

# Uniform consistency of cross-validation estimators for high-dimensional ridge regression

Pratik Patil, Yuting Wei, Alessandro Rinaldo, Ryan J. Tibshirani

Carnegie Mellon University

AISTATS 2021

## Main punchline

- Standard regression with  $n$  data pairs  $(x_i, y_i) \in \mathbb{R}^p \times \mathbb{R}$
- Given a tuning parameter  $\lambda$ , recall that **ridge regression** solves

$$\underset{\beta \in \mathbb{R}^p}{\text{minimize}} \sum_{i=1}^n (y_i - x_i^T \beta)^2 / n + \lambda \|\beta\|_2^2$$

- Choice  $\lambda$  crucially affects the performance of the fitted estimator

Key question: how to **select  $\lambda$**  based on observed data in **high dimensions**

We show: under proportional asymptotics as  $n \rightarrow \infty$ ,  $p/n \rightarrow \gamma \in (0, \infty)$ , the **leave-one-out** and **generalized cross-validation** almost surely,

1. converge to **out-of-sample prediction error** uniformly in  $\lambda$ ;
2. pick optimal  $\lambda$  for prediction error, including when  $\lambda = 0$  or negative

# Outline

Problem setup

Main results

Proof intuitions

## High-dimensional ridge regression

- Let  $X \in \mathbb{R}^{n \times p}$  denote feature matrix,  $y \in \mathbb{R}^n$  denote response vector
- Let  $\hat{\beta}_\lambda := \arg \min_{\beta \in \mathbb{R}^p} \|y - X\beta\|_2^2/n + \lambda\|\beta\|_2^2$  denote ridge estimate
  - if  $\lambda > 0$ , problem convex in  $\beta$  and has an explicit solution:

$$\hat{\beta}_\lambda = (X^T X/n + \lambda I_p)^{-1} X^T y/n$$

- for any  $\lambda \in \mathbb{R}$ , extend using Moore-Penrose inverse:

$$\hat{\beta}_\lambda = (X^T X/n + \lambda I_p)^+ X^T y/n$$

- when  $\lambda = 0$ , this reduces to least squares sol with minimum  $\ell_2$  norm; in particular, when  $\text{rank}(X) = n \leq p$ , the solution interpolates data, i.e.  $X\hat{\beta} = y$ , and has minimum  $\ell_2$  norm among all interpolators

## Prediction error and cross validation

- We measure the performance of fitted models  $\hat{\beta}_\lambda$  by their expected squared **out-of-sample prediction error** defined as

$$\text{err}(\lambda) := \mathbb{E}_{x_0, y_0} [(y_0 - x_0^T \hat{\beta}_\lambda)^2 \mid X, y],$$

where  $(x_0, y_0)$  is test pair sampled from same training distribution

- random (conditional on observed data  $X$  and  $y$ )
  - unknown (depends on characteristics of data generating distribution)
- Several estimators of prediction error:
    - $k$ -fold cross validation (large bias when  $k = 5$  or even when  $k = 10$ )
    - Generalized cross validation
    - Stein unbiased error estimate (in-sample prediction error)

We study the case when  $k = n$  also called **leave-one-out cross-validation**, and **generalized cross-validation**

## Leave-one-out and generalized cross-validation

- Leave-one-out cross-validation (LOOCV):
  - for every  $i$ , train on all data except  $(x_i, y_i)$ , call the estimate  $\widehat{\beta}_\lambda^{-i}$
  - compute test error on the  $i^{\text{th}}$  point and take average

$$\begin{aligned}\text{loo}(\lambda) &= \frac{1}{n} \sum_{i=1}^n \left( y_i - x_i^T \widehat{\beta}_\lambda^{-i} \right)^2 \\ &\stackrel{\text{(shortcut)}}{=} \frac{1}{n} \sum_{i=1}^n \left( \frac{y_i - x_i^T \widehat{\beta}_\lambda}{1 - [L_\lambda]_{ii}} \right)^2\end{aligned}$$

where  $L_\lambda = X(X^T X/n + \lambda I_p)^+ X^T/n$  is the ridge smoothing matrix

