

# Chapter 10. Dimensionality Reduction with PCA



# Chapter 10. Dimensionality Reduction with PCA



- 1 10.1 Problem Setting
- 2 10.2 Maximum Variance Perspective
- 3 10.3 Projection Perspective
- 4 10.4 Eigenvector Computation and Low-Rank Approximations
- 5 10.5 PCA in High Dimensions
- 6 10.6 Key Steps of PCA in Practice

# 10.1 Problem Setting

Consider a data set  $\mathcal{X} = \{x_1, \dots, x_N\}$ ,  $x_n \in \mathbb{R}^D$  with mean 0 and the data covariance matrix:

$$S = \frac{1}{N} \sum_{n=1}^N x_n x_n^T. \quad (1)$$

Assume that there exists a low-dimensional compressed representation (code) of  $x_n$

$$z_n = B^T x_n \in \mathbb{R}^M, \quad (2)$$

where

$$B := [b_1 \dots b_M] \in \mathbb{R}^{D \times M} \text{ the projection matrix is orthonormal.} \quad (3)$$

We seek an  $M$ -dimensional subspace  $U$  of  $\mathbb{R}^D$  onto which we project the data. Denote the projected data by  $\tilde{x}_n \in U$ , and their coordinates by  $z_n$ .

Our aim is to find projections  $\tilde{x}_n$  so that they are as similar to the original data  $x_n$  and minimize the loss due to compression. □

## 10.2 Maximum Variance Perspective

### Direction with Maximal Variance:

Seek a single vector  $b_1 \in \mathbb{R}^D$  that maximizes the variance of the projected data, that is we aim to maximize the variance of the first coordinate  $z_1$  of  $z \in \mathbb{R}^M$  so that

$$\begin{aligned} V_1 &:= V[z_1] = \frac{1}{N} \sum_{n=1}^N z_{1n}^2 = \frac{1}{N} \sum_{n=1}^N (b_1^T x_n)^2 = \frac{1}{N} \sum_{n=1}^N b_1^T x_n x_n^T b_1 \\ &= b_1^T \left( \frac{1}{N} \sum_{n=1}^N x_n x_n^T \right) b_1 = b_1^T S b_1. \end{aligned} \quad (4)$$

We lead to a constrained optimization problem

$$\max_{b_1} b_1^T S b_1 \quad (5)$$

$$\text{subject to } \|b_1\|^2 = 1 \quad (6)$$

which has the Lagrangian

$$\mathcal{L}(b_1, \lambda) = b_1^T S b_1 - \lambda_1 (1 - b_1^T b_1) \quad (7)$$

Setting the partial derivatives of  $\mathcal{L}$  to 0

$$\frac{\partial \mathcal{L}}{\partial b_1} = 2b_1^T S - 2\lambda_1 b_1^T = 0, \quad (8)$$

$$\frac{\partial \mathcal{L}}{\partial \lambda_1} = 1 - b_1^T b_1 = 0, \quad (9)$$

we obtain

$$Sb_1 = \lambda_1 b_1, \quad b_1^T b_1 = 1. \quad (10)$$

Hence  $b_1$  is an eigenvector corresponding to the eigenvalue  $\lambda_1$  of the data covariance matrix  $S$ , and

$$V_1 = b_1^T Sb_1 = \lambda_1. \quad (11)$$

$\lambda_1$  is called **the first principal component**.

# $M$ -dimensional Subspace

Assume that we have found the first  $m - 1$  principal components  $b_1, \dots, b_{m-1}$  as the  $m - 1$  eigenvectors of  $S$  that are associated with the largest  $m - 1$  eigenvalues  $\lambda_1, \dots, \lambda_{m-1}$ .

Since  $S$  is symmetric, these eigenvectors form an orthonormal basis of an  $(m - 1)$ -dimensional subspace of  $\mathbb{R}^D$ . We find principal components that compress the remaining information.

Set  $X := [x_1 \dots x_N] \in \mathbb{R}^{D \times N}$  and  $B_{m-1} := \sum_{i=1}^{m-1} b_i b_i^T$  the projection matrix that projects onto the subspace spanned by  $b_1, \dots, b_{m-1}$ . The new data matrix

$$\hat{X} := X - \sum_{i=1}^{m-1} b_i b_i^T X = X - B_{m-1} X. \quad (12)$$

Let  $\hat{S}$  be the data covariance matrix of the transformed dataset  $\hat{\mathcal{X}} := \{\hat{x}_1, \dots, \hat{x}_N\}$ .

