

## 2. PHẦN MÔ HÌNH TRỘN GAUSSIAN (GAUSSIAN MIXTURE MODEL)

### 2.1. Ví dụ 1. Nhắc lại ví dụ ở bài lý thuyết:

Chúng ta sẽ xét ví dụ nhân tạo với đầu vào là 3 tập ngẫu nhiên, mỗi tập N điểm được khởi tạo theo phân bố Gaussian (phân bố chuẩn) có 3 kỳ vọng ( tâm cụm) xác định trước là means[1:], means[2:] và means[3,:]; có chung ma trận hiệp phương sai là cov[.].

Đoạn mã gọi thư viện và tạo dữ liệu như dưới đây

```
# Gọi các thư viện cần thiết
# Ta tự xây dựng phần k-means nên sẽ không gọi sklearn

from __future__ import print_function
import numpy as np
import matplotlib.pyplot as plt
from scipy.spatial.distance import cdist
np.random.seed(11)

# Kỳ vọng và hiệp phương sai của 3 cụm dữ liệu
means = [[2, 2], [8, 3], [3, 6]]
cov = [[1, 0], [0, 1]]

# Số điểm mỗi cụm dữ liệu
N = 500

# Tạo các cụm dữ liệu qua phân bố chuẩn (Gaussian)
X0 = np.random.multivariate_normal(means[0], cov, N)
X1 = np.random.multivariate_normal(means[1], cov, N)
X2 = np.random.multivariate_normal(means[2], cov, N)
# Tổng hợp dữ liệu từ các cụm
X = np.concatenate((X0, X1, X2), axis = 0)

# Số cụm = 3
K = 3

# Gán nhãn ban đầu cho các cụm, sau đó ta test model và so sánh
original_label = np.asarray([0]*N + [1]*N + [2]*N).T
```

Phương thức để hiển thị dữ liệu X lên mặt phẳng, ở đây sử dụng thông tin nhãn đã được gán ở phần trước trong tham đối label

```
def kmeans_display(X, label):
    K = np.amax(label) + 1
    X0 = X[label == 0, :]
    X1 = X[label == 1, :]
    X2 = X[label == 2, :]

    plt.plot(X0[:, 0], X0[:, 1], 'b^', markersize = 4, alpha = .8)
    plt.plot(X1[:, 0], X1[:, 1], 'go', markersize = 4, alpha = .8)
    plt.plot(X2[:, 0], X2[:, 1], 'rs', markersize = 4, alpha = .8)

    plt.axis('equal')
    plt.plot()
    plt.show()
```

Khởi tạo một bộ tâm cụm trên dữ liệu X với giả thiết có k cụm

```
def kmeans_init_centers(X, k):
    # randomly pick k rows of X as initial centers
    return X[np.random.choice(X.shape[0], k, replace=False)]
```

Ở đây chúng ta chọn các tâm cụm một cách ngẫu nhiên (miễn là các tâm cụm khác nhau). Các bạn hãy thử điều chỉnh bằng cách chọn K điểm xa nhau nhất theo phương pháp đã được trình bày trong phần lý thuyết.

Phương thức để gán cụm cho một điểm dữ liệu bằng cách tính khoảng cách từ điểm đó đến các tâm cụm, khoảng cách đến đâu ngắn nhất thì ta coi điểm hiện tại sẽ thuộc về cụm đó.

```
def kmeans_assign_labels(X, centers):  
    # calculate pairwise distances btw data and centers  
    D = cdist(X, centers)  
    # return index of the closest center  
    return np.argmin(D, axis = 1)
```

Phương thức để cập nhật lại tâm cụm sau mỗi bước lặp: Tâm cụm mới sẽ là trung bình cộng (theo tọa độ) của tất cả các điểm có trong cụm)

```
def kmeans_update_centers(X, labels, K):  
    centers = np.zeros((K, X.shape[1]))  
    for k in range(K):  
        # collect all points assigned to the k-th cluster  
        Xk = X[labels == k, :]  
        # take average  
        centers[k, :] = np.mean(Xk, axis = 0)  
    return centers
```

Hàm để kiểm tra xem thuật toán có thực sự chạy (có hội tụ hay không) thông qua việc tâm cụm sau mỗi bước lặp có thay đổi hay không? Nếu tâm cụm không đổi nghĩa là thuật toán đã dừng (hội tụ) – tức là cần trả về TRUE. Kiểm tra cho tất cả các tâm cụm

```
def has_converged(centers, new_centers):  
    # return True if two sets of centers are the same  
    return (set([tuple(a) for a in centers]) ==  
            set([tuple(a) for a in new_centers]))
```

Vòng lặp để thực hiện tất cả các bước trong thuật toán k-means cho đến khi thuật toán dừng

```
def kmeans(X, K):  
    centers = [kmeans_init_centers(X, K)]  
    labels = []  
    it = 0  
    while True:  
        labels.append(kmeans_assign_labels(X, centers[-1]))  
        new_centers = kmeans_update_centers(X, labels[-1], K)  
        if has_converged(centers[-1], new_centers):  
            break  
        centers.append(new_centers)  
        it += 1  
    return (centers, labels, it)
```

Gọi và thực hiện phương pháp

```
(centers, labels, it) = kmeans(X, K)  
print('Centers found by our algorithm:')  
print(centers[-1])  
  
kmeans_display(X, labels[-1])
```

## 2.2. Ví dụ 2: Thực hiện phân cụm cho bộ dữ liệu chữ số viết tay

Hãy tìm lại phần thực hành đọc dữ liệu từ tệp chứa thông tin hình ảnh các chữ số viết tay và thực hiện việc phân cụm dữ liệu vào 10 cụm (10 chữ số viết tay) theo cách sau: Đọc 5000 mẫu dữ liệu từ phần training, sau

đó thực hiện phân cụm bằng phương pháp k-means, tiếp theo hãy kiểm tra xem trong mỗi cụm, tỷ lệ có nhãn nào (từ 0 đến 9) là cao nhất. Sau đó đếm và in ra tỷ lệ các mẫu không thuộc nhãn đó nhưng được phân vào cùng một cụm với nhãn.