

1.3. Eigenvalues and sums of Hermitian matrices

Let A be a Hermitian $n \times n$ matrix. By the *spectral theorem* for Hermitian matrices (which, for sake of completeness, we prove below), one can diagonalise A using a sequence¹¹

$$\lambda_1(A) \geq \dots \geq \lambda_n(A)$$

of n real eigenvalues, together with an orthonormal basis of eigenvectors $u_1(A), \dots, u_n(A) \in \mathbf{C}^n$. The set $\{\lambda_1(A), \dots, \lambda_n(A)\}$ is known as the *spectrum* of A .

A basic question in linear algebra asks the extent to which the eigenvalues $\lambda_1(A), \dots, \lambda_n(A)$ and $\lambda_1(B), \dots, \lambda_n(B)$ of two Hermitian matrices A, B constrain the eigenvalues $\lambda_1(A+B), \dots, \lambda_n(A+B)$ of the sum. For instance, the linearity of trace

$$\operatorname{tr}(A+B) = \operatorname{tr}(A) + \operatorname{tr}(B),$$

when expressed in terms of eigenvalues, gives the trace constraint

$$(1.52) \quad \lambda_1(A+B) + \dots + \lambda_n(A+B) = \lambda_1(A) + \dots + \lambda_n(A) \\ + \lambda_1(B) + \dots + \lambda_n(B);$$

the identity

$$(1.53) \quad \lambda_1(A) = \sup_{|v|=1} v^*Av$$

(together with the counterparts for B and $A+B$) gives the inequality

$$(1.54) \quad \lambda_1(A+B) \leq \lambda_1(A) + \lambda_1(B),$$

and so forth.

The complete answer to this problem is a fascinating one, requiring a strangely recursive description (once known as *Horn's conjecture*, which is now solved), and connected to a large number of other fields of mathematics, such as geometric invariant theory, intersection theory, and the combinatorics of a certain gadget known as a “honeycomb”. See [KnTa2001] for a survey of this topic.

In typical applications to random matrices, one of the matrices (say, B) is “small” in some sense, so that $A+B$ is a perturbation of A . In this case, one does not need the full strength of the above theory, and instead relies on a simple aspect of it pointed out in [HeRo1995], [To1994], which generates several of the *eigenvalue inequalities* relating A , B , and $A+B$, of which

¹¹The eigenvalues are uniquely determined by A , but the eigenvectors have a little ambiguity to them, particularly if there are repeated eigenvalues; for instance, one could multiply each eigenvector by a complex phase $e^{i\theta}$. In this text we are arranging eigenvalues in descending order; of course, one can also arrange eigenvalues in increasing order, which causes some slight notational changes in the results below.

(1.52) and (1.54) are examples¹². These eigenvalue inequalities can mostly be deduced from a number of *minimax* characterisations of eigenvalues (of which (1.53) is a typical example), together with some basic facts about intersections of subspaces. Examples include the *Weyl inequalities*

$$(1.55) \quad \lambda_{i+j-1}(A+B) \leq \lambda_i(A) + \lambda_j(B),$$

valid whenever $i, j \geq 1$ and $i+j-1 \leq n$, and the *Ky Fan inequality*

$$(1.56) \quad \lambda_1(A+B) + \cdots + \lambda_k(A+B) \leq \lambda_1(A) + \cdots + \lambda_k(A) + \lambda_1(B) + \cdots + \lambda_k(B).$$

One consequence of these inequalities is that the spectrum of a Hermitian matrix is *stable* with respect to small perturbations.

We will also establish some closely related inequalities concerning the relationships between the eigenvalues of a matrix, and the eigenvalues of its minors.

Many of the inequalities here have analogues for the singular values of non-Hermitian matrices (by exploiting the augmented matrix (2.80)). However, the situation is markedly different when dealing with *eigenvalues* of *non-Hermitian* matrices; here, the spectrum can be far more unstable, if *pseudospectrum* is present. Because of this, the theory of the eigenvalues of a random non-Hermitian matrix requires an additional ingredient, namely upper bounds on the prevalence of pseudospectrum, which after recentering the matrix is basically equivalent to establishing lower bounds on least singular values. See Section 2.8.1 for further discussion of this point.

We will work primarily here with Hermitian matrices, which can be viewed as self-adjoint transformations on complex vector spaces such as \mathbf{C}^n . One can of course specialise the discussion to real symmetric matrices, in which case one can restrict these complex vector spaces to their real counterparts \mathbf{R}^n . The specialisation of the complex theory below to the real case is straightforward and is left to the interested reader.

1.3.1. Proof of spectral theorem. To prove the spectral theorem, it is convenient to work more abstractly, in the context of self-adjoint operators on finite-dimensional Hilbert spaces:

Theorem 1.3.1 (Spectral theorem). *Let V be a finite-dimensional complex Hilbert space of some dimension n , and let $T : V \rightarrow V$ be a self-adjoint operator. Then there exists an orthonormal basis $v_1, \dots, v_n \in V$ of V and eigenvalues $\lambda_1, \dots, \lambda_n \in \mathbf{R}$ such that $Tv_i = \lambda_i v_i$ for all $1 \leq i \leq n$.*

¹²Actually, this method eventually generates *all* of the eigenvalue inequalities, but this is a non-trivial fact to prove; see [KnTaWo2004]

The spectral theorem as stated in the introduction then follows by specialising to the case $V = \mathbf{C}^n$ and ordering the eigenvalues.

Proof. We induct on the dimension n . The claim is vacuous for $n = 0$, so suppose that $n \geq 1$ and that the claim has already been proven for $n = 1$.

Let v be a unit vector in V (thus $v^*v = 1$) that maximises the form $\operatorname{Re}(v^*Tv)$; this maximum exists by compactness. By the method of Lagrange multipliers, v is a critical point of $\operatorname{Re}(v^*Tv) - \lambda v^*v$ for some $\lambda \in \mathbf{R}$. Differentiating in an arbitrary direction $w \in V$, we conclude that

$$\operatorname{Re}(v^*Tw + w^*Tv - \lambda v^*w - \lambda w^*v) = 0;$$

this simplifies using self-adjointness to

$$\operatorname{Re}(w^*(Tv - \lambda v)) = 0.$$

Since $w \in V$ was arbitrary, we conclude that $Tv = \lambda v$, thus v is a unit eigenvector of T . By self-adjointness, this implies that the orthogonal complement $v^\perp := \{w \in V : v^*w = 0\}$ of v is preserved by T . Restricting T to this lower-dimensional subspace and applying the induction hypothesis, we can find an orthonormal basis of eigenvectors of T on v^\perp . Adjoining the new unit vector v to the orthonormal basis, we obtain the claim. \square

Suppose we have a self-adjoint transformation $A : \mathbf{C}^n \rightarrow \mathbf{C}^n$, which of course can be identified with a Hermitian matrix. Using the orthogonal eigenbasis provided by the spectral theorem, we can perform an orthonormal change of variables to set that eigenbasis to be the standard basis e_1, \dots, e_n , so that the matrix of A becomes diagonal. This is very useful when dealing with just a single matrix A ; for instance, it makes the task of computing functions of A , such as A^k or $\exp(tA)$, much easier. However, when one has *several* Hermitian matrices in play (e.g., $A, B, A + B$), then it is usually not possible to standardise all the eigenbases simultaneously (i.e., to simultaneously diagonalise all the matrices), except when the matrices all commute. Nevertheless, one can still normalise *one* of the eigenbases to be the standard basis, and this is still useful for several applications, as we shall soon see.

Exercise 1.3.1. Suppose that the eigenvalues $\lambda_1(A) > \dots > \lambda_n(A)$ of an $n \times n$ Hermitian matrix are distinct. Show that the associated eigenbasis $u_1(A), \dots, u_n(A)$ is unique up to rotating each individual eigenvector $u_j(A)$ by a complex phase $e^{i\theta_j}$. In particular, the *spectral projections* $P_j(A) := u_j(A)^*u_j(A)$ are unique. What happens when there is eigenvalue multiplicity?

1.3.2. Minimax formulae. The i^{th} eigenvalue functional $A \mapsto \lambda_i(A)$ is not a linear functional (except in dimension one). It is not even a convex functional (except when $i = 1$) or a concave functional (except when $i = n$). However, it is the next best thing, namely it is a *minimax* expression of linear functionals¹³. More precisely, we have

Theorem 1.3.2 (Courant-Fischer minimax theorem). *Let A be an $n \times n$ Hermitian matrix. Then we have*

$$(1.57) \quad \lambda_i(A) = \sup_{\dim(V)=i} \inf_{v \in V: |v|=1} v^* Av$$

and

$$(1.58) \quad \lambda_i(A) = \inf_{\dim(V)=n-i+1} \sup_{v \in V: |v|=1} v^* Av$$

for all $1 \leq i \leq n$, where V ranges over all subspaces of \mathbf{C}^n with the indicated dimension.

Proof. It suffices to prove (1.57), as (1.58) follows by replacing A by $-A$ (noting that $\lambda_i(-A) = -\lambda_{n-i+1}(A)$).

We first verify the $i = 1$ case, i.e., (1.53). By the spectral theorem, we can assume that A has the standard eigenbasis e_1, \dots, e_n , in which case we have

$$(1.59) \quad v^* Av = \sum_{i=1}^n \lambda_i |v_i|^2$$

whenever $v = (v_1, \dots, v_n)$. The claim (1.53) is then easily verified.

To prove the general case, we may again assume A has the standard eigenbasis. By considering the space V spanned by e_1, \dots, e_i , we easily see the inequality

$$\lambda_i(A) \leq \sup_{\dim(V)=i} \inf_{v \in V: |v|=1} v^* Av,$$

so we only need to prove the reverse inequality. In other words, for every i -dimensional subspace V of \mathbf{C}^n , we have to show that V contains a unit vector v such that

$$v^* Av \leq \lambda_i(A).$$

Let W be the space spanned by e_i, \dots, e_n . This space has codimension $i - 1$, so it must have non-trivial intersection with V . If we let v be a unit vector in $V \cap W$, the claim then follows from (1.59). \square

¹³Note that a convex functional is the same thing as a max of linear functionals, while a concave functional is the same thing as a min of linear functionals.

Remark 1.3.3. By homogeneity, one can replace the restriction $|v| = 1$ with $v \neq 0$ provided that one replaces the quadratic form v^*Av with the Rayleigh quotient v^*Av/v^*v .

A closely related formula is as follows. Given an $n \times n$ Hermitian matrix A and an m -dimensional subspace V of \mathbf{C}^n , we define the *partial trace* $\text{tr}(A \lfloor_V)$ to be the expression

$$\text{tr}(A \lfloor_V) := \sum_{i=1}^m v_i^* A v_i$$

where v_1, \dots, v_m is any orthonormal basis of V . It is easy to see that this expression is independent of the choice of orthonormal basis, and so the partial trace is well-defined.

Proposition 1.3.4 (Extremal partial trace). *Let A be an $n \times n$ Hermitian matrix. Then for any $1 \leq k \leq n$, one has*

$$\lambda_1(A) + \dots + \lambda_k(A) = \sup_{\dim(V)=k} \text{tr}(A \lfloor_V)$$

and

$$\lambda_{n-k+1}(A) + \dots + \lambda_n(A) = \inf_{\dim(V)=k} \text{tr}(A \lfloor_V).$$

As a corollary, we see that $A \mapsto \lambda_1(A) + \dots + \lambda_k(A)$ is a convex function, and $A \mapsto \lambda_{n-k+1}(A) + \dots + \lambda_n(A)$ is a concave function.

Proof. Again, by symmetry it suffices to prove the first formula. As before, we may assume, without loss of generality, that A has the standard eigenbasis e_1, \dots, e_n corresponding to $\lambda_1(A), \dots, \lambda_n(A)$, respectively. By selecting V to be the span of e_1, \dots, e_k we have the inequality

$$\lambda_1(A) + \dots + \lambda_k(A) \leq \sup_{\dim(V)=k} \text{tr}(A \lfloor_V),$$

so it suffices to prove the reverse inequality. For this we induct on the dimension n . If V has dimension k , then it has a $k-1$ -dimensional subspace V' that is contained in the span of e_2, \dots, e_n . By the induction hypothesis applied to the restriction of A to this span (which has eigenvalues $\lambda_2(A), \dots, \lambda_n(A)$), we have

$$\lambda_2(A) + \dots + \lambda_k(A) \geq \text{tr}(A \lfloor_{V'}).$$

On the other hand, if v is a unit vector in the orthogonal complement of V' in V , we see from (1.53) that

$$\lambda_1(A) \geq v^* A v.$$

Adding the two inequalities we obtain the claim. \square

Specialising Proposition 1.3.4 to the case when V is a coordinate subspace (i.e., the span of k of the basis vectors e_1, \dots, e_n), we conclude the *Schur-Horn inequalities*

$$(1.60) \quad \begin{aligned} \lambda_{n-k+1}(A) + \dots + \lambda_n(A) &\leq a_{i_1 i_1} + \dots + a_{i_k i_k} \\ &\leq \lambda_1(A) + \dots + \lambda_k(A) \end{aligned}$$

for any $1 \leq i_1 < \dots < i_k \leq n$, where $a_{11}, a_{22}, \dots, a_{nn}$ are the diagonal entries of A .

Exercise 1.3.2. Show that the inequalities (1.60) are equivalent to the assertion that the diagonal entries $\text{diag}(A) = (a_{11}, a_{22}, \dots, a_{nn})$ lies in the *permutahedron* of $\lambda_1(A), \dots, \lambda_n(A)$, defined as the convex hull of the $n!$ permutations of $(\lambda_1(A), \dots, \lambda_n(A))$ in \mathbf{R}^n .

Remark 1.3.5. It is a theorem of Schur and Horn [Ho1954] that these are the complete set of inequalities connecting the diagonal entries $\text{diag}(A) = (a_{11}, a_{22}, \dots, a_{nn})$ of a Hermitian matrix to its spectrum. To put it another way, the image of any *coadjoint orbit* $\mathcal{O}_A := \{UAU^* : U \in U(n)\}$ of a matrix A with a given spectrum $\lambda_1, \dots, \lambda_n$ under the diagonal map $\text{diag} : A \mapsto \text{diag}(A)$ is the permutahedron of $\lambda_1, \dots, \lambda_n$. Note that the vertices of this permutahedron can be attained by considering the diagonal matrices inside this coadjoint orbit, whose entries are then a permutation of the eigenvalues. One can interpret this diagonal map diag as the *moment map* associated with the conjugation action of the standard maximal torus of $U(n)$ (i.e., the diagonal unitary matrices) on the coadjoint orbit. When viewed in this fashion, the Schur-Horn theorem can be viewed as the special case of the more general *Atiyah convexity theorem* [At1982] (also proven independently by Guillemin and Sternberg [GuSt1982]) in symplectic geometry. Indeed, the topic of eigenvalues of Hermitian matrices turns out to be quite profitably viewed as a question in symplectic geometry (and also in algebraic geometry, particularly when viewed through the machinery of *geometric invariant theory*).

