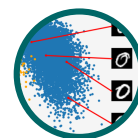


## Dimensionality Reduction with Principal Component Analysis



A  $640 \times 480$  pixels color image is a data point in a million-dimensional space, where every pixel responds to three dimensions, one for each color channel (red, green, blue).

Working directly with high-dimensional data, such as images, comes with some difficulties: it is hard to analyze, interpretation is difficult, visualization is nearly impossible, and (from a practical point of view) storage of the data vectors can be expensive. However, high-dimensional data often has properties that we can exploit. For example, high-dimensional data is often overcomplete, i.e., many dimensions are redundant and can be explained by a combination of other dimensions. Furthermore, dimensions in high-dimensional data are often correlated so that the data possesses an intrinsic lower-dimensional structure. Dimensionality reduction exploits structure and correlation and allows us to work with a more compact representation of the data, ideally without losing information. We can think of dimensionality reduction as a compression technique, similar to jpeg or mp3, which are compression algorithms for images and music.

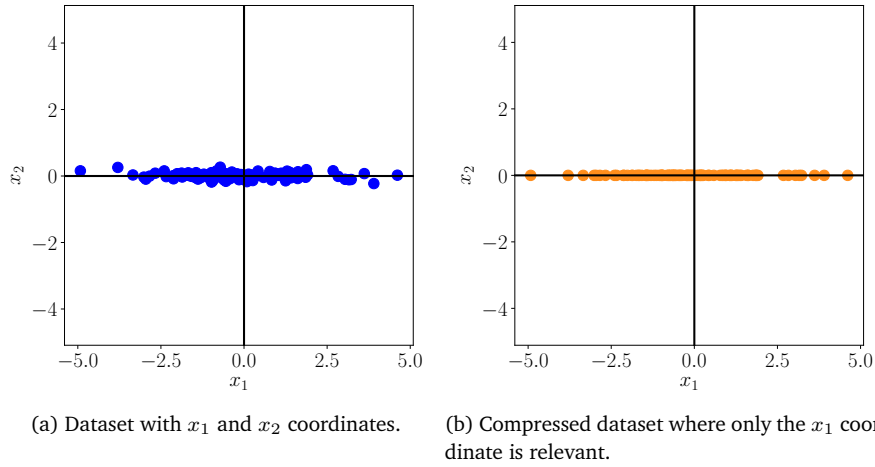
In this chapter, we will discuss *principal component analysis* (PCA), an algorithm for linear *dimensionality reduction*. PCA, proposed by Pearson (1901b) and Hotelling (1933), has been around for more than 100 years and is still one of the most commonly used techniques for data compression and data visualization. It is also used for the identification of simple patterns, latent factors and structures of high-dimensional data. In the signal processing community, PCA is also known as the *Karhunen-Loève transform*. In this chapter, we derive PCA from first principles, drawing on our understanding of basis and basis change (see Sections 2.6.1 and 2.7.2), projections (see Section 3.7), eigenvalues (see Section 4.2), Gaussian distributions (see Section 6.5) and constrained optimization (see Section 7.2).

Dimensionality reduction generally exploits a property of high-dimensional data (e.g., images) that it often lies on a low-dimensional subspace, 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). To describe the data in Figure 10.1(b), only the  $x_1$ -coordinate is required, and the data lies in a one-dimensional subspace of  $\mathbb{R}^2$ .

In the context of Table 1.1, the problem of dimensionality reduction falls

principal component analysis  
dimensionality reduction

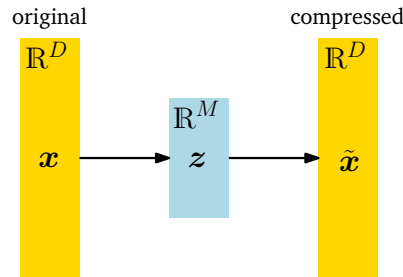
Karhunen-Loève transform



**Figure 10.1**  
Illustration:  
Dimensionality  
reduction. (a) The  
original dataset  
does 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  $z$ .



into the category of an unsupervised learning problem with continuous latent variables.

### 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.

More concretely, we consider an i.i.d. dataset  $\mathcal{X} = \{x_1, \dots, x_N\}$ ,  $x_n \in \mathbb{R}^D$ , with mean  $\mathbf{0}$  that possesses the *data covariance matrix*

$$S = \frac{1}{N} \sum_{n=1}^N x_n x_n^\top. \quad (10.1)$$

Furthermore, we assume there exists a low-dimensional compressed representation (code)

$$z_n = B^\top x_n \in \mathbb{R}^M \quad (10.2)$$

of  $x_n$ , where we define the projection matrix

$$B := [b_1, \dots, b_M] \in \mathbb{R}^{D \times M}. \quad (10.3)$$



**Figure 10.3**  
Examples of  
handwritten digits  
from the MNIST  
dataset. [http:  
//yann.lecun.  
com/exdb/mnist/](http://yann.lecun.com/exdb/mnist/)

We assume that the columns of  $B$  are orthonormal so that  $\mathbf{b}_i^\top \mathbf{b}_j = 0$  if and only if  $i \neq j$ . We seek an  $M$ -dimensional subspace  $U \subseteq \mathbb{R}^D$ ,  $\dim(U) = M < D$  onto which we project the data. We denote the projected data by  $\tilde{\mathbf{x}}_n \in U$ , and their coordinates (with respect to the basis vectors  $\mathbf{b}_1, \dots, \mathbf{b}_M$  of  $U$ ) by  $\mathbf{z}_n$ . Our aim is to find projections  $\tilde{\mathbf{x}}_n \in \mathbb{R}^D$  (or equivalently the codes  $\mathbf{z}_n$  and the basis vectors  $\mathbf{b}_1, \dots, \mathbf{b}_M$ ) so that they are as similar to the original data  $\mathbf{x}_n$  and minimize the loss due to compression.

In Section 10.2, we will find low-dimensional representations that retain as much information as possible and minimize the compression loss. An alternative derivation of PCA is given in Section 10.3, we will be looking at minimizing the squared reconstruction error  $\|\mathbf{x}_n - \tilde{\mathbf{x}}_n\|^2$  between the original data  $\mathbf{x}_n$  and its projection  $\tilde{\mathbf{x}}_n$ .

Figure 10.2 illustrates the setting we consider in PCA, where  $\mathbf{z}$  represents the intrinsic lower dimension of the compressed data  $\tilde{\mathbf{x}}$  and plays the role of a bottleneck, which controls how much information can flow between  $\mathbf{x}$  and  $\tilde{\mathbf{x}}$ . In PCA, we consider a linear relationship between the original data  $\mathbf{x}$  and its low-dimensional code  $\mathbf{z}$  so that  $\mathbf{z} = \mathbf{B}^\top \mathbf{x}$  and  $\tilde{\mathbf{x}} = \mathbf{B}\mathbf{z}$  for a suitable matrix  $\mathbf{B}$ .

#### Example 10.1 (Coordinate Representation/Code)

Consider  $\mathbb{R}^2$  with the canonical basis  $\mathbf{e}_1 = [1, 0]^\top$ ,  $\mathbf{e}_2 = [0, 1]^\top$ . From Chapter 2 we know that  $\mathbf{x} \in \mathbb{R}^2$  can be represented as a linear combination of these basis vectors, e.g.,

$$\begin{bmatrix} 5 \\ 3 \end{bmatrix} = 5\mathbf{e}_1 + 3\mathbf{e}_2. \quad (10.4)$$

However, when we consider vectors of the form

$$\tilde{\mathbf{x}} = \begin{bmatrix} 0 \\ z \end{bmatrix} \in \mathbb{R}^2, \quad z \in \mathbb{R}, \quad (10.5)$$

they can always be written as  $0\mathbf{e}_1 + z\mathbf{e}_2$ . To represent these vectors it is sufficient to remember/store the *coordinate/code*  $z$  of  $\tilde{\mathbf{x}}$  with respect to the  $\mathbf{e}_2$  vector.

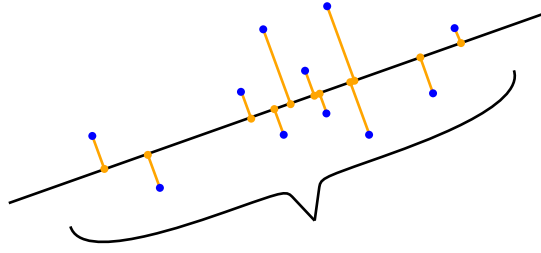
More precisely, the set of  $\tilde{\mathbf{x}}$  vectors (with the standard vector addition and scalar multiplication) forms a vector subspace  $U$  (see Section 2.4) with  $\dim(U) = 1$  because  $U = \text{span}[\mathbf{e}_2]$ .

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  $\tilde{\mathbf{x}} = \mathbf{B}\mathbf{B}^\top \mathbf{x} \in \mathbb{R}^D$  live.

