

## 1 Concentration inequalities

The style of this course will be to consider some phenomenon with real numbers and then to explore what happens when we try to translate it to the world of symmetric matrices. Our first example concerns tail bounds for sums of independent random variables.

**The classical setting.** Let  $X_1, X_2, \dots, X_n$  be independent random variables and denote  $\mathbf{X} := X_1 + \dots + X_n$ . Let us additionally suppose that  $|X_i - \mathbb{E} X_i| \leq L$  holds almost surely for every  $i = 1, 2, \dots, n$ . Then we have the classical exponential concentration inequality: For every  $t \geq 0$ ,

$$\mathbb{P}[\mathbf{X} \geq \mathbb{E} \mathbf{X} + t\sqrt{n}] \leq \exp\left(\frac{-t^2}{2L^2}\right). \quad (1.1)$$

By translation, we may assume that  $\mathbb{E} X_i = 0$  and  $|X_i| \leq L$  for every  $i$ . The standard argument proceeds by applying Markov's inequality to the Laplace transform  $\beta \mapsto e^{\beta \mathbf{X}}$ :

$$\mathbb{P}[\mathbf{X} \geq t\sqrt{n}] = \mathbb{P}[e^{\beta \mathbf{X}} \geq e^{\beta t\sqrt{n}}] \leq e^{-\beta t\sqrt{n}} \mathbb{E}[e^{\beta \mathbf{X}}], \quad (1.2)$$

where  $\beta > 0$  is some parameter we will optimize over later.

The advantage of the exponential representation is that it allows us to use mutual independence:

$$\mathbb{E}[e^{\beta \mathbf{X}}] = \mathbb{E}[e^{\beta(X_1 + \dots + X_n)}] = \prod_{i=1}^n \mathbb{E}[e^{\beta X_i}]. \quad (1.3)$$

Note that

$$1 + x \leq e^x \leq 1 + x + x^2, \quad \forall x \in [-1, 1]. \quad (1.4)$$

Therefore if we ensure that  $\beta \leq 1/L$ , we have

$$\mathbb{E}[e^{\beta X_i}] \leq \mathbb{E}[1 + \beta X_i + \beta^2 X_i^2] = 1 + \beta^2 \mathbb{E}[X_i^2] \leq 1 + \beta^2 L^2 \leq e^{\beta^2 L^2}.$$

Plugging this into (1.2) and (1.3) gives: For  $0 < \beta \leq 1/L$ ,

$$\mathbb{P}[\mathbf{X} \geq t\sqrt{n}] \leq e^{-\beta t\sqrt{n} + \beta^2 L^2 n}. \quad (1.5)$$

Now setting  $\beta := \frac{t}{2L^2\sqrt{n}}$  gives

$$\mathbb{P}[\mathbf{X} \geq t\sqrt{n}] \leq \exp\left(\frac{-t^2}{2L^2}\right),$$

as long as  $t \leq 2L\sqrt{n}$ . But under our assumptions,  $\mathbf{X} \leq nL$  almost surely, hence the bound is true for all  $t \geq 0$ .

**Sums of independent symmetric matrices.** Let us now suppose that  $A_1, \dots, A_n \in \mathbb{M}_d$  is a family of independent random symmetric matrices and  $A := A_1 + \dots + A_n$ . Finding the correct analog of (1.1) might, in general, be driven by the application we have in mind.

A natural approach is to replace numbers by matrices and inequalities by the Loewner order. This would yield the assumptions: Almost surely,

$$-L \cdot I \leq A_i - \mathbb{E} A_i \leq L \cdot I, \quad \forall i = 1, \dots, n, \quad (1.6)$$

and the conclusion

$$\mathbb{P} [A \geq \mathbb{E} A + t\sqrt{n} \cdot I] \leq \exp\left(\frac{-t^2}{2L^2}\right).$$

Actually, this conclusion is substantially weaker than what we might hope for. Indeed, let us assume that  $\mathbb{E} A_i = 0$  for all  $i$ . Then the “bad” event  $\{A \geq t\sqrt{n} \cdot I\}$  entails *all* the eigenvalues of  $A$  being large, whereas we might hope to argue that *none of them* are large.

Our naive translation of (1.1) is due to the fact that  $\leq$  is not a total order:  $A \geq B$  is not the logical negation of  $A < B$ . If we rewrite (1.1) as

$$\mathbb{P} [X \leq \mathbb{E} X + t\sqrt{n}] \geq 1 - \exp\left(\frac{-t^2}{2L^2}\right),$$

then taking again the straightforward translation to matrices yields

$$\mathbb{P} [A \leq \mathbb{E} A + t\sqrt{n} \cdot I] \geq 1 - \exp\left(\frac{-t^2}{2L^2}\right). \quad (1.7)$$

By translation, we may again assume that  $\mathbb{E} A_i = 0$  for each  $i$ , and then (1.6) can be restated as: Almost surely,

$$\|A_i\| \leq L, \quad \forall i = 1, \dots, n,$$

where

$$\|M\| := \sup_{x \neq 0} \frac{\|Mx\|_2}{\|x\|_2}$$

denotes the *operator norm* of  $M \in \mathbb{M}_d$ .

