

Principal component analysis in high dimensions

Principal component analysis (PCA) is a standard technique for exploratory data analysis and dimension reduction. It is based on seeking the maximal variance components of a distribution, or equivalently, a low-dimensional subspace that captures the majority of the variance. Given a finite collection of samples, the empirical form of principal component analysis involves computing some subset of the top eigenvectors of the sample covariance matrix. Of interest is when these eigenvectors provide a good approximation to the subspace spanned by the top eigenvectors of the population covariance matrix. In this chapter, we study these issues in a high-dimensional and non-asymptotic framework, both for classical unstructured forms of PCA as well as for more modern structured variants.

8.1 Principal components and dimension reduction

Let $\mathcal{S}_+^{d \times d}$ denote the space of d -dimensional positive semidefinite matrices, and denote the d -dimensional unit sphere by $\mathbb{S}^{d-1} = \{v \in \mathbb{R}^d \mid \|v\|_2 = 1\}$. Consider a d -dimensional random vector X , say with a zero-mean vector and covariance matrix $\Sigma \in \mathcal{S}_+^{d \times d}$. We use

$$\gamma_1(\Sigma) \geq \gamma_2(\Sigma) \geq \cdots \geq \gamma_d(\Sigma) \geq 0$$

to denote the ordered eigenvalues of the covariance matrix. In its simplest instantiation, principal component analysis asks: along what unit-norm vector $v \in \mathbb{S}^{d-1}$ is the variance of the random variable $\langle v, X \rangle$ maximized? This direction is known as the first principal component at the population level, assumed here for the sake of discussion to be unique. In analytical terms, we have

$$v^* = \arg \max_{v \in \mathbb{S}^{d-1}} \text{var}(\langle v, X \rangle) = \arg \max_{v \in \mathbb{S}^{d-1}} \mathbb{E}[\langle v, X \rangle^2] = \arg \max_{v \in \mathbb{S}^{d-1}} \langle v, \Sigma v \rangle, \quad (8.1)$$

so that by definition, the first principal component is the maximum eigenvector of the covariance matrix Σ . More generally, we can define the top r principal components at the population level by seeking an orthonormal matrix $V \in \mathbb{R}^{d \times r}$, formed with unit-norm and orthogonal columns $\{v_1, \dots, v_r\}$, that maximizes the quantity

$$\mathbb{E}\|V^T X\|_2^2 = \sum_{j=1}^r \mathbb{E}[\langle v_j, X \rangle^2]. \quad (8.2)$$

As we explore in Exercise 8.4, these principal components are simply the top r eigenvectors of the population covariance matrix Σ .

In practice, however, we do not know the covariance matrix, but rather only have access

to a finite collection of samples, say $\{x_i\}_{i=1}^n$, each drawn according to \mathbb{P} . Based on these samples (and using the zero-mean assumption), we can form the sample covariance matrix $\widehat{\Sigma} = \frac{1}{n} \sum_{i=1}^n x_i x_i^T$. The empirical version of PCA is based on the “plug-in” principle, namely replacing the unknown population covariance Σ with this empirical version $\widehat{\Sigma}$. For instance, the empirical analog of the first principal component (8.1) is given by the optimization problem

$$\widehat{v} = \arg \max_{v \in \mathbb{S}^{d-1}} \langle v, \widehat{\Sigma} v \rangle. \quad (8.3)$$

Consequently, from the statistical point of view, we need to understand in what sense the maximizers of these empirically defined problems provide good approximations to their population analogs. Alternatively phrased, we need to determine how the eigenstructures of the population and sample covariance matrices are related.

8.1.1 Interpretations and uses of PCA

Before turning to the analysis of PCA, let us consider some of its interpretations and applications.

Example 8.1 (PCA as matrix approximation) Principal component analysis can be interpreted in terms of low-rank approximation. In particular, given some unitarily invariant¹ matrix norm $\|\cdot\|$, consider the problem of finding the best rank- r approximation to a given matrix Σ —that is,

$$\mathbf{Z}^* = \arg \min_{\text{rank}(\mathbf{Z}) \leq r} \left\{ \|\Sigma - \mathbf{Z}\|^2 \right\}. \quad (8.4)$$

In this interpretation, the matrix Σ need only be symmetric, not necessarily positive semi-definite as it must be when it is a covariance matrix. A classical result known as the *Eckart–Young–Mirsky theorem* guarantees that an optimal solution \mathbf{Z}^* exists, and takes the form of a truncated eigendecomposition, specified in terms of the top r eigenvectors of the matrix Σ . More precisely, recall that the symmetric matrix Σ has an orthonormal basis of eigenvectors, say $\{v_1, \dots, v_d\}$, associated with its ordered eigenvalues $\{\gamma_j(\Sigma)\}_{j=1}^d$. In terms of this notation, the optimal rank- r approximation takes the form

$$\mathbf{Z}^* = \sum_{j=1}^r \gamma_j(\Sigma) (v_j \otimes v_j), \quad (8.5)$$

where $v_j \otimes v_j := v_j v_j^T$ is the rank-one outer product. For the Frobenius matrix norm $\|\mathbf{M}\|_F = \sqrt{\sum_{j,k=1}^d M_{jk}^2}$, the error in the optimal approximation is given by

$$\|\mathbf{Z}^* - \Sigma\|_F^2 = \sum_{j=r+1}^d \gamma_j^2(\Sigma). \quad (8.6)$$

Figure 8.1 provides an illustration of the matrix approximation view of PCA. We first

¹ For a symmetric matrix \mathbf{M} , a matrix norm is unitarily invariant if $\|\mathbf{M}\| = \|\mathbf{V}^T \mathbf{M} \mathbf{V}\|$ for any orthonormal matrix \mathbf{V} . See Exercise 8.2 for further discussion.

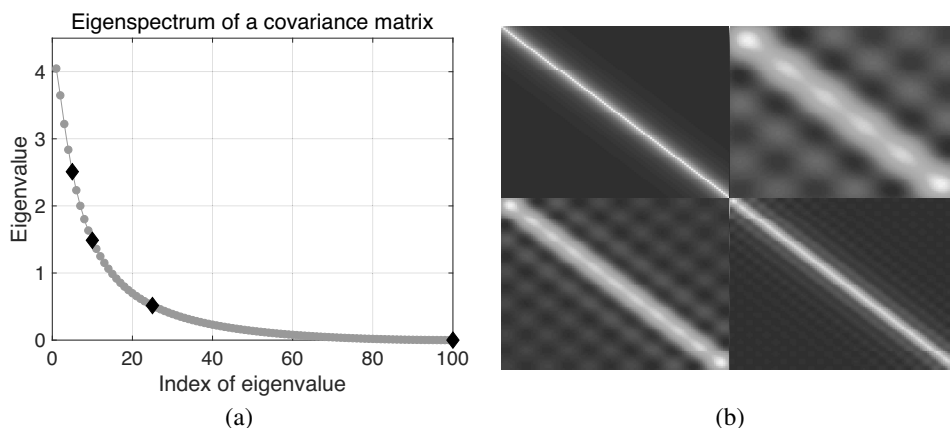


Figure 8.1 Illustration of PCA for low-rank matrix approximation. (a) Eigenspectrum of a matrix $\Sigma \in \mathcal{S}_+^{100 \times 100}$ generated as described in the text. Note the extremely rapid decay of the sorted eigenspectrum. Dark diamonds mark the rank cutoffs $r \in \{5, 10, 25, 100\}$, the first three of which define three approximations to the whole matrix ($r = 100$.) (b) Top left: original matrix. Top right: approximation based on $r = 5$ components. Bottom left: approximation based on $r = 10$ components. Bottom right: approximation based on $r = 25$ components.

generated the Toeplitz matrix $\mathbf{T} \in \mathcal{S}_+^{d \times d}$ with entries $T_{jk} = e^{-\alpha \sqrt{|j-k|}}$ with $\alpha = 0.95$, and then formed the recentered matrix $\Sigma := \mathbf{T} - \gamma_{\min}(\mathbf{T})\mathbf{I}_d$. Figure 8.1(a) shows the eigenspectrum of the matrix Σ : note that the rapid decay of the eigenvalues that renders it amenable to an accurate low-rank approximation. The top left image in Figure 8.1(b) corresponds to the original matrix Σ , whereas the remaining images illustrate approximations with increasing rank ($r = 5$ in top right, $r = 10$ in bottom left and $r = 25$ in bottom right). Although the defects in approximations with rank $r = 5$ or $r = 10$ are readily apparent, the approximation with rank $r = 25$ seems reasonable. ♣

Example 8.2 (PCA for data compression) Principal component analysis can also be interpreted as a linear form of data compression. Given a zero-mean random vector $X \in \mathbb{R}^d$, a simple way in which to compress it is via projection to a lower-dimensional subspace \mathbb{V} —say via a projection operator of the form $\Pi_{\mathbb{V}}(X)$. For a fixed dimension r , how do we choose the subspace \mathbb{V} ? Consider the criterion that chooses \mathbb{V} by minimizing the mean-squared error

$$\mathbb{E}[\|X - \Pi_{\mathbb{V}}(X)\|_2^2].$$

This optimal subspace need not be unique in general, but will be when there is a gap between the eigenvalues $\gamma_r(\Sigma)$ and $\gamma_{r+1}(\Sigma)$. In this case, the optimal subspace \mathbb{V}^* is spanned by the top r eigenvectors of the covariance matrix $\Sigma = \text{cov}(X)$. In particular, the projection operator $\Pi_{\mathbb{V}^*}$ can be written as $\Pi_{\mathbb{V}^*}(x) = \mathbf{V}_r \mathbf{V}_r^T x$, where $\mathbf{V}_r \in \mathbb{R}^{d \times r}$ is an orthonormal matrix with the top r eigenvectors $\{v_1, \dots, v_r\}$ as its columns. Using this optimal projection, the minimal

