

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



BTLT TUẦN 6: KHAI THÁC DỮ LIỆU

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

Ngày 19 tháng 5 năm 2024

MỤC LỤC

1	Thuật toán K-means	1
1.1	Tiền xử lý	1
1.2	Cài đặt thuật toán	2
1.3	Kết quả	7
2	Thuật toán K-medians	9
2.1	Cài đặt thuật toán	9
2.2	Kết quả	10

1. Thuật toán K-means

1. Tiền xử lý

Tạo file `data.csv`

Import thư viện và đọc dữ liệu từ file

```

1 import numpy as np
2 import pandas as pd
3 import matplotlib.pyplot as plt
4 from matplotlib.colors import ListedColormap
5
6 blobs = pd.read_csv('data.csv')
7 colnames = list(blobs.columns[1:-1])

```

Dữ liệu có dạng:

```

1 blobs.head()

```

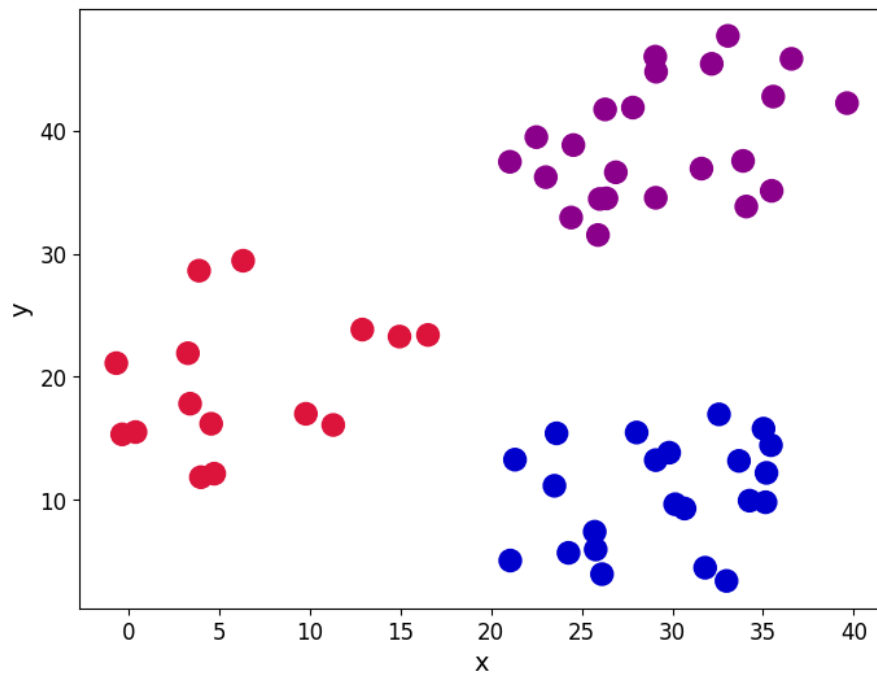
	ID	x	y	cluster
0	0	24.412	32.932	2
1	1	35.190	12.189	1
2	2	26.288	41.718	2
3	3	0.376	15.506	0
4	4	26.116	3.963	1

Chúng ta sẽ sử dụng cột phân cụm để hiển thị các nhóm khác nhau được biểu diễn trong tập dữ liệu

```

1 customcmap = ListedColormap(["crimson", "mediumblue", "darkmagenta"])
2 fig, ax = plt.subplots(figsize=(8,6))
3 plt.scatter(x = blobs['x'], y = blobs['y'], s=150,
4             c=blobs['cluster'].astype('category'),
5             cmap = customcmap)
6 ax.set_xlabel(r'x', fontsize = 14)
7 ax.set_ylabel(r'y', fontsize = 14)
8 plt.xticks(fontsize = 12)
9 plt.yticks(fontsize = 12)
10 plt.show()

```



2. Cài đặt thuật toán

Bước 1 và 2: Xác định k và khởi tạo các tâm

```

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 centroids

