

Precise Asymptotics of Bagging Regularized M-estimators

Takuya Koriyama^{*†} Pratik Patil^{*‡} Jin-Hong Du^{§¶}
tkoriyam@uchicago.edu pratikpatil@berkeley.edu jinhongd@andrew.cmu.edu

Kai Tan^{||} Pierre C. Bellec^{||}
kai.tan@rutgers.edu pierre.bellec@rutgers.edu

Abstract

We characterize the squared prediction risk of ensemble estimators obtained through subagging (subsample bootstrap aggregating) regularized M-estimators and construct a consistent estimator for the risk. Specifically, we consider a heterogeneous collection of $M \geq 1$ regularized M-estimators, each trained with (possibly different) subsample sizes, convex differentiable losses, and convex regularizers. We operate under the proportional asymptotics regime, where the sample size n , feature size p , and subsample sizes k_m for $m \in [M]$ all diverge with fixed limiting ratios n/p and k_m/n . Key to our analysis is a new result on the joint asymptotic behavior of correlations between the estimator and residual errors on overlapping subsamples, governed through a (provably) contractible nonlinear system of equations. Of independent interest, we also establish convergence of trace functionals related to degrees of freedom in the non-ensemble setting (with $M = 1$) along the way, extending previously known cases for square loss and ridge, lasso regularizers.

When specialized to homogeneous ensembles trained with a common loss, regularizer, and subsample size, the risk characterization sheds some light on the implicit regularization effect due to the ensemble and subsample sizes (M, k) . For any ensemble size M , optimally tuning subsample size yields sample-wise monotonic risk. For the full-ensemble estimator (when $M \rightarrow \infty$), the optimal subsample size k^* tends to be in the overparameterized regime ($k^* \leq \min\{n, p\}$), when explicit regularization is vanishing. Finally, joint optimization of subsample size, ensemble size, and regularization can significantly outperform regularizer optimization alone on the full data (without any subagging).

1 Introduction

Ensemble methods combine predictions of multiple models to improve predictive accuracy [HTF09]. Among these methods, bagging (bootstrap aggregating) trains individual models on bootstrapped samples of the dataset and averages their predictions to reduce variance and mitigate overfitting [Bre96]. A popular variant of bagging, known as subagging (subsample bootstrap aggregating), trains models on random subsamples rather than full bootstrapped samples [BY02]. Apart from the computational advantages, subagging can substantially improve predictive performance, especially in the overparameterized regimes and near model interpolation thresholds [PKWR22]. In this paper,

*Corresponding authors.

[†]Booth School of Business, The University of Chicago, Chicago, IL, 60637, USA.

[‡]Department of Statistics, University of California, Berkeley, CA 94720, USA.

[§]Department of Statistics and Data Science, Carnegie Mellon University, Pittsburgh, PA 15213, USA.

[¶]Machine Learning Department, Carnegie Mellon University, Pittsburgh, PA 15213, USA.

^{||}Department of Statistics, Rutgers University, New Brunswick, NJ 08854, USA.

Table 1: **Landscape of subgaging risk analysis in high dimensions.** We summarize the various settings (estimator structure and data structure) of some of the recent works that characterize the prediction risk of subgaging regularized M-estimators trained with loss function `loss` and regularization function `reg` (defined in Section 2). `Convex•` denotes that in addition to being convex, we require the loss function to be differentiable and have Lipschitz continuous derivatives. By RMT features, we refer to features of the form $\mathbf{x} = \Sigma^{1/2}\mathbf{v}$ where \mathbf{v} contains independent entries of bounded moments (of order 4^+) that are common in the random matrix theory literature. `Arbitrary▲` response y refers to no additional modeling assumption on the response other than bounded moments (of order 4^+). Signal and noise refer to $\boldsymbol{\theta}$ and z in the linear data model: $y = \mathbf{x}^\top \boldsymbol{\theta} + z$. `Arbitrary■` signal and noise distributions F_θ and F_z refer to general marginal distribution on the coordinates of signal $\boldsymbol{\theta}$ and noise z , respectively, subject to mild regularity conditions (Assumption D).

Loss (loss)	Penalty (reg)	Features (\mathbf{x})	Covariance (Σ)	Response (y)	Signal ($\boldsymbol{\theta}$)	Signal Dist. (F_θ)	Noise Dist. (F_z)	Reference
Square		Gaussian	Isotropic	Linear	Random	Gaussian	Gaussian	[LJB20]
Square	Ridge	RMT	Anisotropic	Linear	Deterministic		Bnd. Mom.	[PDK23]
Square	Ridge	RMT	Anisotropic	Arbitrary [▲]				[PD23]
Huber		Gaussian	Isotropic	Linear	Deterministic		Arbitrary [■]	[BK24]
Logistic		Gaussian	Isotropic	Logistic	Deterministic			[BK24]
Convex	Ridge	Gaussian	Isotropic	GLM	Random	Gaussian		[CVD ⁺ 24]
Convex [•]	Convex	Gaussian	Isotropic	Linear	Random	Arbitrary [■]	Arbitrary [■]	Corollary 3

we analyze the squared prediction risk of subgaging of regularized M-estimators trained with convex loss and regularizer.

There is growing interest in understanding the prediction risk asymptotics of ensemble methods, particularly subgaging, across different types of predictors and under various data assumptions. For example, [LJB20] study subgaging in the context of ordinary least squares regression without any explicit regularization in the underparameterized regime (where the number of subsamples is higher than the number of features). This is extended to ridge and ridgeless regression by [PDK23] for both the underparameterized and overparameterized regimes (where the number of subsamples is lower than the number of features). [PD23] further generalizes these results and identifies explicit equivalence paths between subsampled estimators and ridge regularized estimators. Beyond linear and ridge regression, [BK24] examines the behavior of subgaging in logistic and Huber regression models without regularization. In addition to subgaging, there has also been considerable work on feature sketching and other ensemble methods. For instance, [LGR⁺22] and [PL24] study feature sketching and ensembling to optimize predictive performance in high-dimensional settings. For other recent developments in the analysis of ensemble methods and related work details, see Section 1.2.

We generalize these previous works by characterizing the prediction risk of subgaging a collection of M regularized M-estimators and constructing a consistent estimator for this risk. We allow the collection to be heterogeneous, consisting of estimators trained with convex differentiable loss function `loss` with Lipschitz continuous derivative and convex regularization function `reg`, which can be non-differentiable and non-strongly convex (that includes ℓ_1 -regularized Huber regression, for instance). We operate under the proportional asymptotics regime, where the sample size n , feature size p , and subsample sizes k_m for $m \in [M]$ diverge while maintaining fixed limiting ratios $n/p \rightarrow \delta \in (0, \infty)$ and $k_m/n \rightarrow c_m \in (0, 1]$.¹ Our results assume a linear response model $y = \mathbf{x}^\top \boldsymbol{\theta} + z$ with isotropic Gaussian features \mathbf{x} , a random signal $\boldsymbol{\theta}$ with independent coordinates drawn from an arbitrary marginal distribution F_θ , and noise variable z with an arbitrary distribution F_z , both

¹Through the paper, we refer to n/p as inverse data aspect ratio and p/n as data aspect ratio of the design $\mathbf{X} \in \mathbb{R}^{n \times p}$, viewing n as the height and p as the width of the rectangular design matrix \mathbf{X} .

subject to a mild regularity condition. We summarize our main results below and situate them within the context of recent related work in Table 1.

1.1 Summary of results and paper outline

A summary of our main results along with an outline for the paper is as follows.²

- (a) **Risk characterization and estimation.** In Section 3, we obtain a precise characterization of the squared risk of subbagging of regularized M-estimators under proportional asymptotics (Theorem 2 and Corollary 3). The asymptotic risk is governed by two nonlinear systems of equations (Systems 1a and 1b) that depend on the loss and regularization functions loss , reg , and the limiting subsample ratios c_m for $m \in [M]$ for the component estimators and the limiting inverse data aspect ratio δ . System 1a is a known system that characterizes the asymptotic risk of regularized M-estimators in non-ensemble settings [TAH18] (see Section 1.2 for more details), while System 1b is a new contribution of this paper. Each scalar unknown in the 2-dimensional System 1b is shown to be the fixed-point equation of a contraction (Theorem 1-(2)). This property plays a crucial role in proving the existence and uniqueness of the solution to System 1b. This contraction property is also crucial in establishing the asymptotic behavior of the inner products between estimator errors and their residuals for estimators trained on overlapping subsamples (Theorem 2), leading to the subagged risk asymptotics (Corollary 3). The contractility, along with a ridge smoothing technique, also allows us to maintain weak assumptions on the regularizer, specifically allowing it to be non-strongly convex. Moreover, we also construct a consistent estimator of the ensemble risk (Theorem 4 and Corollary 5), which can be employed for data-adaptive tuning of hyperparameters such as loss functions, regularizers, and subsample sizes.
- (b) **Homogeneous ensembles.** In Section 4, we consider homogeneous ensembles (where components are trained on the same loss and reg functions and a common subsample size k). In Section 4.1, we analyze (oracle) optimal ensemble optimal risk with optimal ensemble size M^* and subsample size k^* . We first establish the monotonicity of the risk with respect to the ensemble size M (Proposition 6), illustrating the benefits of ensembling, which leads to $M^* = \infty$. We then prove that the risk at the optimal subsample size k_* decreases as the limiting data aspect ratio p/n decreases (Proposition 7). In particular, this implies that the optimally subsampled risk avoids the typical “double (or multiple) descents” observed in regularized M-estimators without subbagging. In Section 4.2, we specialize our main result to convex regularized least squares (including ℓ_q -regularized least squares for $q \geq 1$, such as ridge and lasso estimators) and to general M-estimators (including regularized Huber regression). These recover and generalize various known results in the literature (see Section 4.2 for more details).
- (c) **Subbagging and overparameterization.** In Section 5, we investigate subbagging of estimators with vanishing regularization ($\lambda \rightarrow 0^+$) and also contrast with estimators with optimal explicit regularization (over $\lambda \geq 0$). Our first insight is that when subbagging estimators without any explicit regularization, the optimal subsample size k_* is in the overparameterized regime, regardless of whether the original data aspect ratio p/n is overparameterized. In other words, the optimal subsample size k^* satisfies $k_* \leq \min\{n, p\}$ in such cases. We verify this for the lassoless (minimum ℓ_1 -norm interpolator) ensemble (Figure 5). This highlights the advantages of overparameterization in subbagging in that full-ensemble subsampled lassoless can outper-

²The source code for experimental verifications in this paper is available at the repository [subagging-asymptotics](#). The risk estimator proposed in this paper is also incorporated in the Python library [sklearn-ensemble-cv](#) [DP24b].

form the optimal lasso on the full data (without any subagging). Our second insight is that the joint optimization of the subsample size and explicit regularization parameter can outperform optimizing explicit regularization alone on the full data. We verify this property for ensembles of the lasso, unregularized Huber, and ℓ_1 -regularizer Huber (see Figure 6, 7 and Appendix D). This highlights the benefits of subagging on top of optimal explicit regularization.

Independent results in the non-ensemble setting. In the process of characterizing the risk of the subagging ensemble, we also establish the convergence of certain trace functionals (in particular, see (13) below or the last three rows of Table 2). This implies that the observable adjustments developed in [Bel22] for inference using a single estimator converge to their deterministic counterparts defined as solutions to the System 1a, unifying the mean-field asymptotics featuring System 1a of [TAH18] and the inference results of [Bel22]. This was known only for the lasso [CMW23, Theorem 8] and unregularized M-estimators [BK24]. These convergence results are new for regularized estimators (beyond ridge and lasso) and robust loss functions (beyond squared loss) and are of independent interest even for a single estimator (that is, in the non-ensemble setting with $M = 1$).

1.2 Other related work

Resampling methods, such as bagging and subsampling, are widely used in statistics and machine learning. Given their broad applicability, there is a vast literature on these methods. In this section, we provide an overview of the literature related to the risk analysis of ensemble methods, particularly in high-dimensional regimes that have received considerable recent interest.

Classical work on bagging and subagging includes the work by [Bre96, Bre01, BY02], among others. Beyond bagging, analysis of ensemble methods of different predictors includes smooth weak predictors [BS06, FH07], nonparametric estimators [BY02, LGR⁺22], and classifiers [HS05, Sam12]. Historically, there are also early works by [SK95, KS97] on risk asymptotics for ridge ensembles under Gaussian features. We also mention here some other early work on ensembles, including: [HS90, Per93, SK95, KS97]. For a comprehensive overview of early work on bagging and ensemble methods in general, we refer readers to [PDK23].

Substantial progress has been made in the last decade in understanding the asymptotic behavior of regularized M-estimators in high-dimensional settings, particularly under the proportional asymptotic regime where the number of features scales with the number of observations. Frameworks of Approximate Message Passing (AMP) (developed in a series of papers [DMM09, DMM11, BM11a]), Convex Gaussian Min-Max Theorem (CGMT) (developed in a series of papers [OTH13, OH16, TOH15, TAH18]), and leave-one-out (LOO) and martingale-based analysis common in random matrix theory (used in [EK13, Kar18], for example) have been instrumental in deriving the limiting test risk, often as solutions to (nonlinear) systems of self-consistent equations. More specifically, these include analyses of unregularized estimators [EKBB⁺13, EK13, Kar18, DM16, BBEKY13], ridge estimator [DW18], lasso [BM11b, MM21], bridge estimators [WMZ18], logistic regression [SC19, MLC19, SAH19], convex regularized M-estimators [TOH15, TAH18], among others. Recently, triggered by the empirical success of neural networks that interpolate, these risk analyses have been extended to interpolating estimators with vanishing regularization (in the overparameterized regimes that allow for interpolation), such as ridgeless [HMRT22], lassoless [LW21], max-margin interpolators [MRSY19, DKT22, LS22]; see the survey papers [BMR21, Bel21] for other related references.

Beyond individual regularized M-estimators, there has now been considerable interest over the last few years in the analysis of ensembles of estimators in high-dimensional settings, especially in the

overparameterized regime just mentioned. In particular, [LJB20] consider least squares ensembles obtained by subsampling such that the final subsampled dataset has more observations than the number of features. The work of [PDK23] provides the characterization of the asymptotic risk of ensembles of ridge regression using results from Random Matrix Theory (RMT). Furthermore, recent extensions by [DPK23, PD23] expand the scope of these results by establishing risk equivalences for both optimal and suboptimal risks, considering arbitrary feature covariance and signal structures. Other follow-up works for subbagging of ridge and ridgeless regression include [CZYS23, AK23]. This paper develops tools to study ensembles of regularized estimators with general loss and regularizers, beyond ridge regularization. For instance, our theory accommodates ℓ_1 -regularized Huber regression.

Another line of research focuses on ensemble methods involving random features and feature sketching rather than subsampling. In random features models, the effect of ensembling on various components of the risk has been studied in [AP20, dRBK20, LGR⁺22]. Recently, [PL24] analyze ensembles of ridge regression with sketched features with asymptotically free sketching. There are also analyses of alternative resampling and averaging techniques. For example, in the context of distributed learning, [DS20, DS21, MRRK22] consider the divide-and-conquer approach, or splagging (split aggregating), and investigate their properties for ridge and ridgeless predictors.

Very recently, [CVD⁺24] analyzed the limiting equations of several resampling schemes, including bootstrap and resampling without replacement, and characterized self-consistent equations for the limiting bias and variance functionals of estimators obtained by minimization of the negative log-likelihood plus an additive ridge penalty. This is related to our risk characterization as [CVD⁺24] also covers sampling without replacement, but our nonlinear systems (Systems 1a and 1b) characterizing the subbagging risk do not appear explicitly in their work, which instead focuses on self-consistent equations for bias and variance functionals of the specific resampling scheme. The results of [CVD⁺24] relies on the general AMP analysis and state evolution laid out in [LGR⁺22, Lemmas B.3 and B.5], generalizing [BM11b]. This analysis requires the existence of unique solution to the corresponding limiting system of equations, which is granted under strong convexity (e.g., with a ridge penalty), but was not established until the present paper for the case of ensembling of subsampled regularized estimators.

Finally, complementary to risk characterization, there has also been considerable interest in the cross-validation and model selection of ensemble methods. In particular, [DPK23] study cross-validation for bagging of ridge regression. [BDK⁺24] examine the consistency of generalized cross-validation (GCV) for estimating the prediction risk of arbitrary ensembles of regularized least squares estimators for strongly convex penalties. They show that GCV is inconsistent for any finite ensemble of size greater than one and identify a correction to GCV that is consistent for any finite ensemble size, termed corrected GCV (CGCV). In this paper, we generalize one of the data-dependent estimators proposed in [BDK⁺24] for the general setting of this paper, allowing for general convex losses and heterogeneous component estimators in the ensemble. While we do not attempt to interpret the estimator as a corrected GCV for homogeneous ensembles in the general setting, it may be possible to perform such an analysis further, which we leave for future work.

The proof strategy in this paper extends the approach of [BK24], which studies the bagging of unregularized M-estimators. While their analysis is based on a relatively simple 1-dimensional nonlinear system, the new System 1b below is 2-dimensional, introducing additional complexity to the analysis. The rise in complexity is similar to that from unregularized regression [EK13, DM16, Kar18] and its 2-dimensional system to the 4-dimensional system of regularized M-estimators [TAH18] given in System 1a. One challenge arises from the stochastic control of the trace terms

in (13). In the unregularized case, these trace terms can be approximated by a straightforward product of the norms of the error vector and residuals (see [BK24, Lemma 5.7]), and the stochastic behavior of these norms are well-understood in the existing literature (cf. [TAH18, LGC⁺21]). However, in the regularized case, the trace functional cannot be approximated by such a simple expression, which prevents the direct application of these existing results based on the CGMT. We overcome this by showing that, with high probability, the trace term is a stationary point of a certain (random) strongly convex function. This allows us to control the perturbation of the trace term through the behavior of the convex function (see Appendix A.4).

1.3 Notation

We denote scalars in the regular lower or upper case (e.g., a , A), vectors in bold lower case (e.g., \mathbf{a}), and matrices in bold upper case (e.g., \mathbf{A}). For a natural number n , the shorthand notation $[n]$ denotes the set $\{1, \dots, n\}$. For two real numbers x and y , we use $x \wedge y$ to denote $\min\{x, y\}$. For a vector \mathbf{a} , $\|\mathbf{a}\|_q$ denotes its ℓ_q -norm for $q \geq 1$. If no subscript is present for the norm $\|\mathbf{u}\|$ of a vector \mathbf{u} , then it is assumed to be the ℓ_2 norm of \mathbf{u} . For a univariate function f and a vector $\mathbf{a} \in \mathbb{R}^n$, with a slight overload of notation, we use $f(\mathbf{a}) \in \mathbb{R}^n$ to denote the component-wise application of f to \mathbf{a} . We use $\text{diag}[\mathbf{a}]$ to denote the diagonal matrix whose entries are given by the vector \mathbf{a} . Throughout, we use $\mathbf{0}$, $\mathbf{1}$, and \mathbf{I} to respectively denote the all-zero vector, all-one vector, and identity matrix of varying dimensions, depending on the context.

For any proper, closed, convex function $f: \mathbb{R} \rightarrow \bar{\mathbb{R}}$, the proximal operator and Moreau envelope of f with a parameter $\tau > 0$ at a point $x \in \mathbb{R}$ are, respectively, denoted as:³

$$\text{prox}_f(x; \tau) := \underset{y \in \mathbb{R}}{\text{argmin}} f(y) + \frac{1}{2\tau}(x - y)^2 \quad \text{and} \quad \text{env}_f(x; \tau) := \min_{y \in \mathbb{R}} f(y) + \frac{1}{2\tau}(x - y)^2. \quad (1)$$

For a proper, closed, convex f , the argmin in (1) exists and is unique, and consequently $x \mapsto \text{prox}_f(x; \tau)$ is a well-defined function. Let ∂f denote the subdifferential of f , which is the set of all subgradients of f . We jot down two key relationships between the proximal operator, subdifferential, and Moreau envelopes of f below for the reader's convenience:

$$\frac{\partial}{\partial x} \text{env}_f(x; \tau) = \frac{x - \text{prox}_f(x; \tau)}{\tau} \in \partial f(\text{prox}_f(x; \tau)). \quad (2)$$

For simplicity, we often use $\text{env}'_f(x; \tau)$ to denote the partial derivative $\frac{\partial}{\partial x} \text{env}_f(x; \tau)$.⁴

Finally, we use $\mathcal{O}_{\mathbb{P}}$ and $o_{\mathbb{P}}$ to denote probabilistic big-O and little-o notation, respectively, while the convergences in probability are denoted by $\xrightarrow{\mathbb{P}}$. Most other notation we use is standard and any other non-standard notation is defined inline. For the reader's convenience, we also give a quick overview of the specific notation used in this paper in Table 4.

2 Setup and assumptions

We consider the standard supervised regression setting, in which we observe n data points (\mathbf{x}_i, y_i) for $i \in [n]$. The feature matrix $\mathbf{X} \in \mathbb{R}^{n \times p}$ contains \mathbf{x}_i^\top in its i -th row and the response vector $\mathbf{y} \in \mathbb{R}^n$ contains y_i in its i -th entry. We assume the following distribution on the dataset (\mathbf{X}, \mathbf{y}) :

³Here $\bar{\mathbb{R}}$ is the extended real line (that does two-point $(+\infty$ and $-\infty)$ compactification of the real line).

⁴In general, for a bivariate function $g(\cdot; \cdot)$, we use the notation $g'(\cdot; \cdot)$ to denote the first partial derivative with respect to the *first* argument.

Assumption A (Data distribution). The distribution of (\mathbf{X}, \mathbf{y}) is specified by:

1. The design matrix $\mathbf{X} \in \mathbb{R}^{n \times p}$ has i.i.d. entries drawn from $\mathcal{N}(0, 1/p)$.
2. The response vector $\mathbf{y} = \mathbf{X}\boldsymbol{\theta} + \mathbf{z}$, where $\boldsymbol{\theta} \in \mathbb{R}^p$ is a random signal vector and $\mathbf{z} \in \mathbb{R}^n$ is a random noise vector, both independent of each other and of the feature matrix \mathbf{X} , with:
 - (a) The signal vector $\boldsymbol{\theta} \in \mathbb{R}^p$ has i.i.d. entries drawn from distribution F_θ .
 - (b) The noise vector $\mathbf{z} \in \mathbb{R}^n$ has i.i.d. entries drawn from distribution F_z .

We subsample the dataset (\mathbf{X}, \mathbf{y}) to create M subsampled datasets. Towards that end, define M subsample index subsets $I_m \subset [n]$ of cardinality $k_m = |I_m|$ for $m \in [M]$. The feature matrix and response vector associated with the subsampled dataset (\mathbf{x}_i, y_i) for $i \in I_m$ are denoted as $(\mathbf{X}_{I_m}, \mathbf{y}_{I_m})$. We assume the following sampling strategy for the subsample index sets $\{I_m\}_{m \in [M]}$:

Assumption B (Subsampling strategy). Given deterministic integers $\{k_m \geq 1\}_{m \in [M]}$, the M subsample index sets $\{I_m\}_{m \in [M]}$ are independent of (\mathbf{X}, \mathbf{y}) and are independently sampled from the uniform distribution over subsets of $[n]$ with cardinality k_m for each $m \in [M]$.

It is worth noting that if I_m and I_ℓ (for $m \neq \ell$) are any two independent subsample sets of cardinality k_m and k_ℓ per Assumption B, then the cardinality of intersection $|I_m \cap I_\ell|$ follows a hypergeometric distribution with mean $k_m k_\ell / n$. Using the properties of the hypergeometric distribution, it follows that $|I_m \cap I_\ell| / n \xrightarrow{P} c_m c_\ell$ as both $n, k_m, k_\ell \rightarrow \infty$ with the subsample ratios $k_m / n \rightarrow c_m$ and $k_\ell / n \rightarrow c_\ell$ for some $c_m, c_\ell \in (0, 1]$ (this follows from Chebyshev's inequality and the variance formula of the hypergeometric distribution, see Section S.8.1 of [PDK23] for more details). Intuitively, each sample lands in a subsample I_m with probability c_m (the limiting ratio $|I_m| / n$) and in the overlap of two subsamples with probability $c_m c_\ell$ (the limiting ratio $|I_m \cap I_\ell| / n$), as the subsamples are drawn independently. The overlap between any two subsample sets I_m and I_ℓ is thus of order n with high probability. The randomness in subsampling in Assumption B is not important. Our results can accommodate deterministic sampling where the subsample sets $\{I_m\}_{m \in [M]}$ are selected deterministically, provided that the ratios $|I_m| / n$, $|I_\ell| / n$, and $|I_m \cap I_\ell| / n$ converge to non-zero constants.

For each subsampled dataset $(\mathbf{X}_{I_m}, \mathbf{y}_{I_m})$ for $m \in [M]$, we define the regularized M-estimator $\hat{\boldsymbol{\theta}}_m$ as

$$\hat{\boldsymbol{\theta}}_m(I_m) \in \operatorname{argmin}_{\mathbf{b} \in \mathbb{R}^p} \sum_{i \in I_m} \operatorname{loss}_m(y_i - \mathbf{x}_i^\top \mathbf{b}) + \sum_{j \in [p]} \operatorname{reg}_m(\mathbf{b}_j). \quad (3)$$

When defining (3), we allow the argmin operator to return any one of the minimizers (as emphasized by the element notation in (3)). Here loss_m and reg_m are the loss and regularization functions that satisfy the following assumption for all $m \in [M]$:

Assumption C (Loss and regularizer structure). The loss function $\operatorname{loss}: \mathbb{R} \rightarrow \bar{\mathbb{R}}$ is proper, closed, convex, and differentiable with derivative loss' Lipschitz continuous, and $\min_x \operatorname{loss}(x) = \operatorname{loss}(0)$. The regularizer $\operatorname{reg}: \mathbb{R} \rightarrow \bar{\mathbb{R}}$ is proper, closed, convex and $\min_x \operatorname{reg}(x) = \operatorname{reg}(0)$.

The final ensemble estimator, constructed using the component estimators (3), is defined as:

$$\tilde{\boldsymbol{\theta}}_M(\{I_m\}_{m \in [M]}) := \frac{1}{M} \sum_{m \in [M]} \hat{\boldsymbol{\theta}}_m(I_m). \quad (4)$$

For brevity, we omit the dependency of the component and ensemble estimators on I_m and $\{I_m\}_{m \in [M]}$ and simply write $\hat{\boldsymbol{\theta}}_m$ and $\tilde{\boldsymbol{\theta}}_M$, respectively, when it is clear from the context. We evaluate the performance of the ensemble estimator $\tilde{\boldsymbol{\theta}}_M$ with respect to the signal parameter $\boldsymbol{\theta}$ via:

$$R_M := \frac{1}{p} \|\tilde{\boldsymbol{\theta}}_M - \boldsymbol{\theta}\|_2^2.$$

Note that R_M is the (excess) out-of-sample squared error in our setup because of isotropic features: For an independently sampled test feature $\mathbf{x}_0 \in \mathbb{R}^p$ with i.i.d. entries drawn from $\mathcal{N}(0, 1/p)$, we have $R_M = \mathbb{E}[(\mathbf{x}_0^\top \tilde{\boldsymbol{\theta}}_M - \mathbf{x}_0^\top \boldsymbol{\theta})^2 \mid \mathbf{y}, \mathbf{X}, \{I_m\}_{m \in [M]}]$. We will refer to R_M as the risk of the ensemble. Observe that R_M is a scalar random variable that depends on both the dataset (\mathbf{X}, \mathbf{y}) and the random samples I_m for $m \in [M]$. The goal of the paper is to characterize the asymptotic behavior of this random variable R_M under the proportional asymptotics regime. In this regime, the sample size n , feature size p , and subsample size k_m all diverge while keeping the appropriate ratios fixed: we will assume the inverse data aspect ratio $n/p \rightarrow \delta \in (0, \infty)$ and for each $m \in [M]$, the subsample ratio $k_m/n \rightarrow c_m \in (0, 1]$ as $n, p, k_m \rightarrow \infty$.

3 Risk characterization and estimation

In this section, we will first describe a general technical result on the correlations of the error and residual vectors for regularized M-estimator in Section 3.1. We then state our general result on the risk characterization of the ensemble estimator in Section 3.3 and construct a consistent risk estimator for this risk in Section 3.4.

3.1 Asymptotics of correlations of estimator and residual errors

To state the risk characterization of the ensemble estimator, we first introduce two important nonlinear systems of equations: Systems 1a and 1b. Intuitively, these systems correspond to the corner cases where the ensemble size $M = 1$ and $M = \infty$, respectively. As we shall see in Section 3.3, these systems completely determine the risk asymptotics of the ensemble estimator (4).

