

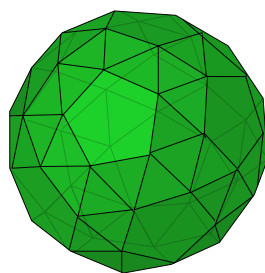
Projects in Mathematics and Applications

PRINCIPAL COMPONENT ANALYSIS

Ngày 27 tháng 8 năm 2018

Trần Thanh Bình *
Lê Quang Kỳ ‡

†Võ Thục Khánh Huyền
§Đỗ Nhật Hoàng



*Trường Phổ Thông Năng Khiếu

†Trường THPT Chuyên Lê Quý Đôn, Quảng Trị

‡Trường THPT Chuyên Hoàng Lê Kha, Tây Ninh

§Trường Phổ Thông Năng Khiếu

Lời cảm ơn

Thời gian được học tập ở Trại hè PiMA 2018 (Project in Mathematics and Applications) là quãng thời gian vô cùng ý nghĩa. Kết quả ngày hôm nay chúng tôi đạt được một phần là nhờ có sự quan tâm hỗ trợ của nhiều cá nhân và các đoàn thể. Trước hết chúng tôi xin gửi lời cảm ơn đến các anh chị Mentor, đặc biệt là anh Cần Trần Thành Trung và Ban Tổ Chức PiMA đã giảng dạy cho chúng tôi về những kiến thức Toán thú vị, chia sẻ cho chúng tôi những cơ hội để tiếp cận với Toán ứng dụng, và giúp chúng tôi hoàn thành được đề tài một cách tốt nhất. Bên cạnh đó, chúng tôi cũng muốn bày tỏ lòng biết ơn đến Trường Đại học Khoa học Tự nhiên Thành phố Hồ Chí Minh đã tạo điều kiện tốt nhất về vật chất và tinh thần để chúng tôi có thể tập trung tối đa vào làm dự án nhóm. Và cuối cùng, chúng tôi muốn cảm ơn các bạn trại sinh đã cùng nhau học tập và vui chơi với chúng tôi trong suốt thời gian qua. Do thời gian làm dự án còn ngắn nên thiếu sót là điều không thể tránh khỏi. Nhóm chúng tôi mong sẽ nhận được đóng góp từ phía độc giả để dự án có thể hoàn thiện tốt hơn.

Tóm tắt nội dung

Những dữ liệu trong thực tế thường có số chiều và số lượng rất lớn, gây khó khăn cho việc lưu trữ và xử lý. Để giải quyết vấn đề này, người ta xây dựng các phương pháp làm giảm chiều dữ liệu nhưng vẫn giữ được những thông tin chính. Phương pháp Principal Component Analysis (PCA), phân tích thành phần chính, là một trong những phương pháp phổ biến nhất. Trong bài viết này, chúng tôi sẽ trình bày và chứng minh lại cơ sở toán học của phương pháp PCA. Do việc tìm ra các thành phần chính của bộ dữ liệu trong PCA tương đương với việc tìm trị riêng và vector riêng của những ma trận có kích thước lớn, bài viết sẽ giới thiệu thêm phương pháp Power để xấp xỉ trị riêng và vector riêng. Sau đó, chúng tôi áp dụng PCA để biểu diễn dữ liệu nhiều chiều và tìm những nét chính của bộ hình ảnh khuôn mặt người.

Mục lục

1	Giới thiệu bài toán Giảm chiều dữ liệu (Dimensionality Reduction)	1
1.1	Bài toán giảm chiều dữ liệu	1
1.2	Các hướng tiếp cận bài toán Giảm chiều dữ liệu	1
2	Cơ sở Toán học	1
2.1	Một số kiến thức quan trọng	1
2.2	Một số kết quả cơ bản	3
3	Mô hình hóa bài toán PCA	4
4	Giải quyết bài toán PCA	5
4.1	Phương pháp nhân tử Lagrange	6
4.2	Phương pháp chuyển hệ cơ sở	7
4.3	Kết luận	8
5	Giới thiệu thuật toán tìm trị riêng	8
5.1	Phương pháp Power	8
5.2	Nhược điểm của thuật toán Power	10
5.3	Tổng kết	10
6	Cài đặt	10
6.1	Các bước thực hiện Phương pháp PCA	10
6.2	Source code trên ngôn ngữ lập trình Python	11
7	Áp dụng mô hình	12
7.1	Hoa Iris	12
7.2	Cơ sở dữ liệu dinh dưỡng quốc gia USDA	12
7.3	EigenFace	13
8	Kết luận đánh giá	16
9	Hướng nghiên cứu trong tương lai	16

1 Giới thiệu bài toán Giảm chiều dữ liệu (Dimensionality Reduction)

1.1 Bài toán giảm chiều dữ liệu

Thông tin trong khoa học dữ liệu nói chung và Machine Learning nói riêng thường được mô tả bằng những vector có số chiều và số phần tử rất lớn, gây khó khăn trong việc lưu trữ và xử lý dữ liệu. Nếu không gian dữ liệu gốc có m tính trạng thì các vector dữ liệu sẽ là những điểm $x \in \mathbb{R}^m$. Giảm chiều dữ liệu (Dimensionality Reduction) thực chất là tìm một hàm số $f: \mathbb{R}^m \rightarrow \mathbb{R}^k$ phù hợp, với $k < m$, để chuyển sang làm việc với các điểm dữ liệu mới $z = f(x)$ có số chiều nhỏ hơn. Những hàm f phổ biến nhất được chọn thường là phép chiếu lên một không gian con thỏa mãn những tính chất nhất định.

1.2 Các hướng tiếp cận bài toán Giảm chiều dữ liệu

Hai hướng tiếp cận chủ yếu trong bài toán Giảm chiều dữ liệu là *Feature Selection* (lựa chọn những tính trạng liên quan) và *Feature Extraction* (xây dựng những tính trạng mới chứa nhiều thông tin hơn). Trong *Feature Extraction*, một số phương pháp thường gặp bao gồm phương pháp *Linear Discriminant Analysis (LDA)*, *Canonical Correlation Analysis (CCA)*, hay *Autoencoder*.

Phép *Phân tích thành phần chính (PCA)* có thể dùng trong cả Selection lẫn Extraction. *Phương pháp PCA* sẽ biểu diễn các vector dữ liệu trong hệ cơ sở trực chuẩn mới (tương ứng với cơ sở riêng), sao cho tầm quan trọng của mỗi chiều dữ liệu là khác nhau. Khi đó, chúng ta có thể chọn và giữ lại k *chi u quan trọng* nhất, tương ứng với việc chọn f là phép chiếu lên k chiều tương ứng. Cơ sở toán của PCA sẽ được trình bày cụ thể ở các phần sau. PCA có rất nhiều ứng dụng như nén ảnh (Image Compression) và nhận dạng khuôn mặt (Facial Recognition).