There is a simultaneous generalisation of Theorem 1.3.2 and Proposition 1.3.4:

Exercise 1.3.3 (Wielandt minimax formula). Let $1 \leq i_1 < \dots < i_k \leq n$ be integers. Define a *partial flag* to be a nested collection $V_1 \subset \dots \subset V_k$ of subspaces of \mathbf{C}^n such that $\dim(V_j) = i_j$ for all $1 \leq j \leq k$. Define the associated *Schubert variety* $X(V_1, \dots, V_k)$ to be the collection of all k -dimensional subspaces W such that $\dim(W \cap V_j) \geq j$. Show that for any $n \times n$ matrix A ,

$$\lambda_{i_1}(A) + \dots + \lambda_{i_k}(A) = \sup_{V_1, \dots, V_k} \inf_{W \in X(V_1, \dots, V_k)} \text{tr}(A|_W).$$

1.3.3. Eigenvalue inequalities. Using the above minimax formulae, we can now quickly prove a variety of eigenvalue inequalities. The basic idea is to exploit the linearity relationship

$$(1.61) \quad v^*(A + B)v = v^*Av + v^*Bv$$

for any unit vector v , and more generally,

$$(1.62) \quad \operatorname{tr}((A + B) \downarrow_V) = \operatorname{tr}(A \downarrow_V) + \operatorname{tr}(B \downarrow_V)$$

for any subspace V .

For instance, as mentioned before, the inequality (1.54) follows immediately from (1.53) and (1.61). Similarly, for the Ky Fan inequality (1.56), one observes from (1.62) and Proposition 1.3.4 that

$$\operatorname{tr}((A + B) \downarrow_W) \leq \operatorname{tr}(A \downarrow_W) + \lambda_1(B) + \cdots + \lambda_k(B)$$

for any k -dimensional subspace W . Substituting this into Proposition 1.3.4 gives the claim. If one uses Exercise 1.3.3 instead of Proposition 1.3.4, one obtains the more general *Lidskii inequality*

$$(1.63) \quad \begin{aligned} & \lambda_{i_1}(A + B) + \cdots + \lambda_{i_k}(A + B) \\ & \leq \lambda_{i_1}(A) + \cdots + \lambda_{i_k}(A) + \lambda_1(B) + \cdots + \lambda_k(B) \end{aligned}$$

for any $1 \leq i_1 < \cdots < i_k \leq n$.

In a similar spirit, using the inequality

$$|v^*Bv| \leq \|B\|_{\text{op}} = \max(|\lambda_1(B)|, |\lambda_n(B)|)$$

for unit vectors v , combined with (1.61) and (1.57), we obtain the eigenvalue stability inequality

$$(1.64) \quad |\lambda_i(A + B) - \lambda_i(A)| \leq \|B\|_{\text{op}},$$

thus the spectrum of $A + B$ is close to that of A if B is small in operator norm. In particular, we see that the map $A \mapsto \lambda_i(A)$ is Lipschitz continuous on the space of Hermitian matrices, for fixed $1 \leq i \leq n$.

More generally, suppose one wants to establish the Weyl inequality (1.55). From (1.57) that it suffices to show that every $i + j - 1$ -dimensional subspace V contains a unit vector v such that

$$v^*(A + B)v \leq \lambda_i(A) + \lambda_j(B).$$

But from (1.57), one can find a subspace U of codimension $i - 1$ such that $v^*Av \leq \lambda_i(A)$ for all unit vectors v in U , and a subspace W of codimension $j - 1$ such that $v^*Bv \leq \lambda_j(B)$ for all unit vectors v in W . The intersection $U \cap W$ has codimension at most $i + j - 2$ and so has a non-trivial intersection with V ; and the claim follows.

Remark 1.3.6. More generally, one can generate an eigenvalue inequality whenever the intersection numbers of three Schubert varieties of compatible dimensions is non-zero; see [HeRo1995]. In fact, this generates a complete set of inequalities; see [Klyachko]. One can in fact restrict attention to those varieties whose intersection number is exactly one; see [KnTaWo2004]. Finally, in those cases, the fact that the intersection is one can be proven by entirely elementary means (based on the standard inequalities relating the dimension of two subspaces V, W to their intersection $V \cap W$ and sum $V + W$); see [BeCoDyLiTi2010]. As a consequence, the methods in this section can, in principle, be used to derive all possible eigenvalue inequalities for sums of Hermitian matrices.

Exercise 1.3.4. Verify the inequalities (1.63) and (1.55) by hand in the case when A and B commute (and are thus simultaneously diagonalisable), without the use of minimax formulae.

Exercise 1.3.5. Establish the dual Lidskii inequality

$$\lambda_{i_1}(A+B) + \cdots + \lambda_{i_k}(A+B) \geq \lambda_{i_1}(A) + \cdots + \lambda_{i_k}(A) \\ + \lambda_{n-k+1}(B) + \cdots + \lambda_n(B)$$

for any $1 \leq i_1 < \cdots < i_k \leq n$ and the dual Weyl inequality

$$\lambda_{i+j-n}(A+B) \geq \lambda_i(A) + \lambda_j(B)$$

whenever $1 \leq i, j, i+j-n \leq n$.

Exercise 1.3.6. Use the Lidskii inequality to establish the more general inequality

$$\sum_{i=1}^n c_i \lambda_i(A+B) \leq \sum_{i=1}^n c_i \lambda_i(A) + \sum_{i=1}^n c_i^* \lambda_i(B)$$

whenever $c_1, \dots, c_n \geq 0$, and $c_1^* \geq \cdots \geq c_n^* \geq 0$ is the decreasing rearrangement of c_1, \dots, c_n . (*Hint:* Express c_i as the integral of $\mathbf{I}(c_i \geq \lambda)$ as λ runs from 0 to infinity. For each fixed λ , apply (1.63).) Combine this with Hölder's inequality to conclude the *p-Weilandt-Hoffman inequality*

$$(1.65) \quad \|(\lambda_i(A+B) - \lambda_i(A))_{i=1}^n\|_{\ell_n^p} \leq \|B\|_{S^p}$$

for any $1 \leq p \leq \infty$, where

$$\|(a_i)_{i=1}^n\|_{\ell_n^p} := \left(\sum_{i=1}^n |a_i|^p \right)^{1/p}$$

is the usual ℓ^p norm (with the usual convention that $\|(a_i)_{i=1}^n\|_{\ell_n^\infty} := \sup_{1 \leq i \leq n} |a_i|$), and

$$(1.66) \quad \|B\|_{S^p} := \|(\lambda_i(B))_{i=1}^n\|_{\ell_n^p}$$

is the *p-Schatten norm* of B .

Exercise 1.3.7. Show that the p -Schatten norms are indeed norms on the space of Hermitian matrices for every $1 \leq p \leq \infty$.

Exercise 1.3.8. Show that for any $1 \leq p \leq \infty$ and any Hermitian matrix $A = (a_{ij})_{1 \leq i, j \leq n}$, one has

$$(1.67) \quad \|(a_{ii})_{i=1}^n\|_{\ell_n^p} \leq \|A\|_{S^p}.$$

Exercise 1.3.9. Establish the *non-commutative Hölder inequality*

$$|\operatorname{tr}(AB)| \leq \|A\|_{S^p} \|B\|_{S^{p'}}$$

whenever $1 \leq p, p' \leq \infty$ with $1/p + 1/p' = 1$, and A, B are $n \times n$ Hermitian matrices. (*Hint:* Diagonalise one of the matrices and use the preceding exercise.)

The most important¹⁴ p -Schatten norms are the ∞ -Schatten norm $\|A\|_{S^\infty} = \|A\|_{\text{op}}$, which is just the operator norm, and the 2-Schatten norm $\|A\|_{S^2} = (\sum_{i=1}^n \lambda_i(A)^2)^{1/2}$, which is also the *Frobenius norm* (or *Hilbert-Schmidt norm*)

$$\|A\|_{S^2} = \|A\|_F := \operatorname{tr}(AA^*)^{1/2} = \left(\sum_{i=1}^n \sum_{j=1}^n |a_{ij}|^2 \right)^{1/2}$$

where a_{ij} are the coefficients of A . Thus we see that the $p = 2$ case of the Weilandt-Hoffman inequality can be written as

$$(1.68) \quad \sum_{i=1}^n |\lambda_i(A+B) - \lambda_i(A)|^2 \leq \|B\|_F^2.$$

We will give an alternate proof of this inequality, based on eigenvalue deformation, in the next section.

1.3.4. Eigenvalue deformation. From the Weyl inequality (1.64), we know that the eigenvalue maps $A \mapsto \lambda_i(A)$ are Lipschitz continuous on Hermitian matrices (and thus also on real symmetric matrices). It turns out that we can obtain better regularity, provided that we avoid repeated eigenvalues. Fortunately, repeated eigenvalues are rare:

Exercise 1.3.10 (Dimension count). Suppose that $n \geq 2$. Show that the space of Hermitian matrices with at least one repeated eigenvalue has codimension 3 in the space of all Hermitian matrices, and the space of real symmetric matrices with at least one repeated eigenvalue has codimension 2 in the space of all real symmetric matrices. (When $n = 1$, repeated eigenvalues of course do not occur.)

¹⁴The 1-Schatten norm S^1 , also known as the *nuclear norm* or *trace class norm*, is important in a number of applications, such as matrix completion, but will not be used in this text.

Let us say that a Hermitian matrix has *simple spectrum* if it has no repeated eigenvalues. We thus see from the above exercise and (1.64) that the set of Hermitian matrices with simple spectrum forms an open dense set in the space of all Hermitian matrices, and similarly for real symmetric matrices; thus simple spectrum is the *generic* behaviour of such matrices. Indeed, the unexpectedly high codimension of the non-simple matrices (naively, one would expect a codimension 1 set for a collision between, say, $\lambda_i(A)$ and $\lambda_{i+1}(A)$) suggests a *repulsion* phenomenon: because it is unexpectedly rare for eigenvalues to be equal, there must be some “force” that “repels” eigenvalues of Hermitian (and to a lesser extent, real symmetric) matrices from getting too close to each other. We now develop some machinery to make this intuition more precise.

We first observe that when A has simple spectrum, the zeroes of the characteristic polynomial $\lambda \mapsto \det(A - \lambda I)$ are simple (i.e., the polynomial has nonzero derivative at those zeroes). From this and the inverse function theorem, we see that each of the eigenvalue maps $A \mapsto \lambda_i(A)$ are smooth on the region where A has simple spectrum. Because the eigenvectors $u_i(A)$ are determined (up to phase) by the equations $(A - \lambda_i(A)I)u_i(A) = 0$ and $u_i(A)^*u_i(A) = 1$, another application of the inverse function theorem tells us that we can (locally¹⁵) select the maps $A \mapsto u_i(A)$ to also be smooth.

Now suppose that $A = A(t)$ depends smoothly on a time variable t , so that (when A has simple spectrum) the eigenvalues $\lambda_i(t) = \lambda_i(A(t))$ and eigenvectors $u_i(t) = u_i(A(t))$ also depend smoothly on t . We can then differentiate the equations

$$(1.69) \quad Au_i = \lambda_i u_i$$

and

$$(1.70) \quad u_i^* u_i = 1$$

to obtain various equations of motion for λ_i and u_i in terms of the derivatives of A .

Let’s see how this works. Taking first derivatives of (1.69), (1.70) using the product rule, we obtain

$$(1.71) \quad \dot{A}u_i + A\dot{u}_i = \dot{\lambda}_i u_i + \lambda_i \dot{u}_i$$

and

$$(1.72) \quad \dot{u}_i^* u_i + u_i^* \dot{u}_i = 0.$$

¹⁵There may be topological obstructions to smoothly selecting these vectors globally, but this will not concern us here as we will be performing a local analysis only. In some applications, it is more convenient not to work with the $u_i(A)$ at all due to their phase ambiguity, and work instead with the *spectral projections* $P_i(A) := u_i(A)u_i(A)^*$, which do not have this ambiguity.

The equation (1.72) simplifies to $\dot{u}_i^* u_i = 0$, thus \dot{u}_i is orthogonal to u_i . Taking inner products of (1.71) with u_i , we conclude the *Hadamard first variation formula*

$$(1.73) \quad \dot{\lambda}_i = u_i^* \dot{A} u_i.$$

This can already be used to give alternate proofs of various eigenvalue identities. For instance, if we apply this to $A(t) := A + tB$, we see that

$$\frac{d}{dt} \lambda_i(A + tB) = u_i(A + tB)^* B u_i(A + tB)$$

whenever $A + tB$ has simple spectrum. The right-hand side can be bounded in magnitude by $\|B\|_{\text{op}}$, and so we see that the map $t \mapsto \lambda_i(A + tB)$ is Lipschitz continuous, with Lipschitz constant $\|B\|_{\text{op}}$ whenever $A + tB$ has simple spectrum, which happens for generic A, B (and all t) by Exercise 1.3.10. By the fundamental theorem of calculus, we thus conclude (1.64).

Exercise 1.3.11. Use a similar argument to the one above to establish (1.68) without using minimax formulae or Lidskii's inequality.

Exercise 1.3.12. Use a similar argument to the one above to deduce Lidskii's inequality (1.63) from Proposition 1.3.4 rather than Exercise 1.3.3.

One can also compute the second derivative of eigenvalues:

Exercise 1.3.13. Suppose that $A = A(t)$ depends smoothly on t . By differentiating (1.71) and (1.72), establish the *Hadamard second variation formula*¹⁶

$$(1.74) \quad \frac{d^2}{dt^2} \lambda_k = u_k^* \ddot{A} u_k + 2 \sum_{j \neq k} \frac{|u_j^* \dot{A} u_k|^2}{\lambda_k - \lambda_j}$$

whenever A has simple spectrum and $1 \leq k \leq n$.

Remark 1.3.7. In the proof of the *four moment theorem* [TaVu2009b] on the fine spacing of Wigner matrices, one also needs the variation formulae for the third, fourth, and fifth derivatives of the eigenvalues (the first four derivatives match up with the four moments mentioned in the theorem, and the fifth derivative is needed to control error terms). Fortunately, one does not need the precise formulae for these derivatives (which, as one can imagine, are quite complicated), but only their general form, and in particular, an upper bound for these derivatives in terms of more easily computable quantities.

¹⁶If one interprets the second derivative of the eigenvalues as being proportional to a “force” on those eigenvalues (in analogy with *Newton's second law*), (1.74) is asserting that each eigenvalue λ_j “repels” the other eigenvalues λ_k by exerting a force that is inversely proportional to their separation (and also proportional to the square of the matrix coefficient of \dot{A} in the eigenbasis). See [Ta2009b, §1.5] for more discussion.

1.3.5. Minors. In the previous sections, we perturbed $n \times n$ Hermitian matrices $A = A_n$ by adding a (small) $n \times n$ Hermitian correction matrix B to them to form a new $n \times n$ Hermitian matrix $A + B$. Another important way to perturb a matrix is to pass to a *principal minor*, for instance to the top left $n - 1 \times n - 1$ minor A_{n-1} of A_n . There is an important relationship between the eigenvalues of the two matrices:

Exercise 1.3.14 (Cauchy interlacing law). For any $n \times n$ Hermitian matrix A_n with top left $n - 1 \times n - 1$ minor A_{n-1} , then

$$(1.75) \quad \lambda_{i+1}(A_n) \leq \lambda_i(A_{n-1}) \leq \lambda_i(A_n)$$

for all $1 \leq i < n$. (*Hint:* Use the Courant-Fischer minimax theorem, Theorem 1.3.2.) Show furthermore that the space of A_n for which equality holds in one of the inequalities in (1.75) has codimension 2 (for Hermitian matrices) or 1 (for real symmetric matrices).

Remark 1.3.8. If one takes successive minors $A_{n-1}, A_{n-2}, \dots, A_1$ of an $n \times n$ Hermitian matrix A_n , and computes their spectra, then (1.75) shows that this triangular array of numbers forms a pattern known as a *Gelfand-Tsetlin pattern*.

