

Dimensionality Reduction with Principal Component Analysis

5209 Data in real life is often high dimensional. For example, if we want to esti-
 5210 mate the price of our house in a year's time, we can use data that helps us
 5211 to do this: the type of house, the size, the number of bedrooms and bath-
 5212 rooms, the value of houses in the neighborhood when they were bought,
 5213 the distance to the next train station and park, the number of crimes com-
 5214 mitted in the neighborhood, the economic climate etc. – there are many
 5215 things that influence the house price, and we collect this information in a
 5216 data set that we can use to estimate the house price. Another example is a
 5217 640×480 pixels color image, which is a data point in a million-dimensional
 5218 space, where every pixel responds to three dimensions - one for each color
 5219 channel (red, green, blue).

5220 Working directly with high-dimensional data comes with some difficul-
 5221 ties: It is hard to analyze, interpretation is difficult, visualization is nearly
 5222 impossible, and (from a practical point of view) storage can be expensive.
 5223 However, high-dimensional data also has some nice properties: For exam-
 5224 ple, high-dimensional data is often overcomplete, i.e., many dimensions
 5225 are redundant and can be explained by a combination of other dimen-
 5226 sions. Dimensionality reduction exploits structure and correlation and al-
 5227 lows us to work with a more compact representation of the data, ideally
 5228 without losing information. We can think of dimensionality reduction as
 5229 a compression technique, similar to jpg or mp3, which are compression
 5230 algorithms for images and music.

5231 In this chapter, we will discuss *principal component analysis* (PCA), an
 5232 algorithm for linear *dimensionality reduction*. PCA, proposed by Pearson
 5233 (1901) and Hotelling (1933), has been around for more than 100 years
 5234 and is still one of the most commonly used techniques for data compres-
 5235 sion, data visualization and the identification of simple patterns, latent
 5236 factors and structures of high-dimensional data. In this chapter, we will
 5237 explore the concept of linear dimensionality reduction with PCA in more
 5238 detail, drawing on our understanding of basis and basis change (see Sec-
 5239 tions 2.6.1 and 2.7.2), projections (see Section 3.6), eigenvalues (see Sec-
 5240 tion 4.2), Gaussian distributions (see Section 6.6) and constrained opti-
 5241 mization (see Section 7.2).

5242 Dimensionality reduction generally exploits the property of high-dimen-
 5243 sional data (e.g., images) that it often lies on a low-dimensional subspace,

principal component
analysis
dimensionality
reduction

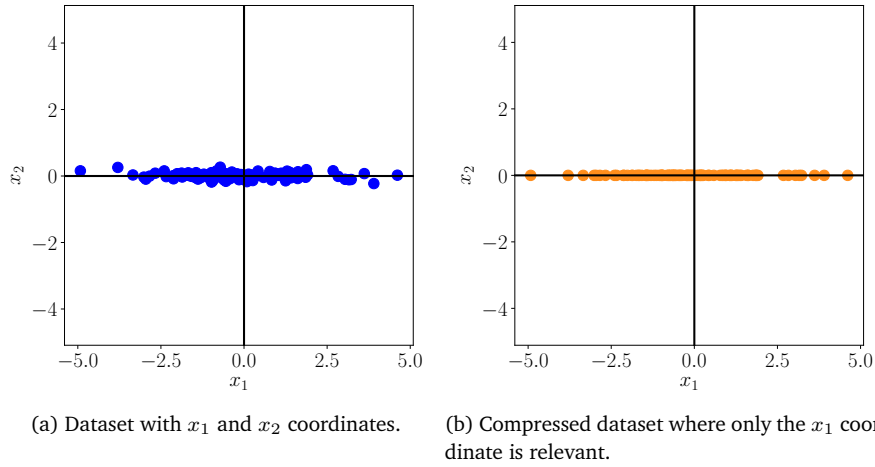
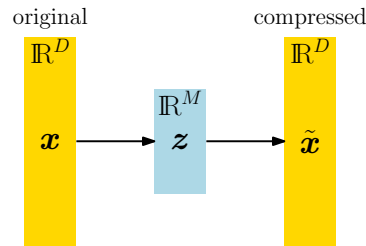


Figure 10.1
Illustration:
Dimensionality
reduction. (a) The
original dataset not
vary much along the
 x_2 direction. (b)
The data from (a)
can be represented
using the
 x_1 -coordinate alone
with nearly no loss.

Figure 10.2
Graphical
illustration of PCA.
In PCA, we find a
compressed version
 \tilde{x} of original data x
that has an intrinsic
lower-dimensional
representation z .



and that many dimensions are highly correlated, redundant or contain little information. Figure 10.1 gives an illustrative example in two dimensions. Although the data in Figure 10.1(a) does not quite lie on a line, the data does not vary much in the x_2 -direction, so that we can express it as if it was on a line – with nearly no loss, see Figure 10.1(b). The data in Figure 10.1(b) requires only the x_1 -coordinate to describe and lies in a one-dimensional subspace of \mathbb{R}^2 .

10.1 Problem Setting

In PCA, we are interested in finding projections \tilde{x}_n of data points x_n that are as similar to the original data points as possible, but which have a significantly lower intrinsic dimensionality. Figure 10.1 gives an illustration what this could look like.

Figure 10.2 illustrates the setting we consider in PCA, where z represents the intrinsic lower dimension of the compressed data \tilde{x} and plays the role of a bottleneck, which controls how much information can flow between x and \tilde{x} .

More concretely, we consider i.i.d. data points $x_1, \dots, x_N \in \mathbb{R}^D$, and we search a low-dimensional, compressed representation (code) z_n of x_n .



Figure 10.3
Examples of
handwritten digits
from the MNIST
dataset.

If our observed data lives in \mathbb{R}^D , we look for an M -dimensional subspace $U \subseteq \mathbb{R}^D$, $\dim(U) = M < D$ onto which we project data. We denote the projected data as $\tilde{x}_n \in U$, and their coordinates (with respect to an appropriate basis in U) with z_n . Our aim is to find \tilde{x}_n so that they are as similar to the original data x_n as possible.

Example 10.1 (Coordinate Representation/Code)

Consider \mathbb{R}^2 with the canonical basis $e_1 = [1, 0]^\top$, $e_2 = [0, 1]^\top$. From Chapter 2 we know that $x \in \mathbb{R}^2$ can be represented as a linear combination of these basis vectors, e.g.,

$$\begin{bmatrix} 5 \\ 3 \end{bmatrix} = 5e_1 + 3e_2. \quad (10.1)$$

However, when we consider the set of vectors

$$\tilde{x} = \begin{bmatrix} 0 \\ z \end{bmatrix} \in \mathbb{R}^2, \quad z \in \mathbb{R}, \quad (10.2)$$

they can always be written as $0e_1 + ze_2$. To represent these vectors it is sufficient to remember/store the *coordinate/code* z of the e_2 vector.

More precisely, the set of \tilde{x} vectors (with the standard vector addition and scalar multiplication) forms a vector subspace U (see Section 2.4) with $\dim(U) = 1$ because $U = \text{span}[e_2]$.

The dimension of a vector space corresponds to the number of its basis vectors (see Section 2.6.1).

In PCA, we consider the relationship between the original data x and its low-dimensional code z to be linear so that $z = B^\top x$ for a suitable matrix B .

Throughout this chapter, we will use the MNIST digits dataset as a re-occurring example, which contains 60,000 examples of handwritten digits 0–9. Each digit is an image of size 28×28 , i.e., it contains 784 pixels so that we can interpret every image in this dataset as a vector $x \in \mathbb{R}^{784}$. Examples of these digits are shown in Figure 10.3.

In the following, we will derive PCA from two different perspectives. First, we derive PCA by maintaining as much variance as possible in the projected space. Second, we will derive PCA by minimizing the average squared reconstruction error, which directly links to many concepts in Chapters 3 and 4.

<http://yann.lecun.com/exdb/mnist/>

10.2 Maximum Variance Perspective

Figure 10.1 gave an example of how a two-dimensional dataset can be represented using a single coordinate. In Figure 10.1(b), we chose to ignore the x_2 -coordinate of the data because it did not add too much information so that the compressed data is similar to the original data in Figure 10.1(a). We could have chosen to ignore the x_1 -coordinate, but then the compressed data had been very dissimilar from the original data, and much information in the data would have been lost.

If we interpret information content in the data as how “space filling” the data set is, then we can describe the information contained in the data by looking at the spread of the data. From Chapter 6 we know that the variance is an indicator of the spread of the data, and it is possible to formulate PCA as a dimensionality reduction algorithm that maximizes the variance in the low-dimensional representation of the data to retain as much information as possible. Now, let us formulate this objective more concretely.

Consider a dataset $\mathbf{x}_1, \dots, \mathbf{x}_N, \mathbf{x}_n \in \mathbb{R}^D$, with mean $\mathbf{0}$ that possesses the *data covariance matrix* (empirical covariance)

data covariance
matrix

$$\mathbf{S} = \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n \mathbf{x}_n^\top. \quad (10.3)$$

Furthermore, we assume a low-dimensional representation $\mathbf{z}_n = \mathbf{B}^\top \mathbf{x}_n \in \mathbb{R}^M$ of \mathbf{x}_n , where $\mathbf{B} \in \mathbb{R}^{D \times M}$.

Our aim is to find a matrix \mathbf{B} that retains as much information as possible when compressing data. We assume that \mathbf{B} is an orthogonal matrix so that $\mathbf{b}_i^\top \mathbf{b}_j = 0$ if and only if $i \neq j$. Retaining most information is formulated as capturing the largest amount of variance in the low-dimensional code (Hotelling, 1933).

The columns
 $\mathbf{b}_1, \dots, \mathbf{b}_M$ of \mathbf{B}
form a basis of the
 M -dimensional
subspace in which
the projected data
 $\hat{\mathbf{x}} = \mathbf{B}\mathbf{B}^\top \mathbf{x} \in \mathbb{R}^D$
live.

Remark. (Centered Data) Let us assume that $\boldsymbol{\mu} = \mathbb{E}_x[\mathbf{x}]$ is the (empirical) mean of the data. Using the properties of the variance, which we discussed in Section 6.4.4 we obtain

$$\mathbb{V}_z[\mathbf{z}] = \mathbb{V}_x[\mathbf{B}^\top (\mathbf{x} - \boldsymbol{\mu})] = \mathbb{V}_x[\mathbf{B}^\top \mathbf{x} - \mathbf{B}^\top \boldsymbol{\mu}] = \mathbb{V}_x[\mathbf{B}^\top \mathbf{x}], \quad (10.4)$$

i.e., the variance of the low-dimensional code does not depend on the mean of the data. Therefore, we assume without loss of generality that the data has mean $\mathbf{0}$ for the remainder of this section. With this assumption the mean of the low-dimensional code is also $\mathbf{0}$ since $\mathbb{E}_z[\mathbf{z}] = \mathbb{E}_x[\mathbf{B}^\top \mathbf{x}] = \mathbf{B}^\top \mathbb{E}_x[\mathbf{x}] = \mathbf{0}$. \diamond

We maximize the variance of the low-dimensional code using a sequential approach. We start by seeking a single vector $\mathbf{b}_1 \in \mathbb{R}^D$ that maximizes the variance of the projected data, i.e., we aim to maximize the first coor-

dinate z_1 of $\mathbf{z} \in \mathbb{R}^M$ so that

$$V_1 := \mathbb{V}[z_1] = \frac{1}{N} \sum_{n=1}^N z_{1n}^2 \quad (10.5)$$

is maximized, where we exploited the i.i.d. assumption of the data and defined z_{1n} as the first coordinate of the low-dimensional representation $\mathbf{z}_n \in \mathbb{R}^M$ of $\mathbf{x}_n \in \mathbb{R}^D$. Note that first component of \mathbf{z}_n is given by

$$z_{1n} = \mathbf{b}_1^\top \mathbf{x}_n. \quad (10.6)$$

We use this relationship now in (10.5), which yields

$$V_1 = \frac{1}{N} \sum_{n=1}^N (\mathbf{b}_1^\top \mathbf{x}_n)^2 = \frac{1}{N} \sum_{n=1}^N \mathbf{b}_1^\top \mathbf{x}_n \mathbf{x}_n^\top \mathbf{b}_1 \quad (10.7a)$$

$$= \mathbf{b}_1^\top \left(\frac{1}{N} \sum_{n=1}^N \mathbf{x}_n \mathbf{x}_n^\top \right) \mathbf{b}_1 = \mathbf{b}_1^\top \mathbf{S} \mathbf{b}_1, \quad (10.7b)$$

where \mathbf{S} is the data covariance matrix defined in (10.3).