System 1a (Error norms of individual regularized M-estimator). Given a triple $(\text{loss}, \text{reg}, c\delta)$ where $c\delta \in (0, \infty)$ and $\text{loss}, \text{reg} : \mathbb{R} \rightarrow \mathbb{R}$ are convex functions, define the following 4-scalar system of equations in variables $(\alpha, \beta, \kappa, \nu)$:

$$\alpha^2 = \mathbb{E}\left[\left(\frac{1}{\nu} \text{env}'_{\text{reg}}\left(\Theta + \frac{\beta}{\nu} H; \frac{1}{\nu}\right) - \frac{\beta}{\nu} H\right)^2\right] \quad (5a)$$

$$\beta^2 = \mathbb{E}\left[\text{env}'_{\text{loss}}(Z + \alpha G; \kappa)^2\right] \cdot c\delta \quad (5b)$$

$$\kappa\beta = \mathbb{E}\left[\left(\frac{1}{\nu} \text{env}'_{\text{reg}}\left(\Theta + \frac{\beta}{\nu} H; \frac{1}{\nu}\right) - \frac{\beta}{\nu} H\right) \cdot (-H)\right] \quad (5c)$$

$$\nu\alpha = \mathbb{E}\left[\text{env}'_{\text{loss}}(Z + \alpha G; \kappa) \cdot G\right] \cdot c\delta \quad (5d)$$

where $H \sim \mathcal{N}(0, 1)$, $G \sim \mathcal{N}(0, 1)$, $\Theta \sim F_\theta$, $Z \sim F_z$, all mutually independent.

System 1a can be found in the literature, specifically in [TAH18, Equation 15]. To be precise, we are applying the result of [TAH18] on the subsample estimator (3) using $k = |I_m|$ observations with $k/n \rightarrow c$, so that the limiting inverse aspect ratio $k/p = k/n \cdot n/p \rightarrow c\delta$. System 1a is known to characterize the limit in probability of the risk of (3) when $|I_m|/p \rightarrow c\delta$: if $(\alpha, \beta, \kappa, \nu)$

is a solution to System 1a, then $p^{-1}\|\widehat{\boldsymbol{\theta}}_m - \boldsymbol{\theta}^*\|^2 \xrightarrow{P} \alpha^2$. The existence and uniqueness of the fixed-point parameters in this system are central to applying results from the Convex Gaussian Min-Max Theorem (CGMT) to derive precise risk characterization for regularized M-estimator (under proportional asymptotics). This is guaranteed under conditions where both **loss** and **reg** are Lipschitz and the problem parameters are such that the perfect signal recovery is not possible, leading to non-zero asymptotic risk. For a detailed discussion on these conditions, see [BK23b]. Next, we describe our second system for risk characterization of the ensemble estimators.

System 1b (Error correlations of overlapped regularized M-estimator). Given $c, \tilde{c} \in (0, 1]$ and convex pairs of functions $(\text{loss}, \text{reg}), (\widetilde{\text{loss}}, \widetilde{\text{reg}})$, let $(\alpha, \beta, \kappa, \nu) \in \mathbb{R}_{>0}^4$ and $(\tilde{\alpha}, \tilde{\beta}, \tilde{\kappa}, \tilde{\nu})$ be parameters that satisfy System 1a with $(\text{loss}, \text{reg}, c\delta)$ and $(\widetilde{\text{loss}}, \widetilde{\text{reg}}, \tilde{c}\delta)$, respectively. Define the following 2-scalar system of equations in variable $(\eta_G, \eta_H) \in [-1, 1]^2$:

$$\begin{bmatrix} \eta_G \\ \eta_H \end{bmatrix} = \begin{bmatrix} F_{\text{reg}}(\eta_H) \\ F_{\text{loss}}(\eta_G) \end{bmatrix} \quad (6)$$

where $F_{\text{loss}}, F_{\text{reg}}: [-1, 1] \rightarrow \mathbb{R}$ are functions defined as:

$$F_{\text{loss}}(\eta_G) := \sqrt{c\tilde{c}} \cdot \frac{\mathbb{E}[\text{env}'_{\text{loss}}(Z + \alpha G; \kappa) \cdot \text{env}'_{\widetilde{\text{loss}}}(Z + \tilde{\alpha}\tilde{G}; \tilde{\kappa})]}{\mathbb{E}[\text{env}'_{\text{loss}}(Z + \alpha G; \kappa)^2]^{1/2} \cdot \mathbb{E}[\text{env}'_{\widetilde{\text{loss}}}(Z + \tilde{\alpha}\tilde{G}; \tilde{\kappa})^2]^{1/2}} \quad (7a)$$

$$F_{\text{reg}}(\eta_H) := \frac{\mathbb{E}[(\frac{1}{\nu} \cdot \text{env}'_{\text{reg}}(\Theta + \frac{\beta}{\nu}H; \frac{1}{\nu}) - \frac{\beta}{\nu}H) \cdot (\frac{1}{\tilde{\nu}} \cdot \text{env}'_{\widetilde{\text{reg}}}(\Theta + \frac{\tilde{\beta}}{\tilde{\nu}}\tilde{H}; \frac{1}{\tilde{\nu}}) - \frac{\tilde{\beta}}{\tilde{\nu}}\tilde{H})]}{\mathbb{E}[(\frac{1}{\nu} \cdot \text{env}'_{\text{reg}}(\Theta + \frac{\beta}{\nu}H; \frac{1}{\nu}) - \frac{\beta}{\nu}H)^2]^{1/2} \cdot \mathbb{E}[(\frac{1}{\tilde{\nu}} \cdot \text{env}'_{\widetilde{\text{reg}}}(\Theta + \frac{\tilde{\beta}}{\tilde{\nu}}\tilde{H}; \frac{1}{\tilde{\nu}}) - \frac{\tilde{\beta}}{\tilde{\nu}}\tilde{H})^2]^{1/2}} \quad (7b)$$

where $\begin{pmatrix} G \\ \tilde{G} \end{pmatrix} \sim \mathcal{N}(\begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 1 & \eta_G \\ \eta_G & 1 \end{bmatrix})$, $\begin{pmatrix} H \\ \tilde{H} \end{pmatrix} \sim \mathcal{N}(\begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 1 & \eta_H \\ \eta_H & 1 \end{bmatrix})$, $\Theta \sim F_\theta$, $Z \sim F_z$, all mutually independent.

System 1b is new and one of the main contributions of this paper. Note that the parameters (η_G, η_H) in the system are correlation parameters (up to scaling factors) of the two random variables visible inside expectations in (7a) and of the two random variables in (7b), respectively. By the Cauchy–Schwarz inequality, the function F_{loss} and F_{reg} are uniformly bounded as $|F_{\text{loss}}(\eta)| \leq \sqrt{c\tilde{c}}$ and $|F_{\text{reg}}(\eta)| \leq 1$ so that any solution (η_G, η_H) to the system (6) lies in the set $[-1, 1] \times [-\sqrt{c\tilde{c}}, \sqrt{c\tilde{c}}]$. As stated, it is not immediately clear if System 1b admits any solution and whether it is unique. Our first result establishes that this is indeed the case:

Theorem 1 (Existence, uniqueness, and sign pattern of solutions to System 1b). The functions F_{loss} and F_{reg} defined in (7) satisfy the following properties:

1. $|F_{\text{loss}}(\eta_G)| < \sqrt{c\tilde{c}}$ for all $|\eta_G| < 1$ and $|F_{\text{reg}}(\eta_H)| < 1$ for all $|\eta_H| < 1$.
2. F_{loss} and F_{reg} are non-decreasing, differentiable, and the compositions $F_{\text{loss}} \circ F_{\text{reg}}$ and $F_{\text{reg}} \circ F_{\text{loss}}$ are $\min\{c, \tilde{c}\}$ -Lipschitz.
3. If $\min\{c, \tilde{c}\} < 1$, then System 1b admits a unique solution $(\eta_G^*, \eta_H^*) \in (-1, 1) \times (-\sqrt{c\tilde{c}}, \sqrt{c\tilde{c}})$.

4. The signs of the solution (η_G^*, η_H^*) are characterized by the following sign pattern:

$$\text{sign} \left(\begin{bmatrix} \eta_G^* \\ \eta_H^* \end{bmatrix} \right) = \text{sign} \left(\begin{bmatrix} F_{\text{reg}} \circ F_{\text{loss}}(0) \\ F_{\text{loss}} \circ F_{\text{reg}}(0) \end{bmatrix} \right) \quad (8)$$

where $\text{sign}(x) := \mathbb{1}_{\{x>0\}} - \mathbb{1}_{\{x<0\}}$ applies component-wise.

Some remarks on Theorem 1 are in order. Among the four properties listed in Theorem 1, the most interesting is the second property: the two maps $(F_{\text{loss}} \circ F_{\text{reg}}, F_{\text{reg}} \circ F_{\text{loss}})$ are strict contractions. Given this property, the third property (the uniqueness and existence of the solution) easily follows from the Banach fixed-point theorem. We briefly explain this next. Indeed, if (η_H^*, η_G^*) is a solution to System 1b, then η_H^* automatically satisfies the following 1-dimensional fixed-point equation:

$$\eta_H^* = F_{\text{loss}}(\eta_G^*) = F_{\text{loss}} \circ F_{\text{reg}}(\eta_H^*).$$

The other direction is also true in the following sense: if η_H^* is a solution to the fixed-point equation $\eta_H = F_{\text{loss}} \circ F_{\text{reg}}(\eta_H)$, then letting $\eta_G^* = F_{\text{reg}}(\eta_H^*)$, we observe that the pair (η_H^*, η_G^*) satisfies System 1b. Since $F_{\text{loss}} \circ F_{\text{reg}}$ is a contraction mapping, the Banach fixed-point theorem implies that such η_H^* uniquely exists. See Appendix A.1 for the full proof details. This contraction property also certifies that the fixed-point iteration algorithm $\eta_H^{(k+1)} = F_{\text{loss}} \circ F_{\text{reg}}(\eta_H^{(k)})$, which we use in our experiments to solve System 1b, numerically converges to the correct solution η_H^* exponentially fast.

Since F_{loss} and F_{reg} are non-decreasing, combined with the fourth property in Theorem 1, we get a simple sufficient condition which determines the sign of (η_H^*, η_G^*) :

$$F_{\text{loss}}(0) \text{ and } F_{\text{reg}}(0) \text{ are non-negative} \Rightarrow \eta_H^* \text{ and } \eta_G^* \text{ are non-negative.}$$

Simplifying the denominators of F_{reg} and F_{loss} by (5a) and (5b), we can write $F_{\text{loss}}(0)$ and $F_{\text{reg}}(0)$ as follows:

$$F_{\text{loss}}(0) = \frac{c\tilde{c}\tilde{\delta}}{\beta\tilde{\beta}} \cdot \mathbb{E} \left[\mathbb{E}[\text{env}'_{\text{loss}}(Z + \alpha G; \kappa) \mid Z] \cdot \mathbb{E}[\text{env}'_{\text{loss}}(Z + \tilde{\alpha} G; \tilde{\kappa}) \mid Z] \right],$$

$$F_{\text{reg}}(0) = \frac{1}{\nu\tilde{\nu}\alpha\tilde{\alpha}} \cdot \mathbb{E} \left[\mathbb{E}[\text{env}'_{\text{reg}}(\Theta + \frac{\beta}{\nu} H; \frac{1}{\nu}) \mid \Theta] \cdot \mathbb{E}[\text{env}'_{\text{reg}}(\Theta + \frac{\tilde{\beta}}{\tilde{\nu}} H; \frac{1}{\tilde{\nu}}) \mid \Theta] \right].$$

In particular, if the same loss and regularizer are used, i.e., $\text{loss} = \widetilde{\text{loss}}$ and $\text{reg} = \widetilde{\text{reg}}$, and the subsample sizes are the same, i.e., $k = \tilde{k}$, then the solutions satisfying System 1b are same, i.e., $(\alpha, \beta, \kappa, \nu) = (\tilde{\alpha}, \tilde{\beta}, \tilde{\kappa}, \tilde{\nu})$, so that it is easy to see from the above formula that $F_{\text{loss}}(0) \geq 0$ and $F_{\text{reg}}(0) \geq 0$. This means that the solutions (η_G^*, η_H^*) are non-negative when the same loss, regularizer, and subsample size are used.

Another case such that the solutions (η_G^*, η_H^*) are positive is when loss and $\widetilde{\text{loss}}$ are least squares and reg and $\widetilde{\text{reg}}$ are ridge (but possibly different regularization parameters). This is because $\text{env}'_f(x; \tau)$ is linear in x for any squared loss (and regularizer) of the form $f(\cdot) = \lambda(\cdot)^2$ so that $F_{\text{loss}}(0) = C_1 \mathbb{E}[Z^2] \geq 0$ and $F_{\text{reg}}(0) = C_2 \mathbb{E}[\Theta^2] \geq 0$ for some positive constants C_1, C_2 .

Remark 1 (Negative estimator error correlation). Let us take $\text{reg}, \widetilde{\text{reg}}$ as the indicator functions

$$\text{reg}(x) = \Pi_{(-\infty, -t]}(x) := \begin{cases} +\infty & x > -t \\ 0 & x \leq -t \end{cases} \quad \text{and} \quad \widetilde{\text{reg}}(x) = \Pi_{[\tilde{t}, +\infty)}(x) := \begin{cases} +\infty & x < \tilde{t} \\ 0 & x \geq \tilde{t} \end{cases}$$

where $t, \tilde{t} \geq 0$ are non-negative constants. Note that the two sets, $(-\infty, -t]$ and $[\tilde{t}, +\infty)$, are disjoint. Noting $\text{prox}_{\text{reg}}(x) = \min\{x, -t\}$ and $\text{prox}_{\widetilde{\text{reg}}}(x) = \max\{x, \tilde{t}\}$, we have

$$F_{\text{reg}}(\eta_H) = \frac{1}{\nu\widetilde{\nu}\alpha\widetilde{\alpha}} \mathbb{E}[(\min\{\Theta + \frac{\beta}{\nu}H, -t\} - \Theta) \cdot (\max\{\Theta + \frac{\widetilde{\beta}}{\widetilde{\nu}}\widetilde{H}, \tilde{t}\} - \Theta)].$$

Thus, if Θ is included in the closed set $[-t, \tilde{t}]$ with probability 1, then we have $F_{\text{reg}}(\eta_H) \leq 0$ for all η_H so that $\eta_G^* = F_{\text{reg}}(\eta_H^*)$ is non-positive. This is intuitive as we will show in Theorem 2 that η_G characterizes the limiting behavior of the correlations between two estimators trained on reg and $\widetilde{\text{reg}}$.

We next show that the correlation parameters (η_G, η_H) from System 1b are the limiting correlations between the estimator and residual errors of estimators trained on overlapped samples. To do so, besides Assumptions A–C, we will need mild regularity conditions that the loss and reg in Assumption C need to satisfy in relation to the distribution F_θ and F_z of the signal and noise coordinates in Assumption A.

Assumption D (Regularity conditions). Let $\Theta \sim F_\theta$ and $Z \sim F_z$ be the signal and noise random variables as in Assumption A, and $G \sim \mathcal{N}(0, 1)$, $H \sim \mathcal{N}(0, 1)$, all mutually independent. In addition to Assumption C, the functions $\text{loss} : \mathbb{R} \rightarrow \overline{\mathbb{R}}$ and $\text{reg} : \mathbb{R} \rightarrow \overline{\mathbb{R}}$ satisfy the following:

1. For all $c \in \mathbb{R}$, we have

$$\mathbb{E}[\text{loss}'_+(cG + Z)^2] < +\infty \quad \text{and} \quad \mathbb{E}[\text{reg}'_+(cH + \Theta)^2] < +\infty$$

where we define $f'_+(x) := \sup_{s \in \partial f(x)} |s|$ for any convex function f .

2. $\mathbb{P}(\Theta \neq 0) > 0$.
3. System 1a admits a unique positive solution $(\alpha, \beta, \kappa, \nu) \in \mathbb{R}_{>0}^4$.
4. There exists interval $\mathcal{I} \subset \mathbb{R}$ where loss' is strictly increasing. For each $z \in \mathbb{R}$, the measure $p_Z(z)$ of Z is either a Dirac delta function or it is continuous.

The conditions in Assumption D are similar to those assumed in [TAH18] when characterizing the asymptotics for the non-ensemble case. The main difference is that we do not require the second moment of Θ to be finite. These conditions ensure that the individual estimator and residual error norms converge, i.e., $\|\widehat{\theta} - \theta\|_2^2/p \xrightarrow{\mathbb{P}} \alpha^2$ and $\|\text{loss}'(\mathbf{y} - \mathbf{X}\widehat{\theta})\|_2^2/p \xrightarrow{\mathbb{P}} \beta^2$ hold, where α and β are solutions to System 1a. (See Appendices A.5 and A.6 for proofs of convergences of error vector and loss gradient norm squared under the relaxation of the conditions.)

For the upcoming statement, recall that when used on a vector, the loss and reg functions are assumed to be operated component-wise. In addition, we denote the feature matrix and response vector associated with the “overlapped” dataset (\mathbf{x}_i, y_i) for $i \in I \cap \widetilde{I}$ using $(\mathbf{X}_{I \cap \widetilde{I}}, \mathbf{y}_{I \cap \widetilde{I}})$.

Theorem 2 (Estimator and residual error correlation characterization). Let $\widehat{\theta}_I$ and $\widehat{\theta}_{\widetilde{I}}$ be component estimators (3) trained on subsamples $(\mathbf{X}_I, \mathbf{y}_I)$ and $(\mathbf{X}_{\widetilde{I}}, \mathbf{y}_{\widetilde{I}})$ corresponding to index sets I and \widetilde{I} with parameters $(\text{loss}, \text{reg})$ and $(\widetilde{\text{loss}}, \widetilde{\text{reg}})$. Under Assumptions A–D, as $n, p, k, \widetilde{k} \rightarrow \infty$ with $n/p \rightarrow \delta \in (0, \infty)$, $k/n \rightarrow c \in (0, 1]$ and $\widetilde{k}/n \rightarrow \widetilde{c} \in (0, 1]$ with $\min\{c, \widetilde{c}\} < 1$, we have

$$\begin{aligned} p^{-1}(\widehat{\theta}_I - \theta)^\top (\widehat{\theta}_{\widetilde{I}} - \theta) &\xrightarrow{\mathbb{P}} \eta_G \alpha \widetilde{\alpha} \\ p^{-1} \text{loss}'(\mathbf{y}_{I \cap \widetilde{I}} - \mathbf{X}_{I \cap \widetilde{I}} \widehat{\theta}_I)^\top \widetilde{\text{loss}}'(\mathbf{y}_{I \cap \widetilde{I}} - \mathbf{X}_{I \cap \widetilde{I}} \widehat{\theta}_{\widetilde{I}}) &\xrightarrow{\mathbb{P}} \eta_H \beta \widetilde{\beta}, \end{aligned} \tag{10}$$

where (η_G, η_H) is the solution to System 1b. Furthermore, for any $i \in I \cap \tilde{I}$ and $j \in [p]$, there exists a jointly normal (G_i, \tilde{G}_i) with mean 0, variance 1 and correlation η_G and (H_j, \tilde{H}_j) with mean 0, variance 1 and correlation η_H (as in System 1b) such that the residuals and estimators are jointly approximated as follows:

$$\max_{j \in [p]} \mathbb{E} \left[\mathbb{1} \wedge \left\| \begin{pmatrix} \mathbf{e}_j^\top \hat{\boldsymbol{\theta}}_I \\ \mathbf{e}_j^\top \hat{\boldsymbol{\theta}}_{\tilde{I}} \end{pmatrix} - \begin{pmatrix} \text{prox}_{\text{reg}}(\theta_j + \frac{\beta}{\nu} H_j; \frac{1}{\nu}) \\ \text{prox}_{\tilde{\text{reg}}}(\theta_j + \frac{\tilde{\beta}}{\tilde{\nu}} \tilde{H}_j; \frac{1}{\tilde{\nu}}) \end{pmatrix} \right\|_2^2 \right] = o(1), \quad (11a)$$

$$\max_{i \in [n]} \mathbb{E} \left[\mathbb{1}_{i \in I \cap \tilde{I}} \wedge \left\| \begin{pmatrix} y_i - \mathbf{x}_i^\top \hat{\boldsymbol{\theta}}_I \\ y_i - \mathbf{x}_i^\top \hat{\boldsymbol{\theta}}_{\tilde{I}} \end{pmatrix} - \begin{pmatrix} \text{prox}_{\text{loss}}(z_i + \alpha G_i; \kappa) \\ \text{prox}_{\tilde{\text{loss}}}(z_i + \tilde{\alpha} \tilde{G}_i; \tilde{\kappa}) \end{pmatrix} \right\|_2^2 \right] = o(1). \quad (11b)$$

The proof is given in Appendix A.2.4. Put in words, η_G and η_H from System 1b encode the cosines of the angles between the estimator errors and loss gradient residual errors of estimators $\hat{\boldsymbol{\theta}}_I$ and $\hat{\boldsymbol{\theta}}_{\tilde{I}}$ respectively. It is worth noting that the proximal representations in (11a) and (11b) allow one to provide limiting behavior of other functionals of the estimator and residual errors by assuming further moments on the signal and error distributions F_θ and F_z ; for example, we can characterize the correlation between the raw residuals (rather than after applying the loss derivative) assuming finite second moment of F_z or consider pseudo-Lipschitz functionals other than squared error.

The main difficulty in showing Theorem 2 is the non-trivial dependence between the two estimators $\hat{\boldsymbol{\theta}}_I$ and $\hat{\boldsymbol{\theta}}_{\tilde{I}}$ as they share the samples $\mathbf{X}_{I \cap \tilde{I}}$. In the case of squared loss and ridge regularizer, the estimators $\hat{\boldsymbol{\theta}}_I$ and $\hat{\boldsymbol{\theta}}_{\tilde{I}}$ have closed-form expressions, and prior work in [PDK23] explicitly analyze the overlapped resolvents by developing conditional calculus of resolvents. However, for general loss and regularizers, the overlapped terms are more challenging to analyze due to the lack of closed-form expressions. Our strategy in this paper is to exploit the recently developed technique in [BK24] to analyze the overlapped terms.

To prove Theorem 2, we show that the left-hand side of (10) concentrate around scalars independent of \mathbf{X} , and that these two scalars are approximate fixed-point of $F_{\text{loss}} \circ F_{\text{reg}}$ and $F_{\text{reg}} \circ F_{\text{loss}}$, respectively. This strategy of first proving the concentration of certain quantities and then obtaining approximate fixed-point equations is reminiscent of the leave-one-out analysis of [EKBB⁺13, Kar18] and was previously used in [BK24] to characterize the ensemble risk in unregularized M-estimators (with no explicit penalty). The setting studied here is significantly more complicated than these works due to the presence of robust loss functions, penalty functions, and shared samples between $\hat{\boldsymbol{\theta}}_I$ and $\hat{\boldsymbol{\theta}}_{\tilde{I}}$.

We believe that once the contractions of $F_{\text{loss}} \circ F_{\text{reg}}$ and $F_{\text{reg}} \circ F_{\text{loss}}$ have been found and the existence and uniqueness of the solution to System 1b has been established, different techniques than those used here and discussed in the previous paragraph could also be used to derive asymptotic results similar to Theorem 2. For instance, after existence and uniqueness of the solution to System 1b is established, there is hope to carry out an AMP analysis for matrix-valued parameters (see for instance [LGR⁺22, Lemmas B.3 and B.5] or [JM13, GB23]), or by using the conditional CGMT technique of [CM24, Appendix F]. However, we emphasize that these alternate techniques would also first require to establish the structure of System 1b (as done in Theorem 1) in order to guarantee the existence and uniqueness of the solution to System 1b, since such existence and uniqueness result is required for both applying CGMT results and ensuring the convergence of AMP to the regularized M-estimator.

Even though we assume the existence and uniqueness of the solution to System 1a, such existence

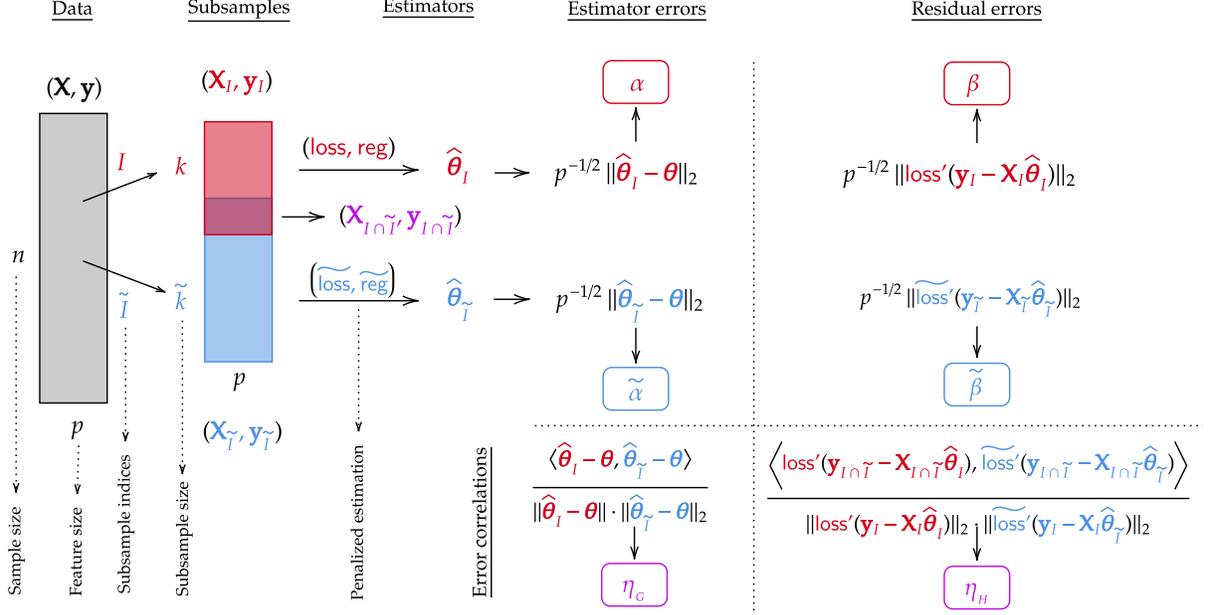


Figure 1: Illustration of the asymptotics of the estimator and residual error norms and correlations of overlapped regularized M-estimator.

and uniqueness have been established under slightly stronger assumptions in [BK23b], namely, for Lipschitz loss and regularizer. We conjecture that the assumptions of Theorem 2 (in particular, Assumption D) are sufficient for existence and uniqueness of a solution to System 1a without requiring a Lipschitz condition on the loss and penalty. We believe that extending the analysis [BK23b] to relax the Lipschitz assumption is a good starting point for this conjecture.

3.2 Interpretation of the parameters in Systems 1a and 1b

As mentioned earlier, the six parameters $(\alpha, \beta, \kappa, \nu)$ and (η_G, η_H) in Systems 1a and 1b essentially characterize the asymptotic risk of the ensemble estimator. These deterministic parameters are limits of various stochastic (observable) quantities that we now describe (see also Figure 1 for a visual illustration).

Here $\hat{\theta}_I$ and $\hat{\theta}_{\tilde{I}}$ are the component estimators (3) trained on subsamples $(\mathbf{X}_I, \mathbf{y}_I)$ and $(\mathbf{X}_{\tilde{I}}, \mathbf{y}_{\tilde{I}})$ corresponding to index sets I and \tilde{I} and parameters $(\text{loss}, \text{reg})$ and $(\widetilde{\text{loss}}, \widetilde{\text{reg}})$, respectively. We further define the scalar df_I and the matrix \mathbf{V}_I by

$$\text{df}_I := \text{tr}[(\partial/\partial \mathbf{y}_I) \mathbf{X}_I \hat{\theta}_I], \quad \mathbf{V}_I := (\partial/\partial \mathbf{y}_I) \text{loss}'(\mathbf{y}_I - \mathbf{X}_I \hat{\theta}_I) \in \mathbb{R}^{|I| \times |I|} \quad (12)$$

and similarly for \tilde{I} . Two scalars of interest, that relates the behavior of $\hat{\theta}_I$ to the scalars (κ, ν) in System 1a, are df_I and $\text{tr}[\mathbf{V}_I]$.

Assuming squared loss, $\text{loss}'(\mathbf{y}_I - \mathbf{X}_I \hat{\theta}_I) = \mathbf{y}_I - \mathbf{X}_I \hat{\theta}_I$ is simply the residual vector, and the matrix \mathbf{V}_I simplifies to $\mathbf{V}_I = \mathbf{I} - (\partial/\partial \mathbf{y}_I) \mathbf{X}_I \hat{\theta}_I$, so that $\text{tr}[\mathbf{V}_I] = n - \text{df}_I$. That is, for the square loss case these quantities can all be related to the usual notion of effective degrees of freedom [Ste81]. The matrix $(\partial/\partial \mathbf{y}_I) \mathbf{X}_I \hat{\theta}_I$ is usually referred to as the “hat” or “smoothing” matrix (for linear smoothers), whose trace is the effective degrees of freedom.

