

University of Pennsylvania
STAT 991: Random Matrix Theory Class Presentation

Asymptotic Risk of

High-Dimensional Regression

Behrad Moniri
bemoniri@seas.upenn.edu

Linear Regression

- **Data Generation:**



High-Dimensional Linear Regression

- **Data Generation:**

$$\left\{ \begin{array}{l} \beta \in \mathbb{R}^d, \Sigma \in \mathbb{R}^{d \times d} \\ \end{array} \right.$$



- **Data Generation:**

$$\begin{cases} \beta \in \mathbb{R}^d, \Sigma \in \mathbb{R}^{d \times d} \\ \{x_i\}_{i=1}^n \sim \mathcal{N}(0, \Sigma) \end{cases}$$



- **Data Generation:**

$$\begin{cases} \beta \in \mathbb{R}^d, \Sigma \in \mathbb{R}^{d \times d} \\ \{x_i\}_{i=1}^n \sim \mathcal{N}(0, \Sigma) \\ y_i = \beta^\top x_i + \varepsilon_i \text{ where } \varepsilon_i \sim \mathcal{N}(0, \sigma^2) \end{cases} \quad (1)$$



- **Data Generation:**

$$\begin{cases} \beta \in \mathbb{R}^d, \Sigma \in \mathbb{R}^{d \times d} \\ \{x_i\}_{i=1}^n \sim \mathcal{N}(0, \Sigma) \\ y_i = \beta^\top x_i + \varepsilon_i \text{ where } \varepsilon_i \sim \mathcal{N}(0, \sigma^2) \end{cases} \quad (1)$$

- **Fit with ridge regression:**



- **Data Generation:**

$$\begin{cases} \beta \in \mathbb{R}^d, \Sigma \in \mathbb{R}^{d \times d} \\ \{x_i\}_{i=1}^n \sim \mathcal{N}(0, \Sigma) \\ y_i = \beta^\top x_i + \varepsilon_i \text{ where } \varepsilon_i \sim \mathcal{N}(0, \sigma^2) \end{cases} \quad (1)$$

- **Fit with ridge regression:**

$$\hat{\beta}_\lambda = \arg \min_{b \in \mathbb{R}^d} \left[\frac{1}{n} \sum_{i=1}^n (y_i - b^\top x_i)^2 + \lambda \|b\|_2^2 \right] \quad (2)$$



- **Data Generation:**

$$\begin{cases} \beta \in \mathbb{R}^d, \Sigma \in \mathbb{R}^{d \times d} \\ \{x_i\}_{i=1}^n \sim \mathcal{N}(0, \Sigma) \\ y_i = \beta^\top x_i + \varepsilon_i \text{ where } \varepsilon_i \sim \mathcal{N}(0, \sigma^2) \end{cases} \quad (1)$$

- **Fit with ridge regression:**

$$\hat{\beta}_\lambda = \arg \min_{b \in \mathbb{R}^d} \left[\frac{1}{n} \sum_{i=1}^n (y_i - b^\top x_i)^2 + \lambda \|b\|_2^2 \right] \quad (2)$$

- **Random design assumption:** β is random with

$$\mathbb{E}[\beta] = 0 \quad \text{and} \quad \text{Cov}(\beta) = \frac{\alpha^2}{d} I_{d \times d} \quad (3)$$



Solution

- Let $X \in \mathbb{R}^{n \times d}$ and $Y \in \mathbb{R}^n$ be the training data.



- Let $X \in \mathbb{R}^{n \times d}$ and $Y \in \mathbb{R}^n$ be the training data.
- Ridge regression has a closed-form solution

$$\hat{\beta}_\lambda = (X^\top X + n\lambda I)^{-1} X^\top Y$$



Risk of the estimator

- **Question:** What is the risk of $\hat{\beta}_\lambda$?



- **Question:** What is the risk of $\hat{\beta}_\lambda$?

$$\begin{aligned}r_\lambda(X) &= \mathbb{E} \left[\left(x^\top \hat{\beta}_\lambda - x^\top \beta - \varepsilon \right)^2 \mid X \right] \\ &= 1 + \mathbb{E} \left[\left\{ x^\top \left(\hat{\beta}_\lambda - \beta \right) \right\}^2 \mid X \right] \\ &= 1 + \mathbb{E} \left[\left(\hat{\beta}_\lambda - \beta \right)^\top \Sigma \left(\hat{\beta}_\lambda - \beta \right) \mid X \right]\end{aligned}$$



- **Question:** What is the risk of $\hat{\beta}_\lambda$?

