

# 泛化理论

## 第一章 传统统计模型

### §1.1.2 岭回归

@ 滕佳烨

## Recall:

Consistency: when we have infinite training samples ( $n \rightarrow \infty$ ), the estimator  $\hat{\beta}$  should converge to the true parameter  $\beta^*$ .

- Basic linear regression:  $y = x^\top \beta^* + \epsilon$
- LSE (least-square estimator): minimize  $MSE$  (min square error)  $\|Y - X\beta\|^2$   
$$\hat{\beta} = (X^\top X)^{-1} X^\top Y, \quad \hat{\beta} | X \sim N(\beta^*, \sigma^2 (X^\top X)^{-1})$$

## Today's topic:

1. Bias-variance tradeoff
2. Ridge regression

## Bias-variance tradeoff

We consider the term  $\mathbb{E}\|\hat{\beta} - \beta^*\|^2$  which measures the distance between  $\hat{\beta}$  and  $\beta^*$ . (The expectation is taken over the randomness of  $Y$ . We still consider fixed design.)

It can be split as the **bias** and **variance** component

$$\mathbb{E}\|\hat{\beta} - \beta^*\|^2 = \|\mathbb{E}\hat{\beta} - \beta^*\|^2 + \mathbb{E}\|\hat{\beta} - \mathbb{E}\hat{\beta}\|^2$$

**For LSE estimator**  $\hat{\beta} = (X^\top X)^{-1}X^\top Y$ .

- It is unbiased:  $\mathbb{E}\hat{\beta}|X = \beta^*$
- Its variance is  $\mathbb{E}\|\hat{\beta} - \mathbb{E}\hat{\beta}\|^2|X = \text{Trace}[\sigma^2(X^\top X)^{-1}]$
- LSE estimator has the smallest variance among all the unbiased estimator (MVUE, minimal variance unbiased estimator).

Question: when the eigenvalues of  $(X^\top X)$  is small, the estimator has large variance!

Can we find a **biased estimator** with **small variance**, while controlling the bias?

## Intuition for Ridge Regression

LSE estimator:  $\hat{\beta} = (X^T X)^{-1} X^T Y$ .

- The variance term  $\sigma^2 (X^T X)^{-1}$  may be too large due to small eigenvalues of  $(X^T X)$ .
- The variance term is stemmed from  $(X^T X)^{-1}$  in  $\hat{\beta}$
- Boost the eigenvalue of  $(X^T X)$ !

Ridge regression estimator:  $\hat{\beta}_r(\rho) = (X^T X + \rho I)^{-1} X^T Y$ .

- The minimal eigenvalue is at least  $\rho$
- The estimator is biased:  $\mathbb{E}\hat{\beta}_r | X \neq \beta^*$
- The variance term can be smaller than  $\hat{\beta}$

It can be proved that under some mild assumptions, there are some  $\rho^*$  such that

$$\mathbb{E}\|\hat{\beta}_r(\rho^*) - \beta^*\|^2 < \mathbb{E}\|\hat{\beta} - \beta^*\|^2$$

Hint: prove that the derivation  $\frac{d\|\hat{\beta}_r(\rho) - \beta^*\|^2}{d\rho} \Big|_{\rho=0} < 0$ . Note that  $\hat{\beta}_r(0) = \hat{\beta}$ .

## Another intuition for Ridge Regression

1. The corresponding loss for Ridge regression:

$$L(\beta) = \frac{1}{n} \|Y - X\beta\|^2 + \rho \|\beta\|^2$$

Penalty on  $\beta$ ! Prior: beta is not too large.

2. Computational stability:  $\hat{\beta} = (X^T X)^{-1} X^T Y$ .

It is not stable to compute  $(X^T X)^{-1}$  with small eigenvalues...

## Take-away messages

- (a) Bias-variance tradeoff.
- (b) Ridge regression: bias but less variance.
- (c) Some intuition behind ridge regression: reduce the variance; prior; stable.

All the slides will be available at [www.tengjiaye.com/generalization](http://www.tengjiaye.com/generalization) soon.

@ 滕佳焯

Thanks!