

Lecture 2: Matrix Chernoff bounds

Nick Harvey

University of British Columbia

Abstract

The purpose of my second and third lectures is to discuss spectral sparsifiers, which are the second key ingredient in most of the fast Laplacian solvers. In this lecture we will discuss concentration bounds for sums of random matrices, which are an important technical tool underlying the simplest sparsifier construction.

1 Introduction

I thought it was a rather trivial lemma, but many things are only trivial once you know them.
Herman Chernoff

The use of randomization in algorithms has been increasingly prevalent since the mid 1970s. The Chernoff bound has been a hugely important tool in randomized algorithms and learning theory since the mid 1980s. It is a concentration inequality for random variables that are the sum of many independent, bounded random variables. A formal statement is:

Theorem 1. Let X_1, \dots, X_k be independent, random, real variables with $0 \leq X_i \leq R$. Let $\mu_{\min} \leq \sum_i \mathbb{E}[X_i] \leq \mu_{\max}$. Then, for all $\delta \geq 0$,

$$\Pr \left[\sum_{i=1}^k X_i \geq (1 + \delta)\mu_{\max} \right] \leq \left(\frac{e^\delta}{(1+\delta)^{1+\delta}} \right)^{\mu_{\max}/R} \leq \begin{cases} e^{-\delta^2 \mu_{\max}/3R} & (\text{if } \delta \leq 1) \\ e^{-\delta \mu_{\max}/3R} & (\text{if } \delta > 1) \end{cases}$$

$$\Pr \left[\sum_{i=1}^k X_i \leq (1 - \delta)\mu_{\min} \right] \leq \left(\frac{e^{-\delta}}{(1-\delta)^{1-\delta}} \right)^{\mu_{\min}/R} \leq e^{-\delta^2 \mu_{\min}/2R} \quad (\text{if } \delta \leq 1).$$

In modern research it has been increasingly useful to study concentration for sums of random *matrices*. Many uses have appeared in compressed sensing, machine learning, and randomized numerical linear algebra. There was a long literature proving matrix concentration bounds, culminating in the following result of Tropp, which shows that the Chernoff bound has a perfect generalization to matrices. The formal statement is syntactically almost identical to the previous theorem.

Theorem 2 (Tropp). Let X_1, \dots, X_k be independent, random, symmetric, real matrices of size $d \times d$ with $0 \preceq X_i \preceq R \cdot I$. Let $\mu_{\min} \cdot I \preceq \sum_i \mathbb{E}[X_i] \preceq \mu_{\max} \cdot I$. Then, for all $\delta \in [0, 1]$,

$$\Pr \left[\lambda_{\max}(\sum_{i=1}^k X_i) \geq (1 + \delta)\mu_{\max} \right] \leq d \cdot \left(\frac{e^\delta}{(1+\delta)^{1+\delta}} \right)^{\mu_{\max}/R} \leq d \cdot \begin{cases} e^{-\delta^2 \mu_{\max}/3R} & (\text{if } \delta \leq 1) \\ e^{-\delta \mu_{\max}/3R} & (\text{if } \delta > 1) \end{cases}$$

$$\Pr \left[\lambda_{\min}(\sum_{i=1}^k X_i) \leq (1 - \delta)\mu_{\min} \right] \leq d \cdot \left(\frac{e^{-\delta}}{(1-\delta)^{1-\delta}} \right)^{\mu_{\min}/R} \leq d \cdot e^{-\delta^2 \mu_{\min}/2R} \quad (\text{if } \delta \leq 1).$$

Portions of these notes are based on scribe notes written by Zachary Drudi.

Here λ_{\max} and λ_{\min} respectively refer to the maximum and minimum eigenvalues of their argument. We are also using the partial ordering \preceq on symmetric matrices defined by $A \preceq B$ if and only if $B - A$ is positive semi-definite (i.e. $\lambda_{\min}(B - A) \geq 0$). This is called the Löwner ordering.

We will use Theorem 2 in the next lecture to construct “graph sparsifiers” — subgraphs which approximate the graph in a very strong sense. Sparsifying the graph is a very useful subroutine for fast algorithms, such as fast Laplacian solvers.

