

Subsample Ridge Ensembles: Equivalences and Generalized Cross-Validation

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

Abstract

We study subsampling-based ridge ensembles in the proportional asymptotics regime, where the feature size grows proportionally with the sample size such that their ratio converges to a constant. By analyzing the squared prediction risk of ridge ensembles as a function of the explicit penalty λ and the limiting subsample aspect ratio ϕ_s (the ratio of the feature size to the subsample size), we characterize contours in the (λ, ϕ_s) -plane at any achievable risk. As a consequence, we prove that the risk of the optimal full ridgeless ensemble (fitted on all possible subsamples) matches that of the optimal ridge predictor. In addition, we prove strong uniform consistency of generalized cross-validation (GCV) over the subsample sizes for estimating the prediction risk of ridge ensembles. This allows for GCV-based tuning of full ridgeless ensembles without sample splitting and yields a predictor whose risk matches optimal ridge risk.

1. Introduction

Ensemble methods (Breiman, 1996) are widely used in various real-world applications in statistics and machine learning. They combine a collection of weak predictors to produce more stable and accurate predictions. One notable example of an ensemble method is bagging (bootstrap aggregating) (Breiman, 1996; Bühlmann & Yu, 2002). Bagging involves averaging base predictors that are fitted on different subsampled datasets and has been shown to stabilize the prediction and reduce the predictive variance (Bühlmann & Yu, 2002). In this paper, we study such a class of ensemble methods that fit each base predictor independently using a different subsampled dataset of the full training data. As a prototypical base predictor, we focus on *ridge regression*

^{*}Equal contribution ¹Department of Statistics and Data Science, Carnegie Mellon University, Pittsburgh, PA 15213, USA. ²Department of Statistics, University of California, Berkeley, CA 94720, USA. Correspondence to: Jin-Hong <jin-hongd@andrew.cmu.edu>, Pratik <pratikpatil@berkeley.edu>.

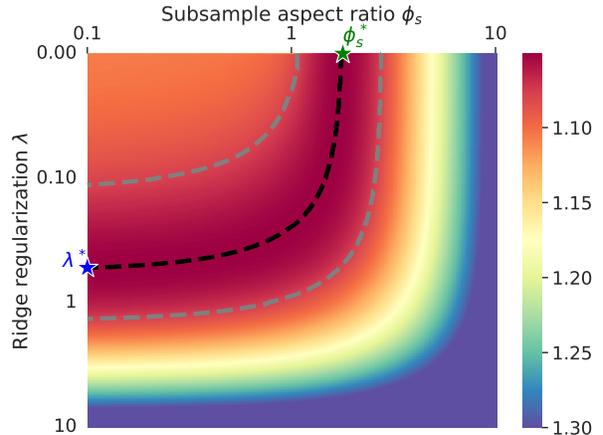


Figure 1. Heat map of the asymptotic prediction risk landscape of full ridge ensembles as the number of observation n , the subsample size k , and the feature dimension p tend to infinity, for varying regularization parameters λ and limiting subsample aspect ratio parameters $\phi_s = \lim p/k$. The data $(\mathbf{x}, y) \in \mathbb{R}^p \times \mathbb{R}$ is generated from a non-isotropic linear model $y = \mathbf{x}^\top \boldsymbol{\beta}_0 + \epsilon$ with $\phi = \lim p/n = 0.1$, where the features, the coefficients, and the residuals are distributed as $\mathbf{x} \sim \mathcal{N}(0, \boldsymbol{\Sigma}_{\text{AR1}})$, $\boldsymbol{\beta}_0 = \frac{1}{5} \sum_{j=1}^5 \mathbf{w}_{(j)}$, and $\epsilon \sim \mathcal{N}(0, 1)$, respectively. Here, the covariance matrix $(\boldsymbol{\Sigma}_{\text{AR1}})_{ij} = 0.5^{|i-j|}$, $\mathbf{w}_{(j)}$ is the top j th eigenvector of $\boldsymbol{\Sigma}_{\text{AR1}}$. The green and blue stars denote the risk of the optimal full-ensemble ridgeless predictor and the optimal ridge predictor without subsampling, respectively. The black dashed line denotes the set of (λ, ϕ_s) pairs that yield the same risk as (λ^*, ϕ) and $(0, \phi_s^*)$, while the gray dashed lines indicate the set of pairs that all result in the same sub-optimal risk.

(Hoerl & Kennard, 1970a;b), one of the most popular statistical methods. We refer readers to the “ridgefest” by Hastie (2020) for the history and review of ridge regression.

Ridge regression has recently attracted great interest, particularly the limiting case of zero regularization (where the regularization parameter tends to zero), termed “ridgeless” regression. In the underparameterized regime, the ridgeless predictor is ordinary least squares. However, in the overparameterized regime, it interpolates the training data and exhibits a peculiar risk behavior (Belkin et al., 2020; Bartlett et al., 2020; Hastie et al., 2022; Muthukumar et al., 2020). LeJeune et al. (2020); Patil et al. (2022a) have recently analyzed the statistical properties of the ensemble ridge and ridgeless predictors under proportional asymp-

otics. Under a linear model with the isotropic Gaussian covariate distribution, LeJeune et al. (2020) prove that the *full ensemble* (ensemble fitted on all possible subsampled datasets) of least squares predictors with optimal subsample size has the same risk as that of ridge predictor with optimal regularization. Under a more general but still isotropic covariate distribution, Patil et al. (2022a) prove similar risk equivalence of the optimized full ridgeless ensemble and the optimized ridge predictor.

These findings inspire two natural avenues to investigate.

(1) *Understanding the extent of risk equivalences.* As a curious experiment, one can empirically observe that a similar phenomenon to the one just mentioned appears to hold under quite general non-isotropic data models, as illustrated in Figure 1. We observe that the optimal ridgeless in the full ensemble (the green star) has the same prediction risk as the optimal ridge on the full data (the blue star). Furthermore, any pair of (λ, ϕ_s) on the black line achieves the same optimal risk. Such a relationship also extends to any other attainable risk value. For example, see the grey lines for (λ, ϕ_s) pairs that all achieve the same sub-optimal risk. This inspires our first investigation to establish risk equivalences between subsampling and ridge regression under general settings.

(2) *Overcoming limitations of split cross-validation.* Apart from its theoretical interest, the risk equivalences also suggest an alternative practical way to tune the ridge regularization parameter by tuning the subsample size. In terms of practical tuning of the ridge and ridgeless ensembles, Patil et al. (2022a) provide a split cross-validation method to estimate the prediction risk of ensembles with a fixed (finite) number of ensemble sizes and further prove that the split cross-validation consistently selects the best subsample size. The split cross-validation procedure has two disadvantages: (a) sample splitting introduces additional external randomness in the predictor; and (b) the reduced sample size, although asymptotically negligible, has significant finite sample effects, especially near the interpolation thresholds. This inspires our second investigation to address these limitations by considering generalized cross-validation (GCV) that does not require any sample splitting. The consideration of GCV as a viable estimator of the prediction risk for ridge ensembles stems from the observation that the ridge ensembles are also in fact linear smoothers.

1.1. Summary of Contributions

Below we provide a brief overview of our main results.

- **General risk equivalences.** We establish general equivalences between the subsample-optimized ridgeless ensemble, the optimal ridge predictor, and the optimal subsample ridge ensemble (see Theorem 2.3). In addition, for any

$\tau \geq 0$, we provide an exact characterization of the sets \mathcal{C}_τ of pairs (λ, ϕ_s) (the regularization parameter and the limiting subsample aspect ratio) such that the risk of the full ridge ensemble with ridge regularization λ and subsample aspect ratio ϕ_s is equal to the risk of the ridge predictor with ridge regularization τ . In essence, this amounts to showing that the implicit regularization of subsampling is the same as additional explicit ridge regularization.

- **Uniform consistency of GCV.** We establish the uniform consistency of GCV across all possible subsample sizes for full ridge ensembles with fixed regularization parameters (see Theorem 3.1). Notably, this result is also applicable to zero explicit regularization and covers the case of ridgeless regression. This finding enables tuning over the subsample size in a data-dependent manner, and in conjunction with Theorem 2.3, it implies that GCV tuning leads to a predictor with the same risk as the optimal ridge predictor (see Corollary 3.2).
- **Finite-ensemble surprises.** Even though GCV is consistent for the non-ensemble ridge and full-ensemble ridge predictors, interestingly, this is the first paper that proves GCV *can* be inconsistent even for ridge ensembles when the ensemble size is two (see Proposition 3.3). This finding is in contrast to other known results of GCV for ridge (see Section 1.2 for more details). Nevertheless, experiments on synthetic data and real-world single-cell multi-omic datasets demonstrate the applicability of GCV for tuning subsample sizes, even with moderate ensemble sizes (roughly of order 10).

1.2. Related Work

Ensembles and risk analysis. Ensemble methods are effective in combining weak predictors to build strong predictors in both regression and classification settings (Hastie et al., 2009). Early work on ensemble methods includes classical papers by Breiman (1996); Bühlmann & Yu (2002). There has been further work on the ensembles of smooth weak predictors (Buja & Stuetzle, 2006; Friedman & Hall, 2007), non-parametric estimators (Bühlmann & Yu, 2002; Loureiro et al., 2022), and classifiers (Hall & Samworth, 2005; Samworth, 2012). Under proportional asymptotics, d’Ascoli et al. (2020); Adlam & Pennington (2020a); Loureiro et al. (2022) study ensemble learning under random feature models. For ridge ensembles, Sollich & Krogh (1995); Krogh & Sollich (1997) derive risk asymptotics under Gaussian features. LeJeune et al. (2020) consider least squares ensembles obtained by subsampling such that the final subsampled dataset has more observations than the number of features. The asymptotic risk characterization for general data models has been derived by Patil et al. (2022a). Both of these works show the equivalence between the subsample optimized full ridgeless ensemble and the optimal ridge under isotropic

models. Our work significantly extends the scope of these results by characterizing risk equivalences for both optimal and suboptimal risks and for arbitrary feature covariance and signal structures. See the remarks after Theorem 2.3 for a detailed comparison,

Cross-validation and consistency. Cross-validation (CV) is arguably the most popular class of methods for model assessment and selection. Classical work on CV include: Allen (1974); Stone (1974; 1977); Geisser (1975), among others. We refer the reader to Arlot & Celisse (2010); Zhang & Yang (2015) for comprehensive surveys of different CV variants. In practice, k -fold CV is widely used with typical k being 5 or 10 (Hastie et al., 2009; Györfi et al., 2006), but such small values of k suffer from bias in high dimensions (Rad & Maleki, 2020). The extreme case of leave-one-out cross-validation (LOOCV) (when $k = n$) alleviates the bias issues in risk estimation, and various statistical consistency properties of LOOCV have been analyzed in recent years; see, e.g., Kale et al. (2011); Kumar et al. (2013); Obuchi & Kabashima (2016); Rad et al. (2020). Except for special cases, LOOCV is computationally expensive, and consequently, various approximations and their theoretical properties have been studied; see, e.g, Wang et al. (2018); Rad & Maleki (2020); Rad et al. (2020); Xu et al. (2019). Generalized cross-validation (GCV) is a sort of approximation for the “shortcut” leave-one-out formula (Hastie et al., 2009), originally studied for the fixed-X design setting for linear smoothers by Golub et al. (1979); Craven & Wahba (1979). The consistency of GCV in such a setting has been investigated in Li (1985; 1986; 1987). More recently, in the random- X setting, GCV has received considerable attention. In particular, consistency of GCV for ridge regression has been established in Adlam & Pennington (2020b); Hastie (2020); Patil et al. (2021; 2022c); Wei et al. (2022) under various data settings. Our work contributes to this body of work by analyzing GCV for subsampled ensemble ridge regression.

2. Subsample and Ridge Equivalences

We consider the standard supervised regression setting. Let $\mathcal{D}_n = \{(\mathbf{x}_j, y_j) : j \in [n]\}$ denote a dataset containing i.i.d. random vectors in $\mathbb{R}^p \times \mathbb{R}$, $\mathbf{X} \in \mathbb{R}^{n \times p}$ denote the feature matrix whose j -th row contains \mathbf{x}_j^\top , and $\mathbf{y} \in \mathbb{R}^n$ denote the response vector whose j -th entry contains y_j . For an index set $I \subseteq [n]$ of size k , let $\mathcal{D}_I = \{(\mathbf{x}_j, y_j) : j \in I\}$ be a subsampled dataset and let $\mathbf{L}_I \in \mathbb{R}^{n \times n}$ denote a diagonal matrix such that its j th diagonal entry is 1 if $j \in I$ and 0 otherwise. Noting that the feature matrix and response vector associated with \mathcal{D}_I are $\mathbf{L}_I \mathbf{X}$ and $\mathbf{L}_I \mathbf{y}$, respectively, the ridge estimator $\hat{\beta}_k^\lambda(\mathcal{D}_I)$ fitted on \mathcal{D}_I (containing k samples)

with regularization parameter $\lambda > 0$ can be expressed as:

$$\begin{aligned} \hat{\beta}_k^\lambda(\mathcal{D}_I) &= \operatorname{argmin}_{\beta \in \mathbb{R}^p} \sum_{j \in I} (y_j - \mathbf{x}_j^\top \beta)^2 / k + \lambda \|\beta\|_2^2 \\ &= (\mathbf{X}^\top \mathbf{L}_I \mathbf{X} / k + \lambda \mathbf{I}_p)^{-1} \mathbf{X}^\top \mathbf{L}_I \mathbf{y} / k. \end{aligned} \quad (1)$$

Letting $\lambda \rightarrow 0^+$, $\hat{\beta}_k^0(\mathcal{D}_I) := (\mathbf{X}^\top \mathbf{L}_I \mathbf{X} / k)^+ \mathbf{X}^\top \mathbf{L}_I \mathbf{y} / k$ becomes the so-called *ridgeless* estimator, where \mathbf{A}^+ denotes the Moore-Penrose inverse of matrix \mathbf{A} .

Ensemble estimator. To introduce the ensemble estimator, it helps to define the set of all k distinct elements from $[n]$ to be $\mathcal{I}_k := \{\{i_1, i_2, \dots, i_k\} : 1 \leq i_1 < i_2 < \dots < i_k \leq n\}$. Note that the cardinality of \mathcal{I}_k is $\binom{n}{k}$. 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 simple random samples from \mathcal{I}_k . The *full-ensemble* ridge estimator is the average of predictors fitted on all possible subsampled datasets:

$$\tilde{\beta}_{k,\infty}^\lambda(\mathcal{D}_n) := \frac{1}{|\mathcal{I}_k|} \sum_{I \in \mathcal{I}_k} \hat{\beta}_k^\lambda(\mathcal{D}_I) = \mathbb{E}[\hat{\beta}_k^\lambda(\mathcal{D}_I) \mid \mathcal{D}_n], \quad (3)$$

where the conditional expectation is taken with respect to the randomness of sampling from \mathcal{I}_k . Lemma A.1 shows that $\tilde{\beta}_{k,\infty}^\lambda(\mathcal{D}_n)$ is also almost surely equivalent to letting the ensemble size M tend to infinity in (2) conditioning on the full dataset \mathcal{D}_n , thus justifying the notation in (3). For simplicity, we drop the dependency on \mathcal{D}_n , $\{I_\ell\}_{\ell=1}^M$ and only write $\tilde{\beta}_{k,M}^\lambda, \tilde{\beta}_{k,\infty}^\lambda$, when it is clear from the context.

Prediction risk. We assess the performance of an M -ensemble predictor via conditional squared prediction risk:

$$R_{k,M}^\lambda := \mathbb{E}_{(\mathbf{x}, y)} [(y - \mathbf{x}^\top \tilde{\beta}_{k,M}^\lambda)^2 \mid \mathcal{D}_n, \{I_\ell\}_{\ell=1}^M], \quad (4)$$

where (\mathbf{x}, y) is an independent test point sampled from the distribution as \mathcal{D}_n . Note that the conditional risk $R_{k,M}^\lambda$ is a random variable that depends on both the dataset \mathcal{D}_n and the random samples I_ℓ , $\ell = 1, \dots, M$. For the full ensemble estimator $\tilde{\beta}_{k,\infty}^\lambda$, the conditional prediction risk is defined analogously, except the risk now only depends on \mathcal{D}_n :

$$R_{k,\infty}^\lambda := \mathbb{E}_{(\mathbf{x}, y)} [(y - \mathbf{x}^\top \tilde{\beta}_{k,\infty}^\lambda)^2 \mid \mathcal{D}_n]. \quad (5)$$

2.1. Data Assumptions

For our theoretical results, we work under a proportional asymptotics regime, in which the original *data aspect ratio* (p/n) converges to $\phi \in (0, \infty)$ as $n, p \rightarrow \infty$, and the *subsample aspect ratio* (p/k) converges to ϕ_s as $k, p \rightarrow \infty$. Note that because $k \leq n$, ϕ_s always lie in $[\phi, \infty]$. In addition, we impose two structural assumptions on the feature matrix and response vector as summarized in Assumptions 2.1 to 2.2, respectively.

Assumption 2.1 (Feature model). The feature matrix decomposes as $\mathbf{X} = \mathbf{Z}\Sigma^{1/2}$, where $\mathbf{Z} \in \mathbb{R}^{n \times p}$ contains i.i.d. entries with mean 0, variance 1, bounded moments of order $4 + \delta$ for some $\delta > 0$, and $\Sigma \in \mathbb{R}^{p \times p}$ is deterministic and symmetric with eigenvalues uniformly bounded between $r_{\min} > 0$ and $r_{\max} < \infty$. Let $\Sigma = \sum_{j=1}^p r_j \mathbf{w}_j \mathbf{w}_j^\top$ denote the eigenvalue decomposition, where $(r_j, \mathbf{w}_j), j \in [p]$, are pairs of associated eigenvalue and normalized eigenvector. We assume there exists a deterministic distribution H such that the empirical spectral distribution of Σ , $H_p(r) := p^{-1} \sum_{i=1}^p \mathbb{1}_{\{r_i \leq r\}}$, weakly converges to H , almost surely (with respect to \mathbf{X}).

Assumption 2.2 (Response model). The response vector decomposes as $\mathbf{y} = \mathbf{X}\beta_0 + \epsilon$, where $\beta_0 \in \mathbb{R}^p$ is an unknown signal vector with ℓ_2 -norm uniformly bounded and $\lim_{p \rightarrow \infty} \|\beta_0\|_2^2 = \rho^2$, and ϵ is an unobserved error vector independent of \mathbf{X} with mean 0, variance σ^2 , and bounded moment of order $4 + \delta$ for some $\delta > 0$. We assume there exists a deterministic distribution G such that the empirical distribution of β_0 's (squared) projection onto Σ 's eigenspace, $G_p(r) := \|\beta_0\|_2^{-2} \sum_{i=1}^p (\beta_0^\top \mathbf{w}_i)^2 \mathbb{1}_{\{r_i \leq r\}}$, weakly converges to G , almost surely (with respect to \mathbf{X}).

Assumptions 2.1 and 2.2 are standard in the study of the ridge and ridgeless regression under proportional asymptotics; see, e.g., [Hastie et al. \(2022\)](#); [Patil et al. \(2022b;a\)](#). It is possible to further relax both of these assumptions. Specifically, one can incorporate other feature models, e.g., random features ([Mei & Montanari, 2022](#)), and can allow for certain non-linearities in the regression function ([Bartlett et al., 2021](#)) for the response model. We leave these for future work.

2.2. Risk Equivalences

Under the above assumptions, Lemma A.2 from [Patil et al. \(2022a\)](#) implies that for every $M \geq 1$, the prediction risk $R_{k,M}^\lambda$ of the ridge and ridgeless predictors in the full ensemble converges to some deterministic limit $\mathcal{R}_M^\lambda(\phi, \phi_s)$ as $k, n, p \rightarrow \infty, p/n \rightarrow \phi$ and $p/k \rightarrow \phi_s$. When $\phi_s = \phi$ (e.g., $k = n$), the asymptotic risk $\mathcal{R}_M^\lambda(\phi, \phi)$ is equal to $\mathcal{R}_1^\lambda(\phi, \phi)$ of the ridge predictor on the full dataset \mathcal{D}_n for all $M \geq 1$, and we denote this risk simply by $\mathcal{R}_\infty^\lambda(\phi, \phi)$. To facilitate our discussion and for simplicity, Table 1 provides pointers to definitions of all important quantities used in the paper.

From a practical point of view, it is important to understand the least attainable risk that could be attained in the full ensemble. For the full ridge ensembles, we found that the explicit ridge regularization is unnecessary when considering optimal bagging and that the implicit regularization of ridgeless and subsampling suffices. The result below formalizes this empirical observation.

Theorem 2.3 (Optimal ridgeless ensemble vs optimal ridge).

Under Assumptions 2.1 and 2.2, for all $\phi \in (0, \infty)$, we have

$$\underbrace{\min_{\phi_s \geq \phi} \mathcal{R}_\infty^0(\phi, \phi_s)}_{\text{opt. ensemble and no ridge}} \stackrel{(a)}{=} \underbrace{\min_{\lambda \geq 0} \mathcal{R}_\infty^\lambda(\phi, \phi)}_{\text{no ensemble and opt. ridge}} \stackrel{(b)}{=} \underbrace{\min_{\substack{\phi_s \geq \phi, \\ \lambda \geq 0}} \mathcal{R}_\infty^\lambda(\phi, \phi_s)}_{\text{opt. ensemble and opt. ridge}}.$$

Further, if ϕ_s^* is the optimal subsample aspect ratio for ridgeless, and λ^* is the optimal ridge regularization with no subsampling, then for any $\theta \in [0, \lambda^*]$, full ridge ensemble with penalty parameter $\lambda = \lambda^* - \theta$ and subsample aspect ratio of $\phi_s = \phi + \theta(\phi_s^* - \phi)/\lambda^*$ also attains the optimal prediction risk.

In words, Theorem 2.3 says that optimizing subsample size (i.e. k) with the full ridgeless ensemble attains the same prediction risk as just optimizing the explicit regularization parameter (i.e., λ) of the ridge predictor. Further, both of them are the same as optimizing both k and λ . If one uses a lesser ridge penalty than needed for optimal prediction (i.e., uses $\lambda < \lambda^*$), then a full ensemble at a specific subsample aspect ratio $\phi_s = \phi + (1 - \lambda/\lambda^*)(\phi_s^* - \phi) > \phi$ can recover the remaining ridge regularization. In this sense, the implicit regularization provided by the ensemble amounts to adding more explicit ridge regularization. Similarly, one can supplement a sub-optimal implicit regularization of subsampling by adding explicit ridge regularization.

A special case of equivalence of (a) in Theorem 2.3 was previously formalized in [LeJeune et al. \(2020\)](#); [Patil et al. \(2022a\)](#) for isotropic covariates. Working with isotropic design helps their proof significantly, as the spectral distributions are the same for all p, n , and the closed-form expression of the asymptotic prediction risk can be derived analytically. However, in the general non-isotropic design, the asymptotic risk does not admit a closed-form expression, and one needs to account for this carefully.

General risk equivalences. Theorem 2.3 proves the risk equivalence of the ridge and full ensemble ridgeless when they attain minimum risk. Appendix A.3 shows a further risk equivalence in the full range, i.e., for any $\bar{\phi}_s \in [\phi, +\infty]$, there exists a $\bar{\lambda} \geq 0$ such that $\mathcal{R}_\infty^0(\phi, \bar{\phi}_s) = \mathcal{R}_\infty^{\bar{\lambda}}(\phi, \phi)$. Further, $\mathcal{R}_\infty^\lambda(\phi, \phi_s)$ remains constant as (λ, ϕ_s) varies on the line segment $(1 - \theta) \cdot (\bar{\lambda}, \phi) + \theta \cdot (0, \bar{\phi}_s)$, for all $\theta \in [0, 1]$.

A remarkable implication of Theorem 2.3 is that for a fixed dataset \mathcal{D}_n , one does not need to tune both the subsample size (i.e., k) and the ridge regularization parameter (i.e., λ), but it suffices to fix for example $\lambda = 0$ and only tune ϕ_s . Alternatively, one can also fix $k = n$ and just tune $\lambda \geq 0$, which was considered in [Patil et al. \(2021\)](#). Performing tuning over $\lambda \geq 0$ requires one to discretize an infinite interval, while tuning the subsample size for a fixed λ only requires searching over a finite grid varying from $k = 1$ to $k = n$. For this reason, we fix λ and focus on tuning over k

Variable	M -ensemble				Full ensemble			
	Finite-sample		Asymptotic		Finite-sample		Asymptotic	
Prediction risk	$R_{k,M}^\lambda$	(4)	$\mathcal{R}_M^\lambda(\phi, \phi_s)$	(21)	$R_{k,\infty}^\lambda$	(5)	$\mathcal{R}_\infty^\lambda(\phi, \phi_s)$	(21)
Test error	$\bar{R}_{k,M}^\lambda$	(7)	$\mathcal{R}_M^\lambda(\phi, \phi_s)$	(21)				
Training error	$T_{k,M}^\lambda$	(6)	$\mathcal{T}_M^\lambda(\phi, \phi_s)$	(32)	$T_{k,\infty}^\lambda$	(8)	$\mathcal{T}_\infty^\lambda(\phi, \phi_s)$	(15)
GCV denominator	$D_{k,M}^\lambda$	(12)			$D_{k,\infty}^\lambda$	(13)	$\mathcal{D}_\infty^\lambda(\phi, \phi_s)$	(26)
GCV estimator	$\text{gcv}_{k,M}^\lambda$	(11)			$\text{gcv}_{k,\infty}^\lambda$	(11)	$\mathcal{G}_\infty^\lambda(\phi, \phi_s)$	(16)

Table 1. Summary of notations and pointers to definitions of important empirical quantities used in this paper and their asymptotic limits.

in this paper. In the next section, we investigate the problem of tuning the subsample size in the full ensemble to achieve the minimum oracle risk via generalized cross-validation.

3. Generalized Cross-Validation

Suppose $\hat{f}(\cdot; \mathcal{D}_n) : \mathbb{R}^p \rightarrow \mathbb{R}$ is a predictor trained on \mathcal{D}_n . We call $\hat{f}(\cdot; \mathcal{D}_n)$ a linear smoother if $\hat{f}(\mathbf{x}; \mathcal{D}_n) = \mathbf{a}_x^\top \mathbf{y}$ for some vector \mathbf{a}_x that only depends on the design \mathbf{X} (and \mathbf{x}). Define the smoothing matrix $\mathbf{S} \in \mathbb{R}^{n \times n}$ with rows $\mathbf{a}_{x_1}^\top, \dots, \mathbf{a}_{x_n}^\top$, which in turn is only a function of \mathbf{X} . For any linear smoother, the generalized cross-validation (GCV) estimator of the prediction risk is defined to be $n^{-1} \|\mathbf{y} - \mathbf{S}\mathbf{y}\|_2^2 / (1 - n^{-1} \text{tr}(\mathbf{S}))^2$; see, e.g., Wasserman (2006, Section 5.3). The numerator of GCV is the training error, which typically is biased downwards, and the denominator attempts to account for such optimism of the predictor.

Ensemble GCV. Before we analyze GCV for the ridge ensemble, we first introduce some notations. Let $I_{1:M} := \cup_{\ell=1}^M I_\ell$ and $I_{1:M}^c := [n] \setminus I_{1:M}$. We define the *in-sample* training error and the *out-of-sample* test error of $\tilde{\beta}_{k,M}^\lambda$ as:

$$T_{k,M}^\lambda := \frac{1}{|I_{1:M}|} \sum_{i \in I_{1:M}} (y_i - \mathbf{x}_i^\top \tilde{\beta}_{k,M}^\lambda)^2, \quad (6)$$

$$\bar{R}_{k,M}^\lambda := \frac{1}{|I_{1:M}^c|} \sum_{i \in I_{1:M}^c} (y_i - \mathbf{x}_i^\top \tilde{\beta}_{k,M}^\lambda)^2. \quad (7)$$

Since the full ensemble estimator $\tilde{\beta}_{k,\infty}^\lambda$ uses all the data \mathcal{D}_n , its training error, denoted by $T_{k,\infty}^\lambda$, is simply:

$$T_{k,\infty}^\lambda := \frac{1}{n} \sum_{i \in [n]} (y_i - \mathbf{x}_i^\top \tilde{\beta}_{k,\infty}^\lambda)^2. \quad (8)$$

Since $I_{1:M} \xrightarrow{\text{a.s.}} [n]$ for any $n \in \mathbb{N}$ as $M \rightarrow \infty$, the notation $T_{k,\infty}^\lambda$ in (8) is justified as a limiting case of (6) (see Appendix A.1 for more details). Now, observe that a ridge ensemble is a linear smoother because $\mathbf{X}_{I_{1:M}} \tilde{\beta}_{k,M}^\lambda =$

$\mathbf{S}_{k,M}^\lambda \mathbf{y}_{I_{1:M}}$, where the smoothing matrix $\mathbf{S}_{k,M}^\lambda$ is given by:

$$\mathbf{S}_{k,M}^\lambda = \frac{1}{M} \sum_{\ell=1}^M \mathbf{X}_{I_\ell} (\mathbf{X}_{I_\ell}^\top \mathbf{X}_{I_\ell} / k + \lambda \mathbf{I}_p)^+ \mathbf{X}_{I_\ell}^\top / k. \quad (9)$$

Analogously, the smoothing matrix for $\tilde{\beta}_{k,\infty}^\lambda$ is given by:

$$\mathbf{S}_{k,\infty}^\lambda = \frac{1}{|\mathcal{I}_k|} \sum_{I \in \mathcal{I}_k} \mathbf{X} (\mathbf{X}^\top \mathbf{L}_I \mathbf{X} / k + \lambda \mathbf{I}_p)^+ \mathbf{X}^\top \mathbf{L}_I / k. \quad (10)$$

Thus, the GCV estimates for ridge predictors in the finite and full ensemble case are respectively given by:

$$\text{gcv}_{k,M}^\lambda = \frac{T_{k,M}^\lambda}{D_{k,M}^\lambda}, \quad \text{gcv}_{k,\infty}^\lambda = \frac{T_{k,\infty}^\lambda}{D_{k,\infty}^\lambda}, \quad (11)$$

where the denominators $D_{k,M}^\lambda$ and $D_{k,\infty}^\lambda$ are as follows:

$$D_{k,M}^\lambda := (1 - |I_{1:M}|^{-1} \text{tr}(\mathbf{S}_{k,M}^\lambda))^2, \quad (12)$$

$$D_{k,\infty}^\lambda := (1 - n^{-1} \text{tr}(\mathbf{S}_{k,\infty}^\lambda))^2. \quad (13)$$

3.1. Full-Ensemble Uniform Consistency

Let $\mathcal{K}_n \subset \{0, 1, \dots, n\}$ be a grid of subsample sizes that covers the full range of $[0, n]$ asymptotically in the sense that $\{k/n : k \in \mathcal{K}_n\}$ “converges” to the set $[0, 1]$ as $n \rightarrow \infty$. One simple choice is to set

$$\mathcal{K}_n = \{0, k_0, 2k_0, \dots, \lfloor n/k_0 \rfloor k_0\},$$

where the increment is $k_0 = \lfloor n^\nu \rfloor$ for some $\nu \in (0, 1)$. Here, we adopt the convention that when $k = 0$, the predictor reduces to a null predictor that always returns zero. Based on the definition above, we now present the uniform consistency results of the GCV estimator (11) for full ensembles when the ridge regularization parameter λ is fixed.

Theorem 3.1 (Uniform consistency of GCV). *Suppose Assumptions 2.1 and 2.2 hold. Then, for all $\lambda \geq 0$, we have*

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

as $n, p \rightarrow \infty$ such that $p/n \rightarrow \phi \in (0, \infty)$.

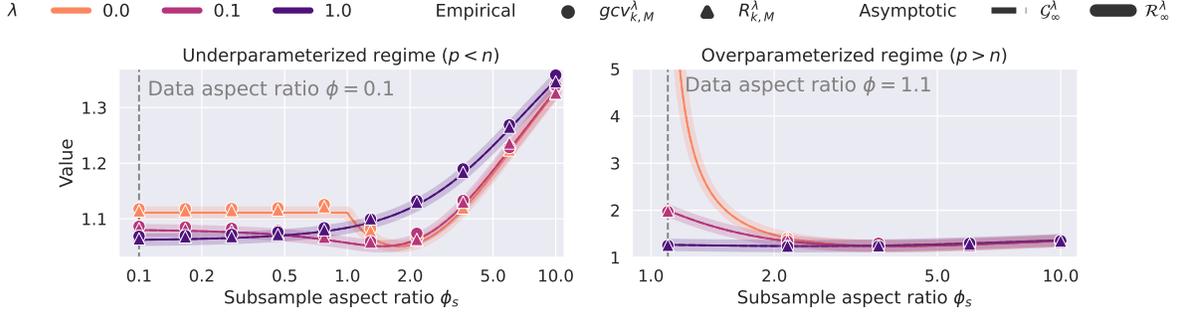


Figure 2. Asymptotic risk and GCV curves for full ridge ensembles, under model (M-AR1) when $\rho_{\text{AR1}} = 0.5$ and $\sigma^2 = 1$ with varying regularization parameters $\lambda \in \{0, 0.1, 1\}$ and subsample sizes $k = \lfloor p/\phi_s \rfloor$. The points denote finite-sample risks averaged over 50 dataset repetitions with an ensemble size of $M = 500$, with $n = \lfloor p/\phi \rfloor$ and $p = 500$. The left and the right panels illustrate the underparameterized and overparameterized cases with the limiting data aspect ratio $\phi = 0.1$ and $\phi = 1.1$, respectively.

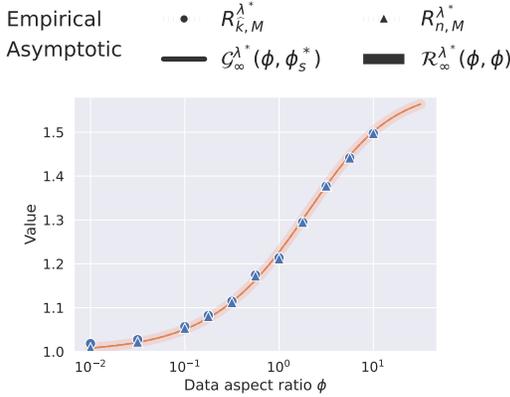


Figure 3. Asymptotic prediction risk curves with optimal tuned parameters λ^* and ϕ_s^* , under model (M-AR1) when $\rho_{\text{AR1}} = 0.5$, $\sigma^2 = 1$, for varying data aspect ratio ϕ . The curves represent the theoretical asymptotic GCV estimate in the full ensemble and the asymptotic risk of the optimal ridge predictors. The points represent the finite-sample risks of the best 500-ensemble ridgeless and the best ridge predictor averaged over 50 dataset repetitions, with $n = \lfloor p/\phi \rfloor$ and $p = 500$.

Theorem 3.1 shows the uniform consistency of GCV in the full ensemble for fixed subsample size k and ridge regularization parameter λ . The almost sure qualification in Theorem 3.1 is with respect to the entire training data (\mathbf{X}, \mathbf{y}) . An implication of Theorem 3.1 is that one can select the optimal subsample size in a data-dependent manner, i.e., selecting $\hat{k}^\lambda \in \operatorname{argmin}_{k \in \mathcal{K}_n} \operatorname{gcv}_{k,\infty}^\lambda$ guarantees to track the minimum prediction risk $\min_{k \in [n]} R_{k,\infty}^\lambda$ asymptotically.

