

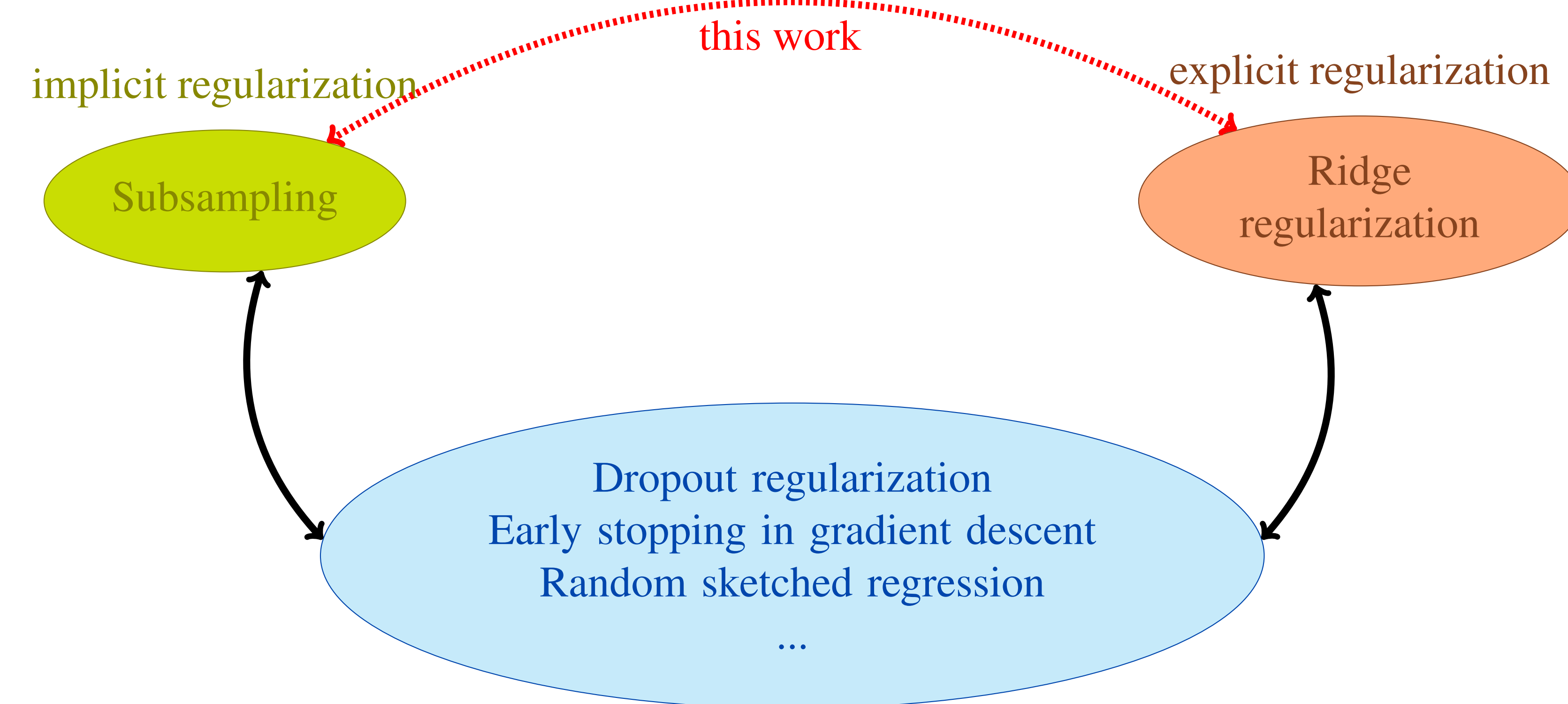
Subsample Ridge Ensembles: Equivalences and Generalized Cross-Validation

Jin-Hong Du^{1*} Pratik Patil^{2*} Arun Kumar Kuchibhotla¹

¹Department of Statistics and Data Science, Carnegie Mellon University

²Department of Statistics, University of California, Berkeley *equal contribution

Background



Ridge estimator: Let $\mathcal{D}_n = \{(\mathbf{x}_j, y_j) : j \in [n]\}$ denote a dataset containing i.i.d. random vectors in $\mathbb{R}^p \times \mathbb{R}$. The ridge estimator fitted on subsampled dataset \mathcal{D}_I is defined as:

$$\hat{\beta}_k^\lambda(\mathcal{D}_I) = \operatorname{argmin}_{\beta \in \mathbb{R}^p} \sum_{j \in I} (y_j - \mathbf{x}_j^\top \beta)^2 / k + \lambda \|\beta\|_2^2, \quad I \subseteq [n], |I| = k \quad (1)$$

