# Chapter 2

# Bounds for all eigenvalues of sums of Hermitian random matrices

## 2.1 Introduction

The classical tools of nonasymptotic random matrix theory can sometimes give quite sharp estimates of the extreme eigenvalues of a Hermitian random matrix, but they are not readily adapted to the study of the interior eigenvalues. This is because, while the extremal eigenvalues are the maxima and minima of a random process, more delicate and challenging minimax problems must be solved to obtain the interior eigenvalues.

This chapter introduces a simple method, based upon the variational characterization of eigenvalues, that parlays bounds on the extreme eigenvalues of sums of random Hermitian matrices into bounds that apply to all the eigenvalues[1]. This technique extends the matrix Laplace transform method detailed in [Tro12]. We combine these ideas to extend several of the inequalities in [Tro12] to address the fluctuations of interior eigenvalues. Specifically, we provide eigenvalue analogs of the classical multiplicative Chernoff bounds and Bennett and Bernstein inequalities.

In this technique, the delicacy of the minimax problems which implicitly define the eigenval-

---

[1]The content of this chapter is adapted from the technical report [GT09] co-authored with Joel Tropp.

ues of Hermitian matrices is encapsulated in terms that reflect the fluctuations of the summands in the appropriate eigenspaces. In particular, we see that the fluctuations of the $k$th eigenvalue of the sum above and below the $k$th eigenvalue of the expected sum are controlled by two different quantities. This satisfies intuition: for instance, given samples from a nondegenerate stationary random process with finite covariance matrix, one expects that the smallest eigenvalue of the sample covariance matrix is more likely to be an underestimate of the smallest eigenvalue of the covariance matrix than it is to be an overestimate.

We provide two illustrative applications of our eigenvalue tail bounds: Theorem 2.14 quantifies the behavior of the singular values of matrices obtained by sampling columns from a short, fat matrix; and Theorem 2.15 quantifies the convergence of the eigenvalues of Wishart matrices.

## 2.2 Notation

We define $\mathbb{M}_{\mathrm{sa}}^n$ to be the set of Hermitian matrices with dimension $n$. We often compare Hermitian matrices using the semidefinite ordering. In this ordering, $\mathbf{A}$ is greater than or equal to $\mathbf{B}$, written $\mathbf{A} \succeq \mathbf{B}$ or $\mathbf{B} \preceq \mathbf{A}$, when $\mathbf{A} - \mathbf{B}$ is positive semidefinite.

The eigenvalues of a matrix $\mathbf{A}$ in $\mathbb{M}_{\mathrm{sa}}^n$ are arranged in weakly decreasing order: $\lambda_{\max}(\mathbf{A}) = \lambda_1(\mathbf{A}) \geq \lambda_2(\mathbf{A}) \geq \cdots \geq \lambda_n(\mathbf{A}) = \lambda_{\min}(\mathbf{A})$. Likewise, the singular values of a rectangular matrix $\mathbf{A}$ with rank $\rho$ are ordered $\sigma_{\max}(\mathbf{A}) = \sigma_1(\mathbf{A}) \geq \sigma_2(\mathbf{A}) \geq \cdots \geq \sigma_\rho(\mathbf{A}) = \sigma_{\min}(\mathbf{A})$. The spectral norm of a matrix $\mathbf{B}$ is written $\|\mathbf{B}\|_2$.

## 2.3 The Courant–Fisher Theorem

In this chapter, we work over the complex field $\mathbb{C}$. One of our central tools is the variational characterization of the eigenvalues of a Hermitian matrix given by the Courant–Fischer Theorem. For integers $d$ and $n$ satisfying $1 \le d \le n$, the complex Stiefel manifold

$$\mathbb{V}_d^n = \{\mathbf{V} \in \mathbb{C}^{n \times d} : \mathbf{V}^* \mathbf{V} = \mathbf{I}\}$$

is the collection of orthonormal bases for the $d$-dimensional subspaces of $\mathbb{C}^n$, or, equivalently, the collection of all isometric embeddings of $\mathbb{C}^d$ into $\mathbb{C}^n$. Let $\mathbf{A}$ be a Hermitian matrix with dimension $n$, and let $\mathbf{V} \in \mathbb{V}_d^n$ be an orthonormal basis for a subspace of $\mathbb{C}^n$. Then the matrix $\mathbf{V}^* \mathbf{A} \mathbf{V}$ can be interpreted as the compression of $\mathbf{A}$ to the space spanned by $\mathbf{V}$.

**Proposition 2.1** (Courant–Fischer ([HJ85, Theorem 4.2.11])). *Let $\mathbf{A}$ be a Hermitian matrix with dimension $n$. Then*

$$\lambda_k(\mathbf{A}) = \min_{\mathbf{V} \in \mathbb{V}_{n-k+1}^n} \lambda_{\max}\left(\mathbf{V}^* \mathbf{A} \mathbf{V}\right) \quad and \tag{2.3.1}$$

$$\lambda_k(\mathbf{A}) = \max_{\mathbf{V} \in \mathbb{V}_k^n} \lambda_{\min}\left(\mathbf{V}^* \mathbf{A} \mathbf{V}\right). \tag{2.3.2}$$

*A matrix $\mathbf{V}_- \in \mathbb{V}_k^n$ achieves equality in (2.3.2) if and only if its columns span a top $k$-dimensional invariant subspace of $\mathbf{A}$. Likewise, a matrix $\mathbf{V}_+ \in \mathbb{V}_{n-k+1}^n$ achieves equality in (2.3.1) if and only if its columns span a bottom $(n - k + 1)$-dimensional invariant subspace of $\mathbf{A}$.*

The $\pm$ subscripts in Proposition 2.1 are chosen to reflect the fact that $\lambda_k(\mathbf{A})$ is the *minimum* eigenvalue of $\mathbf{V}_-^* \mathbf{A} \mathbf{V}_-$ and the *maximum* eigenvalue of $\mathbf{V}_+^* \mathbf{A} \mathbf{V}_+$. As a consequence of Proposition 2.1, when $\mathbf{A}$ is Hermitian, $\lambda_k(-\mathbf{A}) = -\lambda_{n-k+1}(\mathbf{A})$. This fact allows us to use the same

techniques we develop for bounding the eigenvalues from above to bound them from below.

## 2.4 Tail bounds for interior eigenvalues

In this section we develop a generic bound on the tail probabilities of eigenvalues of sums of independent, random, Hermitian matrices. We establish this bound by supplementing the matrix Laplace transform methodology of [Tro12] with Proposition 2.1 and a result, due to Lieb and Seiringer [LS05], on the concavity of a certain trace function on the cone of positive-definite matrices.

First we observe that the Courant–Fischer Theorem allows us to relate the behavior of the $k$th eigenvalue of a matrix to the behavior of the largest eigenvalue of an appropriate compression of the matrix.

**Theorem 2.2.** *Let* **Y** *be a random Hermitian matrix with dimension $n$, and let $k \leq n$ be an integer. Then, for all $t \in \mathbb{R}$,*

$$\mathbb{P}\left\{\lambda_k(\mathbf{Y}) \geq t\right\} \leq \inf_{\theta > 0} \min_{\mathbf{V} \in \mathbb{V}_{n-k+1}^n} \left\{ e^{-\theta t} \cdot \mathbb{E}\operatorname{tr} e^{\theta \mathbf{V}^* \mathbf{Y} \mathbf{V}} \right\}. \tag{2.4.1}$$

*Proof.* Let $\theta$ be a fixed positive number. Then

$$\mathbb{P}\left\{\lambda_k(\mathbf{Y}) \geq t\right\} = \mathbb{P}\left\{\lambda_k(\theta \mathbf{Y}) \geq \theta t\right\} = \mathbb{P}\left\{ e^{\lambda_k(\theta \mathbf{Y})} \geq e^{\theta t}\right\}$$

$$\leq e^{-\theta t} \cdot \mathbb{E}e^{\lambda_k(\theta \mathbf{Y})} = e^{-\theta t} \cdot \mathbb{E}\exp\left\{ \min_{\mathbf{V} \in \mathbb{V}_{n-k+1}^n} \lambda_{\max}\left(\theta \mathbf{V}^* \mathbf{Y} \mathbf{V}\right) \right\}.$$

The first identity follows from the positive homogeneity of eigenvalue maps and the second from the monotonicity of the scalar exponential function. The final two relations are Markov's

inequality and (2.3.1).

To continue, we need to bound the expectation. Use monotonicity to interchange the order of the exponential and the minimum; then apply the spectral mapping theorem to see that

$$\mathbb{E}\exp\left\{\min_{\mathbf{V}\in\mathbb{V}_{n-k+1}^n}\lambda_{\max}\left(\theta\mathbf{V}^*\mathbf{Y}\mathbf{V}\right)\right\}=\mathbb{E}\min_{\mathbf{V}\in\mathbb{V}_{n-k+1}^n}\lambda_{\max}\left(\exp(\theta\mathbf{V}^*\mathbf{Y}\mathbf{V})\right)$$

$$\leq\min_{\mathbf{V}\in\mathbb{V}_{n-k+1}^n}\mathbb{E}\lambda_{\max}\left(\exp(\theta\mathbf{V}^*\mathbf{Y}\mathbf{V})\right)$$

$$\leq\min_{\mathbf{V}\in\mathbb{V}_{n-k+1}^n}\mathbb{E}\operatorname{tr}\exp(\theta\mathbf{V}^*\mathbf{Y}\mathbf{V}).$$

The first inequality is Jensen's. The second inequality follows because the exponential of a Hermitian matrix is positive definite, so its largest eigenvalue is smaller than its trace.

Combine these observations and take the infimum over all positive $\theta$ to complete the argument. □

In most cases it is prohibitively difficult to compute the quantity $\mathbb{E}\operatorname{tr}e^{\theta\mathbf{V}^*\mathbf{Y}\mathbf{V}}$ exactly. The main contribution of [Tro12] is a bound on this quantity, when $\mathbf{V}=\mathbf{I}$, in terms of the cumulant generating functions of the summands. The main tool in the proof is a classical result due to Lieb [Lie73, Thm. 6] that establishes the concavity of the function

$$\mathbf{A}\longmapsto\operatorname{tr}\exp\left(\mathbf{H}+\log(\mathbf{A})\right) \tag{2.4.2}$$

on the positive-definite cone, where $\mathbf{H}$ is Hermitian.

We are interested in the case where $\mathbf{V}\neq\mathbf{I}$ and the matrix $\mathbf{Y}$ in Theorem 2.2 can be expressed as a sum of independent random matrices. In this case, we use the following result to develop the right-hand side of the Laplace transform bound (2.4.1).

**Theorem 2.3.** *Consider a finite sequence $\{\mathbf{X}_j\}$ of independent, random, Hermitian matrices with*

*dimension n and a sequence $\{\mathbf{A}_j\}$ of fixed Hermitian matrices with dimension n that satisfy the relations*

$$\mathbb{E}e^{\mathbf{X}_j} \preceq e^{\mathbf{A}_j}. \qquad (2.4.3)$$

*Let $\mathbf{V} \in \mathbb{V}_k^n$ be an isometric embedding of $\mathbb{C}^k$ into $\mathbb{C}^n$ for some $k \leq n$. Then*

$$\mathbb{E}\operatorname{tr}\exp\left\{\sum_j \mathbf{V}^*\mathbf{X}_j\mathbf{V}\right\} \leq \operatorname{tr}\exp\left\{\sum_j \mathbf{V}^*\mathbf{A}_j\mathbf{V}\right\}. \qquad (2.4.4)$$

*In particular,*

$$\mathbb{E}\operatorname{tr}\exp\left\{\sum_j \mathbf{X}_j\right\} \leq \operatorname{tr}\exp\left\{\sum_j \mathbf{A}_j\right\}. \qquad (2.4.5)$$

Theorem 2.3 is an extension of Lemma 3.4 of [Tro12], which establishes the special case (2.4.5). The proof depends upon a result due to Lieb and Seiringer [LS05, Thm. 3] that extends Lieb's earlier result (2.4.2) by showing that the functional remains concave when the $\log(\mathbf{A})$ term is compressed.

