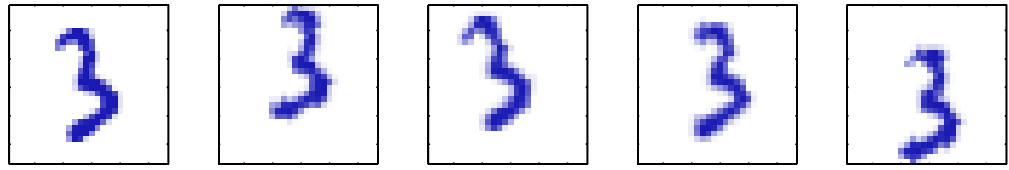


# 12

## Continuous Latent Variables

### Appendix A

In Chapter 9, we discussed probabilistic models having discrete latent variables, such as the mixture of Gaussians. We now explore models in which some, or all, of the latent variables are continuous. An important motivation for such models is that many data sets have the property that the data points all lie close to a manifold of much lower dimensionality than that of the original data space. To see why this might arise, consider an artificial data set constructed by taking one of the off-line digits, represented by a  $64 \times 64$  pixel grey-level image, and embedding it in a larger image of size  $100 \times 100$  by padding with pixels having the value zero (corresponding to white pixels) in which the location and orientation of the digit is varied at random, as illustrated in Figure 12.1. Each of the resulting images is represented by a point in the  $100 \times 100 = 10,000$ -dimensional data space. However, across a data set of such images, there are only three *degrees of freedom* of variability, corresponding to the vertical and horizontal translations and the rotations. The data points will therefore live on a subspace of the data space whose *intrinsic dimensionality* is three. Note



**Figure 12.1** A synthetic data set obtained by taking one of the off-line digit images and creating multiple copies in each of which the digit has undergone a random displacement and rotation within some larger image field. The resulting images each have  $100 \times 100 = 10,000$  pixels.

that the manifold will be nonlinear because, for instance, if we translate the digit past a particular pixel, that pixel value will go from zero (white) to one (black) and back to zero again, which is clearly a nonlinear function of the digit position. In this example, the translation and rotation parameters are latent variables because we observe only the image vectors and are not told which values of the translation or rotation variables were used to create them.

For real digit image data, there will be a further degree of freedom arising from scaling. Moreover there will be multiple additional degrees of freedom associated with more complex deformations due to the variability in an individual's writing as well as the differences in writing styles between individuals. Nevertheless, the number of such degrees of freedom will be small compared to the dimensionality of the data set.

#### Appendix A

Another example is provided by the oil flow data set, in which (for a given geometrical configuration of the gas, water, and oil phases) there are only two degrees of freedom of variability corresponding to the fraction of oil in the pipe and the fraction of water (the fraction of gas then being determined). Although the data space comprises 12 measurements, a data set of points will lie close to a two-dimensional manifold embedded within this space. In this case, the manifold comprises several distinct segments corresponding to different flow regimes, each such segment being a (noisy) continuous two-dimensional manifold. If our goal is data compression, or density modelling, then there can be benefits in exploiting this manifold structure.

In practice, the data points will not be confined precisely to a smooth low-dimensional manifold, and we can interpret the departures of data points from the manifold as 'noise'. This leads naturally to a generative view of such models in which we first select a point within the manifold according to some latent variable distribution and then generate an observed data point by adding noise, drawn from some conditional distribution of the data variables given the latent variables.

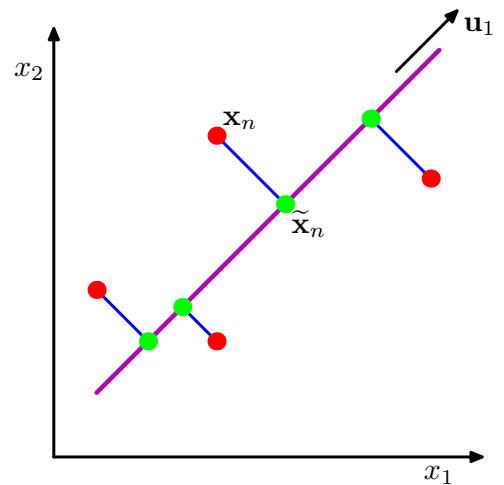
The simplest continuous latent variable model assumes Gaussian distributions for both the latent and observed variables and makes use of a linear-Gaussian dependence of the observed variables on the state of the latent variables. This leads to a probabilistic formulation of the well-known technique of principal component analysis (PCA), as well as to a related model called factor analysis.