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

Throughout this chapter, we will use the MNIST digits dataset as a re-

**Figure 10.4** PCA finds a lower-dimensional subspace (line) that maintains as much variance (spread of the data) as possible when the data (blue) is projected onto this subspace (orange).



occurring example, which contains 60,000 examples of handwritten digits 0–9. Each digit is a grayscale image of size  $28 \times 28$ , 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.

## 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 Section 6.4.1 we know that the variance is an indicator of the spread of the data, and we can derive 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. Figure 10.4 illustrates this.

Considering the setting discussed in Section 10.1, our aim is to find a matrix  $\mathbf{B}$  (see (10.3)) that retains as much information as possible when compressing data by projecting it onto the subspace spanned by the columns  $\mathbf{b}_1, \dots, \mathbf{b}_M$  of  $\mathbf{B}$ . Retaining most information after data compression is equivalent to capturing the largest amount of variance in the low-dimensional code (Hotelling, 1933).

*Remark. (Centered Data)* For the data covariance matrix in (10.1) we assumed centered data. We can make this assumption without loss of generality: Let us assume that  $\boldsymbol{\mu}$  is the mean of the data. Using the properties of the variance, which we discussed in Section 6.4.3 we obtain

$$\mathbb{V}_z[\mathbf{z}] = \mathbb{V}_x[\mathbf{B}^\top(\mathbf{x} - \boldsymbol{\mu})] = \mathbb{V}_x[\mathbf{B}^\top\mathbf{x} - \mathbf{B}^\top\boldsymbol{\mu}] = \mathbb{V}_x[\mathbf{B}^\top\mathbf{x}], \quad (10.6)$$

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  $\mathbb{E}_z[\mathbf{z}] = \mathbb{E}_x[\mathbf{B}^\top \mathbf{x}] = \mathbf{B}^\top \mathbb{E}_x[\mathbf{x}] = \mathbf{0}$ .  $\diamond$

### 10.2.1 Direction with Maximal Variance

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 variance of the first coordinate  $z_1$  of  $\mathbf{z} \in \mathbb{R}^M$  so that

$$V_1 := \mathbb{V}[z_1] = \frac{1}{N} \sum_{n=1}^N z_{1n}^2 \quad (10.7)$$

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

$$z_{1n} = \mathbf{b}_1^\top \mathbf{x}_n, \quad (10.8)$$

i.e., it is the coordinate of the orthogonal projection of  $\mathbf{x}_n$  onto the one-dimensional subspace spanned by  $\mathbf{b}_1$ , see Section 3.7. We substitute (10.8) into (10.7), which yields

$$V_1 = \frac{1}{N} \sum_{n=1}^N (\mathbf{b}_1^\top \mathbf{x}_n)^2 = \frac{1}{N} \sum_{n=1}^N \mathbf{b}_1^\top \mathbf{x}_n \mathbf{x}_n^\top \mathbf{b}_1 \quad (10.9a)$$

$$= \mathbf{b}_1^\top \left( \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n \mathbf{x}_n^\top \right) \mathbf{b}_1 = \mathbf{b}_1^\top \mathbf{S} \mathbf{b}_1, \quad (10.9b)$$

where  $\mathbf{S}$  is the data covariance matrix defined in (10.1). In (10.9a) we have used the fact that the dot product of two vectors is symmetric with respect to its arguments, that is  $\mathbf{b}_1^\top \mathbf{x}_n = \mathbf{x}_n^\top \mathbf{b}_1$ .

Notice that arbitrarily increasing the magnitude of the vector  $\mathbf{b}_1$  increases  $V_1$ , that is, a vector  $\mathbf{b}_1$  that is two times longer can result in  $V_1$  that is potentially four times larger. Therefore, we restrict all solutions to  $\|\mathbf{b}_1\|^2 = 1$ , 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 the vector  $\mathbf{b}_1$  that points in the direction of maximum variance can be found by the constrained optimization problem

$$\begin{aligned} & \max_{\mathbf{b}_1} \mathbf{b}_1^\top \mathbf{S} \mathbf{b}_1 \\ & \text{subject to } \|\mathbf{b}_1\|^2 = 1. \end{aligned} \quad (10.10)$$

Following Section 7.2, we obtain the Lagrangian

$$\mathcal{L}(\mathbf{b}_1, \lambda) = \mathbf{b}_1^\top \mathbf{S} \mathbf{b}_1 + \lambda_1 (1 - \mathbf{b}_1^\top \mathbf{b}_1) \quad (10.11)$$

The vector  $\mathbf{b}_1$  will be the first column of the matrix  $\mathbf{B}$  and therefore the first of  $M$  orthonormal basis vectors that span the lower-dimensional subspace.

$$\|\mathbf{b}_1\|^2 = 1 \iff \|\mathbf{b}_1\| = 1.$$

to solve this constrained optimization problem. The partial derivatives of  $\mathcal{L}$  with respect to  $\mathbf{b}_1$  and  $\lambda_1$  are

$$\frac{\partial \mathcal{L}}{\partial \mathbf{b}_1} = 2\mathbf{b}_1^\top \mathbf{S} - 2\lambda_1 \mathbf{b}_1^\top \quad (10.12)$$

$$\frac{\partial \mathcal{L}}{\partial \lambda_1} = 1 - \mathbf{b}_1^\top \mathbf{b}_1, \quad (10.13)$$

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

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

$$\mathbf{b}_1^\top \mathbf{b}_1 = 1. \quad (10.15)$$

By comparing with the definition of an eigenvalue decomposition (Section 4.4), we see that  $\mathbf{b}_1$  is an eigenvector of the data covariance matrix  $\mathbf{S}$ , and the Lagrange multiplier  $\lambda_1$  plays the role of the corresponding eigenvalue. This eigenvector property (10.14) allows us to rewrite our variance objective (10.10) as

$$V_1 = \mathbf{b}_1^\top \mathbf{S}\mathbf{b}_1 = \lambda_1 \mathbf{b}_1^\top \mathbf{b}_1 = \lambda_1, \quad (10.16)$$

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 associated with 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

$$\tilde{\mathbf{x}}_n = \mathbf{b}_1 z_{1n} = \mathbf{b}_1 \mathbf{b}_1^\top \mathbf{x}_n \in \mathbb{R}^D \quad (10.17)$$

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  $\mathbf{b}_1 \in \mathbb{R}^D$ .  $\diamond$

### 10.2.2 $M$ -dimensional Subspace with Maximal Variance

Assume we have found the first  $m - 1$  principal components as the  $m - 1$  eigenvectors of  $\mathbf{S}$  that are associated with the largest  $m - 1$  eigenvalues. Since  $\mathbf{S}$  is symmetric, these eigenvectors form an ONB of an  $(m - 1)$ -dimensional subspace of  $\mathbb{R}^D$ . Generally, the  $m$ th principal component can be found by subtracting the effect of the first  $m - 1$  principal components  $\mathbf{b}_1, \dots, \mathbf{b}_{m-1}$  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,

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  $\text{span}[\mathbf{b}_1]$ .  
principal component

similar to (10.17), so that we arrive at the new data matrix

$$\hat{\mathbf{X}} := \mathbf{X} - \sum_{i=1}^{m-1} \mathbf{b}_i \mathbf{b}_i^\top \mathbf{X}, \quad (10.18)$$

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  $\hat{\mathbf{X}} := [\hat{\mathbf{x}}_1, \dots, \hat{\mathbf{x}}_N] \in \mathbb{R}^{D \times N}$  in (10.18) 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

$$V_m = \mathbb{V}[z_m] = \frac{1}{N} \sum_{n=1}^N z_{mn}^2 = \frac{1}{N} \sum_{n=1}^N (\mathbf{b}_m^\top \mathbf{x}_n)^2 = \mathbf{b}_m^\top \hat{\mathbf{S}} \mathbf{b}_m, \quad (10.19)$$

subject to  $\|\mathbf{b}_m\|^2 = 1$ , where we followed the same steps as in (10.9b) 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 optimization problem and discover that the optimal solution  $\mathbf{b}_m$  is the eigenvector of  $\hat{\mathbf{S}}$  that is associated with the largest eigenvalue of  $\hat{\mathbf{S}}$ .

