

Principal Component Analysis

Ngoc Hoang Luong

University of Information Technology (UIT), VNU-HCM

March 27, 2023

Phân tích thành phần chính

Lưu ý: Ngọc Hoàng

Trường Đại học Công nghệ Thông tin (UIT), ĐHQG-HCM

Ngày 27 tháng 3 năm 2023

The contents of the slides are from: Gaston Sanchez and Ethan Marzban: *All Models Are Wrong: Concepts of Statistical Learning* - <https://allmodelsarewrong.github.io/pca.html>

Nội dung của các trang trình bày là từ: Gaston Sanchez và Ethan Marzban: Tất cả các mô hình đều sai: Các khái niệm về học tập thống kê - <https://allmodelsarewrong.github.io/pca.html>

Low-dimensional Representations

- Individuals form a cloud of points in a p -dim space. Variables form a cloud of arrows in an n -dim space.
- Suppose some data in which its cloud of points form a mug:



- Is there away to get a low-dimensional representation of this data?

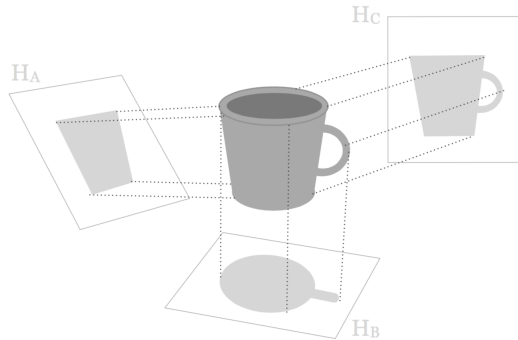
Biểu diễn chiều thấp

- Các cá nhân tạo thành một đám mây điểm trong không gian p -dim. biến dạng một đám mây mũi tên trong không gian n -mở.
- Giả sử một số dữ liệu trong đó đám mây điểm của nó tạo thành một cốc:



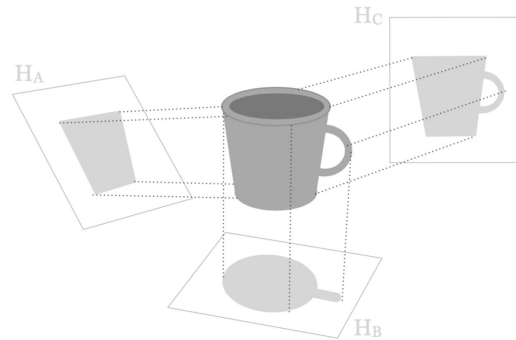
- Có cách nào để có được biểu diễn ít chiều của dữ liệu này không?

Low-dimensional Representations



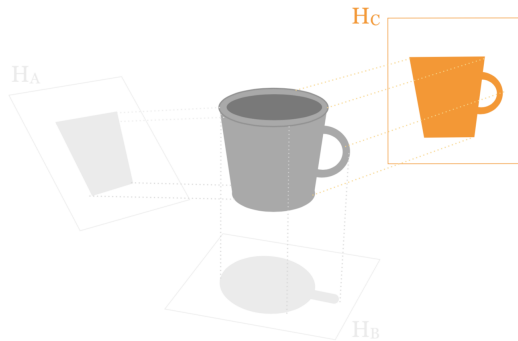
- We can look for projections of the data into sub-spaces of lower dimension.
- Assume we take a photo of the mug from different angles. What is the **best** angle to take a photo to get the images of the mug as similar as possible to the mug?

Biểu diễn chiều thấp



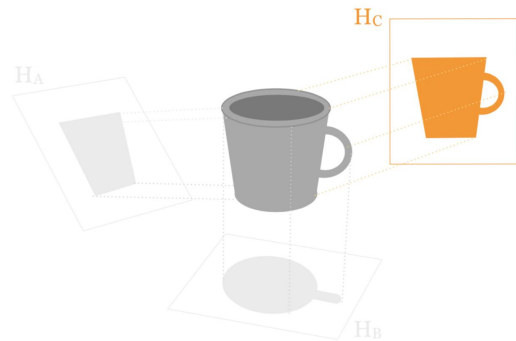
- Chúng ta có thể tìm các phép chiếu của dữ liệu vào các không gian con của cấp dư ới kích thước.
- Giả sử chúng ta chụp ảnh chiếc cốc từ các góc độ khác nhau. Góc tốt nhất để chụp ảnh là gì để có được hình ảnh của chiếc cốc giống với chiếc cốc nhất có thể?

Low-dimensional Representations



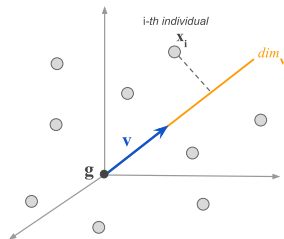
- Among 03 projections $\mathbb{H}_A, \mathbb{H}_B, \mathbb{H}_C$, the subspace \mathbb{H}_C provides the best low-dimensional representation.
- The resulting image in low-dimensional space is not capturing the whole pattern: there is always some loss of information.
- Choosing the right projection, we try to minimize such loss.

Biểu diễn chiều thấp



- Trong 03 phép chiếu H_A, H_B, H_C , không gian con H_C cung cấp đại diện chiều thấp tốt nhất.
- Hình ảnh thu được trong không gian ít chiều không thu được toàn bộ mẫu: luôn có một số thông tin bị mất.
- Chọn phép chiếu phù hợp, chúng tôi cố gắng giảm thiểu tổn thất đó.

