

# Subsample Ridge Ensembles: Equivalences and Generalized Cross-Validation

Jin-Hong Du<sup>1\*</sup> Pratik Patil<sup>2\*</sup> Arun Kumar Kuchibhotla<sup>1</sup>

<sup>1</sup>Department of Statistics and Data Science, Carnegie Mellon University

<sup>2</sup>Department of Statistics, University of California, Berkeley

\*equal contribution

July 2023

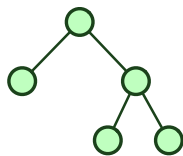
**Carnegie  
Mellon  
University**

**Berkeley**  
UNIVERSITY OF CALIFORNIA

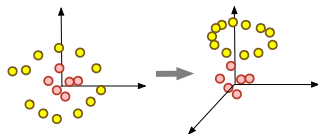
# Regularization

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

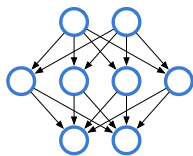
Random forest



Kernel method



Neural network



- I Optimizing such overparameterized models requires different types of regularization.

# Explicit and implicit regularization

implicit regularization



explicit regularization

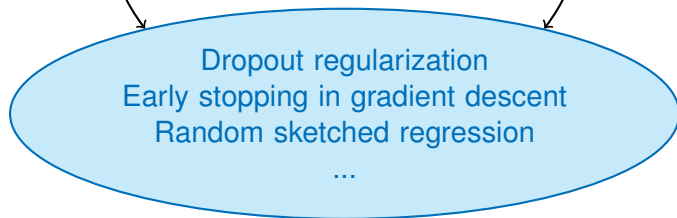


# Explicit and implicit regularization

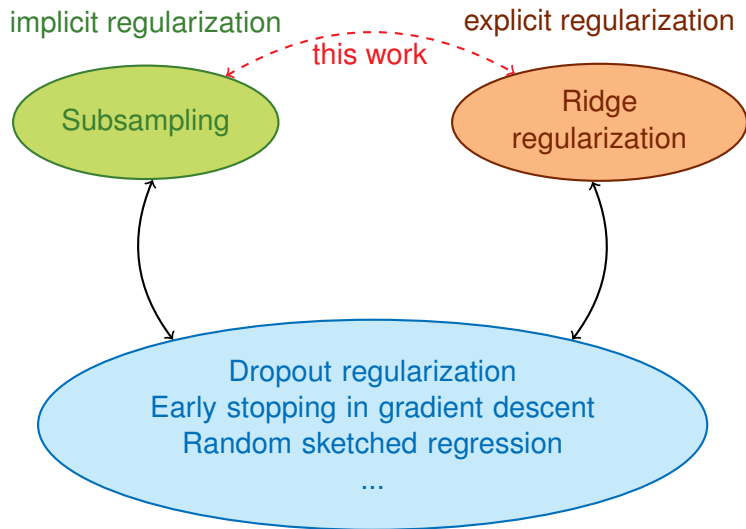
implicit regularization



explicit regularization



# Explicit and implicit regularization



# Ridge ensembles

- I **Ridge estimator:** Let  $D_n = \{(\mathbf{x}_j, y_j) \in \mathbb{R}^p \times \mathbb{R} : j \in [n]\}$  denote a dataset. The ridge estimator fitted on subsampled dataset  $D_I$  with  $I \subseteq [n]; |I| = k$  is defined as:

$$\mathbf{b}_k(D_I) = \underset{\mathbb{R}^p}{\operatorname{argmin}} \frac{1}{k} \sum_{j \in I} (y_j - \langle \mathbf{x}_j, \mathbf{b} \rangle)^2 + k \|\mathbf{b}\|_2^2$$

# Ridge ensembles

- I **Ridge estimator:** Let  $D_n = \{(\mathbf{x}_j; y_j) \in \mathbb{R}^p \times \mathbb{R} : j \in [n]\}$  denote a dataset. The ridge estimator fitted on subsampled dataset  $D_I$  with  $I \subseteq [n]; |I| = k$  is defined as:

$$b_k(D_I) = \underset{\mathbf{x} \in \mathbb{R}^p}{\operatorname{argmin}} \frac{1}{k} \sum_{j \in I} (y_j - \langle \mathbf{x}, \mathbf{x}_j \rangle)^2 + k \lambda^2$$

- I **Ensemble ridge estimator:**

$$e_{k;M}(D_n; \{I_i\}_{i=1}^M) := \frac{1}{M} \sum_{i=1}^M b_k(D_{I_i});$$

with  $I_1, \dots, I_M \subseteq [n]; |I_i| = k; 1 \leq i_1 < \dots < i_k \leq n$ . The *full-ensemble* ridge estimator is defined by letting  $M = \lfloor n/k \rfloor$ .