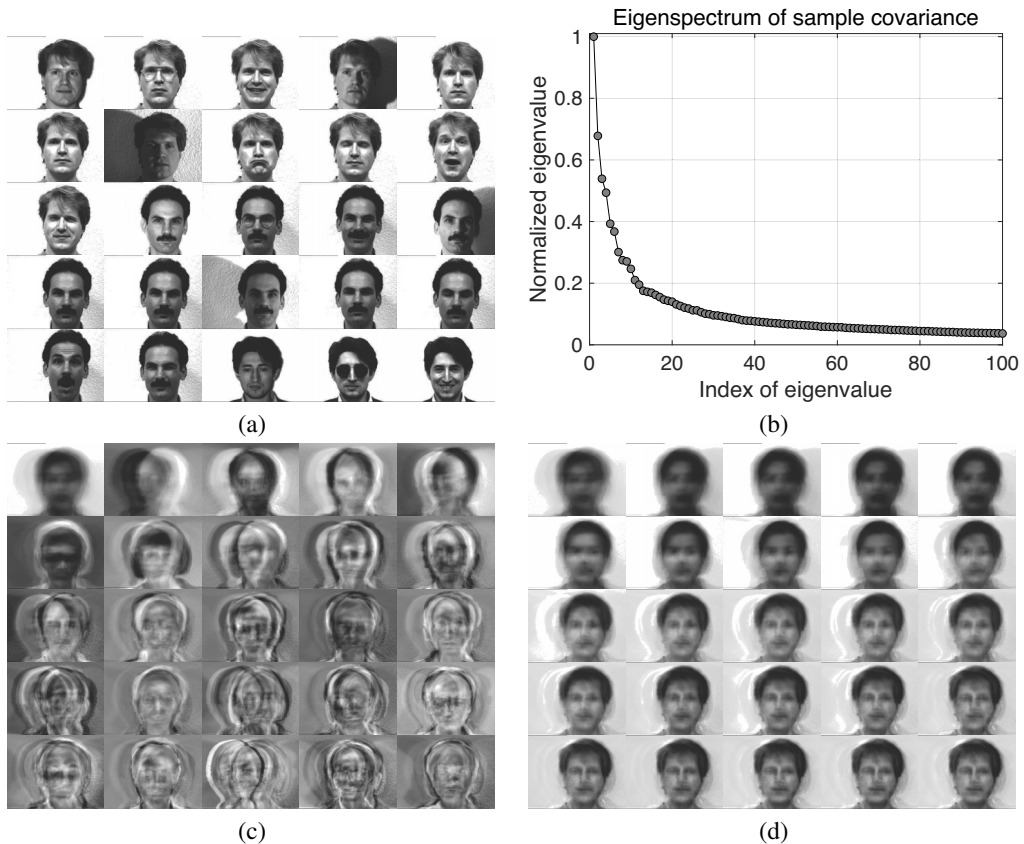


Figure 8.2 (a) Samples of face images from the Yale Face Database. (b) First 100 eigenvalues of the sample covariance matrix. (c) First 25 eigenfaces computed from the sample covariance matrix. (d) Reconstructions based on the first 25 eigenfaces plus the average face.

reconstruction error based on a rank- r projection is given by

$$\mathbb{E}[\|X - \Pi_{V^*}(X)\|_2^2] = \sum_{j=r+1}^d \gamma_j^2(\Sigma), \quad (8.7)$$

where $\{\gamma_j(\Sigma)\}_{j=1}^d$ are the ordered eigenvalues of Σ . See Exercise 8.4 for further exploration of these and other properties.

The problem of face analysis provides an interesting illustration of PCA for data compression. Consider a large database of face images, such as those illustrated in Figure 8.2(a). Taken from the Yale Face Database, each image is gray-scale with dimensions 243×320 . By vectorizing each image, we obtain a vector x in $d = 243 \times 320 = 77\,760$ dimensions. We compute the average image $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$ and the sample covariance matrix $\widehat{\Sigma} = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})(x_i - \bar{x})^T$ based on $n = 165$ samples. Figure 8.2(b) shows the relatively fast decay of the first 100 eigenvalues of this sample covariance matrix. Figure 8.2(c) shows the average face (top left image) along with the first 24 “eigenfaces”, meaning the

top 25 eigenvectors of the sample covariance matrix, each converted back to a 243×320 image. Finally, for a particular sample, Figure 8.2(d) shows a sequence of reconstructions of a given face, starting with the average face (top left image), and followed by the average face in conjunction with principal components 1 through 24. ♣

Principal component analysis can also be used for estimation in mixture models.

Example 8.3 (PCA for Gaussian mixture models) Let $\phi(\cdot; \mu, \Sigma)$ denote the density of a Gaussian random vector with mean vector $\mu \in \mathbb{R}^d$ and covariance matrix $\Sigma \in \mathcal{S}_+^{d \times d}$. A two-component Gaussian mixture model with isotropic covariance structure is a random vector $X \in \mathbb{R}^d$ drawn according to the density

$$f(x; \theta) = \alpha \phi(x; -\theta^*, \sigma^2 \mathbf{I}_d) + (1 - \alpha) \phi(x; \theta^*, \sigma^2 \mathbf{I}_d), \quad (8.8)$$

where $\theta^* \in \mathbb{R}^d$ is a vector parameterizing the means of the two Gaussian components, $\alpha \in (0, 1)$ is a mixture weight and $\sigma > 0$ is a dispersion term. Figure 8.3 provides an illustration of such a mixture model in $d = 2$ dimensions, with mean vector $\theta^* = [0.6 \ -0.6]^T$, standard deviation $\sigma = 0.4$ and weight $\alpha = 0.4$. Given samples $\{x_i\}_{i=1}^n$ drawn from such a model, a natural goal is to estimate the mean vector θ^* . Principal component analysis provides a natural method for doing so. In particular, a straightforward calculation yields that the second-moment matrix

$$\Gamma := \mathbb{E}[X \otimes X] = \theta^* \otimes \theta^* + \sigma^2 \mathbf{I}_d,$$

where $X \otimes X := XX^T$ is the $d \times d$ rank-one outer product matrix. Thus, we see that θ^* is proportional to the maximal eigenvector of Γ . Consequently, a reasonable estimator $\hat{\theta}$ is

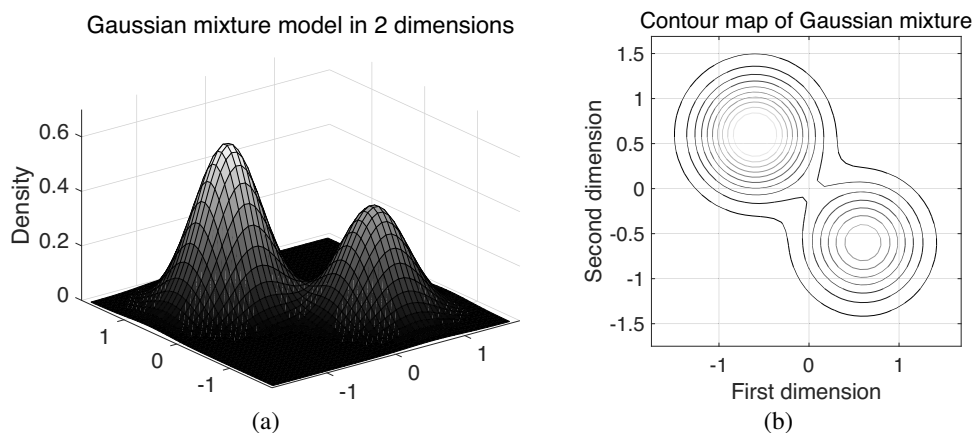


Figure 8.3 Use of PCA for Gaussian mixture models. (a) Density function of a two-component Gaussian mixture (8.8) with mean vector $\theta^* = [0.6 \ -0.6]^T$, standard deviation $\sigma = 0.4$ and weight $\alpha = 0.4$. (b) Contour plots of the density function, which provide intuition as to why PCA should be useful in recovering the mean vector θ^* .

given by the maximal eigenvector of the sample second moment² matrix $\widehat{\mathbf{T}} = \frac{1}{n} \sum_{i=1}^n x_i x_i^T$. We study the properties of this estimator in Exercise 8.6. ♣

8.1.2 Perturbations of eigenvalues and eigenspaces

Thus far, we have seen that the eigenvectors of population and sample covariance matrices are interesting objects with a range of uses. In practice, PCA is always applied to the sample covariance matrix, and the central question of interest is how well the sample-based eigenvectors approximate those of the population covariance.

Before addressing this question, let us make a brief detour into matrix perturbation theory. Let us consider the following general question: given a symmetric matrix \mathbf{R} , how does its eigenstructure relate to the perturbed matrix $\mathbf{Q} = \mathbf{R} + \mathbf{P}$? Here \mathbf{P} is another symmetric matrix, playing the role of the perturbation. It turns out that the eigenvalues of \mathbf{Q} and \mathbf{R} are related in a straightforward manner. Understanding how the eigenspaces change, however, requires some more care.

Let us begin with changes in the eigenvalues. From the standard variational definition of the maximum eigenvalue, we have

$$\gamma_1(\mathbf{Q}) = \max_{v \in \mathbb{S}^{d-1}} \langle v, (\mathbf{R} + \mathbf{P})v \rangle \leq \max_{v \in \mathbb{S}^{d-1}} \langle v, \mathbf{R}v \rangle + \max_{v \in \mathbb{S}^{d-1}} \langle v, \mathbf{P}v \rangle \leq \gamma_1(\mathbf{R}) + \|\mathbf{P}\|_2.$$

Since the same argument holds with the roles of \mathbf{Q} and \mathbf{R} reversed, we conclude that $|\gamma_1(\mathbf{Q}) - \gamma_1(\mathbf{R})| \leq \|\mathbf{Q} - \mathbf{R}\|_2$. Thus, the maximum eigenvalues of \mathbf{Q} and \mathbf{R} can differ by at most the operator norm of their difference. More generally, we have

$$\max_{j=1, \dots, d} |\gamma_j(\mathbf{Q}) - \gamma_j(\mathbf{R})| \leq \|\mathbf{Q} - \mathbf{R}\|_2. \quad (8.9)$$

This bound is a consequence of a more general result known as *Weyl's inequality*; we work through its proof in Exercise 8.3.

Although eigenvalues are generically stable, the same does not hold for eigenvectors and eigenspaces, unless further conditions are imposed. The following example provides an illustration of such instability:

Example 8.4 (Sensitivity of eigenvectors) For a parameter $\epsilon \in [0, 1]$, consider the family of symmetric matrices

$$\mathbf{Q}_\epsilon := \begin{bmatrix} 1 & \epsilon \\ \epsilon & 1.01 \end{bmatrix} = \underbrace{\begin{bmatrix} 1 & 0 \\ 0 & 1.01 \end{bmatrix}}_{\mathbf{Q}_0} + \epsilon \underbrace{\begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix}}_{\mathbf{P}}. \quad (8.10)$$

By construction, the matrix \mathbf{Q}_ϵ is a perturbation of a diagonal matrix \mathbf{Q}_0 by an ϵ -multiple of the fixed matrix \mathbf{P} . Since $\|\mathbf{P}\|_2 = 1$, the magnitude of the perturbation is directly controlled by ϵ . On one hand, the eigenvalues remain stable to this perturbation: in terms of the shorthand $a = 1.01$, we have $\gamma(\mathbf{Q}_0) = \{1, a\}$ and

$$\gamma(\mathbf{Q}_\epsilon) = \left\{ \frac{1}{2}[(a+1) + \sqrt{(a-1)^2 + 4\epsilon^2}], \quad \frac{1}{2}[(a+1) - \sqrt{(a-1)^2 + 4\epsilon^2}] \right\}.$$

² This second-moment matrix coincides with the usual covariance matrix for the special case of an equally weighted mixture pair with $\alpha = 0.5$.

