

Lecture 20

Last time

- ▷ Ordinary least squares
- ▷ Excess risk

Today

- ▷ Ridge Regression
- ▷ Random design
- ▷ Beyond linear features

Ridge Regression

The excess rate that we proved assumed that $d \leq n$ and X is full rank. But in high-dimensional settings $d \gg n$ and our guarantees do not hold.

Even when $\epsilon = 0$ there are many solutions, which one to choose?

One way to do this is by enforcing certain structure with a regularizer, i.e., an additional term in the objective (common choices are the l_1 norm and the l_2 -norm squared).

Def: For any $\lambda > 0$, define ridge least-squares estimator as

Prefers small norm

$$\Theta^\lambda = \underset{\Theta \in \mathbb{R}^d}{\operatorname{argmin}} \frac{1}{n} \|y - X\Theta\|_2^2 + \lambda \|\Theta\|_2^2.$$

Lemma: The ridge least-squares estimator exist and it is unique. Further,

$$\hat{\theta}^\lambda = \frac{1}{n} (\hat{\Sigma} + \lambda I)^{-1} X^T y$$

This matrix doesn't have to be invertible.

Proof: Just as before, the lost is convex and

$$0 = \nabla l_\lambda(\theta) = \frac{2}{n} X^T (X\theta - y) + 2\lambda \theta \\ \Rightarrow (\hat{\Sigma} + \lambda I)\theta = \frac{1}{n} y. \quad \square$$

Note that when $\hat{\Sigma}$ is invertible and $\lambda = 0$, this recovers θ_{OLS} . Indeed, as $\lambda \rightarrow 0$ $\theta^\lambda \rightarrow (X^T X)^+ X^T y$ where $(\cdot)^+$ denotes the pseudo-inverse.

Let's understand what happens to the bias and variance.

Lemma: We have that

$$\text{Bias}_{\hat{\Sigma}}(\theta^\lambda) = \lambda^2 \theta^{*\top} (\hat{\Sigma} + \lambda I)^{-2} \hat{\Sigma} \theta^*$$

$$\text{Var}_{\hat{\Sigma}}(\theta^\lambda) = \frac{\sigma^2}{n} \text{tr}(\hat{\Sigma}^2 (\hat{\Sigma} + \lambda I)^{-2})$$

Proof: We will use the following fact.

Fact: For any symmetric matrix A ,

$$A(A + \lambda I)^{-1} = (A + \lambda I)^{-1} A.$$

This follows since A and $(A + \lambda I)^{-1}$ share the same eigenvectors (Why?).

With this,

$$\begin{aligned}\mathbb{E} \theta^\lambda &= \frac{1}{n} (\hat{\Sigma} + \lambda I)^{-1} X^T \mathbb{E} y \\ &= \frac{1}{n} (\hat{\Sigma} + \lambda I)^{-1} X^T X \theta^* \\ &= (\hat{\Sigma} + \lambda I)^{-1} \hat{\Sigma} \theta^* \\ &= \theta^* + ((\hat{\Sigma} + \lambda I) \hat{\Sigma} - I) \theta^* \\ \text{Add and subtract } \hat{\Sigma} &= \theta^* + ((\hat{\Sigma} + \lambda I)^{-1} (\hat{\Sigma} + \lambda I - \lambda I) - I) \theta^* \\ &= \theta^* + \lambda (\hat{\Sigma} + \lambda I)^{-1} \theta^*.\end{aligned}$$

Thus,

$$\text{Bias}_{\hat{\Sigma}}(\theta^\lambda) = \lambda^2 \| (\hat{\Sigma} + \lambda I)^{-1} \theta^* \|_{\hat{\Sigma}}^2$$

$$\text{The fact above } \rightarrow = \lambda^2 \theta^{*\top} (\hat{\Sigma} + \lambda I)^{-2} \hat{\Sigma} \theta^*.$$

