

ĐẠI HỌC QUỐC GIA THÀNH PHỐ HỒ CHÍ MINH
TRƯỜNG ĐẠI HỌC BÁCH KHOA
KHOA KHOA HỌC - KỸ THUẬT MÁY TÍNH



BÁO CÁO BÀI TẬP LỚN

Môn học: **Học máy và Ứng dụng**

Đề tài:

**Giải thuật K-means với các kỹ thuật
khởi tạo k trung tâm ban đầu và
Ứng dụng trong bài toán phân cụm**

GVHD: Assoc. Prof. Dr. Dương Tuấn Anh

Nhóm lớp: SDH-12

SVTH: Nguyễn Hoàng Phúc
1927030

TP. HỒ CHÍ MINH, THÁNG 6/2021

Mục lục

1	Tìm hiểu về giải thuật K-means	2
1.1	Ý tưởng giải thuật K-means	2
1.2	Tìm hiểu ý tưởng các kỹ thuật khởi tạo k trung tâm ban đầu	4
1.2.1	Phương pháp Ngẫu nhiên	4
1.2.2	Phương pháp K-MEANS++	5
1.2.3	Phương pháp FORGY APPROACH (FA)	6
1.2.4	Phương pháp MACQUEEN APPROACH (MA)	6
1.2.5	Phương pháp KAUFMAN APPROACH (KA)	7
2	Ứng dụng giải thuật K-means với các kỹ thuật khởi tạo k trung tâm ban đầu trong việc phân cụm tập dữ liệu IRIS FLOWER	9
2.1	Mô tả tập dữ liệu	9
2.2	Phân loại dữ liệu với các kỹ thuật khởi tạo k trung tâm ban đầu . .	11
2.2.1	Phương pháp Ngẫu nhiên	11
2.2.2	Phương pháp K-MEANS++	14
2.2.3	Phương pháp FORGY APPROACH (FA)	16
2.2.4	Phương pháp MACQUEEN APPROACH (MA)	19
2.2.5	Phương pháp KAUFMAN APPROACH (KA)	21
2.3	Nhận xét	22
3	Ứng dụng giải thuật K-means với các kỹ thuật khởi tạo k trung tâm ban đầu trong việc phân cụm tập dữ liệu THE MNIST	24
3.1	Mô tả tập dữ liệu	24
3.2	Phân loại dữ liệu với các kỹ thuật khởi tạo k trung tâm ban đầu . .	27
3.2.1	Phương pháp Ngẫu nhiên	27
3.2.2	Phương pháp K-MEANS++	27
3.2.3	Phương pháp FORGY APPROACH (FA)	28
3.2.4	Phương pháp MACQUEEN APPROACH (MA)	29
3.2.5	Phương pháp KAUFMAN APPROACH (KA)	29
3.3	Nhận xét	30
4	Tham khảo	31
5	Phụ lục	32

1 Tìm hiểu về giải thuật K-means

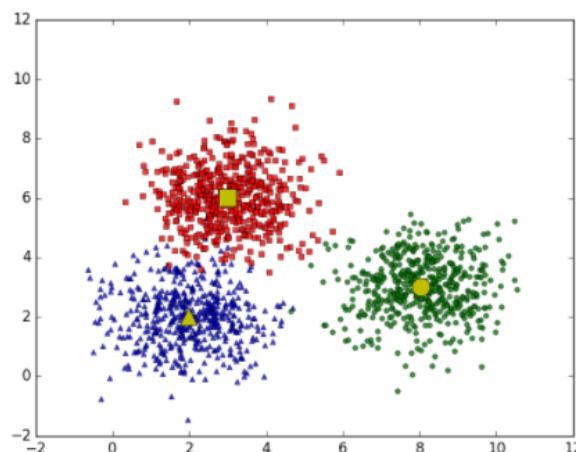
1.1 Ý tưởng giải thuật K-means

Thuật toán K-means clustering (phân cụm K-means) là một trong những thuật toán cơ bản nhất trong Unsupervised learning. Trong thuật toán K-means clustering, chúng ta không biết nhãn (label) của từng điểm dữ liệu. Mục đích là làm thế nào để phân dữ liệu thành các cụm (cluster) khác nhau sao cho dữ liệu trong cùng một cụm có tính chất giống nhau.

Ví dụ: Một công ty muốn tạo ra những chính sách ưu đãi cho những nhóm khách hàng khác nhau dựa trên sự tương tác giữa mỗi khách hàng với công ty đó (số năm là khách hàng; số tiền khách hàng đã chi trả cho công ty; độ tuổi; giới tính; thành phố; nghề nghiệp; ...).

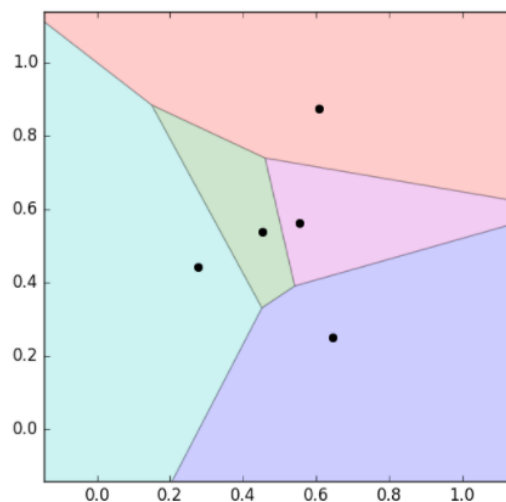
Giả sử công ty đó có rất nhiều dữ liệu của rất nhiều khách hàng nhưng chưa có cách nào chia toàn bộ khách hàng đó thành một số nhóm/cụm khác nhau. Nếu một người biết Machine Learning được đặt câu hỏi này, phương pháp đầu tiên nghĩ đến sẽ là K-means Clustering. Sau khi đã phân ra được từng nhóm, nhân viên công ty đó có thể lựa chọn ra một vài khách hàng trong mỗi nhóm để quyết định xem mỗi nhóm tương ứng với nhóm khách hàng nào. Phần việc cuối cùng này cần sự can thiệp của con người, nhưng lượng công việc đã được rút gọn đi rất nhiều.

Ý tưởng đơn giản nhất về cluster (cụm) là tập hợp các điểm ở gần nhau trong một không gian nào đó (không gian này có thể có rất nhiều chiều trong trường hợp thông tin về một điểm dữ liệu là rất lớn). Hình bên dưới là một ví dụ về 3 cụm dữ liệu).



Hình 1: Bài toán với 3 clusters

Giả sử mỗi cluster có một điểm đại diện (center) màu vàng. Và những điểm xung quanh mỗi center thuộc vào cùng nhóm với center đó. Một cách đơn giản nhất, xét một điểm bất kỳ, ta xét xem điểm đó gần với center nào nhất thì nó thuộc về cùng nhóm với center đó. Tới đây có một bài toán thú vị: Trên một vùng biển hình vuông lớn có ba đảo hình vuông, tam giác, và tròn màu vàng như hình trên. Một điểm trên biển được gọi là thuộc lãnh hải của một đảo nếu nó nằm gần đảo này hơn so với hai đảo kia. Hãy xác định ranh giới lãnh hải của các đảo. Hình dưới đây là một hình minh họa cho việc phân chia lãnh hải nếu có 5 đảo khác nhau được biểu diễn bằng các hình tròn màu đen.



Hình 2: Phân vùng lãnh hải của mỗi đảo

Chúng ta thấy rằng đường phân định giữa các lãnh hải là các đường thẳng (chính xác hơn thì chúng là các đường trung trực của các cặp điểm gần nhau). Vì vậy, lãnh hải của một đảo sẽ là một hình đa giác. Cách phân chia này trong toán học được gọi là Voronoi Diagram.

Trong không gian ba chiều, lấy ví dụ là các hành tinh, thì lãnh không của mỗi hành tinh sẽ là một đa diện. Trong không gian nhiều chiều hơn, sẽ có siêu đa diện (hyperpolygon).

Mặt hạn chế của thuật toán K-mean tiêu chuẩn:

Một nhược điểm của thuật toán K-mean là nó nhạy cảm với việc khởi tạo k trung tâm hoặc các điểm trung bình ban đầu. Vì vậy, nếu một centroid được khởi tạo là một điểm “xa”, nó có thể chỉ kết thúc không có điểm nào được liên kết với nó và đồng thời, nhiều hơn một cụm có thể kết thúc với một centroid duy nhất. Tương tự, nhiều hơn một trung tâm có thể được khởi tạo vào cùng một cụm dẫn đến phân cụm kém. Ví dụ, hãy xem xét các hình ảnh hiển thị bên dưới.



Hình 3: Khởi tạo trung tâm kém dẫn đến phân cụm kém

Để hạn chế nhược điểm trên của giải thuật K-means các nhà khoa học đã đưa ra rất nhiều kỹ thuật khởi tạo k trung tâm cụm ban đầu. Trong phạm vi bài tiểu luận, tác giả trình bày 5 kỹ thuật phổ biến cũng như ứng dụng từng kỹ thuật vào 2 bài toán cụ thể để so sánh hiệu quả của từng phương pháp.

1.2 Tìm hiểu ý tưởng các kỹ thuật khởi tạo k trung tâm ban đầu

1.2.1 Phương pháp Ngẫu nhiên

Chia tập dữ liệu thành k cụm một cách ngẫu nhiên, sau đó tìm tâm của từng cụm để khởi tạo k cụm trung tâm ban đầu. Phương pháp này được sử dụng phổ biến nhất trong việc khởi tạo k cụm ban đầu của giải thuật K-means.

```
1 # RANDOM K CLUSTERS
2 def kmeans_init_randomly_k_clusters(data, k):
3     start = time.time()
4     print('RANDOM K CLUSTERS')
5     centroids = []
6     # randomly separate rows of data as k clusters
7     size_data = len(data)
8     while True:
9         clusters_random = np.random.randint(k, size = size_data)
10        if len(set(clusters_random)) == k:
11            break
12    for i in range(k):
13        cluster = data[clusters_random == i]
14        centroids.append(np.mean(cluster, axis = 0))
15    print(round(time.time() - start, 5))
16    return np.array(centroids)
```

1.2.2 Phương pháp K-MEANS++

Kỹ thuật này đảm bảo khởi tạo centroid thông minh hơn và cải thiện chất lượng của phân cụm.

Thuật toán khởi tạo:

- **Bước 1:** Chọn ngẫu nhiên centroid đầu tiên từ các điểm dữ liệu.
- **Bước 2:** Đối với mỗi điểm dữ liệu, tính toán khoảng cách của nó từ trung tâm gần nhất, đã chọn trước đó.
- **Bước 3:** Chọn centroid tiếp theo từ các điểm dữ liệu sao cho xác suất chọn một điểm làm centroid tỷ lệ thuận với khoảng cách của nó từ centroid gần nhất, đã chọn trước đó. (tức là điểm có khoảng cách tối đa từ tâm gần nhất có nhiều khả năng được chọn tiếp theo làm tâm).
- **Bước 4:** Lặp lại các bước 2 và 3 cho đến khi k centroid được lấy mẫu.