```

Bước 3: tính khoảng cách

Công thức: $Dist(a, b) = \sum (a - b)^2$ với a, b là các điểm dữ liệu

```

1  def rsserr(a, b):
2      '''
3      Calculate sum of squared errors.
4      a and b are numpy arrays
5      '''
6      return np.sum(np.square(a-b))

```

	x	y
0	24.412	32.932
5	25.893	31.515
36	26.878	36.609

Bước 4: gán giá trị các tâm

```

1 def centroid_assignment(dset, centroids):
2     """
3     Given a dataframe `dset` and a set of `centroids`, we assign each
4     data point in `dset` to a centroid.
5     - dset - pandas dataframe with observations
6     - centroids - pandas dataframe with centroids
7     """
8     k = centroids.shape[0]
9     n = dset.shape[0]
10    assignation = []
11    assign_errors = []
12
13    for obs in range(n):
14        # Estimate error
15        all_errors = np.array([])
16        for centroid in range(k):
17            err = rsserr(centroids.iloc[centroid, :], dset.iloc[obs,:])
18            all_errors = np.append(all_errors, err)
19
20        # Get the nearest centroid and the error
21        nearest_centroid = np.where(all_errors==np.amin(all_errors))[0].tolist()[0]
22        nearest_centroid_error = np.amin(all_errors)
23
24        # Add values to corresponding lists
25        assignation.append(nearest_centroid)
26        assign_errors.append(nearest_centroid_error)
27
28    return assignation, assign_errors

```

Thêm cột gán tâm và sai số phát sinh để cập nhật biểu đồ phân tán biểu diễn các trọng tâm

```

1 df['centroid'], df['error'] = centroid_assignment(df, centroids)
2 df.head()

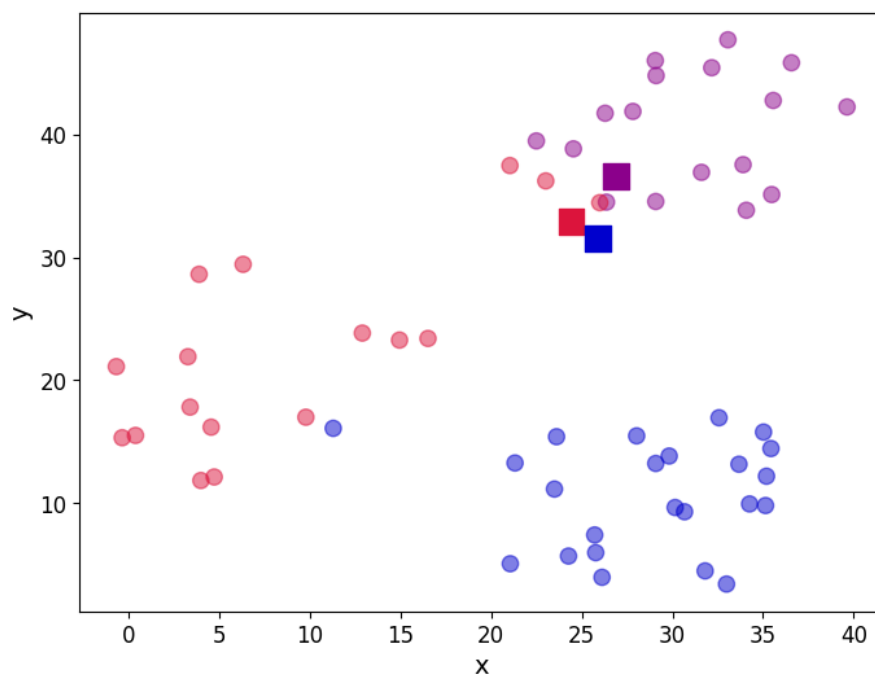
```

	x	y	centroid	error
0	24.412	32.932	0	0.000000
1	35.190	12.189	1	459.928485
2	26.288	41.718	2	26.449981
3	0.376	15.506	0	881.394772
4	26.116	3.963	1	759.162433

```

1 fig, ax = plt.subplots(figsize = (8,6))
2 plt.scatter(df.iloc[:,0], df.iloc[:,1], marker = 'o',
3             c=df['centroid'].astype('category'),
4             cmap = customcmap, s=80, alpha = 0.5)
5 plt.scatter(centroids.iloc[:,0], centroids.iloc[:,1],
6             marker = 's', s= 200, c=[0, 1, 2],
7             cmap = customcmap)
8 ax.set_xlabel(r'x', fontsize = 14)
9 ax.set_ylabel(r'y', fontsize = 14)
10 plt.xticks(fontsize = 12)
11 plt.yticks(fontsize = 12)
12 plt.show()

