

Deck 7: Matrix concentration: Matrix Bernstein

Math 7870: Topics in Randomized Numerical Linear Algebra

Spring 2026

Akil Narayan

Matrix concentration inequalities

A review:

- $\mathbf{X}_j, j \in [n]$ are $m \times m$, independent, positive semi-definite matrices. $\mathbf{S}_n = \sum_{j \in [n]} \mathbf{X}_j$.
- We have master bounds, e.g., $\Pr(\lambda_{\min}(\mathbf{S}_n) \leq t) \leq \inf_{s < 0} e^{-st} \exp \operatorname{tr} \left(\sum_{j \in [n]} \mathbf{K}_{\mathbf{X}_j}(s) \right)$
- We've established that $\operatorname{tr} \exp$ is a monotone function on spd matrices.
- We used this to derive a Chernoff bound for spd matrices, e.g., if $\|\mathbf{X}_j\|_2 \leq B$ with probability 1, then the master bound delivers for an absolute constant $c > 0$:

$$\Pr(\lambda_{\min}(\mathbf{S}_n) \leq (1 - \epsilon)\lambda_{\min}(\mathbb{E}\mathbf{S}_n)) \leq m \exp \left(-c \frac{\lambda_{\min}(\mathbb{E}\mathbf{S}_n)}{B} \epsilon^2 \right),$$

We will now work toward generalizing this result.

A Legendre-Fenchel perspective

We can identify a convex duality argument similar to the scalar case. The benefit here is that we can allow indefinite (Hermitian) matrices.

Theorem 1. *Suppose $\mathbf{X}_j, j \in [n]$ is a sequence of $m \times m$, independent, Hermitian random matrices. Suppose that, for some deterministic spd matrices $\{\mathbf{A}_j\}_{j \in [n]}$, we have,*

$$\mathbf{M}_{\mathbf{X}_j}(s) \leq \exp(f(s)\mathbf{A}_j), \quad f : (0, \infty) \rightarrow [0, \infty],$$

Then:

$$\Pr(\lambda_{\max}(\mathbf{S}_n) \geq t) \leq m \exp(-rf^*(t/r)), \quad r := \lambda_{\max}\left(\sum_{j \in [n]} \mathbf{A}_j\right),$$

where f^* is the convex conjugate of f , and $\mathbf{S}_n = \sum_{j \in [n]} \mathbf{X}_j$.

Compare this to the scalar concentration.

(Here, $f^*(t/r)$ is the rate function, and $r \sim n$ and $t \leftarrow nt$ for the iid case.)

The proof

We “just” have to use the master bounds:

$$\begin{aligned}\Pr(\lambda_{\max}(\mathbf{S}_n) \geq t) &\leq \inf_{s>0} e^{-st} \operatorname{tr} \exp \left(\sum_{j \in [n]} \mathbf{K}_{\mathbf{X}_j}(s) \right) \\ &\leq \inf_{s>0} e^{-st} \operatorname{tr} \exp \left(f(s) \sum_j \mathbf{A}_j \right) \\ &\leq m \inf_{s>0} e^{-st} \lambda_{\max} \exp \left(f(s) \sum_j \mathbf{A}_j \right) \\ &= m \inf_{s>0} \exp(-st + f(s)r) \\ &= m \exp \left(-\sup_{s>0} (st - rf(s)) \right).\end{aligned}$$

For functions f for which f^* is computable, this is a reasonable bound.

However, the assumption $\mathbf{M}_{\mathbf{X}_j(s)} \leq \exp(f(x)\mathbf{A}_j)$ is not very transparent.

On to Matrix Bernstein-type estimates

We'll now consider how to come up with a more practical bound under the assumption that the random summands are bounded.

As is (by now?) usual, the exercise boils down estimating the matrix CGF.