Table 2: **Interpretations of various limiting quantities appearing in Systems 1a and 1b.** See Section 3.2 for definitions and notations.

Interpretation	Stochastic quantity	Limit
Error vector norm squared	$\ \widehat{\boldsymbol{\theta}}_I - \boldsymbol{\theta}\ _2^2/p$	α^2
Loss gradient norm squared	$\ \text{loss}'(\mathbf{y}_I - \mathbf{X}_I \widehat{\boldsymbol{\theta}}_I)\ _2^2/p$	β^2
Inner product of error vectors	$(\widehat{\boldsymbol{\theta}}_I - \boldsymbol{\theta})^\top (\widehat{\boldsymbol{\theta}}_{\bar{I}} - \boldsymbol{\theta})/p$	$\eta_G \alpha \tilde{\alpha}$
Inner product of loss gradients	$\text{loss}'(\mathbf{y}_{I \cap \bar{I}} - \mathbf{X}_{I \cap \bar{I}} \widehat{\boldsymbol{\theta}}_I)^\top \widetilde{\text{loss}}'(\mathbf{y}_{I \cap \bar{I}} - \mathbf{X}_{I \cap \bar{I}} \widehat{\boldsymbol{\theta}}_{\bar{I}})/p$	$\eta_H \beta \tilde{\beta}$
Degrees of freedom	df_I/p where $\text{df}_I := \text{tr}[(\partial/\partial \mathbf{y}_I) \mathbf{X}_I \widehat{\boldsymbol{\theta}}_I]$	$\nu \kappa$
Residual degrees of freedom	$\text{tr}[\mathbf{V}_I]/p$ where $\mathbf{V}_I := (\partial/\partial \mathbf{y}_I) \text{loss}'(\mathbf{y}_I - \mathbf{X}_I \widehat{\boldsymbol{\theta}}_I)$	ν
Generalized resolvent trace	$\text{tr}[(\mathbf{X}_I^\top \text{diag}[\text{loss}''(\mathbf{y}_I - \mathbf{X}_I \widehat{\boldsymbol{\theta}}_I)] \mathbf{X}_I + \text{diag}[\text{reg}''(\widehat{\boldsymbol{\theta}}_I)])^{-1}]$ if reg is twice differentiable and $\text{df}_I/\text{tr}[\mathbf{V}_I]$ if reg is non-smooth	κ

If loss is no the squared loss, but loss' is 1-Lipschitz (as in the Huber loss or several robust regression losses), the quantities $\text{loss}'(\mathbf{y}_I - \mathbf{X}_I \widehat{\boldsymbol{\theta}}_I)$ is still related to a notion of residual vector, and $\text{tr}[\mathbf{V}_I]$ is still related to a notion of degrees of freedom. By [Bel23, Lemma 9.1], the estimator $\widehat{\boldsymbol{\theta}}_I$ is the first part of a solution $(\widehat{\boldsymbol{\theta}}_I, \widehat{\mathbf{u}})$ to the convex optimization problem

$$\min_{\mathbf{b} \in \mathbb{R}^p, \mathbf{u} \in \mathbb{R}^{|I|}} \|\mathbf{y}_I - \mathbf{X}_I \mathbf{b} - \mathbf{u}\|_2^2 + \sum_{j=1}^p \text{reg}(b_j) + \sum_{i \in I} h(u_i),$$

where $h : \mathbb{R} \rightarrow \mathbb{R}$ is a deterministic convex function related to loss . The interpretation of this new optimization problem is that in the presence of heavy-tailed errors or outliers in some components of \mathbf{y}_I , we add additional variables $(u_i)_{i \in I}$ to fit those outliers. As an example, for the Huber loss, $h(\cdot)$ is proportional to the absolute value. The solution satisfies $\text{loss}'(\mathbf{y}_I - \mathbf{X}_I \widehat{\boldsymbol{\theta}}_I) = \mathbf{y}_I - \mathbf{X}_I \widehat{\boldsymbol{\theta}}_I - \widehat{\mathbf{u}} = \mathbf{y}_I - [\mathbf{X} \mid \mathbf{I}_I](\widehat{\boldsymbol{\theta}}_I^\top \mid \widehat{\mathbf{u}}^\top)^\top$. That is, $\text{loss}'(\mathbf{y}_I - \mathbf{X}_I \widehat{\boldsymbol{\theta}}_I)$ is the residual vector of the optimization problem with enlarged design matrix $[\mathbf{X}_I \mid \mathbf{I}_I] \in \mathbb{R}^{|I| \times (|I|+p)}$. Consequently, $\text{tr}[\mathbf{V}_I]$ equals $|I|$ minus the effective degrees-of-freedom of the estimate $(\widehat{\boldsymbol{\theta}}_I^\top, \widehat{\mathbf{u}}^\top)$ fitted using this enlarged design matrix. With this in mind, we refer in Table 2 to $\text{tr}[\mathbf{V}_I]$ as residual degrees of freedom in general, and robust residual degrees of freedom for the special case of the Huber loss.

Another interpretation of the matrix \mathbf{V}_I in (12) is the Hessian, with respect to \mathbf{y}_I , of the objective value (3) at $\widehat{\boldsymbol{\theta}}_I$. More precisely, with

$$F(\mathbf{y}_I) = \sum_{i \in I} \text{loss}(y_i - \mathbf{x}_i^\top \widehat{\boldsymbol{\theta}}_I) + \sum_{j \in [p]} \text{reg}((\widehat{\boldsymbol{\theta}}_I)_j)$$

being the objective value at the minimizer, the envelope theorem gives $(\partial/\partial y_i)F(\mathbf{y}_I) = \text{loss}'(y_i - \mathbf{x}_i^\top \widehat{\boldsymbol{\theta}}_I)$. Differentiating once more reveals that \mathbf{V}_I in (12) is the Hessian of $F(\cdot)$ at \mathbf{y}_I and $\text{tr}[\mathbf{V}_I]$ is the Laplacian. Since partial minimization preserves convexity and the objective function in (3) is jointly convex in (\mathbf{b}, \mathbf{y}) , the function $F(\cdot)$ is convex. This interpretation explains why \mathbf{V}_I is a positive semi-definite matrix in cases where closed form expressions for \mathbf{V}_I are available (see Table 5).

The convergence of the estimator and residual error norms (first two rows of Table 2) is proved in in [TAH18, CMW23, LGC⁺21] using the CGMT. Convergence of the corresponding inner products (third, fourth row of Table 2) is novel and established in Theorem 2. The convergence

$$\text{tr}[\mathbf{V}_I]/p \xrightarrow{P} \nu, \quad \text{df}_I/p \xrightarrow{P} \nu \kappa, \quad \text{df}_I/\text{tr}[\mathbf{V}_I] \xrightarrow{P} \kappa, \quad (13)$$

was so far only known for the lasso [CMW23, Theorem 8] or the square loss [Bel23, Corollary 3.2]. To our knowledge, the present paper is the first to establish the above convergence in probability for regularized estimators and robust loss functions beyond the square loss. The proof is given in Appendix A.4.

Assuming twice differentiable loss and reg functions, the parameter κ is also the limiting trace of a resolvent-like matrix $\mathbf{A}_I := (\mathbf{X}_I^\top \text{diag}[\text{loss}''(\mathbf{y}_I - \mathbf{X}_I \hat{\boldsymbol{\theta}}_I)] \mathbf{X}_I + \text{diag}[\text{reg}''(\hat{\boldsymbol{\theta}}_I)])^{-1}$, which in the further special case of square loss and squared regularizer (with regularization level λ) simplifies to the standard ridge resolvent: $\mathbf{A}_I = (\mathbf{X}_I^\top \mathbf{X}_I + \lambda \mathbf{I})^{-1}$. We refer to \mathbf{A}_I as the generalized resolvent for convenience in Table 2.

3.3 Asymptotics of ensemble risk

Using the parameters in Systems 1a and 1b, we are now ready for our main result on the squared risk asymptotics of the ensemble estimator. Observe that the squared risk of the ensemble estimator $\tilde{\boldsymbol{\theta}}_M = \frac{1}{M} \sum_{m \in [M]} \hat{\boldsymbol{\theta}}_m$ can be decomposed into two terms:

$$\frac{1}{p} \|\tilde{\boldsymbol{\theta}}_M - \boldsymbol{\theta}\|_2^2 = \frac{1}{M^2} \sum_{m \in [M]} \frac{1}{p} \|\hat{\boldsymbol{\theta}}_m - \boldsymbol{\theta}\|_2^2 + \frac{1}{M^2} \sum_{\substack{m, \ell \in [M] \\ m \neq \ell}} \frac{1}{p} (\hat{\boldsymbol{\theta}}_m - \boldsymbol{\theta})^\top (\hat{\boldsymbol{\theta}}_\ell - \boldsymbol{\theta}).$$

Noting $\|\hat{\boldsymbol{\theta}}_m - \boldsymbol{\theta}\|_2^2/p \xrightarrow{P} \alpha_m^2$ and applying Theorem 2 to the cross term $(\hat{\boldsymbol{\theta}}_m - \boldsymbol{\theta})^\top (\hat{\boldsymbol{\theta}}_\ell - \boldsymbol{\theta})/p$ for each $m \neq \ell$, we arrive at the following result:

Corollary 3 (General ensemble risk characterization). Suppose the assumptions of Theorem 2 hold. For $m \in [M]$, let α_m be the parameter satisfying System 1a. For $m, \ell \in [M]$, let $\eta_G(m, \ell)$ be the parameter satisfying System 1b. Then, as $n, p, k \rightarrow \infty$ with $n/p \rightarrow \delta \in (0, \infty)$ and $k_m/n \rightarrow c_m \in (0, 1]$, we have

$$\frac{1}{p} \|\tilde{\boldsymbol{\theta}}_M - \boldsymbol{\theta}\|_2^2 \xrightarrow{P} \mathcal{R}_M := \frac{1}{M^2} \sum_{m \in [M]} \alpha_m^2 + \frac{1}{M^2} \sum_{\substack{m, \ell \in [M] \\ m \neq \ell}} \eta_G(m, \ell) \cdot \alpha_m \alpha_\ell. \quad (14)$$

Since the parameters α_m and $\eta_G(m, \ell)$ implicitly depend on δ and c_m , the asymptotic risk \mathcal{R}_M also implicitly depends on these parameters. For brevity, we will simply write \mathcal{R}_M unless we wish to explicitly point out this dependence. The factor η_G captures the benefit of ensembling. It is worth noting that a negative η_G (which intuitively corresponds to a component that does better in a different direction) will improve the ensemble risk if the components themselves also have small risks. This aligns with the higher level intuition in ensembling that one should ensemble predictors that each does well, preferably on different parts of the input space.

3.4 Risk estimation

The risk characterization in Corollary 3 depends on the population-level characteristics (such as the signal and noise distributions F_θ and F_z) and provides useful theoretical insights into the risk behavior of the ensemble estimator in terms of these quantities. In practical applications, however, the statistician needs to estimate the risk accurately to tune ensemble hyperparameters effectively using the observed data (\mathbf{X}, \mathbf{y}) . These hyperparameters include the choice of component estimators

(through `loss` and `reg`), their level of regularization (regularization level for `reg`), the subsample sizes (k), and the ensemble size (M). For this purpose, we next construct a data-dependent proxy for the squared risk, which one can then tune with respect to various hyperparameters.

Definition 1 (Risk estimator component). Let $\widehat{\boldsymbol{\theta}}_I$ and $\widehat{\boldsymbol{\theta}}_{\widetilde{I}}$ be the component estimators (3) trained on subsamples $(\mathbf{X}_I, \mathbf{y}_I)$ and $(\mathbf{X}_{\widetilde{I}}, \mathbf{y}_{\widetilde{I}})$ corresponding to index sets I and \widetilde{I} with parameters $(\text{loss}, \text{reg})$ and $(\widetilde{\text{loss}}, \widetilde{\text{reg}})$. Corresponding to estimators $\widehat{\boldsymbol{\theta}}_I$ and $\widehat{\boldsymbol{\theta}}_{\widetilde{I}}$, define:

1. Degrees of freedom: $\text{df}_I = \text{tr}[(\partial/\partial \mathbf{y}_I) \mathbf{X}_I \widehat{\boldsymbol{\theta}}_I]$ and $\text{df}_{\widetilde{I}} = \text{tr}[(\partial/\partial \mathbf{y}_{\widetilde{I}}) \mathbf{X}_{\widetilde{I}} \widehat{\boldsymbol{\theta}}_{\widetilde{I}}]$.
2. Residual errors: $\mathbf{r} = \mathbf{y}_I - \mathbf{X}_I \widehat{\boldsymbol{\theta}}_I$ and $\widetilde{\mathbf{r}} = \mathbf{y}_{\widetilde{I}} - \mathbf{X}_{\widetilde{I}} \widehat{\boldsymbol{\theta}}_{\widetilde{I}}$.
3. Residual degrees of freedom: traces of $\mathbf{V}_I = (\partial/\partial \mathbf{y}_I) \text{loss}'(\mathbf{y}_I - \mathbf{X}_I \widehat{\boldsymbol{\theta}}_I)$ and $\mathbf{V}_{\widetilde{I}} = (\partial/\partial \mathbf{y}_{\widetilde{I}}) \widetilde{\text{loss}}'(\mathbf{y}_{\widetilde{I}} - \mathbf{X}_{\widetilde{I}} \widehat{\boldsymbol{\theta}}_{\widetilde{I}})$.

Using these quantities, define an observable quantity $\text{EST}_{I, \widetilde{I}}$ as follows:

$$\text{EST}_{I, \widetilde{I}} := \frac{1}{n} \sum_{i \in [n]} \left(r_i + \mathbb{1}_{\{i \in I\}} \frac{\text{df}_I}{\text{tr}[\mathbf{V}_I]} \text{loss}'(r_i) \right) \left(\widetilde{r}_i + \mathbb{1}_{\{i \in \widetilde{I}\}} \frac{\text{df}_{\widetilde{I}}}{\text{tr}[\mathbf{V}_{\widetilde{I}}]} \widetilde{\text{loss}}'(\widetilde{r}_i) \right) \quad (15)$$

where $\mathbb{1}_\Omega$ denotes the indicator function associated with event Ω .

The quantities $\text{tr}[\mathbf{V}_I]$ and df_I have explicit closed-form expressions for special choices of `loss` and `reg`. Some of these are summarized in Table 5. We show next that $\text{EST}_{I, \widetilde{I}}$ approximates well the component of prediction risk corresponding to the inner product of estimator errors of $\widehat{\boldsymbol{\theta}}_I$ and $\widehat{\boldsymbol{\theta}}_{\widetilde{I}}$. We then naturally construct a risk estimator for the prediction risk of the ensemble estimator.

Theorem 4 (Consistency of risk estimator component). In addition to Assumptions A–D, assume that `reg` and $\widetilde{\text{reg}}$ are strongly convex. Then we have

$$\frac{1}{p} (\widehat{\boldsymbol{\theta}}_I - \boldsymbol{\theta})^\top (\widehat{\boldsymbol{\theta}}_{\widetilde{I}} - \boldsymbol{\theta}) + \frac{\|\mathbf{z}\|_2^2}{n} = \text{EST}_{I, \widetilde{I}} + \mathcal{O}_{\mathbb{P}}(n^{-1/2}) \left(1 + \frac{\|\mathbf{z}\|_2}{\sqrt{n}} \right).$$

The risk estimator $\text{EST}_{I, \widetilde{I}}$ is a generalization of the criterion originally proposed by [BS22] for non-ensemble regularized M-estimator. Although in this paper we focus on the isotropic Gaussian design $\boldsymbol{\Sigma} = \mathbf{I}_p$, the same argument in the proof of Theorem 4 works in the anisotropic design $\boldsymbol{\Sigma} \neq \mathbf{I}_p$. As a result, we can show that the $\text{EST}_{I, \widetilde{I}}$ (without any modification) approximates $p^{-1} (\widehat{\boldsymbol{\theta}}_I - \boldsymbol{\theta})^\top \boldsymbol{\Sigma} (\widehat{\boldsymbol{\theta}}_{\widetilde{I}} - \boldsymbol{\theta}) + n^{-1} \|\mathbf{z}\|_2^2$ under the event that $(\text{tr}[\mathbf{V}_I], \text{tr}[\mathbf{V}_{\widetilde{I}}])$ are bounded from below by a positive constant as in [BS22, Theorem 5.3].

We believe that the strongly convexity assumption on `reg` in Theorem 4 is an artifact of our proof (see Figure 2 for an illustration where `reg` is not strongly convex). Note that this type of assumption of strong convexity has already been assumed in [BS22]. This assumption guarantees for free that the coefficient $\text{df}/\text{tr}[\mathbf{V}]$ does not blow up. However, we emphasize that Theorem 2 and Corollary 3 for risk characterization do not require the strong convexity assumption. The approximation argument (see Appendix A.2.1) used to prove Theorem 2 and Corollary 3 in the non-strongly convex case is again applicable in the context of Theorem 4, although it is not currently sufficient to conclude a version Theorem 4 for non-strongly convex regularizers due to the difficulty of establishing the continuity of df_I and $\text{tr}[\mathbf{V}_I]$ with respect to the perturbation parameter μ as $\mu \rightarrow 0$ in (28).

Equipped with the component risk estimator (15), we can now construct a consistent risk estimator for the ensemble estimator (4):

Corollary 5 (General ensemble risk estimation). Fix $M \geq 1$ and consider the ensemble estimator $\tilde{\boldsymbol{\theta}}_M = \frac{1}{M} \sum_{m \in [M]} \tilde{\boldsymbol{\theta}}_m$ where the component estimator $\tilde{\boldsymbol{\theta}}_m$ is trained with $(\text{loss}_m, \text{reg}_m)$ on subsample I_m for $m \in [M]$ as in (3). Define an estimator EST for the squared prediction risk:

$$\text{EST} := \frac{1}{M^2} \sum_{m, \ell \in [M]} \text{EST}_{m, \ell}$$

where $\text{EST}_{m, \ell}$ is $\text{EST}_{I, \tilde{I}}$ as defined in (15) with $(\text{loss}, \text{reg}, I) = (\text{loss}_m, \text{reg}_m, I_m)$ and $(\widetilde{\text{loss}}, \widetilde{\text{reg}}, \tilde{I}) = (\text{loss}_\ell, \text{reg}_\ell, I_\ell)$. Then we have

$$\frac{1}{p} \|\tilde{\boldsymbol{\theta}}_M - \boldsymbol{\theta}\|_2^2 + \frac{\|\mathbf{z}\|_2^2}{n} = \text{EST} + \mathcal{O}_{\mathbb{P}}(n^{-1/2}) \left(1 + \frac{\|\mathbf{z}\|_2}{\sqrt{n}} \right). \quad (16)$$

If the noise distribution has enough moments, the guarantee (16) implies that EST approximates the (full) prediction risk (that includes the irreducible error) of the ensemble estimator $\tilde{\boldsymbol{\theta}}_M$ under Assumption A:

$$\mathbb{E}[(y_0 - \mathbf{x}_0^\top \tilde{\boldsymbol{\theta}}_M)^2 \mid \mathbf{y}, \mathbf{X}, \{I_m\}_{m \in [M]}] = \frac{1}{p} \|\tilde{\boldsymbol{\theta}}_M - \boldsymbol{\theta}\|_2^2 + \mathbb{E}[Z^2],$$

where $(y_0, \mathbf{x}_0) \in \mathbb{R} \times \mathbb{R}^p$ is an independent test point sampled from the same distribution as the training data (\mathbf{y}, \mathbf{X}) . More precisely, if the noise has a finite second moment, then EST is consistent for the prediction risk. Furthermore, if the noise distribution has a finite fourth moment, by the central limit theorem (on the terms involving noise averages), EST is \sqrt{n} -consistent for the prediction risk. This rate is not improvable because, for the single ordinary least squares (OLS) estimator (with $\text{loss}(x) = x^2/2$, $\text{reg} = 0$, $|I| = |\tilde{I}| = n$) and $F_z = \mathcal{N}(0, 1)$, the risk estimator gives $\text{EST} = \frac{1}{n} \frac{\|\mathbf{y} - \mathbf{X} \hat{\boldsymbol{\theta}}_{\text{ols}}\|_2^2}{(1-p/n)^2} \stackrel{d}{=} \frac{\chi_{n-p}^2}{n(1-p/n)^2}$ and the standard deviation of the χ_{n-p}^2 incurs an unavoidable term of order $n^{-1/2}$.

If the noise distribution F_z does not have a finite second moment but has a finite $(1+\epsilon)$ -moment for $\epsilon \in [0, 1)$, then even if the estimator EST may not track the prediction risk (because the prediction risk does not necessarily converge), minimizing EST is approximately equivalent to minimizing the excess squared risk $p^{-1} \|\tilde{\boldsymbol{\theta}}_M - \boldsymbol{\theta}\|_2^2$ (which does converge). This is because the moment assumption implies $\sum_{i \in [n]} z_i^2 = \mathcal{O}_{\mathbb{P}}(n^{2/(1+\epsilon)})$ (cf. [hp24]) so that (16) yields

$$\frac{1}{p} \|\tilde{\boldsymbol{\theta}}_M - \boldsymbol{\theta}\|_2^2 = \text{EST} - \frac{\|\mathbf{z}\|_2^2}{n} + \mathcal{O}_{\mathbb{P}}(n^{-\frac{\epsilon}{1+\epsilon}}),$$

where the subtraction term $\frac{\|\mathbf{z}\|_2^2}{n}$ is independent of hyperparameters $(\text{loss}_m, \text{reg}_m, I_m)_{m \in [M]}$. We illustrate this in Figure 2 with noise following Student's t_2 distribution (that does not have a finite second moment).

4 Homogeneous ensembles

While Corollary 3 applies for a generic heterogeneous ensemble (where $\text{reg}_m, \text{loss}_m, |I_m|$ are allowed to differ for distinct m), concrete theoretical insights can be obtained for the homogeneous case

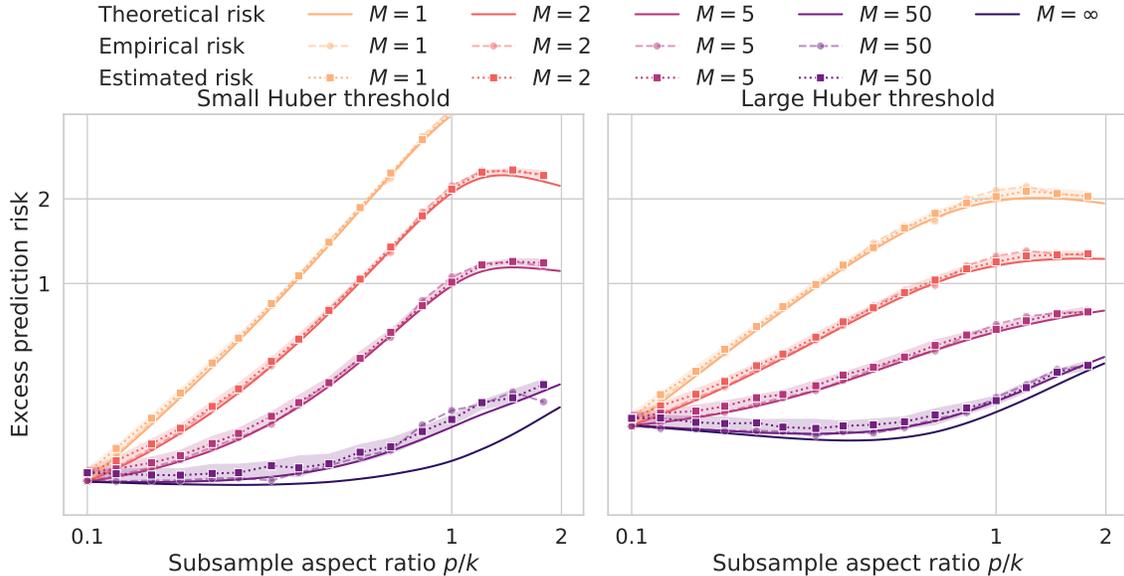


Figure 2: Risk of ℓ_1 -regularized Huber ensemble at different subsample aspect ratios p/k with ℓ_1 -regularization parameter $\lambda = 0.2$ and varying ensemble size M in the underparameterized regime when $p/n = 0.1$ and $n = 5000$. The solid lines represent the theoretical risks, the dashed lines represent the empirical risks averaged over 50 simulations, and the shaded regions represent the standard errors. The data model is given by Appendix D.4.2 where the noise follows Student’s t distribution t_2 . *Left*: Huber threshold parameter 1. *Right*: Huber threshold parameter 5.

(where $\text{reg}_m, \text{loss}_m, |I_m|$ are the same for every m). In the general heterogeneous case, the effect of increasing ensemble size M is not straightforward, as in general we only have:

$$\min_{m \in [M]} \{\alpha_m^2\} \stackrel{?}{\leq} \frac{1}{M^2} \sum_{m \in [M]} \alpha_m^2 + \frac{1}{M^2} \sum_{\substack{m, \ell \in [M] \\ m \neq \ell}} \eta_G(m, \ell) \cdot \alpha_m \alpha_\ell \leq \max_{m \in [M]} \{\alpha_m^2\}.$$

In other words, we will do no worse than the worst component but may not do better than the best component. More ensembles may or may not improve performance depending on the risks of individual estimators.⁵ In contrast, for homogeneous ensembles, we show in the next subsection that increasing ensemble size does indeed help reduce the risk. This uniformity allows for more concrete analytical results and practical insights.

4.1 Risk properties

For the homogeneous ensemble of component estimators trained with the same loss , reg , and subsample size k , as $n, p, k \rightarrow \infty$ with $n/p \rightarrow \delta \in (0, \infty)$ and $k/n \rightarrow c \in (0, 1]$, the limiting risk (14) is given by:

$$\mathcal{R}_M = \frac{1}{M} \mathcal{R}_1 + \left(1 - \frac{1}{M}\right) \mathcal{R}_\infty \quad \text{where} \quad \begin{cases} \mathcal{R}_1 := \alpha^2 & (\text{non-ensemble risk}), \\ \mathcal{R}_\infty := \eta_G \alpha^2 & (\text{full-ensemble risk}). \end{cases} \quad (17)$$

⁵Even if a predictor complements other predictors (in the sense that it has small or negative η_G with the other predictors in the ensemble), it is only “beneficial” for the ensemble if it also has a small risk (α_m) itself.

Observe that the limit \mathcal{R}_M is simply a convex combination of the asymptotic risk of the single estimator \mathcal{R}_1 and of the full-ensemble estimator \mathcal{R}_∞ .⁶ (Note that when $M \rightarrow \infty$, \mathcal{R}_M does indeed converge to \mathcal{R}_∞ , justifying the notation for the limit for the case when $m \neq \ell$.) As a sanity check, note that in the special case when $c = 1$, the above setting corresponds to the non-ensemble case discussed in [TAH18, Section 4].

The following result shows the advantage of ensembling. In the classical bagging and subbagging literature, it is well known that the risk of the ensemble estimator decreases as the ensemble size increases, due to the reduction in variance that comes with having more predictors in the ensemble. In the proportional asymptotic regime that we study in this paper, since subbagging also introduces bias, this is not immediate. A general result along these lines that verifies the monotonicity of the risk of the ensemble itself (not the asymptotic limit) follows from Proposition 3.1 of [PDK23]; see Equation (10) of [PDK23]. Below we verify that the asymptotic risk is strictly monotonic in M by showing that $\eta_G < 1$ in general. This implies that the asymptotic risk is strictly decreasing in M .

Proposition 6 (Improvement due to ensembling). Fix the subsample ratio $c = k/n \in (0, 1)$ and let \mathcal{R}_M be the limiting risk as defined in (17). Then \mathcal{R}_M is strictly decreasing in the number of ensembles M , i.e., $\mathcal{R}_{M+1} < \mathcal{R}_M$ for all $M \in \mathbb{N}$.

The proof follows immediately because of the form of the ensemble risk (17) and the fact that $\eta_G < 1$ when $c < 1$, which follows from Theorem 1 with $c = \tilde{c}$. This monotonicity in the ensemble size M is illustrated by Figure 2 for the ensemble of ℓ_1 -regularized Huber regression⁷.

Because the risk decreases in M , the optimal ensemble size is $M = \infty$.⁸ However, it may not be feasible to use an ensemble size of $M = \infty$. In practice, it suffices to use a large enough M that gives a suboptimal risk close to the full-ensemble risk. For this purpose, a natural idea is to estimate the risk of non-ensemble estimator $M = 1$ and the full estimator $M = \infty$, and obtain an estimate for the risk of M -ensemble using the relationship in (17). This is very similar to the extrapolated cross-validation estimator (ECV) of [DPRK24] which estimates the risk of $M = 1$ and $M = 2$.

