

DDA3020 Machine Learning Assignment 4

Ma Xinxian 121090408@link.cuhk.edu.cn

June 13, 2023

Contents

1	Task Description	2
2	PCA	2
2.1	What And Why	2
2.2	Procedure	2
2.3	Code with Python	2
3	K-means	3
3.1	What And Why	3
3.2	Procedure with Python	3
4	Performance evaluation	6
4.1	Silhouette Coefficient	6
4.2	Rand Index	7

1 TASK DESCRIPTION

Clustering on UCI seed dataset with **PCA** and **K-means**. To contact me, send to my email on the cover page. Welcome any suggestions and requirements!

2 PCA

2.1 WHAT AND WHY

PCA aims to find a low-dimensional data vector to represent the original high-dimensional data vector so that it may help visualization, alleviate overfitting or reduce the computational cost.

2.2 PROCEDURE

Step1: Calculate the empirical covariance matrix [numpy]

$$\Sigma = \frac{1}{N} \sum_{n=1}^N (x^{(n)} - \mu)(x^{(n)} - \mu)^T$$

Step2: Do SVD decomposition of Σ to obtain its D eigenvalues λ_i and eigenvectors q_i by using `np.linalg.eig(cov)` function in numpy, and rank them from large to small according to the eigenvalues.

Step3: Pick the top- K eigenvectors to form the matrix $U = [q_1, \dots, q_K]$ by `argmin` function.

Step4: The new representation of $x^{(n)}$ is $U^T(x^{(n)} - \mu)$.

2.3 CODE WITH PYTHON

```
1 X = np.subtract(X, np.mean(X, axis=0))
2 # Step1
3 COV = np.dot(X.T, X)/X.shape[0]
4 # Step2
5 eigenvalue, eigenvector = np.linalg.eig(COV)
6 print(eigenvalue)
```

```
[1.07419301e+01 2.11931485e+00 7.32794138e-02 1.28261257e-02
 2.73513989e-03 1.56297146e-03 2.95142261e-05]
```

```
1 #Step3 & 4
2 # from above, the first and the second are top two
3 z1 = np.dot(X, eigenvector.T[0])
```



```

4 z2 = np.dot(X, eigenvector.T[1])
5 new_X = np.array([z1, z2]).T
6 print(new_X)

```

3 K-MEANS

3.1 WHAT AND WHY

K-means clustering is a method of vector quantization, that aims to partition n observations/samples into k clusters in which each observation belongs to the cluster with the nearest mean.

3.2 PROCEDURE WITH PYTHON

Step1: First, choose K — the number of clusters. Then you randomly put K feature vectors, called centroids, to the feature space. Here, 3.

Step2: Next, compute the distance from each example x to each centroid c using the Euclidean distance. In the code, this is implemented by *cal_dis()*

```

1 def cal_dis(dataSet, centroids, k):
2     # output's distance corresponding to the index of the center
3     # sort_distance_index outputs a sorted center's index list
4
5     distance_list = []
6     sort_distance_index = []
7
8     for i in range(dataSet.shape[0]):
9         ith_dis = []
10        for center in centroids:
11            ith_dis.append(np.linalg.norm(dataSet[i] - center))
12
13        distance_list.append(ith_dis)
14        sort_distance_index.append(sorted(range(len(ith_dis)), key=
15                                         lambda k: ith_dis[k], reverse=False))
16
17    return distance_list, sort_distance_index

```

Then we assign the closest centroid to each example.

Step3: For each centroid, we calculate the average feature vector of the examples labeled

with it. These average feature vectors become the new locations of the centroids. In the code, this is implemented by `classify()`. New centroids are `newCentroids`.

```
1 def classify(dataSet, centroids, k):
2
3     distance_list, sort_distance_index = cal_dis(dataSet, centroids
4         , k)
5
6     minDistIndices = np.argmin(distance_list, axis=1)
7     newCentroids = pd.DataFrame(dataSet).groupby(minDistIndices).
8         mean()
9
10    newCentroids = newCentroids.values
11    changed = newCentroids - centroids
12
13    return changed, newCentroids
```

Step4: We recompute the distance from each example to each centroid, modify the assignment and repeat the procedure until the assignments don't change after the centroid locations are recomputed. This is implemented by the while loop in the *k-means* function.

```
1 def kmeans(dataSet, k):
2     cluster = [[], [], []]
3     centroids = random.sample(list(dataSet), k)
4
5     # initial assignment
6     changed, newCentroids = classify(dataSet, centroids, k)
7
8     while np.any(changed != 0):
9         changed, newCentroids = classify(dataSet, newCentroids, k)
10
11    centroids = newCentroids.tolist()
12    distance_list, sort_distance_index = cal_dis(dataSet, centroids
13        , k)
14    minDistIndices = np.argmin(distance_list, axis=1)
15
16    for i, j in enumerate(minDistIndices):
17        cluster[j].append(dataSet[i].tolist())
18
19    return centroids, cluster
```

