

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



BTLT TUẦN 7: KHAI THÁC DỮ LIỆU

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

Ngày 26 tháng 5 năm 2024

MỤC LỤC

1	Nhập dữ liệu và tiền xử lý	1
2	Thuật toán gom cụm phân cấp từ thư viện	3
3	Thuật toán gom cụm phân cấp từ cài đặt	6
3.1	Tính hàm tiêu chuẩn liên kết (linkage criteria)	6
3.2	Class Agglomerative Clustering tự cài đặt	7
4	So sánh kết quả	9
4.1	Single Linkage	10
4.2	Complete Linkage	11
4.3	Average Linkage	12
4.4	Ward Linkage	13
5	Kết luận	14

1. Nhập dữ liệu và tiền xử lý

Nhập thư viện

```
1 import numpy as np
2 import pandas as pd
3 import seaborn as sns
4 import scipy.cluster.hierarchy as shc
5 import matplotlib.pyplot as plt
6 from sklearn.decomposition import PCA
7 from sklearn.cluster import AgglomerativeClustering
```

Đọc dữ liệu từ data.csv

```
1 path_to_file = "D:\\Courses HK2 23-24\\Data Mining\\LAB\\LAB 07\\data.csv"
2 customer_data = pd.read_csv(path_to_file)
3
4 customer_data.shape
```

(200,5)

```
1 customer_data.columns
```

Index(['CustomerID', 'Genre', 'Age', 'Annual Income (k\$)', 'Spending Score (1-100)'], dtype='object')

```
1 customer_data.describe().transpose()
```

	count	mean	std	min	25%	50%	75%	max
CustomerID	200.0	100.50	57.879185	1.0	50.75	100.5	150.25	200.0
Age	200.0	38.85	13.969007	18.0	28.75	36.0	49.00	70.0
Annual Income (k\$)	200.0	60.56	26.264721	15.0	41.50	61.5	78.00	137.0
Spending Score (1-100)	200.0	50.20	25.823522	1.0	34.75	50.0	73.00	99.0

```
1 customer_data.head()
```

	CustomerID	Genre	Age	Annual Income (k\$)	Spending Score (1-100)
0	1	Male	19	15	39
1	2	Male	21	15	81
2	3	Female	20	16	6
3	4	Female	23	16	77
4	5	Female	31	17	40

Chia cột Age thành 10 nhóm khác nhau

```

1 intervals = [15, 20, 30, 40, 50, 60, 70]
2 col = customer_data["Age"]
3 customer_data["Age Groups"] = pd.cut(x = col, bins = intervals)
4 customer_data["Age Groups"]

0      (15, 20]
1      (20, 30]
2      (15, 20]
3      (20, 30]
4      (30, 40]
...
195     (30, 40]
196     (40, 50]
197     (30, 40]
198     (30, 40]
199     (20, 30]
Name: Age Groups, Length: 200, dtype: category
Categories (6, interval[int64, right]): [(15, 20] < (20, 30] < (30, 40] < (40, 50] < (50, 60] < (60, 70]]

```

```

1 customer_data.groupby("Age Groups", observed=True)["Age Groups"].count()

```

```

Age Groups
(15, 20]    17
(20, 30]    45
(30, 40]    60
(40, 50]    38
(50, 60]    23
(60, 70]    17
Name: Age Groups, dtype: int64

```

Chuyển 2 cột Age và Genre thành dạng số

```

1 customer_data_oh = pd.get_dummies(customer_data, dtype= int)
2 customer_data_oh

```

	CustomerID	Age	Annual Income (k\$)	Spending Score (1-100)	Genre_Female	Genre_Male	Age Groups_(15, 20]	Age Groups_(20, 30]	Age Groups_(30, 40]	Age Groups_(40, 50]	Age Groups_(50, 60]	Age Groups_(60, 70]
0	1	19	15	39	0	1	1	0	0	0	0	0
1	2	21	15	81	0	1	0	1	0	0	0	0
2	3	20	16	6	1	0	1	0	0	0	0	0
3	4	23	16	77	1	0	0	1	0	0	0	0
4	5	31	17	40	1	0	0	0	1	0	0	0
...
195	196	35	120	79	1	0	0	0	1	0	0	0
196	197	45	126	28	1	0	0	0	0	1	0	0
197	198	32	126	74	0	1	0	0	1	0	0	0
198	199	32	137	18	0	1	0	0	1	0	0	0
199	200	30	137	83	0	1	0	1	0	0	0	0

