SDS 387
Linear Models

Fall 2024

Lecture 20 - Tue, Nov 7, 2024

Instructor: Prof. Ale Rinaldo

📧 Last time: minimax lower bound for OLS $\hat{\beta}$ assuming
a well specified model and deterministic and full-rank
design matrix $\underset{n \times d}{\Phi}$ :

$$Y = \Phi \beta^* + \varepsilon$$

$\longrightarrow$ $n$-dimensional
vector with mean 0
and variance matrix
$\sigma^2 I_n$

We know that, in this setting,

$$\sup_{\beta^*} \underset{\substack{Y = \Phi \beta^* + \varepsilon \\ \varepsilon \sim (0, \sigma^2 I_n)}}{\mathbb{E}} \left[ R(\hat{\beta}) \right] - \sigma^2 = \sigma^2 \frac{d}{n}$$

Now want to establish a lower bound on the largest
possible expected excess risk that holds regardless
of the choice of the estimator.

①

So we are interested in lower-bounding

$$\inf_{A} \sup_{P \in \mathcal{P}} \mathbb{E}_P \left[ R \left( A \left( \underbrace{\Phi \beta^*(P) + \varepsilon}_{Y} \right) \right) \right] - \sigma^2$$

where $\mathcal{P} = \Big\{$ probability distributions for

$$Y = \Phi \beta^* + \varepsilon \qquad \varepsilon \sim (0, \sigma^2 I_n) \Big\}$$
$$\beta^* \in \mathbb{R}^d$$

This expression is, in turn, lower bounded by

$$\inf_{A} \sup_{P \in \mathcal{P}_{Gauss}} \mathbb{E}_P \left[ R \left( A \left( \underbrace{\Phi \beta^*(P) + \varepsilon}_{Y} \right) \right) \right] - \sigma^2$$

where $\mathcal{P}_{Gaussian} = \Big\{ Y \sim N_n \left( \Phi \beta^*, \sigma^2 I_n \right), \ \beta^* \in \mathbb{R}^d \Big\}$

*algorithm that takes as input $Y \in \mathbb{R}^n$ and return an estimator $\tilde{\beta} \in \mathbb{R}^d$*

<u>Remark</u>: $\sigma^2$ is known and $\Phi$ is known and
$\underset{n \times d}{} \qquad$ deterministic

The last expression can be written

$$\inf_{A} \sup_{\beta^* \in \mathbb{R}^d} \mathbb{E}_{\varepsilon \sim N(0, \sigma^2 I_n)} \left[ R \left( A \left( \Phi \beta^* + \varepsilon \right) \right) \right] - \sigma^2$$

only random term

We are going to take a Bayesian approach and lower bound the last expression by

$$\inf_{A} \mathbb{E}_{\beta^* \sim \pi} \mathbb{E}_{\varepsilon \sim N(0, \sigma^2 I_n)} \left[ R \left( A \left( \Phi \beta^* + \varepsilon \right) \right) \right] - \sigma^2$$
$$\hookrightarrow prior$$

②

We can pick any prior $\pi$. We pick a prior $\pi$ that is analytically convenient. Chose as prior $\pi$ the distribution $\beta^* \sim N\left(0, \frac{\sigma^2 I_n}{\lambda n}\right)$ where $\lambda > 0$.

Then $\left(\beta^*, \Phi\beta^* + \varepsilon\right) \in \mathbb{R}^d \times \mathbb{R}^n$ is jointly Gaussian with mean $0$ and variance

$$\frac{\sigma^2}{n\lambda} \begin{bmatrix} I_d & \Phi^T \\ & \\ \Phi & \Phi\Phi^T + n\lambda I_n \end{bmatrix} \begin{array}{c} d \\ \\ n \end{array}$$

$$\phantom{xxxxxxx} d \phantom{xxxxxxxx} n$$

Next, recall that

$$R\left(A(\Phi\beta^* + \varepsilon)\right) - \sigma^2 = \left\| A(\Phi\beta^* + \varepsilon) - \beta^* \right\|_{\hat{\Sigma}}^2$$

where $\hat{\Sigma} = \frac{1}{n}\Phi^T\Phi$.