**Proposition 2.4** (Lieb–Seiringer 2005). *Let $\mathbf{H}$ be a Hermitian matrix with dimension $k$. Let $\mathbf{V} \in \mathbb{V}_k^n$ be an isometric embedding of $\mathbb{C}^k$ into $\mathbb{C}^n$ for some $k \leq n$. Then the function*

$$\mathbf{A} \longmapsto \operatorname{tr}\exp\left(\mathbf{H} + \mathbf{V}^*(\log\mathbf{A})\mathbf{V}\right)$$

*is concave on the cone of positive-definite matrices in $\mathbb{M}_{\mathrm{sa}}^n$.*

*Proof of Theorem 2.3.* First, note that (2.4.3) and the operator monotonicity of the matrix logarithm yield the following inequality for each $k$:

$$\log\mathbb{E}e^{\mathbf{X}_k} \preceq \mathbf{A}_k. \qquad (2.4.6)$$

Let $\mathbb{E}_k$ denote expectation conditioned on the first $k$ summands, $\mathbf{X}_1$ through $\mathbf{X}_k$. Then

$$
\begin{aligned}
\mathbb{E}\operatorname{tr}\exp\left(\sum_{j\le\ell}\mathbf{V}^*\mathbf{X}_j\mathbf{V}\right) &= \mathbb{E}\mathbb{E}_1\cdots\mathbb{E}_{\ell-1}\operatorname{tr}\exp\left(\sum_{j\le\ell-1}\mathbf{V}^*\mathbf{X}_j\mathbf{V}+\mathbf{V}^*\left(\log e^{\mathbf{X}_\ell}\right)\mathbf{V}\right) \\
&\le \mathbb{E}\mathbb{E}_1\cdots\mathbb{E}_{\ell-2}\operatorname{tr}\exp\left(\sum_{j\le\ell-1}\mathbf{V}^*\mathbf{X}_j\mathbf{V}+\mathbf{V}^*\left(\log\mathbb{E}e^{\mathbf{X}_\ell}\right)\mathbf{V}\right) \\
&\le \mathbb{E}\mathbb{E}_1\cdots\mathbb{E}_{\ell-2}\operatorname{tr}\exp\left(\sum_{j\le\ell-1}\mathbf{V}^*\mathbf{X}_j\mathbf{V}+\mathbf{V}^*\left(\log e^{\mathbf{A}_\ell}\right)\mathbf{V}\right) \\
&= \mathbb{E}\mathbb{E}_1\cdots\mathbb{E}_{\ell-2}\operatorname{tr}\exp\left(\sum_{j\le\ell-1}\mathbf{V}^*\mathbf{X}_j\mathbf{V}+\mathbf{V}^*\mathbf{A}_\ell\mathbf{V}\right).
\end{aligned}
$$

The first inequality follows from Proposition 2.4 and Jensen's inequality, and the second depends on (2.4.6) and the monotonicity of the trace exponential. Iterate this argument to complete the proof. $\qquad\square$

Our main result follows from combining Theorem 2.2 and Theorem 2.3.

**Theorem 2.5** (Minimax Laplace Transform). *Consider a finite sequence $\{\mathbf{X}_j\}$ of independent, random, Hermitian matrices with dimension $n$, and let $k\le n$ be an integer.*

*(i) Let $\{\mathbf{A}_j\}$ be a sequence of Hermitian matrices that satisfy the semidefinite relations*

$$
\mathbb{E}e^{\theta\mathbf{X}_j}\preceq e^{g(\theta)\mathbf{A}_j}
$$

*where $g:(0,\infty)\to[0,\infty)$. Then, for all $t\in\mathbb{R}$,*

$$
\mathbb{P}\left\{\lambda_k\left(\sum_j\mathbf{X}_j\right)\ge t\right\}\le\inf_{\theta>0}\min_{\mathbf{V}\in\mathbb{V}_{n-k+1}^n}\left[e^{-\theta t}\cdot\operatorname{tr}\exp\left\{g(\theta)\sum_j\mathbf{V}^*\mathbf{A}_j\mathbf{V}\right\}\right].
$$

*(ii) Let $\{\mathbf{A}_j : \mathbb{V}^n_{n-k+1} \to \mathbb{M}^n_{\mathrm{sa}}\}$ be a sequence of functions that satisfy the semidefinite relations*

$$\mathbb{E}\,\mathrm{e}^{\theta \mathbf{V}^* \mathbf{X}_j \mathbf{V}} \preceq \mathrm{e}^{g(\theta)\mathbf{A}_j(\mathbf{V})}$$

*for all $\mathbf{V} \in \mathbb{V}^n_{n-k+1}$, where $g : (0, \infty) \to [0, \infty)$. Then, for all $t \in \mathbb{R}$,*

$$\mathbb{P}\left\{\lambda_k\left(\sum\nolimits_j \mathbf{X}_j\right) \geq t\right\} \leq \inf_{\theta > 0} \min_{\mathbf{V} \in \mathbb{V}^n_{n-k+1}} \left[\mathrm{e}^{-\theta t} \cdot \mathrm{tr}\exp\left\{g(\theta)\sum\nolimits_j \mathbf{A}_j(\mathbf{V})\right\}\right].$$

The first bound in Theorem 2.5 requires less detailed information on how compression affects the summands but correspondingly does not give as sharp results as the second. For most cases we consider, we use the second inequality because it is straightforward to obtain semidefinite bounds for the compressed summands. The exception occurs in the proof of the subexponential Bernstein inequality (Theorem 2.12 in Section 2.6); here we use the first bound, because in this case there are no nontrivial semidefinite bounds for the compressed summands.

In the following two sections, we use the minimax Laplace transform method to derive Chernoff and Bernstein inequalities for the interior eigenvalues of a sum of independent random matrices. Tail bounds for the eigenvalues of matrix Rademacher and Gaussian series, eigenvalue Hoeffding, and matrix martingale eigenvalue tail bounds can all be derived in a similar manner; see [Tro12] for the details of the arguments leading to such tail bounds for the maximum eigenvalue.

## 2.5  Chernoff bounds

Classical Chernoff bounds establish that the tails of a sum of independent nonnegative random variables decay subexponentially. [Tro12] develops Chernoff bounds for the maximum and minimum eigenvalues of a sum of independent positive semidefinite matrices. We extend this analysis to study the interior eigenvalues.

Intuitively, the eigenvalue tail bounds should depend on how concentrated the summands are; e.g., the maximum eigenvalue of a sum of operators whose ranges are aligned is likely to vary more than that of a sum of operators whose ranges are orthogonal. To measure how much a finite sequence of random summands $\{\mathbf{X}_j\}$ concentrates in a given subspace, we define a function $\Psi : \bigcup_{1 \leq k \leq n} \mathbb{V}_k^n \to \mathbb{R}$ that satisfies

$$\max_j \lambda_{\max} \left( \mathbf{V}^* \mathbf{X}_j \mathbf{V} \right) \leq \Psi(\mathbf{V}) \qquad \text{almost surely for each } \mathbf{V} \in \bigcup_{1 \leq k \leq n} \mathbb{V}_k^n. \tag{2.5.1}$$

The sequence $\{\mathbf{X}_j\}$ associated with $\Psi$ will always be clear from context. We have the following result.

**Theorem 2.6** (Eigenvalue Chernoff Bounds)**.** *Consider a finite sequence $\{\mathbf{X}_j\}$ of independent, random, positive-semidefinite matrices with dimension n. Given an integer $k \leq n$, define*

$$\mu_k = \lambda_k \left( \sum\nolimits_j \mathbb{E} \mathbf{X}_j \right),$$

*and let $\mathbf{V}_+ \in \mathbb{V}_{n-k+1}^n$ and $\mathbf{V}_- \in \mathbb{V}_k^n$ be isometric embeddings that satisfy*

$$\mu_k = \lambda_{\max} \left( \sum\nolimits_j \mathbf{V}_+^* (\mathbb{E} \mathbf{X}_j) \mathbf{V}_+ \right) = \lambda_{\min} \left( \sum\nolimits_j \mathbf{V}_-^* (\mathbb{E} \mathbf{X}_j) \mathbf{V}_- \right).$$

*Then*

$$\mathbb{P}\left\{\lambda_k\left(\sum_j \mathbf{X}_j\right) \geq (1+\delta)\mu_k\right\} \leq (n-k+1)\cdot\left[\frac{e^\delta}{(1+\delta)^{1+\delta}}\right]^{\mu_k/\Psi(\mathbf{V}_+)} \qquad \textit{for } \delta > 0, \textit{ and}$$

$$\mathbb{P}\left\{\lambda_k\left(\sum_j \mathbf{X}_j\right) \leq (1-\delta)\mu_k\right\} \leq k\cdot\left[\frac{e^{-\delta}}{(1-\delta)^{1-\delta}}\right]^{\mu_k/\Psi(\mathbf{V}_-)} \qquad \textit{for } \delta \in [0,1),$$

*where $\Psi$ is a function that satisfies* (2.5.1).

Theorem 2.6 tells us how the tails of the $k$th eigenvalue are controlled by the variation of the random summands in the top and bottom invariant subspaces of $\sum_j \mathbb{E}\mathbf{X}_j$. Up to the dimensional factors $k$ and $n-k+1$, the eigenvalues exhibit binomial-type tails. When $k=1$ (respectively, $k=n$) Theorem 2.6 controls the probability that the largest eigenvalue of the sum is small (respectively, the probability that the smallest eigenvalue of the sum is large), thereby complementing the one-sided Chernoff bounds of [Tro12].

*Remark* 2.7. The results in Theorem 2.6 have the following standard simplifications:

$$\mathbb{P}\left\{\lambda_k\left(\sum_j \mathbf{X}_j\right) \geq t\mu_k\right\} \leq (n-k+1)\cdot\left[\frac{e}{t}\right]^{t\mu_k/\Psi(\mathbf{V}_+)} \qquad \textit{for } t \geq e, \textit{ and}$$

$$\mathbb{P}\left\{\lambda_k\left(\sum_j \mathbf{X}_j\right) \leq t\mu_k\right\} \leq k\cdot e^{-(1-t)^2\mu_k/(2\Psi(\mathbf{V}_-))} \qquad \textit{for } t \in [0,1].$$

*Remark* 2.8. If it is difficult to estimate $\Psi(\mathbf{V}_+)$ or $\Psi(\mathbf{V}_-)$ and the summands are uniformly bounded, one can resort to the weaker estimates

$$\Psi(\mathbf{V}_+) \leq \max_{\mathbf{V}\in\mathbb{V}_{n-k+1}^n}\max_j \left\|\mathbf{V}^*\mathbf{X}_j\mathbf{V}\right\| = \max_j \left\|\mathbf{X}_j\right\| \text{ and}$$

$$\Psi(\mathbf{V}_-) \leq \max_{\mathbf{V}\in\mathbb{V}_k^n}\max_j \left\|\mathbf{V}^*\mathbf{X}_j\mathbf{V}\right\| = \max_j \left\|\mathbf{X}_j\right\|.$$

Theorem 2.6 follows from Theorem 2.5 using an appropriate bound on the matrix moment-

generating functions. The following lemma is due to Ahlswede and Winter [AW02]; see also [Tro12, Lem. 5.8].

