

Random matrices and covariance estimation

2122

2123

Covariance matrices play a central role in statistics, and there exist a variety of methods for estimating them based on data. The problem of covariance estimation dovetails with random matrix theory, since the sample covariance is a particular type of random matrix. A classical framework allows the sample size n to tend to infinity while the matrix dimension d remains fixed; in such a setting, the behavior of the sample covariance matrix is characterized by the usual limit theory. In contrast, for high-dimensional random matrices in which the data dimension is either comparable to the sample size ($d \asymp n$), or possibly much larger than the sample size ($d \gg n$), many new phenomena arise.

2124

2125

2126

2127

2128

2129

2130

2131

2132

High-dimensional random matrices play an important role in many branches of science, mathematics, and engineering, and have been studied extensively. The classical theory is asymptotic in nature, such as the Wigner semi-circle law and the Marcenko-Pastur law for the asymptotic distribution of the eigenvalues of a sample covariance matrix (see Chapter 1 for illustration of the latter). By contrast, this chapter is devoted to an exploration of random matrices in a non-asymptotic setting, with the goal of obtaining explicit deviation inequalities that hold for all sample sizes and matrix dimensions. Beginning with the simplest case—namely ensembles of Gaussian random matrices—we then discuss more general sub-Gaussian ensembles, and then move onwards to ensembles with milder tail conditions. Throughout our development, we bring to bear the techniques from concentration of measure, comparison inequalities, and metric entropy developed previously in Chapters 2 through 5. In addition, this chapter introduces new some techniques, among them a class of matrix tail bounds developed over the past decade (see Section 6.4).

2133

2134

2135

2136

2137

2138

2139

2140

2141

2142

2143

2144

2145

2146

■ 6.1 Some preliminaries

2147

We begin by introducing notation and preliminary results used throughout this chapter, before setting up the problem of covariance estimation more precisely.

2148

2149

■ 6.1.1 Notation and basic facts

2150

The set of symmetric matrices in \mathbb{R}^d is denoted $\mathcal{S}^{d \times d} := \{\mathbf{Q} \in \mathbb{R}^{d \times d} \mid \mathbf{Q} = \mathbf{Q}^T\}$, where the subset of positive semidefinite matrices is given by

$$\mathcal{S}_+^{d \times d} := \{\mathbf{Q} \in \mathcal{S}^{d \times d} \mid \mathbf{Q} \succeq 0\}. \quad (6.1)$$

From standard linear algebra, we recall the facts that any matrix $\mathbf{Q} \in \mathcal{S}^{d \times d}$ is diagonalizable via a unitary transformation, and we use $\gamma(\mathbf{Q}) \in \mathbb{R}^d$ to denote its vector of eigenvalues, ordered as

$$\gamma_{\max}(\mathbf{Q}) = \gamma_1(\mathbf{Q}) \geq \gamma_2(\mathbf{Q}) \geq \cdots \geq \gamma_d(\mathbf{Q}) = \gamma_{\min}(\mathbf{Q}).$$

Note that a matrix is positive semidefinite ($\mathbf{Q} \succeq 0$) if and only if $\gamma_{\min}(\mathbf{Q}) \geq 0$.

2151

Our analysis frequently exploits the Rayleigh-Ritz variational characterization of the minimum and maximum eigenvalues

$$\gamma_{\max}(\mathbf{Q}) = \max_{v \in \mathbb{S}^{d-1}} v^T \mathbf{Q} v, \quad \text{and} \quad \gamma_{\min}(\mathbf{Q}) = \min_{v \in \mathbb{S}^{d-1}} v^T \mathbf{Q} v, \quad (6.2)$$

where $\mathbb{S}^{d-1} := \{v \in \mathbb{R}^d \mid \|v\|_2 = 1\}$ is the Euclidean unit-sphere in \mathbb{R}^d . For any symmetric matrix \mathbf{Q} , the ℓ_2 -operator norm can be written as $\|\mathbf{Q}\|_{\text{op}} = \max\{\gamma_{\max}(\mathbf{Q}), |\gamma_{\min}(\mathbf{Q})|\}$, by virtue of which it inherits the variational representation

$$\|\mathbf{Q}\|_{\text{op}} := \max_{v \in \mathbb{S}^{d-1}} |v^T \mathbf{Q} v|. \quad (6.3)$$

■ 6.1.2 Set-up of covariance estimation

2152

Let us now define the problem of covariance matrix estimation. Let $\{x_1, \dots, x_n\}$ be a collection of n independent and identically distributed samples¹ from a distribution in \mathbb{R}^d with zero mean, and covariance matrix $\Sigma = \text{cov}(x_1) \in \mathcal{S}_+^{d \times d}$. A standard estimator of Σ is the *sample covariance matrix*

$$\hat{\Sigma} := \frac{1}{n} \sum_{i=1}^n x_i x_i^T. \quad (6.4)$$

Since each x_i has zero mean, we are guaranteed that $\mathbb{E}[x_i x_i^T] = \Sigma$, and hence that the random matrix $\hat{\Sigma}$ is an unbiased estimator of the population covariance Σ . Consequently, the error matrix $\hat{\Sigma} - \Sigma$ has mean zero, and our goal in this chapter is to obtain bounds on the error measured in the ℓ_2 -operator norm. By the variational rep-

¹In this chapter, we use a lower case x to denote a random vector, so as it distinguish it from a random matrix.

representation (6.3), a bound of the form $\|\widehat{\Sigma} - \Sigma\|_{\text{op}} \leq \epsilon$ is equivalent to asserting that

$$\max_{v \in \mathbb{S}^{d-1}} \left| \frac{1}{n} \sum_{i=1}^n \langle x_i, v \rangle^2 - v^T \Sigma v \right| \leq \epsilon. \quad (6.5)$$

This representation shows that controlling the deviation $\|\widehat{\Sigma} - \Sigma\|_{\text{op}}$ is equivalent to establishing a uniform law of large numbers for the class of functions $x \mapsto \langle x, v \rangle^2$, indexed by vectors $v \in \mathbb{S}^{d-1}$. See Chapter 4 for further discussion of such uniform laws in a general setting.

Control in the operator norm also guarantees that the eigenvalues of $\widehat{\Sigma}$ are uniformly close to those of Σ . In particular, by a corollary of Weyl's theorem (see the bibliographic section for details), we have

$$\max_{j=1, \dots, d} |\gamma_j(\widehat{\Sigma}) - \gamma_j(\Sigma)| \leq \|\widehat{\Sigma} - \Sigma\|_{\text{op}}. \quad (6.6)$$

A similar type of guarantee can be made for the eigenvectors of the two matrices, but only if one has additional control on the separation between adjacent eigenvalues. See our discussion of principal component analysis in Chapter 8 for more details.

Finally, let us point out the connection to the singular values of the random matrix $\mathbf{X} \in \mathbb{R}^{n \times d}$, which contains the sample x_i^T as its i^{th} row. Since

$$\widehat{\Sigma} = \frac{1}{n} \sum_{i=1}^n x_i x_i^T = \frac{1}{n} \mathbf{X}^T \mathbf{X},$$

it follows that the eigenvalues of $\widehat{\Sigma}$ are the squares of the singular values of \mathbf{X}/\sqrt{n} .

■ 6.2 Wishart matrices and their behavior

We begin by studying the behavior of singular values for random matrices with Gaussian rows. More precisely, let us suppose that each sample x_i is drawn i.i.d. from a multivariate $\mathcal{N}(0, \Sigma)$ distribution, in which case we say that the associated matrix $\mathbf{X} \in \mathbb{R}^{n \times d}$, with x_i^T as its i^{th} row, is drawn from the Σ -Gaussian ensemble. In this case, the sample covariance $\widehat{\Sigma} = \frac{1}{n} \mathbf{X}^T \mathbf{X}$ is said to follow a multivariate Wishart distribution.

Theorem 6.1. Let $\mathbf{X} \in \mathbb{R}^{n \times d}$ be drawn according to the Σ -Gaussian ensemble. Then for all $\delta > 0$, the maximum singular value satisfies the upper deviation inequality

$$\mathbb{P} \left[\frac{\gamma_{\max}(\mathbf{X})}{\sqrt{n}} \geq \gamma_{\max}(\sqrt{\Sigma}) (1 + \delta) + \sqrt{\frac{\text{trace}(\Sigma)}{n}} \right] \leq e^{-n\delta^2/2}. \quad (6.7)$$

Moreover, for $n \geq d$, the minimum singular value satisfies the analogous lower deviation inequality

$$\mathbb{P} \left[\frac{\gamma_{\min}(\mathbf{X})}{\sqrt{n}} \leq \gamma_{\min}(\sqrt{\Sigma}) (1 - \delta) - \sqrt{\frac{\text{trace}(\Sigma)}{n}} \right] \leq e^{-n\delta^2/2}. \quad (6.8)$$


Before proving this result, let us consider some illustrative examples.

Example 6.1 (Operator norm bounds for the standard Gaussian ensemble). Consider a random matrix $\mathbf{W} \in \mathbb{R}^{n \times d}$ generated with i.i.d. $\mathcal{N}(0, 1)$ entries. This choice yields an instance of Σ -Gaussian ensemble, in particular with $\Sigma = \mathbf{I}_d$. By specializing Theorem 6.1, we conclude that for $n \geq d$, we have

$$\frac{\gamma_{\max}(\mathbf{W})}{\sqrt{n}} \leq 1 + \delta + \sqrt{\frac{d}{n}}, \quad \text{and} \quad \frac{\gamma_{\min}(\mathbf{W})}{\sqrt{n}} \geq 1 - \delta - \sqrt{\frac{d}{n}}, \quad (6.9)$$

where both bounds hold with probability greater than $1 - 2e^{-n\delta^2/2}$. These bounds on the singular values of \mathbf{W} imply that

$$\left\| \frac{1}{n} \mathbf{W}^T \mathbf{W} - \mathbf{I}_d \right\|_{\text{op}} \leq 2\epsilon + \epsilon^2, \quad \text{where } \epsilon = \sqrt{\frac{d}{n}} + \delta, \quad (6.10)$$

with the same probability. Consequently, the sample covariance $\hat{\Sigma} = \frac{1}{n} \mathbf{W}^T \mathbf{W}$ is a consistent estimate of the identity matrix \mathbf{I}_d whenever $d/n \rightarrow 0$. 

The preceding example has interesting consequences for the problem of sparse linear regression using standard Gaussian random matrices, as in compressed sensing; in particular, see our discussion of the restricted isometry property in Chapter 7. On the other hand, from the perspective of covariance estimation, estimating the identity matrix is not especially interesting, but a minor modification leads to a more realistic family of problems.

Example 6.2 (Gaussian covariance estimation). Let $\mathbf{X} \in \mathbb{R}^{n \times d}$ be a random matrix from the Σ -Gaussian ensemble. By standard properties of the multivariate Gaussian, we can write $\mathbf{X} = \mathbf{W}\sqrt{\Sigma}$, where $\mathbf{W} \in \mathbb{R}^{n \times d}$ is a standard Gaussian random matrix, and hence

$$\left\| \frac{1}{n} \mathbf{X}^T \mathbf{X} - \Sigma \right\|_{\text{op}} = \left\| \sqrt{\Sigma} \left(\frac{1}{n} \mathbf{W}^T \mathbf{W} - \mathbf{I}_d \right) \sqrt{\Sigma} \right\|_{\text{op}} \leq \left\| \Sigma \right\|_{\text{op}} \left\| \frac{1}{n} \mathbf{W}^T \mathbf{W} - \mathbf{I}_d \right\|_{\text{op}},$$

Consequently, by exploiting the bound (6.10), we are guaranteed that, for all $\delta > 0$

$$\frac{\|\hat{\Sigma} - \Sigma\|_{\text{op}}}{\|\Sigma\|_{\text{op}}} \leq 2\sqrt{\frac{d}{n}} + 2\delta + \left(\sqrt{\frac{d}{n}} + \delta\right)^2, \quad (6.11)$$

with probability at least $1 - 2e^{-n\delta^2/2}$. Consequently, we conclude that the relative error $\|\hat{\Sigma} - \Sigma\|_{\text{op}}/\|\Sigma\|_{\text{op}}$ converges to zero as long the ratio d/n converges to zero. ♣

It is interesting to consider Theorem 6.1 in application to sequences of matrices that satisfy additional structure, one being control on the eigenvalues of the covariance matrix Σ .

Example 6.3 (Faster rates under trace constraints). Recall that $\{\gamma_j(\Sigma)\}_{j=1}^d$ denotes the ordered sequence of eigenvalues of the matrix Σ , with $\gamma_1(\Sigma)$ being the maximum eigenvalue. Now consider a population covariance matrix $\Sigma \succ 0$ that satisfies a “trace constraint” of the form

$$\frac{\text{trace}(\Sigma)}{\|\Sigma\|_{\text{op}}} = \frac{\sum_{j=1}^d \gamma_j(\Sigma)}{\gamma_1(\Sigma)} \leq R, \quad (6.12)$$

where R is some constant independent of dimension. Note that this ratio is a rough measure of the matrix rank, since inequality (6.12) always holds with $R = \text{rank}(\Sigma)$. Perhaps more interesting are matrices that are full-rank but that exhibit a relatively fast eigendecay, with a canonical instance being the Schatten- q -“balls” of matrices. For symmetric matrices, these sets take the form

$$\mathbb{B}_q(R_q) := \left\{ \Sigma \in \mathbb{R}^{d \times d} \mid \sum_{j=1}^d |\gamma_j(\Sigma)|^q \leq R_q \right\}, \quad (6.13)$$

where $q \in [0, 1]$ is a given parameter, and $R_q > 0$ is the radius. Note that these matrix families are nested: the smallest set with $q = 0$ corresponds to the case of matrices with rank at most R_0 , whereas the other extreme $q = 1$ corresponds to an explicit trace constraint. Any one of these families satisfies a bound of the form (6.12) with the parameter R proportional to R_q .

For any matrix class satisfying the bound (6.12), Theorem 6.1 guarantees that, with high probability, the maximum singular value is bounded above as

$$\frac{\gamma_{\max}(\mathbf{X})}{\sqrt{n}} \leq \gamma_{\max}(\sqrt{\Sigma})(1 + \delta + \sqrt{\frac{R}{n}}). \quad (6.14)$$

By comparison to the earlier bound (6.9) for $\Sigma = \mathbf{I}_d$, we conclude that the parameter R plays the role of the *effective dimension*.

We now turn to the proof of Theorem 6.1.

Proof. In order to simplify notation in the proof, let us introduce the convenient short-hand $\bar{\gamma}_{\max} = \gamma_{\max}(\sqrt{\Sigma})$ and $\bar{\gamma}_{\min} = \gamma_{\min}(\sqrt{\Sigma})$. Our proofs of both the upper and lower bounds consist of two steps: first, we use concentration inequalities (see Chapter 2) to argue that the random singular value is close to its expectation with high probability, and second, we use Gaussian comparison inequalities (see Chapter 5) to bound the expected values.

