Generalized equivalences between subsampling and ridge regularization

Pratik Patil¹ Jin-Hong Du²

¹University of California Berkeley ²Carnegie Mellon University

NeurIPS 2023

Over-parameterization and regularization

In the big data era, the success of machine learning and deep learning methods typically have much more parameters than the training samples.



 Optimizing such over-parameterized models requires different types of regularization.

Explicit and implicit regularization



explicit regularization



Explicit and implicit regularization



Explicit and implicit regularization



Ridge ensembles

► Ridge estimator: Consider a dataset D_n = {(x_j, y_j) : j ∈ [n]} containing i.i.d. vectors in ℝ^p × ℝ. The ridge estimator fitted on a subsampled dataset D_I is:

$$\widehat{\boldsymbol{\beta}}_{k}^{\lambda}(\mathcal{D}_{I}) = \operatorname*{argmin}_{\boldsymbol{\beta} \in \mathbb{R}^{p}} \sum_{j \in I} (y_{j} - \boldsymbol{x}_{j}^{\top} \boldsymbol{\beta})^{2} / k + \lambda \|\boldsymbol{\beta}\|_{2}^{2}, I \subseteq [n], |I| = k.$$
(1)

Ridge ensembles

► Ridge estimator: Consider a dataset D_n = {(x_j, y_j) : j ∈ [n]} containing i.i.d. vectors in ℝ^p × ℝ. The ridge estimator fitted on a subsampled dataset D_I is:

$$\widehat{\boldsymbol{\beta}}_{k}^{\boldsymbol{\lambda}}(\mathcal{D}_{I}) = \operatorname*{argmin}_{\boldsymbol{\beta} \in \mathbb{R}^{p}} \sum_{j \in I} (y_{j} - \boldsymbol{x}_{j}^{\top} \boldsymbol{\beta})^{2} / k + \boldsymbol{\lambda} \|\boldsymbol{\beta}\|_{2}^{2}, I \subseteq [n], |I| = k.$$
(1)

Ensemble ridge estimator: For λ ≥ 0, the ensemble estimator is then defined as:

$$\widetilde{\beta}_{k,M}^{\lambda}(\mathcal{D}_n; \{I_\ell\}_{\ell=1}^M) := \frac{1}{M} \sum_{\ell \in [M]} \widehat{\beta}_k^{\lambda}(\mathcal{D}_{I_\ell}),$$
(2)

where I_1, \ldots, I_M are samples from $\mathcal{I}_k := \{\{i_1, i_2, \ldots, i_k\}: 1 \leq i_1 < i_2 < \ldots < i_k \leq n\}$. The *full-ensemble* ridge estimator $\widetilde{\beta}_{k,\infty}^{\lambda}(\mathcal{D}_n)$ is obtained with $M \to \infty$.

Generalized risk

- Let $\beta_0 = \mathbb{E}[xx^{\top}]^{-1}\mathbb{E}[xy]$ be the best linear projection of *y* onto *x*
- Generalized risk. For a linear functional $L_{A,b}(\beta) = A\beta + b$, we study

$$R(\widehat{\boldsymbol{\beta}}; \boldsymbol{A}, \boldsymbol{b}, \boldsymbol{\beta}_0) = \frac{1}{\operatorname{nrow}(\boldsymbol{A})} \| L_{\boldsymbol{A}, \boldsymbol{b}}(\widehat{\boldsymbol{\beta}} - \boldsymbol{\beta}_0) \|_2^2,$$
(3)

under proportional asymptotics where $n, p, k \to \infty$, $p/n \to \phi$ and $p/k \to \psi$. Here, ϕ and ψ are the *data and subsample aspect ratios*, respectively.

Statistical learning problem	$L_{A,b}(\widehat{oldsymbol{eta}}-oldsymbol{eta}_0)$	A	b	nrow(A)
vector coefficient estimation	$\widehat{oldsymbol{eta}} - oldsymbol{eta}_0$	I_p	0	р
projected coefficient estimation	$\boldsymbol{a}^{ op}(\widehat{\boldsymbol{eta}}-\boldsymbol{eta}_0)$	$a^{ op}$	0	1
training error estimation	$X\widehat{eta} - y$	X	$-f_{\scriptscriptstyle NL}$	n
in-sample prediction	$X(\widehat{oldsymbol{eta}}-oldsymbol{eta}_0)$	X	0	n
out-of-sample prediction	$oldsymbol{x}_0^ op \widehat{oldsymbol{eta}} - y_0$	$x_0^ op$	$-\epsilon_0$	1

Motivation and summary of results

Data assumptions. Each feature vector x_i for $i \in [n]$ can be decomposed as $x_i = \Sigma^{1/2} z_i$, where $z_i \in \mathbb{R}^p$ contains i.i.d. entries z_{ij} for $j \in [p]$ with mean 0, variance 1, and bounded $4 + \mu$ moments for some $\mu > 0$. Response distribution: Each response variable y_i for $i \in [n]$ has mean 0, and bounded $4 + \mu$ moments.