# Risk equivalence

**Conditional prediction risk:** The goal is to quantify and estimate the prediction risk:

$$R_{k;M} := \mathbb{E}_{(\mathbf{x};y)}[(y - \mathbf{x}^T \mathbf{e}_{k;M})^2 \mid D_n; fI \cdot \mathcal{G}_{=1}^M]; \quad (1)$$

under proportional asymptotics where  $n; p; k \rightarrow \infty$ ,  $p/n \rightarrow \rho$  and  $p/k \rightarrow s$ . Here,  $\rho$  and  $s$  are the **data** and **subsample** aspect ratios, respectively.



# Risk equivalence

- I As  $p \rightarrow n$  and  $p \rightarrow k$ , the prediction risk in the full ensemble ( $M = 1$ ) converges:

$$R_{k;1} \xrightarrow{\text{a.s.}} R_{k;1}(\cdot; s):$$

- I For  $\alpha = 1$ , the risk profile as a function of  $(\cdot; s)$  is shown in the figure in the log-log scale.

# Risk equivalence

- I As  $p=n$  and  $p=k$ , the prediction risk in the full ensemble ( $M=1$ ) converges:

$$R_{k;1} \stackrel{\text{a.s.}}{\rightarrow} R_{k;1}(\lambda; s):$$

- I For  $\lambda=1$ , the risk profile as a function of  $(\lambda; s)$  is shown in the figure in the log-log scale.

- I Risk equivalence (Theorem 2.3):

$$\min_s R_1^0(\lambda; s) = \min_{\lambda} R_1(\lambda; s) = \min_s R_1(\lambda; s):$$

|-----{z-----}      |-----{z-----}      |-----{z-----}

opt. ridgeless ensemble      opt. ridge predictor      opt. ridge ensemble

# Risk equivalence

- I As  $p=n$  and  $p=k$ , the prediction risk in the full ensemble ( $M=1$ ) converges:

$$R_{k;1} \stackrel{\text{a.s.}}{\rightarrow} R_{k;1}(\cdot; s):$$

- I For  $\lambda=1$ , the risk profile as a function of  $(\cdot; s)$  is shown in the figure in the log-log scale.
- I Implication: the implicit regularization provided by the subsample ensemble (a larger  $s$ , or a smaller  $k$ ) amounts to adding more explicit ridge regularization (a larger  $\lambda$ ).

# Generalized cross-validation for ridge ensembles

- I Beyond quantitative analysis, how can one pick  $(\lambda; s)$  to minimize the prediction risk?

# Generalized cross-validation for ridge ensembles

- I Beyond quantitative analysis, how can one pick  $(\lambda; s)$  to minimize the prediction risk?
- I For ordinary ridge ( $M = 1$  or  $k = n$ ), the **generalized cross-validation (GCV)** estimator is known to be consistent.

# Generalized cross-validation for ridge ensembles

- I Beyond quantitative analysis, how can one pick  $(k; s)$  to minimize the prediction risk?
- I For ordinary ridge ( $M = 1$  or  $k = n$ ), the **generalized cross-validation (GCV)** estimator is known to be consistent.
- I For general  $M$ , the GCV estimator is defined as

$$\text{gcv}_{k;M} = \frac{T_{k;M}}{D_{k;M}}$$

training error  
degree of freedom correction

# Generalized cross-validation for ridge ensembles

- I Beyond quantitative analysis, how can one pick  $(k; s)$  to minimize the prediction risk?
- I For ordinary ridge ( $M = 1$  or  $k = n$ ), the **generalized cross-validation (GCV)** estimator is known to be consistent.
- I For general  $M$ , the GCV estimator is defined as

$$\text{gcv}_{k;M} = \frac{T_{k;M}}{D_{k;M}} = \frac{\sum_{i=1}^M \sum_{j=1}^M (y_i - \mathbf{x}_i^T \mathbf{e}_{k;M})^2}{\left( \sum_{i=1}^M \sum_{j=1}^M \text{tr}(\mathbf{S}_{k;M}) \right)^2};$$

where  $\mathbf{S}_{k;M} = \frac{1}{M} \sum_{i=1}^M \mathbf{X}_i (\mathbf{X}_i^T \mathbf{X}_i + k \mathbf{I}_p)^{-1} \mathbf{X}_i^T$  is the smoothing matrix that represents the degree of freedom.

