

CLUSTERING STATES INTO DIFFERENT COVID-19 ZONES

Hoàng-Nguyên Vũ

1. Mô tả

Bài tập thực hành KMeans nhằm phân vùng bệnh Covid-19 trong bài toán thực tế về dịch bệnh Sar-Cov2 năm 2021 là một ứng dụng quan trọng trong phân tích dữ liệu y tế. Mục tiêu của bài tập là sử dụng thuật toán KMeans để phân nhóm các khu vực hoặc quốc gia dựa trên các chỉ số quan trọng như số ca nhiễm, số ca hồi phục và số ca tử vong. Qua việc phân cụm, ta có thể xác định các khu vực có đặc điểm dịch tễ học tương đồng, từ đó hỗ trợ các nhà quản lý y tế đưa ra quyết định hiệu quả hơn trong việc phân bổ nguồn lực và triển khai các biện pháp phòng chống dịch phù hợp. Đây là một ví dụ điển hình về cách sử dụng machine learning để giải quyết các vấn đề thực tiễn trong y tế cộng đồng.

2. Nội dung

Chúng ta sẽ sử dụng thuật toán K-Means, đây là một thuật toán phân cụm không giám sát. Thuật toán k-means là một thuật toán lặp đi lặp lại, cố gắng chia tập dữ liệu thành K cụm (nhóm con) riêng biệt không chồng chéo, trong đó mỗi điểm dữ liệu chỉ thuộc về một cụm. Thuật toán cố gắng làm cho các điểm dữ liệu trong cùng một cụm càng giống nhau càng tốt, đồng thời giữ cho các cụm càng xa nhau càng tốt. Nó gán các điểm dữ liệu vào một cụm sao cho tổng khoảng cách bình phương giữa các điểm dữ liệu và tâm cụm là nhỏ nhất.

Bước 1: Load dataset: Tải tại đây

```
1 import numpy as np
2 import pandas as pd
3 from sklearn.cluster import KMeans
4 from matplotlib import pyplot as plt
5
6 df_india = # Your coder here to load dataset#
7 print(df_india)
```

	State	Confirmed	Recovered	Deaths
0	Maharashtra	41642	11726	1454
1	Tamil Nadu	13967	6282	95
2	Gujarat	12910	5488	773
3	Delhi	11659	5567	194
4	Rajasthan	6227	3485	151
5	Madhya Pradesh	5981	2844	271

Bước 2: Thực hiện tính tỉ lệ hồi phục và tử vong cho mỗi khu vực:

Chúng ta sẽ thực hiện tính tỉ lệ hồi phục và tử vong từ tập dữ liệu trên theo công thức sau:

$$\text{Tỉ lệ hồi phục} = \frac{\text{Hồi phục}}{\text{Số ca xác nhận}} * 100\% \quad (1)$$

$$\text{Tỉ lệ tử vong} = \frac{\text{Tử vong}}{\text{Số ca xác nhận}} * 100\% \quad (2)$$

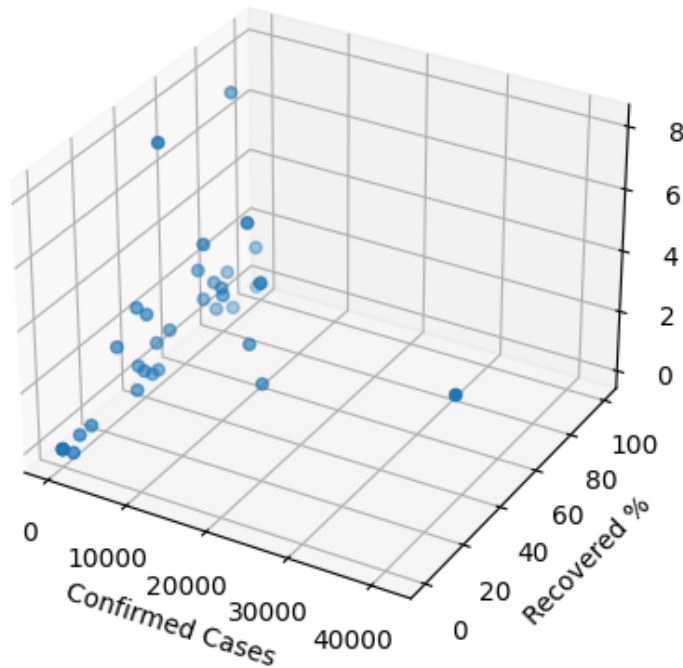
```
1 # Calculate the 'Recovered' and 'Deaths' percentages, ensuring 'Confirmed'
  is not zero
2 df_india['Recovered'] = # Your code here #
3 df_india['Deaths'] = # Your code here #
4 print(df_india)
```

