

ĐẠI HỌC QUỐC GIA THÀNH PHỐ HỒ CHÍ MINH
TRƯỜNG ĐẠI HỌC KHOA HỌC TỰ NHIÊN
KHOA TOÁN - TIN HỌC



BTLT TUẦN 8: KHAI THÁC DỮ LIỆU

Sinh viên thực hiện: Nguyễn Công Hoài Nam
Mã số sinh viên: 21280099

Ngày 2 tháng 6 năm 2024

MỤC LỤC

1	Cài đặt thuật toán Mahalanobis k-means	1
2	Kết quả	3
2.1	Kết quả thuật toán	3
2.2	So sánh với GMM, K-means	4
3	Kết luận	6

1. Cài đặt thuật toán Mahalanobis k-means

Import các thư viện cần thiết

```
1 from sklearn.mixture import GaussianMixture
2 from sklearn.datasets import make_blobs
3 import matplotlib.pyplot as plt
4 import numpy as np
5 import pandas as pd
6 from sklearn.cluster import KMeans
7 import seaborn as sns
```

Các bước để cài đặt thuật toán Mahalanobis K-means

Bước 1: Khởi tạo các phân cụm (clusters) và gán chúng ứng với mỗi điểm dữ liệu

- Mahalanobis k-means cần các điểm dữ liệu được gán với các centroid từ ban đầu bởi vì nó cần chúng để tính ma trận hiệp sai cho các clusters để tính khoảng cách Mahalanobis
- Ở đây em sử dụng khoảng cách Euclidean để khởi tạo (thuật toán k-means)

```
1 def initiate_centroids_clusters(k,dset, random_state = 234):
2     km=KMeans(n_clusters = k,max_iter = 1,random_state = random_state)
3
4     # Lấy iter = 1 , nghĩa là chỉ khởi tạo bước lặp đầu tiên bằng thuật toán kmeans
5     km.fit(dset)
6     clusters=km.labels_
7
8     centroids=pd.DataFrame(km.cluster_centers_,columns=km.feature_names_in_)
9
10    return clusters, centroids
```

Bước 2: tính khoảng cách Mahalanobis từ các centroids đến các điểm dữ liệu và gán mỗi điểm dữ liệu ứng với centroid gần nhất

- Hàm tính khoảng cách Mahalanobis

$$D_{Maha}(X, \mu_r, \Sigma_r) = \sqrt{(X - \mu_r)\Sigma^{-1}(X - \mu_r)^T}$$

- X: mỗi điểm dữ liệu trong dataset (các hàng trong df)
- μ_r : tâm cụm r trong danh sách tâm cụm
- Σ_r : ma trận hiệp phương sai ứng với mỗi tâm cụm r

```
1 def mahalanobis_error(X, mu, sigma):
2     '''
3     D(X,mu,sigma)^2 = (X-mu) * Sigma^-1 * (X - mu).T
4     '''
5     sigma_inv = np.linalg.inv(sigma)
6     return np.sqrt((X - mu) @ sigma_inv @ (X - mu).T)
```

- Để tính được khoảng cách Mahalanobis ta cần ma trận hiệp phương sai, dưới đây là hàm để tính:

```
1 def calculate_covariance_matrices(dset,k):
2     '''
3     Trả về các cov_matrix của các cụm
4     '''
5     cov_matrices = []
6     for centroid in range(k):
7         cov_matrix = np.cov(dset.loc[dset['centroid']==centroid,['x','y']],rowvar=False)
```

```

8     cov_matrices.append(cov_matrix)
9     return cov_matrices

```

- Tính lại ma trận hiệp phương sai và khoảng cách Mahalanobis của các centroids đến từng điểm dữ liệu - Gán lại phân cụm cho các điểm dữ liệu ứng với phân cụm gần nhất

```

1  def centroid_assignment(dset, centroids):
2      # Số centroid
3      k = centroids.shape[0]
4
5      # Số samples trong tập dữ liệu
6      n = dset.shape[0]
7
8      assignation = []
9
10     assign_errors = []
11
12     cov_matrixs = []
13
14     # Tính ma trận hiệp phương sai cho từng cụm
15     cov_matrices = calculate_covariance_matrices(dset, k)
16
17     for obs in range(n):
18         # List error cho từng centroid
19         all_errors = np.array([])
20
21         # Tính error cho từng centroid
22         for centroid in range(k):
23             err = mahalanobis_error(centroids.iloc[centroid, :], dset.iloc[obs, :2],
24                                     cov_matrices[centroid])
25             all_errors = np.append(all_errors, err)
26         # Chọn centroid gần nhất
27         nearest_centroid = np.where(all_errors == np.amin(all_errors))[0].tolist()[0]
28         nearest_centroid_error = np.amin(all_errors)
29
30         assignation.append(nearest_centroid)
31         assign_errors.append(nearest_centroid_error)
32
33     return assignation, assign_errors

