

# Lecture 04: Introduction to High Dimensional Asymptotics and Random Matrix Theory

GEOFF PLEISS

In the previous class, we made an intuitive justification for why the double descent phenomenon occurs in high-dimensional ridgeless regression. We now proceed with a more formal analysis of risk as a function of the number of parameters.

Over the next two lectures, we will derive an asymptotic result for ridge regression, a result which will have implications for neural networks as we will see later in the course. This lecture will focus on a light introduction to the mathematical tools needed to derive an asymptotic result. The following lecture will apply these tools to the ridge risk.

---

## 1) Ridge Regression Problem Setup

To begin, let  $\hat{\theta}_\lambda$  be the ridge regression estimator with regularization parameter  $\lambda$  trained on the dataset  $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$ , which we vectorize as  $\mathbf{X} \in \mathbb{R}^{n \times d}$  and  $\mathbf{y} \in \mathbb{R}^n$ . We assume that

$$\mathbf{x} \stackrel{i.i.d.}{\sim} \mathcal{N}(0, \Sigma), \quad y = \boldsymbol{\theta}^{*\top} \mathbf{x} + \sigma \epsilon, \quad \epsilon \stackrel{i.i.d.}{\sim} \mathcal{N}(0, 1), \quad (1)$$

where  $\boldsymbol{\theta}^* \in \mathbb{R}^d$ ,  $\Sigma \in \mathbb{R}^{d \times d}$ , and  $\sigma^2 > 0$  are fixed.<sup>1</sup> To simplify analysis, we will assume that the amount of regularization scales with  $n$ , i.e.:

$$\hat{\theta}_\lambda = (\mathbf{X}^\top \mathbf{X} + n\lambda \mathbf{I})^{-1} \mathbf{X}^\top \mathbf{y}.$$

We can write  $\mathbf{X}^\top \mathbf{X}$  as a scaled **empirical covariance matrix**:

$$\frac{1}{n} \mathbf{X}^\top \mathbf{X} := \hat{\Sigma}, \quad \mathbb{E}[\hat{\Sigma}] = \Sigma.$$

The risk of  $\hat{f}_\lambda(\mathbf{x}) = \mathbf{x}^\top \hat{\theta}_\lambda$  can be factorized into bias and variance terms (as on the problem set):

$$\mathcal{R}(\hat{\theta}_\lambda) = \underbrace{\mathbb{E} \left[ \left( \boldsymbol{\theta}^* - \mathbb{E}[\hat{\theta}_\lambda] \right)^\top \Sigma \left( \boldsymbol{\theta}^* - \mathbb{E}[\hat{\theta}_\lambda] \right) \right]}_{\mathcal{B}(\hat{\theta}_\lambda) = \text{Bias}^2} + \underbrace{\mathbb{E} \left[ \left( \hat{\theta}_\lambda - \mathbb{E}[\hat{\theta}_\lambda] \right)^\top \Sigma \left( \hat{\theta}_\lambda - \mathbb{E}[\hat{\theta}_\lambda] \right) \right]}_{\mathcal{V}(\hat{\theta}_\lambda) = \text{Var}}.$$

Plugging in  $\hat{\theta}_\lambda = (\mathbf{X}^\top \mathbf{X} + n\lambda \mathbf{I})^{-1} \mathbf{X}^\top (\mathbf{X} \boldsymbol{\theta}^* + \sigma \epsilon) = (\hat{\Sigma} + \lambda \mathbf{I})^{-1} (\hat{\Sigma} \boldsymbol{\theta}^* + \frac{\sigma}{n} \mathbf{X}^\top \epsilon)$ , where  $\sigma \epsilon = \mathbf{y} - \mathbf{X} \boldsymbol{\theta}^*$ , and simplifying, we have