Thus, we find that

$$\max_{j=1,2} |\gamma_j(\mathbf{Q}_0) - \gamma_j(\mathbf{Q}_\epsilon)| = \frac{1}{2} \left| (a-1) - \sqrt{(a-1)^2 + 4\epsilon^2} \right| \leq \epsilon,$$

which confirms the validity of Weyl's inequality (8.9) in this particular case.

On the other hand, the maximal eigenvector of \mathbf{Q}_ϵ is very different from that of \mathbf{Q}_0 , even for relatively small values of ϵ . For $\epsilon = 0$, the matrix \mathbf{Q}_0 has the unique maximal eigenvector $v_0 = [0 \ 1]^T$. However, if we set $\epsilon = 0.01$, a numerical calculation shows that the maximal eigenvector of \mathbf{Q}_ϵ is $v_\epsilon \approx [0.53 \ 0.85]^T$. Note that $\|v - v_\epsilon\|_2 \gg \epsilon$, showing that eigenvectors can be extremely sensitive to perturbations. ♣

What is the underlying problem? The issue is that, while \mathbf{Q}_0 has a unique maximal eigenvector, the gap between the largest eigenvalue $\gamma_1(\mathbf{Q}_0) = 1.01$ and the second largest eigenvalue $\gamma_2(\mathbf{Q}_0) = 1$ is very small. Consequently, even small perturbations of the matrix lead to “mixing” between the spaces spanned by the top and second largest eigenvectors. On the other hand, if this eigengap can be bounded away from zero, then it turns out that we can guarantee stability of the eigenvectors. We now turn to this type of theory.

8.2 Bounds for generic eigenvectors

We begin our exploration of eigenvector bounds with the generic case, in which no additional structure is imposed on the eigenvectors. In later sections, we turn to structured variants of eigenvector estimation.

8.2.1 A general deterministic result

Consider a symmetric positive semidefinite matrix $\mathbf{\Sigma}$ with eigenvalues ordered as

$$\gamma_1(\mathbf{\Sigma}) \geq \gamma_2(\mathbf{\Sigma}) \geq \gamma_3(\mathbf{\Sigma}) \geq \cdots \geq \gamma_d(\mathbf{\Sigma}) \geq 0.$$

Let $\theta^* \in \mathbb{R}^d$ denote its maximal eigenvector, assumed to be unique. Now consider a perturbed version $\widehat{\mathbf{\Sigma}} = \mathbf{\Sigma} + \mathbf{P}$ of the original matrix. As suggested by our notation, in the context of PCA, the original matrix corresponds to the population covariance matrix, whereas the perturbed matrix corresponds to the sample covariance. However, at least for the time being, our theory should be viewed as general.

As should be expected based on Example 8.4, any theory relating the maximum eigenvectors of $\mathbf{\Sigma}$ and $\widehat{\mathbf{\Sigma}}$ should involve the *eigengap* $\nu := \gamma_1(\mathbf{\Sigma}) - \gamma_2(\mathbf{\Sigma})$, assumed to be strictly positive. In addition, the following result involves the transformed perturbation matrix

$$\widetilde{\mathbf{P}} := \mathbf{U}^T \mathbf{P} \mathbf{U} = \begin{bmatrix} \tilde{p}_{11} & \tilde{\mathbf{p}}^T \\ \tilde{\mathbf{p}} & \widetilde{\mathbf{P}}_{22} \end{bmatrix}, \quad (8.11)$$

where $\tilde{p}_{11} \in \mathbb{R}$, $\tilde{\mathbf{p}} \in \mathbb{R}^{d-1}$ and $\widetilde{\mathbf{P}}_{22} \in \mathbb{R}^{(d-1) \times (d-1)}$. Here \mathbf{U} is an orthonormal matrix with the eigenvectors of $\mathbf{\Sigma}$ as its columns.

Theorem 8.5 Consider a positive semidefinite matrix Σ with maximum eigenvector $\theta^* \in \mathbb{S}^{d-1}$ and eigengap $\nu = \gamma_1(\Sigma) - \gamma_2(\Sigma) > 0$. Given any matrix $\mathbf{P} \in \mathbb{S}^{d \times d}$ such that $\|\mathbf{P}\|_2 < \nu/2$, the perturbed matrix $\widehat{\Sigma} := \Sigma + \mathbf{P}$ has a unique maximal eigenvector $\widehat{\theta}$ satisfying the bound

$$\|\widehat{\theta} - \theta^*\|_2 \leq \frac{2\|\mathbf{P}\|_2}{\nu - 2\|\mathbf{P}\|_2}. \quad (8.12)$$

In general, this bound is sharp in the sense that there are problems for which the requirement $\|\mathbf{P}\|_2 < \nu/2$ cannot be loosened. As an example, suppose that $\Sigma = \text{diag}\{2, 1\}$ so that $\nu = 2 - 1 = 1$. Given $\mathbf{P} = \text{diag}\{-\frac{1}{2}, +\frac{1}{2}\}$, the perturbed matrix $\widehat{\Sigma} = \Sigma + \mathbf{P} = \frac{3}{2}\mathbf{I}_2$ no longer has a unique maximal eigenvector. Note that this counterexample lies just at the boundary of our requirement, since $\|\mathbf{P}\|_2 = \frac{1}{2} = \frac{\nu}{2}$.

Proof Our proof is variational in nature, based on the optimization problems that characterize the maximal eigenvectors of the matrices Σ and $\widehat{\Sigma}$, respectively. Define the error vector $\widehat{\Delta} = \widehat{\theta} - \theta^*$, and the function

$$\Psi(\Delta; \mathbf{P}) := \langle \Delta, \mathbf{P}\Delta \rangle + 2\langle \Delta, \mathbf{P}\theta^* \rangle. \quad (8.13)$$

In parallel to our analysis of sparse linear regression from Chapter 7, the first step in our analysis is to prove the *basic inequality for PCA*. For future reference, we state this inequality in a slightly more general form than required for the current proof. In particular, given any subset $C \subseteq \mathbb{S}^{d-1}$, let θ^* and $\widehat{\theta}$ maximize the quadratic objectives

$$\max_{\theta \in C} \langle \theta, \Sigma\theta \rangle \quad \text{and} \quad \max_{\theta \in C} \langle \theta, \widehat{\Sigma}\theta \rangle, \quad (8.14)$$

respectively. The current proof involves the choice $C = \mathbb{S}^{d-1}$.

It is convenient to bound the distance between $\widehat{\theta}$ and θ^* in terms of the inner product $\varrho := \langle \widehat{\theta}, \theta^* \rangle$. Due to the sign ambiguity in eigenvector estimation, we may assume without loss of generality that $\widehat{\theta}$ is chosen such that $\varrho \in [0, 1]$.

Lemma 8.6 (PCA basic inequality) Given a matrix Σ with eigengap $\nu > 0$, the error $\widehat{\Delta} = \widehat{\theta} - \theta^*$ is bounded as

$$\nu(1 - \langle \widehat{\theta}, \theta^* \rangle^2) \leq |\Psi(\widehat{\Delta}; \mathbf{P})|. \quad (8.15)$$

Taking this inequality as given for the moment, the remainder of the proof is straightforward. Recall the transformation $\widetilde{\mathbf{P}} = \mathbf{U}^T \mathbf{P} \mathbf{U}$, or equivalently $\mathbf{P} = \mathbf{U} \widetilde{\mathbf{P}} \mathbf{U}^T$. Substituting this expression into equation (8.13) yields

$$\Psi(\widehat{\Delta}; \mathbf{P}) = \langle \mathbf{U}^T \widehat{\Delta}, \widetilde{\mathbf{P}} \mathbf{U}^T \widehat{\Delta} \rangle + 2\langle \mathbf{U}^T \widehat{\Delta}, \widetilde{\mathbf{P}} \mathbf{U}^T \theta^* \rangle. \quad (8.16)$$

In terms of the inner product $\varrho = \langle \widehat{\theta}, \theta^* \rangle$, we may write $\widehat{\theta} = \varrho \theta^* + \sqrt{1 - \varrho^2} z$, where $z \in \mathbb{R}^d$ is a vector orthogonal to θ^* . Since the matrix \mathbf{U} is orthonormal with its first column given by

θ^* , we have $\mathbf{U}^T \theta^* = e_1$. Letting $\mathbf{U}_2 \in \mathbb{R}^{d \times (d-1)}$ denote the submatrix formed by the remaining $d - 1$ eigenvectors and defining the vector $\tilde{z} = U_2^T z \in \mathbb{R}^{d-1}$, we can write

$$\mathbf{U}^T \widehat{\Delta} = \begin{bmatrix} (\varrho - 1) & (1 - \varrho^2)^{\frac{1}{2}} \tilde{z} \end{bmatrix}^T.$$

Substituting these relations into equation (8.16) yields that

$$\begin{aligned} \Psi(\widehat{\Delta}; \mathbf{P}) &= (\varrho - 1)^2 \tilde{p}_{11} + 2(\varrho - 1) \sqrt{1 - \varrho^2} \langle \tilde{z}, \tilde{p} \rangle + (1 - \varrho^2) \langle \tilde{z}, \tilde{\mathbf{P}}_{22} \tilde{z} \rangle \\ &\quad + 2(\varrho - 1) \tilde{p}_{11} + 2 \sqrt{1 - \varrho^2} \langle \tilde{z}, \tilde{p} \rangle \\ &= (\varrho^2 - 1) \tilde{p}_{11} + 2\varrho \sqrt{1 - \varrho^2} \langle \tilde{z}, \tilde{p} \rangle + (1 - \varrho^2) \langle \tilde{z}, \tilde{\mathbf{P}}_{22} \tilde{z} \rangle. \end{aligned}$$

Putting together the pieces, since $\|\tilde{z}\|_2 \leq 1$ and $|\tilde{p}_{11}| \leq \|\tilde{\mathbf{P}}\|_2$, we have

$$|\Psi(\widehat{\Delta}; \mathbf{P})| \leq 2(1 - \varrho^2) \|\tilde{\mathbf{P}}\|_2 + 2\varrho \sqrt{1 - \varrho^2} \|\tilde{p}\|_2.$$

Combined with the basic inequality (8.15), we find that

$$\nu(1 - \varrho^2) \leq 2(1 - \varrho^2) \|\mathbf{P}\|_2 + 2\varrho \sqrt{1 - \varrho^2} \|\tilde{p}\|_2.$$

Whenever $\nu > 2\|\mathbf{P}\|_2$, this inequality implies that $\sqrt{1 - \varrho^2} \leq \frac{2\varrho\|\tilde{p}\|_2}{\nu - 2\|\mathbf{P}\|_2}$. Noting that $\|\widehat{\Delta}\|_2 = \sqrt{2(1 - \varrho)}$, we thus conclude that

$$\|\widehat{\Delta}\|_2 \leq \frac{\sqrt{2}\varrho}{\sqrt{1 + \varrho}} \left(\frac{2\|\tilde{p}\|_2}{\nu - 2\|\mathbf{P}\|_2} \right) \leq \frac{2\|\tilde{p}\|_2}{\nu - 2\|\mathbf{P}\|_2},$$

where the final step follows since $2\varrho^2 \leq 1 + \varrho$ for all $\varrho \in [0, 1]$.

