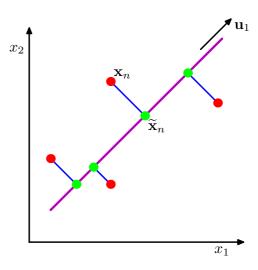
Figure 12.2 Principal component analysis seeks a space of lower dimensionality, known as the principal subspace and denoted by the magenta line, such that the orthogonal projection of the data points (red dots) onto this subspace maximizes the variance of the projected points (green dots). An alternative definition of PCA is based on minimizing the sum-of-squares of the projection errors, indicated by the blue lines.



### Section 12.2

a particular form of linear-Gaussian latent variable model. This probabilistic reformulation brings many advantages, such as the use of EM for parameter estimation, principled extensions to mixtures of PCA models, and Bayesian formulations that allow the number of principal components to be determined automatically from the data. Finally, we discuss briefly several generalizations of the latent variable concept that go beyond the linear-Gaussian assumption including non-Gaussian latent variables, which leads to the framework of *independent component analysis*, as well as models having a nonlinear relationship between latent and observed variables.

#### Section 12.4

# 12.1. Principal Component Analysis

Principal component analysis, or PCA, is a technique that is widely used for applications such as dimensionality reduction, lossy data compression, feature extraction, and data visualization (Jolliffe, 2002). It is also known as the *Karhunen-Loève* transform.

There are two commonly used definitions of PCA that give rise to the same algorithm. PCA can be defined as the orthogonal projection of the data onto a lower dimensional linear space, known as the *principal subspace*, such that the variance of the projected data is maximized (Hotelling, 1933). Equivalently, it can be defined as the linear projection that minimizes the average projection cost, defined as the mean squared distance between the data points and their projections (Pearson, 1901). The process of orthogonal projection is illustrated in Figure 12.2. We consider each of these definitions in turn.

### 12.1.1 Maximum variance formulation

Consider a data set of observations  $\{\mathbf{x}_n\}$  where  $n=1,\ldots,N$ , and  $\mathbf{x}_n$  is a Euclidean variable with dimensionality D. Our goal is to project the data onto a space having dimensionality M < D while maximizing the variance of the projected data. For the moment, we shall assume that the value of M is given. Later in this

chapter, we shall consider techniques to determine an appropriate value of M from the data.

To begin with, consider the projection onto a one-dimensional space (M=1). We can define the direction of this space using a D-dimensional vector  $\mathbf{u}_1$ , which for convenience (and without loss of generality) we shall choose to be a unit vector so that  $\mathbf{u}_1^T\mathbf{u}_1=1$  (note that we are only interested in the direction defined by  $\mathbf{u}_1$ , not in the magnitude of  $\mathbf{u}_1$  itself). Each data point  $\mathbf{x}_n$  is then projected onto a scalar value  $\mathbf{u}_1^T\mathbf{x}_n$ . The mean of the projected data is  $\mathbf{u}_1^T\overline{\mathbf{x}}$  where  $\overline{\mathbf{x}}$  is the sample set mean given by

$$\overline{\mathbf{x}} = \frac{1}{N} \sum_{n=1}^{N} \mathbf{x}_n \tag{12.1}$$

and the variance of the projected data is given by

$$\frac{1}{N} \sum_{n=1}^{N} \left\{ \mathbf{u}_{1}^{\mathrm{T}} \mathbf{x}_{n} - \mathbf{u}_{1}^{\mathrm{T}} \overline{\mathbf{x}} \right\}^{2} = \mathbf{u}_{1}^{\mathrm{T}} \mathbf{S} \mathbf{u}_{1}$$
 (12.2)

where S is the data covariance matrix defined by

$$\mathbf{S} = \frac{1}{N} \sum_{n=1}^{N} (\mathbf{x}_n - \overline{\mathbf{x}}) (\mathbf{x}_n - \overline{\mathbf{x}})^{\mathrm{T}}.$$
 (12.3)

We now maximize the projected variance  $\mathbf{u}_1^T \mathbf{S} \mathbf{u}_1$  with respect to  $\mathbf{u}_1$ . Clearly, this has to be a constrained maximization to prevent  $\|\mathbf{u}_1\| \to \infty$ . The appropriate constraint comes from the normalization condition  $\mathbf{u}_1^T \mathbf{u}_1 = 1$ . To enforce this constraint, we introduce a Lagrange multiplier that we shall denote by  $\lambda_1$ , and then make an unconstrained maximization of

$$\mathbf{u}_{1}^{\mathrm{T}}\mathbf{S}\mathbf{u}_{1} + \lambda_{1} \left(1 - \mathbf{u}_{1}^{\mathrm{T}}\mathbf{u}_{1}\right). \tag{12.4}$$