Bỏ cột CustomerID và vẽ plot

```

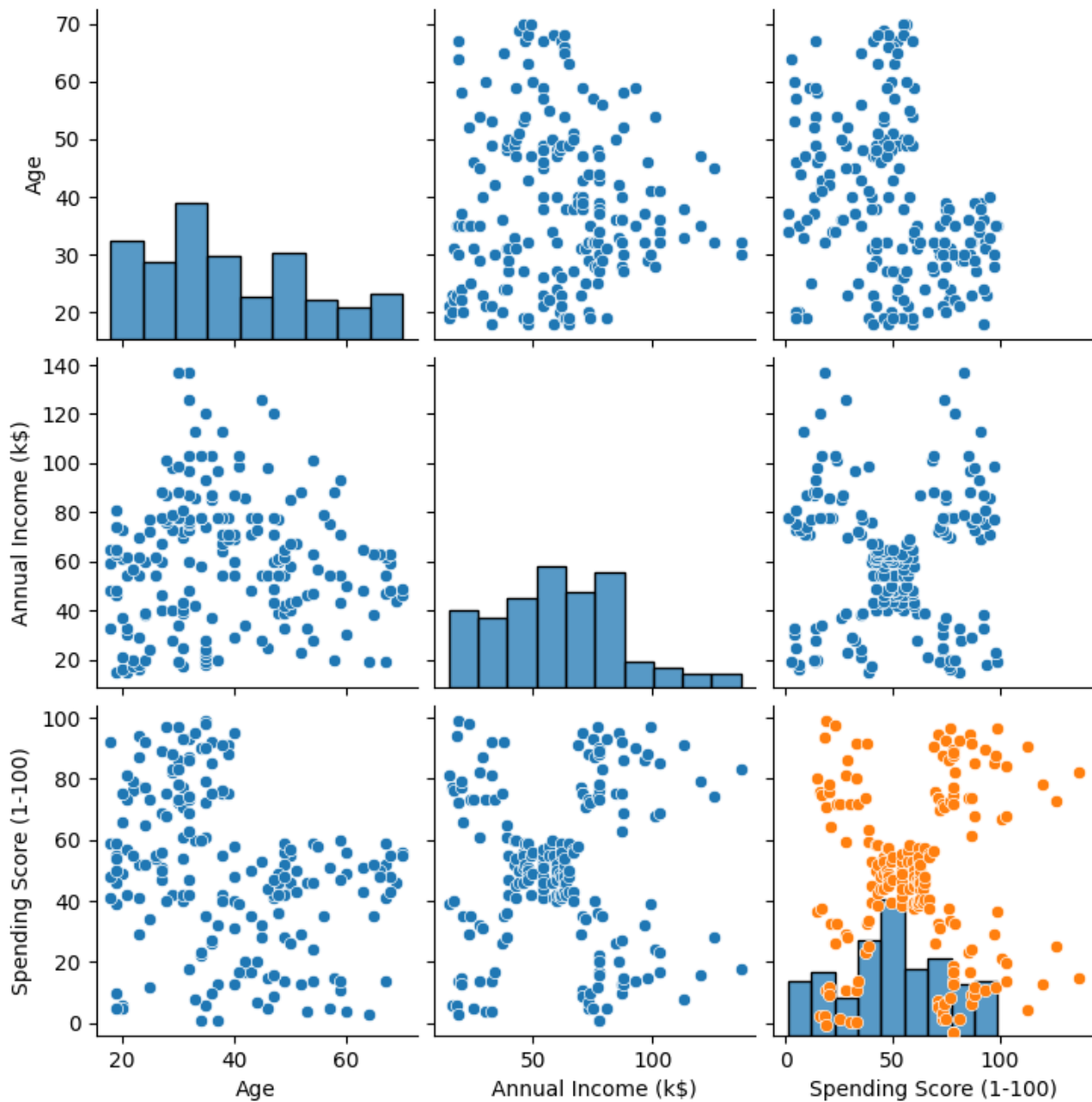
1 customer_data = customer_data.drop("CustomerID", axis= 1)

```

```

1 sns.pairplot(customer_data)
2
3 sns.scatterplot(x =customer_data["Annual Income (k$)"],
4                 y = customer_data["Spending Score (1-100)"])