Let us now return to prove the PCA basic inequality (8.15).

Proof of Lemma 8.6: Since $\widehat{\theta}$ and θ^* are optimal and feasible, respectively, for the programs (8.14), we are guaranteed that $\langle \theta^*, \widehat{\Sigma} \theta^* \rangle \leq \langle \widehat{\theta}, \widehat{\Sigma} \widehat{\theta} \rangle$. Defining the matrix perturbation $\mathbf{P} = \widehat{\Sigma} - \Sigma$, we have

$$\langle \Sigma, \theta^* \otimes \theta^* - \widehat{\theta} \otimes \widehat{\theta} \rangle \leq -\langle \mathbf{P}, \theta^* \otimes \theta^* - \widehat{\theta} \otimes \widehat{\theta} \rangle,$$

where $\langle \mathbf{A}, \mathbf{B} \rangle$ is the trace inner product, and $a \otimes a = aa^T$ denotes the rank-one outer product. Following some simple algebra, the right-hand side is seen to be equal to $-\Psi(\widehat{\Delta}; \mathbf{P})$. The final step is to show that

$$\langle \Sigma, \theta^* \otimes \theta^* - \widehat{\theta} \otimes \widehat{\theta} \rangle \geq \frac{\nu}{2} \|\widehat{\Delta}\|_2^2. \quad (8.17)$$

Recall the representation $\widehat{\theta} = \varrho \theta^* + (\sqrt{1 - \varrho^2}) z$, where the vector $z \in \mathbb{R}^d$ is orthogonal to θ^* , and $\varrho \in [0, 1]$. Using the shorthand notation $\gamma_j \equiv \gamma_j(\Sigma)$ for $j = 1, 2$, define the matrix $\mathbf{\Gamma} = \Sigma - \gamma_1(\theta^* \otimes \theta^*)$, and note that $\mathbf{\Gamma} \theta^* = 0$ and $\|\mathbf{\Gamma}\|_2 \leq \gamma_2$ by construction. Consequently, we can write

$$\begin{aligned} \langle \Sigma, \theta^* \otimes \theta^* - \widehat{\theta} \otimes \widehat{\theta} \rangle &= \gamma_1 \langle \theta^* \otimes \theta^*, \theta^* \otimes \theta^* - \widehat{\theta} \otimes \widehat{\theta} \rangle + \langle \mathbf{\Gamma}, \theta^* \otimes \theta^* - \widehat{\theta} \otimes \widehat{\theta} \rangle \\ &= (1 - \varrho^2) \{ \gamma_1 - \langle \mathbf{\Gamma}, z \otimes z \rangle \}. \end{aligned}$$

Since $\|\Gamma\|_2 \leq \gamma_2$, we have $|\langle \Gamma, z \otimes z \rangle| \leq \gamma_2$. Putting together the pieces, we have shown that

$$\langle \Sigma, \theta^* \otimes \theta^* - \widehat{\theta} \otimes \widehat{\theta} \rangle \geq (1 - \varrho^2) \{\gamma_1 - \gamma_2\} = (1 - \varrho^2) \nu,$$

from which the claim (8.15) follows. \square

8.2.2 Consequences for a spiked ensemble

Theorem 8.5 applies to any form of matrix perturbation. In the context of principal component analysis, this perturbation takes a very specific form—namely, as the difference between the sample and population covariance matrices. More concretely, suppose that we have drawn n i.i.d. samples $\{x_i\}_{i=1}^n$ from a zero-mean random vector with covariance Σ . Principal component analysis is then based on the eigenstructure of the sample covariance matrix $\widehat{\Sigma} = \frac{1}{n} \sum_{i=1}^n x_i x_i^T$, and the goal is to draw conclusions about the eigenstructure of the population matrix.

In order to bring sharper focus to this issue, let us study how PCA behaves for a very simple class of covariance matrices, known as spiked covariance matrices. A sample $x_i \in \mathbb{R}^d$ from the *spiked covariance ensemble* takes the form

$$x_i \stackrel{d}{=} \sqrt{\nu} \xi_i \theta^* + w_i, \quad (8.18)$$

where $\xi_i \in \mathbb{R}$ is a zero-mean random variable with unit variance, and $w_i \in \mathbb{R}^d$ is a random vector independent of ξ_i , with zero mean and covariance matrix \mathbf{I}_d . Overall, the random vector x_i has zero mean, and a covariance matrix of the form

$$\Sigma := \nu \theta^* (\theta^*)^T + \mathbf{I}_d. \quad (8.19)$$

By construction, for any $\nu > 0$, the vector θ^* is the unique maximal eigenvector of Σ with eigenvalue $\gamma_1(\Sigma) = \nu + 1$. All other eigenvalues of Σ are located at 1, so that we have an eigengap $\gamma_1(\Sigma) - \gamma_2(\Sigma) = \nu$.

In the following result, we say that the vector $x_i \in \mathbb{R}^d$ has sub-Gaussian tails if both ξ_i and w_i are sub-Gaussian with parameter at most one.

Corollary 8.7 *Given i.i.d. samples $\{x_i\}_{i=1}^n$ from the spiked covariance ensemble (8.18) with sub-Gaussian tails, suppose that $n > d$ and $\sqrt{\frac{\nu+1}{\nu^2}} \sqrt{\frac{d}{n}} \leq \frac{1}{128}$. Then, with probability at least $1 - c_1 e^{-c_2 n \min\{\sqrt{\nu\delta}, \nu\delta^2\}}$, there is a unique maximal eigenvector $\widehat{\theta}$ of the sample covariance matrix $\widehat{\Sigma} = \frac{1}{n} \sum_{i=1}^n x_i x_i^T$ such that*

$$\|\widehat{\theta} - \theta^*\|_2 \leq c_0 \sqrt{\frac{\nu+1}{\nu^2}} \sqrt{\frac{d}{n}} + \delta. \quad (8.20)$$

Figure 8.4 shows the results of simulations that confirm the qualitative scaling predicted by Corollary 8.7. In each case, we drew $n = 500$ samples from a spiked covariance matrix with the signal-to-noise parameter ν ranging over the interval $[0.75, 5]$. We then computed the ℓ_2 -distance $\|\widehat{\theta} - \theta^*\|_2$ between the maximal eigenvectors of the sample and population

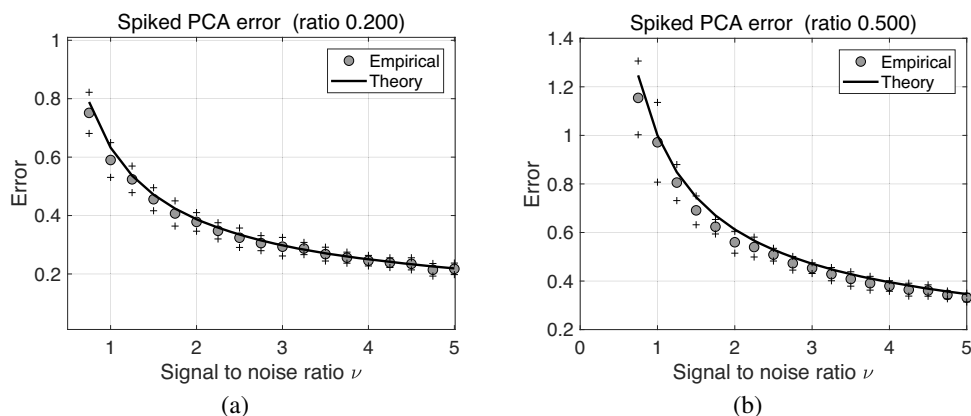


Figure 8.4 Plots of the error $\|\hat{\theta} - \theta^*\|_2$ versus the signal-to-noise ratio, as measured by the eigengap ν . Both plots are based on a sample size $n = 500$. Dots show the average of 100 trials, along with the standard errors (crosses). The full curve shows the theoretical bound $\sqrt{\frac{\nu+1}{\nu^2}} \sqrt{\frac{d}{n}}$. (a) Dimension $d = 100$. (b) Dimension $d = 250$.

covariances, respectively, performing $T = 100$ trials for each setting of ν . The circles in Figure 8.4 show the empirical means, along with standard errors in crosses, whereas the solid curve corresponds to the theoretical prediction $\sqrt{\frac{\nu+1}{\nu^2}} \sqrt{\frac{d}{n}}$. Note that Corollary 8.7 predicts this scaling, but with a looser leading constant ($c_0 > 1$). As shown by Figure 8.4, Corollary 8.7 accurately captures the scaling behavior of the error as a function of the signal-to-noise ratio.