$$\begin{aligned} \mathcal{B}(\hat{\theta}_\lambda) &= \boldsymbol{\theta}^{*\top} \mathbb{E} \left[ \left( \mathbf{I} - (\hat{\Sigma} + \lambda \mathbf{I})^{-1} \hat{\Sigma} \right) \Sigma \left( \mathbf{I} - (\hat{\Sigma} + \lambda \mathbf{I})^{-1} \hat{\Sigma} \right) \right] \boldsymbol{\theta}^* \\ &= \lambda^2 \boldsymbol{\theta}^{*\top} \mathbb{E} \left[ \left( \hat{\Sigma} + \lambda \mathbf{I} \right)^{-1} \Sigma \left( \hat{\Sigma} + \lambda \mathbf{I} \right)^{-1} \right] \boldsymbol{\theta}^*, \quad (\text{Woodbury on } (\mathbf{I} - (\hat{\Sigma} + \lambda \mathbf{I})^{-1} \hat{\Sigma})) \\ \mathcal{V}(\hat{\theta}_\lambda) &= \frac{\sigma^2}{n^2} \mathbb{E} \left[ \epsilon^\top \mathbf{X} \left( \hat{\Sigma} + \lambda \mathbf{I} \right)^{-1} \Sigma \left( \hat{\Sigma} + \lambda \mathbf{I} \right)^{-1} \mathbf{X}^\top \epsilon \right] \\ &= \frac{\sigma^2}{n} \text{Tr} \mathbb{E} \left[ \hat{\Sigma} \left( \hat{\Sigma} + \lambda \mathbf{I} \right)^{-1} \Sigma \left( \hat{\Sigma} + \lambda \mathbf{I} \right)^{-1} \right]. \quad (\text{independence of } \mathbf{X}, \epsilon, \text{ cyclic prop. of Tr}) \end{aligned}$$

---

<sup>1</sup>Technically, we only assume that  $\mathbf{x}$  and  $\epsilon$  are sampled from **sub-Gaussian** distributions with the moments given in Eq. (1).

## 2) The Need for High Dimensional Asymptotics and Random Matrix Theory

Unfortunately, these expressions are about as simple as we can make them in closed form. Both  $\mathcal{B}(\hat{\boldsymbol{\theta}}_\lambda)$  and  $\mathcal{V}(\hat{\boldsymbol{\theta}}_\lambda)$  are expectations of some complex function of the random matrix  $\hat{\boldsymbol{\Sigma}}$ , and the only functions we can compute in closed form are:

1.  $\mathbf{E}[\hat{\boldsymbol{\Sigma}}] = \boldsymbol{\Sigma}$  and
2.  $\mathbf{E}[(\hat{\boldsymbol{\Sigma}})^{-1}] = (\frac{1}{n-d-1})\boldsymbol{\Sigma}^{-1}$ , if  $n > d + 1$  and assuming  $\mathbf{X}$  is Gaussian.

As in most of statistics, we can assume that some “nice simplifying structure” emerges that simplifies these expressions when the problems “get really large.” However, we need to be VERY careful about what we mean by “get really large.” If we simply take  $n \rightarrow \infty$ , then  $\hat{\boldsymbol{\Sigma}} \rightarrow \boldsymbol{\Sigma}$  and our bias and variance both disappear as  $\lambda \rightarrow 0$ . This analysis may be appropriate if we’re trying to model linear regression with  $n \gg d$ , but it’s going to be a horrible model for problem where  $n \approx d$  or  $d > n$ .

The correct framework for asymptotically analyzing these problems requires a conceptual leap. We will examine what happens when  $n, d \rightarrow \infty$  *simultaneously*. In other words, we will assume that  $d$  grows linearly with  $n$ , i.e.  $d = \gamma n$  for some fixed  $\gamma > 0$ , and then we will take  $n$  (and  $d$ )  $\rightarrow \infty$ . This limit is known as the **high-dimensional asymptotic regime** and analyzing it requires tools from **random matrix theory**.

---

## 3) Warm-Up: The Marchenko-Pastur Distribution for Isotropic Data

Imagine for a second that  $\boldsymbol{\Sigma} = \mathbf{I}$ . Then

$$\mathcal{V}(\hat{\boldsymbol{\theta}}_\lambda) = \sigma^2 \text{Tr} \mathbb{E} \left[ \hat{\boldsymbol{\Sigma}}(\hat{\boldsymbol{\Sigma}} + \lambda \mathbf{I})^{-2} \right] = \sigma^2 \sum_{i=1}^d \mathbb{E} \left[ \frac{\hat{s}_i}{(\hat{s}_i + \lambda)^2} \right], \quad (2)$$

where  $\hat{s}_i$  are the eigenvalues of  $\hat{\boldsymbol{\Sigma}}$ . Thus, understanding  $\mathcal{V}(\hat{\boldsymbol{\theta}}_\lambda)$  as  $n, d \rightarrow \infty$  requires understanding what happens to the  $\hat{s}_i$  as  $n, d \rightarrow \infty$ . We will accomplish this understanding by viewing the summation over eigenvalues through a probabilistic lens.