```

Bước 3: lặp lại bước hai cho đến khi các phân cụm không còn thay đổi

- Hàm chính của thuật toán,

```

1  def mahalanobis_kmeans(dset, k = 2, tol = 1e-4, plot = False):
2      # Tạo bản sao để tránh thay đổi data gốc
3      working_dset = dset.copy()
4
5      # Danh sách error cho từng bước lặp
6      err = []
7
8      # Biến điều kiện để tiếp tục vòng lặp
9      goahead = True
10
11     # Biến đếm số lần lặp
12     j = 0
13
14     # Gán giá trị khởi tạo cho centroids
15     working_dset['centroid'], centroids = initiate_centroids_clusters(k, working_dset)

```

```

16 while goahead:
17     # Gán centroid và tính error cho các quan sát
18     working_dset['centroid'], j_err = centroid_assignment(working_dset, centroids)
19
20     # Tính tổng error và vẽ đồ thị
21     err.append(sum(j_err))
22
23     # Tính các centroid mới dựa trên các quan sát đã gán nhãn
24     centroids = working_dset.groupby('centroid').agg('mean').reset_index(drop=True)
25
26     if j > 0:
27         if err[j-1] - err[j] <= tol:
28             goahead = False
29
30     j += 1
31
32 if plot:
33     fig, ax = plt.subplots(1, 1, figsize=(12, 8))
34     show_plot(ax, working_dset, centroids)
35
36 return working_dset['centroid'], j_err, centroids

```

- Để thể hiện kết quả cuối cùng ta viết thêm hàm để trực quan hoá các phân cụm

```

1 def show_plot(ax, df, centroids, title="After Mahalanobis K-means"):
2     k = centroids.shape[0]
3     scatter_points = ax.scatter(df.iloc[:, 0], df.iloc[:, 1], marker='o', c=df['centroid'].
4                                 astype('category'), cmap='viridis', s=80, alpha=0.5)
5     scatter_centroids = ax.scatter(centroids.iloc[:, 0], centroids.iloc[:, 1], marker='X', s
6                                   =250, c=[i for i in range(k)], cmap='viridis', edgecolors='deeppink')
7     ax.set_xlabel('x', fontsize=14)
8     ax.set_ylabel('y', fontsize=14)
9     ax.set_title(title, fontsize=16)
10    ax.tick_params(axis='both', which='major', labelsize=12)
11
12    legend_labels = [f'Centroid {i+1}: ({centroids.iloc[i, 0]:.2f}, {centroids.iloc[i, 1]:.2f
13                                })' for i in range(k)]
14    legend_elements = scatter_centroids.legend_elements()[0]
15    ax.legend(legend_elements, legend_labels, title='Centroids', loc='upper left')