**Lemma 2.9.** *Suppose that $\mathbf{X}$ is a random positive-semidefinite matrix that satisfies $\lambda_{\max}(\mathbf{X}) \leq 1$. Then*

$$\mathbb{E}e^{\theta \mathbf{X}} \preceq \exp\left((e^{\theta} - 1)(\mathbb{E}\mathbf{X})\right) \quad \text{for } \theta \in \mathbb{R}.$$

*Proof of Theorem 2.6, upper bound.* We consider the case where $\Psi(\mathbf{V}_+) = 1$; the general case follows by homogeneity. Define

$$\mathbf{A}_j(\mathbf{V}_+) = \mathbf{V}_+^*(\mathbb{E}\mathbf{X}_j)\mathbf{V}_+ \quad \text{and} \quad g(\theta) = e^{\theta} - 1.$$

Theorem 2.5(ii) and Lemma 2.9 imply that

$$\mathbb{P}\left\{\lambda_k\left(\sum_j \mathbf{X}_j\right) \geq (1+\delta)\mu_k\right\} \leq \inf_{\theta > 0} e^{-\theta(1+\delta)\mu_k} \cdot \operatorname{tr} \exp\left\{g(\theta)\sum_j \mathbf{V}_+^*(\mathbb{E}\mathbf{X}_j)\mathbf{V}_+\right\}.$$

Bound the trace by the maximum eigenvalue, taking into account the reduced dimension of the summands:

$$\operatorname{tr} \exp\left\{g(\theta)\sum_j \mathbf{V}_+^*(\mathbb{E}\mathbf{X}_j)\mathbf{V}_+\right\} \leq (n-k+1) \cdot \lambda_{\max}\left(\exp\left\{g(\theta)\sum_j \mathbf{V}_+^*(\mathbb{E}\mathbf{X}_j)\mathbf{V}_+\right\}\right)$$
$$= (n-k+1) \cdot \exp\left\{g(\theta) \cdot \lambda_{\max}\left(\sum_j \mathbf{V}_+^*(\mathbb{E}\mathbf{X}_j)\mathbf{V}_+\right)\right\}.$$

The equality follows from the spectral mapping theorem. Identify the quantity $\mu_k$; then combine

the last two inequalities to obtain

$$\mathbb{P}\left\{\lambda_k\left(\sum_j \mathbf{X}_j\right) \geq (1+\delta)\mu_k\right\} \leq (n-k+1)\cdot \inf_{\theta>0} \mathrm{e}^{[g(\theta)-\theta(1+\delta)]\mu_k}.$$

The right-hand side is minimized when $\theta = \log(1+\delta)$, which gives the desired upper tail bound. $\qquad\square$

*Proof of Theorem 2.6, lower bound.* As before, we consider the case where $\Psi(\mathbf{V}_-)=1$. Clearly,

$$\mathbb{P}\left\{\lambda_k\left(\sum_j \mathbf{X}_j\right) \leq (1-\delta)\mu_k\right\} = \mathbb{P}\left\{\lambda_{n-k+1}\left(\sum_j -\mathbf{X}_j\right) \geq -(1-\delta)\mu_k\right\}. \qquad (2.5.2)$$

Apply Lemma 2.9 to see that, for $\theta > 0$,

$$\mathbb{E}\mathrm{e}^{\theta(-\mathbf{V}_-^*\mathbf{X}_j\mathbf{V}_-)} = \mathbb{E}\mathrm{e}^{(-\theta)\mathbf{V}_-^*\mathbf{X}_j\mathbf{V}_-} \preceq \exp\left(g(\theta)\cdot \mathbf{V}_-^*(-\mathbb{E}\mathbf{X}_j)\mathbf{V}_-\right),$$

where $g(\theta) = 1-\mathrm{e}^{-\theta}$. Theorem 2.5(ii) thus implies that the latter probability in (2.5.2) is bounded by

$$\inf_{\theta>0} \mathrm{e}^{\theta(1-\delta)\mu_k} \cdot \mathrm{tr}\exp\left(g(\theta)\sum_j \mathbf{V}_-^*(-\mathbb{E}\mathbf{X}_j)\mathbf{V}_-\right).$$

Using reasoning analogous to that in the proof of the upper bound, we justify the first of the following inequalities:

$$\mathrm{tr}\exp\left(g(\theta)\sum_j \mathbf{V}_-^*(-\mathbb{E}\mathbf{X}_j)\mathbf{V}_-\right) \leq k\cdot \exp\left\{\lambda_{\max}\left(g(\theta)\sum_j \mathbf{V}_-^*(-\mathbb{E}\mathbf{X}_j)\mathbf{V}_-\right)\right\}$$

$$= k\cdot \exp\left\{-g(\theta)\cdot\lambda_{\min}\left(\sum_j \mathbf{V}_-^*(\mathbb{E}\mathbf{X}_j)\mathbf{V}_-\right)\right\}$$

$$= k\cdot \exp\left\{-g(\theta)\mu_k\right\}.$$

The remaining equalities follow from the fact that $-g(\theta) < 0$ and the definition of $\mu_k$.

This argument establishes the bound

$$\mathbb{P}\left\{\lambda_k\left(\sum_j \mathbf{X}_j\right) \leq (1-\delta)\mu_k\right\} \leq k \cdot \inf_{\theta > 0} e^{[\theta(1-\delta)-g(\theta)]\mu_k}.$$

The right-hand side is minimized when $\theta = -\log(1-\delta)$, which gives the desired lower tail bound. $\qquad\square$

## 2.6 Bennett and Bernstein inequalities

The classical Bennett and Bernstein inequalities use the variance or knowledge of the moments of the summands to control the probability that a sum of independent random variables deviates from its mean. In [Tro12], matrix Bennett and Bernstein inequalities are developed for the extreme eigenvalues of Hermitian random matrix sums. We establish that the interior eigenvalues satisfy analogous inequalities.

As in the derivation of the Chernoff inequalities of Section 2.5, we need a measure of how concentrated the random summands are in a given subspace. Recall that the function $\Psi : \bigcup_{1 \leq k \leq n} \mathbb{V}_k^n \to \mathbb{R}$ satisfies

$$\max_j \lambda_{\max}\left(\mathbf{V}^*\mathbf{X}_j\mathbf{V}\right) \leq \Psi(\mathbf{V}) \qquad \text{almost surely for each } \mathbf{V} \in \bigcup_{1 \leq k \leq n} \mathbb{V}_k^n. \qquad (2.6.1)$$

The sequence $\{\mathbf{X}_j\}$ associated with $\Psi$ will always be clear from context.

**Theorem 2.10** (Eigenvalue Bennett Inequality)**.** *Consider a finite sequence $\{\mathbf{X}_j\}$ of independent, random, Hermitian matrices with dimension n, all of which have zero mean. Given an integer*

*k ≤ n, define*

$$\sigma_k^2 = \lambda_k \left( \sum_j \mathbb{E}(\mathbf{X}_j^2) \right).$$

*Choose* $\mathbf{V}_+ \in \mathbb{V}_{n-k+1}^n$ *to satisfy*

$$\sigma_k^2 = \lambda_{\max} \left( \sum_j \mathbf{V}_+^* \mathbb{E}(\mathbf{X}_j^2) \mathbf{V}_+ \right).$$

*Then, for all* $t \geq 0$,

$$\mathbb{P}\left\{ \lambda_k \left( \sum_j \mathbf{X}_j \right) \geq t \right\} \leq (n-k+1) \cdot \exp\left\{ -\frac{\sigma_k^2}{\Psi(\mathbf{V}_+)^2} \cdot h\left( \frac{\Psi(\mathbf{V}_+)t}{\sigma_k^2} \right) \right\} \tag{i}$$

$$\leq (n-k+1) \cdot \exp\left\{ \frac{-t^2/2}{\sigma_k^2 + \Psi(\mathbf{V}_+)t/3} \right\} \tag{ii}$$

$$\leq \begin{cases} (n-k+1) \cdot \exp\left\{ -\frac{3}{8}t^2/\sigma_k^2 \right\} & \text{for } t \leq \sigma_k^2/\Psi(\mathbf{V}_+) \\ (n-k+1) \cdot \exp\left\{ -\frac{3}{8}t/\Psi(\mathbf{V}_+) \right\} & \text{for } t \geq \sigma_k^2/\Psi(\mathbf{V}_+), \end{cases} \tag{iii}$$

*where the function* $h(u) = (1+u)\log(1+u) - u$ *for* $u \geq 0$. *The function* $\Psi$ *satisfies* (2.6.1) *above.*

Results (i) and (ii) are, respectively, matrix analogs of the classical Bennett and Bernstein inequalities. As in the scalar case, the Bennett inequality reflects a Poisson-type decay in the tails of the eigenvalues. The Bernstein inequality states that small deviations from the eigenvalues of the expected matrix are roughly normally distributed while larger deviations are subexponential. The split Bernstein inequalities (iii) make explicit the division between these two regimes.

As stated, Theorem 2.10 controls the probability that the eigenvalues of a sum are large. Using the identity

$$\lambda_k \left( -\sum_j \mathbf{X}_j \right) = -\lambda_{n-k+1} \left( \sum_j \mathbf{X}_j \right),$$

Theorem 2.10 can also be applied to control the probability that eigenvalues of a sum are small.

To prove Theorem 2.10, we use the following lemma (Lemma 6.7 in [Tro12]) to control the moment-generating function of a random matrix with bounded maximum eigenvalue.

**Lemma 2.11.** *Let $\mathbf{X}$ be a random Hermitian matrix satisfying $\mathbb{E}\mathbf{X} = \mathbf{0}$ and $\lambda_{\max}(\mathbf{X}) \leq 1$ almost surely. Then*

$$\mathbb{E}e^{\theta\mathbf{X}} \preceq \exp((e^{\theta} - \theta - 1)\cdot\mathbb{E}(\mathbf{X}^2)) \quad \text{for } \theta > 0.$$

*Proof of Theorem 2.10.* Using homogeneity, we assume without loss that $\Psi(\mathbf{V}_+) = 1$. This implies that $\lambda_{\max}(\mathbf{X}_j) \leq 1$ almost surely for all the summands. By Lemma 2.11,

$$\mathbb{E}e^{\theta\mathbf{X}_j} \preceq \exp(g(\theta)\cdot\mathbb{E}(\mathbf{X}_j^2)),$$

with $g(\theta) = e^{\theta} - \theta - 1$.

Theorem 2.5(i) then implies

$$\mathbb{P}\left\{\lambda_k\left(\sum_j\mathbf{X}_j\right) \geq t\right\} \leq \inf_{\theta>0}e^{-\theta t}\cdot\text{tr}\exp\left(g(\theta)\sum_j\mathbf{V}_+^*\mathbb{E}(\mathbf{X}_j^2)\mathbf{V}_+\right)$$

$$\leq (n-k+1)\cdot\inf_{\theta>0}e^{-\theta t}\cdot\lambda_{\max}\left(\exp\left\{g(\theta)\sum_j\mathbf{V}_+^*\mathbb{E}(\mathbf{X}_j^2)\mathbf{V}_+\right\}\right)$$

$$= (n-k+1)\cdot\inf_{\theta>0}e^{-\theta t}\cdot\exp\left\{g(\theta)\cdot\lambda_{\max}\left(\sum_j\mathbf{V}_+^*\mathbb{E}(\mathbf{X}_j^2)\mathbf{V}_+\right)\right\}.$$

The maximum eigenvalue in this expression equals $\sigma_k^2$, thus

$$\mathbb{P}\left\{\lambda_k\left(\sum_j\mathbf{X}_j\right) \geq t\right\} \leq (n-k+1)\cdot\inf_{\theta>0}e^{g(\theta)\sigma_k^2-\theta t}.$$

The Bennett inequality (i) follows by substituting $\theta = \log(1 + t/\sigma_k^2)$ into the right-hand side and simplifying.

