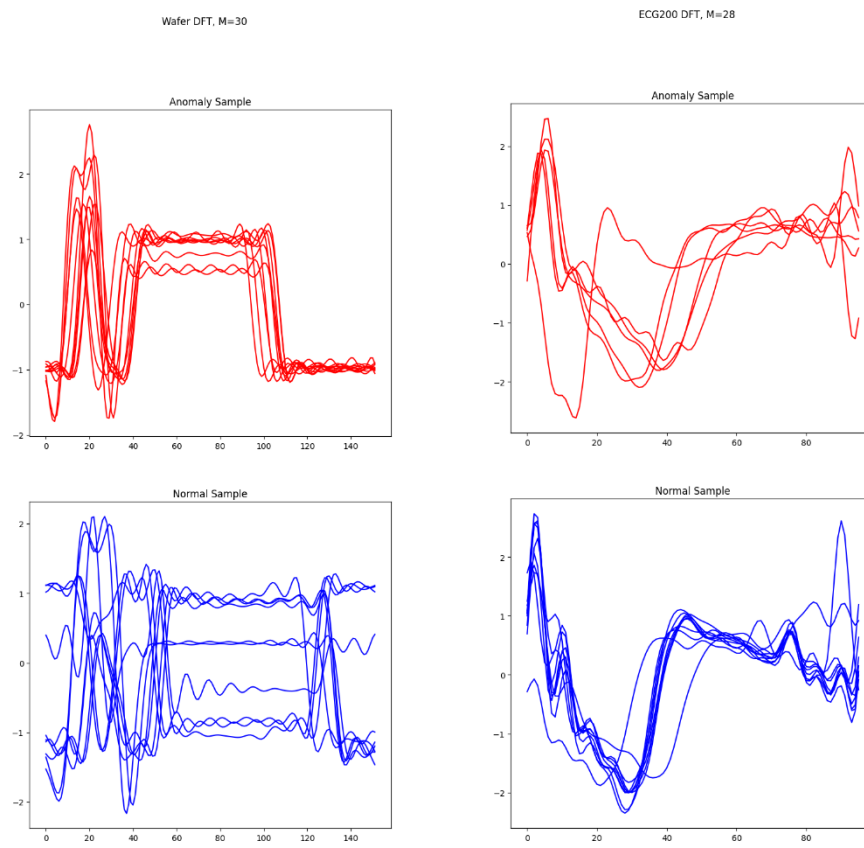
```
DFT_AD = AD(DFT)
DFT_AD(train_data, train_label, test_data, test_label, arg=i)
```

此處 **M** 值嘗試 1 到 data 的 attribute 數的 1/2 。

AD(): 同上。

DFT():

1. 對 data 做 discrete Fourier Transform 取得 coefficients 。選擇 coefficients 中對應最小的 M 個 frequency 的 coefficients 作為 features
2. 用 features 計算 magnitudes 及 reconstruction 。
3. Magnitudes 丟給 KNN 計算 anomaly score , reconstruction 用來 visualization 。



■ Implementation code

```
66 def DFT(X, M=20, use_first_m = False, return_recons = False):
67     ceo      = np.fft.fft(X)
68     mag_idx = np.zeros(X.shape[-1], dtype=bool)
69
70     if use_first_m:      # use first M coes as features
71         mag_idx[:M] = True
72     else:                 # use first M/2 and last M/2 coes as featrues
73         mag_idx[:int(np.ceil(M/2))] = True
74         mag_idx[-int(np.floor(M/2)):] = True
75
76     magnitudes      = np.abs(ceo[:, mag_idx])
77     ceo_selected    = np.zeros_like(ceo)
78     ceo_selected[:, mag_idx] = ceo[:, mag_idx]
79     recons          = np.fft.ifft(ceo_selected)
80
81     if return_recons:
82         return magnitudes, recons
83     else:
84         return magnitudes
```

