Lecture 12: Concentration of sums of independent random matrices



School of Mathematical Sciences, Xiamen University

1. Matrix calculus

• Let us focus on $n \times n$ symmetric matrices.

 $\mathbf{A} \succeq \mathbf{B}$ (of course, $\mathbf{B} \preceq \mathbf{A}$) means $\mathbf{A} - \mathbf{B}$ is positive semidefinite. This defines a *partial order* on the set of $n \times n$ symmetric matrices.

• $\|\mathbf{A}\|_2$ is the smallest non-negative number M such that

$$\|\mathbf{A}\mathbf{x}\|_2 \leqslant M\|\mathbf{x}\|_2$$
 for all $\mathbf{x} \in \mathbb{R}^n$.

Functions of matrices

Let $f : \mathbb{R} \to \mathbb{R}$ be a function and **X** be an $n \times n$ symmetric matrix. We define

$$f(\mathbf{X}) := \sum_{i=1}^{n} f(\lambda_i) \mathbf{u}_i \mathbf{u}_i^{\top}$$

where λ_i are the eigenvalues of **X** and \mathbf{u}_i are the corresponding eigenvectors satisfying

$$\mathbf{X} = \sum_{i=1}^{n} \lambda_i \mathbf{u}_i \mathbf{u}_i^{\top}.$$

2. Matrix Bernstein's inequality

Theorem 1 (Matrix Bernstein's inequality)

Let X_1, \ldots, X_N be independent, mean zero, $n \times n$ symmetric random matrices, such that

$$\|\boldsymbol{X}_i\|_2 \leq K$$

almost surely for all i. Then, for every $t \geq 0$ we have

$$\mathbb{P}\left\{\left\|\sum_{i=1}^{N} \boldsymbol{X}_{i}\right\|_{2} \geqslant t\right\} \leqslant 2n \cdot \exp\left(-\frac{t^{2}/2}{\sigma^{2} + Kt/3}\right).$$

Here

$$\sigma^2 = \left\| \sum_{i=1}^N \mathbb{E} oldsymbol{X}_i^2
ight\|_2$$

is the norm of the sum of the "matrix variance".

- The scalar case, where n = 1, is the classical Bernstein's inequality for random variables.
- A remarkable feature of matrix Bernstein's inequality, which makes it especially powerful, is that it does not require any independence of the entries (or the rows or columns) of X_i ; all is needed is that the random matrices X_i be independent from each other.
- The proof of Theorem 1 is based on bounding the moment generating function (MGF) $\mathbb{E} \exp(\lambda S)$ of the sum $S = \sum_{i=1}^{N} X_i$.
- Golden–Thompson inequality: For any symmetric X and Y,

$$\operatorname{tr}\left(e^{\mathbf{X}+\mathbf{Y}}\right) \leqslant \operatorname{tr}\left(e^{\mathbf{X}}e^{\mathbf{Y}}\right).$$

• Lieb's lemma: For given $n \times n$ symmetric **H**, the function

$$f(\mathbf{X}) = -\operatorname{tr}\exp(\mathbf{H} + \log \mathbf{X})$$

is convex on the space of $n \times n$ SPD matrices.

• Jensen's inequality:

For any convex function f and a random matrix X, it holds

$$f(\mathbb{E}X) \leq \mathbb{E}f(X).$$

• Lieb's inequality for random matrices

Let **H** be a fixed $n \times n$ symmetric matrix and **Z** be an $n \times n$ symmetric random matrix. Then

$$\mathbb{E}\mathrm{tr}\exp(\mathbf{H}+\boldsymbol{Z})\leqslant\mathrm{tr}\exp\left(\mathbf{H}+\log\mathbb{E}\mathrm{e}^{\boldsymbol{Z}}\right).$$

• MGF of a sum of independent random matrices

Let $X_1, ..., X_N$ be independent $n \times n$ symmetric random matrices. Then the sum $S = \sum_{i=1}^{N} X_i$ satisfies

$$\mathbb{E}\mathrm{tr}\exp(\lambda \boldsymbol{S}) \leq \mathrm{tr}\exp\left[\sum_{i=1}^{N}\log\mathbb{E}\mathrm{e}^{\lambda \boldsymbol{X}_{i}}\right].$$

Lemma 2

Let X be an $n \times n$ symmetric random matrix. Assume that

$$\mathbb{E}X = \mathbf{0}$$

and

$$\|\boldsymbol{X}\|_2 \leq K$$

almost surely. Then, for all

$$0<\lambda<\frac{3}{K}$$

we have

$$\mathbb{E}\exp(\lambda \boldsymbol{X}) \leq \exp(g(\lambda)\mathbb{E}\boldsymbol{X}^2)$$

where

$$g(\lambda) = \frac{\lambda^2/2}{1 - \lambda K/3}.$$

Lecture 12

Proof. First, check for $0 < \lambda < 3/K$ and $|x| \le K$, it holds

$$e^{\lambda x} \le 1 + \lambda x + g(\lambda)x^2$$
.

