

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

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

Trong ví dụ này, chúng ta sử dụng dữ liệu là thông tin về chỉ số mua sắm và mức thu nhập có trong tệp dữ liệu shopping-data.csv. Tệp dữ liệu có 4 trường nhưng ta chỉ sử dụng 2 trường liên quan tới thu nhập và chỉ số mua sắm để tìm tương quan, các trường khác chỉ chứa thông tin ID nên ta bỏ qua.

Trong phần này ta sẽ sử dụng thư viện scikit-learn (sklearn) để thực hiện mô hình.

Các bước sẽ như sau:

- (i) Đọc dữ liệu vào, chuẩn hóa (chú ý trong phần lý thuyết ta đã biết nếu kỳ vọng của các phân bố Gaussian là 0 thì tính toán sẽ tốt hơn); Ở trong code ta sẽ in ra một số thông tin của dữ liệu để xem đọc có đúng không.
- (ii) Khởi tạo đối tượng của lớp mô hình trộn Gaussian thông qua hàm dựng  
`gm = GaussianMixture(n_components=K, covariance_type='full', random_state=0)`  
Ở đây số cụm được chia sẽ là `n_component = K`
- (iii) Sau đó ta khớp dữ liệu của ta bằng đối tượng mô hình vừa tạo:  
`gm.fit(X_DATA)`
- (iv) Chúng ta vẽ dữ liệu ra dạng trực quan để có thể quan sát.
- (v) Ở trong code dưới đây, chúng ta bổ sung thêm phần tìm xem K (số cụm nên dùng để chia dữ liệu) bao nhiêu là hợp lý nhất. Các làm của chúng ta ở đây đơn giản chỉ là thử với một dãy các giá trị số cụm, giá trị nào cho kết quả tốt nhất (trong số đó) thì sử dụng.

Dưới đây là các đoạn code:

Bước khai báo thư viện:

```
import numpy as np
import pandas as pd
import seaborn as sns
import itertools
from scipy import linalg
import matplotlib.pyplot as plt
import matplotlib.path_effects as PathEffects
from matplotlib.patches import Ellipse
from sklearn.preprocessing import MinMaxScaler
from sklearn.mixture import GaussianMixture
# Thư viện chứa model Gaussian Mixture
```

Bước (i)

```
data = pd.read_csv("D:\\Teach_n_Train\\Machine
                  Learning\\slide\\Mixture_Model\\code
                  \\shopping-data.csv",
                  header=0,
                  index_col=0)
print(data.shape)
data.head()

# Lấy ra thu nhập và điểm shopping
X = data.iloc[:, 2:4].values
```

```
# Chuẩn hoá dữ liệu
std = MinMaxScaler()
X_std = std.fit_transform(X)
print(X_std.shape)
```

Bước (ii) và bước (iii)

```
# Khởi tạo đối tượng mô hình GaussianMixture
gm = GaussianMixture(n_components=5,
                     covariance_type='full',
                     random_state=0)

gm.fit(X_std)
print('means: \n', gm.means_)
print('covariances: \n ', gm.covariances_)
```