Although the Chernoff bound is quite powerful (and certainly useful), its proof is quite straightforward, as Chernoff himself remarks in the quotation above. The matrix Chernoff bound is also quite comprehensible, especially since its proof has a very similar structure to the scalar bound. The main difficulty is understanding what scalar inequalities have generalizations to matrices. The field that studies such inequalities is called *matrix analysis*.

2 Matrix Analysis

In this section we briefly discuss some results that we will need.

Definition 3 (Spectral mapping). Let $f : \mathbb{R} \rightarrow \mathbb{R}$ be a function. We extend f to a new function $f(A)$ on symmetric matrices by applying f to the eigenvalues of A . That is, let $A = UDU^T$ be the spectral decomposition of A , where U is orthogonal and D is diagonal. Define $f(A) = Uf(D)U^T$, where $f(D)$ is the diagonal matrix with $f(D)_{i,i} = f(D_{i,i})$.

We will primarily be interested in the case $f = \exp$ or $f = \log$.

Definition 4. A function $f : \mathbb{R} \rightarrow \mathbb{R}$ is

- **Operator monotone** if $f(A) \succeq f(B)$ whenever $A \succeq B$.
- **Operator concave** if $f((1-x)A + xB) \succeq (1-x)f(A) + xf(B)$ for all $x \in [0, 1]$ and all A, B .

In many ways, the Löwner ordering interacts nicely with the algebra of matrices and with spectral mapping. Many familiar scalar inequalities generalize to the Löwner ordering. However, many inequalities that one might hope to be true are actually false, so much care is needed!

2.1 The pitfalls of matrix analysis

Hope 1. If $f : \mathbb{R} \rightarrow \mathbb{R}$ is monotone (on a certain interval), is f also operator monotone (for matrices whose eigenvalues lie in that interval)?

This is false. Consider function $f(x) = x^2$, which is monotone on the positive reals. Define $A = \begin{pmatrix} 1 & 1 \\ 1 & 1 \end{pmatrix}$ and $B = \begin{pmatrix} 2 & 1 \\ 1 & 1 \end{pmatrix}$. Then $A \preceq B$ since $B - A = \begin{pmatrix} 1 & 0 \\ 0 & 0 \end{pmatrix}$, but $f(A) \not\preceq f(B)$ since $f(B) - f(A) = \begin{pmatrix} 3 & 1 \\ 1 & 0 \end{pmatrix}$, which has a negative eigenvalue. The same A and B also show that \exp is not operator monotone.

Hope 2. If $f : \mathbb{R} \rightarrow \mathbb{R}$ is concave (on a certain interval), is it also operator concave (for matrices whose eigenvalues lie in that interval)?

Also false. Consider $f(x) = -x^3$, which is concave on the positive reals. The same matrices A and B provide a counterexample with $x = 0.5$.

Before the reader becomes too dismayed, the take-away lesson should not be that matrix analysis

is useless. We will see next that there are many useful and powerful inequalities in matrix analysis. The lesson is that one must carefully determine which matrix inequalities are true before using them!

2.2 Some inequalities of matrix analysis

Proofs for the results of this section may be found in [my 2013 lecture notes](#), with the exception of Theorem 9.

Claim 5. Let $f : \mathbb{R} \rightarrow \mathbb{R}$ and $g : \mathbb{R} \rightarrow \mathbb{R}$ satisfy $f(x) \leq g(x) \forall x \in [l, u]$. Suppose A is symmetric and the eigenvalues of A all lie in $[l, u]$. Then $f(A) \preceq g(A)$.

Claim 6. If X and Y are random matrices and $X \preceq Y$, then $\mathbb{E}[X] \preceq \mathbb{E}[Y]$.

Claim 7 (Weyl's Monotonicity Theorem). Let A and B be symmetric, $n \times n$ matrices. Let $\lambda_i(A)$ denote the i^{th} largest eigenvalue of A . If $A \preceq B$, then $\lambda_i(A) \leq \lambda_i(B)$ for all i .

References: Bhatia page 63.

Corollary 8. If f is monotone, then so is the composition $\text{tr} f$ defined on matrices.

Theorem 9 (Löwner). \log is operator concave.

References: Carlen Theorem 2.6, Horn and Johnson Exercise 6.6.18.

2.3 A commutative multiplication operation

One annoyance with matrix multiplication is that it is not commutative. We will define a new commutative multiplication operation that turns out to be useful.