Lemma. *Assume \mathbf{X} is a random Hermitian matrix, and that $\mathbb{E}\mathbf{X} = \mathbf{0}$ and $\lambda_{\max}(\mathbf{X}) \leq B$ with probability 1.*

Then for $s \in (0, 3/L)$:

$$\mathbf{K}_{\mathbf{X}}(s) \leq g(s)\mathbb{E}\mathbf{X}^2, \quad g(s) := \frac{1}{2} \frac{s^2}{(1 - sL/3)}$$

The proof relies on some technical estimates of g and of a univariate function f . In particular for some $s \in (0, 3/L)$:

$$f(x) := \frac{e^{sx} - sx - 1}{x^2}$$
$$x \leq B \quad \text{implies} \quad f(x) \leq f(B)$$
$$f(x) \leq g(x)$$

Finishing the proof

The rest of the proof uses some familiar properties of the Loewner ordering and operator monotone functions:

$$\exp(s\mathbf{X}) = \mathbf{I} + s\mathbf{X} + (e^{s\mathbf{X}} - s\mathbf{X} - \mathbf{I}) = \mathbf{I} + s\mathbf{X} + \mathbf{X}f(\mathbf{X})\mathbf{X}.$$

Since $\lambda_{\max}(\mathbf{X}) \leq B$, then $\mathbf{X} \leq B\mathbf{I}$, so $f(\mathbf{X}) \leq f(B)\mathbf{I}$. Therefore:

$$\begin{aligned}\exp(s\mathbf{X}) &= \mathbf{I} + s\mathbf{X} + \mathbf{X}f(\mathbf{X})\mathbf{X} \leq \mathbf{I} + s\mathbf{X} + f(B)\mathbf{X}\mathbf{I}\mathbf{X} \\ &= \mathbf{I} + s\mathbf{X} + f(B)\mathbf{X}^2 \\ &\leq \mathbf{I} + s\mathbf{X} + g(B)\mathbf{X}^2\end{aligned}$$

Take expectations on both sides, use the monotone operator property of log.

Toward Hermitian matrix Bernstein

Armed with this inequality, we can prove a matrix Bernstein inequality.

The proof pieces together several results: Suppose \mathbf{X}_j are iid Hermitian, centered, and that $\lambda_{\max}(\mathbf{X}_j) \leq B$ with probability 1 for every $j \in [n]$.

Then, e.g., for the upper probability bound:

$$\begin{aligned}\Pr(\lambda_{\max}(\mathbf{S}_n) \geq t) &\leq \inf_{s>0} e^{-st} \operatorname{tr} \exp \left(\sum_{j \in [n]} \mathbf{K}_{\mathbf{X}_j}(s) \right) \\ &\leq \inf_{0 < s < 3/B} e^{-st} \operatorname{tr} \exp \left(g(s) \sum_{j \in [n]} \mathbb{E} \mathbf{X}_j^2 \right) \\ &\leq m \inf_{0 < s < 3/B} e^{-st} \exp \left(g(s) \lambda_{\max} \left(\sum_{j \in [n]} \mathbb{E} \mathbf{X}_j^2 \right) \right) \\ &= m \inf_{0 < s < 3/B} e^{-st} \exp(g(s) v(\mathbf{S}_n)) \quad (v(\mathbf{S}_n) := \|\mathbf{E} \mathbf{S}_n^2\|_2)\end{aligned}$$

We can then minimize over s .

Hermitian matrix Bernstein

Theorem 2 (Hermitian matrix Bernstein). *Suppose \mathbf{X}_j are iid Hermitian, centered, and that $\lambda_{\max}(\mathbf{X}_j) \leq B$ with probability 1 for every $j \in [n]$. Define,*

$$\mathbf{S}_n = \sum_{j \in [n]} \mathbf{X}_j, \quad v(\mathbf{S}_n) = \|\mathbb{E} \mathbf{S}_n^2\|_2.$$

Then, e.g.,

$$\Pr(\lambda_{\max}(\mathbf{S}_n) \geq t) \leq m \exp\left(-\frac{1}{2} \frac{t^2}{v(\mathbf{S}_n) + Bt/3}\right),$$
$$\mathbb{E} \lambda_{\max}(\mathbf{S}_n) \leq \sqrt{2v(\mathbf{S}_n) \log m} + \frac{B}{3} \log m$$

Bounded rectangular matrix Bernstein, I

With a matrix dilation, we can extend this to non-Hermitian matrices \mathbf{X}_j .

