

KHAI THÁC DỮ LIỆU

Nguyễn Xuân Việt Đức - 22280012

Bài tập thực hành - Lần 6

Bài 6: SỰ PHÂN TÍCH NHÓM

III. Nội dung thực hành:

1. Thuật Toán K-means

K-means là một thuật toán phân cụm lặp, nhằm chia n quan sát thành k cụm sao cho mỗi quan sát thuộc về cụm có trung bình gần nhất (tâm cụm).

Các bước:

- **Khởi tạo:** Chọn ngẫu nhiên k điểm dữ liệu từ tập dữ liệu làm các tâm cụm ban đầu.
- **Gán cụm:** Gán mỗi điểm dữ liệu vào tâm cụm gần nhất. Khoảng cách thường dùng là khoảng cách Euclid và sai số được tính là tổng bình phương sai số (SSE) từ mỗi điểm đến tâm cụm của nó.
- **Cập nhật:** Tính lại các tâm cụm bằng trung bình của tất cả các điểm dữ liệu trong cụm đó.
- **Lặp:** Lặp lại hai bước trên cho đến khi tâm cụm không còn thay đổi đáng kể (nghĩa là sai số tổng thay đổi dưới một ngưỡng xác định) hoặc đạt đến số lần lặp tối đa.

Hàm sai số: Tổng căn bậc hai sai số bình phương:

```
1 rsserr(a, b) = np.sum(np.square(a - b))
```

Listing 1: Hàm sai số

Implementation:

1.1. Khởi tạo ngẫu nhiên centroid:

```
1 def initiate_centroids(k, dset):
2     '''
3     Select k data points as centroids
4     k: number of centroids
5     dset: pandas dataframe
6     '''
7     centroids = dset.sample(k)
8     return centroids
9
10 np.random.seed(42)
11 k = 3
12 df = blobs[['x', 'y']]
13 centroids = initiate_centroids(k, df)
14 print(centroids)
```

	x	y
0	24.412	32.932
5	25.893	31.515
36	26.878	36.609

1.2. Hàm tính sai số:

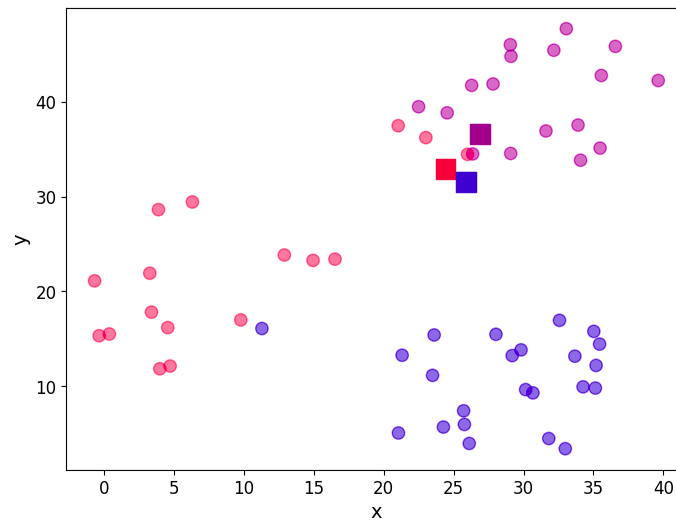
```
1 def rsserr(a, b):
2     '''
3     Calculate the root of sum of squared errors.
4     a and b are numpy arrays
5     '''
6     return np.sum(np.square(a - b))
7
8 for i, centroid in enumerate(range(centroids.shape[0])):
9     err = rsserr(centroids.iloc[centroid, :], df.iloc[36, :])
10    print('Error for centroid {0}: {1:.2f}'.format(i, err))
```

```
Error for centroid 0: 19.60
Error for centroid 1: 26.92
Error for centroid 2: 0.00
```

1.3. Hàm gán tâm cụm:

```
1 def centroid_assignment(dset, centroids):
2
3     k = centroids.shape[0]
4     n = dset.shape[0]
5     assignation = []
6     assign_errors = []
7
8     for obs in range(n):
9         all_errors = np.array([])
10        for centroid in range(k):
11            err = rsserr(centroids.iloc[centroid, :], dset.iloc[obs
12            , :])
13            all_errors = np.append(all_errors, err)
14
15        nearest_centroid = np.where(all_errors == np.amin(
16        all_errors))[0].tolist()[0]
17        nearest_centroid_error = np.amin(all_errors)
18
19        assignation.append(nearest_centroid)
20        assign_errors.append(nearest_centroid_error)
21
22    return assignation, assign_errors
23
24 df.loc[:, 'centroid'], df.loc[:, 'error'] = centroid_assignment(df
25 , centroids)
26 print(df.head())
```

Đồ thị phân tán cùng với cụm ngẫu nhiên:



```
print("The total error is {0:.2f}".format(df['error'].sum()))
```