```

2. Kết quả

2.1. Kết quả thuật toán

Chạy thuật toán với data ban đầu (từ file pdf)

```

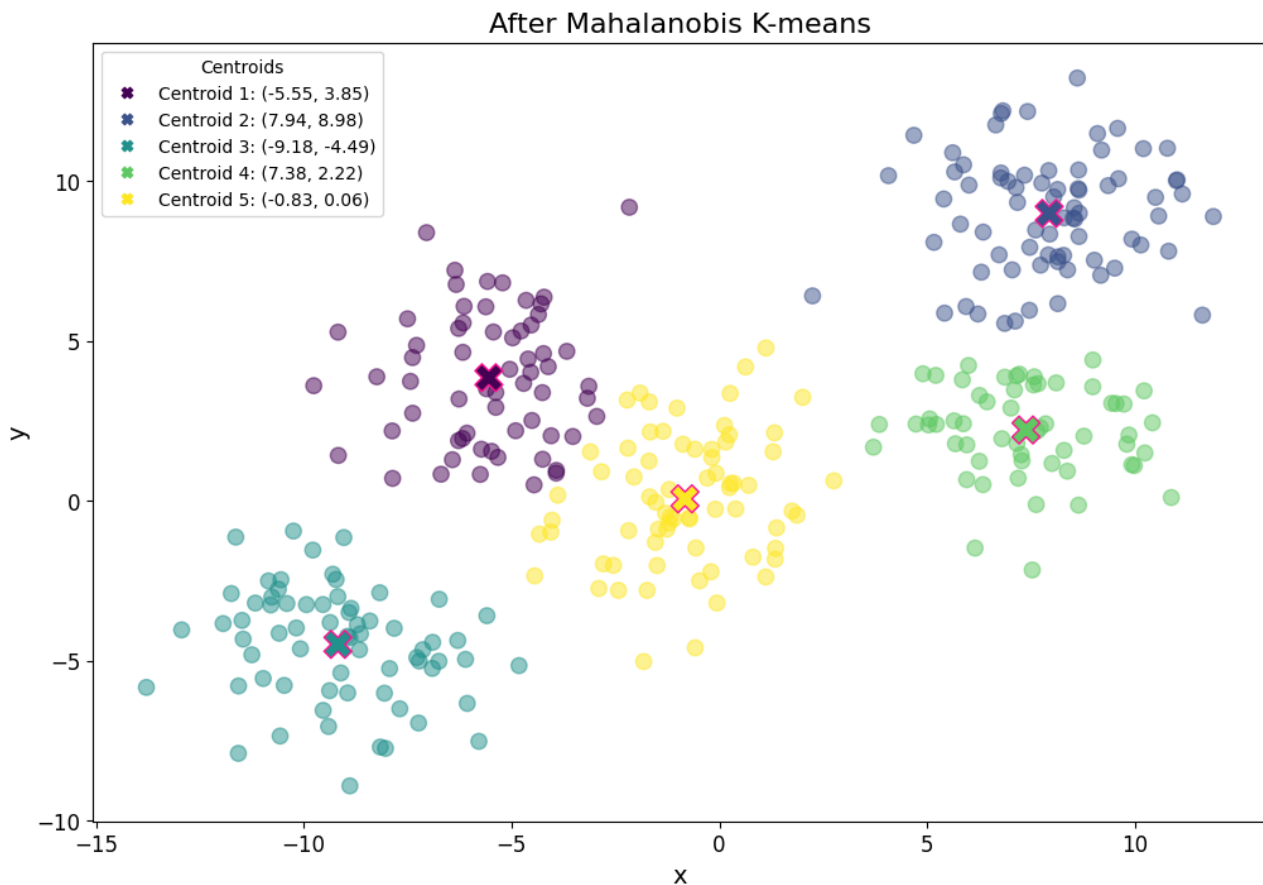
1 np.random.seed(234)
2 x,y=make_blobs(n_samples=330, centers=5,cluster_std=1.84)
3 df = pd.DataFrame({'x':x[:,0], 'y':x[:,1]})

```

```

1 df['centroid'],df['error'],centroids=mahalanobis_kmeans(df,k=5, plot = True)

```



2.2. So sánh với GMM, K-means

Với những bộ dữ liệu có sự tương quan giữa các biến, việc tính toán chúng như những biến độc lập có thể dẫn đến sai số (thuật toán K-means)