The Bernstein inequality (ii) is a consequence of (i) and the fact that

$$h(u) \geq \frac{u^2/2}{1 + u/3} \quad \text{for } u \geq 0,$$

which can be established by comparing derivatives.

The subgaussian and subexponential portions of the split Bernstein inequalities (iii) are verified through algebraic comparisons on the relevant intervals. □

Occasionally, as in the application in Section 2.8 to the problem of covariance matrix estimation, one desires a Bernstein-type tail bound that applies to summands that do not have bounded maximum eigenvalues. In this case, if the moments of the summands satisfy sufficiently strong growth restrictions, one can extend classical scalar arguments to obtain results such as the following Bernstein bound for subexponential matrices.

**Theorem 2.12** (Eigenvalue Bernstein Inequality for Subexponential Matrices)**.** *Consider a finite sequence* $\{\mathbf{X}_j\}$ *of independent, random, Hermitian matrices with dimension n, all of which satisfy the subexponential moment growth condition*

$$\mathbb{E}(\mathbf{X}_j^m) \preceq \frac{m!}{2} B^{m-2} \mathbf{\Sigma}_j^2 \quad \text{for } m = 2, 3, 4, \ldots,$$

*where B is a positive constant and* $\mathbf{\Sigma}_j^2$ *are positive-semidefinite matrices. Given an integer* $k \leq n$, *set*

$$\mu_k = \lambda_k \left( \sum_j \mathbb{E}\mathbf{X}_j \right).$$

*Choose* $\mathbf{V}_+ \in \mathbb{V}_{n-k+1}^n$ *that satisfies*

$$\mu_k = \lambda_{\max} \left( \sum_j \mathbf{V}_+^* (\mathbb{E}\mathbf{X}_j) \mathbf{V}_+ \right),$$

*and define*

$$\sigma_k^2 = \lambda_{\max} \left( \sum_j \mathbf{V}_+^* \mathbf{\Sigma}_j^2 \mathbf{V}_+ \right).$$

*Then, for any $t \geq 0$,*

$$\mathbb{P}\left\{ \lambda_k \left( \sum_j \mathbf{X}_j \right) \geq \mu_k + t \right\} \leq (n - k + 1) \cdot \exp\left\{ -\frac{t^2/2}{\sigma_k^2 + Bt} \right\} \tag{i}$$

$$\leq \begin{cases} (n - k + 1) \cdot \exp\left\{ -\frac{1}{4} t^2 / \sigma_k^2 \right\} & \text{for } t \leq \sigma_k^2 / B \\[2mm] (n - k + 1) \cdot \exp\left\{ -\frac{1}{4} t / B \right\} & \text{for } t \geq \sigma_k^2 / B. \end{cases} \tag{ii}$$

This result is an extension of [Tro12, Theorem 6.2], which, in turn, generalizes a classical scalar argument [DG98].

As with the other matrix inequalities, Theorem 2.12 follows from an application of Theorem 2.5 and appropriate semidefinite bounds on the moment-generating functions of the summands. Thus, the key to the proof lies in exploiting the moment growth conditions of the summands to majorize their moment-generating functions. The following lemma, a trivial extension of Lemma 6.8 in [Tro12], provides what we need.

**Lemma 2.13.** *Let $\mathbf{X}$ be a random Hermitian matrix satisfying the subexponential moment growth conditions*

$$\mathbb{E}(\mathbf{X}^m) \preceq \frac{m!}{2} \mathbf{\Sigma}^2 \quad \text{for } m = 2, 3, 4, \ldots.$$

*Then, for any $\theta$ in $[0, 1)$,*

$$\mathbb{E}\exp(\theta \mathbf{X}) \preceq \exp\left( \theta \mathbb{E}\mathbf{X} + \frac{\theta^2}{2(1 - \theta)} \mathbf{\Sigma}^2 \right).$$

*Proof of Theorem 2.12.* We note that $\mathbf{X}_j$ satisfies the growth condition

$$\mathbb{E}(\mathbf{X}_j^m) \preceq \frac{m!}{2} B^{m-2} \boldsymbol{\Sigma}_j^2 \quad \text{for } m \geq 2$$

if and only if the scaled matrix $\mathbf{X}_j/B$ satisfies

$$\mathbb{E}\left(\frac{\mathbf{X}_j}{B}\right)^m \preceq \frac{m!}{2} \cdot \frac{\boldsymbol{\Sigma}_j^2}{B^2} \quad \text{for } m \geq 2.$$

Thus, by rescaling, it suffices to consider the case $B = 1$.

By Lemma 2.13, the moment-generating functions of the summands satisfy

$$\mathbb{E}\exp(\theta\mathbf{X}_j) \preceq \exp\left(\theta\mathbb{E}\mathbf{X}_j + g(\theta)\boldsymbol{\Sigma}_j^2\right),$$

where $g(\theta) = \theta^2/(2 - 2\theta)$. Now we apply Theorem 2.5(i):

$$\mathbb{P}\left\{\lambda_k\left(\sum_j \mathbf{X}_j\right) \geq \mu_k + t\right\} \leq \inf_{\theta \in [0,1)} e^{-\theta(\mu_k+t)} \cdot \mathrm{tr}\exp\left(\theta\sum_j \mathbf{V}_+^*(\mathbb{E}\mathbf{X}_j)\mathbf{V}_+ + g(\theta)\sum_j \mathbf{V}_+^*\boldsymbol{\Sigma}_j^2\mathbf{V}_+\right)$$

$$\leq \inf_{\theta \in [0,1)} (n - k + 1) \cdot \exp\left\{-\theta(\mu_k + t) + \theta \cdot \lambda_{\max}\left(\sum_j \mathbf{V}_+^*(\mathbb{E}\mathbf{X}_j)\mathbf{V}_+\right)\right.$$

$$\left. + g(\theta) \cdot \lambda_{\max}\left(\sum_j \mathbf{V}_+^*\boldsymbol{\Sigma}_j^2\mathbf{V}_+\right)\right\}$$

$$= \inf_{\theta \in [0,1)} (n - k + 1) \cdot \exp\left(-\theta t + g(\theta)\sigma_k^2\right).$$

To achieve the final simplification, we identified $\mu_k$ and $\sigma_k^2$. Now, select $\theta = t/(t + \sigma_k^2)$. Then simplication gives the Bernstein inequality (i).

Algebraic comparisons on the relevant intervals yield the split Bernstein inequalities (ii). $\qquad\square$

## 2.7 An application to column subsampling

As an application of our Chernoff bounds, we examine how sampling columns from a matrix with orthonormal rows affects the spectrum. This question has applications in numerical linear algebra and compressed sensing. The special cases of the maximum and minimum eigenvalues have been studied in the literature [Tro08, RV07]. The limiting spectral distributions of matrices formed by sampling columns from similarly structured matrices have also been studied: the results of [GH08] apply to matrices formed by sampling columns from any fixed orthogonal matrix, and [Far10] studies matrices formed by sampling columns and rows from the discrete Fourier transform matrix.

Let $\mathbf{U}$ be an $n \times r$ matrix with orthonormal rows. We model the sampling operation using a random diagonal matrix $\mathbf{D}$ whose entries are independent $\mathrm{Bern}(p)$ random variables. Then the random matrix

$$\widehat{\mathbf{U}} = \mathbf{U}\mathbf{D} \tag{2.7.1}$$

can be interpreted as a random column submatrix of $\mathbf{U}$ with an average of $pr$ nonzero columns. Our goal is to study the behavior of the spectrum of $\widehat{\mathbf{U}}$.

Recall that the decay of the Chernoff tail bounds is influenced by the variation of the random summands when compressed to invariant subspaces of the expected sum, as measured by $\Psi(\mathbf{V})$. In this application, the choice of invariant subspace is arbitrary, so we choose that which gives the smallest variations and hence the fastest decay. This gives rise to a coherence-like quantity associated with the matrix $\mathbf{U}$ : Recall that the $j$th column of $\mathbf{U}$ is written $\mathbf{u}_j$. Consider the

following coherence-like quantity associated with $\mathbf{U}$ :

$$\tau_k = \min_{\mathbf{V} \in \mathbb{V}_k^n} \max_j \left\| \mathbf{V}^* \mathbf{u}_j \right\|^2 \quad \text{for } k = 1, \dots, n. \tag{2.7.2}$$

There does not seem to be a simple expression for $\tau_k$. However, by choosing $\mathbf{V}^*$ to be the restriction to an appropriate $k$-dimensional coordinate subspace, we see that $\tau_k$ always satisfies

$$\tau_k \leq \min_{|I| \leq k} \max_j \sum_{i \in I} u_{ij}^2.$$

The following theorem shows that the behavior of $\sigma_k(\widehat{\mathbf{U}})$, the $k$th singular value of $\widehat{\mathbf{U}}$, can be explained in terms of $\tau_k$.

**Theorem 2.14** (Column Subsampling of Matrices with Orthonormal Rows). *Let $\mathbf{U}$ be an $n \times r$ matrix with orthonormal rows, and let $p$ be a sampling probability. Define the sampled matrix $\widehat{\mathbf{U}}$ according to (2.7.1), and the numbers $\{\tau_k\}$ according to (2.7.2). Then, for each $k = 1, \dots, n$,*

$$\mathbb{P}\left\{ \sigma_k(\widehat{\mathbf{U}}) \geq \sqrt{(1+\delta)p} \right\} \leq (n - k + 1) \cdot \left[ \frac{e^{\delta}}{(1+\delta)^{1+\delta}} \right]^{p/\tau_{n-k+1}} \qquad \textit{for } \delta > 0, \textit{ and}$$

$$\mathbb{P}\left\{ \sigma_k(\widehat{\mathbf{U}}) \leq \sqrt{(1-\delta)p} \right\} \leq k \cdot \left[ \frac{e^{-\delta}}{(1-\delta)^{1-\delta}} \right]^{p/\tau_k} \qquad \textit{for } \delta \in [0,1).$$

*Proof.* Observe, using (2.7.1), that

$$\sigma_k(\widehat{\mathbf{U}})^2 = \lambda_k(\mathbf{U}\mathbf{D}^2\mathbf{U}^*) = \lambda_k\left( \sum_j d_j \mathbf{u}_j \mathbf{u}_j^* \right),$$

where $\mathbf{u}_j$ is the $j$th column of $\mathbf{U}$ and $d_j \sim \text{Bern}(p)$. Compute

$$\mu_k = \lambda_k\left( \sum_j \mathbb{E}d_j \mathbf{u}_j \mathbf{u}_j^* \right) = p \cdot \lambda_k(\mathbf{U}\mathbf{U}^*) = p \cdot \lambda_k(\mathbf{I}) = p.$$