```



2. Thuật toán gom cụm phân cấp từ thư viện

Bỏ cột 'Age' và vẽ 'dendrogram'

```

1 customer_data_oh = customer_data_oh.drop(["Age"], axis= 1)
2 customer_data_oh.shape

```

```

1 plt.figure(figsize=(10,7))
2 plt.title("Customer Denrogram")

```

```

3
4 selected_data = customer_data_oh.iloc[:, 1:3]
5 clusters = shc.linkage(selected_data,
6                       method='ward',
7                       metric= 'euclidean')
8 shc.dendrogram(Z=clusters)
9 plt.show()

```

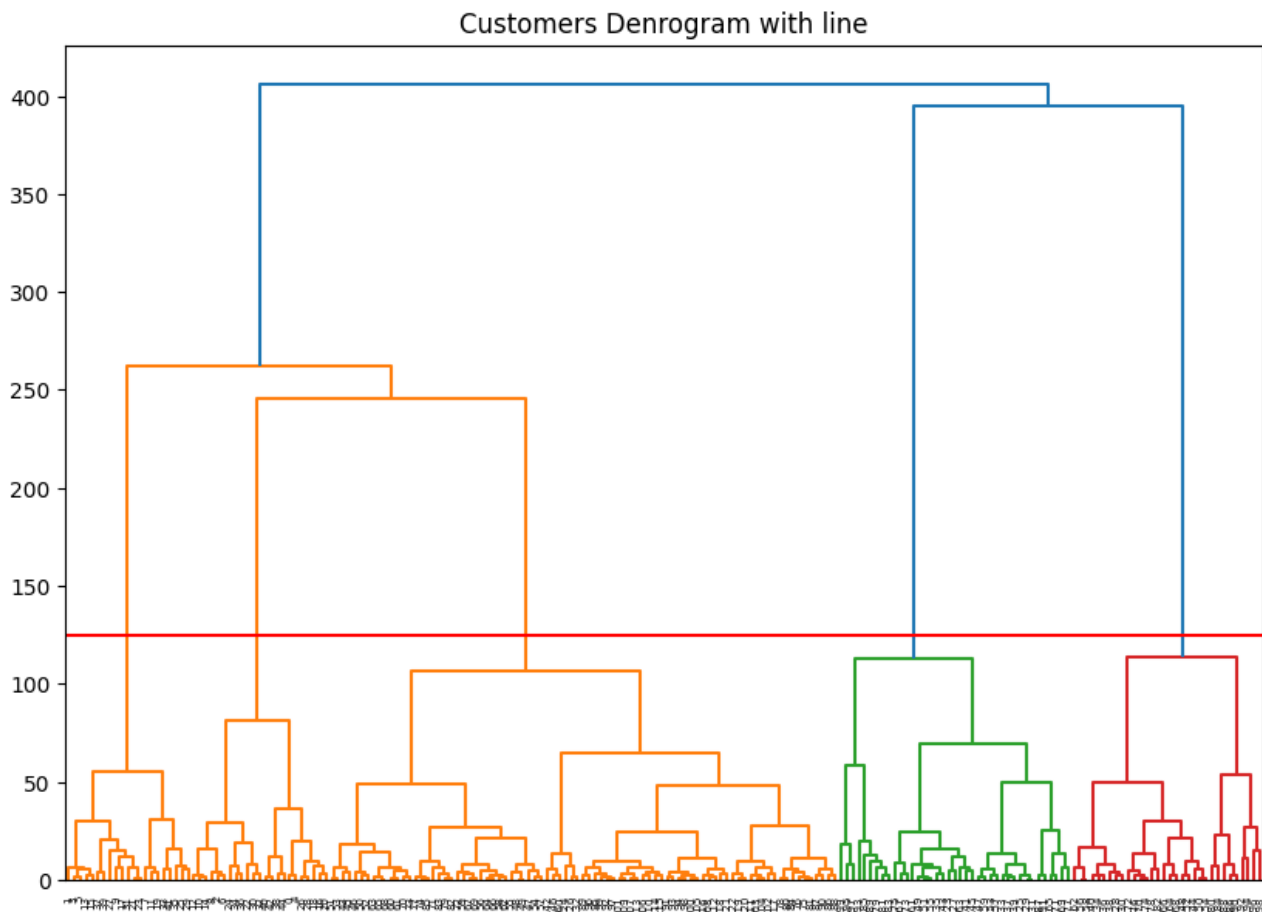


Vẽ đường nằm ngang đi qua khoảng cách dài nhất

```

1 plt.figure(figsize=(10,7))
2 plt.title("Customers Denrogram with line")
3 clusters = shc.linkage(selected_data,
4                       method='ward',
5                       metric= 'euclidean')
6 shc.dendrogram(Z = clusters)
7 plt.axhline(y =125, color = "r", linestyle = "-")
8 plt.show()

