

SDS 387 Linear Models

Fall 2024

Lecture 19 - Tue, Nov 5, 2024

Instructor: Prof. Ale Rinaldo

- Progress report extension: now due on Nov 15.
- Recall we are considering the well-specified, fixed-design linear regression setting:

$$Y = \Phi \beta^* + \varepsilon \quad \varepsilon \sim (0, \sigma^2 I_n)$$

$n \times 1$ $n \times d$ $d \times 1$ $n \times 1$

↓
meaning:
has mean 0
and covariance $\sigma^2 I_n$

- Last time we looked at ridge regression estimator:

$$\hat{\beta}_\lambda = \left(\hat{\Sigma} + \lambda I_d \right) \frac{\Phi^T Y}{n}$$

$\hat{\Sigma} = \frac{\Phi \Phi^T}{n}$

$\lambda > 0$
↓
tuning parameter

- Numerically it can be more stable than OLS and, statistically, it provides a different trade-off between bias and variance

- We saw that the ^{expected} excess risk of $\hat{\beta}_\lambda$ is

$$\mathbb{E} \left[\underbrace{R(\hat{\beta}_\lambda)}_{\sigma^2} \right] - R(\beta^*) = \mathbb{E} \left[\|\hat{\beta}_\lambda - \beta^*\|_{\hat{\Sigma}}^2 \right]$$

$$= d^2 \beta^{*\top} \left(\hat{\Sigma} + dI_d \right)^{-2} \hat{\Sigma} \beta^* + \frac{\sigma^2}{n} \text{tr} \left(\hat{\Sigma}^2 \left(\hat{\Sigma} + dI_d \right)^{-2} \right)$$

bias term $\uparrow d$
variance $\downarrow d$

When $d=0$ this reduces to $\sigma^2 \frac{d}{n}$ the risk of OLS.

Question: what is the optimal d ?

Prop 3.8 in Bach's book

Setting $d_{\text{optimal}} = \frac{\sigma \sqrt{\text{tr}(\hat{\Sigma})}}{\|\beta^*\| \sqrt{n}}$ we get that

$$\mathbb{E} \left[R(\hat{\beta}_d) \right] - \sigma^2 \leq \frac{\sigma \sqrt{\text{tr}(\hat{\Sigma})} \|\beta^*\|}{\sqrt{n}}$$

Remark: compare this with $\sigma^2 \frac{d}{n}$ (the risk of OLS)

i) If d is fixed and $n \rightarrow \infty$ then OLS is better because it vanishes at a rate $\mathcal{O}(1/n)$ while the risk of ridge vanishes at slower rate $\mathcal{O}(1/\sqrt{n})$

When d changes (increases with n) and $\hat{\Sigma}$, σ and $\|\beta^*\|$ also change with n , then ridge can be better than OLS.

ii) this is not necessarily the best choice of d . It is the best choice for a simpler upper bound on the risk.

iii) in practice how do you choose d ?

Use cross-validation.

PA/ First, the eigenvalues of $(\hat{\Sigma} + \lambda I_d)^{-2} \lambda \hat{\Sigma} \leq 1/2$.

To see this, the eigenvalues of this matrix are of the form

$$\frac{\hat{d}_i \lambda}{(\hat{d}_i + \lambda)^2} \quad \text{where } \hat{d}_1, \dots, \hat{d}_d \text{ are eigenvalues of } \hat{\Sigma}$$

which is always $\leq 1/2$ because

$$\frac{2b}{(a+b)^2} \leq 1/2 \quad a, b > 0$$

So, using this fact, the bias term is

$$\lambda \beta^{*\top} (\hat{\Sigma} + \lambda I_d)^{-2} \lambda \hat{\Sigma} \beta^* \leq \frac{\lambda}{2} \|\beta^*\|^2 \quad \text{Exercise}$$

As for the variance,

$$\begin{aligned} \frac{\sigma^2}{n} \text{tr} \left(\hat{\Sigma}^2 (\hat{\Sigma} + \lambda I_d)^{-2} \right) &= \frac{\sigma^2}{n \lambda} \text{tr} \left(\lambda \hat{\Sigma} \hat{\Sigma} (\hat{\Sigma} + \lambda I_d)^{-2} \right) \\ &\leq \frac{\sigma^2}{n \lambda} \underbrace{\left\| \hat{\Sigma} \lambda (\hat{\Sigma} + \lambda I_d)^{-2} \right\|_{\text{op}}}_{\leq 1/2} \text{tr}(\hat{\Sigma}) \\ &\leq \frac{\sigma^2}{n \lambda} \text{tr}(\hat{\Sigma}) \quad \text{Exercise} \end{aligned}$$

So, the expected excess risk of $\hat{\beta}_\lambda$ is upper bounded by:

$$\frac{\lambda}{2} \|\beta^*\|^2 + \frac{\sigma^2 \text{tr}(\hat{\Sigma})}{2n\lambda} := a\lambda + \frac{b}{\lambda}$$

This is minimized at $\lambda = \sqrt{b/a}$, with optimal value $\sqrt{2ab}$. $a, b > 0$

LOWER BOUND ON RISK OF OLS

Section 3.7 in
Bach's book

Recall that the expected excess risk of $\hat{\beta}$ (OLS) is $\sigma^2 \frac{d}{n}$

(assuming fixed covariates and a well-specified model).