We also show that for any ensemble size M , when the subsample ratio c is optimized, the resulting risk decreases in the inverse data aspect ratio $\delta = \lim n/p$. To prepare for the forthcoming statement, let us write the limiting risk \mathcal{R}_M in (17) by $\mathcal{R}_M(\delta, c)$ to make the dependence on the limit $\delta = \lim n/p$ and subsample ratio $c = k/n$ clear. With this notation, we can say the following about the optimally subsampled risk.

Proposition 7 (Monotonicity of optimal subsample risk). The map $\delta \mapsto \inf_{c \in (0,1]} \mathcal{R}_M(\delta, c)$ is non-increasing over $\delta \in (0, \infty)$ for all $M \in \mathbb{N}$.

A consequence of this proposition is that the risk of the optimal ensemble estimator is decreasing in the sample size n for a fixed (large enough) feature size p . Moreover, combined with Proposition 6,

⁶The reason we call this the “full” ensemble is that the ensemble estimator $\hat{\theta}_M$ when $M \rightarrow \infty$ is almost surely equal (coordinate-wise and conditioned on the data) to an ensemble estimator fitted on all possible (and distinct) $\binom{n}{k}$ subsamples of size k ; see Lemma A.1 of [DPK23] for a precise statement and proof.

⁷The Huber loss is defined as $\text{loss}(x) = \frac{x^2}{2\rho} \mathbb{1}(|x| \leq \rho) + (|x| - \frac{\rho}{2}) \mathbb{1}(|x| > \rho)$, where ρ is a positive constant, referred to as the Huber parameter.

⁸In practice, setting $M = \binom{n}{k}$ suffices by only averaging over estimators trained on distinct subsamples (see Appendix A.1 of [DPK23] for more details), but this still can be quite large.

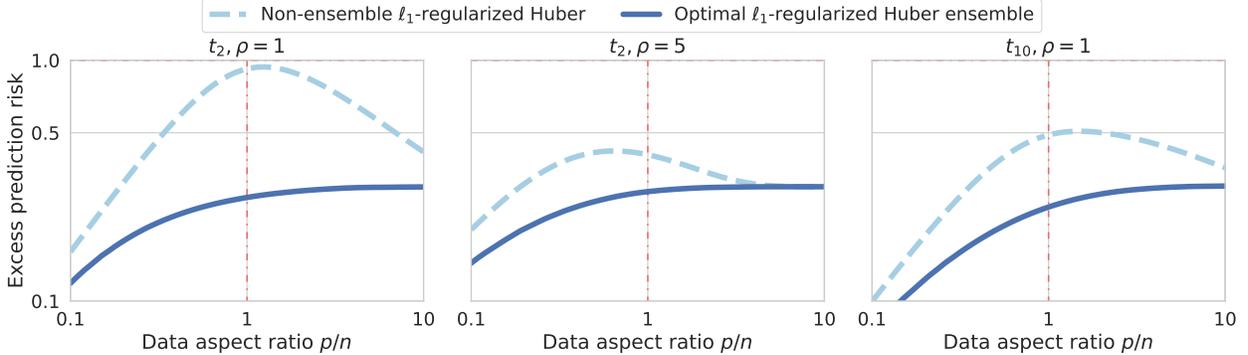


Figure 3: **Optimal subsample risk of the Huber lasso ensemble is monotonic in the data aspect ratio.** Risks of the ℓ_1 -regularized Huber and optimal ℓ_1 -regularized Huber ensemble, for fixed ℓ_1 -regularization parameter $\lambda = 0.5$ and varying Huber parameter ρ , at different data aspect ratios p/n ranging from 0.1 to 10. The data model is given in Appendix D.4.2. *Left:* noise follows Student’s t distribution t_2 and Huber parameter $\rho = 1$. *Middle:* noise follows Student’s t distribution t_2 and Huber parameter $\rho = 5$. *Right:* noise follows Student’s t distribution t_{10} and Huber parameter $\rho = 1$.

we also have that this monotonic risk profile lies below the function $\mathcal{R}_1(\delta, 1)$, the risk profile of the original predictor trained once on the full dataset (\mathbf{X}, \mathbf{y}) with no ensembling (the risk of which, as we discussed above, can be non-monotonic). Such monotonicity in the inverse data aspect ratio is important because it ensures that increasing the amount of data relative to features consistently improves the estimator’s performance. In a sense, a monotonic decrease in risk with the optimal subsample ratio certifies that the model effectively utilizes all the additional data, leading to better performance as the data size grows. This result is illustrated in Figure 3 for the ℓ_1 -regularized Huber regression.

4.2 Examples and connections to literature

In this section, we specialize the general risk characterization in Corollary 3 in several examples of interest. Throughout this section, we will consider ensembles of component predictors trained on the same loss, reg (with a tuning parameter λ), and subsample size k . We begin by considering convex regularized least squares, with further specialization to bridge, ridge, and lasso. We then consider general ridge regularized estimators allowing for non-squared loss in Section 4.2.2 with further specialization to the Huber loss. For the reader’s convenience, the proximal operators, their derivatives, Moreau envelopes, and their derivatives for the ridge and lasso regularization and Huber loss functions are recalled in Table 6.

4.2.1 Ensembles of regularized least squares

In this section, we consider subbagging regularized least squares. Given a common regularizer reg and a regularization parameter $\lambda > 0$, the component estimators are given by:

$$\hat{\boldsymbol{\theta}}_m := \underset{\mathbf{b} \in \mathbb{R}^p}{\text{argmin}} \frac{1}{2} \|\mathbf{y}_{I_m} - \mathbf{X}_{I_m} \mathbf{b}\|_2^2 + \lambda \sum_{j \in [p]} \text{reg}(b_j).$$

Now Assumption D-(1) translates to having a bounded second moment of the noise distribution F_z , which we denote by σ^2 . Using the explicit formula $\text{prox}_{\text{loss}}(x; \tau) = x/(1 + \tau)$ for $\text{loss}(x) = x^2/2$ and performing some change of variables, the risk convergence in (17) holds with $(\mathcal{R}_1, \mathcal{R}_\infty)$ given

by:

$$\mathcal{R}_1 = \tau^2 - \sigma^2 \quad \text{and} \quad \mathcal{R}_\infty = \xi^2 - \sigma^2,$$

where τ and ξ are the solutions to the following systems:

System 2 (Ensembles of regularized least square). Given $\lambda \in (0, \infty)$, $\delta \in (0, \infty)$, $c \in (0, 1]$, $\sigma^2 \in [0, \infty)$, define the following 2-scalar system of equations in variables (τ, a) :

$$\tau^2 = \mathbb{E}[(\text{prox}_{\text{reg}}(\Theta + \frac{\tau}{\sqrt{c\delta}}H; \frac{a\tau}{\sqrt{c\delta}}) - \Theta)^2] + \sigma^2, \quad (18a)$$

$$\frac{\lambda}{\sqrt{c\delta}} = a\tau(1 - \frac{1}{c\delta}\mathbb{E}[\text{prox}'_{\text{reg}}(\Theta + \frac{\tau}{\sqrt{c\delta}}H; \frac{a\tau}{\sqrt{c\delta}})]), \quad (18b)$$

where $H \sim \mathcal{N}(0, 1)$ and $\Theta \sim F_\theta$ are independent. Given $(\tau, a) \in \mathbb{R}_{>0}^2$ that satisfy the above systems, define the following 1-scalar system of equations in variable ξ :

$$\xi^2 = \mathbb{E}[(\text{prox}_{\text{reg}}(\Theta + \frac{\tau}{\sqrt{c\delta}}H; \frac{a\tau}{\sqrt{c\delta}}) - \Theta) \cdot (\text{prox}_{\text{reg}}(\Theta + \frac{\tau}{\sqrt{c\delta}}\tilde{H}; \frac{a\tau}{\sqrt{c\delta}}) - \Theta)] + \sigma^2 \quad (19a)$$

where $\begin{pmatrix} H \\ \tilde{H} \end{pmatrix} \sim \mathcal{N}(\begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{bmatrix} 1 & \eta_H \\ \eta_H & 1 \end{bmatrix})$ with $\eta_H = c\frac{\xi^2}{\tau^2}$, and $\Theta \sim F_\theta$ independent.

Compared to Systems 1a and 1b, the parameterization in System 2 is slightly different. This is done to match the result with some of the existing results for regularized least squares (for $M = 1$). The invertible transformations are given by: $a = \frac{\lambda}{\beta}$, $\tau = \sqrt{c\delta}\frac{\beta}{\nu}$, and $\xi^2 = \eta_C\alpha^2 + \sigma^2$. These parameters also have interpretations as in Section 3.2, summarized below.

Remark 2 (Interpretation of parameters in System 2). The parameter τ^2 is simply the (full) prediction risk of the non-ensemble estimator in the limit, which is $\alpha^2 + \sigma^2$. Note that $\tau^2 \geq \sigma^2$ and is sometimes referred to as effective ‘‘inflated’’ noise variance due to the high dimensionality [DMM09, BM11b]. Moreover, the fact that $\tau^2 = \frac{(c\delta)\beta^2}{\nu^2} = \frac{\beta^2/(c\delta)}{\nu^2/(c\delta)^2}$ is also at the core of consistency of generalized cross-validation (discussed in Section 3.4) for the non-ensemble estimator. Here, the numerator is the asymptotic training error and the denominator is the asymptotic degrees of freedom correction.⁹ The parameter $a\tau$ is the effective threshold parameter at which one applies the proximal operator to the noise-inflated effective observation and appears in approximate message passing (AMP) formulations [DMM09]. The parameter a serves as a proportionality constant between the effective threshold $a\tau$ and standard deviation τ of the inflated effective noise. Finally, the parameter ξ^2 , which is the main contribution of this paper, is the full-ensemble predictor risk (when $M \rightarrow \infty$), which is also $\eta_C\alpha^2 + \sigma^2$.

In the following, we isolate some special cases of regularized M-estimators to compare with existing work.

Remark 3 (Bridge ensembles). Bridge estimators are also known as ℓ_q -regularized least squares and are a popular class of regularized M-estimator [FF93, Fu98]. For ℓ_q regularizer $\text{reg}_q(x) = |x|^q$, the risk of the bridge for $M = 1$ is derived in [WMZ18, WWM20] for general $q \in [1, 2]$. For instance, we recover [WMZ18, Theorem 2.1] by changing of variables $(\lambda', \Theta') = (\lambda/\sqrt{c\delta}, \sqrt{c\delta}\Theta)$ such that the limiting prediction risks match $\tau' = \tau$.¹⁰ Equation (19a) generalizes it for any $M \geq 1$. Further

⁹Observe that the factors of $c\delta$ arise because both the asymptotic training error β^2 and the asymptotic degrees of freedom correction ν are defined with normalization of p in Table 2.

¹⁰The parameter a requires the following change: $a' = (c\delta/\tau^2)^{(1-q)/2}a$ in the two systems.

special boundary cases of $q = 2$ (ridge) and $q = 1$ (lasso) are further isolated in the next two remarks.

Remark 4 (Ridge ensembles). For ridge regression when $q = 2$, System 2 recovers Theorem 4.1 of [PDK23] under isotropic features, using a slight change of variables $(\lambda', \Theta') = (\lambda/\sqrt{c\delta}, \sqrt{c\delta}\Theta)$. The solution v to the (limiting) Stieltjes transform of the spectrum of the sample gram matrix therein satisfies that $v = (a\tau)^{-1}$. For ridge ensembles, there is a deeper connection between \mathcal{R}_1 and \mathcal{R}_∞ . It turns out that $\mathcal{R}_\infty(\lambda, c)$ is exactly equal to $\mathcal{R}_1(\mu, 1)$ for a new level of regularization $\mu = v^{-1} \geq \lambda$ that depends on c (and properties of the data distribution). The lower the subsampling proportion of c , the higher the value of this implicit regularization μ . There are entire paths of equivalences in the (λ, c) plane where not only are the asymptotic squared risks the same, but also the estimators themselves are equivalent. In a sense, one can think of the combined effect of subsampling and ensembling as an additional (implicit) ridge regularization. The prediction risk equivalences are first proved in [DPK23] and later generalized to other risks and estimator equivalences in [PD23].

Remark 5 (Lasso ensembles). For the lasso predictor when $q = 1$, the first set of equations (18) for $M = 1$ in System 2 recovers [BM11b, Theorem 1.5] with a change of variables $\lambda' = \lambda/\sqrt{c\delta}$ and $\Theta' = \sqrt{c\delta}\Theta$. On the other hand, the second set of equations in System 2 for $M = \infty$ is new to the literature. We show empirically in the next section that the optimal full-ensemble subsampled lassoless is not the same as the optimal non-ensemble lasso (on full data). Thus, the subsampling and ensembling of the lasso are qualitatively different than the subsampling and ensembling of the ridge regression, as pointed out in Remark 4. In particular, the effect of subsampling for lasso is not merely additional lasso regularization.

4.2.2 Ensembles of general ridge regularized estimators

In this section, we specialize the results of Corollary 3 for ridge regularized ensembles allowing for general loss functions. Specifically, given $\lambda > 0$, consider component estimators of the form:

$$\hat{\theta}_m(I_m) := \operatorname{argmin}_{\mathbf{b} \in \mathbb{R}^p} \sum_{i \in I_m} \operatorname{loss}(y_i - \mathbf{x}_i^\top \mathbf{b}) + \lambda \|\mathbf{b}\|_2^2.$$

Recall that the risk convergence in (17) holds with $\mathcal{R}_1 = \alpha^2$ and $\mathcal{R}_\infty = \eta_G \alpha^2$ where α and η_G are solutions to System 1a and System 1b respectively. Here, using the explicit formula $\operatorname{prox}_{\operatorname{reg}}(x; \tau) = x/(1 + \lambda\tau)$ for $\operatorname{reg}(x) = \lambda|x|^2/2$, these systems can be simplified as follows:

System 3 (Nonlinear system for general ridge regularized ensembles). Given $\lambda \in (0, \infty)$, $\delta \in (0, \infty)$, and $c \in (0, 1]$, define the following 2-scalar system of equations in variables (α, κ) :

$$\alpha^2 = c\delta \cdot \kappa^2 \mathbb{E}[\operatorname{env}'_{\operatorname{loss}}(Z + \alpha G; \kappa)^2] + \lambda^2 \kappa^2 \mathbb{E}[\Theta^2] \quad \text{and} \quad \alpha = c\delta \cdot \frac{\kappa}{1 - \lambda\kappa} \mathbb{E}[\operatorname{env}'_{\operatorname{loss}}(Z + \alpha G; \kappa) \cdot G], \quad (20)$$

where $G \sim \mathcal{N}(0, 1)$, $\Theta \sim F_\theta$, $Z \sim F_z$, all mutually independent. Let (β, ν) be parameters expressed in terms of (α, κ) as:

$$\beta^2 = \frac{1}{\kappa^2} \alpha^2 - \lambda^2 \mathbb{E}[\Theta^2] \quad \text{and} \quad \nu = \frac{1}{\kappa} - \lambda. \quad (21)$$

Given parameters $(\alpha, \beta, \kappa, \nu)$ that satisfy (20) and (21), define the following 2-scalar system of

equations in variables (η_H, η_G) :

$$\eta_H = \frac{c^2 \delta}{\beta^2} \mathbb{E}[\text{env}'_{\text{loss}}(Z + \alpha G; \kappa) \cdot \text{env}'_{\text{loss}}(Z + \alpha \tilde{G}; \kappa)] \quad \text{and} \quad \eta_G = \frac{\eta_H \beta^2 + \lambda^2 \mathbb{E}[\Theta^2]}{\beta^2 + \lambda^2 \mathbb{E}[\Theta^2]}, \quad (22)$$

where (G, \tilde{G}) are jointly normal with $\mathbb{E}[G\tilde{G}] = \eta_G$.

When $M = 1$, System 3 recovers [EK13, Theorem 2.1]. When $c\delta > 1$ (underparameterized regime) and $\lambda=0$ (unregularized case), we have $\nu = \kappa^{-1}$, $\beta = \kappa^{-1}\alpha$, $\eta_H = \eta_G$. Substituting them to (22), we recover Theorem 2.3 in the recent work of [BK24].

5 Subagging and overparameterization

In Section 4, we discussed subagging of regularized estimators with an explicit regularization level $\lambda > 0$. Triggered by the success of overparameterized neural networks that can (nearly) interpolate, there has been a surge of recent work analyzing the risk behavior of estimators with vanishing regularization, such as the minimum ℓ_2 -norm interpolator (ridgeless), minimum ℓ_1 -norm interpolator (lassoless), and max-margin interpolators, among others. In this section, we discuss subagging of minimum ℓ_q -norm interpolators for $q \in \{1, 2\}$. We will demonstrate some interesting risk properties in Section 5.1, showcase the benefits of subagging in overparameterized regimes in Section 5.2, and contrast with optimal explicit regularization in Section 5.3.

5.1 Subagging of minimum ℓ_q -norm interpolators

We will focus in this section on subagging of bridgeless estimators, that is ℓ_q -norm regularized least squares with $\text{reg}(\mathbf{b}) = \|\mathbf{b}\|_q^q$. Our main cases of interest are the “ridgeless” and “lassoless” estimators, which are the special cases when $q = 1$ [HMRT22] and $q = 2$ [LW21], respectively. The terminology “less” is motivated by the fact that these estimators can be defined in a limiting sense as $\lambda \rightarrow 0^+$ for bridge estimators with regularization level λ . We consider the ensemble of predictors $(\hat{\boldsymbol{\theta}}_m)_{m \in [M]}$ of the form:¹¹

$$\hat{\boldsymbol{\theta}}_m := \lim_{\lambda \rightarrow 0^+} \hat{\boldsymbol{\theta}}_m(\lambda) \quad \text{where} \quad \hat{\boldsymbol{\theta}}_m(\lambda) \in \underset{\mathbf{b} \in \mathbb{R}^p}{\text{argmin}} \frac{1}{2} \|\mathbf{y}_{I_m} - \mathbf{X}_{I_m} \mathbf{b}\|_2^2 + \lambda \|\mathbf{b}\|_q^q.$$

In the underparameterized regime ($p < n$), these are simply the least squares estimators: $\hat{\boldsymbol{\theta}}_m = (\mathbf{X}_{I_m}^\top \mathbf{X}_{I_m})^{-1} \mathbf{X}_{I_m}^\top \mathbf{y}_{I_m}$. In the overparameterized regime ($p > n$), these correspond to the minimum ℓ_q -norm interpolators: $\{\|\hat{\boldsymbol{\theta}}_m\|_q: \mathbf{y}_{I_m} = \mathbf{X}_{I_m} \hat{\boldsymbol{\theta}}_m\}$, when \mathbf{X}_{I_m} has independent rows to allow for interpolation. For $q = 2$, when reg is the ridge penalty, this also has a closed-form expression given by: $\hat{\boldsymbol{\theta}}_m = (\mathbf{X}_{I_m}^\top \mathbf{X}_{I_m})^\dagger \mathbf{X}_{I_m}^\top \mathbf{y}_{I_m}$, where \mathbf{A}^\dagger denotes the Moore-Penrose pseudoinverse of a matrix \mathbf{A} . In other cases, we do not have a closed-form expression for the minimum ℓ_q -norm interpolator. The next system specializes System 2 to convex regularized least squares with vanishing regularization, by taking the limit as $\lambda \rightarrow 0^+$.

System 4 (Ensembles of minimum ℓ_q -norm interpolators). Given $\delta \in (0, \infty)$, $c \in (0, 1]$ such that $c\delta < 1$, and $\sigma^2 \in [0, \infty)$, define the following 2-scalar system of equations in variables

¹¹We refer readers to [Tib13] for details on how the sequence of estimators is defined when the estimators for $\lambda > 0$ are not unique, as in the case of the lasso.

(τ, a) :

$$\tau^2 = \mathbb{E}[(\text{prox}_{\text{reg}}(\Theta + \frac{\tau}{\sqrt{c\delta}}H; \frac{a\tau}{\sqrt{c\delta}}) - \Theta)^2] + \sigma^2 \quad (23a)$$

$$0 = 1 - \frac{1}{c\delta} \mathbb{E}[\text{prox}'_{\text{reg}}(\Theta + \frac{\tau}{\sqrt{c\delta}}H; \frac{a\tau}{\sqrt{c\delta}})] \quad (23b)$$

where $H \sim \mathcal{N}(0, 1)$, and $\Theta \sim F_\theta$ independent. Given (τ, a) that satisfy (23a) and (23b), define the following 1-scalar system of equations in variable ξ :

$$\xi^2 = \mathbb{E}[(\text{prox}_{\text{reg}}(\Theta + \frac{\tau}{\sqrt{c\delta}}H; \frac{a\tau}{\sqrt{c\delta}}) - \Theta) \cdot (\text{prox}_{\text{reg}}(\Theta + \frac{\tau}{\sqrt{c\delta}}\tilde{H}; \frac{a\tau}{\sqrt{c\delta}}) - \Theta)] + \sigma^2 \quad (24a)$$

where $\begin{pmatrix} H \\ \tilde{H} \end{pmatrix} \sim \mathcal{N}(\begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{bmatrix} 1 & \eta_H \\ \eta_H & 1 \end{bmatrix})$ with $\eta_H = c \frac{\xi^2}{\tau^2}$, and $\Theta \sim F_\theta$ independent.

To the best of our knowledge, the existence and uniqueness of the solution (τ, a) to (23) are not fully established in the literature, except for the special cases of $q = 2$ (ridgeless) and $q = 1$ (lassoless). Assuming this is the case, the existence and uniqueness of the solution ξ to (24a) follow from System 1b.

Observe that the equations (23a) and (24a) in System 4 are special cases of (18a) and (19a) in System 2 for ℓ_q penalties. Equation (23b) is the limit of (18b) as $\lambda \rightarrow 0^+$. Indeed, for $q \in \{1, 2\}$, the solution to System 2 with $\lambda > 0$ converges to the solution to System 4 as $\lambda \rightarrow 0^+$ (see Appendix C.2 for a proof). This means that

$$\lim_{\lambda \rightarrow 0^+} \text{p-lim}_{n \rightarrow +\infty} \left\| \frac{1}{M} \sum_{m \in [M]} \hat{\theta}_m(\lambda) - \theta \right\|_2^2 = \mathcal{R}_M := M^{-1} \mathcal{R}_1 + (1 - M^{-1}) \mathcal{R}_\infty$$

where $\mathcal{R}_1 = \tau^2 - \sigma^2$ and $\mathcal{R}_\infty = \xi^2 - \sigma^2$, and by the same argument, the limiting risk \mathcal{R}_M satisfies Proposition 6 and 7 (see also Figure 9 and 10). However, it is challenging to show the above display with the order of \lim_λ and p-lim_n swapped. For the ridgeless estimator ($q = 2$) and any M , this is proved in [PDK23] using a uniform convergence argument. For the lassoless estimator ($q = 1$ and $M = 1$), this is done in [LW21] where the authors directly analyze the interpolator by constructing a suitable AMP algorithm. We conjecture that this is, in general, true at least for bridgeless estimators for any $q \in [1, 2]$ and M . Since this is not the main focus of our paper and is only intended as an illustrative case, we will not work towards this goal in the current paper. We will instead investigate properties and consequences of System 4.

Further special cases of $q = 2$ (ridgeless) and $q = 1$ (lassoless) are isolated in the next two remarks. These will serve as our two main running examples in this section.

Remark 6 (Ridgeless ensembles). For squared loss and ridge regularizer ($q = 2$), when $c\delta < 1$ and $\lambda \rightarrow 0^+$, we get

$$\mathcal{R}_1(\delta, c) + \sigma^2 = \mathbb{E}[\Theta^2](1 - c\delta) + \sigma^2 \frac{1}{1 - c\delta}, \quad \mathcal{R}_\infty(\delta, c) + \sigma^2 = \mathbb{E}[\Theta^2] \frac{(1 - c\delta)^2}{\delta(\delta - (c\delta)^2)} + \sigma^2 \frac{\delta}{\delta - (c\delta)^2}.$$

The result above aligns with the risk ensemble of ridgeless estimators presented in Corollary 6.1 of [PDK23]. This can be seen by substituting $\delta = 1/\phi$ and $c = \phi/\psi$, or equivalently $c\delta = 1/\psi$.

Remark 7 (Lassoless ensembles). For squared loss and lasso regularizer ($q = 1$), when $c\delta < 1$ and $\lambda \rightarrow 0^+$, the total risk τ^2 is the solution to the following equations:

$$\begin{aligned} \tau^2 &= \sigma^2 + \mathbb{E}[(\text{soft}(\frac{\tau}{\sqrt{c\delta}}H + \Theta; \frac{a\tau}{\sqrt{c\delta}}) - \Theta)^2], \\ 1 &= \frac{1}{c\delta} \mathbb{P}(|\frac{\tau}{\sqrt{c\delta}}H + \Theta| > \frac{a\tau}{\sqrt{c\delta}}), \end{aligned} \quad (25)$$

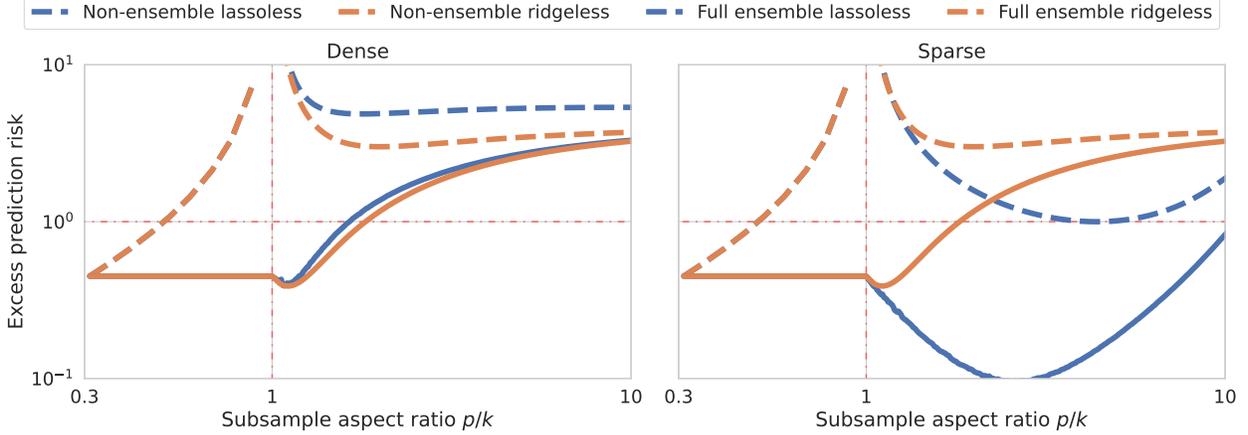


Figure 4: Prediction risks of the full-ensemble lassoless and ridgeless at different subsample aspect ratios p/k ranging from 0.3 to 10. The data model is given by (50) with signal strength $\rho = 2$, noise level $\sigma = 1$, data aspect ratio $p/n = 0.3$, and feature size $p = 500$. The support proportion s varies. *Left*: dense regime with $s = 0.9$. *Right*: sparse regime with $s = 0.01$. We observe that the full-ensemble risk is continuous at the interpolation. Also, the full-ensemble risk has a negative derivative to the right of interpolation. This implies that the optimal subsample size is in the overparameterized regime.

and ξ^2 is the solution to the following equations:

$$\xi^2 = \sigma^2 + \mathbb{E}\left[\left(\text{soft}\left(\frac{\tau}{\sqrt{c\delta}}H + \Theta; \frac{a\tau}{\sqrt{c\delta}}\right) - \Theta\right) \cdot \left(\text{soft}\left(\frac{\tau}{\sqrt{c\delta}}\tilde{H} + \Theta; \frac{a\tau}{\sqrt{c\delta}}\right) - \Theta\right)\right],$$

with $\mathbb{E}[H\tilde{H}] = c\frac{\xi^2}{\tau^2}$. Note that soft is the soft threshold function defined by $\text{soft}(x; \tau) = (|x| - \tau)_+ \text{sign}(x)$. In the ensemble setting, the case of $M = 1$ corresponds to [LW21, Theorem 2] with a slight change of variables $\Theta' = \sqrt{c\delta}\Theta$. The full-ensemble case when $M = \infty$ is new.