2 Cơ sở Toán học

2.1 Một số kiến thức quan trọng

Định nghĩa 2.1. (Biểu diễn trong cơ sở). Cho không gian vector U và một cơ sở hữu hạn $S = \{u_i : i = 1; 2; \dots; n\}$ của U . Với mỗi vector $u \in U$, ký hiệu $[u]_S$ là biểu diễn của u qua cơ sở S :

$$[u]_S = [x_1 \ x_2 \ \dots \ x_n];$$

nếu u là một vector trong \mathbb{R}^n thỏa mãn:

$$u = x_1 u_1 + x_2 u_2 + \dots + x_n u_n. \quad (1)$$

Định nghĩa 2.2. Cơ sở S được gọi là **trực giao** nếu $u_i \cdot u_j = 0$, $\forall i \neq j \in \{1; 2; \dots; n\}$.

Định nghĩa 2.3. Cơ sở S được gọi là **chuẩn** nếu $u_i \cdot u_i = 1$, $\forall i \in \{1; 2; \dots; n\}$.

Định nghĩa 2.4. Cơ sở S được gọi là **trực chuẩn** nếu nó chuẩn và trực giao.

Nhận xét 2.5. Mọi vector $u \in U$ trong cơ sở trực chuẩn S :

$$\begin{aligned} u &= \text{proj}_{u_1}(u) + \text{proj}_{u_2}(u) + \dots + \text{proj}_{u_n}(u) \\ &= (u \cdot u_1)u_1 + (u \cdot u_2)u_2 + \dots + (u \cdot u_n)u_n. \end{aligned}$$

Chứng minh. Theo định nghĩa trên, ta có:

$$u = x_1 u_1 + x_2 u_2 + \dots + x_n u_n \quad (1)$$

Nhân u_i vào 2 vế của phương trình (1) ta được: $u \cdot u_i = x_i$ với $\forall i \in \{1; 2; \dots; n\}$:
Do đó $u = (u \cdot u_1)u_1 + (u \cdot u_2)u_2 + \dots + (u \cdot u_n)u_n$:

Mặt khác, ta có công thức hình chiếu như sau:

$$\text{proj}_{u_i}(u) = \frac{u \cdot u_i}{\|u_i\|^2} u_i = (u \cdot u_i)u_i \text{ với } \forall i \in \{1; 2; \dots; n\}:$$

Do đó $u = \text{proj}_{u_1}(u) + \text{proj}_{u_2}(u) + \dots + \text{proj}_{u_n}(u)$: □

Nhận xét 2.6. Từ các nhận xét trên, ta thấy mọi vector $u \in U$ trong cơ sở trực chuẩn S thì ta có thể bình phương bình phương (1) và rút gọn trở thành:

$$\|u\|^2 = x_1^2 + x_2^2 + \dots + x_n^2:$$

Do đó, ngược lại ta thấy mọi vector d liên thuộc U thì tọa độ của nó trong cơ sở trực chuẩn S thu được từ công thức tính toán và bình phương.

Định nghĩa 2.7. (Trị riêng và vector riêng của ma trận). Cho $A \in \mathbb{R}^{n \times n}$: Ta nói $\lambda \in \mathbb{R}$ là một **trị riêng** (hoặc giá trị riêng) của A nếu tồn tại $v \in \mathbb{R}^n \setminus \{0\}$ sao cho:

$$Av = \lambda v:$$

Lúc đó ta nói v là **vector riêng** ứng với trị riêng λ .

Đa thức đặc trưng của A được kí hiệu $P_A: \mathbb{R} \rightarrow \mathbb{R}$ sao cho $P_A(x) = \det(A - xI_n)$:

Định lý 2.8. Mọi ma trận $n \times n$ có n trị riêng thực.

Định lý 2.9. Trị riêng của ma trận $A \in \mathbb{R}^{n \times n}$ là nghiệm của phương trình $P_A(\lambda) = 0$:

Phương Pháp Nhân Tử Lagrange (Lagrange Multiplier Method)

Ta xét bài toán tối ưu: $\min_{x \in \mathbb{R}^n} f(x)$, thỏa mãn:

$$g_i(x) = 0; i \in \{1; 2; \dots; m\};$$

với các hàm $f; g_i: \mathbb{R}^n \rightarrow \mathbb{R}$ là các hàm khả vi với đạo hàm riêng liên tục.

Để giải quyết bài toán này, ta sẽ phát biểu một điều kiện cần để điểm x^* là một cực trị địa phương của hàm f .

Ta định nghĩa **Lagrangian** $\mathcal{L}: \mathbb{R}^n \times \mathbb{R}^m \rightarrow \mathbb{R}$ của bài toán tối ưu:

$$\mathcal{L}(x; \lambda) = f(x) + \sum_{i=1}^m \lambda_i g_i(x)$$

với $\lambda = [\lambda_1; \lambda_2; \dots; \lambda_m]$ được gọi là **biến đối ngẫu** hay **vector nhân tử Lagrange**.

Nếu x^* là một điểm tối thiểu địa phương của bài toán tối ưu trên, ta có:

$$\exists \lambda^*: \nabla_x \mathcal{L}(x^*; \lambda^*) = 0:$$

Để tìm cực tiểu toàn cục, ta có thể giải hệ phương trình có được từ điều kiện ở trên rồi xét thử tất cả các bộ nghiệm tìm được.

2.2 Một số kết quả cơ bản

Cho $A \in \mathbb{R}^{m \times n}$. Khi đó $AA^T \in \mathbb{R}^{m \times m}$ và $A^T A \in \mathbb{R}^{n \times n}$. Đặt $B = A^T A$.

Tính chất 2.10. Ma trận B đối xứng.

Chứng minh. Vì $B^T = (A^T A)^T = A^T A = B$ nên $B^T = B$. □

Tính chất 2.11. Chúng ta biết rằng ma trận thực đối xứng có trị riêng là các số thực. Giả sử λ_1, λ_2 là 2 trị riêng khác nhau của B với các vector riêng tương ứng v_1, v_2 . Chứng minh rằng $v_1 \cdot v_2 = 0$. Suy ra các vector riêng ứng với trị riêng khác nhau của B trực giao và độc lập tuyến tính.

Chứng minh. Với $B \in \mathbb{R}^{n \times n}$ và $x, y \in \mathbb{R}^n$, ta có:

1. $x \cdot y = x^T y$.
2. $Bx \cdot y = (Bx)^T \cdot y = (x^T B^T) \cdot y = x^T \cdot (B^T y) = x \cdot B^T y$.

Áp dụng với λ_1, λ_2 là 2 trị riêng khác nhau của B với các vector riêng tương ứng v_1, v_2 , ta có:

$$Bv_1 = \lambda_1 v_1 \quad (1) \quad \text{và} \quad Bv_2 = \lambda_2 v_2 \quad (2).$$

Nhân v_2 vào 2 vế của (1) ta có: $v_2 \cdot (Bv_1) = \lambda_1 v_2 \cdot v_1$.

Áp dụng công thức 2 vào biểu thức trên ta có: $(B^T v_2) \cdot v_1 = \lambda_1 v_2 \cdot v_1$.