Ensemble ridge estimator: For $\lambda \geq 0$, the ensemble estimator is then defined as:

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

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

Conditional prediction risk: The goal is to study the prediction risk:

$$R_{k,M}^\lambda := \mathbb{E}_{(\mathbf{x}, y)} [(y - \mathbf{x}^\top \tilde{\beta}_{k,M}^\lambda)^2 \mid \mathcal{D}_n, \{I_\ell\}_{\ell=1}^M], \quad (3)$$

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

Summary

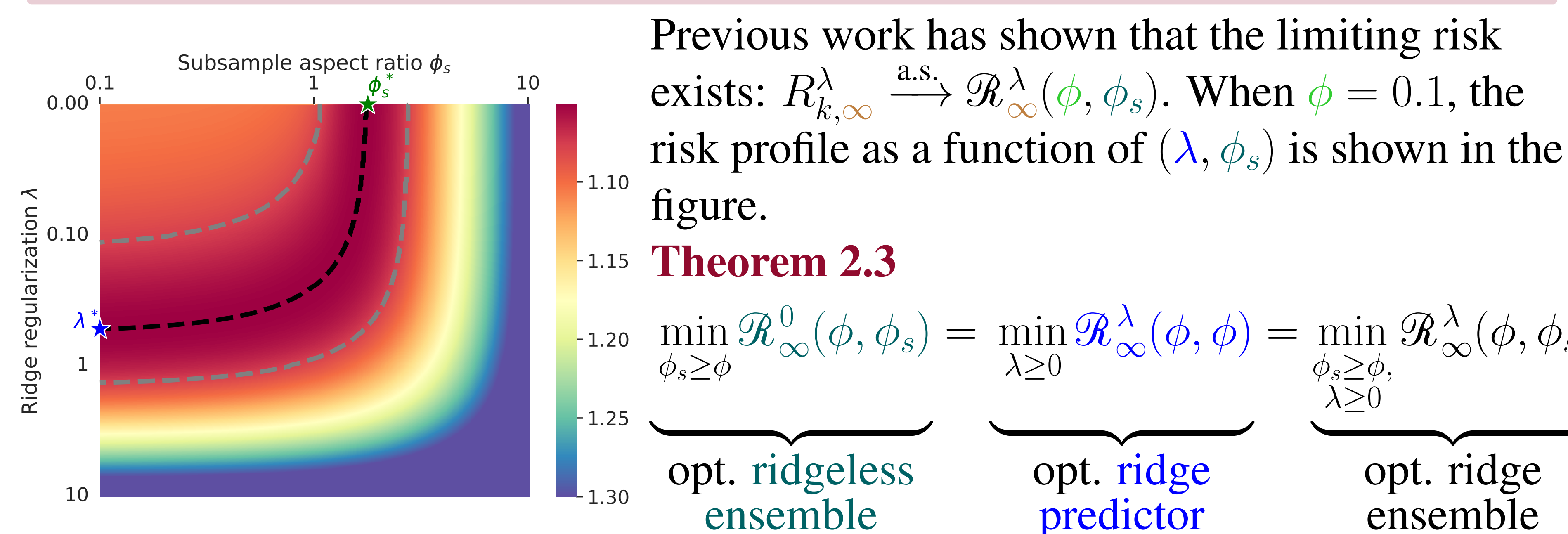
- **General risk equivalences.** We establish prediction risk equivalences between *implicit* regularization of subsampling and *explicit* ridge regularization for the subsample ridge ensemble. For any $\tau \geq 0$, we provide the set \mathcal{C}_τ of pairs (λ, ϕ_s) such that the risk of the full ridge ensemble with ridge penalty λ and subsample aspect ratio ϕ_s is equal to the risk of the ridge predictor with ridge penalty τ .

- **Uniform consistency of GCV.** For full ridge ensembles, we establish the uniform consistency of GCV across all possible subsample sizes k . Notably, this result is also applicable to the ridgeless regression ($\lambda = 0$). This enables tuning the subsample size in a data-dependent manner.
- **Finite-ensemble surprises.** Even though GCV is consistent for $M = 1$ and $M = \infty$, interestingly, this is the first paper that proves GCV *can* be inconsistent even for ridge ensembles when the ensemble size $M = 2$. Nevertheless, GCV is applicable for tuning subsample sizes, even with moderate ensemble sizes in practice.