Remark 8 (Avoiding risk divergence with ensembling). Note that in the underparameterized regime when $c\delta > 1$ is fixed, the estimator of interest is simply the ensemble of least squares. Thus, a standard Stieltjes transform argument or the explicit formula for the expectation of inverse Wishart matrices give

$$\mathcal{R}_1(\delta, c) + \sigma^2 = \sigma^2 \frac{c\delta}{c\delta - 1} \quad \text{and} \quad \mathcal{R}_\infty(\delta, c) + \sigma^2 = \sigma^2 \frac{\delta}{\delta - 1} \quad \text{for all } c > \delta^{-1}.$$

It just so happens that for the full-ensemble least squares estimators, only the inverse aspect ratio δ of the original data matters! In particular, as $c \rightarrow (\delta^{-1})^+$, \mathcal{R}_1 diverges, while the full-ensemble risk \mathcal{R}_∞ is still bounded. Now let us consider the over-parameterized regime $c\delta < 1$. By simple algebra, for any regularizer reg , the solution τ to the sub-system (23a)-(23b) in System 4 is uniformly bounded from below as:

$$\tau^2 \geq (1 - c\delta)^{-1}\sigma^2. \quad (26)$$

(See Appendix C.1 for the proof.) Recalling $\mathcal{R}_1 = \tau^2 - \sigma^2$, this means that the risk of the non-ensemble interpolators blows up as $c \rightarrow (\delta^{-1})^-$. For the minimum ℓ_2 - and ℓ_1 -norm interpolators, this is shown in [HMRT22, LW21]. For the full-ensemble cases, we experimentally observe from Figure 4 that the risk \mathcal{R}_∞ is continuous in c for the full ridgeless and lassoless ensembles. In particular, it does not blow up around $c = \delta^{-1}$. For ridgeless, this claim is easy to verify (and holds more generally, as shown in [PDK23]). For lassoless, given the solution (a, τ) to (25), in Appendix C.3, we identify that the condition $\lim_{c \rightarrow (\delta^{-1})^-} (a\tau) = 0$ is sufficient to obtain the conclusion $\lim_{c \rightarrow (\delta^{-1})^-} (\xi^2) = \frac{\delta}{\delta - 1}\sigma^2$, and we observe experimentally that $a\tau \rightarrow 0$ holds (Figure 8), however we are currently not able to provably establish that $\lim_{c \rightarrow (\delta^{-1})^-} (a\tau) = 0$.

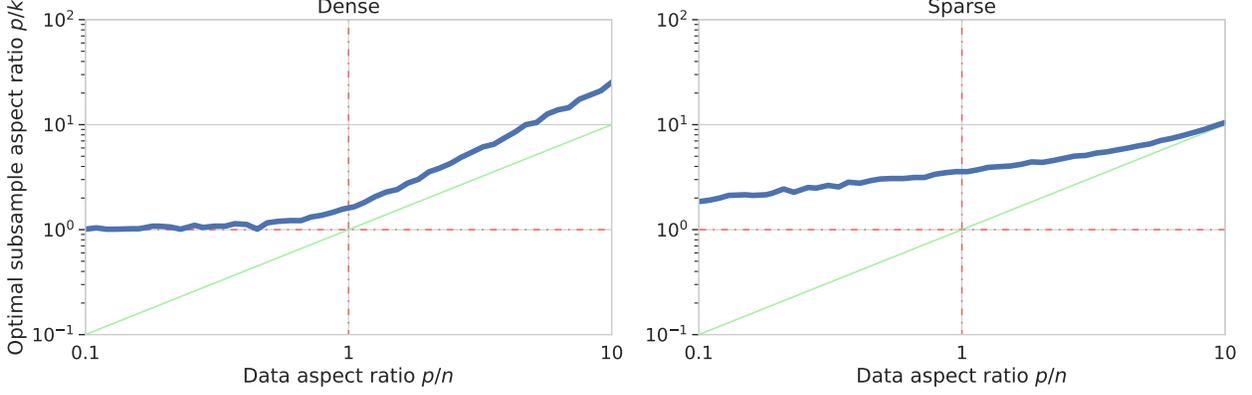


Figure 5: **Optimal subsample size for the lassoless ensemble is always in the overparameterized regime.** Optimal subsample aspect ratio p/k that achieves the optimal risk for the lassoless ensemble at different data aspect ratios p/n ranging from 0.1 to 10. The data model is as in (50) with signal strength $\rho = 2$, noise level $\sigma = 1$, data aspect ratio $p/n = 0.1$, feature size $p = 500$, and varying support proportion s . *Left*: dense regime with $s = 0.9$. *Right*: sparse regime with $s = 0.01$.

5.2 Optimal subsample size

An intriguing observation from Figure 4 concerns the optimal subsample size k_* . When the sample size n and the number of features p are fixed with $n > p$, the optimal subsample size that minimizes the full risk \mathcal{R}_∞ falls below p . This suggests that even when the original sample lies in the underparameterized regime, the optimal subsample size shifts into the overparameterized regime. This phenomenon is proved for ridgeless ensembles ($\text{reg}(x) = x^2$) by [PDK23]. Expanding on this, and utilizing System 4, we empirically show that this behavior extends beyond ridgeless ensembles to lassoless ensembles ($\text{reg}(x) = |x|$) as well (see Figure 5).

5.3 Optimal subbagging versus optimal (explicit) regularization

There are three parameters one can tune to optimize the asymptotic risk $\mathcal{R}_M(\lambda, c)$ of the ensemble estimator, as in Remarks 4 and 5: the regularization level λ , the subsample size c , and the ensemble size M . This hyperparameter optimization is simplified for ridge regression, as shown in Theorem 2.3 of [DPK23]. Minimization with respect to all three parameters is equal to the minimization over λ when $M = 1$ and $c = 1$ (non-ensemble setting), which is the same as minimization over M and c when $\lambda = 0$ (ensemble of ridgeless predictors):

$$\underbrace{\min_{\lambda \in [0, \infty], M \in \mathbb{N}, c \in [0, 1]} \mathcal{R}_M(\lambda, c)}_{\text{opt regularization and opt ensemble}} = \underbrace{\min_{\lambda \in [0, \infty]} \mathcal{R}_1(\lambda, 1)}_{\text{opt regularization but no ensemble}} = \underbrace{\min_{c \in [0, 1]} \mathcal{R}_\infty(0, c)}_{\text{opt ensemble but no regularization}} .$$

In some situations, however, the risk minimization over all three parameters can be strictly better:

$$\underbrace{\min_{\lambda \in [0, \infty], M \in \mathbb{N}, c \in [0, 1]} \mathcal{R}_M(\lambda, c)}_{\text{opt regularization and opt ensemble}} < \underbrace{\min_{\lambda \in [0, \infty]} \mathcal{R}_1(\lambda, 1)}_{\text{opt regularization but no ensemble}} \wedge \underbrace{\min_{c \in [0, 1]} \mathcal{R}_\infty(0, c)}_{\text{opt ensemble but no regularization}} .$$

We illustrate this through a numerical experiment with lasso. We show that the optimal full-ensemble subsampled lassoless is not the same as the optimal non-ensemble lasso (on full data). In Figure 6, we contrast the sparse and dense data settings. In each case, we show the full-ensemble risk heatmap in λ and p/k . In the left panel (the sparse setting), we see that the optimal subsample

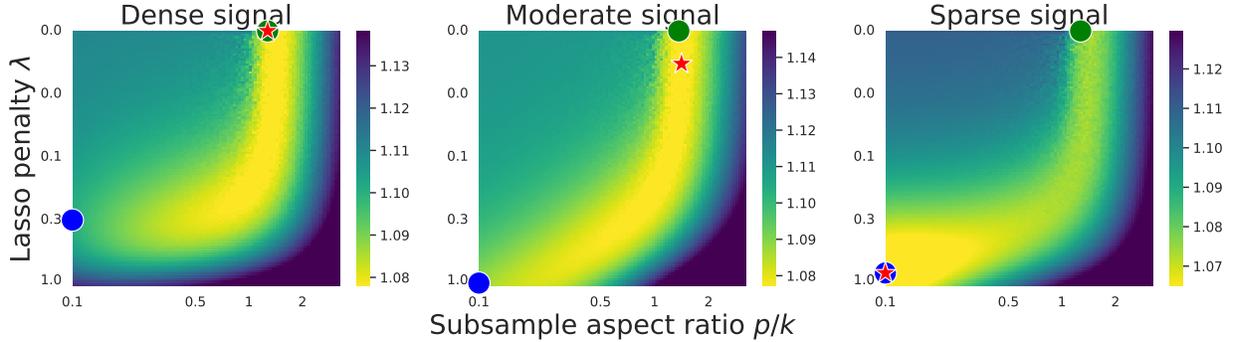


Figure 6: **Optimally subsampled lassoless regression can outperform optimal lasso regression.** Heatmaps of theoretical prediction risk in λ and p/k of full lasso ensemble in the underparameterized regime ($p/n = 0.1$). The data model is given by (50) with signal strength $\rho = 0.5$ and noise level $\sigma = 1$ at different sparsity levels s . *Left*: Dense regime with support proportion $s = 0.9$. Optimal subsample lassoless is better than optimal lasso. *Middle*: Modest sparse regime with support proportion $s = 0.5$. Optimal lasso ensemble is better than the optimal lasso and optimal subsample lassoless. *Right*: Sparse regime with support proportion $s = 0.2$. Optimal lasso is better than optimal subsample lassoless.

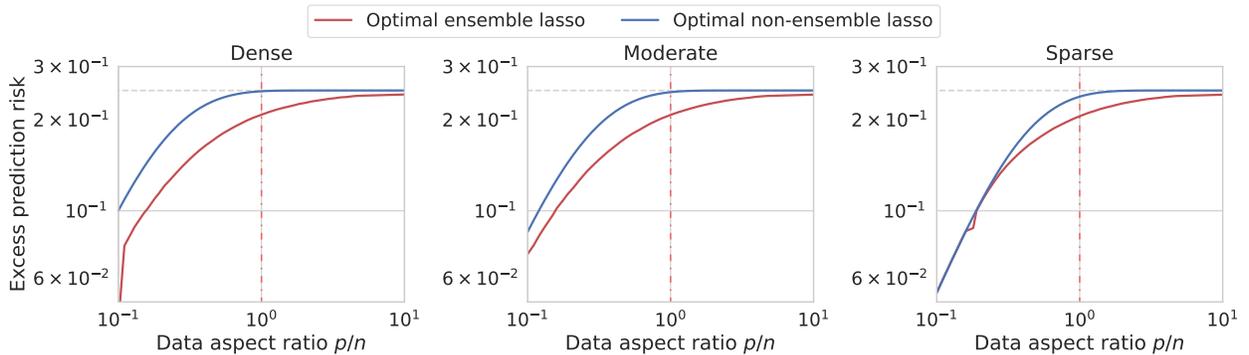


Figure 7: **Optimally subsampled ensemble lasso can uniformly beat optimally tuned non-ensemble lasso across different data aspect ratios.** The theoretical prediction risk of optimal ensemble and non-ensemble lasso at different data aspect ratios p/n is shown. The data model is given by (50) with signal strength $\rho = 0.5$ and noise level $\sigma = 1$ at different sparsity levels s . *Left*: Dense regime with support proportion $s = 0.9$. *Middle*: Modest sparse regime with support proportion $s = 0.5$. *Right*: Sparse regime with support proportion $s = 0.2$.

lassoless is worse than the optimal lasso. This is expected because the lasso is known to perform well in the sparse setting. In the right panel (the dense setting), we see that the optimal subsample lassoless is better than the optimal lasso. This is interesting because it shows that the subsample and ensemble induce an implicit regularization effect. In short, the optimal lassoless ensemble can be better or worse than the optimal lasso. In other words, the full-ensemble lassoless when c is optimized is not the same as the optimized lasso when λ is optimized on the full data.

A similar conclusion holds for overparameterized settings, as shown in Figure 11. In particular, depending on the SNR and δ , the subsample optimized risk may be smaller or larger than the optimized lasso risk. The conclusion is that, in general, it helps to optimize λ , but also to optimize the ensemble size and subsample size. The subsample and ensemble induce an implicit regularization effect. Once optimized, this implicit regularization can improve on the explicit regularization provided by the regularizer.

Finally, Figure 7 shows that the joint optimization of subsampling along with lasso penalty consistently outperforms the optimization of lasso penalty on the full data (without any subsampling) uniformly across varying data aspect ratios. Particularly, in the dense regime ($s = 0.9$) and moderate sparsity ($s = 0.5$) regimes, the ensemble approach leverages implicit regularization to achieve lower prediction risk. In highly sparse scenarios ($s = 0.2$), this effect reduces as one would expect. Overall, we see complementary benefits of optimizing both subsampling and explicit regularization parameters to improve the predictive performance, with the joint optimization being the overall clear winner.

6 Open directions

In this paper, we provide a general risk characterization for the ensemble of regularized M-estimators. The characterization depends on the inverse data aspect ratio $\delta = \lim n/p$, subsample ratio c , and loss and regularizer pair $(\text{loss}, \text{reg})$. We also specialize the results for specific cases of interest, such as the lasso and ridge regression, and analyze various properties related to optimal subsampling and ensembling. Our goal in performing such analysis is to shed light on how subsampling and ensembling influence the risk of the ensemble estimator and the optimal choices of ensemble size and subsample size. The key takeaway is that subsampling and ensembling can be beneficial in terms of reducing the risk of the estimator, and when the subsample size c is optimized, the resulting risk is monotonic in δ for any ensemble size M .

There are several “axes” along which our results on the asymptotics of subagging can be extended. These are all apparent by inspecting the last row of Table 1 (our current results) and contrasting it against the rows above (the known results in specific cases). We briefly mention some of these open directions next to make them explicit. First, our current analysis can handle non-differentiable and non-strongly convex regularizers, but we require a differentiable loss. Extending the analysis to non-differentiable losses through techniques like Moreau smoothing [CMW23, Section B.7] is a promising future direction. In addition to relaxing conditions assumed for risk characterization, we are also interested in relaxing the assumptions for the risk estimation (Theorem 4), particularly the strong convexity assumption on reg . A promising approach in this direction is to apply the Gaussian smoothing technique recently proposed by [BK23a]. Additionally, it is of interest to extend the scope beyond separable regularizers to include non-separable regularizers, which have been studied only for special cases like generalized ridge regression [PD23]. Another potential extension is relaxing the assumption of linear response models. Partial progress in this direction includes the recent work by [BK24] (for logistic models), [CVD⁺24] (for generalized linear models), and [PD23] (for arbitrary response with bounded moments). Furthermore, while we assume isotropic Gaussian features, it is also of interest to extend the results to general feature distributions and anisotropic covariances. This has been studied in special cases of ridge regression ([PDK23], any M), ridgeless regression ([HX23], $M = 1$), lasso regression ([CMW23], $M = 1$), among others; see also, e.g., [PKLS23, HS23], for some progress towards establishing Gaussian universality (for $M = 1$). Finally, broadening the scope of resampling strategies beyond subsampling without replacement to include more general schemes like sampling with replacement or according to a specific distribution, as recently studied in [CVD⁺24, DP24a], is another future direction. We hope our current analysis serves as a step towards these open directions in various axes.

Acknowledgements

We thank Michael Celentano and Ryan Tibshirani for useful discussions. PP and JHD acknowledge the computing support provided by the grant MTH230020 for experiments run on the Bridge-2 system at the Pittsburgh Supercomputing Center. PCB acknowledges partial support from the NSF Grant DMS-1945428.

References

- [AK23] Ryo Ando and Fumiyasu Komaki. On high-dimensional asymptotic properties of model averaging estimators. *arXiv preprint arXiv:2308.09476*, 2023.
- [AP20] Ben Adlam and Jeffrey Pennington. Understanding double descent requires a fine-grained bias-variance decomposition. In *Advances in Neural Information Processing Systems*, 2020.
- [BBEKY13] Derek Bean, Peter J. Bickel, Nouredine El Karoui, and Bin Yu. Optimal M -estimation in high-dimensional regression. *Proceedings of the National Academy of Sciences*, 110(36):14563–14568, 2013.
- [BDK⁺24] Pierre C. Bellec, Jin-Hong Du, Takuya Koriyama, Pratik Patil, and Kai Tan. Corrected generalized cross-validation for finite ensembles of penalized estimators. *Journal of the Royal Statistical Society Series B: Statistical Methodology*, 2024.
- [Bel21] Mikhail Belkin. Fit without fear: remarkable mathematical phenomena of deep learning through the prism of interpolation. *Acta Numerica*, 30:203–248, 2021.
- [Bel22] Pierre C. Bellec. Observable adjustments in single-index models for regularized M -estimators. *arXiv preprint arXiv:2204.06990*, 2022.
- [Bel23] Pierre C. Bellec. Out-of-sample error estimation for M -estimators with convex penalty. *Information and Inference: A Journal of the IMA*, 12(4):2782–2817, 2023.
- [BK23a] Pierre C. Bellec and Takuya Koriyama. Error estimation and adaptive tuning for unregularized robust M -estimator. *arXiv preprint arXiv:2312.13257*, 2023.
- [BK23b] Pierre C. Bellec and Takuya Koriyama. Existence of solutions to the nonlinear equations characterizing the precise error of M -estimators. *arXiv preprint arXiv:2312.13254*, 2023.
- [BK24] Pierre C. Bellec and Takuya Koriyama. Asymptotics of resampling without replacement in robust and logistic regression. *arXiv preprint arXiv:2404.02070*, 2024.
- [BM11a] Mohsen Bayati and Andrea Montanari. The dynamics of message passing on dense graphs, with applications to compressed sensing. *IEEE Transactions on Information Theory*, 57(2):764–785, 2011.
- [BM11b] Mohsen Bayati and Andrea Montanari. The lasso risk for gaussian matrices. *IEEE Transactions on Information Theory*, 58(4):1997–2017, 2011.
- [BMR21] Peter L. Bartlett, Andrea Montanari, and Alexander Rakhlin. Deep learning: a statistical viewpoint. *Acta Numerica*, 30:87–201, 2021.
- [Bre96] Leo Breiman. Bagging predictors. *Machine Learning*, 24(2):123–140, 1996.

- [Bre01] Leo Breiman. Using iterated bagging to debias regressions. *Machine Learning*, 45:261–277, 2001.
- [BS06] Andreas Buja and Werner Stuetzle. Observations on bagging. *Statistica Sinica*, pages 323–351, 2006.
- [BS22] Pierre C. Bellec and Yiwei Shen. Derivatives and residual distribution of regularized M-estimators with application to adaptive tuning. In *Conference on Learning Theory*, 2022.
- [BY02] Peter Bühlmann and Bin Yu. Analyzing bagging. *The Annals of Statistics*, 30(4):927–961, 2002.
- [CM22] Michael Celentano and Andrea Montanari. Fundamental barriers to high-dimensional regression with convex penalties. *The Annals of Statistics*, 50(1):170–196, 2022.
- [CM24] Michael Celentano and Andrea Montanari. Correlation adjusted debiased lasso: debiasing the lasso with inaccurate covariate model. *Journal of the Royal Statistical Society Series B: Statistical Methodology*, 2024.
- [CMW23] Michael Celentano, Andrea Montanari, and Yuting Wei. The lasso with general gaussian designs with applications to hypothesis testing. *The Annals of Statistics*, 51(5):2194–2220, 2023.
- [CVD⁺24] Lucas Clarté, Adrien Vandembroucq, Guillaume Dalle, Bruno Loureiro, Florent Krzakala, and Lenka Zdeborová. Analysis of bootstrap and subsampling in high-dimensional regularized regression. *arXiv preprint arXiv:2402.13622*, 2024.
- [CZYS23] Xin Chen, Yicheng Zeng, Siyue Yang, and Qiang Sun. Sketched ridgeless linear regression: The role of downsampling. In *International Conference on Machine Learning*, 2023.
- [DKT22] Zeyu Deng, Abba Kammoun, and Christos Thrampoulidis. A model of double descent for high-dimensional binary linear classification. *Information and Inference: A Journal of the IMA*, 11(2):435–495, 2022.
- [DM16] David L. Donoho and Andrea Montanari. High dimensional robust M-estimation: Asymptotic variance via approximate message passing. *Probability Theory and Related Fields*, 166(3):935–969, 2016.
- [DMM09] David L. Donoho, Arian Maleki, and Andrea Montanari. Message-passing algorithms for compressed sensing. *Proceedings of the National Academy of Sciences*, 106(45):18914–18919, 2009.
- [DMM11] David L. Donoho, Arian Maleki, and Andrea Montanari. The noise-sensitivity phase transition in compressed sensing. *IEEE Transactions on Information Theory*, 57(10):6920–6941, 2011.
- [DP24a] Jin-Hong Du and Pratik Patil. Implicit regularization paths of weighted neural representations. *arXiv preprint arxiv:2408.15784*, 2024.
- [DP24b] Jin-Hong Du and Pratik Patil. Python package sklearn_ensemble_cv v0.2.3. PyPI, 2024.

- [DPK23] Jin-Hong Du, Pratik Patil, and Arun Kumar Kuchibhotla. Subsample ridge ensembles: Equivalences and generalized cross-validation. In *International Conference on Machine Learning*, 2023.
- [DPRK24] Jin-Hong Du, Pratik Patil, Kathryn Roeder, and Arun Kumar Kuchibhotla. Extrapolated cross-validation for randomized ensembles. *Journal of Computational and Graphical Statistics*, 2024.
- [dRBK20] Stéphane d’Ascoli, Maria Refinetti, Giulio Biroli, and Florent Krzakala. Double trouble in double descent: Bias and variance (s) in the lazy regime. In *International Conference on Machine Learning*, 2020.
- [DS01] Kenneth R. Davidson and Stanislaw J. Szarek. Local operator theory, random matrices and banach spaces. *Handbook of the Geometry of Banach Spaces*, 1, 2001.
- [DS20] Edgar Dobriban and Yue Sheng. Wonder: Weighted one-shot distributed ridge regression in high dimensions. *Journal of Machine Learning Research*, 21(66):1–52, 2020.
- [DS21] Edgar Dobriban and Yue Sheng. Distributed linear regression by averaging. *The Annals of Statistics*, 49(2):918–943, 2021.
- [DW18] Edgar Dobriban and Stefan Wager. High-dimensional asymptotics of prediction: Ridge regression and classification. *The Annals of Statistics*, 46(1):247–279, 2018.
- [EK13] Noureddine El Karoui. Asymptotic behavior of unregularized and ridge-regularized high-dimensional robust regression estimators: rigorous results. *arXiv preprint arXiv:1311.2445*, 2013.
- [EKBB⁺13] Noureddine El Karoui, Derek Bean, Peter J. Bickel, Chinghway Lim, and Bin Yu. On robust regression with high-dimensional predictors. *Proceedings of the National Academy of Sciences*, 110(36):14557–14562, 2013.
- [FF93] Ildiko E. Frank and Jerome H. Friedman. A statistical view of some chemometrics regression tools. *Technometrics*, 35(2):109–135, 1993.
- [FH07] Jerome H. Friedman and Peter Hall. On bagging and nonlinear estimation. *Journal of Statistical Planning and Inference*, 137(3):669–683, 2007.
- [Fu98] Wenjiang J. Fu. Penalized regressions: the bridge versus the lasso. *Journal of Computational and Graphical Statistics*, 7(3):397–416, 1998.
- [GB23] Cédric Gerbelot and Raphaël Berthier. Graph-based approximate message passing iterations. *Information and Inference: A Journal of the IMA*, 12(4):2562–2628, 2023.
- [HMRT22] Trevor Hastie, Andrea Montanari, Saharon Rosset, and Ryan J. Tibshirani. Surprises in high-dimensional ridgeless least squares interpolation. *The Annals of Statistics*, 50(2):949–986, 2022.
- [hp24] Iosif Pinelis (<https://mathoverflow.net/users/36721/iosif-pinelis>). Large deviations: Growth of empirical average of iid non-negative random variables with infinite expectations? MathOverflow, 2024. URL:<https://mathoverflow.net/q/390939> (version: 2024).
- [HS90] L.K. Hansen and P. Salamon. Neural network ensembles. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 12(10):993–1001, 1990.

- [HS05] Peter Hall and Richard J. Samworth. Properties of bagged nearest neighbour classifiers. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 67(3):363–379, 2005.
- [HS23] Qiyang Han and Yandi Shen. Universality of regularized regression estimators in high dimensions. *The Annals of Statistics*, 51(4):1799–1823, 2023.
- [HTF09] Trevor Hastie, Robert Tibshirani, and Jerome H. Friedman. *The Elements of Statistical Learning: Data Mining, Inference, and Prediction*. Springer, 2009. Second edition.
- [Hub64] Peter J. Huber. Robust estimation of a location parameter. *The Annals of Mathematical Statistics*, 35(1):73–101, 1964.
- [HX23] Qiyang Han and Xiacong Xu. The distribution of ridgeless least squares interpolators. *arXiv preprint arXiv:2307.02044*, 2023.
- [JM13] Adel Javanmard and Andrea Montanari. State evolution for general approximate message passing algorithms, with applications to spatial coupling. *Information and Inference: A Journal of the IMA*, 2(2):115–144, 2013.
- [Kar18] Noureddine El Karoui. 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.
- [KS97] Anders Krogh and Peter Sollich. Statistical mechanics of ensemble learning. *Physical Review E*, 55(1):811, 1997.
- [LGC⁺21] Bruno Loureiro, Cedric Gerbelot, Hugo Cui, Sebastian Goldt, Florent Krzakala, Marc Mezard, and Lenka Zdeborová. Learning curves of generic features maps for realistic datasets with a teacher-student model. *Advances in Neural Information Processing Systems*, 34:18137–18151, 2021.
- [LGR⁺22] Bruno Loureiro, Cedric Gerbelot, Maria Refinetti, Gabriele Sicuro, and Florent Krzakala. Fluctuations, bias, variance and ensemble of learners: Exact asymptotics for convex losses in high-dimension. In *International Conference on Machine Learning*, 2022.
- [LJB20] Daniel LeJeune, Hamid Javadi, and Richard Baraniuk. The implicit regularization of ordinary least squares ensembles. In *International Conference on Artificial Intelligence and Statistics*, 2020.
- [LS22] Tengyuan Liang and Pragya Sur. A precise high-dimensional asymptotic theory for boosting and minimum- ℓ_1 -norm interpolated classifiers. *The Annals of Statistics*, 50(3):1669 – 1695, 2022.
- [LW21] Yue Li and Yuting Wei. Minimum ℓ_1 -norm interpolators: Precise asymptotics and multiple descent. *arXiv preprint arXiv:2110.09502*, 2021.
- [MLC19] Xiaoyi Mai, Zhenyu Liao, and Romain Couillet. A large scale analysis of logistic regression: Asymptotic performance and new insights. In *International Conference on Acoustics, Speech and Signal Processing*, 2019.
- [MM21] Léo Miolane and Andrea Montanari. The distribution of the lasso: Uniform control over sparse balls and adaptive parameter tuning. *The Annals of Statistics*, 49(4):2313–2335, 2021.

- [MRRK22] Nicole Mücke, Enrico Reiss, Jonas Rungenhagen, and Markus Klein. Data-splitting improves statistical performance in overparameterized regimes. In *International Conference on Artificial Intelligence and Statistics*, 2022.
- [MRSY19] Andrea Montanari, Feng Ruan, Youngtak Sohn, and Jun Yan. The generalization error of max-margin linear classifiers: High-dimensional asymptotics in the overparametrized regime. *arXiv preprint arXiv:1911.01544*, 2019.
- [OH16] Samet Oymak and Babak Hassibi. Sharp mse bounds for proximal denoising. *Foundations of Computational Mathematics*, 16:965–1029, 2016.
- [OTH13] Samet Oymak, Christos Thrampoulidis, and Babak Hassibi. The squared-error of generalized lasso: A precise analysis. In *Allerton Conference on Communication, Control, and Computing*, 2013.
- [PD23] Pratik Patil and Jin-Hong Du. Generalized equivalences between subsampling and ridge regularization. In *Advances in Neural Information Processing Systems*, 2023.
- [PDK23] Pratik Patil, Jin-Hong Du, and Arun Kumar Kuchibhotla. Bagging in overparameterized learning: Risk characterization and risk monotonicity. *Journal of Machine Learning Research*, 24(319):1–113, 2023.
- [Per93] Michael Perrone. Putting it all together: Methods for combining neural networks. In *Advances in Neural Information Processing Systems*, 1993.
- [PKLS23] Luca Pesce, Florent Krzakala, Bruno Loureiro, and Ludovic Stephan. Are gaussian data all you need? the extents and limits of universality in high-dimensional generalized linear estimation. In *International Conference on Machine Learning*, 2023.
- [PKWR22] Pratik Patil, Arun Kumar Kuchibhotla, Yuting Wei, and Alessandro Rinaldo. Mitigating multiple descents: A model-agnostic framework for risk monotonicity. *arXiv preprint arXiv:2205.12937*, 2022.
- [PL24] Pratik Patil and Daniel LeJeune. Asymptotically free sketched ridge ensembles: Risks, cross-validation, and tuning. In *International Conference on Learning Representations*, 2024.
- [SAH19] Fariborz Salehi, Ehsan Abbasi, and Babak Hassibi. The impact of regularization on high-dimensional logistic regression. In *Advances in Neural Information Processing Systems*, 2019.
- [Sam12] Richard J. Samworth. Optimal weighted nearest neighbour classifiers. *The Annals of Statistics*, 40(5):2733–2763, 2012.
- [SC19] Pragya Sur and Emmanuel J. Candès. A modern maximum-likelihood theory for high-dimensional logistic regression. *Proceedings of the National Academy of Sciences*, 116(29):14516–14525, 2019.
- [SK95] Peter Sollich and Anders Krogh. Learning with ensembles: How overfitting can be useful. In *Advances in Neural Information Processing Systems*, 1995.
- [Ste81] Charles M. Stein. Estimation of the mean of a multivariate normal distribution. *The Annals of Statistics*, 9(6):1135–1151, 1981.