It is clear that arbitrarily increasing the magnitude of the vector \mathbf{b}_1 increases V_1 . Therefore, we restrict all solutions to $\|\mathbf{b}_1\| = 1$, which results in a constrained optimization problem in which we seek the direction along which the data varies most.

With the restriction of the solution space to unit vectors we end up with the constrained optimization problem

$$\max_{\mathbf{b}_1} \mathbf{b}_1^\top \mathbf{S} \mathbf{b}_1 \quad (10.8)$$

$$\text{subject to } \|\mathbf{b}_1\|^2 = 1. \quad (10.9)$$

Following Section 7.2, we obtain the Lagrangian

$$\mathcal{L} = V_1 + \lambda_1(1 - \mathbf{b}_1^\top \mathbf{b}_1) = \mathbf{b}_1^\top \mathbf{S} \mathbf{b}_1 + \lambda_1(1 - \mathbf{b}_1^\top \mathbf{b}_1) \quad (10.10)$$

to solve this constrained optimization problem. The partial derivatives of \mathcal{L} with respect to \mathbf{b}_1 and λ_1 are

$$\frac{\partial \mathcal{L}}{\partial \mathbf{b}_1} = 2\mathbf{b}_1^\top \mathbf{S} - 2\lambda_1 \mathbf{b}_1^\top \quad (10.11)$$

$$\frac{\partial \mathcal{L}}{\partial \lambda_1} = 1 - \mathbf{b}_1^\top \mathbf{b}_1, \quad (10.12)$$

respectively. Setting these partial derivatives to $\mathbf{0}$ gives us the relations

$$\mathbf{b}_1^\top \mathbf{b}_1 = 1, \quad (10.13)$$

$$\mathbf{S} \mathbf{b}_1 = \lambda_1 \mathbf{b}_1, \quad (10.14)$$

i.e., we see that \mathbf{b}_1 is an eigenvector of the data covariance matrix \mathbf{S} , and

the Lagrange multiplier λ_1 plays the role of the corresponding eigenvalue. This eigenvector property allows us to rewrite our variance objective as

$$V_1 = \mathbf{b}_1^\top \mathbf{S} \mathbf{b}_1 = \lambda_1 \mathbf{b}_1^\top \mathbf{b}_1 = \lambda_1, \quad (10.15)$$

i.e., the variance of the data projected onto a one-dimensional subspace equals the eigenvalue that is associated with the basis vector \mathbf{b}_1 that spans this subspace. Therefore, to maximize the variance of the low-dimensional code we choose the basis vector belonging to the largest eigenvalue of the data covariance matrix. This eigenvector is called the first *principal component*. We can determine the effect/contribution of the principal component \mathbf{b}_1 in the original data space by mapping the coordinate z_{1n} back into data space, which gives us the projected data point

$$\tilde{\mathbf{x}}_n = \mathbf{b}_1 z_{1n} = \mathbf{b}_1 \mathbf{b}_1^\top \mathbf{x}_n \in \mathbb{R}^D \quad (10.16)$$

in the original data space.

Remark. Although $\tilde{\mathbf{x}}_n$ is a D -dimensional vector it only requires a single coordinate z_{1n} to represent it with respect to the basis vector $\mathbf{b}_1 \in \mathbb{R}^D$. \diamond

Generally, the m th principal component can be found by subtracting the effect of the first $m-1$ principal components from the data, thereby trying to find principal components that compress the remaining information. We achieve this by first subtracting the contribution of the $m-1$ principal components from the data, similar to (10.16), so that we arrive at the new data matrix

$$\hat{\mathbf{X}} := \mathbf{X} - \sum_{i=1}^{m-1} \mathbf{b}_i \mathbf{b}_i^\top \mathbf{X}, \quad (10.17)$$

where $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_N] \in \mathbb{R}^{D \times N}$ contains the data points as column vectors. The matrix $\hat{\mathbf{X}}$ in (10.17) contains the data that only contains the information that has not yet been compressed.

Remark (Notation). Throughout this chapter, we do not follow the convention of collecting data $\mathbf{x}_1, \dots, \mathbf{x}_N$ as rows of the data matrix, but we define them to be the columns of \mathbf{X} . This means that our data matrix \mathbf{X} is a $D \times N$ matrix instead of the conventional $N \times D$ matrix. The reason for our choice is that the algebra operations work out smoothly without the need to either transpose the matrix or to redefine vectors as row vectors that are left-multiplied onto matrices. \diamond

To find the m th principal component, we maximize the variance

$$V_m = \mathbb{V}[z_m] = \frac{1}{N} \sum_{n=1}^N z_{mn}^2 = \frac{1}{N} \sum_{n=1}^N \mathbf{b}_m^\top \mathbf{x}_n = \mathbf{b}_m^\top \hat{\mathbf{S}} \mathbf{b}_m, \quad (10.18)$$

subject to $\|\mathbf{b}_m\|^2 = 1$, where we followed the same steps as in (10.7b) and defined $\hat{\mathbf{S}}$ as the data covariance matrix of $\hat{\mathbf{X}}$. As previously, when we looked at the first principal component alone, we solve a constrained

principal component

The quantity $\sqrt{\lambda_1}$ is also called the *loading* of the unit vector \mathbf{b}_1 and represents the standard deviation of the data accounted for by the principal subspace $\text{span}[\mathbf{b}_1]$.

optimization problem and discover that the optimal solution \mathbf{b}_m is the eigenvector of $\hat{\mathbf{S}}$ that belongs to the largest eigenvalue of $\hat{\mathbf{S}}$.

However, it also turns out that \mathbf{b}_m is an eigenvector of \mathbf{S} . Since

$$\hat{\mathbf{S}} = \frac{1}{N} \sum_{n=1}^N \hat{\mathbf{x}}_n \hat{\mathbf{x}}_n^\top \stackrel{(10.17)}{=} \frac{1}{N} \sum_{n=1}^N \left(\mathbf{x}_n - \sum_{i=1}^{m-1} \mathbf{b}_i \mathbf{b}_i^\top \mathbf{x}_n \right) \left(\mathbf{x}_n - \sum_{i=1}^{m-1} \mathbf{b}_i \mathbf{b}_i^\top \mathbf{x}_n \right)^\top \quad (10.19a)$$

$$= \frac{1}{N} \sum_{n=1}^N \left(\mathbf{x}_n \mathbf{x}_n^\top - 2 \mathbf{x}_n \mathbf{x}_n^\top \sum_{i=1}^{m-1} \mathbf{b}_i \mathbf{b}_i^\top + \mathbf{x}_n^\top \sum_{i=1}^{m-1} \mathbf{b}_i \mathbf{b}_i^\top \sum_{i=1}^{m-1} \mathbf{b}_i \mathbf{b}_i^\top \right) \quad (10.19b)$$

we can multiply \mathbf{b}_m onto $\hat{\mathbf{S}}$ and obtain

$$\hat{\mathbf{S}} \mathbf{b}_m = \frac{1}{N} \sum_{n=1}^N \hat{\mathbf{x}}_n \hat{\mathbf{x}}_n^\top \mathbf{b}_m = \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n \mathbf{x}_n^\top \mathbf{b}_m = \mathbf{S} \mathbf{b}_m = \lambda_m \mathbf{b}_m. \quad (10.20)$$

Here we applied the orthogonality conditions $\mathbf{b}_i^\top \mathbf{b}_m = 0$ for all $i = 1, \dots, m-1$ (all terms involving sums up to $m-1$ vanish). In the end, we exploited the fact that \mathbf{b}_m is an eigenvector of $\hat{\mathbf{S}}$. Therefore, \mathbf{b}_m is also an eigenvector of the original data covariance matrix \mathbf{S} , and the corresponding eigenvalue is λ_m is the m th largest eigenvalue of \mathbf{S} . Moreover, the variance of the data projected onto the m th principal component

$$V_m = \mathbf{b}_m^\top \mathbf{S} \mathbf{b}_m \stackrel{(10.20)}{=} \lambda_m \mathbf{b}_m^\top \mathbf{b}_m = \lambda_m \quad (10.21)$$

since $\mathbf{b}_m^\top \mathbf{b}_m = 1$. This means that the variance of the data, when projected onto an M -dimensional subspace, equals the sum of the eigenvalues that belong to the corresponding eigenvectors of the data covariance matrix.

In practice, we do not have to compute principal components sequentially, but we can compute all of them at the same time. If we are looking for a projection onto an M -dimensional subspace so that as much variance as possible is retained in the projection, then PCA tells us to choose the columns of \mathbf{B} to be the eigenvectors that belong to the M largest eigenvalues of the data covariance matrix. The maximum amount of variance PCA can capture with the first M principal components is

To maximize the variance of the projected data, we choose the columns of \mathbf{B} to be the eigenvectors that belong to the M largest eigenvalues of the data covariance matrix.

$$V = \sum_{m=1}^M \lambda_m, \quad (10.22)$$

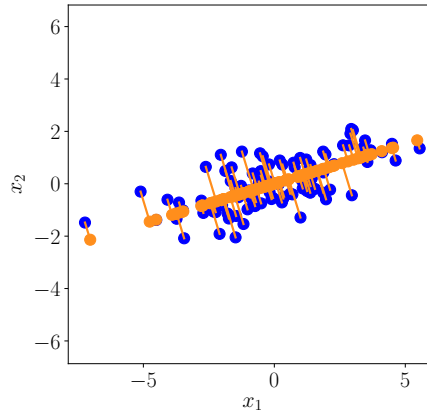
where the λ_m are the M largest eigenvalues of the data covariance matrix \mathbf{S} . Consequently, the variance lost by data compression via PCA is

$$J = \sum_{j=M+1}^D \lambda_j. \quad (10.23)$$

To summarize, to determine the M -dimensional subspace for which an

Figure 10.4

Illustration of the projection approach to PCA. We aim to find a one-dimensional subspace (line) of \mathbb{R}^2 so that the distance vector between projected (orange) and original (blue) data is as small as possible.



orthogonal projection maximizes the variance of the data we need to compute the M eigenvectors that belong to the M largest eigenvalues of the data covariance matrix. In Section 10.4, we will return to this point and discuss how to efficiently compute these M eigenvectors.

10.3 Projection Perspective

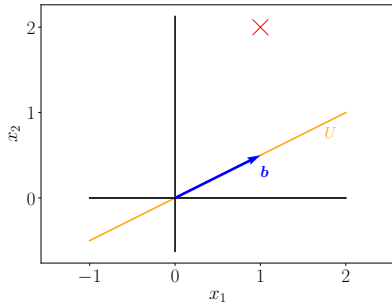
In the following, we will derive PCA as an algorithm for linear dimensionality reduction that minimizes the average projection error. We will draw heavily from Chapters 2 and 3. In the previous section, we derived PCA by maximizing the variance in the projected space to retain as much information as possible. In the following, we will look at the difference vectors between the original data \mathbf{x}_n and their reconstruction $\tilde{\mathbf{x}}_n$ and minimize this distance so that \mathbf{x}_n and $\tilde{\mathbf{x}}_n$ are as close as possible. Figure 10.4 illustrates this setting.