Assumptions

- **Feature model:** $\mathbf{X} = \mathbf{Z}\Sigma^{1/2}$, where $\mathbf{Z} \in \mathbb{R}^{n \times p}$ contains i.i.d. entries with bounded $4 + \delta$ moments, and $\Sigma \in \mathbb{R}^{p \times p}$ has bounded eigenvalues and limiting spectral distribution.
- **Response model:** $\mathbf{y} = \mathbf{X}\beta_0 + \epsilon$, where $\beta_0 \in \mathbb{R}^p$ satisfies $\|\beta_0\|_2^2 \xrightarrow{\text{a.s.}} \rho^2$, and ϵ contains i.i.d. entries with variance σ^2 and bounded $4 + \delta$ moments. The limiting spectral distribution of β_0 's (squared) projection onto Σ exists.

Risk equivalence in the full ensemble



- Implication: the implicit regularization provided by the subsample ensemble (a larger ϕ_s , or a smaller k) amounts to adding more explicit ridge regularization (a larger λ).
- Usage: tuning ridge penalty λ for optimal ridge predictors ($\phi_s = \phi$) by tuning subsample aspect ratio ϕ_s for ridgeless ensembles ($\lambda = 0$).

Generalized Cross-Validation (GCV)

For general M , the GCV estimator is defined as

$$\text{gcv}_{k,M}^\lambda = \frac{T_{k,M}^\lambda}{D_{k,M}^\lambda} = \frac{\frac{1}{|I_{1:M}|} \sum_{i \in I_{1:M}} (y_i - \mathbf{x}_i^\top \tilde{\beta}_{k,M}^\lambda)^2}{(1 - |I_{1:M}|^{-1} \text{tr}(\mathbf{S}_{k,M}^\lambda))^2} \quad \begin{array}{l} \text{training error} \\ \text{degrees-of-freedom correction} \end{array}$$

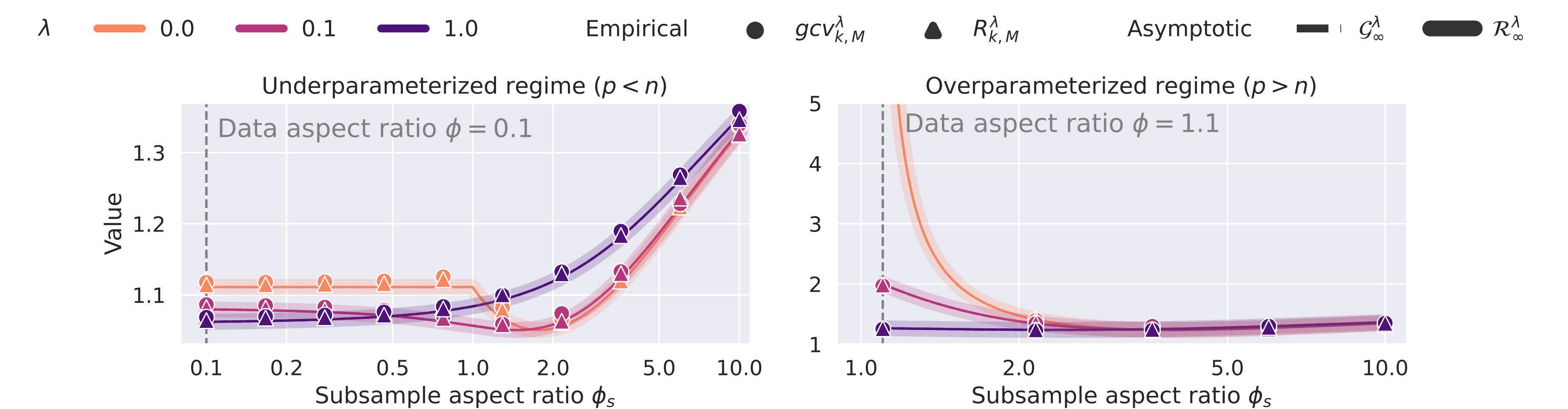
where $\mathbf{S}_{k,M}^\lambda = \frac{1}{M} \sum_{\ell=1}^M \mathbf{X}_{I_\ell} (\mathbf{X}_{I_\ell}^\top \mathbf{X}_{I_\ell} / k + \lambda \mathbf{I}_p)^+ \mathbf{X}_{I_\ell}^\top / k$ is the smoothing matrix that represents the degrees-of-freedom.