We first provide numerical illustrations for Theorem 3.1 under the non-isotropic AR(1) data model, which is the same as the one used for Figure 1; see Appendix I for model details. Figure 2 shows both the GCV estimate and the asymptotic risk for the full ridge ensemble. We observe a close match of the theoretical curves and the GCV estimates.

Combining Theorem 2.3 and Theorem 3.1, we can obtain the following corollary regarding GCV subsample tuning.

Corollary 3.2 (Ridge tuning by GCV subsample tuning). *Suppose Assumptions 2.1 and 2.2 hold. Then, we have*

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

as $k, n, p \rightarrow \infty$ such that $p/n \rightarrow \phi \in (0, \infty)$.

Corollary 3.2 certifies the validity of GCV tuning for achieving the optimal risk over all possible regularization parameters and subsample sizes. In practice, tuning for the ridge parameter λ requires one to determine a grid of λ 's for cross-validation. However, the maximum value for the grid is generally chosen by some ad hoc criteria. For example, there is no default maximum value for ridge tuning in the widely-used package `glmnet` (Friedman et al., 2010). From Theorem 2.3, when the signal-noise ratio ρ^2/σ^2 is small, the subsample size should be small enough (so that ϕ_s is large), and the range of λ 's grid should be large enough to cover its optimal value. On the contrary, the GCV-based method does not need such an upper bound for the grid \mathcal{K}_n of subsample sizes because the sample size provides a natural grid in finite samples, informed by the dataset.

In Corollary 3.2, we fix the ridge regularization parameter λ to be zero. But, one can also use other value of $\lambda < \lambda^*$ and the similar statement still holds with $\operatorname{gcv}_{\hat{k}^0, \infty}$ replaced by $\operatorname{gcv}_{\hat{k}^\lambda, \infty}$ based on Theorem 2.3. Furthermore, one can construct the estimator of λ^* as $\hat{\lambda} = \lambda(n - \hat{k}^0)/(\hat{k}^\lambda - \hat{k}^0)$ by extrapolating the line segment between $(0, \hat{k}^0)$ and $(\lambda, \hat{k}^\lambda)$.

In Figure 3, we numerically compare the optimal subsampled ridgeless ensemble with the optimal ridge predictor to verify Corollary 3.2. As we can see, their theoretical curves exactly match, and the empirical estimates in finite samples are also close to their asymptotic limits.

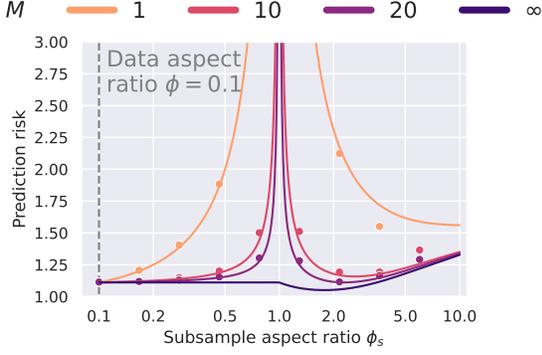


Figure 4. Asymptotic prediction risk curves of ridgeless ensembles, under model (M-AR1) when $\rho_{\text{AR1}} = 0.5$, $\sigma^2 = 1$, and $\phi = 0.1$. The points denote the finite-sample GCV estimates of ridgeless ensembles for varying ensemble sizes $M \in \{1, 10, 20\}$ averaged over 50 dataset repetitions, with $n = \lfloor p/\phi \rfloor$ and $p = 500$.

3.2. A Finite-Ensemble Inconsistency Result

While deriving GCV asymptotics in the proof of Theorem 3.1 for the full ensemble, we also obtain as a byproduct the asymptotic limit of the GCV estimate for finite ensembles. From related work (Patil et al., 2021) and Theorem 3.1, we already know that $\text{gcv}_{k,1}^\lambda$ and $\text{gcv}_{k,\infty}^\lambda$ are consistent estimators of the non-ensemble risk ($R_{k,1}^\lambda$) and the full ensemble risk ($R_{k,\infty}^\lambda$), respectively. However, $\text{gcv}_{k,M}^\lambda$ for $1 < M < \infty$ may not be consistent, which is somewhat surprising. As an example, GCV for $M = 2$ is not a consistent estimator for the prediction risk $R_{k,2}^\lambda$, as shown in the following proposition.

Proposition 3.3 (GCV inconsistency for ridgeless, $M = 2$). *Suppose Assumptions 2.1 and 2.2 hold with $\rho^2, \sigma^2 \in (0, \infty)$. Then, for any $\phi \in (0, \infty)$, we have*

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

as $k, n, p \rightarrow \infty$, $p/n \rightarrow \phi$, $p/k \rightarrow \phi_s \in (1, \infty) \cap (\phi, \infty)$.

Intuitively, the inconsistency for a finite in large part happens because, for a finite M , the residuals computed using the bagged predictor contain non-negligible fractions of out-of-sample and in-sample, and all of them are treated equally. As a result, the GCV estimate for finite ensembles indirectly relates to the original data through the aspect ratios (ϕ, ϕ_s) , even though the GCV estimate is computed only using the training observations. See Section 5 about possible approaches for the corrected GCV estimate for arbitrary ensemble sizes. Though in practice, the correction may not be crucial for a moderate M , because the GCV estimate is close to the underlying target as shown in Figure 4.

3.3. Proof Outline of Theorem 3.1

There are three key steps are involved to prove Theorem 3.1. (1) Deriving the asymptotic limit of the prediction risk $R_{k,M}^\lambda$. (2) Deriving the asymptotic limit of GCV estimate $\text{gcv}_{k,\infty}^\lambda$. (3) Showing pointwise consistency in k by matching the two limits and then lifting to uniform convergence in k . We briefly explain key ideas for showing the three steps below.

(1) Asymptotic limit of risk. We build upon prior results on the risk analysis of ridge ensembles. Under Assumptions 2.1 and 2.2, Lemma A.2 adapted from Patil et al. (2022a) implies that the conditional prediction risks under proportional asymptotics converge to certain deterministic limits:

$$R_{k,M}^\lambda \xrightarrow{\text{a.s.}} \mathcal{R}_M^\lambda(\phi, \phi_s), \quad R_{k,\infty}^\lambda \xrightarrow{\text{a.s.}} \mathcal{R}_\infty^\lambda(\phi, \phi_s), \quad (14)$$

where $\mathcal{R}_M^\lambda(\phi, \phi_s)$, $\mathcal{R}_\infty^\lambda(\phi, \phi_s)$ are as defined in (21).

(2) Asymptotic limit of GCV. To analyze the asymptotic behavior of the GCV estimates, we obtain the asymptotics of the denominator and the numerator of GCV separately. We first show the regular cases when $\phi_s < \infty$ and $\lambda > 0$, and then incorporate boundary cases of $\phi_s = \infty$ and $\lambda = 0$. Our analysis begins with the following lemma that provides asymptotics for the denominator $D_{k,\infty}^\lambda$ (as in (13)) of GCV:

Lemma 3.4 (Asymptotics of the GCV denominator). *Suppose Assumption 2.1 holds. Then, for all $\lambda > 0$,*

$$D_{k,\infty}^\lambda \xrightarrow{\text{a.s.}} \mathcal{D}_\infty^\lambda(\phi, \phi_s),$$

as $k, n, p \rightarrow \infty$, $p/n \rightarrow \phi \in (0, \infty)$, $p/k \rightarrow \phi_s \in [\phi, \infty)$.

It is worth noting that Lemma 3.4 does not require Assumption 2.2 because the smoothing matrix only concerns the design matrix \mathbf{X} and does not depend on the response \mathbf{y} .

Towards obtaining asymptotics for the numerator $T_{k,\infty}^\lambda$ (as in (8)) of GCV, we first decompose $T_{k,\infty}^\lambda$ into simpler components via Lemma D.1. Specifically, the full mean squared training error admits the following decomposition:

$$T_{k,\infty}^\lambda - \sum_{m=1}^2 (c_m T_{k,m}^\lambda + (1 - c_m) \bar{R}_{k,m}^\lambda) \xrightarrow{\text{a.s.}} 0,$$

where $c_1 = \phi/\phi_s$ and $c_2 = 2\phi(2\phi_s - \phi)/\phi_s^2$. Here, $T_{k,m}^\lambda$ and $\bar{R}_{k,m}^\lambda$ are the in-sample training and out-of-sample test errors of the m -ensemble for $m = 1$ and 2, as defined in (6) and (7). This decomposition implies that the full training error is asymptotically simply a linear combination of training and test errors. Therefore, it suffices to obtain the asymptotics of each of these components. As analyzed in Lemma D.2, it is easy to show that the test errors converge to \mathcal{R}_m^λ for $m = 1, 2$. On the other hand, it is more challenging to derive the asymptotic limits for the training errors $T_{k,m}^\lambda$. We first split the $T_{k,m}^\lambda$ into finer components via a

bias-variance decomposition of $T_{k,m}^\lambda$. By developing novel deterministic equivalents of resolvents arising from the decomposition of in-sample errors in Lemma F.8, we are able to show the convergence of the bias and variance components (see Lemma D.3). Combining Lemmas D.1 to D.3 yields the convergence of $T_{k,m}^\lambda$ to a deterministic limit \mathcal{T}_m^λ as summarized in the following lemma:

Lemma 3.5 (Asymptotics of the GCV numerator). *Suppose Assumptions 2.1 and 2.2 hold. Then, for all $\lambda > 0$,*

$$T_{k,\infty}^\lambda \xrightarrow{\text{a.s.}} \mathcal{T}_\infty^\lambda = \sum_{m=1}^2 (c_m \mathcal{T}_m^\lambda + (1 - c_m) \mathcal{R}_m^\lambda), \quad (15)$$

as $k, n, p \rightarrow \infty$, $p/n \rightarrow \phi \in (0, \infty)$, $p/k \rightarrow \phi_s \in [\phi, \infty)$, where $c_1 = \phi/\phi_s$ and $c_2 = 2\phi(2\phi_s - \phi)/\phi_s^2$.

Finally, the boundary cases when $\phi_s = +\infty$ and $\lambda = 0$ are taken care of in succession by Proposition E.1 and Proposition E.2, respectively. Combining the above results provides the asymptotics for the GCV estimate in the full ensemble:

Proposition 3.6 (Asymptotics of GCV for full ensemble). *Suppose Assumptions 2.1 and 2.2 hold. Then, for all $\lambda \geq 0$,*

$$\text{gcv}_{k,\infty}^\lambda \xrightarrow{\text{a.s.}} \mathcal{G}_\infty^\lambda(\phi, \phi_s) := \frac{\mathcal{T}_\infty^\lambda(\phi, \phi_s)}{\mathcal{D}_\infty^\lambda(\phi, \phi_s)}, \quad (16)$$

as $k, n, p \rightarrow \infty$, $p/n \rightarrow \phi \in (0, \infty)$, $p/k \rightarrow \phi_s \in [\phi, \infty)$.

(3) Asymptotics matching and uniform convergence. The asymptotic limits obtained in the first steps can be shown to match with each other by algebraic manipulations. This shows the pointwise consistency in Theorem 3.1. The uniform convergence then follows by applying a certain Cesàro-type mean convergence lemma (see Lemma G.5).

4. Real Data Example: Single-Cell Multiomics

We compare tuning subsample size in the full ridgeless ensemble with tuning the ridge parameter on the full data in a real-world data example from multiomics. This single-cell CITE-seq dataset from Hao et al. (2021) consists of 50,781 human peripheral blood mononuclear cells (PBMCs) originating from eight volunteers post-vaccination (day 3) of an HIV vaccine, which simultaneously measures 20,729 genes and 228 proteins in individual cells.

We follow the standard preprocessing procedure in single-cell data analysis (Hao et al., 2021; Du et al., 2022) to select the top 5,000 highly variable genes and the top 50 highly abundant surface proteins, which exhibit high cell-to-cell variations in the dataset. The gene expression and protein abundance counts for each cell are then divided by the total counts for that cell and multiplied by 10^4 and log-normalized. We randomly hold out half of the cells in each cell type as a test set. The top 500 principal components

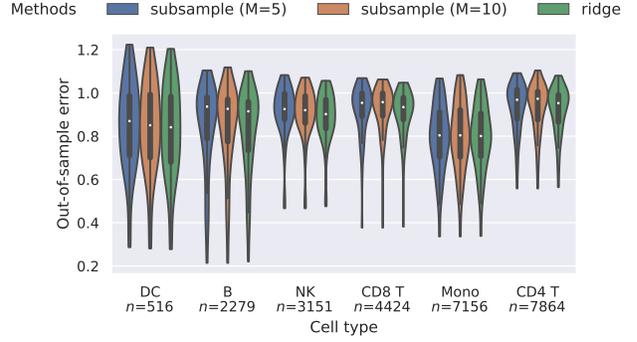


Figure 5. Violin plots of mean squared errors on the randomly held out test sets of different tuning methods for predicting the abundances of 50 proteins in the single-cell CITE-seq dataset. The sizes of the test sets are the same as the sizes of the training sets.

of the standardized gene expressions are used as features to predict protein abundances. The results of using ensembles with subsample size (k) tuning and ridge tuning (λ) without subsampling based on the GCV estimates are compared in Figure 5. For the former, we search over the grid of 25 k 's from n^ν to n spaced evenly on the log scale, with $\nu = 0.5$ and sample size n ranges from 516 to 7864 for different cell types. For the latter, we search over the grid of 100 λ 's from 10^{-2} to 10^2 spaced evenly on log scale. Since different cell types have different sample sizes, this results in different data aspect ratios, presented in increasing order in Figure 5.

From Figure 5, we see that using a moderate ensemble size ($M = 5$ or 10) and tuning the subsample size have a very similar performance to only tuning the ridge regularization parameter in the full dataset. This suggests that the results of Corollary 3.2 also hold even on real data for different data aspect ratios. As discussed after Corollary 3.2, subsample tuning is easier to implement because the dataset provides a natural lower and upper bound of the subsample size. On the other hand, ridge tuning requires one to heuristically pick the upper regularization threshold for the search grid.

5. Discussion and Future Directions

In this work, we provide the risk characterization for the full ridge ensemble and establish the oracle risk equivalences between the full ridgeless ensemble and ridge regression. At a high level, these equivalences show that implicit regularization induced by subsampling matches explicit ridge regularization, i.e., a subsampled ridge predictor with penalty λ_1 has the same risk as another ridge predictor with penalty $\lambda_2 \geq \lambda_1$. Additionally, we prove the uniform consistency of generalized cross-validation for full ridge ensembles, which implies the validity of GCV tuning (that does not require sample splitting) for optimal predictive performance. We describe next some avenues for future work moving forward.

Bias correction for finite ensembles. In Proposition 3.3, we show that the GCV estimate can be inconsistent in the finite ridge and ridgeless ensembles. The inconsistency for $M = 2$ occurs because sampling from the whole dataset induces extra randomness beyond those of the training observations used to compute the GCV estimates. Our analysis of GCV for the full ensemble suggests a way to correct the bias of the GCV estimate. In Appendix H, we outline a possible correction strategy for finite ensembles based on out-of-bag estimates. An intriguing next research direction is to investigate the implementation and uniform consistency of the corrected GCV for finite ensembles in detail.

Extensions to other error metrics. In this paper, we focus on the in-distribution squared prediction risk. It is of interest to extend the equivalences for other error metrics, such as squared estimation risk, general prediction risks, and other functionals of the out-of-sample error distribution, like the quantiles of the error distribution. Additionally, for the purposes of tuning, it is also of interest to extend the GCV analysis to estimate such functionals of the out-of-sample error distribution. Such functional estimation could be valuable in constructing prediction intervals for the unknown response. The technical tools introduced in Patil et al. (2022c) involving leave-one-out perturbation techniques could prove useful for such an extension. Furthermore, this extension would also allow for extending the results presented in this paper to hold under a general non-linear response model.

Extensions to other base predictors. Finally, the focus of this paper is the base ridge predictor. A natural extension of the current work is to consider kernel ridge regression. Going further, it is of much interest to consider other regularized predictors, such as lasso. Whether optimal subsampled lassoless regression still matches with the optimal lasso is an interesting question. There is already empirical evidence along the lines of Figure 1 for such a connection. The results proved in the current paper make us believe that there is a general story quantifying the effect of implicit regularization by subsampling and that provided by explicit regularization. Whether the general story unfolds as neatly as presented here for ridge regression remains an exciting next question!

Acknowledgements

We are grateful to Ryan Tibshirani, Alessandro Rinaldo, Yuting Wei, Matey Neykov, Daniel LeJeune, Shamindra Shrotriya for many helpful conversations surrounding ridge regression, subsampling, and generalized cross-validation. Many thanks are also due to the anonymous reviewers for their encouraging remarks and insightful questions that have informed several new directions for follow-up future work. In particular, special thanks to the reviewer “2t75” for a wonderful review and highlighting other related works on subsampling that have made their way into the manuscript.

References

- Adlam, B. and Pennington, J. Understanding double descent requires a fine-grained bias-variance decomposition. *Advances in neural information processing systems*, 33: 11022–11032, 2020a.
- Adlam, B. and Pennington, J. The neural tangent kernel in high dimensions: Triple descent and a multi-scale theory of generalization. In *International Conference on Machine Learning*. PMLR, 2020b.
- Allen, D. M. The relationship between variable selection and data augmentation and a method for prediction. *Technometrics*, 16(1):125–127, 1974.
- Arlot, S. and Celisse, A. A survey of cross-validation procedures for model selection. *Statistics surveys*, 4:40–79, 2010.
- Bai, Z. and Silverstein, J. W. *Spectral Analysis of Large Dimensional Random Matrices*. Springer, 2010. Second edition.
- Bartlett, P. L., Long, P. M., Lugosi, G., and Tsigler, A. Benign overfitting in linear regression. *Proceedings of the National Academy of Sciences*, 117(48):30063–30070, 2020.
- Bartlett, P. L., Montanari, A., and Rakhlin, A. Deep learning: a statistical viewpoint. *Acta numerica*, 30:87–201, 2021.
- Belkin, M., Hsu, D., and Xu, J. Two models of double descent for weak features. *SIAM Journal on Mathematics of Data Science*, 2(4):1167–1180, 2020.
- Bloemendal, A., Knowles, A., Yau, H.-T., and Yin, J. On the principal components of sample covariance matrices. *Probability theory and Related Fields*, 164(1):459–552, 2016.
- Breiman, L. Bagging predictors. *Machine Learning*, 24(2): 123–140, 1996.
- Bühlmann, P. and Yu, B. Analyzing bagging. *The Annals of Statistics*, 30(4):927–961, 2002.
- Buja, A. and Stuetzle, W. Observations on bagging. *Statistica Sinica*, pp. 323–351, 2006.
- Craven, P. and Wahba, G. Estimating the correct degree of smoothing by the method of generalized cross-validation. *Numerische Mathematik*, 31:377–403, 1979.
- Dobriban, E. and Sheng, Y. Distributed linear regression by averaging. *The Annals of Statistics*, 49(2):918–943, 2021.

- Dobriban, E. and Wager, S. High-dimensional asymptotics of prediction: Ridge regression and classification. *The Annals of Statistics*, 46(1):247–279, 2018.
- Du, J.-H., Cai, Z., and Roeder, K. Robust probabilistic modeling for single-cell multimodal mosaic integration and imputation via scvaeit. *Proceedings of the National Academy of Sciences*, 119(49):e2214414119, 2022.
- d’Ascoli, S., Refinetti, M., Biroli, G., and Krzakala, F. Double trouble in double descent: Bias and variance (s) in the lazy regime. In *International Conference on Machine Learning*, pp. 2280–2290. PMLR, 2020.
- El Karoui, N. Asymptotic behavior of unregularized and ridge-regularized high-dimensional robust regression estimators: rigorous results. *arXiv preprint arXiv:1311.2445*, 2013.
- El Karoui, N. On the impact of predictor geometry on the performance on high-dimensional ridge-regularized generalized robust regression estimators. *Probability Theory and Related Fields*, 170(1):95–175, 2018.
- Erdős, L. and Yau, H.-T. *A Dynamical Approach to Random Matrix Theory*. American Mathematical Society, 2017.
- Friedman, J., Hastie, T., and Tibshirani, R. Regularization paths for generalized linear models via coordinate descent. *Journal of statistical software*, 33(1):1, 2010.
- Friedman, J. H. and Hall, P. On bagging and nonlinear estimation. *Journal of Statistical Planning and Inference*, 137(3):669–683, 2007.
- Geisser, S. The predictive sample reuse method with applications. *Journal of the American statistical Association*, 70(350):320–328, 1975.
- Golub, G. H., Heath, M., and Wahba, G. Generalized cross-validation as a method for choosing a good ridge parameter. *Technometrics*, 21(2):215–223, 1979.
- Greene, E. and Wellner, J. A. Exponential bounds for the hypergeometric distribution. *Bernoulli*, 23(3):1911, 2017.
- Grenander, U. and Szegö, G. *Toeplitz Forms and Their Applications*. University of California Press, 1958. First edition.
- Gut, A. *Probability: A Graduate Course*. Springer, New York, 2005.
- Györfi, L., Kohler, M., Krzyzak, A., and Walk, H. *A Distribution-free Theory of Nonparametric Regression*. Springer Science & Business Media, 2006.
- Hall, P. and Samworth, R. J. Properties of bagged nearest neighbour classifiers. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 67(3):363–379, 2005.
- Hao, Y., Hao, S., Andersen-Nissen, E., Mauck III, W. M., Zheng, S., Butler, A., Lee, M. J., Wilk, A. J., Darby, C., Zager, M., et al. Integrated analysis of multimodal single-cell data. *Cell*, 2021.
- Hastie, T. Ridge regularization: An essential concept in data science. *Technometrics*, 62(4):426–433, 2020.
- Hastie, T., Tibshirani, R., and Friedman, J. *The Elements of Statistical Learning*. Springer Series in Statistics, 2009. Second edition.
- Hastie, T., Montanari, A., Rosset, S., and Tibshirani, R. J. Surprises in high-dimensional ridgeless least squares interpolation. *The Annals of Statistics*, 50(2):949–986, 2022.
- Hoerl, A. E. and Kennard, R. W. Ridge regression: Biased estimation for nonorthogonal problems. *Technometrics*, 12(1):55–67, 1970a.
- Hoerl, A. E. and Kennard, R. W. Ridge regression: applications to nonorthogonal problems. *Technometrics*, 12(1):69–82, 1970b.
- Kale, S., Kumar, R., and Vassilvitskii, S. Cross-validation and mean-square stability. In *In Proceedings of the Second Symposium on Innovations in Computer Science*, 2011.
- Krogh, A. and Sollich, P. Statistical mechanics of ensemble learning. *Physical Review E*, 55(1):811, 1997.
- Kumar, R., Lokshtanov, D., Vassilvitskii, S., and Vattani, A. Near-optimal bounds for cross-validation via loss stability. In *International Conference on Machine Learning*, 2013.
- LeJeune, D., Javadi, H., and Baraniuk, R. The implicit regularization of ordinary least squares ensembles. In *International Conference on Artificial Intelligence and Statistics*, 2020.
- Li, K.-C. From Stein’s unbiased risk estimates to the method of generalized cross validation. *The Annals of Statistics*, pp. 1352–1377, 1985.
- Li, K.-C. Asymptotic optimality of c_l and generalized cross-validation in ridge regression with application to spline smoothing. *The Annals of Statistics*, 14(3):1101–1112, 1986.
- Li, K.-C. Asymptotic optimality for c_p, c_l , cross-validation and generalized cross-validation: Discrete index set. *The Annals of Statistics*, 15(3):958–975, 1987.

- Loureiro, B., Gerbelot, C., Refinetti, M., Sicuro, G., and Krzakala, F. Fluctuations, bias, variance & ensemble of learners: Exact asymptotics for convex losses in high-dimension. In *International Conference on Machine Learning*, pp. 14283–14314. PMLR, 2022.
- Mei, S. and Montanari, A. The generalization error of random features regression: Precise asymptotics and the double descent curve. *Communications on Pure and Applied Mathematics*, 75(4):667–766, 2022.
- Miolane, L. and Montanari, A. The distribution of the lasso: Uniform control over sparse balls and adaptive parameter tuning. *The Annals of Statistics*, 49(4):2313–2335, 2021.
- Muthukumar, V., Vodrahalli, K., Subramanian, V., and Sahai, A. Harmless interpolation of noisy data in regression. *IEEE Journal on Selected Areas in Information Theory*, 1(1):67–83, 2020.
- Obuchi, T. and Kabashima, Y. Cross validation in LASSO and its acceleration. *Journal of Statistical Mechanics: Theory and Experiment*, 2016.
- Patil, P., Wei, Y., Rinaldo, A., and Tibshirani, R. Uniform consistency of cross-validation estimators for high-dimensional ridge regression. In *International Conference on Artificial Intelligence and Statistics*. PMLR, 2021.
- Patil, P., Du, J.-H., and Kuchibhotla, A. K. Bagging in overparameterized learning: Risk characterization and risk monotonicity. *arXiv preprint arXiv:2210.11445*, 2022a.
- Patil, P., Kuchibhotla, A. K., Wei, Y., and Rinaldo, A. Mitigating multiple descents: A model-agnostic framework for risk monotonicity. *arXiv preprint arXiv:2205.12937*, 2022b.
- Patil, P., Rinaldo, A., and Tibshirani, R. Estimating functionals of the out-of-sample error distribution in high-dimensional ridge regression. In *International Conference on Artificial Intelligence and Statistics*. PMLR, 2022c.
- Rad, K. R. and Maleki, A. A scalable estimate of the out-of-sample prediction error via approximate leave-one-out cross-validation. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 82(4):965–996, 2020.
- Rad, K. R., Zhou, W., and Maleki, A. Error bounds in estimating the out-of-sample prediction error using leave-one-out cross validation in high-dimensions. In *International Conference on Artificial Intelligence and Statistics*. PMLR, 2020.
- Rubio, F. and Mestre, X. Spectral convergence for a general class of random matrices. *Statistics & probability letters*, 81(5):592–602, 2011.
- Rudin, W. *Principles of Mathematical Analysis*. McGraw-Hill New York, 1976.
- Samworth, R. J. Optimal weighted nearest neighbour classifiers. *The Annals of Statistics*, 40(5):2733–2763, 2012.
- Sollich, P. and Krogh, A. Learning with ensembles: How overfitting can be useful. *Advances in neural information processing systems*, 8, 1995.
- Stone, M. Cross-validated choice and assessment of statistical predictions. *Journal of the Royal Statistical Society: Series B*, 36(2):111–133, 1974.
- Stone, M. Asymptotics for and against cross-validation. *Biometrika*, 64(1):29–35, 1977.
- Sur, P., Chen, Y., and Candès, E. J. The likelihood ratio test in high-dimensional logistic regression is asymptotically a rescaled chi-square. *Probability Theory and Related Fields*, 175(1):487–558, 2019.
- Thrapoulidis, C., Oymak, S., and Hassibi, B. Regularized linear regression: A precise analysis of the estimation error. In *Conference on Learning Theory*, pp. 1683–1709. PMLR, 2015.
- Thrapoulidis, C., Abbasi, E., and Hassibi, B. Precise error analysis of regularized M -estimators in high dimensions. *IEEE Transactions on Information Theory*, 64(8):5592–5628, 2018.
- Wang, S., Zhou, W., Lu, H., Maleki, A., and Mirrokni, V. Approximate leave-one-out for fast parameter tuning in high dimensions. *arXiv preprint arXiv:1807.02694*, 2018.
- Wasserman, L. *Olive Nonparametric Statistics*. Springer, 2006.
- Wei, A., Hu, W., and Steinhardt, J. More than a toy: Random matrix models predict how real-world neural representations generalize. *arXiv preprint arXiv:2203.06176*, 2022.
- Xu, J., Maleki, A., and Rad, K. R. Consistent risk estimation in high-dimensional linear regression. *arXiv preprint arXiv:1902.01753*, 2019.
- Zhang, Y. and Yang, Y. Cross-validation for selecting a model selection procedure. *Journal of Econometrics*, 187(1):95–112, 2015.

Title in Wasserman (2006) may seem like a typo, but it is not. If you have a moment to chuckle, peek at column 2 of page 266!

Appendix

This serves as an appendix to the paper “Subsample Ridge Ensembles: Equivalences and Generalized Cross-Validation.” Below we provide an outline for the appendix along with a summary of the notation used in the main paper and the appendix.

Organization

The content of the appendix is organized as follows.

- Appendix A presents proofs of results in Section 2.
 - Appendix A.1 relates the ridge estimator in the full ensemble to the M -ensemble ridge estimator as the ensemble size M tends to infinity. Specifically, we provide proof of the fact [mentioned in Section 2] that the estimator $\tilde{\beta}_{k,\infty}^\lambda$ as defined in (3) is almost surely equivalent to letting the ensemble size $M \rightarrow \infty$ for the estimator $\tilde{\beta}_{k,M}^\lambda$ as defined in (2). The main ingredients are results in Appendix G.
 - Appendix A.2 gathers known results from Patil et al. (2022a) in the form of Lemma A.2 that characterize the asymptotic prediction risks of ridge ensembles used in the remaining sections.
 - Appendix A.3 proves Theorem 2.3. The main ingredients are Lemma A.2 and results in Appendix F. See Figure 6.

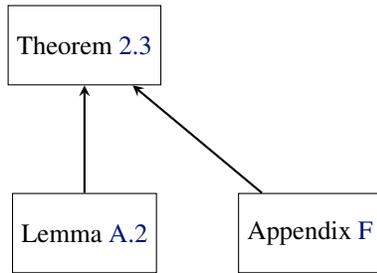


Figure 6. Schematic for the proof of Theorem 2.3.

- Appendix B presents proofs of results in Section 3.
 - Appendix B.1 proves Theorem 3.1. The main ingredients are Lemma B.1 (proved in Appendix B) and results in Appendix G. The main ingredients that prove Lemma B.1 are Proposition 3.6 (proved in Appendix E) and Lemma A.2. See Figure 7.

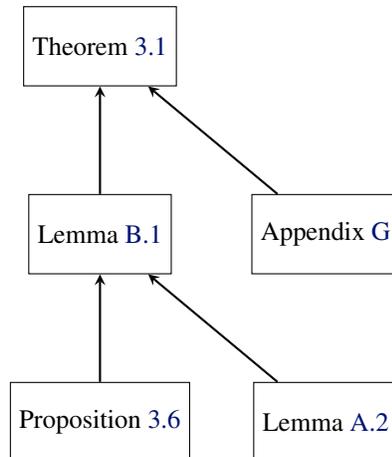


Figure 7. Schematic for the proof of Theorem 3.1.

- Appendix B.2 proves Corollary 3.2. The main ingredient in Theorem 3.1.
- Appendix B.3 proves Proposition 3.3. The main ingredients are Lemma 3.4 and results in Appendix G. See Figure 8.

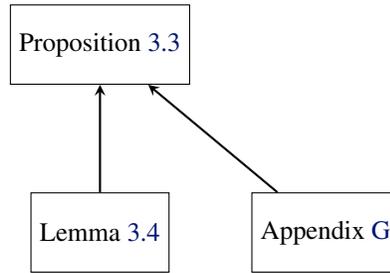


Figure 8. Schematic for the proof of Proposition 3.3.

- Appendix C proves Lemma 3.4. The main ingredient is Lemma C.1 (proved in Appendix C). See Figure 9.

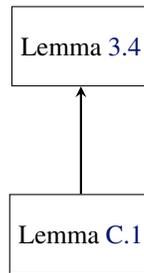


Figure 9. Schematic for the proof of Lemma 3.4.

- Appendix D proves Lemma 3.5. The main ingredients are a series of lemmas, Lemmas D.1 to D.3 (proved in Appendices D.1 to D.3). These lemmas provide structural decompositions for the ensemble train error and obtain the limiting behaviors of the terms in the decompositions. See Figure 10.
 - Appendix D.1 proves Lemma D.1 that shows a certain decomposition of the ensemble train error into in-sample train and out-of-sample test error components.
 - Appendix D.2 proves Lemma D.2 on the convergence of out-of-sample test error components. The main ingredient is Lemma A.2.
 - Appendix D.3 proves Lemma D.3 on the convergence of in-sample train error components. The main ingredients are Lemmas D.4 and D.5 (proved in Appendix D.4) and Lemmas D.6 and D.7 (proved in Appendix D.5).
 - Appendix D.4 proves Lemmas D.4 and D.5 on component concentrations of certain cross and variance terms arising in the aforementioned decompositions. The main ingredients are results in Appendix G.
 - Appendix D.5 proves Lemmas D.6 and D.7 on component deterministic approximations for the concentrated bias and variance functionals in the steps above. The main ingredients are results in Appendix F.
- Appendix E proves Proposition 3.6. The main ingredients are components proved in Lemma 3.4 and Lemma 3.5, Propositions E.1 and E.2 that handle certain boundary cases not covered by Lemma 3.4 and Lemma 3.5 (proved in Appendices E.1 and E.2), and results in Appendix G. See Figure 11.
 - Appendix E.1 proves Proposition E.1 that considers the boundary case as the subsampling ratio $\phi_s \rightarrow \infty$ for ridge regression ($\lambda > 0$).
 - Proposition E.2 proves Proposition E.2 that handles the limiting case of ridgeless regression ($\lambda = 0$), by justifying and taking a suitable limit as $\lambda \rightarrow 0^+$ of the corresponding results for ridge regression.
- Appendix F summarizes auxiliary asymptotic equivalency results used in the proofs throughout.

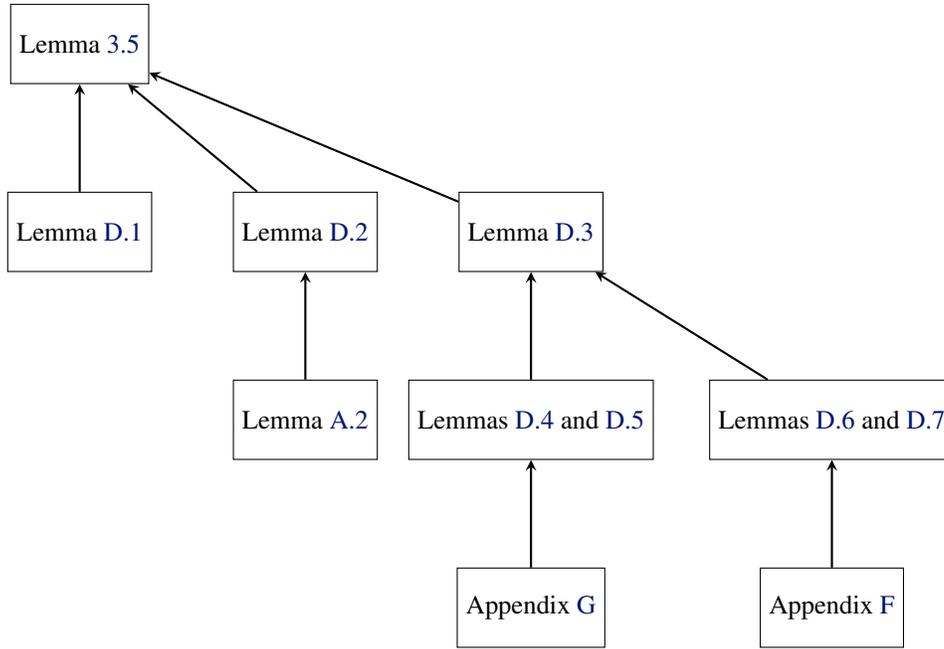


Figure 10. Schematic for the proof of Lemma 3.5.