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

$$\hat{\mathbf{S}} = \frac{1}{N} \sum_{n=1}^N \hat{\mathbf{x}}_n \hat{\mathbf{x}}_n^\top \stackrel{(10.18)}{=} \frac{1}{N} \sum_{n=1}^N \left( \mathbf{x}_n - \sum_{i=1}^{m-1} \mathbf{b}_i \mathbf{b}_i^\top \mathbf{x}_n \right) \left( \mathbf{x}_n - \sum_{i=1}^{m-1} \mathbf{b}_i \mathbf{b}_i^\top \mathbf{x}_n \right)^\top \quad (10.20a)$$

$$= \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n \mathbf{x}_n^\top - 2 \mathbf{x}_n \mathbf{x}_n^\top \sum_{i=1}^{m-1} \mathbf{b}_i \mathbf{b}_i^\top + \sum_{i=1}^{m-1} \mathbf{b}_i \mathbf{b}_i^\top \mathbf{x}_n \mathbf{x}_n^\top \sum_{i=1}^{m-1} \mathbf{b}_i \mathbf{b}_i^\top, \quad (10.20b)$$

where we exploited the symmetries  $\mathbf{x}_n^\top \mathbf{b}_i = \mathbf{b}_i^\top \mathbf{x}_n$  and  $\mathbf{b}_i \mathbf{x}_n^\top = \mathbf{x}_n \mathbf{b}_i^\top$  to summarize

$$-\mathbf{x}_n \mathbf{x}_n^\top \sum_{i=1}^{m-1} \mathbf{b}_i \mathbf{b}_i^\top - \sum_{i=1}^{m-1} \mathbf{b}_i \mathbf{b}_i^\top \mathbf{x}_n \mathbf{x}_n^\top = -2 \mathbf{x}_n \mathbf{x}_n^\top \sum_{i=1}^{m-1} \mathbf{b}_i \mathbf{b}_i^\top. \quad (10.21)$$

If we take a vector  $\mathbf{b}_m$  with  $\|\mathbf{b}_m\| = 1$  that is orthogonal to all  $\mathbf{b}_1, \dots, \mathbf{b}_{m-1}$  and right-multiply  $\mathbf{b}_m$  to  $\hat{\mathbf{S}}$  in (10.20b) we obtain

$$\hat{\mathbf{S}} \mathbf{b}_m = \frac{1}{N} \sum_{n=1}^N \hat{\mathbf{x}}_n \hat{\mathbf{x}}_n^\top \mathbf{b}_m = \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n \mathbf{x}_n^\top \mathbf{b}_m = \mathbf{S} \mathbf{b}_m = \lambda_m \mathbf{b}_m. \quad (10.22)$$

Here we applied the orthogonality property  $\mathbf{b}_i^\top \mathbf{b}_m = 0$  for  $i = 1, \dots, m-1$  (all terms involving sums up to  $m-1$  vanish). Equation (10.22) reveals that  $\mathbf{b}_m$  is an eigenvector of both  $\hat{\mathbf{S}}$  and the original data covariance matrix  $\mathbf{S}$ . In the former case,  $\lambda_m$  is the largest eigenvalue, in latter case,  $\lambda_m$  is the  $m$ th largest eigenvalue. With this the variance of the data projected onto the  $m$ th principal component is

$$V_m = \mathbf{b}_m^\top \mathbf{S} \mathbf{b}_m \stackrel{(10.22)}{=} \lambda_m \mathbf{b}_m^\top \mathbf{b}_m = \lambda_m \quad (10.23)$$

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 is associated with the corresponding eigenvectors of the data covariance matrix.

Overall, 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  $\mathbf{B}$  in (10.3) as the  $M$  eigenvectors of the data covariance matrix  $\mathbf{S}$  that are associated with the  $M$  largest eigenvalues. The maximum amount of variance PCA can capture with the first  $M$  principal components is

$$V_M = \sum_{m=1}^M \lambda_m, \quad (10.24)$$

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

$$J_M := \sum_{j=M+1}^D \lambda_j = V_D - V_M. \quad (10.25)$$

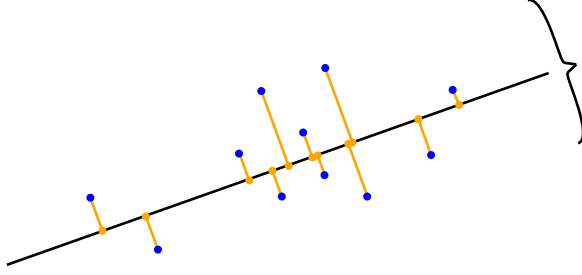
Instead of these absolute quantities, we can also define the relative amount of variance captured as  $\frac{V_M}{V_D}$ , and the relative amount of variance lost by compression as  $1 - \frac{V_M}{V_D}$ .

### 10.3 Projection Perspective

In the following, we will derive PCA as an algorithm for linear dimensionality reduction that directly minimizes the average reconstruction error. This perspective allows us to interpret PCA as an algorithm that implements an optimal linear auto-encoder. 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 10.5 illustrates this setting.





**Figure 10.5**  
Illustration of the projection approach to PCA. We aim to find a lower-dimensional subspace (line) so that the difference vector between projected (orange) and original (blue) data is as short as possible.

### 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.

*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,  $U \cap U^\perp = \{\mathbf{0}\}$  so that any vector  $\mathbf{x} \in V$  can be (uniquely) decomposed into

$$\mathbf{x} = \sum_{m=1}^M \lambda_m \mathbf{b}_m + \sum_{j=1}^{D-M} \psi_j \mathbf{b}_j^\perp, \quad \lambda_m, \psi_j \in \mathbb{R}, \quad (10.26)$$

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  $U^\perp$ .

From Section 2.5 we know that for a basis  $(\mathbf{b}_1, \dots, \mathbf{b}_D)$  of  $\mathbb{R}^D$  any  $\mathbf{x} \in \mathbb{R}^D$  can be written as a linear combination of the basis vectors of  $\mathbb{R}^D$ , i.e.,

$$\mathbf{x} = \sum_{d=1}^D \zeta_d \mathbf{b}_d = \sum_{m=1}^M \zeta_m \mathbf{b}_m + \sum_{j=M+1}^D \zeta_j \mathbf{b}_j \quad (10.27)$$

for suitable coordinates  $\zeta_d \in \mathbb{R}$ .

We are interested in finding vectors  $\tilde{\mathbf{x}} \in \mathbb{R}^D$ , which live in lower-dimensional subspace  $U \subseteq \mathbb{R}^D$ ,  $\dim(U) = M$ , so that

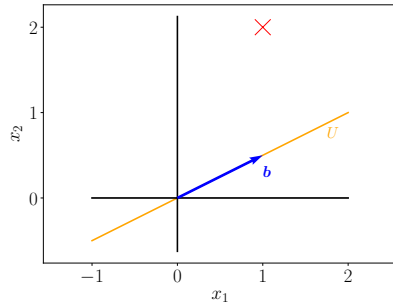
$$\tilde{\mathbf{x}} = \sum_{m=1}^M z_m \mathbf{b}_m \in U \subseteq \mathbb{R}^D \quad (10.28)$$

is as similar to  $\mathbf{x}$  as possible. Note that at this point we need to assume that the coordinates  $z_m$  of  $\tilde{\mathbf{x}}$  and  $\zeta_m$  of  $\mathbf{x}$  are not identical.

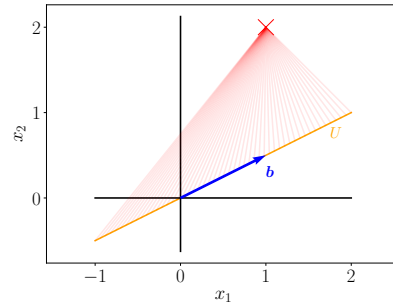
In the following, we use exactly this kind of representation of  $\tilde{\mathbf{x}}$  to find optimal coordinates  $\mathbf{z}$  and basis vectors  $\mathbf{b}_1, \dots, \mathbf{b}_M$  such that  $\tilde{\mathbf{x}}$  is as similar to the original data point  $\mathbf{x}$ , i.e., we aim to minimize the (Euclidean) distance  $\|\mathbf{x} - \tilde{\mathbf{x}}\|$ . Figure 10.6 illustrates this setting. Without loss of generality, we assume that the dataset  $\mathcal{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$ ,  $\mathbf{x}_n \in \mathbb{R}^D$ , is centered at  $\mathbf{0}$ , i.e.,  $\mathbb{E}[\mathcal{X}] = \mathbf{0}$ . Without the zero-mean assumption, we would

orthogonal  
complement

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 with respect to the standard basis of  $\mathbb{R}^3$ .



(a) Setting.

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

**Figure 10.6**  
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  $\tilde{\mathbf{x}}$ .

arrive at exactly the same solution but the notation would be substantially more cluttered.

principal subspace

We are interested in finding the best linear projection of  $\mathcal{X}$  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*. The projections of the data points 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  $\mathbf{z}_n := [z_{1n}, \dots, z_{Mn}]^\top \in \mathbb{R}^M$  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.

reconstruction error

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, 1901b)

$$J_M := \frac{1}{N} \sum_{n=1}^N \|\mathbf{x}_n - \tilde{\mathbf{x}}_n\|^2, \quad (10.30)$$

where we make it explicit that the dimension of the subspace onto which we project the data is  $M$ . In order to find this optimal linear projection, we need to find the orthonormal basis of the principal subspace and the coordinates  $\mathbf{z}_n \in \mathbb{R}^M$  of the projections with respect to this basis.

To find the coordinates  $\mathbf{z}_n$  and the ONB of the principal subspace we follow a two-step approach. First, we optimize the coordinates  $\mathbf{z}_n$  for a given ONB  $(\mathbf{b}_1, \dots, \mathbf{b}_M)$ ; second, we find the optimal ONB.