If we define  $\hat{F}(\cdot)$  as the **empirical distribution over eigenvalues**

$$\hat{F}(s) = \frac{1}{d} \sum_{i=1}^d \mathbf{1}[s = \hat{s}_i],$$

then we can rewrite Eq. (2) as an expectation:

$$\mathcal{V}(\hat{\boldsymbol{\theta}}_\lambda) = \frac{\sigma^2}{d} \int \frac{s}{(s + \lambda)^2} d\hat{F}(s).$$

Just as the central limit theorem tells us that empirical distributions of sums converge to normal distributions, the **Marchenko-Pastur theorem** tells us that the empirical distribution of eigenvalues also converge to a deterministic distribution.

**Theorem 1** (Marchenko and Pastur [1967]). *Let  $\mathbf{X} \in \mathbb{R}^{n \times d}$  have i.i.d. sub-Gaussian entries with mean 0 and variance 1. Assuming the ratio  $\gamma = d/n \in (0, 1]$  is fixed, the empirical distribution of eigenvalues  $\hat{F}(s)$  of  $\hat{\boldsymbol{\Sigma}} = \frac{1}{n} \mathbf{X}^\top \mathbf{X}$  converges (in distribution) to the **Marchenko-Pastur distribution**  $F(s)$  as  $n, d \rightarrow \infty$ :*

$$\lim_{n, d \rightarrow \infty} \hat{F}(s) = F(s), \quad \frac{dF(s)}{ds} = \begin{cases} \frac{1}{2\pi\gamma s} \sqrt{(s_+ - s)(s - s_-^2)} & s \in [s_-, s_+], \\ 0 & \text{otherwise.} \end{cases} \quad (3)$$

where  $s_\pm = (1 \pm \sqrt{\gamma})^2$ .

Note that we can arrive at a density for eigenvalues for the  $\gamma > 1$  (i.e.  $d < n$ ) by recognizing that the non-zero eigenvalues of  $\mathbf{X}^\top \mathbf{X}$  are the same as the non-zero eigenvalues of  $\mathbf{X} \mathbf{X}^\top$ .

It may be strange to think about a continuous distribution over eigenvalues of a matrix! To gain an intuitive understanding, consider a histogram of eigenvalues of  $\frac{1}{n} \mathbf{X}^\top \mathbf{X}$ . For large  $n, d$ , this histogram can be smoothed out with the Marchenko-Pastur density. See Fig. 1 for an illustration.