Bằng cách thực hiện theo quy trình khởi tạo ở trên, chúng ta nhận ra các centroid cách xa nhau. Điều này làm tăng cơ hội ban đầu nhất được các trung tâm nằm trong các cụm khác nhau. Ngoài ra, vì các centroid được chọn từ các điểm dữ liệu, mỗi centroid có một số điểm dữ liệu được liên kết với nó ở cuối.

```
1 # K-MEANS++
2 def kmeans_init_kmeansPlusPlus(data, k):
3     start = time.time()
4     print('K-MEANS++')
5     # initialize the centroids list and add a randomly selected
6     # data point to the list
7     centroids = []
8     centroids.append(data[np.random.randint(data.shape[0]), :])
9     if k < 4:
10         plot(data, np.array(centroids))
11     # compute remaining k - 1 centroids
12     for c_id in range(k - 1):
13         dist = []
14         for centroid in centroids:
15             dist.append(list(map(lambda x: math.sqrt(sum(x)), (data
16 - centroid)**2)))
17         dist = np.array(dist).T
18         # index of the closest center
19         dist = np.array(list(map(lambda x: min(x), dist)))
20         # select data point with maximum distance as our next
21         centroid
22         next_centroid = data[np.argmax(dist)]
23         centroids.append(next_centroid)
```

```
21         if k < 4:
22             plot(data, np.array(centroids))
23     print(round(time.time() - start, 5))
24     return np.array(centroids)
```

1.2.3 Phương pháp FORGY APPROACH (FA)

Được đề xuất bởi Forgy vào năm 1965. Ý tưởng của kỹ thuật này là chọn k điểm một cách ngẫu nhiên trong tập dữ liệu làm k trung tâm cụm ban đầu. Vì là chọn k điểm ngẫu nhiên nên phương pháp này dễ sinh ra rủi ro nhất (chọn các điểm biên và không có điểm nào trong tập dữ liệu được gán vào hoặc 2 tâm cụm được chọn trong cùng 1 cụm thực), nhưng lại là phương pháp khởi tạo nhanh nhất, ít chi phí tính toán nhất.

```
1 # FORGY APPROACH (FA)
2 def kmeans_init_randomly_centers(data, k):
3     start = time.time()
4     print('FORGY APPROACH (FA)')
5     # randomly pick k rows of data as initial centers
6     print(round(time.time() - start, 5))
7     return data[np.random.choice(data.shape[0], k, replace=False)]
```

1.2.4 Phương pháp MACQUEEN APPROACH (MA)

Được đề xuất bởi MacQueen vào năm 1967. Ý tưởng của kỹ thuật này là chọn k điểm một cách ngẫu nhiên trong tập dữ liệu làm k trung tâm cụm ban đầu. Theo thứ tự các điểm dữ liệu trong tập dữ liệu, ta gán lần lượt vào cụm có khoảng cách từ điểm đó đến tâm của cụm gần nhất. Sau mỗi lần gán một điểm dữ liệu, ta tính lại tâm của cụm vừa được gán điểm mới. Thực hiện lần lượt đến khi tất cả các điểm dữ liệu được gán vào k cụm với k trung tâm.

```
1 # MACQUEEN APPROACH (MA)
2 def kmeans_init_MA(data, k):
3     start = time.time()
4     print('MACQUEEN APPROACH (MA)')
5     # randomly pick k rows of data as initial centers
6     choice_index = np.random.choice(data.shape[0], k, replace=False)
7     centroids = np.array(data[choice_index], dtype='f')
8     data = np.delete(data, choice_index, axis=0)
9     # save current mean (centroids) of each cluster and num of
10    # member of k clusters
11    meta_centroids = [1] * k
12    for item in data:
```

```
12     d = np.array(list(map(lambda x: math.sqrt(sum(x)), (
13         centroids - item)**2)))
14     # index of the closest center
15     index_centroid = np.argmin(d)
16     # update centroids and prepare for next iterator
17     centroids[index_centroid] = (centroids[index_centroid] *
18         meta_centroids[index_centroid] + item) / (meta_centroids[
19         index_centroid] + 1)
20     meta_centroids[index_centroid] += 1
21     print(round(time.time() - start, 5))
22     return centroids
```

1.2.5 Phương pháp KAUFMAN APPROACH (KA)

Được đề xuất bởi Kaufman và Rousseeuw vào năm 1990. Trong kỹ thuật này, k trung tâm ban đầu có được bằng cách chọn từng điểm cho đến khi tìm thấy đủ k trung tâm. Trung tâm đầu tiên là medoid của bộ dữ liệu. $K - 1$ trung tâm còn lại được chọn theo quy tắc heuristic rằng điểm được chọn hứa hẹn có giá trị lượng giá cao hơn so với phần còn lại.

Thuật toán khởi tạo:

- **Bước 1:** Chọn trung tâm đầu tiên là medoid của tập dữ liệu.
- **Bước 2:** Với mỗi điểm ω_i chưa được chọn:
 - **Bước 2.1:** Với mỗi điểm ω_j chưa được chọn:
Tính $C_{ij} = \max(D_j - d_{ij}, 0)$, với $d_{ij} = \|\omega_i - \omega_j\|$ và $D_{ij} = \min_s d_{sj}$, với s là một trung tâm cụm đã tìm được.
 - **Bước 2.2:** Tính giá trị lượng giá hứa hẹn của ω_i bằng $\sum_j C_{ij}$.
- **Bước 3:** Chọn ω_i chưa được chọn là trung tâm cụm và có giá trị $\sum_j C_{ij}$ lớn nhất.
- **Bước 4:** Nếu đã tìm đủ k trung tâm cụm thì ngừng giải thuật. Nếu chưa thì quay lại bước 2.
- **Bước 5:** Gán phần còn lại của bộ dữ liệu vào cụm gần nhất với k trung tâm cụm đã tìm được.

Giải thuật này có chi phí tính toán khá lớn. Với bộ dữ liệu lớn giải thuật có thể không chạy đến khi đạt được kết quả.


```
1 # KAUFMAN APPROACH (KA)
2 def kmeans_init_KA(data, k):
3     start = time.time()
4     print('KAUFMAN APPROACH (KA)')
5     dist = []
6     for item in data:
7         dist.append(sum(list(map(lambda x: math.sqrt(sum(x)), (
8 data - item)**2)))))
9     medoid = data[np.argmin(dist)]
10    selected_seed = np.array([medoid])
11    while len(selected_seed) < k:
12        Cji = []
13        for wi in data:
14            sub_Cji = 0
15            for wj in data:
16                Dj = min(list(map(lambda x: math.sqrt(sum(x)), (
17 selected_seed - wj)**2)))
18                sub_Cji += max(Dj - np.linalg.norm(wi-wj), 0)
19            Cji.append(sub_Cji)
20        while True:
21            index = np.argmax(np.array(Cji))
22            prepare_selected_seed = data[index]
23            if list(prepare_selected_seed) not in selected_seed.
24 tolist():
25                break
26            Cji[index] = 0
27            selected_seed = np.append(selected_seed, [
28 prepare_selected_seed], axis=0)
29        print(round(time.time() - start, 5))
30    return selected_seed
```

2 Ứng dụng giải thuật K-means với các kỹ thuật khởi tạo k trung tâm ban đầu trong việc phân cụm tập dữ liệu IRIS FLOWER

2.1 Mô tả tập dữ liệu

Iris flower dataset là một bộ dữ liệu nhỏ đã được gắn nhãn. Bộ dữ liệu này bao gồm thông tin của ba loại hoa Iris (một loài hoa lan) khác nhau: Iris setosa, Iris virginica và Iris versicolor. Mỗi loại có 50 bông hoa được đo với dữ liệu là 4 thông tin: chiều dài, chiều rộng đài hoa (sepal), và chiều dài, chiều rộng cánh hoa (petal). Dưới đây là ví dụ về hình ảnh của ba loại hoa, mỗi điểm dữ liệu trong tập này chỉ là một vector 4 chiều.



Hình 4: Ví dụ về Iris flower dataset

Bài toán phân nhóm giả định:

Bài toán: Giả sử rằng chúng ta không biết nhãn của các dữ liệu hoa này, chúng ta muốn phân nhóm các hoa gần giống nhau về một nhóm. Với thuật toán phân nhóm K-means clustering, chúng ta sẽ giải quyết bài toán này thế nào? Đồng thời sẽ sử dụng Error Sum of Squares (Sai số SSE) để đánh giá hiệu quả phân cụm của bài toán.

Mô tả chung về dữ liệu Iris flower

```
1 #Code
2 from sklearn.datasets import load_iris
3
4 iris = load_iris()
5 iris_X = iris.data
6 iris_y = iris.target
```

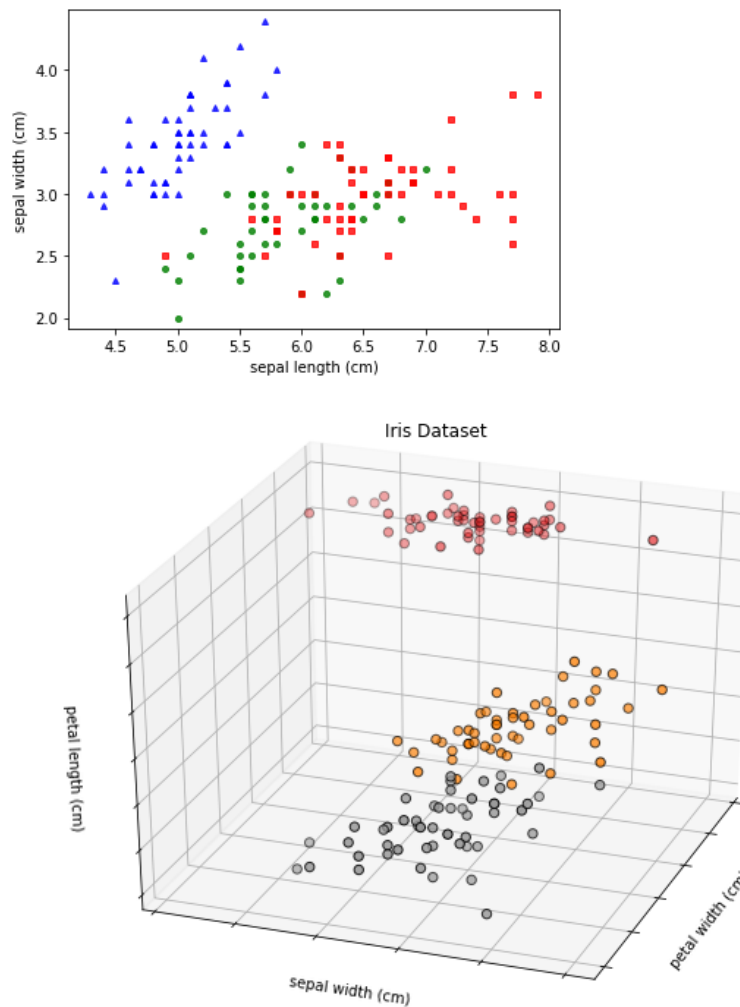
```
7 print('Number of classes: ' , len(np.unique(iris_y)))
8 print('Number of data points: ' , len(iris_y))
9
10 X0 = iris_X[iris_y == 0,:]
11 print('\nSamples from class 0:\n', X0[:5,:])
12
13 X1 = iris_X[iris_y == 1,:]
14 print('\nSamples from class 1:\n', X1[:5,:])
15
16 X2 = iris_X[iris_y == 2,:]
17 print('\nSamples from class 2:\n', X2[:5,:])
18
19 #Out
20 Number of classes: 3
21 Number of data points: 150
22
23 Samples from class 0:
24 [[5.1 3.5 1.4 0.2]
25  [4.9 3.  1.4 0.2]
26  [4.7 3.2 1.3 0.2]
27  [4.6 3.1 1.5 0.2]
28  [5.  3.6 1.4 0.2]]
29
30 Samples from class 1:
31 [[7.  3.2 4.7 1.4]
32  [6.4 3.2 4.5 1.5]
33  [6.9 3.1 4.9 1.5]
34  [5.5 2.3 4.  1.3]
35  [6.5 2.8 4.6 1.5]]
36
37 Samples from class 2:
38 [[6.3 3.3 6.  2.5]
39  [5.8 2.7 5.1 1.9]
40  [7.1 3.  5.9 2.1]
41  [6.3 2.9 5.6 1.8]
42  [6.5 3.  5.8 2.2]]
```

