

Estimating functionals of the out-of-sample error distribution in high-dimensional ridge regression

Pratik Patil, Alessandro Rinaldo, Ryan J. Tibshirani

Carnegie Mellon University

AISTATS 2022

Motivation and punchline of the paper

- Given $\mathcal{D} = \{(x_i, y_i) \in \mathbb{R}^p \times \mathbb{R}, 1 \leq i \leq n\}$, let $\hat{\beta}_\lambda$ be **ridge estimator**:

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

- The **out-of-sample error** of $\hat{\beta}_\lambda$ is $y_0 - x_0^\top \hat{\beta}_\lambda$ for a test point (x_0, y_0)
- Estimating out-of-sample error well is crucial for model assessment
- Prior work shows leave-out-out and generalized cross-validation consistently estimate the expected squared error $\mathbb{E}[(y_0 - x_0^\top \hat{\beta}_\lambda)^2 \mid \mathcal{D}]$

Key question: can we reliably estimate the entire out-of-sample error distribution and its linear and non-linear functionals in high dimensions?

We show, that under proportional asymptotics, almost surely:

- the empirical distributions of re-weighted in-sample errors from **leave-one-out** and **generalized cross-validation** converge weakly to the out-of-sample error distribution, even when $\lambda = 0$
- the plug-in estimators of these empirical distributions consistent for a broad class of linear and non-linear functionals of error distribution

Outline

Problem setup

Distribution estimation

Functional estimation

High-dimensional ridge regression

- Let $X \in \mathbb{R}^{n \times p}$ denote feature matrix, $y \in \mathbb{R}^n$ denote response vector
- Let $\hat{\beta}_\lambda \in \mathbb{R}^p$ denote the ridge estimator at regularization level λ :

$$\hat{\beta}_\lambda := \arg \min_{\beta \in \mathbb{R}^p} \|y - X\beta\|_2^2/n + \lambda \|\beta\|_2^2$$

- if $\lambda > 0$, the problem is convex in β 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 the solution 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 ℓ_2 norm; in particular, when $\text{rank}(X) = n \leq p$, the solution interpolates data, i.e. $X\hat{\beta} = y$, and has minimum ℓ_2 norm among all interpolators

Out-of-sample error distribution and its functionals

- Let P_λ denote distribution of out-of-sample error of $\widehat{\beta}_\lambda$:

$$P_\lambda = \mathcal{L}(y_0 - x_0^\top \widehat{\beta}_\lambda \mid X, y),$$

where (x_0, y_0) is sampled indep from the same training distribution

- a random distribution (conditional on observed data X and y)
- Let ψ denote a functional such that $P \mapsto \psi(P) \in \mathbb{R}$:
 - Linear functional:

$$\psi(P_\lambda) = \int t(z) dP_\lambda(z) = \mathbb{E}[t(y_0 - x_0^\top \widehat{\beta}_\lambda) \mid X, y],$$

where $t: \mathbb{R} \rightarrow \mathbb{R}$ is an error function (e.g., squared or absolute error)

- Nonlinear functional:

$$\psi(P_\lambda) = \text{Quantile}(P_\lambda; \tau) = \inf\{z : F_\lambda(z) \geq \tau\},$$

where F_λ denotes the cumulative distribution function of P_λ

We construct estimators of P_λ and $\psi(P_\lambda)$ by suitably extending [leave-one-out cross-validation](#) and [generalized cross-validation](#) procedures.

Standard 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 $\hat{\beta}_\lambda^{-i}$
 - compute test error on the i^{th} point and take average

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

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

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

- Standard LOOCV and GCV are consistent for the expected squared out-of-sample prediction error

Proposed estimators

Natural estimators for P_λ and $\psi(P_\lambda)$ building off from GCV and LOOCV.