Let

$$\mathbf{Y}_j = \mathcal{D}(\mathbf{X}_j) = \begin{pmatrix} \mathbf{0} & \mathbf{X}_j \\ \mathbf{X}_j^* & \mathbf{0} \end{pmatrix}$$

With

$$\mathbf{S}_n = \sum_j \mathbf{X}_j, \quad \mathbf{T}_n = \sum_j \mathbf{Y}_j$$

From the properties of matrix dilation:

$$\mathbf{T}_n = \mathcal{D}(\mathbf{S}_n), \quad \lambda_{\max}(\mathbf{T}_n) = \|\mathbf{S}_n\|_2.$$

The generic strategy is then to bound the deviation of $\lambda_{\max}(\mathbf{T}_n)$ from $\lambda_{\max}(\mathbb{E}\mathbf{T}_n)$ using Hermitian Bernstein, and then translate quantities involving \mathbf{Y}_j into those depending on \mathbf{X}_j .

Bounded rectangular matrix Bernstein, II

Theorem 3. Let \mathbf{X}_j , $j \in [n]$ be a sequence of $m_1 \times m_2$ independent, centered, random matrices. Assume,

$$\|\mathbf{X}_j\|_2 \leq B, \quad \text{with probability 1 for all } j \in [n].$$

Define the variance statistic

$$v(\mathbf{S}_n) := \max \{ \|\mathbb{E} \mathbf{S}_n \mathbf{S}_n^*\|_2, \|\mathbb{E} \mathbf{S}_n^* \mathbf{S}_n\|_2 \}.$$

Then:

$$\Pr(\|\mathbf{S}_n\|_2 \geq t) \leq (m_1 + m_2) \exp\left(-\frac{1}{2} \frac{t^2}{(v(\mathbf{S}_n) + Bt/3)}\right), \quad t > 0$$

$$\mathbb{E}\|\mathbf{S}_n\|_2 \leq \sqrt{2v(\mathbf{S}_n) \log(m_1 + m_2)} + \frac{B}{3} \log(m_1 + m_2).$$

NB: $d_1 + d_2$ is the size of the dilation $\mathcal{D}(\mathbf{S}_n)$.

Since $\mathbb{E} \mathbf{S}_n = \mathbf{0}$, the above quantifies deviation of \mathbf{S}_n from its mean.

Subexponential matrices?

Can we extend these beyond matrices with bounded norms?

In the scalar case: suppose X is an exponential random variable with rate parameter $\lambda > 0$. By direct computation:

$$f_X(x) \sim e^{-\lambda x}, \quad \mathbb{E}X = 1/\lambda, \quad \mathbb{E}X^p = \frac{p!}{\lambda^p} \sim c^p p!$$

A subexponential random variable colloquially is one with tails that are sufficiently light so the can be dominated by an exponential.

Equivalently, a random variable is subexponential if its moments grow at most as fast as an exponential random variable.

The more common moment bound definition for scalars is not the factorial expression above. Rather, by Stirling's approximation, $p! \sim \sqrt{p} \frac{p^p}{e^p}$, so

$$(\mathbb{E}X^p)^{1/p} \lesssim cp,$$

for some constant c .

However, for the matrix setting, a version of the factorial bound is the more convenient definition.