Thể hiện các điểm dữ liệu lên biểu đồ 2D và 3D

Để thể hiện tốt được sự khác biệt giữa các điểm dữ liệu 4 chiều trên biểu đồ 2D và 3D ta chọn các thuộc tính để vẽ như sau:

- Với biểu đồ 2D, ta chọn 2 thuộc tính độ rộng đài hoa và độ dài đài hoa.
- Với biểu đồ 3D, ta chọn 3 thuộc tính độ rộng cánh hoa, độ dài cánh hoa và độ rộng đài hoa .

Cách biểu diễn này cũng được áp dụng trong phần sau.



Hình 5: Biểu diễn tập dữ liệu Iris flower có nhãn ban đầu

Sai số SSE của tập dữ liệu ban đầu đo được là: 100.39574. Giá trị này giúp so sánh hiệu quả của các phương pháp phân cụm được trình bày ở phần sau.

2.2 Phân loại dữ liệu với các kỹ thuật khởi tạo k trung tâm ban đầu

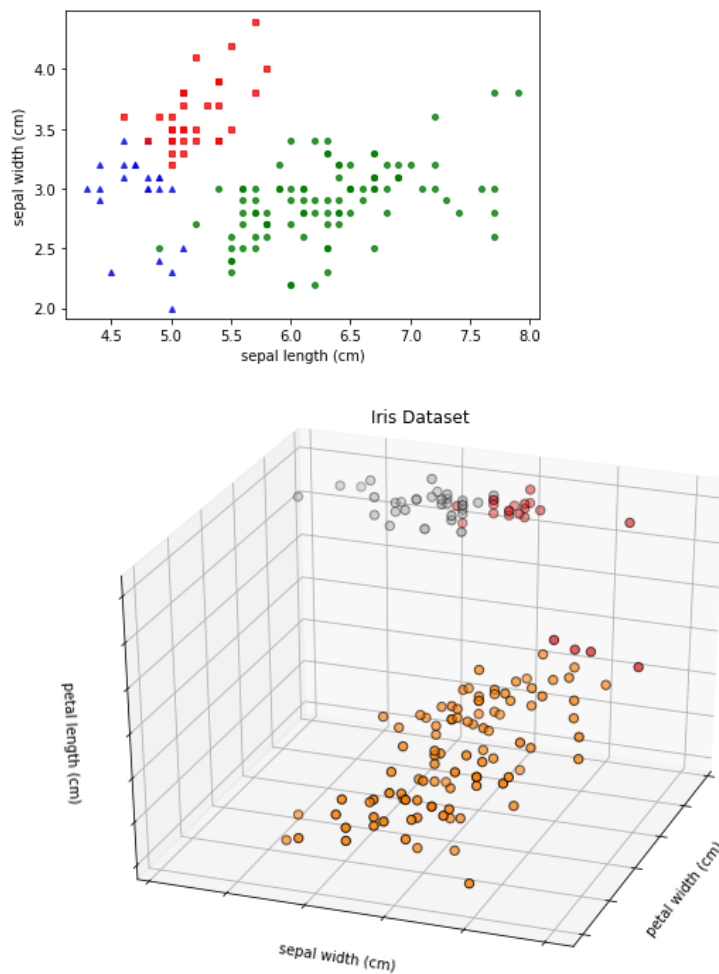
2.2.1 Phương pháp Ngẫu nhiên

Bảng kết quả sau 10 lần chạy thuật toán:

Phương pháp	Thông số	Lần 1	Lần 2	Lần 3	Lần 4	Lần 5	Lần 6	Lần 7	Lần 8	Lần 9	Lần 10	Trung bình
RANDOM K CLUSTERS	Thời gian	0.00085	0.00053	0.00029	0.00542	0.00078	0.00568	0.00039	0.0019	0.00256	0.00337	0.002177
	Vòng lặp	10	11	10	11	11	11	12	8	11	10	10.5
	SSE	97.22487	97.22487	97.22487	97.22487	97.22487	97.22487	124.02237	97.22487	97.22487	97.22487	99.90462

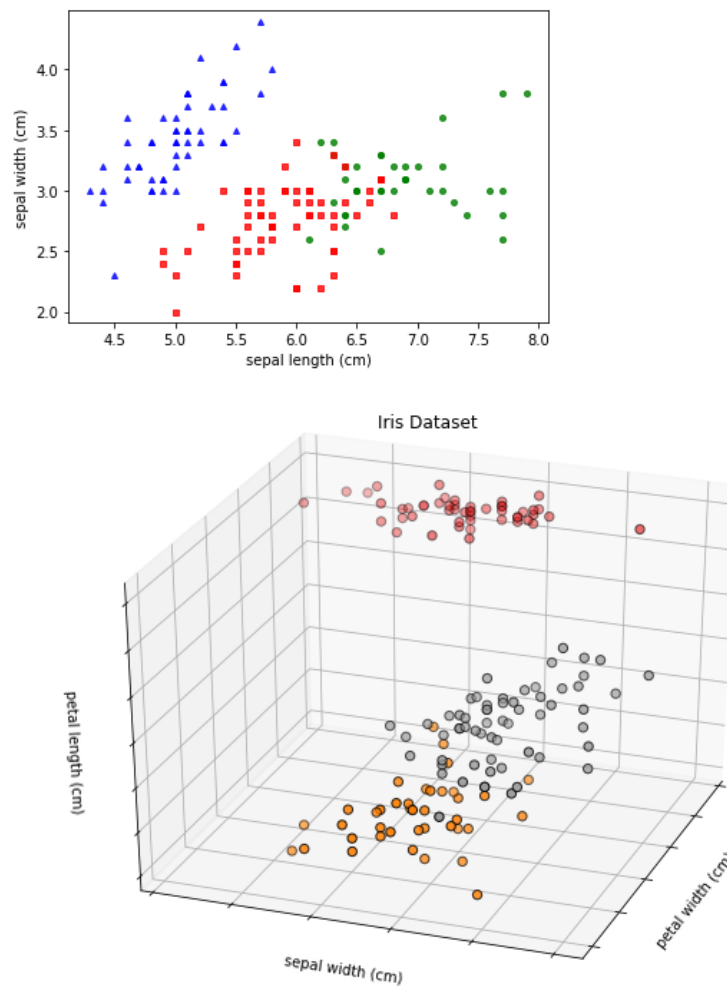
Hình 6: Bảng kết quả với phương pháp Ngẫu nhiên

Kết quả của lần chạy thứ 7, do việc khởi tạo 3 cụm trung tâm ban đầu dẫn đến việc phân chia các cụm không tối ưu (Sai số SSE = 124.02237)



Hình 7: Lần chạy thứ 7 bằng phương pháp Ngẫu nhiên

Kết quả của lần chạy thứ 10, nhờ khởi tạo 3 cụm tốt mà công việc phân cụm tốt (Sai số SSE = 97.22487, xấp xỉ sai số SSE của tập dữ liệu có dán nhãn tương ứng 100.39574).



Hình 8: Lần chạy thứ 10 bằng phương pháp Ngẫu nhiên

Nhìn chung thời gian trung bình 10 lần chạy thuật toán khá bé do chi phí tính toán ít, trái lại số vòng lặp trung bình để hội tụ và sai số SSE trung bình khá lớn.

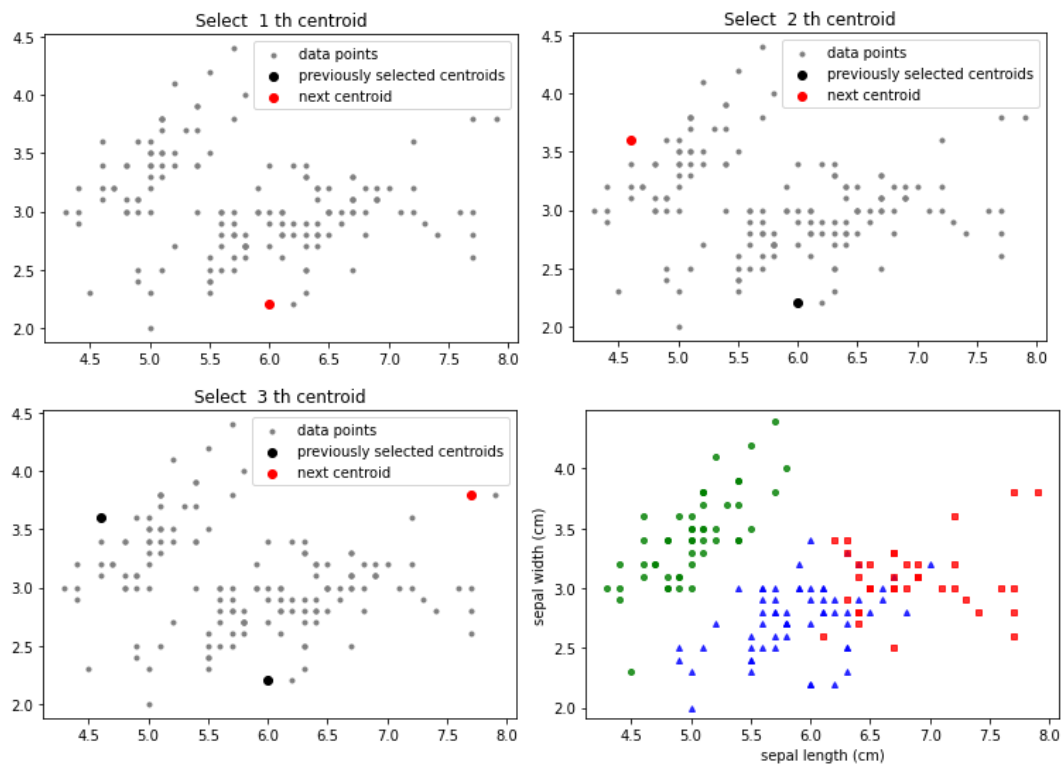
2.2.2 Phương pháp K-MEANS++

Bảng kết quả sau 10 lần chạy thuật toán:

Phương pháp	Thông số	Lần 1	Lần 2	Lần 3	Lần 4	Lần 5	Lần 6	Lần 7	Lần 8	Lần 9	Lần 10	Trung bình
K-MEANS++	Thời gian	0.52934	0.63747	0.56003	0.66934	0.54443	0.53559	0.54364	0.53747	0.62381	0.58109	0.576221
	Vòng lặp	5	3	5	3	10	3	11	3	10	2	5.5
	SSE	97.20457	97.20457	97.20457	97.20457	97.20457	97.20457	97.22487	97.20457	97.22487	97.22487	97.21066

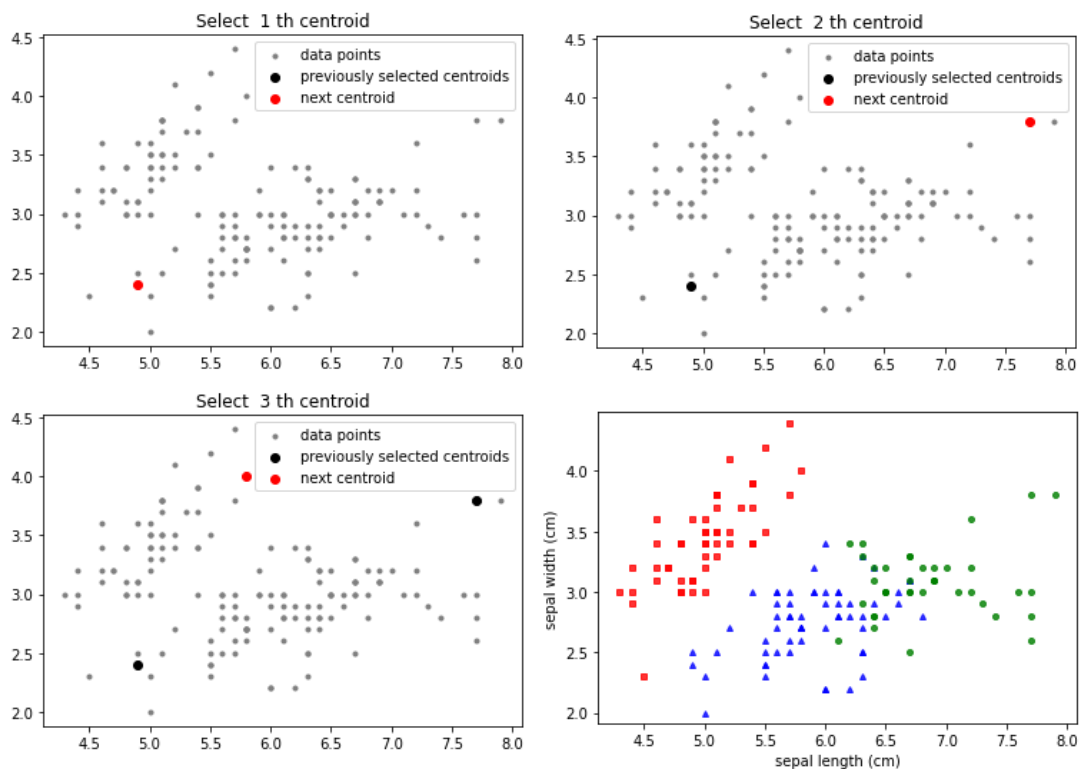
Hình 9: Bảng kết quả với phương pháp K-MEANS++

Lần chạy giải thuật thứ 3 với 5 vòng lặp để hội tụ, cùng sai số SSE là 97.20457.



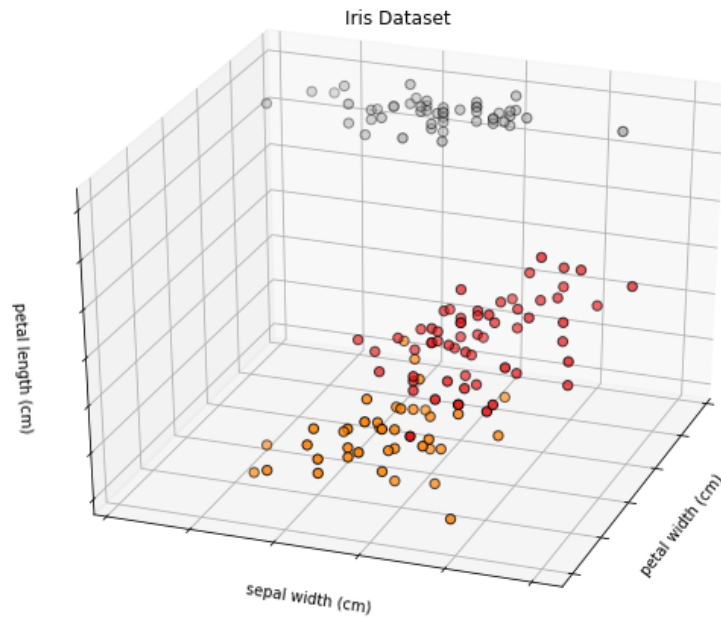
Hình 10: Lần chạy thứ 3 bằng phương pháp K-MEANS++ - các bước khởi tạo từng centroid

Ở lần chạy thứ 10, nhờ việc khởi tạo 3 centroid ở các vị trí thông minh mà chỉ cần 2 vòng lặp đã hội tụ, đồng thời đảm bảo giá trị sai số SSE thấp.



Hình 11: Lần chạy thứ 10 bằng phương pháp K-MEANS++ - các bước khởi tạo từng centroid

Nhờ việc khởi tạo centroid thông minh nên sai số SSE trung bình qua 10 lần chạy khá thấp và ổn định 97.21066, không có trường hợp 2 centroid nằm cùng 1 phân cụm thực tế. Đồng thời số vòng lặp trung bình để hội tụ khá thấp 5.5. Nhưng bù lại thời gian khởi tạo lớn vì cần nhiều chi phí tính toán.



Hình 12: Biểu diễn các điểm dữ liệu đã phân cụm trên biểu đồ 3D

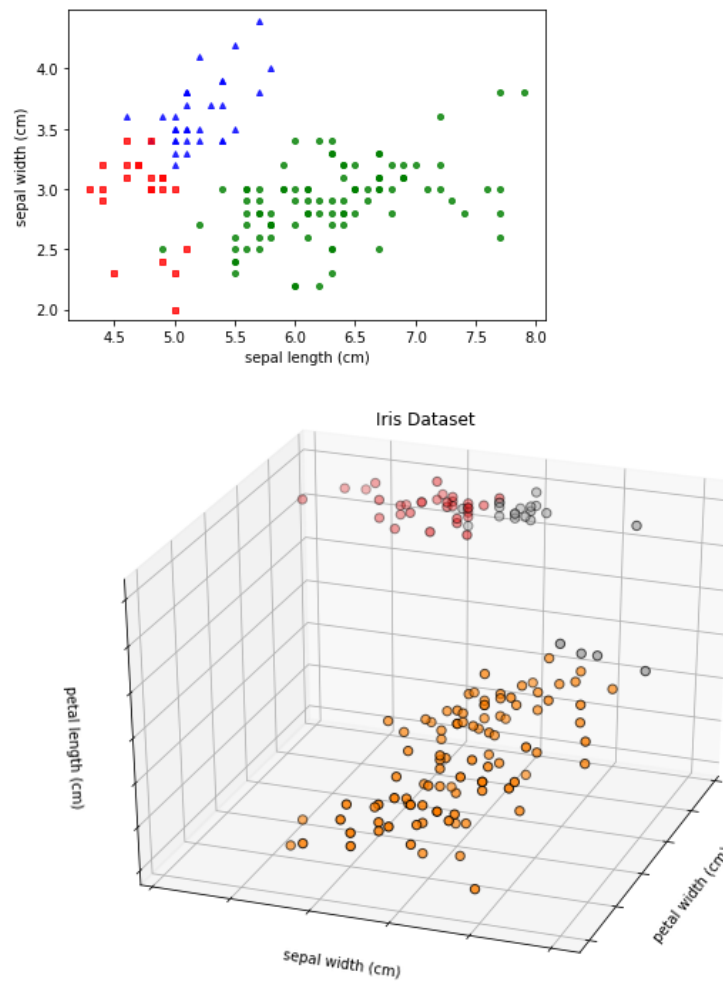
2.2.3 Phương pháp FORGY APPROACH (FA)

Bảng kết quả sau 10 lần chạy thuật toán:

Phương pháp	Thông số	Lần 1	Lần 2	Lần 3	Lần 4	Lần 5	Lần 6	Lần 7	Lần 8	Lần 9	Lần 10	Trung bình
FORGY APPROACH (FA)	Thời gian	0.00357	0.0001	0.00118	0.00034	0.0001	0.00259	0.00013	0.00005	0.00006	0.00008	0.00082
	Vòng lặp	3	3	10	5	11	5	4	6	10	12	6.9
	SSE	97.20457	97.20457	97.22487	97.20457	97.22487	124.0224	122.6001	97.20457	124.0224	97.22487	105.113773

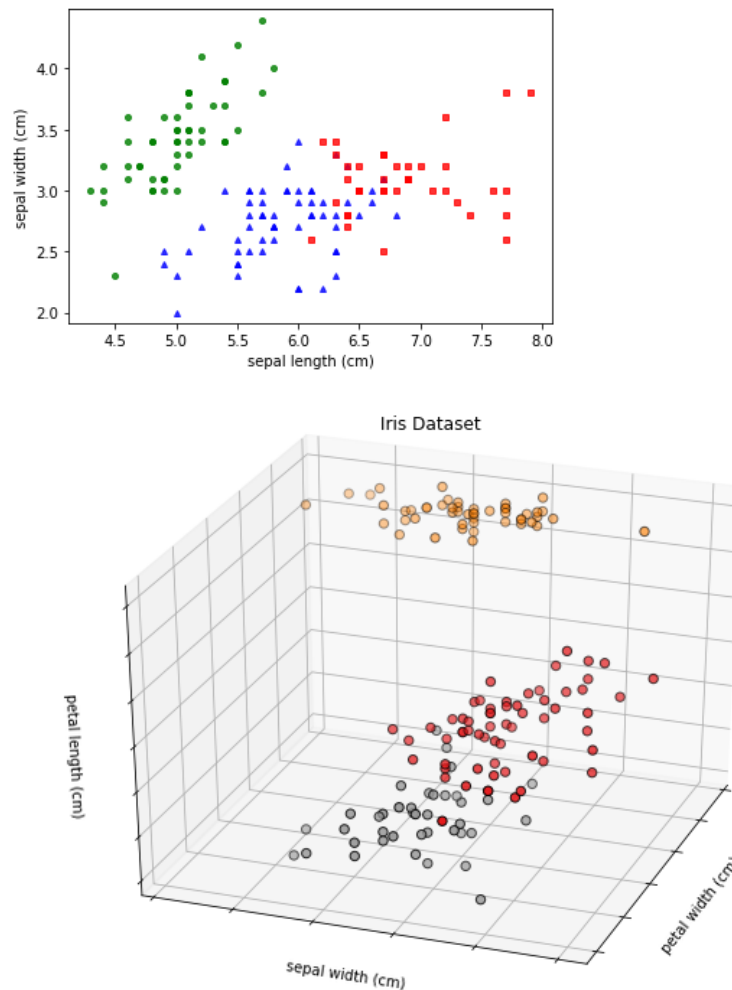
Hình 13: Bảng kết quả với phương pháp FORGY APPROACH

Ở lần chạy thứ 7, thời gian khởi tạo 3 trung tâm cụm rất nhanh 0.00013s và cần 4 vòng lặp để hội tụ nhưng sai số SSE khá lớn 122.6001.



