

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

Pratik Patil Alessandro Rinaldo Ryan J. Tibshirani

Carnegie Mellon University

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 $\hat{\beta}_\lambda$ denote the **ridge estimator**:

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

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

The key question that we ask in this paper is: 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

Out-of-sample error distribution and its functionals

- Let P_λ denote distribution of out-of-sample error of $\hat{\beta}_\lambda$, i.e., $P_\lambda = \mathcal{L}(y_0 - x_0^T \hat{\beta}_\lambda | X, y)$, where (x_0, y_0) is sampled indep from the same training distribution
- Let ψ denote a functional such that $P \mapsto \psi(P) \in \mathbb{R}$:
 - Linear functional**:

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

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

$$\text{loo}(\lambda) = \frac{1}{n} \sum_{i=1}^n (y_i - x_i^T \hat{\beta}_{-i,\lambda})^2$$

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

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

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

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

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

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

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

- Plug-in GCV and LOO estimators:

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

Distribution estimation

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

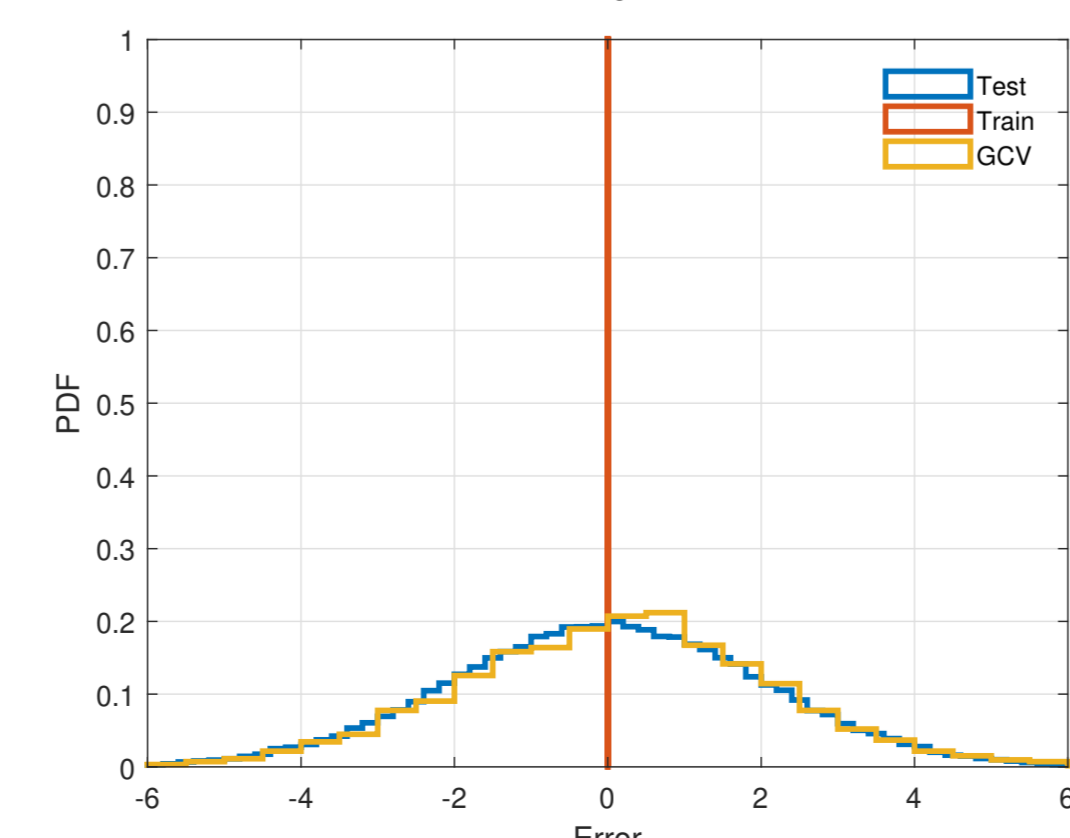
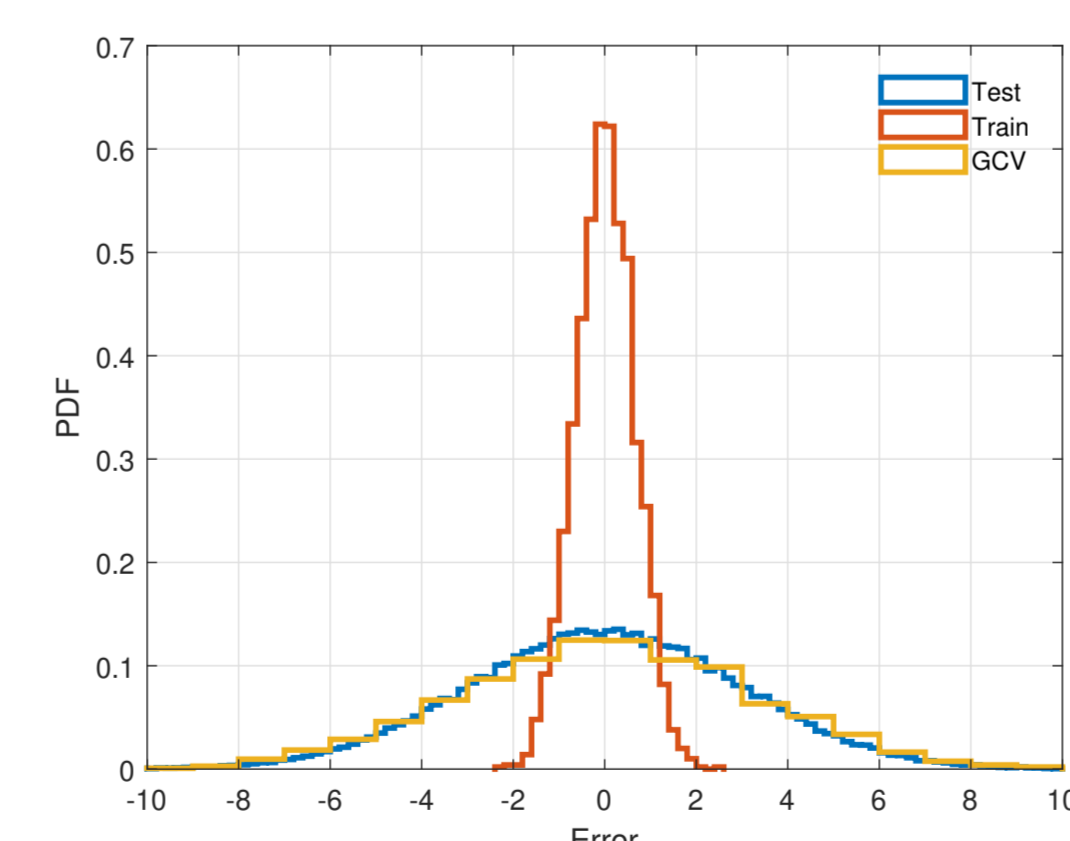
- feature x_i decomposable into $x_i = \Sigma^{1/2} z_i$ where z_i contains i.i.d. entries with mean 0, variance 1 and finite 4+ moment, and max and min eigenvalues of Σ uniformly away from 0 and ∞ ,
- response y_i with bounded 4+ moment,

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

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

- 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: numerical illustration



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

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

Linear functional estimation (pointwise)

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

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

- 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 $t: \mathbb{R} \rightarrow \mathbb{R}$

- that are differentiable,
- 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 ,

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

- 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 work

The main take-away from this work is: 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:

$$y_i - x_i^T \hat{\beta}_{-i,\lambda} = \frac{y_i - x_i^T \hat{\beta}_\lambda}{1 - [L_\lambda]_{ii}} \approx \frac{y_i - x_i^T \hat{\beta}_\lambda}{1 - \text{tr}[L_\lambda] / n}$$

$$y_i - x_i^T \hat{\beta}_{-i,0} = \frac{[(X X^T)^{\dagger} y]_i}{[(X X^T)^{\dagger}]_{ii}} \approx \frac{[(X X^T)^{\dagger} y]_i}{\text{tr}[(X X^T)^{\dagger}] / n}$$

Going beyond ...

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