- Empirical distributions of the GCV, LOO re-weighted errors:

$$\hat{P}_\lambda^{\text{gcv}} = \frac{1}{n} \sum_{i=1}^n \delta\left(\frac{y_i - x_i^\top \hat{\beta}_\lambda}{1 - \text{tr}[L_\lambda]/n}\right) \quad \text{and} \quad \hat{P}_\lambda^{\text{loo}} = \frac{1}{n} \sum_{i=1}^n \delta\left(\frac{y_i - x_i^\top \hat{\beta}_\lambda}{1 - [L_\lambda]_{ii}}\right)$$

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

$$\hat{P}_0^{\text{gcv}} = \frac{1}{n} \sum_{i=1}^n \delta\left(\frac{[(XX^\top)^\dagger y]_i}{\text{tr}[(XX^\top)^\dagger]/n}\right) \quad \text{and} \quad \hat{P}_0^{\text{loo}} = \frac{1}{n} \sum_{i=1}^n \delta\left(\frac{[(XX^\top)^\dagger y]_i}{[(XX^\top)^\dagger]_{ii}}\right)$$

- Plug-in GCV and LOO estimators:

$$\hat{\psi}_\lambda^{\text{gcv}} = \psi(\hat{P}_\lambda^{\text{gcv}}) \quad \text{and} \quad \hat{\psi}_\lambda^{\text{loo}} = \psi(\hat{P}_\lambda^{\text{loo}})$$

Outline

Problem setup

Distribution estimation

Functional estimation

Distribution estimation

Under i.i.d. sampling of (x_i, y_i) , $i = 1, \dots, n$ with

1. feature x_i decomposable into $x_i = \Sigma^{1/2} z_i$ where z_i contains i.i.d. entries with mean 0, variance 1 and finite 4+ moment, and max and min eigenvalues of Σ uniformly away from 0 and ∞ ,
2. response y_i with bounded 4+ moment,

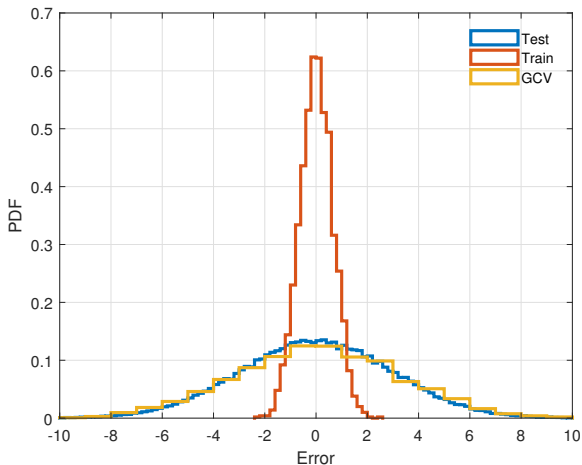
as $n, p \rightarrow \infty$ such that $p/n \rightarrow \gamma \in (0, \infty)$, almost surely

$$\hat{P}_\lambda^{\text{gcv}} \xrightarrow{d} P_\lambda \quad \text{and} \quad \hat{P}_\lambda^{\text{loo}} \xrightarrow{d} P_\lambda.$$

Remarks:

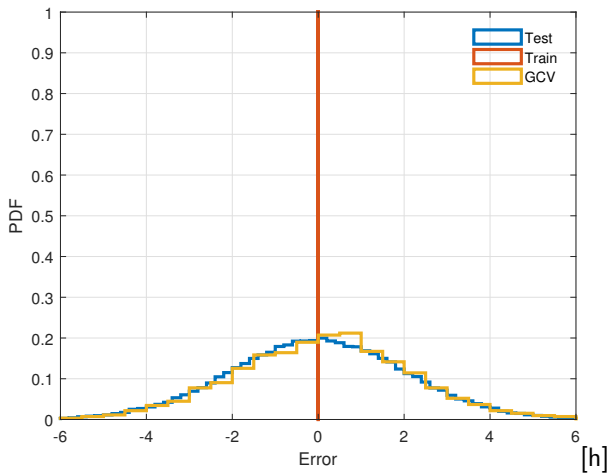
- Almost sure convergence with respect to the training data
- The regression function does not need to be linear in x
- Amazingly, this results also holds when $\lambda = 0$ (min-norm estimator)

