

# Theoretical Understanding the Regularization Effect of Ridge

Xin Bing

Department of Statistical Sciences  
University of Toronto

# A theoretical understanding of the role of regularization

Consider the linear regression

$$Y = \mathbf{X}^\top \boldsymbol{\beta} + \epsilon.$$

Suppose we have i.i.d. observations  $(x_1, y_1), \dots, (x_n, y_n)$ . Further assume the design matrix  $\mathbf{X}$  is deterministic and orthonormal, i.e.

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

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

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  $\beta$  in terms of
  - ▶ bias
  - ▶ variance
  - ▶ mean squared error

# Bias of the OLS and ridge estimators

- OLS: unbiased

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

- 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:

$$\text{Cov}(\hat{\beta}) = \frac{1}{n^2} \mathbf{X}^\top \text{Cov}(\mathbf{y}) \mathbf{X} = \frac{\sigma^2}{n^2} \mathbf{X}^\top \mathbf{X} = \frac{\sigma^2}{n}.$$

- Ridge:

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

# Estimation error of the OLS and ridge estimators

- OLS:

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

- Ridge:

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

**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  $\lambda$ .

# 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  $\lambda$  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.