#### Section 8.1.4

#### Section 12.1

In this chapter we will begin with a standard, nonprobabilistic treatment of PCA, and then we show how PCA arises naturally as the maximum likelihood solution to

**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  $n = 1, \dots, N$ , 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  $\mathbf{u}_1^T \bar{\mathbf{x}}$  where  $\bar{\mathbf{x}}$  is the sample set mean given by

$$\bar{\mathbf{x}} = \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n \quad (12.1)$$

and the variance of the projected data is given by

$$\frac{1}{N} \sum_{n=1}^N \{ \mathbf{u}_1^T \mathbf{x}_n - \mathbf{u}_1^T \bar{\mathbf{x}} \}^2 = \mathbf{u}_1^T \mathbf{S} \mathbf{u}_1 \quad (12.2)$$

where  $\mathbf{S}$  is the data covariance matrix defined by

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

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  $\|\mathbf{u}_1\| \rightarrow \infty$ . 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

$$\mathbf{u}_1^T \mathbf{S} \mathbf{u}_1 + \lambda_1 (1 - \mathbf{u}_1^T \mathbf{u}_1). \quad (12.4)$$

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

$$\mathbf{S} \mathbf{u}_1 = \lambda_1 \mathbf{u}_1 \quad (12.5)$$

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

$$\mathbf{u}_1^T \mathbf{S} \mathbf{u}_1 = \lambda_1 \quad (12.6)$$

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

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  $\mathbf{u}_1, \dots, \mathbf{u}_M$  of the data covariance matrix  $\mathbf{S}$  corresponding to the  $M$  largest eigenvalues  $\lambda_1, \dots, \lambda_M$ . This is easily shown using proof by induction.

To summarize, principal component analysis involves evaluating the mean  $\bar{\mathbf{x}}$  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.

### 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  $i = 1, \dots, D$  that satisfy

$$\mathbf{u}_i^T \mathbf{u}_j = \delta_{ij}. \quad (12.7)$$

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

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

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  $\{x_{n1}, \dots, x_{nD}\}$  are replaced by an equivalent set  $\{\alpha_{n1}, \dots, \alpha_{nD}\}$ . 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

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

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

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

Exercise 12.1

Section 12.2.2

Appendix C



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  $\tilde{\mathbf{x}}_n$ , averaged over the data set, so that our goal is to minimize

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

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

$$z_{nj} = \mathbf{x}_n^T \mathbf{u}_j \quad (12.12)$$

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

$$b_j = \bar{\mathbf{x}}^T \mathbf{u}_j \quad (12.13)$$

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

$$\mathbf{x}_n - \tilde{\mathbf{x}}_n = \sum_{i=M+1}^D \{(\mathbf{x}_n - \bar{\mathbf{x}})^T \mathbf{u}_i\} \mathbf{u}_i \quad (12.14)$$

from which we see that the displacement vector from  $\mathbf{x}_n$  to  $\tilde{\mathbf{x}}_n$  lies in the space orthogonal to the principal subspace, because it is a linear combination of  $\{\mathbf{u}_i\}$  for  $i = M+1, \dots, D$ , as illustrated in Figure 12.2. This is to be expected because the projected points  $\tilde{\mathbf{x}}_n$  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  $\{\mathbf{u}_i\}$  in the form

$$J = \frac{1}{N} \sum_{n=1}^N \sum_{i=M+1}^D (\mathbf{x}_n^T \mathbf{u}_i - \bar{\mathbf{x}}^T \mathbf{u}_i)^2 = \sum_{i=M+1}^D \mathbf{u}_i^T \mathbf{S} \mathbf{u}_i. \quad (12.15)$$

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

$$\tilde{J} = \mathbf{u}_2^T \mathbf{S} \mathbf{u}_2 + \lambda_2 (1 - \mathbf{u}_2^T \mathbf{u}_2). \quad (12.16)$$

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  $J$ .