10.3.1 Setting and Objective

Assume an (ordered) orthonormal basis (ONB) $B = (\mathbf{b}_1, \dots, \mathbf{b}_D)$ of \mathbb{R}^D , i.e., $\mathbf{b}_i^\top \mathbf{b}_j = 1$ if and only if $i = j$ and 0 otherwise. From Section 2.5 we know that every $\mathbf{x} \in \mathbb{R}^D$ can be written as a linear combination of the basis vectors of \mathbb{R}^D , i.e.,

$$\mathbf{x} = \sum_{d=1}^D z_d \mathbf{b}_d \quad (10.24)$$

for $z_d \in \mathbb{R}$. We are interested in finding vectors $\tilde{\mathbf{x}} \in \mathbb{R}^D$, which live in lower-dimensional subspace U of \mathbb{R}^D , so that $\tilde{\mathbf{x}}$ is as similar to \mathbf{x} as possible. As $\tilde{\mathbf{x}} \in U \subseteq \mathbb{R}^D$, we can also express $\tilde{\mathbf{x}}$ as a linear combination



(a) Setting.

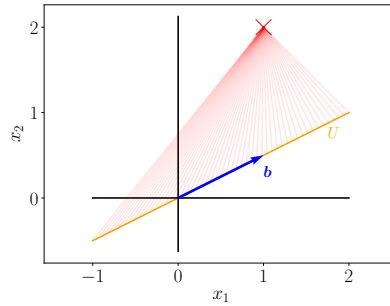
(b) Differences $\mathbf{x} - \tilde{\mathbf{x}}$ for 50 candidates $\tilde{\mathbf{x}}$ are shown by the red lines.

Figure 10.5
Simplified projection setting. (a) A vector $\mathbf{x} \in \mathbb{R}^2$ (red cross) shall be projected onto a one-dimensional subspace $U \subseteq \mathbb{R}^2$ spanned by \mathbf{b} . (b) shows the difference vectors between \mathbf{x} and some candidates $\tilde{\mathbf{x}}$.

of the basis vectors of \mathbb{R}^D so that

$$\tilde{\mathbf{x}} = \sum_{d=1}^D z_d \mathbf{b}_d. \quad (10.25)$$

Let us assume $\dim(U) = M$ where $M < D = \dim(\mathbb{R}^D)$. Then, we can find basis vectors $\mathbf{b}_1, \dots, \mathbf{b}_D$ of \mathbb{R}^D so that at least $D - M$ of the coefficients z_d are equal to 0, and we can rearrange the way we index the basis vectors \mathbf{b}_d such that the coefficients that are zero appear at the end. This allows us to express $\tilde{\mathbf{x}}$ as

$$\tilde{\mathbf{x}} = \sum_{m=1}^M z_m \mathbf{b}_m + \sum_{j=M+1}^D 0 \mathbf{b}_j = \sum_{m=1}^M z_m \mathbf{b}_m = \mathbf{B} \mathbf{z} \in \mathbb{R}^D, \quad (10.26)$$

where we defined

$$\mathbf{B} := [\mathbf{b}_1, \dots, \mathbf{b}_M] \in \mathbb{R}^{D \times M}, \quad (10.27)$$

$$\mathbf{z} := [z_1, \dots, z_M]^\top \in \mathbb{R}^M. \quad (10.28)$$

For example, vectors $\tilde{\mathbf{x}} \in U$ could be vectors on a plane in \mathbb{R}^3 . The dimensionality of the plane is 2, but the vectors still have three coordinates in \mathbb{R}^3 .

In the following, we use exactly this kind of representation of $\tilde{\mathbf{x}}$ to find optimal coordinates \mathbf{z} and basis vectors $\mathbf{b}_1, \dots, \mathbf{b}_M$ such that $\tilde{\mathbf{x}}$ is as similar to the original data point \mathbf{x} , i.e., we aim to minimize the (Euclidean) distance $\|\mathbf{x} - \tilde{\mathbf{x}}\|$. Figure 10.5 illustrates this setting.

Without loss of generality, we assume that the dataset $\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$, $\mathbf{x}_n \in \mathbb{R}^D$, is centered at $\mathbf{0}$, i.e., $\mathbb{E}[\mathbf{X}] = \mathbf{0}$.

Remark. Without the zero-mean assumption, we would arrive at exactly the same solution but the notation would be substantially more cluttered.

◇

We are interested in finding the best linear projection of \mathbf{X} onto a lower-dimensional subspace U of \mathbb{R}^D with $\dim(U) = M$ and orthonormal basis vectors $\mathbf{b}_1, \dots, \mathbf{b}_M$. We will call this subspace U the *principal subspace*, and $(\mathbf{b}_1, \dots, \mathbf{b}_M)$ is an orthonormal basis of the principal sub-

principal subspace

space. The projections are denoted by

$$\tilde{\mathbf{x}}_n := \sum_{m=1}^M z_{mn} \mathbf{b}_m = \mathbf{B} \mathbf{z}_n \in \mathbb{R}^D, \quad (10.29)$$

where $\mathbf{B} \in \mathbb{R}^{D \times M}$ is given in (10.27) and

$$\mathbf{z}_n := [z_{1n}, \dots, z_{Mn}]^\top \in \mathbb{R}^M, \quad n = 1, \dots, N, \quad (10.30)$$

is the coordinate vector of $\tilde{\mathbf{x}}_n$ with respect to the basis $(\mathbf{b}_1, \dots, \mathbf{b}_M)$. More specifically, we are interested in having the $\tilde{\mathbf{x}}_n$ as similar to \mathbf{x}_n as possible. There are many ways to measure similarity.

The similarity measure we use in the following is the squared Euclidean norm $\|\mathbf{x} - \tilde{\mathbf{x}}\|^2$ between \mathbf{x} and $\tilde{\mathbf{x}}$. We therefore define our objective as the minimizing the average squared Euclidean distance (*reconstruction error*) (Pearson, 1901)

$$J := \frac{1}{N} \sum_{n=1}^N \|\mathbf{x}_n - \tilde{\mathbf{x}}_n\|^2. \quad (10.31)$$

In order to find this optimal linear projection, we need to find the orthonormal basis of the principal subspace and the coordinates \mathbf{z}_n of the projections with respect to these basis vectors. All these parameters enter our objective (10.31) through $\tilde{\mathbf{x}}_n$.

In order to find the coordinates \mathbf{z}_n and the ONB of the principal subspace we optimize J by computing the partial derivatives of J with respect to all parameters of interest (i.e., the coordinates and the basis vectors), setting them to $\mathbf{0}$, and solving for the parameters. We detail these steps next. We will first determine the optimal coordinates z_{in} and then the basis vectors $\mathbf{b}_1, \dots, \mathbf{b}_M$ of the principal subspace, i.e., the subspace in which $\tilde{\mathbf{x}}$ lives.

10.3.2 Optimization

Since the parameters we are interested in, i.e., the basis vectors \mathbf{b}_i and the coordinates z_{in} of the projection with respect to the basis of the principal subspace, only enter the objective J through $\tilde{\mathbf{x}}_n$, we obtain

$$\frac{\partial J}{\partial z_{in}} = \frac{\partial J}{\partial \tilde{\mathbf{x}}_n} \frac{\partial \tilde{\mathbf{x}}_n}{\partial z_{in}}, \quad (10.32)$$

$$\frac{\partial J}{\partial \mathbf{b}_i} = \frac{\partial J}{\partial \tilde{\mathbf{x}}_n} \frac{\partial \tilde{\mathbf{x}}_n}{\partial \mathbf{b}_i} \quad (10.33)$$

for $i = 1, \dots, M$ and $n = 1, \dots, N$, where

$$\frac{\partial J}{\partial \tilde{\mathbf{x}}_n} = -\frac{2}{N} (\mathbf{x}_n - \tilde{\mathbf{x}}_n)^\top \in \mathbb{R}^{1 \times D}. \quad (10.34)$$

In the following, we determine the optimal coordinates z_{in} first before finding the ONB of the principal subspace.

Coordinates

Let us start by finding the coordinates z_{1n}, \dots, z_{Mn} of the projections $\tilde{\mathbf{x}}_n$ for $n = 1, \dots, N$. We assume that $(\mathbf{b}_1, \dots, \mathbf{b}_D)$ is an ordered ONB of \mathbb{R}^D . From (10.32) we require the partial derivative

$$\frac{\partial \tilde{\mathbf{x}}_n}{\partial z_{in}} \stackrel{(10.29)}{=} \frac{\partial}{\partial z_{in}} \left(\sum_{m=1}^M z_{mn} \mathbf{b}_m \right) = \mathbf{b}_i \quad (10.35)$$

for $i = 1, \dots, M$, such that we obtain

$$\frac{\partial J}{\partial z_{in}} \stackrel{(10.34)}{=} \stackrel{(10.35)}{=} -\frac{2}{N} (\mathbf{x}_n - \tilde{\mathbf{x}}_n)^\top \mathbf{b}_i \stackrel{(10.29)}{=} -\frac{2}{N} \left(\mathbf{x}_n - \sum_{m=1}^M z_{mn} \mathbf{b}_m \right)^\top \mathbf{b}_i \quad (10.36)$$

$$\stackrel{\text{ONB}}{=} -\frac{2}{N} (\mathbf{x}_n^\top \mathbf{b}_i - z_{in} \underbrace{\mathbf{b}_i^\top \mathbf{b}_i}_{=1}) = -\frac{2}{N} (\mathbf{x}_n^\top \mathbf{b}_i - z_{in}). \quad (10.37)$$

Setting this partial derivative to 0 yields immediately the optimal coordinates

$$z_{in} = \mathbf{x}_n^\top \mathbf{b}_i = \mathbf{b}_i^\top \mathbf{x}_n \quad (10.38)$$

for $i = 1, \dots, M$ and $n = 1, \dots, N$. This means, the optimal coordinates z_{in} of the projection $\tilde{\mathbf{x}}_n$ are the coordinates of the orthogonal projection (see Section 3.6) of the original data point \mathbf{x}_n onto the one-dimensional subspace that is spanned by \mathbf{b}_i . Consequently:

- The optimal projection $\tilde{\mathbf{x}}_n$ of \mathbf{x}_n is an orthogonal projection.
- The coordinates of $\tilde{\mathbf{x}}_n$ with respect to the basis $\mathbf{b}_1, \dots, \mathbf{b}_M$ are the coordinates of the orthogonal projection of \mathbf{x}_n onto the principal subspace.
- An orthogonal projection is the best linear mapping we can find given the objective (10.31).

Remark (Orthogonal Projections with Orthonormal Basis Vectors). Let us briefly recap orthogonal projections from Section 3.6. If $(\mathbf{b}_1, \dots, \mathbf{b}_D)$ is an orthonormal basis of \mathbb{R}^D then

$$\tilde{\mathbf{x}} = \mathbf{b}_j \underbrace{(\mathbf{b}_j^\top \mathbf{b}_j)^{-1}}_{=1} \mathbf{b}_j^\top \mathbf{x} = \mathbf{b}_j \mathbf{b}_j^\top \mathbf{x} \in \mathbb{R}^D \quad (10.39)$$

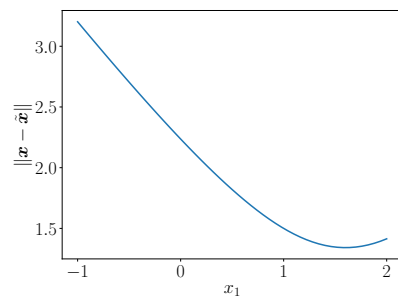
is the orthogonal projection of \mathbf{x} onto the subspace spanned by the j th basis vector, and $z_j = \mathbf{b}_j^\top \mathbf{x}$ is the coordinate of this projection with respect to the basis vector \mathbf{b}_j that spans that subspace since $z_j \mathbf{b}_j = \tilde{\mathbf{x}}$. Figure 10.6 illustrates this setting.

