

# Report of Anomaly Detection HW1

Name: 高煒軒

Student ID: 112062646

## Implementation & explanation

本次作業中的四個 algorithm's implementation code & explanation，會依據四個不同的 algorithm 各自解釋。

### (1) KNN

Function 呼叫順序為：

```
knns.fit(train_data)
knns.prediction(test_data)
```

*fit()*: 記住 input training data。

*score()*: 計算 data 之間的距離。

*prediction()*: 計算出距離後選出前 k 個近的 data，並且這前 k 個近的 data 距離平均即為 anomaly score。

## ■ Implementation code

```
1 from sklearn.metrics import pairwise_distances
2
3 def euclidean(point, data, axis=None):
4     return np.sqrt(np.sum((point - data)**2, axis=axis))
5
6 class KNN_AD:
7     def __init__(self, k=1):
8         """
9         Implement of Anomaly Detection Algorithm with KNN-classifier
10         :param k: "K" nearest neighbor
11         """
12         self.k = k
13
14     def fit(self, X):
15         self.X = X
16
17     def score(self, Y):
18         return pairwise_distances(X=Y, Y=self.X)
19
20     def prediction(self, Y, normal=False):
21         # Get sorted pairwise distances training and testing data
22         pair_dist = self.score(Y)
23         sorted_idx = np.argsort(pair_dist)
24         sorted_idx = sorted_idx[:, :self.k]
25
26         selected = [self.X[index] for index in sorted_idx]
27         selected = np.array(selected)
28
29         output = [np.sum(euclidean(point=Y[idx], data=selected[idx], axis=1)) for idx in range(Y.shape[0])]
30         output = np.array(output)
31         output = output / self.k
32
33         # Mapping output -> [0,1]
34         if normal:
35             Min = min(output)
36             Max = max(output)
37             output = (output - Min)/(Max - Min)
38
39         return output
```

## (2) K-means

Function 呼叫順序為:

```
kms.fit(train_data)
kms.prediction(test_data)
```

*fit():*

1. 選出 k 個 **centroids**
2. 將 **data** 分配到距離最近的 **cluster**，再重新計算 **centroids**(取 **cluster** 中的所有 **data** 的平均)
3. 以上迭代至 **centroids** 不再更新或到達最大迭代次數。

*score():* 計算 **data** 與 **centroids** 之間的距離

*prediction():* 計算出 **data** 與 **centroids** 之間的距離後，將這些距離平均即為該 **data** 的 **anomaly score**。

## ■ Implementation code

```
1 class KMeans_AD:
2     def __init__(self, n_clusters=1, max_iter=300):
3         self.n_clusters = n_clusters
4         self.max_iter = max_iter
5
6     def fit(self, X):
7         idx = np.random.choice(X.shape[0], self.n_clusters, replace=False)
8         self.centroids = np.array([X[i] for i in idx])
9
10        iteration = 0
11        prev_centroids = np.array(None)
12        while (self.centroids != prev_centroids).any() and iteration < self.max_iter:
13
14            dists = pairwise_distances(X, self.centroids)
15            centroids_id = np.argmin(dists, axis=1)
16
17            buckets = []
18            for i in range(self.n_clusters):
19                buckets.append(np.where(centroids_id==i))
20
21            prev_centroids = np.copy(self.centroids)
22            for i in range(self.n_clusters):
23                self.centroids[i] = np.mean(X[buckets[i]], axis=0)
24
25            for i, centroid in enumerate(self.centroids):
26                if np.isnan(centroid).any():
27                    self.centroids[i] = prev_centroids[i]
28            iteration += 1
29
30    def score(self, Y):
31        return pairwise_distances(X=Y, Y=self.centroids)
32
33    def prediction(self, Y, normal=False):
34        pair_dists = self.score(Y)
35        output = np.min(pair_dists, axis=1)
36
37        if normal:
38            Min = min(output)
39            Max = max(output)
40            output = (output - Min)/(Max - Min)
41
42        return output
```

### (3) Distance Base

Function 呼叫順序為:

```
dist_ad.prediction(X=test_data, ref=train_data, metric=metrics[i])
```

*prediction()*: 計算出 data 之間的距離後，將最近的五個 data 距離平均即為該 data 的 anomaly score。

**ref**: input parameter, 用來計算 Mahalanobis Distance 所需之 covariance matrix。

**metric**: input parameter, 用來選擇使用哪個 distance metric。

*CosineDistance (X, Y)*:

i.e.  $(1 - \text{Cosine similarity}(X, Y))$

Cosine similarity definition:

$$\cos(X, Y) = \frac{X \cdot Y}{\|X\| \|Y\|}$$

*L1Distance (X, Y)*:

Minkowski Distance, if  $r=1$ . Also called Manhattan distance.

Definition:

$$d(X, Y) = \sum_i |X_i - Y_i|$$

*L2Distance (X, Y)*:

Minkowski Distance, if  $r=2$ . Also called Euclidean distance.

Definition:

$$d(X, Y) = \sqrt{\sum_i (X_i - Y_i)^2}$$

*ChebyshevDistance (X, Y):*

Minkowski Distance, if  $r=\infty$ .