Table: Comparison with related work. " \checkmark " indicates a partial equivalence result connecting the *optimal* prediction risk of the ridge predictor and the full ridgeless ensemble.

	Туре о	Type of equivalence results			Type of data assumptions			
	pred. risk	gen. risk	estimator	response	feature	lim. spectrum		
Lejeune 2020	√°			linear	isotropic Gaussian	exists		
Patil 2022	√°			linear	isotropic RMT	exists		
Du 2023	\checkmark			linear	anisotropic RMT	exists		
This work	\checkmark	\checkmark	\checkmark	arbitrary	anisotropic RMT	need not exist		

Summary of results

- Risk equivalences. We establish asymptotic equivalences of the full-ensemble ridge estimators at different ridge penalties λ and subsample ratios ψ along specific paths in the (λ, ψ)-plane for a variety of generalized risk functionals.
- Structural equivalences. We establish structural equivalences for linear functionals of the ensemble ridge estimators that hold for all ensemble sizes.
- Equivalence implications. The prediction risk of an optimally tuned ridge estimator is monotonically increasing in p/n under mild regularity conditions.
- Generality of equivalences. The results apply to arbitrary responses with bounded 4 + μ moments, as well as features with general covariance structures.

Asymptotic equivalence

- Let A_p and B_p be sequences of (additively) conformable matrices of arbitrary dimensions (including vectors and scalars).
- ▶ We say that A_p and B_p are *asymptotically equivalent*, denoted as $A_p \simeq B_p$, if $\lim_{p\to\infty} |\operatorname{tr}[C_p(A_p - B_p)]| = 0$ almost surely for any sequence of random matrices C_p with bounded trace norm that are (multiplicatively) conformable and independent of A_p and B_p .
- Note that for sequences of scalar random variables, the definition simply reduces to the typical almost sure convergence of sequences of random variables involved.

Generalized risk equivalences

Equivalence paths. Given $\phi \in (0, \infty)$ and $\bar{\psi} \in [\phi, \infty]$, our statement of equivalences between different ensemble estimators is defined through certain paths characterized by two endpoints $(0, \bar{\psi})$ and $(\bar{\lambda}, \phi)$. Let H_p be the empirical spectral distribution of Σ : $H_p(r) = p^{-1} \sum_{i=1}^p \mathbb{1}_{\{r_i \leq r\}}$, where r_i 's are the eigenvalues of Σ . Consider the following system of equations in $\bar{\lambda}$ and v:

$$\frac{1}{v} = \bar{\lambda} + \phi \int \frac{r}{1 + vr} \, \mathrm{d}H_p(r), \quad \text{and} \quad \frac{1}{v} = \bar{\psi} \int \frac{r}{1 + vr} \, \mathrm{d}H_p(r). \tag{4}$$

Now, define a path $\mathcal{P}(\bar{\lambda}; \phi, \bar{\psi})$ that passes through the endpoints $(0, \bar{\psi})$ and $(\bar{\lambda}, \phi)$:

$$\mathcal{P}(\bar{\lambda};\phi,\bar{\psi}) = \left\{ (1-\theta) \cdot (\bar{\lambda},\phi) + \theta \cdot (0,\bar{\psi}) \mid \theta \in [0,1] \right\}.$$
(5)

Generalized risk equivalences: illustration

Theorem 1. For any $\bar{\psi} \in [\phi, +\infty]$, let $\bar{\lambda}$ be as defined in (4). Then, for any pair of (λ_1, ψ_1) and (λ_2, ψ_2) on the path $\mathcal{P}(\bar{\lambda}; \phi, \bar{\psi})$ as defined in (5), the generalized risk functionals (3) of the full-ensemble estimator are asymptotically equivalent:

$$R(\widehat{\beta}_{\lfloor p/\psi_1 \rfloor,\infty}^{\lambda_1}; \boldsymbol{A}, \boldsymbol{b}, \boldsymbol{\beta}_0) \simeq R(\widehat{\beta}_{\lfloor p/\psi_2 \rfloor,\infty}^{\lambda_2}; \boldsymbol{A}, \boldsymbol{b}, \boldsymbol{\beta}_0).$$
(6)



Structural equivalences

Theorem 3. For any $\bar{\psi} \in [\phi, +\infty]$, let $\bar{\lambda}$ be as in (4). Then, for any $M \in \mathbb{N} \cup \{\infty\}$ and any pair of (λ_1, ψ_1) and (λ_2, ψ_2) on the path (5), the *M*-ensemble estimators are asymptotically equivalent:

 $\widehat{\beta}_{\lfloor p/\psi_1 \rfloor, M}^{\lambda_1} \simeq \widehat{\beta}_{\lfloor p/\psi_2 \rfloor, M}^{\lambda_2}, \qquad \forall (\lambda_1, \psi_1), (\lambda_2, \psi_2) \in \mathcal{P}(\bar{\lambda}; \phi, \bar{\psi}).$ (7)