Proof Let $\mathbf{P} = \widehat{\Sigma} - \Sigma$ be the difference between the sample and population covariance matrices. In order to apply Theorem 8.5, we need to upper bound the quantities $\|\mathbf{P}\|_2$ and $\|\tilde{p}\|_2$. Defining the random vector $\tilde{w} := \frac{1}{n} \sum_{i=1}^n \xi_i w_i$, the perturbation matrix \mathbf{P} can be decomposed as

$$\mathbf{P} = \underbrace{\nu \left(\frac{1}{n} \sum_{i=1}^n \xi_i^2 - 1 \right) \theta^* (\theta^*)^T}_{\mathbf{P}_1} + \underbrace{\sqrt{\nu} (\tilde{w} (\theta^*)^T + \theta^* \tilde{w}^T)}_{\mathbf{P}_2} + \underbrace{\left(\frac{1}{n} \sum_{i=1}^n w_i w_i^T - \mathbf{I}_d \right)}_{\mathbf{P}_3}. \quad (8.21)$$

Since $\|\theta^*\|_2 = 1$, the operator norm of \mathbf{P} can be bounded as

$$\|\mathbf{P}\|_2 \leq \nu \left| \frac{1}{n} \sum_{i=1}^n \xi_i^2 - 1 \right| + 2\sqrt{\nu} \|\tilde{w}\|_2 + \left\| \frac{1}{n} \sum_{i=1}^n w_i w_i^T - \mathbf{I}_d \right\|_2. \quad (8.22a)$$

Let us derive a similar upper bound on $\|\tilde{p}\|_2$ using the decomposition (8.11). Since θ^* is the unique maximal eigenvector of Σ , it forms the first column of the matrix \mathbf{U} . Let $\mathbf{U}_2 \in \mathbb{R}^{d \times (d-1)}$ denote the matrix formed of the remaining $(d-1)$ columns. With this notation, we have $\tilde{p} = \mathbf{U}_2^T \mathbf{P} \theta^*$. Using the decomposition (8.21) of the perturbation matrix and the fact that $\mathbf{U}_2^T \theta^* = 0$, we find that $\tilde{p} = \sqrt{\nu} \mathbf{U}_2^T \tilde{w} + \frac{1}{n} \sum_{i=1}^n \mathbf{U}_2^T w_i \langle w_i, \theta^* \rangle$. Since \mathbf{U}_2 has orthonormal

columns, we have $\|\mathbf{U}_2^T \tilde{\mathbf{w}}\|_2 \leq \|\tilde{\mathbf{w}}\|_2$ and also

$$\left\| \sum_{i=1}^n \mathbf{U}_2^T w_i \langle w_i, \theta^* \rangle \right\|_2 = \sup_{\|v\|_2=1} \left| (\mathbf{U}_2 v)^T \left(\sum_{i=1}^n w_i w_i^T - \mathbf{I}_d \right) \theta^* \right| \leq \left\| \frac{1}{n} \sum_{i=1}^n w_i w_i^T - \mathbf{I}_d \right\|_2.$$

Putting together the pieces, we have shown that

$$\|\tilde{\mathbf{p}}\|_2 \leq \sqrt{\nu} \|\tilde{\mathbf{w}}\|_2 + \left\| \frac{1}{n} \sum_{i=1}^n w_i w_i^T - \mathbf{I}_d \right\|_2. \quad (8.22b)$$

The following lemma allows us to control the quantities appearing the bounds (8.22a) and (8.22b):

Lemma 8.8 *Under the conditions of Corollary 8.7, we have*

$$\mathbb{P} \left[\left| \frac{1}{n} \sum_{i=1}^n \xi_i^2 - 1 \right| \geq \delta_1 \right] \leq 2e^{-c_2 n \min\{\delta_1, \delta_1^2\}}, \quad (8.23a)$$

$$\mathbb{P} [\|\tilde{\mathbf{w}}\|_2 \geq 2\sqrt{\frac{d}{n}} + \delta_2] \leq 2e^{-c_2 n \min\{\delta_2, \delta_2^2\}} \quad (8.23b)$$

and

$$\mathbb{P} \left[\left\| \frac{1}{n} \sum_{i=1}^n w_i w_i^T - \mathbf{I}_d \right\|_2 \geq c_3 \sqrt{\frac{d}{n}} + \delta_3 \right] \leq 2e^{-c_2 n \min\{\delta_3, \delta_3^2\}}. \quad (8.23c)$$

We leave the proof of this claim as an exercise, since it is straightforward application of results and techniques from previous chapters. For future reference, we define

$$\phi(\delta_1, \delta_2, \delta_3) := 2e^{-c_2 n \min\{\delta_1, \delta_1^2\}} + 2e^{-c_2 n \min\{\delta_2, \delta_2^2\}} + 2e^{-c_2 n \min\{\delta_3, \delta_3^2\}}, \quad (8.24)$$

corresponding to the probability with which at least one of the bounds in Lemma 8.8 is violated.

In order to apply Theorem 8.5, we need to first show that $\|\mathbf{P}\|_2 < \frac{\nu}{4}$ with high probability. Beginning with the inequality (8.22a) and applying Lemma 8.8 with $\delta_1 = \frac{1}{16}$, $\delta_2 = \frac{\delta}{4\sqrt{\nu}}$ and $\delta_3 = \delta/16 \in (0, 1)$, we have

$$\|\mathbf{P}\|_2 \leq \frac{\nu}{16} + 8(\sqrt{\nu} + 1) \sqrt{\frac{d}{n}} + \delta \leq \frac{\nu}{16} + 16(\sqrt{\nu} + 1) \sqrt{\frac{d}{n}} + \delta$$

with probability at least $1 - \phi(\frac{1}{4}, \frac{\delta}{3\sqrt{\nu}}, \frac{\delta}{16})$. Consequently, as long as $\sqrt{\frac{\nu+1}{\nu^2}} \sqrt{\frac{d}{n}} \leq \frac{1}{128}$, we have

$$\|\mathbf{P}\|_2 \leq \frac{3}{16} \nu + \delta < \frac{\nu}{4}, \quad \text{for all } \delta \in (0, \frac{\nu}{16}).$$

It remains to bound $\|\tilde{\mathbf{p}}\|_2$. Applying Lemma 8.8 to the inequality (8.22b) with the previously specified choices of $(\delta_1, \delta_2, \delta_3)$, we have

$$\|\tilde{\mathbf{p}}\|_2 \leq 2(\sqrt{\nu} + 1) \sqrt{\frac{d}{n}} + \delta \leq 4\sqrt{\nu+1} \sqrt{\frac{d}{n}} + \delta$$

with probability at least $1 - \phi(\frac{1}{4}, \frac{\delta}{3\sqrt{n}}, \frac{\delta}{16})$. We have shown that conditions of Theorem 8.5 are satisfied, so that the claim (8.20) follows as a consequence of the bound (8.12). \square

8.3 Sparse principal component analysis

Note that Corollary 8.7 requires that the sample size n be larger than the dimension d in order for ordinary PCA to perform well. One might wonder whether this requirement is fundamental: does PCA still perform well in the high-dimensional regime $n < d$?

The answer to this question turns out to be a dramatic “no”. As discussed at more length in the bibliography section, for any fixed signal-to-noise ratio, if the ratio d/n stays suitably bounded away from zero, then the eigenvectors of the sample covariance in a spiked covariance model become *asymptotically orthogonal* to their population analogs. Thus, the classical PCA estimate is no better than ignoring the data, and drawing a vector uniformly at random from the Euclidean sphere. Given this total failure of classical PCA, a next question to ask is whether the eigenvectors might be estimated consistently using a method more sophisticated than PCA. This question also has a negative answer: as we discuss in Chapter 15, for the standard spiked model (8.18), it can be shown via the framework of minimax theory that *no method* can produce consistent estimators of the population eigenvectors when d/n stays bounded away from zero. See Example 15.19 in Chapter 15 for the details of this minimax lower bound.

In practice, however, it is often reasonable to impose structure on eigenvectors, and this structure can be exploited to develop effective estimators even when $n < d$. Perhaps the simplest such structure is that of sparsity in the eigenvectors, which allows for both effective estimation in high-dimensional settings, as well as increased interpretability. Accordingly, this section is devoted to the sparse version of principal component analysis.

Let us illustrate the idea of sparse eigenanalysis by revisiting the eigenfaces from Example 8.2.

Example 8.9 (Sparse eigenfaces) We used the images from the Yale Face Database to set up a PCA problem in $d = 77\,760$ dimensions. In this example, we used an iterative method to approximate sparse eigenvectors with at most $s = \lfloor 0.25d \rfloor = 19\,440$ non-zero coefficients. In particular, we applied a thresholded version of the matrix power method for computing sparse eigenvalues and eigenvectors. (See Exercise 8.5 for exploration of the standard matrix power method.)

Figure 8.5(a) shows the average face (top left image), along with approximations to the first 24 sparse eigenfaces. Each sparse eigenface was restricted to have at most 25% of its pixels non-zero, corresponding to a savings of a factor of 4 in storage. Note that the sparse eigenfaces are more localized than their PCA analogs from Figure 8.2. Figure 8.5(b) shows reconstruction using the average face in conjunction with the first 100 sparse eigenfaces, which require equivalent storage (in terms of pixel values) to the first 25 regular eigenfaces. ♣

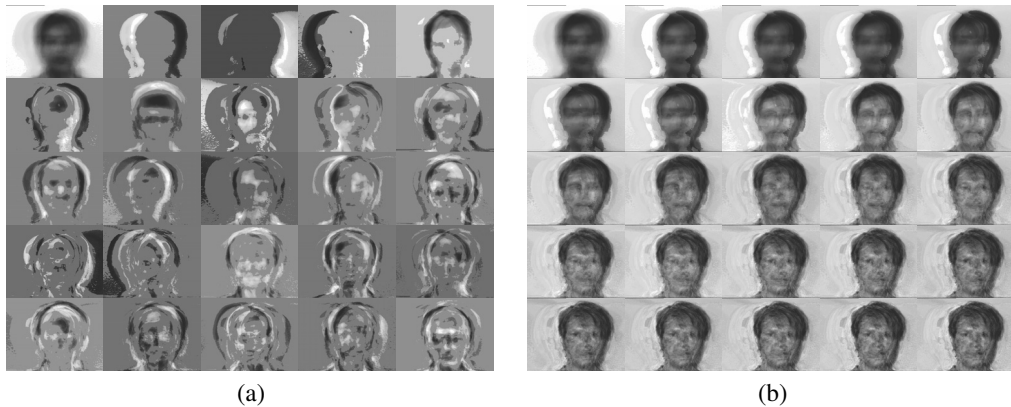


Figure 8.5 Illustration of sparse eigenanalysis for the Yale Face Database. (a) Average face (top left image), and approximations to the first 24 sparse eigenfaces, obtained by a greedy iterative thresholding procedure applied to the eigenvalue power method. Eigenfaces were restricted to have at most 25% of their pixels non-zero, corresponding to a 1/4 reduction in storage. (b) Reconstruction based on sparse eigenfaces.

8.3.1 A general deterministic result

We now turn to the question of how to estimate a maximal eigenvector that is known *a priori* to be sparse. A natural approach is to augment the quadratic objective function underlying classical PCA with an additional sparsity constraint or penalty. More concretely, we analyze both the constrained problem

$$\widehat{\theta} \in \arg \max_{\|\theta\|_2=1} \left\{ \langle \theta, \widehat{\Sigma} \theta \rangle \right\} \quad \text{such that } \|\theta\|_1 \leq R, \quad (8.25a)$$

as well as the penalized variant

$$\widehat{\theta} \in \arg \max_{\|\theta\|_2=1} \left\{ \langle \theta, \widehat{\Sigma} \theta \rangle - \lambda_n \|\theta\|_1 \right\} \quad \text{such that } \|\theta\|_1 \leq \left(\frac{n}{\log d} \right)^{1/4}. \quad (8.25b)$$

In our analysis of the constrained version (8.25a), we set $R = \|\theta^*\|_1$. The advantage of the penalized variant (8.25b) is that the regularization parameter λ_n can be chosen without knowledge of the true eigenvector θ^* . In both formulations, the matrix $\widehat{\Sigma}$ represents some type of approximation to the population covariance matrix Σ , with the sample covariance being a canonical example. Note that neither estimator is convex, since they involve maximization of a positive semidefinite quadratic form. Nonetheless, it is instructive to analyze them in order to understand the statistical behavior of sparse PCA, and in the exercises, we describe some relaxations of these non-convex programs.

