

Lecture 18

Last time

- ▷ Intro to estimation
- ▷ Maximum likelihood estimation

Today

- ▷ Ordinary least squares
- ▷ Excess risk

Least Squares

Suppose we wished to solve the problem

$$(t) \min_{\theta \in \mathbb{R}^d} R_n(\theta) \text{ with } R_n(\theta) = \frac{1}{n} \sum_{i=1}^n (y_i - x_i^\top \theta)^2.$$

where

$y_i = \theta^*{}^\top x_i + \varepsilon_i$ Random or fixed

$\varepsilon_i \stackrel{\text{iid}}{\sim} N(0, \sigma^2)$

As we talked about last time this is known as the empirical risk and it approximates the population risk

$$R(\theta) = \mathbb{E} R_n(\theta).$$

We will use the following matrix notation

$$R_n(\theta) = \frac{1}{n} \|y - X\theta\|_2^2$$

where

$$y = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix} \text{ and } X = \begin{bmatrix} -x_1^T \\ \vdots \\ -x_n^T \end{bmatrix}.$$

We will assume $d \leq n$ and X has rank d . An optimal solution of (★) is called "ordinary least squares estimator" θ_{OLS} .

Lemma: θ_{OLS} exists and is unique. Further, it is given by

$$\theta_{OLS} = (X^T X)^{-1} X^T y.$$

Proof: The first statement follows since $R_n(\cdot)$ is strongly convex. Then, using first order optimality yields

$$0 = \nabla R_n(\theta) = \frac{2}{n} X^T (X\theta - y)$$

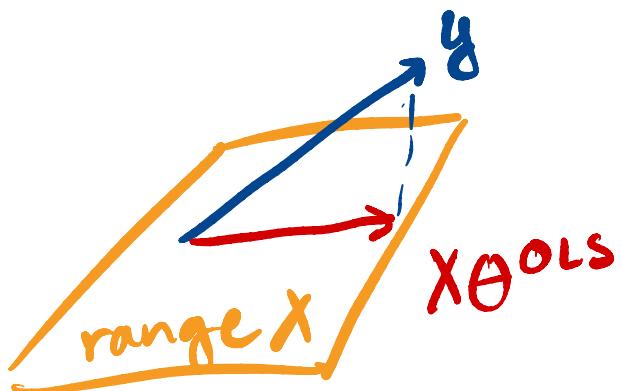
and so the formula for θ_{OLS} follows.

The OLS estimator has a nice geometric interpretation. 4

Lemma: The predictions

$$\hat{X}\theta_{OLS} = X(X^T X)^{-1} X^T y$$

is the orthogonal projection of
 y onto $\text{range}(X) \subseteq \mathbb{R}^n$. \rightarrow



Thus we can see
 θ_{OLS} as solving

- ① $\hat{y} = \text{Proj}_{\text{range}(X)} y$
- ② Solve $X\theta = \hat{y}$.

Excess risk

Natural question:

How close is $R(\theta_{OLS})$
to $\min_{\theta \in \mathbb{R}^d} R(\theta)$?

This is called the excess risk
We will study it in two situations

► Fixed design: Assume that X
is deterministic and we study

$$ER(\theta_{OLS}) = \min_{\theta \in \mathbb{R}^d} R(\theta)$$

↗ Expectation
 w.r.t. ϵ
 only.

R^*

▷ Random design: Assume that X is random and we take the expectation w.r.t. X as well.

We focus on the fixed design setting first. Define

$$\hat{\Sigma} = \frac{1}{n} X^T X$$

which by assumption is invertible. Any positive definite matrix defines an inner product via

$$\langle \theta, \theta' \rangle_{\hat{\Sigma}} := \theta^T \hat{\Sigma} \theta,$$

which induces a norm

$$\|\theta\|_{\hat{\Sigma}}^2 := \langle \theta, \theta \rangle_{\hat{\Sigma}} = \|\hat{\Sigma}^{-1/2} \theta\|_2^2 = \frac{1}{n} \|X\theta\|_2^2$$

Let's characterize the generalization error for any $\theta \in \mathbb{R}^d$.

Lemma(8): Suppose X is fixed. Then,

$$R(\theta) - R^* = \|\theta - \theta^*\|_{\hat{\Sigma}}^2 \quad \forall \theta \in \mathbb{R}^d$$

and, moreover, $R^* = \sigma^2$.