Then extend it to matrix version: for

$$0 < \lambda < 3/K$$
 and $\|\boldsymbol{X}\|_2 \le K$,

it holds

$$e^{\lambda \mathbf{X}} \leq \mathbf{I} + \lambda \mathbf{X} + g(\lambda) \mathbf{X}^2.$$

Finally, take the expectation and recall that

$$\mathbb{E} X = \mathbf{0}$$

to obtain

$$\mathbb{E}\exp(\lambda \boldsymbol{X}) \leq \mathbf{I} + g(\lambda)\mathbb{E}\boldsymbol{X}^2 \leq \exp\left(g(\lambda)\mathbb{E}\boldsymbol{X}^2\right),$$

where the last inequality is from the matrix version of the scalar inequality

$$1 + z \le e^z$$
 for all $z \in \mathbb{R}$. \square

Proof of Matrix Bernstein's inequality (Theorem 1).

$$\mathbb{P}\{\lambda_{\max}(S) \geq t\} = \mathbb{P}\{\exp(\lambda \cdot \lambda_{\max}(S)) \geq \exp(\lambda t)\} \\
\leq e^{-\lambda t} \mathbb{E} e^{\lambda \cdot \lambda_{\max}(S)} \\
= e^{-\lambda t} \mathbb{E} \lambda_{\max}(e^{\lambda S}) \quad \text{(check!)} \\
\leq e^{-\lambda t} \mathbb{E} \text{tr}(e^{\lambda S}) \\
\leq e^{-\lambda t} \text{tr} \exp\left[\sum_{i=1}^{N} \log \mathbb{E} e^{\lambda X_i}\right] \\
\leq e^{-\lambda t} \text{tr} \exp\left[g(\lambda) \mathbf{Z}\right]$$

where $\mathbf{Z} := \sum_{i=1}^{N} \mathbb{E} \mathbf{X}_{i}^{2}$. Consider a special $\lambda = t/(\sigma^{2} + Kt/3)$. With this value of λ , we conclude

$$\mathbb{P}\left\{\lambda_{\max}(\boldsymbol{S}) \geqslant t\right\} \leqslant n \cdot \exp\left(-\frac{t^2/2}{\sigma^2 + Kt/3}\right).$$

Next, repeat the argument for -S to reinstate the absolute value.

DAMC Lecture 12 Spring 2022 8 / 18

Corollary 3

Let $X_1, ..., X_N$ be independent, mean zero, $n \times n$ symmetric random matrices, such that $||X_i||_2 \leq K$ almost surely for all i. Then

$$\mathbb{E} \left\| \sum_{i=1}^{N} \boldsymbol{X}_{i} \right\|_{2} \lesssim \sigma \sqrt{\log n} + K \log n$$

where
$$\sigma = \left\| \sum_{i=1}^{N} \mathbb{E} \boldsymbol{X}_{i}^{2} \right\|_{2}^{1/2}$$

Proof. The link from tail bounds to expectation is provided by the basic identity

$$\mathbb{E}Z = \int_0^\infty \mathbb{P}\left\{Z > t\right\} dt$$

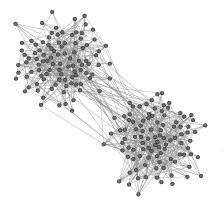
which is valid for any non-negative random variable Z. Integrating the tail bound given by matrix Bernstein's inequality, you will arrive at the expectation bound we claimed.

3. Community recovery in networks

- A network can be mathematically represented by a graph, a set of *n* vertices with edges connecting some of them.
- For simplicity, we will consider undirected graphs where the edges do not have arrows.
- Real world networks often tend to have clusters, or communities subsets of vertices that are connected by unusually many edges.
 (For example, a friendship network where communities form around some common interests.)
- An important problem in data science is to recover communities from a given network.
- Spectral clustering: one of the simplest methods for community recovery.

3.1 Stochastic block model

• Partition a set of n vertices into two subsets ("communities") with n/2 vertices each, and connect every pair of vertices independently with probability p if they belong to the same community and q < p if not. The resulting random graph is said to follow the stochastic block model G(n, p, q).



A network according to stochastic block model

$$n = 200$$

$$p = 1/20$$

$$q = 1/200$$

3.2 Spectral clustering

- Suppose we are shown one instance of a random graph generated according to a stochastic block model G(n, p, q). How can we find which vertices belong to which community?
- The spectral clustering algorithm we are going to explain will do precisely this. It will be based on the spectrum of the adjacency matrix A of the graph, which is the $n \times n$ symmetric matrix whose entries A_{ij} equal 1 if the vertices i and j are connected by an edge, and 0 otherwise.
- The adjacency matrix A is a random matrix. Let us compute its expectation first. This is easy, since the entires of A are Bernoulli random variables. If i and j belong to the same community then $\mathbb{E}A_{ij} = p$ and otherwise $\mathbb{E}A_{ij} = q$.

