

SDS 387 Linear Models

Fall 2024

Lecture 23 - Tue, Nov 19, 2024

Instructor: Prof. Ale Rinaldo

- Last time: linear regression with random design and well-specified model
 $(\Phi_1, Y_1), \dots, (\Phi_n, Y_n) \stackrel{iid}{\sim} P_{\Phi, Y}$

$$Y_i = \Phi_n^T \beta^* + \varepsilon_i$$

$\Sigma_n = \mathbb{E}[\Phi_n \Phi_n^T]$ assumed invertible $\hookrightarrow \varepsilon_1, \dots, \varepsilon_n \mid \Phi_1, \dots, \Phi_n \stackrel{iid}{\sim} (0, \sigma^2)$

- Exercise: what happens when the linear model is not true?
- We saw that the expected excess risk of $\hat{\beta}$ (OLS)

is

$$\sigma^2 \mathbb{E}_{\Phi_1, \dots, \Phi_n} \left[\text{tr} \left(\Sigma_n^{-1} \hat{\Sigma}_n^{-1} \right) \right]$$

and

$$\hat{\Sigma}_n = \frac{1}{n} \sum_{i=1}^n \Phi_i \Phi_i^T$$

which is assumed to be invertible with prob. 1.

- Question: is this optimal?

Theorem 1 of Mourtada (AOS 2022, 20(4), 2157-2178)

i) Assume that either $d > n$ or the distribution of the Φ_i 's is degenerate (supported on an affine subspace of \mathbb{R}^d)
Then the minimax risk is infinity!

ii) If $n \geq d$ and the distribution of the Φ_i 's is not degenerate ($\hat{\Sigma}$ is invertible w.p. 1) then

OLS estimator is minimax optimal

$$\inf_{\tilde{\beta}} \sup_{\beta^*} \mathbb{E} [R(\tilde{\beta})] = \frac{\sigma^2}{n} \mathbb{E} [\text{tr}(\Sigma \hat{\Sigma}^{-1})]$$

all estimators of β^* excess risk

Above, the expectation is w.r.t to the distribution of Φ_1, \dots, Φ_n and $\varepsilon_1, \dots, \varepsilon_n$

~~PP~~ Very similar to the proof of minimax optimality of OLS in fixed design setting.

$$\inf_{\tilde{\beta}} \sup_{\beta^*} \mathbb{E}_{\substack{\Phi_1, \dots, \Phi_n \\ \varepsilon_1, \dots, \varepsilon_n}} [R(\tilde{\beta})] \geq \inf_{\tilde{\beta}} \sup_{\beta^*} \mathbb{E}_{\substack{\Phi_1, \dots, \Phi_n \\ \varepsilon_1, \dots, \varepsilon_n \mid \Phi_1, \dots, \Phi_n \stackrel{i.i.d.}{\sim} N(0, \sigma^2 \mathbf{I}_d)}} [R(\tilde{\beta})]$$

We replace \sup_{β^*} with an average, assuming a prior

distribution for $\beta^* \sim N_d(0, \frac{\sigma^2}{n\lambda} \mathbf{I}_d)$ some $\lambda > 0$ to be chosen later

Using the same calculations in the fixed design setting, the minimax risk is lower bounded by the Bayes risk

$$\frac{\sigma^2}{n} \mathbb{E}_{\Phi_1, \dots, \Phi_n} \left[\text{tr} \left(\left(\hat{\Sigma} + d \text{Id} \right)^{-1} \Sigma_1^* \right) \right] \quad \forall d > 0$$

Case 1): $\hat{\Sigma}$ is not invertible with prob 1 (degenerate case)

Then

largest eigenvalue

eg $d > n$ or the distribution of the Φ_i 's is supported on a affine subspace

$$\text{tr} \left(\Sigma_1^{*1/2} \left(\hat{\Sigma} + d \text{Id} \right)^{-1} \Sigma_1^{*1/2} \right) = d_{\max}(\cdot)$$