Maximum singular value: As noted previously, by standard properties of the multivariate Gaussian distribution, we can write $\mathbf{X} = \mathbf{W}\sqrt{\Sigma}$, where $\mathbf{W} \in \mathbb{R}^{n \times d}$ has i.i.d. $\mathcal{N}(0, 1)$ entries. Now let us view the mapping $\mathbf{W} \mapsto \gamma_{\max}(\mathbf{W}\sqrt{\Sigma})/\sqrt{n}$ as a real-valued function on \mathbb{R}^{nd} . By the argument given in Example 2.16, this function is Lipschitz with respect to the Euclidean norm with constant at most $L = \bar{\gamma}_{\max}/\sqrt{n}$. By concentration of measure for Lipschitz functions of Gaussian random vectors (Theorem 2.4), we conclude that

$$\mathbb{P}[\gamma_{\max}(\mathbf{X}) \geq \mathbb{E}[\gamma_{\max}(\mathbf{X})] + \sqrt{n}\bar{\gamma}_{\max}\delta] \leq e^{-n\delta^2/2}.$$

Consequently, it suffices to show that

$$\mathbb{E}[\gamma_{\max}(\mathbf{X})] \leq \sqrt{n}\bar{\gamma}_{\max} + \sqrt{\text{trace}(\Sigma)}. \quad (6.15)$$

In order to do so, we first write $\gamma_{\max}(\mathbf{X})$ in a variational fashion, as the maximum of a suitably defined Gaussian process. By definition of the maximum singular value, we have $\gamma_{\max}(\mathbf{X}) = \max_{v' \in \mathbb{S}^{d-1}} \|\mathbf{X}v'\|_2$, where \mathbb{S}^{d-1} denotes the Euclidean unit-sphere in \mathbb{R}^d .

Recalling the representation $\mathbf{X} = \mathbf{W}\sqrt{\Sigma}$ and defining $v = \sqrt{\Sigma}v'$, we can write

$$\gamma_{\max}(\mathbf{X}) = \max_{v \in \mathbb{S}^{d-1}(\Sigma^{-1})} \|Wv\|_2 = \max_{u \in \mathbb{S}^{n-1}} \max_{v \in \mathbb{S}^{d-1}(\Sigma^{-1})} \underbrace{u^T W v}_{Z_{u,v}},$$

where $\mathbb{S}^{d-1}(\Sigma^{-1}) := \{v \in \mathbb{R}^d \mid \|\Sigma^{-\frac{1}{2}}v\|_2 = 1\}$ is an ellipse. Consequently, obtaining bounds on the maximum singular value corresponds to controlling the supremum of the zero-mean Gaussian process $\{Z_{u,v}, (u,v) \in \mathbb{T}\}$ indexed by the set $\mathbb{T} := \mathbb{S}^{n-1} \times \mathbb{S}^{d-1}(\Sigma^{-1})$

We upper bound the expected value of this supremum by constructing another Gaussian process $\{Y_{u,v}, (u,v) \in \mathbb{T}\}$ such that $\mathbb{E}[(Z_{u,v} - Z_{\tilde{u},\tilde{v}})^2] \leq \mathbb{E}[(Y_{u,v} - Y_{\tilde{u},\tilde{v}})^2]$ for all pairs (u,v) and (\tilde{u},\tilde{v}) in \mathbb{T} . We can then apply the Sudakov-Fernique comparison (Theorem 5.3) to conclude that

$$\mathbb{E}[\gamma_{\max}(\mathbf{X})] = \mathbb{E}\left[\max_{(u,v) \in \mathbb{T}} Z_{u,v}\right] \leq \mathbb{E}\left[\max_{(u,v) \in \mathbb{T}} Y_{u,v}\right]. \quad (6.16)$$

We begin by computing the metric ρ_Z induced by the Gaussian process $Z_{u,v} = u^T \mathbf{W}v$. Given two pairs (u,v) and (\tilde{u},\tilde{v}) , we may assume without loss of generality that $\|v\|_2 \leq$

$\|\tilde{v}\|_2$. (If not, we simply reverse the roles of (u, v) and (\tilde{u}, \tilde{v}) in the argument to follow.) We begin by observing that $Z_{u,v} = \langle \mathbf{W}, uv^T \rangle$, where $\langle A, B \rangle = \sum_{j,k} A_{j,k} B_{j,k}$ is the trace inner product. Since \mathbf{W} has i.i.d. $\mathcal{N}(0, 1)$ entries, we have

$$\mathbb{E}[(Z_{u,v} - Z_{\tilde{u},\tilde{v}})^2] = \mathbb{E}[(\langle \mathbf{W}, uv^T - \tilde{u}\tilde{v}^T \rangle)^2] = \|uv^T - \tilde{u}\tilde{v}^T\|_F^2.$$

Re-arranging and expanding out this Frobenius norm, we find that

$$\begin{aligned} \|uv^T - \tilde{u}\tilde{v}^T\|_F^2 &= \|u(v - \tilde{v})^T + (u - \tilde{u})\tilde{v}^T\|_F^2 \\ &= \|(u - \tilde{u})\tilde{v}^T\|_F^2 + \|u(v - \tilde{v})^T\|_F^2 + 2\langle u(v - \tilde{v})^T, (u - \tilde{u})\tilde{v}^T \rangle \\ &= \|\tilde{v}\|_2^2 \|u - \tilde{u}\|_2^2 + \|u\|_2^2 + \|v - \tilde{v}\|_2^2 + 2(\|u\|_2^2 - \langle u, \tilde{u} \rangle) (\langle v, \tilde{v} \rangle - \|\tilde{v}\|_2^2) \end{aligned}$$

Now since $\|u\|_2 = \|\tilde{u}\|_2 = 1$ by definition of the set \mathbb{T} , we have $\|u\|_2^2 - \langle u, \tilde{u} \rangle \geq 0$. On the other hand, we have

$$|\langle v, \tilde{v} \rangle| \stackrel{(i)}{\leq} \|v\|_2 \|\tilde{v}\|_2 \stackrel{(ii)}{\leq} \|\tilde{v}\|_2^2,$$

where step (i) follows from the Cauchy-Schwarz inequality, and step (ii) follows from our initial assumption that $\|v\|_2 \leq \|\tilde{v}\|_2$. Combined with our previous bound on $\|u\|_2^2 - \langle u, \tilde{u} \rangle$, we conclude that

$$\underbrace{(\|u\|_2^2 - \langle u, \tilde{u} \rangle)}_{\geq 0} \underbrace{(\langle v, \tilde{v} \rangle - \|\tilde{v}\|_2^2)}_{\leq 0} \leq 0.$$

Putting together the pieces, we conclude that $\|uv^T - \tilde{u}\tilde{v}^T\|_F^2 \leq \|\tilde{v}\|_2^2 \|u - \tilde{u}\|_2^2 + \|v - \tilde{v}\|_2^2$. Finally, by definition of the set $\mathbb{S}^{d-1}(\Sigma^{-1})$, we have $\|\tilde{v}\|_2 \leq \gamma_{\max}(\sqrt{\Sigma})$, and hence

$$\mathbb{E}[(Z_{u,v} - Z_{\tilde{u},\tilde{v}})^2] \leq \sigma^2 \|u - \tilde{u}\|_2^2 + \|v - \tilde{v}\|_2^2, \quad \text{where } \sigma := \gamma_{\max}(\sqrt{\Sigma}).$$

Motivated by this inequality, we define the Gaussian process $Y_{u,v} := \sigma \langle g, u \rangle + \langle h, v \rangle$, where $g \in \mathbb{R}^n$ and $h \in \mathbb{R}^d$ are both standard Gaussian random vectors (i.e., with i.i.d. $\mathcal{N}(0, 1)$ entries), and mutually independent. By construction, we have

$$\mathbb{E}[(Y_\theta - Y_{\hat{\theta}})^2] = \sigma^2 \|u - \tilde{u}\|_2^2 + \|v - \tilde{v}\|_2^2.$$

Thus, we may apply the Sudakov-Fernique bound (6.16) to conclude that

$$\begin{aligned}\mathbb{E}[\gamma_{\max}(\mathbf{X})] &\leq \mathbb{E}\left[\sup_{(u,v) \in \mathbb{T}} Y_{u,v}\right] \\ &= \sigma \mathbb{E}\left[\sup_{u \in \mathbb{S}^{n-1}} \langle g, u \rangle\right] + \mathbb{E}\left[\sup_{u \in \mathbb{S}^{d-1}(\Sigma^{-1})} \langle h, v \rangle\right] \\ &= \sigma \mathbb{E}[\|g\|_2] + \mathbb{E}[\|\sqrt{\Sigma}h\|_2]\end{aligned}$$

By Jensen's inequality, we have $\mathbb{E}[\|g\|_2] \leq \sqrt{n}$, and similarly,

$$\mathbb{E}[\|\sqrt{\Sigma}h\|_2] \leq \sqrt{\mathbb{E}[h^T \Sigma h]} = \sqrt{\text{trace}(\Sigma)},$$

which establishes the claim (6.15). 2201

The lower bound on the minimum singular value is based on a similar argument, 2202
but requires somewhat more technical work, so that we defer it to the Appendix. 2203

□ 2205

■ 6.3 Covariance matrices from sub-Gaussian ensembles 2206

Various aspects of our development thus far have crucially exploited different properties 2207
of the Gaussian distribution, especially our use of the Gaussian comparison inequalities. 2208
In this section, we show a somewhat different approach—namely, discretization and tail 2209
bounds—can be used to establish analogous bounds for general sub-Gaussian random 2210
matrices. 2211

In particular, let us assume that the random vector $x_i \in \mathbb{R}^d$ is zero-mean, and sub-
Gaussian with parameter at most σ , by which we mean that for each fixed $v \in \mathbb{S}^{d-1}$,

$$\mathbb{E}[e^{\lambda \langle v, x_i \rangle}] \leq e^{\frac{\lambda^2 \sigma^2}{2}} \quad \text{for all } \lambda \in \mathbb{R}. \quad (6.17)$$

Equivalently stated, we assume that the scalar random variable $\langle v, x_i \rangle$ is zero-mean 2212
and sub-Gaussian with parameter at most σ . (See Chapter 2 for an in-depth discussion 2213
of sub-Gaussian variables.) Let us consider some examples to illustrate: 2214

- (a) Suppose that the matrix $\mathbf{X} \in \mathbb{R}^{n \times d}$ has i.i.d. entries, where each entry x_{ij} is 2215
zero-mean and sub-Gaussian with parameter $\sigma = 1$. Examples include the stan- 2216
dard Gaussian ensemble ($x_{ij} \sim \mathcal{N}(0, 1)$); the Bernoulli ensemble ($x_{ij} \in \{-1, +1\}$ 2217
equiprobably), and more generally, any zero-mean distribution supported on the 2218
interval $[-1, +1]$. In all of these cases, for any vector $v \in \mathbb{S}^{d-1}$, the random vari- 2219
able $\langle v, x_i \rangle$ is sub-Gaussian with parameter at most σ^2 , using the i.i.d. assumption 2220
on the entries of $x_i \in \mathbb{R}^d$, and standard properties of sub-Gaussian variables. 2221

(b) Now suppose that $x_i \sim \mathcal{N}(0, \Sigma)$. For any $v \in \mathbb{S}^{d-1}$, we have $\langle v, x_i \rangle \sim \mathcal{N}(0, v^T \Sigma v)$. 2222
 Since $v^T \Sigma v \leq \|\Sigma\|_{\text{op}}$, we conclude that x_i is sub-Gaussian with parameter at most 2223
 $\sigma^2 = \|\Sigma\|_{\text{op}}$. 2224

When the random matrix $\mathbf{X} \in \mathbb{R}^{n \times d}$ is formed by drawing each row $x_i \in \mathbb{R}^d$ in an 2225
 i.i.d. manner from a σ -sub-Gaussian distribution, then we say that \mathbf{X} is a sample from 2226
 a σ -sub-Gaussian ensemble. For any random matrix, we have the following result: 2227

Theorem 6.2. Suppose that x_1, \dots, x_n are i.i.d. samples from a zero-mean sub- 2228
 Gaussian distribution with parameter at most σ . Then the sample covariance
 $\hat{\Sigma} = \frac{1}{n} \sum_{i=1}^n x_i x_i^T$ satisfies the bounds

$$\mathbb{E}[e^{\lambda \|\hat{\Sigma} - \Sigma\|_{\text{op}}}] \leq e^{\frac{2\lambda^2 \sigma^2}{n} + 4d} \quad \text{for all } \lambda \in [0, \frac{n}{8\sigma^2}]. \quad (6.18) \quad 2229$$

Moreover, there are universal positive constants c_1 and c_2 such that for all $\delta > 0$

$$\mathbb{P}\left[\|\hat{\Sigma} - \Sigma\|_{\text{op}}/\sigma^2 \geq c_1 \left\{ \sqrt{\frac{d}{n}} + \frac{d}{n} \right\} + \delta\right] \leq 2e^{-c_2 n \min\{\delta, \delta^2\}}. \quad (6.19) \quad 2230$$

Remarks: When $\Sigma = \mathbf{I}_d$ and each x_i is sub-Gaussian with parameter $\sigma = 1$, Theo-
 rem 6.2 implies that

$$\|\hat{\Sigma} - \mathbf{I}_d\|_{\text{op}} \lesssim \sqrt{\frac{d}{n}} + \frac{d}{n}$$

with high probability. For $n \geq d$, this bound implies that the singular values of \mathbf{X}/\sqrt{n}
 satisfy the sandwich relation

$$1 - c' \sqrt{\frac{d}{n}} \leq \frac{\gamma_{\min}(\mathbf{X})}{\sqrt{n}} \leq \frac{\gamma_{\max}(\mathbf{X})}{\sqrt{n}} \leq 1 + c' \sqrt{\frac{d}{n}}, \quad (6.20)$$

for some universal constant $c' > 1$. It is worth comparing this result to the ear- 2231
 lier bounds (6.9), applicable to the special case of a standard Gaussian matrix. The 2232
 bound (6.20) has a qualitatively similar form, except that the constant c' is larger than 2233
 one. 2234

Proof. For notational convenience, we introduce the shorthand $\mathbf{Q} := \hat{\Sigma} - \Sigma$. Recall
 from Section 6.1 the variational representation $\|\mathbf{Q}\|_{\text{op}} = \max_{v \in \mathbb{S}^{d-1}} |\langle v, \mathbf{Q}v \rangle|$. We first reduce
 the supremum to a finite maximum via a discretization argument (see Chapter 5).