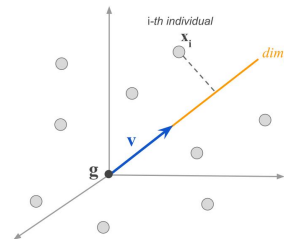
Projections



- Data points are in a p -dimensional space, and the cloud has its centroid g .
- We first try the simplest low-dimensional space: a 1D space, which can be displayed as one axis, denoted as dim_v .

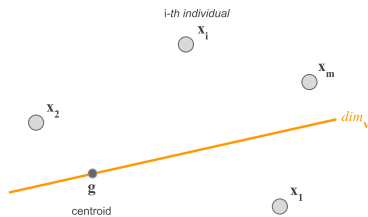
Machine Translated by Google

dự đoán



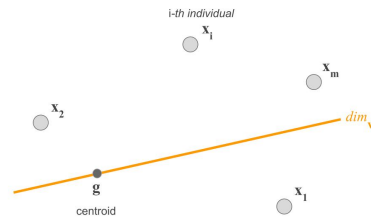
- Các điểm dữ liệu nằm trong không gian p chiều và đám mây có trọng tâm g .
- Trước tiên, chúng tôi thử không gian ít chiều đơn giản nhất: không gian 1D, trong đó có thể được hiển thị dưới dạng một trục, được ký hiệu là dim_v .

Projections



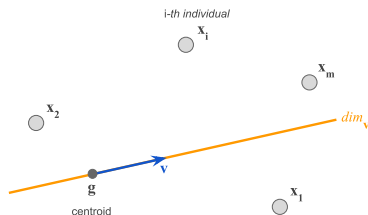
- Data points are in a p -dimensional space, and the cloud has its centroid g .
- We first try the simplest low-dimensional space: a 1D space, which can be displayed as one axis, denoted as dim_v .

dự đoán



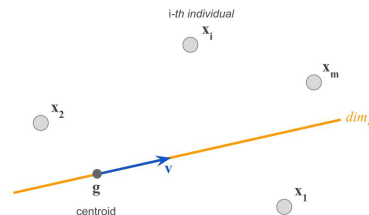
- Các điểm dữ liệu nằm trong không gian p chiều và đám mây có trọng tâm g .
- Trước tiên, chúng tôi thử không gian ít chiều đơn giản nhất: không gian 1D, trong đó có thể được hiển thị dưới dạng một trục, được ký hiệu là dim_v .

Projections



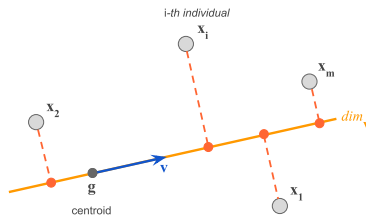
- Data points are in a p -dimensional space, and the cloud has its centroid g .
- We first try the simplest low-dimensional space: a 1D space, which can be displayed as one axis, denoted as dim_v .
- We manipulate dim_v via a vector v along this dimension.

dự đoán



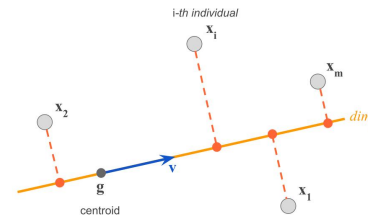
- Các điểm dữ liệu nằm trong không gian p chiều và đám mây có trọng tâm g .
- Trước tiên, chúng tôi thử không gian ít chiều đơn giản nhất: không gian 1D, trong đó có thể được hiển thị dưới dạng một trục, được ký hiệu là dim_v .
- Chúng ta thao tác dim_v thông qua một vectơ v dọc theo chiều này.

Projections



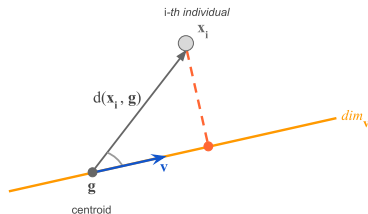
- Data points are in a p -dimensional space, and the cloud has its centroid g .
- We first try the simplest low-dimensional space: a 1D space, which can be displayed as one axis, denoted as dim_v .
- We manipulate dim_v via a vector v along this dimension.
- We want to project orthogonally the individuals onto this dimension.

dự đoán



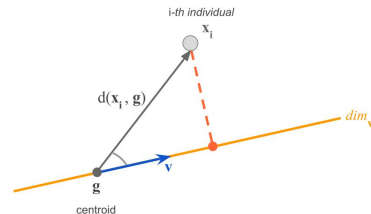
- Các điểm dữ liệu nằm trong không gian p chiều và đám mây có trọng tâm g .
- Trước tiên, chúng tôi thử không gian ít chiều đơn giản nhất: không gian 1D, trong đó có thể được hiển thị dưới dạng một trục, được ký hiệu là dim_v .
- Chúng ta thao tác dim_v thông qua một vectơ v dọc theo chiều này.
- Chúng tôi muốn chiếu trực giao các cá nhân lên điều này kích thước.

Vector and Scalar Projections



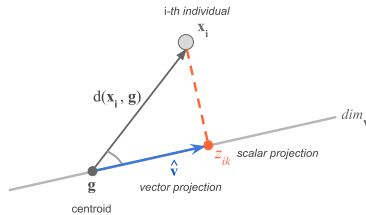
- Take the centroid g as the origin of the clouds of points.
- The dimension that we look for has to pass through the origin.
- Obtain the orthogonal projection of the i -th individual onto dim_v is projecting x_i onto any vector v along this dimension.

Phép chiếu vectơ và vô hướng



- Lấy trọng tâm g làm gốc của đám mây điểm.
- Chiều ta tìm phải đi qua gốc tọa độ.
- Có được hình chiếu trực giao của cá thể thứ i lên dim_v là hình chiếu x_i lên bất kỳ vectơ v nào dọc theo chiều này.