To find the  $m$ th principal component, we solve a constrained optimization problem

$$\max_{b_m} \left[ V_m = V_m[z_m] = \frac{1}{N} \sum_{n=1}^N z_{mn}^2 = \frac{1}{N} \sum_{n=1}^N (b_m^T \hat{x}_n)^2 = b_m^T \hat{S} b_m \right] \quad (13)$$

subject to  $\|b_m\|^2 = 1$ .

The optimal solution  $b_m$  is the eigenvector of  $\hat{S}$  that is associated with the largest eigenvalue of  $\hat{S}$ .

We shall show that the sets of eigenvectors of  $S$  and  $\hat{S}$  are identical.

Consider an eigenvector  $b_i$  of  $S$ , i.e.,  $Sb_i = \lambda_i b_i$ . We have

$$\begin{aligned}\hat{S}b_i &= \frac{1}{N} \hat{X} \hat{X}^T b_i = \frac{1}{N} (X - B_{m-1}X)(X - B_{m-1}X)^T b_i \\ &= (S - SB_{m-1} - B_{m-1}S + B_{m-1}SB_{m-1})b_i.\end{aligned}\quad (14)$$

If  $i \geq m$ , then  $b_i$  is orthogonal to the first  $m-1$  principal components and  $B_{m-1}b_i = 0$ . Hence

$$\hat{S}b_i = (S - B_{m-1}S)b_i = Sb_i = \lambda_i b_i, \quad (15)$$

which shows that  $b_m$  is an eigenvector of  $\hat{S}$  and  $\lambda_m$  is the largest eigenvalue of  $\hat{S}$  and  $\lambda_m$  is the  $m$ th largest eigenvalue of  $S$ , and both have the associated eigenvector  $b_m$ .



If  $i < m$ , then  $b_i$  is a basis vector of the principal subspace onto which  $B_{m-1}$  projects. Since  $b_1, \dots, b_{m-1}$  are an ONB of this principal subspace, we obtain  $B_{m-1}b_i = b_i$ . Thus,

$$\hat{S}b_i = (S - SB_{m-1} - B_{m-1}S + B_{m-1}SB_{m-1})b_i = 0 = 0b_i. \quad (16)$$

This means that  $b_1, \dots, b_{m-1}$  are also eigenvectors of  $\hat{S}$ , but they are associated with eigenvalue 0 so that  $b_1, \dots, b_{m-1}$  span the null space of  $\hat{S}$ .

# $M$ -dimensional Subspace

By (15), the variance of the data projected onto the  $m$ th principal component is

$$V_m = b_m^T S b_m = \lambda_m b_m^T b_m = \lambda_m. \quad (17)$$

To find an  $M$ -dimensional subspace of  $\mathbb{R}^D$  that retains as much information as possible, PCA tells us to choose the columns of the matrix  $B$  in (3) as the  $M$  eigenvectors of the data covariance matrix  $S$  that are associated with the  $M$  largest eigenvalues. The maximum amount of variance PCA can capture with the first  $M$  principal components is

$$V_M = \sum_{m=1}^M \lambda_m \quad (18)$$

where the  $\lambda_m$  are the  $M$  largest eigenvalues of the data covariance matrix  $S$ . Consequently, the variance lost by data compression via PCA is

$$J_M := \sum_{j=M+1}^D \lambda_j = V_D - V_M. \quad (19)$$

**Setting and Objective** We will derive PCA as an algorithm that directly minimizes the average reconstruction error.

Assume that  $B = \{b_1, \dots, b_D\}$  is an orthonormal basis of  $\mathbb{R}^D$ . Any  $x \in \mathbb{R}^D$  can be written as

$$x = \sum_{d=1}^D \xi_d b_d = \sum_{m=1}^M \xi_m b_m + \sum_{j=M+1}^D \xi_j b_j \text{ where } \xi_d \in \mathbb{R}. \quad (20)$$

Find a vectors  $\tilde{x}$  in an  $M$ -dimensional subspace  $U$  of  $\mathbb{R}^D$  so that

$$\tilde{x} = \sum_{m=1}^M z_m b_m \quad (21)$$

is as similar to  $x$  as possible.

Without loss of generality, assume that the dataset  $\mathcal{X} = \{x_1, \dots, x_N\}$  is centered at 0,  $E[\mathcal{X}] = 0$ . We are interested in finding the best linear projection of  $\mathcal{X}$  onto  $U$  the **principal subspace** and ONB  $b_1, \dots, b_M$  of  $U$ . The projections of the data points

$$\tilde{x}_n := \sum_{m=1}^M z_{mn} b_m = B z_n \in \mathbb{R}^D. \quad (22)$$

The **reconstruction error**

$$J_M := \frac{1}{N} \sum_{n=1}^N \|x_n - \tilde{x}_n\|^2. \quad (23)$$

To find the coordinates  $z_n = [z_{1n}, \dots, z_{Mn}]^T \in \mathbb{R}^M$  of  $\tilde{x}_n$  w.r.t the basis  $(b_1, \dots, b_M)$  and the ONB of the principal subspace, we follow two-step. First, we optimize the coordinate  $z_n$  for a given ONB  $(b_1, \dots, b_M)$ ; second, we find the optimal ONB.