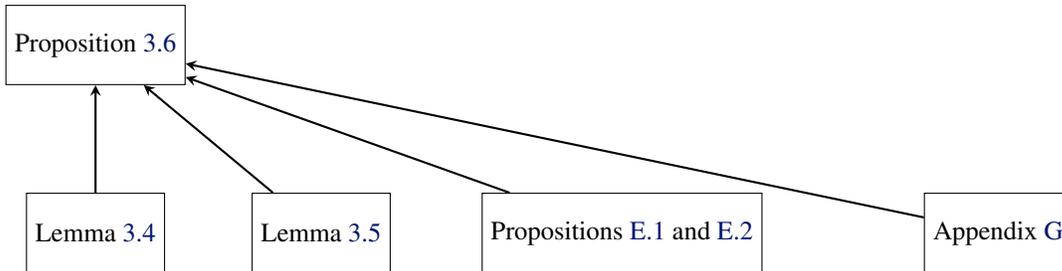


Figure 11. Schematic for the proof of Proposition 3.6.

- Appendix F.1 provides background on the notion of asymptotic matrix equivalents and various calculus rules that such notion of equivalency obeys.
- Appendix F.2 gathers some known asymptotic matrix equivalents and derives some novel asymptotic matrix equivalents that arise in our analysis.
- Appendix F.3 gathers various known analytic properties of certain fixed-point equations and proves some additional properties that arise in our analysis.
- Appendix G collects several helper concentration results used in the proofs throughout.
 - Appendix G.1 provides lemmas deriving the asymptotic proportion of shared observations when subsampling.
 - Appendix G.2 provides lemmas establishing concentrations for linear and quadratic forms of random vectors.
 - Appendix G.3 provides lemmas for lifting original convergences to converges of Ceàro-type mean and max for triangular arrays.
- Appendix H discusses the bias correction of GCV for finite ensembles [mentioned in Section 5].
- Appendix I describes additional numerical details for experiments [mentioned in Section 3].

Notation

An overview of some general notation used in the main paper and the appendix is as follows.

1. **General notation:** We denote scalars in non-bold lower or upper case (e.g., n, λ, C), vectors in bold lower case (e.g., $\mathbf{x}, \boldsymbol{\beta}$), and matrices in bold upper case (e.g., \mathbf{X}). For a real number x , $(x)_+$ denotes its positive part, $\lfloor x \rfloor$ its floor, and $\lceil x \rceil$ its ceiling. For a vector $\boldsymbol{\beta}$, $\|\boldsymbol{\beta}\|_2$ denotes its ℓ_2 norm. For a pair of vectors \mathbf{v} and \mathbf{w} , $\langle \mathbf{v}, \mathbf{w} \rangle$ denotes their inner product. For an event A , $\mathbb{1}_A$ denotes the associated indicator random variable. We denote convergence in probability by “ \xrightarrow{p} ”, almost sure convergence by “ $\xrightarrow{\text{a.s.}}$ ”, and convergence in distribution by “ \xrightarrow{d} ”.
2. **Set notation:** We denote sets using calligraphic letters (e.g., \mathcal{D}), and use blackboard letters to denote some special sets: \mathbb{N} denotes the set of positive integers, \mathbb{R} denotes the set of real numbers, $\mathbb{R}_{\geq 0}$ denotes the set of non-negative real numbers, $\mathbb{R}_{> 0}$ denotes the set of positive real numbers, \mathbb{C} denotes the set of complex numbers, \mathbb{C}^+ denotes the set of complex numbers with positive imaginary part, and \mathbb{C}^- denotes the set of complex numbers with negative imaginary part. For a natural number n , we use $[n]$ to denote the set $\{1, \dots, n\}$.
3. **Matrix notation:** For a matrix $\mathbf{X} \in \mathbb{R}^{n \times p}$, $\mathbf{X}^\top \in \mathbb{R}^{p \times n}$ denotes its transpose, and $\mathbf{X}^+ \in \mathbb{R}^{p \times n}$ denote its Moore-Penrose inverse. For a square matrix $\mathbf{A} \in \mathbb{R}^{p \times p}$, $\text{tr}[\mathbf{A}]$ denotes its trace, and $\mathbf{A}^{-1} \in \mathbb{R}^{p \times p}$ denotes its inverse, provided it is invertible. For a positive semidefinite matrix $\boldsymbol{\Sigma}$, $\boldsymbol{\Sigma}^{1/2}$ denotes its principal square root. A $p \times p$ identity matrix is denoted \mathbf{I}_p , or simply by \mathbf{I} , when it is clear from the context.

For a real matrix \mathbf{X} , its operator norm (or spectral norm) with respect to ℓ_2 vector norm is denoted by $\|\mathbf{X}\|_{\text{op}}$, and its trace norm (or nuclear norm) is denoted by $\|\mathbf{X}\|_{\text{tr}}$ (recall that $\|\mathbf{X}\|_{\text{tr}} = \text{tr}[(\mathbf{X}^\top \mathbf{X})^{1/2}]$). For a positive semidefinite matrix $\mathbf{A} \in \mathbb{R}^{p \times p}$ with eigenvalue decomposition $\mathbf{A} = \mathbf{V} \mathbf{R} \mathbf{V}^{-1}$ for an orthonormal matrix $\mathbf{V} \in \mathbb{R}^{p \times p}$ and a diagonal matrix $\mathbf{R} \in \mathbb{R}^{p \times p}$ with non-negative entries, and a function $f : \mathbb{R}_{\geq 0} \rightarrow \mathbb{R}_{\geq 0}$, we denote by $f(\mathbf{A})$ the $p \times p$ positive semidefinite matrix $\mathbf{V} f(\mathbf{R}) \mathbf{V}^{-1}$. Here, $f(\mathbf{R})$ is a $p \times p$ diagonal matrix obtained by applying the function f to each diagonal entry of \mathbf{R} .

For symmetric matrices \mathbf{A} and \mathbf{B} , $\mathbf{A} \preceq \mathbf{B}$ denotes the Loewner ordering. For sequences of matrices \mathbf{A}_n and \mathbf{B}_n , $\mathbf{A}_n \simeq \mathbf{B}_n$ denotes a certain notion of asymptotic equivalence (see Appendix F).

Finally, in what follows, we will prove the results for n, k, p being a sequence of integers $\{n_m\}_{m=1}^\infty, \{k_m\}_{m=1}^\infty, \{p_m\}_{m=1}^\infty$. One can also view k and p as sequences k_n and p_n that are indexed by n . For simplicity, we drop the subscript when it is clear from the context.

A. Proofs of results in Section 2

A.1. Full-ensemble versus limiting M -ensemble

Lemma A.1 (Almost sure equivalence of full-ensemble and limiting M -ensemble). *For k, n fixed, for the ensemble estimator defined in (2), it holds that*

$$\tilde{\boldsymbol{\beta}}_{k,M}^\lambda(\mathcal{D}_n; \{I_\ell\}_{\ell=1}^M) \xrightarrow{\text{a.s.}} \mathbb{E}[\hat{\boldsymbol{\beta}}_k^\lambda(\mathcal{D}_I) \mid \mathcal{D}_n] = \frac{1}{|\mathcal{I}_k|} \sum_{I \in \mathcal{I}_k} \hat{\boldsymbol{\beta}}_k^\lambda(\mathcal{D}_I),$$

as $M \rightarrow \infty$.

Proof of Lemma A.1. Note that for k, n fixed, the cardinality of \mathcal{I}_k is $\binom{n}{k}$. Thus, we have

$$\tilde{\boldsymbol{\beta}}_{k,M}^\lambda(\mathcal{D}_n; \{I_\ell\}_{\ell=1}^M) = \sum_{I \in \mathcal{I}_k} \frac{n_{M,I}}{M} \hat{\boldsymbol{\beta}}_k^\lambda(\mathcal{D}_I)$$

for random variables $n_{M,i}$'s. Since when sampling with replacement $n_{M,I} \sim \text{Binomial}(M, 1/\binom{n}{k})$ with mean $M/\binom{n}{k}$, from the strong law of large numbers, we have that as $M \rightarrow \infty$,

$$\frac{n_{M,I}}{M} \xrightarrow{\text{a.s.}} \frac{1}{\binom{n}{k}}, \quad \forall I \in \mathcal{I}_k. \quad (17)$$

For sampling without replacement, $n_{M,I} \sim \text{Hypergeometric}(M, 1, \binom{n}{k})$ for $M \leq \binom{n}{k}$ (see Definition G.1) with mean $M/\binom{n}{k}$. When $M = \binom{n}{k}$, $n_{M,I}/M = 1/\binom{n}{k}$. In both cases, we have

$$\tilde{\beta}_{k,M}^\lambda(\mathcal{D}_n; \{I_\ell\}_{\ell=1}^\infty) := \lim_{M \rightarrow \infty} \tilde{\beta}_{k,M}^\lambda(\mathcal{D}_n; \{I_\ell\}_{\ell=1}^M) \stackrel{\text{a.s.}}{=} \frac{1}{\binom{n}{k}} \sum_{I \in \mathcal{I}_k} \hat{\beta}_k^\lambda(\mathcal{D}_I),$$

which concludes the proof. \square

A.2. Risk characterization of ridge ensembles

In the study of ridge ensembles under proportional asymptotics, a key quantity that appears is the solution of a fixed-point equation. For any $\lambda > 0$ and $\theta > 0$, define $v(-\lambda; \theta)$ as the unique nonnegative solution to the fixed-point equation

$$v(-\lambda; \theta)^{-1} = \lambda + \theta \int r(1 + v(-\lambda; \theta)r)^{-1} dH(r). \quad (18)$$

For $\lambda = 0$, define $v(0; \theta) := \lim_{\lambda \rightarrow 0^+} v(-\lambda; \theta)$ for $\theta > 1$ and $+\infty$ otherwise. Such a fixed-point equation has appeared in the literature before. See, e.g., Dobriban & Wager (2018); Hastie et al. (2022); Mei & Montanari (2022) in the context of ridge regression; and more generally, for other M -estimators, see, e.g., El Karoui (2013; 2018); Thrampoulidis et al. (2015; 2018); Sur et al. (2019); Miolane & Montanari (2021), among others. The fact that the fixed-point equation (18) has a unique nonnegative solution follows from Patil et al. (2022b, Lemma S.6.14).

We then define the nonnegative constants $\tilde{v}(-\lambda; \vartheta, \theta)$, and $\tilde{c}(-\lambda; \theta)$ via the following equations:

$$\tilde{v}(-\lambda; \vartheta, \theta) = \frac{\vartheta \int r^2(1 + v(-\lambda; \theta)r)^{-2} dH(r)}{v(-\lambda; \theta)^{-2} - \vartheta \int r^2(1 + v(-\lambda; \theta)r)^{-2} dH(r)}, \quad \tilde{c}(-\lambda; \theta) = \int r(1 + v(-\lambda; \theta)r)^{-2} dG(r). \quad (19)$$

Lemma A.2 (Risk characterization of ridge ensembles, adapted from Patil et al. (2022a)). *Suppose Assumptions 2.1-2.2 hold for the dataset \mathcal{D}_n . Then, as $k, n, p \rightarrow \infty$ such that $p/n \rightarrow \phi \in (0, \infty)$ and $p/k \rightarrow \phi_s \in [\phi, \infty]$ (and $\phi_s \neq 1$ if $\lambda = 0$), there exist deterministic functions $\mathcal{R}_M^\lambda(\phi, \phi_s)$ for $M \in \mathbb{N}$, such that for $I_1, \dots, I_M \stackrel{\text{SRS}}{\sim} \mathcal{I}_k$,*

$$\sup_{M \in \mathbb{N}} |R(\tilde{f}_{M, \mathcal{I}_k}; \mathcal{D}_n, \{I_\ell\}_{\ell=1}^M) - \mathcal{R}_M^\lambda(\phi, \phi_s)| \xrightarrow{P} 0. \quad (20)$$

Furthermore, the function $\mathcal{R}_M^\lambda(\phi, \phi_s)$ decomposes as

$$\mathcal{R}_M^\lambda(\phi, \phi_s) = \sigma^2 + \mathcal{B}_M^\lambda(\phi, \phi_s) + \mathcal{V}_M^\lambda(\phi, \phi_s), \quad (21)$$

where the bias and variance terms are given by

$$\mathcal{B}_M^\lambda(\phi, \phi_s) = M^{-1}B_\lambda(\phi_s, \phi_s) + (1 - M^{-1})B_\lambda(\phi, \phi_s), \quad (22)$$

$$\mathcal{V}_M^\lambda(\phi, \phi_s) = M^{-1}V_\lambda(\phi_s, \phi_s) + (1 - M^{-1})V_\lambda(\phi, \phi_s), \quad (23)$$

and the functions $B_\lambda(\cdot, \cdot)$ and $V_\lambda(\cdot, \cdot)$ are defined as

$$B_\lambda(\vartheta, \theta) = \rho^2(1 + \tilde{v}(-\lambda; \vartheta, \theta))\tilde{c}(-\lambda; \theta), \quad V_\lambda(\vartheta, \theta) = \sigma^2\tilde{v}(-\lambda; \vartheta, \theta), \quad \theta \in (0, \infty], \vartheta \leq \theta. \quad (24)$$

A.3. Proof of Theorem 2.3

Proof of Theorem 2.3. Define $\phi_s^*(\phi) := \operatorname{argmin}_{\phi_s \geq \phi} \mathcal{R}_{\lambda, \infty}(\phi, \phi_s)$ and $\lambda^*(\phi) := \operatorname{argmin}_{\lambda \geq 0} \mathcal{R}_1^\lambda(\phi, \phi)$. We will write ϕ_s^* and λ^* for simplicity and split the proof into different cases.

Part (1) Case of SNR > 0 ($\rho^2 > 0, \sigma^2 > 0$):

From Patil et al. (2022a, Proposition 5.7) we have that $\phi_s^* \in (\phi \vee 1, \infty)$. From Lemma F.11 (1), the function $\phi_s \mapsto v(0; \phi_s)$ is strictly decreasing over $\phi_s \in [1, \infty]$ with range

$$v(0; \phi_s \vee 1) = \begin{cases} v(0; \phi_s), & \phi \in (1, \infty) \\ \lim_{\phi_s \rightarrow 1^+} v(0; \phi_s) = +\infty, & \phi \in (0, 1] \end{cases}, \quad v(0; +\infty) := \lim_{\phi_s \rightarrow +\infty} v(0; \phi_s) = 0.$$

From Lemma F.12 (1), the function $\lambda \mapsto v(-\lambda; \phi)$ is strictly decreasing over $\lambda \in [0, \infty]$ with range

$$v(0; \phi \vee 1) = \begin{cases} v(0; \phi), & \phi \in (1, \infty) \\ \lim_{\lambda \rightarrow 0^+} v(-\lambda; \phi) = +\infty, & \phi \in (0, 1] \end{cases}, \quad v(-\infty; \phi) := \lim_{\lambda \rightarrow +\infty} v(-\lambda; \phi) = 0.$$

By the intermediate value theorem, there exists unique $\lambda_0 \in (0, \infty)$ such that $v(-\lambda_0; \phi) = v(0; \phi_s^*)$. Then we also have $\tilde{c}(-\lambda_0; \phi) = \tilde{c}(0; \phi_s^*)$ and $\tilde{v}(-\lambda_0; \phi, \phi) = \tilde{v}(0; \phi_s^*)$. Substituting this into the optimal ensemble risk, we have

$$\begin{aligned} \min_{\phi_s \geq \phi} \mathcal{R}_\infty^0(\phi, \phi_s) &= \mathcal{R}_\infty^0(\phi, \phi_s^*) \\ &= (\sigma^2 + \rho^2 \tilde{c}(0; \phi_s^*)) (1 + \tilde{v}(0; \phi, \phi_s^*)) \\ &= (\sigma^2 + \rho^2 \tilde{c}(-\lambda_0; \phi)) (1 + \tilde{v}(-\lambda_0; \phi, \phi)) \\ &= \mathcal{R}_1^{\lambda_0}(\phi, \phi) \\ &\leq \min_{\lambda \geq 0} \mathcal{R}_1^\lambda(\phi, \phi). \end{aligned}$$

On the other hand, there exists unique $\phi_0 \in [1, \infty)$ such that $v(-\lambda^*; \phi) = v(0; \phi_0)$, and thus, we have

$$\begin{aligned} \min_{\lambda \geq 0} \mathcal{R}_1^\lambda(\phi, \phi) &= \mathcal{R}_1^{\lambda^*}(\phi, \phi) \\ &= (\sigma^2 + \rho^2 \tilde{c}(-\lambda^*; \phi)) (1 + \tilde{v}(-\lambda^*; \phi, \phi)) \\ &= (\sigma^2 + \rho^2 \tilde{c}(0; \phi_0)) (1 + \tilde{v}(0; \phi, \phi_0)) \\ &= \mathcal{R}_\infty^0(\phi, \phi_0) \\ &\leq \min_{\phi_s \geq \phi} \mathcal{R}_\infty^0(\phi, \phi_s). \end{aligned}$$

Combining the above two inequalities, we have that $\min_{\phi_s \geq \phi} \mathcal{R}_\infty^0(\phi, \phi_s) = \min_{\lambda \geq 0} \mathcal{R}_1^\lambda(\phi, \phi)$.

Part (2) Case of SNR = 0 ($\rho^2 = 0, \sigma^2 > 0$):

From Patil et al. (2022a, Proposition 5.7) we have that $\phi_s^* = +\infty$, which implies that $v(0; \phi_s^*) = 0$. Then, from Lemma F.11 (1) we have $v(0; +\infty) := \lim_{\phi_s \rightarrow +\infty} v(0; \phi_s) = 0$,

$$\begin{aligned} \lim_{\phi_s \rightarrow +\infty} \tilde{v}(0; \phi, +\infty) &= \lim_{\phi_s \rightarrow +\infty} \frac{\phi \int r^2 (1 + v(0; \phi_s) r)^{-2} dH(r)}{v(0; \phi_s)^{-2} - \phi \int r^2 (1 + v(0; \phi_s) r)^{-2} dH(r)} \\ &= \lim_{\phi_s \rightarrow +\infty} \frac{\phi \int (v(0; \phi_s) r)^2 (1 + v(0; \phi_s) r)^{-2} dH(r)}{1 - \phi \int (v(0; \phi_s) r)^2 (1 + v(0; \phi_s) r)^{-2} dH(r)} \\ &= \frac{\phi \int (v(0; +\infty) r)^2 (1 + v(0; +\infty) r)^{-2} dH(r)}{1 - \phi \int (v(0; +\infty) r)^2 (1 + v(0; +\infty) r)^{-2} dH(r)} \\ &= 0, \end{aligned}$$

and thus,

$$\min_{\phi_s \geq \phi} \mathcal{R}_\infty^0(\phi, \phi_s) = \mathcal{R}_\infty^0(\phi, \infty) = \sigma^2 (1 + \tilde{v}(0; \phi, +\infty)) = \sigma^2.$$

On the other hand,

$$\min_{\lambda \geq 0} \mathcal{R}_1^\lambda(\phi, \phi) = \mathcal{R}_1^{\lambda^*}(\phi, \phi) = \sigma^2 \tilde{v}(-\lambda^*; \phi, \phi) \geq \sigma^2$$

where the equality holds when $\lambda^* = +\infty$ because $\tilde{v}(-\lambda^*; \phi, \phi) \geq 0$ from Lemma F.10 (4). Thus, the optimal parameters to the two optimization problems are given by $\phi_s^* = \lambda^* = +\infty$, with $v(0; \phi_s^*) = v(-\lambda^*; \phi) = 0$.

Part (3) Case of $\text{SNR} = \infty$ ($\rho^2 > 0, \sigma^2 = 0$):

When $\phi \leq 1$, from Patil et al. (2022a, Proposition 5.7) we have that any $\phi_s^* \in [\phi, 1]$ minimizes $\min_{\phi_s \geq \phi} \mathcal{R}_\infty^0(\phi, \phi_s)$ and the minimum is 0, which is also the smallest possible prediction risk. As $\mathcal{R}_1^\lambda(\phi, \phi) = 0$ for $\lambda = 0$, the conclusion still holds.

When $\phi \in (1, \infty)$, we know that $\phi_s^* \in (1, \infty)$ from Patil et al. (2022a, Proposition 5.7). Analogous to Part (1), we have that $\min_{\phi_s \geq \phi} \mathcal{R}_\infty^0(\phi, \phi_s) = \min_{\lambda \geq 0} \mathcal{R}_1^\lambda(\phi, \phi)$.

Part (4) Relationship between ϕ^* and λ^* :

Each pair of the optimal solution (ϕ^*, λ^*) satisfies that $v(0; \phi_s^*) = v(-\lambda^*; \phi) =: v^*$, where $v(0; \phi_s^*)$ and $v(-\lambda^*; \phi)$ are non-negative solutions to the following fixed-point equations:

$$\frac{1}{v(0; \phi_s^*)} = \phi_s^* \int \frac{r}{1 + v(0; \phi_s^*)r} dH(r), \quad \frac{1}{v(-\lambda^*; \phi)} = \lambda^* + \phi \int \frac{r}{1 + v(-\lambda^*; \phi)r} dH(r)$$

From the previous parts, if $\text{SNR} = 0$, then $\lambda^* = \phi_s^* = +\infty$ and $v^* = 0$. Otherwise, we have

$$\frac{1}{v^*} = \phi_s^* \int \frac{r}{1 + v^*r} dH(r) = \lambda^* + \phi \int \frac{r}{1 + v^*r} dH(r),$$

which yields that

$$\lambda^* = (\phi_s^* - \phi) \int \frac{r}{1 + v^*r} dH(r).$$

Part (5) Individual and joint optimization:

Note that from Lemma F.10 (2) and Lemma F.11 (1), the function $\phi_s \mapsto v(-\lambda; \phi_s)$ is decreasing with the range $[0, \lambda^{-1}]$ for $\lambda \in [0, \infty]$. Then the function $(\lambda, \phi_s) \mapsto v(-\lambda; \phi_s)$ has the range $[0, +\infty]$, which is the same as $v(0; \phi_s)$. It follows that $\min_{\phi_s \geq \phi} \mathcal{R}_\infty^0(\phi, \phi_s) = \min_{\phi_s \geq \phi, \lambda \geq 0} \mathcal{R}_1^\lambda(\phi, \phi)$ by the analogous argument in Part (1)-(3).

When $\lambda^* = 0$, the curve reduces to a singleton, which is a trivial case. When $\lambda^* > 0$, for any $t \in [0, \lambda^*]$, let $\lambda = \lambda^* - t$ and $\phi_s = \phi + t(\phi_s^* - \phi)/\lambda^*$. Note that

$$\begin{aligned} \frac{1}{v(\lambda; \phi_s)} &= \lambda + \phi_s \int \frac{r}{1 + v^*r} dH(r) \\ &= \lambda^* - t + (\phi + t(\phi_s^* - \phi)/\lambda^*) \int \frac{r}{1 + v(0; \phi_s^*)r} dH(r) \\ &= \lambda^* + \phi \int \frac{r}{1 + v(0; \phi_s^*)r} dH(r) + \frac{t}{\lambda^*} (\phi_s^* - \phi - \lambda) \int \frac{r}{1 + v(0; \phi_s^*)r} dH(r) \\ &= \frac{1}{v^*} + \frac{t}{\lambda^*} \left(\frac{1}{v^*} - \frac{1}{v^*} \right) \\ &= \frac{1}{v^*}, \end{aligned}$$

which implies that $v(\lambda; \phi_s) = v^*$. Then, we have

$$\tilde{c}(-\lambda; \phi_s) = \int \frac{r}{(1 + v(-\lambda; \phi_s)r)^2} dG(r) = \tilde{c}(-\lambda^*; \phi) = \tilde{c}(0; \phi_s^*).$$

and

$$\tilde{v}(-\lambda; \phi, \phi_s) = \frac{\phi \int r^2 (1 + v(-\lambda; \phi_s)r)^{-2} dH(r)}{v(-\lambda; \phi_s)^{-2} - \phi \int r^2 (1 + v(-\lambda; \phi_s)r)^{-2} dH(r)} = \tilde{v}(-\lambda^*; \phi, \phi) = \tilde{v}(0; \phi, \phi_s^*).$$

It then follows that $\mathcal{R}_\infty^\lambda(\phi, \phi_s) = \mathcal{R}_\infty^{\lambda^*}(\phi, \phi) = \mathcal{R}_\infty^0(\phi, \phi_s^*)$, which completes the proof for Theorem 2.3.

Part (6)* Extension of risk equivalence:

Here we extend the results in Theorem 2.3 to a more general equivalence of (λ, ϕ_s) , as indicated following Theorem 2.3 towards the end of Section 2.2.

For any $\bar{\phi}_s \in [\phi, +\infty]$, let $\bar{\lambda} = (\bar{\phi}_s - \phi) \int r(1 + v(0; \phi_s)r)^{-1} dH(r) \geq 0$. Then, we have

$$\frac{1}{v(0; \bar{\phi}_s)} = \bar{\phi}_s \int \frac{r}{1 + v(0; \bar{\phi}_s)r} dH(r) = \bar{\lambda} + \phi \int \frac{r}{1 + v(0; \bar{\phi}_s)r} dH(r),$$

It follows that $v(-\bar{\lambda}; \phi) = v(0; \bar{\phi}_s)$, and consequently, $\mathcal{R}_\infty^0(\phi, \bar{\phi}_s) = \mathcal{R}_\infty^{\bar{\lambda}}(\phi, \phi)$.

□

B. Proofs of results in Section 3

B.1. Proof of Theorem 3.1

To prove Theorem 3.1, we first prove pointwise convergence (over k and λ) as stated in Lemma B.1, which is based on Proposition 3.6 proved in Appendix E.

Lemma B.1 (Consistency of GCV in full ensemble). *Under Assumptions 2.1-2.2, as $k, n, p \rightarrow \infty$, $p/n \rightarrow \phi \in (0, \infty)$ and $p/k \rightarrow \phi_s \in [\phi, +\infty]$, for $\lambda \geq 0$, it holds that*

$$|\text{gcv}_k^\lambda - R_{k, \infty}^\lambda| \xrightarrow{\text{a.s.}} 0. \quad (25)$$

Proof of Lemma B.1. We will first show that proof for $\lambda > 0$ and $\phi_s < \infty$, and then extend the results to these boundary cases.

Recall that from Proposition 3.6, we have

$$\mathcal{G}_\infty^\lambda(\phi, \phi_s) = \frac{\frac{2\phi(2\phi_s - \phi)}{\phi_s^2} \mathcal{T}_2^\lambda(\phi, \phi_s) + \frac{2(\phi_s - \phi)^2}{\phi_s^2} \mathcal{R}_2^\lambda(\phi, \phi_s) - \frac{\phi}{\phi_s} \mathcal{T}_1^\lambda(\phi, \phi_s) - \frac{\phi_s - \phi}{\phi_s} \mathcal{R}_1^\lambda(\phi, \phi_s)}{\mathcal{D}_\infty^\lambda(\phi, \phi_s)}.$$

We next simplify the expression of the numerator:

$$\begin{aligned} & \frac{2\phi(2\phi_s - \phi)}{\phi_s^2} \mathcal{T}_2^\lambda(\phi, \phi_s) + \frac{2(\phi_s - \phi)^2}{\phi_s^2} \mathcal{R}_2^\lambda(\phi, \phi_s) - \frac{\phi}{\phi_s} \mathcal{T}_1^\lambda(\phi, \phi_s) - \frac{\phi_s - \phi}{\phi_s} \mathcal{R}_1^\lambda(\phi, \phi_s) \\ &= \frac{\phi(\phi_s - \phi)}{\phi_s^2} \mathcal{R}_1^\lambda(\phi, \phi_s) + \mathcal{D}^\lambda(\phi_s, \phi_s) \left(\frac{\phi}{\phi_s} \mathcal{R}_1^\lambda(\phi, \phi_s) + \left(\frac{2\phi(\phi_s - \phi)}{\phi_s^2} \frac{1}{\lambda v(-\lambda; \phi_s)} + \frac{\phi^2}{\phi_s^2} \right) \mathcal{R}_\infty^\lambda(\phi, \phi_s) \right) \\ & \quad + \frac{2(\phi_s - \phi)^2}{\phi_s^2} \mathcal{R}_2^\lambda(\phi, \phi_s) - \frac{\phi}{\phi_s} \mathcal{D}^\lambda(\phi_s, \phi_s) \mathcal{R}_1^\lambda(\phi, \phi_s) - \frac{\phi_s - \phi}{\phi_s} \mathcal{R}_1^\lambda(\phi, \phi_s) \\ &= \left(\frac{2\phi(\phi_s - \phi)}{\phi_s^2} \lambda v(-\lambda; \phi_s) + \frac{\phi^2}{\phi_s^2} \lambda^2 v(-\lambda; \phi_s)^2 \right) \mathcal{R}_\infty^\lambda(\phi, \phi_s) + \frac{2(\phi_s - \phi)^2}{\phi_s^2} (\mathcal{R}_2^\lambda(\phi, \phi_s) - \mathcal{R}_1^\lambda(\phi, \phi_s)) \\ &= \left(\frac{(\phi_s - \phi)^2}{\phi_s^2} + \frac{2\phi(\phi_s - \phi)}{\phi_s^2} \lambda v(-\lambda; \phi_s) + \frac{\phi^2}{\phi_s^2} \lambda^2 v(-\lambda; \phi_s)^2 \right) \mathcal{R}_\infty^\lambda(\phi, \phi_s) \\ &= \mathcal{D}_\infty^\lambda(\phi, \phi_s) \mathcal{R}_\infty^\lambda(\phi, \phi_s). \end{aligned}$$

Then, it follows that $\mathcal{G}_\infty^\lambda(\phi, \phi_s) = \mathcal{R}_\infty^\lambda(\phi, \phi_s)$. From Lemma A.2 and Proposition 3.6, we have that $\text{gcv}_k^\lambda \xrightarrow{\text{a.s.}} \mathcal{G}_\infty^\lambda(\phi, \phi_s)$ and $R_{k, \infty}^\lambda \xrightarrow{\text{a.s.}} \mathcal{R}_\infty^\lambda(\phi, \phi_s)$, which finishes the proof. □

We are now ready to prove Theorem 3.1.

Proof of Theorem 3.1. Let $R_{n, k} = \text{gcv}_{k, \infty}^\lambda - R_{k, \infty}^\lambda$ for $n \in \mathbb{N}$ and $k \in \mathcal{K}_n$. From Lemma B.1 we have that $R_{n, k} \xrightarrow{\text{a.s.}} 0$ as $k, n, p \rightarrow \infty$, $p/n \rightarrow \phi \in (0, \infty)$ and $p/k \rightarrow \phi_s \in [\phi, \infty]$. Here we view k and p as k_n and p_n that are indexed by n . Then from Lemma G.5 (1) the conclusion follows. □

B.2. Proof of Corollary 3.2

Proof of Corollary 3.2. From Theorem 3.1, we have

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

This implies that

$$\begin{aligned} \min_{k \in \mathcal{K}_n} \text{gcv}_k^0 &= \min_{k \in \mathcal{K}_n} \mathcal{R}_\infty^0(p/n, p/k) \left(1 + \frac{\text{gcv}_k^0 - \mathcal{R}_\infty^0(p/n, p/k)}{\mathcal{R}_\infty^0(p/n, p/k)} \right) \\ &\leq \min_{k \in \mathcal{K}_n} \mathcal{R}_\infty^0(p/n, p/k) \left(1 \pm \max_{k \in \mathcal{K}_n} \left| \frac{\text{gcv}_k^0 - \mathcal{R}_\infty^0(p/n, p/k)}{\mathcal{R}_\infty^0(p/n, p/k)} \right| \right) \\ &\leq \min_{k \in \mathcal{K}_n} \mathcal{R}_\infty^0(p/n, p/k) \left(1 \pm \frac{1}{\sigma^2} \max_{k \in \mathcal{K}_n} |\text{gcv}_k^0 - \mathcal{R}_\infty^0(p/n, p/k)| \right) \\ &\xrightarrow{\text{a.s.}} \inf_{\phi_s \geq \phi} \mathcal{R}_\infty^0(\phi, \phi_s) \\ &= \inf_{\phi_s \geq \phi, \lambda \geq 0} \mathcal{R}_\infty^\lambda(\phi, \phi_s), \end{aligned}$$

where the last equality is from Theorem 2.3. This finishes the proof. \square

B.3. Proof of Proposition 3.3

Proof of Proposition 3.3. From the proof of Lemma 3.4, we have

$$\frac{1}{k} \text{tr}(\mathbf{M}_m \boldsymbol{\Sigma}_m) \xrightarrow{\text{a.s.}} (1 - \lambda v(-\lambda; \phi_s)).$$

Then, as $k, n, p \rightarrow \infty, p/n \rightarrow \phi$ and $p/k \rightarrow \phi_s$, we have

$$\begin{aligned} D_{k,2}^\lambda &= \left(1 - \frac{1}{|I_1 \cup I_2|} \frac{1}{2} \sum_{m=1}^2 \text{tr}(\mathbf{M}_m \boldsymbol{\Sigma}_m) \right)^2 \\ &= \left(1 - \frac{k}{|I_1 \cup I_2|} \frac{1}{k} \frac{1}{2} \sum_{m=1}^2 \text{tr}(\mathbf{M}_m \boldsymbol{\Sigma}_m) \right)^2 \\ &\xrightarrow{\text{a.s.}} \left(1 - \frac{\phi_s}{2\phi_s - \phi} (1 - \lambda v(-\lambda; \phi_s)) \right)^2 =: \mathcal{D}_2^\lambda(\phi, \phi_s). \end{aligned}$$

where the convergence of $k/|I_1 \cup I_2|$ is from Lemma G.2. It then follows that

$$\text{gcv}_{k,2}^\lambda = \frac{T_{k,2}^\lambda}{D_{k,2}^\lambda} \xrightarrow{\text{a.s.}} \mathcal{G}_2^\lambda(\phi, \phi_s) := \frac{\mathcal{T}_2^\lambda}{\mathcal{D}_2^\lambda},$$

where \mathcal{T}_2^λ defined in (32) has the following expression:

$$\begin{aligned} \mathcal{T}_2^\lambda(\phi, \phi_s) &= \frac{1}{2} \frac{\phi_s - \phi}{2\phi_s - \phi} \mathcal{R}_1^\lambda(\phi, \phi_s) + \frac{1}{2} \mathcal{D}^\lambda(\phi_s, \phi_s) \\ &\left(\frac{\phi_s}{2\phi_s - \phi} \mathcal{R}_1^\lambda(\phi, \phi_s) + \left(\frac{2(\phi_s - \phi)}{2\phi_s - \phi} \frac{1}{\lambda v(-\lambda; \phi_s)} + \frac{\phi}{2\phi_s - \phi} \right) \mathcal{R}_\infty^\lambda(\phi, \phi_s) \right). \end{aligned}$$

On the other hand, we have

$$\mathcal{R}_2^\lambda(\phi, \phi_s) = \frac{1}{2} \mathcal{R}_1^\lambda(\phi, \phi_s) + \frac{1}{2} \mathcal{R}_\infty^\lambda(\phi, \phi_s).$$

Note that $\mathcal{R}_1^\lambda(\phi, \phi_s) > \mathcal{R}_\infty^\lambda(\phi, \phi_s) > \sigma^2$ when $\phi < \phi_s < \infty$ because $\rho^2 > 0$. When $\lambda = 0$ and $\phi_s > 1$, we have $\lambda v(-\lambda; \phi_s) = 0$ and

$$\mathcal{G}_2^0(\phi, \phi_s) = \frac{\frac{\phi_s - \phi}{2\phi_s - \phi} \mathcal{R}_1^\lambda(\phi, \phi_s) + \frac{2(\phi_s - \phi)}{2\phi_s - \phi} \mathcal{R}_\infty^\lambda(\phi, \phi_s)}{2 \left(1 - \frac{\phi_s}{2\phi_s - \phi}\right)} = \frac{1}{2} \cdot \frac{2\phi_s - \phi}{\phi_s - \phi} (\mathcal{R}_1^\lambda(\phi, \phi_s) + 2\mathcal{R}_\infty^\lambda(\phi, \phi_s)).$$

It follows that

$$\mathcal{G}_2^0(\phi, \phi_s) - \mathcal{R}_2^\lambda(\phi, \phi_s) = \frac{1}{2(\phi_s - \phi)} (\phi_s \mathcal{R}_1^\lambda(\phi, \phi_s) + (3\phi_s - \phi) \mathcal{R}_\infty^\lambda(\phi, \phi_s)) =: c,$$

and $\text{gcv}_{k,2}^\lambda - \mathcal{R}_2^\lambda(\phi, \phi_s) \xrightarrow{p} c > 0$, which completes the proof. \square

C. Proof of Lemma 3.4 (convergence of the GCV denominator functional)

Proof of Lemma 3.4. By the definition, the smooth matrix for $M = \infty$ is given by

$$\mathbf{S}_{\lambda, \infty}(\{\mathcal{D}_{I_\ell}\}_{\ell=1}^M) := \lim_{M \rightarrow \infty} \frac{1}{M} \sum_{\ell=1}^M \mathbf{S}_\lambda(\mathcal{D}_{I_\ell}).$$

For the denominator, note that for any fixed $n \in \mathbb{N}$, as $M \rightarrow \infty$, $I_{1:M} \xrightarrow{\text{a.s.}} [n]$. Then from Lemma C.1 (stated and proved below),

$$\frac{1}{|I_{1:M}|} \text{tr}(\mathbf{S}_\lambda(\mathcal{D}_{I_\ell})) \xrightarrow{\text{a.s.}} \frac{1}{n} \text{tr} \left(\mathbf{X} \mathbf{M}_\ell \frac{\mathbf{X}^\top \mathbf{L}_\ell}{k} \right) = \frac{1}{n} \text{tr}(\mathbf{M}_\ell \boldsymbol{\Sigma}_\ell) \xrightarrow{\text{a.s.}} \frac{\phi}{\phi_s} (1 - \lambda v(-\lambda; \phi_s)).$$

By continuous mapping theorem, we have

$$\left(1 - \frac{1}{n} \text{tr}(\mathbf{S}_{\lambda, \infty})\right)^2 \xrightarrow{\text{a.s.}} \mathcal{D}_\infty^\lambda(\phi, \phi_s) := \left(\frac{\phi_s - \phi}{\phi_s} + \frac{\phi}{\phi_s} \lambda v(-\lambda; \phi_s)\right)^2. \quad (26)$$

\square

Lemma C.1 (Deterministic approximation of the denominator functional). *Under Assumption 2.1, for all $m \in [M]$ and $I_m \in \mathcal{I}_k$, let $\widehat{\boldsymbol{\Sigma}}_m = \mathbf{X}^\top \mathbf{L}_m \mathbf{X} / k$, $\mathbf{L}_m \in \mathbb{R}^{n \times n}$ be a diagonal matrix with $(\mathbf{L}_m)_l = 1$ if $l \in I_m$ and 0 otherwise, and $\mathbf{M}_m = (\mathbf{X}^\top \mathbf{L}_m \mathbf{X} / k + \lambda \mathbf{I}_p)^{-1}$. Then, it holds that for all $m \in [M]$ and $I_m \in \mathcal{I}_k$,*

$$\frac{1}{n} \text{tr}(\mathbf{M}_m \widehat{\boldsymbol{\Sigma}}_m) \xrightarrow{\text{a.s.}} \frac{\phi}{\phi_s} (1 - \lambda v(-\lambda; \phi_s)),$$

as $n, k, p \rightarrow \infty$, $p/n \rightarrow \phi \in (0, \infty)$, and $p/k \rightarrow \phi_s \in [\phi, \infty)$, where the nonnegative constant $\tilde{v}(\lambda; \phi, \phi_s)$ is as defined in (19).