$$\begin{aligned}r_\lambda(\mathbf{X}) &= \mathbb{E} \left[\left(\mathbf{x}^\top \hat{\beta}_\lambda - \mathbf{x}^\top \beta - \varepsilon \right)^2 \mid \mathbf{X} \right] \\&= 1 + \mathbb{E} \left[\left\{ \mathbf{x}^\top \left(\hat{\beta}_\lambda - \beta \right) \right\}^2 \mid \mathbf{X} \right] \\&= 1 + \mathbb{E} \left[\left(\hat{\beta}_\lambda - \beta \right)^\top \Sigma \left(\hat{\beta}_\lambda - \beta \right) \mid \mathbf{X} \right]\end{aligned}$$

- We can also write

$$\begin{aligned}\hat{\beta}_\lambda - \beta &= \left(\mathbf{X}^\top \mathbf{X} + \lambda n \mathbf{I} \right)^{-1} \mathbf{X}^\top (\mathbf{X} \beta + \varepsilon) - \beta \\&= -\lambda n \left(\mathbf{X}^\top \mathbf{X} + \lambda n \mathbf{I} \right)^{-1} \beta + \left(\mathbf{X}^\top \mathbf{X} + \lambda n \mathbf{I} \right)^{-1} \mathbf{X}^\top \varepsilon\end{aligned}$$



- If we plug this back to the risk, we get

$$\begin{aligned} r_\lambda(\mathbf{X}) &= 1 + (\lambda n)^2 \mathbb{E} \left[\beta^\top (\mathbf{X}^\top \mathbf{X} + \lambda n \mathbf{I})^{-1} \Sigma (\mathbf{X}^\top \mathbf{X} + \lambda n \mathbf{I})^{-1} \beta \mid \mathbf{X} \right] \\ &\quad + \mathbb{E} \left[\varepsilon^\top \mathbf{X} (\mathbf{X}^\top \mathbf{X} + \lambda n \mathbf{I})^{-1} \Sigma (\mathbf{X}^\top \mathbf{X} + \lambda n \mathbf{I})^{-1} \mathbf{X}^\top \varepsilon \mid \mathbf{X} \right] \end{aligned}$$



- If we plug this back to the risk, we get

$$\begin{aligned}r_{\lambda}(X) &= 1 + (\lambda n)^2 \mathbb{E} \left[\beta^{\top} (X^{\top} X + \lambda n I)^{-1} \Sigma (X^{\top} X + \lambda n I)^{-1} \beta \mid X \right] \\ &\quad + \mathbb{E} \left[\varepsilon^{\top} X (X^{\top} X + \lambda n I)^{-1} \Sigma (X^{\top} X + \lambda n I)^{-1} X^{\top} \varepsilon \mid X \right] \\ &= 1 + \frac{(\lambda n)^2 \alpha^2}{d} \text{Tr} \left(\Sigma (X^{\top} X + \lambda n I)^{-2} \right) \\ &\quad + \text{Tr} \left(\Sigma (X^{\top} X + \lambda n I)^{-1} X^{\top} X (X^{\top} X + \lambda n I)^{-1} \right).\end{aligned}$$



- If we plug this back to the risk, we get

$$\begin{aligned}r_\lambda(X) &= 1 + (\lambda n)^2 \mathbb{E} \left[\beta^\top (X^\top X + \lambda n I)^{-1} \Sigma (X^\top X + \lambda n I)^{-1} \beta \mid X \right] \\ &\quad + \mathbb{E} \left[\varepsilon^\top X (X^\top X + \lambda n I)^{-1} \Sigma (X^\top X + \lambda n I)^{-1} X^\top \varepsilon \mid X \right] \\ &= 1 + \frac{(\lambda n)^2 \alpha^2}{d} \text{Tr} \left(\Sigma (X^\top X + \lambda n I)^{-2} \right) \\ &\quad + \text{Tr} \left(\Sigma (X^\top X + \lambda n I)^{-1} X^\top X (X^\top X + \lambda n I)^{-1} \right).\end{aligned}$$

- Set $\hat{\Sigma} = n^{-1} X^\top X$ and $\gamma_d = d/n$:

$$r_\lambda(X) = 1 + \frac{\gamma_d}{d} \text{Tr} \left(\Sigma \left(\hat{\Sigma} + \lambda I \right)^{-1} \right) + (\lambda \alpha^2 - \gamma_d) \frac{\lambda}{d} \text{Tr} \left(\Sigma \left(\hat{\Sigma} + \lambda I \right)^{-2} \right)$$



So, to analyze the asymptotic risk, we should look at the following two traces:

- **Trace 1:**

$$\frac{\gamma_d}{d} \text{Tr} \left(\Sigma \left(\widehat{\Sigma} + \lambda I \right)^{-1} \right)$$

- **Trace 2:**

$$\left(\lambda \alpha^2 - \gamma_d \right) \frac{\lambda}{d} \text{Tr} \left(\Sigma \left(\widehat{\Sigma} + \lambda I \right)^{-2} \right)$$



Notation and Assumptions

- For a matrix A , define

$$F_A(x) = p^{-1} \sum_{i=1}^p \mathbf{I}(\lambda_i(A) \leq x).$$



Notation and Assumptions

- For a matrix A , define

$$F_A(x) = p^{-1} \sum_{i=1}^p \mathbf{I}(\lambda_i(A) \leq x).$$

- We will take the limit $\gamma_d = d/n \rightarrow \gamma > 0$.