By setting the derivative with respect to  $\mathbf{u}_1$  equal to zero, we see that this quantity will have a stationary point when

$$\mathbf{S}\mathbf{u}_1 = \lambda_1 \mathbf{u}_1 \tag{12.5}$$

which says that  $\mathbf{u}_1$  must be an eigenvector of  $\mathbf{S}$ . If we left-multiply by  $\mathbf{u}_1^T$  and make use of  $\mathbf{u}_1^T\mathbf{u}_1=1$ , we see that the variance is given by

$$\mathbf{u}_1^{\mathrm{T}}\mathbf{S}\mathbf{u}_1 = \lambda_1 \tag{12.6}$$

and so the variance will be a maximum when we set  $\mathbf{u}_1$  equal to the eigenvector having the largest eigenvalue  $\lambda_1$ . This eigenvector is known as the first principal component.

We can define additional principal components in an incremental fashion by choosing each new direction to be that which maximizes the projected variance

### Appendix E

amongst all possible directions orthogonal to those already considered. If we consider the general case of an M-dimensional projection space, the optimal linear projection for which the variance of the projected data is maximized is now defined by the M eigenvectors  $\mathbf{u}_1, \ldots, \mathbf{u}_M$  of the data covariance matrix  $\mathbf{S}$  corresponding to the M largest eigenvalues  $\lambda_1, \ldots, \lambda_M$ . This is easily shown using proof by induction.

To summarize, principal component analysis involves evaluating the mean  $\overline{\mathbf{x}}$  and the covariance matrix  $\mathbf{S}$  of the data set and then finding the M eigenvectors of  $\mathbf{S}$  corresponding to the M largest eigenvalues. Algorithms for finding eigenvectors and eigenvalues, as well as additional theorems related to eigenvector decomposition, can be found in Golub and Van Loan (1996). Note that the computational cost of computing the full eigenvector decomposition for a matrix of size  $D \times D$  is  $O(D^3)$ . If we plan to project our data onto the first M principal components, then we only need to find the first M eigenvalues and eigenvectors. This can be done with more efficient techniques, such as the *power method* (Golub and Van Loan, 1996), that scale like  $O(MD^2)$ , or alternatively we can make use of the EM algorithm.

### Section 12.2.2

Exercise 12.1

## Appendix C

### 12.1.2 Minimum-error formulation

We now discuss an alternative formulation of PCA based on projection error minimization. To do this, we introduce a complete orthonormal set of D-dimensional basis vectors  $\{\mathbf{u}_i\}$  where  $i=1,\ldots,D$  that satisfy

$$\mathbf{u}_i^{\mathrm{T}} \mathbf{u}_j = \delta_{ij}. \tag{12.7}$$

Because this basis is complete, each data point can be represented exactly by a linear combination of the basis vectors

$$\mathbf{x}_n = \sum_{i=1}^D \alpha_{ni} \mathbf{u}_i \tag{12.8}$$

where the coefficients  $\alpha_{ni}$  will be different for different data points. This simply corresponds to a rotation of the coordinate system to a new system defined by the  $\{\mathbf{u}_i\}$ , and the original D components  $\{x_{n1},\ldots,x_{nD}\}$  are replaced by an equivalent set  $\{\alpha_{n1},\ldots,\alpha_{nD}\}$ . Taking the inner product with  $\mathbf{u}_j$ , and making use of the orthonormality property, we obtain  $\alpha_{nj} = \mathbf{x}_n^T \mathbf{u}_j$ , and so without loss of generality we can write

$$\mathbf{x}_n = \sum_{i=1}^{D} \left( \mathbf{x}_n^{\mathrm{T}} \mathbf{u}_i \right) \mathbf{u}_i. \tag{12.9}$$

Our goal, however, is to approximate this data point using a representation involving a restricted number M < D of variables corresponding to a projection onto a lower-dimensional subspace. The M-dimensional linear subspace can be represented, without loss of generality, by the first M of the basis vectors, and so we approximate each data point  $\mathbf{x}_n$  by

$$\widetilde{\mathbf{x}}_n = \sum_{i=1}^M z_{ni} \mathbf{u}_i + \sum_{i=M+1}^D b_i \mathbf{u}_i$$
 (12.10)

where the  $\{z_{ni}\}$  depend on the particular data point, whereas the  $\{b_i\}$  are constants that are the same for all data points. We are free to choose the  $\{\mathbf{u}_i\}$ , the  $\{z_{ni}\}$ , and the  $\{b_i\}$  so as to minimize the distortion introduced by the reduction in dimensionality. As our distortion measure, we shall use the squared distance between the original data point  $\mathbf{x}_n$  and its approximation  $\widetilde{\mathbf{x}}_n$ , averaged over the data set, so that our goal is to minimize

$$J = \frac{1}{N} \sum_{n=1}^{N} \|\mathbf{x}_n - \widetilde{\mathbf{x}}_n\|^2.$$
 (12.11)