When  $A$  is symmetric,  $\|A\| = \max(|\lambda_1|, \dots, |\lambda_d|)$ , where  $\lambda_1, \dots, \lambda_d$  are the eigenvalues of  $A$ . And since  $\mathbb{E} A = 0$ , the bound (1.7) can be rewritten as

$$\mathbb{P} [\lambda_{\max}(A) \geq t\sqrt{n}] \leq \exp\left(\frac{-t^2}{2L^2}\right),$$

where  $\lambda_{\max}(A)$  denotes the maximum eigenvalue of  $A$ .

Note that the RHS is too optimistic: Suppose each  $A_i$  is diagonal with the diagonal containing independent entries. Then the diagonal of  $A$  contains  $d$  independent trials, and we can reasonably expect that in the worst case, applying (1.1) to each entry and taking a union bound will lose a factor of  $d$ .

As we will see, the conclusion is true with this correction:

$$\mathbb{P} [\lambda_{\max}(A) \geq t\sqrt{n}] \leq d \cdot \exp\left(\frac{-t^2}{2L^2}\right). \quad (1.8)$$

**Cautious optimism.** So far everything is proceeding smoothly, so we might as well continue the analogy (cf. (1.2)):

$$\mathbb{P} \left[ \lambda_{\max}(A) \geq t\sqrt{n} \right] \leq \mathbb{P} \left[ e^{\beta\lambda_{\max}(A)} \geq e^{\beta t\sqrt{n}} \right] \leq e^{-\beta t\sqrt{n}} \mathbb{E} \left[ e^{\beta\lambda_{\max}(A)} \right].$$

Here we run into a hiccup: There is no nice formula for  $\lambda_{\max}(A_1 + \dots + A_n)$  in terms of  $\lambda_{\max}(A_1), \dots, \lambda_{\max}(A_n)$ .

But since we willing to lose a factor of  $d$  in the tail bound (1.8), we might consider bounding the maximum eigenvalue by a sum:

$$\mathbb{P} \left[ e^{\beta\lambda_{\max}(A)} \geq e^{\beta t\sqrt{n}} \right] \leq \mathbb{P} \left[ \text{Tr}(e^{\beta A}) \geq e^{\beta t\sqrt{n}} \right], \quad (1.9)$$

where we recall the trace of a matrix:

$$\text{Tr}(A) := \sum_{i=1}^n A_{ii}.$$

Let's take a moment to define  $e^X$  when  $X$  is a matrix.

**Spectral functions.** For a symmetric matrix  $A$ , let  $\text{spec}(A)$  be the set of  $A$ 's eigenvalues. For any function  $f : I \rightarrow \mathbb{R}$ , we can extend  $f$  to symmetric matrices  $A$  with  $\text{spec}(A) \subseteq I$  as follows. Write  $A = \sum_{i=1}^d \lambda_i u_i u_i^T$  where the  $\{u_i\}$  are orthonormal, and define

$$f(A) := \sum_{i=1}^d f(\lambda_i) u_i u_i^T.$$

It is easily checked that  $f(A)$  does not depend on the choice of an eigenvector basis  $\{u_1, \dots, u_d\}$  for  $A$ . When  $f(t) = \sum_{j=1}^m c_j t^j$  is a polynomial, it holds that  $f(A) = \sum_{j=1}^m c_j A^j$ , as one might hope. This continues to hold when  $f$  has a convergent Taylor series. For  $f(t) = e^t$ , it holds that for all symmetric matrices  $A$ ,

$$e^A := f(A) = \sum_{n \geq 0} \frac{A^n}{n!}.$$

For any  $A, B \in \mathbb{M}_d$ , a simple calculation gives

$$\text{Tr}(AB) = \sum_{i,j} A_{ij} B_{ji} = \text{Tr}(BA).$$

This is called the cyclic property of the trace. In particular, this implies that the trace is *invariant under similarities*: For any invertible matrix  $P$ , we have

$$\text{Tr}(PAP^{-1}) = \text{Tr}(A).$$

For a symmetric matrix  $A$  with eigenvalues  $\lambda_1, \dots, \lambda_d$ , we can diagonalize  $A$  and obtain  $\text{Tr}(A) = \lambda_1 + \dots + \lambda_d$ . This justifies (1.9).

Now we can apply Markov's inequality:

$$\mathbb{P} \left[ \text{Tr}(e^{\beta A}) \geq e^{\beta t\sqrt{n}} \right] \leq e^{-\beta t\sqrt{n}} \mathbb{E} \left[ \text{Tr}(e^{\beta A}) \right] = e^{-\beta t\sqrt{n}} \mathbb{E} \left[ \text{Tr} \left( e^{\beta(A_1 + \dots + A_n)} \right) \right]. \quad (1.10)$$

**The trail goes cold.** Unfortunately, non-commutativity now rears its head more seriously: In analogy with the classical setting, we want to utilize independence in evaluating a product. But we don't have a product:  $e^{A+B} \neq e^A e^B$  unless  $A$  and  $B$  commute.