```



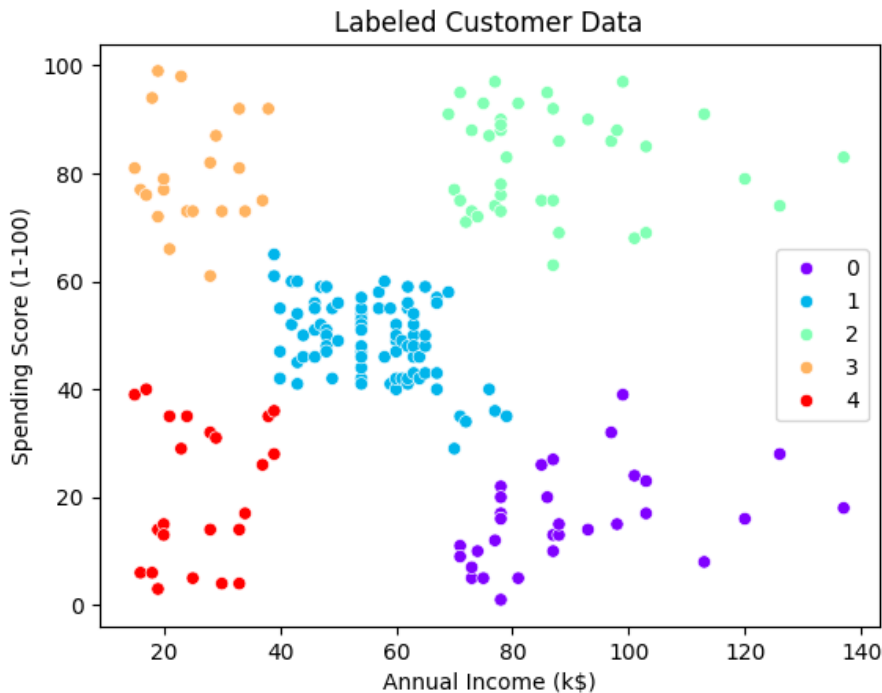
Thực thi phân cụm với dữ liệu ban đầu

```
1 clustering_model = AgglomerativeClustering(n_clusters=5, metric="euclidean", linkage= "ward")
2 clustering_model.fit(selected_data)
3 clustering_model.labels_
```

[illegible]

```
1 data_labels = clustering_model.labels_  
2 sns.scatterplot(x = "Annual Income (k$)",  
3                 y = "Spending Score (1-100)",  
4                 data = selected_data,  
5                 hue= data_labels,  
6                 palette= "rainbow").set_title("Labeled Customer Data")
```

```
plt.show()
```



3. Thuật toán gom cụm phân cấp từ cài đặt

3.1. Tính hàm tiêu chuẩn liên kết (linkage criteria)

Các linkage criteria gồm:

- Complete-linkage:

$$L(r, s) = \max(D(x_{ri}, x_{sj}))$$

- Single-linkage:

$$L(r, s) = \min(D(x_{ri}, x_{sj}))$$

- Average-linkage:

$$L(r, s) = \frac{1}{n_r n_s} \sum_{i=1}^{n_r} \sum_{j=1}^{n_s} D(x_{ri}, x_{sj})$$

- Ward-linkage:

$$L(r, s) = \sum_{x_i \in r \cup s} \|x_i - m\|_2 - \sum_{x_i \in r} \|x_i - m_1\|_2 - \sum_{x_i \in s} \|x_i - m_2\|_2$$

Giả sử m, m_1, m_2 lần lượt là tâm cụm $r \cup s, r, s$.