- Generalized cross-validation (GCV)
  - same as leave-one-out shortcut but a single re-weighting

$$\text{gcv}(\lambda) = \frac{1}{n} \sum_{i=1}^n \left( \frac{y_i - x_i^T \widehat{\beta}_\lambda}{1 - \text{tr}[L_\lambda]/n} \right)^2$$

- When  $\widehat{\beta}_\lambda$  is an interpolator, i.e.  $L_\lambda = I_n$ , both estimates are “0/0”; we then define the estimates as their respective limits as  $\lambda \rightarrow 0$

## Goals of the paper

There are two main questions that we answer in this paper:

1. How do  $\text{gcv}(\lambda)$  and  $\text{loo}(\lambda)$  compare to  $\text{err}(\lambda)$  as functions of  $\lambda$ ?
2. How do  $\text{err}(\hat{\lambda}_I^{\text{gcv}})$  and  $\text{err}(\hat{\lambda}_I^{\text{loo}})$  compare to  $\text{err}(\lambda_i^*)$  where  $\lambda_i^*$  denotes the optimal oracle ridge tuning parameter

$$\lambda_i^* = \arg \min_{\lambda \in I \subseteq \mathbb{R}} \text{err}(\lambda),$$

and  $\hat{\lambda}_I^{\text{gcv}}$  and  $\hat{\lambda}_I^{\text{loo}}$  denote the corresponding tuning parameters that minimize GCV and LOOCV over an interval  $I$ ?

# Outline

Problem setup

**Main results**

Proof intuitions



## Summary of main results

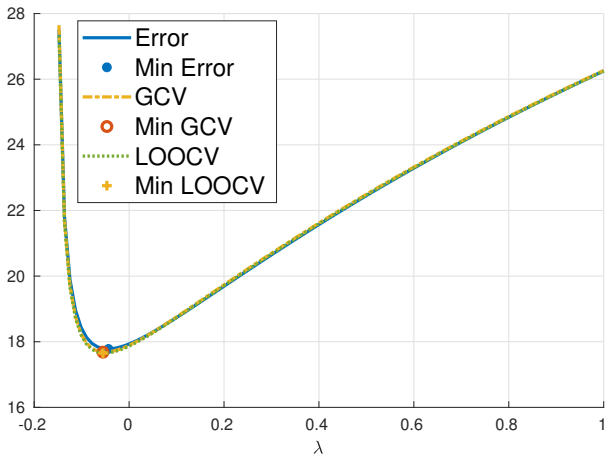
Under i.i.d. sampling with

- a well-specified model  $y = x^T \beta_0 + \varepsilon$  where  $\varepsilon$  is independent of  $x$
- decomposable features  $x = \Sigma^{1/2} z$  where  $z$  contains i.i.d. entries
- certain bnd moment and norm cond. on  $\varepsilon$  and  $z$ , and  $\beta_0$  and  $\Sigma$ , resp.

as  $n \rightarrow \infty$  and  $p/n \rightarrow \gamma \in (0, \infty)$ , we show

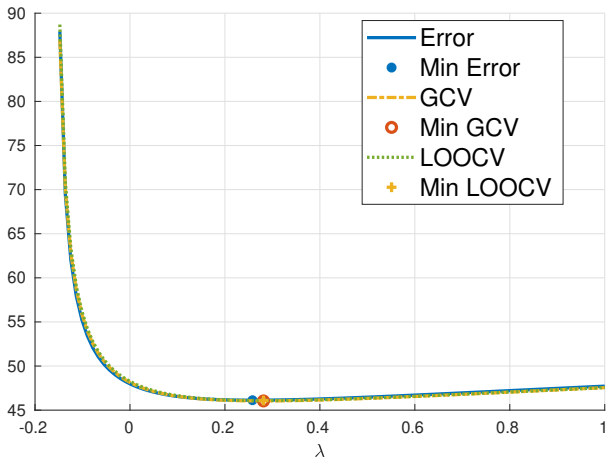
1. GCV pointwise convergence
  - $\text{gcv}(\lambda)$  converges to  $\text{err}(\lambda)$  pointwise in  $\lambda$
2. GCV uniform convergences
  - convergence holds uniformly over compact intervals of  $\lambda$  including 0
3. LOOCV convergences
  - the analogous results hold for  $\text{loo}(\lambda)$  by relating it to  $\text{gcv}(\lambda)$
4. Optimal tuned prediction errors
  - both  $\text{err}(\hat{\lambda}_j^{\text{gcv}})$  and  $\text{err}(\hat{\lambda}_j^{\text{loo}})$  converge to  $\text{err}(\lambda_j^*)$