★

$$= \frac{1}{d_{\min} \left(\Sigma_1^{*1/2} \left(\hat{\Sigma} + d \text{Id} \right) \Sigma_1^{*1/2} \right)}$$

Since $\hat{\Sigma}$ is not invertible w.p. 1, $\exists u \in \mathbb{S}^{d-1}$ (random) possibly

s.t. $u^T \Sigma_1^{*1/2} \left(\hat{\Sigma} + d \text{Id} \right) \Sigma_1^{*1/2} u = 0 + d \| \Sigma_1^{*1/2} u \|^2$

\downarrow
 u is chosen so that $\Sigma_1^{*1/2} u$ is in the null space of $\hat{\Sigma}$

$$\leq d d_{\max} \left(\Sigma_1^{*1/2} \right)$$

$$= \frac{d}{d_{\min} \left(\Sigma_1^* \right)}$$

So $\star \geq \frac{d_{\min} \left(\Sigma_1^* \right)}{d}$

So the minimax risk is lower bounded by

$$\frac{\sigma^2}{n} \frac{\text{dim}(\Sigma_w)}{\lambda}$$

$$\forall \lambda > 0$$

Letting $\lambda \rightarrow 0$ the minimax lower bound $\rightarrow \infty$

in) If $\hat{\Sigma}^c$ is invertible with prob 1 then $\text{tr} \left((\hat{\Sigma}^c + \lambda \text{Id})^{-1} \Sigma^c \right)$ is \downarrow in λ and converges to $\text{tr} \left(\Sigma^{c-1} \Sigma^c \right)$ as $\lambda \rightarrow 0$.

By monotone convergence theorem the minimax lower bound is $\frac{\sigma^2}{n} \mathbb{E} \left[\text{tr} \left(\hat{\Sigma}^c^{-1} \Sigma^c \right) \right]$ ■

Mouradov's Corollary 2: if $n \geq d$ and the distribution of the Φ_i 's is not degenerated you can further lower bound this by

$$\sigma^2 \frac{d}{n-d-1} \quad n > (d+1)$$

If the covariates $\Phi_1, \dots, \Phi_n \stackrel{\text{i.i.d.}}{\sim} N(0, \text{Id})$ this value is the exact minimax risk.

► AN EXACT ANALYSIS OF THE RLSR UNDER GAUSSIAN SETTINGS

Breiman & Freedman (1983) JASA 78 (731)

Assume all the conditions above and further suppose that $\Phi_1, \dots, \Phi_n \stackrel{i.i.d.}{\sim} N(0, I_d)$. Then the risk of \hat{w}

$$\stackrel{15}{=} \frac{\sigma^2}{n} \mathbb{E} \left[\text{tr} \left(\hat{\Sigma}^{-1} \hat{\Sigma} \right) \right] = \sigma^2 \mathbb{E} \left[\left(\sum_{i=1}^n \Phi_i \Phi_i^T \right)^{-1} \right]$$

Wishart distribution with parameter I_d and n degrees of freedom

its inverse is called the inverse Wishart

$$= \begin{cases} \frac{\sigma^2 \text{tr}(I_d)}{n-d-1} = \sigma^2 \frac{d}{n-d-1} & \text{if } n > d+1 \\ \infty & \text{if } n = d \text{ or } d+1 \end{cases}$$

At interpolation ($n=d$ or $n=d+1$) the risk explodes.

- Now, under the same settings, something interesting happens when $d > n$. See Belkin, Hsu and Xu (2020)

SIAM Journal of Mathematics
in Data Science

You would expect "bad" risk behavior.

Surprisingly the risk is stable and in fact it can be smaller than when $d \leq n$.