It follows that, for *any* $\mathbf{V} \in \mathbb{V}^n_{n-k+1}$,

$$\lambda_{\max}\left(\sum_j \mathbf{V}^*(\mathbb{E}d_j\mathbf{u}_j\mathbf{u}_j^*)\mathbf{V}\right) = p \cdot \lambda_{\max}\left(\mathbf{V}^*\mathbf{V}\right) = p = \mu_k,$$

so the choice of $\mathbf{V}_+ \in \mathbb{V}^n_{n-k+1}$ is arbitrary. Similarly, the choice of $\mathbf{V}_- \in \mathbb{V}^n_k$ is arbitrary. We select $\mathbf{V}_+$ to be an isometric embedding that achieves $\tau_{n-k+1}$ and $\mathbf{V}_-$ to be an isometric embedding that achieves $\tau_k$. Accordingly,

$$\Psi(\mathbf{V}_+) = \max_j \|\mathbf{V}_+^*\mathbf{u}_j\mathbf{u}_j^*\mathbf{V}_+\| = \max_j \|\mathbf{V}_+^*\mathbf{u}_j\|^2 = \tau_{n-k+1}, \quad \text{and}$$

$$\Psi(\mathbf{V}_-) = \max_j \|\mathbf{V}_-^*\mathbf{u}_j\mathbf{u}_j^*\mathbf{V}_-\| = \max_j \|\mathbf{V}_-^*\mathbf{u}_j\|^2 = \tau_k.$$

Theorem 2.6 delivers the upper bound

$$\mathbb{P}\left\{\sigma_k(\hat{\mathbf{U}}) \geq \sqrt{(1+\delta)p}\right\} = \mathbb{P}\left\{\lambda_k\left(\sum_j d_j\mathbf{u}_j\mathbf{u}_j^*\right) \geq (1+\delta)p\right\}$$

$$\leq (n-k+1) \cdot \left[\frac{e^\delta}{(1+\delta)^{1+\delta}}\right]^{p/\tau_{n-k+1}}$$

for $\delta > 0$, and the lower bound

$$\mathbb{P}\left\{\sigma_k(\hat{\mathbf{U}}) \leq \sqrt{(1-\delta)p}\right\} = \mathbb{P}\left\{\lambda_k\left(\sum_j d_j\mathbf{u}_j\mathbf{u}_j^*\right) \leq (1-\delta)p\right\} \leq k \cdot \left[\frac{e^{-\delta}}{(1-\delta)^{1-\delta}}\right]^{p/\tau_k}$$

for $\delta \in [0, 1)$. $\qquad\square$

To illustrate the discriminatory power of these bounds, let $\mathbf{U}$ be an $n \times n^2$ matrix consisting of $n$ rows of the $n^2 \times n^2$ Fourier matrix and choose $p = (\log n)/n$ so that, on average, sampling
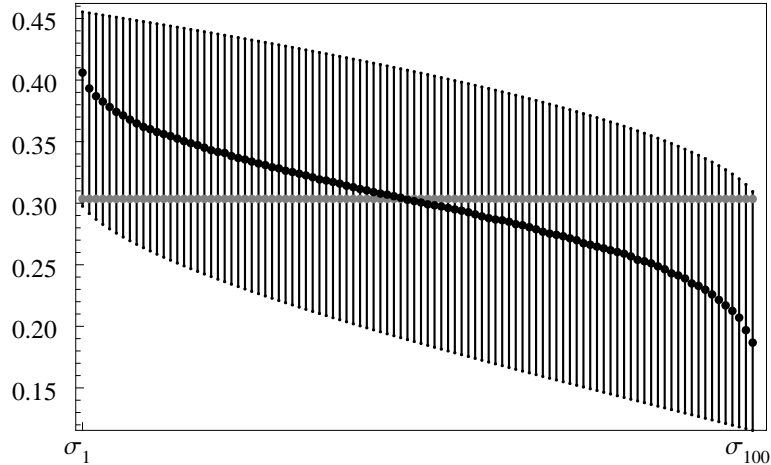
Figure 2.1: SPECTRUM OF A RANDOM SUBMATRIX OF A UNITARY DFT MATRIX. The matrix $\mathbf{U}$ is a $10^2 \times 10^4$ submatrix of the unitary DFT matrix with dimension $10^4$, and the sampling probability $p = 10^{-4} \log(10^4)$. The $k$th vertical bar, calculated using Theorem 2.14, describes an interval containing the median value of the $k$th singular value of the sampled matrix $\widehat{\mathbf{U}}$. The black circles denote the empirical medians of the singular values of $\widehat{\mathbf{U}}$, calculated from 500 trials. The gray circles represent the singular values of $\mathbb{E}\widehat{\mathbf{U}}$.

reduces the aspect ratio from $n$ to $\log n$. For $n = 100$, we determine upper and lower bounds for the median value of $\sigma_k(\widehat{\mathbf{U}})$ by numerically finding the value of $\delta$ where the probability bounds in Theorem 2.14 equal one-half. Figure 2.1 plots the empirical median value along with the computed interval. We see that these ranges reflect the behavior of the singular values more faithfully than the simple estimates $\sigma_k(\mathbb{E}\widehat{\mathbf{U}}) = p$.

## 2.8  Covariance estimation

We conclude with an extended example that illustrates how this circle of ideas allows one to answer interesting statistical questions. Specifically, we investigate the convergence of the individual eigenvalues of sample covariance matrices. Our results establish conditions under which the eigenvalues can be recovered to relative precision, and furthermore reflect the difference in the probabilities of the $k$th eigenvalue of the sample covariance matrix over- or underestimating that of the covariance matrix.

Covariance estimation is a basic and ubiquitious problem that arises in signal processing, graphical modeling, machine learning, and genomics, among other areas. Let $\{\eta_j\}_{j=1}^n \subset \mathbb{R}^p$ be i.i.d. samples drawn from some distribution with zero mean and covariance matrix $\mathbf{C}$. Define the sample covariance matrix

$$\widehat{\mathbf{C}}_n = \frac{1}{n} \sum_{j=1}^n \eta_j \eta_j^*.$$

An important challenge is to determine how many samples are needed to ensure that the empirical covariance estimator has a fixed relative accuracy in the spectral norm. That is, given a fixed $\varepsilon$, how large must $n$ be so that

$$\left\| \widehat{\mathbf{C}}_n - \mathbf{C} \right\|_2 \leq \varepsilon \left\| \mathbf{C} \right\|_2 ? \tag{2.8.1}$$

This estimation problem has been studied extensively. It is now known that for distributions with a finite second moment, $\Omega(p \log p)$ samples suffice [Rud99], and for log-concave distributions, $\Omega(p)$ samples suffice [ALPTJ11]. More broadly, Vershynin [Ver11b] conjectures that, for distributions with finite fourth moment, $\Omega(p)$ samples suffice; he establishes this result to within iterated log factors. In [SV], Srivastava and Vershynin establish that $\Omega(p)$ samples suffice for distributions which have finite $2 + \varepsilon$ moments, for some $\varepsilon > 0$, and satisfy an additional regularity condition.

Inequality (2.8.1) ensures that the difference between the $k$th eigenvalues of $\widehat{\mathbf{C}}_n$ and $\mathbf{C}$ is small, but it requires $O(p)$ samples to obtain estimates of even a few of the eigenvalues. Specifically, letting $\kappa_\ell = \lambda_1(\mathbf{C})/\lambda_\ell(\mathbf{C})$, we see that $O(\varepsilon^{-2} \kappa_\ell^2 p)$ samples are required to obtain relative error estimates of the largest $\ell$ eigenvalues of $\mathbf{C}$ using the results of [ALPTJ11, Ver11b, SV]. However, it is reasonable to expect that when the spectrum of $\mathbf{C}$ exhibits decay and $\ell \ll p$, far fewer than $O(p)$ samples should suffice to ensure relative error recovery of the largest $\ell$

eigenvalues.

In fact, Vershynin shows this is the case when the random vector is subgaussian: in [Ver11a], he defines the effective rank of $\mathbf{C}$ to be $r = \left( \sum_{i=1}^{p} \lambda_i(\mathbf{C}) \right) / \lambda_1(\mathbf{C})$ and uses $r$ to provide bounds of the form (2.8.1). It follow from his arguments that, with high probability, the largest $\ell$ eigenvalues of $\mathbf{C}$ are estimated to relative precision when $n = \mathrm{O}(\varepsilon^{-2} r \kappa_\ell^2 \log p)$ samples are taken. Clearly this result is most of interest when the effective rank is small: e.g. when $r$ is $\mathrm{O}(1)$, we see that $\mathrm{O}(\varepsilon^{-2} \kappa_\ell^2 \log p)$ samples suffice to give relative error accuracy in the largest $\ell$ eigenvalues of $\mathbf{C}$. Note, however, that this result does not supply the rates of convergence of the *individual* eigenvalues, and it requires the effective rank to be small. To the best of the author's knowledge, there are no nonasymptotic estimates of the relative errors of individual eigenvalues that do not require the assumption that $\mathbf{C}$ has low effective rank.

In this section, we derive a relative approximation bound for each eigenvalue of $\mathbf{C}$. For simplicity, we assume the samples are drawn from a $\mathcal{N}(\mathbf{0}, \mathbf{C})$ distribution where $\mathbf{C}$ is full-rank, but we expect that the arguments can be extended to cover other subgaussian distributions.

**Theorem 2.15.** *Assume that $\mathbf{C} \in \mathbb{M}_{\mathrm{sa}}^p$ is positive definite. Let $\{\eta_j\}_{j=1}^n \subset \mathbb{R}^p$ be i.i.d. samples drawn from a $\mathcal{N}(\mathbf{0}, \mathbf{C})$ distribution. Define*

$$\widehat{\mathbf{C}}_n = \frac{1}{n} \sum_{j=1}^n \eta_j \eta_j^*.$$

*Write $\lambda_k$ for the kth eigenvalue of $\mathbf{C}$, and write $\hat{\lambda}_k$ for the kth eigenvalue of $\widehat{\mathbf{C}}_n$. Then for $k = 1, \ldots, p$,*

$$\mathbb{P}\left\{ \hat{\lambda}_k \geq \lambda_k + t \right\} \leq (p - k + 1) \cdot \exp\left( \frac{-nt^2}{32 \lambda_k \sum_{i=k}^p \lambda_i} \right) \quad \text{for } t \leq 4n\lambda_k,$$

*and*

$$\mathbb{P}\left\{\hat{\lambda}_k \leq \lambda_k - t\right\} \leq k \cdot \exp\left(\frac{-3nt^2}{8\lambda_1\left(\lambda_1 + \sum_{i=1}^k \lambda_i\right)}\right) \quad \textit{for } t \leq n\left(\lambda_1 + \sum_{i=1}^k \lambda_i\right).$$

The following corollary provides an answer to our question about relative error estimates.

**Corollary 2.16.** *Let $\lambda_k$ and $\hat{\lambda}_k$ be as in Theorem 2.15. Then*

$$\mathbb{P}\left\{\hat{\lambda}_k \geq (1+\varepsilon)\lambda_k\right\} \leq (p-k+1) \cdot \exp\left(\frac{-cn\varepsilon^2}{\sum_{i=k}^p \frac{\lambda_i}{\lambda_k}}\right) \quad \textit{for } \varepsilon \leq 4n,$$

*and*

$$\mathbb{P}\left\{\hat{\lambda}_k \leq (1-\varepsilon)\lambda_k\right\} \leq k \cdot \exp\left(\frac{-cn\varepsilon^2}{\frac{\lambda_1}{\lambda_k}\left(\sum_{i=1}^k \frac{\lambda_i}{\lambda_k}\right)}\right) \quad \textit{for } \varepsilon \in (0,1],$$

*where the constant $c$ is at least $1/32$.*

The first bound in Corollary 2.16 tells us how many samples are needed to ensure that $\hat{\lambda}_k$ does not overestimate $\lambda_k$. Likewise, the second bound tells us how many samples ensure that $\hat{\lambda}_k$ does not underestimate $\lambda_k$.