(5) Discrete Wavelet Transform

Function 呼叫順序為:

```
DWT_AD = AD(DWT)
DWT_AD(train_data, train_label, test_data, test_label, arg=i)
```

此處 **S** 值嘗試 1 到 data 的 level 值+1。

另外，`level=range([1,ceiling(log2(feature_dim))])`

AD(): 同上。

DWT():

1. 對 Data 做 Discrete Wavelet Transform，並取 last level 的 A coefficient 及倒數 S-1 個 D coefficients 作為 features。
2. 將 features 丟入 KNN 中計算 anomaly score。

■ Implementation code

```
86 def DWT(X, S, return_coes=False):
87     S = 2**S
88     levels = int(np.ceil(np.ma.log2(X.shape[-1])))
89     size = 2**(levels) - 1
90     X_padded = np.pad(X, pad_width=((0,0),(0,2**levels-X.shape[-1])), mode='constant', constant_values=0)
91     A_coes = np.pad(X_padded, ((0,0), (0, size)), mode='constant', constant_values=0)
92     D_coes = np.zeros((X.shape[0], size))
93
94     prev_begin = 0
95     begin = X_padded.shape[-1]
96     D_begin = 0
97
98     for i in range(1, levels+1):
99         len = levels - i
100
101         A_coes[:, begin:begin+2**len] = (A_coes[:, prev_begin+1:begin:2] + A_coes[:, prev_begin:begin:2]) / 2
102         D_coes[:, D_begin:D_begin+2**len] = (A_coes[:, prev_begin+1:begin:2] - A_coes[:, prev_begin:begin:2]) / 2
103
104         prev_begin = begin
105         begin += 2**len
106         D_begin += 2**len
107
108     out_A_coes = np.expand_dims(A_coes[:, -1], axis=1)
109     out_D_coes = D_coes[:, -1:-S:-1]
110     if return_coes:
111         return np.concatenate([out_A_coes, out_D_coes], axis=1), A_coes[:, -size:], D_coes
112
113     return np.concatenate([out_A_coes, out_D_coes], axis=1)
```

Performance

1. Raw Data:

Wafer 及 ECG200 在 raw data 的 knn anomaly detection 上 performance 分別為: Wafer: 0.9884 , ECG200: 0.8750

2. PCA:

Wafer 及 ECG200 在經由 PCA 處理後，performance 最大值分別為:

Wafer: 0.9961, N=49

ECG200: 0.8568, N=11

另外此處 N 值的範圍為 $1 \sim \text{feature_dim}/4$ (實際取到 25 的倍數)

3. DFT:

Wafer 及 ECG200 在經由 DFT 處理後，performance 最大值分別為:

Wafer: 0.9984, M=30

ECG200: 0.8568, M=11

另外此處 M 值的範圍為 $1 \sim \text{feature_dim}/2$ (實際取到 50 的倍數)

4. DWT:

Wafer 及 ECG200 在經由 DWT 處理後，performance 最大值分別為:

Wafer: 0.9986, S=8

ECG200: 0.9479, S=4

Bonus:

此處針對 Raw Data, DFT, DWT 進行以下混合實驗:

KNN : $K \rightarrow [1, 10]$

DFT : $M \rightarrow [1, \text{feature_dim}/4]$ (實際上無條件進位到 25 的倍數)

DWT : $S \rightarrow 2^{[1, \text{ceiling}(\log_2(\text{feature_dim}))+1]}$

最終各個 methods 在兩個 datasets 的最佳 performance 如下:

1. Wafer:

Raw Data, $k=1$, ROC-AUC : 0.9913

DFT, $k=2$, $M=29$, ROC-AUC : 0.9985

DWT, $k=10$, $S=8$, ROC-AUC : 0.9987

2. ECG200:

Raw Data, $k=5$, ROC-AUC : 0.8750

DFT, $k=4$, $M=45$, ROC-AUC : 0.9297

DWT, $k=3$, $S=4$, ROC-AUC : 0.9531