More generally, if we aim to project onto an M -dimensional subspace of \mathbb{R}^D , we obtain the orthogonal projection of \mathbf{x} onto the M -dimensional subspace with orthonormal basis vectors $\mathbf{b}_1, \dots, \mathbf{b}_M$ as

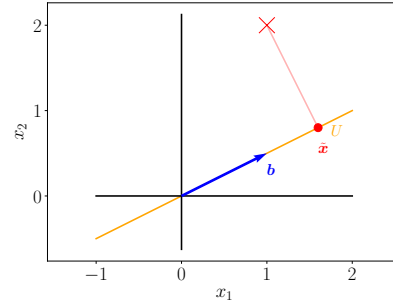
$$\tilde{\mathbf{x}} = \mathbf{B} \underbrace{(\mathbf{B}^\top \mathbf{B})^{-1}}_{=I} \mathbf{B}^\top \mathbf{x} = \mathbf{B} \mathbf{B}^\top \mathbf{x}, \quad (10.40)$$

The coordinates of the optimal projection of \mathbf{x}_n with respect to the basis vectors $\mathbf{b}_1, \dots, \mathbf{b}_M$ are the coordinates of the orthogonal projection of \mathbf{x}_n onto the principal subspace.

$\mathbf{x}^\top \mathbf{b}_j$ is the coordinate of the orthogonal projection of \mathbf{x} onto the one-dimensional subspace spanned by \mathbf{b}_j .



(a) Distances $\|\mathbf{x} - \tilde{\mathbf{x}}\|$ for some $\tilde{\mathbf{x}} \in U$, see panel (b) for the setting.



(b) The vector $\tilde{\mathbf{x}}$ that minimizes the distance in panel (a) is its orthogonal projection onto U . The coordinate of the projection $\tilde{\mathbf{x}}$ with respect to the basis vector \mathbf{b} that spans U is the factor we need to scale \mathbf{b} in order to “reach” $\tilde{\mathbf{x}}$.

Figure 10.6
Optimal projection of a vector $\mathbf{x} \in \mathbb{R}^2$ onto a one-dimensional subspace (continuation from Figure 10.5). (a) Distances $\|\mathbf{x} - \tilde{\mathbf{x}}\|$ for some $\tilde{\mathbf{x}} \in U$. (b) Orthogonal projection and optimal coordinates.

where we defined $\mathbf{B} := [\mathbf{b}_1, \dots, \mathbf{b}_M] \in \mathbb{R}^{D \times M}$. The coordinates of this projection with respect to the ordered basis $(\mathbf{b}_1, \dots, \mathbf{b}_M)$ are $\mathbf{z} := \mathbf{B}^\top \mathbf{x}$ as discussed in Section 3.6.

We can think of the coordinates as a representation of the projected vector in a new coordinate system defined by $(\mathbf{b}_1, \dots, \mathbf{b}_M)$. Note that although $\tilde{\mathbf{x}} \in \mathbb{R}^D$ we only need M coordinates z_1, \dots, z_M to represent this vector; the other $D - M$ coordinates with respect to the basis vectors $(\mathbf{b}_{M+1}, \dots, \mathbf{b}_D)$ are always 0. \diamond

Basis of the Principal Subspace

We already determined the optimal coordinates of the projected data for a given ONB $(\mathbf{b}_1, \dots, \mathbf{b}_D)$ of \mathbb{R}^D , only M of which were non-zero. What remains is to determine the basis vectors that span the principal subspace.

Before we get started, let us briefly introduce the concept of an orthogonal complement.

orthogonal
complement

Remark. (Orthogonal Complement) Consider a D -dimensional vector space V and an M -dimensional subspace $U \subseteq V$. Then its *orthogonal complement* U^\perp is a $(D - M)$ -dimensional subspace of V and contains all vectors in V that are orthogonal to every vector in U . Furthermore, every vector $\mathbf{x} \in V$ can be (uniquely) decomposed into

$$\mathbf{x} = \sum_{m=1}^M \lambda_m \mathbf{b}_m + \sum_{j=1}^{D-M} \psi_j \mathbf{b}_j^\perp, \quad \lambda_i, \psi_j \in \mathbb{R}, \quad (10.41)$$

where $(\mathbf{b}_1, \dots, \mathbf{b}_M)$ is a basis of U and $(\mathbf{b}_1^\perp, \dots, \mathbf{b}_{D-M}^\perp)$ is a basis of U^\perp . \diamond

To determine the basis vectors $\mathbf{b}_1, \dots, \mathbf{b}_M$ of the principal subspace, we rephrase the loss function (10.31) using the results we have so far.

This will make it easier to find the basis vectors. To reformulate the loss function, we exploit our results from before and obtain

$$\tilde{\mathbf{x}}_n = \sum_{m=1}^M z_{mn} \mathbf{b}_m \stackrel{(10.38)}{=} \sum_{m=1}^M (\mathbf{x}_n^\top \mathbf{b}_m) \mathbf{b}_m. \quad (10.42)$$

We now exploit the symmetry of the dot product, which yields

$$\tilde{\mathbf{x}}_n = \left(\sum_{m=1}^M \mathbf{b}_m \mathbf{b}_m^\top \right) \mathbf{x}_n. \quad (10.43)$$

Since we can generally write the original data point \mathbf{x}_n as a linear combination of all basis vectors, we can also write

$$\mathbf{x}_n = \sum_{d=1}^D z_{dn} \mathbf{b}_d \stackrel{(10.38)}{=} \sum_{d=1}^D (\mathbf{x}_n^\top \mathbf{b}_d) \mathbf{b}_d = \left(\sum_{d=1}^D \mathbf{b}_d \mathbf{b}_d^\top \right) \mathbf{x}_n \quad (10.44a)$$

$$= \left(\sum_{m=1}^M \mathbf{b}_m \mathbf{b}_m^\top \right) \mathbf{x}_n + \left(\sum_{j=M+1}^D \mathbf{b}_j \mathbf{b}_j^\top \right) \mathbf{x}_n, \quad (10.44b)$$

where we split the sum with D terms into a sum over M and a sum over $D - M$ terms. With this result, we find that the displacement vector $\mathbf{x}_n - \tilde{\mathbf{x}}_n$, i.e., the difference vector between the original data point and its projection, is

$$\mathbf{x}_n - \tilde{\mathbf{x}}_n = \left(\sum_{j=M+1}^D \mathbf{b}_j \mathbf{b}_j^\top \right) \mathbf{x}_n \quad (10.45)$$

$$= \sum_{j=M+1}^D (\mathbf{x}_n^\top \mathbf{b}_j) \mathbf{b}_j. \quad (10.46)$$

5405 This means the difference is exactly the projection of the data point onto
 5406 the orthogonal complement of the principal subspace: We identify the ma-
 5407 trix $\sum_{j=M+1}^D \mathbf{b}_j \mathbf{b}_j^\top$ in (10.45) as the projection matrix that performs this
 5408 projection. This also means the displacement vector $\mathbf{x}_n - \tilde{\mathbf{x}}_n$ lies in the
 5409 subspace that is orthogonal to the principal subspace as illustrated in Fig-
 5410 ure 10.7.

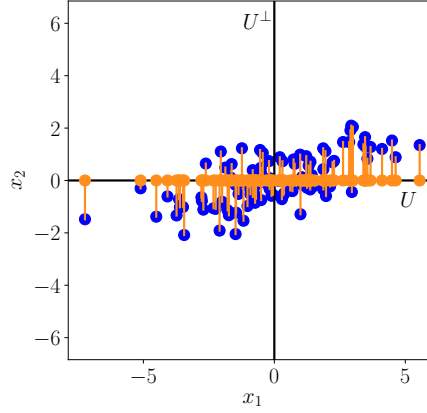
Remark (Low-Rank Approximation). In (10.45), we saw that the projec-
 tion matrix, which projects \mathbf{x} onto $\tilde{\mathbf{x}}$ is given by

$$\sum_{m=1}^M \mathbf{b}_m \mathbf{b}_m^\top = \mathbf{B} \mathbf{B}^\top. \quad (10.47)$$

By construction as a sum of rank-one matrices $\mathbf{b}_m \mathbf{b}_m^\top$ we see that $\mathbf{B} \mathbf{B}^\top$
 is symmetric and has rank M . Therefore, the average reconstruction error

PCA finds the best
 rank- M
 approximation of
 the identity matrix.

Figure 10.7
Orthogonal
projection and
displacement
vectors. When
projecting data
points \mathbf{x}_n (blue)
onto subspace U_1
we obtain $\tilde{\mathbf{x}}_n$
(orange). The
displacement vector
 $\tilde{\mathbf{x}}_n - \mathbf{x}_n$ lies
completely in the
orthogonal
complement U_2 of
 U_1 .



can also be written as

$$\sum_{n=1}^N \|\mathbf{x}_n - \tilde{\mathbf{x}}_n\|^2 = \sum_{n=1}^N \|\mathbf{x}_n - \mathbf{B}\mathbf{B}^\top \mathbf{x}_n\|^2 = \sum_{n=1}^N \|(I - \mathbf{B}\mathbf{B}^\top) \mathbf{x}_n\|^2. \quad (10.48)$$

5411 Finding orthonormal basis vectors $\mathbf{b}_1, \dots, \mathbf{b}_M$ so that the difference be-
5412 tween the original data \mathbf{x}_n and their projections $\tilde{\mathbf{x}}_n$, $n = 1, \dots, N$, is
5413 minimized is equivalent to finding the best rank- M approximation $\mathbf{B}\mathbf{B}^\top$
5414 of the identity matrix \mathbf{I} , see Section 4.6. \diamond

Now, we have all the tools to reformulate the loss function (10.31).

$$J = \frac{1}{N} \sum_{n=1}^N \|\mathbf{x}_n - \tilde{\mathbf{x}}_n\|^2 \stackrel{(10.46)}{=} \frac{1}{N} \sum_{n=1}^N \left\| \sum_{j=M+1}^D (\mathbf{b}_j^\top \mathbf{x}_n) \mathbf{b}_j \right\|^2. \quad (10.49)$$

We now explicitly compute the squared norm and exploit the fact that the \mathbf{b}_j form an ONB, which yields

$$J = \frac{1}{N} \sum_{n=1}^N \sum_{j=M+1}^D (\mathbf{b}_j^\top \mathbf{x}_n)^2 = \frac{1}{N} \sum_{n=1}^N \sum_{j=M+1}^D \mathbf{b}_j^\top \mathbf{x}_n \mathbf{b}_j^\top \mathbf{x}_n \quad (10.50a)$$

$$= \frac{1}{N} \sum_{n=1}^N \sum_{j=M+1}^D \mathbf{b}_j^\top \mathbf{x}_n \mathbf{x}_n^\top \mathbf{b}_j, \quad (10.50b)$$

where we exploited the symmetry of the dot product in the last step to write $\mathbf{b}_j^\top \mathbf{x}_n = \mathbf{x}_n^\top \mathbf{b}_j$. We can now swap the sums and obtain

$$J = \sum_{j=M+1}^D \mathbf{b}_j^\top \underbrace{\left(\frac{1}{N} \sum_{n=1}^N \mathbf{x}_n \mathbf{x}_n^\top \right)}_{=: \mathbf{S}} \mathbf{b}_j = \sum_{j=M+1}^D \mathbf{b}_j^\top \mathbf{S} \mathbf{b}_j \quad (10.51a)$$

$$= \sum_{j=M+1}^D \text{tr}(\mathbf{b}_j^\top \mathbf{S} \mathbf{b}_j) \sum_{j=M+1}^D \text{tr}(\mathbf{S} \mathbf{b}_j \mathbf{b}_j^\top) = \text{tr} \left(\underbrace{\left(\sum_{j=M+1}^D \mathbf{b}_j \mathbf{b}_j^\top \right)}_{\text{projection matrix}} \mathbf{S} \right), \quad (10.51b)$$

where we exploited the property that the trace operator $\text{tr}(\cdot)$, see (4.16), is linear and invariant to cyclic permutations of its arguments. Since we assumed that our dataset is centered, i.e., $\mathbb{E}[\mathbf{X}] = \mathbf{0}$, we identify \mathbf{S} as the data covariance matrix. We see that the projection matrix in (10.51b) is constructed as a sum of rank-one matrices $\mathbf{b}_j \mathbf{b}_j^\top$ so that it itself is of rank $D - M$.