### Exercise 12.2

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

$$\mathbf{S} \mathbf{u}_i = \lambda_i \mathbf{u}_i \quad (12.17)$$

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

$$J = \sum_{i=M+1}^D \lambda_i \quad (12.18)$$

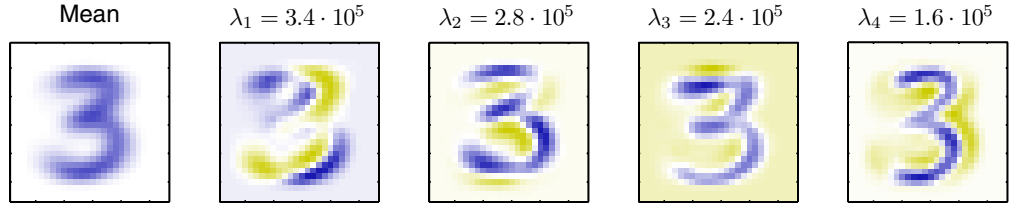
which is simply the sum of the eigenvalues of those eigenvectors that are orthogonal to the principal subspace. We therefore obtain the minimum value of  $J$  by selecting these eigenvectors to be those having the  $D - M$  smallest eigenvalues, and hence the eigenvectors defining the principal subspace 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 off-line digits data set. Because each eigenvector of the covariance matrix is a vector



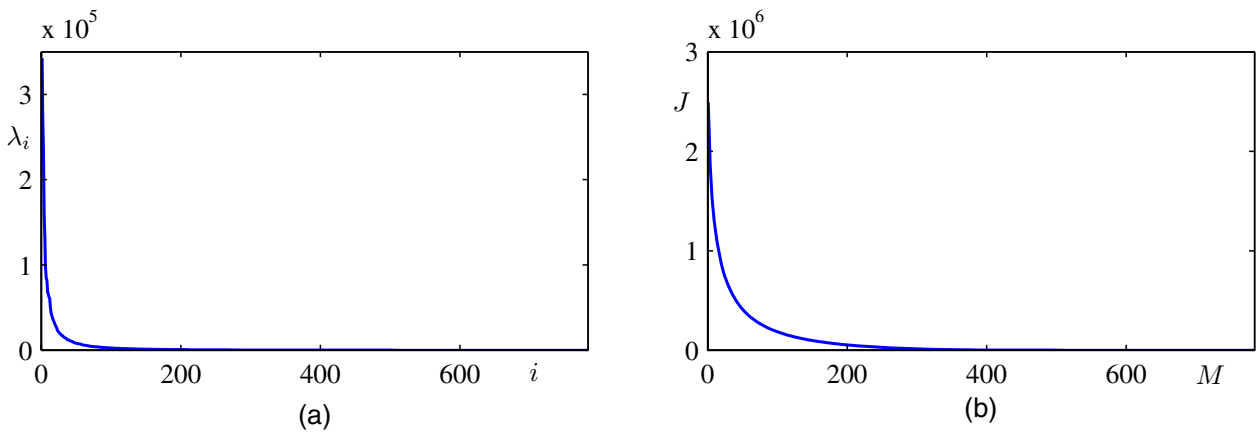
**Figure 12.3** The mean vector  $\bar{\mathbf{x}}$  along with the first four PCA eigenvectors  $\mathbf{u}_1, \dots, \mathbf{u}_4$  for the off-line digits data set, together with the corresponding eigenvalues.

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

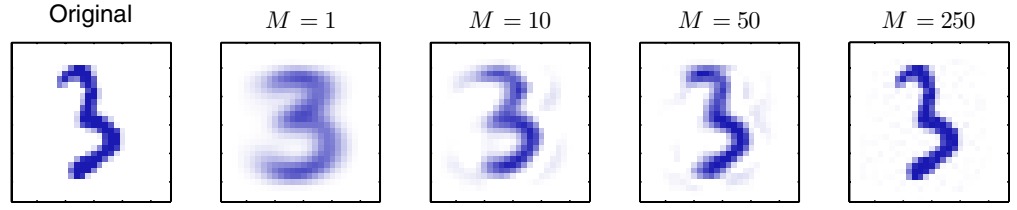
$$\tilde{\mathbf{x}}_n = \sum_{i=1}^M (\mathbf{x}_n^T \mathbf{u}_i) \mathbf{u}_i + \sum_{i=M+1}^D (\bar{\mathbf{x}}^T \mathbf{u}_i) \mathbf{u}_i \quad (12.19)$$

$$= \bar{\mathbf{x}} + \sum_{i=1}^M (\mathbf{x}_n^T \mathbf{u}_i - \bar{\mathbf{x}}^T \mathbf{u}_i) \mathbf{u}_i \quad (12.20)$$



**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  $M$ .





**Figure 12.5** An original example from the off-line digits data set together with its PCA reconstructions obtained by retaining  $M$  principal components for various values of  $M$ . As  $M$  increases the reconstruction becomes more accurate and would become perfect when  $M = D = 28 \times 28 = 784$ .

where we have made use of the relation

$$\bar{\mathbf{x}} = \sum_{i=1}^D (\bar{\mathbf{x}}^T \mathbf{u}_i) \mathbf{u}_i \quad (12.21)$$

which follows from the completeness of the  $\{\mathbf{u}_i\}$ . This represents a compression of the data set, because for each data point we have replaced the  $D$ -dimensional vector  $\mathbf{x}_n$  with an  $M$ -dimensional vector having components  $(\mathbf{x}_n^T \mathbf{u}_i - \bar{\mathbf{x}}^T \mathbf{u}_i)$ . The smaller the value of  $M$ , the greater the degree of compression. Examples of PCA reconstructions of data points for the digits data set are shown in Figure 12.5.