Vector and Scalar Projections



- The **vector projection** of \mathbf{x}_i onto \mathbf{v} is:

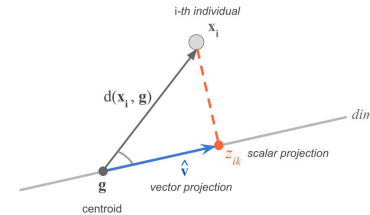
$$\hat{\mathbf{v}} = \frac{\mathbf{v}^T \mathbf{x}_i}{\mathbf{v}^T \mathbf{v}} \mathbf{v}$$

- The **scalar projection** of \mathbf{x}_i onto \mathbf{v} is:

$$z_{ik} = \frac{\mathbf{v}^T \mathbf{x}_i}{\|\mathbf{v}\|}$$

- We would prefer the **scalar projection** to obtain the co-ordinate of \mathbf{x}_i along this axis.

Phép chiếu vectơ và vô hướng



- Hình **chiếu vectơ** của x_i lên v là:

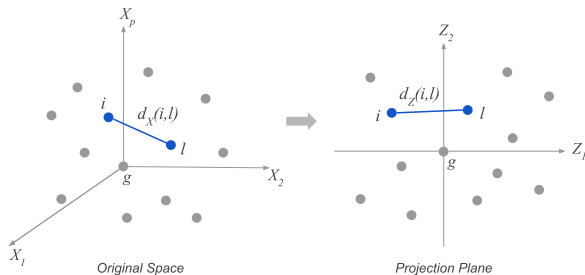
$$\hat{v} = \frac{v^T x_i}{v^T v} v$$

- Hình **chiếu vô hướng** của x_i lên v là: v

$$z_{ik} = \frac{x_i}{v}$$

- Chúng tôi muốn **phép chiếu vô hướng** thu được tọa độ của x_i dọc theo trục này.

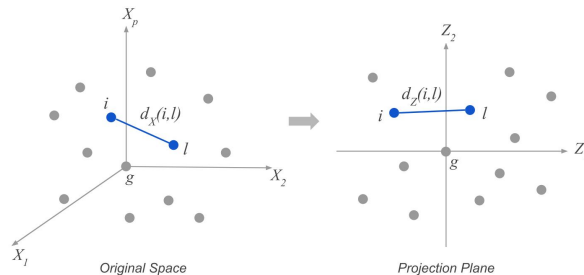
Projected Inertia



- Find the angle that give the best photo of the object \iff Find the subspace that the distances between the points are the most similar to the original points.
- The overall dispersion of the original data is: $\sum_{i=1}^n \sum_{l=1}^n d^2(i,l)$. We try to find a subspace \mathbb{H} such that:

$$\sum_{i=1}^n \sum_{l=1}^n d^2(i,l) \approx \sum_{i=1}^n \sum_{l=1}^n d_{\mathbb{H}}^2(i,l)$$

quán tính dự kiến



- Tìm góc cho ảnh đẹp nhất của vật Tìm không gian con sao cho khoảng cách giữa các điểm bằng nhau nhất so với điểm ban đầu. • Độ phân tán tổng thể của dữ liệu gốc là: cố gắng tìm một không gian con H sao cho:

$$\sum_{i=1}^N \sum_{l=1}^N d_{\text{ngày}}^2(i,l) \approx \sum_{i=1}^N \sum_{l=1}^N d_{\text{ngày xin chào, tôi}}^2(i,l)$$

- The overall dispersion is related to the inertia as:

$$\sum_{i=1}^n \sum_{l=1}^n d^2(i, l) = 2n \sum_{i=1}^n d^2(i, g) = 2n^2 \frac{1}{n} \sum_{i=1}^n d^2(i, g) = 2n^2 \text{Inertia}$$

- Finding the subspace \mathbb{H} that yields similar distances to the original subspace corresponds to maximize the projected inertia:

$$\max_{\mathbb{H}} \left\{ \frac{1}{n} \sum_{i=1}^n d_{\mathbb{H}}^2(i, g) \right\}$$

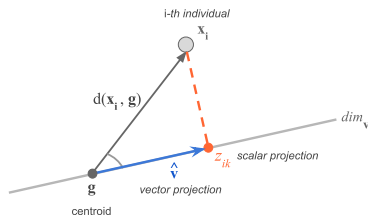
- Độ phân tán tổng thể liên quan đến quán tính như :

$$\sum_{i=1}^N \sum_{l=1}^N d^2(i, l) = 2n \sum_{i=1}^N d^2(i, g) = 2n^2 \frac{1}{n} \sum_{i=1}^N d^2(i, g) = 2n^2 \text{quán tính}$$

- Tìm không gian con H có khoảng cách tứ giác tự với không gian con ban đầu tứ giác ứng để tối đa hóa quán tính dự kiến:

$$\max_{H} \left\{ \frac{1}{N} \sum_{i=1}^N d_H^2(i, g) \right\}$$

Projected Inertia



- We are consider 1D case, $\mathbb{H} \subseteq \mathbb{R}^1$, the projected inertia becomes:

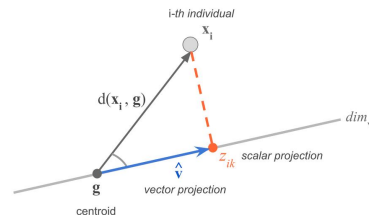
$$\frac{1}{n} \sum_{i=1}^n d_{\mathbb{H}}^2(i, g) = \frac{1}{n} \sum_{i=1}^n (\mathbf{x}_i^T \mathbf{v})^2 = \frac{1}{n} \sum_{i=1}^n z_i^2$$