Notation and Assumptions

- For a matrix A , define

$$F_A(x) = p^{-1} \sum_{i=1}^p \mathbf{I}(\lambda_i(A) \leq x).$$

- We will take the limit $\gamma_d = d/n \rightarrow \gamma > 0$.
- Assume that the spectral distribution F_Σ of Σ converges to a limit H supported on $[0, \infty)$.



Notation and Assumptions

- We defined the Stieltjes transform of a measure G as

$$m_G(z) = \int_{l=0}^{\infty} \frac{dG(l)}{l-z}, \quad z \in \mathbb{C} \setminus \mathbb{R}^+.$$



Notation and Assumptions

- We defined the Stieltjes transform of a measure G as

$$m_G(z) = \int_{l=0}^{\infty} \frac{dG(l)}{l-z}, \quad z \in \mathbb{C} \setminus \mathbb{R}^+.$$

- Stieltjes transform of the spectral measure of $\hat{\Sigma}$

$$m_{\hat{\Sigma}}(z) = \frac{1}{p} \operatorname{tr} \left(\left(\hat{\Sigma} - zI_{p \times p} \right)^{-1} \right) \text{ converges to } m(z)$$



Notation and Assumptions

- We defined the Stieltjes transform of a measure G as

$$m_G(z) = \int_{l=0}^{\infty} \frac{dG(l)}{l-z}, \quad z \in \mathbb{C} \setminus \mathbb{R}^+.$$

- Stieltjes transform of the spectral measure of $\hat{\Sigma}$

$$m_{\hat{\Sigma}}(z) = \frac{1}{p} \operatorname{tr} \left(\left(\hat{\Sigma} - zI_{p \times p} \right)^{-1} \right) \text{ converges to } m(z)$$

- The companion Stieltjes transform is defined as

$$\gamma(m(z) + 1/z) = v(z) + 1/z \text{ for all } z \in \mathbb{C} \setminus \mathbb{R}^+$$



A result of Ledoit and Péché (2011)

The main step to derive the limiting risk, is proving the following theorem:

Theorem (Ledoit and Péché (2011))

Assume that $m_{\widehat{\Sigma}}(z) \rightarrow m(z)$ and let $v(z)$ be the companion transform for $m(z)$. We have

$$\frac{1}{d} \operatorname{tr} \left(\Sigma \left(\widehat{\Sigma} + \lambda I_{d \times d} \right)^{-1} \right) \rightarrow_{a.s.} \frac{1}{\gamma} \left(\frac{1}{\lambda v(-\lambda)} - 1 \right) \quad \text{as } n, d \rightarrow \infty$$



A result of Ledoit and Péché (2011)

The main step to derive the limiting risk, is proving the following theorem:

Theorem (Ledoit and Péché (2011))

Assume that $m_{\widehat{\Sigma}}(z) \rightarrow m(z)$ and let $v(z)$ be the companion transform for $m(z)$. We have

$$\frac{1}{d} \operatorname{tr} \left(\Sigma \left(\widehat{\Sigma} + \lambda I_{d \times d} \right)^{-1} \right) \rightarrow_{a.s.} \frac{1}{\gamma} \left(\frac{1}{\lambda v(-\lambda)} - 1 \right) \quad \text{as } n, d \rightarrow \infty$$

Statements like this are nontrivial. It is clear that these quantities converge, but there is no general theory to tell us what the limit is.



A result of Ledoit and Péché (2011)

The main step to derive the limiting risk, is proving the following theorem:

Theorem (Ledoit and Péché (2011))

Assume that $m_{\widehat{\Sigma}}(z) \rightarrow m(z)$ and let $v(z)$ be the companion transform for $m(z)$. We have

$$\frac{1}{d} \operatorname{tr} \left(\Sigma \left(\widehat{\Sigma} + \lambda I_{d \times d} \right)^{-1} \right) \rightarrow_{a.s.} \frac{1}{\gamma} \left(\frac{1}{\lambda v(-\lambda)} - 1 \right) \quad \text{as } n, d \rightarrow \infty$$

Statements like this are nontrivial. It is clear that these quantities converge, but there is no general theory to tell us what the limit is.

Before we prove it, let's use it to derive the asymptotic risk.