# Finding Optimal Coordinates

We have

$$\begin{aligned}\frac{\partial J_M}{\partial z_{in}} &= \frac{\partial J_M}{\partial \tilde{x}_n} \frac{\partial \tilde{x}_n}{\partial z_{in}}, \\ \frac{\partial J_M}{\partial \tilde{x}_n} &= -\frac{2}{N} (x_n - \tilde{x}_n)^T \in \mathbb{R}^{1 \times D} \\ \frac{\partial \tilde{x}_n}{\partial z_{in}} &= \frac{\partial}{\partial z_{in}} \left( \sum_{m=1}^M z_{mn} b_m \right) = b_i.\end{aligned}\tag{24}$$

Hence

$$\frac{\partial J_M}{\partial z_{in}} = -\frac{2}{N} (x_n - \tilde{x}_n)^T b_i = -\frac{2}{N} \left( x_n - \sum_{m=1}^M z_{mn} b_m \right)^T b_i \tag{25}$$

$$= -\frac{2}{N} (x_n^T b_i - z_{in} b_i^T b_i) = -\frac{2}{N} (x_n^T b_i - z_{in}). \tag{26}$$

Setting this to 0 obtain

$$z_{in} = x_n^T b_i = b_i^T x_n. \tag{27}$$

# Finding the Basis of the Principal Subspace

To determine the basis vectors  $b_1, \dots, b_M$  of the principal subspace, we rephrase the loss function (23) using the results we have so far. To reformulate the loss function, we exploit our results from before

$$\tilde{x}_n = \sum_{m=1}^M z_{mn} b_m = \sum_{m=1}^M (x_n^T b_m) b_m = \left( \sum_{m=1}^M b_m b_m^T \right) x_n. \quad (28)$$

$$\begin{aligned} x_n &= \sum_{d=1}^D z_{dn} b_d \stackrel{(27)}{=} \sum_{d=1}^D (x_n^T b_d) b_d = \left( \sum_{d=1}^D b_d b_d^T \right) x_n \\ &= \left( \sum_{m=1}^M b_m b_m^T \right) x_n + \left( \sum_{j=M+1}^D b_j b_j^T \right) x_n. \end{aligned} \quad (29)$$

Hence the difference vector between the original data point and its projection is

$$x_n - \tilde{x}_n = \left( \sum_{j=M+1}^D b_j b_j^T \right) x_n = \sum_{j=M+1}^D (x_n^T b_j) b_j. \quad (30)$$

Reformulate the loss function

$$J_M = \frac{1}{N} \sum_{n=1}^N \|x_n - \tilde{x}_n\|^2 \stackrel{(30)}{=} \frac{1}{N} \sum_{n=1}^N \left\| \sum_{j=M+1}^D (b_j^T x_n) b_j \right\|^2 \quad (31)$$

$$= \frac{1}{N} \sum_{n=1}^N \sum_{j=M+1}^D (b_j^T x_n)^2 = \frac{1}{N} \sum_{n=1}^N \sum_{j=M+1}^D b_j^T x_n b_j^T x_n \quad (32)$$

$$= \frac{1}{N} \sum_{n=1}^N \sum_{j=M+1}^D b_j^T x_n x_n^T b_j = \sum_{j=M+1}^D b_j^T \left( \frac{1}{N} \sum_{n=1}^N x_n x_n^T \right) b_j = \sum_{j=M+1}^D b_j^T S b_j \quad (33)$$

$$= \sum_{j=M+1}^D \text{tr}(b_j^T S b_j) = \sum_{j=M+1}^D \text{tr}(S b_j b_j^T) = \text{tr} \left( \left( \sum_{j=M+1}^D b_j b_j^T \right) S \right). \quad (34)$$

Hence the followings are equivalent

- Minimize the average squared reconstruction error.
- Minimize the variance of the data when projected onto the orthogonal complement of the principal subspace.
- Maximize the variance of the projection that we retain in the principal subspace which links the projection loss immediately to the maximum-variance formulation of PCA.

The average squared reconstruction error when projecting onto the  $M$ -dimensional principal subspace is

$$J_M = \sum_{j=M+1}^D \lambda_j, \quad (35)$$

where  $\lambda_i$  are the eigenvalues of the data covariance matrix. Therefore, to minimize (35) we need to select the smallest  $D - M$  eigenvalues, which then implies that their corresponding eigenvectors are the basis of the orthogonal complement of the principal subspace. Consequently, the basis of the principal subspace comprises the eigenvectors  $b_1, \dots, b_M$  that are associated with the largest  $M$  eigenvalues of the data covariance matrix.