One can obtain a more precise formula for the eigenvalues of A_n in terms of those for A_{n-1} :

Exercise 1.3.15 (Eigenvalue equation). Let A_n be an $n \times n$ Hermitian matrix with top left $n - 1 \times n - 1$ minor A_{n-1} . Suppose that λ is an eigenvalue of A_n distinct from all the eigenvalues of A_{n-1} (and thus simple, by (1.75)). Show that

$$(1.76) \quad \sum_{j=1}^{n-1} \frac{|u_j(A_{n-1})^* X|^2}{\lambda_j(A_{n-1}) - \lambda} = a_{nn} - \lambda$$

where a_{nn} is the bottom right entry of A , and $X = (a_{nj})_{j=1}^{n-1} \in \mathbf{C}^{n-1}$ is the right column of A (minus the bottom entry). (*Hint:* Expand out the eigenvalue equation $A_n u = \lambda u$ into the \mathbf{C}^{n-1} and \mathbf{C} components.) Note the similarities between (1.76) and (1.74).

Observe that the function $\lambda \rightarrow \sum_{j=1}^{n-1} \frac{|u_j(A_{n-1})^* X|^2}{\lambda_j(A_{n-1}) - \lambda}$ is a rational function of λ which is increasing away from the eigenvalues of A_{n-1} , where it has a pole (except in the rare case when the inner product $u_{j-1}(A_{n-1})^* X$ vanishes, in which case it can have a removable singularity). By graphing this function one can see that the interlacing formula (1.75) can also be interpreted as a manifestation of the intermediate value theorem.

The identity (1.76) suggests that under typical circumstances, an eigenvalue λ of A_n can only get close to an eigenvalue $\lambda_j(A_{n-1})$ if the associated

inner product $u_j(A_{n-1})^*X$ is small. This type of observation is useful to achieve *eigenvalue repulsion*—to show that it is unlikely that the gap between two adjacent eigenvalues is small. We shall see examples of this in later sections.

1.3.6. Singular values. The theory of eigenvalues of $n \times n$ Hermitian matrices has an analogue in the theory of singular values of $p \times n$ non-Hermitian matrices. We first begin with the counterpart to the spectral theorem, namely the *singular value decomposition*.

Theorem 1.3.9 (Singular value decomposition). *Let $0 \leq p \leq n$, and let A be a linear transformation from an n -dimensional complex Hilbert space U to a p -dimensional complex Hilbert space V . (In particular, A could be an $p \times n$ matrix with complex entries, viewed as a linear transformation from \mathbf{C}^n to \mathbf{C}^p .) Then there exist non-negative real numbers*

$$\sigma_1(A) \geq \cdots \geq \sigma_p(A) \geq 0$$

(known as the singular values of A) and orthonormal sets $u_1(A), \dots, u_p(A) \in U$ and $v_1(A), \dots, v_p(A) \in V$ (known as singular vectors of A), such that

$$Au_j = \sigma_j v_j; \quad A^*v_j = \sigma_j u_j$$

for all $1 \leq j \leq p$, where we abbreviate $u_j = u_j(A)$, etc.

Furthermore, $Au = 0$ whenever u is orthogonal to all $u_1(A), \dots, u_p(A)$.

We adopt the convention that $\sigma_i(A) = 0$ for $i > p$. The above theorem only applies to matrices with at least as many rows as columns, but one can also extend the definition to matrices with more columns than rows by adopting the convention $\sigma_i(A^*) := \sigma_i(A)$ (it is easy to check that this extension is consistent on square matrices). All of the results below extend (with minor modifications) to the case when there are more columns than rows, but we have not displayed those extensions here in order to simplify the notation.

Proof. We induct on p . The claim is vacuous for $p = 0$, so suppose that $p \geq 1$ and that the claim has already been proven for $p - 1$.

We follow a similar strategy to the proof of Theorem 1.3.1. We may assume that A is not identically zero, as the claim is obvious otherwise. The function $u \mapsto \|Au\|^2$ is continuous on the unit sphere of U , so there exists a unit vector u_1 which maximises this quantity. If we set $\sigma_1 := \|Au_1\| > 0$, one easily verifies that u_1 is a critical point of the map $u \mapsto \|Au\|^2 - \sigma_1^2 \|u\|^2$, which then implies that $A^*Au_1 = \sigma_1^2 u_1$. Thus, if we set $v_1 := Au_1/\sigma_1$, then $Au_1 = \sigma_1 v_1$ and $A^*v_1 = \sigma_1 u_1$. This implies that A maps the orthogonal complement u_1^\perp of u_1 in U to the orthogonal complement v_1^\perp of v_1 in V . By induction hypothesis, the restriction of A to u_1^\perp (and v_1^\perp) then admits

a singular value decomposition with singular values $\sigma_2 \geq \dots \geq \sigma_p \geq 0$ and singular vectors $u_2, \dots, u_p \in u_1^\perp$, $v_2, \dots, v_p \in v_1^\perp$ with the stated properties. By construction we see that $\sigma_2, \dots, \sigma_p$ are less than or equal to σ_1 . If we now adjoin σ_1, u_1, v_1 to the other singular values and vectors we obtain the claim. \square

Exercise 1.3.16. Show that the singular values $\sigma_1(A) \geq \dots \geq \sigma_p(A) \geq 0$ of a $p \times n$ matrix A are unique. If we have $\sigma_1(A) > \dots > \sigma_p(A) > 0$, show that the singular vectors are unique up to rotation by a complex phase.

By construction (and the above uniqueness claim) we see that $\sigma_i(UAV) = \sigma_i(A)$ whenever A is a $p \times n$ matrix, U is a unitary $p \times p$ matrix, and V is a unitary $n \times n$ matrix. Thus the singular spectrum of a matrix is invariant under left and right unitary transformations.

Exercise 1.3.17. If A is a $p \times n$ complex matrix for some $1 \leq p \leq n$, show that the augmented matrix

$$\tilde{A} := \begin{pmatrix} 0 & A \\ A^* & 0 \end{pmatrix}$$

is a $p+n \times p+n$ Hermitian matrix whose eigenvalues consist of $\pm\sigma_1(A), \dots, \pm\sigma_p(A)$, together with $n-p$ copies of the eigenvalue zero. (This generalises Exercise 2.3.17.) What is the relationship between the singular vectors of A and the eigenvectors of \tilde{A} ?

Exercise 1.3.18. If A is an $n \times n$ Hermitian matrix, show that the singular values $\sigma_1(A), \dots, \sigma_n(A)$ of A are simply the absolute values $|\lambda_1(A)|, \dots, |\lambda_n(A)|$ of A , arranged in descending order. Show that the same claim also holds when A is a *normal matrix* (that is, when A commutes with its adjoint). What is the relationship between the singular vectors and eigenvectors of A ?

Remark 1.3.10. When A is not normal, the relationship between eigenvalues and singular values is more subtle. We will discuss this point in later sections.

Exercise 1.3.19. If A is a $p \times n$ complex matrix for some $1 \leq p \leq n$, show that AA^* has eigenvalues $\sigma_1(A)^2, \dots, \sigma_p(A)^2$, and A^*A has eigenvalues $\sigma_1(A)^2, \dots, \sigma_p(A)^2$ together with $n-p$ copies of the eigenvalue zero. Based on this observation, give an alternate proof of the singular value decomposition theorem using the spectral theorem for (positive semi-definite) Hermitian matrices.

Exercise 1.3.20. Show that the rank of a $p \times n$ matrix is equal to the number of non-zero singular values.

Exercise 1.3.21. Let A be a $p \times n$ complex matrix for some $1 \leq p \leq n$. Establish the Courant-Fischer minimax formula

$$(1.77) \quad \sigma_i(A) = \sup_{\dim(V)=i} \inf_{v \in V; |v|=1} |Av|$$

for all $1 \leq i \leq p$, where the supremum ranges over all subspaces of \mathbf{C}^n of dimension i .

One can use the above exercises to deduce many inequalities about singular values from analogous ones about eigenvalues. We give some examples below.

Exercise 1.3.22. Let A, B be $p \times n$ complex matrices for some $1 \leq p \leq n$.

(i) Establish the Weyl inequality $\sigma_{i+j-1}(A+B) \leq \sigma_i(A) + \sigma_j(B)$ whenever $1 \leq i, j, i+j-1 \leq p$.

(ii) Establish the Lidskii inequality

$$\begin{aligned} \sigma_{i_1}(A+B) + \cdots + \sigma_{i_k}(A+B) &\leq \sigma_{i_1}(A) + \cdots + \sigma_{i_k}(A) \\ &\quad + \sigma_1(B) + \cdots + \sigma_k(B) \end{aligned}$$

whenever $1 \leq i_1 < \cdots < i_k \leq p$.

(iii) Show that for any $1 \leq k \leq p$, the map $A \mapsto \sigma_1(A) + \cdots + \sigma_k(A)$ defines a norm on the space $\mathbf{C}^{p \times n}$ of complex $p \times n$ matrices (this norm is known as the k^{th} Ky Fan norm).

(iv) Establish the Weyl inequality $|\sigma_i(A+B) - \sigma_i(A)| \leq \|B\|_{\text{op}}$ for all $1 \leq i \leq p$.

(v) More generally, establish the q -Weilandt-Hoffman inequality $\|(\sigma_i(A+B) - \sigma_i(A))_{1 \leq i \leq p}\|_{\ell_p^q} \leq \|B\|_{S^q}$ for any $1 \leq q \leq \infty$, where $\|B\|_{S^q} := \|(\sigma_i(B))_{1 \leq i \leq p}\|_{\ell_p^q}$ is the q -Schatten norm of B . (Note that this is consistent with the previous definition of the Schatten norms.)

(vi) Show that the q -Schatten norm is indeed a norm on $\mathbf{C}^{p \times n}$ for any $1 \leq q \leq \infty$.

(vii) If A' is formed by removing one row from A , show that $\lambda_{i+1}(A) \leq \lambda_i(A') \leq \lambda_i(A)$ for all $1 \leq i < p$.

(viii) If $p < n$ and A' is formed by removing one column from A , show that $\lambda_{i+1}(A) \leq \lambda_i(A') \leq \lambda_i(A)$ for all $1 \leq i < p$ and $\lambda_p(A') \leq \lambda_p(A)$. What changes when $p = n$?

Exercise 1.3.23. Let A be a $p \times n$ complex matrix for some $1 \leq p \leq n$. Observe that the linear transformation $A : \mathbf{C}^n \rightarrow \mathbf{C}^p$ naturally induces a linear transformation $A^{\wedge k} : \bigwedge^k \mathbf{C}^n \rightarrow \bigwedge^k \mathbf{C}^p$ from k -forms on \mathbf{C}^n to k -forms on \mathbf{C}^p . We give $\bigwedge^k \mathbf{C}^n$ the structure of a Hilbert space by declaring the basic

forms $e_{i_1} \wedge \dots \wedge e_{i_k}$ for $1 \leq i_1 < \dots < i_k \leq n$ to be orthonormal. For any $1 \leq k \leq p$, show that the operator norm of $A^{\wedge k}$ is equal to $\sigma_1(A) \dots \sigma_k(A)$.

Exercise 1.3.24. Let A be a $p \times n$ matrix for some $1 \leq p \leq n$, let B be a $r \times p$ matrix, and let C be a $n \times s$ matrix for some $r, s \geq 1$. Show that $\sigma_i(BA) \leq \|B\|_{\text{op}} \sigma_i(A)$ and $\sigma_i(AC) \leq \sigma_i(A) \|C\|_{\text{op}}$ for any $1 \leq i \leq p$.

Exercise 1.3.25. Let $A = (a_{ij})_{1 \leq i \leq p; 1 \leq j \leq n}$ be a $p \times n$ matrix for some $1 \leq p \leq n$, let $i_1, \dots, i_k \in \{1, \dots, p\}$ be distinct, and let $j_1, \dots, j_k \in \{1, \dots, n\}$ be distinct. Show that

$$a_{i_1 j_1} + \dots + a_{i_k j_k} \leq \sigma_1(A) + \dots + \sigma_k(A).$$

Using this, show that if $j_1, \dots, j_p \in \{1, \dots, n\}$ are distinct, then

$$\|(a_{ij_i})_{i=1}^p\|_{\ell_p^q} \leq \|A\|_{S^q}$$

for every $1 \leq q \leq \infty$.

Exercise 1.3.26. Establish the Hölder inequality

$$|\text{tr}(AB^*)| \leq \|A\|_{S^q} \|B\|_{S^{q'}}$$

whenever A, B are $p \times n$ complex matrices and $1 \leq q, q' \leq \infty$ are such that $1/q + 1/q' = 1$.

2.3. The operator norm of random matrices

Now that we have developed the basic probabilistic tools that we will need, we now turn to the main subject of this text, namely the study of random matrices. There are many random matrix models (aka matrix ensembles) of interest—far too many to all be discussed here. We will thus focus on just a few simple models. First of all, we shall restrict attention to square matrices $M = (\xi_{ij})_{1 \leq i, j \leq n}$, where n is a (large) integer and the ξ_{ij} are real or complex random variables. (One can certainly study rectangular matrices as well, but for simplicity we will only look at the square case.) Then, we shall restrict to three main models:

- (i) **Iid matrix ensembles**, in which the coefficients ξ_{ij} are iid random variables with a single distribution $\xi_{ij} \equiv \xi$. We will often normalise ξ to have mean zero and unit variance. Examples of iid models include the *Bernoulli ensemble* (aka *random sign matrices*) in which the ξ_{ij} are signed Bernoulli variables, the *real Gaussian matrix ensemble* in which $\xi_{ij} \equiv N(0, 1)_{\mathbf{R}}$, and the *complex Gaussian matrix ensemble* in which $\xi_{ij} \equiv N(0, 1)_{\mathbf{C}}$.
- (ii) **Symmetric Wigner matrix ensembles**, in which the upper triangular coefficients ξ_{ij} , $j \geq i$ are jointly independent and real, but the lower triangular coefficients ξ_{ij} , $j < i$ are constrained to equal their transposes: $\xi_{ij} = \xi_{ji}$. Thus M by construction is always a real symmetric matrix. Typically, the strictly upper triangular coefficients will be iid, as will the diagonal coefficients, but the two classes of coefficients may have a different distribution. One example here is the *symmetric Bernoulli ensemble*, in which both the strictly upper triangular and the diagonal entries are signed Bernoulli variables; another important example is the *Gaussian Orthogonal Ensemble (GOE)*, in which the upper triangular entries have distribution $N(0, 1)_{\mathbf{R}}$ and the diagonal entries have distribution $N(0, 2)_{\mathbf{R}}$. (We will explain the reason for this discrepancy later.)
- (iii) **Hermitian Wigner matrix ensembles**, in which the upper triangular coefficients are jointly independent, with the diagonal entries being real and the strictly upper triangular entries complex, and the lower triangular coefficients ξ_{ij} , $j < i$ are constrained to equal their adjoints: $\xi_{ij} = \overline{\xi_{ji}}$. Thus M by construction is always a Hermitian matrix. This class of ensembles contains the symmetric Wigner ensembles as a subclass. Another very important example is the *Gaussian Unitary Ensemble (GUE)*, in which all

off-diagonal entries have distribution $N(0, 1)_{\mathbf{C}}$, but the diagonal entries have distribution $N(0, 1)_{\mathbf{R}}$.