- [TAH18] Christos Thrampoulidis, Ehsan Abbasi, and Babak Hassibi. Precise error analysis of regularized M-estimators in high dimensions. *IEEE Transactions on Information Theory*, 64(8):5592–5628, 2018.
- [Tib13] Ryan J. Tibshirani. The lasso problem and uniqueness. *Electronic Journal of Statistics*, 7:1456–1490, 2013.
- [TOH15] Christos Thrampoulidis, Samet Oymak, and Babak Hassibi. Regularized linear regression: A precise analysis of the estimation error. In *Conference on Learning Theory*, 2015.
- [VdV00] Aad W. Van der Vaart. *Asymptotic Statistics*, volume 3. Cambridge University Press, 2000.
- [WMZ18] Haolei Weng, Arian Maleki, and Le Zheng. Overcoming the limitations of phase transition by higher order analysis of regularization techniques. *The Annals of Statistics*, 46(6A):3099 – 3129, 2018.
- [WWM20] Shuaiwen Wang, Haolei Weng, and Arian Maleki. Which bridge estimator is the best for variable selection? *The Annals of Statistics*, 48(5):2791 – 2823, 2020.

Supplement

This serves as a supplement to the paper “Precise Asymptotics of Bagging Regularized M-estimators.” Below, we provide an outline of the supplement along with a summary of the notation used in the main paper and the supplement.

Organization

Section	Content
Appendix A	Proofs of Theorems 1, 2 and 4 from Section 3
Appendix B	Proof of Proposition 7 and other miscellaneous details from Section 4
Appendix C	Details of arguments in Remark 8 from Section 5
Appendix D	Additional numerical illustrations and experimental details

Table 3: Outline of the supplement.

Specific notation

Notation	Description
$(\mathbf{x}_i, y_i), i \in [n]$	Observation vector with the feature vector in \mathbb{R}^p and the response variable in \mathbb{R}
\mathbf{X}, \mathbf{y}	Feature matrix in $\mathbb{R}^{n \times p}$ and the response vector in \mathbb{R}^n
$\boldsymbol{\theta}$	Signal vector in the linear model in \mathbb{R}^p with entries drawn i.i.d. from the distribution F_θ
\mathbf{z}	Noise vector in the linear model in \mathbb{R}^n with entries drawn i.i.d. from the distribution F_z
k, M	Size of each subsample and the total number of subsamples (bags) used in the ensemble
$I_m, m \in [M]$	Index set (subset of $[n]$) of the m -th subsample, with $ I_m = k$
$(\mathbf{X}_{I_m}, \mathbf{y}_{I_m}), m \in [M]$	Feature matrix in $\mathbb{R}^{k \times p}$ and response vector in \mathbb{R}^k of the subsampled dataset indexed by I_m
$(\mathbf{X}_{I_m \cap I_\ell}, \mathbf{y}_{I_m \cap I_\ell}), m \neq \ell \in [M]$	Feature matrix in $\mathbb{R}^{ I_m \cap I_\ell \times p}$ and response vector in $\mathbb{R}^{ I_m \cap I_\ell }$ of the overlapped dataset corresponding to the overlap between the subsampled datasets indexed by I_m and I_ℓ
loss, reg, λ	Loss function, regularization function, and regularization level used for the regularized M-estimator
$\hat{\boldsymbol{\theta}}_m, \tilde{\boldsymbol{\theta}}_M$	Component regularized M-estimator fitted on the m -th subsample and the ensemble estimator
$\mathcal{R}_M, \mathcal{R}_M$	Squared prediction risk of the ensemble estimator and its asymptotic limit
$\mathcal{R}_1, \mathcal{R}_\infty$	Asymptotic risk of the non-ensemble estimator ($M = 1$) and the full-ensemble estimator ($M = \infty$)
δ, c	Inverse data aspect ratio n/p and subsample ratio k/n
$(\alpha, \beta, \kappa, \nu)$	Parameters characterizing the ensemble risk asymptotics with $M = 1$ (System 1a)
G, H	Standard normal random variables in the ensemble risk asymptotics with $M = 1$ (System 1a)
Θ, Z	Random variables drawn according to distributions F_θ and F_z , respectively (System 1a)
(η_G, η_H)	Correlation parameters in the ensemble risk asymptotics with $M = \infty$ (System 1b)
$G, \tilde{G}, H, \tilde{H}$	Random variables appearing in the ensemble risk asymptotics with $M = \infty$ (System 1b)
(a, τ)	Parameters in alternate formulation of the ensemble risk asymptotics with $M = 1$ (System 2)
ξ	Parameter alternate formulation of the ensemble risk asymptotics with $M = \infty$ (System 2)

Table 4: Summary of some of the specific notation used in the paper.

A Proofs in Section 3

A.1 Proof of Theorem 1

Part 1. By the Cauchy–Schwarz inequality, if η_G satisfies $|F_{\text{loss}}(\eta_G)| = \sqrt{c\tilde{c}}$, then we must have

$$\text{sign}(F_{\text{loss}}(\eta_G)) \cdot \text{env}'_{\text{loss}}(\alpha G + Z; \kappa) = \text{env}'_{\text{loss}}(\tilde{\alpha}(\eta_G G + \sqrt{1 - \eta_G^2} \bar{G}) + Z; \tilde{\kappa})$$

with probability 1 for any independent Gaussian G, \bar{G} . Multiplying both sides by $\sqrt{1 - \eta_G^2} \bar{G}$ and taking expectation, using independence of (G, \bar{G}) and the third equation in System 1a, we are left with

$$0 = \sqrt{1 - \eta_G^2} \cdot (\tilde{\nu} \tilde{\alpha}) / (c\delta),$$

which gives $|\eta_G| = 1$ since $\tilde{\nu} \tilde{\alpha} > 0$. This means $|F_{\text{loss}}(\eta_G)| < \sqrt{c\tilde{c}}$ for all $|\eta_G| < 1$. By the same argument, if η_H satisfies $|F_{\text{reg}}(\eta_H)| = 1$, then it holds that

$$\text{sign}(F_{\text{reg}}(\eta_H)) \cdot \left[\frac{1}{\tilde{\nu}} \cdot \text{env}'_{\text{reg}}\left(\frac{\beta}{\nu} H + \Theta; \frac{1}{\nu}\right) - \frac{\beta}{\nu} H \right] = \frac{1}{\tilde{\nu}} \cdot \text{env}'_{\text{reg}}\left(\frac{\tilde{\beta}}{\tilde{\nu}}(\eta_H H + \sqrt{1 - \eta_H^2} \bar{H}) + \Theta; \frac{1}{\tilde{\nu}}\right) - \frac{\tilde{\beta}}{\tilde{\nu}} \bar{H}$$

with probability 1 for any independent Gaussian (H, \bar{H}) . Multiplying the the both side by $\sqrt{1 - \eta_H^2} \bar{H}$ and taking the expectation, using the independence of (H, \bar{H}) and the fourth equation in System 1a, we are left with

$$0 = \sqrt{1 - \eta_H^2} (-\tilde{\kappa} \tilde{\beta}),$$

which gives $|\eta_H| = 1$ since $\tilde{\kappa} \tilde{\beta} > 0$. This means $|F_{\text{reg}}(\eta_H)| < 1$ for all $|\eta_H| < 1$.

Part 2. Let us prove the differentiability and contraction of compositions. By System 1a, F_{loss} and F_{reg} can be written as

$$F_{\text{loss}}(\eta_G) = \frac{c\tilde{c}\delta}{\beta\tilde{\beta}} \cdot \mathbb{E}[\text{env}'_{\text{loss}}(\alpha G + Z; \kappa) \cdot \text{env}'_{\text{loss}}(\tilde{\alpha}\tilde{G} + Z; \tilde{\kappa})],$$

$$F_{\text{reg}}(\eta_H) = \frac{1}{\alpha\tilde{\alpha}} \cdot \mathbb{E}\left[\left(\frac{1}{\nu} \cdot \text{env}'_{\text{reg}}\left(\frac{\beta}{\nu} H + \Theta; \frac{1}{\nu}\right) - \frac{\beta}{\nu} H\right) \cdot \left(\frac{1}{\tilde{\nu}} \cdot \text{env}'_{\text{reg}}\left(\frac{\tilde{\beta}}{\tilde{\nu}} \tilde{H} + \Theta; \frac{1}{\tilde{\nu}}\right) - \frac{\tilde{\beta}}{\tilde{\nu}} \tilde{H}\right)\right].$$

We will use the following lemma to argue the differentiability of F_{loss} and F_{reg} .

Lemma 8. *Let G and Z be independent $\mathcal{N}(0, 1)$ random variables. Then, for any Lipschitz functions (f, \tilde{f}) with bounded second moment $\mathbb{E}[f(G)^2], \mathbb{E}[\tilde{f}(G)^2] < +\infty$, the map*

$$\varphi : [-1, 1] \rightarrow \mathbb{R}, \quad \eta \mapsto \mathbb{E}[f(G)\tilde{f}(\eta G + \sqrt{1 - \eta^2} Z)]$$

has the derivative

$$\varphi'(\eta) = \mathbb{E}[f'(G)\tilde{f}'(\eta G + \sqrt{1 - \eta^2} Z)].$$

Proof. Since \tilde{f} is Lipschitz and $\mathcal{N}(0, 1)$ has no point mass, \tilde{f} is differentiable at $G \sim \mathcal{N}(0, 1)$ with probability 1. By the dominated convergence theorem, we have

$$\varphi'(\eta) = \mathbb{E}\left[f(G)\tilde{f}'(\eta G + \sqrt{1 - \eta^2} Z)\left(G - \frac{\eta}{\sqrt{1 - \eta^2}} Z\right)\right].$$

Let us define $A = \eta G + \sqrt{1 - \eta^2} Z$ and $B = \sqrt{1 - \eta^2} G - \eta Z$ so that (A, B) are independent Gaussian $\mathcal{N}(0, 1)$ and $\varphi'(\eta) = (1 - \eta^2)^{-1/2} \mathbb{E}[f(\eta A + \sqrt{1 - \eta^2} B) \tilde{f}'(A) B]$. Using Stein's formula for B conditionally on A , we are left with

$$\varphi'(\eta) = \mathbb{E}[\tilde{f}'(A) f'(\eta A + \sqrt{1 - \eta^2} B)] = \mathbb{E}[\tilde{f}'(\eta G + \sqrt{1 - \eta^2} Z) f'(G)],$$

where we have used $(A, \eta A + \sqrt{1 - \eta^2} B) \stackrel{d}{=} (\eta G + \sqrt{1 - \eta^2} Z, G)$. (Here and throughout $\stackrel{d}{=}$ refers to equality in distribution.) \square

Notice that $\text{env}'_{\text{loss}}(\alpha G + Z; \kappa)$ and $\frac{1}{\nu} \cdot \text{env}'_{\text{reg}}\left(\frac{\beta}{\nu} H + \Theta; \frac{1}{\nu}\right) - \frac{\beta}{\nu} H$ have finite second moments, thanks to the existence of the solution to System 1a. Thus, applying Lemma 8 with $(f(x), \tilde{f}(x)) = (\text{env}'_{\text{loss}}(\alpha x + Z; \kappa), \text{env}'_{\text{loss}}(\tilde{\alpha} x + Z; \tilde{\kappa}))$ and $(f(x), \tilde{f}(x)) = (\frac{\beta}{\nu} x - \frac{1}{\nu} \text{env}_{\text{reg}}(\frac{\beta}{\nu} x + \Theta; \frac{1}{\nu}), \frac{\tilde{\beta}}{\nu} x - \frac{1}{\nu} \text{env}_{\text{reg}}(\frac{\tilde{\beta}}{\nu} x + \Theta; \frac{1}{\nu}))$, we find that F_{loss} and F_{reg} are differentiable and the derivatives are given by:

$$\begin{aligned} F'_{\text{loss}}(\eta_G) &= \frac{c\tilde{c}\tilde{\delta}}{\beta\tilde{\beta}} \cdot \mathbb{E}[\alpha \text{env}''_{\text{loss}}(\alpha G + Z; \kappa) \cdot \tilde{\alpha} \text{env}''_{\text{loss}}(\tilde{\alpha} \tilde{G} + Z; \tilde{\kappa})] \\ F'_{\text{reg}}(\eta_H) &= \frac{1}{\alpha\tilde{\alpha}} \mathbb{E}\left[\frac{\beta}{\nu} \left(1 - \frac{1}{\nu} \cdot \text{env}''_{\text{reg}}\left(\frac{\beta}{\nu} H + \Theta; \frac{1}{\nu}\right)\right) \cdot \frac{\tilde{\beta}}{\nu} \left(1 - \frac{1}{\nu} \cdot \text{env}''_{\text{reg}}\left(\frac{\tilde{\beta}}{\nu} \tilde{H} + \Theta; \frac{1}{\nu}\right)\right)\right] \end{aligned}$$

for all $\eta_G, \eta_H \in [-1, 1]$. Note that the non-expansiveness of the proximal operator implies that the map $x \mapsto \text{env}'_f(x; \tau) = \tau^{-1}(x - \text{prox}_f(x; \tau))$ is τ^{-1} -Lipschitz and non-decreasing for any convex function f . Thus, $F'_{\text{loss}}(\eta_G)$ and $F'_{\text{reg}}(\eta_H)$ are non-negative and uniformly bounded from above as follows:

$$\begin{aligned} 0 \leq F'_{\text{loss}}(\eta_G) &\leq \frac{c\tilde{c}\tilde{\delta}}{\beta\tilde{\beta}} \frac{\alpha}{\kappa} \cdot \mathbb{E}[\tilde{\alpha} \text{env}''_{\text{loss}}(\tilde{\alpha} \tilde{G} + Z; \tilde{\kappa})] && (0 \leq \text{env}''_{\text{loss}}(\alpha G + Z; \kappa) \leq \kappa^{-1}) \\ &= \frac{c\tilde{c}\tilde{\delta}}{\beta\tilde{\beta}} \frac{\alpha}{\kappa} \cdot \mathbb{E}[\tilde{G} \cdot \text{env}'_{\text{loss}}(\tilde{\alpha} \tilde{G} + Z; \tilde{\kappa})] && (\text{by Stein's lemma}) \\ &= \frac{c\tilde{c}\tilde{\delta}}{\beta\tilde{\beta}} \frac{\alpha}{\kappa} \cdot \frac{\tilde{\nu}\tilde{\alpha}}{\tilde{c}\tilde{\delta}} = c \cdot \frac{\alpha\tilde{\alpha}\tilde{\nu}}{\beta\tilde{\beta}\kappa} && (\text{using System 1a}); \\ 0 \leq F'_{\text{reg}}(\eta_H) &\leq \frac{1}{\alpha\tilde{\alpha}} \mathbb{E}\left[\frac{\beta}{\nu} \left(1 - \frac{1}{\nu} \cdot \text{env}''_{\text{reg}}\left(\frac{\beta}{\nu} H + \Theta; \frac{1}{\nu}\right)\right)\right] \cdot \frac{\tilde{\beta}}{\nu} && (0 \leq \text{env}''_{\text{reg}}\left(\frac{\tilde{\beta}}{\nu} \tilde{H} + \Theta; \frac{1}{\nu}\right) \leq \tilde{\nu}) \\ &= \frac{1}{\alpha\tilde{\alpha}} \frac{\tilde{\beta}}{\tilde{\nu}} \mathbb{E}\left[\left(\frac{\beta}{\nu} H - \frac{1}{\nu} \cdot \text{env}'_{\text{reg}}\left(\frac{\beta}{\nu} H + \Theta; \frac{1}{\nu}\right)\right)\right] && (\text{by Stein's lemma}) \\ &= \frac{1}{\alpha\tilde{\alpha}} \frac{\tilde{\beta}}{\tilde{\nu}} \cdot \beta\kappa = \frac{\beta\tilde{\beta}\kappa}{\alpha\tilde{\alpha}\tilde{\nu}} && (\text{using System 1a}). \end{aligned}$$

Thus, by the chain rule, noting that $(\beta\tilde{\beta}\kappa)/(\alpha\tilde{\alpha}\tilde{\nu})$ is cancelled out, we have

$$0 \leq (F_{\text{loss}} \circ F_{\text{reg}})'(\eta_H) \leq c, \quad 0 \leq (F_{\text{reg}} \circ F_{\text{loss}})'(\eta_G) \leq c.$$

By switching the role of $(\alpha, \beta, \kappa, \nu, c)$ and $(\tilde{\alpha}, \tilde{\beta}, \tilde{\kappa}, \tilde{\nu}, \tilde{c})$, it also holds that

$$0 \leq (F_{\text{loss}} \circ F_{\text{reg}})'(\eta_H) \leq \tilde{c}, \quad 0 \leq (F_{\text{reg}} \circ F_{\text{loss}})'(\eta_G) \leq \tilde{c}.$$

Thus, by taking the minimum of (c, \tilde{c}) , we find that the compositions

$$\eta_H \mapsto F_{\text{loss}} \circ F_{\text{reg}}(\eta_H), \quad \eta_G \mapsto F_{\text{reg}} \circ F_{\text{loss}}(\eta_G)$$

are $(c \wedge \tilde{c})$ -Lipschitz.

Part 3. Now let us show the uniqueness and existence of the solution to the system:

$$\eta_H = F_{\text{loss}}(\eta_G), \quad \eta_G = F_{\text{reg}}(\eta_H).$$

By the assumption $\min\{c, \tilde{c}\} < 1$ and the fact that $\eta_H \mapsto F_{\text{loss}} \circ F_{\text{reg}}(\eta_H)$ is $\min\{c, \tilde{c}\}$ -Lipschitz, Banach's fixed-point theorem implies that the fixed-point equation

$$\eta_G = F_{\text{reg}} \circ F_{\text{loss}}(\eta_G), \quad \eta_G \in [-1, 1]$$

admits a unique solution $\eta_G^* \in [-1, 1]$. If we take $\eta_H^* = F_{\text{loss}}(\eta_G^*) \in [-\sqrt{c\tilde{c}}, \sqrt{c\tilde{c}}]$, then we have $F_{\text{reg}}(\eta_H^*) = F_{\text{reg}} \circ F_{\text{loss}}(\eta_G^*) = \eta_G^*$ so that $(\eta_G, \eta_H) = (\eta_G^*, \eta_H^*)$ satisfies the system: $\eta_H = F_{\text{loss}}(\eta_G)$ and $\eta_G = F_{\text{reg}}(\eta_H)$. This proves the existence of the solution to the system. Let us show the uniqueness. Suppose (η_G, η_H) and $(\tilde{\eta}_G, \tilde{\eta}_H)$ satisfy the system. Then, η_G and $\tilde{\eta}_G$ are the solution to the fixed-point equation $\eta = F_{\text{reg}} \circ F_{\text{loss}}(\eta)$:

$$\eta_G = F_{\text{reg}}(\eta_H) = F_{\text{reg}} \circ F_{\text{loss}}(\eta_G), \quad \tilde{\eta}_G = F_{\text{reg}}(\tilde{\eta}_H) = F_{\text{reg}} \circ F_{\text{loss}}(\tilde{\eta}_G).$$

By the uniqueness of the solution to this fixed-point equation, we must have $\eta_G = \tilde{\eta}_G$. By the same argument, the contraction $\|F_{\text{loss}} \circ F_{\text{reg}}\|_{\text{Lip}} \leq c\tilde{c}$ implies that the solution to the fixed-point equation $\eta_H = F_{\text{loss}} \circ F_{\text{reg}}(\eta_H)$ is unique so that we must have $\eta_H = \tilde{\eta}_H$. This proves $(\eta_G, \eta_H) = (\tilde{\eta}_G, \tilde{\eta}_H)$ and hence the systems admit a unique solution.

Let us show $|\eta_G^*| < 1$. We proceed by contradiction. If $|\eta_G^*| = 1$ then noting $\eta_H^* = F_{\text{reg}}(\eta_G^*)$, since we have shown in part 1 that $|F_{\text{reg}}(\eta_H)| < 1$ for all $|\eta_H| < 1$, we must have $|\eta_H^*| = 1$. However, $|\eta_H^*| = |F_{\text{loss}}(\eta_G^*)| \leq \sqrt{c\tilde{c}} < 1$ so that this is a contradiction. Thus, we must have $|\eta_G^*| < 1$. Noting $\eta_H^* = F_{\text{loss}}(\eta_G^*)$ and $|F_{\text{loss}}(\eta_G)| < \sqrt{c\tilde{c}}$ for all $|\eta_G| < 1$, the strict inequality $|\eta_G^*| < 1$ in turn gives $|\eta_H^*| < \sqrt{c\tilde{c}}$.

Part 4. Finally, let us characterize the sign of (η_G^*, η_H^*) . Recall that η_G^* and η_H^* satisfies the fixed-point equations $\eta_G - F_{\text{reg}} \circ F_{\text{loss}}(\eta_G) = 0$ and $\eta_H - F_{\text{loss}} \circ F_{\text{reg}}(\eta_H) = 0$, respectively, and the maps $\eta_G \mapsto \eta_G - F_{\text{reg}} \circ F_{\text{loss}}(\eta_G)$ and $\eta_H \mapsto \eta_H - F_{\text{loss}} \circ F_{\text{reg}}(\eta_H)$ are strictly increasing since $F_{\text{reg}} \circ F_{\text{loss}}$ and $F_{\text{loss}} \circ F_{\text{reg}}$ are $(c \wedge \tilde{c})$ -Lipschitz with $c \wedge \tilde{c} < 1$. Then, the characterization of sign, i.e., $\text{sign}(F_{\text{loss}} \circ F_{\text{reg}}(0)) = \text{sign}(\eta_H^*)$ and $\text{sign}(F_{\text{reg}} \circ F_{\text{loss}}(0)) = \text{sign}(\eta_G^*)$, immediately follows.

A.2 Proof of Theorem 2

In this section, for a vector \mathbf{w} , the norm $\|\mathbf{w}\|$ indicates the ℓ_2 norm unless specified otherwise. By the assumption $0 \in \text{argmin}_x \text{loss}(x) \cap \text{argmin}_x \text{reg}(x)$, taking $\text{loss}_{\text{new}}(x) = \text{loss}(x) - \text{loss}(0)$ and $\text{reg}_{\text{new}}(x) = \text{reg}(x) - \text{reg}(0)$ if necessary, we assume without loss of generality that loss and reg are non-negative and have 0 as a minimizer.

By the change of variable $\mathbf{b} \mapsto \mathbf{h} = (\mathbf{b} - \boldsymbol{\theta})/\sqrt{p}$, denoting $\mathbf{G} = \sqrt{p}\mathbf{X}$ so that \mathbf{G} has i.i.d. $\mathcal{N}(0, 1)$ entries, the regularized M-estimator $\hat{\boldsymbol{\theta}}$ of interest and the residual vector $\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\theta}}$ can be written as

$$\hat{\boldsymbol{\theta}} = \sqrt{p}\hat{\mathbf{h}} + \boldsymbol{\theta}, \quad \mathbf{y} - \mathbf{X}\hat{\boldsymbol{\theta}} = \mathbf{z} - \mathbf{G}\hat{\mathbf{h}}$$

where

$$\hat{\mathbf{h}} \in \underset{\mathbf{h} \in \mathbb{R}^p}{\text{argmin}} \text{obj}(\mathbf{h}) \quad \text{with} \quad \text{obj}(\mathbf{h}) := \sum_{i \in I} \text{loss}(z_i - \mathbf{g}_i^\top \mathbf{h}) + \sum_{j \in [p]} \text{reg}(\sqrt{p}h_j + \theta_j). \quad (27)$$

Throughout this section, we denote $\boldsymbol{\psi} = \sum_{i \in I} \mathbf{e}_i \text{loss}'(z_i - \mathbf{g}_i^\top \mathbf{h}) \in \mathbb{R}^n$ where \mathbf{e}_i are canonical basis of \mathbb{R}^n . Let $\tilde{\boldsymbol{\psi}}, \tilde{\boldsymbol{\theta}}, \tilde{\mathbf{h}}$ be the corresponding notation for another. Then, our goal is to show

$$\hat{\mathbf{h}}^\top \tilde{\mathbf{h}} \xrightarrow{P} \alpha \tilde{\eta}_G \quad \text{and} \quad \boldsymbol{\psi}^\top \tilde{\boldsymbol{\psi}}/p \xrightarrow{P} \beta \tilde{\eta}_H.$$

A.2.1 Smoothing by adding diminishing ridge penalty

For some positive and diminishing scalar $\mu = \mu_n \rightarrow 0$ to be specified later, we define the smoothed regularized M-estimator $\widehat{\mathbf{h}}_\mu$ as

$$\widehat{\mathbf{h}}_\mu \in \operatorname{argmin}_{\mathbf{h} \in \mathbb{R}^p} \operatorname{obj}(\mathbf{h}) + \frac{p\mu}{2} \|\mathbf{h}\|^2 \quad (28)$$

where $\operatorname{obj}(\mathbf{h})$ is the objective function for the original regularized M-estimator $\widehat{\mathbf{h}}$ in (27). We denote by $\boldsymbol{\psi}_\mu = \sum_{i \in I} \mathbf{e}_i \operatorname{loss}'(z_i - \mathbf{g}_i^\top \mathbf{h}_\mu)$ the smoothed version of $\boldsymbol{\psi}$. This strategy of adding and additive ridge penalty for the mathematical analysis as a first step, and then arguing by approximation ($\mu \rightarrow 0$) to study the original $\widehat{\mathbf{h}}_0 = \widehat{\mathbf{h}}$ is ubiquitous; see, for instance, [Kar18, CM22], [LGR⁺22, Appendix B.3], [BK23a]. Here, we use the following lemma to show that (28) approximates well the original estimator for any sequence $\mu = \mu_n$ indexed by n and converging to 0.

Lemma 9 (Ridge smoothing). *Let \mathbf{h}_μ be the smoothed M-estimator and $\widehat{\mathbf{h}}$ be the original M-estimator.*

1. (Monotonicity) *We have $\|\widehat{\mathbf{h}} - \widehat{\mathbf{h}}_\mu\|_2^2 \leq \|\widehat{\mathbf{h}}\|_2^2 - \|\widehat{\mathbf{h}}_\mu\|_2^2$ for any $\mu > 0$.*
2. (Convergence in $\|\cdot\|_2$) *As $n, p \rightarrow +\infty$ with $n/p \rightarrow \delta$, $|I|/n \rightarrow c$, and $\mu = \mu_n$ for any $\mu_n \rightarrow 0$, we have*

$$\|\widehat{\mathbf{h}}_\mu\|_2^2 \xrightarrow{\mathbb{P}} \alpha^2, \quad \frac{\|\boldsymbol{\psi}_\mu\|_2^2}{p} \xrightarrow{\mathbb{P}} \beta^2, \quad \|\widehat{\mathbf{h}}_\mu - \widehat{\mathbf{h}}\|_2^2 = o_{\mathbb{P}}(1), \quad \frac{\|\boldsymbol{\psi}_\mu - \boldsymbol{\psi}\|_2^2}{p} = o_{\mathbb{P}}(1),$$

where α and β are solutions to System 1a.

Proof. By the strong convexity of the ridge term $p\mu^2/2\|\mathbf{h}\|^2$, the smoothed one $\widehat{\mathbf{h}}_\mu$ also minimizes the convex function: $\mathbf{h} \mapsto \operatorname{obj}(\mathbf{h}) + \frac{p\mu}{2}\|\mathbf{h}\|^2 - \frac{p\mu}{2}\|\mathbf{h} - \widehat{\mathbf{h}}_\mu\|^2$. By the optimality of $\widehat{\mathbf{h}}_\mu$ and $\widehat{\mathbf{h}}$, we have

$$\operatorname{obj}(\widehat{\mathbf{h}}_\mu) + \frac{p\mu}{2}\|\widehat{\mathbf{h}}_\mu\|^2 \leq \operatorname{obj}(\widehat{\mathbf{h}}) + \frac{p\mu}{2}\|\widehat{\mathbf{h}}\|^2 - \frac{p\mu}{2}\|\widehat{\mathbf{h}} - \widehat{\mathbf{h}}_\mu\|^2 \leq \operatorname{obj}(\widehat{\mathbf{h}}_\mu) + \frac{p\mu}{2}\|\widehat{\mathbf{h}}\|^2 - \frac{p\mu}{2}\|\widehat{\mathbf{h}} - \widehat{\mathbf{h}}_\mu\|^2$$

so that subtracting $\operatorname{obj}(\widehat{\mathbf{h}}_\mu)$ from the both sides and dividing by $p\mu/2 > 0$, we obtain the first claim.