Subexponential matrix Bernstein

Theorem 4 (Subexponential Matrix Bernstein). Let $\mathbf{X}_j, j \in [n]$ be independent, $m \times m$ self-adjoint matrices, and assume,

$$\mathbb{E}\mathbf{X}_j = \mathbf{0}, \quad \mathbb{E}\mathbf{X}_j^k \leq \frac{1}{2}p! R^{p-2}\mathbf{A}_j^2,$$

for all $j \in [n]$, integer $p \geq 2$. Then by defining,

$$\sigma^2 = \left\| \sum_{j \in [n]} \mathbf{A}_j^2 \right\|_2,$$

we have

$$\Pr(\lambda_{\max}(\mathbf{S}_n) \geq t) \leq m \exp\left(-\frac{1}{2} \frac{t^2}{(\sigma^2 + Rt)}\right).$$

The bound above satisfies:

$$\Pr(\lambda_{\max}(\mathbf{S}_n) \geq t) \leq m \begin{cases} \exp(-t^2/4\sigma^2), & \text{when } t \leq \sigma^2/R \\ \exp(-t/4R), & \text{when } t \geq \sigma^2/R \end{cases}$$

Subexponential matrix Bernstein CGF bound

To prove the subexponential matrix Bernstein, the crucial estimate to prove is a bound on the matrix CGF.

Lemma (Subexponential matrix Bernstein CGF). *Suppose \mathbf{X} is a random centered, self-adjoint matrix satisfying $\mathbb{E}\mathbf{X}^p \leq \frac{p!}{2}\mathbf{A}^2$ for all integers $p \geq 2$. Then,*

$$\mathbf{K}_{\mathbf{X}}(s) \leq \frac{s^2}{2(1-s)}\mathbf{A}^2, \quad s \in (0, 1).$$

Proof: for $0 < s < 1$,

$$\begin{aligned} \mathbf{M}_{\mathbf{X}}(s) &= \sum_{p=0}^{\infty} \frac{s^p \mathbb{E}\mathbf{X}^p}{p!} = \mathbf{I} + s\mathbb{E}\mathbf{X} + \sum_{p=2}^{\infty} \frac{s^p \mathbb{E}\mathbf{X}^p}{p!} \\ &\leq \mathbf{I} + \frac{1}{2}\mathbf{A}^2 \sum_{p=2}^{\infty} s^p \\ &= \mathbf{I} + \frac{s^2}{2(1-s)}\mathbf{A}^2 \leq \exp\left(\frac{s^2}{2(1-s)}\mathbf{A}^2\right). \end{aligned}$$

The rest of the subexponential matrix Bernstein proof is the same as before.

Generalizations

Here are some generalizations:

- Rectangular subexponential matrix Bernstein
- Hermitian matrix Azuma's inequality
(This requires somewhat different machinery than in the scalar case.)
- Matrix Hoeffding
(This turns out to be a consequence of the proof of matrix Azuma)
- Rectangular matrix Azuma
- Hermitian matrix McDiarmid

Application setup: Random matrix approximation

Suppose $\mathbf{A} \in \mathbb{C}^{m_1 \times m_2}$ is a deterministic matrix. We consider the template procedure of approximating \mathbf{A} by sums of iid copies of a random matrix \mathbf{X} . I.e., we consider

$$\mathbf{S}_n = \sum_{j \in [n]} \frac{1}{n} \mathbf{X}_j, \quad \mathbf{X}_j \stackrel{\text{iid}}{\sim} \mathbf{X}, \quad \mathbb{E} \mathbf{X} = \mathbf{A}.$$

Lemma (Matrix approximation via sampling). *With the above setup, assume $\|\mathbf{X}\|_2 \leq B$ with probability 1, and define,*

$$v_2(\mathbf{X}) = \max \{ \|\mathbb{E}(\mathbf{X}\mathbf{X}^*)\|_2, \|\mathbf{X}^*\mathbf{X}\|_2 \}.$$

Then:

$$\mathbb{E} \|\mathbf{S}_n - \mathbf{A}\|_2 \leq \sqrt{\frac{v_2(\mathbf{X}) \log(m_1 + m_2)}{n}} + \frac{2B \log(m_1 + m_2)}{3n},$$
$$\Pr(\|\mathbf{S}_n - \mathbf{A}\|_2 \geq t) \leq (m_1 + m_2) \exp\left(-\frac{1}{2} \frac{nt^2}{v_2(\mathbf{X}) + 2Bt/3}\right), \quad t > 0.$$

Proof steps: (i) Write $\mathbf{S}_n - \mathbf{A}$ as an iid sum of centered summands, (ii) compute uniform bounds on the centered summand ($2B/n$), (iii) compute the variance statistic v of the centered summands ($v_2(\mathbf{X})/n$). Apply (bounded) matrix Bernstein.