- Our maximization problem becomes:

$$\max_{\mathbf{v}} \left\{ \frac{1}{n} \sum_{i=1}^n (\mathbf{x}_i^T \mathbf{v})^2 \right\} \quad \text{s.t.} \quad \mathbf{v}^T \mathbf{v} = 1$$

- We constraint \mathbf{v} to be a unit vector; otherwise, the maximization objective is unbounded.

quán tính dự kiến



- Ta xét trường hợp 1D, $\mathbb{H} \subseteq \mathbb{R}^1$, quán tính dự kiến trở thành:

$$\frac{1}{n} \sum_{i=1}^n d_{\mathbb{H}}^2(i, g) = \frac{1}{n} \sum_{i=1}^n (\mathbf{x}_i^T \mathbf{v})^2 = \frac{1}{n} \sum_{i=1}^n z_i^2$$

- Bài toán tối đa hóa của chúng ta trở thành:

$$\text{tối đa} \quad \frac{1}{n} \sum_{i=1}^n (\mathbf{x}_i^T \mathbf{v})^2 \quad \text{s.t.} \quad \mathbf{v}^T \mathbf{v} = 1$$

- Ta ràng buộc \mathbf{v} là véc tơ đơn vị; mặt khác, việc tối đa hóa mục tiêu là không giới hạn.

Maximization Problem

- Assume mean-centered data, the centroid \mathbf{g} of the cloud of points is the origin $\mathbf{g} = \mathbf{0}$.
- We are projecting onto a line spanned by a unit-vector \mathbf{v} , the projected inertia $I_{\mathbb{H}}$ is the variance of the projected data points:

$$\begin{aligned} I_{\mathbb{H}} &= \frac{1}{n} \sum_{i=1}^n d_{\mathbb{H}}^2(i, 0) = \frac{1}{n} \sum_{i=1}^n (\mathbf{x}_i^{\top} \mathbf{v})^2 \\ &= \frac{1}{n} \sum_{i=1}^n z_i^2 = \frac{1}{n} \mathbf{z}^{\top} \mathbf{z} = \frac{1}{n} \mathbf{v}^{\top} \mathbf{X}^{\top} \mathbf{X} \mathbf{v} \end{aligned}$$

where

$$\mathbf{z} = \begin{pmatrix} z_1 \\ z_2 \\ \vdots \\ z_n \end{pmatrix} = \mathbf{X} \mathbf{v} = \begin{bmatrix} - & - & - & \mathbf{x}_1^{\top} & - & - & - \\ - & - & - & \mathbf{x}_2^{\top} & - & - & - \\ - & - & - & - & - & - & - \\ - & - & - & \mathbf{x}_n^{\top} & - & - & - \end{bmatrix} \begin{pmatrix} v_1 \\ v_2 \\ \vdots \\ v_p \end{pmatrix}$$

Vấn đề tối đa hóa

- Giả sử dữ liệu lấy trung bình làm trung tâm, tâm \mathbf{g} của đám mây điểm là gốc tọa độ $\mathbf{g} = \mathbf{0}$.

Chúng ta đang chiếu lên một đường được kéo dài bởi một vectơ đơn vị \mathbf{v} , quán tính dự kiến $I_{\mathbb{H}}$ là phương sai của các điểm dữ liệu dự kiến:

$$\begin{aligned} I_{\mathbb{H}} &= \frac{1}{n} \sum_{i=1}^n d_{\mathbb{H}}^2(i, 0) = \frac{1}{n} \sum_{i=1}^n (\mathbf{x}_i^{\top} \mathbf{v})^2 \\ &= \frac{1}{n} \sum_{i=1}^n z_i^2 = \frac{1}{n} \mathbf{z}^{\top} \mathbf{z} = \frac{1}{n} \mathbf{v}^{\top} \mathbf{X}^{\top} \mathbf{X} \mathbf{v} \end{aligned}$$

Ở đây

$$\begin{aligned} \begin{pmatrix} z_1 \\ z_2 \\ \vdots \\ z_n \end{pmatrix} &= \mathbf{X} \mathbf{v} = \begin{bmatrix} \mathbf{x}_1^{\top} & \mathbf{x}_2^{\top} & \dots & \mathbf{x}_n^{\top} \end{bmatrix} \begin{pmatrix} v_1 \\ v_2 \\ \vdots \\ v_p \end{pmatrix} \end{aligned}$$

Maximization Problem

- The maximization problem becomes:

$$\max_{\mathbf{v}} \left\{ \frac{1}{n} \mathbf{v}^T \mathbf{X}^T \mathbf{X} \mathbf{v} \right\} \quad \text{s.t.} \quad \mathbf{v}^T \mathbf{v} = 1$$

- To solve this maximization, problem, we use Lagrange multipliers.

$$\mathcal{L} = \frac{1}{n} \mathbf{v}^T \mathbf{X}^T \mathbf{X} \mathbf{v} - \lambda (\mathbf{v}^T \mathbf{v} - 1)$$

- Set the derivative of the Lagrangian \mathcal{L} wrt \mathbf{v} to 0:

$$\frac{\partial \mathcal{L}}{\partial \mathbf{v}} = \frac{2}{n} \mathbf{X}^T \mathbf{X} \mathbf{v} - 2\lambda \mathbf{v} = \mathbf{0} \Rightarrow \underbrace{\frac{1}{n} \mathbf{X}^T \mathbf{X} \mathbf{v}}_{\mathbf{S} \in \mathbb{R}^{p \times p}} = \lambda \mathbf{v} \Rightarrow \mathbf{S} \mathbf{v} = \lambda \mathbf{v}$$

- This means that \mathbf{v} is an eigenvector (with eigenvalue λ) of \mathbf{S} .
- λ is the value of the projected inertia $I_{\mathbb{H}}$ that we want to maximize.