So the expression becomes

$$\mathbb{E}_{(\beta^*, \, y)} \left[ \left\| A(Y) - \beta^* \right\|_{\hat{\Sigma}}^2 \right] =$$

$$\underset{\uparrow}{\phantom{xx}} \Phi\beta^* + \varepsilon$$

$$= \int_{\mathbb{R}^n} \int_{\mathbb{R}^d} \left\| A(y) - \beta^* \right\|_{\hat{\Sigma}}^2 \, d\,P(\beta^*|y) \, d\,P(y)$$

$$\underbrace{\phantom{d\,P(\beta^*|y)}}$$

posterior of $\beta^*$ given $Y$

$$(3)$$

A standard calculation gives that the posterior is

$$\beta^* | Y \sim N_d \left( \hat{\beta}_\lambda, \frac{\sigma^2}{n} \left( \hat{\Sigma} + \lambda I_d \right)^{-1} \right)$$

where

$$\hat{\beta}_\lambda = \left( \hat{\Sigma} + \lambda I_d \right)^{-1} \frac{\Phi^T}{n} \varepsilon .$$

Next,

$$\int_{\mathbb{R}^d} \| A(Y) - \beta^* \|_{\hat{\Sigma}}^2 \, dP(\beta^* | Y) = \mathop{\mathbb{E}}_{\beta^* | Y} \left[ \| A(Y) - \beta^* \|_{\hat{\Sigma}}^2 \right]$$

$$\geq \inf_A \mathop{\mathbb{E}}_{\beta^* | Y} \left[ \| A(Y) - \beta^* \|_{\hat{\Sigma}}^2 \right]$$

$$= \mathop{\mathbb{E}}_{\beta^* | Y} \left[ \| \hat{\beta}_\lambda - \beta^* \|_{\hat{\Sigma}}^2 \right]$$

because $\mathop{\mathbb{E}}_{\beta^* | Y} \left[ \| A(Y) - \beta^* \|_{\hat{\Sigma}}^2 \right]$ is minimized when

$A(Y) = \hat{\beta}_\lambda$, the posterior mean

Putting everything together, we have found the following lower bound for the minimax risk:

$$\mathop{\mathbb{E}}_{(\beta^*, Y)} \left[ \| \hat{\beta}_\lambda - \beta^* \|_{\hat{\Sigma}}^2 \right] =$$

$$\mathop{\mathbb{E}}_{\beta^* \sim N\left(0, \frac{\sigma^2 I_d}{n \lambda}\right)} \mathop{\mathbb{E}}_{\varepsilon \sim N(0, \sigma^2 I_n)} \left[ \| \left( \Phi^T \Phi + n \lambda I_d \right)^{-1} \Phi^T \left( \Phi \beta^* + \varepsilon \right) - \beta^* \|_{\hat{\Sigma}}^2 \right]$$

$\underbrace{\phantom{xxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxx}}$
☆

④

We have that

$$\cancel{A} = \left(\underline{\Phi}^T \underline{\Phi} + n\lambda Id\right)^{-1} \underline{\Phi}^T \varepsilon - n\lambda \left(\underline{\Phi}^T \underline{\Phi} + n\lambda Id\right)^{-1} \beta^*$$

Because $\beta^* \perp\!\!\!\perp \varepsilon$ the expression is

$$\mathbb{E}_{\varepsilon \sim N(0, \sigma^2 I_n)} \left[ \left\| \left(\hat{\Sigma} + \lambda\right)^{-1} \frac{\underline{\Phi}^T \varepsilon}{n} \right\|^2_{\hat{\Sigma}} \right] + \mathbb{E}_{\beta^* \sim N\left(0, \frac{\sigma^2 Id}{n\lambda}\right)} \left[ \left\| \lambda \left(\hat{\Sigma} + \lambda Id\right)^{-1} \beta^* \right\|^2_{\hat{\Sigma}} \right]$$

$$= \quad T_1 \quad + \quad T_2$$

We have that
$$T_1 = \frac{\sigma^2}{n} \, \mathrm{tr}\left( \left(\hat{\Sigma} + \lambda Id\right)^{-2} \hat{\Sigma}^2 \right)$$

$$T_2 = \lambda^2 \, \mathbb{E}_{\beta^*} \left[ \beta^{*T} \left(\hat{\Sigma} + \lambda Id\right)^{-1} \hat{\Sigma} \left(\hat{\Sigma} + \lambda Id\right)^{-1} \beta^* \right]$$

$$= \frac{\lambda^2 \sigma^2}{n\lambda} \, \mathrm{tr}\left( \left(\hat{\Sigma} + \lambda Id\right)^{-2} \hat{\Sigma} \right) \qquad \text{}$$

$$\hookrightarrow \qquad T_1 + T_2 = \boxed{\frac{\sigma^2}{n} \, \mathrm{tr}\left( \left(\hat{\Sigma} + \lambda\right)^{-1} \hat{\Sigma} \right)} \qquad \lambda > 0$$

$$\sum_{j=1}^{d} \frac{\hat{\lambda}_j}{\hat{\lambda}_j + \lambda}$$
$$\hat{\lambda}_j \quad j\text{-th eigenvalue of } \hat{\Sigma}$$

$\hookrightarrow$ lower bound on the minimax risk.