Equation (10.51a) implies that we can formulate the average squared reconstruction error equivalently as the covariance matrix of the data, projected onto the orthogonal complement of the principal subspace. Minimizing the average squared reconstruction error is therefore equivalent to minimizing the variance of the data when projected onto the subspace we ignore, i.e., the orthogonal complement of the principal subspace. Equivalently, we 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 discussed in Section 10.2. But this then also means that we will obtain the same solution that we obtained for the maximum-variance perspective. Therefore, we skip the slightly lengthy derivation here and summarize the results from earlier in the light of the projection perspective.

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

$$J = \sum_{j=M+1}^D \lambda_j, \quad (10.52)$$

where λ_j are the eigenvalues of the data covariance matrix. Therefore, to minimize (10.52) 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, this means that the basis of the principal subspace are the eigenvectors $\mathbf{b}_1, \dots, \mathbf{b}_M$ that belong to the largest M eigenvalues of the data covariance matrix.

Minimizing the average squared reconstruction error is equivalent to minimizing the projection of the data covariance matrix onto the orthogonal complement of the principal subspace.

Minimizing the average squared reconstruction error is equivalent to maximizing the variance of the projected data.

Example 10.2 (MNIST Digits Embedding)

Figure 10.8
Embedding of
MNIST digits 0
(blue) and 1
(orange) in a
two-dimensional
principal subspace
using PCA. Four
examples
embeddings of the
digits ‘0’ and ‘1’ in
the principal
subspace are
highlighted in red
with their
corresponding
original digit.

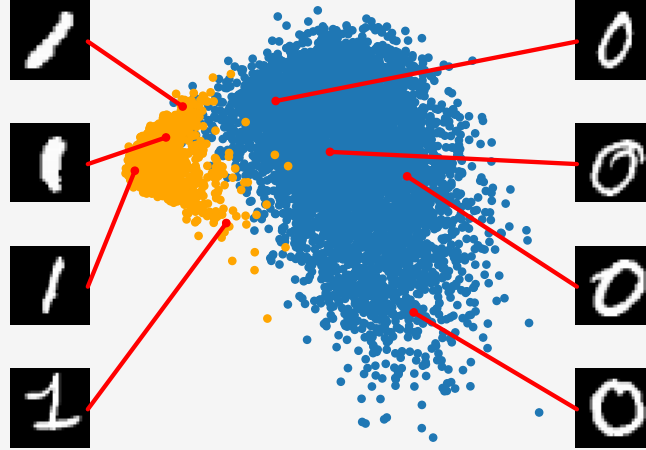


Figure 10.8 visualizes the training data of the MNIST digits ‘0’ and ‘1’ embedded in the vector subspace spanned by the first two principal components. We can see a relatively clear separation between ‘0’s (blue dots) and ‘1’s (orange dots), and we can see the variation within each individual cluster.

10.4 Eigenvector Computation

In the previous sections, we obtained the basis of the principal subspace as the eigenvectors that belong to the largest eigenvalues of the data covariance matrix

$$\mathbf{S} = \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n \mathbf{x}_n^\top = \frac{1}{N} \mathbf{X} \mathbf{X}^\top, \quad (10.53)$$

$$\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_N] \in \mathbb{R}^{D \times N}. \quad (10.54)$$

To get the eigenvalues (and the corresponding eigenvectors) of \mathbf{S} , we can follow two approaches:

- We perform an eigendecomposition (see Section 4.2) and compute the eigenvalues and eigenvectors of \mathbf{S} directly.
- We use a singular value decomposition (see Section 4.5). Since \mathbf{S} is symmetric and factorizes into $\mathbf{X} \mathbf{X}^\top$ (ignoring the factor $\frac{1}{N}$), the eigenvalues of \mathbf{S} are the squared singular values of \mathbf{X} . More specifically, if

the SVD of \mathbf{X} is given by

$$\mathbf{X} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^\top, \quad (10.55)$$

where $\mathbf{U} \in \mathbb{R}^{D \times D}$ and $\mathbf{V}^\top \in \mathbb{R}^{D \times N}$ are orthogonal matrices and $\mathbf{\Sigma} \in \mathbb{R}^{D \times N}$ is a matrix whose only non-zero entries are the singular values $\sigma_{ii} \geq 0$. Then it follows that

$$\mathbf{S} = \frac{1}{N} \mathbf{X} \mathbf{X}^\top = \frac{1}{N} \mathbf{U} \mathbf{\Sigma} \mathbf{V}^\top \mathbf{V} \mathbf{\Sigma}^\top \mathbf{U}^\top = \frac{1}{N} \mathbf{U} \mathbf{\Sigma} \mathbf{\Sigma}^\top \mathbf{U}^\top. \quad (10.56)$$

With the results from Section 4.5 we get that the columns of \mathbf{U} are the eigenvectors of $\mathbf{X} \mathbf{X}^\top$ (and therefore \mathbf{S}). Furthermore, the eigenvalues of \mathbf{S} are related to the singular values of \mathbf{X} via

$$\lambda_i = \frac{\sigma_i^2}{N}. \quad (10.57)$$

Practical aspects Finding eigenvalues and eigenvectors is also important in other fundamental machine learning methods that require matrix decompositions. In theory, as we discussed in Section 4.2, we can solve for the eigenvalues as roots of the characteristic polynomial. However, for matrices larger than 4×4 this is not possible because we would need to find the roots of a polynomial of degree 5 or higher. However, the Abel-Ruffini theorem (Ruffini, 1799; Abel, 1826) states that there exists no algebraic solution to this problem for polynomials of degree 5 or more. Therefore, in practice, we solve for eigenvalues or singular values using iterative methods, which are implemented in all modern packages for linear algebra.

`np.linalg.eigh`
or
`np.linalg.svd`

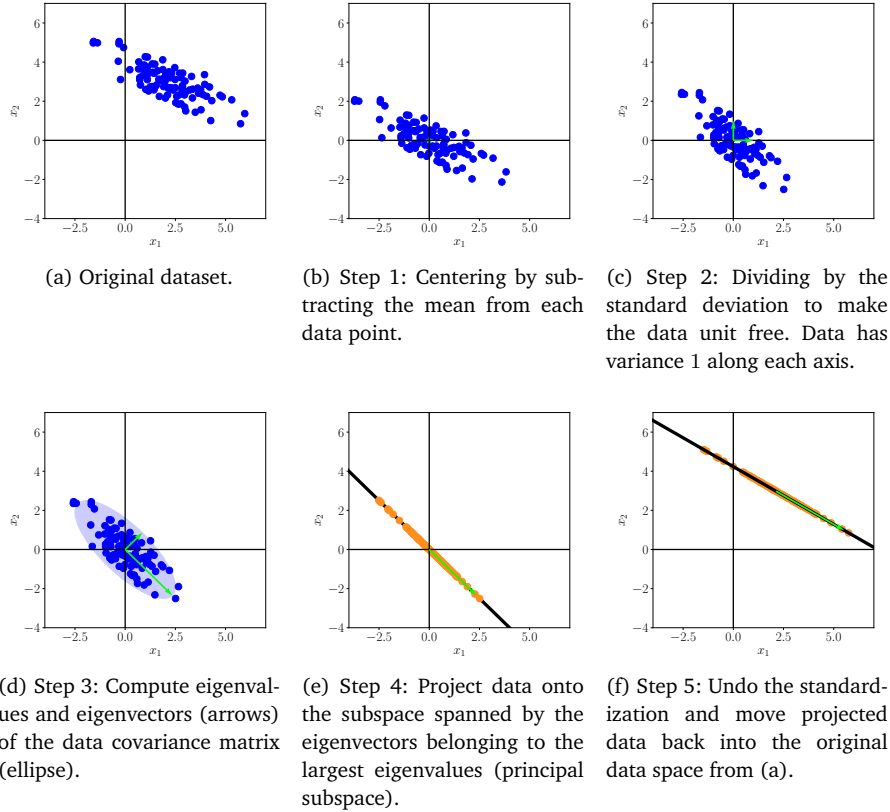
In many applications (such as PCA presented in this chapter), we only require a few eigenvectors. It would be wasteful to compute the full decomposition, and then discard all eigenvectors with eigenvalues that are beyond the first few. It turns out that if we are interested in only the first few eigenvectors (with the largest eigenvalues) 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* is very efficient. Power iteration chooses a random vector \mathbf{x}_0 and follows the iteration

power iteration

$$\mathbf{x}_{k+1} = \frac{\mathbf{S} \mathbf{x}_k}{\|\mathbf{S} \mathbf{x}_k\|}, \quad k = 0, 1, \dots \quad (10.58)$$

This means the vector \mathbf{x}_k is multiplied by \mathbf{S} in every iteration and then normalized, i.e., we always have $\|\mathbf{x}_k\| = 1$. This sequence of vectors converges to the eigenvector associated with the largest eigenvalue of \mathbf{S} . The original Google PageRank algorithm (Page et al., 1999) uses such an algorithm for ranking web pages based on their hyperlinks.

Figure 10.9 Steps of PCA.



10.5 PCA Algorithm

In the following, we will go through the individual steps of PCA using a running example, which is summarized in Figure 10.9. We are given a two-dimensional data set (Figure 10.9(a)), and we want to use PCA to project it onto a one-dimensional subspace.

1. **Mean subtraction** We start by centering the data by computing the mean μ of the dataset and subtracting it from every single data point. This ensures that the data set has mean 0 (Figure 10.9(b)). Mean subtraction is not strictly necessary but reduces the risk of numerical problems.
2. **Standardization** Divide the data points by the standard deviation σ_d of the dataset for every dimension $d = 1, \dots, D$. Now the data is unit free, and it has variance 1 along each axis, which is indicated by the two arrows in Figure 10.9(c). This step completes the *standardization* of the data.
3. **Eigendecomposition of the covariance matrix** Compute the data covariance matrix and its eigenvalues and corresponding eigenvectors. In Figure 10.9(d), the eigenvectors are scaled by the magnitude of the

standardization

corresponding eigenvalue. The longer vector spans the principal subspace, which we denote by U . The data covariance matrix is represented by the ellipse.