Naturally, the proximity of $\widehat{\theta}$ to the maximum eigenvector θ^* of Σ depends on the perturbation matrix $\mathbf{P} := \widehat{\Sigma} - \Sigma$. How to measure the effect of the perturbation? As will become clear, much of our analysis of ordinary PCA can be modified in a relatively straightforward way so as to obtain results for the sparse version. In particular, a central object in our analysis of ordinary PCA was the basic inequality stated in Lemma 8.6: it shows that the perturbation

matrix enters via the function

$$\Psi(\Delta; \mathbf{P}) := \langle \Delta, \mathbf{P}\Delta \rangle + 2 \langle \Delta, \mathbf{P}\theta^* \rangle.$$

As with our analysis of PCA, our general deterministic theorem for sparse PCA involves imposing a form of uniform control on $\Psi(\Delta; \mathbf{P})$ as Δ ranges over all vectors of the form $\theta - \theta^*$ with $\theta \in \mathbb{S}^{d-1}$. The sparsity constraint enters in the form of this uniform bound that we assume. More precisely, letting $\varphi_\nu(n, d)$ and $\psi_\nu(n, d)$ be non-negative functions of the eigengap ν , sample size and dimension, we assume that there exists a universal constant $c_0 > 0$ such that

$$\sup_{\substack{\Delta = \theta - \theta^* \\ \|\theta\|_2 = 1}} |\Psi(\Delta; \mathbf{P})| \leq c_0 \nu \|\Delta\|_2^2 + \varphi_\nu(n, d) \|\Delta\|_1 + \psi_\nu^2(n, d) \|\Delta\|_1^2. \quad (8.26)$$

As a concrete example, for a sparse version of the spiked PCA ensemble (8.18) with sub-Gaussian tails, this condition is satisfied with high probability with $\varphi_\nu^2(n, d) \asymp (\nu + 1) \frac{\log d}{n}$ and $\psi_\nu^2(n, d) \asymp \frac{1}{\nu} \frac{\log d}{n}$. This fact will be established in the proof of Corollary 8.12 to follow.

Theorem 8.10 *Given a matrix Σ with a unique, unit-norm, s -sparse maximal eigenvector θ^* with eigengap ν , let $\widehat{\Sigma}$ be any symmetric matrix satisfying the uniform deviation condition (8.26) with constant $c_0 < \frac{1}{6}$, and $16s\psi_\nu^2(n, d) \leq c_0\nu$.*

(a) *For any optimal solution $\widehat{\theta}$ to the constrained program (8.25a) with $R = \|\theta^*\|_1$,*

$$\min \left\{ \|\widehat{\theta} - \theta^*\|_2, \|\widehat{\theta} + \theta^*\|_2 \right\} \leq \frac{8}{\nu(1 - 4c_0)} \sqrt{s} \varphi_\nu(n, d). \quad (8.27)$$

(b) *Consider the penalized program (8.25b) with the regularization parameter lower bounded as $\lambda_n \geq 4 \left(\frac{n}{\log d} \right)^{1/4} \psi_\nu^2(n, d) + 2\varphi_\nu(n, d)$. Then any optimal solution $\widehat{\theta}$ satisfies the bound*

$$\min \left\{ \|\widehat{\theta} - \theta^*\|_2, \|\widehat{\theta} + \theta^*\|_2 \right\} \leq \frac{2 \left(\frac{\lambda_n}{\varphi_\nu(n, d)} + 4 \right)}{\nu(1 - 4c_0)} \sqrt{s} \varphi_\nu(n, d). \quad (8.28)$$

Proof We begin by analyzing the constrained estimator, and then describe the modifications necessary for the regularized version.

Argument for constrained estimator: Note that $\|\widehat{\theta}\|_1 \leq R = \|\theta^*\|_1$ by construction of the estimator, and moreover $\theta_{S^c}^* = 0$ by assumption. By splitting the ℓ_1 -norm into two components, indexed by S and S^c , respectively, it can be shown³ that the error $\widehat{\Delta} = \widehat{\theta} - \theta^*$ satisfies the inequality $\|\widehat{\Delta}_{S^c}\|_1 \leq \|\widehat{\Delta}_S\|_1$. So as to simplify our treatment of the regularized estimator, let us proceed by assuming only the weaker inequality $\|\widehat{\Delta}_{S^c}\|_1 \leq 3\|\widehat{\Delta}_S\|_1$, which implies that $\|\widehat{\Delta}\|_1 \leq 4\sqrt{s}\|\widehat{\Delta}\|_2$. Combining this inequality with the uniform bound (8.26) on Ψ , we find

³ We leave this calculation as an exercise for the reader: helpful details can be found in Chapter 7.

that

$$|\Psi(\hat{\Delta}; \mathbf{P})| \leq c_0 \nu \|\hat{\Delta}\|_2^2 + 4 \sqrt{s} \varphi_\nu(n, d) \|\hat{\Delta}\|_2 + 16 s \psi_\nu^2(n, d) \|\hat{\Delta}\|_2^2. \quad (8.29)$$

Substituting back into the basic inequality (8.15) and performing some algebra yields

$$\underbrace{\nu \left\{ \frac{1}{2} - c_0 - 16 \frac{s}{\nu} \psi_\nu^2(n, d) \right\}}_{\kappa} \|\hat{\Delta}\|_2^2 \leq 4 \sqrt{s} \varphi_\nu(n, d) \|\hat{\Delta}\|_2.$$

Note that our assumptions imply that $\kappa > \frac{1}{2}(1 - 4c_0) > 0$, so that the bound (8.27) follows after canceling a term $\|\hat{\Delta}\|_2$ and rearranging.

Argument for regularized estimator: We now turn to the regularized estimator (8.25b). With the addition of the regularizer, the basic inequality (8.15) now takes the slightly modified form

$$\frac{\nu}{2} \|\hat{\Delta}\|_2^2 - |\Psi(\hat{\Delta}; \mathbf{P})| \leq \lambda_n \{ \|\theta^*\|_1 - \|\hat{\theta}\|_1 \} \leq \lambda_n \{ \|\hat{\Delta}_S\|_1 - \|\hat{\Delta}_{S^c}\|_1 \}, \quad (8.30)$$

where the second inequality follows by the S -sparsity of θ^* and the triangle inequality (see Chapter 7 for details).

We claim that the error vector $\hat{\Delta}$ still satisfies a form of the cone inequality. Let us state this claim as a separate lemma.

Lemma 8.11 *Under the conditions of Theorem 8.10, the error vector $\hat{\Delta} = \hat{\theta} - \theta^*$ satisfies the cone inequality*

$$\|\hat{\Delta}_{S^c}\|_1 \leq 3 \|\hat{\Delta}_S\|_1 \quad \text{and hence} \quad \|\hat{\Delta}\|_1 \leq 4 \sqrt{s} \|\hat{\Delta}\|_2. \quad (8.31)$$

Taking this lemma as given, let us complete the proof of the theorem. Given Lemma 8.11, the previously derived upper bound (8.29) on $|\Psi(\hat{\Delta}; \mathbf{P})|$ is also applicable to the regularized estimator. Substituting this bound into our basic inequality, we find that

$$\underbrace{\nu \left\{ \frac{1}{2} - c_0 - \frac{16}{\nu} s \psi_\nu^2(n, d) \right\}}_{\kappa} \|\hat{\Delta}\|_2^2 \leq \sqrt{s} (\lambda_n + 4 \varphi_\nu(n, d)) \|\hat{\Delta}\|_2.$$

Our assumptions imply that $\kappa \geq \frac{1}{2}(1 - 4c_0) > 0$, from which claim (8.28) follows.

It remains to prove Lemma 8.11. Combining the uniform bound with the basic inequality (8.30)

$$0 \leq \underbrace{\nu \left(\frac{1}{2} - c_0 \right)}_{>0} \|\Delta\|_2^2 \leq \varphi_\nu(n, d) \|\Delta\|_1 + \psi_\nu^2(n, d) \|\Delta\|_1^2 + \lambda_n \{ \|\hat{\Delta}_S\|_1 - \|\hat{\Delta}_{S^c}\|_1 \}.$$

Introducing the shorthand $R = \left(\frac{n}{\log d} \right)^{1/4}$, the feasibility of $\hat{\theta}$ and θ^* implies that $\|\hat{\Delta}\|_1 \leq 2R$,

and hence

$$\begin{aligned} 0 &\leq \underbrace{\left\{ \varphi_v(n, d) + 2R\psi_v^2(n, d) \right\}}_{\leq \frac{4n}{2}} \|\hat{\Delta}\|_1 + \lambda_n \left\{ \|\hat{\Delta}_S\|_1 - \|\hat{\Delta}_{S^c}\|_1 \right\} \\ &\leq \lambda_n \left\{ \frac{3}{2} \|\hat{\Delta}_S\|_1 - \frac{1}{2} \|\hat{\Delta}_{S^c}\|_1 \right\}, \end{aligned}$$

and rearranging yields the claim. \square

8.3.2 Consequences for the spiked model with sparsity

Theorem 8.10 is a general deterministic guarantee that applies to any matrix with a sparse maximal eigenvector. In order to obtain more concrete results in a particular case, let us return to the spiked covariance model previously introduced in equation (8.18), and analyze a sparse variant of it. More precisely, consider a random vector $x_i \in \mathbb{R}^d$ generated from the usual spiked ensemble—namely, as $x_i \stackrel{d}{=} \sqrt{v} \xi_i \theta^* + w_i$, where $\theta^* \in \mathbb{S}^{d-1}$ is an s -sparse vector, corresponding to the maximal eigenvector of $\Sigma = \text{cov}(x_i)$. As before, we assume that both the random variable ξ_i and the random vector $w_i \in \mathbb{R}^d$ are independent, each sub-Gaussian with parameter 1, in which case we say that the random vector $x_i \in \mathbb{R}^d$ has sub-Gaussian tails.

