

K-means Algorithm Visualization with SDL2/C++

Nguyen Dinh Hoang, 20021361
University of Engineering and Technology - VNU

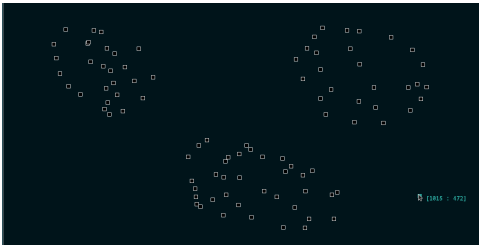
Tóm tắt nội dung— Bài toán phân cụm có nhiệm vụ phân chia một tập hợp các đối tượng (còn gọi là các thành viên) thành các nhóm khác nhau (gọi là các cụm) dựa trên các đặc điểm của đối tượng. Các thành viên của một nhóm sẽ có nhiều điểm tương đồng hơn so với những thành viên trong nhóm khác. Bài viết này nói về một phương pháp phân cụm được gọi là thuật toán K-Means dưới góc độ toán học và quá trình tạo lên nó từ thư viện SDL2 và ngôn ngữ C++.

I. GIỚI THIỆU

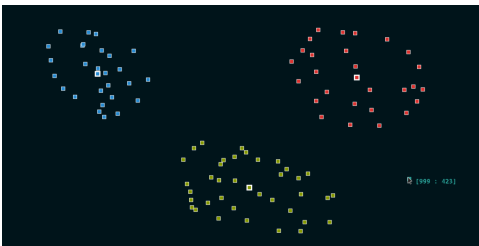
BÀI toán phân cụm có nhiều ứng dụng khác nhau: khai thác dữ liệu học máy và khám phá kiến thức, nén dữ liệu, nhận dạng mẫu và phân loại mẫu.

Mục tiêu của Thuật toán K-Means là phân tách chính xác các đối tượng trong tập dữ liệu thành các nhóm dựa trên thuộc tính của đối tượng. Ví dụ, các đối tượng có thể là ngôi nhà và thuộc tính của chúng là kích thước, số tầng, vị trí, mức tiêu thụ điện năng mỗi năm, v.v. Mục đích là phân loại tập dữ liệu nhà thành các nhóm giàu, trung bình, nghèo. Trong trường hợp đó, tất cả các thuộc tính của các ngôi nhà phải được xử lý để chuyển thành số để tạo ra một vector, quá trình này được gọi là vectơ hóa.

Một ví dụ khác, lấy các điểm trong bảng làm đối tượng và mỗi đối tượng có hai thuộc tính là vị trí trục x và trục y, với đầu vào $K = 3$.



Hình 1. Các đối tượng chưa được chia nhóm.



Hình 2. Các đối tượng sau khi được chia nhóm.

II. GIẢI THÍCH THUẬT TOÁN

Thuật toán nhận vào một tập hợp:

$$X = [x_1, x_2, \dots, x_n] \in \mathbf{R}^{d \times n}$$

d là chiều của mỗi vector, n là số vector. Cuối cùng nhận đầu vào K là số nhóm muốn chia. Thuật toán sẽ xuất ra K điểm trung tâm $[m_1, m_2, \dots, m_K] \in \mathbf{R}^{d \times K}$ của K nhóm và gắn mác(nhãn, nhóm thuộc về) cho từng điểm đầu vào.

Thuật toán có thể áp dụng cho tập hợp các dữ liệu vector nhiều chiều, vì phần mềm mô phỏng trên mặt phẳng Oxy hai chiều nên bài viết sẽ nói về hai chiều để có một cách tiếp cận dễ nhất, đối với nhiều chiều các bước làm cũng tương tự.

Tập hợp nhận vào là các vector có tọa độ (x, y) là tọa độ của chuột tại vị trí nhấp chuột.

Bước đầu chúng ta cần tạo một mảng, vector hai chiều để lưu trữ các điểm đầu vào, sau đó sẽ nhận vào vị trí của chuột và input vào mảng, vector hai chiều đó.

```
1 // Khai bao mang luu tru cac diem dau vao
2 vector<vector<double>> points;
3 ...
4 while(window.isRunning) {
5     ...
6
7     int mouse_x, mouse_y;
8     SDL_GetMouseState(&mouse_x, &mouse_y);
9
10    ...
11    if (event.type == SDL_MOUSEBUTTONDOWN) {
12        if ([mouse in points space]) {
13            labels.clear();
14            vector<double> tmp;
15            int x = mouse_x - 11; // nhan vao toa
16
17            do x
18                int y = mouse_y - 11; // nhan vao toa
19
20            do y
21                tmp.push_back(x);
22                tmp.push_back(y);
23                points.push_back(tmp);
24                tmp.clear();
25            }
26        }
27    }
```