Vấn đề tối đa hóa

- Bài toán cực đại hóa trở thành:

$$\max_{\mathbf{v}} \left\{ \frac{1}{N} \mathbf{v}^T \mathbf{X}^T \mathbf{X} \mathbf{v} \right\} \quad \text{s.t.} \quad \mathbf{v}^T \mathbf{v} = 1$$

- Để giải bài toán cực đại hóa này, ta sử dụng hệ số nhân Lagrange.

$$L = \frac{1}{N} \mathbf{v}^T \mathbf{X}^T \mathbf{X} \mathbf{v} - \lambda (\mathbf{v}^T \mathbf{v} - 1)$$

- Đặt đạo hàm của Lagrangian L wrt \mathbf{v} thành 0:

$$\frac{\partial L}{\partial \mathbf{v}} = \frac{2}{N} \mathbf{X}^T \mathbf{X} \mathbf{v} - 2\lambda \mathbf{v} = \mathbf{0} \Rightarrow \underbrace{\frac{1}{N} \mathbf{X}^T \mathbf{X} \mathbf{v}}_{\mathbf{S} \in \mathbb{R}^{p \times p}} = \lambda \mathbf{v} \Rightarrow \mathbf{S} \mathbf{v} = \lambda \mathbf{v}$$

- Điều này có nghĩa là \mathbf{v} là một vectơ riêng (với giá trị riêng λ) của \mathbf{S} .
- λ là giá trị của $I_{\mathbb{H}}$ quán tính dự kiến mà chúng ta muốn cực đại hóa.

Eigenvectors of S

- Assume \mathbf{X} is full rank ($\text{rank}(\mathbf{X}) = p$). We have p eigenvectors:

$$\mathbf{V} = [\mathbf{v}_1 \quad \mathbf{v}_2 \quad \dots \quad \mathbf{v}_k \quad \dots \quad \mathbf{v}_p]$$

- We also have the matrix of eigenvalues $\mathbf{\Lambda} = \text{diag}\{\lambda_i\}_{i=1}^n$

$$\mathbf{\Lambda} = \begin{bmatrix} \lambda_1 & 0 & \dots & 0 \\ 0 & \lambda_2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \lambda_p \end{bmatrix}$$

- We then have the matrix of projected points \mathbf{Z} (also known as **the matrix of principal components (PC's)**):

$$\mathbf{Z} = [\mathbf{z}_1 \quad \mathbf{z}_2 \quad \dots \quad \mathbf{z}_k \quad \dots \quad \mathbf{z}_p]$$

where the k -th principal component \mathbf{z}_k is:

$$\mathbf{z}_k = \mathbf{X}\mathbf{v}_k = v_{1k}\mathbf{x}_1 + v_{2k}\mathbf{x}_2 + \dots + v_{pk}\mathbf{x}_p$$

with \mathbf{x}_k denotes columns of \mathbf{X} .

Véc tơ riêng của S .

Giả sử \mathbf{X} là hạng đầy đủ ($\text{xếp hạng}(\mathbf{X}) = p$). Ta có p vectơ riêng:

$$\mathbf{V} = [\mathbf{v}_1 \quad \mathbf{v}_2 \quad \dots \quad \mathbf{v}_k \quad \dots \quad \text{phó chủ tịch}]$$

- Ta cũng có ma trận các giá trị riêng $\mathbf{\Lambda} = \text{diag}\{\lambda_i\}_{i=1}^n$

$$= \begin{bmatrix} \lambda_1 & 0 & \dots & 0 & 0 \\ 0 & \lambda_2 & \dots & 0 & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & \dots & \lambda_p & 0 \end{bmatrix}$$

- Khi đó ta có ma trận các điểm hình chiếu \mathbf{Z} (còn được gọi là **ma trận các thành phần chính (PC's)**):

$$\mathbf{Z} = [\mathbf{z}_1 \quad \mathbf{z}_2 \quad \dots \quad \mathbf{z}_k \quad \dots \quad \mathbf{z}_p]$$

trong đó thành phần chính thứ k \mathbf{z}_k là:

$$\mathbf{z}_k = \mathbf{X}\mathbf{v}_k = v_{1k}\mathbf{x}_1 + v_{2k}\mathbf{x}_2 + \dots + v_{pk}\mathbf{x}_p$$

với \mathbf{x}_k biểu thị các cột của \mathbf{X} .

Eigenvalues of S

- Because the data is mean-centered, we have $\text{mean}(\mathbf{x}_i) = 0$. Then, $\text{mean}(\mathbf{z}_k) = 0$.
- How about the variance of \mathbf{z}_k ?

$$\begin{aligned} \text{Var}(\mathbf{z}_k) &= \frac{1}{n} \mathbf{z}^\top \mathbf{z} = \frac{1}{n} (\mathbf{X} \mathbf{v}_k)^\top (\mathbf{X} \mathbf{v}_k) = \frac{1}{n} \mathbf{v}_k^\top \mathbf{X}^\top \mathbf{X} \mathbf{v}_k \\ &= \mathbf{v}_k^\top \mathbf{S} \mathbf{v}_k = \mathbf{v}_k^\top (\lambda_k \mathbf{v}_k) = \lambda_k (\mathbf{v}_k^\top \mathbf{v}_k) = \lambda_k \end{aligned}$$

- The k -th eigenvalue of \mathbf{S} is the variance of the k -th principal component.
- If \mathbf{X} is mean centered, $\mathbf{S} = \frac{1}{n} \mathbf{X}^\top \mathbf{X}$ is the covariance matrix of data.
- If \mathbf{X} is standardized (mean-centered and scaled by the variance), then \mathbf{S} is the correlation matrix.