Step5: Finally, we output the centroids and the cluster (too long, see in the code file). And we also plot the classification.

```
1 centroids, cluster = kmeans(new_X, 3)
2 print('Centroids{}'.format(centroids))
3
4 for i in range(3):
5     print('Cluster{}:{}'.format(i, cluster[i]))
6
7 color_list = ['green', 'blue', 'purple']
8 for i in range(3):
9     clus = cluster[i]
10    color = color_list[i]
11    for point in clus:
12        plt.scatter(point[0],point[1], marker = 'o',color = color ,
13                    s = 40, label = 'original point')
14
15 for i in range(len(centroids)):
16     plt.scatter(centroids[i][0],centroids[i][1],marker='x',color='
17                 red', s=50,label='centroids')
```

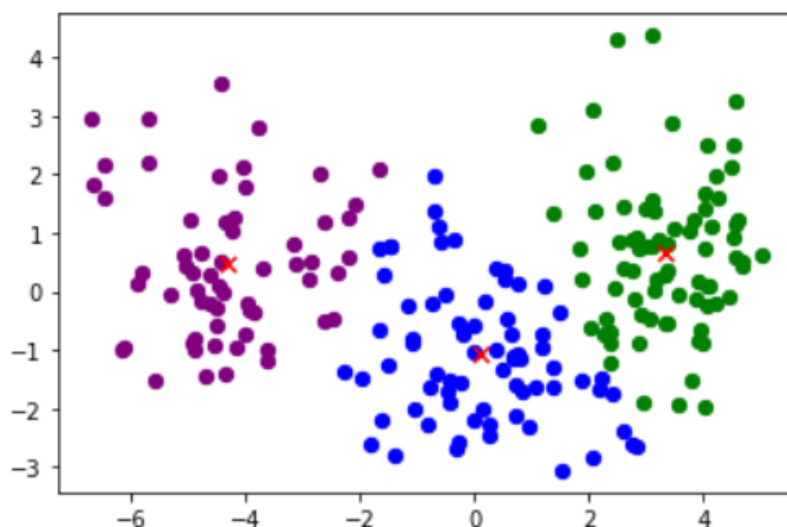


FIGURE 1: CLASSIFICATION.

4 PERFORMANCE EVALUATION

4.1 SILHOUETTE COEFFICIENT

Given a clustering, we define

- **a**: The mean distance between a point and all other points in the same cluster.
- **b**: The mean distance between a point and all other points in the next nearest cluster.

Then, Silhouette coefficients for a single sample is formulated as:

$$s = \frac{b - a}{\max(a, b)}$$

And Silhouette coefficient s for a set of samples is defined as the mean of the Silhouette Coefficient for each sample. Thus, we implement it as following:

```
1 # Silhouette Coefficient
2 distance_list, sort_distance_index = cal_dis(new_X, centroids, k)
3 s_sum = 0
4 for i in range(new_X.shape[0]):
5     s_data = new_X[i]
6     smallest_index = sort_distance_index[i][0]
7     second_index = sort_distance_index[i][1]
8     a_s_sum = 0
9     b_s_sum = 0
10    #a
11    for oth_point in cluster[smallest_index]:
12        a_s_sum += np.linalg.norm(s_data - oth_point)
13    a_s = a_s_sum/len(cluster[smallest_index])
14    #b
15    for oth_point in cluster[second_index]:
16        b_s_sum += np.linalg.norm(s_data - oth_point)
17    b_s = b_s_sum/len(cluster[second_index])
18
19    s_s = (b_s - a_s)/max(a_s, b_s)
20    s_sum += s_s
21
22 s = s_sum/new_X.shape[0]
23 print(s)
```

0.4875992369330414

It is easy to know that $s \in (-1, 1)$, and larger s value indicates better clustering performance. So this classification is relatively good.

4.2 RAND INDEX

Given a set of n samples $S = o_1, o_2, \dots, o_n$, there are **two clusterings of S to compare**, including: X with r clusters and Y with s clusters. Let us define:

- **a**: The number of pairs of elements in S that are in the same subset in X and Y
 - **b**: The number of pairs of elements in S that are in the different subset in X and Y
 - **c**: The number of pairs of elements in S that are in the same subset in X but different in Y
 - **d**: The number of pairs of elements in S that are in the different subset in X but same in Y
- The rand index (RI) can be computed as follows:

$$RI = \frac{a + b}{a + b + c + d} = \frac{a + b}{\frac{(1+n)n}{2}}$$

Note that $RI \in [0, 1]$, and **higher score corresponds higher similarity**. In this task, we do not have two clustering, but the true label and PCA & K-means clustering. If we do rand index to compare them, we can evaluate the performance of our clustering (similarity with true label). Thus, we implement the rand index as followings:

```
1 a = 0
2 b = 0
3 label = txtDF[txtDF.columns[-1]]
4 for i in range(new_X.shape[0]):
5     for j in range(i+1, new_X.shape[0]):
6         if sort_distance_index[i][0] == sort_distance_index[j][0]
7             and label[i] == label[j]:
8             a +=1
9         elif sort_distance_index[i][0] != sort_distance_index[j][0]
10            and label[i] != label[j]:
11            b +=1
12 RI = 2* (a + b) / (new_X.shape[0] * (new_X.shape[0] - 1))
13 print(RI)
```

0.8743677375256322

Which has a high score, so the classification is good.

<END>