Corollary 8.12 Consider n i.i.d. samples $\{x_i\}_{i=1}^n$ from an s -sparse spiked covariance matrix with eigengap $v > 0$ and suppose that $\frac{s \log d}{n} \leq c \min \left\{ 1, \frac{v^2}{v+1} \right\}$ for a sufficiently small constant $c > 0$. Then for any $\delta \in (0, 1)$, any optimal solution $\hat{\theta}$ to the constrained program (8.25a) with $R = \|\theta^*\|_1$, or to the penalized program (8.25b) with $\lambda_n = c_3 \sqrt{v+1} \left\{ \sqrt{\frac{\log d}{n}} + \delta \right\}$, satisfies the bound

$$\min \left\{ \|\hat{\theta} - \theta^*\|_2, \|\hat{\theta} + \theta^*\|_2 \right\} \leq c_4 \sqrt{\frac{v+1}{v^2}} \left\{ \sqrt{\frac{s \log d}{n}} + \delta \right\} \quad \text{for all } \delta \in (0, 1) \quad (8.32)$$

with probability at least $1 - c_1 e^{-c_2(n/s) \min\{\delta^2, v^2, v\}}$.

Proof Letting $\mathbf{P} = \hat{\Sigma} - \Sigma$ be the deviation between the sample and population covariance matrices, our goal is to show that $\Psi(\cdot, \mathbf{P})$ satisfies the uniform deviation condition (8.26). In particular, we claim that, uniformly over $\Delta \in \mathbb{R}^d$, we have

$$\underbrace{|\Psi(\Delta; \mathbf{P})|}_{c_0} \leq \underbrace{\frac{1}{8} v \|\Delta\|_2^2 + 16 \sqrt{v+1} \left\{ \sqrt{\frac{\log d}{n}} + \delta \right\} \|\Delta\|_1}_{\varphi_v(n, d)} + \underbrace{\frac{c'_3 \log d}{v n} \|\Delta\|_1^2}_{\psi_v^2(n, d)}, \quad (8.33)$$

with probability at least $1 - c_1 e^{-c_2 n \min\{\delta^2, v^2\}}$. Here (c_1, c_2, c'_3) are universal constants. Taking this intermediate claim as given, let us verify that the bound (8.32) follows as a consequence

of Theorem 8.10. We have

$$\frac{9s\psi_v^2(n, d)}{c_0} = \frac{72c'_3}{v} \frac{s \log d}{n} \leq v \left\{ 72c'_3 \frac{v+1}{v^2} \frac{s \log d}{n} \right\} \leq v,$$

using the assumed upper bound on the ratio $\frac{s \log d}{n}$ for a sufficiently small constant c . Consequently, the bound for the constrained estimator follows from Theorem 8.10. For the penalized estimator, there are a few other conditions to be verified: let us first check that $\|\theta^*\|_1 \leq v \sqrt{\frac{n}{\log d}}$. Since θ^* is s -sparse with $\|\theta^*\|_2 = 1$, it suffices to have $\sqrt{s} \leq v \sqrt{\frac{n}{\log d}}$, or equivalently $\frac{1}{v^2} \frac{s \log d}{n} \leq 1$, which follows from our assumptions. Finally, we need to check that λ_n satisfies the lower bound requirement in Theorem 8.10. We have

$$\begin{aligned} 4R\psi_v^2(n, d) + 2\varphi_v(n, d) &\leq 4v \sqrt{\frac{n}{\log d}} \frac{c'_3 \log d}{v} \frac{1}{n} + 24 \sqrt{v+1} \left\{ \sqrt{\frac{\log d}{n}} + \delta \right\} \\ &\leq \underbrace{c_3 \sqrt{v+1} \left\{ \sqrt{\frac{\log d}{n}} + \delta \right\}}_{\lambda_n} \end{aligned}$$

as required.

It remains to prove the uniform bound (8.33). Recall the decomposition $\mathbf{P} = \sum_{j=1}^3 \mathbf{P}_j$ given in equation (8.21). By linearity of the function Ψ in its second argument, this decomposition implies that $\Psi(\Delta; \mathbf{P}) = \sum_{j=1}^3 \Psi(\Delta; \mathbf{P}_j)$. We control each of these terms in turn.

Control of first component: Lemma 8.8 guarantees that $\left| \frac{1}{n} \sum_{i=1}^n \xi_i^2 - 1 \right| \leq \frac{1}{16}$ with probability at least $1 - 2e^{-cn}$. Conditioned on this bound, for any vector of the form $\Delta = \theta - \theta^*$ with $\theta \in \mathbb{S}^{d-1}$, we have

$$|\Psi(\Delta; \mathbf{P}_1)| \leq \frac{v}{16} \langle \Delta, \theta^* \rangle^2 = \frac{v}{16} (1 - \langle \theta^*, \theta \rangle)^2 \leq \frac{v}{32} \|\Delta\|_2^2, \quad (8.34)$$

where we have used the fact that $2(1 - \langle \theta^*, \theta \rangle)^2 \leq 2(1 - \langle \theta^*, \theta \rangle) = \|\Delta\|_2^2$.

Control of second component: We have

$$\begin{aligned} |\Psi(\Delta; \mathbf{P}_2)| &\leq 2\sqrt{v} \left\{ \langle \Delta, \bar{w} \rangle \langle \Delta, \theta^* \rangle + \langle \bar{w}, \Delta \rangle + \langle \theta^*, \bar{w} \rangle \langle \Delta, \theta^* \rangle \right\} \\ &\leq 4\sqrt{v} \|\Delta\|_1 \|\bar{w}\|_\infty + 2\sqrt{v} |\langle \theta^*, \bar{w} \rangle| \frac{\|\Delta\|_2^2}{2}. \end{aligned} \quad (8.35)$$

The following lemma provides control on the two terms in this upper bound:

Lemma 8.13 *Under the conditions of Corollary 8.12, we have*

$$\mathbb{P} \left[\|\bar{w}\|_\infty \geq 2\sqrt{\frac{\log d}{n}} + \delta \right] \leq c_1 e^{-c_2 n \delta^2} \quad \text{for all } \delta \in (0, 1), \text{ and} \quad (8.36a)$$

$$\mathbb{P} \left[|\langle \theta^*, \bar{w} \rangle| \geq \frac{\sqrt{\nu}}{32} \right] \leq c_1 e^{-c_2 n \nu}. \quad (8.36b)$$

We leave the proof of these bounds as an exercise for the reader, since they follow from standard results in Chapter 2. Combining Lemma 8.13 with the bound (8.35) yields

$$|\Psi(\Delta; \mathbf{P}_2)| \leq \frac{\nu}{32} \|\Delta\|_2^2 + 8\sqrt{\nu+1} \left\{ \sqrt{\frac{\log d}{n}} + \delta \right\} \|\Delta\|_1. \quad (8.37)$$

Control of third term: Recalling that $\mathbf{P}_3 = \frac{1}{n} \mathbf{W}^T \mathbf{W} - \mathbf{I}_d$, we have

$$|\Psi(\Delta; \mathbf{P}_3)| \leq |\langle \Delta, \mathbf{P}_3 \Delta \rangle| + 2\|\mathbf{P}_3 \theta^*\|_\infty \|\Delta\|_1. \quad (8.38)$$

Our final lemma controls the two terms in this bound:

Lemma 8.14 *Under the conditions of Corollary 8.12, for all $\delta \in (0, 1)$, we have*

$$\|\mathbf{P}_3 \theta^*\|_\infty \leq 2\sqrt{\frac{\log d}{n}} + \delta \quad (8.39a)$$

and

$$\sup_{\Delta \in \mathbb{R}^d} |\langle \Delta, \mathbf{P}_3 \Delta \rangle| \leq \frac{\nu}{16} \|\Delta\|_2^2 + \frac{c'_3}{\nu} \frac{\log d}{n} \|\Delta\|_1^2, \quad (8.39b)$$

where both inequalities hold with probability greater than $1 - c_1 e^{-c_2 n \min\{\nu, \nu^2, \delta^2\}}$.

Combining this lemma with our earlier inequality (8.38) yields the bound

$$|\Psi(\Delta; \mathbf{P}_3)| \leq \frac{\nu}{16} \|\Delta\|_2^2 + 8 \left\{ \sqrt{\frac{\log d}{n}} + \delta \right\} \|\Delta\|_1 + \frac{c'_3}{\nu} \frac{\log d}{n} \|\Delta\|_1^2. \quad (8.40)$$

Finally, combining the bounds (8.34), (8.37) and (8.40) yields the claim (8.33).

The only remaining detail is the proof of Lemma 8.14. The proof of the tail bound (8.39a) is a simple exercise, using the sub-exponential tail bounds from Chapter 2. The proof of the bound (8.39b) requires more involved argument, one that makes use of both Exercise 7.10 and our previous results on estimation of sample covariances from Chapter 6.

For a constant $\xi > 0$ to be chosen, consider the positive integer $k := \lceil \xi \nu^2 \frac{n}{\log d} \rceil$, and the

collection of submatrices $\{(\mathbf{P}_3)_{SS}, |S| = k\}$. Given a parameter $\alpha \in (0, 1)$ to be chosen, a combination of the union bound and Theorem 6.5 imply that there are universal constants c_1 and c_2 such that

$$\mathbb{P}\left[\max_{|S|=k} \|(\mathbf{P}_3)_{SS}\|_2 \geq c_1 \sqrt{\frac{k}{n}} + \alpha v\right] \leq 2e^{-c_2 n \alpha^2 v^2 + \log \binom{d}{k}}.$$

Since $\log \binom{d}{k} \leq 2k \log(d) \leq 4\xi v^2 n$, this probability is at most $e^{-c_2 n v^2 (\alpha^2 - 4\xi)} = e^{-c_2 n v^2 \alpha^2 / 2}$, as long as we set $\xi = \alpha^2 / 8$. The result of Exercise 7.10 then implies that

$$|\langle \Delta, \mathbf{P}_3 \Delta \rangle| \leq 27c'_1 \alpha v \left\{ \|\Delta\|_2^2 + \frac{8}{\alpha^2 v^2} \frac{\log d}{n} \|\Delta\|_1^2 \right\} \quad \text{for all } \Delta \in \mathbb{R}^d,$$

with the previously stated probability. Setting $\alpha = \frac{1}{(16 \times 27)c'_1}$ yields the claim (8.39b) with $c'_3 = (2\alpha^2)^{-1}$. \square

8.4 Bibliographic details and background

Further details on PCA and its applications can be found in books by Anderson (1984) (cf. chapter 11), Jolliffe (2004) and Muirhead (2008). See the two-volume set by Horn and Johnson (1985; 1991) for background on matrix analysis, as well as the book by Bhatia (1997) for a general operator-theoretic viewpoint. The book by Stewart and Sun (1980) is more specifically focused on matrix perturbation theory, whereas Stewart (1971) provides perturbation theory in the more general setting of closed linear operators.