Giá trị riêng của S

- Vì dữ liệu lấy trung bình làm trung tâm nên chúng ta có $\text{mean}(\mathbf{x}_i) = 0$. Khi đó, $\text{mean}(\mathbf{z}_k) = 0$.

• Còn phương sai của \mathbf{z}_k thì sao?

$$\begin{aligned} \text{Var}(\mathbf{z}_k) &= \frac{1}{n} \mathbf{z}^\top \mathbf{z} = \frac{1}{n} (\mathbf{X} \mathbf{v}_k)^\top (\mathbf{X} \mathbf{v}_k) = \frac{1}{n} \mathbf{v}_k^\top \mathbf{X}^\top \mathbf{X} \mathbf{v}_k \\ &= \mathbf{v}_k^\top \mathbf{S} \mathbf{v}_k = \mathbf{v}_k^\top (\lambda_k \mathbf{v}_k) = \lambda_k (\mathbf{v}_k^\top \mathbf{v}_k) = \lambda_k \end{aligned}$$

- Giá trị riêng thứ k của \mathbf{S} là phương sai của giá trị gốc thứ k thành phần.

liệu. • Nếu \mathbf{X} là trung bình, $\mathbf{S} = \frac{1}{n} \mathbf{X}^\top \mathbf{X}$ là ma trận hiệp phương sai của dữ

- Nếu \mathbf{X} được chuẩn hóa (trung bình là trung tâm và được chia tỷ lệ theo phương sai), thì \mathbf{S} là ma trận tương quan.

Eigenvalues of S

$$\text{Inertia} = \frac{1}{n} \sum_{i=1}^n d^2(i, g) = \sum_k \lambda_k = \text{tr} \left(\frac{1}{n} \mathbf{X}^\top \mathbf{X} \right)$$

- $\sum_{k=1}^p \lambda_k$ relates to the total amount of variability in the data.
- The principal components capture different parts of the variability in the data.

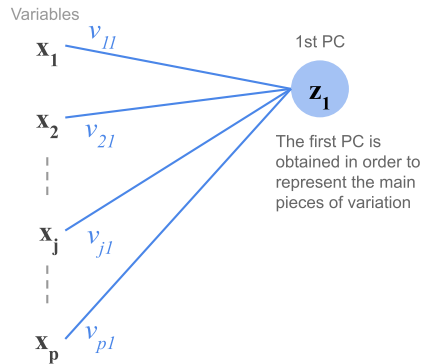
Giá trị riêng của S

$$\text{quán tính} = \frac{1}{n} \sum_{i=1}^n d^2(i, g) = \sum_k \lambda_k = \text{tr} \left(\frac{1}{n} \mathbf{X}^\top \mathbf{X} \right)$$

- $\sum_{k=1}^p \lambda_k$ liên quan đến tổng lượng biến đổi trong dữ liệu.
 - Các thành phần chính nắm bắt các phần khác nhau của sự thay đổi trong dữ liệu.

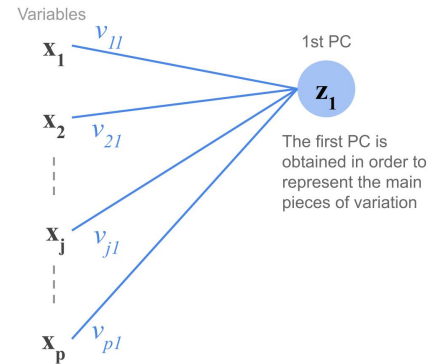
Principal Component Analysis (PCA)

- Given a set of p variables $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_p$, we want to obtain new k variables $\mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_k$, called the **Principal Components (PCs)**.
- A principal component is a **linear combination** of the p variables: $\mathbf{z} = \mathbf{X}\mathbf{v}$.
- The first PC is a linear combination:



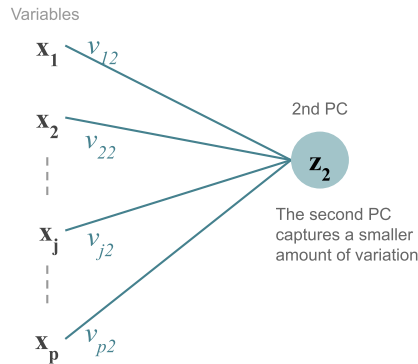
Phân tích thành phần chính (PCA)

- Cho tập p biến x_1, x_2, \dots, x_p , chúng tôi muốn lấy k các biến z_1, z_2, \dots, z_k , mới, đư ợc gọi là **Thành phần chính (PC)**.
- Thành phần chính là sự kết hợp tuyến tính của p biến: $\mathbf{z} = \mathbf{X}\mathbf{v}$.
- PC đầu tiên là tổ hợp tuyến tính:



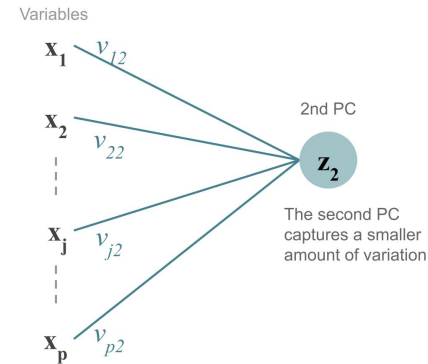
Principal Component Analysis (PCA)