The total error is 20606.95

1.4. Hàm tính K-means:

```

1 def kmeans(dset, k=2, tol=1e-4):
2     '''
3     K-means implementation for a
4     'dset': DataFrame with observations
5     'k': number of clusters, default k=2
6     'tol': tolerance=1E-4
7     '''
8     # Let us work in a copy, so we don't mess the original
9     working_dset = dset.copy()
10    # We define some variables to hold the error, the
11    # stopping signal and a counter for the iterations
12    err = []
13    goahead = True
14    j = 0
15
16    # Step 2: Initiate clusters by defining centroids
17    centroids = initiate_centroids(k, dset)
18
19    while(goahead):
20        # Step 3 and 4 - Assign centroids and calculate error
21        working_dset['centroid'], j_err = centroid_assignment(
22            working_dset, centroids)
23        err.append(sum(j_err))
24
25        # Step 5 - Update centroid position
26        centroids = working_dset.groupby('centroid').agg('mean').
27            reset_index(drop=True)
28
29        # Step 6 - Restart the iteration
30        if j > 0:
31            # Is the error less than a tolerance (1E-4)
32            if err[j-1] - err[j] <= tol:
33                goahead = False
34
35        j += 1
36
37    return working_dset['centroid'], j_err, centroids

```

```

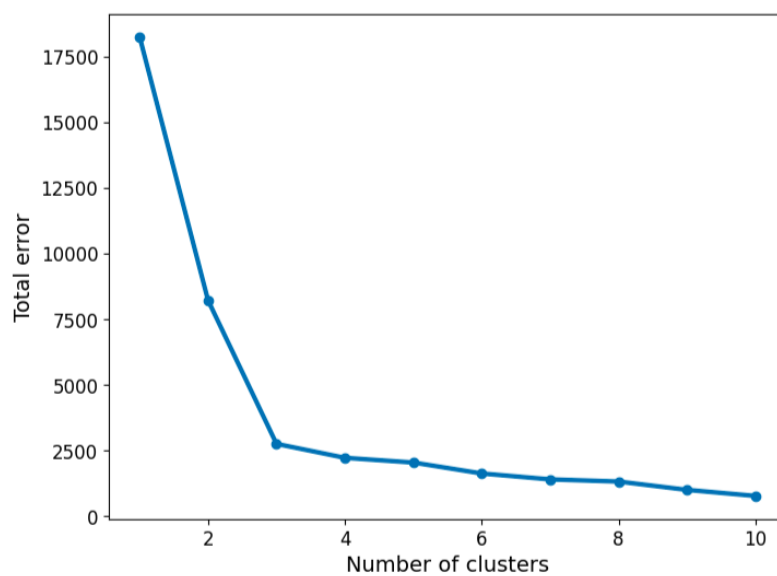
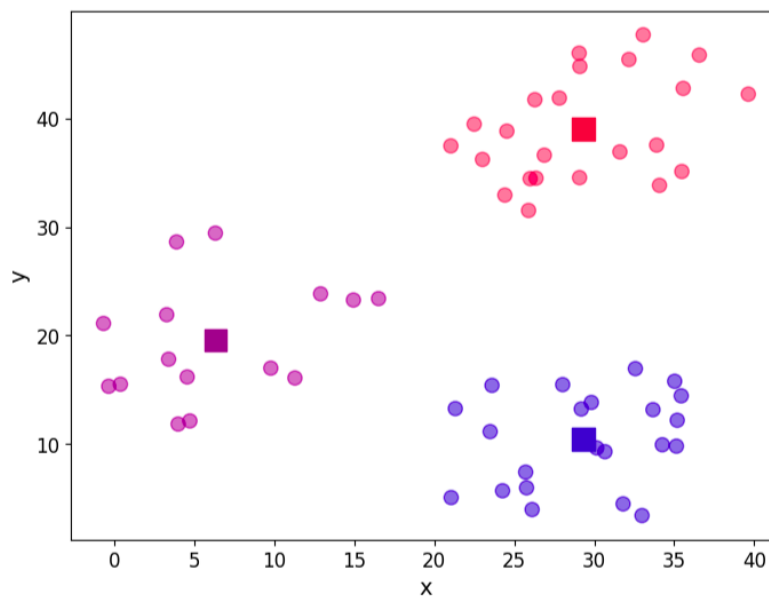
1 np.random.seed(432)
2 df['centroid'], df['error'], centroids = kmeans(df[['x', 'y']], 3)
3 print(df.head())

```

	x	y	centroid	error
0	24.412	32.932	0	61.380524
1	35.190	12.189	1	37.415091
2	26.288	41.718	0	16.216075
3	0.376	15.506	2	51.798518
4	26.116	3.963	1	52.188602

	x	y
0	29.304957	39.050783
1	29.330864	10.432409
2	6.322867	19.559800

1.5. Kết quả sau áp dụng K-means Biểu đồ khủ tay của k-means:



Từ biểu đồ trên ta thấy được khủ tay được xác định tại vị trí số lượng cụm là $3 \Rightarrow k = 3$, điều này hoàn toàn hợp lý với dữ liệu ban đầu.