Recall $\|\mathbf{h}\|^2 \xrightarrow{\mathbb{P}} \alpha^2$ and $\|\boldsymbol{\psi}\|^2/p \xrightarrow{\mathbb{P}} \beta^2$. Then, by Lemma 9-(1) and the Lipschitz condition of loss' , it suffices to show the convergence $\|\widehat{\mathbf{h}}_\mu\|^2 \xrightarrow{\mathbb{P}} \alpha^2$ for the smoothed estimator $\widehat{\mathbf{h}}_\mu$. Note that $\widehat{\mathbf{h}}_\mu$ minimizes the function below

$$\mathbf{h} \mapsto \sum_{i \in I} \operatorname{loss}(z_i - \mathbf{g}_i^\top \mathbf{h}) + \sum_{j \in [p]} \operatorname{reg}_j^{\mu_n}(\sqrt{p}h_j) \quad \text{where} \quad \operatorname{reg}_j^\mu(x) := \operatorname{reg}(x + \theta_j) + \mu \frac{x^2}{2}.$$

Now we suppose that for any standard normal $\mathbf{g} = (g_j)_{j=1}^p \sim \mathcal{N}(\mathbf{0}_p, \mathbf{I}_p)$, the convergence of the Moreau envelope

$$\frac{1}{p} \sum_{j \in [p]} \operatorname{env}_{\operatorname{reg}_j^{\mu_n}}(cg_j; \tau) - \operatorname{reg}_j^{\mu_n}(0) \xrightarrow{\mathbb{P}} \mathbb{E}[\operatorname{env}_{\operatorname{reg}}(cH + \Theta; \tau) - \operatorname{reg}(\Theta)]$$

holds for all $c \in \mathbb{R}$ and $\tau > 0$. Then, by [TAH18, Theorem 3.1], we have $\|\widehat{\mathbf{h}}_\mu\|^2 \xrightarrow{\mathbb{P}} \alpha^2$ and complete the proof. By the weak law of large numbers, the above display holds with $\mu = 0$ (see [TAH18, Lemma 4.1] for details):

$$\frac{1}{p} \sum_{j \in [p]} \operatorname{env}_{\operatorname{reg}_j^{\mu=0}}(cg_j; \tau) - \operatorname{reg}_j^{\mu=0}(0) \xrightarrow{\mathbb{P}} \mathbb{E}[\operatorname{env}_{\operatorname{reg}}(cH + \Theta; \tau) - \operatorname{reg}(\Theta)].$$

Then, noting $\text{reg}_j^{\mu=0}(0) = \text{reg}(\theta_j) = \text{reg}_j^\mu(0)$ for any $\mu \geq 0$, it suffices to show

$$\frac{1}{p} \sum_{j \in [p]} \text{env}_{\text{reg}_j^{\mu_n}}(cg_j; \tau) - \frac{1}{p} \sum_{j \in [p]} \text{env}_{\text{reg}_j^{\mu=0}}(cg_j; \tau) = o_{\mathbb{P}}(1).$$

For each j , the monotonicity $\text{env}_{\text{reg}_j^{\mu=0}}(cg_j; \tau) \leq \text{env}_{\text{reg}_j^{\mu_n}}(cg_j; \tau)$ holds since the objective function is monotone in the sense of $\text{reg}_j^{\mu=0}(x) \leq \text{reg}_j^\mu(x)$ for all $x \in \mathbb{R}$ and all $\mu \geq 0$. On the other hand, by the optimality of $\text{prox}_{\text{reg}_j^{\mu=0}}(cg_j; \tau)$ and $\text{prox}_{\text{reg}_j^{\mu_n}}(cg_j; \tau)$, we find

$$\begin{aligned} \text{env}_{\text{reg}_j^{\mu_n}}(cg_j; \tau) &\leq \frac{1}{2\tau} (cg_j - \text{prox}_{\text{reg}_j^{\mu=0}}(cg_j; \tau))^2 + \text{reg}_j^{\mu_n}(\text{prox}_{\text{reg}_j^{\mu=0}}(cg_j; \tau)) \\ &= \text{env}_{\text{reg}_j^{\mu=0}}(cg_j; \tau) + \frac{\mu_n}{2} (\text{prox}_{\text{reg}_j^{\mu=0}}(cg_j; \tau))^2 \end{aligned}$$

for each $j \in [p]$. Thus, it holds that

$$0 \leq \frac{1}{p} \sum_{j \in [p]} \text{env}_{\text{reg}_j^{\mu_n}}(cg_j; \tau) - \frac{1}{p} \sum_{j \in [p]} \text{env}_{\text{reg}_j^{\mu=0}}(cg_j; \tau) \leq \frac{\mu_n}{2} \cdot \frac{1}{p} \sum_{j \in [p]} (\text{prox}_{\text{reg}_j^{\mu=0}}(cg_j; \tau))^2,$$

where $\text{prox}_{\text{reg}_j^{\mu=0}}(cg_j; \tau) = \text{prox}_{\text{reg}}(cg_j + \theta_j; \tau) - \theta_j$ by the definition of reg_j^μ . Here, under Assumption D-(1), the expectation $\mathbb{E}[(\text{prox}_{\text{reg}}(cg_j + \theta_j; \tau) - \theta_j)^2]$ under $g_j \sim \mathcal{N}(0, 1)$ and $\theta_j \sim \Theta$ is finite for all $c \in \mathbb{R}$ and $\tau > 9$ (cf. [TAH18, equation (123)]). Therefore, by the weak law of large numbers, we have $\frac{1}{p} \sum_{j \in [p]} (\text{prox}_{\text{reg}_j^{\mu=0}}(cg_j; \tau))^2 \xrightarrow{P} \mathbb{E}[(\text{prox}_{\text{reg}}(cH + \Theta; \tau) - \Theta)^2] < +\infty$. This means that for any $\mu = \mu_n \rightarrow 0$, the RHS of the previous display is $o_{\mathbb{P}}(1)$. \square

A.2.2 Bounding norm of regularized M-estimators

For some positive scalar K to be specified later, we add another regularization term; we define $\hat{\mathbf{h}}_{\mu, K}$ as

$$\hat{\mathbf{h}}_{\mu, K} \in \underset{\mathbf{h} \in \mathbb{R}^p}{\text{argmin}} \text{obj}(\mathbf{h}) + \frac{p\mu}{2} \|\mathbf{h}\|^2 + \underbrace{\frac{\hat{\lambda}}{2} \mathbb{F}\left(\frac{\|\mathbf{h}\|^2 - K}{2}\right)}_{\text{additional term}} \quad \text{with} \quad \hat{\lambda} := \text{obj}(\mathbf{0})$$

where $\mathbb{F} : \mathbb{R} \rightarrow \mathbb{R}$ is convex, non-negative, and non-decreasing with $\lim_{x \rightarrow +\infty} \mathbb{F}(x) = +\infty$, as well as differentiable with $\mathbb{F}'(u) = 0$ if $u \leq 0$ and $\mathbb{F}'(u) = 1$ if $u \geq 1$. For instance, we may take \mathbb{F} as an integral of the smoothed step function as follows:

$$\mathbb{F}(x) := \int_{-\infty}^x f(u) \, du \quad \text{with} \quad f(u) := \begin{cases} 1 & u \geq 1 \\ 3u^2 - 2u & u \in (0, 1) \\ 0 & u \leq 0 \end{cases}$$

Note in passing that the regularization parameter $\hat{\lambda} = \text{obj}(\mathbf{0}) = \sum_{i \in I} \text{loss}(z_i) + \sum_{j \in [p]} \text{reg}(\theta_j)$ is independent of the design matrix \mathbf{G} and non-negative since loss and reg are non-negative.

Now we claim that for sufficiently large $K > 0$ this modified estimator $\hat{\mathbf{h}}_{\mu, K}$ coincides with the smoothed estimator $\hat{\mathbf{h}}_\mu$. Furthermore, thanks to the additional regularizer $\frac{\hat{\lambda}}{2} \mathbb{F}\left(\frac{\|\mathbf{h}\|^2 - K}{2}\right)$, the norm of $\mathbf{h}_{\mu, K}$ and $\boldsymbol{\psi}_{\mu, K}$ are suitably bounded as follows:

Lemma 10. *The following convergences hold.*

1. *If we set $K = 2\alpha^2$ where $\alpha > 0$ is the solution to System 1a, we have*

$$\mathbb{P}(\widehat{\mathbf{h}}_{\mu,K} = \widehat{\mathbf{h}}_{\mu}) \xrightarrow{\mathbb{P}} 1.$$

2. *For any $K \geq 0$, there exists a positive constant C_K that only depends on K such that*

$$\|\widehat{\mathbf{h}}_{\mu,K}\|^2 \leq C_K, \quad \|\boldsymbol{\psi}_{\mu,K}\|^2 \leq C_K(1 + \|\text{loss}'\|_{\text{Lip}}^2)(\|\text{loss}'(\mathbf{z})\|^2 + \|\mathbf{G}\|_{\text{op}}^2).$$

Thus, the constant $C^ = 1.1C_K(1 + \|\text{loss}'\|_{\text{Lip}}^2)(\mathbb{E}[\text{loss}'(Z)^2] + (1 + \delta^{-1/2}))^2$ satisfies*

$$\mathbb{P}\left(\|\widehat{\mathbf{h}}_{\mu,K}\|^2 \vee (n^{-1}\|\boldsymbol{\psi}_{\mu,K}\|^2) \vee \sup_{m \geq 1} \mathbb{E}\left[\|\widehat{\mathbf{h}}_{\mu,K}\|^{2m} \vee (n^{-1}\|\boldsymbol{\psi}_{\mu,K}\|^2)^m \mid \mathbf{z}, \boldsymbol{\theta}\right]^{1/m} \leq C^*\right) \rightarrow 1.$$

Proof. Let us consider the event $\Omega := \{\|\mathbf{h}\|^2 \leq K\}$ with $K = 2\alpha^2$, which holds with high probability, since $\|\mathbf{h}\|^2 \xrightarrow{\mathbb{P}} \alpha^2 > 0$. Combined with the monotonicity $\|\mathbf{h}_{\mu}\|^2 \leq \|\mathbf{h}\|^2$ by Lemma 9-(1), we have $\|\widehat{\mathbf{h}}_{\mu}\|^2 \leq K$ under the event Ω . Since $\mathbf{F}'(u) = 0$ for all $u \leq 0$, combined with the Karush-Kuhn-Tucker (KKT) condition $-p\mu\widehat{\mathbf{h}}_{\mu} \in \partial\text{obj}(\widehat{\mathbf{h}}_{\mu})$ for the smoothed estimator $\widehat{\mathbf{h}}_{\mu}$, we observe that $\widehat{\mathbf{h}}_{\mu}$ satisfies the KKT condition for $\widehat{\mathbf{h}}_{\mu,K}$ under the event Ω :

$$-p\mu\widehat{\mathbf{h}}_{\mu} - \widehat{\lambda}\mathbf{F}'\left(\frac{\|\widehat{\mathbf{h}}_{\mu}\|^2 - K}{2}\right)\widehat{\mathbf{h}}_{\mu} = -p\mu\widehat{\mathbf{h}}_{\mu} - 0 \cdot \widehat{\mathbf{h}}_{\mu} \in \partial\text{obj}(\widehat{\mathbf{h}}_{\mu}).$$

This implies $\mathbb{P}(\widehat{\mathbf{h}}_{\mu} = \widehat{\mathbf{h}}_{\mu,K}) \geq \mathbb{P}(\Omega) \rightarrow 1$.

By the non-negativity of $(\text{obj}(\cdot), p\mu\|\cdot\|^2, \frac{\widehat{\lambda}}{2}\mathbf{F}(\cdot))$ and the optimality of $\widehat{\mathbf{h}}_{\mu,K}$, it holds that

$$\begin{aligned} 0 &\leq \frac{\widehat{\lambda}}{2}\mathbf{F}\left(\frac{\|\widehat{\mathbf{h}}_{\mu,K}\|^2 - K}{2}\right) \leq \text{obj}(\widehat{\mathbf{h}}_{\mu,K}) + \frac{p\mu}{2}\|\mathbf{h}_{\mu,K}\|^2 + \frac{\widehat{\lambda}}{2}\mathbf{F}\left(\frac{\|\mathbf{h}_{\mu,K}\|^2 - K}{2}\right) \\ &\leq \text{obj}(\mathbf{0}) + \frac{p\mu}{2}\|\mathbf{0}\|^2 + \frac{\widehat{\lambda}}{2}\mathbf{F}\left(\frac{\|\mathbf{0}\|^2 - K}{2}\right) \\ &= \widehat{\lambda} + 0 + 0. \end{aligned}$$

When $\widehat{\lambda} = 0$ then all inequalities above holds with equality, which means that $\mathbf{0}$ minimizes the objective function $\text{obj}(\mathbf{h}) + \frac{p\mu}{2}\|\mathbf{h}\|^2 + \frac{\widehat{\lambda}}{2}\mathbf{F}\left(\frac{\|\mathbf{h}\|^2 - K}{2}\right)$ for $\widehat{\mathbf{h}}_{\mu,K}$. Since this objective function is strongly convex due to the ridge term, the minimizer is unique. This means $\widehat{\mathbf{h}}_{\mu,K} = \mathbf{0}$ when $\widehat{\lambda} = 0$. On the other hand, if $\widehat{\lambda} > 0$ then dividing the above display by $\widehat{\lambda} > 0$ we are left with $\mathbf{F}\left(\frac{\|\widehat{\mathbf{h}}_{\mu,K}\|^2 - K}{2}\right) \leq 2$. Since \mathbf{F} is non-decreasing and coercive on the positive side, i.e., $\lim_{x \rightarrow +\infty} \mathbf{F}(x) \rightarrow +\infty$, this gives $\|\widehat{\mathbf{h}}_{\mu,K}\|^2 \leq C(K)$ for a constant $C(K) > 0$ depending on K only. Combined with the Lipschitz condition of loss' , the norm of $\boldsymbol{\psi}_{\mu,K} = \sum_{i \in I} \mathbf{e}_i \text{loss}'(z_i - \mathbf{g}_i^{\top} \mathbf{h}_{\mu,K})$ is bounded as

$$\|\boldsymbol{\psi}_{\mu,K}\| \leq \|\text{loss}'(\mathbf{z})\| + \|\text{loss}'(\mathbf{z} - \mathbf{G}\widehat{\mathbf{h}}_{\mu,K}) - \text{loss}'(\mathbf{z})\| \leq \|\text{loss}'(\mathbf{z})\| + \|\text{loss}'\|_{\text{Lip}}\|\mathbf{G}\|_{\text{op}}\|\widehat{\mathbf{h}}_{\mu,K}\|.$$

Since $\text{loss}'(z_i)^2$ has a finite second moment by Assumption D-(1), the weak law of large numbers gives $n^{-1}\|\text{loss}'(\mathbf{z})\|^2 \xrightarrow{\mathbb{P}} \mathbb{E}[\text{loss}'(Z)^2] < +\infty$. Since $\mathbf{G} \in \mathbb{R}^{n \times p}$ is a Gaussian matrix with i.i.d. $\mathcal{N}(0, 1)$ entries, we also have $\|\mathbf{G}\|_{\text{op}}^2/n \xrightarrow{\mathbb{P}} (1 + \delta^{-1/2})^2$ by standard results of the maximal singular value of a Gaussian matrix, e.g., [DS01, Theorem II.13]. This completes the proof. \square

A.2.3 Derivative formulae for strongly convex regularizer

Lemma 11 ([BS22]). *Let $\hat{\mathbf{h}} \in \mathbb{R}^p$ be the regularized M -estimator of the form*

$$\hat{\mathbf{h}} \in \operatorname{argmin}_{\mathbf{h} \in \mathbb{R}^p} \sum_{i \in I} \operatorname{loss}(z_i - \mathbf{g}_i^\top \mathbf{h}) + \mathbf{R}(\mathbf{h}),$$

where $I \subset [n]$ is a subset independent of $(\mathbf{g}_i)_{i \in [n]}$, $\operatorname{loss} : \mathbb{R} \rightarrow \mathbb{R}$ is a convex and differentiable function with Lipschitz derivative, $\mathbf{R} : \mathbb{R}^p \rightarrow \mathbb{R}$ is a τ -strongly convex regularizer for some $\tau > 0$, and \mathbf{R} is independent from (\mathbf{g}_i) . Then, there exists a matrix $\mathbf{A} \in \mathbb{R}^{p \times p}$ satisfying $\|\mathbf{A}\|_{\text{op}} \leq \tau^{-1}$ and $\operatorname{tr}[\mathbf{A}] \geq 0$ such that $\hat{\mathbf{h}} \in \mathbb{R}^p$ and $\boldsymbol{\psi} = \sum_{i \in I} \mathbf{e}_i \operatorname{loss}'(z_i - \mathbf{g}_i^\top \hat{\mathbf{h}}) \in \mathbb{R}^n$ are both differentiable as functions of the design $(\mathbf{G}) = (g_{ij})$ with the derivative given by

$$\forall i \in [n], \forall j \in [p], \quad \frac{\partial \hat{\mathbf{h}}}{\partial g_{ij}} = \mathbf{A} \left(\mathbf{e}_j \mathbf{e}_i^\top \boldsymbol{\psi} - \mathbf{G}^\top \mathbf{D} \mathbf{e}_i \mathbf{e}_j^\top \hat{\mathbf{h}} \right), \quad \frac{\partial \boldsymbol{\psi}}{\partial g_{ij}} = -\mathbf{D} \mathbf{G} \mathbf{A} \mathbf{e}_j \mathbf{e}_i^\top \boldsymbol{\psi} - \mathbf{V} \mathbf{e}_i \mathbf{e}_j^\top \hat{\mathbf{h}}. \quad (29)$$

where \mathbf{D} and \mathbf{V} are $n \times n$ matrices defined by $\mathbf{D} = \sum_{i \in I} \mathbf{e}_i \mathbf{e}_i^\top \operatorname{loss}''(z_i - \mathbf{g}_i^\top \hat{\mathbf{h}})$ and $\mathbf{V} = \mathbf{D} - \mathbf{D} \mathbf{G} \mathbf{A} \mathbf{G}^\top \mathbf{D}$. Furthermore, the matrix \mathbf{V} defined above is positive semidefinite with its operator norm bounded as $\|\mathbf{V}\|_{\text{op}} \leq \|\operatorname{loss}'\|_{\text{Lip}}$.

Recall that the modified estimator $\hat{\mathbf{h}}_{\mu, K}$ in Appendix A.2.2 minimizes the objective function $\sum_{i \in I} \operatorname{loss}(z_i - \mathbf{g}_i^\top \mathbf{h}) + \mathbf{R}(\mathbf{h})$ where \mathbf{R} is the regularizer of the form

$$\mathbf{R}(\mathbf{h}) := \sum_{j \in [p]} \operatorname{reg}(\sqrt{p} h_j + \theta_j) + \frac{p\mu}{2} \|\mathbf{h}\|^2 + \frac{\hat{\lambda}}{2} \mathbf{F} \left(\frac{\|\mathbf{h}\|^2 - K}{2} \right),$$

which is $(p\mu)$ -strongly convex. Thus, we can apply Lemma 11; there exists a matrix $\mathbf{A}_{\mu, K} \in \mathbb{R}^{p \times p}$ such that

$$\|\mathbf{A}_{\mu, K}\|_{\text{op}} \leq (p\mu)^{-1}, \quad \operatorname{tr}[\mathbf{A}_{\mu, K}] \geq 0, \quad (30)$$

and $\hat{\mathbf{h}}_{\mu, K}$ and $\boldsymbol{\psi}_{\mu, K} = \sum_{i \in I} \mathbf{e}_i \operatorname{loss}'(z_i - \mathbf{g}_i^\top \hat{\mathbf{h}}_{\mu, K})$ are differentiable with respect to the design matrix $\mathbf{G} = (g_{ij})$ as in (29) with

$$\mathbf{D}_{\mu, K} = \sum_{i \in I} \mathbf{e}_i \mathbf{e}_i^\top \operatorname{loss}''(z_i - \mathbf{g}_i^\top \hat{\mathbf{h}}_{\mu, K}) \quad \mathbf{V}_{\mu, K} = \mathbf{D}_{\mu, K} - \mathbf{D}_{\mu, K} \mathbf{G} \mathbf{A}_{\mu, K} \mathbf{G}^\top \mathbf{D}_{\mu, K}.$$

Now we claim that the trace of $\mathbf{V}_{\mu, K}$ and $\mathbf{A}_{\mu, K}$ are empirical quantities that converge to the remaining solution ν and κ :

Lemma 12. *For any $\mu \rightarrow 0$ such that $\mu^{-1} = O(n^{1/8})$, we have*

$$\operatorname{tr}[\mathbf{V}_{\mu, K}]/p \xrightarrow{P} \nu \quad \text{and} \quad \operatorname{tr}[\mathbf{A}_{\mu, K}] \xrightarrow{P} \kappa.$$

Proof. See Appendix A.4. □

A.2.4 Proof of Theorem 2

Below we will take the diminishing constant $\mu_n \rightarrow 0$ as in Lemma 12. Dropping the dependence on (μ, K) for simplicity, we denote $\mathbf{h} = \hat{\mathbf{h}}_{\mu, K}$, $\boldsymbol{\psi} = \boldsymbol{\psi}_{\mu, K}$, and $\mathbf{V} = \mathbf{V}_{\mu, K}$, $\mathbf{A} = \mathbf{A}_{\mu, K}$. In the same way, we use the notation $(\tilde{\mathbf{h}}, \tilde{\boldsymbol{\psi}}, \tilde{\mathbf{V}}, \tilde{\mathbf{A}})$ for the other estimator. Then, by the convergence in $\|\cdot\|_2$ from Lemma 9 and Lemma 10, it suffices to show the convergence of correlation for the modified ones:

$$\mathbf{h}^\top \tilde{\mathbf{h}} / (\alpha \tilde{\alpha}) \xrightarrow{P} \eta_G \quad \text{and} \quad \boldsymbol{\psi}^\top \tilde{\boldsymbol{\psi}} / (p\beta\tilde{\beta}) \xrightarrow{P} \eta_H.$$

First, we will argue by the Gaussian Poincaré inequality that $\mathbf{h}^\top \tilde{\mathbf{h}} / (\alpha \tilde{\alpha})$ and $\boldsymbol{\psi}^\top \tilde{\boldsymbol{\psi}} / (p\beta\tilde{\beta})$ concentrate on random quantities that are independent of the scaled design matrix $\mathbf{G} = \sqrt{p}\mathbf{X}$. Throughout this section, we denote by $\bar{\mathbb{E}}[\cdot] = \mathbb{E}[\cdot | \mathbf{z}, \boldsymbol{\theta}, I, \tilde{I}]$ the conditional expectation with respect to the design matrix \mathbf{G} given signal $\boldsymbol{\theta}$, noise \mathbf{z} and subsample index I, \tilde{I} . Since \mathbf{G} is independent of $(\boldsymbol{\theta}, \mathbf{z}, I, \tilde{I})$, the conditional distribution of \mathbf{G} given $(\boldsymbol{\theta}, \mathbf{z}, I, \tilde{I})$ is still that of a matrix with i.i.d. $\mathcal{N}(0, 1)$ entries.

Lemma 13. *For any $\mu_n > 0$ such that $\mu_n^{-1} = o(n^{1/2})$, we have*

$$\begin{aligned} \bar{\mathbb{E}}[(\mathbf{h}^\top \tilde{\mathbf{h}} - \bar{\mathbb{E}}[\mathbf{h}^\top \tilde{\mathbf{h}}])^2] &= \mathcal{O}_{\mathbb{P}}(n^{-1}\mu^{-2}) = o_{\mathbb{P}}(1). \\ \bar{\mathbb{E}}[(\boldsymbol{\psi}^\top \tilde{\boldsymbol{\psi}} - \bar{\mathbb{E}}[\boldsymbol{\psi}^\top \tilde{\boldsymbol{\psi}}])^2] &= \mathcal{O}_{\mathbb{P}}(n\mu^{-2}) = o_{\mathbb{P}}(n^2). \end{aligned}$$

Proof. By the Gaussian Poincaré inequality, noting $\partial_{ij}(\mathbf{h}^\top \tilde{\mathbf{h}}) = \tilde{\mathbf{h}}^\top \partial_{ij}\mathbf{h} + \mathbf{h}^\top \partial_{ij}\tilde{\mathbf{h}}$ with $\partial_{ij} = \frac{\partial}{\partial g_{ij}}$, the conditional variance of $\mathbf{h}^\top \tilde{\mathbf{h}}$ is bounded from above as

$$\bar{\mathbb{E}}[(\mathbf{h}^\top \tilde{\mathbf{h}} - \bar{\mathbb{E}}[\mathbf{h}^\top \tilde{\mathbf{h}}])^2] \leq \sum_{ij} \bar{\mathbb{E}}[(\partial_{ij}(\mathbf{h}^\top \tilde{\mathbf{h}}))^2] \leq 2 \sum_{ij} \bar{\mathbb{E}}[(\tilde{\mathbf{h}}^\top \partial_{ij}\mathbf{h})^2] + 2 \sum_{ij} \bar{\mathbb{E}}[(\mathbf{h}^\top \partial_{ij}\tilde{\mathbf{h}})^2],$$

where we denote by $\sum_{ij} = \sum_{i=1}^n \sum_{j=1}^p$ for brevity. By the derivative formula (29), the first term on the RHS is bounded as

$$\begin{aligned} \sum_{ij} (\tilde{\mathbf{h}}^\top \partial_{ij}\mathbf{h})^2 &= \sum_{i,j} (\tilde{\mathbf{h}}^\top \mathbf{A}(\mathbf{e}_j \psi_i - \mathbf{G}^\top \mathbf{D} \mathbf{e}_i h_j))^2 \\ &\leq 2\|\mathbf{A}\|_{\text{op}}^2 \|\mathbf{h}\|^2 \|\boldsymbol{\psi}\|^2 + 2\|\mathbf{A}\|_{\text{op}}^2 \|\mathbf{G}\|_{\text{op}}^2 \|\mathbf{D}\|_{\text{op}}^2 \|\tilde{\mathbf{h}}\|^2 \|\mathbf{h}\|^2. \end{aligned}$$

Using the upper bound $\|\mathbf{A}\|_{\text{op}} \leq (p\mu_n)^{-1}$ from (30) and the moment bound in Lemma 10-(2), the conditional expectation (with respect to $\bar{\mathbb{E}}$) of the RHS is $\mathcal{O}_{\mathbb{P}}(n^{-1}\mu_n^{-2})$. By symmetry, we also get $\bar{\mathbb{E}}[\sum_{ij} \bar{\mathbb{E}}[(\mathbf{h}^\top \partial_{ij}\tilde{\mathbf{h}})^2]] = \mathcal{O}_{\mathbb{P}}(n^{-1}\mu_n^{-2})$.