Given a matrix ensemble M , there are many statistics of M that one may wish to consider, e.g., the eigenvalues or singular values of M , the trace and determinant, etc. In this section we will focus on a basic statistic, namely the *operator norm*

$$(2.57) \quad \|M\|_{\text{op}} := \sup_{x \in \mathbf{C}^n: |x|=1} |Mx|$$

of the matrix M . This is an interesting quantity in its own right, but also serves as a basic upper bound on many other quantities. (For instance, $\|M\|_{\text{op}}$ is also the largest singular value $\sigma_1(M)$ of M and thus dominates the other singular values; similarly, all eigenvalues $\lambda_i(M)$ of M clearly have magnitude at most $\|M\|_{\text{op}}$.) Because of this, it is particularly important to get good *upper tail bounds*,

$$\mathbf{P}(\|M\|_{\text{op}} \geq \lambda) \leq \dots,$$

on this quantity, for various thresholds λ . (Lower tail bounds are also of interest, of course; for instance, they give us confidence that the upper tail bounds are sharp.) Also, as we shall see, the problem of upper bounding $\|M\|_{\text{op}}$ can be viewed as a non-commutative analogue¹⁴ of upper bounding the quantity $|S_n|$ studied in Section 2.1.

An $n \times n$ matrix consisting entirely of 1s has an operator norm of exactly n , as can for instance be seen from the Cauchy-Schwarz inequality. More generally, any matrix whose entries are all uniformly $O(1)$ will have an operator norm of $O(n)$ (which can again be seen from Cauchy-Schwarz, or alternatively from *Schur's test* (see e.g. [Ta2010, §1.11]), or from a computation of the *Frobenius norm* (see (2.63))). However, this argument does not take advantage of possible cancellations in M . Indeed, from analogy with concentration of measure, when the entries of the matrix M are independent, bounded and have mean zero, we expect the operator norm to be of size $O(\sqrt{n})$ rather than $O(n)$. We shall see shortly that this intuition is indeed correct¹⁵.

As mentioned before, there is an analogy here with the concentration of measure¹⁶ phenomenon, and many of the tools used in the latter (e.g., the moment method) will also appear here. Similarly, just as many of the tools

¹⁴The analogue of the central limit theorem studied in Section 2.2 is the *Wigner semicircular law*, which will be studied in Section 2.4.

¹⁵One can see, though, that the mean zero hypothesis is important; from the triangle inequality we see that if we add the all-ones matrix (for instance) to a random matrix with mean zero, to obtain a random matrix whose coefficients all have mean 1, then at least one of the two random matrices necessarily has operator norm at least $n/2$.

¹⁶Indeed, we will be able to use some of the concentration inequalities from Section 2.1 directly to help control $\|M\|_{\text{op}}$ and related quantities.

from concentration of measure could be adapted to help prove the central limit theorem, several of the tools seen here will be of use in deriving the semicircular law in Section 2.4.

The most advanced knowledge we have on the operator norm is given by the *Tracy-Widom law*, which not only tells us where the operator norm is concentrated in (it turns out, for instance, that for a Wigner matrix (with some additional technical assumptions), it is concentrated in the range $[2\sqrt{n} - O(n^{-1/6}), 2\sqrt{n} + O(n^{-1/6})]$), but what its distribution in that range is. While the methods in this section can eventually be pushed to establish this result, this is far from trivial, and will only be briefly discussed here. We will, however, discuss the Tracy-Widom law at several later points in the text.

2.3.1. The epsilon net argument. The slickest way to control $\|M\|_{\text{op}}$ is via the moment method. But let us defer using this method for the moment, and work with a more “naive” way to control the operator norm, namely by working with the definition (2.57). From that definition, we see that we can view the upper tail event $\|M\|_{\text{op}} > \lambda$ as a union of many simpler events:

$$(2.58) \quad \mathbf{P}(\|M\|_{\text{op}} > \lambda) \leq \mathbf{P}\left(\bigvee_{x \in S} |Mx| > \lambda\right)$$

where $S := \{x \in \mathbf{C}^d : |x| = 1\}$ is the unit sphere in the complex space \mathbf{C}^d .

The point of doing this is that the event $|Mx| > \lambda$ is easier to control than the event $\|M\|_{\text{op}} > \lambda$, and can in fact be handled by the concentration of measure estimates we already have. For instance:

Lemma 2.3.1. *Suppose that the coefficients ξ_{ij} of M are independent, have mean zero, and are uniformly bounded in magnitude by 1. Let x be a unit vector in \mathbf{C}^n . Then for sufficiently large A (larger than some absolute constant), one has*

$$\mathbf{P}(|Mx| \geq A\sqrt{n}) \leq C \exp(-cAn)$$

for some absolute constants $C, c > 0$.

Proof. Let X_1, \dots, X_n be the n rows of M , then the column vector Mx has coefficients $X_i \cdot x$ for $i = 1, \dots, n$. If we let x_1, \dots, x_n be the coefficients of x , so that $\sum_{j=1}^n |x_j|^2 = 1$, then $X_i \cdot x$ is just $\sum_{j=1}^n \xi_{ij} x_j$. Applying standard concentration of measure results (e.g., Exercise 2.1.4, Exercise 2.1.5, or Theorem 2.1.13), we see that each $X_i \cdot x$ is uniformly sub-Gaussian, thus

$$\mathbf{P}(|X_i \cdot x| \geq \lambda) \leq C \exp(-c\lambda^2)$$

for some absolute constants $C, c > 0$. In particular, we have

$$\mathbf{E}e^{c|X_i \cdot x|^2} \leq C$$

for some (slightly different) absolute constants $C, c > 0$. Multiplying these inequalities together for all i , we obtain

$$\mathbf{E}e^{c|Mx|^2} \leq C^n$$

and the claim then follows from Markov's inequality (1.14). \square

Thus (with the hypotheses of Lemma 2.3.1), we see that for each individual unit vector x , we have $|Mx| = O(\sqrt{n})$ with overwhelming probability. It is then tempting to apply the union bound and try to conclude that $\|M\|_{\text{op}} = O(\sqrt{n})$ with overwhelming probability also. However, we encounter a difficulty: the unit sphere S is uncountable, and so we are taking the union over an uncountable number of events. Even though each event occurs with exponentially small probability, the union could well be everything.

Of course, it is extremely wasteful to apply the union bound to an uncountable union. One can pass to a countable union just by working with a countable dense subset of the unit sphere S instead of the sphere itself, since the map $x \mapsto |Mx|$ is continuous. Of course, this is still an infinite set and so we still cannot usefully apply the union bound. However, the map $x \mapsto |Mx|$ is not just continuous; it is *Lipschitz* continuous, with a Lipschitz constant of $\|M\|_{\text{op}}$. Now, of course there is some circularity here because $\|M\|_{\text{op}}$ is precisely the quantity we are trying to bound. Nevertheless, we can use this stronger continuity to refine the countable dense subset further, to a *finite* (but still quite dense) subset of S , at the slight cost of modifying the threshold λ by a constant factor. Namely:

Lemma 2.3.2. *Let Σ be a maximal $1/2$ -net of the sphere S , i.e., a set of points in S that are separated from each other by a distance of at least $1/2$, and which is maximal with respect to set inclusion. Then for any $n \times n$ matrix M with complex coefficients, and any $\lambda > 0$, we have*

$$\mathbf{P}(\|M\|_{\text{op}} > \lambda) \leq \mathbf{P}\left(\bigvee_{y \in \Sigma} |My| > \lambda/2\right).$$

Proof. By (2.57) (and compactness) we can find $x \in S$ such that

$$|Mx| = \|M\|_{\text{op}}.$$

This point x need not lie in Σ . However, as Σ is a maximal $1/2$ -net of S , we know that x lies within $1/2$ of some point y in Σ (since otherwise we could add x to Σ and contradict maximality). Since $|x - y| \leq 1/2$, we have

$$|M(x - y)| \leq \|M\|_{\text{op}}/2.$$

By the triangle inequality we conclude that

$$|My| \geq \|M\|_{\text{op}}/2.$$

In particular, if $\|M\|_{\text{op}} > \lambda$, then $|My| > \lambda/2$ for some $y \in \Sigma$, and the claim follows. \square

Remark 2.3.3. Clearly, if one replaces the maximal $1/2$ -net here with a maximal ε -net for some other $0 < \varepsilon < 1$ (defined in the obvious manner), then we get the same conclusion, but with $\lambda/2$ replaced by $\lambda(1 - \varepsilon)$.

Now that we have discretised the range of points y to be finite, the union bound becomes viable again. We first make the following basic observation:

Lemma 2.3.4 (Volume packing argument). *Let $0 < \varepsilon < 1$, and let Σ be an ε -net of the sphere S . Then Σ has cardinality at most $(C/\varepsilon)^n$ for some absolute constant $C > 0$.*

Proof. Consider the balls of radius $\varepsilon/2$ centred around each point in Σ ; by hypothesis, these are disjoint. On the other hand, by the triangle inequality, they are all contained in the ball of radius $3/2$ centred at the origin. The volume of the latter ball is at most $(C/\varepsilon)^n$ the volume of any of the small balls, and the claim follows. \square

Exercise 2.3.1. Improve the bound $(C/\varepsilon)^n$ to C^n/ε^{n-1} . In the converse direction, if Σ is a *maximal* ε -net, show that Σ has cardinality *at least* c^n/ε^{n-1} for some absolute constant $c > 0$.

And now we get an upper tail estimate:

Corollary 2.3.5 (Upper tail estimate for iid ensembles). *Suppose that the coefficients ξ_{ij} of M are independent, have mean zero, and uniformly bounded in magnitude by 1. Then there exist absolute constants $C, c > 0$ such that*

$$\mathbf{P}(\|M\|_{\text{op}} > A\sqrt{n}) \leq C \exp(-cAn)$$

for all $A \geq C$. In particular, we have $\|M\|_{\text{op}} = O(\sqrt{n})$ with overwhelming probability.

Proof. From Lemma 2.3.2 and the union bound, we have

$$\mathbf{P}(\|M\|_{\text{op}} > A\sqrt{n}) \leq \sum_{y \in \Sigma} \mathbf{P}(|My| > A\sqrt{n}/2)$$

where Σ is a maximal $1/2$ -net of S . By Lemma 2.3.1, each of the probabilities $\mathbf{P}(|My| > A\sqrt{n}/2)$ is bounded by $C \exp(-cAn)$ if A is large enough. Meanwhile, from Lemma 2.3.4, Σ has cardinality $O(1)^n$. If A is large enough, the *entropy loss*¹⁷ of $O(1)^n$ can be absorbed into the exponential gain of $\exp(-cAn)$ by modifying c slightly, and the claim follows. \square

¹⁷Roughly speaking, the *entropy* of a configuration is the logarithm of the number of possible states that configuration can be in. When applying the union bound to control all possible configurations at once, one often loses a factor proportional to the number of such states; this factor is sometimes referred to as the *entropy factor* or *entropy loss* in one's argument.

Exercise 2.3.2. If Σ is a maximal $1/4$ -net instead of a maximal $1/2$ -net, establish the following variant

$$\mathbf{P}(\|M\|_{\text{op}} > \lambda) \leq \mathbf{P}\left(\bigvee_{x,y \in \Sigma} |x^* M y| > \lambda/4\right)$$

of Lemma 2.3.2. Use this to provide an alternate proof of Corollary 2.3.5.

The above result was for matrices with independent entries, but it easily extends to the Wigner case:

Corollary 2.3.6 (Upper tail estimate for Wigner ensembles). *Suppose that the coefficients ξ_{ij} of M are independent for $j \geq i$, mean zero, and uniformly bounded in magnitude by 1, and let $\xi_{ij} := \overline{\xi_{ji}}$ for $j < i$. Then there exist absolute constants $C, c > 0$ such that*

$$\mathbf{P}(\|M\|_{\text{op}} > A\sqrt{n}) \leq C \exp(-cAn)$$

for all $A \geq C$. In particular, we have $\|M\|_{\text{op}} = O(\sqrt{n})$ with overwhelming probability.

Proof. From Corollary 2.3.5, the claim already holds for the upper-triangular portion of M , as well as for the strict lower-triangular portion of M . The claim then follows from the triangle inequality (adjusting the constants C, c appropriately). \square

Exercise 2.3.3. Generalise Corollary 2.3.5 and Corollary 2.3.6 to the case where the coefficients ξ_{ij} have uniform sub-Gaussian tails, rather than being uniformly bounded in magnitude by 1.

Remark 2.3.7. What we have just seen is a simple example of an *epsilon net argument*, which is useful when controlling a supremum of random variables $\sup_{x \in S} X_x$ such as (2.57), where each individual random variable X_x is known to obey a large deviation inequality (in this case, Lemma 2.3.1). The idea is to use metric arguments (e.g., the triangle inequality, see Lemma 2.3.2) to refine the set of parameters S to take the supremum over to an ε -net $\Sigma = \Sigma_\varepsilon$ for some suitable ε , and then apply the union bound. One takes a loss based on the cardinality of the ε -net (which is basically the *covering number* of the original parameter space at scale ε), but one can hope that the bounds from the large deviation inequality are strong enough (and the metric entropy bounds sufficiently accurate) to overcome this entropy loss.

There is of course the question of what scale ε to use. In this simple example, the scale $\varepsilon = 1/2$ sufficed. In other contexts, one has to choose the scale ε more carefully. In more complicated examples with no natural preferred scale, it often makes sense to take a large range of scales (e.g., $\varepsilon = 2^{-j}$ for $j = 1, \dots, J$) and *chain* them together by using telescoping series such as $X_x = X_{x_1} + \sum_{j=1}^J X_{x_{j+1}} - X_{x_j}$ (where x_j is the nearest point in Σ_j to

x for $j = 1, \dots, J$, and x_{J+1} is x by convention) to estimate the supremum, the point being that one can hope to exploit cancellations between adjacent elements of the sequence X_{x_j} . This is known as the method of *chaining*. There is an even more powerful refinement of this method, known as the method of *generic chaining*, which has a large number of applications; see [Ta2005] for a beautiful and systematic treatment of the subject. However, we will not use this method in this text.

2.3.2. A symmetrisation argument (optional). We pause here to record an elegant *symmetrisation argument* that exploits convexity to allow us to reduce, without loss of generality, to the symmetric case $M \equiv -M$, albeit at the cost of losing a factor of 2. We will not use this type of argument directly in this text, but it is often used elsewhere in the literature.

Let M be any random matrix with mean zero, and let \tilde{M} be an independent copy of M . Then, conditioning on M , we have

$$\mathbf{E}(M - \tilde{M} | M) = M.$$

As the operator norm $M \mapsto \|M\|_{\text{op}}$ is convex, we can then apply Jensen's inequality (Exercise 1.1.8) to conclude that

$$\mathbf{E}(\|M - \tilde{M}\|_{\text{op}} | M) \geq \|M\|_{\text{op}}.$$

Undoing the conditioning over M , we conclude that

$$(2.59) \quad \mathbf{E}\|M - \tilde{M}\|_{\text{op}} \geq \mathbf{E}\|M\|_{\text{op}}.$$

Thus, to upper bound the expected operator norm of M , it suffices to upper bound the expected operator norm of $M - \tilde{M}$. The point is that even if M is not symmetric ($M \not\equiv -M$), $M - \tilde{M}$ is automatically symmetric.