2. Thuật Toán K-medians

K-medians tương tự K-means nhưng khác ở cách cập nhật tâm cụm và cách tính khoảng cách (và sai số). Thuật toán này tối thiểu hóa tổng khoảng cách từ các điểm đến trung vị của cụm. Nó thường bền vững hơn với các giá trị ngoại lai so với K-means.

Các bước:

- **Khởi tạo:** Chọn ngẫu nhiên k điểm dữ liệu làm tâm cụm ban đầu (giống phương pháp của K-means).
- **Gán cụm:** Gán mỗi điểm vào tâm cụm gần nhất. Khoảng cách dùng là khoảng cách Manhattan (Tổng độ lệch tuyệt đối - SAD).
- **Cập nhật:** Tính lại các tâm cụm bằng trung vị của tất cả các điểm trong cụm.
- **Lặp:** Lặp lại quá trình gán cụm và cập nhật cho đến khi các tâm cụm không còn thay đổi đáng kể hoặc đạt đến số lần lặp tối đa.

Hàm sai số: Tổng độ lệch tuyệt đối:

```
1 sad_err(a, b) = np.sum(np.abs(a - b))
```

Listing 2: Hàm sai số

Implementation:

2.1. Hàm tính sai số cho K-medians:

```
1 def sad_err(a, b):  
2     '''  
3     Calculate the sum of absolute differences (Manhattan distance).  
4     a and b are numpy arrays or pandas series  
5     '''  
6     return np.sum(np.abs(a - b))
```

2.2. Hàm gán tâm cụm cho thuật toán K-medians:

```
1  
2 def centroid_assignment_median(dset, centroids):  
3     '''  
4     Given a dataframe 'dset' and a set of 'centroids', we assign  
5     each  
6     data point in 'dset' to a centroid based on Sum of Absolute  
7     Differences.  
8     - dset: pandas dataframe with observations  
9     - centroids: pandas dataframe with centroids  
10    '''  
11    k = centroids.shape[0]  
12    n = dset.shape[0]  
13    assignation = []  
14    assign_errors = []
```

```

14     for obs_idx in range(n):
15         # Estimate error
16         all_errors = np.array([])
17         for centroid_idx in range(k):
18             err = sad_err(centroids.iloc[centroid_idx, :], dset.
19                             iloc[obs_idx, :])
20             all_errors = np.append(all_errors, err)
21
22         # Get the nearest centroid and the error
23         nearest_centroid = np.where(all_errors == np.amin(
24             all_errors))[0].tolist()[0]
25         nearest_centroid_error = np.amin(all_errors)
26
27         # Add values to corresponding lists
28         assignation.append(nearest_centroid)
29         assign_errors.append(nearest_centroid_error)
30
31     return assignation, assign_errors

```

2.3. Hàm thuật toán K-medians:

```

1 def kmedians(dset, k=2, tol=1e-4):
2     '''
3     K-medians implementation for a
4     'dset': DataFrame with observations
5     'k': number of clusters, default k=2
6     'tol': tolerance=1E-4
7     '''
8     working_dset = dset.copy()
9     err = []
10    goahead = True
11    j = 0
12    centroids = initiate_centroids(k, dset)
13
14    while(goahead):
15        working_dset['centroid'], j_err =
16        centroid_assignment_median(working_dset, centroids)
17        err.append(sum(j_err))
18
19        centroids = working_dset.groupby('centroid').median().
20        reset_index(drop=True)
21
22        if j > 0:
23            if err[j-1] - err[j] <= tol:
24                goahead = False
25            j += 1
26
27    return working_dset['centroid'], j_err, centroids
28
29 np.random.seed(432)
30 df_kmedians = blobs[['x', 'y']].copy()
31 df_kmedians['centroid'], df_kmedians['error'], centroids_kmedians =
32 kmedians(df_kmedians[['x', 'y']], 3)
33 print("Data points with assigned centroids (K-medians):")
34 print(df_kmedians.head())
35 print("\nFinal centroids (K-medians):")
36 print(centroids_kmedians)

```

2.4. Kết quả

Data points with assigned centroids (K-medians):

	x	y	centroid	error
0	24.412	32.932	0	9.2470
1	35.190	12.189	1	6.8565
2	26.288	41.718	0	6.9510
3	0.376	15.506	2	6.4780
4	26.116	3.963	1	10.4435

Final centroids (K-medians):

	x	y
0	29.056	37.5350
1	29.986	10.5365
2	4.550	17.8100