Limiting Risk of Ridge Regression

- **First Trace:** Ledoit and P ech e (2011) proves that

$$\frac{1}{d} \text{tr} \left(\Sigma \left(\widehat{\Sigma} + \lambda I_{d \times d} \right)^{-1} \right) \rightarrow_{a.s.} \frac{1}{\gamma} \left(\frac{1}{\lambda v(-\lambda)} - 1 \right) := \kappa(\lambda)$$



Limiting Risk of Ridge Regression

- **First Trace:** Ledoit and P  ch   (2011) proves that

$$\frac{1}{d} \operatorname{tr} \left(\Sigma \left(\widehat{\Sigma} + \lambda I_{d \times d} \right)^{-1} \right) \rightarrow_{a.s.} \frac{1}{\gamma} \left(\frac{1}{\lambda v(-\lambda)} - 1 \right) := \kappa(\lambda)$$

- **Second Functional:**

$$\left(\lambda \alpha^2 - \gamma_d \right) \frac{\lambda}{d} \operatorname{tr} \left(\Sigma \left(\widehat{\Sigma} + \lambda I_{d \times d} \right)^{-2} \right)$$



Limiting Risk of Ridge Regression

- **First Trace:** Ledoit and P  ch   (2011) proves that

$$\frac{1}{d} \text{tr} \left(\Sigma \left(\widehat{\Sigma} + \lambda I_{d \times d} \right)^{-1} \right) \rightarrow_{a.s.} \frac{1}{\gamma} \left(\frac{1}{\lambda v(-\lambda)} - 1 \right) := \kappa(\lambda)$$

- **Second Functional:**

$$\left(\lambda \alpha^2 - \gamma_d \right) \frac{\lambda}{d} \text{tr} \left(\Sigma \left(\widehat{\Sigma} + \lambda I_{d \times d} \right)^{-2} \right) \rightarrow_{a.s.} - \left(\lambda \alpha^2 - \gamma_d \right) \lambda \kappa'(\lambda)$$



Limiting Risk of Ridge Regression

- **First Trace:** Ledoit and Péché (2011) proves that

$$\frac{1}{d} \operatorname{tr} \left(\Sigma \left(\widehat{\Sigma} + \lambda I_{d \times d} \right)^{-1} \right) \rightarrow_{a.s.} \frac{1}{\gamma} \left(\frac{1}{\lambda v(-\lambda)} - 1 \right) := \kappa(\lambda)$$

- **Second Functional:**

$$\left(\lambda \alpha^2 - \gamma_d \right) \frac{\lambda}{d} \operatorname{tr} \left(\Sigma \left(\widehat{\Sigma} + \lambda I_{d \times d} \right)^{-2} \right) \rightarrow_{a.s.} - \left(\lambda \alpha^2 - \gamma_d \right) \lambda \kappa'(\lambda)$$

- After simplification, the limiting risk converges almost surely

$$r_\lambda(X) \rightarrow_{a.s.} \frac{1}{\lambda v(-\lambda)} \left\{ 1 + \left(\frac{\lambda \alpha^2}{\gamma} - 1 \right) \left(1 - \frac{\lambda v^\top(-\lambda)}{v(-\lambda)} \right) \right\}.$$



Properties of the Solution

- From the formula, the optimal ridge parameter λ is

$$\lambda^* = \gamma\alpha^{-2},$$

and we have

$$r_{\lambda^*}(X) \xrightarrow{a.s.} \frac{1}{\lambda^*v(-\lambda^*)}.$$



Properties of the Solution

- From the formula, the optimal ridge parameter λ is

$$\lambda^* = \gamma\alpha^{-2},$$

and we have

$$r_{\lambda^*}(X) \xrightarrow{a.s.} \frac{1}{\lambda^*v(-\lambda^*)}.$$

- With $\Sigma = I_{d \times d}$, by the Marchenko-Pastur theorem we have

$$r_{\lambda}(X) \rightarrow 1 + \gamma m_I(-\lambda; \gamma) + \lambda \left(\lambda\alpha^2 - \gamma \right) m_I^{\top}(-\lambda; \gamma)$$

where $m_I(\cdot; \gamma)$ is the Stieltjes transform of the Marchenko Pastur law.



Properties of the Solution

- From the formula, the optimal ridge parameter λ is

$$\lambda^* = \gamma\alpha^{-2},$$

and we have

$$r_{\lambda^*}(X) \xrightarrow{a.s.} \frac{1}{\lambda^*v(-\lambda^*)}.$$

- With $\Sigma = I_{d \times d}$, by the Marchenko-Pastur theorem we have

$$r_{\lambda}(X) \rightarrow 1 + \gamma m_I(-\lambda; \gamma) + \lambda \left(\lambda\alpha^2 - \gamma \right) m_I^{\top}(-\lambda; \gamma)$$

where $m_I(\cdot; \gamma)$ is the Stieltjes transform of the Marchenko Pastur law.