### 10.3.2 Finding Optimal Coordinates

Let us start by finding the optimal coordinates  $z_{1n}, \dots, z_{Mn}$  of the projections  $\tilde{\mathbf{x}}_n$  for  $n = 1, \dots, N$ . Consider Figure 10.6(b) where the principal

subspace is spanned by a single vector  $\mathbf{b}$ . Geometrically speaking, finding the optimal coordinates  $z$  corresponds to finding the representation of the linear projection  $\tilde{\mathbf{x}}$  with respect to  $\mathbf{b}$  that minimizes the distance between  $\tilde{\mathbf{x}} - \mathbf{x}$ . From Figure 10.6(b) it is clear that this will be the orthogonal projection, and in the following we will show exactly this.

We assume an ONB  $(\mathbf{b}_1, \dots, \mathbf{b}_M)$  of  $U \subseteq \mathbb{R}^D$ . To find the optimal coordinates  $z_m$  with respect to this basis, we require the partial derivatives

$$\frac{\partial J_M}{\partial z_{in}} = \frac{\partial J_M}{\partial \tilde{\mathbf{x}}_n} \frac{\partial \tilde{\mathbf{x}}_n}{\partial z_{in}}, \quad (10.31a)$$

$$\frac{\partial J_M}{\partial \tilde{\mathbf{x}}_n} = -\frac{2}{N}(\mathbf{x}_n - \tilde{\mathbf{x}}_n)^\top \in \mathbb{R}^{1 \times D}, \quad (10.31b)$$

$$\frac{\partial \tilde{\mathbf{x}}_n}{\partial z_{in}} \stackrel{(10.29)}{=} \frac{\partial}{\partial z_{in}} \left( \sum_{m=1}^M z_{mn} \mathbf{b}_m \right) = \mathbf{b}_i \quad (10.31c)$$

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

$$\frac{\partial J_M}{\partial z_{in}} \stackrel{(10.31b)}{\stackrel{(10.31c)}}{=} -\frac{2}{N}(\mathbf{x}_n - \tilde{\mathbf{x}}_n)^\top \mathbf{b}_i \stackrel{(10.29)}{=} -\frac{2}{N} \left( \mathbf{x}_n - \sum_{m=1}^M z_{mn} \mathbf{b}_m \right)^\top \mathbf{b}_i \quad (10.32a)$$

$$\stackrel{\text{ONB}}{=} -\frac{2}{N}(\mathbf{x}_n^\top \mathbf{b}_i - z_{in} \underbrace{\mathbf{b}_i^\top \mathbf{b}_i}_{=1}) = -\frac{2}{N}(\mathbf{x}_n^\top \mathbf{b}_i - z_{in}). \quad (10.32b)$$

Setting this partial derivative to 0 yields immediately the optimal coordinates

$$z_{in} = \mathbf{x}_n^\top \mathbf{b}_i = \mathbf{b}_i^\top \mathbf{x}_n \quad (10.33)$$

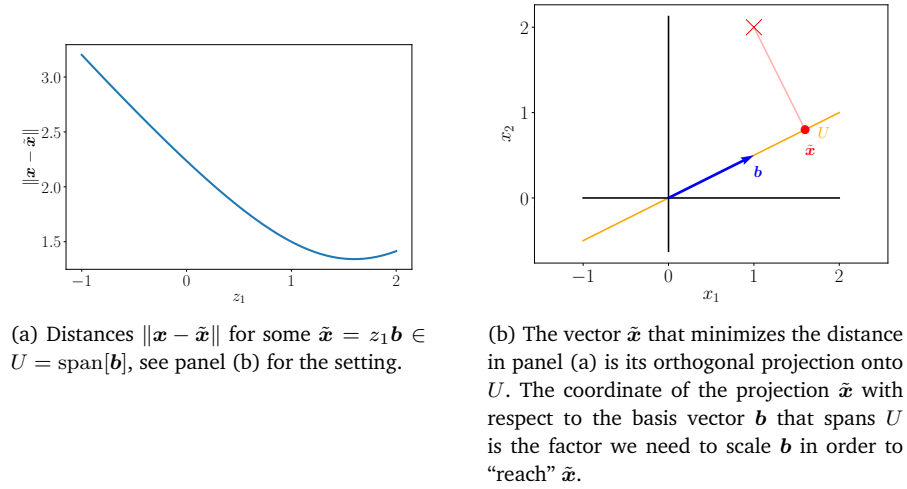
for  $i = 1, \dots, M$  and  $n = 1, \dots, N$ . This means, the optimal coordinates  $z_{in}$  of the projection  $\tilde{\mathbf{x}}_n$  are the coordinates of the orthogonal projection (see Section 3.7) of the original data point  $\mathbf{x}_n$  onto the one-dimensional subspace that is spanned by  $\mathbf{b}_i$ . Consequently:

- The optimal linear 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 (10.30).
- The coordinates  $\zeta_m$  of  $\mathbf{x}$  in (10.27) and the coordinates  $z_m$  of  $\tilde{\mathbf{x}}$  in (10.28) must be identical for  $m = 1, \dots, M$  in PCA since  $U^\perp = \text{span}[\mathbf{b}_{M+1}, \dots, \mathbf{b}_D]$  is the orthogonal complement of  $U = \text{span}[\mathbf{b}_1, \dots, \mathbf{b}_M]$ .

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.

*Remark* (Orthogonal Projections with Orthonormal Basis Vectors). Let us briefly recap orthogonal projections from Section 3.7. If  $(\mathbf{b}_1, \dots, \mathbf{b}_D)$  is an

**Figure 10.7**  
Optimal projection  
of a vector  $\mathbf{x} \in \mathbb{R}^2$   
onto a  
one-dimensional  
subspace  
(continuation from  
Figure 10.6).  
(a) Distances  
 $\|\mathbf{x} - \tilde{\mathbf{x}}\|$  for some  
 $\tilde{\mathbf{x}} \in U$ .  
(b) Orthogonal  
projection and  
optimal coordinates.



orthonormal basis of  $\mathbb{R}^D$  then

$$\tilde{\mathbf{x}} = \mathbf{b}_j \underbrace{(\mathbf{b}_j^\top \mathbf{b}_j)^{-1}}_{=1} \mathbf{b}_j^\top \mathbf{x} = \mathbf{b}_j \mathbf{b}_j^\top \mathbf{x} \in \mathbb{R}^D \quad (10.34)$$

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

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 10.7 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

$$\tilde{\mathbf{x}} = \mathbf{B} \underbrace{(\mathbf{B}^\top \mathbf{B})^{-1}}_{=\mathbf{I}} \mathbf{B}^\top \mathbf{x} = \mathbf{B} \mathbf{B}^\top \mathbf{x}, \quad (10.35)$$

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 3.7.

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.  $\diamond$

So far, we showed that for a given ONB we can find the optimal coordinates of  $\tilde{\mathbf{x}}$  by an orthogonal projection onto the principal subspace. In the following, we will determine what the best basis is.

### 10.3.3 Finding the Basis of the Principal Subspace

To determine the basis vectors  $\mathbf{b}_1, \dots, \mathbf{b}_M$  of the principal subspace, we rephrase the loss function (10.30) 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

$$\tilde{\mathbf{x}}_n = \sum_{m=1}^M z_{mn} \mathbf{b}_m \stackrel{(10.33)}{=} \sum_{m=1}^M (\mathbf{x}_n^\top \mathbf{b}_m) \mathbf{b}_m. \quad (10.36)$$

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

$$\tilde{\mathbf{x}}_n = \left( \sum_{m=1}^M \mathbf{b}_m \mathbf{b}_m^\top \right) \mathbf{x}_n. \quad (10.37)$$

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

$$\mathbf{x}_n = \sum_{d=1}^D z_{dn} \mathbf{b}_d \stackrel{(10.33)}{=} \sum_{d=1}^D (\mathbf{x}_n^\top \mathbf{b}_d) \mathbf{b}_d = \left( \sum_{d=1}^D \mathbf{b}_d \mathbf{b}_d^\top \right) \mathbf{x}_n \quad (10.38a)$$

$$= \left( \sum_{m=1}^M \mathbf{b}_m \mathbf{b}_m^\top \right) \mathbf{x}_n + \left( \sum_{j=M+1}^D \mathbf{b}_j \mathbf{b}_j^\top \right) \mathbf{x}_n, \quad (10.38b)$$

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

$$\mathbf{x}_n - \tilde{\mathbf{x}}_n = \left( \sum_{j=M+1}^D \mathbf{b}_j \mathbf{b}_j^\top \right) \mathbf{x}_n \quad (10.39a)$$

$$= \sum_{j=M+1}^D (\mathbf{x}_n^\top \mathbf{b}_j) \mathbf{b}_j. \quad (10.39b)$$

This means the difference is exactly the projection of the data point onto the orthogonal complement of the principal subspace: We identify the matrix  $\sum_{j=M+1}^D \mathbf{b}_j \mathbf{b}_j^\top$  in (10.39a) as the projection matrix that performs this projection. This also means the displacement vector  $\mathbf{x}_n - \tilde{\mathbf{x}}_n$  lies in the subspace that is orthogonal to the principal subspace as illustrated in Figure 10.8.