Definition:

$$d(X, Y) = \max_i |X_i - Y_i|$$

*MahalanobisDistance (X, Y, VI):*

Mahalanobis distance, 其中 VI 為 inverse covariance matrix 以 training data 計算求得。

Definition:

$$VI = V^{-1}$$

$$d(X, Y, V^{-1}) = \sqrt{(X - Y)^T V^{-1} (X - Y)}$$

## ■ Implementation code

```
1 def CosineDistance(X, Y):
2     X_dot_Y = np.dot(X,Y)
3     X_norm = np.linalg.norm(X)
4     Y_norm = np.linalg.norm(Y)
5     return 1 - (X_dot_Y / (X_norm * Y_norm))
6
7 def L1Distance(X, Y):
8     return np.sum(np.abs(X - Y))
9
10 def L2Distance(X, Y):
11     return np.linalg.norm(X - Y)
12
13 def ChebyshevDistance(X, Y):
14     return np.max(np.abs(X-Y))
15
16 def MahalanobisDistance(X, Y, VI):
17     diff = X - Y
18     return np.sqrt(np.matmul(np.matmul(diff, VI), diff.T))
19
20 class Distance_AD:
21     def __init__(self, k=5):
22         self.k = k
23         self.metrics = {"CosineDistance": CosineDistance,
24                         "L1Distance": L1Distance,
25                         "L2Distance": L2Distance,
26                         "ChebyshevDistance": ChebyshevDistance,
27                         "MahalanobisDistance": MahalanobisDistance}
28
29     def predict(self, X, ref=None, metric="L2Distance"):
30         if (ref != None).any() and metric=="MahalanobisDistance":
31             self.VI = np.linalg.inv(np.cov(ref.T))
32             dists = pairwise_distances(X=X, metric=self.metrics[metric], VI=self.VI)
33         else:
34             dists = pairwise_distances(X=X, metric=self.metrics[metric])
35         output = np.sort(dists, axis=1)[: , self.k]
36         return output
```

## (4) LOF

Function 呼叫順序為:

```
Lof.predict(test_data)
```

*prediction()*: 依序呼叫 *\_KNN()*, *\_reach\_distance()*, *\_LRD()*, *\_LOF\_score()* 後輸出 data 的 anomaly score。

*\_KNN()*: 計算出 data 之間的距離(dist)，並以此記錄各個 data 的 k neighbor(即 k\_near)和第 k 個 neighbor 的距離(即 k\_dist)。

k\_near 及 k\_dist 請參考 implementation code。

*\_reach\_distance()*: 計算各個 data 的 reachable distance。

※為加速 algorithm，我不使用 loop 來計算 reachable distance，而是將 k\_dist 擴展到與 dists 相同維度並比較大小。

*\_LRD()*: 計算 local reachable density of data，此處 output 為 lrd 的 inverse。

*\_LOF\_score()*: 計算 local outlier factor of data。

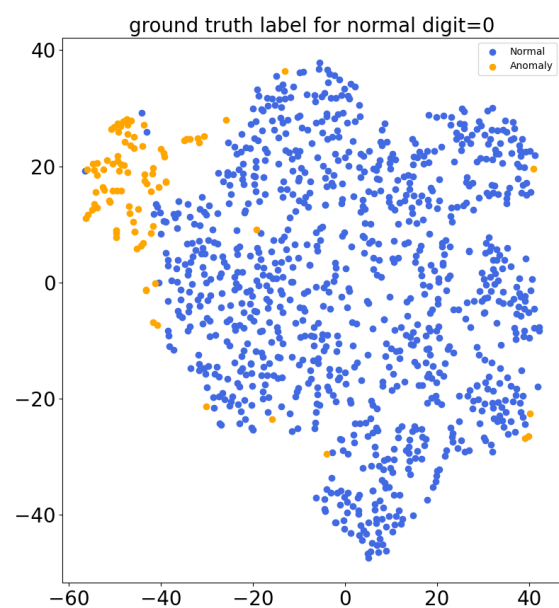
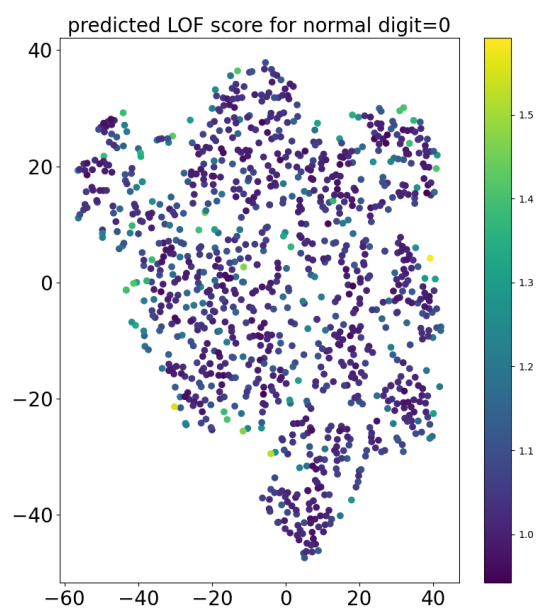