Application I: Randomized matrix multiplication, I

Let's tie up a loose end: randomized matrix multiplication.

Recall: Let \mathbf{A}, \mathbf{B} be two matrices of sizes $\mathbf{A} \in \mathbb{C}^{m_1 \times k}$, $\mathbf{B} \in \mathbb{C}^{k \times m_2}$.

If $\mathbf{a}_j, \mathbf{b}_j, j \in [k]$, are the rows and columns of \mathbf{A} and \mathbf{B} , respectively, then we define a rank-1 matrix \mathbf{X} as:

$$\mathbf{X} = \frac{1}{p_j} \mathbf{a}_j \mathbf{b}_j^* \quad \text{w/ probability } p_j, j \in [k].$$

Then $\mathbb{E}\mathbf{X} = \mathbf{AB}$.

Also recall: the choice,

$$p_j = \frac{\|\mathbf{a}_j\|_2 \|\mathbf{b}_j\|_2}{\sum_{q \in [k]} \|\mathbf{a}_q\|_2 \|\mathbf{b}_q\|_2},$$

minimizes a variance. For computational convenience, we'll choose the slightly different distribution:

$$p_j = \frac{\|\mathbf{a}_j\|_2^2 + \|\mathbf{b}_j\|_2^2}{\|\mathbf{A}\|_F^2 + \|\mathbf{B}\|_F^2}$$

(The latter p_j are comparable to the former.)

Application I: Randomized matrix multiplication, II

Our estimator is $\mathbf{S}_n = \sum_{j \in [k]} (\frac{1}{n} \mathbf{X}_j)$, with $\mathbf{X}_j \stackrel{\text{iid}}{\sim} \mathbf{X}$.

Some computations reveal:

$$\begin{aligned}\|\mathbf{X}\|_2 &\leq \frac{1}{2} (\|\mathbf{A}\|_F^2 + \|\mathbf{B}\|_F^2) = B. \\ v_2(\mathbf{X}) &\leq (\|\mathbf{A}\|_F^2 + \|\mathbf{B}\|_F^2) \max\{\|\mathbf{A}\|_2^2, \|\mathbf{B}\|_2^2\}.\end{aligned}$$

Then, e.g., in expectation we obtain,

$$\frac{1}{\|\mathbf{A}\|_2^2 + \|\mathbf{B}\|_2^2} \mathbb{E} \|\mathbf{S}_n - \mathbf{AB}\|_2 \lesssim \sqrt{\frac{2\tau}{n}} + \frac{\tau}{3n}, \quad \tau = (\text{srank}(\mathbf{A}) + \text{srank}(\mathbf{B})) \log(m_1 + m_2),$$

where, srank is the *stable* or numerical rank, defined as,

$$\text{srank}(\mathbf{A}) = \frac{\|\mathbf{A}\|_F^2}{\|\mathbf{A}\|_2^2} = \frac{\sum_j \sigma_j(\mathbf{A})^2}{\sigma_1^2(\mathbf{A})} \in [1, \min\{k, m_1\}].$$

(The expected statements hold in probability: sub-Gaussian tails for $\|\mathbf{S}_n - \mathbf{AB}\|_2$ close to 0, sub-exponential tails when it's far from 0.)

Application II: Random sparsification of a matrix, I

We can now consider a slightly more complicated, and useful, numerical linear algebra application: matrix sparsification.