Theorem 3.1 For all $\lambda \geq 0$, we have

$$\max_{k \in \mathcal{K}_n} |\text{gcv}_{k,\infty}^\lambda - R_{k,\infty}^\lambda| \xrightarrow{\text{a.s.}} 0. \quad (4)$$

Coupled with the risk equivalence, we further have data-dependent tuning:

Corollary 3.2 $\text{gcv}_{k,\infty}^0 \xrightarrow{\text{a.s.}} \min_{\phi_s \geq \phi, \lambda \geq 0} \mathcal{R}_\infty^\lambda(\phi, \phi_s)$.

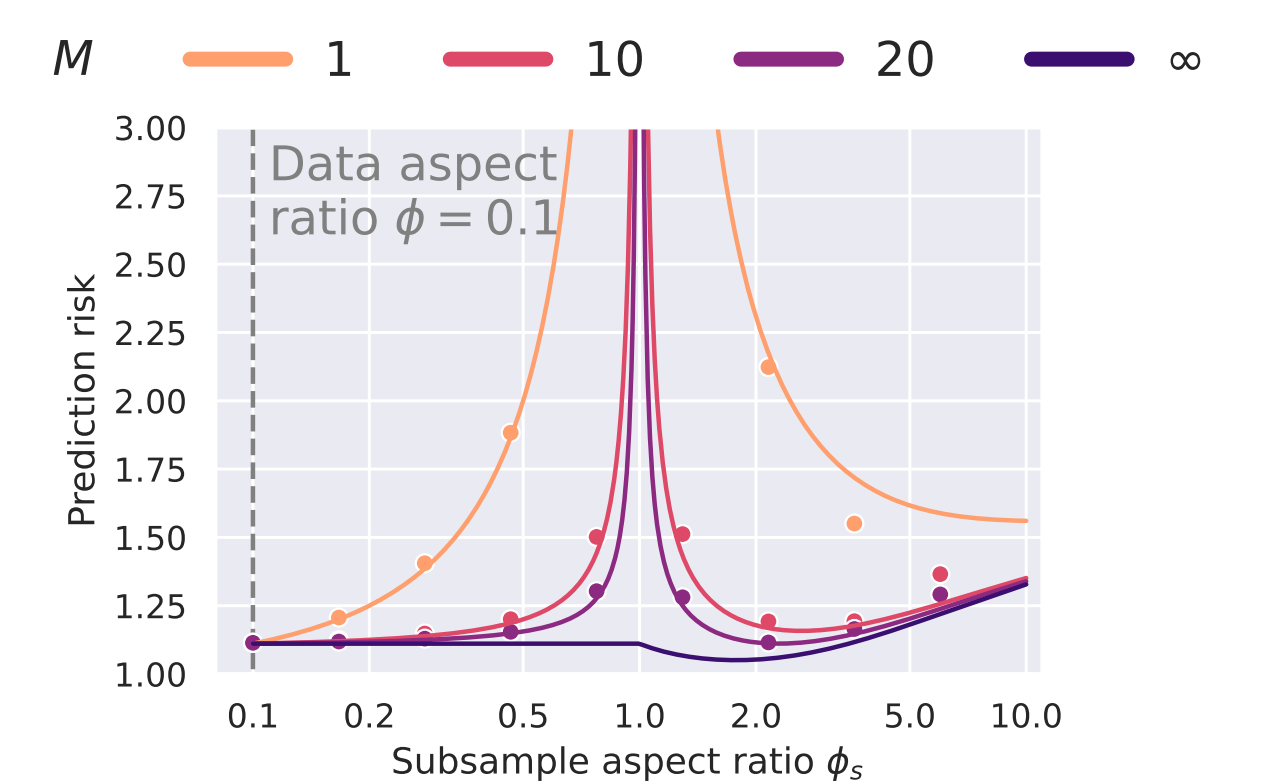


Inconsistency

Proposition 3.3 For ensemble size $M = 2$, ridge penalty $\lambda = 0$, and any $\phi \in (0, \infty)$,

$$|\text{gcv}_{k,2}^0 - R_{k,2}^0| \not\xrightarrow{P} 0.$$

The bias scales as $1/M$ and is negligible for large M .



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).