```



Tổng các sai số

```

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

```

The total error is **20606.25**

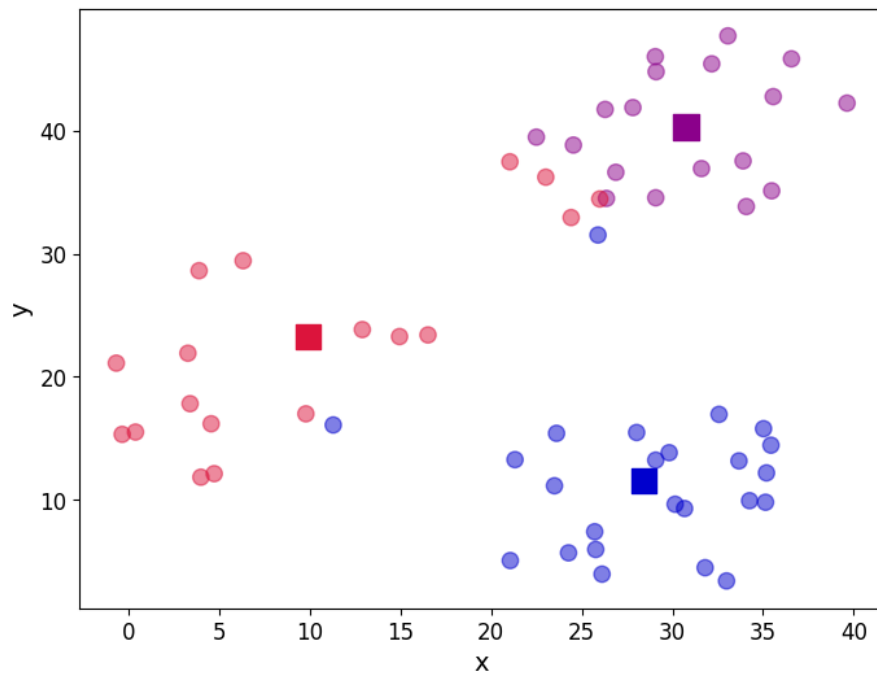
Bước 5: cập nhật vị trí của k tâm bằng việc tính giá trị trung bình của các quan sát được gán cho mỗi tâm

```
1 centroids = df.groupby('centroid').agg('mean').loc[:, colnames].reset_index(drop = True)
2 centroids
```

	x	y
0	9.889444	23.242611
1	28.431250	11.546250
2	30.759333	40.311167

Xem lại biểu đồ phân tán với vị trí của k tâm đã được cập nhật

```
1 fig, ax = plt.subplots(figsize = (8,6))
2 plt.scatter(df.iloc[:,0], df.iloc[:,1], marker = 'o',
3             c=df['centroid'].astype('category'),
4             cmap = customcmap, s=80, alpha = 0.5)
5 plt.scatter(centroids.iloc[:,0], centroids.iloc[:,1],
6             marker = 's', s= 200, c=[0, 1, 2],
7             cmap = customcmap)
8 ax.set_xlabel(r'x', fontsize = 14)
9 ax.set_ylabel(r'y', fontsize = 14)
10 plt.xticks(fontsize = 12)
11 plt.yticks(fontsize = 12)
12 plt.show()
```