Với $D(a, b) = \sqrt{(a - b)^2} = \|a - b\|_2$ (Norm 2 hay khoảng cách Euclidean).

```
1 import numpy as np
2
3 def euclid_distance(a, b):
4     # Tính khoảng cách Euclid giữa hai điểm a và b
5     return np.sqrt(np.sum((np.array(a) - np.array(b)) ** 2))
6
7 def single_linkage(r, s):
8     # Tính khoảng cách liên kết đơn giữa hai cụm r và s
9     # Khoảng cách liên kết đơn là khoảng cách ngắn nhất giữa bất kỳ cặp điểm nào từ hai cụm
10    return np.min([euclid_distance(xi, xj) for xi in r for xj in s])
```



```

11
12 def complete_linkage(r, s):
13     # Tính khoảng cách liên kết đầy đủ giữa hai cụm r và s
14     # Khoảng cách liên kết đầy đủ là khoảng cách lớn nhất giữa bất kỳ cặp điểm nào từ hai cụm
15     return np.max([euclid_distance(xi, xj) for xi in r for xj in s])
16
17 def average_linkage(r, s):
18     # Tính khoảng cách liên kết trung bình giữa hai cụm r và s
19     # Khoảng cách liên kết trung bình là khoảng cách trung bình giữa tất cả các cặp điểm từ
20     # hai cụm
21     return np.mean([euclid_distance(xi, xj) for xi in r for xj in s])
22
23 def ward_linkage(r, s):
24     # Tính khoảng cách liên kết Ward giữa hai cụm r và s
25     # Phương pháp liên kết Ward tối thiểu hóa tổng phương sai tăng thêm khi gộp hai cụm
26     r_s = r + s # Gộp hai cụm r và s
27     centroid_r = np.mean(r, axis=0) # Tính tâm cụm r
28     centroid_s = np.mean(s, axis=0) # Tính tâm cụm s
29     centroid_rs = np.mean(r_s, axis=0) # Tính tâm của cụm gộp r_s
30
31     # Tính tổng bình phương khoảng cách từ các điểm trong cụm r đến tâm cụm r
32     ss_r = np.sum([euclid_distance(xi, centroid_r) ** 2 for xi in r])
33
34     # Tính tổng bình phương khoảng cách từ các điểm trong cụm s đến tâm cụm s
35     ss_s = np.sum([euclid_distance(xj, centroid_s) ** 2 for xj in s])
36
37     # Tính tổng bình phương khoảng cách từ các điểm trong cụm gộp r_s đến tâm cụm r_s
38     ss_rs = np.sum([euclid_distance(xij, centroid_rs) ** 2 for xij in r_s])
39
40     # Tính khoảng cách liên kết Ward
41     return ss_rs - (ss_r + ss_s)
42
43 def get_linkage_method(linkage):
44     # Chọn phương pháp liên kết dựa trên tham số đầu vào
45     if linkage == 'single':
46         return single_linkage
47     elif linkage == 'complete':
48         return complete_linkage
49     elif linkage == 'average':
50         return average_linkage
51     elif linkage == 'ward':
52         return ward_linkage
53     else:
54         # Nếu tham số không hợp lệ, báo lỗi
55         raise ValueError("Invalid linkage method. Choose from 'single', 'complete', 'average', 'ward'.")

```

3.2. Class Agglomerative Clustering tự cài đặt

```

1 import matplotlib.pyplot as plt
2 import seaborn as sns
3 import pandas as pd
4
5 class AgglomerativeClusteringFromScratch:
6     def __init__(self, n_clusters, linkage='ward'):
7         self.n_clusters = n_clusters # Số cụm
8         self.linkage = get_linkage_method(linkage) # Lấy hàm linkage ứng với tham số linkage
9         self.method = linkage # Lưu lại linkage method để trực quan
10
11     # Hàm khởi tạo các phân cụm