Hình 14: Lần chạy thứ 7 bằng phương pháp FORGY APPROACH

Với lần chạy thứ 10, Thời gian khởi tạo các trung tâm cụm vẫn rất nhanh 0.00008s, với 12 vòng lặp để hội tụ và sai số SSE lúc này tốt hơn 97.22487 (xấp xỉ sai số SSE của tập dữ liệu có dán nhãn tương ứng 100.39574).



Hình 15: Lần chạy thứ 10 bằng phương pháp FORGY APPROACH

Dễ thấy rằng thời gian khởi tạo 3 trung tâm cụm của phương pháp FORGY APPROACH ổn định ở mức rất thấp do chỉ chọn ngẫu nhiên 3 trong số các điểm trong điểm dữ liệu, nhưng đồng thời cũng mang lại nhiều rủi ro khác như sai số SSE trung bình cao 105.113773, số vòng lặp để hội tụ cũng không ổn định lúc lớn lúc nhỏ (trung bình 6.9) và đây chỉ là tập dữ liệu gồm 150 điểm dữ liệu nên nếu ở bộ dữ liệu lớn thì độ sai biệt này càng rõ rệt hơn.

Đồng thời phương pháp này cũng mang nặng yếu tố ngẫu nhiên và cục bộ, chưa xét đến ảnh hưởng của toàn bộ tập dữ liệu khi khởi tạo nên hiệu quả của phương pháp này không ổn định hay có thể nói là không tốt.

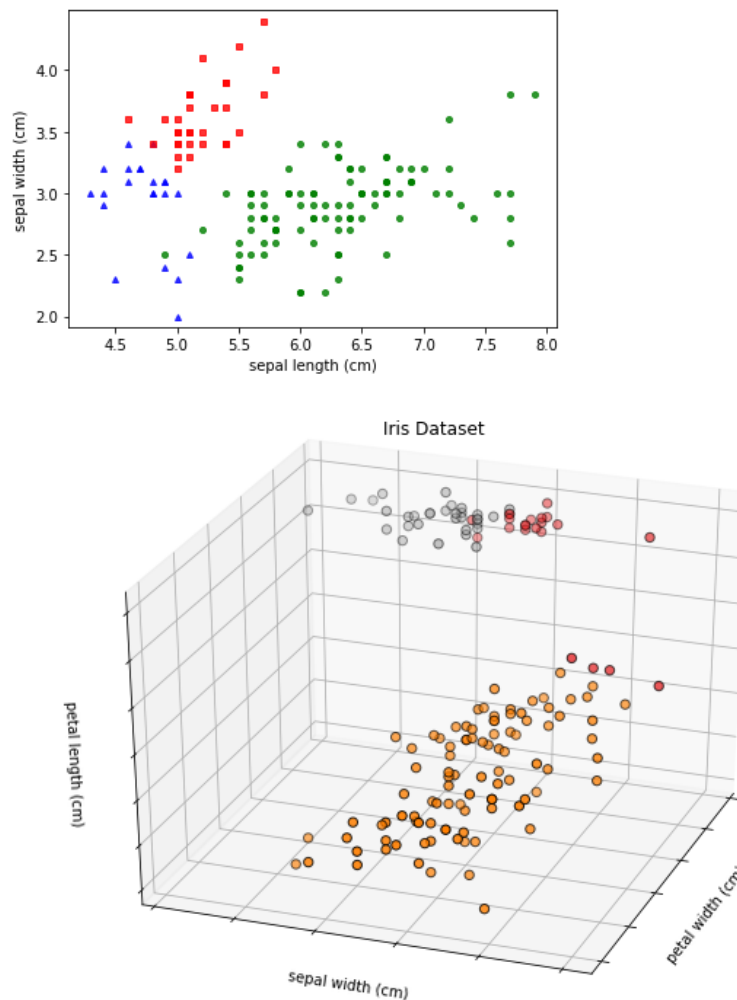
2.2.4 Phương pháp MACQUEEN APPROACH (MA)

Bảng kết quả sau 10 lần chạy thuật toán:

Phương pháp	Thông số	Lần 1	Lần 2	Lần 3	Lần 4	Lần 5	Lần 6	Lần 7	Lần 8	Lần 9	Lần 10	Trung bình
MACQUEEN APPROACH (MA)	Thời gian	0.00928	0.00567	0.00796	0.00918	0.00922	0.01149	0.0126	0.01705	0.01126	0.01147	0.010518
	Vòng lặp	1	8	1	2	9	9	7	6	1	5	4.9
	SSE	97.20457	97.22487	97.20457	97.20457	124.02237	97.22487	124.02237	124.02237	97.20457	124.0224	107.93575

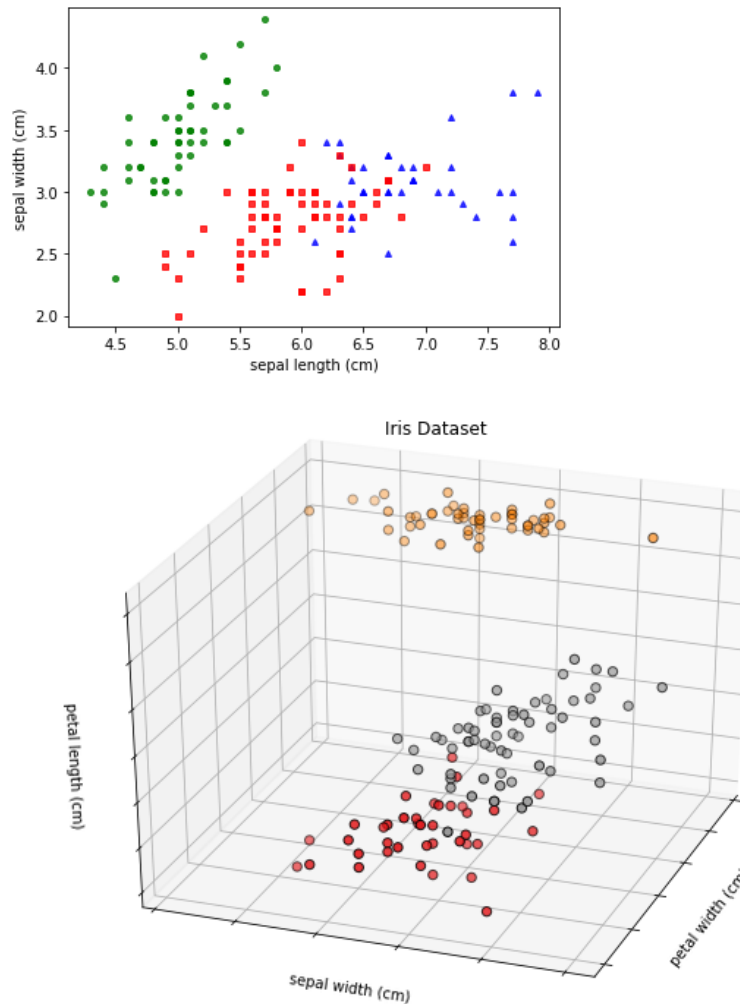
Hình 16: Bảng kết quả với phương pháp MACQUEEN APPROACH

Kết quả ở lần chạy thứ 5 với sai số SSE là 124.02237, mất 0.00922s khởi tạo 3 cụm trung tâm và 9 lần lặp để hội tụ.



Hình 17: Lần chạy thứ 5 bằng phương pháp MACQUEEN APPROACH

Kết quả ở lần chạy thứ 9 với sai số SSE là 97.20457, mất 0.01126s khởi tạo 3 cụm trung tâm và chỉ cần 1 lần lặp để hội tụ.



Hình 18: Lần chạy thứ 9 bằng phương pháp MACQUEEN APPROACH

Với phương pháp MACQUEEN APPROACH số vòng lặp trung bình để hội tụ khá thấp, nhưng giá trị quan trọng nhất là sai số SSE trung bình lại rất cao đồng thời yêu cầu tính toán trên toàn bộ tập dữ liệu làm giải thuật hoạt động khá chậm.

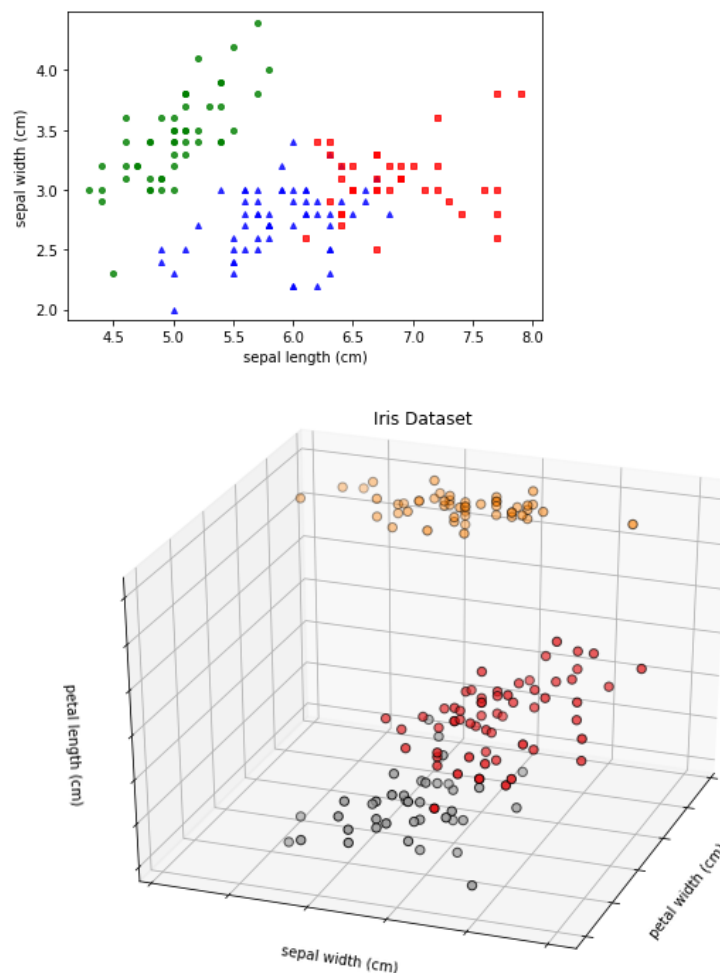
2.2.5 Phương pháp KAUFMAN APPROACH (KA)

Bảng kết quả sau 10 lần chạy thuật toán:

Phương pháp	Thông số	Lần 1	Lần 2	Lần 3	Lần 4	Lần 5	Lần 6	Lần 7	Lần 8	Lần 9	Lần 10	Trung bình
KAUFMAN APPROACH (KA)	Thời gian	0.90235	0.92017	0.9026	0.98473	1.14646	0.95677	1.00863	1.04121	1.05406	0.9537	0.987068
	Vòng lặp	3	3	3	3	3	3	3	3	3	3	3
	SSE	97.22487	97.22487	97.22487	97.22487	97.22487	97.22487	97.22487	97.22487	97.22487	97.22487	97.22487

Hình 19: Bảng kết quả với phương pháp KAUFMAN APPROACH

Kết quả ở lần chạy th 9 với sai số SSE là 97.22487, mất 1.05406 khởi tạo 3 cụm trung tâm và chỉ cần 3 lần lặp để hội tụ.