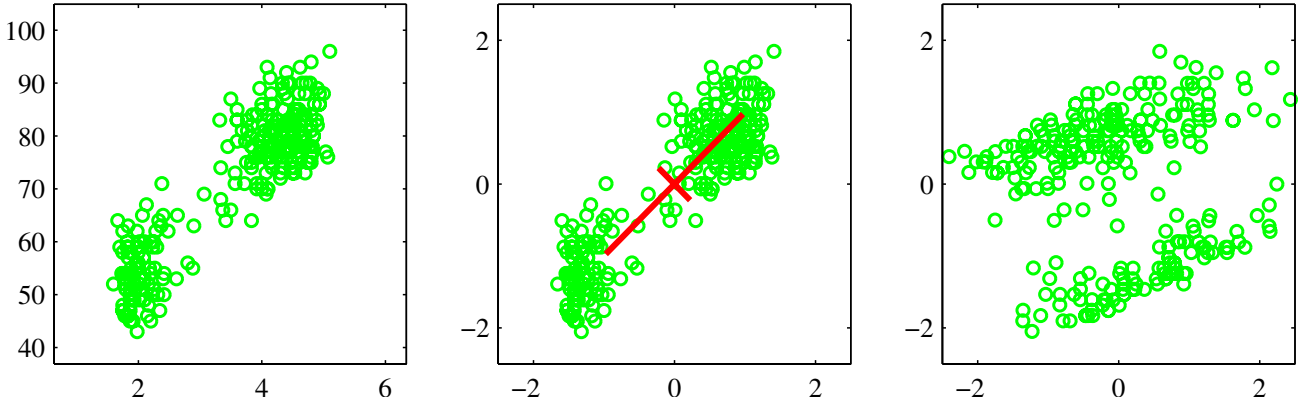
Another application of principal component analysis is to data pre-processing. In this case, the goal is not dimensionality reduction but rather the transformation of a data set in order to standardize certain of its properties. This can be important in allowing subsequent pattern recognition algorithms to be applied successfully to the data set. Typically, it is done when the original variables are measured in various different units or have significantly different variability. For instance in the Old Faithful data set, the time between eruptions is typically an order of magnitude greater than the duration of an eruption. When we applied the  $K$ -means algorithm to this data set, we first made a separate linear re-scaling of the individual variables such that each variable had zero mean and unit variance. This is known as *standardizing* the data, and the covariance matrix for the standardized data has components

$$\rho_{ij} = \frac{1}{N} \sum_{n=1}^N \frac{(x_{ni} - \bar{x}_i)}{\sigma_i} \frac{(x_{nj} - \bar{x}_j)}{\sigma_j} \quad (12.22)$$

where  $\sigma_i$  is the variance of  $x_i$ . This is known as the *correlation* matrix of the original data and has the property that if two components  $x_i$  and  $x_j$  of the data are perfectly correlated, then  $\rho_{ij} = 1$ , and if they are uncorrelated, then  $\rho_{ij} = 0$ .

However, using PCA we can make a more substantial normalization of the data to give it zero mean and unit covariance, so that different variables become decorrelated. To do this, we first write the eigenvector equation (12.17) in the form

$$\mathbf{S}\mathbf{U} = \mathbf{U}\mathbf{L} \quad (12.23)$$



**Figure 12.6** Illustration of the effects of linear pre-processing applied to the Old Faithful data set. The plot on the left shows the original data. The centre plot shows the result of standardizing the individual variables to zero mean and unit variance. Also shown are the principal axes of this normalized data set, plotted over the range  $\pm\lambda_i^{1/2}$ . The plot on the right shows the result of whitening of the data to give it zero mean and unit covariance.