Here we show that this value is optimal in a minimax sense.

- Suppose we are interested in estimate a parameter θ^* , generally defined as a functional of a probability distribution, say P . To highlight this fact, we write $\theta^*(P)$ [Note: θ^* does not need to identify or fully specify P .] We also specify a collection, say \mathcal{P} , of probability distributions whose parameter θ^* is of interest. We observe data (an iid sequence of length, say, n) and construct an estimator, $\hat{\theta}_n$ of θ^* [we do not know which P in \mathcal{P} has generated the data, so we do not know $\theta^* = \theta^*(P)$].

We measure the quality of $\hat{\theta}_n$ by its risk

$$R(\theta^*(P), \hat{\theta}_n)$$

Example: i) \mathcal{P} : set of distributions of Y , where

$$\text{For each } P \in \mathcal{P} \quad Y = \underbrace{\beta^*}_{\substack{\text{fixed and} \\ \text{known}}} + \underbrace{\varepsilon}_{\text{unknown}} \rightarrow \varepsilon \sim (0, \sigma^2) \quad \downarrow \text{known}$$

$$\theta^*(P) = \beta^*(P)$$

$$\text{iii) } \mathcal{P}_{\text{Gauss}}: Y \sim N_n(\Phi \beta^*, \sigma^2 I_n)$$

$$\text{For } P \in \mathcal{P}_{\text{Gauss}} \quad \theta^*(P) = \beta^*(P) \quad \downarrow \quad \text{unknown } \beta^* \in \mathbb{R}^d$$

Of course $\mathcal{P}_{\text{Gauss}} \subset \mathcal{P}$

Of course for any estimator $\tilde{\beta}$ of β^* the risk is

$$\begin{aligned} R(\tilde{\beta}, \beta^*) &= \mathbb{E}[R(\tilde{\beta})] - \sigma^2 \\ &= \mathbb{E}\left[\|\tilde{\beta} - \beta^*\|_{\Sigma}^2\right] \end{aligned}$$

Remark think of the risk $R(\theta^*, \hat{\theta}_n)$ as a function of θ^* .

- How do we use this setting to evaluate whether an estimator is good or optimal?
- The minimax approach requires you to evaluate the minimax risk and find an estimator that at least asymptotically, achieves this risk value.
- The minimax risk is

$$\inf_{\hat{\theta}_n} \sup_{P \in \mathcal{P}} R(\hat{\theta}_n, \theta^*(P))$$

where $\inf_{\hat{\theta}_n}$ is the infimum over all estimators (functions of the data)

Minimax risk measures intrinsic statistical hardness of an estimation task (or any statistical task)

For regression, this translates into

$$\inf_{\hat{\beta}} \sup_{P \in \mathcal{P}} \left(\mathbb{E}_P [R(\hat{\beta})] - \sigma^2 \right)$$

- An estimator $\hat{\theta}_n$ is **minimax rate-optimal** if asymptotically as $n \rightarrow \infty$, its risk is of the same order as the minimax risk. Formally, let

$R_{n, \text{minimax}}$ the value of the minimax risk and for an estimator $\hat{\theta}_n$, let $R_n^{\text{sup}}(\hat{\theta}_n) \geq \sup_{\theta^*} R(\hat{\theta}_n, \theta^*)$

Then $\hat{\theta}_n$ is minimax rate optimal if

$$\limsup_n \frac{R_n^{\text{sup}}(\hat{\theta}_n)}{R_{n, \text{minimax}}} \leq c$$

$\underbrace{\hspace{10em}}_{\geq 1}$

- $\hat{\theta}_n$ is sharp minimax rate optimal if $c = 1$

- $\hat{\theta}_n$ is exact minimax optimal if $R_n^{\text{sup}}(\hat{\theta}_n) = R_{n, \text{minimax}}$ all n .

Thm $\hat{\beta}$ (OLS) is exact minimax optimal for $\mathcal{P}_{\text{Gauss}}$

$\hookrightarrow \hat{\beta}$ (OLS) is minimax optimal for \mathcal{P}

- For regression, we are interested in computing a lower bound on this quantity:

$$\inf_A \sup_{\beta^* \in \mathbb{R}^d} \mathbb{E}_{\substack{\varepsilon \in \mathbb{R}^d: \\ \varepsilon \sim \mathcal{N}(0, \sigma^2 \mathbf{I}_n)}} \left[R(\mathcal{A}(\Phi \beta^* + \varepsilon)) \right] - \sigma^2$$

algorithm
taking as input

Y and Φ
→ fixed

$$\geq \inf_A \sup_{\beta^* \in \mathbb{R}^d} \mathbb{E}_{\varepsilon \sim \mathcal{N}_n(0, \sigma^2 \mathbf{I}_n)} \left[R(\mathcal{A}(\Phi \beta^* + \varepsilon)) \right] - \sigma^2$$