Hình 20: Lần chạy thứ 5 bằng phương pháp KAUFMAN APPROACH

Đễ dàng nhận thấy sai số SSE và số vòng lặp hội tụ không đổi qua 10 lần chạy do không có yếu tố ngẫu nhiên trong giải thuật. Các trung tâm cụm tìm được qua n lần chạy vẫn không đổi, chỉ khác nhau thời gian chạy do ảnh hưởng của hệ điều hành trong việc sắp xếp tiến trình. Dù mang lại kết quả tốt về số vòng lặp hội tụ cũng như sai số SSE nhưng hạn chế lớn nhất của phương pháp KAUFMAN APPROACH chính là chi phí tính toán quá lớn (hàm mũ) sẽ không hoạt động tốt ở bộ dữ liệu lớn. Phần sau sẽ thực hiện trên bộ dữ liệu lớn The Mnist để thấy rõ hạn chế này.

2.3 Nhận xét

Phương pháp	Thông số	Lần 1	Lần 2	Lần 3	Lần 4	Lần 5	Lần 6	Lần 7	Lần 8	Lần 9	Lần 10	Trung bình
RANDOM K CLUSTERS	Thời gian	0.00085	0.00053	0.00029	0.00542	0.00078	0.00568	0.00039	0.0019	0.00256	0.00337	0.002177
	Vòng lặp	10	11	10	11	11	11	12	8	11	10	10.5
	SSE	97.22487	97.22487	97.22487	97.22487	97.22487	97.22487	124.02237	97.22487	97.22487	97.22487	99.90462
FORGY APPROACH (FA)	Thời gian	0.00357	0.0001	0.00118	0.00034	0.0001	0.00259	0.00013	0.00005	0.00006	0.00008	0.00082
	Vòng lặp	3	3	10	5	11	5	4	6	10	12	6.9
	SSE	97.20457	97.20457	97.22487	97.20457	97.22487	124.0224	122.6001	97.20457	124.0224	97.22487	105.113773
K-MEANS++	Thời gian	0.52934	0.63747	0.56003	0.66934	0.54443	0.53559	0.54364	0.53747	0.62381	0.58109	0.576221
	Vòng lặp	5	3	5	3	10	3	11	3	10	2	5.5
	SSE	97.20457	97.20457	97.20457	97.20457	97.20457	97.20457	97.22487	97.20457	97.22487	97.22487	97.21066
MACQUEEN APPROACH (MA)	Thời gian	0.00928	0.00567	0.00796	0.00918	0.00922	0.01149	0.0126	0.01705	0.01126	0.01147	0.010518
	Vòng lặp	1	8	1	2	9	9	7	6	1	5	4.9
	SSE	97.20457	97.22487	97.20457	97.20457	124.02237	97.22487	124.02237	124.02237	97.20457	124.0224	107.93575
KAUFMAN APPROACH (KA)	Thời gian	0.90235	0.92017	0.9026	0.98473	1.14646	0.95677	1.00863	1.04121	1.05406	0.9537	0.987068
	Vòng lặp	3	3	3	3	3	3	3	3	3	3	3
	SSE	97.22487	97.22487	97.22487	97.22487	97.22487	97.22487	97.22487	97.22487	97.22487	97.22487	97.22487

Hình 21: Bảng thống kê kết quả của 5 kỹ thuật khởi tạo trung tâm cụm đầu tiên qua 10 lần chạy

Nhận xét về hiệu quả của 5 phương pháp khởi tạo trung tâm cụm ban đầu với bộ dữ liệu nhỏ Iris Flower:

- **Phương pháp Ngẫu nhiên:** dựa vào yếu tố ngẫu nhiên trên toàn bộ tập dữ liệu để chọn ra các phần tử thuộc k cụm ban đầu nhưng các trung tâm cụm này tạo ra có xu hướng nằm gần nhau nên quá trình hội tụ lâu hơn, đồng thời hiệu quả của giải thuật thông qua sai số SSE cũng không tối ưu so với các phương pháp khác.
- **Phương pháp K-MEANS++:** nhờ việc khởi tạo centroid thông minh (các trung tâm cụm cố gắng tách rời nhau, và ít bị ảnh hưởng bởi yếu tố ngẫu nhiên) nên sai số SSE thấp và ổn định, số vòng lặp để hội tụ cũng khá thấp, nhưng bù lại thời gian khởi tạo lớn vì cần nhiều chi phí tính toán (dù vậy độ phức tạp giải thuật chỉ là $O(n)$ nên có thể áp dụng tốt cho các bộ dữ liệu vừa và nhỏ).

- **Phương pháp FORGY APPROACH:** mang nặng yếu tố ngẫu nhiên và cục bộ, chưa xét đến ảnh hưởng của toàn bộ tập dữ liệu khi khởi tạo nên chứa nhiều rủi ro như sai số SSE cao, số vòng lặp để hội tụ cũng không ổn định (với bộ dữ liệu lớn thì độ sai biệt này càng rõ rệt hơn), bù lại vì đơn giản nên thời gian khởi tạo rất nhanh.
Hiệu quả của phương pháp này không ổn định hay có thể nói là không tốt.
- **Phương pháp MACQUEEN APPROACH:** bị ảnh hưởng bởi yếu tố ngẫu nhiên (khởi tạo ngẫu nhiên k trung tâm ban đầu), tuy có xét đến ảnh hưởng của toàn bộ tập dữ liệu xong vẫn có xu hướng hội tụ cục bộ. Kết quả là với chi phí tính toán khá lớn nhưng hiệu quả vẫn không cao (sai số SSE lớn).
- **Phương pháp KAUFMAN APPROACH:** không bị ảnh hưởng bởi các yếu tố ngẫu nhiên, số vòng lặp để hội tụ ít và sai số SSE thấp, bù lại thì chi phí tính toán quá lớn (hàm mũ theo số lượng dữ liệu) nên chỉ thích hợp với bộ dữ liệu nhỏ, số chiều ít). Thời gian khởi tạo trung bình của phương pháp này cao gấp 1204 lần phương pháp đơn giản nhất (FORGY APPROACH).

Vậy với bộ dữ liệu nhỏ như Iris Flower phương pháp khởi tạo K-MEANS++ đem lại hiệu quả tốt cả về mặt thời gian cũng như sai số SSE của toàn bộ tập dữ liệu khi đã phân cụm.

3 Ứng dụng giải thuật K-means với các kỹ thuật khởi tạo k trung tâm ban đầu trong việc phân cụm tập dữ liệu THE MNIST

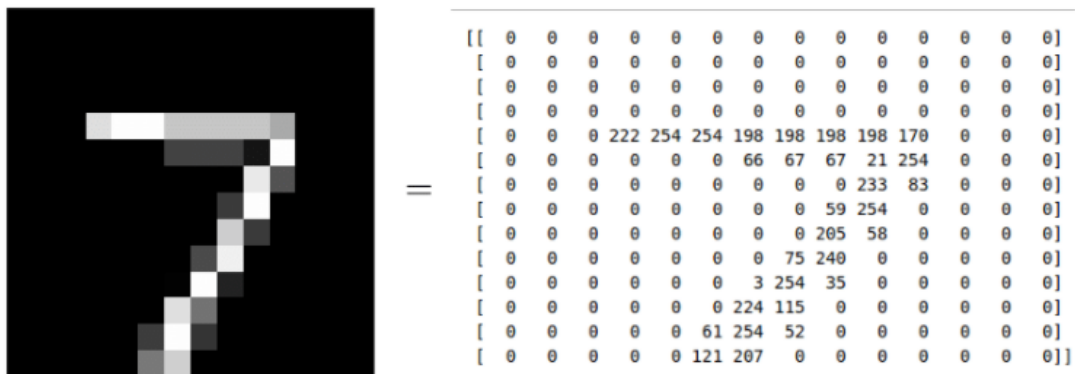
3.1 Mô tả tập dữ liệu

Bộ cơ sở dữ liệu MNIST là bộ cơ sở dữ liệu lớn nhất về chữ số viết tay và được sử dụng trong hầu hết các thuật toán nhận dạng hình ảnh (Image Classification). MNIST bao gồm hai tập con: tập dữ liệu huấn luyện (training set) có tổng cộng 60k ví dụ khác nhau về chữ số viết tay từ 0 đến 9, tập dữ liệu kiểm tra (test set) có 10k ví dụ khác nhau. Tất cả đều đã được gán nhãn. Hình dưới đây là ví dụ về một số hình ảnh được trích ra từ MNIST.



Hình 22: MNIST: bộ cơ sở dữ liệu của chữ số viết tay

Mỗi bức ảnh là một ảnh đen trắng (có 1 channel), có kích thước 28x28 pixel (tổng cộng 784 pixels). Mỗi pixel mang một giá trị là một số tự nhiên từ 0 đến 255. Các pixel màu đen có giá trị bằng 0, các pixel càng trắng thì có giá trị càng cao (nhưng không quá 255). Dưới đây là một ví dụ về chữ số 7 và giá trị các pixel của nó.



Hình 23: Mô tả chữ số 7

Việc áp dụng các kỹ thuật khởi tạo k trung tâm cụm ban đầu của giải thuật K-means lên bộ dữ liệu lớn The Mnist này giúp thấy rõ các mặt hạn chế của các kỹ thuật mà Iris Flower chưa thấy được.

Bài toán phân nhóm giả định:

Bài toán: Giả sử rằng chúng ta không biết nhãn của các chữ số này, chúng ta muốn phân nhóm các bức ảnh gần giống nhau về một nhóm. Với thuật toán phân nhóm K-means clustering, chúng ta sẽ giải quyết bài toán này thế nào?

Trước khi áp dụng thuật toán K-means clustering, chúng ta cần coi mỗi bức ảnh là một điểm dữ liệu. Và vì mỗi điểm dữ liệu là 1 vector (hàng hoặc cột) chứ không phải ma trận như số 7 ở trên, chúng ta phải làm thêm một bước đơn giản trung gian gọi là vectorization (vector hóa). Nghĩa là, để có được 1 vector, ta có thể tách các hàng của ma trận pixel ra, sau đó đặt chúng cạnh nhau, và chúng ta được một vector hàng rất dài biểu diễn 1 bức ảnh chữ số. Cách làm này chỉ là cách đơn giản nhất để mô tả dữ liệu ảnh bằng 1 vector. Trên thực tế có rất nhiều kỹ thuật khác nhau để có thể tạo ra các vector đặc trưng (feature vector) giúp các thuật toán có được kết quả tốt hơn.

