

Assignment 3

Random Matrix Theory in Data Science and Statistics

(EN.553.796, Fall 2024)

Assigned: November 4, 2024 Due: November 22, 2024, 11:59pm

Solve Problem 1, and any three out of the four remaining problems. If you solve more, we will grade the first three solutions (past Problem 1) that you include. Each problem is worth an equal amount towards your grade.

Submit solutions in \LaTeX . Write in complete sentences. Include and justify all steps of your arguments, but try to avoid writing excessive explanation that is not contributing to our understanding your solution. You are welcome to include images if you think that will help explain your solutions.

Problem 1 (Project update). Give an update, in one or two paragraphs, on how your project is going. What have you read or done? Has your goal changed, or have you found that you will need to read something additional or different to understand what you originally planned to? What do you plan to do in the remaining time? Do you have any big-picture questions for us about your reading or work?

Problem 2. This problem is a continuation of Homework 2, Problem 5, Part 3. You may use that result even if you did not solve that problem. Recall the statement: if $\mathbf{g} \sim \mathcal{N}(\mathbf{0}, \Sigma)$ and $\mathbf{h} \sim \mathcal{N}(\mathbf{0}, \Lambda)$ are arbitrary N -dimensional centered Gaussian vectors such that $\mathbb{E}(\mathbf{g}_i - \mathbf{g}_j)^2 \leq \mathbb{E}(\mathbf{h}_i - \mathbf{h}_j)^2$ for all $i, j \in [N]$, then $\mathbb{E} \max_i \mathbf{g}_i \leq \mathbb{E} \max_i \mathbf{h}_i$.

1. The *Gaussian width* of a compact set $\mathcal{X} \subset \mathbb{R}^d$ is

$$\omega(\mathcal{X}) := \mathbb{E} \max_{\mathbf{x} \in \mathcal{X}} \langle \mathbf{g}, \mathbf{x} \rangle.$$

By the rotational invariance of \mathbf{g} , you may view this as measuring the width of \mathcal{X} in the direction of a random ray through the origin in \mathbb{R}^d —hence the name. Let $\mathbf{W} \sim \text{GOE}(d)$ (recall this means $W_{ij} = W_{ji} \sim \mathcal{N}(0, 1 + \mathbb{1}\{i = j\})$) independently for all $i \leq j$). Show that, for any compact $\mathcal{X} \subseteq \mathbb{S}^{d-1}(1)$,

$$\mathbb{E} \max_{\mathbf{x} \in \mathcal{X}} \mathbf{x}^\top \mathbf{W} \mathbf{x} \leq 2\omega(\mathcal{X}).$$

(**HINT:** Use ϵ -nets to reduce to the case of finite \mathcal{X} and use the inequality from Homework 2 cited above. You will find it useful to prove the linear-algebraic inequality $\|\mathbf{w}\mathbf{x}^\top - \mathbf{y}\mathbf{z}^\top\|_F^2 \leq \|\mathbf{w} - \mathbf{y}\|^2 + \|\mathbf{x} - \mathbf{z}\|^2$ for all vectors $\mathbf{w}, \mathbf{x}, \mathbf{y}, \mathbf{z}$ of compatible sizes and each of unit norm.)

- Specializing the above result, prove the following. Note that there are no further “error terms” in any of these statements, and they hold non-asymptotically for all d .

$$\mathbb{E} \lambda_1(\mathbf{W}) = \mathbb{E} \max_{\mathbf{x} \in \mathbb{S}^{d-1}(1)} \mathbf{x}^\top \mathbf{W} \mathbf{x} \leq 2 \cdot \sqrt{d},$$

$$\mathbb{E} \max_{\mathbf{x} \in \{\pm 1/\sqrt{d}\}^d} \mathbf{x}^\top \mathbf{W} \mathbf{x} \leq 2\sqrt{\frac{2}{\pi}} \cdot \sqrt{d},$$

$$\mathbb{E} \max_{\substack{\mathbf{x} \in \mathbb{S}^{d-1}(1) \\ x_i \geq 0 \text{ for all } i \in [d]}} \mathbf{x}^\top \mathbf{W} \mathbf{x} \leq \sqrt{2} \cdot \sqrt{d}.$$

