

Dimensionality Reduction

Yin-Yun Li

December 24, 2025

1 Principal Component Analysis

Suppose we have a dataset with n observations and k variables, represented as a matrix:

$$X = \begin{bmatrix} x_{11} & x_{12} & \cdots & x_{1k} \\ x_{21} & x_{22} & \cdots & x_{2k} \\ \vdots & \vdots & \ddots & \vdots \\ x_{n1} & x_{n2} & \cdots & x_{nk} \end{bmatrix}.$$

For convenience, we usually center the data by subtracting the sample mean across individual for each variable. Let each column vector $\mathbf{x}_i \in \mathbb{R}^k$ be i -th observation with k dimensions, we can also write X as:

$$\begin{bmatrix} \text{---} & \mathbf{x}_1^\top & \text{---} \\ \text{---} & \mathbf{x}_2^\top & \text{---} \\ \vdots & & \vdots \\ \text{---} & \mathbf{x}_n^\top & \text{---} \end{bmatrix}.$$

Intuitively, PCA aims to find a new coordinate system such that the first direction (principal component) captures the largest variance, and the second one captures the maximum remaining variance unrelated to the previous one, and so forth.

Let the first new variable $z_{i1} = \sum_{p=1}^k \omega_{p1} x_{ip}$ for each $i = 1, 2, \dots, n$, or equivalently in matrix form:

$$\underbrace{\begin{bmatrix} z_{11} \\ z_{21} \\ \vdots \\ z_{n1} \end{bmatrix}}_{\mathbf{z}_1} = \underbrace{\begin{bmatrix} x_{11} & x_{12} & \cdots & x_{1k} \\ x_{21} & x_{22} & \cdots & x_{2k} \\ \vdots & \vdots & \ddots & \vdots \\ x_{n1} & x_{n2} & \cdots & x_{nk} \end{bmatrix}}_X \underbrace{\begin{bmatrix} \omega_{11} \\ \omega_{21} \\ \vdots \\ \omega_{k1} \end{bmatrix}}_{\boldsymbol{\omega}_1}$$

Our goal is to choose ω_1 such that the variance of \mathbf{z}_1 is maximized, subject to the constraint that ω_1 is a unit vector, so we can write down the maximization problem as:

$$\begin{aligned} \max_{\omega_1} \quad & \text{Var}(\mathbf{z}_1) = \frac{1}{n} (\mathbf{X}\omega_1)^\top \mathbf{X}\omega_1 \\ \text{subject to} \quad & \omega_1^\top \omega_1 = \|\omega_1\| = 1, \end{aligned} \quad (1)$$

where

$$\mathbf{X}\omega_1 = \begin{bmatrix} \mathbf{x}_1^\top \omega_1 \\ \mathbf{x}_2^\top \omega_1 \\ \vdots \\ \mathbf{x}_n^\top \omega_1 \end{bmatrix}, \quad (\mathbf{X}\omega_1)^\top \mathbf{X}\omega_1 = \sum_{i=1}^n (\mathbf{x}_i^\top \omega_1)^2 = \sum_{i=1}^n (\mathbf{x}_i^\top \omega_1)^\top (\mathbf{x}_i^\top \omega_1) = \omega_1^\top \left(\sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^\top \right) \omega_1.$$

Note that we first focus on the sum of squares of \mathbf{z}_1 . The Lagrangian function is given by:

$$\mathcal{L}(\omega_1, \lambda_1) = \omega_1^\top \mathbf{X}^\top \mathbf{X}\omega_1 - \lambda_1(\omega_1^\top \omega_1 - 1) \quad (2)$$

The First-order condition yields:

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial \omega_1} &= 2\mathbf{X}^\top \mathbf{X}\omega_1 - 2\lambda_1\omega_1 = \mathbf{0} \\ \mathbf{X}^\top \mathbf{X}\omega_1 &= \lambda_1\omega_1. \end{aligned} \quad (3)$$

The above equation shows that ω_1 is an eigenvector of $\mathbf{X}^\top \mathbf{X}$ w.r.t. the eigenvalue λ_1 .

Next, we try to find the second principal component $\mathbf{z}_2 = \mathbf{X}\omega_2$, which is uncorrelated with the first one. Thus, we need to solve the following maximization problem:

$$\begin{aligned} \max_{\omega_2} \quad & \text{Var}(\mathbf{z}_2) = \frac{1}{n} (\mathbf{X}\omega_2)^\top \mathbf{X}\omega_2 \\ \text{subject to} \quad & \omega_2^\top \omega_2 = \|\omega_2\| = 1; \quad \omega_1^\top \omega_2 = 0. \end{aligned} \quad (4)$$

Apply the Lagrangian method again:

$$\mathcal{L}(\omega_2, \lambda_2, \gamma) = \omega_2^\top \mathbf{X}^\top \mathbf{X}\omega_2 - \lambda_2(\omega_2^\top \omega_2 - 1) - \gamma(\omega_1^\top \omega_2). \quad (5)$$

The First-order condition yields:

$$\frac{\partial \mathcal{L}}{\partial \omega_2} = 2\mathbf{X}^\top \mathbf{X}\omega_2 - 2\lambda_2\omega_2 - \gamma\omega_1 = \mathbf{0} \quad (6)$$

Multiplying both sides of Equation (6) by ω_1^\top gives:

$$2\omega_1^\top X^\top X \omega_2 - \gamma \omega_1^\top \omega_1 = \mathbf{0}. \quad (7)$$

Note that $\omega_1^\top X^\top X = \lambda_1 \omega_1^\top$ from (3), and this forces $\gamma = 0$ and hence

$$X^\top X \omega_2 = \lambda_2 \omega_2. \quad (8)$$

Now we construct the matrix of eigenvectors as $\Omega = [\omega_1, \omega_2, \dots, \omega_k]$. Inductively, we can regard the system of equations as $Z = X\Omega$, where $Z = [\mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_k]$. Since $X^\top X$ is a symmetric matrix, there exists an invertible matrix P such that: $P^{-1}X^\top X P$ is diagonal, i.e., $X^\top X$ is diagonalizable¹.

Moreover, the diagonal matrix takes the form of:

$$Z^\top Z = \Omega^\top X^\top X \Omega = \begin{bmatrix} \lambda_1 & 0 & \cdots & 0 \\ 0 & \lambda_2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \lambda_k \end{bmatrix},$$

which shows that $P = \Omega$ is an orthonormal matrix and $\text{tr}(Z^\top Z) = \text{tr}(X^\top X) = \sum_{i=1}^k \lambda_i$. For i -th principal component, the contribution to total variance is given by the proportion $\frac{\lambda_i}{\sum_{i=1}^k \lambda_i}$. Suppose, WLOG, $n \geq k$. If $\text{rank}(X) < k$, then some eigenvalues are exactly zero, indicating that those components contribute no variance and can be discarded. If $\text{rank}(X) = k$, we ideally hope that the trailing eigenvalues are sufficiently small, such that their contribution to the total variance is negligible.

¹For now, I have not learned the spectral theorem, I can only write these problems in terms of diagonalization theorem.