When $d > n$ we will fit the min-norm estimator:

$$\hat{\beta}_{MN} = \Phi^+ Y = (\Phi^T \Phi)^+ \Phi^T Y = \Phi^T (\Phi \Phi^T)^{-1} Y$$

Φ with rows Φ^T $Y = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}$ $n \times n$ invertible

Let's compute the expected excess risk $\mathbb{E} \left[\|\hat{\beta}_{MN} - \beta^*\|_{\mathbb{I}_d}^2 \right]$

The first thing to notice there is a bias term:

$$\begin{aligned} \beta^* - \hat{\beta}_{MN} &= \beta^* - \Phi^T (\Phi \Phi^T)^{-1} \Phi \beta^* - \Phi^T (\Phi \Phi^T)^{-1} \varepsilon \\ &= \underbrace{(\mathbb{I}_d - \Phi^T (\Phi \Phi^T)^{-1} \Phi)}_{\Pi} \beta^* - \Phi^T (\Phi \Phi^T)^{-1} \varepsilon \end{aligned}$$

$\mathbb{I} - \Pi$ orthogonal projection
onto the null space of
 Φ
 $n \times d$

So

$$\begin{aligned} \mathbb{E} \left[\|\beta^* - \hat{\beta}_{MN}\|^2 \right] &= \mathbb{E}_{\Phi} \left[\|\underbrace{(\mathbb{I} - \Pi)}_{\text{random projection matrix}} \beta^*\|^2 \right] + \\ &\quad \mathbb{E}_{\Phi} \left[\text{tr} \left((\Phi \Phi^T)^{-1} \Phi \Phi^T (\Phi \Phi^T)^{-1} \mathbb{E} \left[\underbrace{\varepsilon \varepsilon^T}_{\sigma^2 \mathbb{I}_n} \right] \right) \right] \end{aligned}$$

$$= \mathbb{E}_{\Phi} \left[\|\underbrace{(\mathbb{I} - \Pi)}_{\beta^*}\|^2 \right] + \mathbb{E}_{\Phi} \left[(\Phi \Phi^T)^{-1} \right]$$

$$= T_1 + T_2$$

(6)

Next,

$$T_1 = \|\beta^*\|^2 - \mathbb{E} \left[\|\Pi \beta^*\|^2 \right]$$

To compute $\mathbb{E} \left[\|\Pi \beta^*\|^2 \right]$ we will use the fact that if $Z \sim N(0, I)$ then $UZ \sim N(0, I)$ where U is an orthogonal matrix. Let U_1, U_2, \dots, U_d be $d \times d$ orthogonal matrices s.t. $U_i \beta^* = \|\beta^*\| e_i$

So, $\forall i=1, \dots, d$

\downarrow i th standard basis vector

$$\begin{aligned} \|\Pi \beta^*\|^2 &= \beta^{*\top} \Phi^\top (\Phi \Phi^\top)^{-1} \Phi \beta^* \stackrel{d}{=} \beta^{*\top} U_i^\top \Phi^\top (\Phi U_i U_i^\top \Phi^\top)^{-1} \Phi U_i \beta^* \\ &= \|\beta^*\|^2 e_i^\top \Phi^\top (\Phi \Phi^\top)^{-1} \Phi e_i \\ &= \|\beta^*\|^2 \text{tr} \left(\Phi^\top (\Phi \Phi^\top)^{-1} \Phi e_i e_i^\top \right) \end{aligned}$$

So

$$\begin{aligned} \mathbb{E} \left[\|\Pi \beta^*\|^2 \right] &= \frac{1}{d} \sum_{i=1}^d \mathbb{E} \left[\|\beta^*\|^2 \text{tr} \left(\Phi^\top (\Phi \Phi^\top)^{-1} \Phi e_i e_i^\top \right) \right] \\ &= \|\beta^*\|^2 \frac{1}{d} \mathbb{E} \left[\text{tr} \left(\Phi^\top (\Phi \Phi^\top)^{-1} \Phi \underbrace{\sum_{i=1}^d e_i e_i^\top}_{I_d} \right) \right] \\ &= \|\beta^*\|^2 \frac{1}{d} \mathbb{E} \left[\text{tr} \left(\underbrace{\Phi \Phi^\top}_{I_n} \Phi^\top (\Phi \Phi^\top)^{-1} \Phi \right) \right] \\ &= \|\beta^*\|^2 \frac{n}{d} \end{aligned}$$