Let $\{v^1, \dots, v^N\}$ be a $\frac{1}{8}$ -covering of the sphere \mathbb{S}^{d-1} in the Euclidean norm; from
 Example 5.4, there exists such a covering with $N \leq 17^d$ vectors. Given any $v \in \mathbb{S}^{d-1}$,

we can write $v = v^j + \Delta$ for some v^j in the cover, and an error $\|\Delta\|_2 \leq \frac{1}{8}$, and hence

$$\langle v, \mathbf{Q}v \rangle = \langle v^j, \mathbf{Q}v^j \rangle + 2\langle \Delta, \mathbf{Q}v^j \rangle + \langle \Delta, \mathbf{Q}\Delta \rangle.$$

Applying the triangle inequality and the definition of operator norm yields

$$\begin{aligned} |\langle v, \mathbf{Q}v \rangle| &\leq |\langle v^j, \mathbf{Q}v^j \rangle| + 2\|\Delta\|_2 \|\mathbf{Q}\|_{\text{op}} \|v^j\|_2 + \|\mathbf{Q}\|_{\text{op}} \|\Delta\|_2^2 \\ &\leq |\langle v^j, \mathbf{Q}v^j \rangle| + \frac{1}{4} \|\mathbf{Q}\|_{\text{op}} + \frac{1}{64} \|\mathbf{Q}\|_{\text{op}} \\ &\leq |\langle v^j, \mathbf{Q}v^j \rangle| + \frac{1}{2} \|\mathbf{Q}\|_{\text{op}}. \end{aligned}$$

Re-arranging and then taking the supremum over $v \in \mathbb{S}^{d-1}$, and the associated maximum over $j \in \{1, 2, \dots, N\}$, we obtain

$$\|\mathbf{Q}\|_{\text{op}} = \max_{v \in \mathbb{S}^{d-1}} |\langle v, \mathbf{Q}v \rangle| \leq 2 \max_{j=1, \dots, N} |\langle v^j, \mathbf{Q}v^j \rangle|.$$

Our next step is to control the moment generating function of the random variable $|\langle u, \mathbf{Q}u \rangle|$, where $u \in \mathbb{S}^{d-1}$ is any fixed vector. For any $\lambda > 0$, by the definition of \mathbf{Q} and independence, we have

$$\mathbb{E}[e^{2\lambda \langle u, \mathbf{Q}u \rangle}] = \prod_{i=1}^n \mathbb{E}[e^{\frac{2\lambda}{n} \{ \langle x_i, u \rangle^2 - \langle u, \Sigma u \rangle \}}]$$

Since $z_i = \langle x_i, u \rangle$ is zero-mean and sub-Gaussian, Theorem 2.1 implies that

$$\mathbb{E}[e^{\frac{t}{2\gamma_i^2} (z_i^2 - \gamma_i^2)}] \leq \frac{e^{-t/2}}{\sqrt{1-t}} \leq e^{-t^2/2} \quad \text{for all } t \in [-\frac{1}{2}, \frac{1}{2}],$$

where $\gamma_i^2 = \mathbb{E}[z_i^2] \leq \sigma^2$. Setting $\lambda = \frac{nt}{4\gamma_i^2}$, we find that

$$\mathbb{E}[e^{2\lambda \langle u, \mathbf{Q}u \rangle}] \leq e^{8\frac{\lambda^2}{n^2} \sum_{i=1}^n \gamma_i^2} \leq e^{\frac{8\lambda^2 \sigma^2}{n}}, \quad \text{valid for all } \lambda \in [-\frac{n}{8\sigma^2}, \frac{n}{8\sigma^2}].$$

Lastly, dealing with the absolute value, we have

$$\mathbb{E}[e^{2\lambda |\langle u, \mathbf{Q}u \rangle|}] \leq \mathbb{E}[e^{2\lambda \langle u, \mathbf{Q}u \rangle}] + \mathbb{E}[e^{-2\lambda \langle u, \mathbf{Q}u \rangle}] \leq 2e^{\frac{8\lambda^2 \sigma^2}{n}}$$

for all $\lambda \in [0, \frac{nt}{4\sigma^2}]$. Since this bound holds for each choice of u , we have

$$\mathbb{E}[e^{\lambda \|\mathbf{Q}\|_{\text{op}}}] \leq \mathbb{E}[e^{2\lambda \max_{j=1, \dots, N} |\langle v^j, \mathbf{Q}v^j \rangle|}] \leq 2N e^{\frac{8\lambda^2 \sigma^2}{n}}.$$

Since $2(17^d) \leq e^{4d}$, the first bound (6.18) follows. The tail bound (6.19) follows as a 2235

consequence of Proposition 2.2.

2236

□ 2237

■ 6.4 Bounds for general matrices

2238

The preceding sections were devoted to bounds applicable to sample covariances under Gaussian or sub-Gaussian tail conditions. This section is devoted to developing extensions to more general tail conditions. In order to do so, it is convenient to introduce some more general methodology, which applies not only to sample covariance matrices, but also to more general random matrices. The main results in this section are Theorems 6.3 and 6.4, which are (essentially) matrix-based analogs of our earlier Hoeffding and Bernstein bounds for random variables. Before proving these results, we develop some useful matrix-theoretic generalizations of ideas from Chapter 2, including various types of tail conditions, as well as decompositions for the cumulant function for independent random matrices.

2239

2240

2241

2242

2243

2244

2245

2246

2247

2248

■ 6.4.1 Background on matrix analysis

2249

We begin by introducing some additional background on matrix-valued functions. Recall the class $\mathcal{S}^{d \times d}$ of symmetric $d \times d$ matrices. Any function $f : \mathbb{R}^d \rightarrow \mathbb{R}^d$ can be extended to a map from the set $\mathcal{S}^{d \times d}$ to itself in the following way. We begin with the eigendecomposition $\mathbf{W} = \mathbf{U}^T \Gamma \mathbf{U}$, where $\mathbf{U} \in \mathbb{R}^{d \times d}$ is a unitary matrix, satisfy the relation $\mathbf{U}^T \mathbf{U} = \mathbf{I}_d$, whereas $\Gamma = \text{diag}(\gamma(\mathbf{Q}))$ is a diagonal matrix specified by the vector of eigenvalues $\gamma(\mathbf{Q}) \in \mathbb{R}^d$. Using this notation, we consider the mapping on $\mathcal{S}^{d \times d}$ defined via $\mathbf{Q} \mapsto \mathbf{U}^T \text{diag}(f(\gamma(\mathbf{Q}))) \mathbf{U}$. In words, we apply the original function f to the vector of eigenvalues $\gamma(\mathbf{Q})$, and then rotate the resulting matrix $\text{diag}(f(\gamma(\mathbf{Q})))$ back to the original co-ordinate system defined by the eigenvectors of \mathbf{Q} . With a slight abuse of notation, we write

$$f(\mathbf{Q}) := \mathbf{U}^T \text{diag}(f(\gamma(\mathbf{Q}))) \mathbf{U} \quad (6.21)$$

for this induced mapping $f : \mathcal{S}^{d \times d} \rightarrow \mathcal{S}^{d \times d}$. By construction, this mapping is unitarily invariant, meaning that

$$f(\mathbf{V}^T \mathbf{Q} \mathbf{V}) = \mathbf{V}^T f(\mathbf{Q}) \mathbf{V} \quad \text{for all unitary matrices } \mathbf{V} \in \mathbb{R}^{d \times d},$$

since it affects only the eigenvalues (but not the eigenvectors) of \mathbf{Q} . Moreover, the eigenvalues of $f(\mathbf{Q})$ transform in a simple way, since we have

$$\gamma(f(\mathbf{Q})) = f(\gamma(\mathbf{Q})). \quad (6.22)$$

In words, the eigenvalues of the $f(\mathbf{Q})$ are simply the eigenvalues of \mathbf{Q} transformed by f , a result often referred to as the *spectral mapping property*.

Two functions that play a central role in our development of matrix tails bounds are the matrix exponential and the matrix logarithm. As a particular case of our construction, the matrix exponential has the power series expansion $e^{\mathbf{Q}} = \sum_{k=0}^{\infty} \frac{\mathbf{Q}^k}{k!}$. By the spectral mapping property, the eigenvalues of $e^{\mathbf{Q}}$ are positive, so that it is a positive definite matrix for any choice of \mathbf{Q} . Parts of our analysis also involve the matrix logarithm; when restricted to the cone of strictly positive definite matrices, as suffices for our purposes, the matrix logarithm corresponds to the inverse of the matrix exponential.

■ 6.4.2 Tail conditions for matrices

Given a symmetric random matrix $\mathbf{Q} \in \mathcal{S}^{d \times d}$, its polynomial moments, assuming that they exist, are the matrices defined by $\mathbb{E}[\mathbf{Q}^j]$. As shown in Exercise 6.2, the variance of \mathbf{Q} is a positive semidefinite matrix given by $\text{var}(\mathbf{Q}) = \mathbb{E}[\mathbf{Q}^2] - (\mathbb{E}[\mathbf{Q}])^2$. If \mathbf{Q} has polynomial moments of all orders, then its cumulant generating function $\Phi_{\mathbf{Q}} : \mathbb{R} \rightarrow \mathcal{S}^{d \times d}$ is given by

$$\Phi_{\mathbf{Q}}(\lambda) := \log \mathbb{E}[e^{\lambda \mathbf{Q}}], \quad (6.23)$$

and is guaranteed to be finite for all λ in an interval centered at zero. In parallel with our discussion in Chapter 2, various tail conditions are based on imposing bounds on this cumulant function. We begin with the simplest case:

Definition 6.1. A zero-mean symmetric random matrix $\mathbf{Q} \in \mathcal{S}^{d \times d}$ is sub-Gaussian with matrix parameter $\mathbf{V} \in \mathcal{S}_+^{d \times d}$ if

$$\Phi_{\mathbf{Q}}(\lambda) \preceq \frac{\lambda^2 \mathbf{V}}{2} \quad \text{for all } \lambda \in \mathbb{R}. \quad (6.24)$$

This definition is best understood by working through some simple examples.

Example 6.4. Suppose that $\mathbf{Q} = \varepsilon \mathbf{B}$ where $\varepsilon \in \{-1, +1\}$ is a Rademacher variable, and $\mathbf{B} \in \mathcal{S}^{d \times d}$ is a fixed matrix. Random matrices of this form frequently arise as the result of symmetrization arguments, as discussed at more length in the sequel. For the given random matrix, we have $\mathbb{E}[\mathbf{Q}^k] = 0$ for all odd k , and $\mathbb{E}[\mathbf{Q}^k] = \mathbf{B}^k$ for all even k , and hence

$$\mathbb{E}[e^{\lambda \mathbf{Q}}] = \sum_{k=0}^{\infty} \frac{\lambda^{2k}}{(2k)!} \mathbf{B}^{2k} \preceq \sum_{k=1}^{\infty} \frac{1}{k!} \left(\frac{\lambda^2 \mathbf{B}^2}{2} \right)^k = e^{\frac{\lambda^2 \mathbf{B}^2}{2}}$$

showing that the sub-Gaussian condition (6.24) holds with $\mathbf{V} = \mathbf{B}^2 = \text{var}(\mathbf{Q})$. ♣ 2269

As we show in Exercise 6.3, more generally, a random matrix of the form $\mathbf{Q} = g\mathbf{B}$, where $g \in \mathbb{R}$ is a σ -sub-Gaussian variable with distribution symmetric around zero, satisfies the condition (6.24) with matrix parameter $\mathbf{V} = \sigma^2\mathbf{B}^2$. 2270
2271
2272

Example 6.5. As an extension of the previous example, consider a random matrix of the form $\mathbf{Q} = \varepsilon\mathbf{C}$, where ε is a Rademacher variable as before, and \mathbf{C} is now a random matrix, independent of ε , that satisfies the bound $\|\mathbf{C}\|_{\text{op}} \leq b$. First fixing \mathbf{C} and taking expectations over the Rademacher variable, the previous example yields $\mathbb{E}_\varepsilon[e^{\lambda\varepsilon\mathbf{C}}] \preceq e^{\frac{\lambda^2}{2}\mathbf{C}^2}$. Since $\|\mathbf{C}\|_{\text{op}} \leq b$, we have $e^{\frac{\lambda^2}{2}\mathbf{C}^2} \preceq e^{\frac{\lambda^2}{2}b^2\mathbf{I}_d}$, and hence

$$\Phi_{\mathbf{Q}}(\lambda) \preceq \frac{\lambda^2}{2}b^2\mathbf{I}_d \quad \text{for all } \lambda \in \mathbb{R},$$

showing that \mathbf{Q} is sub-Gaussian with matrix parameter $\mathbf{V} = b^2\mathbf{I}_d$. 2273

♣ 2274

In parallel with our treatment of scalar random variables in Chapter 2, it is natural to consider various weakenings of the sub-Gaussian requirement. 2275
2276

Definition 6.2 (Sub-exponential random matrices). A zero-mean random matrix is sub-exponential with parameters (\mathbf{V}, b) if the cumulant function $\Phi_{\mathbf{Q}}(\lambda)$ is finite for all $|\lambda| < \frac{1}{b}$. 2277
2278
2279

Thus, any sub-Gaussian random matrix is also sub-exponential with parameters (\mathbf{V}, ∞) . However, there also exist sub-exponential random matrices that are not sub-Gaussian. One example is the zero-mean random matrix $\mathbf{M} = \varepsilon g^2\mathbf{B}$, where $\varepsilon \in \{-1, +1\}$ is a Rademacher, the variable $g \sim \mathcal{N}(0, 1)$ is independent of ε , and \mathbf{B} is a fixed symmetric matrix. 2280
2281
2282
2283
2284
2285

The Bernstein condition for random matrices provides one useful way of certifying the sub-exponential condition: 2286
2287

Definition 6.3 (Bernstein condition for matrices). A zero-mean symmetric random matrix \mathbf{Q} satisfies a Bernstein condition with parameter $b > 0$ if

$$\mathbb{E}[\mathbf{Q}^j] \preceq \frac{1}{2}j!b^{j-2}\text{var}(\mathbf{Q}) \quad \text{for } j = 3, 4, \dots \quad (6.25)$$

We note that (a stronger form of) the Bernstein condition holds whenever the matrix \mathbf{Q} has a bounded operator norm—say $\|\mathbf{Q}\|_{\text{op}} \leq b$ almost surely. In this case, it can be 2288
2289
2290

shown (see Exercise 6.5) that

$$\mathbb{E}[\mathbf{Q}^j] \preceq b^{j-2} \text{var}(\mathbf{Q}) \quad \text{for all } j = 3, 4, \dots \quad (6.26)$$

Exercise 6.7 gives an example of a random matrix with unbounded operator norm for which the Bernstein condition holds.

The following lemma shows how the general Bernstein condition (6.25) implies the sub-exponential condition. More generally, the argument given here provides an explicit bound on the cumulant generating function:

Lemma 6.1. For any symmetric zero-mean random matrix satisfying the Bernstein condition (6.25), we have

$$\Phi_{\mathbf{Q}}(\lambda) \preceq \frac{\lambda^2 \text{var}(\mathbf{Q})}{1 - b|\lambda|} \quad \text{for all } |\lambda| < \frac{1}{b}. \quad (6.27)$$

Proof. Since $\mathbb{E}[\mathbf{Q}] = 0$, applying the definition of the matrix exponential for a suitably small $\lambda \in \mathbb{R}$ yields

$$\begin{aligned} \mathbb{E}[e^{\lambda \mathbf{Q}}] &= \mathbf{I}_d + \frac{\lambda^2 \text{var}(\mathbf{Q})}{2} + \sum_{j=3}^{\infty} \frac{\lambda^j \mathbb{E}[\mathbf{Q}^j]}{j!} \\ &\stackrel{(i)}{\preceq} \mathbf{I}_d + \frac{\lambda^2 \text{var}(\mathbf{Q})}{2} \left\{ \sum_{j=0}^{\infty} |\lambda|^j b^j \right\} \\ &\stackrel{(ii)}{\preceq} \mathbf{I}_d + \frac{\lambda^2 \text{var}(\mathbf{Q})}{2(1 - b|\lambda|)}, \\ &\stackrel{(iii)}{\preceq} \exp\left(\frac{\lambda^2 \text{var}(\mathbf{Q})}{2(1 - b|\lambda|)}\right), \end{aligned}$$

where step (i) applies the Bernstein condition; step (ii) is valid for any $|\lambda| < 1/b$, a choice for which the geometric series is summable; and step (iii) follows from the spectral theorem, and the elementary inequality $1 + v \leq e^v$. Since taking matrix logarithms preserves the positive semidefinite order, this is equivalent to the claim (6.27). \square

■ 6.4.3 Matrix-Chernoff approach and independent decompositions

The Chernoff approach to tail bounds, as discussed in Chapter 2, is based on controlling the cumulant generating function of a random variable. In this section, we begin by showing that the trace of the matrix cumulant generating function (6.23) plays a similar role in bounding the operator norm of random matrices.