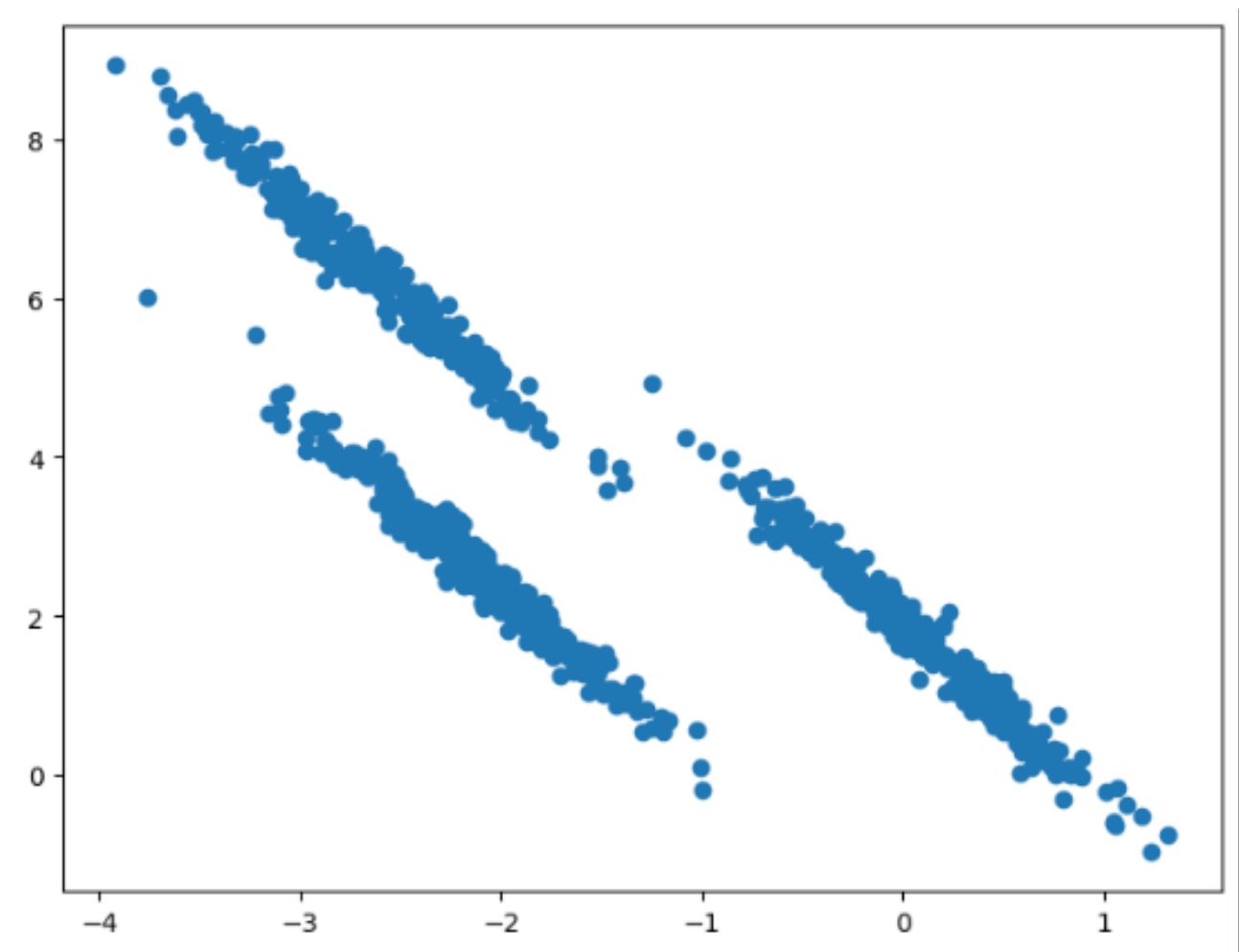
- Do đó chúng ta cần những thuật toán quan tâm đến sự tương quan của các biến
- Đó chính là GMM và Mahalanobis K-means

Ta khởi tạo một bộ data có sự tương quan để thấy rõ sự khác biệt giữa các thuật toán

```

1 x, y = make_blobs(n_samples=1000, random_state=234, centers=3)
2 x = np.dot(x, [[0.2, -0.6], [-0.4, 0.8]])
3 df=pd.DataFrame({'x':x[:,0], 'y':x[:,1]})
4
5 plt.figure(figsize=(8,6))
6 plt.scatter(x[:,0],x[:,1])
7 plt.show()

```



Tính toán GMM và K-means của thư viện Sklearn

```

1 # Apply standard KMeans
2 kmeans = KMeans(n_clusters=3, random_state=234)
3 df['centroid'] = kmeans.fit_predict(df[['x', 'y']])
4 centroids_kmeans = pd.DataFrame(kmeans.cluster_centers_, columns=['x', 'y'])
5
6 # Apply GMM
7 gmm = GaussianMixture(n_components=3, random_state=234)
8 df['centroid_gmm'] = gmm.fit_predict(df[['x', 'y']])
9 centroids_gmm = pd.DataFrame(gmm.means_, columns=['x', 'y'])