Definition 10. If A, B are positive definite, define $A \odot B = \exp(\log(A) + \log(B))$.

References: Warmuth-Kuzmin Section 4.

Note that if A and B commute then $A \odot B$ is the usual product AB .

Theorem 11 (Lieb). Fix any symmetric H . The map $A \mapsto \text{tr} \exp(\log(A) + H)$ is concave on the set of positive definite matrices.

References: Lieb 1973, Epstein 1973, Ohya-Petz 1993 Theorem 3.7.

Corollary 12. For any fixed B , the map $A \mapsto \text{tr}(A \odot B)$ is concave.

PROOF: $\text{tr}(A \odot B) = \text{tr} \exp(\log A + \log B)$. Apply Lieb's theorem with $H = \log B$. \square

Corollary 13. Let B be fixed, and A a random matrix. Then $\mathbb{E}[\text{tr}(A \odot B)] \leq \text{tr}(\mathbb{E}[A] \odot B)$.

PROOF: Apply Jensen's inequality. \square

Corollary 14. Let A_1, \dots, A_k be independent random positive definite matrices. Then

$$\mathbb{E}[\text{tr}(A_1 \odot \dots \odot A_k)] \leq \text{tr}(\mathbb{E}[A_1] \odot \dots \odot \mathbb{E}[A_k]).$$

PROOF: By symmetry Corollary 12 and Corollary 13 are also true if the roles of A and B are swapped. Thus we may inductively apply Corollary 13 to obtain the desired result. \square

3 The Chernoff Bound

We now prove only the first inequality of Theorem 1.

Claim 15.

$$\Pr \left[\sum_{i=1}^k X_i \geq t \right] \leq \inf_{\theta > 0} e^{-\theta t} \cdot \prod_{i=1}^k \mathbb{E} \left[e^{\theta X_i} \right].$$

PROOF: Fix $\theta > 0$.

$$\begin{aligned} \Pr \left[\sum_i X_i \geq t \right] &= \Pr \left[\sum_i \theta X_i \geq \theta t \right] \\ &= \Pr \left[\exp(\sum_i \theta X_i) \geq \exp(\theta t) \right] \quad (\text{monotonicity of } e^x) \\ &\leq e^{-\theta t} \cdot \mathbb{E} \left[\exp(\sum_i \theta X_i) \right] \quad (\text{Markov's inequality}) \end{aligned}$$

This expectation can be simplified:

$$\mathbb{E} \left[\exp(\sum_i \theta X_i) \right] = \mathbb{E} \left[\prod_i e^{\theta X_i} \right] = \prod_i \mathbb{E} \left[e^{\theta X_i} \right] \quad (\text{by independence}).$$

Combining these proves the claim. \square

Claim 16. Let X be a random variable with $0 \leq X \leq 1$. Then

$$\mathbb{E} \left[e^{\theta X} \right] \leq 1 + (e^\theta - 1) \cdot \mathbb{E} [X].$$

PROOF: For $x \in [0, 1]$ we have $e^{\theta x} \leq 1 + (e^\theta - 1) \cdot x$, by convexity of the left-hand side. Since $X \in [0, 1]$,

$$\begin{aligned} e^{\theta X} &\leq 1 + (e^\theta - 1) \cdot X \\ \implies \mathbb{E} \left[e^{\theta X} \right] &\leq 1 + (e^\theta - 1) \cdot \mathbb{E} [X], \end{aligned}$$

since inequalities are preserved under taking expectation. \square

Proof (of Chernoff Upper Bound). Without loss of generality $R = 1$.

$$\begin{aligned} \prod_{i=1}^k \mathbb{E} \left[e^{\theta X_i} \right] &\leq \prod_{i=1}^k (1 + (e^\theta - 1) \cdot \mathbb{E} [X_i]) \quad (\text{by Claim 16}) \\ &= \exp \left(\sum_{i=1}^k \log (1 + (e^\theta - 1) \cdot \mathbb{E} [X_i]) \right) \\ &\leq \exp \left(\sum_{i=1}^k (e^\theta - 1) \cdot \mathbb{E} [X_i] \right) \quad (\text{using } \log(1+x) \leq x) \\ &\leq \exp \left((e^\theta - 1) \mu_{\max} \right) \end{aligned}$$

Applying Claim 15 with $t = (1 + \delta) \mu_{\max}$ and $\theta = \ln(1 + \delta)$

$$\begin{aligned} \Pr \left[\sum_i X_i \geq (1 + \delta) \mu_{\max} \right] &\leq \exp \left(- \ln(1 + \delta) \cdot (1 + \delta) \mu_{\max} \right) \cdot \exp(\delta \cdot \mu_{\max}) \\ &= \left(\frac{e^\delta}{(1 + \delta)^{1 + \delta}} \right)^{\mu_{\max}} \quad \blacksquare \end{aligned}$$

4 Tropp's Matrix Chernoff Bound

We now prove only the first inequality of Theorem 2.