Sau khi vector hóa một ảnh 28x28 thành vector 784 chiều nên để đơn giản tác giả chỉ chọn 1000 ảnh của tập testing để thực hiện tiểu luận này (giảm quá trình tính toán khi hội tụ các trung tâm cụm).

Mô tả chung về dữ liệu The Mnist

```
1 #Code
2 mndata = MNIST('/content/gdrive/MyDrive/Colab Notebooks/MNIST')
3 mndata.load_testing()
4 y = np.asarray(mndata.test_labels)[:1000]
5 X = np.asarray(mndata.test_images)[:1000,:]/256.0
```

```
6 print(X.shape)
7 print(y.shape)
8 SSE(X, y)
9 centroids = []
10 for i in range(10):
11     cluster = X[y == i]
12     centroids.append(np.mean(cluster, axis = 0))
13 display_MNIST(X, np.array(centroids), y)
14
15 #Out
16 (1000, 784)
17 (1000,)
18 Error Sum of Squares: 6093.73051
```

Thể hiện các điểm dữ liệu có nhãn ban đầu

Để thể hiện tốt được hiệu quả phân cụm của các chữ số viết tay, ta vẽ 2 hình như sau:

- Hình bên trái: là hình vẽ giá trị trung tâm của các cụm ứng với từng hàng bên phải.
- Hình bên phải: 10 hình vẽ đầu tiên của từng cụm đã được phân loại.

Cách biểu diễn này cũng được áp dụng trong phần sau.



Hình 24: Biểu diễn tập dữ liệu The Mnist có nhãn ban đầu

Sai số SSE của tập dữ liệu ban đầu đo được là: 6093.73051. Giá trị này giúp so sánh hiệu quả của các phương pháp phân cụm được trình bày ở phần sau.

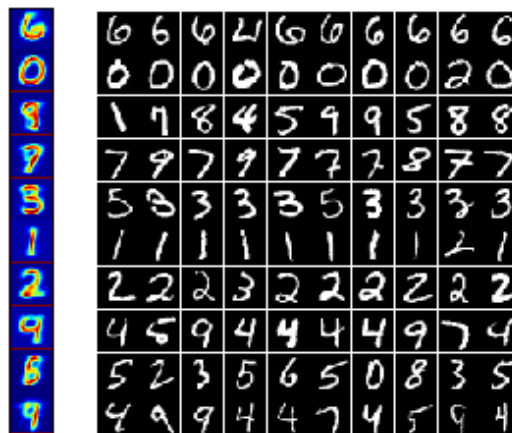
3.2 Phân loại dữ liệu với các kỹ thuật khởi tạo k trung tâm ban đầu

3.2.1 Phương pháp Ngẫu nhiên

Bảng kết quả sau 10 lần chạy thuật toán:

Phương pháp	Thông số	Lần 1	Lần 2	Lần 3	Lần 4	Lần 5	Lần 6	Lần 7	Lần 8	Lần 9	Lần 10	Trung bình
RANDOM K CLUSTERS	Thời gian	0.00855	0.00532	0.00556	0.00305	0.0051	0.00479	0.0051	0.00383	0.00299	0.00415	0.004844
	Vòng lặp	17	32	19	19	36	22	36	24	14	14	23.3
	SSE	5930.5494	5939.1930	5941.8433	5961.6496	5935.9937	5927.3018	5964.1246	5941.1095	5949.7992	5948.9486	5944.0513

Hình 25: Bảng kết quả với phương pháp Ngẫu nhiên



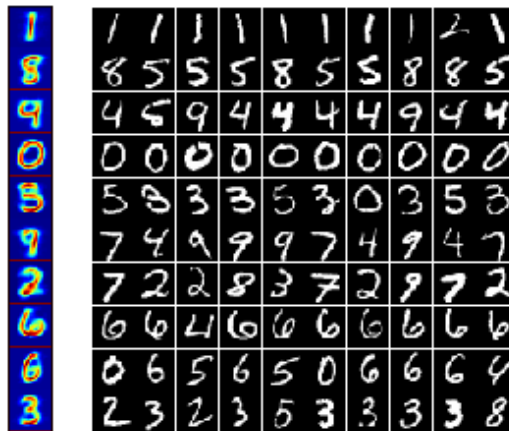
Hình 26: Biểu diễn tập dữ liệu được phân loại với phương pháp Ngẫu nhiên

3.2.2 Phương pháp K-MEANS++

Bảng kết quả sau 10 lần chạy thuật toán:

Phương pháp	Thông số	Lần 1	Lần 2	Lần 3	Lần 4	Lần 5	Lần 6	Lần 7	Lần 8	Lần 9	Lần 10	Trung bình
K-MEANS++	Thời gian	5.80787	5.8044	5.81957	5.80238	5.82219	5.77378	5.79221	5.79412	5.82603	5.8484	5.809095
	Vòng lặp	14	20	17	36	20	18	20	20	15	26	20.6
	SSE	5970.2789	5949.6351	5960.8560	5948.4690	5950.7241	5930.0117	5959.1173	5949.6351	5956.6539	5954.3616	5952.9743

Hình 27: Bảng kết quả với phương pháp K-MEANS++



Hình 28: Biểu diễn tập dữ liệu được phân loại với phương pháp K-MEANS++

3.2.3 Phương pháp FORGY APPROACH (FA)

Bảng kết quả sau 10 lần chạy thuật toán:

Phương pháp	Thông số	Lần 1	Lần 2	Lần 3	Lần 4	Lần 5	Lần 6	Lần 7	Lần 8	Lần 9	Lần 10	Trung bình
FORGY	Thời gian	0.00007	0.00003	0.00003	0.00002	0.00024	0.00005	0.00013	0.00047	0.00003	0.00023	0.00013
APPROACH	Vòng lặp	28	35	23	11	28	36	21	26	13	25	24.6
(FA)	SSE	5948.3647	5930.8015	5986.3016	5944.6580	5949.7843	5939.3301	5944.7821	5950.6811	5940.1967	5948.6813	5948.3581

Hình 29: Bảng kết quả với phương pháp FORGY APPROACH



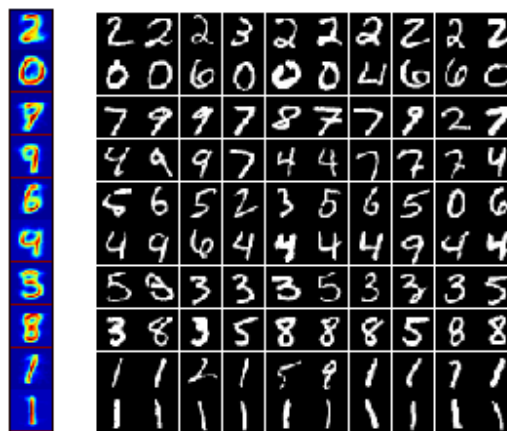
Hình 30: Biểu diễn tập dữ liệu được phân loại với phương pháp FORGY APPROACH

3.2.4 Phương pháp MACQUEEN APPROACH (MA)

Bảng kết quả sau 10 lần chạy thuật toán:

Phương pháp	Thông số	Lần 1	Lần 2	Lần 3	Lần 4	Lần 5	Lần 6	Lần 7	Lần 8	Lần 9	Lần 10	Trung bình
MACQUEEN	1.37322	1.35456	1.28899	1.29595	1.32568	1.30472	1.32406	1.3608	1.29773	1.33477	1.326048	1.321331
APPROACH	26	27	39	19	24	16	22	24	8	22	22.7	22.37
(MA)	5953.0583	6031.5797	5982.3435	5924.9864	6026.1666	5965.8175	5945.8815	5985.3529	5997.7372	5972.4951	5978.5419	5981.09

Hình 31: Bảng kết quả với phương pháp MACQUEEN APPROACH



Hình 32: Biểu diễn tập dữ liệu được phân loại với phương pháp MACQUEEN APPROACH

3.2.5 Phương pháp KAUFMAN APPROACH (KA)

Như đã phân tích về độ phức tạp hàm mũ của kỹ thuật KAUFMAN APPROACH khi thực hiện với bộ dữ liệu Iris flower thì với số lượng dữ liệu cũng như số chiều lớn của The Mnist, việc phân loại với số lượng dữ liệu là 1000 như với các phương pháp khác thì không thể trả về kết quả, vì vậy tác giả chỉ dùng 200 dữ liệu để phân loại.

Bảng kết quả như sau::

Phương pháp	Thông số	N	Trung bình
KAUFMAN	Thời gian	200	265.96389
APPROACH	Vòng lặp		5
(KA)	SSE		1139.26754

Hình 33: Bảng kết quả với phương pháp KAUFMAN APPROACH



Hình 34: Biểu diễn tập dữ liệu được phân loại với phương pháp KAUFMAN APPROACH

3.3 Nhận xét

Phương pháp	Thông số	Lần 1	Lần 2	Lần 3	Lần 4	Lần 5	Lần 6	Lần 7	Lần 8	Lần 9	Lần 10	Trung bình
RANDOM K CLUSTERS	Thời gian	0.00855	0.00532	0.00556	0.00305	0.0051	0.00479	0.0051	0.00383	0.00299	0.00415	0.004844
	Vòng lặp	17	32	19	19	36	22	36	24	14	14	23.3
	SSE	5930.5494	5939.1930	5941.8433	5961.6496	5935.9937	5927.3018	5964.1246	5941.1095	5949.7992	5948.9486	5944.0513
FORGY APPROACH (FA)	Thời gian	0.00007	0.00003	0.00003	0.00002	0.00024	0.00005	0.00013	0.00047	0.00003	0.00023	0.00013
	Vòng lặp	28	35	23	11	28	36	21	26	13	25	24.6
	SSE	5948.3647	5930.8015	5986.3016	5944.6580	5949.7843	5939.3301	5944.7821	5950.6811	5940.1967	5948.6813	5948.3581
K-MEANS++	Thời gian	5.80787	5.8044	5.81957	5.80238	5.82219	5.77378	5.79221	5.79412	5.82603	5.8484	5.809095
	Vòng lặp	14	20	17	36	20	18	20	20	15	26	20.6
	SSE	5970.2789	5949.6351	5960.8560	5948.4690	5950.7241	5930.0117	5959.1173	5949.6351	5956.6539	5954.3616	5952.9743
MACQUEEN APPROACH (MA)	Thời gian	1.37322	1.35456	1.28899	1.29595	1.32568	1.30472	1.32406	1.3608	1.29773	1.33477	1.326048
	Vòng lặp	26	27	39	19	24	16	22	24	8	22	22.7
	SSE	5953.0583	6031.5797	5982.3435	5924.9864	6026.1666	5965.8175	5945.8815	5985.3529	5997.7372	5972.4951	5978.5419
KAUFMAN APPROACH (KA)	Thời gian	N = 200	265.9639									
	Vòng lặp		5									
	SSE		1139.268									

Hình 35: Bảng thống kê kết quả của 5 kỹ thuật khởi tạo trung tâm cụm đầu tiên qua 10 lần chạy

Các thông số đánh giá trung bình qua 10 lần chạy của 4 kỹ thuật đều tương đồng nhau với sai số SSE chênh lệch với nhãn ban đầu rất ít (6093.73051). Với bộ dữ liệu lớn và nhiều chiều, phương pháp K-MEANS++ vẫn cho kết quả rất tốt mặc dù chênh lệch với 3 phương pháp còn lại là không nhiều. Vậy với bộ dữ liệu lớn và nhiều chiều như The Mnist thì ngoại trừ phương pháp KAUFMAN APPROACH đều đem lại hiệu quả phân cụm tốt.