```

```

12     def init_clusters(self, data):
13         return {data_id: [data_point] for data_id, data_point in enumerate(data)} # Mỗi điểm
           dữ liệu là một cụm riêng biệt
14
15     # Hàm tìm cặp cụm gần nhất để hợp nhất
16     def find_closest_clusters(self):
17         # Lấy danh sách các ID của các cụm hiện tại
18         clusters_ids = list(self.clusters.keys())
19         # Tạo tất cả các cặp có thể từ các cụm hiện tại
20         cluster_combinations = [(clusters_ids[i], clusters_ids[j])
21                                 for i in range(len(clusters_ids)-1)
22                                 for j in range(i+1, len(clusters_ids))]
23         # Tính khoảng cách giữa các cụm và chọn cặp gần nhất
24         distances = [self.linkage(self.clusters[xi], self.clusters[xj])
25                     for xi, xj in cluster_combinations]
26         closest_index = np.argmin(distances) # Tìm chỉ số của cặp cụm có khoảng cách nhỏ
           nhất
27         closest_clusters = cluster_combinations[closest_index] # Lấy cặp cụm gần nhất
28         return closest_clusters # Trả về cặp cụm gần nhất
29
30     # Hàm hợp nhất các cụm và tạo cụm mới
31     def merge_and_form_new_clusters(self, xi_id, xj_id):
32         new_clusters = {0: self.clusters[xi_id] + self.clusters[xj_id]} # Tạo một cụm mới
           bằng cách hợp nhất hai cụm gần nhất
33         new_cluster_id = 1 # Khởi tạo ID cho các cụm mới
34         # Duyệt qua tất cả các cụm hiện tại và hợp nhất cụm nếu không phải là cụm được chọn
           để hợp nhất
35         for cluster_id, points in self.clusters.items():
36             if cluster_id != xi_id and cluster_id != xj_id:
37                 new_clusters[new_cluster_id] = points # Thêm các cụm không bị hợp nhất vào
           dictionary mới
38                 new_cluster_id += 1 # Tăng ID cụm mới
39         self.clusters = new_clusters # Cập nhật lại các cụm mới
40
41     # Hàm huấn luyện mô hình
42     def fit(self, data):
43         self.data = data.values.tolist() # Chuyển đổi DataFrame thành danh sách các điểm dữ
           liệu
44         self.col_names = data.columns # Lưu tên của các cột dữ liệu
45         self.clusters = self.init_clusters(self.data) # Khởi tạo các cụm ban đầu từ dữ liệu
46
47         # Tiếp tục hợp nhất các cụm cho đến khi đạt được số lượng cụm mong muốn
48         while len(self.clusters.keys()) > self.n_clusters:
49             closest_clusters = self.find_closest_clusters() # Tìm cặp cụm gần nhất
50             self.merge_and_form_new_clusters(*closest_clusters) # Hợp nhất cặp cụm gần nhất
51
52         # Gán nhãn cho từng điểm dữ liệu dựa trên các cụm
53         self.labels = np.array([None] * len(data)) # Khởi tạo mảng nhãn cho các điểm dữ liệu
54         for i, d in enumerate(self.data):
55             for cluster_id, points in self.clusters.items():
56                 if d in points:
57                     self.labels[i] = cluster_id # Gán nhãn cho điểm dữ liệu
58                     break # Thoát khỏi vòng lặp khi đã tìm thấy cụm chứa điểm dữ liệu
59
60     # Hàm trực quan hóa dữ liệu
61     def visualize(self, show_dendrogram=False):
62         data = pd.DataFrame(self.data, columns=self.col_names) # Tạo DataFrame từ dữ liệu đã
           chuyển đổi
63
64         if show_dendrogram: # 2 plot
65             fig, axes = plt.subplots(1, 2, figsize=(24, 7)) # Tạo subplot cho cả scatter

```

```

66     plot và dendrogram
67     fig.suptitle(f'Cluster Analysis From Scratch Using {self.method.capitalize()}
68     Linkage', fontsize=20, fontweight='bold') # Đặt tiêu đề chung cho cả hai
69     subplot
70
71     # Trực quan hóa dữ liệu cụm
72     sns.scatterplot(ax=axes[0], x=self.col_names[0], y=self.col_names[1],
73                     data=data, hue=self.labels, palette="rainbow").set_title("
74                     Visualization of Clustered Data") # Vẽ scatter plot của dữ
75                     liệu cụm
76     axes[0].legend(title="Clusters", loc="upper right") # Thêm chú thích cho các
77                     nhãn cụm
78
79     # Trực quan hóa dendrogram
80     axes[1].set_title("Dendrogram Plot with Line") # Đặt tiêu đề cho dendrogram
81     clusters = shc.linkage(data, method=self.method, metric='euclidean') # Tạo liên
82                     kết cho dendrogram
83     shc.dendrogram(Z=clusters, ax=axes[1]) # Vẽ dendrogram trên subplot thứ hai
84     max_y = max(shc.maxdists(clusters)) # Tìm giá trị y lớn nhất trên dendrogram
85     axes[1].axhline(y=0.7 * max_y, color="black", linestyle="--") # Vẽ đường ngang
86                     trên dendrogram
87
88     plt.show() # Hiển thị plot
89
90     else: # 1 plot
91         # Trực quan hóa dữ liệu cụm
92         sns.scatterplot(x=self.col_names[0], y=self.col_names[1],
93                         data=data, hue=self.labels, palette="rainbow").set_title("
94                         Visualization of Clustered Data From Scratch") # Vẽ scatter
95                         plot của dữ liệu cụm
96         plt.show() # Hiển thị plot