**Bước 6:** lặp lại bước 3-5

```

1  def kmeans(dset, k=2, tol=1e-4):
2      '''
3      K-means implementationd 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(working_dset, centroids)
22         err.append(sum(j_err))
23
24         # Step 5 - Update centroid position
25         centroids = working_dset.groupby('centroid').agg('mean').reset_index(drop = True)
26
27         # Step 6 - Restart the iteration
28         if j>0:
29             # Is the error less than a tolerance (1E-4)
30             if err[j-1]-err[j]<=tol:
31                 goahead = False
32             j+=1
33
34     working_dset['centroid'], j_err = centroid_assignment(working_dset, centroids)
35     centroids = working_dset.groupby('centroid').agg('mean').reset_index(drop = True)
36     return working_dset['centroid'], j_err, centroids

```


3. Kết quả

Thực thi các hàm trên

```
1 np.random.seed(42)
2 df['centroid'], df['error'], centroid = kmeans(df[['x', 'y']], 3)
3 df.head()
```

	x	y	centroid	error
0	24.412	32.932	2	61.380524
1	35.190	12.189	1	37.472641
2	26.288	41.718	2	16.216075
3	0.376	15.506	0	51.798518
4	26.116	3.963	1	52.157062

Vị trí của các tâm cuối cùng

```
1 centroids
```

	x	y
0	9.889444	23.242611
1	28.431250	11.546250
2	30.759333	40.311167

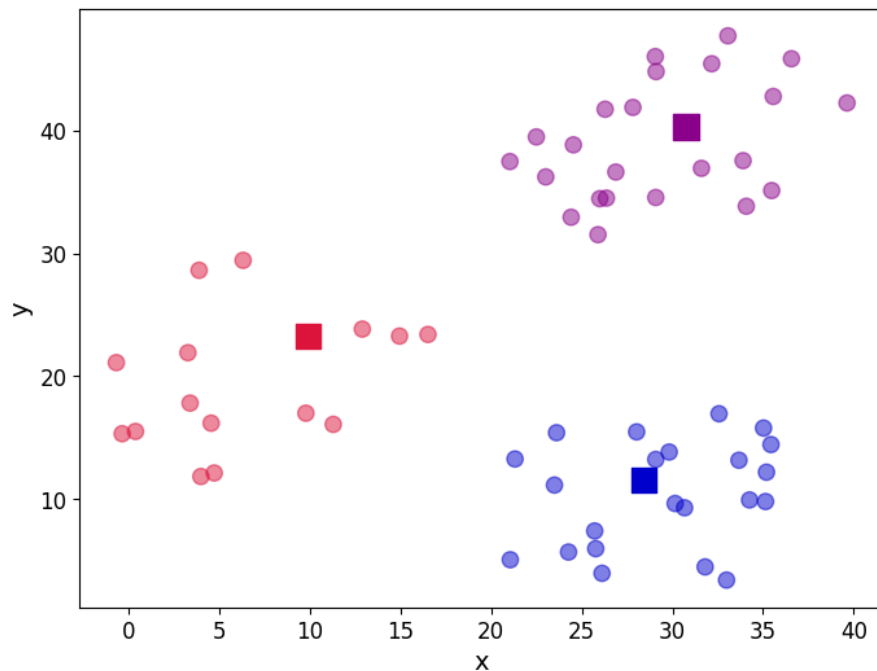
Xem lại biểu đồ phân tán

```
1 fig, ax = plt.subplots(figsize = (8,6))
2 plt.scatter(df.iloc[:,0], df.iloc[:,1], marker = 'o',
3             c=df['centroid'].astype('category'),
4             cmap = customcmap, s=80, alpha = 0.5)
5 plt.scatter(centroids.iloc[:,0], centroids.iloc[:,1],
6             marker = 's', s= 200, c=[0, 1, 2],
7             cmap = customcmap)
8 ax.set_xlabel(r'x', fontsize = 14)
```

```

9  ax.set_ylabel(r'y', fontsize = 14)
10 plt.xticks(fontsize = 12)
11 plt.yticks(fontsize = 12)
12 plt.show()