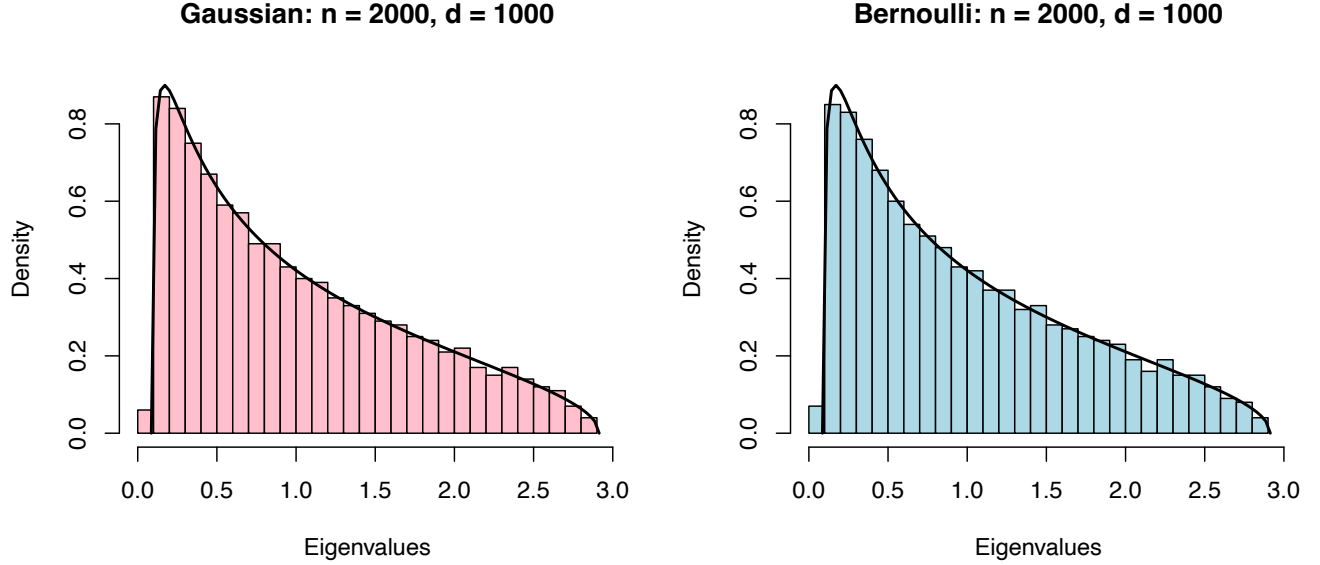


Figure 1: Histogram of eigenvalues of  $\frac{1}{n} \mathbf{X}^\top \mathbf{X}$  for  $n = 2000$ ,  $d = 1000$ , when the entries of  $\mathbf{X}$  are drawn i.i.d. from a Gaussian distribution (left) versus a Bernoulli distribution (right) with zero mean and unit variance. The values of the histogram roughly correspond to the Marchenko-Pastur density. (Figure reproduced from Tibishriani [2023].)

Therefore, the variance for isotropic ridge regression can be approximated by its high-dimensional asymptotic limit:

$$\mathcal{V}(\hat{\theta}_\lambda) = \frac{\sigma^2}{d} \int \frac{s}{(s + \lambda)^2} d\hat{F}(s) \approx \frac{\sigma^2}{d} \int \frac{s}{(s + \lambda)^2} \frac{dF(s)}{ds} ds,$$

which is now an analytic one-dimensional integral that we can numerically solve!

### 3.1 Gaussian Universality

Notably, from Fig. 1 we see that the distribution of eigenvalues doesn't really depend on the distribution of the entries of  $\mathbf{X}$ . The left and right histograms, corresponding to Gaussian and Bernoulli entries (respectively) for  $\mathbf{Z}$ , are nearly identical. This **universality** is a key feature of random matrix theory: in the high-dimensional asymptotic limit, the properties of a random matrix often only depends on its first and second moments and not on the specific distribution of the entries. Therefore, we can often analyze the behaviour of random matrices with Gaussian entries without much loss of generality.

#### 4) A Hand-Wavy Derivation of the Marchenko-Pastur Theorem with the Stieltjes Transform

How do we arrive at the Marchenko-Pastur distribution for isotropic data? We will need to introduce a concept from probability known as the **Stieltjes transform**. The Stieltjes transform of a symmetric matrix  $\mathbf{A}$  is a  $\mathbb{R} \rightarrow \mathbb{R}$  function defined as

$$m_{\mathbf{A}}(-\lambda) = \frac{1}{d} \text{Tr} \left( (\mathbf{A} + \lambda \mathbf{I})^{-1} \right), \quad \lambda > 0, \quad \mathbf{A} \in \mathbb{R}^{d \times d}.$$

More generally, the Stieltjes transform of a probability distribution  $F(s)$  is defined as  $m_F(-\lambda) = \int \frac{1}{s+\lambda} dF(s)$ , and so the Stieltjes transform of a matrix is just the Stieltjes transform of its empirical eigenvalue distribution. Like characteristic functions, Stieltjes transforms are an alternative characterization of a probability measure.<sup>2</sup> Moreover, convergence in Stieltjes transform implies convergence in distribution—i.e. given a sequence of probability distributions  $F_n(s)$  and some limiting distribution  $F(s)$ :

$$m_{F_n}(-\lambda) \rightarrow m_F(-\lambda) \iff F_n(s) \rightarrow F(s).$$

If we assume that there is limiting matrix  $\mathbf{A}$  so that, for all  $\lambda > 0$ ,

$$\frac{1}{d} \text{Tr} \left( (\hat{\Sigma} + \lambda \mathbf{I})^{-1} \right) \rightarrow \frac{1}{d} \text{Tr} \left( (\mathbf{A} + \lambda \mathbf{I})^{-1} \right), \quad (4)$$