4 Tham khảo

Tài liệu

- [1] *K-means Clustering*, <https://machinelearningcoban.com/2017/01/01/kmeans/>
- [2] *K-means Clustering: Simple Applications*, <https://machinelearningcoban.com/2017/01/04/kmeans2/>
- [3] *An empirical comparison of four initialization methods for the K-means algorithm*

5 Phụ lục

```
1 #Import module
2 !pip install python-mnist
3 from google.colab import drive
4 drive.mount("/content/gdrive")
5 from __future__ import print_function
6 import numpy as np
7 import matplotlib.pyplot as plt
8 import math
9 import time
10 np.random.seed(11)
11 from mpl_toolkits.mplot3d import Axes3D
12 from mnist import MNIST
13 from sklearn.neighbors import NearestNeighbors
14 from sklearn.preprocessing import normalize
15
16 #K-means and other functions
17 # function for plot
18 def plot(data, centroids):
19     plt.scatter(data[:, 0], data[:, 1], marker = '.',
20                color = 'gray', label = 'data points')
21     plt.scatter(centroids[:-1, 0], centroids[:-1, 1],
22                color = 'black', label = 'previously selected
23                centroids')
24     plt.scatter(centroids[-1, 0], centroids[-1, 1],
25                color = 'red', label = 'next centroid')
26     plt.title('Select % d th centroid'%(centroids.shape[0]))
27     plt.legend()
28     plt.show()
29
30 def kmeans_display(X, label):
31     K = np.amax(label) + 1
32     X0 = X[label == 0, :]
33     X1 = X[label == 1, :]
34     X2 = X[label == 2, :]
35
36     plt.plot(X0[:, 0], X0[:, 1], 'b~', markersize = 4, alpha = .8)
37     plt.plot(X1[:, 0], X1[:, 1], 'go', markersize = 4, alpha = .8)
38     plt.plot(X2[:, 0], X2[:, 1], 'rs', markersize = 4, alpha = .8)
39     plt.xlabel(iris.feature_names[0])
40     plt.ylabel(iris.feature_names[1])
41     plt.axis('equal')
42     plt.plot()
43     plt.show()
44
45 def plot_Iris_Dataset_2D_3D(X, labels):
46     # plot 2D
```

```
46 kmeans_display(X, labels)
47
48 # plot 3D
49 fig = plt.figure(1, figsize=(8, 6))
50 ax = Axes3D(fig, elev=-150, azimuth=160)
51 ax.scatter(X[:, 3], X[:, 1], X[:, 2], c=labels, cmap=plt.cm.
Set1, edgecolor='k', s=40)
52 ax.set_title("Iris Dataset")
53 ax.set_xlabel(iris.feature_names[3])
54 ax.w_xaxis.set_ticklabels([])
55 ax.set_ylabel(iris.feature_names[1])
56 ax.w_yaxis.set_ticklabels([])
57 ax.set_zlabel(iris.feature_names[2])
58 ax.w_zaxis.set_ticklabels([])
59 plt.show()
60
61 # function support for k-means
62 def kmeans_assign_labels(init_function, X, centers):
63     # calculate distances btw data and centers
64     k = len(centers)
65     while True:
66         d = []
67         for center in centers:
68             d.append(list(map(lambda x: math.sqrt(sum(x -
center)**2)), (X -
center)**2)))
69         d = np.array(d).T
70         # index of the closest center
71         d = np.array(list(map(lambda x: np.argmin(x), d)))
72         # sometime some centers not have any instances so init
again with same init k-means
73         if len(set(d)) == k:
74             break
75         centers = init_function(X, k)
76     return d
77
78 def kmeans_update_centers(X, labels, K):
79     centers = np.zeros((K, X.shape[1]))
80     for k in range(K):
81         # collect all points assigned to the k-th cluster
82         Xk = X[labels == k, :]
83         # take average
84         centers[k, :] = np.mean(Xk, axis = 0)
85     return centers
86
87 def has_converged(centers, new_centers):
88     # return True if two sets of centers are the same
89     return (set([tuple(a) for a in centers]) ==
90             set([tuple(a) for a in new_centers]))
91
```

```
92 def SSE(X, labels):
93     SSErr = 0
94     for i in range(len(set(labels))):
95         Xi = X[labels == i, :]
96         centeri = np.mean(Xi, axis = 0)
97         SSErr += sum(list(map(lambda x: math.sqrt(sum(x)), (Xi -
98             centeri)**2)))
99         print('Error Sum of Squares: ', round(SSErr, 5))
100
101 def kmeans(init_function, X, K):
102     centers = init_function(X, K)
103     labels = []
104     it = 0
105     while True:
106         labels = kmeans_assign_labels(init_function, X, centers)
107         new_centers = kmeans_update_centers(X, labels, K)
108         if has_converged(centers, new_centers):
109             break
110         centers = new_centers
111         it += 1
112     return (centers, labels, it)
113
114 #IRIS FLOWER dataset
115 plot_Iris_Dataset_2D_3D(iris_X, iris_y)
116 SSE(iris_X, iris_y)
117
118 #Run with different init methods
119 #FORGY APPROACH (FA)
120 centers, labels, it = kmeans(kmeans_init_randomly_centers, iris_X,
121     K=3)
122 print('iterator:', it)
123 plot_Iris_Dataset_2D_3D(iris_X, labels)
124 SSE(iris_X, labels)
125
126 #K-MEANS++
127 centers, labels, it = kmeans(kmeans_init_kmeansPlusPlus, iris_X, K
128     =3)
129 print('iterator:', it)
130 plot_Iris_Dataset_2D_3D(iris_X, labels)
131 SSE(iris_X, labels)
132
133 #RANDOM K CLUSTERS
134 centers, labels, it = kmeans(kmeans_init_randomly_k_clusters,
135     iris_X, K=3)
136 print('iterator:', it)
137 plot_Iris_Dataset_2D_3D(iris_X, labels)
138 SSE(iris_X, labels)
139
140 #MACQUEEN APPROACH (MA)
```

```
137 centers, labels, it = kmeans(kmeans_init_MA, iris_X, K=3)
138 print('iterator:', it)
139 plot_Iris_Dataset_2D_3D(iris_X, labels)
140 SSE(iris_X, labels)
141
142 #KAUFMAN APPROACH (KA)
143 centers, labels, it = kmeans(kmeans_init_KA, iris_X, K=3)
144 print('iterator:', it)
145 plot_Iris_Dataset_2D_3D(iris_X, labels)
146 SSE(iris_X, labels)
147
148 #THE MNIST dataset
149 #Load THE MNIST dataset and functions
150 def display_network(A, m = -1, n = -1):
151     opt_normalize = True
152     opt_graycolor = True
153
154     # Rescale
155     A = A - np.average(A)
156
157     # Compute rows & cols
158     (row, col) = A.shape
159     sz = int(np.ceil(np.sqrt(row)))
160     buf = 1
161     if m < 0 or n < 0:
162         n = np.ceil(np.sqrt(col))
163         m = np.ceil(col / n)
164
165
166     image = np.ones(shape=(buf + m * (sz + buf), buf + n * (sz +
167 buf)))
168
169     if not opt_graycolor:
170         image *= 0.1
171
172     k = 0
173
174     for i in range(int(m)):
175         for j in range(int(n)):
176             if k >= col:
177                 continue
178
179             clim = np.max(np.abs(A[:, k]))
180
181             if opt_normalize:
182                 image[buf + i * (sz + buf):buf + i * (sz + buf) +
183 sz, buf + j * (sz + buf):buf + j * (sz + buf) + sz] = \
184                 A[:, k].reshape(sz, sz) / clim
185             else:
```

```
184         image[buf + i * (sz + buf):buf + i * (sz + buf) +
185             sz, buf + j * (sz + buf):buf + j * (sz + buf) + sz] = \
186             A[:, k].reshape(sz, sz) / np.max(np.abs(A))
187         k += 1
188     return image
189
190 def display_MNIST(X, centers, labels):
191     A = display_network(centers.T, 10, 1)
192     f1 = plt.imshow(A, interpolation='nearest', cmap = "jet")
193     f1.axes.get_xaxis().set_visible(False)
194     f1.axes.get_yaxis().set_visible(False)
195     plt.show()
196
197     K = len(centers)
198     NO = 10;
199     X1 = np.zeros((NO*K, 784))
200     X2 = np.zeros((NO*K, 784))
201     for k in range(K):
202         Xk = X[labels == k, :]
203         center_k = [centers[k]]
204         neigh = NearestNeighbors(NO).fit(Xk)
205         dist, nearest_id = neigh.kneighbors(center_k, NO)
206         X1[NO*k: NO*k + NO, :] = Xk[nearest_id, :]
207         X2[NO*k: NO*k + NO, :] = Xk[:NO, :]
208     plt.axis('off')
209     A = display_network(X2.T, K, NO)
210     f2 = plt.imshow(A, interpolation='nearest' )
211     plt.gray()
212     plt.show()
213
214 mndata = MNIST('/content/gdrive/MyDrive/Colab Notebooks/MNIST')
215 mndata.load_testing()
216 X = np.asarray(mndata.test_images)[:1000,:]/256.0
217 print(X.shape)
218
219 #Run with different init methods
220 #FORGY APPROACH (FA)
221 centers, labels, it = kmeans(kmeans_init_randomly_centers,X, K=10)
222 print('iterator:', it)
223 SSE(X, labels)
224 display_MNIST(X, centers, labels)
225
226 #K-MEANS++
227 centers, labels, it = kmeans(kmeans_init_kmeansPlusPlus,X, K=10)
228 print('iterator:', it)
229 SSE(X, labels)
230 display_MNIST(X, centers, labels)
231
```

```
232 #RANDOM K CLUSTERS
233 centers, labels, it = kmeans(kmeans_init_randomly_k_clusters,X, K
    =10)
234 print('iterator:', it)
235 SSE(X, labels)
236 display_MNIST(X, centers, labels)
237
238 #MACQUEEN APPROACH (MA)
239 centers, labels, it = kmeans(kmeans_init_MA,X, K=10)
240 print('iterator:', it)
241 SSE(X, labels)
242 display_MNIST(X, centers, labels)
243
244 #KAUFMAN APPROACH (KA)
245 centers, labels, it = kmeans(kmeans_init_KA,X[:200,:], K=10)
246 print('iterator:', it)
247 SSE(X[:200,:], labels)
248 display_MNIST(X[:200,:], centers, labels)
```