- Let $\mathbf{G} \sim \mathcal{N}(0, 1)^{\otimes d \times m}$, i.e., a random rectangular $d \times m$ matrix with i.i.d. standard Gaussian entries. Adapt your argument from above to show that

$$\mathbb{E} \|\mathbf{G}\| \leq \sqrt{d} + \sqrt{m}.$$

- Let $\mathbf{g} \sim \mathcal{N}(\mathbf{0}, \Sigma)$ be N -dimensional and suppose that $\text{rank}(\Sigma) = N$. Show that the function $D(i, j) := \sqrt{\mathbb{E}(g_i - g_j)^2}$ defines a metric on the set $[N] = \{1, \dots, N\}$. Fix $\delta > 0$, and suppose that there exists a δ -packing of P points in this space: a subset $S \subset [N]$ with $|S| = P$ and such that, for all $i, j \in S$ distinct, $D(i, j) > \delta$. Show that, for some absolute constant $c > 0$,

$$\mathbb{E} \max_{i \in [N]} g_i \geq c\delta \sqrt{\log P}.$$

(**HINT:** Use the above Homework 2 inequality again. Find a very simple Gaussian vector with which to compare the restriction of \mathbf{g} to the indices in S .)

- Derive an inequality relating the Gaussian width of a compact $\mathcal{X} \subset \mathbb{R}^d$ to the *packing number* $P(\delta)$, the maximum number of points in \mathcal{X} at pairwise (Euclidean) distance at least δ . Note that this is a very geometric statement that you have derived by probabilistic reasoning!

Problem 3. This problem, using a technique in a somewhat similar spirit to Homework 2, Problem 5, explores some general aspects of concentration of functions of Gaussian random variables in a different way than we will see in class. Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be a smooth function and let $\mathbf{g} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_n)$ be a standard n -dimensional Gaussian vector.

- Let $\mathbf{g}, \mathbf{h} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_n)$ independently, so that \mathbf{h} is an independent copy of \mathbf{g} . Write $\mathbf{g}(t) = \sin(t)\mathbf{g} + \cos(t)\mathbf{h}$ for $t \in [0, \pi/2]$. This is a kind of interpolation in the style of

Homework 2, Problem 5, but between two independent copies of the *same* Gaussian random vector. Show that

$$f(\mathbf{g}) - f(\mathbf{h}) = \int_0^{\pi/2} \langle \nabla f(\mathbf{g}(t)), \mathbf{g}'(t) \rangle dt,$$

and that $\text{Law}((\mathbf{g}(t), \mathbf{g}'(t))) = \text{Law}(\mathbf{g}, \mathbf{h}) = \mathcal{N}(\mathbf{0}, \mathbf{I}_{2n})$ for all $t \in [0, \pi/2]$, whereby the integrand above has the same law at each $t \in [0, \pi/2]$.

2. Let $\Psi : \mathbb{R} \rightarrow \mathbb{R}$ be a convex function. Show that

$$\mathbb{E}_{\mathbf{g}, \mathbf{h}} \Psi(f(\mathbf{g}) - f(\mathbf{h})) \leq \mathbb{E}_{\mathbf{g}, \mathbf{h}} \Psi\left(\frac{\pi}{2} \langle \nabla f(\mathbf{g}), \mathbf{h} \rangle\right),$$

and therefore also that

$$\mathbb{E}_{\mathbf{g}} \Psi(f(\mathbf{g}) - \mathbb{E}f(\mathbf{g})) \leq \mathbb{E}_{\mathbf{g}, \mathbf{h}} \Psi\left(\frac{\pi}{2} \langle \nabla f(\mathbf{g}), \mathbf{h} \rangle\right).$$

(HINT: Jensen's inequality, repeatedly.)

3. Suppose now that f is also L -Lipschitz, i.e. that $\|\nabla f(\mathbf{x})\| \leq L$ for all $\mathbf{x} \in \mathbb{R}^n$ (for a smooth function this is equivalent to the usual definition). Prove that

$$\mathbb{P}[|f(\mathbf{g}) - \mathbb{E}f(\mathbf{g})| > t] \leq 2 \exp\left(-\frac{2}{\pi^2} \frac{t^2}{L^2}\right).$$

This is a slightly weaker (in the constant in the exponential) version of the Gaussian Lipschitz concentration inequality we discussed in class. In fact relatively simple further arguments show that the same holds even if f is not smooth, and can be Lipschitz in the weaker sense that $|f(\mathbf{x}) - f(\mathbf{y})| \leq L\|\mathbf{x} - \mathbf{y}\|$ for all $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$. You may use this without proof below.