Mặt khác $B = B^T$ và kết hợp (2) nên biểu thức trên được viết lại thành: $v_2 \cdot v_1 = \lambda_1 v_2 \cdot v_1$.

Hay $(\lambda_1 - \lambda_2) v_2 \cdot v_1 = 0$, mà λ_1 khác λ_2 nên $v_2 \cdot v_1 = 0$.

Do đó, ta suy ra được các vector riêng ứng với trị riêng khác nhau của B **trực giao**. □

Ma trận đối xứng $B \in \mathbb{R}^{n \times n}$ thì có đúng n trị riêng thực, gọi $v_1; v_2; \dots; v_n$ là n vector riêng ứng với n trị riêng đó. Giả sử tồn tại các số $c_1; c_2; \dots; c_n \in \mathbb{R}$ thỏa mãn:

$$c_1 v_1 + c_2 v_2 + \dots + c_n v_n = \theta:$$

Bình phương hai vế của biểu thức trên và sử dụng điều kiện $v_i \cdot v_j = 0$ với $i \neq j; i, j \in \{1; 2; \dots; n\}$, ta có:

$$\sum_{i=1}^n c_i^2 \|v_i\|^2 = 0;$$

nên $c_i = 0$ với $\forall i \in \{1; 2; \dots; n\}$.

Do đó, ta suy ra được các vector riêng ứng với trị riêng khác nhau của B **độc lập tuyến tính**, hay chúng tạo thành một hệ cơ sở.

Ý nghĩa: Ta sẽ biểu diễn bộ dữ liệu dưới cơ sở S trực giao mới mà $S = \{v_1; v_2; \dots; v_n\}$. Cơ sở trực chuẩn giúp việc tính toán biểu diễn của một vector và tích vô hướng giữa các vector dễ dàng hơn.

Tính chất 2.12. Chứng minh rằng $A^T A$ và AA^T có cùng giá trị riêng (Lưu ý: có thể khác nhau về số lần lặp của trị riêng 0)

Chứng minh. Gọi v là vector riêng của $A^T A$ ứng với trị riêng λ ($v \neq \theta, \lambda \neq 0$). Ta có:

$$(A^T A)v = \lambda v \quad \text{hay} \quad A^T (Av) = \lambda v:$$

Nhân hai vế của biểu thức cho A và gộp nhóm lại như sau:

$$AA^T(Av) = (Av):$$

Điều này cho chúng ta thấy Av là vector riêng của AA^T ứng với trị riêng $\neq 0$. Bây giờ chúng ta cần đi chứng minh Av khác vector 0. Nếu $Av = 0$ thì theo biểu thức trên $v = 0$, vô lý vì $v \neq 0, \neq 0$. Do đó những trị riêng $\neq 0$ của $A^T A$ cũng là những trị riêng của AA^T . \square

Ý nghĩa: Kết quả này thực sự hiệu quả, nếu ta xét ma trận $A \in \mathbb{R}^{100 \times 3}$. Ta sẽ mất rất nhiều thời gian để đi tìm trị riêng của ma trận $AA^T \in \mathbb{R}^{100 \times 100}$, thay vào đó ta sẽ đi tìm trị riêng của ma trận $A^T A \in \mathbb{R}^{3 \times 3}$. Đó cũng chính là trị riêng của ma trận AA^T , 97 trị riêng còn lại bằng 0.

3 Mô hình hóa bài toán PCA

Ta sẽ xét một bộ dữ liệu gồm n điểm dữ liệu, mỗi điểm dữ liệu có m tính trạng. Đối với bài toán PCA, ta sẽ biểu diễn mỗi điểm dữ liệu thành một vector m chiều, mỗi chiều sẽ ứng với một tính trạng. Ở đây, chúng ta thống nhất nếu không nói gì thêm thì X là vector cột:

$$X_i = \begin{bmatrix} X_{1,i} \\ X_{2,i} \\ \vdots \\ X_{m,i} \end{bmatrix}; \text{ với } i = 1; 2; \dots; n:$$

Ta xét ma trận sau:

$$X_{m \times n} = [X_1 \quad X_2 \quad \dots \quad X_n]:$$

Viết đầy đủ, ta có:

$$X_{m \times n} = \begin{bmatrix} X_{1,1} & X_{1,2} & \dots & X_{1,n} \\ X_{2,1} & X_{2,2} & \dots & X_{2,n} \\ \vdots & \vdots & \ddots & \vdots \\ X_{m,1} & X_{m,2} & \dots & X_{m,n} \end{bmatrix}:$$

Tiếp theo, chúng ta sẽ định nghĩa thế nào là một *chi u quan tr ng* của dữ liệu. *Chi u quan tr ng* của dữ liệu có thể hiểu là chiều của một vector w trong không gian \mathbb{R}^m mà theo phương là giá của vector đó, *phân tán* của các tính trạng là lớn nhất. Và để xác định *phân tán*, ta sẽ dùng đến phương sai. Như vậy, chiều quan trọng nhất sẽ được định nghĩa là chiều mà trên đó phương sai của bộ dữ liệu là lớn nhất.

Ta có công thức phương sai:

$$\text{Var}_{w_i}(X) = \frac{2}{w_i}(X) = \frac{1}{n} \sum_{j=1}^n (X_{i,j} - \bar{x}_i)^2; \text{ với } i = 1; 2; \dots; n.$$

Trong đó, $\bar{x}_i = \frac{X_{i,1} + X_{i,2} + \dots + X_{i,n}}{n}$ là giá trị trung bình hay giá trị kì vọng của mỗi tính trạng.

Nhận thấy rằng việc chuẩn hóa ma trận này về quanh gốc tọa độ không làm thay đổi bản chất về *phân tán* của bộ dữ liệu, nhưng lại làm cho việc tính toán trở nên dễ dàng hơn. Vì vậy, ta có thể xét việc biến đổi trên ma trận $\hat{X}_{m \times n}$ như sau:

$$\hat{X}_{m \times n} = [x_{i,j} - \bar{x}_i] = [x_1 \quad x_2 \quad \dots \quad x_n] :$$

Với bộ dữ liệu $\hat{X}_{m \times n}$, ta có $\hat{x}_{m \times n} = \theta$, ta có định nghĩa sau:

$$\text{Var}_{w_i}(X) = \frac{1}{n} \sum_{i=1}^n \|\text{proj}_{w_i} x_i\|^2.$$

Để tìm *chi u quan trọng* nhất, ta sẽ tìm w_i sao cho $\sum_{i=1}^n \|\text{proj}_{w_i} x_i\|^2$ đạt giá trị lớn nhất.