- In this case, the optimal risk can be written in a closed form.



Silverstein and Choi (1995) show that $v(z)$ is the unique solution with positive imaginary part of the Silverstein equation:

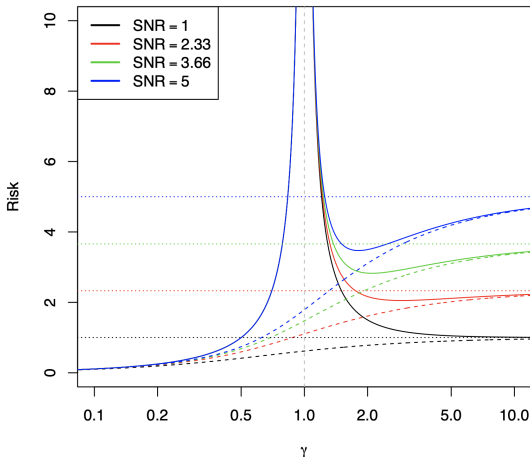
$$-\frac{1}{v(z)} = z - \gamma \int \frac{tdH(t)}{1 + tv(z)}, z \in \mathbb{C}^+.$$

This equation can be solved by a fixed-point algorithm to compute $v(z)$ for all $z \in \mathbb{C}^+$.



Plots!

Assume $\Sigma = I$. This is the risk plot as a function of γ for ridgeless $\lambda \rightarrow 0$ (dashed) and optimal ridge (solid), for different SNRs.





Proof of Ledoit and Péché (2011)

[We will use bold symbols for matrices from now on]

- Remember that $\mathbf{X} \in \mathbb{R}^{n \times d}$ and define $\hat{\Sigma} = \frac{1}{n} \mathbf{X}^\top \mathbf{X}$. Let $\mathbf{G} = (\lambda \mathbf{I} + \hat{\Sigma})^{-1}$ be the resolvent.



Proof of Ledoit and Péché (2011)

[We will use bold symbols for matrices from now on]

- Remember that $\mathbf{X} \in \mathbb{R}^{n \times d}$ and define $\hat{\Sigma} = \frac{1}{n} \mathbf{X}^\top \mathbf{X}$. Let $\mathbf{G} = (\lambda \mathbf{I} + \hat{\Sigma})^{-1}$ be the resolvent.
- We are interested in computing traces of the form

$$\frac{1}{d} \text{Tr}(\mathbf{G}\Sigma) \quad \text{and} \quad \frac{1}{d} \text{Tr}(\mathbf{G})$$



Proof of Ledoit and Péché (2011)

[We will use bold symbols for matrices from now on]

- Remember that $\mathbf{X} \in \mathbb{R}^{n \times d}$ and define $\hat{\Sigma} = \frac{1}{n} \mathbf{X}^\top \mathbf{X}$. Let $\mathbf{G} = (\lambda \mathbf{I} + \hat{\Sigma})^{-1}$ be the resolvent.
- We are interested in computing traces of the form

$$\frac{1}{d} \text{Tr}(\mathbf{G}\Sigma) \quad \text{and} \quad \frac{1}{d} \text{Tr}(\mathbf{G})$$

- There are many ways to do it; e.g., leave-one-out analysis, free probability, etc. Here, we use a method based on Stein's formula.



Let's first remind the Stein's formula. Hong Hu talked about it briefly last time:

Lemma (Stein's Formula)

Let $X \sim \mathcal{N}(0, \Sigma)$ and $f : \mathbb{R}^d \rightarrow \mathbb{R}$ have gradient with at most polynomial growth at infinity. Then for all $i_0 = 1, \dots, d$:

$$\mathbb{E} X_{i_0} f(X_1, \dots, X_d) = \sum_{k=1}^d \Sigma_{i_0 k} \mathbb{E} (\partial_k f)(X_1, \dots, X_d)$$



Matrix version of Stein's formula

Theorem

For $\mathbf{F} : \mathbb{R}^d \rightarrow \mathbb{R}^{d \times d}$ a good behaving matrix-valued function, we have

$$\mathbb{E} [X^\top \mathbf{F}(X) X] = \mathbb{E} [\text{Tr} \Sigma \mathbf{F}(X)] + \mathbb{E} \sum_{k=1}^d (\Sigma (\partial_k \mathbf{F})(X) X)_k$$



Matrix version of Stein's formula

Theorem

For $\mathbf{F} : \mathbb{R}^d \rightarrow \mathbb{R}^{d \times d}$ a good behaving matrix-valued function, we have

$$\mathbb{E} [X^\top \mathbf{F}(X) X] = \mathbb{E} [\text{Tr} \Sigma \mathbf{F}(X)] + \mathbb{E} \sum_{k=1}^d (\Sigma (\partial_k \mathbf{F})(X) X)_k$$