One can modify (2.59) in a few ways, given some more hypotheses on M . Suppose now that $M = (\xi_{ij})_{1 \leq i, j \leq n}$ is a matrix with independent entries, thus $M - \tilde{M}$ has coefficients $\xi_{ij} - \tilde{\xi}_{ij}$ where $\tilde{\xi}_{ij}$ is an independent copy of ξ_{ij} . Introduce a random sign matrix $E = (\varepsilon_{ij})_{1 \leq i, j \leq n}$ which is (jointly) independent of M, \tilde{M} . Observe that as the distribution of $\xi_{ij} - \tilde{\xi}_{ij}$ is symmetric, that

$$(\xi_{ij} - \tilde{\xi}_{ij}) \equiv (\xi_{ij} - \tilde{\xi}_{ij})\varepsilon_{ij},$$

and thus

$$(M - \tilde{M}) \equiv (M - \tilde{M}) \cdot E$$

where $A \cdot B := (a_{ij}b_{ij})_{1 \leq i, j \leq n}$ is the *Hadamard product* of $A = (a_{ij})_{1 \leq i, j \leq n}$ and $B = (b_{ij})_{1 \leq i, j \leq n}$. We conclude from (2.59) that

$$\mathbf{E}\|M\|_{\text{op}} \leq \mathbf{E}\|(M - \tilde{M}) \cdot E\|_{\text{op}}.$$

By the distributive law and the triangle inequality we have

$$\|(M - \tilde{M}) \cdot E\|_{\text{op}} \leq \|M \cdot E\|_{\text{op}} + \|\tilde{M} \cdot E\|_{\text{op}}.$$

But as $M \cdot E \equiv \tilde{M} \cdot E$, the quantities $\|M \cdot E\|_{\text{op}}$ and $\|\tilde{M} \cdot E\|_{\text{op}}$ have the same expectation. We conclude the *symmetrisation inequality*

$$(2.60) \quad \mathbf{E}\|M\|_{\text{op}} \leq 2\mathbf{E}\|M \cdot E\|_{\text{op}}.$$

Thus, if one does not mind losing a factor of two, one has the freedom to randomise the sign of each entry of M independently (assuming that the entries were already independent). Thus, in proving Corollary 2.3.5, one could have reduced to the case when the ξ_{ij} were symmetric, though in this case this would not have made the argument that much simpler.

Sometimes it is preferable to multiply the coefficients by a Gaussian rather than by a random sign. Again, let $M = (\xi_{ij})_{1 \leq i, j \leq n}$ have independent entries with mean zero. Let $G = (g_{ij})_{1 \leq i, j \leq n}$ be a real Gaussian matrix independent of M , thus the $g_{ij} \equiv N(0, 1)_{\mathbf{R}}$ are iid. We can split $G = E \cdot |G|$, where $E := (\text{sgn}(g_{ij}))_{1 \leq i, j \leq n}$ and $|G| = (|g_{ij}|)_{1 \leq i, j \leq n}$. Note that E , M , $|G|$ are independent, and E is a random sign matrix. In particular, (2.60) holds. We now use

Exercise 2.3.4. If $g \equiv N(0, 1)_{\mathbf{R}}$, show that $\mathbf{E}|g| = \sqrt{\frac{2}{\pi}}$.

From this exercise we see that

$$\mathbf{E}(M \cdot E \cdot |G| | M, E) = \sqrt{\frac{2}{\pi}} M \cdot E$$

and hence by Jensen's inequality (Exercise 1.1.8) again

$$\mathbf{E}(\|M \cdot E \cdot |G|\|_{\text{op}} | M, E) \geq \sqrt{\frac{2}{\pi}} \|M \cdot E\|_{\text{op}}.$$

Undoing the conditional expectation in M, E and applying (2.60) we conclude the *Gaussian symmetrisation inequality*

$$(2.61) \quad \mathbf{E}\|M\|_{\text{op}} \leq \sqrt{2\pi} \mathbf{E}\|M \cdot G\|_{\text{op}}.$$

Thus, for instance, when proving Corollary 2.3.5, one could have inserted a random Gaussian in front of each coefficient. This would have made the proof of Lemma 2.3.1 marginally simpler (as one could compute directly with Gaussians, and reduce the number of appeals to concentration of measure results) but in this case the improvement is negligible. In other situations though it can be quite helpful to have the additional random sign or random Gaussian factor present. For instance, we have the following result of Latala [La2005]:

Theorem 2.3.8. *Let $M = (\xi_{ij})_{1 \leq i, j \leq n}$ be a matrix with independent mean zero entries, obeying the second moment bounds*

$$\sup_i \sum_{j=1}^n \mathbf{E}|\xi_{ij}|^2 \leq K^2 n,$$

$$\sup_j \sum_{i=1}^n \mathbf{E}|\xi_{ij}|^2 \leq K^2 n,$$

and the fourth moment bound

$$\sum_{i=1}^n \sum_{j=1}^n \mathbf{E}|\xi_{ij}|^4 \leq K^4 n^2$$

for some $K > 0$. Then $\mathbf{E}\|M\|_{\text{op}} = O(K\sqrt{n})$.

Proof (Sketch only). Using (2.61) one can replace ξ_{ij} by $\xi_{ij} \cdot g_{ij}$ without much penalty. One then runs the epsilon-net argument with an explicit net, and uses concentration of measure results for Gaussians (such as Theorem 2.1.12) to obtain the analogue of Lemma 2.3.1. The details are rather intricate, and we refer the interested reader to [La2005]. \square

As a corollary of Theorem 2.3.8, we see that if we have an iid matrix (or Wigner matrix) of mean zero whose entries have a fourth moment of $O(1)$, then the expected operator norm is $O(\sqrt{n})$. The fourth moment hypothesis is sharp. To see this, we make the trivial observation that the operator norm of a matrix $M = (\xi_{ij})_{1 \leq i, j \leq n}$ bounds the magnitude of any of its coefficients, thus

$$\sup_{1 \leq i, j \leq n} |\xi_{ij}| \leq \|M\|_{\text{op}}$$

or, equivalently, that

$$\mathbf{P}(\|M\|_{\text{op}} \leq \lambda) \leq \mathbf{P}\left(\bigvee_{1 \leq i, j \leq n} |\xi_{ij}| \leq \lambda\right).$$

In the iid case $\xi_{ij} \equiv \xi$, and setting $\lambda = A\sqrt{n}$ for some fixed A independent of n , we thus have

$$(2.62) \quad \mathbf{P}(\|M\|_{\text{op}} \leq A\sqrt{n}) \leq \mathbf{P}(|\xi| \leq A\sqrt{n})^{n^2}.$$

With the fourth moment hypothesis, one has from dominated convergence that

$$\mathbf{P}(|\xi| \leq A\sqrt{n}) \geq 1 - o_A(1/n^2),$$

and so the right-hand side of (2.62) is asymptotically trivial. But with weaker hypotheses than the fourth moment hypothesis, the rate of convergence of $\mathbf{P}(|\xi| \leq A\sqrt{n})$ to 1 can be slower, and one can easily build

examples for which the right-hand side of (2.62) is $o_A(1)$ for every A , which forces $\|M\|_{\text{op}}$ to typically be much larger than \sqrt{n} on the average.

Remark 2.3.9. The symmetrisation inequalities remain valid with the operator norm replaced by any other convex norm on the space of matrices. The results are also just as valid for rectangular matrices as for square ones.

2.3.3. Concentration of measure. Consider a random matrix M of the type considered in Corollary 2.3.5 (e.g., a random sign matrix). We now know that the operator norm $\|M\|_{\text{op}}$ is of size $O(\sqrt{n})$ with overwhelming probability. But there is much more that can be said. For instance, by taking advantage of the convexity and Lipschitz properties of $\|M\|_{\text{op}}$, we have the following quick application of Talagrand's inequality (Theorem 2.1.13):

Proposition 2.3.10. *Let M be as in Corollary 2.3.5. Then for any $\lambda > 0$, one has*

$$\mathbf{P}(\|\|M\|_{\text{op}} - \mathbf{M}\|M\|_{\text{op}}\| \geq \lambda) \leq C \exp(-c\lambda^2)$$

for some absolute constants $C, c > 0$, where $\mathbf{M}\|M\|_{\text{op}}$ is a median value for $\|M\|_{\text{op}}$. The same result also holds with $\mathbf{M}\|M\|_{\text{op}}$ replaced by the expectation $\mathbf{E}\|M\|_{\text{op}}$.

Proof. We view $\|M\|_{\text{op}}$ as a function $F((\xi_{ij})_{1 \leq i, j \leq n})$ of the independent complex variables ξ_{ij} , thus F is a function from \mathbf{C}^{n^2} to \mathbf{R} . The convexity of the operator norm tells us that F is convex. The triangle inequality, together with the elementary bound

$$(2.63) \quad \|M\|_{\text{op}} \leq \|M\|_F$$

(easily proven by Cauchy-Schwarz), where

$$(2.64) \quad \|M\|_F := \left(\sum_{i=1}^n \sum_{j=1}^n |\xi_{ij}|^2 \right)^{1/2}$$

is the *Frobenius norm* (also known as the *Hilbert-Schmidt norm* or *2-Schatten norm*), tells us that F is Lipschitz with constant 1. The claim then follows directly from Talagrand's inequality (Theorem 2.1.13). \square

Exercise 2.3.5. Establish a similar result for the matrices in Corollary 2.3.6.

From Corollary 2.3.5 we know that the median or expectation of $\|M\|_{\text{op}}$ is of size $O(\sqrt{n})$; we now know that $\|M\|_{\text{op}}$ concentrates around this median to width at most $O(1)$. (This turns out to be non-optimal; the Tracy-Widom law actually gives a concentration of $O(n^{-1/6})$, under some additional assumptions on M . Nevertheless, this level of concentration is already non-trivial.)

However, this argument does not tell us much about what the median or expected value of $\|M\|_{\text{op}}$ actually *is*. For this, we will need to use other methods, such as the moment method which we turn to next.

Remark 2.3.11. Talagrand’s inequality, as formulated in Theorem 2.1.13, relies heavily on convexity. Because of this, we cannot apply this argument directly to non-convex matrix statistics, such as singular values $\sigma_j(M)$ other than the largest singular value $\sigma_1(M)$. Nevertheless, one can still use this inequality to obtain good concentration results, by using the convexity of related quantities, such as the partial sums $\sum_{j=1}^J \sigma_j(M)$; see [Me2004]. Other approaches include the use of alternate large deviation inequalities, such as those arising from log-Sobolev inequalities (see e.g., [Gu2009]), or by using more abstract versions of Talagrand’s inequality (see [AlKrVu2002], [GuZe2000]).

2.3.4. The moment method. We now bring the moment method to bear on the problem, starting with the easy moments and working one’s way up to the more sophisticated moments. It turns out that it is easier to work first with the case when M is symmetric or Hermitian; we will discuss the non-symmetric case near the end of this section.

The starting point for the moment method is the observation that for symmetric or Hermitian M , the operator norm $\|M\|_{\text{op}}$ is equal to the ℓ^∞ norm

$$(2.65) \quad \|M\|_{\text{op}} = \max_{1 \leq i \leq n} |\lambda_i|$$

of the eigenvalues $\lambda_1, \dots, \lambda_n \in \mathbf{R}$ of M . On the other hand, we have the standard linear algebra identity

$$\text{tr}(M) = \sum_{i=1}^n \lambda_i$$

and more generally

$$\text{tr}(M^k) = \sum_{i=1}^n \lambda_i^k.$$

In particular, if $k = 2, 4, \dots$ is an even integer, then $\text{tr}(M^k)^{1/k}$ is just the ℓ^k norm of these eigenvalues, and we have the inequalities

$$(2.66) \quad \|M\|_{\text{op}}^k \leq \text{tr}(M^k) \leq n \|M\|_{\text{op}}^k.$$

To put this another way, knowledge of the k^{th} moment $\text{tr}(M^k)$ controls the operator norm up to a multiplicative factor of $n^{1/k}$. Taking larger and larger k , we should thus obtain more accurate control on the operator norm¹⁸.

¹⁸This is also the philosophy underlying the *power method* in numerical linear algebra.

Remark 2.3.12. In most cases, one expects the eigenvalues to be reasonably uniformly distributed, in which case the upper bound in (2.66) is closer to the truth than the lower bound. One scenario in which this can be rigorously established is if it is known that the eigenvalues of M all come with a high multiplicity. This is often the case for matrices associated with group actions (particularly those which are *quasirandom* in the sense of Gowers [Go2008]). However, this is usually not the case with most random matrix ensembles, and we must instead proceed by increasing k as described above.

Let's see how this method works in practice. The simplest case is that of the second moment $\text{tr}(M^2)$, which in the Hermitian case works out to

$$\text{tr}(M^2) = \sum_{i=1}^n \sum_{j=1}^n |\xi_{ij}|^2 = \|M\|_F^2.$$

Note that (2.63) is just the $k = 2$ case of the lower inequality in (2.66), at least in the Hermitian case.

The expression $\sum_{i=1}^n \sum_{j=1}^n |\xi_{ij}|^2$ is easy to compute in practice. For instance, for the symmetric Bernoulli ensemble, this expression is exactly equal to n^2 . More generally, if we have a Wigner matrix in which all off-diagonal entries have mean zero and unit variance, and the diagonal entries have mean zero and bounded variance (this is the case for instance for GOE), then the off-diagonal entries have mean 1, and by the law of large numbers¹⁹ we see that this expression is almost surely asymptotic to n^2 .

From the weak law of large numbers, we see, in particular, that one has

$$(2.67) \quad \sum_{i=1}^n \sum_{j=1}^n |\xi_{ij}|^2 = (1 + o(1))n^2$$

asymptotically almost surely.

Exercise 2.3.6. If the ξ_{ij} have uniformly sub-exponential tail, show that we in fact have (2.67) with overwhelming probability.

Applying (2.66), we obtain the bounds

$$(2.68) \quad (1 + o(1))\sqrt{n} \leq \|M\|_{\text{op}} \leq (1 + o(1))n$$

asymptotically almost surely. This is already enough to show that the median of $\|M\|_{\text{op}}$ is at least $(1 + o(1))\sqrt{n}$, which complements (up to constants) the upper bound of $O(\sqrt{n})$ obtained from the epsilon net argument. But the upper bound here is terrible; we will need to move to higher moments to improve it.

¹⁹There is of course a dependence between the upper triangular and lower triangular entries, but this is easy to deal with by folding the sum into twice the upper triangular portion (plus the diagonal portion, which is lower order).

Accordingly, we now turn to the fourth moment. For simplicity let us assume that all entries ξ_{ij} have zero mean and unit variance. To control moments beyond the second moment, we will also assume that all entries are bounded in magnitude by some K . We expand

$$\mathrm{tr}(M^4) = \sum_{1 \leq i_1, i_2, i_3, i_4 \leq n} \xi_{i_1 i_2} \xi_{i_2 i_3} \xi_{i_3 i_4} \xi_{i_4 i_1}.$$

To understand this expression, we take expectations:

$$\mathbf{E} \mathrm{tr}(M^4) = \sum_{1 \leq i_1, i_2, i_3, i_4 \leq n} \mathbf{E} \xi_{i_1 i_2} \xi_{i_2 i_3} \xi_{i_3 i_4} \xi_{i_4 i_1}.$$

One can view this sum graphically, as a sum over length four cycles in the vertex set $\{1, \dots, n\}$; note that the four edges $\{i_1, i_2\}, \{i_2, i_3\}, \{i_3, i_4\}, \{i_4, i_1\}$ are allowed to be degenerate if two adjacent ξ_i are equal. The value of each term

$$(2.69) \quad \mathbf{E} \xi_{i_1 i_2} \xi_{i_2 i_3} \xi_{i_3 i_4} \xi_{i_4 i_1}$$

in this sum depends on what the cycle does.

First, there is the case when all of the four edges $\{i_1, i_2\}, \{i_2, i_3\}, \{i_3, i_4\}, \{i_4, i_1\}$ are distinct. Then the four factors $\xi_{i_1 i_2}, \dots, \xi_{i_4 i_1}$ are independent; since we are assuming them to have mean zero, the term (2.69) vanishes. Indeed, the same argument shows that the only terms that do not vanish are those in which each edge is repeated at least twice. A short combinatorial case check then shows that, up to cyclic permutations of the i_1, i_2, i_3, i_4 indices there are now only a few types of cycles in which the term (2.69) does not automatically vanish:

- (i) $i_1 = i_3$, but i_2, i_4 are distinct from each other and from i_1 .
- (ii) $i_1 = i_3$ and $i_2 = i_4$.
- (iii) $i_1 = i_2 = i_3$, but i_4 is distinct from i_1 .
- (iv) $i_1 = i_2 = i_3 = i_4$.