Nhớ lại công thức tính $\|\text{proj}_{w_i} x_i\|$:

$$\|\text{proj}_{w_i} x_i\| = \left\| \frac{x_i \cdot w_i}{\|w_i\|^2} w_i \right\|.$$

Như vậy, nếu ta chuẩn hóa w_i sao cho $\|w_i\| = 1$ thì công thức tính $\text{proj}_{w_i} x_i$ có thể được viết lại:

$$\|\text{proj}_{w_i} x_i\| = |x_i \cdot w_i|.$$

Khi đó, công thức tính w_1 có thể được viết lại thành:

$$w_1 = \underset{\|w_1\|=1}{\text{argmax}} \sum_{i=1}^n (w_1 \cdot x_i)^2 :$$

Hay có thể được viết gọn lại thành:

$$w_1 = \underset{\|w_1\|=1}{\text{argmax}} [(\hat{X}_{m \times n}^T w_1) \cdot (\hat{X}_{m \times n}^T w_1)] :$$

Áp dụng công thức biến đổi $x \cdot y = x^T y$, ta có:

$$w_1 = \underset{\|w_1\|=1}{\text{argmax}} (w_1^T \hat{X} \hat{X}^T w_1) :$$

Vậy, việc xác định thành phần quan trọng nhất của bộ dữ liệu đã được mô hình hóa thành việc giải quyết bài toán tối ưu trên.

Sau khi đã có thành phần quan trọng nhất, ta sẽ tiếp tục đi tìm các thành phần còn lại. Để thuận tiện cho việc quan sát và tính toán, ta sẽ tìm thành phần quan trọng thứ i đôi một vuông góc với những thành phần trước nó. Hay nói khác đi, ta sẽ tìm w_i sao cho $\text{Var}_{w_i}(X)$ đạt max và w_i thỏa $w_i \cdot w_j = 0; \forall j < i$. Biến đổi tương tự như đối với chiều thứ nhất, ta sẽ có:

$$w_i = \underset{\substack{\|w_i\|=1 \\ w_i \cdot w_j = 0; \forall j < i}}{\text{argmax}} (w_i^T \hat{X} \hat{X}^T w_i) :$$

4 Giải quyết bài toán PCA

Bài toán PCA đã được mô hình hóa về việc tìm những chiều dữ liệu. Bây giờ ta sẽ giải bài toán tối ưu đã được nhắc đến ở phần trên. Vì đây là một bài toán tối ưu hóa có điều kiện, phương pháp đầu tiên chúng ta nghĩ đến là *phương pháp nhân t Lagrange*.

4.1 Phương pháp nhân tử Lagrange

Ta có hàm cần tối ưu:

$$w_1 = \operatorname{argmax}_{\|w_1\|=1} (w_1^T \hat{X} \hat{X}^T w_1);$$

và có Lagrangian:

$$\mathcal{L}(w_1; \lambda) = w_1^T \hat{X} \hat{X}^T w_1 + \lambda (\|w_1\|^2 - 1);$$

Đặt $A = \hat{X} \hat{X}^T$, ta có $A \in \mathbb{R}^{m \times m}$ và $A = A^T$: Viết lại công thức:

$$\mathcal{L}(w_1; \lambda) = w_1^T A w_1 + \lambda (\|w_1\|^2 - 1);$$

Đặt $w_1 = \begin{bmatrix} w_{1,1} \\ w_{1,2} \\ \vdots \\ w_{1,m} \end{bmatrix}$, khi đó: $w_1^T A w_1 = \sum_{i=1}^m \sum_{j=1}^m w_{1,i} w_{1,j} a_{ij}$:

Vậy ta có:

$$\mathcal{L}(w_1; \lambda) = \sum_{i=1}^m \sum_{j=1}^m w_{1,i} w_{1,j} a_{ij} + \lambda (\|w_1\|^2 - 1);$$

Điều kiện cần để \mathcal{L} đạt giá trị lớn nhất: $\nabla \mathcal{L}(w_1; \lambda) = \begin{bmatrix} \frac{\partial \mathcal{L}}{\partial \lambda} \\ \frac{\partial \mathcal{L}}{\partial w_{1,1}} \\ \frac{\partial \mathcal{L}}{\partial w_{1,2}} \\ \vdots \\ \frac{\partial \mathcal{L}}{\partial w_{1,m}} \end{bmatrix} = 0$:

Như vậy, nếu \mathcal{L} đạt \max thì w_1 thỏa:

$$\begin{cases} \frac{\partial \mathcal{L}}{\partial w_{1,i}} = 0, \forall i = 1; 2; \dots; m \\ \frac{\partial \mathcal{L}}{\partial \lambda} = 0 \end{cases}$$

Ta có:

$$\frac{\partial \mathcal{L}}{\partial w_{1,i}} = 2 \cdot \sum_{j=1}^m a_{ij} w_{1,j} + 2 \lambda w_{1,i};$$

Như vậy, hệ có thể viết lại thành:

$$\begin{aligned} (A + \lambda I_n) w_1 &= 0 \\ \Leftrightarrow A w_1 &= -\lambda w_1 \end{aligned}$$

Do đó, w_1 là một vector riêng của A và $-\lambda$ là một trị riêng của A .

Thế điều kiện cần vào hàm nhân tử Lagrange:

$$\begin{aligned}\mathcal{L}(w_1; \lambda) &= w_1^T A w_1 \\ &= w_1^T (-\lambda w_1) \\ &= -\lambda \|w_1\|^2 \\ &= -\lambda.\end{aligned}$$

Như vậy, $\max(w_1^T X X^T w_1) = -\lambda_{\max}$ khi và chỉ khi w_1 là vector riêng có độ dài bằng 1 ứng với $-\lambda_{\max}$ (trị riêng lớn nhất của ma trận $X X^T$).

Giải các bài toán tối ưu kia tương tự bằng *phương pháp nhân tử Lagrange*, ta sẽ tìm được kết quả là hệ cơ sở gồm các vector riêng có độ dài 1.

Sau khi đã giải ra hệ nghiệm vector bằng *phương pháp nhân tử Lagrange*, ta nghĩ đến một phương pháp thứ hai đẹp hơn như sau.

4.2 Phương pháp chuyển hệ cơ sở

Với chiều đầu tiên, bài toán của chúng ta là việc tối ưu hóa hàm:

$$w_1 = \operatorname{argmax}_{\|w_1\|=1} (w_1^T \hat{X} \hat{X}^T w_1):$$

Đặt $A = \hat{X} \hat{X}^T \in \mathbb{R}^{m \times m}$. A là một ma trận đối xứng nên sẽ có m vector riêng ứng với m trị riêng. Gọi $v_1; v_2; \dots; v_m$ là m vector riêng của ma trận A ứng với m trị riêng $\lambda_1; \lambda_2; \dots; \lambda_m$ thỏa $\|v_i\| = 1$ ($\forall i = 1; 2; \dots; m$). Ta thêm điều kiện $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_m$. Như chứng minh trên, ta có $V = \{v_1; v_2; \dots; v_m\}$ là một cơ sở trực chuẩn trong không gian $\mathbb{R}^{m \times m}$. Ta sẽ biểu diễn w_1 theo cơ sở V :

$$w_1 = c_{1,1} v_1 + c_{1,2} v_2 + \dots + c_{1,m} v_m:$$

Vì $\|w_1\| = 1$ nên ta có điều kiện $\sum_{i=1}^m c_{1,i}^2 = 1$. Đồng thời, theo định nghĩa của trị riêng và vector riêng, ta có:

$$A v_i = \lambda_i v_i:$$

Suy ra:

$$\begin{aligned}w_1 &= \operatorname{argmax}_{\|w_1\|=1} (w_1^T A w_1) \\ &= \operatorname{argmax}_{\|w_1\|=1} [w_1 \cdot (A \sum_{i=1}^m c_{1,i} v_i)] \\ &= \operatorname{argmax}_{\|w_1\|=1} (w_1 \cdot \sum_{i=1}^m \lambda_i c_{1,i} v_i):\end{aligned}$$