where  $\mathbf{L}$  is a  $D \times D$  diagonal matrix with elements  $\lambda_i$ , and  $\mathbf{U}$  is a  $D \times D$  orthogonal matrix with columns given by  $\mathbf{u}_i$ . Then we define, for each data point  $\mathbf{x}_n$ , a transformed value given by

$$\mathbf{y}_n = \mathbf{L}^{-1/2} \mathbf{U}^T (\mathbf{x}_n - \bar{\mathbf{x}}) \quad (12.24)$$

where  $\bar{\mathbf{x}}$  is the sample mean defined by (12.1). Clearly, the set  $\{\mathbf{y}_n\}$  has zero mean, and its covariance is given by the identity matrix because

$$\begin{aligned} \frac{1}{N} \sum_{n=1}^N \mathbf{y}_n \mathbf{y}_n^T &= \frac{1}{N} \sum_{n=1}^N \mathbf{L}^{-1/2} \mathbf{U}^T (\mathbf{x}_n - \bar{\mathbf{x}}) (\mathbf{x}_n - \bar{\mathbf{x}})^T \mathbf{U} \mathbf{L}^{-1/2} \\ &= \mathbf{L}^{-1/2} \mathbf{U}^T \mathbf{S} \mathbf{U} \mathbf{L}^{-1/2} = \mathbf{L}^{-1/2} \mathbf{L} \mathbf{L}^{-1/2} = \mathbf{I}. \end{aligned} \quad (12.25)$$

This operation is known as *whitening* or *sphereing* the data and is illustrated for the Old Faithful data set in Figure 12.6.

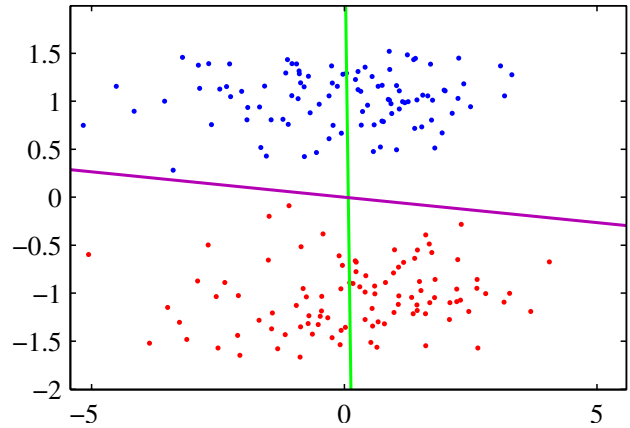
It is interesting to compare PCA with the Fisher linear discriminant which was discussed in Section 4.1.4. Both methods can be viewed as techniques for linear dimensionality reduction. However, PCA is unsupervised and depends only on the values  $\mathbf{x}_n$  whereas Fisher linear discriminant also uses class-label information. This difference is highlighted by the example in Figure 12.7.

Another common application of principal component analysis is to data visualization. Here each data point is projected onto a two-dimensional ( $M = 2$ ) principal subspace, so that a data point  $\mathbf{x}_n$  is plotted at Cartesian coordinates given by  $\mathbf{x}_n^T \mathbf{u}_1$  and  $\mathbf{x}_n^T \mathbf{u}_2$ , where  $\mathbf{u}_1$  and  $\mathbf{u}_2$  are the eigenvectors corresponding to the largest and second largest eigenvalues. An example of such a plot, for the oil flow data set, is shown in Figure 12.8.

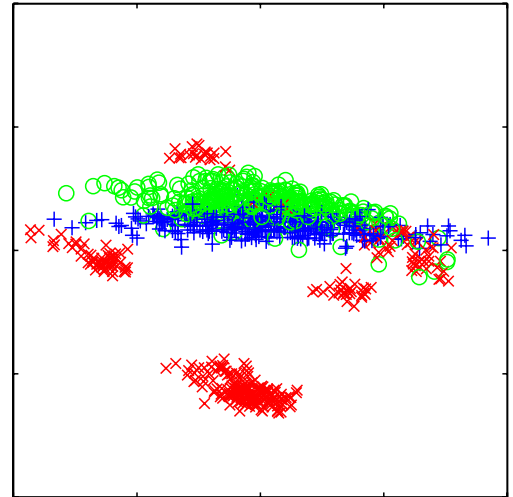
Appendix A

Appendix A

**Figure 12.7** A comparison of principal component analysis with Fisher's linear discriminant for linear dimensionality reduction. Here the data in two dimensions, belonging to two classes shown in red and blue, is to be projected onto a single dimension. PCA chooses the direction of maximum variance, shown by the magenta curve, which leads to strong class overlap, whereas the Fisher linear discriminant takes account of the class labels and leads to a projection onto the green curve giving much better class separation.