+

Proof: Expanding

$$\begin{aligned} R(\theta) &= \frac{1}{n} \mathbb{E}_{\epsilon} \|y - X\theta\|_2^2 \\ &= \frac{1}{n} \mathbb{E}_{\epsilon} \|X(\theta^* - \theta) + \epsilon\|_2^2 \\ &= \|\theta^* - \theta\|_{\hat{\Sigma}}^2 + \frac{1}{n} \mathbb{E} \|\epsilon\|_2^2 \\ &= \|\theta^* - \theta\|_{\hat{\Sigma}}^2 + \sigma^2. \end{aligned}$$

Thus, $R^* = R(\theta^*) = \sigma^2$ and the result follows. \square

Recall from our computation last class that

$$\mathbb{E} \|\hat{\theta} - \theta^*\|_{\hat{\Sigma}}^2 = \underbrace{\|\mathbb{E}\hat{\theta} - \theta^*\|_{\hat{\Sigma}}^2}_{\text{Bias}_{\hat{\Sigma}}(\hat{\theta})} + \underbrace{\mathbb{E} \|\hat{\theta} - \mathbb{E}\hat{\theta}\|_{\hat{\Sigma}}^2}_{\text{Var}_{\hat{\Sigma}}(\hat{\theta})}.$$

Next we study these two for the OLS estimator.

Lemma (ii): Suppose that X is fixed.

Then,

$$\text{Bias}_{\hat{\theta}}(\hat{\theta}) = 0 \quad \text{and} \quad E(\theta^{\text{as}} - \theta^*) (\theta^{\text{OLS}} - \theta^*)^T = \frac{\sigma^2}{n} \hat{\Sigma}.$$

Proof: Recall

$$\theta^{\text{OLS}} = (X^T X)^{-1} X^T y = (X^T X)^{-1} X^T (X \theta^* + \varepsilon) = \theta^* + (X^T X)^{-1} X^T \varepsilon.$$

So by linearity of expectation, $E \theta^{\text{OLS}} = \theta^*$.

Further,

$$\begin{aligned} E(\theta^{\text{OLS}} - \theta^*)(\theta^{\text{OLS}} - \theta^*)^T &= (X^T X)^{-1} X^T E \underbrace{\varepsilon \varepsilon^T}_{=\sigma^2 I} X (X^T X)^{-1} \\ &= \sigma^2 (X^T X)^{-1} \\ &= \frac{\sigma^2}{n} \hat{\Sigma}^{-1}. \end{aligned}$$

□

As a direct corollary of Lemmas (i) and (ii) we obtain a characterization of the excess risk.

Corollary: Suppose X is fixed. Then

$$E R(\theta^{\text{OLS}}) - R^* = \frac{d}{n} \sigma^2.$$

Proof: We know that since θ^{OLS} is

unbiased, the excess risk is $\text{Var}_{\hat{\Sigma}}(\hat{\theta})$. Thus

$$\text{Var}_{\hat{\Sigma}}(\theta^{\text{OLS}}) = \mathbb{E} \| \theta^{\text{OLS}} - \theta^* \|_2^2$$

$$\begin{aligned} & \xrightarrow{\substack{\text{Cyclic} \\ \text{invariance} \\ \text{of the trace}}} = \mathbb{E} \text{tr} (\hat{\Sigma} (\theta^{\text{OLS}} - \theta^*) (\theta^{\text{OLS}} - \theta^*)^\top) \\ & = \text{tr} (\hat{\Sigma} \mathbb{E} (\theta^{\text{OLS}} - \theta^*) (\theta^{\text{OLS}} - \theta^*)^\top) \\ & = \frac{\sigma^2}{n} \text{tr}(I) = \frac{d}{n} \sigma^2. \end{aligned}$$

□

Another natural question given that we don't have access to R , how close is $R_n(\theta^{\text{OLS}})$ to $R(\theta^{\text{OLS}})$?

Lemma: Suppose that X is fixed. Then,

$$\mathbb{E} R_n(\theta^{\text{OLS}}) = \sigma^2 - \frac{d}{n} \sigma^2.$$

Proof: Expanding

$$\begin{aligned} \mathbb{E} R_n(\theta^{\text{OLS}}) &= \frac{1}{n} \mathbb{E}_\epsilon \| X \theta^{\text{OLS}} - y \|_2^2 \\ &= \frac{1}{n} \mathbb{E}_\epsilon \| (X(X^\top X)^{-1} X - I) y \|_2^2 \\ &= \frac{1}{n} \mathbb{E}_\epsilon \| (X(X^\top X)^{-1} X - I)(X \theta^* + \epsilon) \|_2^2 \end{aligned}$$

$$= \frac{1}{n} \mathbb{E} \| \underbrace{(X(X^T X)^{-1} X - I) \varepsilon}_{P \text{ (projection matrix)}} \|_F^2$$

$$= \frac{1}{n} \text{tr}((P - I) \mathbb{E} \varepsilon \varepsilon^T (P - I)^T)$$

$$= \frac{\sigma^2}{n} \text{tr}((P - I) (P - I)^T)$$

↑
Projection onto a $(n-d)$ -dim
subspace.

$$= \frac{\sigma^2}{n} (n - d).$$

□