## ■ Implementation code

```
1 class LOF:
2     def __init__(self, k=5):
3         self.k = k
4
5     def _KNN(self, X):
6         dists = pairwise_distances(X=X)
7         k_near = np.argsort(dists, axis=1)[: , 1:self.k+1]
8         k_dist = np.sort(dists, axis=1)[: , self.k]
9
10        return dists, k_near, k_dist
11
12    def _reach_distance(self, dists, k_dist):
13        k_dist_ex = np.expand_dims(k_dist, axis=0)
14        k_dist_rep = np.repeat(k_dist_ex, self.size, axis=0)
15        reach_dist = np.maximum(k_dist_rep, dists)
16        assert (reach_dist > 0).all(), "Error:: Reachable Distance expect to be > 0, but receive <= 0 value"
17
18        return reach_dist
19
20    def _LRD(self, reach_dist, k_near, (variable) idx: int):
21        # Local reachability dist
22        k_reachDist = [reach_dist[idx, k_near[i]] for idx in range(self.size)]
23        lrd = np.array([len(k_near[i])/(k_reachDist[i].mean()) for i in range(self.size)])
24
25        return lrd
26
27    def _LOF_score(self, lrd, k_near):
28        k_lrds = [lrd[k_near[idx]] for idx in range(self.size)]
29        lof = np.array([k_lrds[idx].sum()/lrd[idx] for idx in range(self.size)])
30
31        return lof / self.k
32
33    def predict(self, X):
34        self.size = X.shape[0]
35        dists, k_near, k_dist = self._KNN(X)
36        reach_dist = self._reach_distance(dists, k_dist)
37        lrd = self._LRD(reach_dist, k_near)
38        output = self._LOF_score(lrd, k_near)
39
40        return output
```

## ■ Visualization

```
1 # The TSNE.png of LOF
2 test_data, test_label = resample(orig_test_data, orig_test_label, target_label=0, outlier_ratio=0.1)
3 test_label = np.where(test_label==0, 0, 1)
4 Lof_score = Lof.predict(test_data)
5
6 t_sne = TSNE()
7 t_sne_out = t_sne.fit_transform(test_data)
8 normal_idx = test_label == 0
9 anomaly_idx = test_label == 1
10
11 fig, axs = plt.subplots(1, 2, figsize=(20, 10))
12
13 im = axs[0].scatter(t_sne_out[:, 0], t_sne_out[:, 1], c=Lof_score)
14 axs[0].set_title("predicted LOF score for normal digit=0", fontsize=20)
15 axs[0].tick_params(axis='both', which='major', labelsize=20)
16
17 axs[1].scatter(t_sne_out[normal_idx, 0],
18               t_sne_out[normal_idx, 1],
19               c='royalblue',
20               label='Normal')
21 axs[1].scatter(t_sne_out[anomaly_idx, 0],
22               t_sne_out[anomaly_idx, 1],
23               c='orange',
24               label='Anomaly')
25 axs[1].legend()
26 axs[1].set_title("ground truth label for normal digit=0", fontsize=20)
27 axs[1].tick_params(axis='both', which='major', labelsize=20)
28
29 plt.colorbar(im, ax=axs[0])
30 plt.savefig("TSNE.png", format="png")
31 plt.show()
```



# Performance

```
Average ROC-AUC of KNN:      (k=1)    0.9658, (k=5) 0.9683, (k=10) 0.9669
Average ROC-AUC of K-means:   (k=1)    0.9032, (k=5) 0.9487, (k=10) 0.9611
Average ROC-AUC of Distance-Base: (Cosine) 0.9744, (r=1) 0.9483, (r=2) 0.9517, (r=inf) 0.9532, (mahalanobis) 0.9794
Average ROC-AUC of LOF: 0.7644
```

## 1. KNN:

根據 **ROC-AUC**，可以看到 **KNN** 在這次的 **anomaly detection task** 的效能不會

因為 **k** 值的改變而有太多的效能增減。此現象可能與 **MNIST** 的不同的 **label** 對

應的圖像經過 **PCA** 降維後分布差距較大有關。

## 2. K-means:

相對於上述之 **KNN**，我們可以看到 **K-means** 的效能根據 **k** 值而有所改變。這

可能是因為當 **k** 值設定過小時，會導致過多的 **normal data** 被判斷為 **anomaly**。

## 3. Distance Base:

在我的作業撰寫過程中，我使用了 **sklearn** 的 **api : pairwise\_distance** 來幫助計

算距離，並使用自己寫的 **metrics function** 作為此 **api** 的 **callback function**。

並且我發現不同的 **metrics function** 除了在 **ROC-AUC** 的結果上有差異之外，

使用自定義的 **metrics function** 還會降低 **pairwise\_distance** 的執行速度(相較

於直接使用 **sklearn** 提供的 **metric** 選項)，我猜測這可能與 **sklearn** 的執行速度

優化有關。

## 4. LOF:

在這項作業中 **LOF** 的效能是最差的，並且根據 **visualize** 的結果，我感覺

**density** 與是否為 **anomaly** 並非完全相關。