In the first case, the independence and unit variance assumptions tell us that (2.69) is 1, and there are $O(n^3)$ such terms, so the total contribution here to $\mathbf{E} \mathrm{tr}(M^4)$ is at most $O(n^3)$. In the second case, the unit variance and bound by K tells us that the term is $O(K^2)$, and there are $O(n^2)$ such terms, so the contribution here is $O(n^2 K^2)$. Similarly, the contribution of the third type of cycle is $O(n^2)$, and the fourth type of cycle is $O(n K^2)$, so we can put it all together to get

$$\mathbf{E} \mathrm{tr}(M^4) \leq O(n^3) + O(n^2 K^2).$$

In particular, if we make the hypothesis $K = O(\sqrt{n})$, then we have

$$\mathbf{E} \mathrm{tr}(M^4) \leq O(n^3),$$

and thus by Markov's inequality (1.13) we see that for any $\varepsilon > 0$, $\text{tr}(M^4) \leq O_\varepsilon(n^3)$ with probability at least $1 - \varepsilon$. Applying (2.66), this leads to the upper bound

$$\|M\|_{\text{op}} \leq O_\varepsilon(n^{3/4})$$

with probability at least $1 - \varepsilon$; a similar argument shows that for any fixed $\varepsilon > 0$, one has

$$\|M\|_{\text{op}} \leq n^{3/4+\varepsilon}$$

with high probability. This is better than the upper bound obtained from the second moment method, but still non-optimal.

Exercise 2.3.7. If $K = o(\sqrt{n})$, use the above argument to show that

$$(\mathbf{E}\|M\|_{\text{op}}^4)^{1/4} \geq (2^{1/4} + o(1))\sqrt{n},$$

which in some sense improves upon (2.68) by a factor of $2^{1/4}$. In particular, if $K = O(1)$, conclude that the median of $\|M\|_{\text{op}}$ is at least $(2^{1/4} + o(1))\sqrt{n}$.

Now let us take a quick look at the sixth moment, again with the running assumption of a Wigner matrix in which all entries have mean zero, unit variance, and bounded in magnitude by K . We have

$$\mathbf{E} \text{tr}(M^6) = \sum_{1 \leq i_1, \dots, i_6 \leq n} \mathbf{E} \xi_{i_1 i_2} \cdots \xi_{i_5 i_6} \xi_{i_6 i_1},$$

a sum over cycles of length 6 in $\{1, \dots, n\}$. Again, most of the summands here vanish; the only ones which do not are those cycles in which each edge occurs at least twice (so in particular, there are at most three distinct edges).

Classifying all the types of cycles that could occur here is somewhat tedious, but it is clear that there are going to be $O(1)$ different types of cycles. But we can organise things by the multiplicity of each edge, leaving us with four classes of cycles to deal with:

- (i) Cycles in which there are three distinct edges, each occurring two times.
- (ii) Cycles in which there are two distinct edges, one occurring twice and one occurring four times.
- (iii) Cycles in which there are two distinct edges, each occurring three times²⁰.
- (iv) Cycles in which a single edge occurs six times.

It is not hard to see that summands coming from the first type of cycle give a contribution of 1, and there are $O(n^4)$ of these (because such cycles span at most four vertices). Similarly, the second and third types of cycles

²⁰Actually, this case ends up being impossible, due to a “bridges of Königsberg” type of obstruction, but we will retain it for this discussion.

give a contribution of $O(K^2)$ per summand, and there are $O(n^3)$ summands; finally, the fourth type of cycle gives a contribution of $O(K^4)$, with $O(n^2)$ summands. Putting this together we see that

$$\mathbf{E} \operatorname{tr}(M^6) \leq O(n^4) + O(n^3 K^2) + O(n^2 K^4);$$

so, in particular, if we assume $K = O(\sqrt{n})$ as before, we have

$$\mathbf{E} \operatorname{tr}(M^6) \leq O(n^4)$$

and if we then use (2.66) as before we see that

$$\|M\|_{\text{op}} \leq O_\varepsilon(n^{2/3})$$

with probability $1 - \varepsilon$, for any $\varepsilon > 0$; so we are continuing to make progress towards what we suspect (from the epsilon net argument) to be the correct bound of $n^{1/2}$.

Exercise 2.3.8. If $K = o(\sqrt{n})$, use the above argument to show that

$$(\mathbf{E} \|M\|_{\text{op}}^6)^{1/6} \geq (5^{1/6} + o(1))\sqrt{n}.$$

In particular, if $K = O(1)$, conclude that the median of $\|M\|_{\text{op}}$ is at least $(5^{1/6} + o(1))\sqrt{n}$. Thus this is a (slight) improvement over Exercise 2.3.7.

Let us now consider the general k^{th} moment computation under the same hypotheses as before, with k an even integer, and make some modest attempt to track the dependency of the constants on k . Again, we have

$$(2.70) \quad \mathbf{E} \operatorname{tr}(M^k) = \sum_{1 \leq i_1, \dots, i_k \leq n} \mathbf{E} \xi_{i_1 i_2} \dots \xi_{i_k i_1},$$

which is a sum over cycles of length k . Again, the only non-vanishing expectations are those for which each edge occurs twice; in particular, there are at most $k/2$ edges, and thus at most $k/2 + 1$ vertices.

We divide the cycles into various classes, depending on which edges are equal to each other. (More formally, a class is an equivalence relation \sim on a set of k labels, say $\{1, \dots, k\}$ in which each equivalence class contains at least two elements, and a cycle of k edges $\{i_1, i_2\}, \dots, \{i_k, i_1\}$ lies in the class associated to \sim when we have that $\{i_j, i_{j+1}\} = \{i_{j'}, i_{j'+1}\}$ iff $j \sim j'$, where we adopt the cyclic notation $i_{k+1} := i_1$.)

How many different classes could there be? We have to assign up to $k/2$ labels to k edges, so a crude upper bound here is $(k/2)^k$.

Now consider a given class of cycle. It has j edges e_1, \dots, e_j for some $1 \leq j \leq k/2$, with multiplicities a_1, \dots, a_j , where a_1, \dots, a_j are at least 2 and add up to k . The j edges span at most $j + 1$ vertices; indeed, in addition to the first vertex i_1 , one can specify all the other vertices by looking at the first appearance of each of the j edges e_1, \dots, e_j in the path from i_1 to i_k , and recording the final vertex of each such edge. From this, we see that the

total number of cycles in this particular class is at most n^{j+1} . On the other hand, because each ξ_{ij} has mean zero, unit variance, and is bounded by K , the a^{th} moment of this coefficient is at most K^{a-2} for any $a \geq 2$. Thus each summand in (2.70) coming from a cycle in this class has magnitude at most

$$K^{a_1-2} \dots K^{a_j-2} = K^{a_1+\dots+a_j-2j} = K^{k-2j}.$$

Thus the total contribution of this class to (2.70) is $n^{j+1}K^{k-2j}$, which we can upper bound by

$$\max(n^{\frac{k}{2}+1}, n^2 K^{k-2}) = n^{k/2+1} \max(1, K/\sqrt{n})^{k-2}.$$

Summing up over all classes, we obtain the (somewhat crude) bound

$$\mathbf{E} \operatorname{tr}(M^k) \leq (k/2)^k n^{k/2+1} \max(1, K/\sqrt{n})^{k-2}$$

and thus by (2.66),

$$\mathbf{E} \|M\|_{\text{op}}^k \leq (k/2)^k n^{k/2+1} \max(1, K/\sqrt{n})^{k-2}$$

and so by Markov's inequality (1.13) we have

$$\mathbf{P}(\|M\|_{\text{op}} \geq \lambda) \leq \lambda^{-k} (k/2)^k n^{k/2+1} \max(1, K/\sqrt{n})^{k-2}$$

for all $\lambda > 0$. This, for instance, places the median of $\|M\|_{\text{op}}$ at

$$O(n^{1/k} k \sqrt{n} \max(1, K/\sqrt{n})).$$

We can optimise this in k by choosing k to be comparable to $\log n$, and so we obtain an upper bound of $O(\sqrt{n} \log n \max(1, K/\sqrt{n}))$ for the median; indeed, a slight tweaking of the constants tells us that $\|M\|_{\text{op}} = O(\sqrt{n} \log n \max(1, K/\sqrt{n}))$ with high probability.

The same argument works if the entries have at most unit variance rather than unit variance, thus we have shown

Proposition 2.3.13 (Weak upper bound). *Let M be a random Hermitian matrix, with the upper triangular entries ξ_{ij} , $i \leq j$ being independent with mean zero and variance at most 1, and bounded in magnitude by K . Then $\|M\|_{\text{op}} = O(\sqrt{n} \log n \max(1, K/\sqrt{n}))$ with high probability.*

When $K \leq \sqrt{n}$, this gives an upper bound of $O(\sqrt{n} \log n)$, which is still off by a logarithmic factor from the expected bound of $O(\sqrt{n})$. We will remove this logarithmic loss later in this section.

2.3.5. Computing the moment to top order. Now let us consider the case when $K = o(\sqrt{n})$, and each entry has variance exactly 1. We have an upper bound

$$\mathbf{E} \operatorname{tr}(M^k) \leq (k/2)^k n^{k/2+1};$$

let us try to get a more precise answer here (as in Exercises 2.3.7, 2.3.8). Recall that each class of cycles contributed a bound of $n^{j+1}K^{k-2j}$ to this

expression. If $K = o(\sqrt{n})$, we see that such expressions are $o_k(n^{k/2+1})$ whenever $j < k/2$, where the $o_k()$ notation means that the decay rate as $n \rightarrow \infty$ can depend on k . So the total contribution of all such classes is $o_k(n^{k/2+1})$.

Now we consider the remaining classes with $j = k/2$. For such classes, each equivalence class of edges contains exactly two representatives, thus each edge is repeated exactly once. The contribution of each such cycle to (2.70) is exactly 1, thanks to the unit variance and independence hypothesis. Thus, the total contribution of these classes to $\mathbf{E} \operatorname{tr}(M^k)$ is equal to a purely combinatorial quantity, namely the number of cycles of length k on $\{1, \dots, n\}$ in which each edge is repeated exactly once, yielding $k/2$ unique edges. We are thus faced with the enumerative combinatorics problem of bounding this quantity as precisely as possible.

With $k/2$ edges, there are at most $k/2 + 1$ vertices traversed by the cycle. If there are fewer than $k/2 + 1$ vertices traversed, then there are at most $O_k(n^{k/2}) = o_k(n^{k/2+1})$ cycles of this type, since one can specify such cycles by identifying up to $k/2$ vertices in $\{1, \dots, n\}$ and then matching those coordinates with the k vertices of the cycle. So we set aside these cycles, and only consider those cycles which traverse exactly $k/2 + 1$ vertices. Let us call such cycles (i.e., cycles of length k with each edge repeated exactly once, and traversing exactly $k/2 + 1$ vertices) *non-crossing cycles* of length k in $\{1, \dots, n\}$. Our remaining task is then to count the number of non-crossing cycles.

Example 2.3.14. Let a, b, c, d be distinct elements of $\{1, \dots, n\}$. Then $(i_1, \dots, i_6) = (a, b, c, d, c, b)$ is a non-crossing cycle of length 6, as is (a, b, a, c, a, d) . Any cyclic permutation of a non-crossing cycle is again a non-crossing cycle.

Exercise 2.3.9. Show that a cycle of length k is non-crossing if and only if there exists a tree²¹ in $\{1, \dots, n\}$ of $k/2$ edges and $k/2 + 1$ vertices, such that the cycle lies in the tree and traverses each edge in the tree exactly twice.

Exercise 2.3.10. Let i_1, \dots, i_k be a cycle of length k . Arrange the integers $1, \dots, k$ around a circle. Whenever $1 \leq a < b \leq k$ are such that $i_a = i_b$, with no c between a and b for which $i_a = i_c = i_b$, draw a line segment between a and b . Show that the cycle is non-crossing if and only if the number of line segments is exactly $k/2 - 1$, and the line segments do not cross each other. This may help explain the terminology “non-crossing”.

²¹In graph theory, a *tree* is a finite collection of vertices and (undirected) edges between vertices, which do not contain any cycles.

Now we can complete the count. If k is a positive even integer, define a *Dyck word*²² of length k to be the number of words consisting of left and right parentheses $(,)$ of length k , such that when one reads from left to right, there are always at least as many left parentheses as right parentheses (or in other words, the parentheses define a valid nesting). For instance, the only Dyck word of length 2 is $()$, the two Dyck words of length 4 are $(())$ and $()()$, and the five Dyck words of length 6 are

$$()()(), ((()))(), ()(()), ((())), ((())),$$

and so forth.

Lemma 2.3.15. *The number of non-crossing cycles of length k in $\{1, \dots, n\}$ is equal to $C_{k/2}n(n-1)\dots(n-k/2)$, where $C_{k/2}$ is the number of Dyck words of length k . (The number $C_{k/2}$ is also known as the $(k/2)^{\text{th}}$ Catalan number.)*

Proof. We will give a *bijective proof*. Namely, we will find a way to store a non-crossing cycle as a Dyck word, together with an (ordered) sequence of $k/2 + 1$ distinct elements from $\{1, \dots, n\}$, in such a way that any such pair of a Dyck word and ordered sequence generates exactly one non-crossing cycle. This will clearly give the claim.

So, let us take a non-crossing cycle i_1, \dots, i_k . We imagine traversing this cycle from i_1 to i_2 , then from i_2 to i_3 , and so forth until we finally return to i_1 from i_k . On each leg of this journey, say from i_j to i_{j+1} , we either use an edge that we have not seen before, or else we are using an edge for the second time. Let us say that the leg from i_j to i_{j+1} is an *innovative* leg if it is in the first category, and a *returning* leg otherwise. Thus there are $k/2$ innovative legs and $k/2$ returning legs. Clearly, it is only the innovative legs that can bring us to vertices that we have not seen before. Since we have to visit $k/2 + 1$ distinct vertices (including the vertex i_1 we start at), we conclude that each innovative leg must take us to a new vertex. We thus record, in order, each of the new vertices we visit, starting at i_1 and adding another vertex for each innovative leg; this is an ordered sequence of $k/2 + 1$ distinct elements of $\{1, \dots, n\}$. Next, traversing the cycle again, we write a $($ whenever we traverse an innovative leg, and a $)$ otherwise. This is clearly a Dyck word. For instance, using the examples in Example 2.3.14, the non-crossing cycle (a, b, c, d, c, b) gives us the ordered sequence (a, b, c, d) and the Dyck word $((()))$, while (a, b, a, c, a, d) gives us the ordered sequence (a, b, c, d) and the Dyck word $()()()$.