## Numerical illustration (negative optimal regularization)



- Overparametrized regime ( $p = 12000$ ,  $n = 6000$ )
- Autoregressive  $\Sigma$
- $\beta_0$  aligned with the top eigendirection of  $\Sigma$

## Numerical illustration (positive optimal regularization)



- Overparametrized regime ( $p = 12000$ ,  $n = 6000$ )
- Autoregressive  $\Sigma$
- $\beta_0$  aligned with the bottom eigendirection of  $\Sigma$

# Outline

Problem setup

Main results

Proof intuitions

## GCV versus prediction error: two key proof steps

Step 1: bias and variance decompositions of prediction error and GCV

Let  $\widehat{\Sigma} := X^T X/n$  denote the sample covariance matrix.

- limiting bias-like components:

– prediction error

$$\text{err}_b(\lambda) := \lambda^2 \beta_0^T (\widehat{\Sigma} + \lambda I)^+ \Sigma (\widehat{\Sigma} + \lambda I)^+ \beta_0$$

– gcv

$$\text{gcv}_b(\lambda) := \frac{\lambda^2 \beta_0^T (\widehat{\Sigma} + \lambda I)^+ \widehat{\Sigma} (\widehat{\Sigma} + \lambda I)^+ \beta_0}{(1 - \text{tr}[(\widehat{\Sigma} + \lambda I_p)^+ \widehat{\Sigma}]/n)^2}$$

- limiting variance-like components:

– prediction error

$$\text{err}_v(\lambda) := \sigma^2 \left[ 1 + \text{tr}[(\widehat{\Sigma} + \lambda I_p)^+ \Sigma]/n \right] - \sigma^2 \text{tr}[(\widehat{\Sigma} + \lambda I_p)^+ \Sigma (\widehat{\Sigma} + \lambda I_p)^+]/n$$

– gcv

$$\text{gcv}_v(\lambda) := \sigma^2 \left[ \frac{1}{1 - \text{tr}[(\widehat{\Sigma} + \lambda I_p)^+ \widehat{\Sigma}]/n} \right] - \frac{\sigma^2 \text{tr}[(\widehat{\Sigma} + \lambda I_p)^+ \widehat{\Sigma} (\widehat{\Sigma} + \lambda I_p)^+]/n}{(1 - \text{tr}[(\widehat{\Sigma} + \lambda I_p)^+ \widehat{\Sigma}]/n)^2}$$

## GCV versus prediction error: two key proof steps

Step 2: bias and variance equivalences for prediction error and GCV

- bias component equivalence:

$$\lambda^2 \beta_0^T (\widehat{\Sigma} + \lambda I)^+ \Sigma (\widehat{\Sigma} + \lambda I)^+ \beta_0 - \frac{\lambda^2 \beta_0^T (\widehat{\Sigma} + \lambda I)^+ \widehat{\Sigma} (\widehat{\Sigma} + \lambda I)^+ \beta_0}{(1 - \text{tr} [(\widehat{\Sigma} + \lambda I_p)^+ \widehat{\Sigma}] / n)^2} \xrightarrow{\text{a.s.}} 0$$