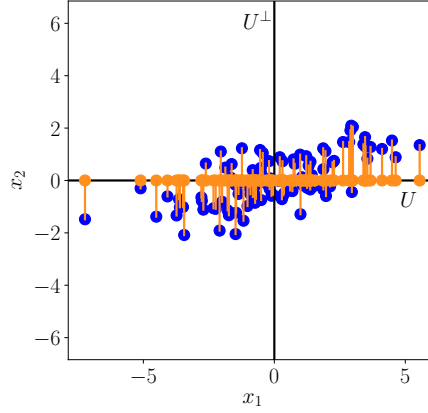
*Remark (Low-Rank Approximation).* In (10.39a), we saw that the projection matrix, which projects  $\mathbf{x}$  onto  $\tilde{\mathbf{x}}$ , is given by

$$\sum_{m=1}^M \mathbf{b}_m \mathbf{b}_m^\top = \mathbf{B} \mathbf{B}^\top. \quad (10.40)$$

By construction as a sum of rank-one matrices  $\mathbf{b}_m \mathbf{b}_m^\top$  we see that  $\mathbf{B} \mathbf{B}^\top$

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

**Figure 10.8**  
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  $U_1$ .



is symmetric and has rank  $M$ . Therefore, the average squared reconstruction error can also be written as

$$\frac{1}{N} \sum_{n=1}^N \|\mathbf{x}_n - \tilde{\mathbf{x}}_n\|^2 = \frac{1}{N} \sum_{n=1}^N \|\mathbf{x}_n - \mathbf{B}\mathbf{B}^\top \mathbf{x}_n\|^2 \quad (10.41a)$$

$$= \frac{1}{N} \sum_{n=1}^N \|(I - \mathbf{B}\mathbf{B}^\top) \mathbf{x}_n\|^2. \quad (10.41b)$$

5785 Finding orthonormal basis vectors  $\mathbf{b}_1, \dots, \mathbf{b}_M$  so that the difference be-  
 5786 tween the original data  $\mathbf{x}_n$  and their projections  $\tilde{\mathbf{x}}_n$ ,  $n = 1, \dots, N$ , is  
 5787 minimized is equivalent to finding the best rank- $M$  approximation  $\mathbf{B}\mathbf{B}^\top$   
 5788 of the identity matrix  $I$ , see Section 4.6.  $\diamond$

Now, we have all the tools to reformulate the loss function (10.30).

$$J_M = \frac{1}{N} \sum_{n=1}^N \|\mathbf{x}_n - \tilde{\mathbf{x}}_n\|^2 \stackrel{(10.39b)}{=} \frac{1}{N} \sum_{n=1}^N \left\| \sum_{j=M+1}^D (\mathbf{b}_j^\top \mathbf{x}_n) \mathbf{b}_j \right\|^2. \quad (10.42)$$

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

$$J_M = \frac{1}{N} \sum_{n=1}^N \sum_{j=M+1}^D (\mathbf{b}_j^\top \mathbf{x}_n)^2 = \frac{1}{N} \sum_{n=1}^N \sum_{j=M+1}^D \mathbf{b}_j^\top \mathbf{x}_n \mathbf{b}_j^\top \mathbf{x}_n \quad (10.43a)$$

$$= \frac{1}{N} \sum_{n=1}^N \sum_{j=M+1}^D \mathbf{b}_j^\top \mathbf{x}_n \mathbf{x}_n^\top \mathbf{b}_j, \quad (10.43b)$$

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

$$J_M = \sum_{j=M+1}^D \mathbf{b}_j^\top \underbrace{\left( \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n \mathbf{x}_n^\top \right)}_{=: \mathbf{S}} \mathbf{b}_j = \sum_{j=M+1}^D \mathbf{b}_j^\top \mathbf{S} \mathbf{b}_j \quad (10.44a)$$

$$= \sum_{j=M+1}^D \text{tr}(\mathbf{b}_j^\top \mathbf{S} \mathbf{b}_j) = \sum_{j=M+1}^D \text{tr}(\mathbf{S} \mathbf{b}_j \mathbf{b}_j^\top) = \text{tr} \left( \underbrace{\left( \sum_{j=M+1}^D \mathbf{b}_j \mathbf{b}_j^\top \right)}_{\text{projection matrix}} \mathbf{S} \right), \quad (10.44b)$$

where we exploited the property that the trace operator  $\text{tr}(\cdot)$ , see (4.18), is linear and invariant to cyclic permutations of its arguments. Since we assumed that our dataset is centered, i.e.,  $\mathbb{E}[\mathcal{X}] = \mathbf{0}$ , we identify  $\mathbf{S}$  as the data covariance matrix. We see that the projection matrix in (10.44b) 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 (10.44a) 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 omit a derivation that is identical to the one in Section 10.2 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

$$J_M = \sum_{j=M+1}^D \lambda_j, \quad (10.45)$$

where  $\lambda_j$  are the eigenvalues of the data covariance matrix. Therefore, to minimize (10.45) 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 are associated with 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 10.9**  
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.

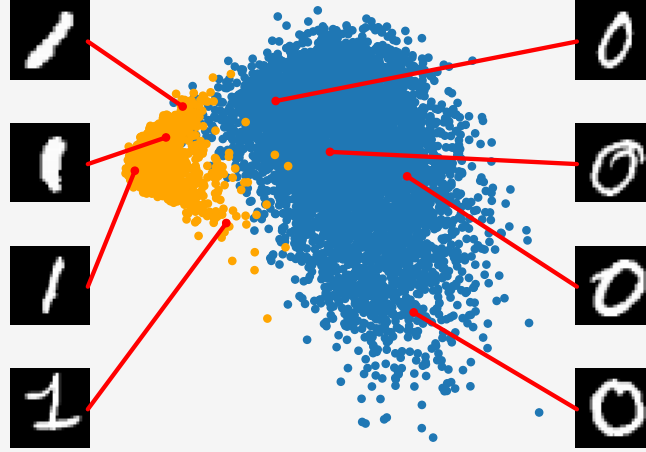


Figure 10.9 visualizes the training data of the MNIST digits ‘0’ and ‘1’ embedded in the vector subspace spanned by the first two principal components. We can see a relatively clear separation between ‘0’s (blue dots) and ‘1’s (orange dots), and we can see the variation within each individual cluster.

**10.4 Eigenvector Computation and Low-Rank Approximations**

obtained the basis of the principal subspace as the eigenvectors that are associated with the largest eigenvalues of the data covariance matrix

$$\mathbf{S} = \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n \mathbf{x}_n^\top = \frac{1}{N} \mathbf{X} \mathbf{X}^\top, \quad (10.46)$$

$$\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_N] \in \mathbb{R}^{D \times N}. \quad (10.47)$$

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

$$\underbrace{\mathbf{X}}_{D \times N} = \underbrace{\mathbf{U}}_{D \times D} \underbrace{\mathbf{\Sigma}}_{D \times N} \underbrace{\mathbf{V}^\top}_{N \times N}, \quad (10.48)$$

where  $\mathbf{U} \in \mathbb{R}^{D \times D}$  and  $\mathbf{V}^\top \in \mathbb{R}^{N \times N}$  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

$$\mathbf{S} = \frac{1}{N} \mathbf{X} \mathbf{X}^\top = \frac{1}{N} \mathbf{U} \mathbf{\Sigma} \underbrace{\mathbf{V}^\top \mathbf{V}}_{=\mathbf{I}_N} \mathbf{\Sigma}^\top \mathbf{U}^\top = \frac{1}{N} \mathbf{U} \mathbf{\Sigma} \mathbf{\Sigma}^\top \mathbf{U}^\top. \quad (10.49)$$

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  $\lambda_d$  of  $\mathbf{S}$  are related to the singular values of  $\mathbf{X}$  via

$$\lambda_d = \frac{\sigma_d^2}{N}. \quad (10.50)$$

The columns of  $\mathbf{U}$  are the eigenvectors of  $\mathbf{S}$ .