(i.e.  $m_{\hat{\Sigma}}(-\lambda) \rightarrow m_{\mathbf{A}}(-\lambda)$ ), then it should also be true that, for all  $i \in [1, n]$ ,

$$\frac{1}{d} \text{Tr} \left( \left( (\hat{\Sigma} + \lambda \mathbf{I}) - \frac{1}{n} \mathbf{x}_i \mathbf{x}_i^\top \right)^{-1} \right) \rightarrow \frac{1}{d} \text{Tr} \left( (\mathbf{A} + \lambda \mathbf{I})^{-1} \right), \quad (5)$$

i.e. that  $m_{\hat{\Sigma}}(-\lambda) \rightarrow m_{\mathbf{A}}(-\lambda)$  even if we removed one  $\mathbf{x}_i$  from the empirical covariance matrix, since each data point contributes infinitesimally to  $\hat{\Sigma}$  as  $n, d \rightarrow \infty$ . Following this argument<sup>3</sup>, coupled with copious amounts of measure theory, the fact that both Eqs. (4) and (5) hold implies that  $\lim_{n,d \rightarrow \infty} m_{\hat{\Sigma}}(-\lambda)$  must adhere to the **self-consistency equation**:

$$\frac{1}{d} \text{Tr} \left( (\hat{\Sigma} (\hat{\Sigma} + \lambda \mathbf{I})^{-1}) \right) \longrightarrow \frac{1}{d} \text{Tr} \left( (\Sigma (\Sigma + \kappa(\lambda) \mathbf{I})^{-1}) \right), \quad 1/\kappa(\lambda) := \lim_{n,d \rightarrow \infty} \frac{1}{n} \text{Tr} \left( \left( \frac{1}{n} \mathbf{X} \mathbf{X}^\top + \lambda \mathbf{I} \right)^{-1} \right), \quad (6)$$

where we recognize  $\frac{1}{n} \text{Tr} \left( \left( \frac{1}{n} \mathbf{X} \mathbf{X}^\top + \lambda \mathbf{I} \right)^{-1} \right)$  as the Steiltjes transform of the empirical distribution of the spectrum of  $\frac{1}{n} \mathbf{X} \mathbf{X}^\top$ .

**Aside.** It is worth thinking about what  $\lim_{d \rightarrow \infty} \frac{1}{d} \text{Tr}(\Sigma(\Sigma + \lambda \mathbf{I})^{-1})$  even means.  $\Sigma \in \mathbb{R}^{d \times d}$  is a fixed matrix, so what does it to consider a limit over  $d$ ? There are two points worth considering:

- Rigorously defining what  $\lim_{d \rightarrow \infty} \frac{1}{d} \text{Tr}(\Sigma(\Sigma + \lambda \mathbf{I})^{-1})$  means depends on the particular problem setup. Most of the “work” in theory papers using random matrix theory involves setting up a mathematically valid interpretation Eq. (6). For the purposes of this course, we can just assume the hand-wavy interpretation that  $\frac{1}{d} \text{Tr}(\Sigma(\Sigma + \lambda \mathbf{I})^{-1}) \approx \frac{1}{d} \text{Tr}(\hat{\Sigma}(\hat{\Sigma} + \lambda \mathbf{I})^{-1})$  when  $n, d$  are “sufficiently large.”
- As  $d \rightarrow \infty$  and  $\Sigma$  becomes infinitely large, we can go back to our notion of kernels. If  $\Sigma = \mathbb{E}[\mathbf{x} \mathbf{x}^\top]$ , then  $k(\mathbf{x}, \mathbf{x}') = \lim_{d \rightarrow \infty} \mathbf{x}^\top \mathbf{x}'$  remains valid and the spectrum of  $k(\cdot, \cdot)$  matches the limiting spectrum of  $\Sigma$ .

