- Today finish the proof of minimax optimality of the OLS estimator when is the model is linear and the covariates are fixed. That is, our data consist of

$$Y = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix} \quad \text{s.t.}$$

True unknown coeff.

$$Y = \Phi \beta^* + \varepsilon \qquad \varepsilon_1 \dots \varepsilon_n \overset{iid}{\sim} (0, \sigma^2)$$

$\downarrow$ $n \times d$ fixed matrix

known!

- <u>Remark</u> : the extension to the random design case can be found in Mourtada's paper. The result is the same the OLS estimator is minimax estimator of the model is linear.

- Last time we established the following lower bound

$$\inf_{A} \quad \sup_{\beta \in \mathbb{R}^d} \quad \underset{\varepsilon = \begin{bmatrix} \varepsilon_1 \\ \vdots \\ \varepsilon_n \end{bmatrix}}{\mathbb{E}} \left[ R \left( A(\Phi\beta + \varepsilon) \right) \right] - \sigma^2$$

$$\underset{\text{iid } \sim (0, \sigma^2)}{}$$

$$\underbrace{\qquad\qquad\qquad\qquad\qquad}_{\text{excess risk of } A(\cdot)}$$

$$\geq \quad \underset{(\beta, Y)}{\mathbb{E}} \left[ \| \hat{\beta}_\lambda - \beta \|^2_{\hat{\Sigma}} \right] \qquad\qquad (\cancel{A})$$

ridge estimator ← (pointing to $\hat{\beta}_\lambda$)

$$\hookrightarrow \frac{\hat{\Phi^T \Phi}}{n}$$

$$\text{where} \qquad \beta \sim N_d \left( 0, \frac{\sigma^2}{n\lambda} I_d \right) \quad \text{prior}$$

$$Y | \beta \sim N_n \left( \Phi\beta, \sigma^2 I_n \right)$$

This is a standard argument: lower bound the maximal possible risk of any estimator $A(\cdot)$ over all $\beta \in \mathbb{R}^d$ by the average risk of $A(\cdot)$ with respect to a carefully chosen distribution for $\beta$ (a prior).

- Next, we have that $(\cancel{A})$ is equal

$$\underset{\beta \sim N_d \left( 0, \frac{\sigma^2}{n\lambda} I_d \right)}{\mathbb{E}} \underset{\varepsilon \sim N_n(0, \sigma^2 I_n)}{\mathbb{E}} \left[ \| \underbrace{(\Phi^T\Phi + n\lambda I_d)^{-1} \Phi^T (\Phi\beta + \varepsilon)}_{T_0} - \beta \|^2_{\hat{\Sigma}} \right]$$

Next, we have that $\xrightarrow{\text{exercise}}$

$$T_0 = \left( \Phi^T\Phi + n\lambda I_d \right)^{-1} \Phi^T \varepsilon - n\lambda \left( \Phi^T\Phi + n\lambda I_d \right) \beta$$

Because $\varepsilon \perp \beta$ the expression reduces to

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

$$= T_1 + T_2$$

It can be seen to see that

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

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

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

$\hookrightarrow$ on hw

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

This is a lower bound on the minimax risk.

EXERCISE

Notice that $\mathrm{tr}\left((\hat{\Sigma} + \lambda I_d)^{-1}\hat{\Sigma}\right) = \sum_{j=1}^{d}\frac{\hat{\lambda}_j}{\hat{\lambda}_j + \lambda}$

where $\hat{\lambda}_j$ is the $j^{th}$ eigenvalue of $\hat{\Sigma}$.

This is decreasing in $\lambda$. So

$$\sup_\lambda T_1 + T_2 = \frac{\sigma^2}{n}\lim_{\lambda \downarrow 0}\mathrm{tr}\left((\hat{\Sigma} + \lambda I_d)^{-1}\hat{\Sigma}\right)$$

③

$$= \frac{\sigma^2}{n} \ tr \left( \hat{\Sigma}^{-1} \Sigma \right)$$

$$= \boxed{\frac{\sigma^2}{n} d}$$

$\downarrow$

This is the excess risk of $\hat{\beta}$ OLS

- About minimaxity for estimation: let $\mathcal{P} = \{ P_\theta, \theta \in \Theta \}$ (parameter space)

  be a parametric family of prob. distributions

  ( a parametric statistical model ). We are interested

  in estimating $\theta^*$, the true parameter, s.t.