#### 10.4.1 PCA using Low-rank Matrix Approximations

To maximize the variance of the projected data (or minimize the average squared reconstruction error), PCA chooses the columns of  $\mathbf{U}$  in (10.49) to be the eigenvectors that are associated with the  $M$  largest eigenvalues of the data covariance matrix  $\mathbf{S}$  so that we identify  $\mathbf{U}$  as the projection matrix  $\mathbf{B}$  in (10.3), 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

$$\tilde{\mathbf{X}}_M := \operatorname{argmin}_{\operatorname{rk}(\mathbf{A}) \leq M} \|\mathbf{X} - \mathbf{A}\|_2 \in \mathbb{R}^{D \times N} \quad (10.51)$$

of  $\mathbf{X}$ , where  $\|\cdot\|_2$  is the spectral norm defined in (4.110). The Eckart-Young Theorem states that  $\tilde{\mathbf{X}}_M$  is given by truncating the SVD at the top- $M$  singular value. In other words, we obtain

$$\tilde{\mathbf{X}}_M = \underbrace{\mathbf{U}_M}_{D \times M} \underbrace{\mathbf{\Sigma}_M}_{M \times M} \underbrace{\mathbf{V}_M^\top}_{M \times N} \in \mathbb{R}^{D \times N} \quad (10.52)$$

with orthogonal matrices  $\mathbf{U}_M := [\mathbf{u}_1, \dots, \mathbf{u}_M] \in \mathbb{R}^{D \times M}$  and  $\mathbf{V}_M := [\mathbf{v}_1, \dots, \mathbf{v}_M] \in \mathbb{R}^{N \times M}$  and a diagonal matrix  $\mathbf{\Sigma}_M \in \mathbb{R}^{M \times M}$  whose diagonal entries are the  $M$  largest singular values of  $\mathbf{X}$ .

Eckart-Young Theorem

#### 10.4.2 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

5829  $4 \times 4$  this is not possible because we would need to find the roots of a poly-  
 5830 nomial of degree 5 or higher. However, the Abel-Ruffini theorem (Ruffini,  
 5831 1799; Abel, 1826) states that there exists no algebraic solution to this  
 5832 problem for polynomials of degree 5 or more. Therefore, in practice, we  
 np.linalg.eigh 5833 solve for eigenvalues or singular values using iterative methods, which are  
 or 5834 implemented in all modern packages for linear algebra.  
 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  $\mathbf{x}_0$  that is not in the null space of  $\mathbf{S}$  and follows the iteration

power iteration

If  $\mathbf{S}$  is invertible, it is sufficient to ensure that  $\mathbf{x}_0 \neq \mathbf{0}$ .

$$\mathbf{x}_{k+1} = \frac{\mathbf{S}\mathbf{x}_k}{\|\mathbf{S}\mathbf{x}_k\|}, \quad k = 0, 1, \dots \quad (10.53)$$

5835 This means the vector  $\mathbf{x}_k$  is multiplied by  $\mathbf{S}$  in every iteration and then  
 5836 normalized, i.e., we always have  $\|\mathbf{x}_k\| = 1$ . This sequence of vectors con-  
 5837 verges to the eigenvector associated with the largest eigenvalue of  $\mathbf{S}$ . The  
 5838 original Google PageRank algorithm (Page et al., 1999) uses such an al-  
 5839 gorithm for ranking web pages based on their hyperlinks.

## 10.5 PCA in High Dimensions

5840

5841 In order to do PCA, we need to compute the data covariance matrix. In  $D$   
 5842 dimensions, the data covariance matrix is a  $D \times D$  matrix. Computing the  
 5843 eigenvalues and eigenvectors of this matrix is computationally expensive  
 5844 as it scales cubically in  $D$ . Therefore, PCA, as we discussed earlier, will be  
 5845 infeasible in very high dimensions. For example, if our  $\mathbf{x}_n$  are images with  
 5846 10,000 pixels (e.g.,  $100 \times 100$  pixel images), we would need to compute  
 5847 the eigendecomposition of a  $10,000 \times 10,000$  covariance matrix. In the  
 5848 following, we provide a solution to this problem for the case that we have  
 5849 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

$$\mathbf{S} = \frac{1}{N} \mathbf{X} \mathbf{X}^\top \in \mathbb{R}^{D \times D}, \quad (10.54)$$

5850 where  $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_N]$  is a  $D \times N$  matrix whose columns are the data  
 5851 points.

5852 We now assume that  $N \ll D$ , i.e., the number of data points is smaller  
 5853 than the dimensionality of the data. If there are no duplicate data points

the rank of the covariance matrix  $\mathbf{S}$  is  $N$ , so 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

$$\mathbf{S}\mathbf{b}_m = \lambda_m \mathbf{b}_m, \quad m = 1, \dots, M, \quad (10.55)$$

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 (10.54), we obtain

$$\mathbf{S}\mathbf{b}_m = \frac{1}{N} \mathbf{X} \mathbf{X}^\top \mathbf{b}_m = \lambda_m \mathbf{b}_m. \quad (10.56)$$

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

$$\frac{1}{N} \underbrace{\mathbf{X}^\top \mathbf{X}}_{N \times N} \underbrace{\mathbf{X}^\top \mathbf{b}_m}_{=: \mathbf{c}_m} = \lambda_m \mathbf{X}^\top \mathbf{b}_m \iff \frac{1}{N} \mathbf{X}^\top \mathbf{X} \mathbf{c}_m = \lambda_m \mathbf{c}_m, \quad (10.57)$$

and we get a new eigenvector/eigenvalue equation:  $\lambda_m$  remains eigenvalue, which confirms our results from Section 4.5.3 that the non-zero eigenvalues of  $\mathbf{X} \mathbf{X}^\top$  equal the non-zero eigenvalues of  $\mathbf{X}^\top \mathbf{X}$ . We obtain the eigenvector of the matrix  $\frac{1}{N} \mathbf{X}^\top \mathbf{X} \in \mathbb{R}^{N \times N}$  associated with  $\lambda_m$  as  $\mathbf{c}_m := \mathbf{X}^\top \mathbf{b}_m$ . 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

$$\underbrace{\frac{1}{N} \mathbf{X} \mathbf{X}^\top}_{\mathbf{S}} \mathbf{X} \mathbf{c}_m = \lambda_m \mathbf{X} \mathbf{c}_m \quad (10.58)$$

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 10.6 we need to normalize the eigenvectors  $\mathbf{X} \mathbf{c}_m$  of  $\mathbf{S}$  so that they have norm 1.  $\diamond$

## 10.6 Key Steps of PCA in Practice

In the following, we will go through the individual steps of PCA using a running example, which is summarized in Figure 10.10. We are given a

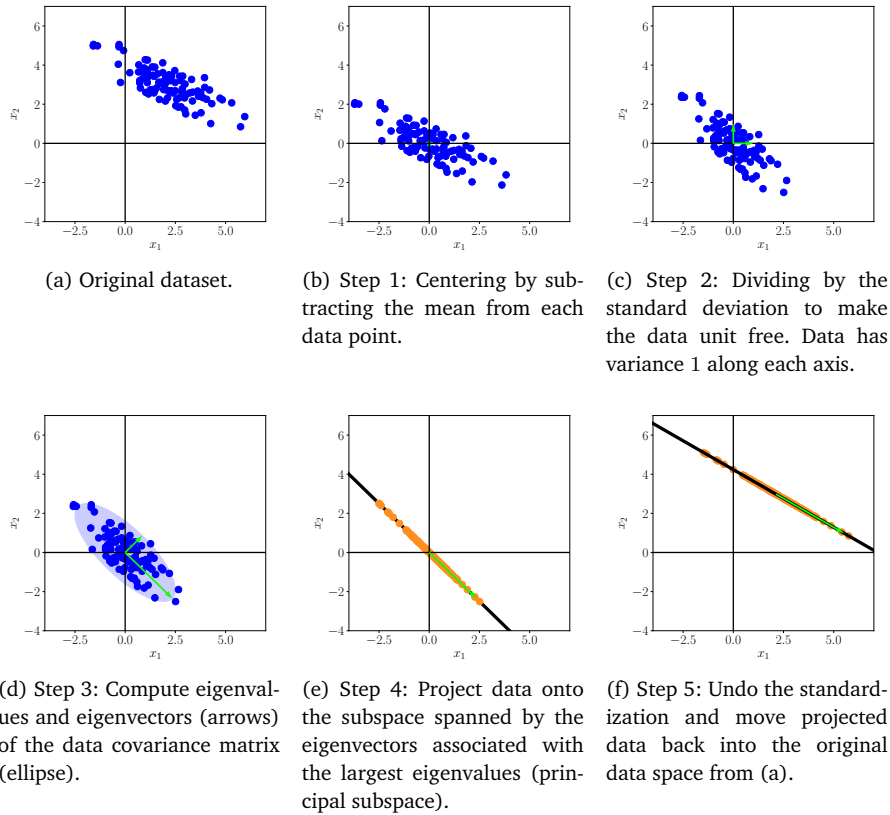


Figure 10.10 Steps of PCA.

two-dimensional data set (Figure 10.10(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  $\mu$  of the dataset and subtracting it from every single data point. This ensures that the data set has mean 0 (Figure 10.10(b)). Mean subtraction is not strictly necessary but reduces the risk of numerical problems.
2. **Standardization** Divide the data points by the standard deviation  $\sigma_d$  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 10.10(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. Since the covariance matrix is symmetric, the eigenvectors form an orthogonal basis. In Figure 10.10(d), the eigenvectors are scaled by the magnitude of the 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  $\mu_d$  and standard deviation  $\sigma_d$  of the training data in the  $d$ th dimension, respectively, so that

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

where  $x_*^{(d)}$  is the  $d$ th component of  $\mathbf{x}_*$ . We obtain the projection 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 are associated with 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