Bước (v): Ta nêu phân chọn số K tốt nhất trước:

```
lowest_bic = np.infty
bic = []
n_components_range = range(1, 7)
# cv_types = ['spherical', 'tied', 'diag', 'full']
cv_types = ['full', 'tied']
for cv_type in cv_types:
    for n_components in n_components_range:
        # Fit Gaussian mixture theo phương pháp huấn luyện EM
        gmm = GaussianMixture(n_components=n_components,
                              covariance_type=cv_type)

        gmm.fit(X_std)
        bic.append(gmm.bic(X_std))
        # Gán model có BIC scores thấp nhất là model tốt nhất
        if bic[-1] < lowest_bic:
            lowest_bic = bic[-1]
            best_gmm = gmm

bic = np.array(bic)
color_iter = itertools.cycle(['navy', 'turquoise'])
clf = best_gmm
bars = []

# Vẽ biểu đồ BIC scores
plt.figure(figsize=(12, 8))
for i, (cv_type, color) in enumerate(zip(cv_types, color_iter)):
    xpos = np.array(n_components_range) + .2 * (i - 2)
    bars.append(plt.bar(xpos, bic[i * len(n_components_range):
                               (i + 1) * len(n_components_range)],
                        width=.2, color=color))

plt.xticks(n_components_range)
plt.ylim([bic.min() * 1.01 - .01 * bic.max(), bic.max()])
plt.title('BIC score per model')
xpos = np.mod(bic.argmin(), len(n_components_range)) + .65 + \
    .2 * np.floor(bic.argmin() / len(n_components_range))
plt.text(xpos, bic.min() * 0.97 + .03 * bic.max(), '*', fontsize=14)
plt.xlabel('Number of components')
plt.legend([b[0] for b in bars], cv_types)
```

Bước (iv) – Hiển thị kết quả

```
def _plot_kmean_scatter(X, labels):
    """
    X: dữ liệu đầu vào
    labels: nhãn dự báo
    """
    # lựa chọn màu sắc
    num_classes = len(np.unique(labels))
```

```

palette = np.array(sns.color_palette("hls", num_classes))

# vẽ biểu đồ scatter
fig = plt.figure(figsize=(12, 8))
ax = plt.subplot()
sc = ax.scatter(X[:,0], X[:,1], lw=0, s=40,
                c=palette[labels.astype(np.int)])

# thêm nhãn cho mỗi cluster
txts = []

for i in range(num_classes):
    # Vẽ text tên cụm tại trung vị của mỗi cụm
    xtext, ytext = np.median(X[labels == i, :], axis=0)
    txt = ax.text(xtext, ytext, str(i), fontsize=24)
    txt.set_path_effects([
        PathEffects.Stroke(linewidth=5, foreground="w"),
        PathEffects.Normal()])
    txts.append(txt)
plt.title('t-sne visualization')

```

Cuối cùng: Gọi và thực hiện 2 phương thức chính: Thực hiện mô hình trộn Gaussian và phương thức hiển thị dữ liệu ra mặt phẳng 2 chiều để dễ quan sát.

```

labels = best_gmm.predict(X_std)
plot_kmean_scatter(X_std, labels)

```

Chú ý với các ví dụ tiếp theo, chúng ta cần đọc dữ liệu đầy đủ và thiết lập vector dữ liệu X phù hợp (vì số chiều có thể lớn hơn). Sau đó trong phần hiển thị dữ liệu, chúng ta cần sử dụng mô hình phân tích thành phần chính (PCA) để giảm số chiều dữ liệu xuống còn 2 chiều để có thể hiển thị lên mặt phẳng.

**1.2.Ví dụ mở rộng.** Trong ví dụ này chúng ta sử dụng dữ liệu cho trong tệp đính kèm [Sales Transactions Dataset Weekly.csv](#) hoặc tại link <https://archive.ics.uci.edu/ml/machine-learning-databases/00396/>.

**Mô tả dữ liệu:** Dữ liệu chứa thông tin giao dịch bán hàng của 819 mã sản phẩm (Product ID) trong 819 dòng dữ liệu. Mỗi mã được đặc trưng bởi các thông tin: Mã (cột đầu); lượng giao dịch trong 52 tuần (week 0 – week 51), sau đó có thông tin chuẩn hóa của giao dịch các tuần (ta sử dụng thông tin này hoặc thông tin theo tuần – chỉ chọn 1 trong 2).

**Yêu cầu:** Hãy dựa vào ví dụ 1, mở rộng áp dụng với dữ liệu trên đây. Biết rằng số cụm tối ưu sẽ nằm trong khoảng từ 5 đến 15, hãy xác định số cụm và sau đó phân các mã sản phẩm thành các cụm dựa trên tiêu chí lượng giao dịch đã cho trong dữ liệu.