4. **Projection** We can project any data point $\mathbf{x}_* \in \mathbb{R}^D$ onto the principal subspace: To get this right, we need to standardize \mathbf{x}_* using the mean and standard deviation of the data set that we used to compute the data covariance matrix, i.e.,

$$x_*^{(d)} \leftarrow \frac{x_*^{(d)} - \mu^{(d)}}{\sigma_d}, \quad d = 1, \dots, D, \quad (10.59)$$

where $x^{(d)}$ is the d th component of \mathbf{x} . Then, we obtain the projected data point as

$$\tilde{\mathbf{x}}_* = \mathbf{B}\mathbf{B}^\top \mathbf{x}_* \quad (10.60)$$

with coordinates $\mathbf{z}_* = \mathbf{B}^\top \mathbf{x}_*$ with respect to the basis of the principal subspace. Here, \mathbf{B} is the matrix that contains the eigenvectors that belong to the largest eigenvalues of the data covariance matrix as columns.

5. **Moving back to data space** To see our projection in the original data format (i.e., before standardization), we need to undo the standardization (10.59) and multiply by the standard deviation before adding the mean so that we obtain

$$\tilde{x}_*^{(d)} \leftarrow \tilde{x}_*^{(d)} \sigma_d + \mu^{(d)}, \quad d = 1, \dots, D, \quad (10.61)$$

where $\mu^{(d)}$ and σ_d are the mean and standard deviation of the training data in the d th dimension, respectively. Figure 10.9(f) illustrates the projection in the original data format.

Example 10.3 (MNIST Digits: Reconstruction)

In the following, we will apply PCA to the MNIST digits dataset, which contains 60,000 examples of handwritten digits 0–9. Each digit is an image of size 28×28 , i.e., it contains 784 pixels so that we can interpret every image in this dataset as a vector $\mathbf{x} \in \mathbb{R}^{784}$. Examples of these digits are shown in Figure 10.3. For illustration purposes, we apply PCA to a subset of the MNIST digits, and we focus on the digit ‘8’. We used 5,389 training images of the digit ‘8’ and determined the principal subspace as detailed in this chapter. We then used the learned projection matrix to reconstruct a set of test images, which is illustrated in Figure 10.10. The first row of Figure 10.10 shows a set of four original digits from the test set. The following rows show reconstructions of exactly these digits when using a principal subspace of dimensions 1, 10, 100, 500, respectively. We can

<http://yann.lecun.com/exdb/mnist/>

see that even with a single-dimensional principal subspace we get a half-way decent reconstruction of the original digits, which, however, is blurry and generic. With an increasing number of principal components (PCs) the reconstructions become sharper and more details can be accounted for. With 500 principal components, we effectively obtain a near-perfect reconstruction. If we were to choose 784 PCs we would recover the exact digit without any compression loss.

Figure 10.10 Effect of increasing number of principal components on reconstruction.

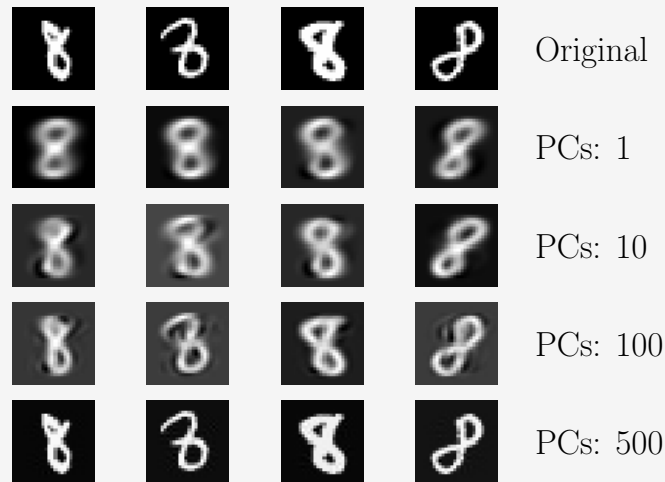


Figure 10.11 Average reconstruction error as a function of the number of principal components.

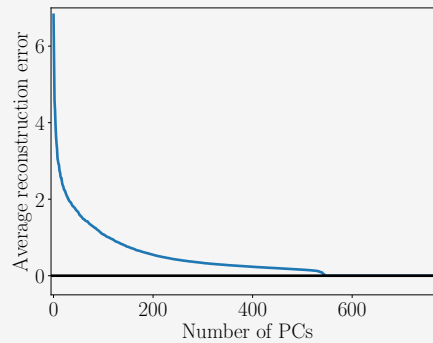


Figure 10.11 shows the average reconstruction error, which is

$$\frac{1}{N} \sum_{n=1}^N \|\mathbf{x}_n - \tilde{\mathbf{x}}_n\| = \sum_{d=1}^D \sqrt{\lambda_d}, \quad (10.62)$$

as a function of the number of principal components. We can see that the importance of the principal components drops off rapidly, and only

marginal gains can be achieved by adding more PCs. With about 550 PCs, we can essentially fully reconstruct the training data that contains the digit ‘8’.

10.6 PCA in High Dimensions

In order to do PCA, we need to compute the data covariance matrix. In D dimensions, the data covariance matrix is a $D \times D$ matrix. Computing the eigenvalues and eigenvectors of this matrix is computationally expensive as it scales cubically in D . Therefore, PCA, as we discussed earlier, will be infeasible in very high dimensions. For example, if our \mathbf{x}_n are images with 10,000 pixels (e.g., 100×100 pixel images), we would need to compute the eigendecomposition of a $10,000 \times 10,000$ covariance matrix. In the following, we provide a solution to this problem for the case that we have substantially fewer data points than dimensions, i.e., $N \ll D$.

Assume we have a data set $\mathbf{x}_1, \dots, \mathbf{x}_N, \mathbf{x}_n \in \mathbb{R}^D$. Assuming the data is centered, the data covariance matrix is given as

$$\mathbf{S} = \frac{1}{N} \mathbf{X} \mathbf{X}^\top \in \mathbb{R}^{D \times D}, \quad (10.63)$$

where $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_N]$ is a $D \times N$ matrix whose columns are the data points.

We now assume that $N \ll D$, i.e., the number of data points is smaller than the dimensionality of the data. Then the rank of the covariance matrix \mathbf{S} is N , and it has $D - N + 1$ many eigenvalues that are 0. Intuitively, this means that there are some redundancies.

In the following, we will exploit this and turn the $D \times D$ covariance matrix into an $N \times N$ covariance matrix whose eigenvalues are all greater than 0.

In PCA, we ended up with the eigenvector equation

$$\mathbf{S} \mathbf{b}_m = \lambda_m \mathbf{b}_m, \quad m = 1, \dots, M, \quad (10.64)$$

where \mathbf{b}_m is a basis vector of the principal subspace. Let us re-write this equation a bit: With \mathbf{S} defined in (10.63), we obtain

$$\mathbf{S} \mathbf{b}_m = \frac{1}{N} \mathbf{X} \mathbf{X}^\top \mathbf{b}_m = \lambda_m \mathbf{b}_m. \quad (10.65)$$

We now multiply $\mathbf{X}^\top \in \mathbb{R}^{N \times D}$ from the left-hand side, which yields

$$\frac{1}{N} \underbrace{\mathbf{X}^\top \mathbf{X}}_{N \times N} \underbrace{\mathbf{X}^\top \mathbf{b}_m}_{=: \mathbf{c}_m} = \lambda_m \mathbf{X}^\top \mathbf{b}_m \iff \frac{1}{N} \mathbf{X}^\top \mathbf{X} \mathbf{c}_m = \lambda_m \mathbf{c}_m, \quad (10.66)$$

and we get a new eigenvector/eigenvalue equation: λ_m is still the eigenvalue, but the eigenvector is now $\mathbf{c}_m := \mathbf{X}^\top \mathbf{b}_m$ of the matrix $\frac{1}{N} \mathbf{X}^\top \mathbf{X} \in$

$\mathbb{R}^{N \times N}$. Assuming we have no duplicate data points, this matrix has rank N and is invertible. This also implies that $\frac{1}{N} \mathbf{X}^\top \mathbf{X}$ has the same (non-zero) eigenvalues as the data covariance matrix \mathbf{S} . But this is now an $N \times N$ matrix, so that we can compute the eigenvalues and eigenvectors much more efficiently than for the original $D \times D$ data covariance matrix.

Now, that we have the eigenvectors of $\frac{1}{N} \mathbf{X}^\top \mathbf{X}$, we are going to recover the original eigenvectors, which we still need for PCA. Currently, we know the eigenvectors of $\frac{1}{N} \mathbf{X}^\top \mathbf{X}$. If we left-multiply our eigenvalue/eigenvector equation with \mathbf{X} , we get

$$\underbrace{\frac{1}{N} \mathbf{X} \mathbf{X}^\top}_{\mathbf{S}} \mathbf{X} \mathbf{c}_m = \lambda_m \mathbf{X} \mathbf{c}_m \quad (10.67)$$

and we recover the data covariance matrix again. This now also means that we recover $\mathbf{X} \mathbf{c}_m$ as an eigenvector of \mathbf{S} .

Remark. If we want to apply the PCA algorithm that we discussed in Section 10.5 we need to normalize the eigenvectors $\mathbf{X} \mathbf{c}_m$ of \mathbf{S} so that they have norm 1. \diamond

10.7 Probabilistic Principal Component Analysis

In the previous sections, we derived PCA without any notion of a probabilistic model using the maximum-variance and the projection perspectives. On the one hand this approach may be appealing as it allows us to sidestep all the mathematical difficulties that come with probability theory, on the other hand a probabilistic model would offer us more flexibility and useful insights. More specifically, a probabilistic model would

- come with a likelihood function, and we can explicitly deal with noisy observations (which we did not even discuss earlier),
- allow us to do Bayesian model comparison via the marginal likelihood as discussed in Section 8.4,
- view PCA as a generative model, which allows us to simulate new data,
- allow us to make straightforward connections to related algorithms
- deal with data dimensions that are missing at random by applying Bayes' theorem,
- give us a notion of the novelty of a new data point,
- allow us to extend the model fairly straightforwardly, e.g., to a mixture of PCA models,
- have the PCA we derived in earlier sections as a special case,
- allow for a fully Bayesian treatment by marginalizing out the model parameters.

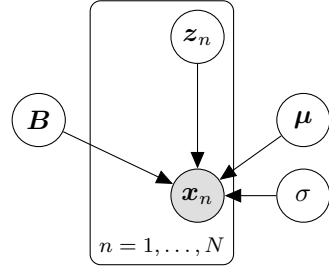


Figure 10.12
Graphical model for probabilistic PCA. The observations \mathbf{x}_n explicitly depend on corresponding latent variables $\mathbf{z}_n \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$. The model parameters \mathbf{B} , $\boldsymbol{\mu}$ and the likelihood parameter σ are shared across the dataset.

Tipping and Bishop (1999) proposed *Probabilistic PCA* (PPCA) to address most of these issues, and the PCA solution that we obtained by maximizing the variance in the projected space or by minimizing the reconstruction error is obtained as the special case of maximum likelihood in a noise-free setting.