- Given a set of p variables $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_p$, we want to obtain new k variables $\mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_k$, called the **Principal Components (PCs)**.
- A principal component is a **linear combination** of the p variables: $\mathbf{z} = \mathbf{X}\mathbf{v}$.
- The second PC is another linear combination:



Phân tích thành phần chính (PCA)

- Cho tập p biến x_1, x_2, \dots, x_p , chúng tôi muốn lấy k các biến z_1, \dots, z_k , mới, đợc gọi là **Thành phần chính (PC)**.
- Thành phần chính là sự kết hợp tuyến tính của p biến: $\mathbf{z} = \mathbf{X}\mathbf{v}$.
- PC thứ hai là một tổ hợp tuyến tính khác:



Principal Component Analysis (PCA)

- Given a set of p variables $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_p$, we want to obtain new k variables $\mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_k$, called the **Principal Components (PCs)**.
- A principal component is a **linear combination** of the p variables: $\mathbf{z} = \mathbf{X}\mathbf{v}$.
- We compute PCs as linear combinations of original variables:

$$\mathbf{z}_1 = v_{11}\mathbf{x}_1 + v_{21}\mathbf{x}_2 + \dots + v_{p1}\mathbf{x}_p$$

$$\mathbf{z}_2 = v_{12}\mathbf{x}_1 + v_{22}\mathbf{x}_2 + \dots + v_{p2}\mathbf{x}_p$$

$$\vdots = \vdots$$

$$\mathbf{z}_k = v_{1k}\mathbf{x}_1 + v_{2k}\mathbf{x}_2 + \dots + v_{pk}\mathbf{x}_p$$

Or:

$$\mathbf{Z} = \mathbf{X}\mathbf{V}$$

where \mathbf{Z} is an $n \times k$ matrix of principal components, and \mathbf{V} is a $p \times k$ matrix of weights (directional vectors of the principal axes).

Phân tích thành phần chính (PCA)

- Cho tập p biến x_1, x_2, \dots, x_p , chúng tôi muốn lấy k mới các biến z_1, z_2, \dots, z_k , đư ợc gọi là **Thành phần chính (PC)**.
- Thành phần chính là sự kết hợp tuyến tính của p biến: $z = \mathbf{X}\mathbf{v}$.
- Chúng tôi tính toán PC dư ới dạng tổ hợp tuyến tính của các biến ban đầu:

$$z_1 = v_{11}x_1 + v_{21}x_2 + \dots + v_{p1}x_p$$

$$z_2 = v_{12}x_1 + v_{22}x_2 + \dots + v_{p2}x_p$$

$$=$$

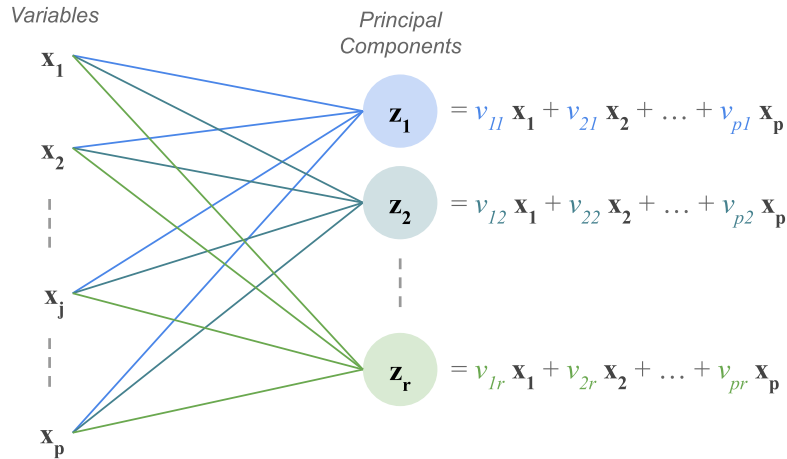
$$z_k = v_{1k}x_1 + v_{2k}x_2 + \dots + v_{pk}x_p$$

Hoặc:

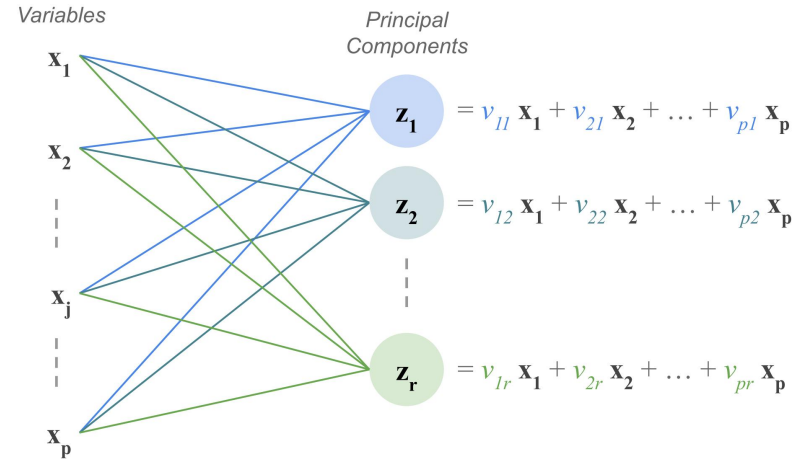
$$\mathbf{Z} = \mathbf{X}\mathbf{V}$$

trong đó \mathbf{Z} là ma trận $n \times k$ của các thành phần chính và \mathbf{V} là một ma trận trọng số $p \times k$ (vectơ chỉ phương của các trục chính).