Claim 17.

$$\Pr \left[\lambda_{\max} \left(\sum_{i=1}^k X_i \right) \geq t \right] \leq \inf_{\theta > 0} e^{-\theta t} \cdot \text{tr} \left(\bigodot_{i=1}^k \mathbb{E} \left[e^{\theta X_i} \right] \right).$$

PROOF: Fix $\theta > 0$.

$$\begin{aligned} \Pr \left[\lambda_{\max}(\sum_i X_i) \geq t \right] &= \Pr \left[\lambda_{\max}(\sum_i \theta X_i) \geq \theta t \right] && \text{(homogeneity of max eigenvalue)} \\ &= \Pr \left[\exp(\lambda_{\max}(\sum_i \theta X_i)) \geq \exp(\theta t) \right] && \text{(monotonicity of } e^x \text{)} \\ &\leq e^{-\theta t} \cdot \mathbb{E} \left[\exp(\lambda_{\max}(\sum_i \theta X_i)) \right] && \text{(Markov's inequality)} \end{aligned}$$

We can bound the maximum eigenvalue by a trace:

$$\begin{aligned} \exp(\lambda_{\max}(\sum_i \theta X_i)) &= \lambda_{\max}(\exp(\sum_i \theta X_i)) && \text{(definition of matrix exponentiation)} \\ &\leq \text{tr}(\exp(\sum_i \theta X_i)) && \text{(max eigenvalue } \leq \text{ sum of eigenvalues)} \end{aligned}$$

Taking the expectation gives the bound:

$$\Pr \left[\lambda_{\max}(\sum_i X_i) \geq t \right] \leq e^{-\theta t} \cdot \mathbb{E} \left[\text{tr}(\exp(\sum_i \theta X_i)) \right].$$

This expectation can be bounded:

$$\begin{aligned} \mathbb{E} \left[\text{tr}(\exp(\sum_i \theta X_i)) \right] &= \mathbb{E} \left[\text{tr}(\exp(\sum_i \log A_i)) \right] && \text{(let } A_i = e^{\theta X_i} \text{)} \\ &= \mathbb{E} \left[\text{tr}(A_1 \odot \cdots \odot A_k) \right] && \text{(definition of } \odot \text{)} \\ &\leq \text{tr}(\mathbb{E}[A_1] \odot \cdots \odot \mathbb{E}[A_k]) && \text{(by Corollary 14)} \end{aligned}$$

Combining these inequalities proves the claim. \square

Claim 18. Let X be a random symmetric $d \times d$ matrix with $0 \preceq X \preceq I$. Then

$$\mathbb{E} \left[e^{\theta X} \right] \preceq I + (e^\theta - 1) \cdot \mathbb{E}[X].$$

PROOF: For $x \in [0, 1]$ we have $e^{\theta x} \leq 1 + (e^\theta - 1) \cdot x$, by convexity of the left-hand side. Since X has all eigenvalues in $[0, 1]$, Claim 5 gives

$$\begin{aligned} e^{\theta X} &\preceq I + (e^\theta - 1) \cdot X \\ \implies \mathbb{E} \left[e^{\theta X} \right] &\preceq I + (e^\theta - 1) \cdot \mathbb{E}[X], \end{aligned}$$

since the Löwner ordering is preserved under taking expectation by Claim 6. \square

Proof (of Matrix Chernoff Upper Bound). Without loss of generality $R = 1$. Our first observation is a bound for a sum of logs:

$$\begin{aligned} \sum_{i=1}^k \log \mathbb{E} \left[e^{\theta X_i} \right] &= k \cdot \sum_{i=1}^k \frac{1}{k} \log \mathbb{E} \left[e^{\theta X_i} \right] \\ &\leq k \cdot \log \left(\sum_{i=1}^k \frac{1}{k} \mathbb{E} \left[e^{\theta X_i} \right] \right) && \text{(by Theorem 9)} \end{aligned} \tag{1}$$