```



Sử dụng **elbow** để chỉ ra số cụm tối ưu

```

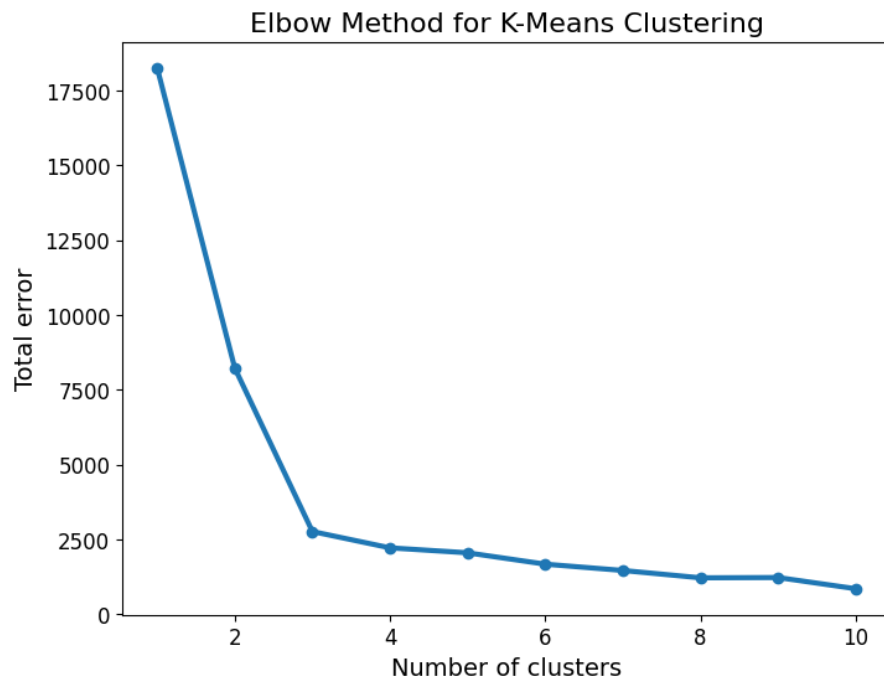
1  err_total = []
2  n = 10
3  df_elbow = blobs[['x', 'y']]
4  for i in range(n):
5      _, my_errs, _ = kmeans(df_elbow, i+1)
6      err_total.append(sum(my_errs))

```

```

1  fig, ax = plt.subplots(figsize = (8,6))
2  plt.plot(range(1,n+1), err_total, linewidth=3, marker = 'o')
3  ax.set_xlabel(r'Number of clusters', fontsize = 14)
4  ax.set_ylabel(r'Total error', fontsize = 14)
5  ax.set_title('Elbow Method for K-Means Clustering', fontsize=16)
6  plt.xticks(fontsize = 12)
7  plt.yticks(fontsize = 12)
8  plt.show()

```



Số cụm tối ưu là 3 (chỗ khuỷu tay)

2. Thuật toán K-medians

Về cơ bản, K-medians chỉ khác K-means ở:

- Cách tính khoảng cách (sử dụng khoảng cách Mahattan thay vì Euclid)
- Cách cập nhật trọng tâm (dùng median thay cho mean)

1. Cài đặt thuật toán

Lấy ngẫu nhiên trọng tâm ban đầu, tương tự như với hàm `inititate_means` của `kmeans`

```

1 def initiate_medians(k, dset):
2     '''
3     select k data points as medians
4     k: number of medians
5     dset: pandas dataframe
6     '''
7     centroids = dset.sample(k)
8     return centroids

```

Hàm tính khoảng cách, khác với `kmeans`, `kmedians` sử dụng khoảng cách Manhattan

Công thức: $Dist(a, b) = \sum |a - b|$ với a, b là các điểm dữ liệu

```

1 def rsserr(a,b):
2     '''
3     Calculate sum of squared errors.
4     a and b are numpy arrays
5     '''
6     return np.sum(np.abs((a-b)))

```

Gán các giá trị tâm, tương tự `kmeans`

```

1 def median_assignment(dset, centroids):