10.7.1 Generative Process and Probabilistic Model

In PPCA, we explicitly write down the probabilistic model for linear dimensionality reduction. For this we assume a continuous latent variable $\mathbf{z} \in \mathbb{R}^M$ with a standard-Normal prior $p(\mathbf{z}) = \mathcal{N}(\mathbf{0}, \mathbf{I})$ and a linear relationship between the latent variables and the observed \mathbf{x} data where

$$\mathbf{x} = \mathbf{B}\mathbf{z} + \boldsymbol{\mu} + \boldsymbol{\epsilon} \in \mathbb{R}^D, \quad (10.68)$$

where $\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$ is Gaussian observation noise, $\mathbf{B} \in \mathbb{R}^{D \times M}$ and $\boldsymbol{\mu} \in \mathbb{R}^D$ describe the linear/affine mapping from latent to observed variables. Therefore, PPCA links latent and observed variables via

$$p(\mathbf{x} | \mathbf{z}, \mathbf{B}, \boldsymbol{\mu}, \sigma^2) = \mathcal{N}(\mathbf{x} | \mathbf{B}\mathbf{z} + \boldsymbol{\mu}, \sigma^2 \mathbf{I}). \quad (10.69)$$

Overall, PPCA induces the following generative process:

$$\mathbf{z}_n \sim \mathcal{N}(\mathbf{z} | \mathbf{0}, \mathbf{I}) \quad (10.70)$$

$$\mathbf{x}_n | \mathbf{z}_n \sim \mathcal{N}(\mathbf{x} | \mathbf{B}\mathbf{z}_n + \boldsymbol{\mu}, \sigma^2 \mathbf{I}) \quad (10.71)$$

To generate a data point that is typical given the model parameters, we follow an *ancestral sampling* scheme: We first sample a latent variable \mathbf{z}_n from $p(\mathbf{z})$. Then, we use \mathbf{z}_n in (10.69) to sample a data point conditioned on the sampled \mathbf{z}_n , i.e., $\mathbf{x}_n \sim p(\mathbf{x} | \mathbf{z}_n, \mathbf{B}, \boldsymbol{\mu}, \sigma^2)$.

ancestral sampling

This generative process allows us to write down the probabilistic model (i.e., the joint distribution of all random variables) as

$$p(\mathbf{x}, \mathbf{z} | \mathbf{B}, \boldsymbol{\mu}, \sigma^2) = p(\mathbf{x} | \mathbf{z}, \mathbf{B}, \boldsymbol{\mu}, \sigma^2) p(\mathbf{z}), \quad (10.72)$$

which immediately gives rise to the graphical model in Figure 10.12 using the results from Section 8.5.

5553 *Remark.* Note the direction of the arrow that connects the latent variables
 5554 z and the observed data x : The arrow points from z to x , which means
 5555 that the PPCA model assumes a lower-dimensional latent cause z for high-
 5556 dimensional observations x . In the end, we are obviously interested in
 5557 finding something out about z given some observations. To get there we
 5558 will apply Bayesian inference to “invert” the arrow implicitly and go from
 5559 observations to latent variables. \diamond

Example 10.4 (Generating Data from Latent Variables)

Figure 10.13
 Generating new
 MNIST digits. The
 latent variables z
 can be used to
 generate new data
 $\tilde{x} = Bz$. The closer
 we stay to the
 training data the
 more realistic the
 generated data.

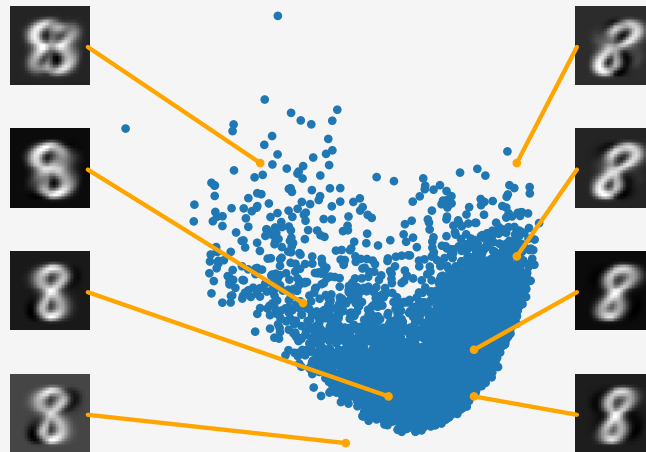


Figure 10.13 shows the latent coordinates of the MNIST digits ‘8’ found by PCA when using a two-dimensional principal subspace (blue dots). We can query any vector z_* in this latent space and generate an image $\tilde{x}_* = Bz_*$ that resembles the digit ‘8’. We show eight of such generated images with their corresponding latent space representation. Depending on where we query the latent space, the generated images look different (shape, rotation, size, ...). If we query away from the training data, we see more and more artefacts, e.g., the top-left and top-right digits. Note that the intrinsic dimensionality of these generated images is only two.

10.7.2 Likelihood and Joint Distribution

Using the results from Chapter 6, we obtain the marginal distribution of the data \mathbf{x} by integrating out the latent variable \mathbf{z} so that

$$\begin{aligned} p(\mathbf{x} | \mathbf{B}, \boldsymbol{\mu}, \sigma^2) &= \int p(\mathbf{x} | \mathbf{z}, \boldsymbol{\mu}, \sigma^2) p(\mathbf{z}) d\mathbf{z} \\ &= \int \mathcal{N}(\mathbf{x} | \mathbf{B}\mathbf{z} + \boldsymbol{\mu}, \sigma^2 \mathbf{I}) \mathcal{N}(\mathbf{z} | \mathbf{0}, \mathbf{I}) d\mathbf{z}. \end{aligned} \quad (10.73)$$

From Section 6.6, we know that the solution to this integral is a Gaussian distribution with mean

$$\mathbb{E}[\mathbf{x}] = \mathbb{E}_z[\mathbf{B}\mathbf{z} + \boldsymbol{\mu}] + \mathbb{E}_\epsilon[\boldsymbol{\epsilon}] = \boldsymbol{\mu} \quad (10.74)$$

and with covariance matrix

$$\begin{aligned} \mathbb{V}[\mathbf{x}] &= \mathbb{V}_z[\mathbf{B}\mathbf{z} + \boldsymbol{\mu}] + \mathbb{V}_\epsilon[\boldsymbol{\epsilon}] = \mathbb{V}_z[\mathbf{B}\mathbf{z}] \\ &= \mathbf{B} \mathbb{V}_z[\mathbf{z}] \mathbf{B}^\top + \sigma^2 \mathbf{I} = \mathbf{B} \mathbf{B}^\top + \sigma^2 \mathbf{I}. \end{aligned} \quad (10.75)$$

The marginal distribution in (10.73) is the *PPCA likelihood*, which we can use for maximum likelihood or MAP estimation of the model parameters.

PPCA likelihood

Remark. Although the conditional distribution in (10.69) is also a likelihood, we cannot use it for maximum likelihood estimation as it still depends on the latent variables. The likelihood function we require for maximum likelihood (or MAP) estimation should only be a function of the data \mathbf{x} and the model parameters, but not on the latent variables. \diamond

From Section 6.6 we also know that the joint distribution of a Gaussian random variable \mathbf{z} and a linear/affine transformation $\mathbf{x} = \mathbf{B}\mathbf{z}$ of it are jointly Gaussian distributed. We already know the marginals $p(\mathbf{z}) = \mathcal{N}(\mathbf{z} | \mathbf{0}, \mathbf{I})$ and $p(\mathbf{x}) = \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}, \mathbf{B}\mathbf{B}^\top + \sigma^2 \mathbf{I})$. The missing cross-covariance is given as

$$\text{Cov}[\mathbf{x}, \mathbf{z}] = \text{Cov}_z[\mathbf{B}\mathbf{z} + \boldsymbol{\mu}, \mathbf{z}] = \mathbf{B} \text{Cov}_z[\mathbf{z}, \mathbf{z}] = \mathbf{B}. \quad (10.76)$$

Therefore, the probabilistic model of PPCA, i.e., the joint distribution of latent and observed random variables is explicitly given by

$$p(\mathbf{x}, \mathbf{z} | \mathbf{B}, \boldsymbol{\mu}, \sigma^2) = \mathcal{N} \left(\begin{bmatrix} \mathbf{x} \\ \mathbf{z} \end{bmatrix} \middle| \begin{bmatrix} \boldsymbol{\mu} \\ \mathbf{0} \end{bmatrix}, \begin{bmatrix} \mathbf{B}\mathbf{B}^\top + \sigma^2 \mathbf{I} & \mathbf{B} \\ \mathbf{B}^\top & \mathbf{I} \end{bmatrix} \right), \quad (10.77)$$

with a mean vector of length $D + M$ and a covariance matrix of size $(D + M) \times (D + M)$.

10.7.3 Posterior Distribution

The joint Gaussian distribution $p(\mathbf{x}, \mathbf{z} | \mathbf{B}, \boldsymbol{\mu}, \sigma^2)$ in (10.77) allows us to determine the posterior distribution $p(\mathbf{z} | \mathbf{x})$ immediately by applying the

rules of Gaussian conditioning from Section 6.6.1. The posterior distribution of the latent variable given an observation \mathbf{x} is then

$$p(\mathbf{z} | \mathbf{x}) = \mathcal{N}(\mathbf{z} | \mathbf{m}, \mathbf{C}), \quad (10.78a)$$

$$\mathbf{m} = \mathbf{B}^\top (\mathbf{B}\mathbf{B}^\top + \sigma^2 \mathbf{I})^{-1} (\mathbf{x} - \boldsymbol{\mu}), \quad (10.78b)$$

$$\mathbf{C} = \mathbf{I} - \mathbf{B}^\top (\mathbf{B}\mathbf{B}^\top + \sigma^2 \mathbf{I})^{-1} \mathbf{B}. \quad (10.78c)$$

5571 The posterior distribution in (10.78) reveals a few things. First, the poste-
 5572 rior mean \mathbf{m} is effectively an orthogonal projection of the mean-centered
 5573 data $\mathbf{x} - \boldsymbol{\mu}$ onto the vector subspace spanned by the columns of \mathbf{B} . If we
 5574 ignore the measurement noise contribution we immediately recover the
 5575 projection equation (3.55) from Section 3.6. Second, the posterior covari-
 5576 ance matrix \mathbf{C} does not directly depend on the observed data \mathbf{x} .

5577 If we now have a new observation \mathbf{x}_* in data space, we can use (10.78)
 5578 to determine the posterior distribution of the corresponding latent vari-
 5579 able \mathbf{z}_* . The covariance matrix \mathbf{C} allows us to assess how confident the
 5580 embedding is. A covariance matrix \mathbf{C} with a small determinant (which
 5581 measures volumes) tells us that the latent embedding \mathbf{z}_* is fairly certain.
 5582 However, if we obtain a posterior distribution $p(\mathbf{z}_* | \mathbf{x}_*)$ that is uncertain,
 5583 we may be faced with an outlier. However, we can explore this posterior
 5584 distribution to understand what other data points \mathbf{x} are plausible under
 5585 this posterior. To do this, we can exploit PPCA's generative process. The
 5586 generative process underlying PPCA allows us to explore the posterior dis-
 5587 tribution on the latent variables by generating new data that are plausible
 5588 under this posterior. This can be achieved as follows:

- 5589 1. Sample a latent variable $\mathbf{z}_* \sim p(\mathbf{z} | \mathbf{x}_*)$ from the posterior distribution
 5590 over the latent variables (10.78)
- 5591 2. Sample a reconstructed vector $\tilde{\mathbf{x}}_* \sim p(\mathbf{x} | \mathbf{z}_*, \mathbf{B}, \boldsymbol{\mu}, \sigma^2)$ from (10.69)

5592 If we repeat this process many times, we can explore the posterior dis-
 5593 tribution (10.78) on the latent variables \mathbf{z}_* and its implications on the
 5594 observed data. The sampling process effectively hypothesizes data, which
 5595 is plausible under the posterior distribution.