Lemma 6.2 (Matrix Chernoff technique). Let \mathbf{Q} be a zero-mean symmetric random matrix whose cumulant function $\Phi_{\mathbf{Q}}$ exists in an open interval $(-a, a)$. Then for any $\delta > 0$, we have

$$\mathbb{P}[\gamma_{\max}(\mathbf{Q}) \geq \delta] \leq \text{tr}(e^{\Phi_{\mathbf{Q}}(\lambda)}) e^{-\lambda\delta} \quad \text{for all } \lambda \in [0, a), \quad (6.28) \quad 2309$$

where $\text{tr}(\cdot)$ denotes the trace operator on matrices. Similarly, we have

$$\mathbb{P}[\|\mathbf{Q}\|_{\text{op}} \geq \delta] \leq 2 \text{tr}(e^{\Phi_{\mathbf{Q}}(\lambda)}) e^{-\lambda\delta} \quad \text{for all } \lambda \in [0, a). \quad (6.29) \quad 2310$$

Proof. For each $\lambda \in [0, a)$, we have

$$\mathbb{P}[\gamma_{\max}(\mathbf{Q}) \geq \delta] = \mathbb{P}[e^{\gamma_{\max}(\lambda\mathbf{Q})} \geq e^{\lambda\delta}] \stackrel{(i)}{=} \mathbb{P}[\gamma_{\max}(e^{\lambda\mathbf{Q}}) \geq e^{\lambda\delta}],$$

where step (i) uses the functional calculus relating the eigenvalues of $\lambda\mathbf{Q}$ to those of $e^{\lambda\mathbf{Q}}$. Applying Markov's inequality yields

$$\mathbb{P}[\gamma_{\max}(e^{\lambda\mathbf{Q}}) \geq e^{\lambda\delta}] \leq \mathbb{E}[\gamma_{\max}(e^{\lambda\mathbf{Q}})] e^{-\lambda\delta} \stackrel{(i)}{\leq} \mathbb{E}[\text{tr}(e^{\lambda\mathbf{Q}})] e^{-\lambda\delta}. \quad (6.30)$$

Here inequality (i) uses the upper bound $\gamma_{\max}(e^{\lambda\mathbf{Q}}) \leq \text{tr}(e^{\lambda\mathbf{Q}})$, which holds since $e^{\lambda\mathbf{Q}}$ is positive definite. Finally, since trace and expectation commute, we have

$$\mathbb{E}[\text{tr}(e^{\lambda\mathbf{Q}})] = \text{tr}(\mathbb{E}[e^{\lambda\mathbf{Q}}]) \stackrel{(ii)}{=} \text{tr}(e^{\Phi_{\mathbf{Q}}(\lambda)}),$$

where equality (ii) uses the fact that the matrix logarithm and exponential are inverses, and the definition (6.23) of the matrix cumulant generating function. 2311
2312

Note that the same argument can be applied to bound the event $\gamma_{\max}(-\mathbf{Q}) \geq \delta$, or equivalently the event $\gamma_{\min}(\mathbf{Q}) \leq -\delta$. Since $\|\mathbf{Q}\|_{\text{op}} = \max\{\gamma_{\max}(\mathbf{Q}), |\gamma_{\min}(\mathbf{Q})|\}$, the tail bound on the operator norm (6.29) follows. 2313
2314
2315

An important property of independent random variables is that the cumulant function of their sum also decomposes additively. For random matrices, this type of decomposition is no longer guaranteed to hold with equality, essentially because matrix products need not commute. However, for independent random matrices, it is nonetheless possible to establish an upper bound in terms of the trace of the cumulant generating functions, as we now show. 2316
2317
2318
2319
2320
2321

Lemma 6.3. Let $\mathbf{Q}_1, \dots, \mathbf{Q}_n$ be independent symmetric random matrices whose cumulant functions exist for all $\lambda \in I$, and define the sum $\mathbf{S}_n = \sum_{i=1}^n \mathbf{Q}_i$. Then 2322

$$\text{tr}(e^{\Phi_{\mathbf{S}_n}(\lambda)}) \leq \text{tr}(e^{\sum_{i=1}^n \Phi_{\mathbf{Q}_i}(\lambda)}) \quad \text{for all } \lambda \in I. \quad (6.31) \quad 2323$$

Remark: In conjunction with Lemma 6.2, this lemma provides an avenue for obtaining tail bounds on the operator norm of sums of independent random matrices. In particular, if we apply the upper bound (6.29) to the random matrix \mathbf{S}_n/n , we find that

$$\mathbb{P}\left[\left\|\frac{1}{n}\sum_{i=1}^n \mathbf{Q}_i\right\|_{\text{op}} \geq \delta\right] \leq 2 \operatorname{tr}\left(e^{\sum_{i=1}^n \Phi_{\mathbf{Q}_i}(\lambda)}\right) e^{-\lambda n \delta} \quad \text{for all } \lambda \in [0, a). \quad (6.32)$$

Proof. In order to prove this lemma, we require the following result due to Lieb [Lie73]: for any fixed matrix $\mathbf{H} \in \mathcal{S}^{d \times d}$, the function $f : \mathcal{S}_+^{d \times d} \rightarrow \mathbb{R}$ given by

$$f(\mathbf{A}) := \operatorname{tr}\left(e^{\mathbf{H} + \log(\mathbf{A})}\right)$$

is concave. Introducing the shorthand notation $G(\lambda) := \operatorname{tr}\left(e^{\Phi_{\mathbf{S}_n}(\lambda)}\right)$, we note that, by linearity of trace and expectation, we have

$$G(\lambda) = \operatorname{tr}\left(\mathbb{E}\left[e^{\lambda \mathbf{S}_{n-1} + \log \exp(\lambda \mathbf{Q}_n)}\right]\right) = \mathbb{E}_{\mathbf{S}_{n-1}} \left[\operatorname{tr}\left(e^{\lambda \mathbf{S}_{n-1} + \log \exp(\lambda \mathbf{Q}_n)}\right) \right],$$

Using concavity of the function f with $\mathbf{H} = \lambda \mathbf{S}_{n-1}$ and $\mathbf{A} = e^{\lambda \mathbf{Q}_n}$, Jensen's inequality implies that

$$\begin{aligned} \mathbb{E}_{\mathbf{Q}_n} \left[\operatorname{tr}\left(e^{\lambda \mathbf{S}_{n-1} + \log \exp(\lambda \mathbf{Q}_n)}\right) \right] &\leq \operatorname{tr}\left(e^{\lambda \mathbf{S}_{n-1} + \log \mathbb{E}_{\mathbf{Q}_n} \exp(\lambda \mathbf{Q}_n)}\right) \\ &= \operatorname{tr}\left(e^{\lambda \mathbf{S}_{n-1} + \Phi_{\mathbf{Q}_n}(\lambda)}\right), \end{aligned}$$

so that we have shown that $G(\lambda) \leq \mathbb{E}_{\mathbf{S}_{n-1}} \left[\operatorname{tr}\left(e^{\lambda \mathbf{S}_{n-1} + \Phi_{\mathbf{Q}_n}(\lambda)}\right) \right]$. 2325

We now recurse this argument, in particular peeling off the term involving \mathbf{Q}_{n-1} , so that we have

$$G(\lambda) \leq \mathbb{E}_{\mathbf{S}_{n-2}} \mathbb{E}_{\mathbf{Q}_{n-1}} \left[\operatorname{tr}\left(e^{\lambda \mathbf{S}_{n-2} + \Phi_{\mathbf{Q}_n}(\lambda) + \log \exp(\lambda \mathbf{Q}_{n-1})}\right) \right].$$

We again exploit the concavity of the function f , this time with $\mathbf{H} = \lambda \mathbf{S}_{n-2} + \Phi_{\mathbf{Q}_n}(\lambda)$ and $\mathbf{A} = e^{\lambda \mathbf{Q}_{n-1}}$ to conclude that

$$G(\lambda) \leq \mathbb{E}_{\mathbf{S}_{n-2}} \left[\operatorname{tr}\left(e^{\lambda \mathbf{S}_{n-2} + \Phi_{\mathbf{Q}_{n-1}}(\lambda) + \Phi_{\mathbf{Q}_n}(\lambda)}\right) \right],$$

and continuing on in this manner yields the claim. □ 2326

■ 6.4.4 Upper tail bounds for random matrices 2327

We now have collected the ingredients necessary for stating and proving various tail 2328
bounds for the deviations of sums of zero-mean independent random matrices. 2329

Sub-Gaussian case:

We begin with a tail bound for sub-Gaussian random matrices. It provides an approximate analog of the Hoeffding-type tail bound for random variables (Proposition 2.1).

Theorem 6.3 (Hoeffding bound for random matrices). Let $\{\mathbf{Q}_i\}_{i=1}^n$ be a sequence of zero-mean independent symmetric random matrices that satisfy the sub-Gaussian condition with parameters $\{\mathbf{V}_i\}_{i=1}^n$. Then for all $\delta > 0$, we have the upper tail bound

$$\mathbb{P}\left[\left\|\frac{1}{n}\sum_{i=1}^n \mathbf{Q}_i\right\|_{\text{op}} \geq \delta\right] \leq 2 \min\{n, d\} e^{-\frac{n\delta^2}{2\sigma^2}} \quad \text{where } \sigma^2 = \left\|\frac{1}{n}\sum_{i=1}^n \mathbf{V}_i\right\|_{\text{op}}. \quad (6.33)$$

Proof. From Definition 6.1 and Lemma 6.3, we have $\text{tr}(e^{\sum_{i=1}^n \Phi_{\mathbf{Q}_i}(\lambda)}) \leq \text{tr}(e^{\frac{\lambda^2}{2} \sum_{i=1}^n \mathbf{V}_i})$. This upper bound, when combined with the matrix Chernoff bound (6.32), yields

$$\mathbb{P}\left[\left\|\frac{1}{n}\sum_{i=1}^n \mathbf{Q}_i\right\|_{\text{op}} \geq \delta\right] \leq 2 \text{tr}\left(e^{\frac{\lambda^2}{2} \sum_{i=1}^n \mathbf{V}_i}\right) e^{-\lambda n \delta}$$

For any symmetric matrix \mathbf{Q} , we have $\text{tr}(e^{\mathbf{Q}}) \leq \text{rank}(\mathbf{Q}) e^{\|\mathbf{Q}\|_{\text{op}}}$. Applying this inequality to the matrix $\mathbf{Q} = \frac{\lambda^2}{2} \sum_{i=1}^n \mathbf{V}_i$, we have $\text{rank}(\mathbf{Q}) \leq \min\{n, d\}$, whence

$$\mathbb{P}\left[\left\|\frac{1}{n}\sum_{i=1}^n \mathbf{Q}_i\right\|_{\text{op}} \geq \delta\right] \leq 2 \min\{n, d\} e^{\frac{\lambda^2}{2} n \sigma^2 - \lambda n \delta}.$$

This upper bound holds for all $\lambda \geq 0$; the optimal choice of $\lambda = \delta/\sigma^2$ yields the claim. \square

An important fact is that inequality (6.33) also implies an analogous bound for general independent but potentially non-symmetric and/or non-square matrices, with d replaced by $(d_1 + d_2)$. More specifically, a problem involving general zero-mean random matrices $\mathbf{A}_i \in \mathbb{R}^{d_1 \times d_2}$ can be transformed to a symmetric version by defining the matrices

$$\mathbf{Q}_i := \begin{bmatrix} \mathbf{0}_{d_1 \times d_1} & \mathbf{A}_i \\ \mathbf{A}_i^T & \mathbf{0}_{d_2 \times d_2} \end{bmatrix} \quad (6.34)$$

and imposing a sub-Gaussian condition (6.24) on the symmetric matrices \mathbf{Q}_i . See Exercise 6.6 for further details.

A significant feature of the tail bound (6.33) is the appearance of the factor $\min\{n, d\}$ in front of the exponent. In certain cases, this factor is superfluous, and leads to sub-optimal bounds. However, it cannot be avoided in general. The following example

illustrates these two extremes.

2344

Example 6.6 (Looseness/sharpness of Theorem 6.3). For simplicity, let us consider 2345
examples with $n = d$. For each $i = 1, 2, \dots, d$, let $\mathbf{E}_i \in \mathcal{S}^{d \times d}$ denote the diagonal 2346
matrix with 1 in position (i, i) , and zeroes elsewhere. Define $\mathbf{Q}_i = y_i \mathbf{E}_i$, where $\{y_i\}_{i=1}^n$ 2347
is an i.i.d. sequence of 1-sub-Gaussian variables. Two specific cases to keep in mind are 2348
Rademacher variables $\{\varepsilon_i\}_{i=1}^n$, and $\mathcal{N}(0, 1)$ variables $\{g_i\}_{i=1}^n$. 2349

For any such choice of sub-Gaussian variables, a calculation similar to that of Ex-
ample 6.4 shows that each \mathbf{Q}_i satisfies the sub-Gaussian bound (6.24) with $\mathbf{V}_i = \mathbf{E}_i$,
and hence $\sigma^2 = \|\frac{1}{d} \sum_{i=1}^d \mathbf{V}_i\|_{\text{op}} = 1/d$. Consequently, an application of Theorem 6.3
yields the tail bound

$$\mathbb{P}\left[\left\|\frac{1}{d} \sum_{i=1}^d \mathbf{Q}_i\right\|_{\text{op}} \geq \delta\right] \leq 2 d e^{-\frac{d^2 \delta^2}{2}} \quad \text{for all } \delta > 0, \quad (6.35)$$

which implies that $\|\frac{1}{d} \sum_{j=1}^d \mathbf{Q}_j\|_{\text{op}} \lesssim \frac{\sqrt{2 \log(2d)}}{d}$ with high probability. On the other
hand, an explicit calculation shows that

$$\left\|\frac{1}{d} \sum_{i=1}^n \mathbf{Q}_i\right\|_{\text{op}} = \max_{i=1, \dots, d} \frac{|y_i|}{d}. \quad (6.36)$$

Comparing the exact result (6.36) with the bound (6.35) yields a range of behav-
ior. At one extreme, for i.i.d. Rademacher variables $y_i = \varepsilon_i \in \{-1, +1\}$, we have
 $\|\frac{1}{d} \sum_{i=1}^n \mathbf{Q}_i\|_{\text{op}} = 1/d$, showing that the bound (6.35) is off by the order $\sqrt{\log d}$. On the
other hand, for i.i.d. Gaussian variables $y_i = g_i \sim \mathcal{N}(0, 1)$, we have

$$\left\|\frac{1}{d} \sum_{i=1}^d \mathbf{Q}_i\right\|_{\text{op}} = \max_{i=1, \dots, d} \frac{|g_i|}{d} \simeq \frac{\sqrt{2 \log d}}{d},$$

using the fact that the maximum of d i.i.d. $\mathcal{N}(0, 1)$ variables scales as $\sqrt{2 \log d}$. Conse- 2350
quently, Theorem 6.3 cannot be improved for this class of random matrices. ♣ 2351

Bernstein-type bounds for random matrices

2352

We now turn to bounds on random matrices that satisfy sub-exponential tail conditions, 2353
in particular of the Bernstein form (6.25). 2354

2355

Theorem 6.4 (Bernstein bound for random matrices). Let $\mathbf{Q}_1, \dots, \mathbf{Q}_n$ be a sequence of independent, zero-mean, symmetric random matrices that satisfy the Bernstein condition (6.25) with parameter $b > 0$. Then for all $\delta \geq 0$, the operator norm satisfies the tail bound

$$\mathbb{P} \left[\left\| \frac{1}{n} \sum_{i=1}^n \mathbf{Q}_i \right\|_{\text{op}} \geq \delta \right] \leq 2 \min\{d, n\} \exp \left\{ - \frac{n\delta^2}{2(\sigma^2 + b\delta)} \right\}, \quad (6.37)$$

where $\sigma^2 := \frac{1}{n} \left\| \sum_{j=1}^n \text{var}(\mathbf{Q}_j) \right\|_{\text{op}}$.