## 10.4 Eigenvector Computation and Low-Rank Approximations

The SVD of data matrix  $X = [x_1 \ \cdots \ x_N] \in \mathbb{R}^{D \times N}$  is given by

$$X = U\Sigma V^T, \quad (36)$$

where  $U \in \mathbb{R}^{D \times D}$ ,  $V \in \mathbb{R}^{N \times N}$  are orthogonal matrices and  $\Sigma \in \mathbb{R}^{D \times N}$  is a matrix whose only nonzero entries are singular values  $\sigma_{ii} \geq 0$ . The data covariance matrix

$$S = \frac{1}{N}XX^T = \frac{1}{N}U\Sigma V^T V\Sigma^T U^T = \frac{1}{N}U\Sigma\Sigma^T U^T \quad (37)$$

The columns of  $U$  are the eigenvectors of  $XX^T$  (see Section 4.5).

Futhermore, the eigenvalues  $\lambda_d$  of  $S$  are related to the singular values of  $X$  via

$$\lambda_d = \frac{\sigma_d^2}{N} \quad (38)$$

This relation provides the connection between the maximum variance view and the singular value decomposition.

To maximize the variance of the projected data, PCA chooses the columns of  $U$  in (37) to be the eigenvectors that are associated with the  $M$  largest eigenvalues of the data covariance matrix  $S$  so that we identify  $U$  as the projection matrix  $B$  in (3), which projects the original data onto a lower-dimensional subspace of dimension  $M$ . The Eckart-Young theorem (Section 4.6) offers a direct way to estimate the low-dimensional representation. Consider the best rank- $M$  approximation of  $X$

$$\tilde{X}_M := \arg \min_{rk(A) \leq M} \|X - A\|_2 \in \mathbb{R}^{D \times N}, \quad (39)$$

where  $\|\cdot\|_2$  is the spectral norm defined in Section 4.6.

For large size matrices, this is not possible to find all eigenvalues and eigenvectors. In practice, we solve for eigenvalues or singular values using iterative methods, which are implemented in all modern packages for linear algebra.

In PCA, we only require a few eigenvectors. If we are interested in only the first few eigenvectors (with the largest eigenvalues), then iterative processes, which directly optimize these eigenvectors, are computationally more efficient than a full eigendecomposition (or SVD). In the extreme case of only needing the first eigenvector, a simple method called the **power iteration**: chooses a random vector  $x_0$  not in the null space of  $S$  and follows the iteration

$$x_{k+1} = \frac{Sx_k}{\|Sx_k\|} \quad (40)$$

This sequence of vectors converges to the eigenvector associated with the largest eigenvalue of  $S$ .

We will provide a method for PCA when fewer data points than dimensions  $N \ll D$ .

Assume we have a centered dataset  $x_1, \dots, x_N \in \mathbb{R}^D$  and there are no duplicate data points. So  $rk(S) = N$ ,  $S$  has  $D - N + 1$  many eigenvalues that are 0. We will exploit this and turn the  $D \times D$  covariance matrix into an  $N \times N$  covariance matrix whose eigenvalues are all positive.

Let  $(b_m, m = 1, \dots, M)$  is a basis vector of the principal subspace.

$$Sb_m = \frac{1}{N}XX^T b_m = \lambda_m b_m \quad (41)$$

$$\Rightarrow \frac{1}{N}X^T XX^T b_m = \lambda_m X^T b_m \Leftrightarrow \frac{1}{N}X^T Xc_m = \lambda_m c_m. \quad (42)$$

The eigenvector of the  $N \times N$  matrix  $\frac{1}{N}X^T X \in \mathbb{R}^{N \times N}$  associated with  $\lambda_m$  as  $c_m := X^T b_m$  ( $\frac{1}{N}X^T X$  and  $S$  have the same nonzero eigenvalues).

Recover the original eigenvectors via the eigenvectors of  $\frac{1}{N}X^T X$

$$\frac{1}{N}XX^T Xc_m = \lambda_m Xc_m, \text{ which means } Xc_m \text{ as an eigenvector of } S. \quad (43)$$

# 10.6 Key Steps of PCA in Practice

1. Mean subtraction
2. Standardization
3. Eigendecomposition of the covariance matrix
4. Projection

We have studied:

- Maximum variance perspective of PCA;
- PCA as an algorithm that directly minimizes the average reconstruction error;
- PCA using low-rank matrix approximations;
- PCA in high dimensions.