```

Vẽ biểu đồ trực quan kết quả

```

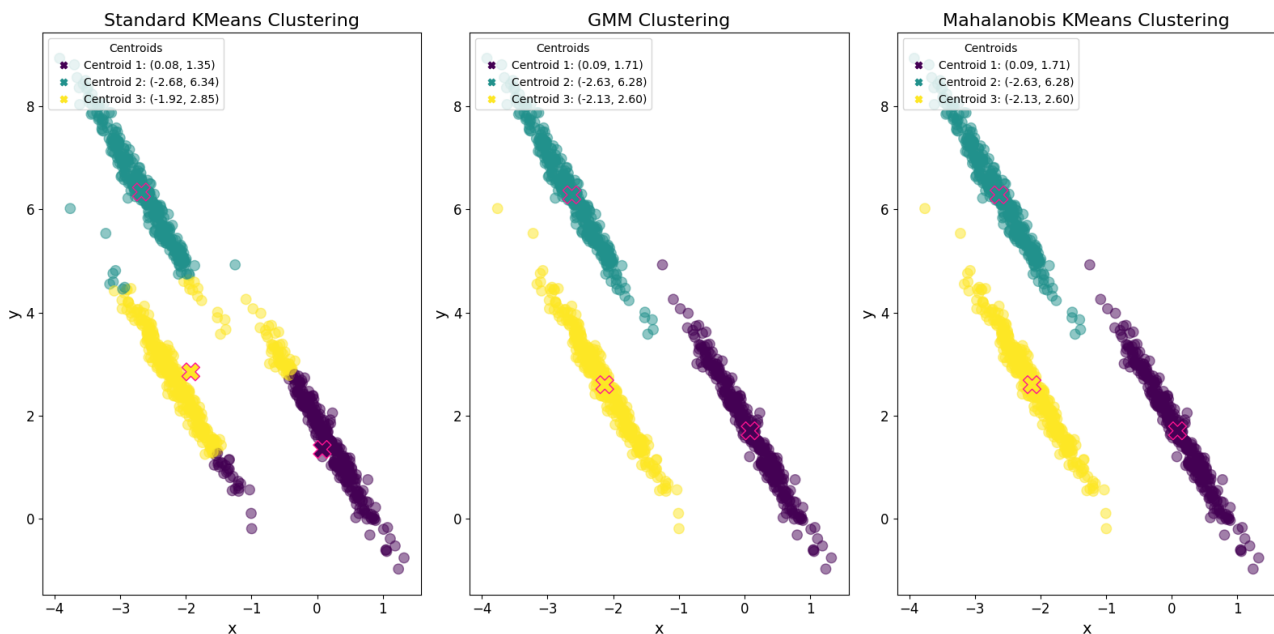
1 # Mahalanobis KMeans (assuming mahalanobis_kmeans is defined)
2 df['centroid_maha'], _, centroids_maha_kmeans = mahalanobis_kmeans(df[['x', 'y']], k=3, plot=
    False)
3
4 fig, axs = plt.subplots(1, 3, figsize=(16, 8))
5
6 # Plot for Standard KMeans
7 df['centroid'] = kmeans.fit_predict(df[['x', 'y']])
8 show_plot(axs[0], df[['x', 'y', 'centroid']], centroids_kmeans, title='Standard KMeans
    Clustering')
9
10 # Plot for GMM

```

```

11 df['centroid'] = df['centroid_gmm'] # Assign centroid_gmm to centroid for consistent
    plotting
12 show_plot(axes[1], df[['x', 'y', 'centroid']], centroids_gmm, title='GMM Clustering')
13
14 # Plot for Mahalanobis KMeans
15 df['centroid'] = df['centroid_maha'] # Assign centroid_maha to centroid for consistent
    plotting
16 show_plot(axes[2], df[['x', 'y', 'centroid']], centroids_maha_kmeans, title='Mahalanobis
    KMeans Clustering')
17
18 plt.tight_layout()
19 plt.show()

```



Có thể thấy GMM và Mahalanobis phân cụm tốt hơn nhiều so với K-means

3. Kết luận

Khi dữ liệu có tương quan giữa các biến, các phương pháp phân cụm như GMM (Gaussian Mixture Models) và Mahalanobis K-means thường vượt trội so với thuật toán K-means truyền thống.

GMM có khả năng mô hình hóa các cụm với hình dạng và kích thước khác nhau nhờ việc sử dụng các phân phối Gaussian. Hơn nữa, GMM có thể tính toán và mô hình hóa ma trận hiệp phương sai của dữ liệu, do đó có thể phát hiện các cấu trúc tương quan giữa các biến. Điều này cho phép GMM phân chia dữ liệu thành các cụm phức tạp hơn so với K-means.

Mahalanobis K-means sử dụng khoảng cách Mahalanobis tính đến cả phương sai và tương quan của các biến, do đó có thể phát hiện các cụm có hình dạng và kích thước khác nhau một cách hiệu quả hơn.

Tuy nhiên, với dữ liệu có các biến độc lập, K-means vẫn là sự lựa chọn hợp lý nhờ sự đơn giản và hiệu quả của thuật toán.

Vì vậy, việc lựa chọn thuật toán phân cụm phụ thuộc vào đặc điểm của dữ liệu, cần phải đánh giá ưu và nhược điểm của từng phương pháp để quyết định sử dụng phù hợp.