Proof of Lemma C.1. Note that $\mathbf{M}_m \widehat{\boldsymbol{\Sigma}}_m = \mathbf{I}_p - \lambda \mathbf{M}_m$. From Corollary F.5, we have that $\lambda \mathbf{M}_m \simeq (v(-\lambda; \phi_s) \boldsymbol{\Sigma} + \mathbf{I}_p)^{-1}$. Then by Lemma F.3 (4), it follows that

$$\begin{aligned} \frac{1}{n} \text{tr}(\mathbf{M}_m \boldsymbol{\Sigma}_m) &\xrightarrow{\text{a.s.}} \phi \lim_{p \rightarrow \infty} \frac{1}{p} \text{tr}(\mathbf{I}_p - (v(-\lambda; \phi_s) \boldsymbol{\Sigma} + \mathbf{I}_p)^{-1}) \\ &= \phi \lim_{p \rightarrow \infty} \left(1 - \int \frac{1}{1 + v(-\lambda; \phi_s) r} dH_p(r)\right) \\ &= \phi \int \frac{v(-\lambda; \phi_s) r}{1 + v(-\lambda; \phi_s) r} dH(r) \\ &= \frac{\phi}{\phi_s} (1 - \lambda v(-\lambda; \phi_s)) \end{aligned}$$

where in the second last line, we used the fact that H_p and H have compact supports and Assumption 2.2 and the last equality is due to the definition of $v(-\lambda; \phi_s)$ in (18). \square

D. Proof of Lemma 3.5 (convergence of the GCV numerator functional)

Proof of Lemma 3.5. For any $m \in [M]$, let I_m be a sample from \mathcal{I}_k , and $\mathbf{L}_{I_m} \in \mathbb{R}^{n \times n}$ be a diagonal matrix with $(\mathbf{L}_{I_m})_{ll} = 1$ if $l \in I_m$ and 0 otherwise. The ingredient estimator takes the form:

$$\begin{aligned} \tilde{\beta}_{k,M}^\lambda(\{\mathcal{D}_{I_\ell}\}_{\ell=1}^M) &= \frac{1}{M} \sum_{m=1}^M \hat{\beta}^\lambda(\mathcal{D}_{I_m}) \\ &= \frac{1}{M} \sum_{m=1}^M (\mathbf{X}^\top \mathbf{L}_{I_m} \mathbf{X} / k + \lambda \mathbf{I}_p)^{-1} (\mathbf{X}^\top \mathbf{L}_{I_m} \mathbf{y} / k) \\ &= \frac{1}{M} \sum_{m=1}^M \left[\left(\frac{\mathbf{X}^\top \mathbf{L}_{I_m} \mathbf{X}}{k} + \lambda \mathbf{I}_p \right)^{-1} \frac{\mathbf{X}^\top \mathbf{L}_{I_m}}{k} \beta_0 + \left(\frac{\mathbf{X}^\top \mathbf{L}_{I_m} \mathbf{X}}{k} + \lambda \mathbf{I}_p \right)^{-1} \frac{\mathbf{X}^\top \mathbf{L}_{I_m}}{k} \boldsymbol{\epsilon} \right]. \end{aligned}$$

We will write $\tilde{\beta}_{\lambda,M} = \tilde{\beta}_{k,M}^\lambda$ and $\mathbf{L}_m = \mathbf{L}_{I_m}$ for simplicity when they are clear from the context. The set operation will be propagated to such notations, e.g., $\mathbf{L}_{m \cup l} = \mathbf{L}_{I_m \cup I_l}$, $\mathbf{L}_{m \cap l} = \mathbf{L}_{I_m \cap I_l}$, $\mathbf{L}_{m \setminus l} = \mathbf{L}_{I_m \setminus I_l}$, etc. Let $\mathbf{M}_m = (\mathbf{X}^\top \mathbf{L}_{I_m} \mathbf{X} / k + \lambda \mathbf{I}_p)^{-1}$ for $m \in [M]$, we have

$$\tilde{\beta}_{\lambda,M} = \frac{1}{M} \sum_{m=1}^M (\mathbf{I}_p - \lambda \mathbf{M}_m) \beta_0 + \frac{1}{M} \sum_{m=1}^M \mathbf{M}_m (\mathbf{X}^\top \mathbf{L}_m / k) \boldsymbol{\epsilon}. \quad (27)$$

The proof follows by combing the squared error decomposition in Lemma D.1, with the component convergence of test errors in Lemma D.2 and of train errors in Lemma D.3. To prove Lemma D.3, we further make of the component concentration results presented in Appendices D.4 and D.5, and component deterministic approximation results presented in Appendix D.5. \square

D.1. Decomposition of the mean squared error (Lemma D.1)

Lemma D.1 (Decomposition of the mean squared error for the M -ensemble estimator). *For a dataset \mathcal{D}_n , let $\mathbf{X} \in \mathbb{R}^{n \times p}$ and $\mathbf{y} \in \mathbb{R}^n$ be the design matrix and response vector. Let $\mathbf{L}_I \in \mathbb{R}^{n \times n}$ be a diagonal matrix with $(\mathbf{L}_I)_{ll} = 1$ if $l \in I$ and 0 otherwise. Then the mean squared error evaluated on \mathcal{D}_n decomposes as*

$$\begin{aligned} \|\mathbf{y} - \mathbf{X} \tilde{\beta}_k^\lambda\|_2^2 &= -\mathbb{E} \left[\text{Err}_{\text{train}}(\hat{\beta}_k^\lambda(\mathcal{D}_I)) + \text{Err}_{\text{test}}(\hat{\beta}_k^\lambda(\{\mathcal{D}_I\})) \mid \mathcal{D}_n \right] \\ &\quad + 2\mathbb{E} \left[\text{Err}_{\text{train}}(\tilde{\beta}_{k,2}^\lambda(\{\mathcal{D}_I, \mathcal{D}_J\})) + \text{Err}_{\text{test}}(\tilde{\beta}_{k,2}^\lambda(\{\mathcal{D}_I, \mathcal{D}_J\})) \mid \mathcal{D}_n \right] \end{aligned} \quad (28)$$

where the training and test errors are defined by

$$\begin{aligned} \text{Err}_{\text{train}}(\hat{\beta}_k^\lambda(\mathcal{D}_{I_\ell})) &= \|\mathbf{L}_{I_\ell} (\mathbf{y} - \mathbf{X} \hat{\beta}_k^\lambda(\mathcal{D}_{I_\ell}))\|_2^2 \\ \text{Err}_{\text{test}}(\hat{\beta}_k^\lambda(\mathcal{D}_{I_\ell})) &= \|\mathbf{L}_{I_\ell^c} (\mathbf{y} - \mathbf{X} \hat{\beta}_k^\lambda(\mathcal{D}_{I_\ell}))\|_2^2 \\ \text{Err}_{\text{train}}(\tilde{\beta}_{k,2}^\lambda(\{\mathcal{D}_{I_m}, \mathcal{D}_{I_\ell}\})) &= \|\mathbf{L}_{I_m \cup I_\ell} (\mathbf{y} - \mathbf{X} \tilde{\beta}_{k,2}^\lambda(\{\mathcal{D}_{I_m}, \mathcal{D}_{I_\ell}\}))\|_2^2 \\ \text{Err}_{\text{test}}(\tilde{\beta}_{k,2}^\lambda(\{\mathcal{D}_{I_m}, \mathcal{D}_{I_\ell}\})) &= \|\mathbf{L}_{I_m^c \cap I_\ell^c} (\mathbf{y} - \mathbf{X} \tilde{\beta}_{k,2}^\lambda(\{\mathcal{D}_{I_m}, \mathcal{D}_{I_\ell}\}))\|_2^2. \end{aligned} \quad (29)$$

From Lemma D.1, the numerator of the GCV estimate for a M -ensemble estimator decomposes into a linear combination of the training and test error of all possible 1-ensemble and 2-ensemble estimators. Then the asymptotics of the numerator can be obtained, if we can show that the limits of all components exist and their linear combination remains invariable when M goes off to infinity.

Proof of Lemma D.1. We first decompose the training error into the linear combination of the mean squared errors (evaluated on \mathcal{D}_n) for 1-ensemble and 2-ensemble estimators:

$$\frac{1}{n} \|\mathbf{y} - \mathbf{X} \tilde{\beta}_{k,M}^\lambda\|_2^2$$

$$\begin{aligned}
 &= \frac{1}{n} \sum_{i=1}^n \left(\frac{1}{M} \sum_{\ell=1}^M (y_i - \mathbf{x}_i \widehat{\boldsymbol{\beta}}_k^\lambda(\mathcal{D}_{I_\ell})) \right)^2 \\
 &= \frac{1}{n} \sum_{i=1}^n \frac{1}{M^2} \sum_{\ell=1}^M (y_i - \mathbf{x}_i \widehat{\boldsymbol{\beta}}_k^\lambda(\mathcal{D}_{I_\ell}))^2 + \frac{1}{n} \sum_{i=1}^n \frac{1}{M^2} \sum_{\substack{m, \ell \in [M] \\ i \neq j}} (y_i - \mathbf{x}_i \widehat{\boldsymbol{\beta}}_k^\lambda(\mathcal{D}_{I_m})) (y_i - \mathbf{x}_i \widehat{\boldsymbol{\beta}}_k^\lambda(\mathcal{D}_{I_\ell})) \\
 &= \frac{1}{n} \frac{1}{M^2} \sum_{\ell=1}^M \|\mathbf{y} - \mathbf{X} \widehat{\boldsymbol{\beta}}_k^\lambda(\mathcal{D}_{I_\ell})\|_2^2 \\
 &\quad + \frac{1}{n} \sum_{i=1}^n \frac{1}{M^2} \sum_{\substack{m, \ell \in [M] \\ i \neq j}} \frac{1}{2} [4(y_i - \mathbf{x}_i \widetilde{\boldsymbol{\beta}}_{\lambda,2}(\{\mathcal{D}_{I_m}, \mathcal{D}_{I_\ell}\}))^2 - (y_i - \mathbf{x}_i \widehat{\boldsymbol{\beta}}_k^\lambda(\mathcal{D}_{I_m}))^2 - (y_i - \mathbf{x}_i \widehat{\boldsymbol{\beta}}_k^\lambda(\mathcal{D}_{I_\ell}))^2] \\
 &= - \left(\frac{1}{M} - \frac{2}{M^2} \right) \sum_{\ell=1}^M \frac{1}{n} \|\mathbf{y} - \mathbf{X} \widehat{\boldsymbol{\beta}}_k^\lambda(\mathcal{D}_{I_\ell})\|_2^2 + \frac{2}{M^2} \sum_{\substack{m, \ell \in [M] \\ i \neq j}} \frac{1}{n} \|\mathbf{y} - \mathbf{X} \widetilde{\boldsymbol{\beta}}_{\lambda,2}(\{\mathcal{D}_{I_m}, \mathcal{D}_{I_\ell}\})\|_2^2.
 \end{aligned}$$

Next, we further decompose the MSE into training and test errors for 1-ensemble and 2-ensemble estimators:

$$\begin{aligned}
 \frac{1}{n} \|\mathbf{y} - \mathbf{X} \widehat{\boldsymbol{\beta}}_k^\lambda(\mathcal{D}_{I_\ell})\|_2^2 &= \frac{1}{n} \|\mathbf{L}_{I_\ell}(\mathbf{y} - \mathbf{X} \widehat{\boldsymbol{\beta}}_k^\lambda(\mathcal{D}_{I_\ell}))\|_2^2 + \frac{1}{n} \|\mathbf{L}_{I_\ell^c}(\mathbf{y} - \mathbf{X} \widehat{\boldsymbol{\beta}}_k^\lambda(\mathcal{D}_{I_\ell}))\|_2^2, \\
 \frac{1}{n} \|\mathbf{y} - \mathbf{X} \widetilde{\boldsymbol{\beta}}_{\lambda,2}(\{\mathcal{D}_{I_m}, \mathcal{D}_{I_\ell}\})\|_2^2 &= \frac{1}{n} \|\mathbf{L}_{I_m \cup I_\ell}(\mathbf{y} - \mathbf{X} \widetilde{\boldsymbol{\beta}}_{\lambda,2}(\{\mathcal{D}_{I_m}, \mathcal{D}_{I_\ell}\}))\|_2^2 + \frac{1}{n} \|\mathbf{L}_{I_m^c \cap I_\ell^c}(\mathbf{y} - \mathbf{X} \widetilde{\boldsymbol{\beta}}_{\lambda,2}(\{\mathcal{D}_{I_m}, \mathcal{D}_{I_\ell}\}))\|_2^2
 \end{aligned}$$

The conclusion then readily follows. \square

D.2. Convergence of test errors (Lemma D.2)

Lemma D.2 (Convergence of test errors). *Under Assumptions 2.1-2.2, for the test error defined in (29) with $I_1, I_2 \stackrel{\text{SRS}}{\sim} \mathcal{I}_k$, we have that as $k, n, p \rightarrow \infty$, $p/n \rightarrow \phi \in (0, \infty)$ and $p/k \rightarrow \phi_s \in [\phi, +\infty]$,*

$$\frac{\text{Err}_{\text{test}}(\widehat{\boldsymbol{\beta}}_k^\lambda(\mathcal{D}_{I_\ell}))}{n - k} \xrightarrow{\text{a.s.}} \mathcal{R}_1^\lambda(\phi, \phi_s) \quad (30)$$

$$\frac{\text{Err}_{\text{test}}(\widetilde{\boldsymbol{\beta}}_{k,2}^\lambda(\{\mathcal{D}_{I_m}, \mathcal{D}_{I_\ell}\}))}{|I_m^c \cap I_\ell^c|} \xrightarrow{\text{a.s.}} \mathcal{R}_2^\lambda(\phi, \phi_s), \quad (31)$$

where the deterministic functions \mathcal{R}_M is defined in Lemma A.2.

Proof of Lemma D.2. From the strong law of large numbers, we have

$$\begin{aligned}
 &\frac{1}{k} \text{Err}_{\text{test}}(\widehat{\boldsymbol{\beta}}_k^\lambda(\mathcal{D}_{I_\ell})) \xrightarrow{\text{a.s.}} \mathbb{E} \left[(y_0 - \mathbf{x}_0^\top \widehat{\boldsymbol{\beta}}_k^\lambda(\mathcal{D}_{I_1}))^2 \mid \mathcal{D}_n \right] \\
 &\frac{1}{|I_m^c \cap I_\ell^c|} \text{Err}_{\text{test}}(\widetilde{\boldsymbol{\beta}}_{k,2}^\lambda(\{\mathcal{D}_{I_m}, \mathcal{D}_{I_\ell}\})) \xrightarrow{\text{a.s.}} \mathbb{E} \left[(y_0 - \mathbf{x}_0^\top \widetilde{\boldsymbol{\beta}}_{k,2}^\lambda(\{\mathcal{D}_{I_m}, \mathcal{D}_{I_\ell}\}))^2 \mid \mathcal{D}_n \right].
 \end{aligned}$$

From Lemma A.2 (Patil et al., 2022a, Theorem 4.1), the condition prediction risks converge in the sense that

$$\begin{aligned}
 &\mathbb{E} \left[(y_0 - \mathbf{x}_0^\top \widehat{\boldsymbol{\beta}}_k^\lambda(\mathcal{D}_{I_1}))^2 \mid \mathcal{D}_n \right] \xrightarrow{\text{a.s.}} \mathcal{R}_1(\phi, \phi_s) \\
 &\mathbb{E} \left[(y_0 - \mathbf{x}_0^\top \widetilde{\boldsymbol{\beta}}_{k,2}^\lambda(\{\mathcal{D}_{I_m}, \mathcal{D}_{I_\ell}\}))^2 \mid \mathcal{D}_n \right] \xrightarrow{\text{a.s.}} \mathcal{R}_2(\phi, \phi_s),
 \end{aligned}$$

and the conclusions follow. \square

D.3. Convergence of train errors (Lemma D.3)

Lemma D.3 (Convergence of train errors). *Under Assumptions 2.1-2.2, for the train error defined in (29) with $I_1, I_2 \stackrel{\text{SRS}}{\sim} \mathcal{I}_k$, we have that as $k, n, p \rightarrow \infty$, $p/n \rightarrow \phi \in (0, \infty)$ and $p/k \rightarrow \phi_s \in [\phi, +\infty)$,*

$$\begin{aligned} \frac{1}{k} \text{Err}_{\text{train}}(\widehat{\beta}_k^\lambda(\mathcal{D}_{I_\ell})) &\xrightarrow{\text{a.s.}} \mathcal{F}_1^\lambda(\phi, \phi_s) := \mathcal{D}^\lambda(\phi_s, \phi_s) \mathcal{R}_1^\lambda(\phi, \phi_s) \\ \frac{1}{|I_m \cup I_\ell|} \text{Err}_{\text{train}}(\widetilde{\beta}_{k,1}^\lambda(\{\mathcal{D}_{I_m}, \mathcal{D}_{I_\ell}\})) &\xrightarrow{\text{a.s.}} \mathcal{F}_2^\lambda(\phi, \phi_s) := \frac{1}{2} \frac{\phi_s - \phi}{2\phi_s - \phi} \mathcal{R}_1^\lambda(\phi, \phi_s) + \frac{1}{2} \mathcal{D}^\lambda(\phi_s, \phi_s) \\ &\quad \left(\frac{\phi_s}{2\phi_s - \phi} \mathcal{R}_1^\lambda(\phi, \phi_s) + \left(\frac{2(\phi_s - \phi)}{2\phi_s - \phi} \frac{1}{\lambda v(-\lambda; \phi_s)} + \frac{\phi}{2\phi_s - \phi} \right) (2\mathcal{R}_2^\lambda(\phi, \phi_s) - \mathcal{R}_1^\lambda(\phi, \phi_s)) \right), \end{aligned} \quad (32)$$

where the deterministic function \mathcal{R}_M is defined in Lemma A.2.

Proof of Lemma D.3. From (27), we have

$$\beta_0 - \widetilde{\beta}_{\lambda, M} = \frac{1}{M} \sum_{m=1}^M \lambda \mathbf{M}_m \beta_0 - \frac{1}{M} \sum_{m=1}^M \mathbf{M}_m (\mathbf{X}^\top \mathbf{L}_m / k) \epsilon. \quad (33)$$

Part (1) Case of $M = 1$:

From (33), the train error can be decomposed as follows:

$$\begin{aligned} \frac{1}{k} \text{Err}_{\text{train}}(\widehat{\beta}_k^\lambda(\mathcal{D}_{I_\ell})) &= \|\mathbf{L}_\ell \epsilon + \mathbf{L}_\ell \mathbf{X} (\beta_0 - \widehat{\beta}_k^\lambda(\mathcal{D}_{I_\ell}))\|_2^2 / k \\ &= \|(\mathbf{L}_\ell - \mathbf{L}_\ell \mathbf{X} \mathbf{M}_\ell \mathbf{X}^\top \mathbf{L}_\ell / k) \epsilon + \lambda \mathbf{L}_\ell \mathbf{X} \mathbf{M}_\ell \beta_0\|_2^2 / k \\ &= T_C + T_B + T_V, \end{aligned}$$

where the constant term T_C , bias term T_B , and the variance term T_V are given by

$$T_C = \frac{2\lambda}{k} \epsilon^\top \mathbf{L}_\ell \left(\mathbf{I}_n - \mathbf{X} \mathbf{M}_\ell \frac{\mathbf{X}^\top \mathbf{L}_\ell}{k} \right)^\top \mathbf{L}_\ell \mathbf{X} \mathbf{M}_\ell \beta_0, \quad (34)$$

$$T_B = \lambda^2 \beta_0^\top \mathbf{M}_\ell \widehat{\Sigma}_\ell \mathbf{M}_\ell \beta_0, \quad (35)$$

$$T_V = \frac{1}{k} \epsilon^\top \mathbf{L}_\ell \left(\mathbf{I}_n - \mathbf{X} \mathbf{M}_\ell \frac{\mathbf{X}^\top \mathbf{L}_\ell}{k} \right)^\top \mathbf{L}_\ell \left(\mathbf{I}_n - \mathbf{X} \mathbf{M}_\ell \frac{\mathbf{X}^\top \mathbf{L}_\ell}{k} \right) \mathbf{L}_\ell \epsilon. \quad (36)$$

Next, we analyze the three terms separately. From Lemmas D.4 and D.5 with $n = k$, we have that $T_C \xrightarrow{\text{a.s.}} 0$, and

$$T_V \xrightarrow{\text{a.s.}} \sigma^2 \left(1 - \frac{2}{k} \text{tr}(\mathbf{M}_m \widehat{\Sigma}_m) + \frac{1}{k} \text{tr}(\mathbf{M}_m \widehat{\Sigma}_m \mathbf{M}_m \widehat{\Sigma}_m) \right) := T_{VT}.$$

Thus, it remains to obtain the deterministic equivalent for the bias term T_B and the trace term T_{VT} .

From Lemma D.6 and Lemma D.7, we have that for all $I_1 \in \mathcal{I}_k$,

$$\begin{aligned} T_B &\xrightarrow{\text{a.s.}} \rho^2 \mathcal{D}^\lambda(\phi_s, \phi_s) (1 + \widetilde{v}(-\lambda; \phi_s, \phi_s)) \widetilde{c}(-\lambda; \phi_s) \\ T_{VT} &\xrightarrow{\text{a.s.}} \sigma^2 \mathcal{D}^\lambda(\phi_s, \phi_s) (1 + \widetilde{v}(-\lambda; \phi_s, \phi_s)). \end{aligned}$$

Then, we have

$$\frac{1}{k} \text{Err}_{\text{train}}(\widehat{\beta}_k^\lambda(\mathcal{D}_{I_\ell})) \xrightarrow{\text{a.s.}} \mathcal{R}_{\lambda, 1}(\phi, \phi_s) \mathcal{D}^\lambda(\phi_s, \phi_s)$$

where $\mathcal{R}_{\lambda, 1}$ and \mathcal{D}^λ are defined in (21) and (26), respectively.

Part (2) Case of $M = 2$:

From (33), the train error can be analogously decomposed as follows:

$$\frac{1}{|I_m \cup I_l|} \text{Err}_{\text{train}}(\tilde{\beta}_{k,2}^\lambda(\{\mathcal{D}_{I_1}, \mathcal{D}_{I_2}\})) = \frac{1}{|I_m \cup I_l|} \|\mathbf{L}_{m \cup l}(\boldsymbol{\epsilon} + \mathbf{X}(\boldsymbol{\beta}_0 - \hat{\boldsymbol{\beta}}_{\lambda,2}))\|_2^2 = T'_C + T'_B + T'_V,$$

where the constant term T_C , bias term T_B , and the variance term T_V are given by

$$T'_C = \frac{\lambda}{2|I_m \cup I_l|} \sum_{i,j \in \{m,l\}} \boldsymbol{\epsilon}^\top \left(\mathbf{I}_n - \mathbf{X} \mathbf{M}_i \frac{\mathbf{X}^\top \mathbf{L}_i}{k} \right)^\top \mathbf{L}_{m \cup l} \mathbf{X} \mathbf{M}_j \boldsymbol{\beta}_0, \quad (37)$$

$$\begin{aligned} T'_B &= \frac{\lambda^2}{4|I_m \cup I_l|} \sum_{i,j \in \{m,l\}} \boldsymbol{\beta}_0^\top \mathbf{M}_i \hat{\boldsymbol{\Sigma}}_{m \cup l} \mathbf{M}_j \boldsymbol{\beta}_0 \\ &= \frac{\lambda^2 k}{4|I_m \cup I_l|} \sum_{i \in \{m,l\}} \boldsymbol{\beta}_0^\top \mathbf{M}_i \hat{\boldsymbol{\Sigma}}_i \mathbf{M}_i \boldsymbol{\beta}_0 + \frac{\lambda^2}{4|I_m \cup I_l|} \sum_{i \in \{m,l\}} |I_{m+l-i} \setminus I_i| \boldsymbol{\beta}_0^\top \mathbf{M}_i \hat{\boldsymbol{\Sigma}}_{(m+l-i) \setminus i} \mathbf{M}_i \boldsymbol{\beta}_0 \\ &\quad + \frac{\lambda^2}{4} \sum_{i \in \{m,l\}} \boldsymbol{\beta}_0^\top \mathbf{M}_i \hat{\boldsymbol{\Sigma}}_{m \cup l} \mathbf{M}_{m+l-i} \boldsymbol{\beta}_0, \end{aligned} \quad (38)$$

$$T'_V = \frac{1}{4|I_m \cup I_l|} \sum_{i,j \in \{m,l\}} \boldsymbol{\epsilon}^\top \left(\mathbf{I}_n - \mathbf{X} \mathbf{M}_i \frac{\mathbf{X}^\top \mathbf{L}_i}{k} \right)^\top \mathbf{L}_{m \cup l} \left(\mathbf{I}_n - \mathbf{X} \mathbf{M}_j \frac{\mathbf{X}^\top \mathbf{L}_j}{k} \right) \boldsymbol{\epsilon}. \quad (39)$$

Next, we analyze the three terms separately. From Lemmas D.4 and D.5, we have that $T_C \xrightarrow{\text{a.s.}} 0$, and

$$\begin{aligned} T'_V &\xrightarrow{\text{a.s.}} \frac{\sigma^2}{4} \sum_{i \in \{m,l\}} \left(1 - \frac{2}{|I_m \cup I_l|} \text{tr}(\mathbf{M}_i \hat{\boldsymbol{\Sigma}}_i) + \frac{1}{k} \text{tr}(\mathbf{M}_i \hat{\boldsymbol{\Sigma}}_i \mathbf{M}_i \hat{\boldsymbol{\Sigma}}_{m \cup l}) \right) \\ &\quad + \frac{\sigma^2}{2} \left(1 - \frac{1}{|I_m \cup I_l|} \sum_{j \in \{m,l\}} \text{tr}(\mathbf{M}_j \hat{\boldsymbol{\Sigma}}_j) + \frac{1}{n} \text{tr}(\mathbf{M}_l \hat{\boldsymbol{\Sigma}}_{m \cap l} \mathbf{M}_l \hat{\boldsymbol{\Sigma}}_{m \cup l}) \right) \\ &= \frac{\phi_s}{2\phi_s - \phi} \frac{T_{VT}}{2} + \frac{\sigma^2}{4} \frac{\phi_s - \phi}{2\phi_s - \phi} \left(2 + \frac{1}{k} \sum_{i \in \{m,l\}} \text{tr}(\mathbf{M}_i \hat{\boldsymbol{\Sigma}}_i \mathbf{M}_i \hat{\boldsymbol{\Sigma}}_{(m+l-i) \setminus i}) \right) \\ &\quad + \frac{\sigma^2}{2} \left(1 - \frac{1}{|I_m \cup I_l|} \sum_{j \in \{m,l\}} \text{tr}(\mathbf{M}_j \hat{\boldsymbol{\Sigma}}_j) + \frac{1}{n} \text{tr}(\mathbf{M}_l \hat{\boldsymbol{\Sigma}}_{m \cap l} \mathbf{M}_m \hat{\boldsymbol{\Sigma}}_{m \cup l}) \right) := T'_{VT}. \end{aligned}$$

Thus, it remains to obtain the deterministic equivalent for the bias term T'_B and the trace term T'_{VT} .

From Lemma D.6 and Lemma D.7, and the convergence of the cardinality from Lemma G.2, we have that for all $I_m, I_l \stackrel{\text{SRS}}{\sim} \mathcal{I}_k$,

$$T'_B \xrightarrow{\text{a.s.}} \frac{\rho^2}{2} \tilde{t}(\phi, \phi_s) \tilde{c}(-\lambda; \phi_s), \quad \text{and} \quad T'_V \xrightarrow{\text{a.s.}} \frac{\sigma^2}{2} \tilde{t}(\phi, \phi_s),$$

where

$$\begin{aligned} \tilde{t}(\phi, \phi_s) &= \frac{\phi_s}{2\phi_s - \phi} \mathcal{D}^\lambda(\phi_s, \phi_s) (1 + \tilde{v}(-\lambda; \phi_s, \phi_s)) + \frac{\phi_s - \phi}{2\phi_s - \phi} (1 + \tilde{v}(-\lambda; \phi_s, \phi_s)) \\ &\quad + \mathcal{D}^\lambda(\phi_s, \phi_s) \left(\frac{2(\phi_s - \phi)}{2\phi_s - \phi} \frac{1}{\lambda v(-\lambda; \phi_s)} + \frac{\phi}{2\phi_s - \phi} \right) (1 + \tilde{v}(-\lambda; \phi, \phi_s)). \end{aligned}$$

Then, we have

$$\frac{1}{|I_m \cup I_l|} \text{Err}_{\text{train}}(\tilde{\beta}_{k,2}^\lambda(\{\mathcal{D}_{I_1}, \mathcal{D}_{I_2}\}))$$

$$\begin{aligned}
 & \xrightarrow{\text{a.s.}} \frac{\rho^2}{2} \tilde{t}(\phi, \phi_s) \tilde{c}(-\lambda; \phi_s) + \frac{\sigma^2}{2} \tilde{t}(\phi, \phi_s) \\
 & = \frac{1}{2} \mathcal{D}^\lambda(\phi_s, \phi_s) \left(\frac{\phi_s}{2\phi_s - \phi} \mathcal{R}_{\lambda,1}(\phi, \phi_s) + \left(\frac{2(\phi_s - \phi)}{2\phi_s - \phi} \frac{1}{\lambda v(-\lambda; \phi_s)} + \frac{\phi}{2\phi_s - \phi} \right) (2\mathcal{R}_{\lambda,2}(\phi, \phi_s) - \mathcal{R}_{\lambda,2}(\phi, \phi_s)) \right) \\
 & \quad + \frac{1}{2} \frac{\phi_s - \phi}{2\phi_s - \phi} \mathcal{R}_1(\phi, \phi_s),
 \end{aligned}$$

which finishes the proof. \square

D.4. Component concentrations

In this subsection, we will show that the cross-term T_C converges to zero and the variance term T_V converges to its corresponding trace expectation.