<sup>2</sup>To gain intuition for why the Steiltjes transform characterizes a probability measure, note that the Taylor expansion suggests that  $\mathbb{E}_s[1/(s - \lambda)] = -\frac{1}{\lambda} \sum_{j=1}^{\infty} \mathbb{E}_s[(s/\lambda)^j]$  and thus the Steiltjes transform determines the distributions’ moments.

<sup>3</sup>A hand-wavy version of this argument can be shown by applying the Woodbury formula to  $((\hat{\Sigma} + \lambda \mathbf{I}) - \frac{1}{n} \mathbf{x}_i \mathbf{x}_i^\top)^{-1}$  and recognizing that  $\mathbf{x}_i^\top \mathbf{B} \mathbf{x}_i \approx \text{Tr} \mathbf{B}$  for a matrix  $\mathbf{B}$  that is independent of  $\mathbf{x}_i$ . See [Pedregosa et al., 2021, Part 2] for a simple introduction to the argument.

Using the Woodbury formula, the cyclic property of the trace, and rewriting  $\hat{\Sigma} = \frac{1}{n}\mathbf{X}\mathbf{X}^\top$  we have that

$$\begin{aligned} \frac{1}{d} \left[ \text{Tr}(\hat{\Sigma}(\hat{\Sigma} + \lambda\mathbf{I})^{-1}) \right] &= \frac{1}{d} \left[ \text{Tr}(\mathbf{X}(\mathbf{X}^\top\mathbf{X} + n\lambda\mathbf{I})^{-1}\mathbf{X}^\top) \right] = \frac{1}{d} \left[ \text{Tr}(\mathbf{I}_{n \times n} - \lambda \text{Tr}(\frac{1}{n}\mathbf{X}\mathbf{X}^\top + \lambda\mathbf{I})) \right] \\ &= \frac{n}{d} \left[ 1 - \lambda \frac{1}{n} \text{Tr}(\frac{1}{n}\mathbf{X}\mathbf{X}^\top + \lambda\mathbf{I}) \right] \\ &\rightarrow \gamma [1 - \lambda/\kappa(\lambda)]. \end{aligned}$$

Thus Eq. (6) can be rearranged to be written as

$$\underbrace{\gamma \frac{1}{d} \text{Tr}(\Sigma(\Sigma + \kappa(\lambda)\mathbf{I})^{-1})}_{\frac{1}{n} \sum_{i=1}^d \frac{s_i}{s_i + \kappa(\lambda)}} + \frac{\lambda}{\kappa(\lambda)} = 1, \quad (7)$$

which is known as the **Silverstein equation** [Silverstein, 1995]. In the deep learning literature, this equation is often referred to as the implicit regularization equation [Jacot et al., 2020] for reasons that we will see in next lecture.

We now have a recipe to derive the limiting spectrum of  $\hat{\Sigma}$ :

1. Given a  $\Sigma$  and a  $\gamma$ , solve for the  $\kappa(\lambda)$  function that satisfies Eq. (7).
2. Recognizing from Eq. (6) that  $1/\kappa(\lambda)$  is the limiting Steiltjes transform of the spectrum of  $\frac{1}{n}\mathbf{X}\mathbf{X}^\top$ , and that limiting Steiltjes transformations characterize limiting distributions, find the distribution  $F(s)$  with a Steiltjes transform that matches  $1/\kappa(\lambda)$ .
3.  $F(s)$  is the limiting distribution of the eigenvalues of  $\frac{1}{n}\mathbf{X}\mathbf{X}^\top$ , which is also the limiting distribution of the non-zero eigenvalues of  $\frac{1}{n}\mathbf{X}^\top\mathbf{X} = \hat{\Sigma}$ .

When  $\Sigma = \mathbf{I}$ ,  $\kappa(\lambda)$  in Eq. (7) admits a closed-form analytic solution, where  $1/\kappa(\lambda)$  is exactly equal to the Steiltjes transformation of the Marchenko-Pastur distribution given in Eq. (3).