```

4. So sánh kết quả

Chuẩn bị data

```

1 df = selected_data.copy()
2 df

```

	Annual Income (k\$)	Spending Score (1-100)
0	15	39
1	15	81
2	16	6
3	16	77
4	17	40
...
195	120	79
196	126	28
197	126	74
198	137	18
199	137	83

200 rows × 2 columns

Tạo hàm tính Agglomerative Clustering từ thư viện Sklearn.

Sau đó tạo hàm vẽ biểu đồ trực quan cho kết quả (gồm 2 plot scatter cluster và dendrogram giống hàm tự code). Dùng để so sánh với hàm tự cài đặt ở trên

```

1  def sklearn_agglomerative_clustering(df, n_clusters, linkage):
2      # Khởi tạo mô hình phân cụm phân cấp từ thư viện sklearn
3      clustering_model = AgglomerativeClustering(n_clusters=n_clusters, metric="euclidean",
4          linkage=linkage)
5      # Huấn luyện mô hình với dữ liệu đầu vào
6      clustering_model.fit(df)
7      # Trả về mô hình đã được huấn luyện
8      return clustering_model
9
10 def show_plot(df, labels, linkage):
11     # Tạo một figure với 2 subplot, kích thước là 24x7 inch
12     fig, axes = plt.subplots(1, 2, figsize=(24, 7))
13     # Đặt tiêu đề chính cho figure
14     fig.suptitle(f'Cluster Analysis From Library Using {linkage.capitalize()} Linkage',
15         fontsize=20, fontweight='bold')
16
17     # Vẽ biểu đồ scatter plot cho dữ liệu đã được phân cụm
18     sns.scatterplot(
19         ax=axes[0], # Vẽ trên subplot đầu tiên
20         x="Annual Income (k$)", # Trục x là thu nhập hàng năm
21         y="Spending Score (1-100)", # Trục y là điểm chỉ tiêu
22         data=df, # Dữ liệu đầu vào
23         hue=labels, # Nhãn cụm
24         palette="rainbow" # Sử dụng bảng màu rainbow
25     ).set_title("Visualization of Clustered Data") # Đặt tiêu đề cho biểu đồ
26     # Thêm chú thích cho biểu đồ scatter plot
27     axes[0].legend(title="Clusters", loc="upper right")
28
29     # Vẽ biểu đồ dendrogram
30     axes[1].set_title("Dendrogram Plot with Line") # Đặt tiêu đề cho dendrogram
31     clusters = shc.linkage(df, method=linkage, metric='euclidean') # Tính toán liên kết giữa
32     # các điểm dữ liệu
33     shc.dendrogram(Z=clusters, ax=axes[1]) # Vẽ dendrogram
34     max_y = max(shc.maxdists(clusters)) # Tìm giá trị y lớn nhất trên dendrogram
35     axes[1].axhline(y=0.7 * max_y, color="black", linestyle="--") # Vẽ đường ngang trên
36     # dendrogram tại vị trí 70% giá trị y lớn nhất
37
38     # Hiển thị figure với 2 biểu đồ
39     plt.show()

```

4.1. Single Linkage

Hàm thư viện

```

1  model = sklearn_agglomerative_clustering(df, n_clusters = 5, linkage = "single")
2  labels = model.labels_
3  labels

```

Hàm tự code

```

1  agg_model = AgglomerativeClusteringFromScratch(n_clusters=5, linkage="single")
2  agg_model.fit(df)
3  agg_model.labels

```

Hai kết quả hoàn toàn giống nhau (tên cụm khác nhau nhưng ý nghĩa như nhau)
Trực quan của thư viện