D.4.1. CONVERGENCE OF THE CROSS TERM

Lemma D.4 (Convergence of the cross term). *Under Assumptions 2.1-2.2, for T_C and T'_C as defined in (34) and (37), we have $T_C \xrightarrow{\text{a.s.}} 0$ and $T'_C \xrightarrow{\text{a.s.}} 0$ as $k, p \rightarrow \infty$ and $p/k \rightarrow \phi_s$.*

Proof of Lemma D.4. We first prove the result for T'_C . Note that

$$T'_C = -\frac{\lambda}{M^2} \cdot \frac{1}{|I_1 \cup I_2|} \left\langle \sum_{m=1}^2 \left(\mathbf{I}_n - \mathbf{X} \frac{\mathbf{M}_m \mathbf{X}^\top \mathbf{L}_m}{k} \right)^\top \mathbf{L}_{m \cup l} \mathbf{X} \sum_{m=1}^2 \mathbf{M}_m \boldsymbol{\beta}_0, \boldsymbol{\epsilon} \right\rangle.$$

We next bound the squared norm:

$$\begin{aligned}
 & \frac{1}{|I_m \cup I_l|} \left\| \frac{1}{2} \sum_{m=1}^2 \left(\mathbf{I}_n - \frac{\mathbf{X} \mathbf{M}_m \mathbf{X}^\top \mathbf{L}_m}{k} \right)^\top \mathbf{L}_{m \cup l} \mathbf{X} \sum_{m=1}^2 \mathbf{M}_m \boldsymbol{\beta}_0 \right\|_2^2 \\
 & \leq \sum_{j=1}^2 \sum_{l=1}^2 \left[\frac{|I_m \cup I_l|}{4k^2} \left\| (\mathbf{M}_j \mathbf{X}^\top \mathbf{L}_j)^\top \widehat{\boldsymbol{\Sigma}}_{m \cup l} \mathbf{M}_l \boldsymbol{\beta}_0 \right\|_2^2 + \frac{1}{4|I_1 \cup I_2|} \left\| \mathbf{L}_j \mathbf{X} \mathbf{M}_l \boldsymbol{\beta}_0 \right\|_2^2 \right] \\
 & \leq \frac{\|\boldsymbol{\beta}_0\|_2^2}{4} \cdot \sum_{j=1}^2 \sum_{l=1}^2 \left[\frac{|I_m \cup I_l|}{k^2} \left\| \mathbf{M}_l \widehat{\boldsymbol{\Sigma}}_{m \cup l} \mathbf{M}_j \mathbf{X}^\top \mathbf{L}_j \mathbf{X} \mathbf{M}_j \widehat{\boldsymbol{\Sigma}}_{m \cup l} \mathbf{M}_l \right\|_{\text{op}} + \frac{k}{|I_1 \cup I_2|} \left\| \widehat{\boldsymbol{\Sigma}}_j \right\|_{\text{op}} \left\| \mathbf{M}_l \right\|_{\text{op}} \right] \\
 & \leq \frac{\|\boldsymbol{\beta}_0\|_2^2}{4} \cdot \sum_{j=1}^2 \sum_{l=1}^2 \left[\frac{|I_m \cup I_l|}{k} \left\| \mathbf{M}_l \right\|_{\text{op}}^2 \left\| \widehat{\boldsymbol{\Sigma}}_l \right\|_{\text{op}}^2 \left\| \mathbf{M}_j (\mathbf{X}^\top \mathbf{L}_j \mathbf{X} / k) \mathbf{M}_j \right\|_{\text{op}} + \frac{k}{|I_1 \cup I_2|} \left\| \mathbf{M}_l \right\|_{\text{op}} \left\| \widehat{\boldsymbol{\Sigma}}_l \right\|_{\text{op}} \right] \\
 & = \frac{\|\boldsymbol{\beta}_0\|_2^2}{4} \cdot \sum_{j=1}^2 \sum_{l=1}^2 \left[\frac{|I_m \cup I_l|}{k} \left\| \mathbf{M}_l \right\|_{\text{op}}^2 \left\| \widehat{\boldsymbol{\Sigma}}_l \right\|_{\text{op}}^2 \left\| \mathbf{M}_j \right\|_{\text{op}} \left\| \mathbf{I}_p - \lambda \mathbf{M}_j \right\|_{\text{op}} + \frac{k}{|I_1 \cup I_2|} \left\| \mathbf{M}_l \right\|_{\text{op}} \left\| \widehat{\boldsymbol{\Sigma}}_l \right\|_{\text{op}} \right] \\
 & \leq \frac{\|\boldsymbol{\beta}_0\|_2^2}{\lambda} \left(\frac{|I_m \cup I_l|}{k\lambda^2} + \frac{k}{|I_1 \cup I_2|} \right) \left\| \widehat{\boldsymbol{\Sigma}}_l \right\|_{\text{op}}^2,
 \end{aligned}$$

where the last inequality is due to the fact that $\|\mathbf{M}_j\|_{\text{op}} \leq 1/\lambda$ and $\|\mathbf{I}_p - \lambda \mathbf{M}_j\|_{\text{op}} \leq 1$. By Assumption 2.2, $\|\boldsymbol{\beta}_0\|_2^2$ is uniformly bounded in p . From Bai & Silverstein (2010), we have $\limsup \left\| \widehat{\boldsymbol{\Sigma}} \right\|_{\text{op}} \leq \limsup \max_{1 \leq i \leq p} s_i^2 \leq r_{\max}(1 + \sqrt{\phi_s})^2$ almost surely as $k, p \rightarrow \infty$ and $p/k \rightarrow \phi_s \in (0, \infty)$. From Lemma G.2, we have that $k/|I_1 \cup I_2| \xrightarrow{\text{a.s.}} k/(2k - k^2/n)$, which is $\phi_s/(2\phi_s - \phi)$ almost surely. Then we have that the square norm is almost surely upper bounded by some constant. Applying Lemma G.3, we thus have that $T'_C \xrightarrow{\text{a.s.}} 0$.

Note that when $I_1 = I_2$, T'_C reduces to T_C , and thus the conclusion for T_C also holds. \square

D.4.2. CONVERGENCE OF THE VARIANCE TERM

Lemma D.5 (Convergence of the variance term). *Under Assumptions 2.1-2.2, let $M \in \mathbb{N}$ and $\widehat{\boldsymbol{\Sigma}} = \mathbf{X}^\top \mathbf{X} / n$. For all $m \in [M]$ and $I_m \in \mathcal{I}_k$, let $\widehat{\boldsymbol{\Sigma}}_m = \mathbf{X}^\top \mathbf{L}_m \mathbf{X} / k$, $\mathbf{L}_m \in \mathbb{R}^{n \times n}$ be a diagonal matrix with $(\mathbf{L}_m)_{ll} = 1$ if $l \in I_m$ and 0*

otherwise, and $\mathbf{M}_m = (\mathbf{X}^\top \mathbf{L}_m \mathbf{X} / k + \lambda \mathbf{I}_p)^{-1}$. Then, for all $m, l \in [M]$ and $m \neq l$, it holds that

$$\begin{aligned} & \frac{1}{k} \boldsymbol{\epsilon}^\top \left(\mathbf{I}_n - \mathbf{X} \mathbf{M}_m \frac{\mathbf{X}^\top \mathbf{L}_m}{k} \right)^\top \mathbf{L}_m \left(\mathbf{I}_n - \mathbf{X} \mathbf{M}_m \frac{\mathbf{X}^\top \mathbf{L}_m}{k} \right) \boldsymbol{\epsilon} \\ & - \sigma^2 \left(1 - \frac{2}{k} \text{tr}(\mathbf{M}_m \widehat{\boldsymbol{\Sigma}}_m) + \frac{1}{k} \text{tr}(\mathbf{M}_m \widehat{\boldsymbol{\Sigma}}_m \mathbf{M}_m \widehat{\boldsymbol{\Sigma}}_m) \right) \xrightarrow{\text{a.s.}} 0, \end{aligned} \quad (40)$$

$$\begin{aligned} & \frac{1}{|I_m \cup I_l|} \boldsymbol{\epsilon}^\top \left(\mathbf{I}_n - \mathbf{X} \mathbf{M}_i \frac{\mathbf{X}^\top \mathbf{L}_i}{k} \right)^\top \mathbf{L}_{m \cup l} \left(\mathbf{I}_n - \mathbf{X} \mathbf{M}_i \frac{\mathbf{X}^\top \mathbf{L}_i}{k} \right) \boldsymbol{\epsilon} \\ & - \sigma^2 \left(1 - \frac{2}{|I_m \cup I_l|} \text{tr}(\mathbf{M}_i \widehat{\boldsymbol{\Sigma}}_i) + \frac{1}{k} \text{tr}(\mathbf{M}_i \widehat{\boldsymbol{\Sigma}}_i \mathbf{M}_i \widehat{\boldsymbol{\Sigma}}_{m \cup l}) \right) \xrightarrow{\text{a.s.}} 0, \end{aligned} \quad (41)$$

$$\begin{aligned} & \frac{1}{|I_m \cup I_l|} \boldsymbol{\epsilon}^\top \left(\mathbf{I}_n - \mathbf{X} \mathbf{M}_m \frac{\mathbf{X}^\top \mathbf{L}_m}{k} \right)^\top \mathbf{L}_{m \cup l} \left(\mathbf{I}_n - \mathbf{X} \mathbf{M}_l \frac{\mathbf{X}^\top \mathbf{L}_l}{k} \right) \boldsymbol{\epsilon} \\ & - \sigma^2 \left(1 - \frac{1}{|I_m \cup I_l|} \sum_{\ell \in \{i, j\}} \text{tr}(\mathbf{M}_\ell \widehat{\boldsymbol{\Sigma}}_\ell) + \frac{1}{n} \text{tr}(\mathbf{M}_i \widehat{\boldsymbol{\Sigma}}_{i \cap j} \mathbf{M}_j \widehat{\boldsymbol{\Sigma}}_{m \cup l}) \right) \xrightarrow{\text{a.s.}} 0, \end{aligned} \quad (42)$$

where $i, j \in \{m, l\}$, $i \neq j$, and $\widehat{\boldsymbol{\Sigma}}_{m \cup l} = \mathbf{X}^\top \mathbf{L}_{m \cup l} \mathbf{X} / |I_m \cup I_l|$, as $n, k, p \rightarrow \infty$, $p/n \rightarrow \phi \in (0, \infty)$, and $p/k \rightarrow \phi_s \in [\phi, \infty)$.

Proof of Lemma D.5. We first prove the last convergence result. Note that

$$\begin{aligned} & \left\| \left(\mathbf{I}_n - \mathbf{X} \mathbf{M}_m \frac{\mathbf{X}^\top \mathbf{L}_m}{k} \right)^\top \mathbf{L}_{m \cup l} \left(\mathbf{I}_n - \mathbf{X} \mathbf{M}_l \frac{\mathbf{X}^\top \mathbf{L}_l}{k} \right) \right\|_{\text{op}} \\ & = \left\| \mathbf{L}_{m \cup l} - \frac{1}{k} \mathbf{L}_m \mathbf{X} \mathbf{M}_m \mathbf{X}^\top \mathbf{L}_{m \cup l} - \frac{1}{k} \mathbf{L}_{m \cup l} \mathbf{X} \mathbf{M}_l \mathbf{X}^\top \mathbf{L}_l + \frac{1}{k^2} \mathbf{L}_m \mathbf{X} \mathbf{M}_m \mathbf{X}^\top \mathbf{L}_{m \cup l} \mathbf{X} \mathbf{M}_l \mathbf{X}^\top \mathbf{L}_l \right\|_{\text{op}} \\ & \leq 1 + \sqrt{\frac{|I_m \cup I_l|}{k}} \sum_{j \in \{m, l\}} \left\| \widehat{\boldsymbol{\Sigma}}_j \right\|_{\text{op}}^{\frac{1}{2}} \left\| \mathbf{M}_j \right\|_{\text{op}} \left\| \widehat{\boldsymbol{\Sigma}}_{m \cup l} \right\|_{\text{op}}^{\frac{1}{2}} + \frac{|I_m \cup I_l|}{k} \left\| \widehat{\boldsymbol{\Sigma}}_m \right\|_{\text{op}}^{\frac{1}{2}} \left\| \mathbf{M}_m \right\|_{\text{op}} \left\| \widehat{\boldsymbol{\Sigma}}_{m \cup l} \right\|_{\text{op}} \left\| \mathbf{M}_l \right\|_{\text{op}} \left\| \widehat{\boldsymbol{\Sigma}}_l \right\|_{\text{op}}^{\frac{1}{2}} \\ & \leq 1 + \frac{2}{\lambda} \sqrt{\frac{|I_m \cup I_l|}{k}} \left\| \widehat{\boldsymbol{\Sigma}}_m \right\|_{\text{op}}^{\frac{1}{2}} \left\| \widehat{\boldsymbol{\Sigma}}_l \right\|_{\text{op}}^{\frac{1}{2}} + \frac{1}{\lambda^2} \frac{|I_m \cup I_l|}{k} \left\| \widehat{\boldsymbol{\Sigma}}_m \right\|_{\text{op}}^{\frac{1}{2}} \left\| \widehat{\boldsymbol{\Sigma}}_{m \cup l} \right\|_{\text{op}} \left\| \widehat{\boldsymbol{\Sigma}}_l \right\|_{\text{op}}^{\frac{1}{2}}. \end{aligned}$$

Now, we have $\limsup \|\widehat{\boldsymbol{\Sigma}}\|_{\text{op}} \leq \limsup \max_{1 \leq i \leq p} s_i^2 \leq r_{\max}(1 + \sqrt{\phi})^2$ almost surely as $n, p \rightarrow \infty$ and $p/n \rightarrow \phi \in (0, \infty)$ from Bai & Silverstein (2010). Similarly, $\limsup \|\widehat{\boldsymbol{\Sigma}}_m\|_{\text{op}} \leq r_{\max}(1 + \sqrt{\phi_s})^2$ almost surely. From Lemma G.2, $|I_m \cup I_l|/k \xrightarrow{\text{a.s.}} (2\phi_s - \phi)/\phi_s$. Then the above quantity is asymptotically upper bounded by some constant as $n, k, p \rightarrow \infty$, $p/n \rightarrow \phi \in (0, \infty)$ and $p/k \rightarrow \phi_s \in [\phi, \infty)$. From Lemma G.4, it follows that

$$\begin{aligned} & \frac{1}{|I_m \cup I_l|} \boldsymbol{\epsilon}^\top \left(\mathbf{I}_n - \mathbf{X} \mathbf{M}_i \frac{\mathbf{X}^\top \mathbf{L}_i}{k} \right)^\top \mathbf{L}_{m \cup l} \left(\mathbf{I}_n - \mathbf{X} \mathbf{M}_j \frac{\mathbf{X}^\top \mathbf{L}_j}{k} \right) \boldsymbol{\epsilon} \\ & - \frac{\sigma^2}{|I_m \cup I_l|} \text{tr} \left(\left(\mathbf{I}_n - \mathbf{X} \mathbf{M}_i \frac{\mathbf{X}^\top \mathbf{L}_i}{k} \right)^\top \mathbf{L}_{m \cup l} \left(\mathbf{I}_n - \mathbf{X} \mathbf{M}_j \frac{\mathbf{X}^\top \mathbf{L}_j}{k} \right) \right) \xrightarrow{\text{a.s.}} 0. \end{aligned}$$

Expanding the trace term above, we have

$$\begin{aligned} & \frac{\sigma^2}{|I_m \cup I_l|} \text{tr} \left(\left(\mathbf{I}_n - \mathbf{X} \mathbf{M}_i \frac{\mathbf{X}^\top \mathbf{L}_i}{k} \right)^\top \mathbf{L}_{m \cup l} \left(\mathbf{I}_n - \mathbf{X} \mathbf{M}_j \frac{\mathbf{X}^\top \mathbf{L}_j}{k} \right) \right) \\ & = \sigma^2 \left(1 - \frac{1}{|I_m \cup I_l|} \sum_{\ell \in \{i, j\}} \text{tr}(\mathbf{M}_\ell \widehat{\boldsymbol{\Sigma}}_\ell) + \frac{|I_i \cap I_j|}{k^2} \text{tr}(\mathbf{M}_i \widehat{\boldsymbol{\Sigma}}_{i \cup j} \mathbf{M}_j \widehat{\boldsymbol{\Sigma}}_{m \cap l}) \right). \end{aligned} \quad (43)$$

Since $I_m, I_l \stackrel{\text{SRS}}{\sim} \mathcal{I}_k$, from Lemma G.2 we have that $|I_m \cap I_l|/k \xrightarrow{\text{a.s.}} k/n$. Then, we have

$$\begin{aligned} & \frac{1}{|I_m \cup I_l|} \boldsymbol{\epsilon}^\top \left(\mathbf{I}_n - \mathbf{X} \mathbf{M}_i \frac{\mathbf{X}^\top \mathbf{L}_i}{k} \right)^\top \mathbf{L}_{m \cup l} \left(\mathbf{I}_n - \mathbf{X} \mathbf{M}_j \frac{\mathbf{X}^\top \mathbf{L}_j}{k} \right) \boldsymbol{\epsilon} \\ & - \sigma^2 \left(1 - \frac{1}{|I_m \cup I_l|} \sum_{\ell \in \{i, j\}} \text{tr}(\mathbf{M}_\ell \widehat{\boldsymbol{\Sigma}}_\ell) + \frac{1}{n} \text{tr}(\mathbf{M}_i \widehat{\boldsymbol{\Sigma}}_{i \cap j} \mathbf{M}_j \widehat{\boldsymbol{\Sigma}}_{m \cup l}) \right) \xrightarrow{\text{a.s.}} 0, \end{aligned}$$

and thus (42) follows.

Setting $i = j$ in (43) yields (41).

Finally, setting $i = j = l = m$ in (43) finishes the proof for (40). \square

D.5. Component deterministic approximations

D.5.1. DETERMINISTIC APPROXIMATION OF THE BIAS FUNCTIONAL

Lemma D.6 (Deterministic approximation of the bias functional). *Under Assumptions 2.1-2.2, for all $m \in [M]$ and $I_m \in \mathcal{I}_k$, let $\widehat{\boldsymbol{\Sigma}}_m = \mathbf{X}^\top \mathbf{L}_m \mathbf{X}/k$, $\mathbf{L}_m \in \mathbb{R}^{n \times n}$ be a diagonal matrix with $(\mathbf{L}_m)_{ll} = 1$ if $l \in I_m$ and 0 otherwise, and $\mathbf{M}_m = (\mathbf{X}^\top \mathbf{L}_m \mathbf{X}/k + \lambda \mathbf{I}_p)^{-1}$. Then, it holds that*

1. For all $m \in [M]$,

$$\lambda^2 \boldsymbol{\beta}_0^\top \mathbf{M}_m \widehat{\boldsymbol{\Sigma}}_m \mathbf{M}_m \boldsymbol{\beta}_0 \xrightarrow{\text{a.s.}} \rho^2 \lambda^2 v(-\lambda; \phi_s)^2 (1 + \tilde{v}(-\lambda; \phi_s, \phi_s)) \tilde{c}(-\lambda; \phi_s).$$

2. For all $m, l \in [M]$, $m \neq l$ and $I_m, I_l \stackrel{\text{SRSWR}}{\sim} \mathcal{I}_k$,

$$\lambda^2 \boldsymbol{\beta}_0^\top \mathbf{M}_l \widehat{\boldsymbol{\Sigma}}_{m \setminus l} \mathbf{M}_l \boldsymbol{\beta}_0 \xrightarrow{\text{a.s.}} \rho^2 (1 + \tilde{v}(-\lambda; \phi_s, \phi_s)) \tilde{c}(-\lambda; \phi_s). \quad (44)$$

3. For all $m, l \in [M]$, $m \neq l$ and $I_m, I_l \stackrel{\text{SRSWR}}{\sim} \mathcal{I}_k$,

$$\lambda^2 \boldsymbol{\beta}_0^\top \mathbf{M}_m \widehat{\boldsymbol{\Sigma}}_{m \cup l} \mathbf{M}_l \boldsymbol{\beta}_0 \xrightarrow{\text{a.s.}} \rho^2 \lambda^2 \left(\frac{2(\phi_s - \phi)}{2\phi_s - \phi} \frac{1}{\lambda v(-\lambda; \phi_s)} + \frac{\phi}{2\phi_s - \phi} \right) v(-\lambda; \phi_s)^2 (1 + \tilde{v}(-\lambda; \phi, \phi_s)) \tilde{c}(-\lambda; \phi_s), \quad (45)$$

where $\widehat{\boldsymbol{\Sigma}}_{m \cup l} = |I_m \cup I_l|^{-1} \mathbf{X}^\top \mathbf{L}_{m \cup l} \mathbf{X}$, as $n, k, p \rightarrow \infty$, $p/n \rightarrow \phi \in (0, \infty)$, and $p/k \rightarrow \phi_s \in [\phi, \infty)$, where $\phi_0 = \phi_s^2/\phi$, T_B is as defined in (35), and the nonnegative constants $\tilde{v}(-\lambda; \phi, \phi_s)$ and $\tilde{c}(-\lambda; \phi_s)$ are as defined in (19).

Proof of Lemma D.6. We split the proof into different parts.

Part (1) From Lemma F.6 (2) (with $\mathbf{A} = \mathbf{I}_p$), we have that

$$\lambda^2 \mathbf{M}_m \widehat{\boldsymbol{\Sigma}}_m \mathbf{M}_m \simeq v(-\lambda; \phi_s)^2 (1 + \tilde{v}(-\lambda; \phi_s, \phi_s)) \cdot (v(-\lambda; \phi_s) \boldsymbol{\Sigma} + \mathbf{I}_p)^{-1} \boldsymbol{\Sigma} (v(-\lambda; \phi_s) \boldsymbol{\Sigma} + \mathbf{I}_p)^{-1}. \quad (46)$$

By the definition of deterministic equivalent, we have

$$\begin{aligned} \lambda^2 \boldsymbol{\beta}_0^\top \mathbf{M}_m \widehat{\boldsymbol{\Sigma}}_m \mathbf{M}_m \boldsymbol{\beta}_0 & \xrightarrow{\text{a.s.}} \lim_{p \rightarrow \infty} v(-\lambda; \phi_s)^2 (1 + \tilde{v}(-\lambda; \phi_s, \phi_s)) \sum_{i=1}^p \frac{r_i}{(1 + r_i v(-\lambda; \phi_s))^2} (\boldsymbol{\beta}_0^\top w_i)^2 \\ & = \lim_{p \rightarrow \infty} \|\boldsymbol{\beta}_0\|_2^2 v(-\lambda; \phi_s)^2 (1 + \tilde{v}(-\lambda; \phi_s, \phi_s)) \int \frac{r}{(1 + v(-\lambda; \phi_s) r)^2} dG_p(r) \\ & = \rho^2 v(-\lambda; \phi_s)^2 (1 + \tilde{v}(-\lambda; \phi_s, \phi_s)) \int \frac{r}{(1 + v(-\lambda; \phi_s) r)^2} dG(r), \end{aligned} \quad (47)$$

where the last equality holds since G_p and G have compact supports and invoking Assumption 2.2.

Part (2) From Lemma F.8 (1), we have

$$M_l \widehat{\Sigma}_{m \setminus l} M_l \simeq M_l \Sigma M_l.$$

Then, from Patil et al. (2022a, Lemma S.2.4), the conclusion follows.

Part (3) For the cross term, it suffices to derive the deterministic equivalent of $\beta_0^\top M_1 \widehat{\Sigma}_{1 \cup 2} M_2 \beta_0$. We begin with analyzing the deterministic equivalent of $M_1 \widehat{\Sigma}_{1 \cup 2} M_2$. Let $i_0 = \text{tr}(\mathbf{L}_1 \mathbf{L}_2)$ be the number of shared samples between \mathcal{D}_{I_1} and \mathcal{D}_{I_2} , we use the decomposition

$$M_j^{-1} = \frac{i_0}{k} (\widehat{\Sigma}_0 + \lambda \mathbf{I}_p) + \frac{k - i_0}{k} (\widehat{\Sigma}_j^{\text{ind}} + \lambda \mathbf{I}_p), \quad j = 1, 2,$$

where $\widehat{\Sigma}_0 = \mathbf{X}^\top \mathbf{L}_1 \mathbf{L}_2 \mathbf{X} / i_0$ and $\widehat{\Sigma}_j^{\text{ind}} = \mathbf{X}^\top (\mathbf{L}_j - \mathbf{L}_1 \mathbf{L}_2) \mathbf{X} / (k - i_0)$ are the common and individual covariance estimators of the two datasets. Let $\mathbf{N}_0 = (\widehat{\Sigma}_0 + \lambda \mathbf{I}_p)^{-1}$ and $\mathbf{N}_j = (\widehat{\Sigma}_j^{\text{ind}} + \lambda \mathbf{I}_p)^{-1}$ for $j = 1, 2$. Then

$$M_j = \left(\frac{i_0}{k} \mathbf{N}_0^{-1} + \frac{k - i_0}{k} \mathbf{N}_j^{-1} \right)^{-1}, \quad j = 1, 2, \quad (48)$$

where the equalities hold because \mathbf{N}_0 is invertible when $\lambda > 0$. Note that

$$\begin{aligned} \widehat{\Sigma}_{1 \cup 2} &= \frac{k}{2k - i_0} \widehat{\Sigma}_1 + \frac{k}{2k - i_0} \widehat{\Sigma}_2 - \frac{i_0}{2k - i_0} \widehat{\Sigma}_0, \\ &= \frac{k}{2k - i_0} \sum_{j=1}^2 (M_j^{-1} - \lambda \mathbf{I}_p) - \frac{i_0}{2k - i_0} \widehat{\Sigma}_0. \end{aligned}$$

We have that

$$\lambda^2 M_1 \widehat{\Sigma}_{1 \cup 2} M_2 = \frac{k}{2k - i_0} \lambda^2 \sum_{j=1}^2 M_j - \frac{2k}{2k - i_0} \lambda^3 M_1 M_2 - \frac{i_0}{2k - i_0} \lambda^2 M_1 \widehat{\Sigma}_0 M_2. \quad (49)$$

Next, we derive the deterministic equivalents for the three terms in (49). From Corollary F.5, the first term admits

$$\lambda M_j \simeq (v(-\lambda; \phi_s) \Sigma + \mathbf{I}_p)^{-1}. \quad (50)$$

Note that

$$\lambda^2 M_1 M_2 \simeq (v(-\lambda; \phi_s) \Sigma + \mathbf{I}_p)^{-1} (\tilde{v}(-\lambda; \phi, \phi_s, \mathbf{I}_p) \Sigma + \mathbf{I}_p) (v(-\lambda; \phi_s) \Sigma + \mathbf{I}_p)^{-1} \quad (51)$$

$$\simeq (v(-\lambda; \phi_s) \Sigma + \mathbf{I}_p)^{-2} (\tilde{v}(-\lambda; \phi, \phi_s, \mathbf{I}_p) \Sigma + \mathbf{I}_p), \quad (52)$$

where

$$\tilde{v}(-\lambda; \phi, \phi_s, \mathbf{I}_p) = \frac{\phi \int \frac{r}{(1 + v(-\lambda; \phi_s) r)^2} dH(r)}{v(-\lambda; \phi_s)^{-2} - \phi \int \frac{r^2}{(1 + v(-\lambda; \phi_s) r)^2} dH(r)}.$$

For the third term,

$$M_1 \widehat{\Sigma}_0 M_2 \simeq \tilde{v}_v(-\lambda; \phi, \phi_s) (v(-\lambda; \phi_s) \Sigma + \mathbf{I}_p)^{-2} \Sigma, \quad (53)$$

where

$$\tilde{v}_v(-\lambda; \phi, \phi_s) := \frac{1}{v(-\lambda; \phi_s)^{-2} - \phi \int \frac{r^2}{(1 + v(-\lambda; \phi_s) r)^2} dH(r)}.$$

Combining (50)-(53), we get

$$\begin{aligned} \lambda^2 \mathbf{M}_1 \widehat{\Sigma}_{1 \cup 2} \mathbf{M}_2 &\simeq \left(\frac{2\phi_s}{2\phi_s - \phi} (v(-\lambda; \phi_s) - \tilde{v}(-\lambda; \phi, \phi_s, \mathbf{I}_p)) - \frac{\phi}{2\phi_s - \phi} \lambda \tilde{v}_v(-\lambda; \phi, \phi_s) \right) \lambda (v(-\lambda; \phi_s) \Sigma + \mathbf{I}_p)^{-2} \Sigma \\ &= \lambda^2 v(-\lambda; \phi_s)^2 (1 + \tilde{v}(-\lambda; \phi, \phi_s)) \left(\frac{2(\phi_s - \phi)}{2\phi_s - \phi} \frac{1}{\lambda v(-\lambda; \phi_s)} + \frac{\phi}{2\phi_s - \phi} \right) (v(-\lambda; \phi_s) \Sigma + \mathbf{I}_p)^{-2} \Sigma. \end{aligned} \quad (54)$$

The last conclusion then follows analogously as in (47). \square

D.5.2. DETERMINISTIC APPROXIMATION OF THE VARIANCE FUNCTIONAL

Lemma D.7 (Deterministic approximation of the variance functional). *Under Assumptions 2.1-2.2, for all $m \in [M]$ and $I_m \in \mathcal{I}_k$, let $\widehat{\Sigma}_m = \mathbf{X}^\top \mathbf{L}_m \mathbf{X} / k$, $\mathbf{L}_m \in \mathbb{R}^{n \times n}$ be a diagonal matrix with $(\mathbf{L}_m)_{ll} = 1$ if $l \in I_m$ and 0 otherwise, and $\mathbf{M}_m = (\mathbf{X}^\top \mathbf{L}_m \mathbf{X} / k + \lambda \mathbf{I}_p)^{-1}$. Then, it holds that*

1. For all $m \in [M]$ and $I_m \in \mathcal{I}_k$,

$$1 - \frac{2}{k} \text{tr}(\mathbf{M}_m \widehat{\Sigma}_m) + \frac{1}{k} \text{tr}(\mathbf{M}_m \widehat{\Sigma}_m \mathbf{M}_m \widehat{\Sigma}_m) \xrightarrow{\text{a.s.}} \lambda^2 v(-\lambda; \phi_s)^2 (1 + \tilde{v}(-\lambda; \phi_s, \phi_s)). \quad (55)$$

2. For all $m, l \in [M]$, $m \neq l$ and $I_m, I_l \stackrel{\text{SRSWR}}{\sim} \mathcal{I}_k$,

$$\frac{1}{k} \text{tr}(\mathbf{M}_m \widehat{\Sigma}_m \mathbf{M}_m \widehat{\Sigma}_{l \setminus m}) \xrightarrow{\text{a.s.}} \tilde{v}(-\lambda; \phi_s, \phi_s). \quad (56)$$

3. For all $m, l \in [M]$, $m \neq l$ and $I_m, I_l \stackrel{\text{SRSWR}}{\sim} \mathcal{I}_k$,

$$\begin{aligned} 1 - \frac{1}{|I_m \cup I_l|} \sum_{j \in \{m, l\}} \text{tr}(\mathbf{M}_j \widehat{\Sigma}_j) + \frac{1}{n} \text{tr}(\mathbf{M}_l \widehat{\Sigma}_{m \cap l} \mathbf{M}_m \widehat{\Sigma}_{m \cup l}) &\xrightarrow{\text{a.s.}} \\ \mathcal{D}^\lambda(\phi_s, \phi_s) \left(\frac{2(\phi_s - \phi)}{2\phi_s - \phi} \frac{1}{\lambda v(-\lambda; \phi_s)} + \frac{\phi}{2\phi_s - \phi} \right) &(1 + \tilde{v}(-\lambda; \phi, \phi_s)), \end{aligned} \quad (57)$$

where $\widehat{\Sigma}_{l \setminus m} = |I_l \setminus I_m|^{-1} \mathbf{X}^\top \mathbf{L}_{l \setminus m} \mathbf{X}$ and $\widehat{\Sigma}_{m \cup l} = |I_m \cup I_l|^{-1} \mathbf{X}^\top \mathbf{L}_{m \cup l} \mathbf{X}$, as $n, k, p \rightarrow \infty$, $p/n \rightarrow \phi \in (0, \infty)$, and $p/k \rightarrow \phi_s \in [\phi, \infty)$, where the nonnegative constant $\tilde{v}(\lambda; \phi, \phi_s)$ is as defined in (19).

Proof of Lemma D.7. We split the proof into different parts.

Part (1) Note that

$$\text{tr}(\mathbf{M}_m \widehat{\Sigma}_m \mathbf{M}_m \widehat{\Sigma}_m) = \text{tr}(\mathbf{M}_m \widehat{\Sigma}_m) - \lambda \text{tr}(\mathbf{M}_m^2 \widehat{\Sigma}_m).$$

We now have

$$\begin{aligned} 1 - \frac{2}{k} \text{tr}(\mathbf{M}_m \widehat{\Sigma}_m) + \frac{1}{k} \text{tr}(\mathbf{M}_m \widehat{\Sigma}_m \mathbf{M}_m \widehat{\Sigma}_m) &= 1 - \frac{1}{k} \text{tr}(\mathbf{M}_m \widehat{\Sigma}_m) - \frac{\lambda}{k} \text{tr}(\mathbf{M}_m^2 \widehat{\Sigma}_m) \\ &= 1 - \frac{p}{k} + \frac{\lambda}{k} \text{tr}(\mathbf{M}_m) - \frac{\lambda}{k} \text{tr}(\mathbf{M}_m^2 \widehat{\Sigma}_m). \end{aligned}$$

From Corollary F.5 we have that $\lambda \mathbf{M}_m \simeq (v(-\lambda; \phi_s) \Sigma + \mathbf{I}_p)^{-1}$. From Lemma F.6 (2) (with $\mathbf{A} = \mathbf{I}$), we have that for $j \in [M]$,

$$\mathbf{M}_m \widehat{\Sigma}_m \mathbf{M}_m \simeq \tilde{v}_v(-\lambda; \phi_s) (v(-\lambda; \phi_s) \Sigma + \mathbf{I}_p)^{-2} \Sigma. \quad (58)$$

By the trace rule Lemma F.3 (4), we have

$$\begin{aligned}
 \frac{\lambda}{k} \operatorname{tr}(\mathbf{M}_m) - \frac{\lambda}{k} \operatorname{tr}(\mathbf{M}_m^2 \widehat{\Sigma}_m) &\xrightarrow{\text{a.s.}} \lim_{p \rightarrow \infty} \frac{p}{k} \cdot \frac{1}{p} \left(\operatorname{tr}((v(-\lambda; \phi_s) \Sigma + \mathbf{I}_p)^{-1}) - \operatorname{tr}(\lambda \tilde{v}_v(-\lambda; \phi_s)(v(-\lambda; \phi_s) \Sigma + \mathbf{I}_p)^{-2} \Sigma) \right) \\
 &= \phi_s \lim_{p \rightarrow \infty} \frac{1}{p} \sum_{i=1}^p \frac{1 + v(-\lambda; \phi_s) r_i - \lambda \tilde{v}_v(-\lambda; \phi_s) r_i}{(v(-\lambda; \phi_s) r_i + 1)^2} \\
 &= \phi_s \lim_{p \rightarrow \infty} \int \frac{1 + v(-\lambda; \phi_s) r - \lambda \tilde{v}_v(-\lambda; \phi_s) r}{(1 + v(-\lambda; \phi_s) r)^2} dH_p(r) \\
 &= \phi_s \int \frac{1 + v(-\lambda; \phi_s) r - \lambda \tilde{v}_v(-\lambda; \phi_s) r}{(1 + v(-\lambda; \phi_s) r)^2} dH(r), \quad j = 1, 2,
 \end{aligned} \tag{59}$$

where in the last line we used the fact that H_p and H have compact supports and Assumption 2.2. Then, we have

$$\begin{aligned}
 1 - \frac{2}{k} \operatorname{tr}(\mathbf{M}_m \widehat{\Sigma}_m) + \frac{1}{k} \operatorname{tr}(\mathbf{M}_m \widehat{\Sigma}_m \mathbf{M}_m \widehat{\Sigma}_m) \\
 &\xrightarrow{\text{a.s.}} 1 - \phi_s + \phi_s \int \frac{1 + v(-\lambda; \phi_s) r - \lambda \tilde{v}_v(-\lambda; \phi_s) r}{(1 + v(-\lambda; \phi_s) r)^2} dH(r) \\
 &= 1 - \phi_s \int \frac{v(-\lambda; \phi_s) r}{1 + v(-\lambda; \phi_s) r} dH(r) - \phi_s \int \frac{\lambda \tilde{v}_v(-\lambda; \phi_s) r}{(1 + v(-\lambda; \phi_s) r)^2} dH(r) \\
 &= \lambda v(-\lambda; \phi_s) - \phi_s \int \frac{\lambda \tilde{v}_v(-\lambda; \phi_s) r}{(1 + v(-\lambda; \phi_s) r)^2} dH(r) \\
 &= \lambda \tilde{v}_v(-\lambda; \phi_s) \left(v(-\lambda; \phi_s)^{-1} - \phi_s \int \frac{v(-\lambda; \phi_s) r^2}{(1 + v(-\lambda; \phi_s) r)^2} dH(r) - \phi_s \int \frac{r}{(1 + v(-\lambda; \phi_s) r)^2} dH(r) \right) \\
 &= \lambda \tilde{v}_v(-\lambda; \phi_s) \left(v(-\lambda; \phi_s)^{-1} - \phi_s \int \frac{r}{1 + v(-\lambda; \phi_s) r} dH(r) \right) \\
 &= \lambda^2 \tilde{v}_v(-\lambda; \phi_s) \\
 &= \lambda^2 v(-\lambda; \phi_s)^2 (1 + \tilde{v}(-\lambda; \phi_s, \phi_s)),
 \end{aligned}$$

and thus, (55) follows.