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

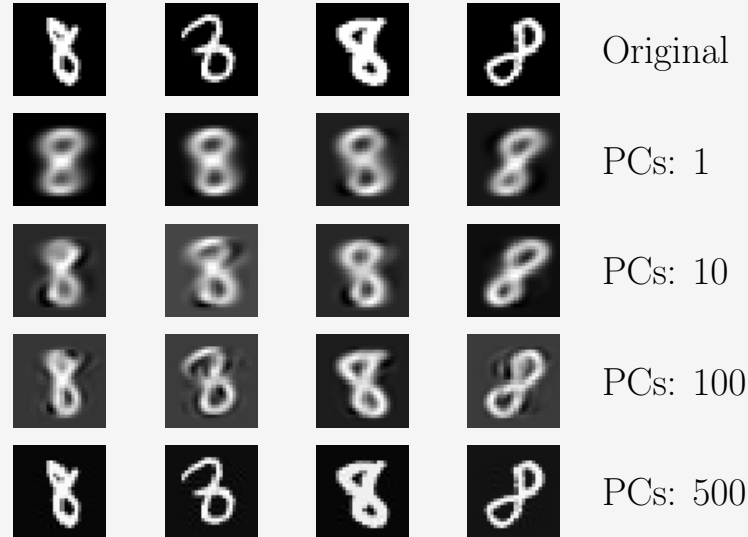
Figure 10.10(f) illustrates the projection in the original data format.

### 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  $28 \times 28$ , 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 10.11. The first row of Figure 10.11 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 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.

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

**Figure 10.11** Effect of increasing number of principal components on reconstruction.



**Figure 10.12** Average squared reconstruction error as a function of the number of principal components.

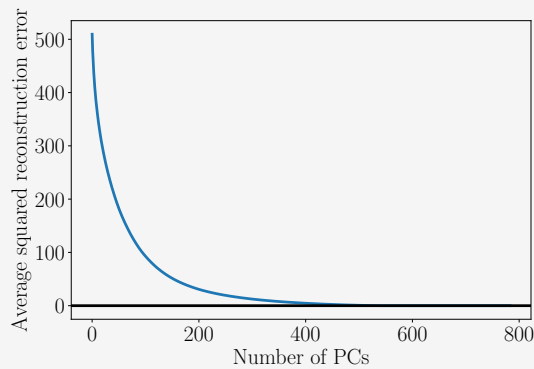


Figure 10.12 shows the average squared reconstruction error, which is

$$\frac{1}{N} \sum_{n=1}^N \|\mathbf{x}_n - \tilde{\mathbf{x}}_n\|^2 = \sum_{i=M+1}^D \lambda_i, \quad (10.62)$$

as a function of the number  $M$  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’ (some pixels around the boundaries show no variation across the dataset as they are always black).

## 10.7 Latent Variable Perspective

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 8.5,
- 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,
- give us a principled way to extend the model, 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.

By introducing a continuous-valued latent variable  $\mathbf{z} \in \mathbb{R}^M$  it is possible to phrase PCA as a probabilistic latent-variable model. Tipping and Bishop (1999) proposed this latent-variable model as *Probabilistic PCA* (PPCA). PPCA addresses most of the issues above, 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 estimation in a noise-free setting.

Probabilistic PCA

### 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

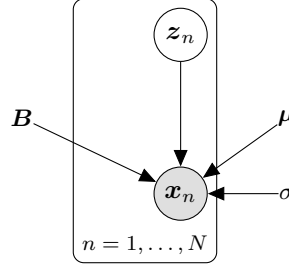
$$\mathbf{x} = \mathbf{B}\mathbf{z} + \boldsymbol{\mu} + \boldsymbol{\epsilon} \in \mathbb{R}^D, \quad (10.63)$$

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

$$p(\mathbf{x} | \mathbf{z}, \mathbf{B}, \boldsymbol{\mu}, \sigma^2) = \mathcal{N}(\mathbf{x} | \mathbf{B}\mathbf{z} + \boldsymbol{\mu}, \sigma^2 \mathbf{I}). \quad (10.64)$$

**Figure 10.13**

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  $\sigma$  are shared across the dataset.



Overall, PPCA induces the following generative process:

$$\mathbf{z}_n \sim \mathcal{N}(\mathbf{z} | \mathbf{0}, \mathbf{I}) \quad (10.65)$$

$$\mathbf{x}_n | \mathbf{z}_n \sim \mathcal{N}(\mathbf{x} | \mathbf{B}\mathbf{z}_n + \boldsymbol{\mu}, \sigma^2 \mathbf{I}) \quad (10.66)$$

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 (10.64) to sample a data point conditioned on the sampled  $\mathbf{z}_n$ , i.e.,  $\mathbf{x}_n \sim p(\mathbf{x} | \mathbf{z}_n, \mathbf{B}, \boldsymbol{\mu}, \sigma^2)$ .

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

$$p(\mathbf{x}, \mathbf{z} | \mathbf{B}, \boldsymbol{\mu}, \sigma^2) = p(\mathbf{x} | \mathbf{z}, \mathbf{B}, \boldsymbol{\mu}, \sigma^2) p(\mathbf{z}), \quad (10.67)$$

which immediately gives rise to the graphical model in Figure 10.13 using the results from Section 8.4.

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

#### Example 10.4 (Generating Data from Latent Variables)

Figure 10.14 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  $\mathbf{z}_*$  in this latent space and generate an image  $\tilde{\mathbf{x}}_* = \mathbf{B}\mathbf{z}_*$  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 artefacts, e.g., the top-left and top-right digits. Note that the intrinsic dimensionality of these generated images is only two.





**Figure 10.14**  
Generating new MNIST digits. The latent variables  $\mathbf{z}$  can be used to generate new data  $\tilde{\mathbf{x}} = \mathbf{B}\mathbf{z}$ . The closer we stay to the training data the more realistic the generated data.

### 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

$$\begin{aligned} p(\mathbf{x} | \mathbf{B}, \boldsymbol{\mu}, \sigma^2) &= \int p(\mathbf{x} | \mathbf{z}, \boldsymbol{\mu}, \sigma^2) p(\mathbf{z}) d\mathbf{z} \\ &= \int \mathcal{N}(\mathbf{x} | \mathbf{B}\mathbf{z} + \boldsymbol{\mu}, \sigma^2 \mathbf{I}) \mathcal{N}(\mathbf{z} | \mathbf{0}, \mathbf{I}) d\mathbf{z}. \end{aligned} \quad (10.68)$$

From Section 6.5, we know that the solution to this integral is a Gaussian distribution with mean

$$\mathbb{E}[\mathbf{x}] = \mathbb{E}_{\mathbf{z}}[\mathbf{B}\mathbf{z} + \boldsymbol{\mu}] + \mathbb{E}_{\epsilon}[\epsilon] = \boldsymbol{\mu} \quad (10.69)$$

and with covariance matrix

$$\begin{aligned} \mathbb{V}[\mathbf{x}] &= \mathbb{V}_{\mathbf{z}}[\mathbf{B}\mathbf{z} + \boldsymbol{\mu}] + \mathbb{V}_{\epsilon}[\epsilon] = \mathbb{V}_{\mathbf{z}}[\mathbf{B}\mathbf{z}] + \sigma^2 \mathbf{I} \\ &= \mathbf{B} \mathbb{V}_{\mathbf{z}}[\mathbf{z}] \mathbf{B}^{\top} + \sigma^2 \mathbf{I} = \mathbf{B} \mathbf{B}^{\top} + \sigma^2 \mathbf{I}. \end{aligned} \quad (10.70)$$

The marginal distribution in (10.68) 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 (10.64) 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 6.5 we also know that the joint distribution of a Gaussian random variable  $z$  and a linear/affine transformation  $x = Bz$  of it are jointly Gaussian distributed. We already know the marginals  $p(z) = \mathcal{N}(z | \mathbf{0}, \mathbf{I})$  and  $p(x) = \mathcal{N}(x | \mu, BB^\top + \sigma^2 \mathbf{I})$ . The missing cross-covariance is given as

$$\text{Cov}[x, z] = \text{Cov}_z[Bz + \mu] = B \text{Cov}_z[z, z] = B. \quad (10.71)$$

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

$$p(x, z | B, \mu, \sigma^2) = \mathcal{N} \left( \begin{bmatrix} x \\ z \end{bmatrix} \middle| \begin{bmatrix} \mu \\ \mathbf{0} \end{bmatrix}, \begin{bmatrix} BB^\top + \sigma^2 \mathbf{I} & B \\ B^\top & \mathbf{I} \end{bmatrix} \right), \quad (10.72)$$

5951 with a mean vector of length  $D + M$  and a covariance matrix of size  
5952  $(D + M) \times (D + M)$ .