Next:

$$\begin{aligned}
& \operatorname{tr} \left(\mathbb{E} \left[e^{\theta X_1} \right] \odot \dots \odot \mathbb{E} \left[e^{\theta X_k} \right] \right) \\
&= \operatorname{tr} \exp \left(\sum_{i=1}^k \log \mathbb{E} \left[e^{\theta X_i} \right] \right) && \text{(definition of } \odot \text{)} \\
&\stackrel{(a)}{\leq} \operatorname{tr} \exp \left(k \cdot \log \left(\sum_{i=1}^k \frac{1}{k} \mathbb{E} \left[e^{\theta X_i} \right] \right) \right) && \text{(by (1) and Corollary 8)} \\
&\leq d \cdot \lambda_{\max} \left(\exp \left(k \cdot \log \left(\sum_{i=1}^k \frac{1}{k} \mathbb{E} \left[e^{\theta X_i} \right] \right) \right) \right) && \text{(sum of eigenvalues } \leq d \text{ times maximum)} \\
&= d \cdot \exp \left(k \cdot \log \lambda_{\max} \left(\sum_{i=1}^k \frac{1}{k} \mathbb{E} \left[e^{\theta X_i} \right] \right) \right) && \text{(definition of spectral mapping)} \\
&\leq d \cdot \exp \left(k \cdot \log \lambda_{\max} \left(I + \sum_{i=1}^k \frac{1}{k} (e^\theta - 1) \mathbb{E} \left[X_i \right] \right) \right) && \text{(by Claim 18 and Claim 7)} \\
&= d \cdot \exp \left(k \cdot \log \left(1 + \frac{e^\theta - 1}{k} \lambda_{\max} \left(\sum_{i=1}^k \mathbb{E} \left[X_i \right] \right) \right) \right) \\
&\leq d \cdot \exp \left((e^\theta - 1) \cdot \lambda_{\max} \left(\sum_{i=1}^k \mathbb{E} \left[X_i \right] \right) \right) && \text{(using } \log(1+x) \leq x \text{)} \\
&\leq d \cdot \exp \left((e^\theta - 1) \cdot \mu_{\max} \right)
\end{aligned}$$

Apply Claim 17 with $t = (1 + \delta)\mu_{\max}$ and $\theta = \ln(1 + \delta)$:

$$\begin{aligned}
\Pr \left[\lambda_{\max} \left(\sum_i X_i \right) \geq (1 + \delta)\mu_{\max} \right] &\leq \exp \left(-\ln(1 + \delta) \cdot (1 + \delta)\mu_{\max} \right) \cdot \left(d \cdot \exp(\delta \cdot \mu_{\max}) \right) \\
&= d \cdot \left(\frac{e^\delta}{(1 + \delta)^{1+\delta}} \right)^{\mu_{\max}}
\end{aligned}$$

■

Remark. In inequality (a), it is *not* true that \exp is operator monotone, but it is true (by Corollary 8) that $\operatorname{tr} \exp$ is monotone.

References

- [1] R. Bhatia. *Matrix Analysis*. Springer, 1997.
- [2] E. Carlen. Trace inequalities and quantum entropy: An introductory course. *Contemporary Mathematics*, 529, 2009.
- [3] H. Epstein. Remarks on Two Theorems of E. Lieb. *Communications in Mathematical Physics*, 31:317–325, 1973.
- [4] R. Horn and C. Johnson. *Topics in Matrix Analysis*. Cambridge University Press, 1991.
- [5] E. H. Lieb. Convex trace functions and the Wigner-Yanase-Dyson conjecture. *Advances in Mathematics*, 11:267–288, 1973.
- [6] M. Ohya and D. Petz. *Quantum Entropy and Its Use*. Springer-Verlag, 1993.
- [7] J. A. Tropp. User-friendly tail bounds for sums of random matrices. *Foundations of Computational Mathematics*, 2011.
- [8] M. K. Warmuth and D. Kuzmin. Bayesian generalized probability calculus for density matrices. *Machine Learning*, 78(1), 2010.