**Figure 12.8** Visualization of the oil flow data set obtained by projecting the data onto the first two principal components. The red, blue, and green points correspond to the 'laminar', 'homogeneous', and 'annular' flow configurations respectively.



#### 12.1.4 PCA for high-dimensional data

In some applications of principal component analysis, the number of data points is smaller than the dimensionality of the data space. For example, we might want to apply PCA to a data set of a few hundred images, each of which corresponds to a vector in a space of potentially several million dimensions (corresponding to three colour values for each of the pixels in the image). Note that in a  $D$ -dimensional space a set of  $N$  points, where  $N < D$ , defines a linear subspace whose dimensionality is at most  $N - 1$ , and so there is little point in applying PCA for values of  $M$  that are greater than  $N - 1$ . Indeed, if we perform PCA we will find that at least  $D - N + 1$  of the eigenvalues are zero, corresponding to eigenvectors along whose directions the data set has zero variance. Furthermore, typical algorithms for finding the eigenvectors of a  $D \times D$  matrix have a computational cost that scales like  $O(D^3)$ , and so for applications such as the image example, a direct application of PCA will be computationally infeasible.

We can resolve this problem as follows. First, let us define  $\mathbf{X}$  to be the  $(N \times D)$ -

dimensional centred data matrix, whose  $n^{\text{th}}$  row is given by  $(\mathbf{x}_n - \bar{\mathbf{x}})^T$ . The covariance matrix (12.3) can then be written as  $\mathbf{S} = N^{-1}\mathbf{X}^T\mathbf{X}$ , and the corresponding eigenvector equation becomes

$$\frac{1}{N}\mathbf{X}^T\mathbf{X}\mathbf{u}_i = \lambda_i\mathbf{u}_i. \quad (12.26)$$

Now pre-multiply both sides by  $\mathbf{X}$  to give

$$\frac{1}{N}\mathbf{X}\mathbf{X}^T(\mathbf{X}\mathbf{u}_i) = \lambda_i(\mathbf{X}\mathbf{u}_i). \quad (12.27)$$

If we now define  $\mathbf{v}_i = \mathbf{X}\mathbf{u}_i$ , we obtain

$$\frac{1}{N}\mathbf{X}\mathbf{X}^T\mathbf{v}_i = \lambda_i\mathbf{v}_i \quad (12.28)$$

which is an eigenvector equation for the  $N \times N$  matrix  $N^{-1}\mathbf{X}\mathbf{X}^T$ . We see that this has the same  $N - 1$  eigenvalues as the original covariance matrix (which itself has an additional  $D - N + 1$  eigenvalues of value zero). Thus we can solve the eigenvector problem in spaces of lower dimensionality with computational cost  $O(N^3)$  instead of  $O(D^3)$ . In order to determine the eigenvectors, we multiply both sides of (12.28) by  $\mathbf{X}^T$  to give

$$\left(\frac{1}{N}\mathbf{X}^T\mathbf{X}\right)(\mathbf{X}^T\mathbf{v}_i) = \lambda_i(\mathbf{X}^T\mathbf{v}_i) \quad (12.29)$$

from which we see that  $(\mathbf{X}^T\mathbf{v}_i)$  is an eigenvector of  $\mathbf{S}$  with eigenvalue  $\lambda_i$ . Note, however, that these eigenvectors need not be normalized. To determine the appropriate normalization, we re-scale  $\mathbf{u}_i \propto \mathbf{X}^T\mathbf{v}_i$  by a constant such that  $\|\mathbf{u}_i\| = 1$ , which, assuming  $\mathbf{v}_i$  has been normalized to unit length, gives

$$\mathbf{u}_i = \frac{1}{(N\lambda_i)^{1/2}}\mathbf{X}^T\mathbf{v}_i. \quad (12.30)$$

In summary, to apply this approach we first evaluate  $\mathbf{X}\mathbf{X}^T$  and then find its eigenvectors and eigenvalues and then compute the eigenvectors in the original data space using (12.30).

## 12.2. Probabilistic PCA

The formulation of PCA discussed in the previous section was based on a linear projection of the data onto a subspace of lower dimensionality than the original data space. We now show that PCA can also be expressed as the maximum likelihood solution of a probabilistic latent variable model. This reformulation of PCA, known as *probabilistic PCA*, brings several advantages compared with conventional PCA:

- Probabilistic PCA represents a constrained form of the Gaussian distribution in which the number of free parameters can be restricted while still allowing the model to capture the dominant correlations in a data set.