### 5953 10.7.3 Posterior Distribution

The joint Gaussian distribution  $p(x, z | B, \mu, \sigma^2)$  in (10.72) allows us to determine the posterior distribution  $p(z | x)$  immediately by applying the rules of Gaussian conditioning from Section 6.5.1. The posterior distribution of the latent variable given an observation  $x$  is then

$$p(z | x) = \mathcal{N}(z | m, C), \quad (10.73a)$$

$$m = B^\top (BB^\top + \sigma^2 \mathbf{I})^{-1} (x - \mu), \quad (10.73b)$$

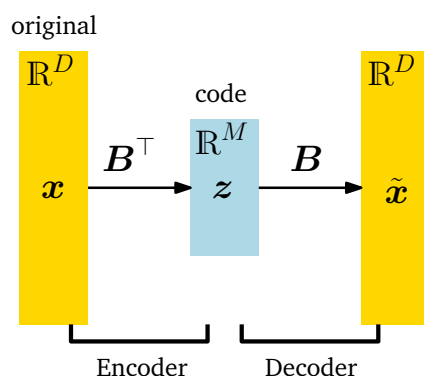
$$C = \mathbf{I} - B^\top (BB^\top + \sigma^2 \mathbf{I})^{-1} B. \quad (10.73c)$$

5954 Note that the posterior covariance does not depend on the observation  $x$ .

5955 For a new observation  $x_*$  in data space, we can use (10.73) to deter-  
5956 mine the posterior distribution of the corresponding latent variable  $z_*$ .  
5957 The covariance matrix  $C$  allows us to assess how confident the embed-  
5958 ding is. A covariance matrix  $C$  with a small determinant (which measures  
5959 volumes) tells us that the latent embedding  $z_*$  is fairly certain. If we ob-  
5960 tain a posterior distribution  $p(z_* | x_*)$  with much variance, we may be  
5961 faced with an outlier. However, we can explore this posterior distribution  
5962 to understand what other data points  $x$  are plausible under this posterior.  
5963 To do this, we can exploit PPCA's generative process. The generative pro-  
5964 cess underlying PPCA allows us to explore the posterior distribution on  
5965 the latent variables by generating new data that are plausible under this  
5966 posterior. This can be achieved as follows:

- 5967 1. Sample a latent variable  $z_* \sim p(z | x_*)$  from the posterior distribution  
5968 over the latent variables (10.73)
- 5969 2. Sample a reconstructed vector  $\tilde{x}_* \sim p(x | z_*, B, \mu, \sigma^2)$  from (10.64)

5970 If we repeat this process many times, we can explore the posterior dis-  
5971 tribution (10.73) on the latent variables  $z_*$  and its implications on the



**Figure 10.15** 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.

observed data. The sampling process effectively hypothesizes data, which is plausible under the posterior distribution.

## 10.8 Further Reading

We derived PCA from two perspectives: a) maximizing the variance in the projected space; b) minimizing the average reconstruction error. However, PCA can also be interpreted from different perspectives. Let us re-cap what we have done: We took high-dimensional data  $x \in \mathbb{R}^D$  and used a matrix  $B^\top$  to find a lower-dimensional representation  $z \in \mathbb{R}^M$ . The columns of  $B$  are the eigenvectors of the data covariance matrix  $S$  that are associated with the largest eigenvalues. Once we have a low-dimensional representation  $z$ , we can get a high-dimensional version of it (in the original data 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 10.15. 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

$$\frac{1}{N} \sum_{n=1}^N \|x_n - \tilde{x}_n\|^2 = \frac{1}{N} \sum_{n=1}^N \|x_n - B^\top B x_n\|^2. \quad (10.74)$$

This means we end up with the same objective function as in (10.30) 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

auto-encoder  
code  
  
encoder  
decoder

recognition network  
inference network  
generator  
The code is a  
compressed version  
of the original data

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, 2003a).

In our discussion on PPCA, we assumed that the parameters of the model, i.e.,  $\mathbf{B}$ ,  $\boldsymbol{\mu}$  and the likelihood parameter  $\sigma^2$  are known. Tipping and Bishop (1999) describe how to derive maximum likelihood estimates for these parameters 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

The matrix  $\mathbf{\Lambda} - \sigma^2 \mathbf{I}$  in (10.76) is guaranteed to be positive semi-definite as the smallest eigenvalue of the data covariance matrix is bounded from below by the noise variance  $\sigma^2$ .

$$\boldsymbol{\mu}_{\text{ML}} = \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n, \quad (10.75)$$

$$\mathbf{B}_{\text{ML}} = \mathbf{T}(\mathbf{\Lambda} - \sigma^2 \mathbf{I})^{\frac{1}{2}} \mathbf{R}, \quad (10.76)$$

$$\sigma_{\text{ML}}^2 = \frac{1}{D - M} \sum_{j=M+1}^D \lambda_j, \quad (10.77)$$

where  $\mathbf{T} \in \mathbb{R}^{D \times M}$  contains  $M$  eigenvectors of the data covariance matrix,  $\mathbf{\Lambda} = \text{diag}(\lambda_1, \dots, \lambda_M) \in \mathbb{R}^{M \times M}$  is a diagonal matrix with the eigenvalues associated with the principal axes on its diagonal, and  $\mathbf{R} \in \mathbb{R}^{M \times M}$  is an arbitrary orthogonal matrix. The maximum likelihood solution  $\mathbf{B}_{\text{ML}}$  is unique up to an arbitrary orthogonal transformation, e.g., we can right-multiply  $\mathbf{B}_{\text{ML}}$  with any rotation matrix  $\mathbf{R}$  so that (10.76) 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 (10.75) is the sample mean of the data. The maximum likelihood estimator for the observation noise variance  $\sigma^2$  given in (10.77) is the average variance in the orthogonal complement of the principal subspace, i.e., the average 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

$$\mathbf{S} = \mathbf{T}\mathbf{\Lambda}\mathbf{T}^{-1}. \quad (10.78)$$

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 (10.70). 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 (10.78)) so that

$$\text{Cov}[\mathbf{X}] = \mathbf{T}\mathbf{\Lambda}\mathbf{T}^{-1} = \mathbf{B}\mathbf{B}^\top \iff \mathbf{B} = \mathbf{T}\mathbf{\Lambda}^{\frac{1}{2}}\mathbf{R}, \quad (10.79)$$

so that we obtain exactly the maximum likelihood estimate in (10.76) for  $\sigma = 0$ .

From (10.76) and (10.78) it becomes clear that (P)PCA performs a decomposition of the data covariance matrix.

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).

To determine the dimensionality of the latent variables (the length of the code, the dimensionality of the lower-dimensional subspace onto which we project the data) Gavish and Donoho (2014) suggest the heuristic that, if we can estimate the noise variance  $\sigma^2$  of the data, we should discard all singular values smaller than  $\frac{4\sigma\sqrt{D}}{\sqrt{3}}$ . Alternatively, we can use cross validation or the Bayesian model selection criteria (discussed in Section 8.5.2) to determine the true dimensionality of the data (Minka, 2001).

Similar to our discussion on linear regression in Chapter 9, we can place a prior distribution on the parameters of the model and integrate them out, thereby avoiding a) point estimates of the parameters and the issues that come with these point estimates (see Section 8.5) 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  $\sigma_d^2$  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

factor analysis

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

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.

**Independent Component Analysis** *Independent Component Analysis* (ICA) is also closely related to PCA. Starting again with the model  $\mathbf{x}_n = \mathbf{B}\mathbf{z}_n + \epsilon$  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 Murphy (2012) for more details on ICA.

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.

The (P)PCA model we discussed here allows for several important extensions. In Section 10.5, we explained how to do PCA when the input dimensionality  $D$  is significantly greater than the number  $N$  of data points. By exploiting the insight that PCA can be performed by computing (many) inner products, this idea can be pushed to the extreme by considering infinite-dimensional features. The *kernel trick* is the basis of *kernel PCA* and allows us to implicitly compute inner products between infinite-dimensional features (Schölkopf et al., 1998; Schölkopf and Smola, 2002).

There are nonlinear dimensionality reduction techniques that are derived from PCA (Burges (2010) provide a good overview). The auto-encoder perspective of PCA that we discussed above can be used to render PCA as a special case of a *deep auto-encoder*. In the deep auto-encoder, both the encoder and the decoder are represented by multi-layer feedforward neural networks, which themselves are nonlinear mappings. If we set the activation functions in these neural networks to be the identity, the model becomes equivalent to PCA. A different approach to nonlinear dimensionality reduction is the *Gaussian Process Latent Variable Model* (GP-LVM) proposed by Lawrence (2005). The GP-LVM starts off with the latent-variable perspective that we used to derive PPCA and replaces the linear relationship between the latent variables  $\mathbf{z}$  and the observations  $\mathbf{x}$  with a Gaussian process (GP). Instead of estimating the parameters of the mapping (as we do in PPCA), the GP-LVM marginalizes out the model parameters and makes point estimates of the latent variables  $\mathbf{z}$ . Similar to Bayesian PCA, the *Bayesian GP-LVM* proposed by Titsias and Lawrence (2010) maintains

6097 a distribution on the latent variables  $z$  and uses approximate inference to  
6098 integrate them out as well.