(HINT: Choose the same Ψ that you would use to show a Chernoff bound.)

4. Let $\mathbf{W} \sim \text{GOE}(d)$. Prove that, for an absolute constant $c > 0$,

$$\mathbb{P}[|\lambda_1(\mathbf{W}) - \mathbb{E}\lambda_1(\mathbf{W})| > t] \leq 2 \exp(-ct^2).$$

That is, the norm of a Gaussian random matrix has Gaussian tails. Note that, remarkably, the argument requires no knowledge whatsoever of $\mathbb{E}\lambda_1(\mathbf{W})$, though we have seen separately (including partly in the previous problem) that $\mathbb{E}\lambda_1(\mathbf{W}) \approx 2\sqrt{d}$, so this shows quite strong concentration since the inequality does not depend on d .

5. Prove that the map $\lambda : \mathbb{R}_{\text{sym}}^{d \times d} \rightarrow \mathbb{R}^d$ mapping a matrix to its eigenvalues in descending order is 1-Lipschitz when matrices are given the Frobenius norm and vectors the ℓ^2 norm. Using that, prove that there is an absolute constant $c > 0$ such that, for $\mathbf{W} \sim \text{GOE}(d)$, $\widehat{\mathbf{W}} := \frac{1}{\sqrt{d}}\mathbf{W}$, any L -Lipschitz function $f : \mathbb{R} \rightarrow \mathbb{R}$, and any $i \in [d]$,

$$\mathbb{P}\left[\left|\frac{1}{d} \sum_{i=1}^d f(\lambda_i(\widehat{\mathbf{W}})) - \mathbb{E} \frac{1}{d} \sum_{i=1}^d f(\lambda_i(\widehat{\mathbf{W}}))\right| > t\right] \leq 2 \exp\left(-c \frac{t^2}{L^2} d^2\right),$$

$$\mathbb{P}\left[|f(\lambda_i(\widehat{\mathbf{W}})) - \mathbb{E}f(\lambda_i(\widehat{\mathbf{W}}))| > t\right] \leq 2 \exp\left(-c \frac{t^2}{L^2} d\right).$$

The first result easily upgrades weak convergence in expectation to the semicircle law to various stronger kinds of weak convergence (weak convergence in probability, L^2 or any L^p , almost surely, and so forth).

(HINT: For the part about the λ map, look through your old homework problems.)

Problem 4. I mentioned early in the class that the semicircle distribution is entirely alien to the limit theorems of probability theory outside of random matrix theory. In this problem, you will see that I was lying.

1. Define a sequence of polynomials

$$\begin{aligned} H_0(x) &= 1, \\ H_1(x) &= x, \\ H_{d+1}(x) &= xH_d(x) - H'_d(x). \end{aligned}$$

Show that these are *orthogonal polynomials* for the standard Gaussian measure:

$$\mathbb{E}_{g \sim \mathcal{N}(0,1)} H_a(g)H_b(g) = \begin{cases} 0 & \text{if } a \neq b \\ a! & \text{if } a = b \end{cases}.$$

(HINT: Work with $\mathbb{E}H_a(g)f(g)$ for a general (nice) f .)

2. Show that $H'_d(x) = dH_{d-1}(x)$. Let $\mathbf{W} \sim \text{GOE}(d)$. Recall that the *characteristic polynomial* of \mathbf{W} is the polynomial $p_{\mathbf{W}}(t) = \det(t\mathbf{I}_d - \mathbf{W})$, whose roots are the eigenvalues of \mathbf{W} . Using the above identities, show that, for all d ,

$$\mathbb{E}_{\mathbf{W} \sim \text{GOE}(d)} p_{\mathbf{W}}(t) = H_d(t).$$

Make a guess for what you think the limiting empirical distribution of the collection of roots of $H_d(\sqrt{d} \cdot t)$ might be. Be bold, or, if you can't, try a numerical experiment.

(HINT: Use the Laplace expansion of the determinant.)

3. Show that the H_d are equivalently defined by the identity

$$H_a(x) \exp\left(-\frac{x^2}{2}\right) = (-1)^a \frac{d^a}{dx^a} \exp\left(-\frac{x^2}{2}\right).$$

Conclude that all roots of H_d are real.