Let $\mathbf{A} \in \mathbb{C}^{m_1 \times m_2}$ be a deterministic matrix. We seek to compute an approximation to \mathbf{A} via sparsification, i.e., if we write,

$$\mathbf{A} = \sum_{(i,j) \in [m_1] \times [m_2]} A_{i,j} \mathbf{e}_i \mathbf{e}_j^*.$$

Given a set $S \subset [m_1] \times [m_2]$, with $|S| = n$, we consider,

$$\mathbf{A}_S = \sum_S A_{i,j} \mathbf{e}_i \mathbf{e}_j^*, \quad \text{Is } \|\mathbf{A} - \mathbf{A}_S\|_2 \sim \text{“small”?}$$

If $n = |S| \ll m_1 m_2$, this is a sparsification of \mathbf{A} .

We'll construct S randomly.

Note that best possible fixed-rank approximation in the 2-norm is given by the truncated SVD, and the trailing singular values provide an estimate of error.

Therefore it's reasonable to believe that the trailing singular values play a role in estimating error. In particular, we'll use the stable rank $\text{srnk}(\mathbf{A})$.

Application II: Random sparsification of a matrix, II

We construct the subset S randomly by selecting indices. This is similar to how we constructed matrix multiplication. Let

$$\mathbf{X} = \frac{1}{p_{i,j}} A_{i,j} \mathbf{e}_i \mathbf{e}_j^*, \quad (i,j) = (k,\ell) \text{ with probability } p_{k,\ell},$$

where we choose the probabilities to be of the form,

$$2p_{i,j} = \frac{|A_{i,j}|^2}{\|\mathbf{A}\|_F^2} + \frac{|A_{i,j}|}{\|\mathbf{A}\|_1}, \quad (i,j) \in [m_1] \times [m_2].$$

The estimator is $\mathbf{S}_n = \frac{1}{n} \sum_{j \in [n]} \mathbf{X}_j$, with $\mathbf{X}_j \stackrel{\text{iid}}{\sim} \mathbf{X}$.

We follow a similar trajectory:

$$\mathbb{E} \mathbf{X} = \mathbf{A}, \quad B = \max_{i,j} \|\mathbf{X}\|_2 \leq 2\|\mathbf{A}\|_1, \quad v_2(\mathbf{X}) \leq 2 \max\{m_1, m_2\}.$$

Then, e.g.,

$$\frac{1}{\|\mathbf{A}\|_2} \mathbb{E} \|\mathbf{A} - \mathbf{S}_n\|_2 \lesssim \frac{\|\mathbf{A}\|_F}{\|\mathbf{A}\|_2} \sqrt{\frac{\max\{m_1, m_2\} \log(m_1 + m_2)}{n}}$$

For example, this implies that $n \gtrsim \text{srnk}(\mathbf{A})/\epsilon^2$ guarantees a relative expected error of 4ϵ .

Application III: Connectedness of a random graph, I

An Erdős-Rényi-Gilbert (ER) graph is one where, given a fixed number of vertices, every edge is present according to an iid Bernoulli draw.

Given $n \in \mathbb{N}$ vertices, and a density parameter $p \in [0, 1]$, $G(n, p)$ is the ER model where a graph is constructed over the n vertices by sampling an independent Bernoulli(p) random variable on each edge to determine if it's present or not.

Erdős-Rényi-Gilbert Random Graphs $G(n, p)$

$n = 8, p = 0.2$



sparse (6 edges)

$n = 8, p = 0.5$



medium (14 edges)

$n = 8, p = 0.8$



dense (22 edges)

Application III: Connectedness of a random graph, II

There are $\sim n^2/2$ possible edges over n vertices.

It's well-known that if $p \sim (\log n)/n$, then the probability that the graph is connected is high. (A graph is connected if edges are present to allow you to travel to one vertex from any other one.)

We can show this result using our matrix concentration results.

We'll need some background from (spectral) graph theory.

Given a graph G with vertices $[n]$, the *degree* d_j of vertex $j \in [n]$ is the number of vertices directly connected to it. (Equivalently, the number of edges connected to j .)