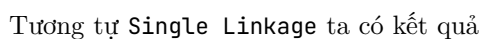
```
array([1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1,  
       1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1,  
       1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1,  
       1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1,  
       1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1,  
       1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1,  
       1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1,  
       1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1,  
       1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 0, 1, 0, 2, 0, 2,  
       3, 4], dtype=object)
```

Hình 2: Scratch

1

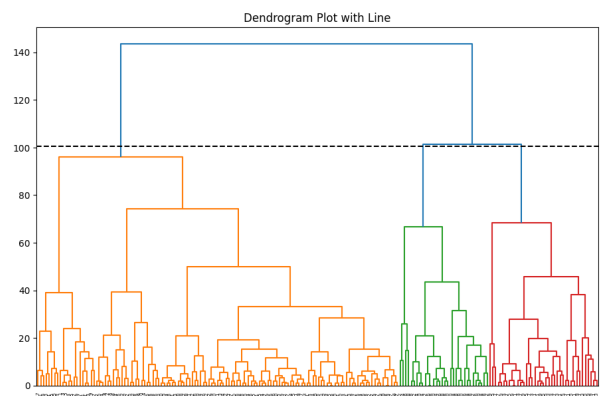
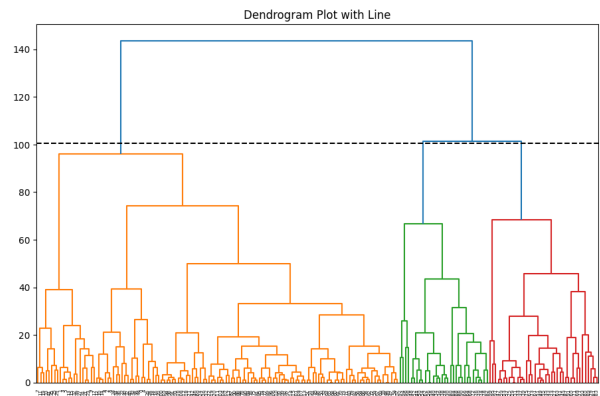


1

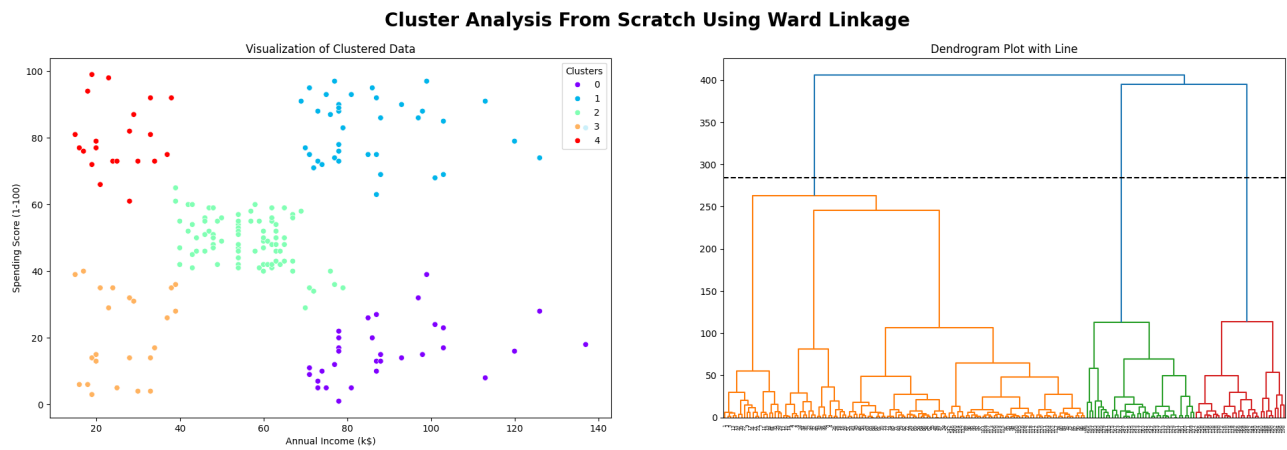


[illegible]

Hình 4: Scratch

[illegible]

Hình 6: Scratch



5. Kết luận

Như vậy kết quả em tự code hoàn toàn giống với kết quả của thư viện.

Cả hai phương pháp đều cho phép phân cụm dữ liệu và trực quan hóa kết quả phân cụm.

Việc trực quan hóa giúp chúng ta hiểu rõ hơn về cấu trúc của dữ liệu và quá trình phân cụm.