Vì $v_i \cdot v_j = 0, \forall i \neq j$ nên:

$$w_1 = \operatorname{argmax}_{\|w_1\|=1} \sum_{i=1}^m \lambda_i c_{1,i}^2 v_i^2:$$

Đến đây, vì ta đã chuẩn hóa cho $\|v_i\| = 1$ nên ta có thể rút gọn được công thức như sau:

$$w_1 = \operatorname{argmax}_{\|w_1\|=1} \sum_{i=1}^m c_{1,i}^2$$

Áp dụng điều kiện $\sum_{i=1}^m c_{1,i}^2 = 1$, và $c_{1,i} = \max\{c_{i,j}\}$, ta có:

$$\sum_{i=1}^m c_{1,i}^2 \leq c_{1,i} \cdot \sum_{i=1}^m c_{1,i} = c_{1,i}$$

Suy ra:

$$\max \sum_{i=1}^m c_{1,i}^2 = c_{1,i}$$

Đẳng thức xảy ra khi và chỉ khi $c_{1,1} = 1$ và $c_{1,i} = 0 \forall i \neq 1$, hay ta có $w_1 \equiv v_1$.

Vậy ta đã tìm được thành phần có **phân tán** lớn nhất. Để đi tìm thành phần có **phân tán** lớn tiếp theo, ta giải lần lượt các bài toán đã nêu ở trên:

$$w_i = \operatorname{argmax}_{\substack{\|w_i\|=1 \\ w_i \cdot w_j = 0, \forall j < i}} (w_i^T A w_i)$$

Từ điều kiện $w_i \cdot w_j = 0 \forall j < i$, ta thấy rằng $\forall i \geq 2: [w_i]_v = [0 \quad \dots \quad 0 \quad c_{i,i} \quad c_{i,i+1} \quad \dots \quad c_{i,m}]$:

Giải lần lượt các bài toán này một cách tương tự như chiều thứ nhất, ta tìm ra được m **chi u quan trọng** của bộ dữ liệu, cũng chính là m vector riêng ứng với m trị riêng của ma trận A . Ở đây giá trị riêng càng lớn thì chiều của vector riêng đó càng quan trọng.

4.3 Kết luận

Như vậy, **bài toán PCA** thực chất là **bài toán tìm các vector riêng của ma trận $A = \hat{X}\hat{X}^T$** . Ma trận A còn được gọi là **ma trận hiệp phương sai** của tập dữ liệu X .

5 Giới thiệu thuật toán tìm trị riêng

Như đã phân tích trong các phần trên, bài toán của chúng ta được đưa về bài toán tìm các **trị riêng** và các **vector riêng** của **ma trận hiệp phương sai $\hat{X}\hat{X}^T$** . Nhưng việc tìm trị riêng chính là giải phương trình $\det(A - \lambda I_n) = 0$ ($A \in \mathbb{R}^{n \times n}$). Phương trình này sẽ đưa về phương trình bậc n . Nhưng ta đã biết việc giải phương trình bậc cao sẽ rất khó khăn và với $n \geq 5$ thì ta không có công thức nghiệm tổng quát.

Tuy nhiên **ma trận hiệp phương sai** của chúng ta là ma trận **đối xứng** và **hữu hạn**, cho nên chúng ta sẽ có một vài thuật toán có thể giải được nghiệm của định thức này.

5.1 Phương pháp Power

Đây là phương pháp tìm trị riêng và vector riêng của một ma trận **đối xứng** $A \in \mathbb{R}^{n \times n}$

Thuật toán: Trước tiên, ta sẽ tìm vector riêng ứng với trị riêng lớn nhất λ_1 :

- Bước 1: Chọn bất kì một vector $q_0 \in \mathbb{R}^{n \times n} : \|q_0\| = 1$:
- Bước 2: Với $k = 1; 2; 3; \dots$, ta tính $q_k = Aq_{k-1}$:
- Bước 3: Chuẩn hóa $q_k = \frac{q_k}{\|q_k\|}$:
- Bước 4: Kiểm tra nếu $\|q_k - q_{k-1}\|$ đủ nhỏ (các vector q hội tụ) thì ta dừng lại. Nếu không, tăng k lên 1 rồi quay lại Bước 2.
- Bước 5: q_k chính là vector riêng ứng với trị riêng lớn nhất λ_1 suy ra $\lambda_1 = q_k^T A q_k$:

Ta chứng minh được rằng q sẽ hội tụ khá nhanh.

Chứng minh. Gọi S là hệ cơ sở tạo bởi các vector riêng x của ma trận A trong \mathbb{R}^n . Khi đó vector q_0 có thể được biểu diễn dạng:

$$q_0 = a_1 x_1 + a_2 x_2 + a_3 x_3 + \dots + a_n x_n$$

với $a_1; a_2; a_3; \dots$ là tọa độ của q_0 trong S .

Sau k lần nhân như trên, ta có:

$$\begin{aligned} A^k q_0 &= A^k (a_1 x_1 + a_2 x_2 + a_3 x_3 + \dots + a_n x_n) \\ &= a_1 A^k x_1 + a_2 A^k x_2 + a_3 A^k x_3 + \dots + a_n A^k x_n \\ &= a_1 \lambda_1^k x_1 + a_2 \lambda_2^k x_2 + a_3 \lambda_3^k x_3 + \dots + a_n \lambda_n^k x_n \\ &= a_1 \lambda_1^k \left(x_1 + \sum_{j=2}^n \frac{a_j}{a_1} \left(\frac{\lambda_j}{\lambda_1} \right)^k x_j \right) \end{aligned}$$

Nếu $|\lambda_1| > |\lambda_2| > |\lambda_3| > \dots > |\lambda_n|$ thì $\left| \frac{\lambda_j}{\lambda_1} \right| < 1$ với $j > 1$:

Suy ra $\left(\frac{\lambda_j}{\lambda_1} \right)^k$ hội tụ về 0 khi k tiến đến ∞ .

Vậy $A^k q_0$ hội tụ về $a_1 \lambda_1^k x_1$. Nhưng do ta đã thực hiện chuẩn hóa liên tục Aq^{k-1} tại bước thứ k , nên ta suy ra được dãy $A^k q_0$ sẽ hội tụ về x_1 . \square

Vậy ta tính được trị riêng lớn nhất λ_1 và vector riêng ứng với nó, ta sẽ tìm được các trị riêng và vector riêng còn lại theo định lý sau.