```

```

2      '''
3      Given a dataframe `dset` and a set of `medians`, we assign each
4      data point in `dset` to a median.
5      - dset - pandas dataframe with observations
6      - medians - pandas dataframe with medians
7      '''
8      k = centroids.shape[0]
9      n = dset.shape[0]
10     assignation = []
11     assign_errors = []
12
13     for obs in range(n):
14         all_errors = np.array([])
15         for centroid in range(k):
16             err = rsserr(centroids.iloc[centroid, :], dset.iloc[obs,:])
17             all_errors = np.append(all_errors, err)
18
19         nearest_centroid = np.where(all_errors==np.amin(all_errors))[0].tolist()[0]
20         nearest_centroid_error = np.amin(all_errors)
21
22         assignation.append(nearest_centroid)
23         assign_errors.append(nearest_centroid_error)
24
25     return assignation, assign_errors

```

Khác với kmeans, kmedians sử dụng trung vị (median) để tính lại các centroid

```

1  def kmedians(dset, k=2, tol=1e-4):
2      working_dset = dset.copy()
3      err = []
4      goahead = True
5      j = 0
6      centroids = initiate_medians(k, working_dset)
7
8      while goahead:
9          working_dset['centroid'], j_err = median_assignment(working_dset.iloc[:, :2],
10                                                             centroids)
11          err.append(sum(j_err))
12          centroids = working_dset.groupby('centroid').agg('median').reset_index(drop=True)
13
14          if j > 0:
15              if err[j-1] - err[j] <= tol:
16                  goahead = False
17              j += 1
18
19      return working_dset['centroid'], j_err, centroids

```

2. Kết quả

Kết quả hàm

```

1  np.random.seed(42)
2  df['centroid'], df['error'], centroids = kmedians(df[['x','y']], 3)
3  print('The total error is {:.2f}'.format(df['error'].sum()))
4  df.head()

```

The total error is 486.99

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

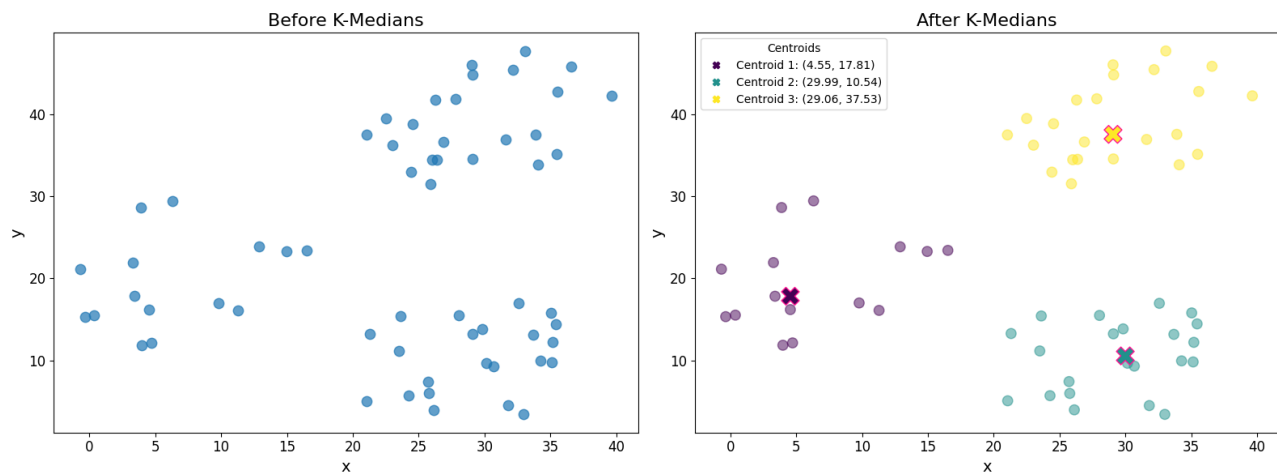
Vẽ biểu đồ để trực quan kết quả vừa tìm được

```

1  def plot_kmedians(df, centroids):
2      k = centroids.shape[0]
3      plt.figure(figsize=(16, 6))
4
5      plt.subplot(1, 2, 1)
6      plt.scatter(df.iloc[:, 0], df.iloc[:, 1], marker='o', s=80, alpha=0.7)
7      plt.xlabel('x', fontsize=14)
8      plt.ylabel('y', fontsize=14)
9      plt.title('Before K-Medians', fontsize=16)
10     plt.xticks(fontsize=12)
11     plt.yticks(fontsize=12)
12
13     plt.subplot(1, 2, 2)
14     plt.scatter(df.iloc[:, 0], df.iloc[:, 1], marker='o', c=df['centroid'].astype('category')
15                 , cmap='viridis', s=80, alpha=0.5)
16     scatter_centroids = plt.scatter(centroids.iloc[:, 0], centroids.iloc[:, 1], marker='X', s
17                                     =250, c=[i for i in range(k)], cmap='viridis', edgecolors='deeppink')
18     plt.xlabel('x', fontsize=14)
19     plt.ylabel('y', fontsize=14)
20     plt.title('After K-Medians', fontsize=16)
21     plt.xticks(fontsize=12)
22     plt.yticks(fontsize=12)
23
24     legend_labels = [f'Centroid {i+1}: ({centroids.iloc[i, 0]:.2f}, {centroids.iloc[i, 1]:.2f
25                     })' for i in range(k)]
26     legend_elements = scatter_centroids.legend_elements()[0]
27     plt.legend(legend_elements, legend_labels, title='Centroids', loc='upper left')
28
29     plt.tight_layout()
30     plt.show()