Corollary 2.16 suggests that the relationship of $\hat{\lambda}_k$ to $\lambda_k$ is determined by the spectrum of **C** in the following manner. When the eigenvalues below $\lambda_k$ are small compared with $\lambda_k$, the quantity

$$\sum_{i=k}^p \lambda_i / \lambda_k$$

is small (viz., it is no larger than $p - k + 1$), and so $\hat{\lambda}_k$ is not likely to overestimate $\lambda_k$. Similarly,

when the eigenvalues above $\lambda_k$ are comparable with $\lambda_k$, the quantity

$$\frac{\lambda_1}{\lambda_k}\left(\sum_{i=1}^{k}\lambda_i/\lambda_k\right)$$

is small (viz., it is no larger than $k \cdot \kappa_k^2$), and so $\hat{\lambda}_k$ is not likely to underestimate $\lambda_k$.

*Remark* 2.17. The results in Theorem 2.15 and Corollary 2.16 also apply when $\mathbf{C}$ is rank-deficient: simply replace each occurence of the dimension $p$ in the bounds with rank($\mathbf{C}$).

Indeed, assume that $\mathbf{C}$ is rank-deficient and take its truncated eigenvalue decomposition to be $\mathbf{C} = \mathbf{U}\mathbf{\Lambda}\mathbf{U}^*$. If $\eta_j \sim \mathcal{N}(\mathbf{0}, \mathbf{C})$, then $\eta_j$ lies in the span of $\mathbf{C}$. It follows that $\hat{\lambda}_k = \lambda_k = 0$ for all $k > \text{rank}(\mathbf{C})$. When $k \leq \text{rank}(\mathbf{C})$, we observe that

$$\lambda_k(\mathbf{C}) = \lambda_k(\mathbf{\Lambda}) \quad \text{and} \quad \lambda_k\left(\sum_j \eta_j\eta_j^*\right) = \lambda_k\left(\sum_j \xi_j\xi_j^*\right),$$

where $\xi_j = \mathbf{U}^*\eta_j$ is distributed $\mathcal{N}(\mathbf{0}, \mathbf{\Lambda})$. Thus,

$$\left|\lambda_k\left(\sum_j \eta_j\eta_j^*\right) - \lambda_k(\mathbf{C})\right| = \left|\lambda_k\left(\sum_j \xi_j\xi_j^*\right) - \lambda_k(\mathbf{\Lambda})\right|.$$

Consequently, the problem of estimating the eigenvalues of $\mathbf{C}$ to relative error using the samples $\{\eta_j\}$ is equivalent to that of estimating the eigenvalues of the full-rank covariance matrix $\mathbf{\Lambda}$ to relative error using the samples $\{\xi_j\}$.

It is reasonable to expect that one should be able to use Corollary 2.16 to recover Vershynin's result in [Ver11a] for Wishart matrices: that $\Omega(\varepsilon^{-2} r \kappa_\ell^2 \log p)$ samples suffice to estimate the eigenvalues of the covariance matrix of a Gaussian random variable to within a relative precision of $1 \pm \varepsilon$. Indeed, this result follows from Corollary 2.16 and a simple union bound argument.

**Corollary 2.18.** *Assume* $\mathbf{C}$ *is positive semidefinite. Let* $\{\eta_j\}_{j=1}^{n} \subset \mathbb{R}^p$ *be i.i.d. samples drawn from*

*a* $\mathcal{N}(\mathbf{0}, \mathbf{C})$ *distribution. If* $n = \Omega(\varepsilon^{-2} r \kappa_\ell^2 \log p)$*, then with high probability*

$$|\lambda_k(\widehat{\mathbf{C}}_n) - \lambda_k(\mathbf{C})| \leq \varepsilon \lambda_k(\mathbf{C}) \quad \text{for } k = 1, \dots, \ell.$$

*Proof.* From Corollary 2.16, we see that

$$\mathbb{P}\left\{\lambda_k(\widehat{\mathbf{C}}_n) \leq (1 - \varepsilon)\lambda_k\right\} \leq p^{-\beta} \quad \text{when } n \geq 32\varepsilon^{-2}\left(\frac{\lambda_1}{\lambda_k}\sum_{i \leq k}\frac{\lambda_i}{\lambda_k}\right)(\log k + \beta \log p).$$

Recall that $\kappa_k = \lambda_1(\mathbf{C})/\lambda_k(\mathbf{C})$ and $r = \left(\sum_{i=1}^{p}\lambda_i(\mathbf{C})\right)/\lambda_1(\mathbf{C})$, so

$$\left(\frac{\lambda_1}{\lambda_k}\sum_{i \leq k}\frac{\lambda_i}{\lambda_k}\right) \leq \kappa_k^2 r.$$

Clearly, taking $n = \Omega(\varepsilon^{-2} r \kappa_\ell^2 \log p)$ samples ensures that, with high probability, each of the top $\ell$ eigenvalues of the sample covariance matrix satisfies $\lambda_k(\widehat{\mathbf{C}}_n) > (1 - \varepsilon)\lambda_k$.

Likewise,

$$\mathbb{P}\left\{\lambda_k(\widehat{\mathbf{C}}_n) \geq (1 - \varepsilon)\lambda_k\right\} \leq p^{-\beta} \quad \text{when } n \geq 32\varepsilon^{-2}\left(\sum_{i \geq k}\frac{\lambda_i}{\lambda_k}\right)(\log(p - k + 1) + \beta \log p)$$

and

$$\sum_{i \geq k}\frac{\lambda_i}{\lambda_k} = \frac{\lambda_1}{\lambda_k}\frac{\left(\sum_{i \geq k}\lambda_i\right)}{\lambda_1} \leq \kappa_k\frac{\left(\sum_{i=1}^{p}\lambda_i\right)}{\lambda_1} = \kappa_k r,$$

so we see that taking $n = \Omega(\varepsilon^{-2} r \kappa_\ell \log p)$ samples ensures that, with high probability, each of the top $\ell$ eigenvalues of the sample covariance matrix satisfies $\lambda_k(\widehat{\mathbf{C}}_n) < (1 + \varepsilon)\lambda_k$.

Combining these two results, we conclude that $n = \Omega(\varepsilon^{-2} r \kappa_\ell^2 \log p)$ ensures that the top $\ell$ eigenvalues of $\mathbf{C}$ are estimated to within relative precision $1 \pm \varepsilon$ with probability at least

$1 - 2\ell p^{-\beta}$. □

### 2.8.1  Proof of Theorem 2.15

We now prove Theorem 2.15. This result requires a number of supporting lemmas; we defer their proofs until after a discussion of extensions to Theorem 2.15.

We study the error $|\lambda_k(\widehat{\mathbf{C}}_n) - \lambda_k(\mathbf{C})|$. To apply the methods developed in this chapter, we pass to a question about the eigenvalues of a difference of two matrices. The first lemma accomplishes this goal by compressing both the population covariance matrix and the sample covariance matrix to a fixed invariant subspace of the population covariance matrix.

**Lemma 2.19.** *Let $\mathbf{X}$ be a random Hermitian matrix with dimension $p$, and let $\mathbf{A}$ be a fixed Hermitian matrix with dimension $p$. Choose $\mathbf{W}_+ \in \mathbb{V}_{p-k+1}^p$ and $\mathbf{W}_- \in \mathbb{V}_k^p$ for which*

$$\lambda_k(\mathbf{A}) = \lambda_{\max}\left(\mathbf{W}_+^* \mathbf{A} \mathbf{W}_+\right) = \lambda_{\min}\left(\mathbf{W}_-^* \mathbf{A} \mathbf{W}_-\right).$$

*Then, for all $t > 0$,*

$$\mathbb{P}\left\{\lambda_k(\mathbf{X}) \geq \lambda_k(\mathbf{A}) + t\right\} \leq \mathbb{P}\left\{\lambda_{\max}\left(\mathbf{W}_+^* \mathbf{X} \mathbf{W}_+\right) \geq \lambda_k(\mathbf{A}) + t\right\} \tag{2.8.2}$$

*and*

$$\mathbb{P}\left\{\lambda_k(\mathbf{X}) \leq \lambda_k(\mathbf{A}) - t\right\} \leq \mathbb{P}\left\{\lambda_{\max}\left(\mathbf{W}_-^*(\mathbf{A} - \mathbf{X})\mathbf{W}_-\right) \geq t\right\}. \tag{2.8.3}$$

We apply this result with $\mathbf{A} = \mathbf{C}$ and $\mathbf{X} = \widehat{\mathbf{C}}_n$. The first estimate (2.8.2) and the second estimate (2.8.3) are handled using different arguments. The second estimate is easier because

the maximum eigenvalue of the matrix $\mathbf{C} - \widehat{\mathbf{C}}_n$ is bounded. Indeed,

$$\lambda_{\max}\left(\mathbf{W}_+^*(\mathbf{C} - \widehat{\mathbf{C}}_n)\mathbf{W}_+\right) \leq \lambda_{\max}\left(\mathbf{W}_+^*\mathbf{C}\mathbf{W}_+\right).$$

Thus, we may use Theorem 2.10 to complete the second estimate. The next lemma gives the matrix variances that we need to apply this theorem.

**Lemma 2.20.** *Let $\xi \sim \mathcal{N}(\mathbf{0}, \mathbf{G})$. Then*

$$\mathbb{E}(\xi\xi^* - \mathbf{G})^2 = \mathbf{G}^2 + \operatorname{tr}(\mathbf{G}) \cdot \mathbf{G}.$$

The first inequality (2.8.2) is harder because $\widehat{\mathbf{C}}_n$ is unbounded. In this case, we may apply Theorem 2.12. To use this theorem, we need the following moment growth estimate for rank-one Wishart matrices.

**Lemma 2.21.** *Let $\xi \sim \mathcal{N}(\mathbf{0}, \mathbf{G})$. Then for any integer $m \geq 2$,*

$$\mathbb{E}\left(\xi\xi^*\right)^m \preceq 2^m m!(\operatorname{tr}\mathbf{G})^{m-1} \cdot \mathbf{G}.$$

With these preliminaries addressed, we prove Theorem 2.15.