Expanding for Var gives

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

$$\begin{aligned}
&= \mathbb{E} \| \frac{1}{n} (\hat{\Sigma} + \lambda I)^{-1} X^T y - \\
&\quad \frac{1}{n} (\hat{\Sigma} + \lambda I)^{-1} X^T X \theta^* \|_{\hat{\Sigma}}^2 \\
&= \mathbb{E} \| \frac{1}{n} (\hat{\Sigma} + \lambda I)^{-1} X^T \varepsilon \|_{\hat{\Sigma}}^2 \\
&= \frac{1}{n} \mathbb{E} \text{tr}(\varepsilon^T X (\hat{\Sigma} + \lambda I)^{-1} \hat{\Sigma} (\hat{\Sigma} + \lambda I) X^T \varepsilon)
\end{aligned}$$

Cyclic invariance of tr \Rightarrow

$$= \frac{\sigma^2}{n} \text{tr}(\hat{\Sigma} (\hat{\Sigma} + \lambda I)^{-1} \hat{\Sigma} (\hat{\Sigma} + \lambda I))$$

Fact \Rightarrow

$$= \frac{\sigma^2}{n} \text{tr}(\hat{\Sigma}^2 (\Sigma + \lambda I)^{-2}).$$

Notice that

eigenvalues of $\hat{\Sigma}$.

$$\text{tr}(\hat{\Sigma}^2 (\hat{\Sigma} + \lambda I)^{-2}) = \sum_{i=1}^d \frac{\lambda_i^3}{(\lambda_i + \lambda)^2} < d$$

quantifies the intrinsic dimension or degrees of freedom of the problem. It provides a soft count on how many eigenvalues $\lambda_i \gg \lambda$.

Altogether we have proven:

Corollary:

$$\mathbb{E} R(\hat{\theta}_\lambda) - R^* = \lambda^2 \theta^* (\hat{\Sigma} + \lambda I)^{-2} \Sigma \theta^*$$

$$+ \frac{\sigma^2}{n} \text{tr}(\hat{\Sigma}^2 (\hat{\Sigma} + \lambda I)^{-2}).$$

Remarks:

- ▷ No dimension dependence! This is key for learning in infinite dim.
- ▷ This estimator is not unbiased any more (unless $\lambda = 0$).
- ▷ We can optimize over λ !

Theorem: Take $\bar{\lambda} = \frac{\sigma}{\|\theta^*\|_2} \sqrt{\frac{\text{tr}(\hat{\Sigma})}{n}}$. Then,

$$\mathbb{E} R(\theta^{\bar{\lambda}}) - R^* \leq \sigma \|\theta^*\|_2 \sqrt{\frac{\text{tr}(\hat{\Sigma})}{n}}.$$

Proof: Note that we can upper bound

$$\lambda^2 \theta^* (\hat{\Sigma} + \lambda I)^{-2} \hat{\Sigma} \theta^* \leq \lambda \underbrace{\|(\hat{\Sigma} + \lambda I)^{-2} \lambda \hat{\Sigma}\|_{\text{op}}}_{(\lambda_i + \lambda)^{-2} \lambda \lambda_i \leq \frac{1}{2}} \|\theta^*\|_2^2$$

$$(\lambda_i + \lambda)^{-2} \lambda \lambda_i \leq \frac{1}{2} \Leftrightarrow 2 \lambda \lambda_i \leq (\lambda_i + \lambda)^2 \forall i.$$

$$\leq \frac{\lambda}{2} \|\theta^*\|_2^2.$$

Further,

$$\frac{\sigma^2}{n} \text{tr}(\hat{\Sigma}^2 (\hat{\Sigma} + \lambda I)^{-2}) \leq \frac{\sigma^2}{\lambda n} \text{tr}(\hat{\Sigma} \lambda \hat{\Sigma} (\hat{\Sigma} + \lambda I)^{-2})$$

Show that $\text{tr}(AM) \leq \text{tr}(A) \lambda_{\max}(M)$

$$\leq \frac{\sigma^2}{2\lambda n} \text{tr}(\hat{\Sigma}).$$

Altogether

$$ER(\theta^\lambda) - R^* \leq \frac{\lambda}{2} \|\theta^*\|_2^2 + \frac{\sigma^2 \text{tr}(\hat{\Sigma})}{2\lambda n}.$$

Notice that this has the form

$$a\lambda + \frac{b}{\lambda} \quad \text{with } a, b > 0,$$

Check!

which is always minimized at $\lambda = \sqrt{b/a}$.