Định lý 5.1. Nếu ma trận $n \times n$ xác định A có các trị riêng $\lambda_1 > \lambda_2 > \lambda_3 > \dots > \lambda_n > 0$ và các vector riêng tương ứng $v_1; v_2; v_3; \dots; v_n$ thì ma trận $n \times n$ các vector riêng này tạo thành 1 họ cơ sở của \mathbb{R}^n , thì ma trận $n \times n$:

$$B = A - \lambda_1 v_1 v_1^T$$

có các trị riêng $\lambda_2 > \lambda_3 > \dots > \lambda_n > 0$ và các vector riêng tương ứng là $v_2; v_3; \dots; v_n$.

Chứng minh. Ta có:

- Với $i = 1$:

$$Bv_1 = (A - \lambda_1 v_1 v_1^T) v_1 = Av_1 - \lambda_1 v_1 = 0:$$

Suy ra: λ_1 không là trị riêng của B :

- Với $i > 1$:

$$\begin{aligned} Bv_i &= (A - \lambda_1 v_1 v_1^T) v_i \\ &= Av_i - \lambda_1 v_1 (v_1^T v_i) \\ &= Av_i \text{ (vì } v_1 \text{ và } v_i \text{ trực giao)} \\ &= \lambda_i v_i \end{aligned}$$

Vậy khi $i > 1$ thì λ_i là trị riêng của B . Vậy ta tìm được trị riêng lớn nhất λ_2 và vector riêng v_2 tương ứng của B . Tiếp tục làm như vậy, ta sẽ tìm được n trị riêng và vector riêng của ma trận A lúc đầu. \square

5.2 Nhược điểm của thuật toán Power

Thuật toán Power cho ta một cách xấp xỉ trị riêng, vector riêng khá tốt và đơn giản, song, phương pháp này cũng có những hạn chế khiến nó không thường được sử dụng trong thực tế cho việc này. Trong số đó có thể kể đến như:

- Trong trường hợp không may mắn, có xác suất (rất nhỏ) việc vector q khởi tạo vuông góc với vector riêng ứng với trị riêng lớn nhất. Khi đó, thuật toán Power sẽ cho ra kết quả là vector riêng ứng với trị riêng lớn thứ 2 trong số các vector riêng của ma trận A .
- Nếu một ma trận có nhiều vector riêng ứng với một trị riêng thì thuật toán Power sẽ chỉ cho ra một vector riêng ứng với trị riêng đó.
- Nếu các trị riêng của ma trận A có giá trị xấp xỉ bằng nhau, tốc độ hội tụ của thuật toán Power sẽ khá chậm.

Tuy nhiên, thuật toán Power vẫn có những ứng dụng thực tế trong một số trường hợp cụ thể ví dụ như trong Google Pagerank để sắp xếp các website theo độ phổ biến giảm dần của chúng.

5.3 Tổng kết

Như vậy, phương pháp Power tuy là một thuật toán xấp xỉ đơn giản nhưng nó lại có những trường hợp riêng mà thuật toán hoạt động không hiệu quả. Chính vì vậy, người ta thường dùng những thuật toán khác như thuật toán **QR** hay **SVD** để thay thế.

6 Cài đặt

6.1 Các bước thực hiện Phương pháp PCA

Sau khi phân tích về PCA, chúng ta rút ra được các bước thực hiện chính như sau:

1. Tìm vector trung bình của toàn bộ dữ liệu:

$$\forall i \leq m : \bar{x}_i = \frac{1}{n} \sum_{j=1}^n x_{ij}$$

2. Chuẩn hóa lại toàn bộ dữ liệu:

$$\forall j \leq n : \text{col}_j(\hat{X}) = \text{col}_j(X) - \bar{x}$$

3. Tính ma trận hiệp phương sai:

$$\text{Cov}(X) = \frac{1}{n} \hat{X} \hat{X}^T:$$

Nhưng thực tế khi lập trình, ta có thể không cần chia n (do các ma trận tương đương nhau nên không ảnh hưởng kết quả cuối cùng của ta).

4. Tính các trị riêng và vector riêng có *norm* bằng 1 của ma trận $\text{Cov}(X)$, sắp xếp chúng theo thứ tự giảm dần của trị riêng.
5. Chọn k vector để giữ lại ứng với k trị riêng lớn nhất của $\text{Cov}(X)$ để tạo thành hệ cơ sở S có k chiều. Hệ cơ sở S này tạo thành không gian con mà ta cần chiếu bộ dữ liệu của mình lên.
6. Thực hiện phép chiếu \hat{X} vào hệ cơ sở S , ta được dữ liệu mới chính là tọa độ của các điểm của \hat{X} trong S :

$$Z = S\hat{X}:$$

6.2 Source code trên ngôn ngữ lập trình Python

Để dễ dàng hơn khi lập trình, ta quy ước ngược lại với phân tích ở trên một tí là ma trận X nằm trong $\mathbb{R}^{n \times m}$ với n hàng là các bộ dữ liệu, m cột là các đặc tính của dữ liệu.

6.2.1 Code lại thủ công toàn bộ PCA

```
1 def PCA(data, k):
2     import numpy as np
3     X = np.matrix(data, dtype = float)
4     n, m = X.shape # n là số dữ liệu, m là số đặc tính
5     # Chuẩn hóa theo vector trung bình
6     X -= X.mean(axis=0)
7     # Tính ma trận hiệp phương sai
8     Cov = np.cov(X, rowvar=False)
9     # tính trị riêng và vector riêng
10    eVector = np.linalg.eigh(Cov)[1]
11    # np.linalg.eigh(Cov)[0] trả ra các trị riêng sắp theo thứ tự tăng dần
12    # np.linalg.eigh(Cov)[1] trả ra ma trận các vector riêng ở dạng cột
13    # theo thứ tự tương ứng như trên
14    eVector = (eVector.T)[::-1]
15    # đổi thành ma trận vector hàng và sắp lại
16    # theo thứ tự trị riêng giảm dần
17    S = np.matrix(eVector[:k]) # tạo hệ cơ sở mới với k chiều đầu tiên
18    return (S*(X.T)).A
19    # trả ra kết quả ở dạng np.ndarray để dễ xử lý hơn thay np.matrix
```

6.2.2 Sử dụng thư viện sklearn

```
1 import numpy as np
2 import scipy
3 from sklearn.decomposition import PCA
4
5 X = # Bộ dữ liệu
6 k = # Số chiều muốn giữ lại
7 pca = PCA(n_components = k)
8 pca.fit(X.T)
9 res = pca.components_.T
```