Distribution estimation: illustration ($p < n$)



- $n = 2500$, $p = 2000$, $p/n = 0.8$
- $\lambda = 0$, i.e., least squares

Distribution estimation: illustration ($p > n$)



- $n = 2500$, $p = 5000$, $p/n = 2$
- $\lambda = 0$, i.e., the min-norm estimator, zero in-sample errors

Outline

Problem setup

Distribution estimation

Functional estimation

Linear functional estimation (pointwise)

- Let T_λ be a linear functional of the out-of-sample error distribution:

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

- Let $\hat{T}_\lambda^{\text{gcv}}$ and $\hat{T}_\lambda^{\text{loo}}$ be plug-in estimators from GCV and LOOCV:

$$\hat{T}_\lambda^{\text{gcv}} = \frac{1}{n} \sum_{i=1}^n t \left(\frac{y_i - x_i^T \hat{\beta}_\lambda}{1 - \text{tr}[L_\lambda]/n} \right) \quad \text{and} \quad \hat{T}_\lambda^{\text{loo}} = \frac{1}{n} \sum_{i=1}^n t \left(\frac{y_i - x_i^T \hat{\beta}_\lambda}{1 - [L_\lambda]_{ii}} \right)$$

For error functions $t : \mathbb{R} \rightarrow \mathbb{R}$

- that are continuous,
- have quadratic growth, i.e., there exist constants $a, b, c > 0$ such that $|t(z)| \leq az^2 + b|z| + c$ for any $z \in \mathbb{R}$,

as $n, p \rightarrow \infty$ with $p/n \rightarrow \gamma \in (0, \infty)$, almost surely

$$\hat{T}_\lambda^{\text{gcv}} \rightarrow T_\lambda \quad \text{and} \quad \hat{T}_\lambda^{\text{loo}} \rightarrow T_\lambda.$$

Linear functional estimation (uniform)

For error functions $t : \mathbb{R} \rightarrow \mathbb{R}$

1. that are differentiable,
2. have derivative with linear growth rate, i.e., there exist constants $g, h > 0$ such that $|t'(z)| \leq g|z| + h$ for any $z \in \mathbb{R}$

as $n, p \rightarrow \infty$ with $p/n \rightarrow \gamma \in (0, \infty)$ for any compact set Λ ,

$$\sup_{\lambda \in \Lambda} |\widehat{T}_{\lambda}^{\text{gcv}} - T_{\lambda}| \rightarrow 0 \quad \text{and} \quad \sup_{\lambda \in \Lambda} |\widehat{T}_{\lambda}^{\text{loo}} - T_{\lambda}| \rightarrow 0.$$

Remarks:

- Special case of $t(r) = r^2$ exploits bias-variance decomposition
- No bias-variance decomposition for general error functions and result requires a different proof technique via leave-one-out arguments
- Using uniformity arguments, the result can be extended for non-linear variational functionals (see paper for more details)

Discussion and future directions

Take-away from this work: empirical distributions of GCV and LOOCV track out-of-sample error distribution and a wide class of its functionals for ridge regression under proportional asymptotics framework

Key relation that we exploit:

$$y_i - x_i^\top \widehat{\beta}_{-i,\lambda} = \frac{y_i - x_i^\top \widehat{\beta}_\lambda}{1 - [L_\lambda]_{ii}} \approx \frac{y_i - x_i^\top \widehat{\beta}_\lambda}{1 - \text{tr}[L_\lambda]/n}$$
$$y_i - x_i^\top \widehat{\beta}_{-i,0} = \frac{[(XX^\top)^\dagger y]_i}{[(XX^\top)^\dagger]_{ii}} \approx \frac{[(XX^\top)^\dagger y]_i}{\text{tr}[(XX^\top)^\dagger]/n}$$

Going beyond ...

- Equivalences for ridge variants and other smoothers
- Finite sample analysis and rates of convergence
-

Thanks for listening!

Questions/comments/thoughts?