5596 10.8 Further Reading

5597 We derived PCA from two perspectives: a) maximizing the variance in the
 5598 projected space; b) minimizing the average reconstruction error. However,
 5599 PCA can also be interpreted from different perspectives. Let us re-cap what
 5600 we have done: We took high-dimensional data $\mathbf{x} \in \mathbb{R}^D$ and used a matrix
 5601 \mathbf{B}^\top to find a lower-dimensional representation $\mathbf{z} \in \mathbb{R}^M$. The columns of
 5602 \mathbf{B} are the eigenvectors of the data covariance matrix \mathbf{S} that are associated
 5603 with the largest eigenvalues. Once we have a low-dimensional represen-
 5604 tation \mathbf{z} , we can get a high-dimensional version of it (in the original data

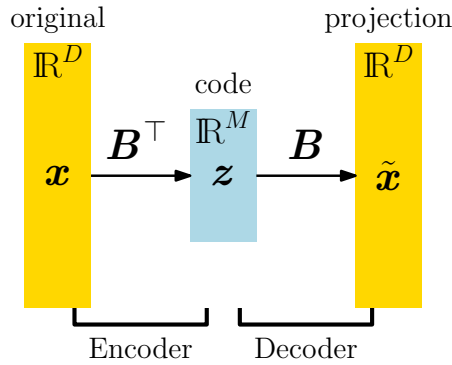


Figure 10.14 PCA can be viewed as a linear auto-encoder. It encodes the high-dimensional data x into a lower-dimensional representation (code) $z \in \mathbb{R}^M$ and decodes z using a decoder. The decoded vector \tilde{x} is the orthogonal projection of the original data x onto the M -dimensional principal subspace.

space) as $x \approx \tilde{x} = Bz = BB^\top x \in \mathbb{R}^D$, where BB^\top is a projection matrix.

We can also think of PCA as a linear *auto-encoder* as illustrated in Figure 10.14. An auto-encoder encodes the data $x_n \in \mathbb{R}^D$ to a *code* $z_n \in \mathbb{R}^M$ and tries to decode it to a \tilde{x}_n similar to x_n . The mapping from the data to the code is called the *encoder*, the mapping from the code back to the original data space is called the *decoder*. If we consider linear mappings where the code is given by $z_n = B^\top x_n \in \mathbb{R}^M$ and we are interested in minimizing the average squared error between the data x_n and its reconstruction $\tilde{x}_n = Bz_n$, $n = 1, \dots, N$, we obtain

$$\frac{1}{N} \sum_{n=1}^N \|x_n - \tilde{x}_n\|^2 = \frac{1}{N} \sum_{n=1}^N \|x_n - B^\top B x_n\|^2. \quad (10.79)$$

This means we end up with the same objective function as in (10.31) that we discussed in Section 10.3 so that we obtain the PCA solution when we minimize the squared auto-encoding loss. If we replace the linear mapping of PCA with a nonlinear mapping, we get a nonlinear auto-encoder. A prominent example of this is a deep auto-encoder where the linear functions are replaced with deep neural networks. In this context, the encoder is also known as *recognition network* or *inference network*, whereas the decoder is also called a *generator*.

Another interpretation of PCA is related to information theory. We can think of the code as a smaller or compressed version of the original data point. When we reconstruct our original data using the code, we do not get the exact data point back, but a slightly distorted or noisy version of it. This means that our compression is “lossy”. Intuitively we want to maximize the correlation between the original data and the lower-dimensional code. More formally, this is related to the mutual information. We would then get the same solution to PCA we discussed in Section 10.3 by maximizing the mutual information, a core concept in information theory (MacKay, 2003).

In our discussion on PPCA, we assumed that the parameters of the

auto-encoder
code

encoder
decoder

recognition network
inference network
generator

The code is a
compressed version
of the original data.

model, i.e., \mathbf{B} , $\boldsymbol{\mu}$ and the likelihood parameter σ^2 are known. Tipping and Bishop (1999) describe how to derive maximum likelihood estimates for these parameter in the PPCA setting (note that we use a different notation in this chapter). The maximum likelihood parameters, when projecting D -dimensional data onto an M -dimensional subspace, are given by

$$\boldsymbol{\mu}_{\text{ML}} = \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n, \quad (10.80)$$

$$\mathbf{B}_{\text{ML}} = \mathbf{T}(\boldsymbol{\Lambda} - \sigma^2 \mathbf{I})^{\frac{1}{2}}, \quad (10.81)$$

$$\sigma_{\text{ML}}^2 = \frac{1}{D - M} \sum_{j=M+1}^D \lambda_j, \quad (10.82)$$

where $\mathbf{T} \in \mathbb{R}^{D \times M}$ contains M eigenvectors of the data covariance matrix, and $\boldsymbol{\Lambda} = \text{diag}(\lambda_1, \dots, \lambda_M) \in \mathbb{R}^{M \times M}$ is a diagonal matrix with the eigenvalues belonging to the principal axes on its diagonal. The maximum likelihood solution \mathbf{B}_{ML} is unique up to a arbitrary rotations, i.e., we can right-multiply \mathbf{B}_{ML} with any rotation matrix $\mathbf{R} \in \mathbb{R}^{M \times M}$ so that (10.81) essentially is a singular value decomposition (see Section 4.5). An outline of the proof is given by Tipping and Bishop (1999).

The maximum likelihood estimate for $\boldsymbol{\mu}$ given in (10.80) is the sample mean of the data. The maximum likelihood estimator for the observation noise variance σ^2 given in (10.82) is the sum of all variances in the orthogonal complement of the principal subspace, i.e., the leftover variance that we cannot capture with the first M principal components are treated as observation noise.

In the noise-free limit where $\sigma \rightarrow 0$, PPCA and PCA provide identical solutions: Since the data covariance matrix \mathbf{S} is symmetric, it can be diagonalized (see Section 4.4), i.e., there exists a matrix \mathbf{T} of eigenvectors of \mathbf{S} so that

$$\mathbf{S} = \mathbf{T} \boldsymbol{\Lambda} \mathbf{T}^{-1}. \quad (10.83)$$

In the PPCA model, the data covariance matrix is the covariance matrix of the likelihood $p(\mathbf{X} | \mathbf{B}, \boldsymbol{\mu}, \sigma^2)$, which is $\mathbf{B} \mathbf{B}^\top + \sigma^2 \mathbf{I}$, see (10.75). For $\sigma \rightarrow 0$, we obtain $\mathbf{B} \mathbf{B}^\top$ so that this data covariance must equal the PCA data covariance (and its factorization given in (10.83)) so that

$$\text{Cov}[\mathbf{X}] = \mathbf{T} \boldsymbol{\Lambda} \mathbf{T}^{-1} = \mathbf{B} \mathbf{B}^\top \iff \mathbf{B} = \mathbf{T} \boldsymbol{\Lambda}^{\frac{1}{2}}, \quad (10.84)$$

which is exactly the maximum likelihood estimate in (10.81) for $\sigma = 0$.

From (10.81) and (10.83) it becomes clear that (P)PCA performs a decomposition of the data covariance matrix. We can also think of PCA as a method for finding the best rank- M approximation of the data covariance matrix so that we can apply the Eckart-Young Theorem from Section 4.6, which states that the optimal solution can be determined using a singular value decomposition.

In a streaming setting, where data arrives sequentially, it is recommended to use the iterative Expectation Maximization (EM) algorithm for maximum likelihood estimation (Roweis, 1998).

Similar to our discussion on linear regression in Chapter 9, we can place a prior distribution on the parameters of the model and to integrate them out, thereby avoiding a) point estimates of the parameters and the issues that come with these point estimates (see Section 8.4) and b) allowing for an automatic selection of the appropriate dimensionality M of the latent space. In this *Bayesian PCA*, which was proposed by Bishop (1999), places a (hierarchical) prior $p(\boldsymbol{\mu}, \mathbf{B}, \sigma^2)$ on the model parameters. The generative process allows us to integrate the model parameters out instead of conditioning on them, which addresses overfitting issues. Since this integration is analytically intractable, Bishop (1999) proposes to use approximate inference methods, such as MCMC or variational inference. We refer to the work by Gilks et al. (1996) and Blei et al. (2017) for more details on these approximate inference techniques.

Bayesian PCA

In PPCA, we considered the linear model $\mathbf{x}_n = \mathbf{B}\mathbf{z}_n + \boldsymbol{\epsilon}$ with $p(\mathbf{z}_n) = \mathcal{N}(\mathbf{0}, \mathbf{I})$ and $\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$, i.e., all observation dimensions are affected by the same amount of noise. If we allow each observation dimension d to have a different variance σ_d^2 we obtain *factor analysis* (FA) (Spearman, 1904; Bartholomew et al., 2011). This means, FA gives the likelihood some more flexibility than PPCA, but still forces the data to be explained by the model parameters \mathbf{B} , $\boldsymbol{\mu}$. However, FA no longer allows for a closed-form solution to maximum likelihood so that we need to use an iterative scheme, such as the EM algorithm, to estimate the model parameters. While in PPCA all stationary points are global optima, this no longer holds for FA. Compared to PPCA, FA does not change if we scale the data, but it does return different solutions if we rotate the data.

factor analysis

An overly flexible likelihood would be able to explain more than just the noise.

Independent Component Analysis (ICA) is also closely related to PCA. Starting again with the model $\mathbf{x}_n = \mathbf{B}\mathbf{z}_n + \boldsymbol{\epsilon}$ we now change the prior on \mathbf{z}_n to non-Gaussian distributions. ICA can be used for *blind-source separation*. Imagine you are in a busy train station with many people talking. Your ears play the role of microphones, and they linearly mix different speech signals in the train station. The goal of blind-source separation is to identify the constituent parts of the mixed signals. As discussed above in the context of maximum likelihood estimation for PPCA, the original PCA solution is invariant to any rotation. Therefore, PCA can identify the best lower-dimensional subspace in which the signals live, but not the signals themselves (Murphy, 2012). ICA addresses this issue by modifying the prior distribution $p(\mathbf{z})$ on the latent sources to require non-Gaussian priors $p(\mathbf{z})$. We refer to the book by Murphy2012 for more details on ICA.

Independent Component Analysis

blind-source separation

PCA, factor analysis and ICA are three examples for dimensionality reduction with linear models. Cunningham and Ghahramani (2015) provide a broader survey of linear dimensionality reduction.

Murphy2012

The (P)PCA model we discussed here allows for several important ex-

5690 tensions. In Section 10.6, we explained how to do PCA when the in-
 5691 put dimensionality D is significantly greater than the number N of data
 5692 points. By exploiting the insight that PCA can be performed by computing
 5693 (many) inner products, this idea can be pushed to the extreme by consid-
 5694 ering infinite-dimensional features. The *kernel trick* is the basis of *kernel*
 5695 *PCA* and allows us to implicitly compute inner products between infinite-
 5696 dimensional features (Schölkopf et al., 1998; Schölkopf and Smola, 2002).
 5697 There are nonlinear dimensionality reduction techniques that are de-
 5698 rived from PCA. The auto-encoder perspective of PCA that we discussed
 5699 above can be used to render PCA as a special case of a *deep auto-encoder*.
 5700 In the deep auto-encoder, both the encoder and the decoder are repre-
 5701 sented by multi-layer feedforward neural networks, which themselves are
 5702 nonlinear mappings. If we set the activation functions in these neural net-
 5703 works to be the identity, the model becomes equivalent to PCA. A different
 5704 approach to nonlinear dimensionality reduction is the *Gaussian Process La-*
 5705 *tent Variable Model* (GP-LVM) proposed by Lawrence (2005). The GP-LVM
 5706 starts off with the latent-variable perspective that we used to derive PPCA
 5707 and replaces the linear relationship between the latent variables z and the
 5708 observations x with a Gaussian process (GP). Instead of estimating the pa-
 5709 rameters of the mapping (as we do in PPCA), the GP-LVM marginalizes out
 5710 the model parameters and makes point estimates of the latent variables
 5711 z . Similar to Bayesian PCA, the *Bayesian GP-LVM* proposed by Titsias and
 5712 Lawrence (2010) maintains a distribution on the latent variables z and
 5713 uses approximate inference to integrate them out as well.