(HINT: Show that all derivatives of $\exp(-x^2/2)$ decay to zero as $x \rightarrow \pm\infty$. What can you say about the relationship between the roots of the n th and $(n+1)$ th derivative of such a function?)

4. Show that the roots of H_d are the eigenvalues of the $d \times d$ matrix

$$\mathbf{T}^{(d)} := \begin{bmatrix} 0 & \sqrt{1} & 0 & 0 & \cdots & 0 & 0 \\ \sqrt{1} & 0 & \sqrt{2} & 0 & \cdots & 0 & 0 \\ 0 & \sqrt{2} & 0 & \sqrt{3} & \cdots & 0 & 0 \\ 0 & 0 & \sqrt{3} & 0 & \cdots & 0 & 0 \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & 0 & 0 & \ddots & 0 & \sqrt{d-1} \\ 0 & 0 & 0 & 0 & \ddots & \sqrt{d-1} & 0 \end{bmatrix}$$

(**HINT:** Use again the alternate form of the defining recursion $H_{d+1}(x) = xH_d(x) - dH_{d-1}(x)$.)

5. Show that your prediction in Part 2 is correct: the empirical distribution of the roots of $H_d(\sqrt{d} \cdot t)$ converges weakly to what you guessed. You may be somewhat heuristic since this is a bit of a challenging computation to formalize fully, but do your best to make a convincing calculation.

(**HINT:** Calculate traces of powers of $\mathbf{T}^{(d)}$.)

Problem 5. You will numerically study several variants of the spiked matrix model we saw in class. Recall that that model considered the top eigenpair of $\mathbf{Y} = \mathbf{W} + \beta\sqrt{d}\mathbf{x}\mathbf{x}^\top$ for $\mathbf{W} \sim \text{GOE}(d)$ and some \mathbf{x} with $\|\mathbf{x}\| = 1$. As always, include all plots, tables, etc. needed for us to understand and believe the conclusions you claim to come to based on numerical experiments.

1. Consider a *quantized spiked matrix model*, where instead of \mathbf{Y} you observe $\text{sgn}(\mathbf{Y})$, the matrix with entries $\text{sgn}(Y_{ij}) \in \{\pm 1\}$. Perform experiments. At what value of β does an outlier eigenvalue appear? Consider rounding the values of \mathbf{Y} instead to points on a grid with a given width. How do the results depend on the grid width?
2. Consider a *censored spiked matrix model*, where each entry of \mathbf{Y} is hidden from you with some probability $\delta \in (0, 1)$. A natural way to fill in this “missing data” is to set those entries to zero. At what value of β does an outlier eigenvalue appear in this matrix with random entries set to zero? Make a prediction about the critical β as a function of δ .
3. Take the spiked matrix model with a random spike $\mathbf{x} \sim \text{Unif}(\{\pm 1/\sqrt{d}\}^d)$. Consider attempting to recover \mathbf{x} by fixing some initial guess $\hat{\mathbf{x}}^{(0)}$ and then iterating the map $\hat{\mathbf{x}}^{(t+1)} = \tanh(c\mathbf{Y}\hat{\mathbf{x}}^{(t)})$ for some $c > 0$. Why is this a sensible idea? What is the role of c ? Try it numerically: try starting from $\hat{\mathbf{x}}^{(0)}$ uniformly random, and also from $\hat{\mathbf{x}}^{(0)}$ some scaling of the top eigenvector of \mathbf{Y} . Vary c and β . Can you find a setting where you can design an algorithm that convincingly beats the top eigenvector estimator? (In the sense of $\langle \hat{\mathbf{x}}^{(T)} / \|\hat{\mathbf{x}}^{(T)}\|, \mathbf{x} \rangle^2$ for some large T being typically substantially larger than $\langle \mathbf{v}_1(\mathbf{Y}), \mathbf{x} \rangle^2$.)

4. Suppose instead you choose a random spike by first drawing $\tilde{x} \sim \text{Unif}(\mathbb{S}^{d-1}(1))$, and then taking $x_i := |\tilde{x}_i|$ (i.e., the entrywise absolute value). Think of a way to adapt the strategy from Part 3 to the case where you are promised the spike was generated in this way. Can you again beat the top eigenvector estimator?