# Generalized cross-validation for ridge ensembles

- I Beyond quantitative analysis, how can one pick  $(k; s)$  to minimize the prediction risk?
- I For ordinary ridge ( $M = 1$  or  $k = n$ ), the **generalized cross-validation (GCV)** estimator is known to be consistent.
- I For general  $M$ , the GCV estimator is defined as

$$\text{gcv}_{k;M} = \frac{T_{k;M}}{D_{k;M}} = \frac{\frac{1}{j_{1:M}^j} \sum_{i=1}^M (y_i - \mathbf{x}_i^T \mathbf{e}_{k;M})^2}{(1 - \frac{1}{j_{1:M}^j} \text{tr}(\mathbf{S}_{k;M}))^2};$$

where  $\mathbf{S}_{k;M} = \frac{1}{M} \sum_{i=1}^M \mathbf{X}_i (\mathbf{X}_i^T \mathbf{X}_i + \mathbf{I}_p)^{-1} \mathbf{X}_i^T$  is the smoothing matrix that represents the degree of freedom.

- I The GCV for full ensemble is defined by letting  $M$  tend to infinity.

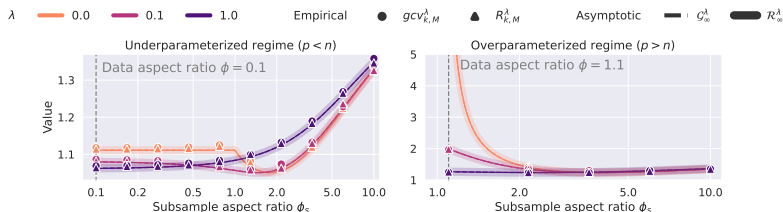


# Uniform consistency of GCV for full-ensemble ridge

I (Theorem 3.1, informal) For all  $\lambda > 0$ , we have

$$\max_{k \geq K_n} j \text{gcv}_{k;1} = R_{k;1} \text{ a.s. } 0:$$

I This allows selecting the optimal ensemble and subsample sizes in a data-dependent manner:



Coupled with the risk equivalence result, it suffices to fix  $\lambda$  and only tune the subsample size  $k$  or subsample aspect ratio  $\phi_s$ .

# Inconsistency on finite ensembles

- I (Proposition 3.3, informal) For ensemble size  $M = 2$ , ridge penalty  $\lambda = 0$ , and any  $\mathcal{Z}(\mathbf{0}; 1)$ ,

$$j \text{gcv}_{k;2}^0 \neq R_{k;2j}^0 \neq 0:$$

# Inconsistency on finite ensembles

- I (Proposition 3.3, informal) For ensemble size  $M = 2$ , ridge penalty  $\lambda = 0$ , and any  $\mathcal{L}(\mathbf{0}; 1)$ ,

$$j \text{gcv}_{k;2}^0 \quad R_{k;2j}^0 \quad \mathcal{E} \quad 0:$$

- I The bias scales as  $1/M$ , which is negligible for large  $M$ :

point - empirical GCV

line - theoretical risk

# Summary

- I This work [1] reveals the connections between the *implicit regularization* induced by subsampling and *explicit ridge regularization* for subsample ridge ensembles.
- I We establish the *uniform consistency* of GCV for full ridge ensembles.
- I We show that GCV can be *inconsistent* even for ridge ensembles when  $M = 2$ .
- I Future directions: bias correction of GCV for finite  $M$ ; extension to other metrics [2]; extension to other base predictors.

[1] Jin-Hong Du, Pratik Patil, and Arun Kumar Kuchibhotla. "Subsample Ridge Ensembles: Equivalences and Generalized Cross-Validation". In: *International Conference on Machine Learning* (2023)

[2] Pratik Patil and Jin-Hong Du. "Generalized equivalences between subsampling and ridge regularization". In: *arXiv preprint arXiv:2305.18496* (2023)