Proof. By Lemma 6.3, we have $\text{tr}(e^{\Phi \mathbf{s}_n(\lambda)}) \leq \text{tr}(e^{\sum_{i=1}^n \Phi \mathbf{Q}_i(\lambda)})$. By Lemma 6.1, the Bernstein condition implies that $\Phi \mathbf{Q}_i(\lambda) \preceq \frac{\lambda^2 \text{var}(\mathbf{Q}_i)}{1-b|\lambda|}$ for any $|\lambda| < \frac{1}{b}$. Putting together the pieces yields

$$\text{tr}(e^{\sum_{i=1}^n \Phi \mathbf{Q}_i(\lambda)}) \leq \text{tr} \left(\exp \left(\frac{\lambda^2 \sum_{i=1}^n \text{var}(\mathbf{Q}_i)}{1-b|\lambda|} \right) \right) \leq \min\{n, d\} e^{\frac{n\lambda^2 \sigma^2}{1-b|\lambda|}}, \quad (6.38)$$

where the final inequality uses the same argument as the proof of Theorem 6.3. Combined with the upper bound (6.32), we find that

$$\mathbb{P} \left[\left\| \frac{1}{n} \sum_{i=1}^n \mathbf{Q}_i \right\|_{\text{op}} \geq \delta \right] \leq 2 \min\{n, d\} e^{\frac{n\sigma^2 \lambda^2}{1-b|\lambda|} - \lambda n \delta},$$

valid for all $\lambda \in [0, 1/b)$. Setting $\lambda = \frac{\delta^2}{\sigma^2 + b\delta} \in (0, \frac{1}{b})$ and simplifying yields the claim (6.37). \square

Remarks: Note that the tail bound (6.37) is of the sub-exponential type, with two regimes of behavior depending on the relative sizes of the parameters σ^2 and b . Thus, it is a natural generalization of the classical Bernstein bound for scalar random variables. As with Theorem 6.3, Theorem 6.4 can also be generalized to non-symmetric (and potentially non-square) matrices $\{\mathbf{A}_i\}_{i=1}^n$, as long as we adopt the new definition

$$\sigma^2 := \max \left\{ \left\| \frac{1}{n} \sum_{i=1}^n \mathbb{E}[\mathbf{A}_i \mathbf{A}_i^T] \right\|_{\text{op}}, \left\| \frac{1}{n} \sum_{i=1}^n \mathbb{E}[\mathbf{A}_i^T \mathbf{A}_i] \right\|_{\text{op}} \right\}, \quad (6.39)$$

and replace d by $(d_1 + d_2)$. Doing so involves defining the symmetrized analogues (6.34), and analyzing their properties. We provide an instance of this operation in the next example.

The problem of matrix completion provides an interesting class of examples in which Theorem 6.4 can be fruitfully applied. See Chapter 10 for a detailed description of the

underlying problem, which motivates the following discussion.

Example 6.7 (Tail bounds in matrix completion). Consider an i.i.d. sequence of matrices of the form $\mathbf{A}_i = \xi_i \mathbf{X}_i \in \mathbb{R}^{d \times d}$. Here ξ_i is a zero-mean sub-exponential variable that satisfies the Bernstein condition with parameter b and variance ν^2 . For the moment, we assume that its distribution is symmetric around zero (meaning that $-\xi_i$ has the same distribution as ξ_i). Suppose that \mathbf{X}_i is a random “mask matrix”, independent from ξ_i , with a single entry equal to d in a position chosen uniformly at random from all d^2 entries, and zeros elsewhere. With this particular scaling, for any fixed matrix $\Theta \in \mathbb{R}^{d \times d}$, we have $\mathbb{E}[\langle \mathbf{A}_i, \Theta \rangle^2] = \|\Theta\|_F^2$ — a property that plays an important role in our later analysis of matrix completion.

This is a sequence of non-symmetric matrices, but as discussed following Theorem 6.4, we can bound the operator norm $\|\frac{1}{n} \sum_{i=1}^n \mathbf{A}_i\|_{\text{op}}$ in terms of the operator norm $\|\frac{1}{n} \sum_{i=1}^n \mathbf{Q}_i\|_{\text{op}}$, where the symmetrized version $\mathbf{Q}_i \in \mathbb{R}^{2d \times 2d}$ was defined in equation (6.34). By the independence between ξ_i and \mathbf{A}_i , we have $\mathbb{E}[\mathbf{Q}_i^{2m+1}] = 0$ for all $m = 0, 1, 2, \dots$. Turning to the even moments, suppose that entry (a, b) is the only non-zero in the mask matrix \mathbf{X}_i . We then have

$$\mathbf{Q}_i^{2m} = (\xi_i)^{2m} d^{2m} \begin{bmatrix} \mathbf{D}_a & 0 \\ 0 & \mathbf{D}_b \end{bmatrix} \quad \text{for all } m = 1, 2, \dots, \quad (6.40)$$

where $\mathbf{D}_a \in \mathbb{R}^{d \times d}$ is the diagonal matrix with a single one in entry (a, a) , with \mathbf{D}_b defined analogously. By the Bernstein condition, we have $\mathbb{E}[\xi_i^{2m}] \leq \frac{1}{2}(2m)! b^{2m-2} \nu^2$ for all $m = 1, 2, \dots$

On the other hand, $\mathbb{E}[\mathbf{D}_a] = \frac{1}{d} \mathbf{I}_d$ since the probability of choosing a in the first co-ordinate is $1/d$. We thus see that $\text{var}(\mathbf{Q}_i) = \nu^2 d \mathbf{I}_{2d}$. Putting together the pieces, we have shown that

$$\mathbb{E}[\mathbf{Q}_i^{2m}] \preceq \frac{1}{2}(2m)! b^{2m-2} \sigma^2 d^{2m} \frac{1}{d} \mathbf{I}_{2d} = \frac{1}{2}(2m)! (bd)^{2m-2} \text{var}(\mathbf{Q}_i),$$

showing that \mathbf{Q}_i satisfies the Bernstein condition with parameter bd , and that

$$\sigma^2 = \left\| \frac{1}{n} \sum_{i=1}^n \text{var}(\mathbf{Q}_i) \right\|_{\text{op}} \leq \nu^2 d.$$

Consequently, Theorem 6.4 implies that

$$\mathbb{P} \left[\left\| \frac{1}{n} \sum_{i=1}^n \mathbf{A}_i \right\|_{\text{op}} \geq \delta \right] \leq 4d e^{-\frac{n\delta^2}{2d(\nu^2 + b\delta)}}. \quad (6.41)$$

♣ 2379

In many problems, it is convenient to deal with random matrices \mathbf{Q}_i that have a

distribution symmetric around zero (meaning that $-\mathbf{Q}_i$ and \mathbf{Q}_i follow the same distribution). In the previous example, we enforced this distributional property by assuming that ξ_i had a symmetric distribution. However, in general, it is always possible to reduce to the case of symmetric distributions, as shown in the following example:

Example 6.8 (Symmetrization and operator norm bounds). By Markov's inequality, we have

$$\mathbb{P}[\gamma_{\max}(\sum_{i=1}^n \mathbf{Q}_i) \geq \delta] \leq \mathbb{E}[e^{\gamma_{\max}(\sum_{i=1}^n \mathbf{Q}_i)}] e^{-\lambda \delta}.$$

By the variational representation of the maximum eigenvalue, we have

$$\begin{aligned} \mathbb{E}[e^{\lambda \gamma_{\max}(\sum_{i=1}^n \mathbf{Q}_i)}] &= \mathbb{E}[\exp(\lambda \sup_{\|u\|_2=1} u^T (\sum_{i=1}^n \mathbf{Q}_i) u)] \\ &\stackrel{(i)}{\leq} \mathbb{E}[\exp(2\lambda \sup_{\|u\|_2=1} u^T (\sum_{i=1}^n \varepsilon_i \mathbf{Q}_i) u)] \\ &= \mathbb{E}[\exp(2\lambda \gamma_{\max}(\sum_{i=1}^n \varepsilon_i \mathbf{Q}_i))], \end{aligned}$$

where inequality (i) makes use of the symmetrization inequality from Proposition 4.1(b) with $\Phi(t) = e^{t}$. From this point, the argument proceeds as before: in particular, upper bounding the operator norm by the trace and applying Lemma 6.3, we find that

$$\mathbb{E}[e^{\lambda \gamma_{\max}(\sum_{i=1}^n \mathbf{Q}_i)}] \leq \text{tr}(\mathbb{E}[\exp(2\lambda \sum_{i=1}^n \varepsilon_i \mathbf{Q}_i)]) \leq \text{tr}(e^{\sum_{i=1}^n \Phi_{\tilde{\mathbf{Q}}_i}(\lambda)}),$$

where we have defined the symmetrized and rescaled versions $\tilde{\mathbf{Q}}_i = 2\varepsilon_i \mathbf{Q}_i$. Consequently, apart from the factor of two, we may assume without loss of generality that our matrices have a distribution symmetric around zero. ♣

■ 6.4.5 Consequences for covariance matrices

We conclude with a useful corollary of Theorem 6.4 for the estimation of covariance matrices.

Corollary 6.1. Let x_1, \dots, x_n be i.i.d. zero-mean random vectors with covariance Σ such that $\|x_j\|_2 \leq \sqrt{b}$ almost surely. Then for all $\delta > 0$, the sample covariance matrix $\hat{\Sigma} = \frac{1}{n} \sum_{i=1}^n x_i x_i^T$ satisfies

$$\mathbb{P}[\|\hat{\Sigma} - \Sigma\|_{\text{op}} \geq \delta] \leq 2 \min\{d, n\} \exp\left(-\frac{n\delta^2}{2b(\|\Sigma\|_{\text{op}} + \delta)}\right). \quad (6.42)$$

2393

Proof. We apply Theorem 6.4 to the zero-mean random matrices $\mathbf{Q}_i := x_i x_i^T - \mathbf{\Sigma}$. These matrices have controlled operator norm: indeed, by triangle inequality, we have

$$\|\mathbf{Q}_i\|_{\text{op}} \leq \|x_i\|_2^2 + \|\mathbf{\Sigma}\|_{\text{op}} \leq b + \|\mathbf{\Sigma}\|_{\text{op}}.$$

Since $\mathbf{\Sigma} = \mathbb{E}[x_i x_i^T]$, we have $\|\mathbf{\Sigma}\|_{\text{op}} = \max_{v \in \mathbb{S}^{d-1}} \mathbb{E}[\langle v, x_i \rangle^2] \leq b$, and hence $\|\mathbf{Q}_i\|_{\text{op}} \leq 2b$. Turning to the variance of \mathbf{Q}_i , we have

$$\text{var}(\mathbf{Q}_i) = \mathbb{E}[(x_i x_i^T)^2] - \mathbf{\Sigma}^2 \preceq \mathbb{E}[\|x_i\|_2^2 x_i x_i^T] \preceq b \mathbf{\Sigma},$$

so that $\|\text{var}(\mathbf{Q}_i)\|_{\text{op}} \leq b \|\mathbf{\Sigma}\|_{\text{op}}$. Substituting into the tail bound (6.37) yields the claim. 2394

□ 2395

Let us illustrate some consequences of this corollary with some examples. 2396

2397

Example 6.9 (Random vectors uniform on sphere). Suppose that the random vectors x_i are chosen uniformly from the sphere $\mathbb{S}^{d-1}(\sqrt{d})$, so that $\|x_i\|_2 = \sqrt{d}$ for all $i = 1, \dots, n$. By construction, we have $\mathbb{E}[x_i x_i^T] = \mathbf{\Sigma} = \mathbf{I}_d$, and hence $\|\mathbf{\Sigma}\|_{\text{op}} = 1$. Applying Corollary 6.1 yields

$$\mathbb{P}[\|\widehat{\mathbf{\Sigma}} - \mathbf{I}_d\|_{\text{op}} \geq \delta] \leq 2 \min\{n, d\} e^{-\frac{n\delta^2}{2d+2\delta}} \quad \text{for all } \delta \geq 0. \quad (6.43)$$

This bound implies that

$$\|\widehat{\mathbf{\Sigma}} - \mathbf{I}_d\|_{\text{op}} \lesssim \sqrt{\frac{d \log d}{n}} + \frac{d \log d}{n} \quad (6.44)$$

with high probability, so that the sample covariance is a consistent estimate as long 2398
as $\frac{d \log d}{n} \rightarrow 0$. This result is close to optimal, with only the extra logarithmic factor 2399
being superfluous. For instance, it can be removed by noting that x_i is a sub-Gaussian 2400
random vector, and then applying Theorem 6.2. ♣ 2401

Example 6.10 (“Spiked” random vectors). We now consider an ensemble of random 2402
vectors that are rather different than the previous example, but still satisfy the same 2403
bound. In particular, consider a random vector of the form $x_i = \sqrt{d} e_{a(i)}$, where $a(i)$ is 2404
an index chosen uniformly at random from $\{1, \dots, d\}$, and $e_{a(i)} \in \mathbb{R}^d$ is the canonical 2405
basis vector with 1 in position $a(i)$. As before, we have $\|x_i\|_2 = \sqrt{d}$, and $\mathbb{E}[x_i x_i^T] = \mathbf{I}_d$ 2406
so that the tail bound (6.43) also applies to this ensemble. An interesting fact is that 2407
for this particular ensemble, the bound (6.44) is sharp, meaning it cannot be improved 2408
beyond constant factors. ♣ 2409

■ 6.5 Bounds for structured covariance matrices

2410

In the preceding sections, our primary focus has been estimation of general unstructured covariance matrices via the sample covariance. When a covariance matrix is equipped with additional structure, faster rates of estimation are possible using different estimators than the sample covariance matrix. We have already seen instances of this phenomenon: for example, when Theorem 6.1 is applied to matrices with trace norm (6.12), it is this trace constraint, as opposed to the dimension d , that enters the rates. In this section, we explore other forms of structure, such as those involving sparsity or graph structure.

In the simplest setting, the covariance matrix is known to be sparse, and the positions of the non-zero entries are known. In such settings, it is natural to consider matrix estimators that are non-zero only in these known positions. For instance, if we are given a priori knowledge that the covariance matrix is diagonal, then it would be natural to use the estimate $\hat{\mathbf{D}} := \text{diag}\{\hat{\Sigma}_{11}, \hat{\Sigma}_{22}, \dots, \hat{\Sigma}_{dd}\}$, corresponding to the diagonal entries of the sample covariance matrix $\hat{\Sigma}$. As we explore in Exercise 6.11, the performance of this estimator can be substantially better: in particular, for sub-Gaussian variables, it achieves an estimation error of the order $\sqrt{\frac{\log d}{n}}$, as opposed to the order $\sqrt{\frac{d}{n}}$ rates in the unstructured setting. Similar statements apply to other forms of known sparsity.