*Proof of lower estimate in Theorem 2.15.* First we consider the probability that $\hat{\lambda}_k$ underestimates $\lambda_k$. Let $\mathbf{W}_- \in \mathbb{V}_k^p$ satisfy

$$\lambda_k(\mathbf{C}) = \lambda_{\min}\left(\mathbf{W}_-^*\mathbf{C}\mathbf{W}_-\right).$$

Then Lemma 2.19 implies

$$\mathbb{P}\left\{\lambda_k(\widehat{\mathbf{C}}_n) \leq \lambda_k(\mathbf{C}) - t\right\} \leq \mathbb{P}\left\{\lambda_{\max}\left(\mathbf{W}_-^*(\mathbf{C} - \widehat{\mathbf{C}}_n)\mathbf{W}_-\right) \geq t\right\}$$

$$= \mathbb{P}\left\{\lambda_{\max}\left(\sum_j \mathbf{W}_-^*(\mathbf{C} - \eta_j\eta_j^*)\mathbf{W}_-\right) \geq nt\right\}.$$

The factor $n$ comes from the normalization of the sample covariance matrix. Each term in the sum is zero mean and bounded above by $\mathbf{W}_-^*\mathbf{C}\mathbf{W}_-$ in the semidefinite order, so Theorem 2.10 applies. As we desire a bound on the maximum eigenvalue of the sum, we take $\mathbf{V}_+ = \mathbf{I}$ when we invoke Theorem 2.10. Then

$$\sigma_1^2 = \lambda_{\max}\left(\sum_j \mathbb{E}\left[\mathbf{W}_-^*(\mathbf{C} - \eta_j\eta_j^*)\mathbf{W}_-\right]^2\right) = n\lambda_{\max}\left(\mathbb{E}\left[\mathbf{W}_-^*(\mathbf{C} - \eta_1\eta_1^*)\mathbf{W}_-\right]^2\right).$$

The covariance matrix of $\eta_1$ is $\mathbf{C}$, so that of $\mathbf{W}_-^*\eta_1$ is $\mathbf{W}_-^*\mathbf{C}\mathbf{W}_-$. It follows from Lemma 2.20 that

$$\mathbb{E}\left[\mathbf{W}_-^*(\mathbf{C} - \eta_1\eta_1^*)\mathbf{W}_-\right]^2 = (\mathbf{W}_-^*\mathbf{C}\mathbf{W}_-)^2 + \mathrm{tr}(\mathbf{W}_-^*\mathbf{C}\mathbf{W}_-) \cdot \mathbf{W}_-^*\mathbf{C}\mathbf{W}_-.$$

Observe that $\mathbf{W}_-^*\mathbf{C}\mathbf{W}_-$ is the restriction of $\mathbf{C}$ to its top $k$-dimensional invariant subspace, so

$$\sigma_1^2 = n\lambda_{\max}\left(\mathbb{E}\left[\mathbf{W}_-^*(\mathbf{C} - \eta_1\eta_1^*)\mathbf{W}_-\right]^2\right) = n\lambda_1(\mathbf{C})\left(\lambda_1(\mathbf{C}) + \sum_{i=1}^{k}\lambda_i(\mathbf{C})\right)$$

and we can take $\Psi(\mathbf{V}_+) = \lambda_{\max}(\mathbf{C})$.

The subgaussian branch of the split Bernstein inequality of Theorem 2.10 shows that

$$\mathbb{P}\left\{\lambda_{\max}\left(\sum_j \mathbf{W}_-^*(\mathbf{C} - \eta_j\eta_j^*)\mathbf{W}_-\right) \geq nt\right\} \leq k \cdot \exp\left(\frac{-3nt^2}{8\lambda_1(\mathbf{C})(\lambda_1(\mathbf{C}) + \sum_{i=1}^{k}\lambda_i(\mathbf{C}))}\right)$$

when $t \leq n\left(\lambda_1(\mathbf{C}) + \sum_{i=1}^{k}\lambda_i(\mathbf{C})\right)$. This inequality provides the desired bound on the probability

that $\lambda_k(\widehat{\mathbf{C}}_n)$ underestimates $\lambda_k(\mathbf{C})$. □

*Proof of upper estimate in Theorem 2.15.* Now we consider the probability that $\hat{\lambda}_k$ overestimates $\lambda_k$. Let $\mathbf{W}_+ \in \mathbb{V}^p_{p-k+1}$ satisfy

$$\lambda_k(\mathbf{C}) = \lambda_{\max}\left(\mathbf{W}^*_+ \mathbf{C} \mathbf{W}_+\right).$$

Then Lemma 2.19 implies

$$\mathbb{P}\left\{\lambda_k(\widehat{\mathbf{C}}_n) \geq \lambda_k(\mathbf{C}) + t\right\} \leq \mathbb{P}\left\{\lambda_{\max}\left(\mathbf{W}^*_+ \widehat{\mathbf{C}}_n \mathbf{W}_+\right) \geq \lambda_k(\mathbf{C}) + t\right\}$$

$$= \mathbb{P}\left\{\lambda_{\max}\left(\sum_j \mathbf{W}^*_+(\eta_j \eta_j^*)\mathbf{W}_+\right) \geq n\lambda_k(\mathbf{C}) + nt\right\}. \qquad (2.8.4)$$

The factor $n$ comes from the normalization of the sample covariance matrix.

The covariance matrix of $\eta_j$ is $\mathbf{C}$, so that of $\mathbf{W}^*_+ \eta_j$ is $\mathbf{W}^*_+ \mathbf{C} \mathbf{W}_+$. Apply Lemma 2.21 to verify that $\mathbf{W}^*_+ \eta_j$ satisfies the subexponential moment growth bound required by Theorem 2.12 with

$$B = 2\operatorname{tr}(\mathbf{W}^*_+ \mathbf{C} \mathbf{W}_+) \quad \text{and} \quad \boldsymbol{\Sigma}^2_j = 8\operatorname{tr}(\mathbf{W}^*_+ \mathbf{C} \mathbf{W}) \cdot \mathbf{W}^*_+ \mathbf{C} \mathbf{W}_+.$$

In fact, $\mathbf{W}^*_+ \mathbf{C} \mathbf{W}_+$ is the compression of $\mathbf{C}$ to the invariant subspace corresponding with its bottom $p - k + 1$ eigenvalues, so

$$B = 2\sum_{i=k}^p \lambda_i(\mathbf{C}) \quad \text{and} \quad \lambda_{\max}\left(\boldsymbol{\Sigma}^2_j\right) = 8\lambda_k(\mathbf{C})\sum_{i=k}^p \lambda_i(\mathbf{C}).$$

We are concerned with the maximum eigenvalue of the sum in (2.8.4), so we take $\mathbf{V}_+ = \mathbf{I}$ in the

statement of Theorem 2.12 to find that

$$\sigma_1^2 = \lambda_{\max}\left(\sum_j \Sigma_j^2\right) = n\lambda_{\max}\left(\Sigma_1^2\right) = 8n\lambda_k(\mathbf{C})\sum_{i=k}^{p}\lambda_i(\mathbf{C}) \quad \text{and}$$

$$\mu_1 = \lambda_{\max}\left(\sum_j \mathbf{W}_+^* \mathbb{E}(\eta_j\eta_j^*)\mathbf{W}_+\right) = n\lambda_{\max}\left(\mathbf{W}_+^*\mathbf{C}\mathbf{W}_+\right) = n\lambda_k(\mathbf{C}).$$

It follows from the subgaussian branch of the split Bernstein inequality of Theorem 2.12 that

$$\mathbb{P}\left\{\lambda_k\left(\sum_j \mathbf{W}_+^*(\eta_j\eta_j^*)\mathbf{W}_+\right) \geq n\lambda_k(\mathbf{C}) + nt\right\} \leq (p-k+1)\cdot\exp\left(\frac{-nt^2}{32\lambda_k(\mathbf{C})\sum_{i=k}^{p}\lambda_i(\mathbf{C})}\right)$$

when $t \leq 4n\lambda_k(\mathbf{C})$. This provides the desired bound on the probability that $\lambda_k(\widehat{\mathbf{C}}_n)$ overestimates $\lambda_k(\mathbf{C})$. $\qquad\square$

### 2.8.2 Extensions of Theorem 2.15

Results analogous to Theorem 2.15 can be established for other distributions. If the distribution is bounded, the possibility that $\hat{\lambda}_k$ deviates above or below $\lambda_k$ can be controlled using the Bernstein inequality of Theorem 2.10. If the distribution is unbounded but has matrix moments that satisfy a sufficiently nice growth condition, the probability that $\hat{\lambda}_k$ deviates below $\lambda_k$ can be controlled with the Bernstein inequality of Theorem 2.10 and the probability that it deviates above $\lambda_k$ can be bounded using a Bernstein inequality analogous to that in Theorem 2.12.

We established Theorem 2.15 using this technique to demonstrate the simplicity of the Laplace transform machinery. However, the results of [ALPTJ11] on the convergence of empirical covariance matrices of isotropic log-concave random vectors lead to tighter bounds on the probability that $\hat{\lambda}_k$ overestimates $\lambda_k$. There does not seem to be an analogous reduction for handling the probability that $\hat{\lambda}_k$ is an underestimate.

To see the relevance of the results in [ALPTJ11], first observe the following consequence of the subadditivity of the maximum eigenvalue mapping:

$$\lambda_{\max}\left(\mathbf{W}_+^*(\mathbf{X}-\mathbf{A})\mathbf{W}_+\right) \geq \lambda_{\max}\left(\mathbf{W}_+^*\mathbf{X}\mathbf{W}_+\right) - \lambda_{\max}\left(\mathbf{W}_+^*\mathbf{A}\mathbf{W}_+\right)$$

$$= \lambda_{\max}\left(\mathbf{W}_+^*\mathbf{X}\mathbf{W}_+\right) - \lambda_k(\mathbf{A}).$$

In conjunction with (2.8.2), this gives us the following control on the probability that $\lambda_k(\mathbf{X})$ overestimates $\lambda_k(\mathbf{A})$ :

$$\mathbb{P}\left\{\lambda_k(\mathbf{X}) \geq \lambda_k(\mathbf{A}) + t\right\} \leq \mathbb{P}\left\{\lambda_{\max}\left(\mathbf{W}_+^*(\mathbf{X}-\mathbf{A})\mathbf{W}_+\right) \geq t\right\}.$$

In our application, $\mathbf{X}$ is the empirical covariance matrix and $\mathbf{A}$ is the actual covariance matrix. The spectral norm dominates the maximum eigenvalue, so

$$\mathbb{P}\left\{\lambda_k(\widehat{\mathbf{C}}_n) \geq \lambda_k(\mathbf{C}) + t\right\} \leq \mathbb{P}\left\{\lambda_{\max}\left(\mathbf{W}_+^*(\widehat{\mathbf{C}}_n - \mathbf{C})\mathbf{W}_+\right) \geq t\right\}$$

$$\leq \mathbb{P}\left\{\|\mathbf{W}_+^*(\widehat{\mathbf{C}}_n - \mathbf{C})\mathbf{W}_+\| \geq t\right\} = \mathbb{P}\left\{\|\mathbf{W}_+^*\widehat{\mathbf{C}}_n\mathbf{W}_+ - \mathbf{S}^2\| \geq t\right\},$$

where $\mathbf{S}$ is the square root of $\mathbf{W}_+^*\mathbf{C}\mathbf{W}_+$. Now factor out $\mathbf{S}^2$ and identify $\lambda_k(\mathbf{C}) = \|\mathbf{S}^2\|$ to obtain

$$\mathbb{P}\left\{\lambda_k(\widehat{\mathbf{C}}) \geq \lambda_k(\mathbf{C}) + t\right\} \leq \mathbb{P}\left\{\|\mathbf{S}^{-1}\mathbf{W}_+^*\widehat{\mathbf{C}}_n\mathbf{W}_+\mathbf{S}^{-1} - \mathbf{I}\|\|\mathbf{S}^2\| \geq t\right\}$$

$$= \mathbb{P}\left\{\|\mathbf{S}^{-1}\mathbf{W}_+^*\widehat{\mathbf{C}}_n\mathbf{W}_+\mathbf{S}^{-1} - \mathbf{I}\| \geq t/\lambda_k(\mathbf{C})\right\}.$$

Note that if $\eta$ is drawn from a $\mathcal{N}(\mathbf{0}, \mathbf{C})$ distribution, then the covariance matrix of the transformed

sample $\mathbf{S}^{-1}\mathbf{W}_+^*\eta$ is the identity:

$$\mathbb{E}\left(\mathbf{S}^{-1}\mathbf{W}_+^*\eta\eta^*\mathbf{W}_+\mathbf{S}^{-1}\right) = \mathbf{S}^{-1}\mathbf{W}_+^*\mathbf{C}\mathbf{W}_+\mathbf{S}^{-1} = \mathbf{I}.$$

Thus $\mathbf{S}^{-1}\mathbf{W}_+^*\widehat{\mathbf{C}}_n\mathbf{W}_+\mathbf{S}^{-1}$ is the empirical covariance matrix of a standard Gaussian vector in $\mathbb{R}^{p-k+1}$. By Theorem 1 of [ALPTJ11], it follows that $\hat{\lambda}_k$ is unlikely to overestimate $\lambda_k$ in relative error when the number $n$ of samples is $\Omega(p - k + 1)$.