• Thus $\mathbb{E}\boldsymbol{A}$ has block structure: for example, if n=4 then $\mathbb{E}\boldsymbol{A}$ looks like this:

$$\mathbb{E}\boldsymbol{A} = \begin{bmatrix} p & p & q & q \\ p & p & q & q \\ \hline q & q & p & p \\ q & q & p & p \end{bmatrix}$$

For illustration purposes, we grouped the vertices from each community together.

ullet EA has rank 2, and the non-zero eigenvalues and the corresponding eigenvectors are

$$\lambda_1 = \left(\frac{p+q}{2}\right)n, \ \mathbf{v}_1 = \begin{bmatrix} \frac{1}{1} \\ \frac{1}{1} \end{bmatrix}; \quad \lambda_2 = \left(\frac{p-q}{2}\right)n, \ \mathbf{v}_2 = \begin{bmatrix} \frac{1}{1} \\ \frac{-1}{-1} \end{bmatrix}.$$

DAMC Lecture 12 Spring 2022 13 / 18

ullet The eigenvalues and eigenvectors of $\mathbb{E} A$ tell us a lot about the community structure of the underlying graph.

The first (larger) eigenvalue,

$$d := \left(\frac{p+q}{2}\right)n,$$

is the *expected degree* of any vertex of the graph.

The second eigenvalue

$$\left(\frac{p-q}{2}\right)n,$$

tells us whether there is any community structure at all (which happens when $p \neq q$ and thus $\lambda_2 \neq 0$).

The signs of the components of the second eigenvector \mathbf{v}_2 tell us exactly how to separate the vertices into the two communities.

• The problem is that we do not know $\mathbb{E}A$. Instead, we know the adjacency matrix A. So is it true that $A \approx \mathbb{E}A$?

Theorem 4 (Concentration of the stochastic block model)

Let **A** be the adjacency matrix of a G(n, p, q) random graph. Then

$$\mathbb{E}\|\boldsymbol{A} - \mathbb{E}\boldsymbol{A}\|_2 \lesssim \sqrt{d\log n} + \log n.$$

Here d = (p+q)n/2 is the expected degree.

Proof. To use matrix Bernstein's inequality, let us break \boldsymbol{A} into a sum of independent random matrices

$$oldsymbol{A} = \sum_{i,j:i\leqslant j} oldsymbol{X}_{ij}$$

where each matrix X_{ij} contains a pair of symmetric entries of A, or one diagonal entry. Matrix Bernstein's inequality obviously applies for the sum

$$oldsymbol{A} - \mathbb{E}oldsymbol{A} = \sum_{i \leqslant j} \left(oldsymbol{X}_{ij} - \mathbb{E}oldsymbol{X}_{ij}
ight).$$

DAMC Lecture 12 Spring 2022 15 / 18

Corollary 3 gives

$$\mathbb{E}\|\boldsymbol{A} - \mathbb{E}\boldsymbol{A}\|_2 \lesssim \sigma \sqrt{\log n} + K \log n$$

where

$$\sigma^2 = \left\| \sum_{i \leqslant j} \mathbb{E} \left(oldsymbol{X}_{ij} - \mathbb{E} oldsymbol{X}_{ij}
ight)^2
ight\|_2$$

and

$$K = \max_{i,j} \|\boldsymbol{X}_{i,j} - \mathbb{E}\boldsymbol{X}_{i,j}\|_{2}.$$

It is a good exercise to check that

$$\sigma^2 \lesssim d$$

and

$$K \leqslant 2$$
.

This completes the proof.

DAMC Lecture 12 Spring 2022 16 / 18

• How useful is Theorem 4 for community recovery? Suppose that the network is not too sparse, namely

$$d \gg \log n$$
.

Then (with high probability)

$$\|\boldsymbol{A} - \mathbb{E}\boldsymbol{A}\|_2 \lesssim \sqrt{d \log n}$$

while

$$\|\mathbb{E}\boldsymbol{A}\|_2 = \lambda_1(\mathbb{E}\boldsymbol{A}) = d,$$

which implies that

$$\|\boldsymbol{A} - \mathbb{E}\boldsymbol{A}\|_2 \ll \|\mathbb{E}\boldsymbol{A}\|_2$$
.

In other words, A nicely approximates $\mathbb{E}A$: the relative error or approximation is small in the operator norm.

DAMC Lecture 12 Spring 2022 17 / 18

- Classical results from the perturbation theory for matrices state that if \boldsymbol{A} and $\mathbb{E}\boldsymbol{A}$ are close, then their eigenvalues and eigenvectors must also be close. The relevant perturbation results are Weyl's inequality for eigenvalues and Davis–Kahan's inequality for eigenvectors. (See the bible NLA book MC)
- So we can use $\mathbf{v}_2(\mathbf{A})$ to approximately recover the communities. This method is called spectral clustering:

Spectral Clustering Algorithm.

Compute $\mathbf{v}_2(\mathbf{A})$, the eigenvector corresponding to the second largest eigenvalue of the adjacency matrix \mathbf{A} of the network. Use the signs of the coefficients of $\mathbf{v}_2(\mathbf{A})$ to predict the community membership of the vertices.

DAMC Lecture 12 Spring 2022 18 / 18