7 Áp dụng mô hình

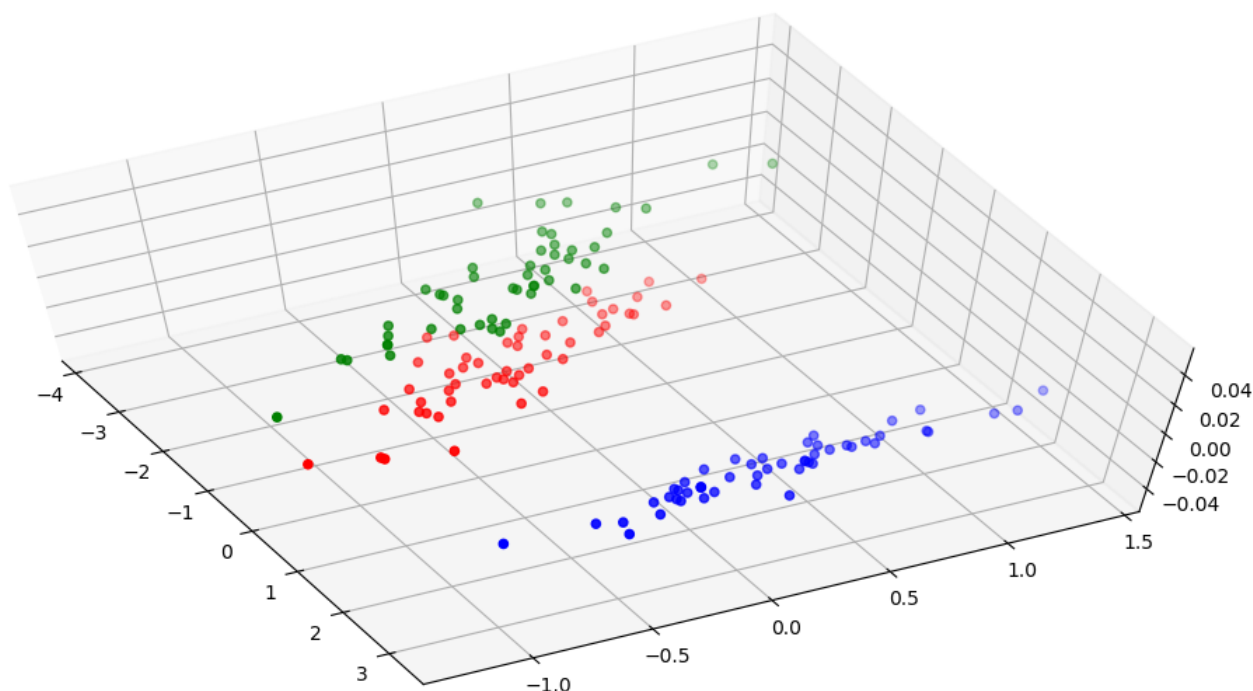
Trong thực tế, dữ liệu thường sẽ có rất nhiều tính trạng và vector biểu diễn chúng sẽ có số chiều tương ứng. Tuy vậy, con người chỉ có thể cảm nhận được vị trí và hình ảnh trong không gian nhỏ hơn hay bằng 3 chiều. Chính vì thế, nếu ta có thể áp dụng PCA để giảm chiều dữ liệu xuống còn 3 hay thấp hơn, ta có thể biểu diễn những điểm dữ liệu này trên đồ thị. Bây giờ, ta sẽ thực hiện PCA lên hai bộ dữ liệu mẫu được tải free trên mạng.

7.1 Hoa Iris

Bộ dữ liệu này gồm dữ liệu về độ dài và độ rộng của đài hoa và cánh hoa của 150 bông hoa Iris thuộc 3 chủng loại Iris Setosa, Iris Versicolour, Iris Virginica. Ứng với mỗi chủng loại hoa sẽ có 50 bông hoa được khảo sát.

1	sepal length	sepal width	petal length	petal width	class
2	5,1	3,5	1,4	0,2	Iris-setosa
3	4,9	3	1,4	0,2	Iris-setosa
4	4,7	3,2	1,3	0,2	Iris-setosa
5	4,6	3,1	1,5	0,2	Iris-setosa
6	5	3,6	1,4	0,2	Iris-setosa

Kết quả sau khi PCA là:

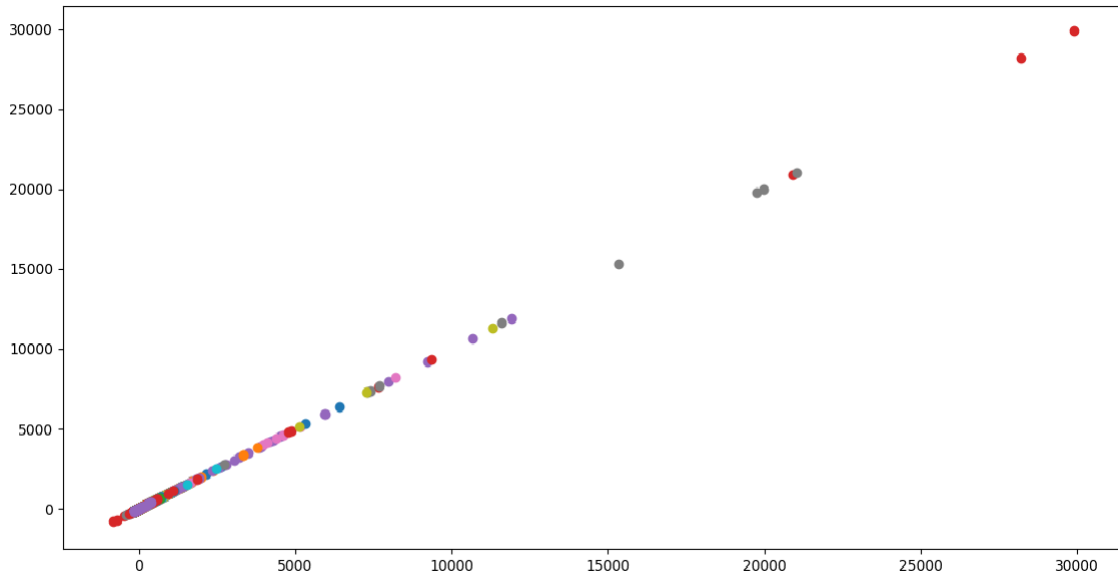


7.2 Cơ sở dữ liệu dinh dưỡng quốc gia USDA

Bộ dữ liệu này gồm danh sách các loại thức ăn và những thông số dinh dưỡng của nó như: energy kcal, protein, fat, ...