Similarly, for more general distributions, the bounds on the probability of $\hat{\lambda}_k$ exceeding $\lambda_k$ can be tightened beyond those suggested in Theorem 2.15 by using the results in [ALPTJ11] or [Ver11b].

Finally, we note that the techniques developed in the proof of Theorem 2.15 can be used to investigate the spectrum of the error matrices $\widehat{\mathbf{C}}_n - \mathbf{C}$.

### 2.8.3 Proofs of the supporting lemmas

We now establish the lemmas used in the proof of Theorem 2.15.

*Proof of Lemma 2.19.* The probability that $\lambda_k(\mathbf{X})$ overestimates $\lambda_k(\mathbf{A})$ is controlled with the sequence of inequalities

$$\mathbb{P}\left\{\lambda_k(\mathbf{X}) \geq \lambda_k(\mathbf{A}) + t\right\} = \mathbb{P}\left\{\inf_{\mathbf{W}\in\mathbb{V}_{p-k+1}^p} \lambda_{\max}\left(\mathbf{W}^*\mathbf{X}\mathbf{W}\right) \geq \lambda_k(\mathbf{A}) + t\right\}$$

$$\leq \mathbb{P}\left\{\lambda_{\max}\left(\mathbf{W}_+^*\mathbf{X}\mathbf{W}_+\right) \geq \lambda_k(\mathbf{A}) + t\right\}.$$

We use a related approach to study the probability that $\lambda_k(\mathbf{X})$ underestimates $\lambda_k(\mathbf{A})$. Our

choice of $\mathbf{W}_-$ implies that

$$\lambda_{p-k+1}(-\mathbf{A}) = -\lambda_k(\mathbf{A}) = -\lambda_{\min}\left(\mathbf{W}_-^* \mathbf{A} \mathbf{W}_-\right) = \lambda_{\max}\left(\mathbf{W}_-^*(-\mathbf{A})\mathbf{W}_-\right).$$

It follows that

$$\begin{aligned}
\mathbb{P}\left\{\lambda_k(\mathbf{X}) \leq \lambda_k(\mathbf{A}) - t\right\} &= \mathbb{P}\left\{\lambda_{p-k+1}(-\mathbf{X}) \geq \lambda_{p-k+1}(-\mathbf{A}) + t\right\} \\
&= \mathbb{P}\left\{\inf_{\mathbf{W} \in \mathbb{V}_k^p} \lambda_{\max}\left(\mathbf{W}^*(-\mathbf{X})\mathbf{W}\right) \geq \lambda_{\max}\left(\mathbf{W}_-^*(-\mathbf{A})\mathbf{W}_-\right) + t\right\} \\
&\leq \mathbb{P}\left\{\lambda_{\max}\left(\mathbf{W}_-^*(-\mathbf{X})\mathbf{W}_-\right) - \lambda_{\max}\left(\mathbf{W}_-^*(-\mathbf{A})\mathbf{W}_-\right) \geq t\right\} \\
&\leq \mathbb{P}\left\{\lambda_{\max}\left(\mathbf{W}_-^*(\mathbf{A} - \mathbf{X})\mathbf{W}_-\right) \geq t\right\}.
\end{aligned}$$

The final inequality follows from the subadditivity of the maximum eigenvalue mapping. $\quad\square$

*Proof of Lemma 2.20.* We begin by taking $\mathbf{S}$ to be the positive-semidefinite square root of $\mathbf{G}$. Let $\mathbf{S} = \mathbf{U}\boldsymbol{\Lambda}\mathbf{U}^*$ be the eigenvalue decomposition of $\mathbf{S}$, and let $\gamma$ be a $\mathcal{N}(\mathbf{0}, \mathbf{I}_p)$ random variable. Recalling that $\mathbf{G}$ is the covariance matrix of $\xi$, we see that $\xi$ and $\mathbf{U}\boldsymbol{\Lambda}\gamma$ are identically distributed. Thus,

$$\begin{aligned}
\mathbb{E}(\xi\xi^* - \mathbf{G})^2 &= \mathbb{E}(\mathbf{U}\boldsymbol{\Lambda}\gamma\gamma^*\boldsymbol{\Lambda}\mathbf{U}^* - \mathbf{U}\boldsymbol{\Lambda}^2\mathbf{U}^*)^2 \\
&= \mathbf{U}\boldsymbol{\Lambda}\mathbb{E}(\gamma\gamma^*\boldsymbol{\Lambda}^2\gamma\gamma^*)\boldsymbol{\Lambda}\mathbf{U}^* - \mathbf{G}^2. \quad\quad\quad (2.8.5)
\end{aligned}$$

Consider the $(i, j)$ entry of the matrix being averaged:

$$\mathbb{E}(\gamma\gamma^*\boldsymbol{\Lambda}^2\gamma\gamma^*)_{ij} = \sum_k \mathbb{E}(\gamma_i\gamma_j\gamma_k^2)\lambda_k^2.$$

The $(i, j)$ entry of this matrix is zero because the entries of $\gamma$ are independent and symmetric.

Furthermore, the $(i, i)$ entry satisfies

$$\mathbb{E}(\gamma\gamma^*\boldsymbol{\Lambda}^2\gamma\gamma^*)_{ii} = \mathbb{E}(\gamma_i^4)\lambda_i^2 + \sum_{k \neq i} \mathbb{E}(\gamma_k^2)\lambda_k^2 = 2\lambda_i^2 + \text{tr}(\boldsymbol{\Lambda}^2).$$

We have shown

$$\mathbb{E}(\gamma\gamma^*\boldsymbol{\Lambda}^2\gamma\gamma^*) = 2\boldsymbol{\Lambda}^2 + \text{tr}(\mathbf{G}) \cdot \mathbf{I}.$$

This equality and (2.8.5) imply the desired result. $\qquad\square$

*Proof of Lemma 2.21.* Factor the covariance matrix of $\xi$ as $\mathbf{G} = \mathbf{U}\boldsymbol{\Lambda}\mathbf{U}^*$ where $\mathbf{U}$ is orthogonal and $\boldsymbol{\Lambda} = \text{diag}(\lambda_1, \ldots, \lambda_p)$ is the matrix of eigenvalues of $\mathbf{G}$. Let $\gamma$ be a $\mathcal{N}(\mathbf{0}, \mathbf{I}_p)$ random variable. Then $\xi$ and $\mathbf{U}\boldsymbol{\Lambda}^{1/2}\gamma$ are identically distributed, so

$$\mathbb{E}(\xi\xi^*)^m = \mathbb{E}\left[(\xi^*\xi)^{m-1}\xi\xi^*\right] = \mathbb{E}\left[(\gamma^*\boldsymbol{\Lambda}\gamma)^{m-1}\mathbf{U}\boldsymbol{\Lambda}^{1/2}\gamma\gamma^*\boldsymbol{\Lambda}^{1/2}\mathbf{U}^*\right]$$

$$= \mathbf{U}\boldsymbol{\Lambda}^{1/2}\mathbb{E}\left[(\gamma^*\boldsymbol{\Lambda}\gamma)^{m-1}\gamma\gamma^*\right]\boldsymbol{\Lambda}^{1/2}\mathbf{U}^*. \tag{2.8.6}$$

Consider the $(i, j)$ entry of the bracketed matrix in (2.8.6):

$$\mathbb{E}\left[(\gamma^*\boldsymbol{\Lambda}\gamma)^{m-1}\gamma_i\gamma_j\right] = \mathbb{E}\left[\left(\sum_{\ell=1}^p \lambda_\ell\gamma_\ell^2\right)^{m-1}\gamma_i\gamma_j\right]. \tag{2.8.7}$$

From this expression, and the independence and symmetry of the Gaussian variables $\{\gamma_i\}$, we see that this matrix is diagonal.

To bound the diagonal entries, use a multinomial expansion to further develop the sum in

(2.8.7) for the $(i, i)$ entry:

$$\mathbb{E}\left[(\gamma^*\Lambda\gamma)^{m-1}\gamma_i^2\right] = \sum_{\ell_1+\cdots+\ell_p=m-1} \binom{m-1}{\ell_1,\ldots,\ell_p} \lambda_1^{\ell_1}\cdots\lambda_p^{\ell_p} \mathbb{E}\left[\gamma_1^{2\ell_1}\cdots\gamma_p^{2\ell_p}\gamma_i^2\right].$$

Now we use the generalized AM–GM inequality to replace the expectation of the product of Gaussians with the $2m$th moment of a single standard Gaussian $g$. Denote the $L_r$ norm of a random variable $X$ by

$$\|X\|_{L_r} = (\mathbb{E}|X|^r)^{1/r}.$$

Since $\ell_1,\ldots,\ell_p$ are nonnegative integers summing to $m-1$, the generalized AM-GM inequality justifies the first of the following inequalities:

$$\mathbb{E}\gamma_1^{2\ell_1}\cdots\gamma_p^{2\ell_p}\gamma_i^2 \le \mathbb{E}\left(\frac{\ell_1|\gamma_1|+\cdots+\ell_p|\gamma_p|+|\gamma_i|}{m}\right)^{2m} = \left\|\frac{1}{m}\left(|\gamma_i|+\sum_{j=1}^p \ell_j|\gamma_j|\right)\right\|_{L_{2m}}^{2m}$$

$$\le \left(\frac{1}{m}\left(\|\gamma_i\|_{L_{2m}}+\sum_{j=1}^p \ell_j\|\gamma_j\|_{L_{2m}}\right)\right)^{2m}$$

$$= \left(\frac{1+\ell_1+\ldots+\ell_p}{m}\right)^{2m}\|g\|_{L_{2m}}^{2m} = \mathbb{E}(g^{2m}).$$

The second inequality is the triangle inequality for $L_r$ norms. Now we reverse the multinomial expansion to see that the diagonal terms satisfy the inequality

$$\mathbb{E}\left[(\gamma^*\Lambda\gamma)^{m-1}\gamma_i^2\right] \le \sum_{\ell_1+\cdots+\ell_p=m-1} \binom{m-1}{\ell_1,\ldots,\ell_p} \lambda_1^{\ell_1}\cdots\lambda_p^{\ell_p}\mathbb{E}(g^{2m})$$

$$= (\lambda_1+\ldots+\lambda_p)^{m-1}\mathbb{E}(g^{2m}) = \operatorname{tr}(\mathbf{G})^{m-1}\mathbb{E}(g^{2m}). \qquad (2.8.8)$$

Estimate $\mathbb{E}(g^{2m})$ using the fact that $\Gamma(x)$ is increasing for $x \geq 1$:

$$\mathbb{E}\left(g^{2m}\right) = \frac{2^m}{\sqrt{\pi}}\Gamma(m+1/2) < \frac{2^m}{\sqrt{\pi}}\Gamma(m+1) = \frac{2^m}{\sqrt{\pi}}m! \quad \text{for } m \geq 1.$$

Combine this result with (2.8.8) to see that

$$\mathbb{E}\left[(\gamma^*\Lambda\gamma)^{m-1}\gamma\gamma^*\right] \preceq \frac{2^m}{\sqrt{\pi}}m!\,\text{tr}(\mathbf{G})^{m-1}\cdot\mathbf{I}.$$

Complete the proof by using this estimate in (2.8.6). $\quad\square$