## 5) Strategy for Analyzing High-Dimensional Ridge Risk: Deterministic Equivalents

Unfortunately, the recipe above will only work in very special cases, since most  $\Sigma \neq \mathbf{I}$  don't afford closed-form solutions to  $\kappa(\lambda)$  in Eq. (7). Luckily, for our purposes, we won't need to actually solve for the limiting spectrum of  $\hat{\Sigma}$ ; we only care about specific *reductions* of the limiting spectrum. Recall that for the variance we need to compute  $\text{Tr}(\hat{\Sigma}(\hat{\Sigma} + \lambda\mathbf{I})^{-1}\Sigma(\hat{\Sigma} + \lambda\mathbf{I})^{-1})$ . This computation does not require knowing the exact eigenvalues of  $\hat{\Sigma}$ , but instead requires knowing the summation (trace) of some function of the eigenvalues.

As a starting point, we know an asymptotic limit of the reduction  $\text{Tr}(\hat{\Sigma}(\hat{\Sigma} + \lambda\mathbf{I})^{-1})$ ; it is approximately equal to  $\text{Tr}(\Sigma(\Sigma + \kappa(\lambda)\mathbf{I})^{-1})$ . By adding and subtracting  $\mathbf{I}$  from both sides of Eq. (6), we have that

$$\lambda \text{Tr}(\hat{\Sigma} + \lambda\mathbf{I})^{-1} \approx \kappa(\lambda) \text{Tr}(\Sigma + \kappa(\lambda)\mathbf{I})^{-1}.$$

More rigorously, we have that

$$\frac{1}{d} \left[ \lambda \text{Tr}(\mathbf{B}(\hat{\Sigma} + \lambda\mathbf{I})^{-1}) - \kappa(\lambda) \text{Tr}(\mathbf{B}(\Sigma + \kappa(\lambda)\mathbf{I})^{-1}) \right] \rightarrow 0, \quad (8)$$

for any  $\mathbf{B} \in \mathbb{R}^{d \times d}$  under certain regularity conditions [Rubio and Mestre, 2011]. We will denote this asymptotic convergence by the symbol

$$\lambda(\hat{\Sigma} + \lambda\mathbf{I})^{-1} \asymp \kappa(\lambda)(\Sigma + \kappa(\lambda)\mathbf{I})^{-1} \iff \text{Eq. (8) holds for all "regular" } \mathbf{B} \in \mathbb{R}^{d \times d}.$$

Eq. (8) gives us a powerful mechanism for analyzing equations involving random matrices. Every time we come across a term that looks like  $\lambda \text{Tr}(\mathbf{B}(\hat{\Sigma} + \lambda \mathbf{I})^{-1})$ , we can “swap it out”<sup>4</sup> with its **deterministic equivalent**  $\kappa(\lambda) \text{Tr}(\mathbf{B}(\Sigma + \kappa(\lambda) \mathbf{I})^{-1})$ . Once we have arrived at a final expression that involves only deterministic matrices, we can then think about bounding (or numerically solving) for  $\kappa(\lambda)$ .

We will apply this strategy to the high-dimensional ridge risk in the next lecture.

---

## References

- A. Jacot, B. Simsek, F. Spadaro, C. Hongler, and F. Gabriel. Implicit regularization of random feature models. In *International Conference on Machine Learning*, pages 4631–4640, 2020.
- V. A. Marchenko and L. A. Pastur. Distribution of eigenvalues for some sets of random matrices. *Matematicheskii Sbornik*, 114(4):507–536, 1967.
- F. Pedregosa, C. Paquette, T. Trogon, and J. Pennington. Random matrix theory and machine learning tutorial, 2021. URL <https://random-matrix-learning.github.io/>.
- F. Rubio and X. Mestre. Spectral convergence for a general class of random matrices. *Statistics and Probability Letters*, 81(5):592–602, 2011.
- J. W. Silverstein. Strong convergence of the empirical distribution of eigenvalues of large dimensional random matrices. *Journal of Multivariate Analysis*, 55(2):331–339, 1995.
- R. Tibishriani. High-dimensional regression: Ridge, 2023. URL <https://www.stat.berkeley.edu/~ryantibs/statlearn-s23/lectures/ridge.pdf>.

---

<sup>4</sup>For mathematically rigorous results, this “swap” requires appropriate conditions and analysis of limits. However, in this course we will assume that there is some reasonable notion of approximation that allows us to swap these terms.