Proof.

$$\begin{aligned} \mathbb{E} X^\top \mathbf{F}(X) X &= \mathbb{E} \sum_{ij} X_i X_j F(X)_{ij} \\ &= \mathbb{E} \sum_{ijk} \Sigma_{ik} \frac{\partial}{\partial X_k} X_j F(X)_{ij} \\ &= \mathbb{E} \sum_{ijk} \Sigma_{ik} [\delta_{j=k} F(X)_{ij} + X_j (\partial_k \mathbf{F})(X)_{ij}] \\ &= \mathbb{E} [\text{Tr} \Sigma \mathbf{F}(X)] + \mathbb{E} \sum_{ijk} \Sigma_{ik} X_j (\partial_k \mathbf{F})(X)_{ij} \\ &= \mathbb{E} [\text{Tr} \Sigma \mathbf{F}(X)] + \mathbb{E} \sum_k [\Sigma (\partial_k \mathbf{F})(X) X]_k \end{aligned}$$



Proof of Ledoit and Péché (2011)

- We want to write $\frac{1}{n}\text{Tr}(\mathbf{G}\Sigma)$ in terms of $\frac{1}{n}\text{Tr}(\mathbf{G})$.



Proof of Ledoit and Péché (2011)

- We want to write $\frac{1}{n}\text{Tr}(\mathbf{G}\Sigma)$ in terms of $\frac{1}{n}\text{Tr}(\mathbf{G})$.
- Note that we have

$$\frac{1}{n}\text{Tr}(\mathbf{G}\hat{\Sigma}) = \frac{1}{n}\text{Tr}\left(\mathbf{G}(\hat{\Sigma} + \lambda I - \lambda I)\right) = \gamma - \frac{\lambda}{n}\text{Tr}(\mathbf{G}).$$



Proof of Ledoit and P ech e (2011)

- We want to write $\frac{1}{n}\text{Tr}(\mathbf{G}\Sigma)$ in terms of $\frac{1}{n}\text{Tr}(\mathbf{G})$.
- Note that we have

$$\frac{1}{n}\text{Tr}(\mathbf{G}\hat{\Sigma}) = \frac{1}{n}\text{Tr}\left(\mathbf{G}(\hat{\Sigma} + \lambda I - \lambda I)\right) = \gamma - \frac{\lambda}{n}\text{Tr}(\mathbf{G}).$$

- We will start with $\frac{1}{n}\text{Tr}(\mathbf{G}\hat{\Sigma})$.



Proof of Ledoit and Péché (2011)

$$\mathbb{E} [X^\top \mathbf{F}(X)X] = \mathbb{E} [\text{Tr} \Sigma \mathbf{F}(X)] + \mathbb{E} \sum_{k=1}^d (\Sigma (\partial_k \mathbf{F})(X)X)_k$$



Proof of Ledoit and P ech e (2011)

$$\mathbb{E} [X^\top \mathbf{F}(X)X] = \mathbb{E} [\text{Tr} \boldsymbol{\Sigma} \mathbf{F}(X)] + \mathbb{E} \sum_{k=1}^d (\boldsymbol{\Sigma} (\partial_k \mathbf{F})(X)X)_k$$

$$\begin{aligned} \mathbb{E} \text{Tr} \mathbf{G} \hat{\boldsymbol{\Sigma}} &= \mathbb{E} \left[\frac{1}{n} \sum_{i=1}^n \text{Tr} \mathbf{G} X(i) X(i)^\top \right] \\ &= \frac{1}{n} \sum_{i=1}^n \mathbb{E} [X(i)^\top \mathbf{G} X(i)] \\ &= \frac{1}{n} \sum_{i=1}^n \mathbb{E} [\text{Tr} \mathbf{G} \boldsymbol{\Sigma}] + \frac{1}{n} \sum_{i=1}^n \sum_{k=1}^d \mathbb{E} \left[e_k^\top \boldsymbol{\Sigma} \left(\frac{\partial}{\partial X(i)_k} \mathbf{G} \right) X(i) \right] \end{aligned}$$