We have seen that every non-crossing cycle gives rise to an ordered sequence and a Dyck word. A little thought shows that the cycle can be

²²Dyck words are also closely related to *Dyck paths* in enumerative combinatorics.

uniquely reconstructed from this ordered sequence and Dyck word (the key point being that whenever one is performing a returning leg from a vertex v , one is forced to return along the unique innovative leg that discovered v). A slight variant of this thought also shows that every Dyck word of length k and ordered sequence of $k/2+1$ distinct elements gives rise to a non-crossing cycle. This gives the required bijection, and the claim follows. \square

Next, we recall the classical formula for the Catalan number:

Exercise 2.3.11. Establish the recurrence

$$C_{n+1} = \sum_{i=0}^n C_i C_{n-i}$$

for any $n \geq 1$ (with the convention $C_0 = 1$), and use this to deduce that

$$(2.71) \quad C_{k/2} := \frac{k!}{(k/2+1)!(k/2)!}$$

for all $k = 2, 4, 6, \dots$

Exercise 2.3.12. Let k be a positive even integer. Given a string of $k/2$ left parentheses and $k/2$ right parentheses which is *not* a Dyck word, define the *reflection* of this string by taking the first right parenthesis which does not have a matching left parenthesis, and then reversing all the parentheses after that right parenthesis. Thus, for instance, the reflection of $()()()$ is $())))$. Show that there is a bijection between non-Dyck words with $k/2$ left parentheses and $k/2$ right parentheses, and arbitrary words with $k/2 - 1$ left parentheses and $k/2 + 1$ right parentheses. Use this to give an alternate proof of (2.71).

Note that $n(n-1)\dots(n-k/2) = (1+o_k(1))n^{k/2+1}$. Putting all the above computations together, we conclude

Theorem 2.3.16 (Moment computation). *Let M be a real symmetric random matrix, with the upper triangular elements ξ_{ij} , $i \leq j$ jointly independent with mean zero and variance one, and bounded in magnitude by $o(\sqrt{n})$. Let k be a positive even integer. Then we have*

$$\mathbf{E} \operatorname{tr}(M^k) = (C_{k/2} + o_k(1))n^{k/2+1}$$

where $C_{k/2}$ is given by (2.71).

Remark 2.3.17. An inspection of the proof also shows that if we allow the ξ_{ij} to have variance at most one, rather than equal to one, we obtain the upper bound

$$\mathbf{E} \operatorname{tr}(M^k) \leq (C_{k/2} + o_k(1))n^{k/2+1}.$$

Exercise 2.3.13. Show that Theorem 2.3.16 also holds for Hermitian random matrices. (*Hint:* The main point is that with non-crossing cycles, each non-innovative leg goes in the reverse direction to the corresponding innovative leg—why?)

Remark 2.3.18. Theorem 2.3.16 can be compared with the formula

$$\mathbf{E}S^k = (C'_{k/2} + o_k(1))n^{k/2}$$

derived in Section 2.1, where $S = X_1 + \cdots + X_n$ is the sum of n iid random variables of mean zero and variance one, and

$$C'_{k/2} := \frac{k!}{2^{k/2}(k/2)!}.$$

Exercise 2.3.10 shows that $C_{k/2}$ can be interpreted as the number of ways to join k points on the circle by $k/2 - 1$ non-crossing chords. In a similar vein, $C'_{k/2}$ can be interpreted as the number of ways to join k points on the circle by $k/2$ chords which are allowed to cross each other (except at the endpoints). Thus moments of Wigner-type matrices are in some sense the “non-crossing” version of moments of sums of random variables. We will discuss this phenomenon more when we turn to free probability in Section 2.5.

Combining Theorem 2.3.16 with (2.66) we obtain a lower bound

$$\mathbf{E}\|M\|_{\text{op}}^k \geq (C_{k/2} + o_k(1))n^{k/2}.$$

In the bounded case $K = O(1)$, we can combine this with Exercise 2.3.5 to conclude that the median (or mean) of $\|M\|_{\text{op}}$ is at least $(C_{k/2}^{1/k} + o_k(1))\sqrt{n}$. On the other hand, from Stirling’s formula (Section 1.2) we see that $C_{k/2}^{1/k}$ converges to 2 as $k \rightarrow \infty$. Taking k to be a slowly growing function of n , we conclude

Proposition 2.3.19 (Lower Bai-Yin theorem). *Let M be a real symmetric random matrix, with the upper triangular elements ξ_{ij} , $i \leq j$ jointly independent with mean zero and variance one, and bounded in magnitude by $O(1)$. Then the median (or mean) of $\|M\|_{\text{op}}$ is at least $(2 - o(1))\sqrt{n}$.*

Remark 2.3.20. One can in fact obtain an exact asymptotic expansion of the moments $\mathbf{E} \text{tr}(M^k)$ as a polynomial in n , known as the *genus expansion* of the moments. This expansion is, however, somewhat difficult to work with from a combinatorial perspective (except at top order) and will not be used here.

2.3.6. Removing the logarithm. The upper bound in Proposition 2.3.13 loses a logarithm in comparison to the lower bound coming from Theorem 2.3.16. We now discuss how to remove this logarithm.

Suppose that we could eliminate the $o_k(1)$ error in Theorem 2.3.16. Then from (2.66) we would have

$$\mathbf{E}\|M\|_{\text{op}}^k \leq C_{k/2} n^{k/2+1}$$

and hence by Markov's inequality (1.13),

$$\mathbf{P}(\|M\|_{\text{op}} > \lambda) \leq \lambda^{-k} C_{k/2} n^{k/2+1}.$$

Applying this with $\lambda = (2 + \varepsilon)\sqrt{n}$ for some fixed $\varepsilon > 0$, and setting k to be a large multiple of $\log n$, we see that $\|M\|_{\text{op}} \leq (2 + O(\varepsilon))\sqrt{n}$ asymptotically almost surely, which on selecting ε to go to zero slowly in n gives in fact that $\|M\|_{\text{op}} \leq (2 + o(1))\sqrt{n}$ asymptotically almost surely, thus complementing the lower bound in Proposition 2.3.19.

This argument was not rigorous because it did not address the $o_k(1)$ error. Without a more quantitative accounting of this error, one cannot set k as large as $\log n$ without losing control of the error terms; and indeed, a crude accounting of this nature will lose factors of k^k which are unacceptable. Nevertheless, by tightening the hypotheses a little bit and arguing more carefully, we can get a good bound, for k in the region of interest:

Theorem 2.3.21 (Improved moment bound). *Let M be a real symmetric random matrix, with the upper triangular elements ξ_{ij} , $i \leq j$ jointly independent with mean zero and variance one, and bounded in magnitude by $O(n^{0.49})$ (say). Let k be a positive even integer of size $k = O(\log^2 n)$ (say). Then we have*

$$\mathbf{E} \operatorname{tr}(M^k) = C_{k/2} n^{k/2+1} + O(k^{O(1)} 2^k n^{k/2+0.98})$$

where $C_{k/2}$ is given by (2.71). In particular, from the trivial bound $C_{k/2} \leq 2^k$ (which is obvious from the Dyck words definition) one has

$$(2.72) \quad \mathbf{E} \operatorname{tr}(M^k) \leq (2 + o(1))^k n^{k/2+1}.$$

One can of course adjust the parameters $n^{0.49}$ and $\log^2 n$ in the above theorem, but we have tailored these parameters for our application to simplify the exposition slightly.

Proof. We may assume n large, as the claim is vacuous for bounded n .

We again expand using (2.70), and discard all the cycles in which there is an edge that only appears once. The contribution of the non-crossing cycles was already computed in the previous section to be

$$C_{k/2} n(n-1) \dots (n-k/2),$$

which can easily be computed (e.g., by taking logarithms, or using Stirling's formula) to be $(C_{k/2} + o(1))n^{k/2+1}$. So the only task is to show that the net contribution of the remaining cycles is $O(k^{O(1)}2^k n^{k/2+0.98})$.

Consider one of these cycles (i_1, \dots, i_k) ; it has j distinct edges for some $1 \leq j \leq k/2$ (with each edge repeated at least once).

We order the j distinct edges e_1, \dots, e_j by their first appearance in the cycle. Let a_1, \dots, a_j be the multiplicities of these edges, thus the a_1, \dots, a_j are all at least 2 and add up to k . Observe from the moment hypotheses that the moment $\mathbf{E}|\xi_{ij}|^a$ is bounded by $O(n^{0.49})^{a-2}$ for $a \geq 2$. Since $a_1 + \dots + a_j = k$, we conclude that the expression

$$\mathbf{E}\xi_{i_1 i_2} \dots \xi_{i_k i_1}$$

in (2.70) has magnitude at most $O(n^{0.49})^{k-2j}$, and so the net contribution of the cycles that are not non-crossing is bounded in magnitude by

$$(2.73) \quad \sum_{j=1}^{k/2} O(n^{0.49})^{k-2j} \sum_{a_1, \dots, a_j} N_{a_1, \dots, a_j}$$

where a_1, \dots, a_j range over integers that are at least 2 and which add up to k , and N_{a_1, \dots, a_j} is the number of cycles that are not non-crossing and have j distinct edges with multiplicity a_1, \dots, a_j (in order of appearance). It thus suffices to show that (2.73) is $O(k^{O(1)}2^k n^{k/2+0.98})$.

Next, we estimate N_{a_1, \dots, a_j} for a fixed a_1, \dots, a_j . Given a cycle (i_1, \dots, i_k) , we traverse its k legs (which each traverse one of the edges e_1, \dots, e_j) one at a time and classify them into various categories:

- (i) *High-multiplicity legs*, which use an edge e_i whose multiplicity a_i is larger than two.
- (ii) *Fresh legs*, which use an edge e_i with $a_i = 2$ for the first time.
- (iii) *Return legs*, which use an edge e_i with $a_i = 2$ that has already been traversed by a previous fresh leg.

We also subdivide fresh legs into *innovative* legs, which take one to a vertex one has not visited before, and *non-innovative* legs, which take one to a vertex that one has visited before.

At any given point in time when traversing this cycle, we define an *available* edge to be an edge e_i of multiplicity $a_i = 2$ that has already been traversed by its fresh leg, but not by its return leg. Thus, at any given point in time, one travels along either a high-multiplicity leg, a fresh leg (thus creating a new available edge), or one returns along an available edge (thus removing that edge from availability).

Call a return leg starting from a vertex v *forced* if, at the time one is performing that leg, there is only one available edge from v , and *unforced* otherwise (i.e., there are two or more available edges to choose from).

We suppose that there are $l := \#\{1 \leq i \leq j : a_i > 2\}$ high-multiplicity edges among the e_1, \dots, e_j , leading to $j - l$ fresh legs and their $j - l$ return leg counterparts. In particular, the total number of high-multiplicity legs is

$$(2.74) \quad \sum_{a_i > 2} a_i = k - 2(j - l).$$

Since $\sum_{a_i > 2} a_i \geq 3l$, we conclude the bound

$$(2.75) \quad l \leq k - 2j.$$

We assume that there are m non-innovative legs among the $j - l$ fresh legs, leaving $j - l - m$ innovative legs. As the cycle is not non-crossing, we either have $j < k/2$ or $m > 0$.

Similarly, we assume that there are r unforced return legs among the $j - l$ total return legs. We have an important estimate:

Lemma 2.3.22 (Not too many unforced return legs). *We have*

$$r \leq 2(m + \sum_{a_i > 2} a_i).$$

In particular, from (2.74), (2.75), we have

$$r \leq O(k - 2j) + O(m).$$

Proof. Let v be a vertex visited by the cycle which is not the initial vertex i_1 . Then the very first arrival at v comes from a fresh leg, which immediately becomes available. Each departure from v may create another available edge from v , but each subsequent arrival at v will delete an available leg from v , unless the arrival is along a non-innovative or high-multiplicity edge²³. Finally, any returning leg that departs from v will also delete an available edge from v .

This has two consequences. First, if there are no non-innovative or high-multiplicity edges arriving at v , then whenever one arrives at v , there is at most one available edge from v , and so every return leg from v is forced (and there will be only one such return leg). If, instead, there are non-innovative or high-multiplicity edges arriving at v , then we see that the total number of return legs from v is at most one plus the number of such edges. In both cases, we conclude that the number of unforced return legs from v is bounded by twice the number of non-innovative or high-multiplicity edges arriving at v . Summing over v , one obtains the claim. \square

²³Note that one can loop from v to itself and create an available edge, but this is along a non-innovative edge and so is not inconsistent with the previous statements.

Now we return to the task of counting N_{a_1, \dots, a_j} , by recording various data associated to any given cycle (i_1, \dots, i_k) contributing to this number. First, fix m, r . We record the initial vertex i_1 of the cycle, for which there are n possibilities. Next, for each high-multiplicity edge e_i (in increasing order of i), we record all the a_i locations in the cycle where this edge is used; the total number of ways this can occur for each such edge can be bounded above by k^{a_i} , so the total entropy cost here is $k^{\sum_{a_i > 2} a_i} = k^{k-2(j-l)}$. We also record the final endpoint of the first occurrence of the edge e_i for each such i ; this list of l vertices in $\{1, \dots, n\}$ has at most n^l possibilities.

For each innovative leg, we record the final endpoint of that leg, leading to an additional list of $j-l-m$ vertices with at most n^{j-l-m} possibilities.

For each non-innovative leg, we record the position of that leg, leading to a list of m numbers from $\{1, \dots, k\}$, which has at most k^m possibilities.

For each unforced return leg, we record the position of the corresponding fresh leg, leading to a list of r numbers from $\{1, \dots, k\}$, which has at most k^r possibilities.

Finally, we record a Dyck-like word of length k , in which we place a (whenever the leg is innovative, and) otherwise (the brackets need not match here). The total entropy cost here can be bounded above by 2^k .

We now observe that all this data (together with l, m, r) can be used to completely reconstruct the original cycle. Indeed, as one traverses the cycle, the data already tells us which edges are high-multiplicity, which ones are innovative, which ones are non-innovative, and which ones are return legs. In all edges in which one could possibly visit a new vertex, the location of that vertex has been recorded. For all unforced returns, the data tells us which fresh leg to backtrack upon to return to. Finally, for forced returns, there is only one available leg to backtrack to, and so one can reconstruct the entire cycle from this data.

As a consequence, for fixed l, m and r , there are at most

$$nk^{k-2(j-l)}n^l n^{j-l-m} k^m k^r 2^k$$

contributions to N_{a_1, \dots, a_j} ; using (2.75), (2.3.22) we can bound this by

$$k^{O(k-2j)+O(m)} n^{j-m+1} 2^k.$$

Summing over the possible values of m, r (recalling that we either have $j < k/2$ or $m > 0$, and also that $k = O(\log^2 n)$) we obtain

$$N_{a_1, \dots, a_j} \leq k^{O(k-2j)+O(1)} n^{\max(j+1, k/2)} 2^k.$$

The expression (2.73) can then be bounded by

$$2^k \sum_{j=1}^{k/2} O(n^{0.49})^{k-2j} k^{O(k-2j)+O(1)} n^{\max(j+1, k/2)} \sum_{a_1, \dots, a_j} 1.$$

When j is exactly $k/2$, then all the a_1, \dots, a_j must equal 2, and so the contribution of this case simplifies to $2^k k^{O(1)} n^{k/2}$. For $j < k/2$, the numbers $a_1 - 2, \dots, a_j - 2$ are non-negative and add up to $k - 2j$, and so the total number of possible values for these numbers (for fixed j) can be bounded crudely by $j^{k-2j} \leq k^{k-2j}$ (for instance). Putting all this together, we can bound (2.73) by

$$2^k [k^{O(1)} n^{k/2} + \sum_{j=1}^{k/2-1} O(n^{0.49})^{k-2j} k^{O(k-2j)+O(1)} n^{j+1} k^{k-2j}],$$

which simplifies by the geometric series formula (and the hypothesis $k = O(\log^2 n)$) to

$$O(2^k k^{O(1)} n^{k/2+0.98}),$$

as required. \square

We can use this to conclude the following matching upper bound to Proposition 2.3.19, due to Bai and Yin [BaYi1988]:

Theorem 2.3.23 (Weak Bai-Yin theorem, upper bound). *Let $M = (\xi_{ij})_{1 \leq i, j \leq n}$ be a real symmetric matrix whose entries all have the same distribution ξ , with mean zero, variance one, and fourth moment $O(1)$. Then for every $\varepsilon > 0$ independent of n , one has $\|M\|_{\text{op}} \leq (2 + \varepsilon)\sqrt{n}$ asymptotically almost surely. In particular, $\|M\|_{\text{op}} \leq (2 + o(1))\sqrt{n}$ asymptotically almost surely; as another consequence, the median of $\|M\|_{\text{op}}$ is at most $(2 + o(1))\sqrt{n}$. (If ξ is bounded, we see, in particular, from Proposition 2.3.19 that the median is in fact equal to $(2 + o(1))\sqrt{n}$.)*

The fourth moment hypothesis is best possible, as seen in the discussion after Theorem 2.3.8. We will discuss some generalisations and improvements of this theorem in other directions below.