The above bound holds for any $\lambda > 0$. Of course
$\mathrm{tr}\left( \left(\hat{\Sigma} + \lambda Id\right)^{-1} \hat{\Sigma} \right)$ is $\downarrow$ in $\lambda$.

⑤

So the final lower bound:

$$\sup_{\lambda > 0} \frac{\sigma^2}{n} \, \text{tr}\left( \left(\hat{\Sigma} + \lambda \, I_d\right)^{-1} \hat{\Sigma} \right) =$$

$$\frac{\sigma^2}{n} \lim_{\lambda \downarrow 0} \left(\hat{\Sigma} + \lambda \, I_d\right)^{-1} \hat{\Sigma} =$$

Recall $\hat{\Sigma}$ is invertible by assumption

$$\frac{\sigma^2}{n} \, \text{tr}\left( I_d \right) = \boxed{\sigma^2 \, \frac{d}{n}}$$

↓

expected excess risk of $\hat{\beta}$ (OLS)

So $\hat{\beta}$ (OLS) is the minimax estimator

## ⊟ STATISTICAL INFERENCE FOR $\beta^*$

- As usual the model is

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

↙ full column rank
and deterministic

<u>Goal</u> : statistical inference for $\beta^*$

- Is $\hat{\beta}$ (OLS) <u>consistent</u>, meaning $\hat{\beta} \xrightarrow{P} \beta^*$ ?

↓

when $\hat{\beta}$ is computed using data $Y = \Phi \beta^* + \varepsilon$

⑥

Yes !   To see this :

$$\hat{\beta} = (\hat{\Phi}^T \Phi)^{-1} \Phi^T Y = \beta^* + \hat{\Sigma}^{-1} \frac{\Phi^T \varepsilon}{n}$$

<u>Claim</u> :  if  $\boxed{\hat{\Sigma} = \frac{\Phi^T \Phi}{n} \xrightarrow[d \times d]{} \Sigma}$  then

$$\hat{\Sigma}^{-1} \frac{\Phi^T \varepsilon}{n} \xrightarrow{P} 0 \qquad \text{or} \qquad \hat{\Sigma}^{-1} \frac{\Phi^T \varepsilon}{n} = o_P(1)$$

<u>Pf</u>/

By  WLLN,  $\frac{\Phi^T \varepsilon}{n} \xrightarrow{P} 0$  because

$$\frac{\Phi^T \varepsilon}{n} = \sum_{n=1}^{n} \frac{\Phi_n \varepsilon_n}{n} \qquad \text{when } \Phi_n \text{ is transpose of } n\text{th row of } \Phi$$

and $\qquad \Phi_i \varepsilon_i \sim (0, \sigma^2 \Phi_n \Phi_n^T)$

So $\qquad Var\left[\frac{\Phi^T \varepsilon}{n}\right] = \frac{\sigma^2}{n} \Phi^T \Phi$

So  by  chebyshev  $\frac{\Phi^T \varepsilon}{n} \xrightarrow{P} 0$

Next $\qquad \hat{\Sigma}^{-1} \to \Sigma^{-1}$  so  by  slutsky's theorem

$$\hat{\Sigma}^{-1} \frac{\Phi^T \varepsilon}{n} \xrightarrow{P} \Sigma^{-1} \times 0 = 0$$

Also $\qquad \hat{\beta}$  is asymptotically normal :

$$\sqrt{n} \left(\hat{\beta} - \beta^*\right) \xrightarrow{d} N\left(0, \sigma^2 \Sigma^{-1}\right)$$

$$E[\hat{\sigma}]$$

We have that $\sqrt{n}\left(\hat{\beta}-\beta^{*}\right) = \hat{\Sigma}_{\cdot}^{-1} \sqrt{n} \frac{\Phi^{T}\varepsilon}{n}$

$$\downarrow$$

$$\to \Sigma^{-1}$$

So we only need to show that $\sqrt{n} \frac{\Phi^{T}\varepsilon}{n} \xrightarrow{d} N(0, \sigma^2 \Sigma)$

and the result will follow by Slutsky

- Next time we will going to verify this using the
  LF conditions, which holds if

$$\max_{i} \frac{\|\Phi_{i}\|}{\sqrt{n}} \to 0$$

$$\boxed{8}$$