Listing 1. Nhận đầu vào

Bước tiếp theo chúng ta sẽ khai báo và nhận đầu vào K , và tạo một mảng, vector hai chiều rồi random ra K điểm có tọa độ (x, y) bất kỳ và nhập nó vào trong mảng, vector hai chiều vừa tạo.

```
1 #include <time.h> // ham random
2 int K = 0;
3 vector<vector<double>> clusters;
4 ...
5 while(window.isRunning()) {
6     ...
7
8     if ([mouse in random space]) {
```

```

9     clusters.clear(); // make sure = []
10    labels.clear();
11    for (int i = 0; i < K; i++) {
12        // random 10 < r1 < 1178 - 10
13        // random 10 < r2 < 598 - 10
14        int r1 = 10 + rand() % (1178 - 10 + 1
15        - 10);
16        int r2 = 10 + rand() % (598 - 10 + 1
17        - 10);
18        vector<double> rand;
19        rand.push_back(r1);
20        rand.push_back(r2);
21        clusters.push_back(rand);
22        rand.clear();
23    }
24    cerr << "Random exc" << endl;

```

Listing 2. Random các điểm

Bước tiếp theo là chạy thuật toán. Nói một cách đơn giản chúng ta đang có $points = [p_1, p_2, \dots, p_n] \in \mathbf{R}^{2 \times n}$ là tập hợp các điểm ta nhập lên màn hình. $clusters = [c_1, c_2, \dots, c_K] \in \mathbf{R}^{2 \times K}$ là số điểm mà chúng ta vừa random. Bây giờ chúng ta sẽ tính khoảng cách của từng điểm trong $points$ tới các điểm trong $clusters$ rồi lưu nó vào một list mới $distancetocl$.

$$distancetocl = [d_1, d_2, \dots, d_K] \in \mathbf{R}^{2 \times K}$$

Sau khi có list $distancetocl$ ta sẽ tìm ra (giá trị nhỏ nhất) khoảng cách nhỏ nhất của điểm đầu (points[0]) đến từng điểm trong $clusters$ ghi nhớ vị trí (index) của nó rồi lưu vào list $labels$. Tiếp tục tính tương tự đối với các phần tử tiếp theo trong list $points$.

$$labels = [l_1, l_2, \dots, l_n] \in \mathbf{R}^{1 \times n}$$

Mỗi giá trị trong $labels$ sẽ đánh dấu các điểm K. Mục đích của việc này là để gắn nhãn cho các điểm.

Giờ hãy khai báo mảng và tạo hàm tính khoảng cách:

```

1 double dis(vector<double> vect1, vector<double>
2 vect2) {
3     return double(sqrt(pow((vect1[0] - vect2[0]),
4     2) + pow((vect1[1] - vect2[1]), 2)));
5 }

```

Listing 3. Hàm tính khoảng cách

```

1 if ([mouse in run space]) {
2     // change points color
3     labels.clear();
4     for (int i = 0; i < points.size(); i++) {
5         vector<double> distance_to_cl;
6         for (int j = 0; j < clusters.size(); j++)
7         {
8             double tmp1 = dis(points[i], clusters
9             [j]);
10            distance_to_cl.push_back(tmp1);
11        }
12        double min_dis = *min_element(
13        distance_to_cl.begin(), distance_to_cl.end())
14        ;
15        labels.push_back(getIndex(
16        distance_to_cl, min_dis));
17    }
18    ...
19 }

```

Listing 4. Chạy thuật toán

Tiếp theo chúng ta sẽ cập nhật tọa độ của từng điểm K trong $clusters$, cứ tiếp tục vòng lặp cho đến khi phân chia được kết quả.

```

1 ...
2 // change clusters pos
3 for (int i = 0; i < clusters.size(); i++) {
4     double sum_x = 0;
5     double sum_y = 0;
6     double count = 0;
7     for (int j = 0; j < points.size(); j++) {
8         if (labels[j] == i) {
9             sum_x = sum_x + points[j][0];
10            sum_y = sum_y + points[j][1];
11            count = count + 1;
12        }
13    }
14    if (count != 0) {
15        clusters[i][0] = sum_x/count;
16        clusters[i][1] = sum_y/count;
17    }
18 }
19 ...

```

Listing 5. Chạy thuật toán

Tọa độ mới của từng điểm K là trung bình cộng của các điểm trong $points$ được gán vào nhóm K đó. Sau khi có vị trí K mới ta sẽ cập nhật lại cho đến khi vị trí của K được cố định.