■ 6.5.1 Unknown sparsity and thresholding

2428

More generally, suppose that the covariance matrix Σ is known to be relatively sparse, but that the positions of the non-zero entries are no longer known. It is then natural to consider estimators based on thresholding. Given a parameter $\lambda > 0$, the *hard thresholding operator* is a function $T_\lambda : \mathbb{R} \rightarrow \mathbb{R}$ such that

$$T_\lambda(u) := u \mathbb{I}[|u| > \lambda] = \begin{cases} u & \text{if } |u| > \lambda \\ 0 & \text{otherwise.} \end{cases} \quad (6.45)$$

With a minor abuse of notation, for a matrix M , we write $T_\lambda(M)$ for the matrix obtained by applying the thresholding operator to each element of M . In this section, we study the performance of the estimator $T_{\lambda_n}(\hat{\Sigma})$, where the parameter $\lambda_n > 0$ is suitably chosen as a function of the sample size n and matrix dimension d . The sparsity of the covariance matrix can be measured in various ways. The zero-pattern of the covariance matrix is captured by the adjacency matrix $\mathbf{A} \in \mathbb{R}^{d \times d}$ with entries $A_{j\ell} = \mathbb{I}[\Sigma_{j\ell} \neq 0]$. This adjacency matrix defines the edge structure of an undirected graph G on the vertices $\{1, 2, \dots, d\}$, with edge (j, ℓ) included in the graph if and only if $\Sigma_{j\ell} \neq 0$. The operator norm $\|\mathbf{A}\|_{\text{op}}$ of the adjacency matrix provides a natural measure of sparsity. In particular, it can be verified that $\|\mathbf{A}\|_{\text{op}} \leq d$, with equality holding when G is fully connected, meaning that Σ has no zero entries. More generally, it can be verified (see

Exercise 6.9) that $\|\mathbf{A}\|_{\text{op}} \leq s$ whenever $\mathbf{\Sigma}$ has at most s non-zero entries per row, or equivalently when the graph G has maximum degree at most s . The following result provides a guarantee for the thresholded sample covariance matrix that involves the graph adjacency matrix \mathbf{A} defined by $\mathbf{\Sigma}$.

Theorem 6.5 (Thresholding-based covariance estimation). Let $\{x_i\}_{i=1}^n$ be an i.i.d. sequence of zero-mean random vectors with covariance matrix $\mathbf{\Sigma}$, and suppose that each x_{ij} is sub-Gaussian with parameter at most σ . If $n > \log d$, then for any $\delta > 0$, the thresholded sample covariance matrix $T_{\lambda_n}(\hat{\mathbf{\Sigma}})$ with $\lambda_n/\sigma^2 = 8\sqrt{\frac{\log d}{n}} + \delta$ satisfies

$$\mathbb{P}\left[\|T_{\lambda_n}(\hat{\mathbf{\Sigma}}) - \mathbf{\Sigma}\|_{\text{op}} \geq 2\|\mathbf{A}\|_{\text{op}}\lambda_n\right] \leq 8e^{-\frac{n}{16}\min\{\delta, \delta^2\}}. \quad (6.46)$$

Underlying the proof of Theorem 6.5 is the following (deterministic) result: for any choice of λ_n such that $\|\hat{\mathbf{\Sigma}} - \mathbf{\Sigma}\|_{\text{max}} \leq \lambda_n$, we are guaranteed that

$$\|T_{\lambda_n}(\hat{\mathbf{\Sigma}}) - \mathbf{\Sigma}\|_{\text{op}} \leq 2\|\mathbf{A}\|_{\text{op}}\lambda_n. \quad (6.47)$$

The proof of this intermediate claim is straightforward. First, for any index pair (j, ℓ) such that $\Sigma_{j\ell} = 0$, the bound $\|\hat{\mathbf{\Sigma}} - \mathbf{\Sigma}\|_{\text{max}} \leq \lambda_n$ guarantees that $|\hat{\Sigma}_{j\ell}| \leq \lambda_n$, and hence that $T_{\lambda_n}(\hat{\Sigma}_{j\ell}) = 0$ by definition of the thresholding operator. On the other hand, for any pair (j, ℓ) for which $\Sigma_{j\ell} \neq 0$, we have

$$|T_{\lambda_n}(\hat{\Sigma}_{j\ell}) - \Sigma_{j\ell}| \stackrel{(i)}{\leq} |T_{\lambda_n}(\hat{\Sigma}_{j\ell}) - \hat{\Sigma}_{j\ell}| + |\hat{\Sigma}_{j\ell} - \Sigma_{j\ell}| \stackrel{(ii)}{\leq} 2\lambda_n,$$

where step (i) follows from the triangle inequality, and step (ii) follows from the fact that $|T_{\lambda_n}(\hat{\Sigma}_{j\ell}) - \hat{\Sigma}_{j\ell}| \leq \lambda_n$, and a second application of the assumption $\|\hat{\mathbf{\Sigma}} - \mathbf{\Sigma}\|_{\text{max}} \leq \lambda_n$. Consequently, we have shown that the matrix $\mathbf{B} := |T_{\lambda_n}(\hat{\mathbf{\Sigma}}) - \mathbf{\Sigma}|$ satisfies the elementwise inequality $\mathbf{B} \leq 2\lambda_n \mathbf{A}$. Since both \mathbf{B} and \mathbf{A} have non-negative entries, we are guaranteed that $\|\mathbf{B}\|_{\text{op}} \leq 2\lambda_n \|\mathbf{A}\|_{\text{op}}$, and hence that $\|T_{\lambda_n}(\hat{\mathbf{\Sigma}}) - \mathbf{\Sigma}\|_{\text{op}} \leq 2\lambda_n \|\mathbf{A}\|_{\text{op}}$ as claimed. (See Exercise 6.10 for the details of these last steps.)

Theorem 6.5 has a number of interesting corollaries for particular classes of covariance matrices.

Corollary 6.2. Suppose that, in addition to the conditions of Theorem 6.5, the covariance matrix Σ has at most s non-zeros entries per row. Then with $\lambda_n/\sigma^2 = 8\sqrt{\frac{\log d}{n}} + \delta$ for some $\delta > 0$, we have

$$\mathbb{P}\left[\|T_{\lambda_n}(\hat{\Sigma}) - \Sigma\|_{\text{op}}/\sigma^2 \geq 16s\sqrt{\frac{\log d}{n}} + 2\delta\right] \leq 8e^{-\frac{n}{16}\min\{\delta, \delta^2\}}. \quad (6.48)$$

In order to establish these claims from Theorem 6.5, it suffices to show that $\|\mathbf{A}\|_{\text{op}} \leq s$. Since A has at most s ones per row (with the remaining entries equal to zero), this claim follows from the result of Exercise 6.9.

Example 6.11 (Sparsity and adjacency matrices). In certain ways, the bound (6.48) is more appealing than the bound (6.46), since it is based on a local quantity—namely, the maximum degree of the graph defined by the covariance matrix, as opposed to the spectral norm $\|\mathbf{A}\|_{\text{op}}$. In certain cases, these two bounds coincide. As an example, consider any graph with maximum degree $s - 1$ that contains a s -clique (i.e., a subset of s nodes that are all joined by edges). As we explore in Exercise 6.12, for any such graph, we have $\|\mathbf{A}\|_{\text{op}} = s - 1$, so that the two bounds are equivalent.

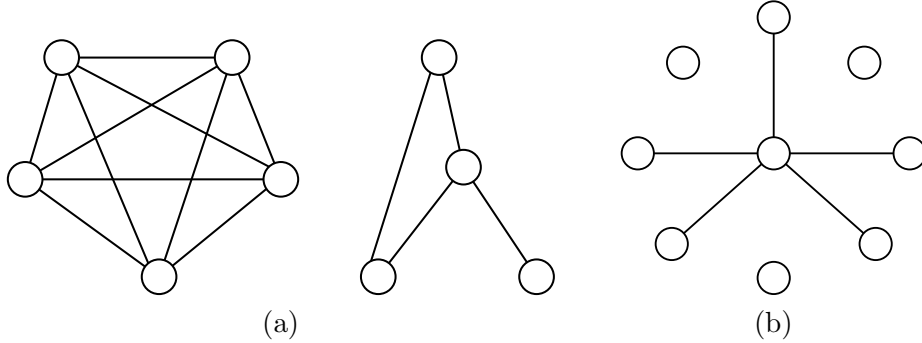


Figure 6-1. (a) An instance of a graph on $d = 9$ nodes containing a $s = 5$ clique. For this class of graphs, the bounds (6.46) and (6.48) coincide. (b) A hub-and-spoke graph on $d = 9$ nodes with maximum degree $s = 5$. For this class of graphs, the bounds differ by a factor of \sqrt{s} .

However, in general, the bound (6.46) can be substantially sharper than the bound (6.48). As an example, consider a hub-and-spoke graph, in which one central node known as the hub is connected to s of the remaining $d - 1$ nodes, as illustrated in Figure 6-1(a). For such a graph, we have $\|\mathbf{A}\|_{\text{op}} = 1 + \sqrt{s - 1}$, so that in this case, Theorem 6.5 guarantees that

$$\|T_{\lambda_n}(\hat{\Sigma}) - \Sigma\|_{\text{op}} \lesssim \sqrt{\frac{s \log d}{n}},$$

with high probability, a bound that is sharper by a factor of order \sqrt{s} compared to the

bound (6.48) from Corollary 6.2.

♣ 2469

We now turn to the proof of the remainder of Theorem 6.5. Based on the reasoning leading to equation (6.47), it suffices to establish a high probability bound on the elementwise infinity norm of the error matrix $\hat{\Delta} := \hat{\Sigma} - \Sigma$.

Lemma 6.4. Under the conditions of Theorem 6.5, we have

$$\mathbb{P}[\|\hat{\Delta}\|_{\max}/\sigma^2 \geq t] \leq 8e^{-\frac{n}{16}\min\{t, t^2\} + 2\log d} \quad \text{for all } t > 0. \quad (6.49)$$

Setting $t = \lambda_n/\sigma^2 = 8\sqrt{\frac{\log d}{n}} + \delta$ in the bound (6.49) yields

$$\mathbb{P}[\|\hat{\Delta}\|_{\max} \geq \lambda_n] \leq 8e^{-\frac{n}{16}\min\{\delta, \delta^2\}},$$

where we have used the fact that $n > \log d$ by assumption.

It remains to prove Lemma 6.4. Note that the rescaled vector x_i/σ is sub-Gaussian with parameter at most 1. Consequently, we may assume without loss of generality that $\sigma = 1$, and then rescale at the end. First considering a diagonal entry, the result of Exercise 6.11(a) guarantees that there are universal positive constants c_1, c_2 such that

$$\mathbb{P}[|\Delta_{jj}| \geq c_1\delta] \leq 2e^{-c_2n\delta^2} \quad \text{for all } \delta \in (0, 1). \quad (6.50)$$

Turning to the non-diagonal entries, for any $j \neq \ell$, we have

$$2\Delta_{j\ell} = \frac{2}{n} \sum_{i=1}^n x_{ij}x_{i\ell} - 2\Sigma_{j\ell} = \frac{1}{n} \sum_{i=1}^n (x_{ij} + x_{i\ell})^2 - (\Sigma_{jj} + \Sigma_{\ell\ell} - 2\Sigma_{j\ell}) - \Delta_{jj} - \Delta_{\ell\ell}.$$

Since x_{ij} and $x_{i\ell}$ are both zero-mean and sub-Gaussian with parameter σ , the sum $x_{ij} + x_{i\ell}$ is zero-mean and sub-Gaussian with parameter at most $2\sqrt{2}\sigma$ (see part (c) of Exercise 2.13). Consequently, there are universal constants c_2, c_3 such that for all $\delta \in (0, 1)$, we have

$$\mathbb{P}\left[\left|\frac{1}{n} \sum_{i=1}^n (x_{ij} + x_{i\ell})^2 - (\Sigma_{jj} + \Sigma_{\ell\ell} - 2\Sigma_{j\ell})\right| \geq c_3\delta\right] \leq 2e^{-c_2n\delta^2},$$

and hence $\mathbb{P}[|\hat{\Delta}_{j\ell}| \geq c_1\delta] \leq 6e^{-c_2n\delta^2}$. By combining this bound with the earlier inequality (6.50) and taking a union bound over all d^2 entries of the matrix, we obtain the stated claim (6.49).

■ 6.5.2 Approximate sparsity

Given a covariance matrix Σ with no entries that are exactly zero, the bounds of Theorem 6.5 are very poor. In particular, for a completely dense matrix, the associated adjacency matrix \mathbf{A} is simply the all-ones matrix, so that $\|\mathbf{A}\|_{\text{op}} = d$. Intuitively, one might expect that these bounds could be improved if Σ had a large number of non-zero entries, but many of them were “near-zero”.

Recall that one way in which to measure the sparsity of Σ is in terms of the maximum number of non-zero entries per row. A generalization of this idea is to measure the ℓ_q -“norm” of each row. More specifically, given a parameter $q \in [0, 1]$ and a radius R_q , we impose the constraint

$$\max_{j=1,\dots,d} \sum_{\ell=1}^d |\Sigma_{j\ell}|^q \leq R_q. \quad (6.51)$$

See Figure 7-1 in Chapter 7 for an illustration of these types of sets. In the special case $q = 0$, this constraint is equivalent to requiring that each row of Σ have at most R_0 non-zero entries. For intermediate values $q \in (0, 1]$, it allows for many non-zero entries but requires that their absolute magnitudes (if ordered from largest to smallest) drop off relatively quickly.

Theorem 6.6 (Covariance estimation under ℓ_q -sparsity). Suppose that the covariance matrix Σ satisfies the ℓ_q -sparsity constraint (6.51). Then for any λ_n such that $\|\hat{\Sigma} - \Sigma\|_{\max} \leq \lambda_n/2$, we are guaranteed that

$$\|T_{\lambda_n}(\hat{\Sigma}) - \Sigma\|_{\text{op}} \leq 2R_q \lambda_n^{1-q}. \quad (6.52)$$

Consequently, if the sample covariance is formed using i.i.d. samples $\{x_i\}_{i=1}^n$ that are zero-mean with sub-Gaussian parameter at most σ , then with $\lambda_n/\sigma^2 = 8\sqrt{\frac{\log d}{n}} + \delta$, we have

$$\mathbb{P}[\|T_{\lambda_n}(\hat{\Sigma}) - \Sigma\|_{\text{op}} \geq 2R_q \lambda_n^{1-q}] \leq 8e^{-\frac{n}{16} \min\{\delta, \delta^2\}} \quad \text{for all } \delta > 0. \quad (6.53)$$

Proof. It suffices to verify the deterministic claim, which is based on the assumption that $\|\hat{\Sigma} - \Sigma\|_{\max} \leq \lambda_n/2$. We next observe that the operator norm can be upper bounded as

$$\|T_{\lambda_n}(\hat{\Sigma}) - \Sigma\|_{\text{op}} \leq \max_{j=1,\dots,d} \sum_{\ell=1}^d |T_{\lambda_n}(\hat{\Sigma}_{j\ell}) - \Sigma_{j\ell}|$$

(See Exercise 6.9 for details of this claim.) We now fix an index $j \in \{1, 2, \dots, d\}$, and

define the index set $S_j(\lambda_n/2) = \{\ell \in \{1, \dots, d\} \mid |\Sigma_{j\ell}| > \lambda_n/2\}$. We then have

$$\sum_{\ell=1}^d |T_{\lambda_n}(\widehat{\Sigma}_{j\ell}) - \Sigma_{j\ell}| \leq |S_j(\lambda_n/2)|\lambda_n + \sum_{\ell \notin S_j(\lambda_n)} |T_{\lambda_n}(\widehat{\Sigma}_{j\ell}) - \Sigma_{j\ell}|. \quad (6.54)$$

For any index $\ell \notin S_j(\ell)$, we have $|\Sigma_{j\ell}| \leq \lambda_n/2$, and hence

$$|\widehat{\Sigma}_{j\ell}| \leq |\widehat{\Sigma}_{j\ell} - \Sigma_{j\ell}| + |\Sigma_{j\ell}| \leq \lambda_n/2 + \lambda_n/2 = \lambda_n.$$

By definition of the thresholding operator, we then have $T_{\lambda_n}(\widehat{\Sigma}_{j\ell}) = 0$ for all $\ell \notin S_j(\lambda_n/2)$, and hence from our earlier bound (6.54),

$$\sum_{\ell=1}^d |T_{\lambda_n}(\widehat{\Sigma}_{j\ell}) - \Sigma_{j\ell}| \leq |S_j(\lambda_n/2)|\lambda_n + \sum_{\ell \notin S_j(\lambda_n)} |\Sigma_{j\ell}|. \quad (6.55)$$