Data-dependent paths. For any $M \in \mathbb{N} \cup \{\infty\}$, let $\overline{\lambda}_n$ be the value that satisfies the following equation in ensemble ridgeless and ridge gram matrices:

$$\frac{1}{M}\sum_{\ell=1}^{M}\frac{1}{k}\operatorname{tr}\left[\left(\frac{1}{k}\boldsymbol{L}_{I_{\ell}}\boldsymbol{X}\boldsymbol{X}^{\top}\boldsymbol{L}_{I_{\ell}}\right)^{+}\right] = \frac{1}{n}\operatorname{tr}\left[\left(\frac{1}{n}\boldsymbol{X}\boldsymbol{X}^{\top} + \bar{\boldsymbol{\lambda}}_{n}\boldsymbol{I}_{n}\right)^{-1}\right].$$
(8)

Define the data-dependent path $\mathcal{P}_n = \mathcal{P}(\overline{\lambda}_n; \phi_n, \overline{\psi}_n)$. Theorems 1 & 3 hold with \mathcal{P}_n .

Structural equivalences: illustration



Implications: Monotonicity of optimal ridge

- Many common methods, such as ridgeless or lassoless predictors, exhibit non-monotonic behavior in the sample size or the limiting aspect ratio.
- An open problem raised by Nakkiran et al. (2021) asks whether the prediction risk of ridge regression with optimal ridge penalty λ* is monotonically increasing in the data aspect ratio φ = p/n.
- Our equivalences imply that the prediction risk of an optimally-tuned ridge estimator is monotonically increasing in the data aspect ratio under mild regularity conditions.
- Under proportional asymptotics, our result settles a recent open question raised by Conjecture 1 of Nakkiran et al. (2021) concerning the monotonicity of optimal ridge regression under anisotropic features and general data models while maintaining a regularity condition that preserves the linearized signal-to-noise ratios across regression problems.

Implications of equivalences: illustration

Theorem 6. Let $k, n, p \to \infty$ such that $p/n \to \phi \in (0, \infty)$ and $p/k \to \psi \in [\phi, \infty]$. Then, for $A = \Sigma^{1/2}$ and b = 0, the optimal risk of the ridgeless ensemble, $\min_{\psi \ge \phi} \mathscr{R}(0; \phi, \psi)$, is monotonically increasing in ϕ . Consequently, the optimal risk of the ridge predictor, $\min_{\ge 0} \mathscr{R}(;\phi,\phi)$, is also monotonically increasing in ϕ .



Equivalences for random features

Conjecture 7. Define $\phi_n = p/n$. Let $k \le n$ be the subsample size and denote by $\overline{\psi}_n = p/k$. Suppose φ satisfies certain regularity conditions. For any $M \in \mathbb{N} \cup \{\infty\}$, let $\overline{\lambda}_n$ be the value that satisfies

$$\frac{1}{M}\sum_{\ell=1}^{M}\frac{1}{k}\operatorname{tr}\left[\left(\frac{1}{k}\varphi(\boldsymbol{L}_{I_{\ell}}\boldsymbol{X}\boldsymbol{F}^{\top})\varphi(\boldsymbol{L}_{I_{\ell}}\boldsymbol{X}\boldsymbol{F}^{\top})^{\top}\right)^{+}\right] = \frac{1}{n}\operatorname{tr}\left[\left(\frac{1}{n}\varphi(\boldsymbol{X}\boldsymbol{F}^{\top})\varphi(\boldsymbol{X}\boldsymbol{F}^{\top})^{\top} + \bar{\lambda}_{n}\boldsymbol{I}_{n}\right)^{-1}\right].$$

Define the data-dependent path $\mathcal{P}_n = \mathcal{P}(\bar{\lambda}_n; \phi_n, \bar{\psi}_n)$. Then similar equivalences continue to hold along \mathcal{P}_n .



Equivalences for kernel features

Conjecture 8. Define $\phi_n = p/n$. Suppose the kernel *K* satisfies certain regularity conditions. Let $k \le n$ be the subsample size and denote by $\overline{\psi}_n = p/k$. For any $M \in \mathbb{N} \cup \{\infty\}$, let $\overline{\lambda}_n$ be a solution to

$$\frac{1}{M}\sum_{\ell=1}^{M} \operatorname{tr}\left[\boldsymbol{K}_{I_{\ell}}^{+}\right] = \operatorname{tr}\left[\left(\boldsymbol{K}_{[n]} + \frac{n}{p}\bar{\lambda}_{n}\boldsymbol{I}_{n}\right)^{-1}\right]$$

Define the data-dependent path $\mathcal{P}_n = \mathcal{P}(\bar{\lambda}_n; \phi_n, \bar{\psi}_n)$. Then similar equivalences continue to hold along \mathcal{P}_n .