	State	Confirmed	Recovered	Deaths
0	Maharashtra	41642	28.159070	3.491667
1	Tamil Nadu	13967	44.977447	0.680175
2	Gujarat	12910	42.509682	5.987607
3	Delhi	11659	47.748520	1.663951
4	Rajasthan	6227	55.965955	2.424924
5	Madhya Pradesh	5981	47.550577	4.531015

Hình 1: Kết quả sau khi tính tỉ lệ

Bước 3: Trực quan hóa dữ liệu trước khi áp dụng K-Means:

```
1 from mpl_toolkits.mplot3d import Axes3D
2
3 fig = plt.figure()
4 ax = fig.add_subplot(111, projection='3d')
5
6 ax.scatter(df_india['Confirmed'], df_india['Recovered'], df_india['Deaths']
7            ])
8 ax.set_xlabel('Confirmed Cases')
9 ax.set_ylabel('Recovered %')
10 ax.set_zlabel('Deaths %')
11 plt.show()
```

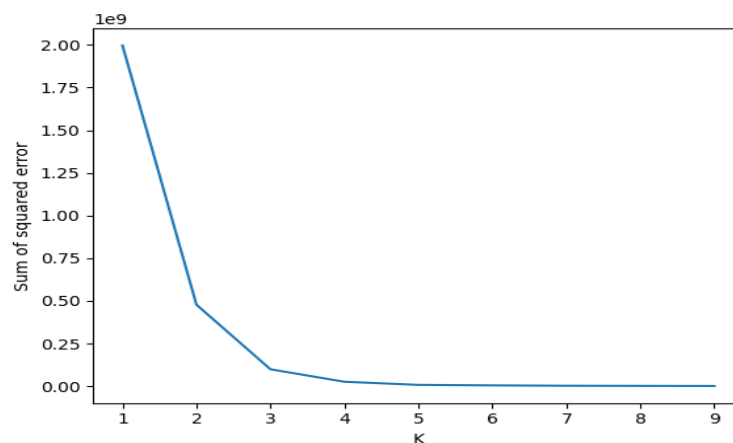


Hình 2: Phân bố data covid trước khi phân cụm

Bước 4: Cài đặt K-Means:

Chúng ta sẽ áp dụng mô hình K-Means và áp dụng phương pháp Elbow để chọn K tốt nhất trong tập dữ liệu này:

```
1 sse = []  
2 k_rng = range(1,10)  
3 # Your code here #
```



Hình 3: Kết quả Elbow - Qua đồ thị này dễ thấy K = 4 là kết quả tối ưu

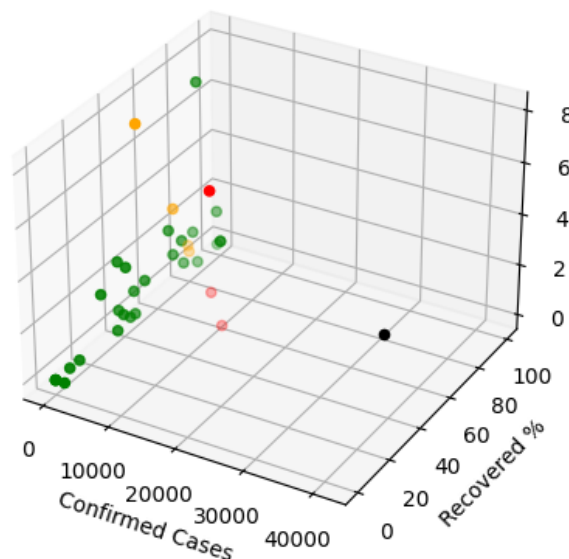
Bước 5: Mapping kết quả cluster vào dữ liệu để trực quan hóa trên bản đồ thế giới:

Chúng ta sẽ mapping kết quả cluster trên với $K = 4$ vào dữ liệu với cột mới có tên là **Zone**, giá trị cột này là giá trị clustering từ mô hình K-Mean với $K = 4$

```
1 # Your code here #
2 print(df_india)
```

	State	Confirmed	Recovered	Death	Zone
0	Maharashtra	41642	28.159070	3.491667	3
1	Tamil Nadu	13967	44.977447	0.680175	2
2	Gujarat	12910	42.509682	5.987607	2
3	Delhi	11659	47.748520	1.663951	2
4	Rajasthan	6227	55.965955	2.424924	1
5	Madhya Pradesh	5981	47.550577	4.531015	1

Hình 4: Kết quả mapping zone



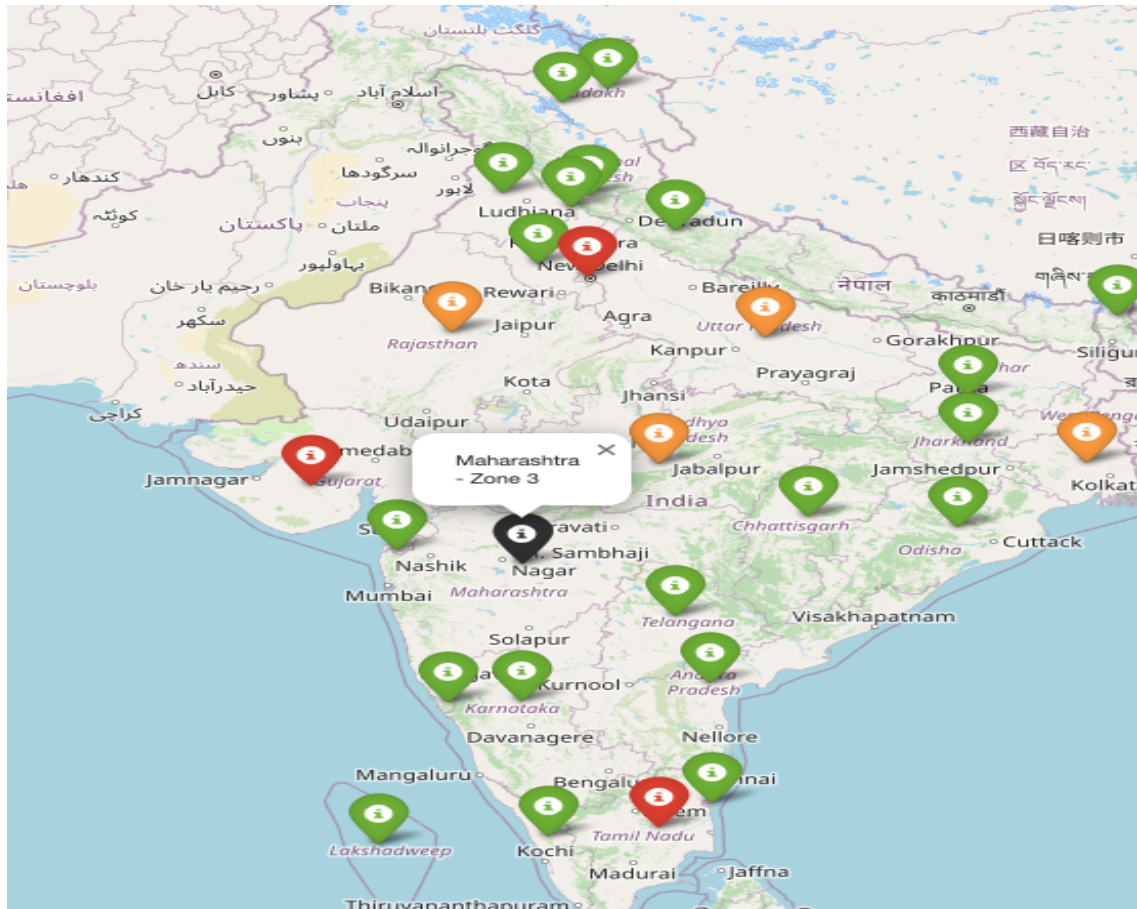
Hình 5: Kết quả phân bố sau khi thực hiện phân cụm