The **degree matrix** of G is the diagonal matrix:

$$\mathbf{D} = \text{diag}(d_1, \dots, d_n).$$

Application III: Connectedness of a random graph, III

With the vertices $[n]$ of G , the **adjacency matrix** \mathbf{A} of G encodes the edges:

$$\mathbf{A} \in \mathbb{R}^{n \times n}, \quad (\mathbf{A})_{ij} = \begin{cases} 0, & i \text{ not connected to } j \\ 1, & i \text{ connected to } j \end{cases}$$

The matrix \mathbf{A} is Hermitian (like \mathbf{D}).

Finally, we need the (unnormalized) *Graph Laplacian*:

$$\mathbf{L} = \mathbf{D} - \mathbf{A} \geq \mathbf{0}.$$

Since \mathbf{L} is spd, its smallest possible eigenvalue is 0.

It turns out: any connected component in G results in a unique eigenvector with eigenvalue 0.

I.e., G is fully connected iff $\lambda = 0$ is a simple eigenvalue, or if the *second* smallest eigenvalue of \mathbf{L} is strictly positive.

Application III: Connectedness of a random graph, IV

By explicit computation, the pair $(0, \mathbf{1})$ is *always* an eigenpair of any graph Laplacian.

With this knowledge, we can consider the graph Laplacian matrix corresponding to a compression that's orthogonal to $\mathbf{1}$.

$$\mathbf{S} := \mathbf{Q}^* \mathbf{L} \mathbf{Q}, \quad \mathbf{Q}^* \mathbf{Q} = \mathbf{I}_{n-1}, \quad \mathbf{Q}^* \mathbf{1} = \mathbf{0}.$$

The smallest eigenvalue of \mathbf{S} is equal to the second smallest eigenvalue of \mathbf{L} .

What's the probability that the smallest eigenvalue of \mathbf{S} is larger than 0?

To connect this to random matrices, suppose that G is from the $G(n, p)$ ERG model.

In this model, given a unique pair $\{i, j\}$ of distinct vertices:

- An edge connecting $\{i, j\}$ is activated according to the outcome of an independent RV $\gamma_{i,j} \sim \text{Bernoulli}(p)$.
- Both d_i and d_j are incremented by 1
- The (i, j) and (j, i) entries of the adjacency matrix take value 1.

Therefore:

$$\mathbf{L} = \sum_{i=1}^n \sum_{j < i} \gamma_{i,j} (\mathbf{e}_i^* \mathbf{e}_i + \mathbf{e}_j^* \mathbf{e}_j - \mathbf{e}_i^* \mathbf{e}_j - \mathbf{e}_j^* \mathbf{e}_i).$$

Application III: Connectedness of a random graph, V

$$\mathbf{S} := \mathbf{Q}^* \mathbf{L} \mathbf{Q} = \sum_{i \in [n]} \sum_{j < i} \mathbf{X}_{i,j}, \quad \mathbf{L} = \sum_{i=1}^n \sum_{j < i} \gamma_{i,j} (\mathbf{e}_i^* \mathbf{e}_i + \mathbf{e}_j^* \mathbf{e}_j - \mathbf{e}_i^* \mathbf{e}_j - \mathbf{e}_j^* \mathbf{e}_i).$$

Therefore, \mathbf{S} is an iid sum of $n(n-1)/2$ random (spd) matrices. One can compute:

$$\mathbb{E} \mathbf{S} = pn \mathbf{I}_{n-1}, \quad \lambda_{\min}(\mathbb{E} \mathbf{S}) = pn, \quad \|\mathbf{X}_{i,j}\|_2 \leq 2.$$

This is enough to employ the multiplicative spd Chernoff estimate:

$$\Pr(\lambda_{\min}(\mathbf{S}) \leq t pn) \leq (n-1) e^{-pn/2} \underbrace{\left(\frac{e}{t}\right)^{tpn/2}}_{\sim 1 \text{ for } t \ll 1}$$

We want this probability to be small, which translates into $p \gtrsim \frac{2 \log(n-1)}{n}$.
(This ensures that the second smallest eigenvalue of \mathbf{L} is positive with large probability.)