□

Remarks:

- ▷ The dependency on n is worst than for OLS. But there is no d and dependency on σ is better.
- ▷ It's hard to find λ (cross-validation in practice).
- ▷ This is extremely relevant for learning with Kernels (see Project topics).

Random Design

So far we assumed X was fixed, next we suppose (X_i, y_i) are drawn iid with

$$y_i = \theta^T X_i + \epsilon_i$$

iid $\mathbb{E}\epsilon_i = 0, \mathbb{E}\epsilon_i^2 = 1$

Lemma: Under the random design model.

$$\forall \theta, R(\theta) - R^* = \|\theta - \theta^*\|_{\Sigma}^2 \text{ with } \Sigma = \mathbb{E}XX^T$$

Proof: Just as before

$$\begin{aligned} R(\theta) &= \underset{\in \mathbb{R}}{\mathbb{E}} (y - X^T \theta)^2 \\ &= \underset{\in \mathbb{R}^n}{\mathbb{E}} (X^T (\theta^* - \theta) + \epsilon)^2 \\ &= (\theta^* - \theta)^T \mathbb{E} XX^T (\theta^* - \theta)^T + \sigma^2 \quad \square \end{aligned}$$

Lemma: Under the random design model and assuming $\hat{\Sigma}$ is invertible,

$$\mathbb{E} R(\theta^{OLS}) - R^* = \frac{\sigma^2}{n} \mathbb{E} \text{tr}(\Sigma \hat{\Sigma}^{-1})$$

Population $\xrightarrow{\text{Empirical}}$

Proof: Recall $\theta^{OLS} = \theta^* + \frac{1}{n} \hat{\Sigma}^{-1} X^T \epsilon$. Thus,

$$\mathbb{E} R(\theta^{OLS}) - R^* = \mathbb{E} \|\theta^{OLS} - \theta^*\|_{\Sigma}^2$$

$$\begin{aligned}
&= \frac{1}{n^2} \mathbb{E} \epsilon^\top X \hat{\Sigma}^{-1} \Sigma^{-1} \hat{\Sigma}^{-1} X^\top \epsilon \\
&= \frac{1}{n} \mathbb{E} \text{tr}(\epsilon \epsilon^\top \hat{\Sigma} \hat{\Sigma}^{-1} \Sigma \hat{\Sigma}^{-1}) \\
&= \frac{1}{n} \mathbb{E} \text{tr}(\mathbb{E}_{\epsilon} \epsilon \epsilon^\top \Sigma \hat{\Sigma}^{-1}) \\
&= \frac{\sigma^2}{n} \mathbb{E} \text{tr}(\Sigma \hat{\Sigma}^{-1}). \quad \square
\end{aligned}$$

Note that unlike before $\Sigma \hat{\Sigma}^{-1} \neq I$ and the bound depends on how fast does $\mathbb{E} \text{tr}(\Sigma \hat{\Sigma}^{-1}) \rightarrow d$.

For instance if $z \sim N(0, \Sigma)$ one can show

$$\mathbb{E} \text{tr}(\Sigma \hat{\Sigma}^{-1}) = \frac{d}{1 - \frac{d+1}{n}}$$

Try to prove it.

Hint: Inverse Wishart dist.

Beyond linear features.

Everything we developed holds

if we imagine that

$$y_i = \theta^*{}^\top \psi(x_i) + \varepsilon_i$$

where $\psi: \mathbb{R}^k \rightarrow \mathbb{R}^d$. Then we substitute
to the matrix X by

$$\tilde{\Phi} = \begin{pmatrix} -\psi(x_1) - \\ \vdots \\ -\psi(x_n) - \end{pmatrix},$$

and $\hat{\Sigma}$ by $\frac{1}{n} \tilde{\Phi}^\top \tilde{\Phi}$ and every-
thing follows just as before.