```

plot_kmedians(df, centroids)

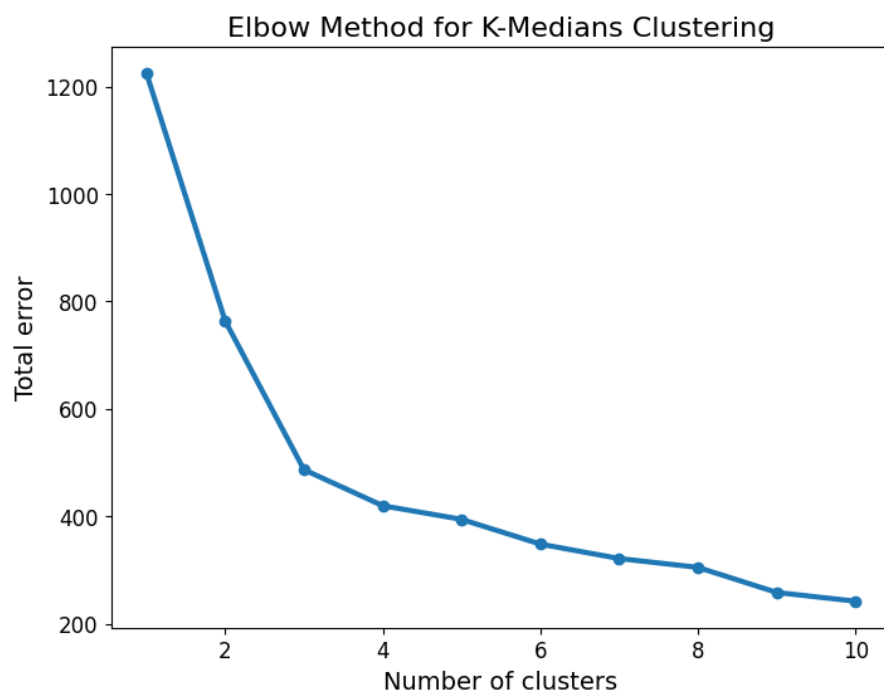


Cùng khám phá biểu đồ elbow để chọn số cluster thích hợp

```

1 err_total=[]
2 n = 10
3 df_elbow=blobs[['x','y']]
4 for i in range(n):
5     _,my_errs,_= kmedians(df_elbow,k=i+1)
6     err_total.append(sum(my_errs))
7
8 fig,ax=plt.subplots(figsize=(8,6))
9 plt.plot(range(1,n+1),err_total,linewidth=3,marker='o')
10 ax.set_xlabel(r'Number of clusters',fontsize=14)
11 ax.set_ylabel(r'Total error',fontsize=14)
12 ax.set_title('Elbow Method for K-Medians Clustering', fontsize=16)
13 plt.xticks(fontsize=12)
14 plt.yticks(fontsize=12)
15 plt.show()

```



Số cluster phù hợp là 3.