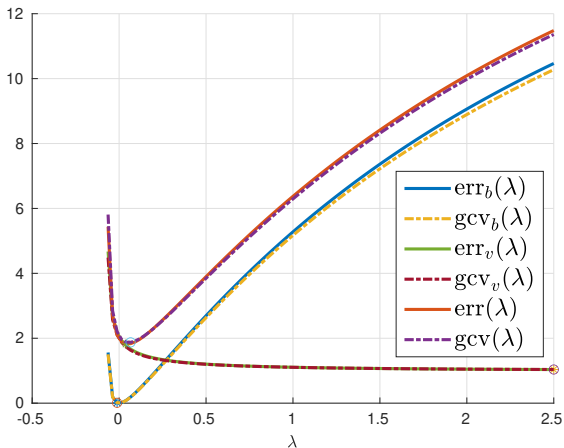
- variance component equivalences:

$$\sigma^2 \text{tr} [(\widehat{\Sigma} + \lambda I_p)^+ \Sigma (\widehat{\Sigma} + \lambda I_p)^+] / n - \frac{\sigma^2 \text{tr} [(\widehat{\Sigma} + \lambda I_p)^+ \widehat{\Sigma} (\widehat{\Sigma} + \lambda I_p)^+] / n}{(1 - \text{tr} [(\widehat{\Sigma} + \lambda I_p)^+ \widehat{\Sigma}] / n)^2} \xrightarrow{\text{a.s.}} 0$$

$$\sigma^2 \text{tr} [(\widehat{\Sigma} + \lambda I_p)^+ \Sigma] / n - \frac{\sigma^2 \text{tr} [(\widehat{\Sigma} + \lambda I_p)^+ \widehat{\Sigma}] / n}{1 - \text{tr} [(\widehat{\Sigma} + \lambda I_p)^+ \widehat{\Sigma}] / n} \xrightarrow{\text{a.s.}} 0$$

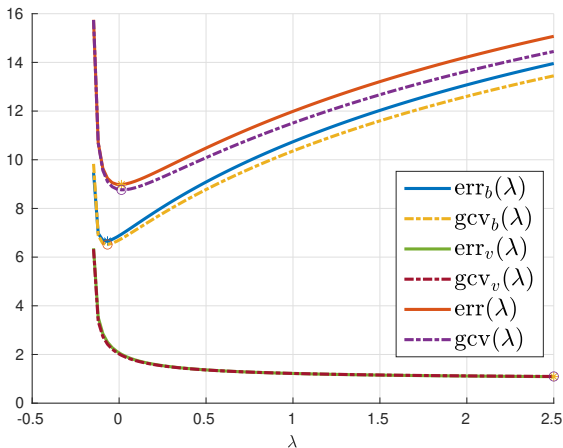
Main message: the GCV denominator proves to be the right correction for the excess optimism in the biased GCV numerator of training error

## Bias and variance equivalence numerical illustration



- Underparametrized ( $n = 6000, p = 3000$ )
- Bias minimized at  $\lambda = 0$  and variance decreases as  $\lambda$  increases
- Optimal  $\lambda$  always positive is underparametrized regime!

## Bias and variance equivalence numerical illustration



- Overparametrized ( $p = 12000$ ,  $n = 6000$ )
- Bias no longer minimized  $\lambda = 0$  and variance still decreasing in  $\lambda$
- Optimal  $\lambda$  may be negative in overparametrized regime!



## Discussion and future directions

This work shows GCV and LOOCV uniformly track squared out-of-sample prediction error for ridge regression under proportional asymptotics.

Main tool:

$$(\widehat{\Sigma} + \lambda I_p)^+ \Sigma \asymp \frac{(\widehat{\Sigma} + \lambda I_p)^+ \widehat{\Sigma}}{1 - \text{tr}[(\widehat{\Sigma} + \lambda I_p)^+ \widehat{\Sigma}]/n}$$

where for any two sequences of matrices  $A_p$  and  $B_p$ ,  $A_p \asymp B_p$  is used to mean  $\text{tr}[C_p(A_p - B_p)] \xrightarrow{\text{a.s.}} 0$  for any deterministic seq of matrices  $C_p$  of bnd trace norm

Going beyond ...

- Equivalences for general functionals of out-of-sample distributions
- Equivalences for general estimators
- Finite sample analysis and rates of convergence
-

Thanks for listening!

Questions/comments/thoughts?