Part (2) Since $\mathbf{M}_m \widehat{\Sigma}_m \mathbf{M}_m$ and $\widehat{\Sigma}_{l \setminus m}$ are independent, from Lemma F.8 (1), we have

$$\mathbf{M}_m \widehat{\Sigma}_m \mathbf{M}_m \widehat{\Sigma}_{l \setminus m} \simeq \mathbf{M}_m \widehat{\Sigma}_m \mathbf{M}_m \Sigma.$$

Then, by the definition of deterministic equivalents, it follows that

$$\frac{1}{k} \operatorname{tr}(\mathbf{M}_m \widehat{\Sigma}_m \mathbf{M}_m \widehat{\Sigma}_{l \setminus m}) = \frac{1}{k} \operatorname{tr}(\mathbf{M}_m \widehat{\Sigma}_m \mathbf{M}_m \Sigma) \xrightarrow{\text{a.s.}} \tilde{v}(-\lambda; \phi_s, \phi_s),$$

where the convergence is due to Patil et al. (2022a, Lemma S.2.5).

Part (3) Let $i_0 = |I_m \cap I_l|$. The first two terms in (59) satisfy that

$$1 - \frac{1}{|I_m \cup I_l|} \sum_{j \in \{m, l\}} \operatorname{tr}(\mathbf{M}_j \widehat{\Sigma}_j) = 1 - \frac{p}{2k - i_0} \sum_{j \in \{m, l\}} \frac{1}{p} \operatorname{tr}(\mathbf{M}_j \widehat{\Sigma}_j) \xrightarrow{\text{a.s.}} 1 - \frac{2\phi_s^2}{2\phi_s - \phi} \int \frac{v(-\lambda; \phi_s) r}{1 + v(-\lambda; \phi_s) r} dH(r),$$

where the convergence is from Part 1. The last term can be further decomposed because

$$\begin{aligned}
 \operatorname{tr}(\mathbf{M}_l \widehat{\Sigma}_{m \cap l} \mathbf{M}_m \widehat{\Sigma}_{m \cup l}) &= \frac{k}{2k - i_0} \sum_{j \in \{m, l\}} \operatorname{tr}(\mathbf{M}_l \widehat{\Sigma}_{m \cap l} \mathbf{M}_m \widehat{\Sigma}_j) - \frac{i_0}{2k - i_0} \operatorname{tr}(\mathbf{M}_l \widehat{\Sigma}_{m \cap l} \mathbf{M}_m \widehat{\Sigma}_{m \cap l}) \\
 &= \frac{k}{2k - i_0} \sum_{j \in \{m, l\}} [\operatorname{tr}(\mathbf{M}_j \widehat{\Sigma}_{m \cap l}) - \lambda \operatorname{tr}(\mathbf{M}_m \widehat{\Sigma}_{m \cap l} \mathbf{M}_l)] - \frac{i_0}{2k - i_0} \operatorname{tr}(\mathbf{M}_l \widehat{\Sigma}_{m \cap l} \mathbf{M}_m \widehat{\Sigma}_{m \cap l}).
 \end{aligned}$$

From Lemma F.8 (4), (3), and (5), we have

$$\begin{aligned}
 \mathbf{M}_j \widehat{\boldsymbol{\Sigma}}_{m \cap l} &\simeq \mathbf{I}_p - (v(-\lambda; \phi_s) \boldsymbol{\Sigma} + \mathbf{I}_p)^{-1} \\
 \mathbf{M}_m \widehat{\boldsymbol{\Sigma}}_{m \cap l} \mathbf{M}_l &\simeq \tilde{v}_v(-\lambda; \phi, \phi_s) (v(-\lambda; \phi_s) \boldsymbol{\Sigma} + \mathbf{I}_p)^{-2} \boldsymbol{\Sigma} \\
 \mathbf{M}_l \widehat{\boldsymbol{\Sigma}}_{m \cap l} \mathbf{M}_m \widehat{\boldsymbol{\Sigma}}_{m \cap l} &\simeq \left(\frac{\phi_s}{\phi} v(-\lambda; \phi_s) - \frac{\phi_s - \phi}{\phi} \lambda \tilde{v}_v(-\lambda; \phi, \phi_s) \right) (v(-\lambda; \phi_s) \boldsymbol{\Sigma} + \mathbf{I}_p)^{-1} \boldsymbol{\Sigma} \\
 &\quad - \lambda \tilde{v}_v(-\lambda; \phi, \phi_s) (v(-\lambda; \phi_s) \boldsymbol{\Sigma} + \mathbf{I}_p)^{-2} \boldsymbol{\Sigma}.
 \end{aligned}$$

Combining the above terms and by Assumption 2.2, we have

$$\begin{aligned}
 &1 - \frac{1}{|I_m \cup I_l|} \sum_{j \in \{m, l\}} \text{tr}(\mathbf{M}_j \widehat{\boldsymbol{\Sigma}}_j) + \frac{1}{n} \text{tr}(\mathbf{M}_l \widehat{\boldsymbol{\Sigma}}_{m \cap l} \mathbf{M}_m \widehat{\boldsymbol{\Sigma}}_{m \cup l}) \\
 &\stackrel{\text{a.s.}}{\rightarrow} 1 - \frac{2\phi_s^2}{2\phi_s - \phi} \int \frac{v(-\lambda; \phi_s)r}{1 + v(-\lambda; \phi_s)r} dH(r) \\
 &\quad + \phi \frac{2\phi_s}{2\phi_s - \phi} \left(\int \frac{v(-\lambda; \phi_s)r}{1 + v(-\lambda; \phi_s)r} dH(r) - \lambda \tilde{v}_v(-\lambda; \phi, \phi_s) \int \frac{r}{(1 + v(-\lambda; \phi_s)r)^2} dH(r) \right) \\
 &\quad - \phi \frac{\phi}{2\phi_s - \phi} \left(\frac{\phi_s}{\phi} \int \frac{r}{1 + v(-\lambda; \phi_s)r} dH(r) - \frac{\phi_s - \phi}{\phi} \lambda \tilde{v}_v(-\lambda; \phi, \phi_s) \int \frac{r}{1 + v(-\lambda; \phi_s)r} dH(r) \right. \\
 &\quad \quad \left. - \lambda \tilde{v}_v(-\lambda; \phi, \phi_s) \int \frac{r}{(1 + v(-\lambda; \phi_s)r)^2} dH(r) \right) \\
 &= 1 - \phi_s \int \frac{v(-\lambda; \phi_s)r}{1 + v(-\lambda; \phi_s)r} dH(r) + \frac{\phi(\phi_s - \phi)}{2\phi_s - \phi} \lambda \tilde{v}_v(-\lambda; \phi, \phi_s) \int \frac{r}{1 + v(-\lambda; \phi_s)r} dH(r) \\
 &\quad - \phi \lambda \tilde{v}_v(-\lambda; \phi, \phi_s) \int \frac{r}{(1 + v(-\lambda; \phi_s)r)^2} dH(r) \\
 &= \lambda v(-\lambda; \phi_s) + \frac{\phi(\phi_s - \phi)}{2\phi_s - \phi} \int \frac{\lambda \tilde{v}_v(-\lambda; \phi, \phi_s)r}{1 + v(-\lambda; \phi_s)r} dH(r) - \phi \int \frac{\lambda \tilde{v}_v(-\lambda; \phi, \phi_s)r}{(1 + v(-\lambda; \phi_s)r)^2} dH(r) \\
 &= \frac{\lambda^2 \tilde{v}_v(-\lambda; \phi, \phi_s)}{2\phi_s - \phi} \left(\frac{2\phi_s - \phi}{\lambda v(-\lambda; \phi_s)} - \frac{\phi(2\phi_s - \phi)}{\lambda} \int \frac{v(-\lambda; \phi_s)r^2}{(1 + v(-\lambda; \phi_s)r)^2} dH(r) \right. \\
 &\quad \left. + \frac{\phi(\phi_s - \phi)}{\lambda} \int \frac{r}{1 + v(-\lambda; \phi_s)r} dH(r) - \frac{\phi}{\lambda} \int \frac{r}{(1 + v(-\lambda; \phi_s)r)^2} dH(r) \right) \\
 &= \frac{\lambda^2 \tilde{v}_v(-\lambda; \phi, \phi_s)}{2\phi_s - \phi} \left(\frac{2\phi_s - \phi}{\lambda v(-\lambda; \phi_s)} - \frac{\phi(2\phi_s - \phi)}{\lambda} \int \frac{r}{1 + v(-\lambda; \phi_s)r} dH(r) \right. \\
 &\quad \left. + \frac{\phi(\phi_s - \phi)}{\lambda} \int \frac{r}{1 + v(-\lambda; \phi_s)r} dH(r) \right) \\
 &= \lambda^2 \tilde{v}_v(-\lambda; \phi, \phi_s)^2 \left(\frac{2(\phi_s - \phi)}{2\phi_s - \phi} \frac{1}{\lambda v(-\lambda; \phi_s)} + \frac{\phi}{2\phi_s - \phi} \right) \\
 &= \mathcal{D}^\lambda(\phi_s, \phi_s) \left(\frac{2(\phi_s - \phi)}{2\phi_s - \phi} \frac{1}{\lambda v(-\lambda; \phi_s)} + \frac{\phi}{2\phi_s - \phi} \right) (1 + \tilde{v}(-\lambda; \phi, \phi_s)).
 \end{aligned}$$

□

E. Proof of Proposition 3.6

Proof of Proposition 3.6. From Lemma D.1, we have

$$\begin{aligned}
 \frac{1}{n} \|\mathbf{y} - \mathbf{X} \tilde{\boldsymbol{\beta}}_k^\lambda\|_2^2 &= \frac{1}{n} \mathbb{E}_{I \stackrel{\text{SRS}}{\sim} \mathcal{I}_k} \left[\text{Err}_{\text{train}}(\tilde{\boldsymbol{\beta}}_k^\lambda(\mathcal{D}_I)) + \text{Err}_{\text{test}}(\tilde{\boldsymbol{\beta}}_k^\lambda(\{\mathcal{D}_I\})) \right] \\
 &\quad + \frac{2}{n} \mathbb{E}_{(I_m, I_\ell) \stackrel{\text{SRS}}{\sim} \mathcal{I}_k} \left[\text{Err}_{\text{train}}(\tilde{\boldsymbol{\beta}}_k^\lambda(\{\mathcal{D}_{I_m}, \mathcal{D}_{I_\ell}\})) + \text{Err}_{\text{test}}(\tilde{\boldsymbol{\beta}}_k^\lambda(\{\mathcal{D}_{I_m}, \mathcal{D}_{I_\ell}\})) \right].
 \end{aligned}$$

Since by Lemma D.2, Lemma D.3 and Lemma G.6, each expectation converges, we have that

$$\frac{1}{n} \|\mathbf{y} - \mathbf{X} \tilde{\boldsymbol{\beta}}_k^\lambda\|_2^2 \xrightarrow{\text{a.s.}} \frac{2\phi(2\phi_s - \phi)}{\phi_s^2} \mathcal{T}_2^\lambda(\phi, \phi_s) + \frac{2(\phi_s - \phi)^2}{\phi_s^2} \mathcal{R}_2^\lambda(\phi, \phi_s) - \frac{\phi}{\phi_s} \mathcal{T}_1^\lambda(\phi, \phi_s) - \frac{\phi_s - \phi}{\phi_s} \mathcal{R}_1^\lambda(\phi, \phi_s),$$

where the convergence of the averages is from Lemma G.5 and the convergence of coefficients is from Lemma G.2. Since the denominator converges from Lemma 3.4, we further have

$$\text{gcv}_{k,\infty}^\lambda \xrightarrow{\text{a.s.}} \mathcal{G}_\infty^\lambda(\phi, \phi_s) = \frac{\frac{2\phi(2\phi_s - \phi)}{\phi_s^2} \mathcal{T}_2^\lambda(\phi, \phi_s) + \frac{2(\phi_s - \phi)^2}{\phi_s^2} \mathcal{R}_2^\lambda(\phi, \phi_s) - \frac{\phi}{\phi_s} \mathcal{T}_1^\lambda(\phi, \phi_s) - \frac{\phi_s - \phi}{\phi_s} \mathcal{R}_1^\lambda(\phi, \phi_s)}{\mathcal{D}_\infty^\lambda(\phi, \phi_s)},$$

for $\lambda > 0$ and $\phi_s \in [\phi, +\infty)$.

For the boundary case when $\lambda > 0$ but $\phi_s = +\infty$, we require Proposition E.1; for the boundary case when $\lambda = 0$, we require Proposition E.2. Applying Proposition E.1 and Proposition E.2 finishes the proof. \square

E.1. Boundary case: diverging subsample aspect ratio for the ridge predictor

Proposition E.1 (Risk approximation when $\phi_s \rightarrow +\infty$). *Under Assumptions 2.1-2.2, it holds for all $\lambda > 0$,*

$$\text{gcv}_k^\lambda \xrightarrow{\text{a.s.}} \mathcal{G}_\infty^\lambda(\phi, \infty),$$

as $k, n, p \rightarrow \infty$, $p/n \rightarrow \phi \in (0, \infty)$ and $p/k \rightarrow \infty$, where $\mathcal{G}_\infty^\lambda(\cdot, \cdot)$ is defined in Proposition 3.6.

Proof of Proposition E.1. Recall that

$$\frac{1}{n} \|\mathbf{y} - \mathbf{X} \tilde{\boldsymbol{\beta}}_k^\lambda\|_2^2 = \lim_{M \rightarrow \infty} \frac{1}{n} \|\mathbf{y} - \mathbf{X} \tilde{\boldsymbol{\beta}}_{k,M}^\lambda\|_2^2 = (\boldsymbol{\beta}_0 - \tilde{\boldsymbol{\beta}}_k^\lambda)^\top \widehat{\boldsymbol{\Sigma}} (\boldsymbol{\beta}_0 - \tilde{\boldsymbol{\beta}}_k^\lambda) + \frac{1}{n} \boldsymbol{\epsilon}^\top \boldsymbol{\epsilon} + \frac{2}{n} (\boldsymbol{\beta}_0 - \tilde{\boldsymbol{\beta}}_k^\lambda)^\top \widehat{\boldsymbol{\Sigma}} \boldsymbol{\epsilon}$$

From Lemma G.3 and Lemma G.4, we have that $\boldsymbol{\beta}_0^\top \mathbf{X}^\top \boldsymbol{\epsilon}/n \xrightarrow{\text{a.s.}} 0$ and $\boldsymbol{\epsilon}^\top \boldsymbol{\epsilon}/n \xrightarrow{\text{a.s.}} \sigma^2$ as $n \rightarrow \infty$. For the other term, note that for any $(I_1, \dots, I_M) \stackrel{\text{SRs}}{\sim} \mathcal{I}_k$,

$$\begin{aligned} \|\tilde{\boldsymbol{\beta}}_k^\lambda\|_2 &\leq \lim_{M \rightarrow \infty} \mathbb{E}_{(I_1, \dots, I_M) \stackrel{\text{SRs}}{\sim} \mathcal{I}_k} \left[\frac{1}{M} \sum_{m=1}^M \|(\mathbf{X}^\top \mathbf{L}_m \mathbf{X}/k + \lambda \mathbf{I}_p)^{-1} (\mathbf{X}^\top \mathbf{L}_m \mathbf{y}/k)\|_2 \right] \\ &\leq \lim_{M \rightarrow \infty} \mathbb{E}_{(I_1, \dots, I_M) \stackrel{\text{SRs}}{\sim} \mathcal{I}_k} \left[\frac{1}{M} \sum_{m=1}^M \|(\mathbf{X}^\top \mathbf{L}_m \mathbf{X}/k + \lambda \mathbf{I}_p)^{-1} \mathbf{X}^\top \mathbf{L}_m / \sqrt{k}\| \cdot \|\mathbf{L}_m \mathbf{y} / \sqrt{k}\|_2 \right] \\ &\leq C \sqrt{\rho^2 + \sigma^2} \cdot \max_{I_m \in \mathcal{I}_k} \|(\mathbf{X}^\top \mathbf{L}_m \mathbf{X}/k + \lambda \mathbf{I}_p)^{-1} \mathbf{X}^\top \mathbf{L}_m / \sqrt{k}\|_{\text{op}}, \end{aligned}$$

where the last inequality holds eventually almost surely since Assumptions 2.1-2.2 imply that the entries of \mathbf{y} have bounded 4-th moment, and thus from the strong law of large numbers, $\|\mathbf{L}_m \mathbf{y} / \sqrt{k}\|_2$ is eventually almost surely bounded above by $C \sqrt{\mathbb{E}[y_1^2]} = C \sqrt{\rho^2 + \sigma^2}$ for some constant C . Observe that operator norm of the matrix $(\mathbf{X}^\top \mathbf{L}_m \mathbf{X}/k + \lambda \mathbf{I}_p)^{-1} \mathbf{X}^\top \mathbf{L}_m / \sqrt{k}$ is upper bounded $\max_i s_i / (s_i^2 + \lambda) \leq 1/s_{\min}$ where s_i 's are the singular values of \mathbf{X} and s_{\min} is the smallest nonzero singular value. As $k, p \rightarrow \infty$ such that $p/k \rightarrow \infty$, $s_{\min} \rightarrow \infty$ almost surely (e.g., from results of Bloemendal et al. (2016)) and therefore, $\|\tilde{\boldsymbol{\beta}}_k^\lambda\|_2 \rightarrow 0$ almost surely. Because $\|\widehat{\boldsymbol{\Sigma}}\|_{\text{op}}$ is upper bounded almost surely, we further have $\tilde{\boldsymbol{\beta}}_k^\lambda \mathbf{X}^\top \boldsymbol{\epsilon}/n \xrightarrow{\text{a.s.}} 0$. Consequently we have $(\boldsymbol{\beta}_0 - \tilde{\boldsymbol{\beta}}_k^\lambda)^\top \mathbf{X}^\top \boldsymbol{\epsilon}/n \xrightarrow{\text{a.s.}} 0$ and

$$\frac{1}{n} \|\mathbf{y} - \mathbf{X} \tilde{\boldsymbol{\beta}}_k^\lambda\|_2^2 \xrightarrow{\text{a.s.}} \boldsymbol{\beta}_0^\top \widehat{\boldsymbol{\Sigma}} \boldsymbol{\beta}_0 + \sigma^2.$$

Finally, from Lemma F.8 (1) $\boldsymbol{\beta}_0^\top \widehat{\boldsymbol{\Sigma}} \boldsymbol{\beta}_0 \xrightarrow{\text{a.s.}} \boldsymbol{\beta}_0^\top \boldsymbol{\Sigma} \boldsymbol{\beta}_0$ and from Assumption 2.2, we have

$$\frac{1}{n} \|\mathbf{y} - \mathbf{X} \tilde{\boldsymbol{\beta}}_k^\lambda\|_2^2 \xrightarrow{\text{a.s.}} \sigma^2 + \rho^2 \int r \, dG(r).$$

Since $\mathbf{S}_k^\lambda = \mathbf{X} \tilde{\boldsymbol{\beta}}_k^\lambda$, we have that $\text{tr}(\mathbf{S}_k^\lambda)/n \xrightarrow{\text{a.s.}} 0$. So the denominator converges to 1 almost surely.

From Lemma F.10, we have $\mathcal{G}_\infty^\lambda(\phi, \infty) := \lim_{\phi_s \rightarrow +\infty} \mathcal{G}_\infty^\lambda(\phi, \phi_s) = \sigma^2 + \rho^2 \int r \, dG(r)$, which is also the limit of the GCV estimate. Thus, $\mathcal{G}_\infty^\lambda(\phi, \infty)$ is well defined and $\mathcal{G}_\infty^\lambda(\phi, \phi_s)$ is right continuous at $\phi_s = +\infty$. \square

E.2. Boundary case: the ridgeless predictor

Proposition E.2 (Risk approximation when $\lambda = 0$). *Under Assumptions 2.1-2.2, suppose that the conclusion of Proposition 3.6 holds for $\lambda > 0$. Then it holds that,*

$$\text{gcv}_k^0 \xrightarrow{\text{a.s.}} \mathcal{G}_\infty^0(\phi, \phi_s) := \lim_{\lambda \rightarrow 0^+} \mathcal{G}_\infty^\lambda(\phi, \phi_s),$$

as $k, n, p \rightarrow \infty$, $p/n \rightarrow \phi \in (0, \infty)$ and $p/k \rightarrow [\phi, +\infty]$, where $\mathcal{G}_\infty^\lambda(\cdot, \cdot)$ is defined in Proposition 3.6.

Proof of Proposition E.2. We analyze the numerator and the denominator separately.

Part (1) For the denominator, note that

$$P_{n,\lambda} := (1 - \text{tr}(\mathbf{S}_k^\lambda)/n)^2 = \lim_{M \rightarrow \infty} (1 - \text{tr}(\mathbf{S}_{k,M}^\lambda)/n)^2,$$

where $\mathbf{S}_k^\lambda = \lim_{M \rightarrow \infty} \mathbf{S}_{k,M}^\lambda$ is the smoothing matrix. Since $\mathbf{S}_k^\lambda \succeq \mathbf{0}_{n \times n}$ and

$$\|\mathbf{S}_k^\lambda\|_{\text{op}} \leq \max_{I_m \in \mathcal{L}_k} \|\mathbf{X}(\mathbf{X}^\top \mathbf{L}_m \mathbf{X}/k + \lambda \mathbf{I}_p)^{-1} \mathbf{X}^\top \mathbf{L}_m / \sqrt{k}\|_{\text{op}}, \quad (60)$$

which is also upper bounded almost surely from the proof in Proposition E.1 (when $\lambda = 0$, the inverse in the above display is replaced by pseudo-inverse). Thus, we have $P_{n,\lambda}$ is almost surely upper bounded $\lambda \in \Lambda := [0, \lambda_{\max}]$ for any $\lambda_{\max} \in (0, \infty)$ fixed.

Next we inspect the boundedness of the derivative of $P_{n,\lambda}$:

$$\frac{\partial}{\partial \lambda} P_{n,\lambda} = \frac{\partial}{\partial \lambda} \lim_{M \rightarrow \infty} \lim_{M \rightarrow \infty} (1 - \text{tr}(\mathbf{S}_{k,M}^\lambda)/n)^2 =: \frac{\partial}{\partial \lambda} \lim_{M \rightarrow \infty} Q_{M,\lambda}.$$

We claim that

$$\frac{\partial}{\partial \lambda} \lim_{M \rightarrow \infty} Q_{M,\lambda} = \lim_{M \rightarrow \infty} \frac{\partial}{\partial \lambda} Q_{M,\lambda}.$$

To see this, we need to show that $Q_{M,\lambda}$ is equicontinuous in λ over Λ . First we know that $Q_{M,\lambda}$ is differentiable in λ . From (60), we have that $Q_{M,\lambda}$ is uniformly upper bounded over $\lambda \in \Lambda$ almost surely. Note that

$$\frac{\partial}{\partial \lambda} Q_{M,\lambda} = 2(1 - \text{tr}(\mathbf{S}_{k,M}^\lambda)) \text{tr} \left(\frac{\partial}{\partial \lambda} \mathbf{S}_{k,M}^\lambda \right),$$

where

$$\frac{\partial}{\partial \lambda} \mathbf{S}_{k,M}^\lambda = \frac{1}{M} \sum_{m=1}^M \mathbf{X} \left(\frac{\mathbf{X}^\top \mathbf{L}_m \mathbf{X}}{k} + \lambda \mathbf{I} \right)^{-2} \frac{\mathbf{X}^\top \mathbf{L}_m}{k}.$$

By the similar arguments as in Proposition E.1, we have that $\left\| \frac{\partial \mathbf{S}_{k,M}^\lambda}{\partial \lambda} \right\|_{\text{op}}$, and $\|\mathbf{S}_{k,M}^\lambda\|_2^2$ are uniformly upper bounded almost surely over Λ , the equicontinuity conclusion follows. Then by Moore-Osgood theorem, we have

$$\frac{\partial}{\partial \lambda} P_{n,\lambda} = \lim_{M \rightarrow \infty} 2(1 - \text{tr}(\mathbf{S}_{k,M}^\lambda)) \text{tr} \left(\frac{\partial}{\partial \lambda} \mathbf{S}_{k,M}^\lambda \right)$$

is uniformly upper bounded almost surely over $[0, +\infty]$ independent of λ and M . Therefore, we conclude that $|\partial P_{n,\lambda} / \partial \lambda|$ is upper bounded almost surely over $\lambda \in \Lambda$.

On the other hand, we know that $P_{n,\lambda} \xrightarrow{\text{a.s.}} \mathcal{D}_\infty^\lambda(\phi, \phi_s)$ for $\lambda > 0$. Define $\mathcal{D}^0(\phi, \phi_s) := \lim_{\lambda \rightarrow 0^+} \mathcal{D}_\infty^\lambda(\phi, \phi_s)$. When $\lambda = 0$ and $\phi_s > 1$, we know that $\mathcal{D}^0(\phi, \phi_s)$ is well-defined because $v(-\lambda; \phi_s)$ is finite and continuous from Lemma F.12. When $\lambda = 0$ and $\phi_s \in (0, 1]$, from the definition of fixed-point solution (18), we have

$$1 = v(-\lambda; \phi_s) \lambda + \phi_s \int \frac{v(-\lambda; \phi_s) r}{1 + v(-\lambda; \phi_s) r} dH(r).$$

In this case, $v(0; \phi_s) = +\infty$ from Lemma F.12. Let $\lambda \rightarrow 0^+$, we have

$$1 = \lim_{\lambda \rightarrow 0^+} v(-\lambda; \phi_s)\lambda + \phi_s \lim_{\lambda \rightarrow 0^+} \int \frac{v(-\lambda; \phi_s)r}{1 + v(-\lambda; \phi_s)r} dH(r) = \lim_{\lambda \rightarrow 0^+} v(-\lambda; \phi_s)\lambda + \phi_s.$$

Then we have $\lim_{\lambda \rightarrow 0^+} v(-\lambda; \phi_s)\lambda = 1 - \phi_s$ and $\mathcal{D}_\infty^\lambda(\phi, \phi_s) = (1 - \phi_s)^2$. Thus, $\mathcal{D}^0(\phi, \phi_s)$ is always well-defined.

From Lemma F.12, there exists $M' > 0$ such that the magnitudes of $v(-\lambda; \phi_s)$ and its derivative with respect to λ are continuous and bounded by M' for all $\lambda \in [0, +\infty]$. It follows that $|\mathcal{D}_\infty^\lambda(\phi, \phi_s)|$ and $|\partial \mathcal{D}_\infty^\lambda(\phi, \phi_s)/\partial \lambda|$ are uniformly upper bounded almost surely. From Moore-Osgood theorem and the continuity property from Lemma F.12, we have

$$\lim_{n \rightarrow \infty} \lim_{\lambda \rightarrow 0^+} P_{n,\lambda} = \lim_{\lambda \rightarrow 0^+} \lim_{n \rightarrow \infty} P_{n,\lambda} = \lim_{\lambda \rightarrow 0^+} \mathcal{D}_\infty^\lambda(\phi, \phi_s) = \mathcal{D}^0(\phi, \phi_s).$$

Part (2) For the numerator, note that

$$P'_{n,\lambda} := \frac{1}{n} \|\mathbf{y} - \mathbf{X}\tilde{\boldsymbol{\beta}}_k^\lambda\|_2^2 = \frac{1}{n} \|(\mathbf{I}_n - \mathbf{S}_k^\lambda)\mathbf{y}\|_2^2.$$

Assumptions 2.1-2.2 imply that the entries of \mathbf{y} have bounded 4-th moment, and thus from the strong law of large numbers, $\|\mathbf{y}/\sqrt{n}\|_2$ is eventually almost surely bounded above by $C\sqrt{\mathbb{E}[y_1^2]} = C\sqrt{\rho^2 + \sigma^2}$ for some constant C . On the other hand, $\mathbf{S}_k^\lambda \succeq \mathbf{0}_{n \times n}$ and $\|\mathbf{S}_k^\lambda\|_{\text{op}}$ is also upper bounded almost surely from Part (1). Thus, we have $P'_{n,\lambda}$ is almost surely upper bounded $\lambda \in \Lambda$.

Next we inspect the boundedness of the derivative of $P'_{n,\lambda}$:

$$\begin{aligned} \frac{\partial}{\partial \lambda} P'_{n,\lambda} &= \frac{2}{n} (\mathbf{y} - \mathbf{X}\tilde{\boldsymbol{\beta}}_k^\lambda)^\top \frac{\partial}{\partial \lambda} \mathbf{S}_k^\lambda \mathbf{y} \\ &= \frac{\partial}{\partial \lambda} \lim_{M \rightarrow \infty} \frac{2}{n} (\mathbf{y} - \mathbf{X}\tilde{\boldsymbol{\beta}}_{k,M}^\lambda)^\top \mathbf{S}_{k,M}^\lambda \mathbf{y} \\ &= \frac{\partial}{\partial \lambda} \lim_{M \rightarrow \infty} \frac{2}{n} \mathbf{y}^\top (\mathbf{I}_n - \mathbf{S}_{k,M}^\lambda) \mathbf{S}_{k,M}^\lambda \mathbf{y} =: \frac{\partial}{\partial \lambda} \lim_{M \rightarrow \infty} Q'_{M,\lambda}. \end{aligned}$$

We claim that

$$\frac{\partial}{\partial \lambda} \lim_{M \rightarrow \infty} Q'_{M,\lambda} = \lim_{M \rightarrow \infty} \frac{\partial}{\partial \lambda} Q'_{M,\lambda}.$$

To see this, we need to show that $Q'_{M,\lambda}$ is equicontinuous in λ over Λ . First we know that $Q'_{M,\lambda}$ is differentiable in λ . From (60), we have that $Q'_{M,\lambda}$ is uniformly upper bounded over $\lambda \in \Lambda$ almost surely. Similarly, we have

$$\frac{\partial}{\partial \lambda} Q'_{M,\lambda} = \frac{2}{n} \mathbf{y}^\top (\mathbf{I}_n - 2\mathbf{S}_{k,M}^\lambda) \frac{\partial}{\partial \lambda} \mathbf{S}_{k,M}^\lambda \mathbf{y},$$

and

$$\left| \frac{\partial}{\partial \lambda} Q'_{M,\lambda} \right| \leq \|\mathbf{I}_n - 2\mathbf{S}_{k,M}^\lambda\|_{\text{op}} \left\| \frac{\partial}{\partial \lambda} \mathbf{S}_{k,M}^\lambda \right\|_{\text{op}} \frac{1}{n} \|\mathbf{y}\|_2^2.$$

and the equicontinuity conclusion follows analogously as in Part (1). Therefore, we conclude that $|\partial P'_{n,\lambda}/\partial \lambda|$ is upper bounded almost surely over $\lambda \in [0, +\infty]$.

On the other hand, we know that $P'_{n,\lambda} \xrightarrow{\text{a.s.}} \mathcal{G}_\infty^\lambda(\phi, \phi_s)$ for $\lambda > 0$. Define $\mathcal{D}^0(\phi, \phi_s) \mathcal{G}_\infty^0(\phi, \phi_s) := \lim_{\lambda \rightarrow 0^+} (\mathcal{D}_\infty^\lambda(\phi, \phi_s) \mathcal{G}_\infty^\lambda(\phi, \phi_s)) = \mathcal{D}^0(\phi, \phi_s) \mathcal{R}_M^0(\phi, \phi_s)$, which is well defined from Part (1) and Lemma A.2. From Lemma F.12, there exists $M' > 0$ such that the magnitudes of $v(-\lambda; \phi_s)$, $\tilde{v}(\lambda; \phi_s, \phi)$ and $\tilde{c}(\lambda; \phi_s)$, and their derivatives with respect to λ are continuous and bounded by M' for all $\lambda \in [0, +\infty]$. It follows that $|\mathcal{G}_\infty^\lambda(\phi, \phi_s)|$ is upper bounded almost surely. Analogously, we have that $|\partial(\mathcal{D}_\infty^\lambda(\phi, \phi_s) \mathcal{G}_\infty^\lambda(\phi, \phi_s))/\partial \lambda|$ is also upper bounded almost surely on $\lambda \in \Lambda$. From Moore-Osgood theorem and the continuity property from Lemma F.12, we have

$$\lim_{n \rightarrow \infty} \lim_{\lambda \rightarrow 0^+} P'_{n,\lambda} = \lim_{\lambda \rightarrow 0^+} \lim_{n \rightarrow \infty} P'_{n,\lambda} = \lim_{\lambda \rightarrow 0^+} \mathcal{G}_\infty^\lambda(\phi, \phi_s) = \mathcal{G}_\infty^0(\phi, \phi_s).$$

□

F. Auxiliary results on asymptotic equivalents

F.1. Preliminary background

We use the notion of asymptotic equivalence of sequences of random matrices in various proofs. In this section, we provide a basic review of the related definitions and corresponding calculus rules. See [Dobriban & Wager \(2018\)](#); [Dobriban & Sheng \(2021\)](#); [Patil et al. \(2022b;a\)](#) for more details.

Definition F.1 (Asymptotic equivalence). Consider sequences $\{\mathbf{A}_p\}_{p \geq 1}$ and $\{\mathbf{B}_p\}_{p \geq 1}$ of (random or deterministic) matrices of growing dimensions. We say that \mathbf{A}_p and \mathbf{B}_p are equivalent and write $\mathbf{A}_p \simeq \mathbf{B}_p$ if $\lim_{p \rightarrow \infty} |\text{tr}[\mathbf{C}_p(\mathbf{A}_p - \mathbf{B}_p)]| = 0$ almost surely for any sequence of random matrices \mathbf{C}_p independent to \mathbf{A}_p and \mathbf{B}_p , with bounded trace norm such that $\limsup_{p \rightarrow \infty} \|\mathbf{C}_p\|_{\text{tr}} < \infty$ almost surely.