  $X_1, \dots, X_n \overset{iid}{\sim} P_{\theta^*}$ $\longrightarrow$ function of data

  For any estimator $\hat{\theta}$ of $\theta^*$ ( where $\hat{\theta} = \hat{\theta}(X_1, \dots, X_n)$ )

  let $L(\hat{\theta}, \theta^*)$ be the loss function

  associated to $\hat{\theta}$ ( e.g. $L(\hat{\theta}, \theta^*) = \| \hat{\theta} - \theta^* \|^2$ ).

  The risk of $\hat{\theta}$ is the function

  $$\theta \in \Theta \longmapsto \mathbb{E}_{X_1, \dots, X_n \overset{iid}{\sim} P_\theta} \left[ L(\hat{\theta}, \theta) \right] = R(\hat{\theta}, \theta)$$

- The minimax risk for this problem

  $$\inf_{\hat{\theta}} \ \sup_{\theta \in \Theta} \ R(\hat{\theta}, \theta)$$

  $\downarrow$
  over all estimators

(4)

A natural lower bound on the minimax risk is the Bayes risk.

$$\inf_{\hat{\theta}} \mathbb{E}_{\theta \sim \pi} \left[ R(\hat{\theta}, \theta) \right] = R(\pi)$$

$$\downarrow$$
$$\text{prior}$$

a procedure attaining that infimum is called a Bayes procedure wrt $\pi$.

When the loss function is e.g. quadratic, the Bayes procedure is the posterior mean of $\theta$.

- If $\{\pi_k\}$ is a sequence of priors s.t.

$$R(\pi_k) \longrightarrow r \qquad \text{as} \qquad k \to \infty$$

and $\hat{\theta}$ is a procedure s.t

$$\sup_{\theta \subseteq \Theta} R(\hat{\theta}, \theta) = r$$

Then $\hat{\theta}$ is minimax $\longrightarrow$ HW

- Remark: if the covariate (i.e. the rows of $\Phi$) are random, the risk of OLS $\hat{\beta}$ is

$$\frac{\sigma^2}{n} \mathbb{E}_{\Phi} \quad \text{tr}\left( \hat{\Sigma}^{-1} \Sigma \right) \qquad \text{where}$$

⑤

$$\Sigma_i' = \mathbb{E}\left[\Phi_i \Phi_i^T\right]$$

$\quad\quad\quad\quad\quad\quad$ $\hookrightarrow$ $i$st row of $\Phi$.

This is also the minimax risk.

See Mourtarta's paper.

◼ STATISTICAL INFERENCE FOR $\beta^*$.

As before we assume a well-specified linear model and fixed covariates, ie.

$\quad\quad\quad\quad\quad\quad\quad\quad$ unknown

$$Y = \underset{\underset{\text{fixed}}{\downarrow}}{\Phi}\, \beta^* + \varepsilon$$

$\quad\quad n \times 1 \quad\quad\quad\quad\quad \hookrightarrow \overset{iid}{\sim} (0, \sigma^2)$

Goal : to estimate and carry out statistical inference for $\beta^*$, in fixed dimensions (i.e. $d$ is fixed)

• Is the OLS consistent ?

$$\hat{\beta} \overset{p}{\longrightarrow} \beta^*$$

• Yes ! Assume that $\hat{\Sigma_i'} = \dfrac{\Phi^T \Phi}{n} \longrightarrow \underset{d \times d}{\Sigma_i'}$ p.d.

$\quad\quad$ Then

$$\hat{\Sigma_i'}^{-1}\, \dfrac{\Phi^T \varepsilon}{n} \overset{p}{\longrightarrow} 0$$

$\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad$ transpose of $i$th
$\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad$ $\hookrightarrow$ row of $\Phi$

Pf/ By WLLN $\quad \dfrac{\Phi^T \varepsilon}{n} = \dfrac{1}{n}\sum_{i=1}^{n} \Phi_i\, \varepsilon_i \overset{p}{\longrightarrow} 0$ ⑥

To   see   this   $\Phi_i \varepsilon_i \sim (0, \sigma^2 \Phi_n \Phi_n^T)$

So   because   the   $\varepsilon_i$'s are indep.

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

$$= \frac{\sigma^2}{n} \left( \frac{\Phi^T \Phi}{n} - \Sigma + \Sigma_t \right)$$

$$= \frac{\sigma^2}{n} \Sigma_t + \frac{\sigma^2}{n} \left( \frac{\Phi^T \Phi}{n} - \Sigma_t \right) \to 0$$

$$\text{of } n \to \infty$$

Then   $\hat{\Sigma}_t^{-1} \frac{\Phi^T \varepsilon}{n} \xrightarrow{P} 0$   by   slutsky's theorem

$\hat{\beta}$   is   asymptotically   normal

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

$$\Downarrow$$

$$\mathbb{E}[\hat{\beta}]$$