Application IV: Random features, I

Let $\{\mathbf{p}_j\}_{j \in [N]} \subset \mathbb{R}^P$ be a set of N points in P -dimensional space. (Formally $P = \infty$ is allowable.)

These points correspond to observed data, and the goal is to develop pairwise comparisons for downstream tasks (classification, identification, compression, etc.)

The *kernel trick* embeds such pairwise comparisons as scalars through a *kernel function*:

$$A_{i,j} = K(\mathbf{p}_i, \mathbf{p}_j), \quad K : \mathbb{R}^P \times \mathbb{R}^P \rightarrow \mathbb{C}$$

The function K is symmetric (e.g., $K(\mathbf{p}, \mathbf{q}) \sim \exp(-\|\mathbf{p} - \mathbf{q}\|_2^2)$), so that \mathbf{A} is Hermitian.

The problem: Evaluating $\mathbf{A} \in \mathbb{R}^{N \times N}$ is expensive, requiring say $\mathcal{O}(N^2P)$.

(P -dependent complexity arises when evaluating function over P inputs.)

Application IV: Random features, II

A particular class of kernels of interest are translation-invariant kernels on \mathbb{R}^P , which are given by,

$$K(\mathbf{p}, \mathbf{q}) = \phi(\mathbf{p} - \mathbf{q}), \quad \phi : \mathbb{R}^P \rightarrow \mathbb{C}.$$

Through Bochner's Theorem, these translation-invariant kernels can be represented as,

$$\phi(\mathbf{p} - \mathbf{q}) = c \int \exp(i \langle \mathbf{p} - \mathbf{q}, \mathbf{w} \rangle) d\mu(\mathbf{w}),$$

for a constant c and a measure μ on \mathbb{R}^P . Written another way, if \mathbf{w} is a random vector whose density function is $d\mu(\mathbf{w})$, then

$$A_{i,j} = \mathbb{E} (\psi(\mathbf{p}_i, \mathbf{w}) \psi(\mathbf{p}_j, \mathbf{w})^*), \quad \psi(\mathbf{p}, \mathbf{w}) = \sqrt{c} \exp(i \langle \mathbf{p}, \mathbf{w} \rangle).$$

Application IV: Random features, III

To demonstrate that this characterization has concrete applicability:

$$K(\mathbf{p}, \mathbf{q}) = \exp(-a\|\mathbf{p} - \mathbf{q}\|_2^2/2) \implies \psi(\mathbf{p}, (\mathbf{w}, U)) = \sqrt{2} \cos(\langle \mathbf{p}, \mathbf{w} \rangle + U),$$

where $\mathbf{w} \sim \mathcal{N}(\mathbf{0}, \alpha \mathbf{I})$ and $U \sim \text{uniform}([0, 2\pi))$.

Hence, we know that \mathbf{A} is the expectation of a random variable \mathbf{X} :

$$\mathbf{X} = \mathbf{x}\mathbf{x}^*, \quad x_j = \psi(\mathbf{p}_j; \mathbf{w}), \quad \mathbb{E}\mathbf{X} = \mathbf{A}.$$

Therefore, we develop a familiar strategy and perform some computations:

$$\mathbf{S}_n = \frac{1}{n} \sum_{j \in [n]} \mathbf{X}_j, \quad \mathbf{X}_j \stackrel{\text{iid}}{\sim} \mathbf{X}, \quad \|\mathbf{X}\|_2 \leq CN, \quad v_2(\mathbf{X}) \leq CN\|\mathbf{A}\|_2,$$

where $|\psi(\mathbf{p}, \mathbf{w})| \leq C$.

These imply that $\mathbb{E}\|\mathbf{S}_n - \mathbf{A}\|_2 \lesssim \sqrt{CN\|\mathbf{A}\|_2 \log(2N)}$.

This translates into $n \gtrsim C \text{idim}(\mathbf{A})/\epsilon^2$ implies $\frac{1}{\|\mathbf{A}\|_2} \mathbb{E}\|\mathbf{S}_n - \mathbf{A}\|_2 \lesssim \epsilon$.