Now we have

$$\sum_{\ell \notin S_j(\lambda_n/2)} |\Sigma_{j\ell}| = \frac{\lambda_n}{2} \sum_{\ell \notin S_j(\lambda_n/2)} \frac{|\Sigma_{j\ell}|}{\lambda_n/2} \stackrel{(i)}{\leq} \frac{\lambda_n}{2} \sum_{\ell \notin S_j(\lambda_n/2)} \left(\frac{|\Sigma_{j\ell}|}{\lambda_n/2}\right)^q \stackrel{(ii)}{\leq} \lambda_n^{1-q} R_q,$$

where step (i) follows since $|\Sigma_{j\ell}| \leq \lambda/2$ for all $\ell \notin S_j(\lambda_n/2)$ and $q \in [0, 1]$, and step (ii) follows by the assumption (6.51). On the other hand, we have

$$R_q = \sum_{\ell=1}^d |\Sigma_{j\ell}|^q \geq |S_j(\lambda_n/2)| \left(\frac{\lambda_n}{2}\right)^q,$$

whence $|S_j(\lambda_n/2)| \leq R_q \lambda_n^{-q}$. Combining these ingredients with the inequality (6.55), we find that

$$\sum_{\ell=1}^d |T_{\lambda_n}(\widehat{\Sigma}_{j\ell}) - \Sigma_{j\ell}| \leq R_q \lambda_n^{1-q} + R_q \lambda_n^{1-q} = 2R_q \lambda_n^{1-q}.$$

Since this same argument holds for each index $j = 1, \dots, d$, the claim (6.52) follows. \square 2495

■ 6.6 Appendix: Proof of Theorem 6.1

2496

It remains to prove the lower bound (6.8) on the minimal singular value. In order to do so, we follow an argument similar to that used to upper bound the maximal singular value. We begin by lower bounding the expectation using a Gaussian comparison principle due to Gordon [Gor85]. By definition, the minimum singular value has the

variational representation $\gamma_{\min}(\mathbf{X}) = \min_{v' \in \mathbb{S}^{d-1}} \|\mathbf{X}v'\|_2$. Let us reformulate this representation slightly for later theoretical convenience. Recalling the shorthand notation $\bar{\gamma}_{\min} = \gamma_{\min}(\sqrt{\Sigma})$, let us define the radius $R = 1/\bar{\gamma}_{\min}$, and then consider the set

$$\mathcal{V}(R) := \{z \in \mathbb{S}^{d-1} \mid \|\sqrt{\Sigma}z\|_2 = 1, \|z\|_2 \leq R\}. \quad (6.56)$$

It suffices to show that for any $\delta > 0$, a lower bound of the form

$$\min_{z \in \mathcal{V}(R)} \frac{\|\mathbf{X}z\|_2}{\sqrt{n}} \geq 1 - \delta - R \sqrt{\frac{\text{trace}(\Sigma)}{n}} \quad (6.57)$$

holds with probability at least $1 - e^{-n\delta^2/2}$. Indeed, suppose that inequality (6.57) holds. Then for any $v' \in \mathbb{S}^{d-1}$, we can define the rescaled vector $z := \frac{v'}{\|\sqrt{\Sigma}v'\|_2}$. By construction, we have

$$\|\sqrt{\Sigma}z\|_2 = 1, \quad \text{and} \quad \|z\|_2 = \frac{1}{\|\sqrt{\Sigma}v'\|_2} \leq \frac{1}{\gamma_{\min}(\sqrt{\Sigma})} = R,$$

so that $z \in \mathcal{V}(R)$. We now observe that

$$\frac{\|\mathbf{X}v'\|_2}{\sqrt{n}} = \|\sqrt{\Sigma}v'\|_2 \frac{\|\mathbf{X}z\|_2}{\sqrt{n}} \geq \gamma_{\min}(\sqrt{\Sigma}) \min_{z \in \mathcal{V}(R)} \frac{\|\mathbf{X}z\|_2}{\sqrt{n}}.$$

Since this bound holds for all $v' \in \mathbb{S}^{d-1}$, we can take the infimum on the left-hand side, thereby obtaining

$$\begin{aligned} \min_{v' \in \mathbb{S}^{d-1}} \frac{\|\mathbf{X}v'\|_2}{\sqrt{n}} &\geq \bar{\gamma}_{\min} \min_{z \in \mathcal{V}(R)} \frac{\|\mathbf{X}z\|_2}{\sqrt{n}} \\ &\stackrel{(i)}{\geq} \bar{\gamma}_{\min} \left\{ 1 - R \sqrt{\frac{\text{trace}(\Sigma)}{n}} - \delta \right\} \\ &= (1 - \delta) \bar{\gamma}_{\min} - \sqrt{\frac{\text{trace}(\Sigma)}{n}}, \end{aligned}$$

where step (i) follows from the bound (6.57).

2497

It remains to prove the lower bound (6.57). We begin by showing concentration of the random variable $\min_{v \in \mathcal{V}(R)} \|\mathbf{X}v\|_2/\sqrt{n}$ around its expected value. Since the matrix $\mathbf{X} \in \mathbb{R}^{n \times d}$ has i.i.d. rows, each drawn from the $\mathcal{N}(0, \Sigma)$ distribution, we can write $\mathbf{X} = \mathbf{W}\sqrt{\Sigma}$, where the random matrix \mathbf{W} is standard Gaussian. Using the fact that $\|\sqrt{\Sigma}v\|_2 = 1$ for all $v \in \mathcal{V}(R)$, it follows that the function $\mathbf{W} \mapsto \min_{v \in \mathcal{V}(R)} \|\mathbf{W}\sqrt{\Sigma}v\|_2/\sqrt{n}$ is Lipschitz with parameter $L = 1/\sqrt{n}$. Applying Theo-

rem 2.4, we conclude that

$$\min_{v \in \mathcal{V}(R)} \frac{\|\mathbf{X}v\|_2}{\sqrt{n}} \geq \mathbb{E} \left[\min_{v \in \mathcal{V}(R)} \frac{\|\mathbf{X}v\|_2}{\sqrt{n}} \right] - \delta$$

with probability at least $1 - e^{-n\delta^2/2}$.

2498

Consequently, the proof will be complete if we can show that

$$\mathbb{E} \left[\min_{v \in \mathcal{V}(R)} \frac{\|\mathbf{X}v\|_2}{\sqrt{n}} \right] \geq 1 - R \sqrt{\frac{\text{trace}(\mathbf{\Sigma})}{n}}. \quad (6.58)$$

In order to do so, we make use of an extension of the Sudakov-Fernique inequality, known as Gordon's inequality, which we now state. Let $\{Z_{u,v}\}$ and $\{Y_{u,v}\}$ be a pair of zero-mean Gaussian processes indexed by a non-empty index set $\mathbb{T} = U \times V$. Suppose that

$$\mathbb{E}[(Z_{u,v} - Z_{\tilde{u},\tilde{v}})^2] \leq \mathbb{E}[(Y_{u,v} - Y_{\tilde{u},\tilde{v}})^2] \quad \text{for all pairs } (u,v) \text{ and } (\tilde{u},\tilde{v}) \in \mathbb{T}, \quad (6.59)$$

and moreover, this inequality holds with *equality* whenever $v = \tilde{v}$. Under these conditions, Gordon's inequality guarantees that

$$\mathbb{E} \left[\max_{v \in V} \min_{u \in U} Z_{u,v} \right] \leq \mathbb{E} \left[\max_{v \in V} \min_{u \in U} Y_{u,v} \right]. \quad (6.60)$$

In order to exploit this result, we first observe that

$$- \min_{z \in \mathcal{V}(R)} \|\mathbf{X}z\|_2 = \max_{z \in \mathcal{V}(R)} \{ -\|\mathbf{X}z\|_2 \} = \max_{z \in \mathcal{V}(R)} \min_{u \in \mathbb{S}^{n-1}} u^T \mathbf{X}z.$$

As before, if we introduce the standard Gaussian random matrix $\mathbf{W} \in \mathbb{R}^{n \times d}$, then for any $z \in \mathcal{V}(R)$, we can write $u^T \mathbf{X}z = u^T \mathbf{W}v$, where $v := \sqrt{\mathbf{\Sigma}}z$. Whenever $z \in \mathcal{V}(R)$, then the vector v must belong to the set $\mathcal{V}'(R) := \{v \in \mathbb{S}^{d-1} \mid \|\mathbf{\Sigma}^{-\frac{1}{2}}v\|_2 \leq R\}$, and we have shown that

$$\min_{v \in \mathcal{V}(R)} \|\mathbf{X}v\|_2 = \max_{v \in \mathcal{V}'(R)} \min_{u \in \mathbb{S}^{n-1}} \underbrace{u^T \mathbf{W}v}_{Z_{u,v}}.$$

Let $\theta = (u,v)$ and $\tilde{\theta} = (\tilde{u},\tilde{v})$ be any two members of the Cartesian product space $\mathbb{S}^{n-1} \times \mathcal{V}'(R)$. Since $\|u\|_2 = \|\tilde{u}\|_2 = \|v\|_2 = \|\tilde{v}\|_2 = 1$, following the same argument as in bounding the maximal singular value shows that

$$\rho_Z^2(\theta, \tilde{\theta}) \leq \|u - \tilde{u}\|_2^2 + \|v - \tilde{v}\|_2^2, \quad (6.61)$$

with equality holding when $v = \tilde{v}$. Consequently, if we define the Gaussian process

$Y_{u,v} := \langle g, u \rangle + \langle h, v \rangle$, where $g \in \mathbb{R}^n$ and $h \in \mathbb{R}^d$ are standard Gaussian vectors and mutually independent, then we have $\rho_Y^2(\theta, \tilde{\theta}) = \|u - \tilde{u}\|_2^2 + \|v - \tilde{v}\|_2^2$, so that the Sudakov-Fernique increment condition (6.59) holds. In addition, whenever $v = \tilde{v}$, holds in the upper bound (6.61), which guarantees that $\rho_Z((u, v), (\tilde{u}, v)) = \rho_Y((u, v), (\tilde{u}, v))$. Consequently, we may apply Gordon's inequality (6.60) to conclude that

$$\begin{aligned} \mathbb{E} \left[- \min_{z \in \mathcal{V}(R)} \|\mathbf{X}z\|_2 \right] &\leq \mathbb{E} \left[\max_{v \in \mathcal{V}'(R)} \min_{u \in S^{n-1}} Y_{u,v} \right] \\ &= \mathbb{E} \left[\min_{u \in S^{n-1}} \langle g, u \rangle \right] + \mathbb{E} \left[\max_{v \in \mathcal{V}'(R)} \langle h, v \rangle \right] \\ &\leq -\mathbb{E} \|g\|_2 + \mathbb{E} [\|\sqrt{\Sigma}h\|_2] R, \end{aligned}$$

where we have used the upper bound $|\langle h, v \rangle| = |\langle \sqrt{\Sigma}h, \Sigma^{-\frac{1}{2}}v \rangle| \leq \|\sqrt{\Sigma}h\|_2 R$, by definition of the set $\mathcal{V}'(R)$. 2499
2500

We now claim that

$$\frac{\mathbb{E} [\|\sqrt{\Sigma}h\|_2]}{\sqrt{\text{trace}(\Sigma)}} \leq \frac{\mathbb{E} [\|h\|_2]}{\sqrt{d}}. \quad (6.62)$$

Indeed, by the rotation invariance of the Gaussian distribution, we may assume that Σ is diagonal, with non-negative entries $\{\gamma_k\}_{k=1}^d$, and the claim is equivalent to showing that the function $F(\gamma) := \mathbb{E}[(\sum_{j=1}^d \gamma_j h_j^2)^{1/2}]$ achieves its maximum over the probability simplex at $\gamma_j = 1/d$. Since F is continuous and the probability simplex is compact, the maximum is achieved. Moreover, since F is concave and permutation invariant—meaning that $F(\gamma) = F(\Pi(\gamma))$ for all permutation matrices Π —the maximum must be achieved at $\gamma_j = 1/d$, which establishes the inequality (6.62). 2501
2502
2503
2504
2505
2506
2507

Recalling that $R = \frac{1}{\gamma_{\min}}$, we then have

$$\begin{aligned} -\mathbb{E}[\|g\|_2] + R \mathbb{E}[\|\sqrt{\Sigma}h\|_2] &\leq -\mathbb{E}[\|g\|_2] + \sqrt{\text{trace}(\Sigma)} \tilde{\gamma}_{\min} - \frac{\mathbb{E}[\|h\|_2]}{\sqrt{d}} \\ &= \underbrace{\{-\mathbb{E}[\|g\|_2] + \mathbb{E}[\|h\|_2]\}}_{T_1} + \underbrace{\left\{ \sqrt{\frac{\text{trace}(\Sigma)}{\tilde{\gamma}_{\min}^2 d}} - 1 \right\} \mathbb{E}[\|h\|_2]}_{T_2} \end{aligned}$$

By Jensen's inequality, we have $\mathbb{E}[\|h\|_2] \leq \sqrt{\mathbb{E}[\|h\|_2^2]} = \sqrt{d}$. Since $\frac{\text{trace}(\Sigma)}{\tilde{\gamma}_{\min}^2 d} \geq 1$, we conclude that $T_2 \leq \left\{ \sqrt{\frac{\text{trace}(\Sigma)}{\tilde{\gamma}_{\min}^2 d}} - 1 \right\} \sqrt{d}$. On the other hand, a direct calculation, using our assumption that $n \geq d$, shows that $T_1 \leq -\sqrt{n} + \sqrt{d}$. Combining the pieces, we

conclude that

$$\begin{aligned}\mathbb{E}\left[-\min_{z \in \mathcal{V}(R)} \|\mathbf{X}z\|_2\right] &\leq -\sqrt{n} + \sqrt{d} + \left\{\sqrt{\frac{\text{trace}(\mathbf{\Sigma})}{\tilde{\gamma}_{\min}^2 d}} - 1\right\}\sqrt{d} \\ &= -\sqrt{n} + \frac{\sqrt{\text{trace}(\mathbf{\Sigma})}}{\tilde{\gamma}_{\min}},\end{aligned}$$

which establishes the initial claim (6.57), thereby completing the proof.