Consider first of all the minimization with respect to the quantities  $\{z_{ni}\}$ . Substituting for  $\widetilde{\mathbf{x}}_n$ , setting the derivative with respect to  $z_{nj}$  to zero, and making use of the orthonormality conditions, we obtain

$$z_{nj} = \mathbf{x}_n^{\mathrm{T}} \mathbf{u}_j \tag{12.12}$$

where j = 1, ..., M. Similarly, setting the derivative of J with respect to  $b_i$  to zero, and again making use of the orthonormality relations, gives

$$b_j = \overline{\mathbf{x}}^{\mathrm{T}} \mathbf{u}_j \tag{12.13}$$

where  $j = M+1, \ldots, D$ . If we substitute for  $z_{ni}$  and  $b_i$ , and make use of the general expansion (12.9), we obtain

$$\mathbf{x}_n - \widetilde{\mathbf{x}}_n = \sum_{i=M+1}^D \left\{ (\mathbf{x}_n - \overline{\mathbf{x}})^{\mathrm{T}} \mathbf{u}_i \right\} \mathbf{u}_i$$
 (12.14)

from which we see that the displacement vector from  $\mathbf{x}_n$  to  $\widetilde{\mathbf{x}}_n$  lies in the space orthogonal to the principal subspace, because it is a linear combination of  $\{\mathbf{u}_i\}$  for  $i=M+1,\ldots,D$ , as illustrated in Figure 12.2. This is to be expected because the projected points  $\widetilde{\mathbf{x}}_n$  must lie within the principal subspace, but we can move them freely within that subspace, and so the minimum error is given by the orthogonal projection.

We therefore obtain an expression for the distortion measure J as a function purely of the  $\{u_i\}$  in the form

$$J = \frac{1}{N} \sum_{n=1}^{N} \sum_{i=M+1}^{D} \left( \mathbf{x}_{n}^{\mathrm{T}} \mathbf{u}_{i} - \overline{\mathbf{x}}^{\mathrm{T}} \mathbf{u}_{i} \right)^{2} = \sum_{i=M+1}^{D} \mathbf{u}_{i}^{\mathrm{T}} \mathbf{S} \mathbf{u}_{i}.$$
 (12.15)

There remains the task of minimizing J with respect to the  $\{\mathbf{u}_i\}$ , which must be a constrained minimization otherwise we will obtain the vacuous result  $\mathbf{u}_i=0$ . The constraints arise from the orthonormality conditions and, as we shall see, the solution will be expressed in terms of the eigenvector expansion of the covariance matrix. Before considering a formal solution, let us try to obtain some intuition about the result by considering the case of a two-dimensional data space D=2 and a one-dimensional principal subspace M=1. We have to choose a direction  $\mathbf{u}_2$  so as to

minimize  $J = \mathbf{u}_2^T \mathbf{S} \mathbf{u}_2$ , subject to the normalization constraint  $\mathbf{u}_2^T \mathbf{u}_2 = 1$ . Using a Lagrange multiplier  $\lambda_2$  to enforce the constraint, we consider the minimization of

$$\widetilde{J} = \mathbf{u}_2^{\mathrm{T}} \mathbf{S} \mathbf{u}_2 + \lambda_2 \left( 1 - \mathbf{u}_2^{\mathrm{T}} \mathbf{u}_2 \right). \tag{12.16}$$

Setting the derivative with respect to  $\mathbf{u}_2$  to zero, we obtain  $\mathbf{S}\mathbf{u}_2 = \lambda_2\mathbf{u}_2$  so that  $\mathbf{u}_2$  is an eigenvector of  $\mathbf{S}$  with eigenvalue  $\lambda_2$ . Thus any eigenvector will define a stationary point of the distortion measure. To find the value of J at the minimum, we back-substitute the solution for  $\mathbf{u}_2$  into the distortion measure to give  $J = \lambda_2$ . We therefore obtain the minimum value of J by choosing  $\mathbf{u}_2$  to be the eigenvector corresponding to the smaller of the two eigenvalues. Thus we should choose the principal subspace to be aligned with the eigenvector having the *larger* eigenvalue. This result accords with our intuition that, in order to minimize the average squared projection distance, we should choose the principal component subspace to pass through the mean of the data points and to be aligned with the directions of maximum variance. For the case when the eigenvalues are equal, any choice of principal direction will give rise to the same value of J.