Proof of Ledoit and Péché (2011)

$$\mathbb{E} [X^\top \mathbf{F}(X)X] = \mathbb{E} [\text{Tr} \boldsymbol{\Sigma} \mathbf{F}(X)] + \mathbb{E} \sum_{k=1}^d (\boldsymbol{\Sigma} (\partial_k \mathbf{F})(X)X)_k$$

$$\begin{aligned} \mathbb{E} \text{Tr} \mathbf{G} \hat{\boldsymbol{\Sigma}} &= \mathbb{E} \left[\frac{1}{n} \sum_{i=1}^n \text{Tr} \mathbf{G} X(i) X(i)^\top \right] \\ &= \frac{1}{n} \sum_{i=1}^n \mathbb{E} [X(i)^\top \mathbf{G} X(i)] \\ &= \frac{1}{n} \sum_{i=1}^n \mathbb{E} [\text{Tr} \mathbf{G} \boldsymbol{\Sigma}] + \frac{1}{n} \sum_{i=1}^n \sum_{k=1}^d \mathbb{E} \left[e_k^\top \boldsymbol{\Sigma} \left(\frac{\partial}{\partial X(i)_k} \mathbf{G} \right) X(i) \right] \end{aligned}$$

What is $\frac{\partial \mathbf{G}}{\partial X(i)_k} = \frac{\partial}{\partial X(i)_k} [(\lambda \mathbf{I} + \hat{\boldsymbol{\Sigma}})^{-1}]$?

$$\frac{\partial \mathbf{G}}{\partial X(i)_k} = \frac{1}{n} \left[\mathbf{G} \left(e_k X(i)^\top + X(i) e_k^\top \right) \mathbf{G} \right]$$



Proof of Ledoit and Péché (2011)

$$\begin{aligned} & \mathbb{E} \operatorname{Tr} \left[\mathbf{G} \hat{\Sigma} \right] \\ &= \frac{1}{n} \sum_{i=1}^n \mathbb{E} [\operatorname{Tr} \mathbf{G} \Sigma] + \frac{1}{n^2} \sum_{i=1}^n \sum_{k=1}^d \mathbb{E} [e_k^\top \Sigma (\mathbf{G} (e_k X(i)^\top + X(i) e_k^\top) \mathbf{G}) X(i)] \\ &= \mathbb{E} [\operatorname{Tr} \mathbf{G} \Sigma] + \frac{1}{n^2} \sum_{i=1}^n \sum_{k=1}^d \mathbb{E} [e_k^\top \Sigma \mathbf{G} e_k X(i)^\top \mathbf{G} X(i) + e_k^\top \Sigma \mathbf{G} X(i) e_k^\top \mathbf{G} X(i)] \\ &= \mathbb{E} [\operatorname{Tr} \mathbf{G} \Sigma] + \frac{1}{n^2} \sum_{i=1}^n \mathbb{E} (\operatorname{Tr}(\Sigma \mathbf{G}) X(i)^\top \mathbf{G} X(i) + X(i)^\top \mathbf{G} \Sigma \mathbf{G} X(i)) \\ &= \mathbb{E} [\operatorname{Tr} \mathbf{G} \Sigma] + \frac{1}{n} \mathbb{E} \left[\operatorname{Tr}(\Sigma \mathbf{G}) \operatorname{Tr}(\mathbf{G} \hat{\Sigma}) \right] + \frac{1}{n} \mathbb{E} \left[\operatorname{Tr}(\mathbf{G} \Sigma \mathbf{G} \hat{\Sigma}) \right] \end{aligned}$$



Proof of Ledoit and Péché (2011)

$$\begin{aligned}
& \mathbb{E} \operatorname{Tr} \left[\mathbf{G} \hat{\Sigma} \right] \\
&= \frac{1}{n} \sum_{i=1}^n \mathbb{E} [\operatorname{Tr} \mathbf{G} \Sigma] + \frac{1}{n^2} \sum_{i=1}^n \sum_{k=1}^d \mathbb{E} [e_k^\top \Sigma (\mathbf{G} (e_k X(i)^\top + X(i) e_k^\top) \mathbf{G}) X(i)] \\
&= \mathbb{E} [\operatorname{Tr} \mathbf{G} \Sigma] + \frac{1}{n^2} \sum_{i=1}^n \sum_{k=1}^d \mathbb{E} [e_k^\top \Sigma \mathbf{G} e_k X(i)^\top \mathbf{G} X(i) + e_k^\top \Sigma \mathbf{G} X(i) e_k^\top \mathbf{G} X(i)] \\
&= \mathbb{E} [\operatorname{Tr} \mathbf{G} \Sigma] + \frac{1}{n^2} \sum_{i=1}^n \mathbb{E} (\operatorname{Tr}(\Sigma \mathbf{G}) X(i)^\top \mathbf{G} X(i) + X(i)^\top \mathbf{G} \Sigma \mathbf{G} X(i)) \\
&= \mathbb{E} [\operatorname{Tr} \mathbf{G} \Sigma] + \frac{1}{n} \mathbb{E} \left[\operatorname{Tr}(\Sigma \mathbf{G}) \operatorname{Tr}(\mathbf{G} \hat{\Sigma}) \right] + \frac{1}{n} \mathbb{E} \left[\operatorname{Tr}(\mathbf{G} \Sigma \mathbf{G} \hat{\Sigma}) \right]
\end{aligned}$$