## 1.1 A sprinkle of magic

It's not surprising that we got stuck; many weaker, more difficult approaches were tried before the work of Ahlswede and Winter (2002). They employed the Golden-Thompson inequality.

**Lemma 1.1** (Golden-Thompson). *For all symmetric matrices  $A, B \in \mathbb{M}_d$ , it holds that*

$$\mathrm{Tr}(e^{A+B}) \leq \mathrm{Tr}(e^A e^B).$$

And, indeed, it is not too difficult to establish our desired tail bound with this in hand.

$$\begin{aligned} \mathbb{E} \left[ \mathrm{Tr} \left( e^{\beta(A_1 + \dots + A_n)} \right) \right] &\leq \mathbb{E} \left[ \mathrm{Tr} \left( e^{\beta(A_1 + \dots + A_{n-1})} e^{\beta A_n} \right) \right] \\ &= \mathbb{E} \left[ \mathrm{Tr} \left( e^{\beta(A_1 + \dots + A_{n-1})} \mathbb{E}[e^{\beta A_n}] \right) \right], \end{aligned} \quad (1.11)$$

where the last equality follows from independence (if  $A, B$  are independent random matrices, it holds that  $\mathbb{E}[AB] = \mathbb{E}[A] \mathbb{E}[B]$ ).

For  $A, B$  both symmetric, it holds that

$$|\mathrm{Tr}(AB)| \leq \|A\| \cdot \|B\|_1. \quad (1.12)$$

where  $\|A\| = \|(\lambda_1(A), \dots, \lambda_d(A))\|_\infty$  denotes the operator norm, and  $\|B\|_1 = \|(\lambda_1(B), \dots, \lambda_d(B))\|_1$  is the  $\ell_1$ -norm of the eigenvalues values of  $B$ . The proof of (1.12) is straightforward: Write  $B = \sum_i \lambda_i u_i u_i^T$  where  $\{u_i\}$  is an orthonormal basis. Then,

$$|\mathrm{Tr}(AB)| = \left| \sum_i \langle u_i, ABu_i \rangle \right| = \left| \sum_i \lambda_i \langle u_i, Au_i \rangle \right| \leq \|A\| \sum_i |\lambda_i| = \|A\| \cdot \|B\|_1.$$

When  $B$  is PSD, we have  $\|B\|_1 = \mathrm{Tr}(B)$ , hence (1.12) gives us

$$\mathbb{E} \left[ \mathrm{Tr} \left( e^{\beta(A_1 + \dots + A_{n-1})} \mathbb{E}[e^{\beta A_n}] \right) \right] \leq \|\mathbb{E}[e^{\beta A_n}]\| \mathbb{E} \left[ \mathrm{Tr} \left( e^{\beta(A_1 + \dots + A_{n-1})} \right) \right].$$

Continuing inductively, we end up with

$$\mathbb{E} \left[ \mathrm{Tr} \left( e^{\beta(A_1 + \dots + A_n)} \right) \right] \leq \mathrm{Tr}(I) \prod_{i=1}^n \|\mathbb{E}[e^{\beta A_i}]\| = d \prod_{i=1}^n \|\mathbb{E}[e^{\beta A_i}]\| \quad (1.13)$$

Note that since  $e^A = \sum_{n \geq 0} \frac{A^n}{n!}$ , it holds that  $e^A$  commutes with powers of  $A$ . Since  $I, A, A^2, e^A$  all commute, they are simultaneously diagonalizable, meaning that the inequality (1.4) for real numbers extends:

$$I + A \leq e^A \leq I + A + A^2, \quad \forall \|A\| \leq 1.$$

Indeed, for diagonal matrices, this is just a pointwise inequality on the matrix entries: If  $D, D'$  are diagonal, then  $D \leq D' \iff D_{ii} \leq D'_{ii} \forall i$ .

Therefore we have, for each  $i$ ,

$$\|\mathbb{E}[e^{\beta A_i}]\| \leq \|\mathbb{E}[1 + \beta A_i + \beta^2 A_i^2]\| = \|\mathbb{E}[1 + \beta^2 A_i^2]\| \leq 1 + \beta^2 L^2 \leq e^{\beta^2 L^2}$$

where we have used the assumption  $\|A_i\| \leq L$ . So combining (1.10) and (1.13) gives us

$$\mathbb{P}[\lambda_{\max}(\mathbf{A}) \geq t\sqrt{n}] \leq d e^{-\beta t\sqrt{n}} e^{\beta^2 L^2 n},$$

which is analogous to the upper bound (1.5) we obtained in the scalar setting, except for the leading factor of  $d$ .