ID	FoodGroup	ShortDescrip	Descrip	Energy_kcal	Protein_g	Fat_g	Carb_g	Sugar_g
01001	Dairy and Egg Products	BUTTER,WITH SALT	Butter, salted	717	0,85	81,11	0,06	0,06
01002	Dairy and Egg Products	BUTTER,WHIPPED,WITH SALT	Butter, whipped, with salt	717	0,85	81,11	0,06	0,06
01003	Dairy and Egg Products	BUTTER OIL,ANHYDROUS	Butter oil, anhydrous	876	0,28	99,48	0	0
01004	Dairy and Egg Products	CHEESE,BLUE	Cheese, blue	353	21,4	28,74	2,34	0,5
01005	Dairy and Egg Products	CHEESE,BRICK	Cheese, brick	371	23,24	29,68	2,79	0,51

Kết quả sau khi PCA là:



7.3 EigenFace

EigenFaces là một trong những phương pháp phổ biến trong bài toán nhận diện khuôn mặt. Ý tưởng của phương pháp này là từ một bộ ảnh khuôn mặt người có kích thước lớn và số lượng nhiều, ta sẽ tìm một không gian nhỏ hơn để biểu diễn các khuôn mặt này.

EigenFace chính là PCA: các Eigen Faces (*khuôn m t riêng*) chính là các vector riêng ứng với các trị riêng lớn nhất. Ta có thể hiểu là các *khuôn m t riêng* sẽ mang phần nào thông tin của các khuôn mặt ban đầu. Vì thế, khuôn mặt ban đầu có thể xấp xỉ bằng tổng có trọng số của các *khuôn m t riêng* này.

Bây giờ chúng ta tìm EigenFaces trên bộ dữ liệu Yale face database. Bộ dữ liệu này gồm hình trắng đen khuôn mặt của 15 người khác nhau, mỗi người có 11 tấm ảnh được chụp ở điều kiện ánh sáng, cảm xúc khác nhau, gồm: *centerlight, glasses, happy, leftlight, noglasses, normal, rightlight, sad, sleepy, surprised* và *wink*. Các ảnh trong bộ dữ liệu này đã được chuẩn hóa cùng kích thước 116×98 và khuôn mặt nằm gọn trong ảnh.

Source code:

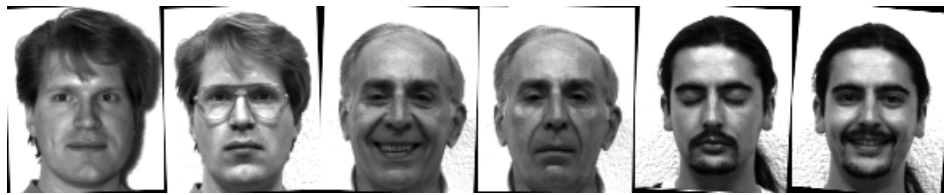
```
1 import numpy as np
2 from matplotlib import pyplot as plt
3 from sklearn.decomposition import PCA
4 import imageio
5
6 # read data
7 states = ["centerlight", "glasses", "happy", "leftlight", "noglasses",
8           "normal", "rightlight", "sad", "sleepy", "surprised", "wink"]
9 h = 116
10 w = 98
11 D = h * w
12 n = len(states) * 15
13 X = np.zeros((n, D))
14
15 cnt = 0
16 for person_id in range(1, 16):
17     for state in states:
18         path = "Yale/subject" + str(person_id).zfill(2) + '.' + state + ".pgm"
19         X[cnt] = imageio.imread(path).reshape(D)
20         cnt += 1
21
22 # Do PCA
23 k = int(input())
24 pca = PCA(n_components=k)
25 pca.fit(X)
26
27 # projection matrix
28 X_proj = pca.components_
29
30 for i in range(X_proj.shape[0]):
31     plt.axis('off')
32     plt.imshow(X_proj[i].reshape(h, w), interpolation="nearest")
33     plt.gray()
34     path = "Eigen Faces/" + str(i).zfill(2) + ".png"
35     plt.savefig(path, bbox_inches='tight', pad_inches=0)
36
37 # See reconstruction of first 6 persons
38 for person_id in range(1, 16):
39     for state in states:
40         path = "Yale/subject" + str(person_id).zfill(2) + '.' + state + ".pgm"
41         img = imageio.imread(path)
42         # subtract mean
43         x = img.reshape(D, 1) - pca.mean_.reshape(D, 1)
44         # encode
45         z = X_proj.dot(x)
46         # decode
47         x_tilde = X_proj.T.dot(z) + pca.mean_.reshape(D, 1)
48         img_tilde = x_tilde.reshape(h, w)
49         plt.axis('off')
50         plt.imshow(img_tilde, interpolation='nearest')
51         plt.gray()
52         path = "reconstruction/" + str(person_id).zfill(2) + '.' + state + ".png"
53         plt.savefig(path, bbox_inches='tight', pad_inches=0)
```

Dưới đây là ảnh của 9 Eigen Faces đầu tiên:



Để đánh giá độ hiệu quả của EigenFace, ta biểu diễn lại một vài ảnh để so sánh với ảnh ban đầu. Bộ dữ liệu gồm 165 ảnh này được biểu diễn lại bằng k khuôn $m \times t$ riêng, với k lần lượt bằng 75, 100, và 130

- Ảnh gốc:



- Tái tạo với $k = 75$:



- Tái tạo với $k = 100$:



- Tái tạo với $k = 130$:



8 Kết luận đánh giá

PCA là một cách tiếp cận phổ biến cho bài toán giảm chiều dữ liệu vì có nhiều ưu điểm như:

- PCA là một cách tiếp cận trực quan và dễ hiểu.
- PCA giúp giảm chiều của dữ liệu đi một cách đáng kể, từ đó làm cho việc xử lý và lưu trữ trở nên nhanh chóng và dễ dàng hơn. Vì vậy, trong thực tế, người ta thường thực hiện PCA kết hợp với những thuật toán khác nhằm tăng hiệu quả của thuật toán đó lên.

Tuy vậy, phương pháp này vẫn có những nhược điểm của nó:

- Độ phức tạp tính toán của PCA thường rơi vào khoảng $O(\min(n^3; m^3))$. Vì vậy, trong nhiều trường hợp, việc áp dụng PCA sẽ không hiệu quả.
- Nếu *phân tán* trên các chiều dữ liệu là xấp xỉ nhau, việc PCA sẽ làm cho lượng thông tin bị mất đi là khá lớn gây ảnh hưởng đến việc xử lý về sau.

9 Hướng nghiên cứu trong tương lai

Trong tương lai, nhóm sẽ đi sâu vào những kĩ thuật thực hiện như những thuật toán tìm trị riêng tối ưu hơn, như phương pháp QR factorization hay SVM factorization. Bên cạnh đó, các hướng mở rộng của PCA như Sparse PCA, Kernel PCA hay Robust PCA cũng là những hướng tìm hiểu thú vị.

Tài liệu

[1] *Tài liệu Tr i hệ PiMA 2018*.

[2] *Wikipedia*.

[3] *Yale face database*. <http://vismod.media.mit.edu/vismod/classes/mas622-00/datasets/>, 2016.

[4] data.world. *USDA National Nutrient Database data*. 2016.

[5] R.A. Fisher. *Iris Data Set*. <https://archive.ics.uci.edu/ml/datasets/iris>, 1936.

[6] Jeff Jauregui. *Principal component analysis with linear algebra*. 2012.

[7] Vũ Hữu Tiệp. *Machine Learning c b n*.