Johnstone (2001) introduced the spiked covariance model (8.18), and investigated the high-dimensional asymptotics of its eigenstructure; see also the papers by Baik and Silverstein (2006) and Paul (2007) for high-dimensional asymptotics. Johnstone and Lu (2009) introduced the sparse variant of the spiked ensemble, and proved consistency results for a simple estimator based on thresholding the diagonal entries of the sample covariance matrix. Amini and Wainwright (2009) provided a more refined analysis of this same estimator, as well as of a semidefinite programming (SDP) relaxation proposed by d'Asprémont et al. (2007). See Exercise 8.8 for the derivation of this latter SDP relaxation. The non-convex estimator (8.25a) was first proposed by Jolliffe et al. (2003), and called the SCOTLASS criterion; Witten et al. (2009) derive an alternating algorithm for finding a local optimum of this criterion. Other authors, including Ma (2010; 2013) and Yuan and Zhang (2013), have studied iterative algorithms for sparse PCA based on combining the power method with soft or hard thresholding.

Minimax lower bounds for estimating principal components in various types of spiked ensembles can be derived using techniques discussed in Chapter 15. These lower bounds show that the upper bounds obtained in Corollaries 8.7 and 8.12 for ordinary and sparse PCA, respectively, are essentially optimal. See Birnbaum et al. (2012) and Vu and Lei (2012) for lower bounds on the ℓ_2 -norm error in sparse PCA. Amini and Wainwright (2009) derived lower bounds for the problem of variable selection in sparse PCA. Some of these lower bounds are covered in this book: in particular, see Example 15.19 for minimax lower bounds on ℓ_2 -error in ordinary PCA, Example 15.20 for lower bounds on variable selection in sparse PCA, and Exercise 15.16 for ℓ_2 -error lower bounds on sparse PCA. Berthet

and Rigollet (2013) derived certain hardness results for the problem of sparse PCA detection, based on relating it to the (conjectured) average-case hardness of the planted k -clique problem in Erdős–Rényi random graphs. Ma and Wu (2013) developed a related but distinct reduction, one which applies to a Gaussian detection problem over a family of sparse-plus-low-rank matrices. See also the papers (Wang et al., 2014; Cai et al., 2015; Gao et al., 2015) for related results using the conjectured hardness of the k -clique problem.

8.5 Exercises

Exercise 8.1 (Courant–Fischer variational representation) For a given integer $j \in \{2, \dots, d\}$, let \mathcal{V}_{j-1} denote the collection of all subspaces of dimension $j-1$. For any symmetric matrix \mathbf{Q} , show that the j th largest eigenvalue is given by

$$\gamma_j(\mathbf{Q}) = \min_{\mathbb{V} \in \mathcal{V}_{j-1}} \max_{u \in \mathbb{V}^\perp \cap S^{d-1}} \langle u, \mathbf{Q}u \rangle, \quad (8.41)$$

where \mathbb{V}^\perp denotes the orthogonal subspace to \mathbb{V} .

Exercise 8.2 (Unitarily invariant matrix norms) For positive integers $d_1 \leq d_2$, a matrix norm on $\mathbb{R}^{d_1 \times d_2}$ is *unitarily invariant* if $\|\mathbf{M}\| = \|\mathbf{V}\mathbf{M}\mathbf{U}\|$ for all orthonormal matrices $\mathbf{V} \in \mathbb{R}^{d_1 \times d_1}$ and $\mathbf{U} \in \mathbb{R}^{d_2 \times d_2}$.

(a) Which of the following matrix norms are unitarily invariant?

- (i) The Frobenium norm $\|\mathbf{M}\|_F$.
- (ii) The nuclear norm $\|\mathbf{M}\|_{\text{nuc}}$.
- (iii) The ℓ_2 -operator norm $\|\mathbf{M}\|_2 = \sup_{\|u\|_2=1} \|\mathbf{M}u\|_2$.
- (iv) The ℓ_∞ -operator norm $\|\mathbf{M}\|_\infty = \sup_{\|u\|_\infty=1} \|\mathbf{M}u\|_\infty$.

(b) Let ρ be a norm on \mathbb{R}^{d_1} that is invariant to permutations and sign changes—that is

$$\rho(x_1, \dots, x_{d_1}) = \rho(z_1 x_{\pi(1)}, \dots, z_{d_1} x_{\pi(d_1)})$$

for all binary strings $z \in \{-1, 1\}^{d_1}$ and permutations π on $\{1, \dots, d_1\}$. Such a function is known as a *symmetric gauge function*. Letting $\{\sigma_j(\mathbf{M})\}_{j=1}^{d_1}$ denote the singular values of \mathbf{M} , show that

$$\|\mathbf{M}\|_\rho := \rho(\underbrace{\sigma_1(\mathbf{M}), \dots, \sigma_{d_1}(\mathbf{M})}_{\sigma(\mathbf{M}) \in \mathbb{R}^{d_1}})$$

defines a matrix norm. (*Hint:* For any pair of $d_1 \times d_2$ matrices \mathbf{M} and \mathbf{N} , we have $\text{trace}(\mathbf{N}^T \mathbf{M}) \leq \langle \sigma(\mathbf{N}), \sigma(\mathbf{M}) \rangle$, where $\sigma(\mathbf{M})$ denotes the ordered vector of singular values.)

(c) Show that all matrix norms in the family from part (b) are unitarily invariant.

Exercise 8.3 (Weyl’s inequality) Prove Weyl’s inequality (8.9). (*Hint:* Exercise 8.1 may be useful.)

Exercise 8.4 (Variational characterization of eigenvectors) Show that the orthogonal matrix $\mathbf{V} \in \mathbb{R}^{d \times r}$ maximizing the criterion (8.2) has columns formed by the top r eigenvectors of $\Sigma = \text{cov}(X)$.

Exercise 8.5 (Matrix power method) Let $\mathbf{Q} \in \mathcal{S}^{d \times d}$ be a strictly positive definite symmetric matrix with a unique maximal eigenvector θ^* . Given some non-zero initial vector $\theta^0 \in \mathbb{R}^d$, consider the sequence $\{\theta^t\}_{t=0}^\infty$,

$$\theta^{t+1} = \frac{\mathbf{Q}\theta^t}{\|\mathbf{Q}\theta^t\|_2}. \quad (8.42)$$

- Prove that there is a large set of initial vectors θ^0 for which the sequence $\{\theta^t\}_{t=0}^\infty$ converges to θ^* .
- Give a “bad” initialization for which this convergence does not take place.
- Based on part (b), specify a procedure to compute the second largest eigenvector, assuming it is also unique.

Exercise 8.6 (PCA for Gaussian mixture models) Consider an instance of the Gaussian mixture model from Example 8.3 with equal mixture weights ($\alpha = 0.5$) and unit-norm mean vector ($\|\theta^*\|_2 = 1$), and suppose that we implement the PCA-based estimator $\widehat{\theta}$ for the mean vector θ^* .

- Prove that if the sample size is lower bounded as $n > c_1 \sigma^2 (1 + \sigma^2)d$ for a sufficiently large constant c_1 , this estimator satisfies a bound of the form

$$\|\widehat{\theta} - \theta^*\|_2 \leq c_2 \sigma \sqrt{1 + \sigma^2} \sqrt{\frac{d}{n}}$$

with high probability.

- Explain how to use your estimator to build a classification rule—that is, a mapping $x \mapsto \psi(x) \in \{-1, +1\}$, where the binary labels code whether sample x has mean $-\theta^*$ or $+\theta^*$.
- Does your method still work if the shared covariance matrix is *not* a multiple of the identity?

Exercise 8.7 (PCA for retrieval from absolute values) Suppose that our goal is to estimate an unknown vector $\theta^* \in \mathbb{R}^d$ based on n i.i.d. samples $\{(x_i, y_i)\}_{i=1}^n$ of the form $y_i = |\langle x_i, \theta^* \rangle|$, where $x_i \sim \mathcal{N}(0, \mathbf{I}_d)$. This model is a real-valued idealization of the problem of phase retrieval, to be discussed at more length in Chapter 10. Suggest a PCA-based method for estimating θ^* that is consistent in the limit of infinite data. (*Hint*: Using the pair (x, y) , try to construct a random matrix \mathbf{Z} such that $\mathbb{E}[\mathbf{Z}] = \sqrt{\frac{2}{\pi}}(\theta^* \otimes \theta^* + \mathbf{I}_d)$.)

Exercise 8.8 (Semidefinite relaxation of sparse PCA) Recall the non-convex problem (8.25a), also known as the SCOTLASS estimator. In this exercise, we derive a convex relaxation of the objective, due to d’Aspremont et al. (2007).

- Show that the non-convex problem (8.25a) is equivalent to the optimization problem

$$\max_{\Theta \in \mathcal{S}_+^{d \times d}} \text{trace}(\widehat{\Sigma}\Theta) \quad \text{such that } \text{trace}(\Theta) = 1, \sum_{j,k=1}^d |\Theta_{jk}| \leq R^2 \text{ and } \text{rank}(\Theta) = 1,$$

where $\mathcal{S}_+^{d \times d}$ denotes the cone of symmetric, positive semidefinite matrices.

(b) Dropping the rank constraint yields the convex program

$$\max_{\Theta \in S_+^{d \times d}} \text{trace}(\widehat{\Sigma}\Theta) \quad \text{such that } \text{trace}(\Theta) = 1 \text{ and } \sum_{j,k=1}^d |\Theta_{jk}| \leq R^2.$$

What happens when its optimum is achieved at a rank-one matrix?

Exercise 8.9 (Primal–dual witness for sparse PCA) The SDP relaxation from Exercise 8.8(b) can be written in the equivalent Lagrangian form

$$\max_{\substack{\Theta \in S_+^{d \times d} \\ \text{trace}(\Theta)=1}} \left\{ \text{trace}(\widehat{\Sigma}\Theta) - \lambda_n \sum_{j,k=1}^d |\Theta_{jk}| \right\}. \quad (8.43)$$

Suppose that there exists a vector $\widehat{\theta} \in \mathbb{R}^d$ and a matrix $\widehat{U} \in \mathbb{R}^{d \times d}$ such that

$$\widehat{U}_{jk} = \begin{cases} \text{sign}(\widehat{\theta}_j \widehat{\theta}_k) & \text{if } \widehat{\theta}_j \widehat{\theta}_k \neq 0, \\ \in [-1, 1] & \text{otherwise,} \end{cases}$$

and moreover such that $\widehat{\theta}$ is a maximal eigenvector of the matrix $\widehat{\Sigma} - \lambda_n \widehat{U}$. Prove that the rank-one matrix $\widehat{\Theta} = \widehat{\theta} \otimes \widehat{\theta}$ is an optimal solution to the SDP relaxation (8.43).