Bước 6: Trực quan hóa phân cụm trên bản đồ thế giới:

Để thực hiện trực quan hóa trên bản đồ thế giới, chúng ta cần cài đặt thêm thư viện Folium để giúp vẽ bản đồ thế giới. Cũng như các có dữ liệu kinh tuyến và vĩ tuyến của các bang của Ấn Độ. Trong project này, chúng ta đã được cung cấp dữ liệu các tọa độ của các bang, tuy nhiên đối với bài toán thực tế các bạn có thể lấy dữ liệu tọa độ theo dữ liệu GeoJSON Tại đây. Chúng ta sẽ thực hiện trực quan hóa toàn bộ phân cụm lên bản đồ thế giới với Folium như sau:

```
1 !pip install folium
2 import folium
3
4 # Assuming you have a dataframe called df_india with 'State' and 'Zone'
  columns
5 # and a dictionary called state_coords with state names as keys and
  latitude, longitude tuples as values.
6
7 # Example state_coords dictionary (replace with your actual data)
8 state_coords = {
9     "Andhra Pradesh": (15.9129, 79.7399),
10    "Arunachal Pradesh": (28.2180, 94.7278),
11    "Assam": (26.2006, 92.9376),
12    "Bihar": (25.0961, 85.3131),
13    "Chhattisgarh": (21.2787, 81.8661),
14    "Goa": (15.2993, 74.1240),
15    "Gujarat": (22.2587, 71.1924),
16    "Haryana": (29.0588, 76.0856),
17    "Himachal Pradesh": (31.1048, 77.1734),
18    "Jharkhand": (23.6102, 85.2799),
19    "Karnataka": (15.3173, 75.7139),
20    "Kerala": (10.8505, 76.2711),
21    "Madhya Pradesh": (22.9734, 78.6569),
22    "Maharashtra": (19.7515, 75.7139),
23    "Manipur": (24.6637, 93.9063),
24    "Meghalaya": (25.4670, 91.3662),
25    "Mizoram": (23.1645, 92.9376),
26    "Nagaland": (26.1584, 94.5624),
27    "Odisha": (20.9517, 85.0985),
28    "Punjab": (31.1471, 75.3412),
29    "Rajasthan": (27.0238, 74.2179),
30    "Sikkim": (27.5330, 88.5122),
31    "Tamil Nadu": (11.1271, 78.6569),
32    "Telangana": (18.1124, 79.0193),
33    "Tripura": (23.9408, 91.9882),
34    "Uttar Pradesh": (26.8467, 80.9462),
35    "Uttarakhand": (30.0668, 79.0193),
36    "West Bengal": (22.9868, 87.8550),
37    "Andaman and Nicobar Islands": (11.7401, 92.6586),
38    "Chandigarh": (30.7333, 76.7794),
39    "Dadra and Nagar Haveli and Daman and Diu": (20.2270, 73.0169),
40    "Delhi": (28.7041, 77.1025),
41    "Jammu and Kashmir": (33.7782, 76.5762),
42    "Ladakh": (34.1526, 77.5806),
```

```
43     "Lakshadweep": (10.5593, 72.6358),
44     "Puducherry": (11.9416, 79.8083),
45     "Orissa": (20.9517, 85.0985) # Added Orissa for consistency
46 }
47
48
49 # Create a map centered on India
50 map_india = folium.Map(location=[20.5937, 78.9629], zoom_start=4)
51
52 # Add markers for each state with color based on cluster
53 for index, row in df_india.iterrows():
54     state = row['State']
55     zone = row['Zone']
56     if state in state_coords:
57         lat, lon = state_coords[state]
58         if zone == 0:
59             color = "green"
60         elif zone == 1:
61             color = "orange"
62         elif zone == 2:
63             color = "red"
64         else:
65             color = "black"
66         folium.Marker(
67             location=[lat, lon],
68             popup=f"{state} - Zone {zone}",
69             icon=folium.Icon(color=color)
70         ).add_to(map_india)
71
72 # Display the map
73 map_india
```



Hình 6: Kết quả phân vùng các mức độ bệnh Covid ở Ấn Độ

- Hết -