

Comparison between OLS and Ridge in terms of ℓ_2 estimation error

Xin Bing

Department of Statistical Sciences
University of Toronto

A theoretical understanding of the role of regularization

Consider the linear regression

$$Y = X^T \beta + \epsilon.$$

Suppose we have i.i.d. observations $(x_1, y_1), \dots, (x_n, y_n)$. Further assume the design matrix $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_n)^T \in \mathbb{R}^{n \times p}$ is deterministic and orthonormal, i.e.

$$\frac{1}{n} \mathbf{X}^T \mathbf{X} = \mathbf{I}_p.$$

Consider the ridge estimator $\hat{\beta}_\lambda^R$ of β for any given regularization parameter $\lambda \geq 0$. Let $\hat{\beta}$ be the OLS estimator of β .

We now contrast the behaviour of the ridge estimator with that of the OLS estimator side by side.

- Criteria:

$$\hat{\beta} = \underset{\beta}{\operatorname{argmin}} \|\mathbf{y} - \mathbf{X}\beta\|_2^2$$

$$\hat{\beta}_\lambda^R = \underset{\beta}{\operatorname{argmin}} \|\mathbf{y} - \mathbf{X}\beta\|_2^2 + \lambda \|\beta\|_2^2.$$

- Closed-form solutions:

$$\hat{\beta} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y} = \frac{1}{n} \mathbf{X}^\top \mathbf{y}$$

$$\hat{\beta}_\lambda^R = (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I}_p)^{-1} \mathbf{X}^\top \mathbf{y} = \frac{1}{n + \lambda} \mathbf{X}^\top \mathbf{y}.$$

- We examine their statistical properties of estimating β in terms of
 - ▶ bias
 - ▶ variance
 - ▶ mean squared error

- OLS: unbiased

$$\mathbb{E}[\hat{\beta}] = \mathbb{E}\left[\frac{1}{n}\mathbf{X}^\top \mathbf{y}\right]$$

$$= \mathbb{E}\left[\frac{1}{n}\mathbf{X}^\top (\mathbf{X}\beta + \epsilon)\right]$$

$$= \frac{1}{n}\mathbf{X}^\top \mathbf{X}\beta + \mathbb{E}\left[\frac{1}{n}\mathbf{X}^\top \epsilon\right]$$

$$= \frac{1}{n}\mathbf{X}^\top \mathbf{X}\beta + \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbb{E}[\epsilon_i]$$

$$= \beta.$$

$$\text{by } \frac{1}{n}\mathbf{X}^\top \mathbf{X} = \mathbf{I}_p$$

$$\text{by } \mathbf{y} = \mathbf{X}\beta + \epsilon$$

\mathbf{X} and β are deterministic

$$\text{by } \mathbb{E}[\epsilon_i] = 0$$

Bias of the ridge estimator

By repeating similar arguments as before:

- Ridge: biased

$$\begin{aligned}\mathbb{E}[\hat{\beta}_\lambda^R] &= \mathbb{E}\left[\frac{1}{n+\lambda}\mathbf{X}^\top\mathbf{y}\right] \\ &= \mathbb{E}\left[\frac{1}{n+\lambda}\mathbf{X}^\top(\mathbf{X}\beta + \epsilon)\right] \\ &= \frac{1}{n+\lambda}\mathbf{X}^\top\mathbf{X}\beta + \mathbb{E}\left[\frac{1}{n+\lambda}\mathbf{X}^\top\epsilon\right] \\ &= \frac{n}{n+\lambda}\beta \\ &= \beta - \frac{\lambda}{n+\lambda}\beta.\end{aligned}$$

Variance of the OLS and ridge estimators

- OLS:

$$\begin{aligned}\text{Cov}(\hat{\beta}) &= \frac{1}{n^2} \mathbf{X}^\top \text{Cov}(\mathbf{y}) \mathbf{X} \\ &= \frac{\sigma^2}{n^2} \mathbf{X}^\top \mathbf{X} && \text{by } \text{Cov}(\mathbf{y}) = \sigma^2 \mathbf{I}_n \\ &= \frac{\sigma^2}{n} \mathbf{I}_p.\end{aligned}$$

- Ridge:

$$\text{Cov}(\hat{\beta}_\lambda^R) = \frac{1}{(n + \lambda)^2} \mathbf{X}^\top \text{Cov}(\mathbf{y}) \mathbf{X} = \frac{\sigma^2 n}{(n + \lambda)^2} \mathbf{I}_p.$$

- OLS:

$$\begin{aligned} & \mathbb{E}[\|\hat{\beta} - \beta\|_2^2] \\ &= \mathbb{E}[(\hat{\beta} - \beta)^\top (\hat{\beta} - \beta)] \\ &= \mathbb{E}[(\hat{\beta} - \mathbb{E}[\hat{\beta}] + \mathbb{E}[\hat{\beta}] - \beta)^\top (\hat{\beta} - \mathbb{E}[\hat{\beta}] + \mathbb{E}[\hat{\beta}] - \beta)] \\ &= \mathbb{E}[(\hat{\beta} - \mathbb{E}[\hat{\beta}])^\top (\hat{\beta} - \mathbb{E}[\hat{\beta}])] + \mathbb{E}[(\mathbb{E}[\hat{\beta}] - \beta)^\top (\mathbb{E}[\hat{\beta}] - \beta)] \\ &= \text{trace} \mathbb{E}[(\hat{\beta} - \mathbb{E}[\hat{\beta}]) (\hat{\beta} - \mathbb{E}[\hat{\beta}])^\top] + (\mathbb{E}[\hat{\beta}] - \beta)^\top (\mathbb{E}[\hat{\beta}] - \beta) \\ &= \text{trace}[\text{Cov}(\hat{\beta})] + \|\mathbb{E}[\hat{\beta}] - \beta\|_2^2 \\ &= \underbrace{\frac{\sigma^2 p}{n}}_{\text{Variance}} + \underbrace{0}_{\text{Bias}}. \end{aligned}$$

ℓ_2 estimation error of the ridge predictor

- Ridge:

$$\begin{aligned}\mathbb{E}[\|\hat{\beta}_\lambda^R - \beta\|_2^2] &= \text{trace}[\text{Cov}(\hat{\beta}_\lambda^R)] + \left\| \mathbb{E}[\hat{\beta}_\lambda^R] - \beta \right\|_2^2 \\ &= \underbrace{\frac{\sigma^2 pn}{(n + \lambda)^2}}_{\text{Variance}} + \underbrace{\left(\frac{\lambda}{n + \lambda} \right)^2 \|\beta\|_2^2}_{\text{Bias}}.\end{aligned}$$

- Recall for OLS:

$$\mathbb{E}[\|\hat{\beta} - \beta\|_2^2] = \underbrace{\frac{\sigma^2 p}{n}}_{\text{Variance}} + \underbrace{0}_{\text{Bias}}.$$

Remark: Ridge estimator has smaller variance by paying extra bias as the price. **This is the essential idea of regularization!** The balance between variance and bias of ridge is controlled by the magnitude of λ .

Same phenomenon for prediction

Since we predict $X = x$ by

- OLS:

$$\hat{y} = x^T \hat{\beta}$$

- Ridge:

$$\hat{y}_\lambda^R = x^T \hat{\beta}_\lambda^R$$

Regularization controlled by λ has the same effects on prediction MSE.

Same phenomenon for the Lasso

The same idea holds for the Lasso. But the analysis of the MSE estimation error of the Lasso is less straightforward than that of Ridge.