The general solution to the minimization of J for arbitrary D and arbitrary M < D is obtained by choosing the  $\{\mathbf{u}_i\}$  to be eigenvectors of the covariance matrix given by

$$\mathbf{S}\mathbf{u}_i = \lambda_i \mathbf{u}_i \tag{12.17}$$

where i = 1, ..., D, and as usual the eigenvectors  $\{\mathbf{u}_i\}$  are chosen to be orthonormal. The corresponding value of the distortion measure is then given by

$$J = \sum_{i=M+1}^{D} \lambda_i \tag{12.18}$$

which is simply the sum of the eigenvalues of those eigenvectors that are orthogonal to the principal subspace. We therefore obtain the minimum value of J by selecting these eigenvectors to be those having the D-M smallest eigenvalues, and hence the eigenvectors defining the principal subspace are those corresponding to the M largest eigenvalues.

Although we have considered M < D, the PCA analysis still holds if M = D, in which case there is no dimensionality reduction but simply a rotation of the coordinate axes to align with principal components.

Finally, it is worth noting that there exists a closely related linear dimensionality reduction technique called *canonical correlation analysis*, or *CCA* (Hotelling, 1936; Bach and Jordan, 2002). Whereas PCA works with a single random variable, CCA considers two (or more) variables and tries to find a corresponding pair of linear subspaces that have high cross-correlation, so that each component within one of the subspaces is correlated with a single component from the other subspace. Its solution can be expressed in terms of a generalized eigenvector problem.

## 12.1.3 Applications of PCA

We can illustrate the use of PCA for data compression by considering the offline digits data set. Because each eigenvector of the covariance matrix is a vector

### Exercise 12.2

### Appendix A

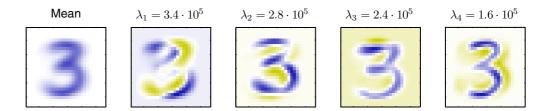


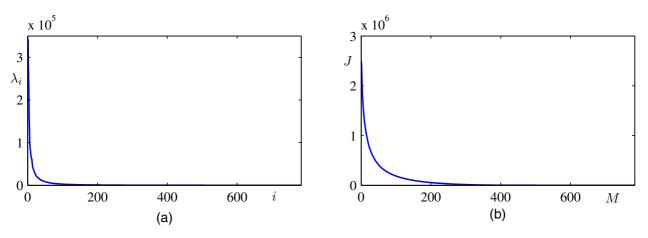
Figure 12.3 The mean vector  $\overline{\mathbf{x}}$  along with the first four PCA eigenvectors  $\mathbf{u}_1, \dots, \mathbf{u}_4$  for the off-line digits data set, together with the corresponding eigenvalues.

in the original D-dimensional space, we can represent the eigenvectors as images of the same size as the data points. The first five eigenvectors, along with the corresponding eigenvalues, are shown in Figure 12.3. A plot of the complete spectrum of eigenvalues, sorted into decreasing order, is shown in Figure 12.4(a). The distortion measure J associated with choosing a particular value of M is given by the sum of the eigenvalues from M+1 up to D and is plotted for different values of M in Figure 12.4(b).

If we substitute (12.12) and (12.13) into (12.10), we can write the PCA approximation to a data vector  $\mathbf{x}_n$  in the form

$$\widetilde{\mathbf{x}}_n = \sum_{i=1}^{M} (\mathbf{x}_n^{\mathrm{T}} \mathbf{u}_i) \mathbf{u}_i + \sum_{i=M+1}^{D} (\overline{\mathbf{x}}^{\mathrm{T}} \mathbf{u}_i) \mathbf{u}_i$$
 (12.19)

$$= \overline{\mathbf{x}} + \sum_{i=1}^{M} (\mathbf{x}_{n}^{\mathrm{T}} \mathbf{u}_{i} - \overline{\mathbf{x}}^{\mathrm{T}} \mathbf{u}_{i}) \mathbf{u}_{i}$$
 (12.20)



**Figure 12.4** (a) Plot of the eigenvalue spectrum for the off-line digits data set. (b) Plot of the sum of the discarded eigenvalues, which represents the sum-of-squares distortion J introduced by projecting the data onto a principal component subspace of dimensionality M.