2508

■ 6.7 Bibliographic details and background

2509

The books by Horn and Johnson [HJ85, HJ91] are standard references on linear algebra. 2510
A statement of Weyl's theorem and its corollaries can be found in Section 4.3 of the 2511
first volume [HJ85]. The monograph by Bhatia [Bha97] is more advanced in nature, 2512
taking a functional-analytic perspective, and includes discussion of Lidskii's theorem 2513
(see Section III.4). 2514

Some classical papers on asymptotic random matrix theory (RMT) include those 2515
by Wigner [Wig55], Pastur [Pas72], and Marcenko and Pastur [MP67]. Mehta [Meh91] 2516
provides an overview of asymptotic RMT, primarily from the physicist's perspective, 2517
whereas the book by Bai and Silverstein [BS10] takes a more statistical perspective. 2518
The lecture notes of Vershynin [Ver11] focus on the non-asymptotic aspects of random 2519
matrix theory, as partially covered here. 2520

Davidson and Szarek [DS01] describe the use of Sudakov-Fernique (Slepian) and 2521
Gordon inequalities in bounding expectations of random matrices; see also the earlier 2522
papers by Gordon [Gor85] and Szarek [Sza91]. The results in Davidson and Szarek [DS01] 2523
are for the special case of the standard Gaussian ensemble ($\mathbf{\Sigma} = \mathbf{I}_d$), but the underlying 2524
arguments are easily extended to the general case, as given here. 2525

The proof of Theorem 6.2 is based on the lecture notes of Vershynin [Ver11]. The 2526
underlying discretization argument is classical, used extensively in early work on ran- 2527
dom constructions in Banach space geometry (e.g., see the book by Pisier [Pis89] and 2528
references therein). Note that the discretization is a one-step version of the more so- 2529
phisticated chaining methods described in Chapter 5. 2530

Bounds for the expected operator norm of random matrices follow from non-commutative 2531
Bernstein inequalities, as derived initially by Rudelson [Rud99]. Ahlswede and Win- 2532
ter [AW02] developed techniques for matrix tail bounds based on controlling the ma- 2533
trix moment generating function, and exploiting the Golden-Thompson inequality. 2534
Other authors, among them Gross [Gro11], Recht [Rec11] and Oliveira [Oli10], de- 2535
veloped various extensions and refinements of the original Ahlswede-Winter approach. 2536
Tropp [Tro10] introduced the idea of controlling the matrix cumulant function directly, 2537
and developed the argument that underlies Lemma 6.3. Controlling the cumulant func- 2538

tion leads to tail bounds involving the variance parameter $\sigma^2 := \frac{1}{n} \|\sum_{i=1}^n \text{var}(\mathbf{Q}_i)\|_{\text{op}}$ as opposed to the quantity $\tilde{\sigma}^2 := \frac{1}{n} \sum_{i=1}^n \|\text{var}(\mathbf{Q}_i)\|_{\text{op}}$ that follows from the original Ahlswede-Winter argument. By the triangle inequality for the operator norm, we have $\sigma^2 \leq \tilde{\sigma}^2$, and the latter quantity can be substantially larger. Independent work by Oliveira [Oli10] also derived bounds involving the variance parameter σ^2 , using a technique that sharpened the original Ahlswede-Winter approach. Tropp [Tro10] also provides various extensions of the basic Bernstein bound, among them results for matrix martingales as opposed to the independent random matrices considered here. Mackey et al. [MJC⁺12] show how to derive matrix concentration bounds with sharp constants using the method of exchangeable pairs introduced by Chatterjee [Cha07].

For covariance estimation, Adamczak et al. [ALPTJ10] provide sharp results on the deviation $\|\hat{\Sigma} - \Sigma\|_{\text{op}}$ for distributions with sub-exponential tails. These results remove the superfluous logarithmic factor that arises from an application of Corollary 6.1 to a sub-exponential ensemble. For thresholded sample covariances, the first high-dimensional analyses were undertaken in independent work by Bickel and Levina [BL08a], and El Karoui [El 08]. Bickel and Levina studied the problem under sub-Gaussian tail conditions, and introduced the row-wise sparsity model, defined in terms of the maximum ℓ_q -“norm” taken over the rows. In contrast, El Karoui imposed a milder set of moment conditions, and measured sparsity in terms of the growth rates of path lengths in the graph; this approach is essentially equivalent to controlling the operator norm $\|\mathbf{A}\|_{\text{op}}$ of the adjacency matrix, as in Theorem 6.5. The star graph is an interesting example that illustrates the difference between the row-wise sparsity model, and the operator norm approach.

An alternative model for covariance matrices is a banded decay model, in which entries decay according to their distance from the diagonal. Bickel and Levina [BL08b] introduced this model in the covariance setting, and proposed a certain kind of tapering estimator. Cai et al. [CZZ10] analyzed the minimax optimal rates associated with this class of covariance matrices, and provided a modified estimator that achieves these optimal rates.

■ 6.8 Exercises

Exercise 6.1 (Bounds on eigenvalues). Given two symmetric matrices \mathbf{A} and \mathbf{B} , show directly, without citing any other theorems, that

$$|\gamma_{\max}(\mathbf{A}) - \gamma_{\max}(\mathbf{B})| \leq \|\mathbf{A} - \mathbf{B}\|_{\text{op}}, \quad \text{and} \quad |\gamma_{\min}(\mathbf{A}) - \gamma_{\min}(\mathbf{B})| \leq \|\mathbf{A} - \mathbf{B}\|_{\text{op}}.$$

Exercise 6.2 (Variance and positive semidefiniteness). Recall that the variance of symmetric random matrix \mathbf{Q} is given by $\text{var}(\mathbf{Q}) = \mathbb{E}[\mathbf{Q}^2] - (\mathbb{E}[\mathbf{Q}])^2$. Show that $\text{var}(\mathbf{Q}) \succeq 0$.

2572

Exercise 6.3 (Sub-Gaussian random matrices). Consider the random matrix $\mathbf{Q} = g\mathbf{B}$, where $g \in \mathbb{R}$ is a zero-mean σ sub-Gaussian variable.

2574

(a) Assume that g has a distribution symmetric around zero, and $\mathbf{B} \in \mathcal{S}^{d \times d}$ is a deterministic matrix. Show that \mathbf{Q} is sub-Gaussian with matrix parameter $\mathbf{V} = c^2 \sigma^2 \mathbf{B}^2$, for some universal constant c .

2577

(b) Now assume that $\mathbf{B} \in \mathcal{S}^{d \times d}$ is random and independent of g , with $\|\mathbf{B}\|_{\text{op}} \leq b$ almost surely. Now show that \mathbf{Q} is sub-Gaussian with matrix parameter $\mathbf{V} = c^2 \sigma^2 b^2 \mathbf{I}_d$.

2578

2579

Exercise 6.4 (Sub-Gaussian matrices and mean bounds). Consider a sequence of independent, zero-mean random matrices $\{\mathbf{Q}_i\}_{i=1}^n$ in $\mathcal{S}^{d \times d}$, each sub-Gaussian with matrix parameter \mathbf{V}_i . In this exercise, we provide bounds on the expected value of eigenvalues and operator norm of $\mathbf{S}_n = \frac{1}{n} \sum_{i=1}^n \mathbf{Q}_i$.

2580

2581

2582

2583

(a) Show that $\mathbb{E}[\gamma_{\max}(\mathbf{S}_n)] \leq \sqrt{\frac{2\sigma^2 \log d}{n}}$, where $\sigma^2 = \left\| \frac{1}{n} \sum_{i=1}^n \mathbf{V}_i \right\|_{\text{op}}$.

2584

(Hint: Start by showing that $\mathbb{E}[e^{\lambda \gamma_{\max}(\mathbf{S}_n)}] \leq de^{\frac{\lambda^2 \sigma^2}{2n}}$.)

2585

(b) Show that

$$\mathbb{E}\left[\left\| \frac{1}{n} \sum_{i=1}^n \mathbf{Q}_i \right\|_{\text{op}}\right] \leq \sqrt{\frac{2\sigma^2 \log(2d)}{n}}, \quad (6.63)$$

Exercise 6.5 (Bounded matrices and Bernstein condition). Let $\mathbf{Q} \in \mathcal{S}^{d \times d}$ be an arbitrary symmetric matrix.

2586

2587

(a) Show that the bound $\|\mathbf{Q}\|_{\text{op}} \leq b$ implies that $\mathbf{Q}^{j-2} \preceq b^{j-2} \mathbf{I}_d$.

2588

(b) Show that the positive semidefinite order is preserved under left-right multiplication, meaning that if $\mathbf{A} \preceq \mathbf{B}$, then we also have $\mathbf{Q}\mathbf{A}\mathbf{Q} \preceq \mathbf{Q}\mathbf{B}\mathbf{Q}$ for any matrix $\mathbf{Q} \in \mathcal{S}^{d \times d}$.

2589

2590

2591

(c) Use parts (a) and (b) to prove the inequality (6.26).

2592

Exercise 6.6 (Tail bounds for non-symmetric matrices). In this exercise, we prove that a version of the tail bound (6.37) holds for general independent zero-mean matrices $\{\mathbf{A}_i\}_{i=1}^n$, as long as we adopt new definition (6.39) of σ^2 .

2593

2594

2595

- (a) Given a general matrix $\mathbf{A}_i \in \mathbb{R}^{d_1 \times d_2}$, define a symmetric matrix of dimension $(d_1 + d_2)$ via

$$\mathbf{Q}_i := \begin{bmatrix} \mathbf{0}_{d_1 \times d_2} & \mathbf{A}_i \\ \mathbf{A}_i^T & \mathbf{0}_{d_2 \times d_1} \end{bmatrix}$$

Prove that $\|\mathbf{Q}_i\|_{\text{op}} = \|\mathbf{A}_i\|_{\text{op}}$.

2596

- (b) Prove that $\|\frac{1}{n} \sum_{i=1}^n \text{var}(\mathbf{Q}_i)\|_{\text{op}} \leq \sigma^2$ where σ^2 is defined in equation (6.39).

2597

- (c) Conclude that

$$\mathbb{P}\left[\left\|\frac{1}{n} \sum_{i=1}^n \mathbf{A}_i\right\|_{\text{op}} \geq \delta\right] \leq 2(d_1 + d_2)e^{-\frac{n\delta^2}{2(\sigma^2 + b\delta)}}. \quad (6.64)$$

Exercise 6.7 (Unbounded matrices and Bernstein bounds). Consider an independent sequence of random matrices $\{\mathbf{A}_i\}_{i=1}^n$ in $\mathbb{R}^{d_1 \times d_2}$, each of the form $\mathbf{A}_i = g_i \mathbf{B}_i$ where $g_i \in \mathbb{R}$ is a zero-mean scalar random variable, and \mathbf{B}_i is an independent random matrix. Suppose that $\mathbb{E}[g_i^j] \leq \frac{j!}{2} b_1^{j-2} \sigma^2$ for $j = 2, 3, \dots$, and that $\|\mathbf{B}_i\|_{\text{op}} \leq b_2$ almost surely.

2598

2599

2600

2601

- (a) For any $\delta > 0$, show that

$$\mathbb{P}\left[\left\|\frac{1}{n} \sum_{i=1}^n \mathbf{A}_i\right\|_{\text{op}} \geq \delta\right] \leq (d_1 + d_2)e^{-\frac{n\delta^2}{2(\sigma^2 b_2^2 + b_1 b_2 \delta)}}.$$

(Hint: The result of Exercise 6.6(a) could be useful.)

2602

- (b) Show that

$$\mathbb{E}\left[\left\|\frac{1}{n} \sum_{i=1}^n \mathbf{A}_i\right\|_{\text{op}}\right] \leq \frac{2\sigma b_2}{\sqrt{n}} \left\{ \sqrt{\log(d_1 + d_2)} + \sqrt{\pi} \right\} + \frac{4b_1 b_2}{n} \left\{ \log(d_1 + d_2) + 1 \right\}.$$

(Hint: The result of Exercise 2.8 could be useful.)

2603

Exercise 6.8 (Random packings). Prove that there exists a subset $\mathcal{P} = \{\theta^1, \dots, \theta^M\}$ of the sphere \mathbb{S}^{d-1} such that

2604

2605

- (a) The set \mathcal{P} forms a $1/2$ -packing.

2606

- (b) The set \mathcal{P} has cardinality $M \geq e^{c_0 d}$ for some universal constant c_0 .

2607

- (c) The inequality $\left\|\frac{1}{M} \sum_{j=1}^M (\theta^j \otimes \theta^j)\right\|_{\text{op}} \leq \frac{2}{d}$ holds.

2608

(Note: You may assume that d is larger than some universal constant so as to avoid annoying subcases.)

Exercise 6.9 (Relations between matrix operator norms). For a matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$ and $q \in [1, \infty]$, the $(\ell_q \rightarrow \ell_q)$ -operator norms are given by

$$\|\mathbf{A}\|_q = \sup_{\|x\|_q=1} \|\mathbf{A}x\|_q.$$

- (a) Derive explicit expressions for the operator norms $\|\mathbf{A}\|_2$, $\|\mathbf{A}\|_1$ and $\|\mathbf{A}\|_\infty$ in terms of elements and/or singular values of \mathbf{A} .
- (b) Prove that $\|\mathbf{A}\mathbf{B}\|_q \leq \|\mathbf{A}\|_q \|\mathbf{B}\|_q$ for any size-compatible matrices \mathbf{A} and \mathbf{B} .
- (c) Prove that $\|\mathbf{A}\|_2^2 \leq \|\mathbf{A}\|_1 \|\mathbf{A}\|_\infty$. What happens when \mathbf{A} is symmetric?

Exercise 6.10 (Non-negative matrices and operator norms). Given two d -dimensional symmetric matrices \mathbf{A} and \mathbf{B} , suppose that $0 \preceq \mathbf{A} \preceq \mathbf{B}$ in an elementwise sense (i.e., $0 \leq A_{j\ell} \leq B_{j\ell}$ for all $j, \ell = 1, \dots, d$.)

- (a) Show that $0 \preceq \mathbf{A}^m \preceq \mathbf{B}^m$ for all integers $m = 1, 2, \dots$.
- (b) Use part (a) to show that $\|\mathbf{A}\|_{\text{op}} \leq \|\mathbf{B}\|_{\text{op}}$.
- (c) Use a similar argument to show that $\|\mathbf{C}\|_{\text{op}} \leq \|\mathbf{C}\|_{\text{op}}$ for any symmetric matrix \mathbf{C} .

Exercise 6.11 (Estimation of diagonal covariances). Let $\{x_i\}_{i=1}^n$ be an i.i.d. sequence of d -dimensional vectors, drawn from a zero-mean distribution with diagonal covariance matrix $\Sigma = \mathbf{D}$. Consider the estimate $\hat{\mathbf{D}} = \text{diag}(\hat{\Sigma})$, where $\hat{\Sigma}$ is the usual sample covariance matrix.

- (a) When each vector x_i is sub-Gaussian with parameter at most σ , show that there are universal positive constants c_j such that

$$\mathbb{P}\left[\|\hat{\mathbf{D}} - \mathbf{D}\|_{\text{op}}/\sigma^2 \geq c_0 \sqrt{\frac{\log d}{n}} + \delta\right] \leq c_1 e^{-c_2 n \min\{\delta, \delta^2\}}. \quad \text{for all } \delta > 0.$$

- (b) Instead of sub-Gaussianity, suppose that for some even integer $m \geq 2$, there is a universal constant K_m such that

$$\underbrace{\mathbb{E}[(x_{ij}^2 - \Sigma_{jj})^m]}_{\|x_{ij}^2 - \Sigma_{jj}\|_m^m} \leq K_m,$$

for each $i = 1, \dots, n$ and $j = 1, \dots, d$. Show that

$$\mathbb{P}[\|\hat{\mathbf{D}} - \mathbf{D}\|_{\text{op}} \geq 4\delta \sqrt{\frac{d^{2/m}}{n}}] \leq K'_m \left(\frac{1}{2\delta}\right)^m \quad \text{for all } \delta > 0,$$

where K'_m is another universal constant. *Hint:* You may find Rosenthal's inequality useful: given zero-mean independent random variables Z_i such that $\|Z_i\|_m < +\infty$, there is a universal constant C_m such that

$$\left\| \sum_{i=1}^n Z_i \right\|_m \leq C_m \left\{ \left(\sum_{i=1}^n \mathbb{E}[Z_i^2] \right)^{1/2} + \left(\sum_{i=1}^n \mathbb{E}[|Z_i|^m] \right)^{1/m} \right\}.$$

Exercise 6.12 (Graphs and adjacency matrices). Let G be a graph with maximum degree $s - 1$ that contains a s -clique. Letting \mathbf{A} denote its adjacency matrix, show that $\|\mathbf{A}\|_{\text{op}} = s - 1$.