Principal Component Analysis (PCA)



Phân tích thành phần chính (PCA)



Finding Principal Components

- The components $\mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_k$ are required to capture most of the variation in data \mathbf{X} .
- We look for a vector \mathbf{v}_h such that a component $\mathbf{z}_h = \mathbf{X}\mathbf{v}_h$ has maximum variance:

$$\max_{\mathbf{v}_h} \text{var}(\mathbf{z}_h) \Rightarrow \max_{\mathbf{v}_h} \text{var}(\mathbf{X}\mathbf{v}_h) \Rightarrow \max_{\mathbf{v}_h} \frac{1}{n} \mathbf{v}_h^\top \mathbf{X}^\top \mathbf{X} \mathbf{v}_h$$

- If \mathbf{v}_h can be arbitrarily big, the problem is unbounded. We need to restrict \mathbf{v}_h to be of unit norm:

$$\|\mathbf{v}_h\| = 1 \Rightarrow \mathbf{v}_h^\top \mathbf{v}_h = 1$$

- If we denote the covariance matrix $\mathbf{S} = (1/n)\mathbf{X}^\top \mathbf{X}$, then

$$\max_{\mathbf{v}_h} \mathbf{v}_h^\top \mathbf{S} \mathbf{v}_h \quad \text{s.t.} \quad \mathbf{v}_h^\top \mathbf{v}_h = 1$$

- To avoid redundancy, we require $\mathbf{z}_h^\top \mathbf{z}_l = 0$ mutually orthogonal if $h \neq l$.

Tìm các thành phần chính • Các thành phần \mathbf{z}_1 ,

$\mathbf{z}_2, \dots, \mathbf{z}_k$ được yêu cầu để nắm bắt hầu hết các biến thể trong dữ liệu \mathbf{X} .

- Ta tìm véc tơ \mathbf{v}_h sao cho thành phần $\mathbf{z}_h = \mathbf{X}\mathbf{v}_h$ có phương sai lớn nhất:

$$\text{tối đa}_{\mathbf{v}_h} \text{var}(\mathbf{z}_h) = \max_{\mathbf{v}_h} \text{var}(\mathbf{X}\mathbf{v}_h) = \max_{\mathbf{v}_h} \frac{1}{n} \mathbf{v}_h^\top \mathbf{X}^\top \mathbf{X} \mathbf{v}_h$$

- Nếu \mathbf{v}_h có thể lớn tùy ý thì bài toán là vô giới hạn. Chúng ta cần hạn chế \mathbf{v}_h ở dạng chuẩn đơn vị:

$$\mathbf{v}_h^\top \mathbf{v}_h = 1$$

- Nếu chúng ta ký hiệu ma trận hiệp phương sai $\mathbf{S} = (1/n)\mathbf{X}^\top \mathbf{X}$, thì

$$\max_{\mathbf{v}_h} \mathbf{v}_h^\top \mathbf{S} \mathbf{v}_h \quad \text{s.t.} \quad \mathbf{v}_h^\top \mathbf{v}_h = 1$$

- Để tránh dư thừa, chúng tôi yêu cầu $\mathbf{z}_h^\top \mathbf{z}_l = 0$ trực giao với nhau nếu $h \neq l$.

Finding Principal Components

All PCs can be found by **diagonalizing** $S = (1/n)X^T X$.

$$S = V\Lambda V^T$$

- Λ is a diagonal matrix. The diagonal elements of Λ are the eigenvalues of S .
- The columns of V are orthonormal: $V^T V = I$
- The columns of V are the eigenvectors of S .
- $V^T = V^{-1}$

Because S is a $p \times p$ symmetric matrix, we have:

- S has p real eigenvalues.
- The eigenvectors corresponding to different eigenvalues are orthogonal. S is orthogonally diagonalizable ($S = V\Lambda V^T$).
- The set of eigenvalues of S is called the **spectrum of S** .
- The PCA is obtained via an Eigenvalue Decomposition of S .

Tìm các thành phần chính Tất cả các PC

có thể được tìm thấy bằng cách **chéo hóa** $S = (1/n)X^T X$.

$$S = V\Lambda V$$

- Λ là ma trận đường chéo. Các phần tử đường chéo của Λ là các giá trị riêng của S .
- Các cột của V là trực giao: $V^T V = I$
- Các cột của V là vectơ riêng của S .
- $V^T = V^{-1}$

Vì S là ma trận đối xứng $p \times p$ nên ta có:

- S có p giá trị riêng thực.
- Các vectơ riêng tương ứng với các giá trị riêng khác nhau là trực giao. S có thể chéo hóa trực giao ($S = V\Lambda V^T$).
- Tập hợp các giá trị riêng của S được gọi là **phổ** của S .
- PCA thu được thông qua Phân tích giá trị riêng của S .

Examples

- Principal Component Analysis - Intuitions:
<https://fmin.xyz/docs/applications/pca/>
- Principal Component Analysis - Explained Visually:
<https://setosa.io/ev/principal-component-analysis/>
- Principal Component Analysis (PCA): Iris data: https://www.math.umd.edu/~petersd/666/html/iris_pca.html
- Face Recognition using Principal Component Analysis:
<https://machinelearningmastery.com/face-recognition-using-principal-component-analysis/>

ví dụ

- Phân tích thành phần chính - Trực giác:
<https://fmin.xyz/docs/applications/pca/>
- Phân tích thành phần chính - Giải thích trực quan:
<https://setosa.io/ev/principal-component-analysis/>
- Phân tích thành phần chính (PCA): Dữ liệu mống mắt: https://www.math.umd.edu/~petersd/666/html/iris_pca.html
- Nhận dạng khuôn mặt bằng cách sử dụng Phân tích thành phần chính : <https://machinelearningmastery.com/face-recognition-using-principal-component-analysis/>