Dividing both sides by n and by Gaussian Lipschitz concentration inequality, we have

$$\frac{1}{n} \operatorname{Tr} (\mathbf{G} \hat{\Sigma}) = \frac{1}{n} \operatorname{Tr}(\mathbf{G} \Sigma) + \frac{1}{n} \operatorname{Tr}(\Sigma \mathbf{G}) \cdot \frac{1}{n} \operatorname{Tr}(\mathbf{G} \hat{\Sigma}) + o(1).$$



Proof of Ledoit and P ech e (2011)

We proved that

$$\frac{1}{n} \text{Tr}(\mathbf{G}\hat{\Sigma}) = \frac{1}{n} \text{Tr}(\mathbf{G}\Sigma) + \frac{1}{n} \text{Tr}(\Sigma\mathbf{G}) \cdot \frac{1}{n} \text{Tr}(\mathbf{G}\hat{\Sigma}) + o(1).$$



Proof of Ledoit and Péché (2011)

We proved that

$$\frac{1}{n} \text{Tr}(\mathbf{G}\hat{\Sigma}) = \frac{1}{n} \text{Tr}(\mathbf{G}\Sigma) + \frac{1}{n} \text{Tr}(\Sigma\mathbf{G}) \cdot \frac{1}{n} \text{Tr}(\mathbf{G}\hat{\Sigma}) + o(1).$$

This implies that

$$\frac{1}{n} \text{Tr}(\mathbf{G}\Sigma) \approx \frac{\frac{1}{n} \text{Tr}(\mathbf{G}\hat{\Sigma})}{1 + \frac{1}{n} \text{Tr}(\mathbf{G}\hat{\Sigma})} + o(1).$$



Proof of Ledoit and Péché (2011)

We proved that

$$\frac{1}{n} \text{Tr}(\mathbf{G}\hat{\Sigma}) = \frac{1}{n} \text{Tr}(\mathbf{G}\Sigma) + \frac{1}{n} \text{Tr}(\Sigma\mathbf{G}) \cdot \frac{1}{n} \text{Tr}(\mathbf{G}\hat{\Sigma}) + o(1).$$

This implies that

$$\frac{1}{n} \text{Tr}(\mathbf{G}\Sigma) \approx \frac{\frac{1}{n} \text{Tr}(\mathbf{G}\hat{\Sigma})}{1 + \frac{1}{n} \text{Tr}(\mathbf{G}\hat{\Sigma})} + o(1).$$

In other words, we have

$$\begin{aligned} \frac{1}{n} \text{Tr}((\lambda\mathbf{I} + \hat{\Sigma})^{-1}\Sigma) &\approx \frac{\gamma - \frac{\lambda\gamma}{d} \text{Tr}((\lambda\mathbf{I} + \hat{\Sigma})^{-1})}{1 + \gamma - \frac{\lambda\gamma}{d} \text{Tr}((\lambda\mathbf{I} + \hat{\Sigma})^{-1})} + o(1) \\ &= \frac{1}{\gamma} \left(\frac{1}{\lambda v_{\hat{\Sigma}}(-\lambda)} - 1 \right) + o(1) \rightarrow \frac{1}{\gamma} \left(\frac{1}{\lambda v_{\Sigma}(-\lambda)} - 1 \right) \quad \square \end{aligned}$$

More General Results



- Eigenvalue Decomposition: $\Sigma \rightsquigarrow (s_1, v_1), \dots, (s_d, v_d)$.
- Assume that

$$\widehat{H}_n(s) := \frac{1}{d} \sum_{i=1}^d \mathbf{1}_{\{s \geq s_i\}} \rightarrow H(s)$$

$$\widehat{G}_n(s) = \frac{1}{\|\beta\|_2^2} \sum_{i=1}^d \langle \beta, v_i \rangle^2 \mathbf{1}_{\{s \geq s_i\}} \rightarrow G(s)$$

- Hastie, Montanari, Rosset, and Tibshirani (2020) derive the asymptotic risk in this case.

References



- **Dobriban and Wager (2016)**
High-Dimensional Asymptotics of Prediction: Ridge Regression and Classification
Annals of Statistics.
- **Hastie, Montanari, Rosset, and Tibshirani (2020)**
Surprises in High-Dimensional Ridgeless Least Squares Interpolation
Annals of Statistics.
- **Ledoit and Péché (2011)**
Eigenvectors of some large sample covariance matrix ensembles
Probability Theory and Related Fields.
- **Benaych-Georges (2022)**
A short proof of Ledoit-Péché's RIE formula for covariance matrices
Technical Report.

Thank You!