Proof. To obtain Theorem 2.3.23 from Theorem 2.3.21 we use the truncation method. We split each ξ_{ij} as $\xi_{ij, \leq n^{0.49}} + \xi_{ij, > n^{0.49}}$ in the usual manner, and split $M = M_{\leq n^{0.49}} + M_{> n^{0.49}}$ accordingly. We would like to apply Theorem 2.3.21 to $M_{\leq n^{0.49}}$, but unfortunately the truncation causes some slight adjustment to the mean and variance of the $\xi_{ij, \leq n^{0.49}}$. The variance is not much of a problem; since ξ_{ij} had variance 1, it is clear that $\xi_{ij, \leq n^{0.49}}$ has variance at most 1, and it is easy to see that reducing the variance only

serves to improve the bound (2.72). As for the mean, we use the mean zero nature of ξ_{ij} to write

$$\mathbf{E}\xi_{ij,\leq n^{0.49}} = -\mathbf{E}\xi_{ij,>n^{0.49}}.$$

To control the right-hand side, we use the trivial inequality $|\xi_{ij,\leq n^{0.49}}| \leq n^{-3 \times 0.49} |\xi_{ij}|^4$ and the bounded fourth moment hypothesis to conclude that

$$\mathbf{E}\xi_{ij,\leq n^{0.49}} = O(n^{-1.47}).$$

Thus we can write $M_{\leq n^{0.49}} = \tilde{M}_{\leq n^{0.49}} + \mathbf{E}M_{\leq n^{0.49}}$, where $\tilde{M}_{\leq n^{0.49}}$ is the random matrix with coefficients

$$\tilde{\xi}_{ij,\leq n^{0.49}} := \xi_{ij,\leq n^{0.49}} - \mathbf{E}\xi_{ij,\leq n^{0.49}}$$

and $\mathbf{E}M_{\leq n^{0.49}}$ is a matrix whose entries have magnitude $O(n^{-1.47})$. In particular, by Schur's test this matrix has operator norm $O(n^{-0.47})$, and so by the triangle inequality

$$\|M\|_{\text{op}} \leq \|\tilde{M}_{\leq n^{0.49}}\|_{\text{op}} + \|M_{>n^{0.49}}\|_{\text{op}} + O(n^{-0.47}).$$

The error term $O(n^{-0.47})$ is clearly negligible for n large, and it will suffice to show that

$$(2.76) \quad \|\tilde{M}_{\leq n^{0.49}}\|_{\text{op}} \leq (2 + \varepsilon/3)\sqrt{n}$$

and

$$(2.77) \quad \|M_{>n^{0.49}}\|_{\text{op}} \leq \frac{\varepsilon}{3}\sqrt{n}$$

asymptotically almost surely.

We first show (2.76). We can now apply Theorem 2.3.21 to conclude that

$$\mathbf{E}\|\tilde{M}_{\leq n^{0.49}}\|_{\text{op}}^k \leq (2 + o(1))^k n^{k/2+1}$$

for any $k = O(\log^2 n)$. In particular, we see from Markov's inequality (1.13) that (2.76) holds with probability at most

$$\left(\frac{2 + o(1)}{2 + \varepsilon/3}\right)^k n.$$

Setting k to be a large enough multiple of $\log n$ (depending on ε), we thus see that this event (2.76) indeed holds asymptotically almost surely²⁴.

Now we turn to (2.77). The idea here is to exploit the sparseness of the matrix $M_{>n^{0.49}}$. First let us dispose of the event that one of the entries ξ_{ij} has magnitude larger than $\frac{\varepsilon}{3}\sqrt{n}$ (which would certainly cause (2.77) to fail). By the union bound, the probability of this event is at most

$$n^2 \mathbf{P}\left(|\xi| \geq \frac{\varepsilon}{3}\sqrt{n}\right).$$

²⁴Indeed, one can ensure it happens with overwhelming probability, by letting $k/\log n$ grow slowly to infinity.

By the fourth moment bound on ξ and dominated convergence, this expression goes to zero as $n \rightarrow \infty$. Thus, asymptotically almost surely, all entries are less than $\frac{\varepsilon}{3}\sqrt{n}$.

Now let us see how many non-zero entries there are in $M_{>n^{0.49}}$. By Markov's inequality (1.13) and the fourth moment hypothesis, each entry has a probability $O(n^{-4 \times 0.49}) = O(n^{-1.96})$ of being non-zero; by the first moment method, we see that the expected number of entries is $O(n^{0.04})$. As this is much less than n , we expect it to be unlikely that any row or column has more than one entry. Indeed, from the union bound and independence, we see that the probability that any given row and column has at least two non-zero entries is at most

$$n^2 \times O(n^{-1.96})^2 = O(n^{-1.92})$$

and so by the union bound again, we see that with probability at least $1 - O(n^{-0.92})$ (and in particular, asymptotically almost surely), none of the rows or columns have more than one non-zero entry. As the entries have magnitude at most $\frac{\varepsilon}{3}\sqrt{n}$, the bound (2.77) now follows from Schur's test, and the claim follows. \square

We can upgrade the asymptotic almost sure bound to almost sure boundedness:

Theorem 2.3.24 (Strong Bai-Yin theorem, upper bound). *Let ξ be a real random variable with mean zero, variance 1, and finite fourth moment, and for all $1 \leq i \leq j$, let ξ_{ij} be an iid sequence with distribution ξ , and set $\xi_{ji} := \xi_{ij}$. Let $M_n := (\xi_{ij})_{1 \leq i, j \leq n}$ be the random matrix formed by the top left $n \times n$ block. Then almost surely one has $\limsup_{n \rightarrow \infty} \|M_n\|_{\text{op}}/\sqrt{n} \leq 2$.*

Exercise 2.3.14. By combining the above results with Proposition 2.3.19 and Exercise 2.3.5, show that with the hypotheses of Theorem 2.3.24 with ξ bounded, one has $\lim_{n \rightarrow \infty} \|M_n\|_{\text{op}}/\sqrt{n} = 2$ almost surely²⁵.

Proof. We first give ourselves an epsilon of room (cf. [Ta2010, §2.7]). It suffices to show that for each $\varepsilon > 0$, one has

$$(2.78) \quad \limsup_{n \rightarrow \infty} \|M_n\|_{\text{op}}/\sqrt{n} \leq 2 + \varepsilon$$

almost surely.

Next, we perform dyadic sparsification (as was done in the proof of the strong law of large numbers, Theorem 2.1.8). Observe that any minor of a matrix has its operator norm bounded by that of the larger matrix, and so $\|M_n\|_{\text{op}}$ is increasing in n . Because of this, it will suffice to show (2.78) almost surely for n restricted to a *lacunary* sequence, such as $n = n_m :=$

²⁵The same claim is true without the boundedness hypothesis; we will see this in Section 2.4.

$\lfloor (1 + \varepsilon)^m \rfloor$ for $m = 1, 2, \dots$, as the general case then follows by rounding n upwards to the nearest n_m (and adjusting ε a little bit as necessary).

Once we sparsified, it is now safe to apply the Borel-Cantelli lemma (Exercise 1.1.1), and it will suffice to show that

$$\sum_{m=1}^{\infty} \mathbf{P}(\|M_{n_m}\|_{\text{op}} \geq (2 + \varepsilon)\sqrt{n_m}) < \infty.$$

To bound the probabilities $\mathbf{P}(\|M_{n_m}\|_{\text{op}} \geq (2 + \varepsilon)\sqrt{n_m})$, we inspect the proof of Theorem 2.3.23. Most of the contributions to this probability decay polynomially in n_m (i.e., are of the form $O(n_m^{-c})$ for some $c > 0$) and so are summable. The only contribution which can cause difficulty is the contribution of the event that one of the entries of M_{n_m} exceeds $\frac{\varepsilon}{3}\sqrt{n_m}$ in magnitude; this event was bounded by

$$n_m^2 \mathbf{P}(|\xi| \geq \frac{\varepsilon}{3}\sqrt{n_m}).$$

But if one sums over m using Fubini's theorem and the geometric series formula, we see that this expression is bounded by $O_\varepsilon(\mathbf{E}|\xi|^4)$, which is finite by hypothesis, and the claim follows. \square

Now we discuss some variants and generalisations of the Bai-Yin result.

First, we note that the results stated above require the diagonal and off-diagonal terms to have the same distribution. This is not the case for important ensembles such as the Gaussian Orthogonal Ensemble (GOE), in which the diagonal entries have twice as much variance as the off-diagonal ones. But this can easily be handled by considering the diagonal separately. For instance, consider a diagonal matrix $D = \text{diag}(\xi_{11}, \dots, \xi_{nn})$ where the $\xi_{ii} \equiv \xi$ are identically distributed with finite second moment. The operator norm of this matrix is just $\sup_{1 \leq i \leq n} |\xi_{ii}|$, and so by the union bound

$$\mathbf{P}(\|D\|_{\text{op}} > \varepsilon\sqrt{n}) \leq n\mathbf{P}(|\xi| > \varepsilon\sqrt{n}).$$

From the finite second moment and dominated convergence, the right-hand side is $o_\varepsilon(1)$, and so we conclude that for every fixed $\varepsilon > 0$, $\|D\|_{\text{op}} \leq \varepsilon\sqrt{n}$ asymptotically almost surely; diagonalising, we conclude that $\|D\|_{\text{op}} = o(\sqrt{n})$ asymptotically almost surely. Because of this and the triangle inequality, we can modify the diagonal by any amount with identical distribution and bounded second moment (a similar argument also works for non-identical distributions if one has uniform control of some moment beyond the second, such as the fourth moment) while only affecting all operator norms by $o(\sqrt{n})$.

Exercise 2.3.15. Modify this observation to extend the weak and strong Bai-Yin theorems to the case where the diagonal entries are allowed to have

different distribution than the off-diagonal terms, and need not be independent of each other or of the off-diagonal terms, but have uniformly bounded fourth moment.

Second, it is a routine matter to generalise the Bai-Yin result from real symmetric matrices to Hermitian matrices, basically for the same reasons that Exercise 2.3.13 works. We leave the details to the interested reader.

The Bai-Yin results also hold for iid random matrices, where $\xi_{ij} \equiv \xi$ has mean zero, unit variance, and bounded fourth moment; this is a result of Yin, Bai, and Krishnaiah [YiBaKr1988], building upon the earlier work of Geman [Ge1980]. Because of the lack of symmetry, the eigenvalues need not be real, and the bounds (2.66) no longer apply. However, there is a substitute, namely the bound

$$(2.79) \quad \|M\|_{\text{op}}^k \leq \text{tr}((MM^*)^{k/2}) \leq n\|M\|_{\text{op}}^k,$$

valid for any $n \times n$ matrix M with complex entries and every even positive integer k .

Exercise 2.3.16. Prove (2.79).

It is possible to adapt all of the above moment calculations for $\text{tr}(M^k)$ in the symmetric or Hermitian cases to give analogous results for $\text{tr}((MM^*)^{k/2})$ in the non-symmetric cases; we do not give the details here, but mention that the cycles now go back and forth along a bipartite graph with n vertices in each class, rather than in the complete graph on n vertices, although this ends up not affecting the enumerative combinatorics significantly. Another way of viewing this is through the simple observation that the operator norm of a non-symmetric matrix M is equal to the operator norm of the *augmented matrix*

$$(2.80) \quad \tilde{M} := \begin{pmatrix} 0 & M \\ M^* & 0 \end{pmatrix},$$

which is a $2n \times 2n$ Hermitian matrix. Thus, one can to some extent identify an $n \times n$ iid matrix M with a $2n \times 2n$ Wigner-type matrix \tilde{M} , in which two $n \times n$ blocks of that matrix are set to zero.

Exercise 2.3.17. If M has singular values $\sigma_1, \dots, \sigma_n$, show that \tilde{M} has eigenvalues $\pm\sigma_1, \dots, \pm\sigma_n$. This suggests that the theory of the singular values of an iid matrix should resemble to some extent the theory of eigenvalues of a Wigner matrix; we will see several examples of this phenomenon in later sections.

When one assumes more moment conditions on ξ than bounded fourth moment, one can obtain substantially more precise asymptotics on $\text{tr}(M^k)$ than given by results such as Theorem 2.3.21, particularly if one also assumes

that the underlying random variable ξ is symmetric (i.e., $\xi \equiv -\xi$). At a practical level, the advantage of symmetry is that it allows one to assume that the high-multiplicity edges in a cycle are traversed an *even* number of times; see the following exercise.

Exercise 2.3.18. Let X be a bounded real random variable. Show that X is symmetric if and only if $\mathbf{E}X^k = 0$ for all positive odd integers k .

Next, extend the previous result to the case when X is sub-Gaussian rather than bounded. (*Hint:* The slickest way to do this is via the characteristic function e^{itX} and analytic continuation; it is also instructive to find a “real-variable” proof that avoids the use of this function.)

By using these methods, it is in fact possible to show that under various hypotheses, $\|M\|_{\text{op}}$ is concentrated in the range $[2\sqrt{n} - O(n^{-1/6}), 2\sqrt{n} + O(n^{-1/6})]$, and even to get a universal distribution for the normalised expression $(\|M\|_{\text{op}} - 2\sqrt{n})n^{1/6}$, known as the *Tracy-Widom law*. See [So1999] for details. There have also been a number of subsequent variants and refinements of this result (as well as counterexamples when not enough moment hypotheses are assumed); see²⁶ [So2004, SoFy2005, Ru2007, Pe2006, Vu2007, PeSo2007, Pe2009, Kh2009, TaVu2009c].

2.4. The semicircular law

We can now turn attention to one of the centerpiece universality results in random matrix theory, namely the *Wigner semicircle law* for Wigner matrices. Recall from Section 2.3 that a *Wigner Hermitian matrix ensemble* is a random matrix ensemble $M_n = (\xi_{ij})_{1 \leq i, j \leq n}$ of Hermitian matrices (thus $\xi_{ij} = \overline{\xi_{ji}}$; this includes *real symmetric matrices* as an important special case), in which the upper-triangular entries ξ_{ij} , $i > j$ are iid complex random variables with mean zero and unit variance, and the diagonal entries ξ_{ii} are iid real variables, independent of the upper-triangular entries, with bounded mean and variance. Particular special cases of interest include the *Gaussian Orthogonal Ensemble (GOE)*, the *symmetric random sign matrices* (aka *symmetric Bernoulli ensemble*), and the *Gaussian Unitary Ensemble (GUE)*.

In Section 2.3 we saw that the operator norm of M_n was typically of size $O(\sqrt{n})$, so it is natural to work with the normalised matrix $\frac{1}{\sqrt{n}}M_n$. Accordingly, given any $n \times n$ Hermitian matrix M_n , we can form the (normalised)

²⁶Similar results for some non-independent distributions are also available, see e.g., the paper [DeGi2007], which (like many of the other references cited above) builds upon the original work of Tracy and Widom [TrWi2002] that handled special ensembles such as GOE and GUE.