The notion of asymptotic equivalence of two sequences of random matrices from Definition F.1 can be further extended to incorporate conditioning on another sequence of random matrices.

Definition F.2 (Conditional asymptotic equivalence). Consider sequences $\{\mathbf{A}_p\}_{p \geq 1}$, $\{\mathbf{B}_p\}_{p \geq 1}$ and $\{\mathbf{D}_p\}_{p \geq 1}$ of (random or deterministic) matrices of growing dimensions. We say that \mathbf{A}_p and \mathbf{B}_p are equivalent given \mathbf{D}_p and write $\mathbf{A}_p \simeq \mathbf{B}_p \mid \mathbf{D}_p$ if $\lim_{p \rightarrow \infty} |\text{tr}[\mathbf{C}_p(\mathbf{A}_p - \mathbf{B}_p)]| = 0$ almost surely conditional on $\{\mathbf{D}_p\}_{p \geq 1}$, i.e.,

$$\mathbb{P} \left(\lim_{p \rightarrow \infty} |\text{tr}[\mathbf{C}_p(\mathbf{A}_p - \mathbf{B}_p)]| = 0 \mid \{\mathbf{D}_p\}_{p \geq 1} \right) = 1,$$

for any sequence of random matrices \mathbf{C}_p independent to \mathbf{A}_p and \mathbf{B}_p conditional on \mathbf{D}_p , with bounded trace norm such that $\limsup \|\mathbf{C}_p\|_{\text{tr}} < \infty$ as $p \rightarrow \infty$.

Below we summarize the calculus rules for conditional asymptotic equivalence Definition F.2 adapted from [Patil et al. \(2022a\)](#), Lemma S.7.4 and S.7.6).

Lemma F.3 (Calculus of deterministic equivalents). *Let \mathbf{A}_p , \mathbf{B}_p , \mathbf{C}_p and \mathbf{D}_p be sequences of random matrices. The calculus of deterministic equivalents (\simeq_D and \simeq_R) satisfies the following properties:*

- (1) *Equivalence: The relation \simeq is an equivalence relation.*
- (2) *Sum: If $\mathbf{A}_p \simeq \mathbf{B}_p \mid \mathbf{E}_p$ and $\mathbf{C}_p \simeq \mathbf{D}_p \mid \mathbf{E}_p$, then $\mathbf{A}_p + \mathbf{C}_p \simeq \mathbf{B}_p + \mathbf{D}_p \mid \mathbf{E}_p$.*
- (3) *Product: If \mathbf{A}_p has bounded operator norms such that $\limsup_{p \rightarrow \infty} \|\mathbf{A}_p\|_{\text{op}} < \infty$, \mathbf{A}_p is conditional independent to \mathbf{B}_p and \mathbf{C}_p given \mathbf{E}_p for $p \geq 1$, and $\mathbf{B}_p \simeq \mathbf{C}_p \mid \mathbf{E}_p$, then $\mathbf{A}_p \mathbf{B}_p \simeq \mathbf{A}_p \mathbf{C}_p \mid \mathbf{E}_p$.*
- (4) *Trace: If $\mathbf{A}_p \simeq \mathbf{B}_p \mid \mathbf{E}_p$, then $\text{tr}[\mathbf{A}_p]/p - \text{tr}[\mathbf{B}_p]/p \rightarrow 0$ almost surely when conditioning on \mathbf{E}_p .*
- (5) *Differentiation: Suppose $f(z, \mathbf{A}_p) \simeq g(z, \mathbf{B}_p) \mid \mathbf{E}_p$ where the entries of f and g are analytic functions in $z \in S$ and S is an open connected subset of \mathbb{C} . Suppose for any sequence \mathbf{C}_p of deterministic matrices with bounded trace norm we have $|\text{tr}[\mathbf{C}_p(f(z, \mathbf{A}_p) - g(z, \mathbf{B}_p))]| \leq M$ for every p and $z \in S$. Then we have $f'(z, \mathbf{A}_p) \simeq g'(z, \mathbf{B}_p) \mid \mathbf{E}_p$ for every $z \in S$, where the derivatives are taken entrywise with respect to z .*
- (6) *Unconditioning: If $\mathbf{A}_p \simeq \mathbf{B}_p \mid \mathbf{E}_p$, then $\mathbf{A}_p \simeq \mathbf{B}_p$.*
- (7) *Substitution: Let $v : \mathbb{R}^{p \times p} \rightarrow \mathbb{R}$ and $f(v(\mathbf{C}), \mathbf{C}) : \mathbb{R}^{p \times p} \rightarrow \mathbb{R}^{p \times p}$ be a matrix function for matrix $\mathbf{C} \in \mathbb{R}^{p \times p}$ and $p \in \mathbb{N}$, that is continuous in the first argument with respect to operator norm. If $v(\mathbf{C}) \stackrel{\text{a.s.}}{=} v(\mathbf{D})$ such that \mathbf{C} is independent to \mathbf{D} , then $f(v(\mathbf{C}), \mathbf{C}) \simeq f(v(\mathbf{D}), \mathbf{C}) \mid \mathbf{C}$.*

F.2. Standard ridge resolvents and various extensions

In this section, we collect various asymptotic equivalents. Appendix F.2.1 introduces the basic concepts and definitions. The extended equivalents developed in the work of [Patil et al. \(2022a\)](#) are summarized in Appendix F.2.2. Based on results in Appendices F.2.1 and F.2.2, we prove some useful deterministic equivalent relations in Appendix F.2.3, which are subsequently used in the proof of Lemma D.3 (Lemmas D.6 and D.7).

F.2.1. STANDARD RIDGE RESOLVENTS

The following lemma provides a deterministic equivalent for the standard ridge resolvent and implies Corollary F.5. It is adapted from Theorem 1 of Rubio & Mestre (2011). See also Theorem 3 of Dobriban & Sheng (2021).

Lemma F.4 (Deterministic equivalent for standard ridge resolvent). *Suppose $\mathbf{x}_i \in \mathbb{R}^p$, $1 \leq i \leq n$, are i.i.d. random vectors such that each $\mathbf{x}_i = \mathbf{z}_i \boldsymbol{\Sigma}^{1/2}$, where \mathbf{z}_i is a random vector consisting of i.i.d. entries z_{ij} , $1 \leq j \leq p$, satisfying $\mathbb{E}[z_{ij}] = 0$, $\mathbb{E}[z_{ij}^2] = 1$, and $\mathbb{E}[|z_{ij}|^{8+\alpha}] \leq M_\alpha$ for some constants $\alpha > 0$ and $M_\alpha < \infty$, and $\boldsymbol{\Sigma} \in \mathbb{R}^{p \times p}$ is a positive semidefinite matrix satisfying $0 \preceq \boldsymbol{\Sigma} \preceq r_{\max} \mathbf{I}_p$ for some constant $r_{\max} < \infty$ (independent of p). Let $\mathbf{X} \in \mathbb{R}^{n \times p}$ the concatenated matrix with \mathbf{x}_i^\top , $1 \leq i \leq n$, as rows, and let $\widehat{\boldsymbol{\Sigma}} \in \mathbb{R}^{p \times p}$ denote the random matrix $\mathbf{X}^\top \mathbf{X}/n$. Let $\gamma := p/n$. Then, for $z \in \mathbb{C}^+$, as $n, p \rightarrow \infty$ such that $0 < \liminf \gamma \leq \limsup \gamma < \infty$, we have*

$$(\widehat{\boldsymbol{\Sigma}} - z \mathbf{I}_p)^{-1} \simeq (c(e(z; \gamma)) \boldsymbol{\Sigma} - z \mathbf{I}_p)^{-1}, \quad (61)$$

where the scalar $c(e(z; \gamma))$ is defined in terms of another scalar $e(z; \gamma)$ by the equation

$$c(e(z; \gamma)) = \frac{1}{1 + \gamma e(z; \gamma)}, \quad (62)$$

and $e(z; \gamma)$ is the unique solution in \mathbb{C}^+ to the fixed-point equation

$$e(z; \gamma) = \text{tr}[\boldsymbol{\Sigma}(c(e(z; \gamma)) \boldsymbol{\Sigma} - z \mathbf{I}_p)^{-1}]/p. \quad (63)$$

Note that both the scalars $c(e(z; \gamma))$ and $e(z; \gamma)$ also implicitly depend on $\boldsymbol{\Sigma}$. For notation brevity, we do not always explicitly indicate this dependence. However, we will be explicit in such dependence for certain extensions to follow. Additionally, observe that one can eliminate $e(z; \gamma)$ in the statement of Lemma F.4 by combining (62) and (63) so that for $z \in \mathbb{C}^+$, one has

$$(\widehat{\boldsymbol{\Sigma}} - z \mathbf{I}_p)^{-1} \simeq (c(z; \gamma) \boldsymbol{\Sigma} - z \mathbf{I}_p)^{-1},$$

where $c(z)$ is the unique solution in \mathbb{C}^- to the fixed-point equation

$$\frac{1}{c(z; \gamma)} = 1 + \gamma \text{tr}[\boldsymbol{\Sigma}(c(z; \gamma) \boldsymbol{\Sigma} - z \mathbf{I}_p)^{-1}]/p.$$

The following corollary is a simple consequence of Lemma F.4, which supplies a deterministic equivalent for the (regularization) scaled ridge resolvent. We will work with a real regularization parameter λ from here on.

Corollary F.5 (Deterministic equivalent for scaled ridge resolvent). *Assume the setting of Lemma F.4. For $\lambda > 0$, we have*

$$\lambda(\widehat{\boldsymbol{\Sigma}} + \lambda \mathbf{I}_p)^{-1} \simeq (v(-\lambda; \gamma) \boldsymbol{\Sigma} + \mathbf{I}_p)^{-1},$$

where $v(-\lambda; \gamma) > 0$ is the unique solution to the fixed-point equation

$$\frac{1}{v(-\lambda; \gamma)} = \lambda + \gamma \int \frac{r}{1 + v(-\lambda; \gamma)r} dH_n(r). \quad (64)$$

Here H_n is the empirical distribution (supported on $\mathbb{R}_{\geq 0}$) of the eigenvalues of $\boldsymbol{\Sigma}$.

As a side note, the parameter $v(-\lambda; \gamma)$ in Corollary F.5 is also the companion Stieltjes transform of the spectral distribution of the sample covariance matrix $\widehat{\boldsymbol{\Sigma}}$, which is also the Stieltjes transform of the spectral distribution of the gram matrix $\mathbf{X} \mathbf{X}^\top/n$.

The following lemma uses Corollary F.5 along with calculus of deterministic equivalents (from Lemma F.3), and provides deterministic equivalents for resolvents needed to obtain asymptotic bias and variance of standard ridge regression. It is adapted from Lemma S.6.10 of Patil et al. (2022b).

Lemma F.6 (Deterministic equivalents for ridge resolvents associated with generalized bias and variance). *Suppose $\mathbf{x}_i \in \mathbb{R}^p$, $1 \leq i \leq n$, are i.i.d. random vectors with each $\mathbf{x}_i = \mathbf{z}_i \boldsymbol{\Sigma}^{1/2}$, where $\mathbf{z}_i \in \mathbb{R}^p$ is a random vector that contains i.i.d. random variables z_{ij} , $1 \leq j \leq p$, each with $\mathbb{E}[z_{ij}] = 0$, $\mathbb{E}[z_{ij}^2] = 1$, and $\mathbb{E}[|z_{ij}|^{8+\alpha}] \leq M_\alpha$ for some constants $\alpha > 0$ and $M_\alpha < \infty$, and $\boldsymbol{\Sigma} \in \mathbb{R}^{p \times p}$ is a positive semidefinite matrix with $r_{\min} \mathbf{I}_p \preceq \boldsymbol{\Sigma} \preceq r_{\max} \mathbf{I}_p$ for some constants $r_{\min} > 0$ and $r_{\max} < \infty$ (independent of p). Let $\mathbf{X} \in \mathbb{R}^{n \times p}$ be the concatenated random matrix with \mathbf{x}_i , $1 \leq i \leq n$, as its rows, and define $\widehat{\boldsymbol{\Sigma}} := \mathbf{X}^\top \mathbf{X}/n \in \mathbb{R}^{p \times p}$. Let $\gamma := p/n$. Then, for $\lambda > 0$, as $n, p \rightarrow \infty$ with $0 < \liminf \gamma \leq \limsup \gamma < \infty$, the following statements hold:*

(1) *Bias of ridge regression:*

$$\lambda^2(\widehat{\Sigma} + \lambda\mathbf{I}_p)^{-1}\mathbf{A}(\widehat{\Sigma} + \lambda\mathbf{I}_p)^{-1} \simeq (v(-\lambda; \gamma, \Sigma)\Sigma + \mathbf{I}_p)^{-1}(\tilde{v}_b(-\lambda; \gamma, \Sigma, \mathbf{A})\Sigma + \mathbf{A})(v(-\lambda; \gamma, \Sigma)\Sigma + \mathbf{I}_p)^{-1}. \quad (65)$$

(2) *Variance of ridge regression:*

$$(\widehat{\Sigma} + \lambda\mathbf{I}_p)^{-2}\widehat{\Sigma}\mathbf{A} \simeq \tilde{v}_v(-\lambda; \gamma, \Sigma)(v(-\lambda; \gamma, \Sigma)\Sigma + \mathbf{I}_p)^{-2}\Sigma\mathbf{A}. \quad (66)$$

Here $v(-\lambda; \gamma, \Sigma) > 0$ is the unique solution to the fixed-point equation

$$\frac{1}{v(-\lambda; \gamma, \Sigma)} = \lambda + \int \frac{\gamma r}{1 + v(-\lambda; \gamma, \Sigma)r} dH_n(r; \Sigma), \quad (67)$$

and $\tilde{v}_b(-\lambda; \gamma, \Sigma)$ and $\tilde{v}_v(-\lambda; \gamma, \Sigma)$ are defined through $v(-\lambda; \gamma, \Sigma)$ by the following equations:

$$\tilde{v}_b(-\lambda; \gamma, \Sigma, \mathbf{A}) = \frac{\gamma \operatorname{tr}[\mathbf{A}\Sigma(v(-\lambda; \gamma, \Sigma)\Sigma + \mathbf{I}_p)^{-2}]/p}{v(-\lambda; \gamma, \Sigma)^{-2} - \int \gamma r^2(1 + v(-\lambda; \gamma, \Sigma)r)^{-2} dH_n(r; \Sigma)}, \quad (68)$$

$$\tilde{v}_v(-\lambda; \gamma, \Sigma)^{-1} = v(-\lambda; \gamma, \Sigma)^{-2} - \int \gamma r^2(1 + v(-\lambda; \gamma, \Sigma)r)^{-2} dH_n(r; \Sigma), \quad (69)$$

where $H_n(\cdot; \Sigma)$ is the empirical distribution (supported on $[r_{\min}, r_{\max}]$) of the eigenvalues of Σ .

Though Lemma F.6 states the dependency explicitly, we will simply write $H_n(r)$, $v(-\lambda; \gamma)$, $\tilde{v}_b(-\lambda; \gamma, \mathbf{A})$, and $\tilde{v}_v(-\lambda; \gamma)$ to denote $H_n(r; \Sigma)$, $v(-\lambda; \gamma, \Sigma)$, $\tilde{v}_b(-\lambda; \gamma, \Sigma, \mathbf{A})$, and $\tilde{v}_v(-\lambda; \gamma, \Sigma)$, respectively, for simplifying notations when it is clear from the context. When $\mathbf{A} = \Sigma$, we simply write $\tilde{v}_b(-\lambda; \gamma) = \tilde{v}_b(-\lambda; \gamma, \mathbf{A})$. The moment assumption of order $8 + \alpha$ for some $\alpha > 0$ in the above lemma can be relaxed to only requiring the existence of moments of order $4 + \alpha$ by a truncation argument as in the proof of Theorem 6 of Hastie et al. (2022) (in Appendix A.4 therein). We omit the details and refer the readers to Hastie et al. (2022).

F.2.2. EXTENDED RIDGE RESOLVENTS

The lemma below extends the deterministic equivalents of the ridge resolvents in Lemma F.6 to provide deterministic equivalents for Tikhonov resolvents, where the regularization matrix $\lambda\mathbf{I}_p$ is replaced with $\lambda(\mathbf{I}_p + \mathbf{C})$ and $\mathbf{C} \in \mathbb{R}^{p \times p}$ is an arbitrary positive semidefinite random matrix.

Lemma F.7 (Tikhonov resolvents, adapted from Patil et al. (2022a)). *Suppose the conditions in Lemma F.6 holds. Let $\mathbf{C} \in \mathbb{R}^{p \times p}$ be any symmetric and positive semidefinite random matrix with uniformly bounded operator norm in p that is independent to \mathbf{X} for all $n, p \in \mathbb{N}$, and let $\mathbf{N} = (\widehat{\Sigma} + \lambda\mathbf{I}_p)^{-1}$. Then the following statements hold:*

(1) *Tikhonov resolvent:*

$$\lambda(\mathbf{N}^{-1} + \lambda\mathbf{C})^{-1} \simeq \widetilde{\Sigma}_C^{-1}. \quad (70)$$

(2) *Bias of Tikhonov regression:*

$$\lambda^2(\mathbf{N}^{-1} + \lambda\mathbf{C})^{-1}\Sigma(\mathbf{N}^{-1} + \lambda\mathbf{C})^{-1} \simeq \widetilde{\Sigma}_C^{-1}(\tilde{v}_b(-\lambda; \gamma, \Sigma_C)\Sigma + \Sigma)\widetilde{\Sigma}_C^{-1}. \quad (71)$$

(3) *Variance of Tikhonov regression:*

$$(\mathbf{N}^{-1} + \lambda\mathbf{C})^{-1}\widehat{\Sigma}(\mathbf{N}^{-1} + \lambda\mathbf{C})^{-1}\Sigma \simeq \tilde{v}_v(-\lambda; \gamma, \Sigma_C)\widetilde{\Sigma}_C^{-1}\Sigma\widetilde{\Sigma}_C^{-1}\Sigma, \quad (72)$$

where $\Sigma_C = (\mathbf{I}_p + \mathbf{C})^{-\frac{1}{2}}\Sigma(\mathbf{I}_p + \mathbf{C})^{-\frac{1}{2}}$, $\widetilde{\Sigma}_C = v(-\lambda; \gamma, \Sigma_C)\Sigma + \mathbf{I}_p + \mathbf{C}$. Here, $v(-\lambda; \gamma, \Sigma_C)$, $\tilde{v}_b(-\lambda; \gamma, \Sigma_C)$, and $\tilde{v}_v(-\lambda; \gamma, \Sigma_C)$ defined by (67)-(69) simplify to

$$\frac{1}{v(-\lambda; \gamma, \Sigma_C)} = \lambda + \gamma \operatorname{tr}[(v(-\lambda; \gamma, \Sigma_C)\Sigma + \mathbf{I}_p + \mathbf{C})^{-1}\Sigma]/p, \quad (73)$$

$$\frac{1}{\tilde{v}_v(-\lambda; \gamma, \Sigma_C)} = \frac{1}{v(-\lambda; \gamma, \Sigma_C)^2} - \gamma \operatorname{tr}[(v(-\lambda; \gamma, \Sigma_C)\Sigma + \mathbf{I}_p + \mathbf{C})^{-2}\Sigma^2]/p, \quad (74)$$

$$\tilde{v}_b(-\lambda; \gamma, \Sigma_C) = \gamma \operatorname{tr}[(v(-\lambda; \gamma, \Sigma_C)\Sigma + \mathbf{I}_p + \mathbf{C})^{-2}\Sigma^2]/p \cdot \tilde{v}_v(-\lambda; \gamma, \Sigma_C). \quad (75)$$

If $\gamma \rightarrow \phi \in (0, \infty)$, then γ in (1)-(3) can be replaced by ϕ , with the empirical distribution H_n of eigenvalues replaced by the limiting distribution H .

F.2.3. RESOLVENTS FOR TRAINING ERROR

The following lemma concerns the deterministic equivalents of quantities that arise in the proof for Lemma D.3.

Lemma F.8 (Resolvents for in-sample error). *Suppose the conditions in Lemma F.6 holds. Let $\mathbf{C} \in \mathbb{R}^{p \times p}$ be any symmetric and positive semidefinite random matrix with uniformly bounded operator norm in p that is independent to \mathbf{X} for all $n, p \in \mathbb{N}$. Let $I_1, I_2 \stackrel{\text{SRS}}{\sim} \mathcal{I}_k$ and $\widehat{\Sigma}_j$ be the sample covariance matrix computed using k observations of \mathbf{X} indexed by I_j ($j = 0, 1$). For $j = 1, 2$, let $\mathbf{M}_j = (\widehat{\Sigma}_j + \lambda \mathbf{I}_p)^{-1}$ be the resolvent for $\widehat{\Sigma}_j$. Then, as $k, n, p \rightarrow \infty$ such that $p/n \rightarrow \phi \in (0, \infty)$ and $p/k \rightarrow \phi_s \in [\phi, \infty)$, the following statements hold:*

(1) *Independent product with sample covariance:*

$$\mathbf{C} \widehat{\Sigma}_j \simeq \mathbf{C} \Sigma.$$

(2) *Bias term 1:*

$$\lambda^2 \mathbf{M}_1 \mathbf{C} \mathbf{M}_2 \simeq (v(-\lambda; \phi_s) \Sigma + \mathbf{I}_p)^{-1} (\tilde{v}(-\lambda; \phi, \phi_s, \mathbf{C}) \Sigma + \mathbf{C}) (v(-\lambda; \phi_s) \Sigma + \mathbf{I}_p)^{-1}. \quad (76)$$

(3) *Bias term 2:*

$$\mathbf{M}_1 \widehat{\Sigma}_{1 \cap 2} \mathbf{M}_2 \mathbf{C} \simeq \tilde{v}_v(-\lambda; \phi, \phi_s) (v(-\lambda; \phi_s) \Sigma + \mathbf{I}_p)^{-2} \Sigma \mathbf{C}, \quad (77)$$

(4) *Variance term 1:*

$$\mathbf{M}_1 \widehat{\Sigma}_{1 \cap 2} \simeq \mathbf{I}_p - (v(-\lambda; \phi_s) \Sigma + \mathbf{I}_p)^{-1}, \quad (78)$$

(5) *Variance term 2:*

$$\begin{aligned} \mathbf{M}_1 \widehat{\Sigma}_{1 \cap 2} \mathbf{M}_2 \widehat{\Sigma}_{1 \cap 2} &\simeq \frac{\phi_s}{\phi} \left(v(-\lambda; \phi_s) - \frac{\phi_s - \phi}{\phi_s} \lambda \tilde{v}_v(-\lambda; \phi, \phi_s) \right) (v(-\lambda; \phi_s) \Sigma + \mathbf{I}_p)^{-1} \Sigma \\ &\quad - \lambda \tilde{v}_v(-\lambda; \phi, \phi_s) (v(-\lambda; \phi_s) \Sigma + \mathbf{I}_p)^{-2} \Sigma, \end{aligned} \quad (79)$$

where

$$\begin{aligned} \tilde{v}(-\lambda; \phi, \phi_s, \mathbf{C}) &= \frac{\lim_{k, n, p} \phi \operatorname{tr}[\mathbf{C} \Sigma (v(-\lambda; \phi_s) \Sigma + \mathbf{I}_p)^{-2}] / p}{v(-\lambda; \phi_s)^{-2} - \phi \int \frac{r^2}{(1 + v(-\lambda; \phi_s) r)^2} dH(r)}, \\ \tilde{v}_v(-\lambda; \phi, \phi_s) &:= \frac{1}{v(-\lambda; \phi_s)^{-2} - \phi \int \frac{r^2}{(1 + v(-\lambda; \phi_s) r)^2} dH(r)}. \end{aligned}$$

Proof of Lemma F.8. We split the proof into different parts.

Part (1) Note that $\operatorname{tr}(\widehat{\Sigma}_j) = \sum_{i \in I_j} \|\mathbf{x}_i\|_2^2 / k$ and $\mathbf{x}_i = \mathbf{z}_i^\top \Sigma \mathbf{z}_i$. By Lemma G.4, we have that $\operatorname{tr}(\widehat{\Sigma}_j) / p - \operatorname{tr}(\Sigma) / p \xrightarrow{\text{a.s.}} 0$. Since $\|\mathbf{C}\|_{\text{op}}$ is uniformly upper bounded and

$$\left| \frac{1}{p} \operatorname{tr}(\mathbf{C} \widehat{\Sigma}_j) - \frac{1}{p} \operatorname{tr}(\mathbf{C} \Sigma) \right| \leq \frac{1}{p} |\operatorname{tr}(\mathbf{C}(\widehat{\Sigma}_j - \Sigma))| \leq \frac{1}{p} \|\mathbf{C}\|_{\text{op}} |\operatorname{tr}(\widehat{\Sigma}_j - \Sigma)|,$$

it follows that $\frac{1}{p} \operatorname{tr}(\mathbf{C} \widehat{\Sigma}_j) - \frac{1}{p} \operatorname{tr}(\mathbf{C} \Sigma) \xrightarrow{\text{a.s.}} 0$, which implies that $\mathbf{C} \widehat{\Sigma}_j \simeq \mathbf{C} \Sigma$

Part (2) This is a direct consequence of Patil et al. (2022a, Part (c) of the proof for Lemma S.24).

Part (3) This is a direct consequence of Patil et al. (2022a, Part (c) of the proof for Lemma S.25).

Part (4) Let $i_0 = |I_1 \cap I_2|$. Conditioning on $\widehat{\Sigma}_{1 \cap 2}$ and i_0 , from Definition F.2 and Lemma F.7 (1) we have

$$\lambda \mathbf{M}_1 \simeq \mathbf{M}_{\mathbf{M}_{1 \cap 2}, i_0}^{\det} := \frac{k}{k - i_0} (v_1 \Sigma + \mathbf{I}_p + \mathbf{C}_1)^{-1} \Big| i_0,$$

where $v_1 = v(-\lambda; \gamma_1, \Sigma_{\mathbf{C}_1})$, $\Sigma_{\mathbf{C}_1} = (\mathbf{I}_p + \mathbf{C}_1)^{-\frac{1}{2}} \Sigma (\mathbf{I}_p + \mathbf{C}_1)^{-\frac{1}{2}}$, $\mathbf{C}_1 = i_0 (\lambda(k - i_0))^{-1} \mathbf{M}_{1 \cap 2}^{-1}$, and $\gamma_1 = p/(k - i_0)$. Here the subscripts of v_1 and \mathbf{C}_1 are related to the aspect ratio γ_1 . Because

$$\limsup \left\| \widehat{\Sigma}_{1 \cap 2} \right\|_{\text{op}} \leq r_{\max} (1 + \sqrt{\phi_s^2 / \phi})^2,$$

almost surely as $k, n, p \rightarrow \infty$ such that $p/n \rightarrow \phi$ and $p/k \rightarrow \phi_s$, by Lemma F.3 (3), we have

$$\mathbf{M}_1 \widehat{\Sigma}_{1 \cap 2} \simeq \lambda^{-1} \mathbf{M}_{\mathbf{M}_{1 \cap 2}, i_0}^{\det} \widehat{\Sigma}_{1 \cap 2} \Big| i_0.$$

That is,

$$\mathbf{M}_1 \widehat{\Sigma}_{1 \cap 2} \simeq \frac{k}{i_0} (\mathbf{M}_{1 \cap 2}^{-1} + \lambda \mathbf{C}_0)^{-1} \widehat{\Sigma}_{1 \cap 2} \Big| i_0,$$

where $\mathbf{C}_0 = (k - i_0)/i_0 \cdot (v_1 \Sigma + \mathbf{I}_p)$. Define $\Sigma_{\mathbf{C}_0} = (\mathbf{I} + \mathbf{C}_0)^{-\frac{1}{2}} \Sigma (\mathbf{I} + \mathbf{C}_0)^{-\frac{1}{2}}$. Conditioning on i_0 , by Lemma F.7 (1), we have

$$\begin{aligned} \text{tr}[\Sigma_{\mathbf{C}_1} (v_1 \Sigma_{\mathbf{C}_1} + \mathbf{I}_p)^{-1}] &= \text{tr}[\Sigma (v_1 \Sigma + \mathbf{I}_p + \mathbf{C}_1)^{-1}] \\ &= \frac{\lambda(k - i_0)}{i_0} \text{tr} \left[\Sigma \left(\mathbf{M}_{1 \cap 2}^{-1} + \frac{\lambda(k - i_0)}{i_0} (v_1 \Sigma + \mathbf{I}_p) \right)^{-1} \right] \\ &\stackrel{\text{a.s.}}{=} \frac{k - i_0}{i_0} \text{tr} \left[\Sigma \left(v_0 \Sigma + \mathbf{I}_p + \frac{k - i_0}{i_0} (v_1 \Sigma + \mathbf{I}_p) \right)^{-1} \right] \\ &= \text{tr} \left[\Sigma \left(\left(\frac{i_0}{k - i_0} v_0 + v_1 \right) \Sigma + \frac{k}{k - i_0} \mathbf{I}_p \right)^{-1} \right], \end{aligned}$$

where $v_0 = v(-\lambda; \gamma_0, \Sigma_{\mathbf{C}_0})$ and $\gamma_0 = p/i_0$. Note that the fixed-point solution v_0 depends on v_1 . The fixed-point equations reduce to

$$\begin{aligned} \frac{1}{v_0} &= \lambda + \gamma_0 \text{tr}[\Sigma_{\mathbf{C}_0} (v_0 \Sigma_{\mathbf{C}_0} + \mathbf{I}_p)^{-1}] / p = \lambda + \frac{p}{k} \text{tr} \left[\Sigma \left(\left(\frac{i_0}{k} v_0 + \frac{k - i_0}{k} v_1 \right) \Sigma + \mathbf{I}_p \right)^{-1} \right] / p \\ \frac{1}{v_1} &= \lambda + \gamma_1 \text{tr}[\Sigma_{\mathbf{C}_1} (v_1 \Sigma_{\mathbf{C}_1} + \mathbf{I}_p)^{-1}] / p = \lambda + \frac{p}{k} \text{tr} \left[\Sigma \left(\left(\frac{i_0}{k} v_0 + \frac{k - i_0}{k} v_1 \right) \Sigma + \mathbf{I}_p \right)^{-1} \right] / p \end{aligned}$$

almost surely. Note that the solution (v_0, v_1) to the above equations is a pair of positive numbers and does not depend on samples. If (v_0, v_1) is a solution to the above system, then (v_1, v_0) is also a solution. Thus, any solution to the above equations must be unique. On the other hand, since $v_0 = v_1 = v(-\lambda; p/k)$ satisfies the above equations, it is the unique solution. By Lemma F.3 (7), we can replace $v(-\lambda; \gamma_1, \Sigma_{\mathbf{C}_1})$ by the solution $v_0 = v_1 = v(-\lambda; p/k)$ of the above system, which does not depend on samples. Thus,

$$\mathbf{M}_1 \widehat{\Sigma}_{1 \cap 2} \simeq \frac{k}{i_0} (\mathbf{M}_{1 \cap 2}^{-1} + \lambda \mathbf{C}^*)^{-1} \widehat{\Sigma}_{1 \cap 2} \Big| i_0, \quad (80)$$

where $\mathbf{C}^* = (k - i_0)/i_0 \cdot (v(-\lambda; p/k) \Sigma + \mathbf{I}_p)$. Again from Lemma F.7 (1) we have

$$\begin{aligned} (\mathbf{M}_{1 \cap 2}^{-1} + \lambda \mathbf{C}^*)^{-1} \widehat{\Sigma}_{1 \cap 2} &= \mathbf{I}_p - \lambda (\mathbf{M}_{1 \cap 2}^{-1} + \lambda \mathbf{C}^*)^{-1} (\mathbf{I}_p + \mathbf{C}^*) \\ &\simeq \mathbf{I}_p - (v(-\lambda; p/k) \Sigma + \mathbf{I}_p + \mathbf{C}^*)^{-1} (\mathbf{I}_p + \mathbf{C}^*) \Big| i_0 \\ &= \frac{i_0}{k} (\mathbf{I}_p - (v(-\lambda; p/k) \Sigma + \mathbf{I}_p)^{-1}). \end{aligned}$$

Finally, from Lemma F.3 (6), we have

$$\mathbf{M}_1 \widehat{\Sigma}_{1 \cap 2} \simeq \mathbf{I}_p - (v(-\lambda; p/k) \Sigma + \mathbf{I}_p)^{-1} \simeq \mathbf{I}_p - (v(-\lambda; \phi_s) \Sigma + \mathbf{I}_p)^{-1}.$$

Part (5) From Patil et al. (2022a, Part (c) of the proof for Lemma S.2.5), we have that

$$\mathbf{M}_1 \widehat{\boldsymbol{\Sigma}}_{1\cap 2} \mathbf{M}_2 \widehat{\boldsymbol{\Sigma}}_{1\cap 2} \simeq \frac{k^2}{i_0^2} (\mathbf{M}_{1\cap 2}^{-1} + \lambda \mathbf{C}^*)^{-1} \widehat{\boldsymbol{\Sigma}}_{1\cap 2} (\mathbf{M}_{1\cap 2}^{-1} + \lambda \mathbf{C}^*)^{-1} \widehat{\boldsymbol{\Sigma}}_{1\cap 2},$$

where $\mathbf{M}_{1\cap 2} = (\widehat{\boldsymbol{\Sigma}}_{1\cap 2} + \lambda \mathbf{I}_p)^{-1}$ and $\mathbf{C}^* = (k - i_0)/i_0(v(-\lambda; \phi_s)\boldsymbol{\Sigma} + \mathbf{I}_p)$. Since

$$(\mathbf{M}_{1\cap 2}^{-1} + \lambda \mathbf{C}^*)^{-1} \widehat{\boldsymbol{\Sigma}}_{1\cap 2} = \mathbf{I}_p - \lambda (\mathbf{M}_{1\cap 2}^{-1} + \lambda \mathbf{C}^*)^{-1} (\mathbf{I}_p + \mathbf{C}^*),$$

we have

$$\begin{aligned} & \mathbf{M}_1 \widehat{\boldsymbol{\Sigma}}_{1\cap 2} \mathbf{M}_2 \widehat{\boldsymbol{\Sigma}}_{1\cap 2} \\ & \simeq \frac{k^2}{i_0^2} (\mathbf{M}_{1\cap 2}^{-1} + \lambda \mathbf{C}^*)^{-1} \widehat{\boldsymbol{\Sigma}}_{1\cap 2} - \lambda \frac{k^2}{i_0^2} (\mathbf{M}_{1\cap 2}^{-1} + \lambda \mathbf{C}^*)^{-1} \widehat{\boldsymbol{\Sigma}}_{1\cap 2} (\mathbf{M}_{1\cap 2}^{-1} + \lambda \mathbf{C}^*)^{-1} (\mathbf{I}_p + \mathbf{C}^*) \\ & = \frac{k^2}{i_0^2} (\mathbf{I}_p - \lambda (\mathbf{M}_{1\cap 2}^{-1} + \lambda \mathbf{C}^*)^{-1} (\mathbf{I}_p + \mathbf{C}^*)) - \lambda \frac{k^2}{i_0^2} (\mathbf{M}_{1\cap 2}^{-1} + \lambda \mathbf{C}^*)^{-1} \widehat{\boldsymbol{\Sigma}}_{1\cap 2} (\mathbf{M}_{1\cap 2}^{-1} + \lambda \mathbf{C}^*)^{-1} (\mathbf{I}_p + \mathbf{C}^*) \end{aligned} \quad (81)$$

From Lemma F.7 (1) and (3), we have that

$$\begin{aligned} & \lambda (\mathbf{M}_{1\cap 2}^{-1} + \lambda \mathbf{C}^*)^{-1} \simeq \frac{\phi}{\phi_s} (v(-\lambda; \phi_s)\boldsymbol{\Sigma} + \mathbf{I}_p)^{-1} \\ & (\mathbf{M}_{1\cap 2}^{-1} + \lambda \mathbf{C}^*)^{-1} \widehat{\boldsymbol{\Sigma}}_{1\cap 2} (\mathbf{M}_{1\cap 2}^{-1} + \lambda \mathbf{C}^*)^{-1} \simeq \frac{\phi^2}{\phi_s^2} \tilde{v}_v(-\lambda; \phi, \phi_s) (v(-\lambda; \phi_s)\boldsymbol{\Sigma} + \mathbf{I}_p)^{-2} \boldsymbol{\Sigma}. \end{aligned}$$

Combing the above two equivalents, the expression in (81) can be further simplified as:

$$\begin{aligned} \mathbf{M}_1 \widehat{\boldsymbol{\Sigma}}_{1\cap 2} \mathbf{M}_2 \widehat{\boldsymbol{\Sigma}}_{1\cap 2} & \simeq \frac{\phi_s}{\phi} \left(v(-\lambda; \phi_s) - \frac{\phi_s - \phi}{\phi_s} \lambda \tilde{v}_v(-\lambda; \phi, \phi_s) \right) (v(-\lambda; \phi_s)\boldsymbol{\Sigma} + \mathbf{I}_p)^{-1} \boldsymbol{\Sigma} \\ & \quad - \lambda \tilde{v}_v(-\lambda; \phi, \phi_s) (v(-\lambda; \phi_s)\boldsymbol{\Sigma} + \mathbf{I}_p)^{-2} \boldsymbol{\Sigma}. \end{aligned}$$

□

F.3. Analytic properties of associated fixed-point equations

In this section, we gather results on the properties of the fixed-point solution $v(-\lambda; \phi)$ defined in (64).

The following lemma provides the existence and uniqueness of the solution $v(-\lambda; \phi)$. The properties of the derivatives in Lemma F.9 are related to the properties of $\tilde{v}_v(-\lambda; \phi)$ defined in Lemma F.10, which equals $-f'(x)$, where the function f is defined in (82).

Lemma F.9 (Properties of the solution to the fixed-point equation, adapted from Patil et al. (2022a)). *Let $\lambda, \phi, a > 0$ and $b < \infty$ be real numbers. Let P be a probability measure supported on $[a, b]$. Define the function f such that*

$$f(x) = \frac{1}{x} - \phi \int \frac{r}{1+rx} dP(r) - \lambda. \quad (82)$$

Then, the following properties hold:

- (1) For $\lambda = 0$ and $\phi \in (1, \infty)$, there is a unique $x_0 \in (0, \infty)$ such that $f(x_0) = 0$. The function f is positive and strictly decreasing over $(0, x_0)$ and negative over (x_0, ∞) , with $\lim_{x \rightarrow 0^+} f(x) = \infty$ and $\lim_{x \rightarrow \infty} f(x) = 0$.
- (2) For $\lambda > 0$ and $\phi \in (0, \infty)$, there is a unique $x_0^\lambda \in (0, \infty)$ such that $f(x_0^\lambda) = 0$. The function f is positive and strictly decreasing over $(0, x_0^\lambda)$ and negative over (x_0^λ, ∞) , with $\lim_{x \rightarrow 0^+} f(x) = \infty$ and $\lim_{x \rightarrow \infty} f(x) = -\lambda$.
- (3) For $\lambda = 0$ and $\phi \in (1, \infty)$, f is differentiable on $(0, \infty)$ and its derivative f' is strictly increasing over $(0, x_0)$, with $\lim_{x \rightarrow 0^+} f'(x) = -\infty$ and $f'(x_0) < 0$.
- (4) For $\lambda > 0$ and $\phi \in (0, \infty)$, f is differentiable on $(0, \infty)$ and its derivative f' is strictly increasing over $(0, \infty)$, with $\lim_{x \rightarrow 0^+} f'(x) = -\infty$ and $f'(x_0^\lambda) < 0$.

The continuity and limiting behavior of the function $\phi \mapsto v(-\lambda; \phi)$ is given for ridge regression ($\lambda > 0$) in Lemma F.10 and for ridgeless regression ($\lambda = 0$) in Lemma F.11.

Lemma F.10 (Continuity in the aspect ratio for ridge regression, adapted from Patil et al. (2022a)). *Let $\lambda, a > 0$ and $b < \infty$ be real numbers. Let P be a probability measure supported on $[a, b]$. Consider the function $v(-\lambda; \cdot) : \phi \mapsto v(-\lambda; \phi)$, over $(0, \infty)$, where $v(-\lambda; \phi) > 0$ is the unique solution to the fixed-point equation*

$$\frac{1}{v(-\lambda; \phi)} = \lambda + \phi \int \frac{r}{1 + rv(-\lambda; \phi)} dP(r) \quad (83)$$

Then, the following properties hold:

- (1) The range of the function $v(-\lambda; \cdot)$ is a subset of $(0, \lambda^{-1})$.
- (2) The function $v(-\lambda; \cdot)$ is continuous and strictly decreasing over $(0, \infty)$. Furthermore, $\lim_{\phi \rightarrow 0^+} v(-\lambda; \phi) = \lambda^{-1}$, and $\lim_{\phi \rightarrow \infty} v(-\lambda; \phi) = 0$.
- (3) The function $\tilde{v}_v(-\lambda; \cdot) : \phi \mapsto \tilde{v}_v(-\lambda; \phi)$, where

$$\tilde{v}_v(-\lambda; \phi) = \left(v(-\lambda; \phi)^{-2} - \int \phi r^2 (1 + rv(-\lambda; \phi))^{-2} dP(r) \right)^{-1},$$

is positive and continuous over $(0, \infty)$. Furthermore, $\lim_{\phi \rightarrow 0^+} \tilde{v}_v(-\lambda; \phi) = \lambda^{-2}$, and $\lim_{\phi \rightarrow \infty} \tilde{v}_v(-\lambda; \phi) = 0$.

- (4) The function $\tilde{v}_b(-\lambda; \cdot) : \phi \mapsto \tilde{v}_b(-\lambda; \phi)$, where

$$\tilde{v}_b(-\lambda; \phi) = \tilde{v}_v(-\lambda; \phi) \int \phi r^2 (1 + v(-\lambda; \phi)r)^{-2} dP(r),$$

is positive and continuous over $(0, \infty)$. Furthermore, $\lim_{\phi \rightarrow 0^+} \tilde{v}_b(-\lambda; \phi) = \lim_{\phi \rightarrow \infty} \tilde{v}_b(-\lambda; \phi) = 0$.

Lemma F.11 (Continuity in the aspect ratio for ridgeless regression, adapted from Patil et al. (2022b)). *Let $a > 0$ and $b < \infty$ be real numbers. Let P be a probability measure supported on $[a, b]$. Consider the function $v(0; \cdot) : \phi \mapsto v(0; \phi)$, over $(1, \infty)$, where $v(0; \phi) > 0$ is the unique solution to the fixed-point equation*

$$\frac{1}{\phi} = \int \frac{v(0; \phi)r}{1 + v(0; \phi)r} dP(r). \quad (84)$$

Then, the following properties hold:

- (1) The function $v(0; \cdot)$ is continuous and strictly decreasing over $(1, \infty)$. Furthermore, $\lim_{\phi \rightarrow 1^+} v(0; \phi) = \infty$, and $\lim_{\phi \rightarrow \infty} v(0; \phi) = 0$.
- (2) The function $\phi \mapsto (\phi v(0; \phi))^{-1}$ is strictly increasing over $(1, \infty)$. Furthermore, $\lim_{\phi \rightarrow 1^+} (\phi v(0; \phi))^{-1} = 0$ and $\lim_{\phi \rightarrow \infty} (\phi v(0; \phi))^{-1} = 1$.
- (3) The function $\tilde{v}_v(0; \cdot) : \phi \mapsto \tilde{v}_v(0; \phi)$, where

$$\tilde{v}_v(0; \phi) = \left(v(0; \phi)^{-2} - \phi \int r^2 (1 + rv(0; \phi))^{-2} dP(r) \right)^{-1},$$

is positive and continuous over $(1, \infty)$. Furthermore, $\lim_{\phi \rightarrow 1^+} \tilde{v}_v(0; \phi) = \infty$, and $\lim_{\phi \rightarrow \infty} \tilde{v}_v(0; \phi) = 0$.

- (4) The function $\tilde{v}_b(0; \cdot) : \phi \mapsto \tilde{v}_b(0; \phi)$, where

$$\tilde{v}_b(0; \phi) = \tilde{v}_v(0; \phi) \int r^2 (1 + v(0; \phi)r)^{-2} dP(r),$$

is positive and continuous over $(1, \infty)$. Furthermore, $\lim_{\phi \rightarrow 1^+} \tilde{v}_b(0; \phi) = \infty$, and $\lim_{\phi \rightarrow \infty} \tilde{v}_b(0; \phi) = 0$.

The continuity and differentiability of the function $\lambda \mapsto v(-\lambda; \phi)$ on a closed interval $[0, \lambda_{\max}]$ for some constant λ_{\max} is given for $\phi \in (1, \infty)$ in Lemma F.12 adapted from Patil et al. (2022b). This ensures that $v(0; \phi) = \lim_{\lambda \rightarrow 0^+} v(-\lambda; \phi)$ is well-defined for $\phi > 1$ and also implies that related functions are bounded.

Lemma F.12 (Differentiability in the regularization parameter). *Let $0 < a \leq b < \infty$ be real numbers. Let P be a probability measure supported on $[a, b]$. Let $\phi > 0$ be a real number. Let $\Lambda = [0, \lambda_{\max}]$ for some constant $\lambda_{\max} \in (0, \infty)$. For $\lambda \in \Lambda$, let $v(-\lambda; \phi) > 0$ denote the solution to the fixed-point equation*

$$\frac{1}{v(-\lambda; \phi)} = \lambda + \phi \int \frac{r}{v(-\lambda; \phi)r + 1} dP(r).$$

When $\lambda = 0$ and $\phi \in (0, 1]$, $v(-\lambda; \phi) := +\infty$. Then, the following properties hold:

- (1) (Monotonicity) For $\phi \in (0, \infty)$, the function $\lambda \mapsto v(-\lambda; \phi)$ is strictly decreasing in $\lambda \in [0, \infty)$.
- (2) (Differentiability) For $\phi \in (1, \infty)$, the function $\lambda \mapsto v(-\lambda; \phi)$ is twice differentiable over Λ .
- (3) (Boundedness of the second derivative) For $\phi \in (1, \infty)$, $v(-\lambda; \phi)$, $\partial/\partial\lambda[v(-\lambda; \phi)]$, and $\partial^2/\partial\lambda^2[v(-\lambda; \phi)]$ are bounded over Λ .

Proof of Lemma F.12. Start by re-writing the fixed-point equation as

$$\lambda = \frac{1}{v(-\lambda; \phi)} - \phi \int \frac{r}{v(-\lambda; \phi)r + 1} dP(r).$$

Define a function f by

$$f(x) = \frac{1}{x} - \phi \int \frac{r}{xr + 1} dP(r).$$

Observe that $v(-\lambda; \phi) = f^{-1}(\lambda)$. We next study various properties of f and prove the different parts in the statement.

Part (1) Properties of f and f^{-1} :

Observe that

$$f(x) = \frac{1}{x} - \phi \int \frac{r}{xr + 1} dP(r) = \frac{1}{x} \left(1 - \phi \int \frac{xr}{xr + 1} dP(r) \right).$$

The function $g : x \mapsto 1/x$ is positive and strictly decreasing over $(0, \infty)$ with $\lim_{x \rightarrow 0^+} g(x) = \infty$ and $\lim_{x \rightarrow \infty} g(x) = 0$, while the function

$$h : x \mapsto 1 - \phi \int \frac{xr}{xr + 1} dP(r)$$

is strictly decreasing over $(0, \infty)$ with $h(0) = 1$ and $\lim_{x \rightarrow \infty} h(x) = 1 - \phi$.

Thus, there is a unique $0 < x_0 < \infty$ when $\phi > 1$ such that $h(x_0) = 0$, and consequently $f(x_0) = 0$; and $x_0 = +\infty$ when $\phi \in (0, 1]$ such that $g(x_0) = 0$, and consequently $f(x_0) = 0$. Because h and g are positive over $[0, x_0)$, f , a product of two positive strictly decreasing functions, is strictly decreasing over $(0, x_0)$, with $\lim_{x \rightarrow 0^+} f(x) = \infty$ and $f(x_0) = 0$.

Because f is strictly decreasing over $(0, x_0)$, f^{-1} is strictly decreasing (see, e.g., Problem 2, Chapter 5 of Rudin (1976)). Since $f(x_0) = 0$, $f^{-1}(0) = x_0$, and since $\lim_{x \rightarrow 0^+} f(x) = \infty$, $\lim_{y \rightarrow \infty} f^{-1}(y) = 0$. Hence, f^{-1} is strictly decreasing over $[0, \infty)$ for all $\phi > 0$ and bounded above by $x_0 < \infty$ for all $\phi > 1$.

Parts (2) and (3) We will prove the remaining two parts together.

Properties of f' and $(f^{-1})'$:

The derivative f' at x is given by

$$f'(x) = -\frac{1}{x^2} + \phi \int \frac{r^2}{(xr + 1)^2} dP(r) = -\frac{1}{x^2} \left(1 - \phi \int \left(\frac{xr}{xr + 1} \right)^2 dP(r) \right).$$

The function $g : x \mapsto 1/x^2$ is positive and strictly decreasing over $(0, \infty)$ with $\lim_{x \rightarrow 0^+} g(x) = \infty$ and $\lim_{x \rightarrow \infty} g(x) = 0$. On the other hand, the function

$$h : x \mapsto 1 - \phi \int \left(\frac{xr}{xr+1} \right)^2 dP(r)$$

is strictly decreasing over $(0, \infty)$ with $h(0) = 1$ and $h(x_0) > 0$. This follows because for $x \in [0, x_0]$,

$$\begin{aligned} \phi \int \left(\frac{xr}{xr+1} \right)^2 dP(r) &\leq \left(\frac{x_0 b}{x_0 b + 1} \right) \phi \int \left(\frac{xr}{xr+1} \right) dP(r) \\ &< \phi \int \frac{xr}{xr+1} dP(r) \leq \phi \int \frac{x_0 r}{x_0 r + 1} dP(r) = 1, \end{aligned} \quad (85)$$

where the first inequality in the chain above follows as the support of P is $[a, b]$, and the last inequality follows since $f(x_0) = 0$ and $x_0 > 0$, which implies that

$$\frac{1}{x_0} = \phi \int \frac{r}{x_0 r + 1} dP(r), \quad \text{or equivalently that} \quad 1 = \phi \int \frac{x_0 r}{x_0 r + 1} dP(r).$$

Thus, $-f'$, a product of two positive strictly decreasing functions, is strictly decreasing, and in turn, f' is strictly increasing. Moreover, $\lim_{x \rightarrow 0^+} f'(x) = -\infty$; when $\phi > 1$, $f'(x_0) < 0$ and when $\phi \in (0, 1]$, $f'(x)$ approaches zero from below as $x \rightarrow +\infty$.

When $\phi > 1$, because $f'(x) \neq 0$ over $(0, x_0)$, by the inverse function theorem, $(f^{-1})'$, we have

$$|(f^{-1})'(f(x))| = \left| \frac{1}{f'(x)} \right| < \left| \frac{1}{f'(x_0)} \right| = \frac{1}{\frac{1}{x_0^2} \left(1 - \phi \int \left(\frac{xr}{xr+1} \right)^2 dP(r) \right)} < \infty,$$

where the first inequality uses the fact that $|f'(x_0)| < |f'(x)|$ for $x \in (0, x_0]$ from Part 1, and the last inequality uses the bound from (85).

Properties of f'' and $(f^{-1})''$:

The second derivative f'' at x is given by

$$f''(x) = \frac{2}{x^3} - 2\phi \int \frac{r^3}{(xr+1)^3} dP(r) = \frac{2}{x^3} \left(1 - \phi \int \left(\frac{xr}{xr+1} \right)^3 dP(r) \right).$$

The rest of the arguments are similar to those in Part 2. The function $g : x \mapsto 1/x^3$ is positive and strictly decreasing over $(0, \infty)$ with $\lim_{x \rightarrow 0^+} g(x) = \infty$ and $\lim_{x \rightarrow \infty} g(x) = 0$, while the function

$$h : x \mapsto 1 - \phi \int \left(\frac{xr}{xr+1} \right)^3 dP(r)$$

is strictly decreasing over $(0, \infty)$ with $h(0) = 1$ and $h(x_0) > 0$ as

$$\begin{aligned} \phi \int \left(\frac{xr}{xr+1} \right)^3 dP(r) &\leq \left(\frac{x_0 b}{x_0 b + 1} \right)^2 \phi \int \left(\frac{xr}{xr+1} \right) dP(r) \\ &< \phi \int \frac{xr}{xr+1} dP(r) \leq \phi \int \frac{x_0 r}{x_0 r + 1} dP(r) = 1. \end{aligned} \quad (86)$$

It then follows that f'' is strictly decreasing, with $\lim_{x \rightarrow 0^+} f''(x) = \infty$; when $\phi > 1$, $f''(x_0) > 0$ and when $\phi \in (0, 1]$, $f''(x)$ approaches zero from above as $x \rightarrow +\infty$.

When $\phi > 1$, by inverse function theorem, we have

$$|(f^{-1})''(f(x))| = \left| \frac{f''(x)}{f'(x)^3} \right| = \frac{\frac{2}{x^3} \left(1 - \phi \int \left(\frac{xr}{xr+1} \right)^3 dP(r) \right)}{\left(\frac{1}{x_0^2} \left(1 - \phi \int \left(\frac{xr}{xr+1} \right)^2 dP(r) \right) \right)^3} \leq \frac{2x_0^3}{\left(1 - \phi \int \left(\frac{xr}{xr+1} \right)^2 dP(r) \right)^3} < \infty,$$

where the first inequality uses the bound from (86), and the second inequality uses the bound from (85).

This finishes all the parts and concludes the proof. \square

G. Helper concentration results

G.1. Size of the intersection of randomly sampled datasets

In this section, we collect various helper results concerned with concentrations and convergences. Below we recall the definition of a hypergeometric random variable, along with its mean and variance. See, e.g., [Greene & Wellner \(2017\)](#) for more related details.

Definition G.1 (Hypergeometric random variable). A random variable X follows the hypergeometric distribution $X \sim \text{Hypergeometric}(n, K, N)$ if its probability mass function is given by

$$\mathbb{P}(X = k) = \frac{\binom{K}{k} \binom{N-K}{n-k}}{\binom{N}{n}}, \quad \max\{0, n + K - N\} \leq k \leq \min\{n, K\}.$$

The expectation and variance of X are given by

$$\mathbb{E}[X] = \frac{nK}{N}, \quad \text{Var}(X) = \frac{nK(N-K)(N-n)}{N^2(N-1)}.$$

The following lemma characterizes the limiting proportions of shared observations in two simple random samples under proportional asymptotics, when both the subsample size and the full data size tend to infinity, which is adapted from [Patil et al. \(2022a\)](#).

Lemma G.2 (Asymptotic proportions of shared observations). For $n \in \mathbb{N}$, define $\mathcal{I}_k := \{\{i_1, i_2, \dots, i_k\} : 1 \leq i_1 < i_2 < \dots < i_k \leq n\}$. Let $I_1, I_2 \stackrel{\text{SRSSWR}}{\sim} \mathcal{I}_k$, define the random variable $i_0^{\text{SRSSWR}} := |I_1 \cap I_2|$ to be the number of shared samples, and define i_0^{SRSSWOR} accordingly. Let $\{k_m\}_{m=1}^\infty$ and $\{n_m\}_{m=1}^\infty$ be two sequences of positive integers such that n_m is strictly increasing in m , $n'_m \leq k_m \leq n_m$ for some constant $\nu \in (0, 1)$. Then, $i_0^{\text{SRSSWR}}/k_m - k_m/n_m \xrightarrow{\text{a.s.}} 0$, and $i_0^{\text{SRSSWOR}}/k_m - k_m/n_m \xrightarrow{\text{a.s.}} 0$.

G.2. Convergence of random linear and quadratic forms

In this section, we collect helper lemmas on the concentration of linear and quadratic forms of random vectors.

The following lemma provides the concentration of a linear form of a random vector with independent components. It follows from a moment bound from Lemma 7.8 of [Erdős & Yau \(2017\)](#), along with the Borel-Cantelli lemma, and is adapted from Lemma S.8.5 of [Patil et al. \(2022b\)](#).

Lemma G.3 (Concentration of linear form with independent components). Let $\mathbf{z}_p \in \mathbb{R}^p$ be a sequence of random vector with i.i.d. entries z_{pi} , $i = 1, \dots, p$ such that for each i , $\mathbb{E}[z_{pi}] = 0$, $\mathbb{E}[z_{pi}^2] = 1$, $\mathbb{E}[|z_{pi}|^{4+\alpha}] \leq M_\alpha$ for some $\alpha > 0$ and constant $M_\alpha < \infty$. Let $\mathbf{a}_p \in \mathbb{R}^p$ be a sequence of random vectors independent of \mathbf{z}_p such that $\limsup_p \|\mathbf{a}_p\|^2/p \leq M_0$ almost surely for a constant $M_0 < \infty$. Then, $\mathbf{a}_p^\top \mathbf{z}_p/p \rightarrow 0$ almost surely as $p \rightarrow \infty$.

The following lemma provides the concentration of a quadratic form of a random vector with independent components. It follows from a moment bound from Lemma B.26 of [Bai & Silverstein \(2010\)](#), along with the Borel-Cantelli lemma, and is adapted from Lemma S.8.6 of [Patil et al. \(2022b\)](#).

Lemma G.4 (Concentration of quadratic form with independent components). Let $\mathbf{z}_p \in \mathbb{R}^p$ be a sequence of random vector with i.i.d. entries z_{pi} , $i = 1, \dots, p$ such that for each i , $\mathbb{E}[z_{pi}] = 0$, $\mathbb{E}[z_{pi}^2] = 1$, $\mathbb{E}[|z_{pi}|^{4+\alpha}] \leq M_\alpha$ for some $\alpha > 0$ and constant $M_\alpha < \infty$. Let $\mathbf{D}_p \in \mathbb{R}^{p \times p}$ be a sequence of random matrix such that $\limsup \|\mathbf{D}_p\|_{\text{op}} \leq M_0$ almost surely as $p \rightarrow \infty$ for some constant $M_0 < \infty$. Then, $\mathbf{z}_p^\top \mathbf{D}_p \mathbf{z}_p/p - \text{tr}[\mathbf{D}_p]/p \rightarrow 0$ almost surely as $p \rightarrow \infty$.

G.3. Convergence of Cesàro-type mean and max for triangular array

In this section, we collect a helper lemma on deducing almost sure convergence of a Cesàro-type mean from almost sure convergence of the original sequence, which is adapted from [Patil et al. \(2022a\)](#).

Lemma G.5 (Convergence of conditional expectation). *For $n \in \mathbb{N}$, suppose $\{R_{n,\ell}\}_{\ell=1}^{N_n}$ is a set of N_n random variables defined over the probability space $(\Omega, \mathcal{F}, \mathbb{P})$, with $1 < N_n < \infty$ almost surely. If there exists a constant c such that $R_{n,p_n} \xrightarrow{\text{a.s.}} c$ for all deterministic sequences $\{p_n \in [N_n]\}_{n=1}^\infty$, then the following statements hold:*

$$(1) \max_{\ell \in [N_n]} |R_{n,\ell}(\omega) - c| \xrightarrow{\text{a.s.}} 0.$$

$$(2) N_n^{-1} \sum_{\ell=1}^{N_n} R_{n,\ell} \xrightarrow{\text{a.s.}} c.$$

Lemma G.6 (Convergence of conditional expectation over simple random sampling). *For $n \in \mathbb{N}$ and $k = k_n \in \mathcal{K}_n$, let $M_n = |\mathcal{I}_k|$ and suppose $\{R_{n,1}(I_\ell)\}_{\ell \in [M_n]}$ and $\{R_{n,2}(I_m, I_\ell)\}_{m, \ell \in [M_n], m \neq \ell}$ are sets of M_n and $M_n(M_n - 1)$ random variables, such that $R_{n,2}(I_m, I_\ell) \leq (R_{n,1}(I_m) + R_{n,2}(I_\ell))/2$. Then, the following holds:*

$$(1) \text{ If there exists a constant } c_1 \text{ such that } R_{n,1}(I_{\ell_n}) \xrightarrow{\text{a.s.}} c_1 \text{ for all deterministic sequences } \{\ell_n \in [M_n]\}_{n=1}^\infty, \text{ then } \max_{\ell \in [M_n]} |R_{n,\ell}(I_\ell) - c| \xrightarrow{\text{a.s.}} 0 \text{ and } \mathbb{E}_{I_\ell \sim \mathcal{I}_k}^{\text{SRS}} [|R_{n,\ell}(I_\ell) - c|] \xrightarrow{\text{a.s.}} 0.$$

$$(2) \text{ Further, if there exists a constant } c_2 \text{ such that } R_{n,2}(I_{m_n}, I_{\ell_n}) \xrightarrow{\text{a.s.}} c_2 \text{ for all sequences of simple random samples } \{(I_{m_n}, I_{\ell_n})\}_{n=1}^\infty \stackrel{\text{SRS}}{\sim} \mathcal{I}_{k_n}, \text{ then } \max_{(I_m, I_\ell) \sim \mathcal{I}_k}^{\text{SRS}} |R_{n,2}(I_m, I_\ell) - c_2| \xrightarrow{\text{a.s.}} 0 \text{ and } \mathbb{E}_{(I_m, I_\ell) \sim \mathcal{I}_k}^{\text{SRS}} [|R_{n,2}(I_m, I_\ell) - c_2|] \xrightarrow{\text{a.s.}} 0.$$

Proof of Lemma G.6. We split the proof into two cases.

Part (1) The conclusion directly follows from Lemma G.5.

Part (2) Observe that

$$R_{n,2}(I_m, I_\ell) \leq \frac{1}{2}(R_{n,1}(I_m) + R_{n,2}(I_\ell)). \quad (87)$$

From (1), we have that $\mathbb{E}_I[R_{n,1}(I)] \xrightarrow{\text{a.s.}} c_1$, where the expectation is taken with respect to the uniform distribution over \mathcal{I}_k . From the condition, we have $R_{n,2}(I_m, I_\ell) \xrightarrow{\text{a.s.}} c_2$ for any $I_m, I_\ell \stackrel{\text{SRS}}{\sim} \mathcal{I}_k$. Then, by Pratt's lemma (see, e.g., Gut, 2005, Theorem 5.5), the conclusion follows. \square

H. GCV correction for arbitrary M

Note that the asymptotic limit of the training error for arbitrary $M \in \mathbb{N}$ is given by

$$\mathcal{F}_M^\lambda = 2\mathcal{E}_{k,2}^\lambda - \mathcal{E}_{k,1}^\lambda + \frac{2}{M}(\mathcal{E}_{k,1}^\lambda - \mathcal{E}_{k,2}^\lambda),$$

where $\mathcal{E}_{k,j}^\lambda = c_{k,M,j} \mathcal{F}_{k,j}^\lambda + (1 - c_{k,M,j}) \mathcal{B}_{k,j}^\lambda$. Here, $c_{k,M,j}$ is the limiting proportion of the distinct number of observations from j simple random samples to the distinct number of observations from M simple random samples of size k . Roughly speaking, the proportion of unseen observations from M simple random samples of size k is $(n - k)^M / n^M$ and thus

$$c_{k,M,j} = \lim \frac{1 - (n - k)^j / n^j}{1 - (n - k)^M / n^M} = \frac{1 - (1 - \phi / \phi_s)^j}{1 - (1 - \phi / \phi_s)^M}.$$

From the expression, one knows for certain that the GCV asymptotics will not match the risk of the estimator in general. In addition, the form of the expression also leads to an approach to correct the GCV estimator for general M that we will discuss below. What we prove in Theorem 3.1 is that the difference between the two asymptotics vanishes as $M \rightarrow \infty$. We expect the difference to scale as $1/M$. The explicit analysis of the finite-ensemble effect requires carefully analyzing the coefficients $c_{k,M,j}$, and even for the isotropic design, the expression for the GCV asymptotics for general appears to be very involved. It is in principle possible to perform this analysis, but we did not pursue it further in the paper given our primary focus on the full-ensemble estimator. Numerically, we observe that the bias is small for a moderate M (e.g., for $M = 10$) and a reasonable data model with SNR (SNR = 0.6) from Figure 4. In general, we expect this to be the case for either moderate k or M and typical real-world SNR ranges. We will consider adding more numerical illustrations of the finite-ensemble effect in the revision under different settings.

We aim to define the corrected GCV as

$$\overline{\text{gcv}}_{k,M}^\lambda := \frac{a_1 T_{k,M}^\lambda + a_2 \bar{R}_{k,M}^\lambda}{D_{k,M}^\lambda},$$

where a_1 and a_2 are two unknown parameters to be determined. To determine the unknown parameters, we need to match the limiting GCV with the true risk. Since

$$\mathcal{T}_1^\lambda(\phi, \phi_s) = \mathcal{D}_1^\lambda(\phi, \phi_s) \mathcal{R}_1^\lambda(\phi, \phi_s), \quad \text{and} \quad \mathcal{T}_2^\lambda(\phi, \phi_s) = b_1 \mathcal{R}_1^\lambda(\phi, \phi_s) + b_2 \mathcal{R}_2^\lambda(\phi, \phi_s),$$

for some known constants b_1 and b_2 which can be derived in the proof of Proposition 3.3, the adjustment is given by

$$\begin{aligned} & a_1 \left[- \left(1 - \frac{2}{M} \right) (c_{k,M,1} \mathcal{D}_1^\lambda(\phi, \phi_s) + 1 - c_{k,M,1}) \mathcal{R}_1^\lambda(\phi, \phi_s) \right. \\ & \quad \left. + 2 \left(1 + \frac{1}{M} \right) (c_{k,M,2} b_1 \mathcal{R}_1^\lambda(\phi, \phi_s) + (1 - c_{k,M,2} + b_2) \mathcal{R}_2^\lambda(\phi, \phi_s)) \right] \\ & \quad + \mathcal{D}_M^\lambda(\phi, \phi_s) a_2 \left(- \left(1 - \frac{2}{M} \right) \mathcal{R}_1^\lambda(\phi, \phi_s) + 2 \left(1 + \frac{1}{M} \right) \mathcal{R}_2^\lambda(\phi, \phi_s) \right) \\ & = - \left(1 - \frac{2}{M} \right) [a_1 (c_{k,M,1} \mathcal{D}_1^\lambda(\phi, \phi_s) + 1 - c_{k,M,1}) + a_2 \mathcal{D}_M^\lambda(\phi, \phi_s)] \mathcal{R}_1^\lambda(\phi, \phi_s) \\ & \quad + 2 \left(1 + \frac{1}{M} \right) [a_1 (c_{k,M,2} b_1 \mathcal{R}_1^\lambda(\phi, \phi_s) + (1 - c_{k,M,2} + b_2) \mathcal{R}_2^\lambda(\phi, \phi_s)) + a_2 \mathcal{D}_M^\lambda(\phi, \phi_s)] \mathcal{R}_2^\lambda(\phi, \phi_s), \end{aligned}$$

which implies that

$$\begin{aligned} & a_1 (c_{k,M,1} \mathcal{D}_1^\lambda(\phi, \phi_s) + 1 - c_{k,M,1}) + a_2 \mathcal{D}_M^\lambda(\phi, \phi_s) = \mathcal{D}_M^\lambda(\phi, \phi_s) \\ & a_1 (c_{k,M,2} b_1 \mathcal{R}_1^\lambda(\phi, \phi_s) + (1 - c_{k,M,2} + b_2) \mathcal{R}_2^\lambda(\phi, \phi_s)) + a_2 \mathcal{D}_M^\lambda(\phi, \phi_s) \mathcal{R}_2^\lambda(\phi, \phi_s) = \mathcal{D}_M^\lambda(\phi, \phi_s) \mathcal{R}_2^\lambda(\phi, \phi_s). \end{aligned}$$

Solving the above linear system for $a_1 > 0$ and $a_2 \in \mathbb{R}$ gives the correct weights for defining a consistent GCV estimate. The solutions will depend on $\mathcal{D}_M^\lambda(\phi, \phi_s)$, $\mathcal{R}_1^\lambda(\phi, \phi_s)$ and $\mathcal{R}_2^\lambda(\phi, \phi_s)$. For the denominator, (12) is a consistent estimate for $\mathcal{D}_M^\lambda(\phi, \phi_s)$. For the prediction risks of $M = 1, 2$, we can use out-of-bag observations to estimate $\mathcal{R}_1^\lambda(\phi, \phi_s)$ and $\mathcal{R}_2^\lambda(\phi, \phi_s)$.

I. Additional details for numerical experiments

The covariance matrix of an auto-regressive process of order 1 (AR(1)) is given by Σ_{ar1} , where $(\Sigma_{\text{ar1}})_{ij} = \rho_{\text{AR1}}^{|i-j|}$ for some parameter $\rho_{\text{AR1}} \in (0, 1)$, and the AR(1) data model is defined as:

$$\begin{aligned} & y_i = \mathbf{x}_i^\top \boldsymbol{\beta}_0 + \epsilon_i, \quad \mathbf{x}_i \sim \mathcal{N}(0, \Sigma_{\text{ar1}}), \\ & \boldsymbol{\beta}_0 = \frac{1}{5} \sum_{j=1}^5 \mathbf{w}_{(j)}, \quad \epsilon_i \sim \mathcal{N}(0, \sigma^2), \end{aligned} \tag{M-AR1}$$

where $\mathbf{w}_{(j)}$ is the eigenvector of Σ_{ar1} associated with the top j th eigenvalue $r_{(j)}$. From Grenander & Szegő (1958, pp. 69-70), the top j -th eigenvalue can be written as $r_{(j)} = (1 - \rho_{\text{AR1}}^2) / (1 - 2\rho_{\text{AR1}} \cos \theta_{jp} + \rho_{\text{AR1}}^2)$ for some $\theta_{jp} \in ((j-1)\pi/(p+1), j\pi/(p+1))$. Then, under model (M-AR1), the signal strength ρ^2 defined in Assumption 2.2 is $5^{-1}(1 - \rho_{\text{AR1}}^2) / (1 - \rho_{\text{AR1}})^2$, which is the limit of $25^{-1} \sum_{j=1}^5 r_{(j)}$. Thus, model (M-AR1) parameterized by two parameters ρ_{AR1} and σ^2 satisfies Assumption 2.1-2.2.