We use a similar argument for the variance of $\boldsymbol{\psi}^\top \tilde{\boldsymbol{\psi}}$. The Gaussian Poincaré gives $\bar{\mathbb{E}}[(\boldsymbol{\psi}^\top \tilde{\boldsymbol{\psi}} - \bar{\mathbb{E}}[\boldsymbol{\psi}^\top \tilde{\boldsymbol{\psi}}])^2] \leq 2 \sum_{ij} \bar{\mathbb{E}}[(\tilde{\boldsymbol{\psi}}^\top \partial_{ij}\boldsymbol{\psi})^2 + (\boldsymbol{\psi}^\top \partial_{ij}\tilde{\boldsymbol{\psi}})^2]$, where

$$\begin{aligned} \sum_{ij} (\tilde{\boldsymbol{\psi}}^\top \partial_{ij}\boldsymbol{\psi})^2 &= \sum_{i,j} (\tilde{\boldsymbol{\psi}}^\top (-\mathbf{D}\mathbf{G}\mathbf{A}\mathbf{e}_j \psi_i - \mathbf{V}\mathbf{e}_i w_j))^2 \\ &\lesssim \|\mathbf{A}^\top \mathbf{G}^\top \mathbf{D}^\top\|_{\text{op}}^2 \|\tilde{\boldsymbol{\psi}}\|^2 \|\boldsymbol{\psi}\|^2 + \|\mathbf{V}\|_{\text{op}}^2 \|\tilde{\boldsymbol{\psi}}\|^2 \|\mathbf{h}\|^2. \end{aligned}$$

by the derivative formula. Using $\|\mathbf{A}\|_{\text{op}} \leq (p\mu_n)^{-1}$ and the moment bound in Lemma 10-(2), the conditional expectation $\bar{\mathbb{E}}[\cdot]$ of the RHS is $\mathcal{O}_{\mathbb{P}}(n\mu_n^{-2} + n) = \mathcal{O}_{\mathbb{P}}(n\mu_n^{-2})$. \square

Let us define $\widehat{\eta}_G$ and $\widehat{\eta}_H$ as the “truncated” values of the conditional expectation of $\mathbf{h}^\top \widetilde{\mathbf{h}}$ and $\boldsymbol{\psi}^\top \widetilde{\boldsymbol{\psi}}$, respectively:

$$\widehat{\eta}_G := \Pi_{[-1,1]}(\overline{\mathbb{E}}[\mathbf{h}^\top \widetilde{\mathbf{h}}]/(\alpha \widetilde{\alpha})), \quad \widehat{\eta}_H := \Pi_{[-1,1]}(\overline{\mathbb{E}}[\boldsymbol{\psi}^\top \widetilde{\boldsymbol{\psi}}]/(p\beta \widetilde{\beta}))$$

where $\Pi_{[-1,1]}(x) := \max(\min(x, 1), -1)$ is the projection map onto $[-1, 1]$. Here, we emphasize that $(\widehat{\eta}_H, \widehat{\eta}_G)$ is independent from the design matrix \mathbf{G} since \mathbf{G} is integrated out. Furthermore, the absolute values of $\widehat{\eta}_G$ and $\widehat{\eta}_H$ are less than 1 due to the truncation, and hence the two matrices

$$\begin{pmatrix} 1 & \widehat{\eta}_H \\ \widehat{\eta}_H & 1 \end{pmatrix} \text{ and } \begin{pmatrix} 1 & \widehat{\eta}_G \\ \widehat{\eta}_G & 1 \end{pmatrix}$$

are both positive semidefinite. By the concentration from Lemma 13 and the convergence $\|\mathbf{h}\|^2 \xrightarrow{\mathbb{P}} \alpha^2 > 0$, $\|\boldsymbol{\psi}\|^2/p \xrightarrow{\mathbb{P}} \beta^2 > 0$, noting the truncation $x \mapsto \Pi_{[-1,1]}(x)$ is continuous, we find that the random variables $(\widehat{\eta}_H, \widehat{\eta}_G)$ defined above still capture the correlations, that is,

$$\begin{aligned} \widehat{\eta}_H &= \Pi_{[-1,1]}\left(\frac{\boldsymbol{\psi}^\top \widetilde{\boldsymbol{\psi}}}{\|\boldsymbol{\psi}\| \|\widetilde{\boldsymbol{\psi}}\|}\right) + o_{\mathbb{P}}(1) = \frac{\boldsymbol{\psi}^\top \widetilde{\boldsymbol{\psi}}}{\|\boldsymbol{\psi}\| \|\widetilde{\boldsymbol{\psi}}\|} + o_{\mathbb{P}}(1). \\ \widehat{\eta}_G &= \Pi_{[-1,1]}\left(\frac{\mathbf{h}^\top \widetilde{\mathbf{h}}}{\|\mathbf{h}\| \|\widetilde{\mathbf{h}}\|}\right) + o_{\mathbb{P}}(1) = \frac{\mathbf{h}^\top \widetilde{\mathbf{h}}}{\|\mathbf{h}\| \|\widetilde{\mathbf{h}}\|} + o_{\mathbb{P}}(1). \end{aligned}$$

where the second equation follows from the fact that the correlations are less than 1 in absolute values by the Cauchy–Schwarz inequality.

Now, we use the multivariate normal approximation below to invite System 1b.

Lemma 14 (Proposition 5.1 in [BK24]). *Let $\mathbf{z} \sim \mathcal{N}(\mathbf{0}_q, \mathbf{I}_q)$ and let $\mathbf{F} : \mathbb{R}^q \rightarrow \mathbb{R}^{q \times M}$ be a locally Lipschitz function with $M \leq q$. Then there exists a standard normal vector $\mathbf{w} \sim \mathcal{N}(\mathbf{0}_M, \mathbf{I}_M)$ such that*

$$\mathbb{E}\left[\left\|\mathbf{F}(\mathbf{z})^\top \mathbf{z} - \sum_{l \in [q]} \frac{\partial \mathbf{F}(\mathbf{z})^\top \mathbf{e}_l}{\partial z_l} - \left\{\mathbf{F}(\mathbf{z})^\top \mathbf{F}(\mathbf{z})\right\}^{1/2} \mathbf{w}\right\|^2\right] \leq C_3 \sum_{l \in [q]} \mathbb{E}\left[\left\|\frac{\partial \mathbf{F}(\mathbf{z})}{\partial z_l}\right\|_F^2\right],$$

where $\{\cdot\}^{1/2}$ is the square root of the positive semidefinite matrix.

For each $j \in [p]$, applying Lemma 14 with $\mathbf{F} = \begin{bmatrix} \boldsymbol{\psi} & \widetilde{\boldsymbol{\psi}} \\ \sqrt{p\beta} & \sqrt{p\beta} \end{bmatrix} \in \mathbb{R}^{n \times 2}$ and $\mathbf{z} = \mathbf{G}\mathbf{e}_j \in \mathbb{R}^n$, using the derivative formula (29), we find that there exists a random vector $\mathbf{w}_j \in \mathbb{R}^2$ such that

$$\mathbf{w}_j \mid (\boldsymbol{\theta}, \boldsymbol{\epsilon}, \mathbf{G}^{-j}, I, \widetilde{I}) \stackrel{d}{=} \mathcal{N}(\mathbf{0}_2, \mathbf{I}_2)$$

and

$$\begin{aligned} &\mathbb{E}\left[\left\|\frac{1}{\sqrt{p}} \begin{pmatrix} \boldsymbol{\psi}^\top \mathbf{G}\mathbf{e}_j^\top + \text{tr}[\mathbf{V}]h_j + \boldsymbol{\psi}^\top \mathbf{D}\mathbf{G}\mathbf{A}\mathbf{e}_j/\beta \\ \widetilde{\boldsymbol{\psi}}^\top \mathbf{G}\mathbf{e}_j^\top + \text{tr}[\widetilde{\mathbf{V}}]\widetilde{h}_j + \widetilde{\boldsymbol{\psi}}^\top \widetilde{\mathbf{D}}\mathbf{G}\widetilde{\mathbf{A}}\mathbf{e}_j/\widetilde{\beta} \end{pmatrix} - \begin{pmatrix} \|\boldsymbol{\psi}\|^2/(p\beta^2) & \boldsymbol{\psi}^\top \widetilde{\boldsymbol{\psi}}/(p\beta \widetilde{\beta}) \\ \boldsymbol{\psi}^\top \widetilde{\boldsymbol{\psi}}/(p\beta \widetilde{\beta}) & \|\widetilde{\boldsymbol{\psi}}\|^2/(p\beta^2) \end{pmatrix}^{1/2} \mathbf{w}_j\right\|^2\right] \\ &\leq C_4 \sum_{i \in [n]} \mathbb{E}\left[\frac{1}{p\beta^2} \left\|\frac{\partial \boldsymbol{\psi}}{\partial g_{ij}}\right\|^2 + \frac{1}{p\widetilde{\beta}^2} \left\|\frac{\partial \widetilde{\boldsymbol{\psi}}}{\partial g_{ij}}\right\|^2\right]. \end{aligned}$$

Using $\|\mathbf{A}\|_{\text{op}} \leq (p\mu_n)^{-1}$ from (30) and Lemma 10, we have

$$\mathbb{E}\left[\|\boldsymbol{\psi}^\top \mathbf{D}\mathbf{G}\mathbf{A}\|^2\right] = \mathcal{O}_{\mathbb{P}}(\mu^{-2}),$$

$$\sum_{ij} \mathbb{E} \left[\left\| \frac{\partial \boldsymbol{\psi}}{\partial g_{ij}} \right\|^2 \right] \leq \mathbb{E} \left[2 \|\mathbf{DGA}\|_{\text{F}}^2 \|\boldsymbol{\psi}\|^2 + 2 \|\mathbf{V}\|_{\text{F}}^2 \|\mathbf{h}\|^2 \right] = \mathcal{O}_{\mathbb{P}}(n\mu^{-2}).$$

Thus, summing over $j \in [p]$ in the previous display, noting $\mu^{-2} = o(n)$, we are left with

$$\sum_{j \in [p]} \left\| \frac{1}{\sqrt{p}} \begin{pmatrix} (\boldsymbol{\psi}^{\top} \mathbf{G} \mathbf{e}_j^{\top} + \text{tr}[\mathbf{V}]h_j)/\beta \\ (\tilde{\boldsymbol{\psi}}^{\top} \mathbf{G} \mathbf{e}_j^{\top} + \text{tr}[\tilde{\mathbf{V}}]h_j)/\tilde{\beta} \end{pmatrix} - \begin{pmatrix} \|\boldsymbol{\psi}\|^2/(p\beta^2) & \boldsymbol{\psi}^{\top} \tilde{\boldsymbol{\psi}}/(p\beta\tilde{\beta}) \\ \boldsymbol{\psi}^{\top} \tilde{\boldsymbol{\psi}}/(p\beta\tilde{\beta}) & \|\tilde{\boldsymbol{\psi}}\|^2/(p\tilde{\beta}^2) \end{pmatrix}^{1/2} \mathbf{w}_j \right\|^2 = o_{\mathbb{P}}(n).$$

By the (1/2)-Hölder continuity for the matrix square root: $\|\mathbf{M}^{1/2} - \mathbf{N}^{1/2}\|_{\text{op}} \leq \|\mathbf{M} - \mathbf{N}\|_{\text{op}}^{1/2}$ for positive semidefinite matrices \mathbf{M}, \mathbf{N} , combined with the convergence $\|\boldsymbol{\psi}\|^2/p \rightarrow \beta^2$, $\|\tilde{\boldsymbol{\psi}}\|^2/p \rightarrow \tilde{\beta}^2$ and the concentration $\hat{\eta}_H = \boldsymbol{\psi}^{\top} \tilde{\boldsymbol{\psi}}/(p\beta\tilde{\beta}) + o_{\mathbb{P}}(1)$, we have

$$\left\| \begin{pmatrix} \|\boldsymbol{\psi}\|^2/(p\beta^2) & \boldsymbol{\psi}^{\top} \tilde{\boldsymbol{\psi}}/(p\beta\tilde{\beta}) \\ \boldsymbol{\psi}^{\top} \tilde{\boldsymbol{\psi}}/(p\beta\tilde{\beta}) & \|\tilde{\boldsymbol{\psi}}\|^2/(p\tilde{\beta}^2) \end{pmatrix}^{1/2} - \begin{pmatrix} 1 & \hat{\eta}_H \\ \hat{\eta}_H & 1 \end{pmatrix}^{1/2} \right\|_{\text{op}} = o_{\mathbb{P}}(1).$$

Combined with the previous display, we get

$$\sum_{j \in [p]} \left\| \frac{1}{\sqrt{p}} \begin{pmatrix} (\boldsymbol{\psi}^{\top} \mathbf{G} \mathbf{e}_j^{\top} + \text{tr}[\mathbf{V}]h_j)/\beta \\ (\tilde{\boldsymbol{\psi}}^{\top} \mathbf{G} \mathbf{e}_j^{\top} + \text{tr}[\tilde{\mathbf{V}}]h_j)/\tilde{\beta} \end{pmatrix} - \begin{pmatrix} 1 & \hat{\eta}_H \\ \hat{\eta}_H & 1 \end{pmatrix}^{1/2} \mathbf{w}_j \right\|^2 = o_{\mathbb{P}}(n) + o_{\mathbb{P}}(1) \sum_{j \in [p]} \|\mathbf{w}_j\|^2 = o_{\mathbb{P}}(n),$$

where the last equation follows from $\sum_{j \in [p]} \mathbb{E}[\|\mathbf{w}_j\|^2] = 2p$, which follows from the fact that the marginal distribution of \mathbf{w}_j given $(\mathbf{z}, \boldsymbol{\theta}, I, \tilde{I})$ is $\mathcal{N}(0, 1)$ (note, however, that we do not establish or take for granted that $(\mathbf{w}_j)_{j \in [p]}$ are i.i.d.). Using $\text{tr}[\mathbf{V}]/p \xrightarrow{\mathbb{P}} \nu$ (Lemma 12) and $\|\mathbf{h}\|^2 = \mathcal{O}_{\mathbb{P}}(1)$ from Lemma 10, we can also replace $\text{tr}[\mathbf{V}]h_j$ by $p\nu h_j$. As a consequence, we get

$$\sum_{j \in [p]} \left\| \begin{pmatrix} \frac{1}{\sqrt{p\beta}} (\boldsymbol{\psi}^{\top} \mathbf{G} \mathbf{e}_j^{\top} + p\nu h_j) \\ \frac{1}{\sqrt{p\tilde{\beta}}} (\tilde{\boldsymbol{\psi}}^{\top} \mathbf{G} \mathbf{e}_j^{\top} + p\tilde{\nu} h_j) \end{pmatrix} - \begin{pmatrix} \hat{w}_j \\ \hat{w}_j \end{pmatrix} \right\|^2 = o_{\mathbb{P}}(n) \text{ where } \begin{pmatrix} \hat{w}_j \\ \hat{w}_j \end{pmatrix} := \begin{pmatrix} 1 & \hat{\eta}_H \\ \hat{\eta}_H & 1 \end{pmatrix}^{1/2} \mathbf{w}_j. \quad (31)$$

Here, we emphasize that the conditional distribution of (\hat{w}_j, \tilde{w}_j) is given by

$$\begin{pmatrix} \hat{w}_j \\ \tilde{w}_j \end{pmatrix} \mid (\mathbf{z}, \boldsymbol{\theta}, \mathbf{G}^{-j}, I, \tilde{I}) \stackrel{\text{d}}{=} \mathcal{N}\left(\mathbf{0}_2, \begin{pmatrix} 1 & \hat{\eta}_H \\ \hat{\eta}_H & 1 \end{pmatrix}\right)$$

since the conditional distribution of \mathbf{w}_j given $(\mathbf{z}, \boldsymbol{\theta}, \mathbf{G}^{-j}, I, \tilde{I})$ is standard normal $\mathcal{N}(0_2, I_2)$ and $\hat{\eta}_H$ is $\sigma(\mathbf{z}, \boldsymbol{\theta}, I, \tilde{I})$ -measurable. For each $j \in [p]$, let us define Ξ_j and $\tilde{\Xi}_j$ as

$$\Xi_j = \frac{\boldsymbol{\psi}^{\top} \mathbf{G} \mathbf{e}_j^{\top} + p\nu h_j}{\sqrt{p\beta}} - \hat{w}_j, \quad \tilde{\Xi}_j = \frac{\tilde{\boldsymbol{\psi}}^{\top} \mathbf{G} \mathbf{e}_j^{\top} + p\tilde{\nu} h_j}{\sqrt{p\tilde{\beta}}} - \tilde{w}_j$$

so that the bound (31) reads $\sum_{j \in [p]} [\Xi_j^2 + \tilde{\Xi}_j^2] = o_{\mathbb{P}}(n)$. By the KKT conditions $\mathbf{G}^{\top} \boldsymbol{\psi} \in \sqrt{p} \partial \text{reg}(\sqrt{p}\mathbf{h} + \boldsymbol{\theta})$ and $\mathbf{G}^{\top} \tilde{\boldsymbol{\psi}} \in \sqrt{p} \partial \tilde{\text{reg}}(\sqrt{p}\tilde{\mathbf{h}} + \boldsymbol{\theta})$ for \mathbf{h} and $\tilde{\mathbf{h}}$, respectively, $\sqrt{p}h_j$ and $\sqrt{p}\tilde{h}_j$ can be written as

$$\begin{pmatrix} \sqrt{p}h_j \\ \sqrt{p}\tilde{h}_j \end{pmatrix} = \begin{pmatrix} \text{prox}_{\text{reg}}(\theta_j + \frac{\beta}{\nu}(\hat{w}_j + \Xi_j); \nu^{-1}) - \theta_j \\ \text{prox}_{\tilde{\text{reg}}}(\theta_j + \frac{\tilde{\beta}}{\tilde{\nu}}(\tilde{w}_j + \tilde{\Xi}_j); \tilde{\nu}^{-1}) - \theta_j \end{pmatrix}.$$

Now we claim that the bound $\sum_{j \in [p]} \Xi_j^2 = o_{\mathbb{P}}(n)$ also holds in the conditional expectation $\bar{\mathbb{E}}$, i.e., $\bar{\mathbb{E}}[\sum_{j \in [p]} \Xi_j^2] = o_{\mathbb{P}}(n)$. To this end, it suffices to show $\bar{\mathbb{E}}[(p^{-1} \sum_{j \in [p]} \Xi_j^2)^2] \leq C$ with high probability for a constant C . Using the upper estimate

$$0 \leq \frac{1}{p} \sum_{j \in [p]} \Xi_j^2 \leq 3 \left(\frac{\|\mathbf{G}^\top \boldsymbol{\psi}\|^2}{\beta^2 p^2} + \frac{\nu^2}{\beta^2} \|\mathbf{h}\|^2 + \frac{1}{p} \sum_{j \in [p]} (\hat{w}_j)^2 \right)$$

and the moment bound from Lemma 10, we get

$$\bar{\mathbb{E}} \left[\left(\frac{1}{p} \sum_{j \in [p]} \Xi_j^2 \right)^2 \right] \leq C_5 \left(1 + \bar{\mathbb{E}} \left[\left(\frac{1}{p} \sum_{j \in [p]} (\hat{w}_j)^2 \right)^2 \right] \right)$$

with high probability. For the second term, expanding the square of the summation,

$$\bar{\mathbb{E}} \left[\left(p^{-1} \sum_{j \in [p]} (\hat{w}_j)^2 \right)^2 \right] = \frac{1}{p^2} \sum_{i,j} \bar{\mathbb{E}}[(\hat{w}_i)^2 (\hat{w}_j)^2] \leq \frac{1}{p^2} \sum_{i,j} \sqrt{\bar{\mathbb{E}}[(\hat{w}_i)^4]} \sqrt{\bar{\mathbb{E}}[(\hat{w}_j)^4]} = 6.$$

where we have used $\bar{\mathbb{E}}[(\hat{w}_j)^4] = 6$ for all j , which follows from the fact that the marginal law of \hat{w}_j is $\mathcal{N}(0, 1)$. This gives $\bar{\mathbb{E}}[(\frac{1}{p} \sum_{j \in [p]} \Xi_j^2)^2] \leq C$ with high probability for a constant C , and hence we get the estimate $\bar{\mathbb{E}}[\frac{1}{p} \sum_{j \in [p]} \Xi_j^2] = o_{\mathbb{P}}(1)$. The same argument yields $\bar{\mathbb{E}}[\frac{1}{p} \sum_{j \in [p]} \tilde{\Xi}_j^2] = o_{\mathbb{P}}(1)$.

Combined with the proximal representation of $(\sqrt{p}h_j, \sqrt{p}\tilde{h}_j)$ using $(\Xi_j, \tilde{\Xi}_j)$, since $\text{prox}_f(\cdot)$ is 1-Lipschitz for any convex function f , the upper bounds of $\bar{\mathbb{E}}[\frac{1}{p} \sum_{j \in [p]} \Xi_j^2]$ and $\bar{\mathbb{E}}[\frac{1}{p} \sum_{j \in [p]} \tilde{\Xi}_j^2]$ yield the following simple proximal approximation of $(\sqrt{p}h_j, \sqrt{p}\tilde{h}_j)$:

$$\frac{1}{p} \sum_{j \in [p]} \bar{\mathbb{E}} \left[\left\| \begin{pmatrix} \sqrt{p}h_j \\ \sqrt{p}\tilde{h}_j \end{pmatrix} - \begin{pmatrix} \text{prox}_{\text{reg}}(\theta_j + (\beta/\nu)\hat{w}_j; \nu^{-1}) - \theta_j \\ \text{prox}_{\tilde{\text{reg}}}(\theta_j + (\tilde{\beta}/\tilde{\nu})\tilde{w}_j; \tilde{\nu}^{-1}) - \theta_j \end{pmatrix} \right\|^2 \right] = o_{\mathbb{P}}(1).$$

Noting $\bar{\mathbb{E}}[\|\mathbf{h}\|^2] = \mathcal{O}_{\mathbb{P}}(1)$, this lets us approximate $\bar{\mathbb{E}}[\mathbf{h}^\top \tilde{\mathbf{h}}]/\alpha\tilde{\alpha}$ by the inner product of proximal operators:

$$\begin{aligned} \frac{\bar{\mathbb{E}}[\mathbf{h}^\top \tilde{\mathbf{h}}]}{\alpha\tilde{\alpha}} + o_{\mathbb{P}}(1) &= \frac{1}{p} \sum_{j \in [p]} \frac{1}{\alpha\tilde{\alpha}} \bar{\mathbb{E}} \left[\left(\text{prox}_{\text{reg}}(\theta_j + (\beta/\nu)\hat{w}_j; \nu^{-1}) - \theta_j \right) \left(\text{prox}_{\tilde{\text{reg}}}(\theta_j + (\tilde{\beta}/\tilde{\nu})\tilde{w}_j; \tilde{\nu}^{-1}) - \theta_j \right) \right] \\ &= \frac{1}{p} \sum_{j \in [p]} F_{\text{reg}}(\hat{\eta}_H; \theta_j), \end{aligned}$$

where for the second inequality, we used the fact that the marginal law of (\hat{w}_j, \tilde{w}_j) given $(\boldsymbol{\theta}, \mathbf{z}, I, \tilde{I})$ is jointly Gaussian, with zero mean, unit variance, and correlation $\hat{\eta}_H$, and where $F_{\text{reg}}(\cdot; \theta_j) : [-1, 1] \rightarrow \mathbb{R}$ is the function defined by

$$F_{\text{reg}}(\eta; \theta_j) = \iint \frac{(\text{prox}_{\text{reg}}(\theta_j + \frac{\beta}{\nu}x; \frac{1}{\nu}) - \theta_j) (\text{prox}_{\tilde{\text{reg}}}(\theta_j + \frac{\tilde{\beta}}{\tilde{\nu}}(\eta x + \sqrt{1-\eta^2}y); \frac{1}{\tilde{\nu}}) - \theta_j)}{\alpha\tilde{\alpha}} \varphi(x)\varphi(y) dx dy$$

with φ being the standard normal probability function and $\iint = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty}$. Notice that for each $\eta \in [-1, 1]$, the sequence $(F_{\text{reg}}(\eta; \theta_j))_{j=1}^p$ are i.i.d. random variables with mean $\mathbb{E}[F_{\text{reg}}(\eta; \theta_j)] = F_{\text{reg}}(\eta)$. Furthermore, by the Jensen's inequality and the Cauchy-Schwarz inequality, the expectation of the absolute value is finite:

$$\mathbb{E} \left[|F_{\text{reg}}(\eta; \theta_j)| \right] \leq \frac{1}{\alpha\tilde{\alpha}} \mathbb{E} \left[\left| (\text{prox}_{\text{reg}}(\Theta + \frac{\beta}{\nu}H; \frac{1}{\nu}) - \Theta) \cdot (\text{prox}_{\tilde{\text{reg}}}(\Theta + \frac{\tilde{\beta}}{\tilde{\nu}}\tilde{H}; \frac{1}{\tilde{\nu}}) - \Theta) \right| \right]$$

$$\begin{aligned}
&\leq \frac{1}{\alpha\tilde{\alpha}} \mathbb{E} \left[\left(\text{prox}_{\text{reg}} \left(\Theta + \frac{\beta}{\nu} H; \frac{1}{\nu} \right) - \Theta \right)^2 \right]^{1/2} \cdot \mathbb{E} \left[\left(\text{prox}_{\text{reg}} \left(\Theta + \frac{\tilde{\beta}}{\tilde{\nu}} \tilde{H}; \frac{1}{\tilde{\nu}} \right) - \Theta \right)^2 \right]^{1/2} \\
&= 1 \quad (\text{by (5a) in System 1a}).
\end{aligned}$$

Thus, the weak law of large numbers gives $p^{-1} \sum_{j \in [p]} F_{\text{reg}}(\eta; \theta_j) \xrightarrow{\mathbb{P}} F_{\text{reg}}(\eta)$ for each $\eta \in [-1, 1]$. Next, we claim that this convergence holds uniformly over $\eta \in [-1, 1]$.

Lemma 15. *Let I be a bounded and closed interval of \mathbb{R} and let $(f_n)_{n \geq 1} : I \rightarrow \mathbb{R}$ be a sequence of non-decreasing functions such that $f_n(x) \xrightarrow{\mathbb{P}} f(x)$ pointwise for some function $f : I \rightarrow \mathbb{R}$. If f is non-decreasing and continuous, then the uniform convergence $\sup_{x \in I} |f_n(x) - f(x)| \xrightarrow{\mathbb{P}} 0$ holds.*

Note that Lemma 15 is a probabilistic analogue of a similar statement for deterministic functions: if a sequence of real-valued monotone functions (on \mathbb{R}) converges pointwise to a continuous function on a compact set $I \subset \mathbb{R}$, then the convergence is uniform on the set I . The probabilistic version is known but in a more general setting, so to keep the treatment self-contained we give a basic proof below which is similar to proof of Glivenko-Cantelli theorem (cf. [VdV00, Theorem 19.1]).

Proof. Let us write $I = [a, b]$. Since f is continuous and I is compact, f is uniformly continuous on I . For any $\epsilon > 0$, there exists some $\delta_\epsilon > 0$ such that $|f(x) - f(y)| < \epsilon/2$ for all $x, y \in I$ such that $|x - y| \leq \delta_\epsilon$. Now for sufficiently large integer $k = k_\epsilon \in \mathbb{N}$ such that $(b - a)/k < \delta_\epsilon$, let us take equally spaced grids $(x_i)_{i=0}^k$ over $[a, b]$ such that $a = x_0 < x_1 < \dots < x_k = b$ and $(x_i - x_{i-1}) = (b - a)/k$ for all $i \in \{0, \dots, k\}$. Let $\Omega_\epsilon = \cap_{i=0}^k \{|f_n(x_i) - f(x_i)| \leq \epsilon/2\}$ be the event under which f_n and f are sufficiently close at the finite grids. Note that this event holds with probability converging to 1, thanks to the pointwise convergence $f_n(x_i) \xrightarrow{\mathbb{P}} f(x_i)$ at finitely many x_i . Then, since f_n is non-decreasing while f does not move more than $\epsilon/2$ in $[x_{i-1}, x_i]$, for all $i \in \{0, 1, \dots, k\}$ and for all $x \in [x_{i-1}, x_i]$ we have

$$\begin{aligned}
f_n(x_{i-1}) &\leq f_n(x) \leq f_n(x_i), \\
-\epsilon/2 - f(x_{i-1}) &\leq -f(x) \leq -f(x_i) + \epsilon/2.
\end{aligned}$$

In the event Ω_ϵ , summing the two lines it holds that $-\epsilon \leq f_n(x) - f(x) \leq \epsilon$ for all $i \in \{0, 1, \dots, k\}$ and for all $x \in [x_{i-1}, x_i]$. \square

Recall that we have shown in the proof of Theorem 1 that F_{reg} is differentiable with a non-negative derivative. By the same argument, the map $\eta \mapsto F_{\text{reg}}(\eta; \theta_j)$ is differentiable with a non-negative derivative. Thus, applying Lemma 15 with $f_n(\cdot) = p^{-1} \sum_{j \in [p]} F_{\text{reg}}(\cdot; \theta_j)$, $f(\cdot) = F_{\text{reg}}(\cdot)$ and $I = [-1, 1]$, we get the uniform convergence:

$$\sup_{\eta \in [-1, 1]} \left| \frac{1}{p} \sum_{j \in [p]} F_{\text{reg}}(\eta; \theta_j) - F_{\text{reg}}(\eta) \right| \xrightarrow{\mathbb{P}} 0.$$

Combined with $\frac{\mathbb{E}[\mathbf{h}^\top \tilde{\mathbf{h}}]}{\alpha\tilde{\alpha}} = \frac{1}{p} \sum_{j \in [p]} \widehat{F}_{\text{reg}}(\widehat{\eta}_H; \theta_j) + o_{\mathbb{P}}(1)$ and $\widehat{\eta}_H \in [-1, 1]$, we are left with

$$\frac{\mathbb{E}[\mathbf{h}^\top \tilde{\mathbf{h}}]}{\alpha\tilde{\alpha}} = \frac{1}{p} \sum_{j \in [p]} F_{\text{reg}}(\widehat{\eta}_H; \theta_j) + o_{\mathbb{P}}(1) = F_{\text{reg}}(\widehat{\eta}_H) + o_{\mathbb{P}}(1).$$

Recall the definition $\hat{\eta}_G = \Pi_{[-1,1]}(\frac{\mathbb{E}[\mathbf{h}^\top \tilde{\mathbf{h}}]}{\alpha \tilde{\alpha}})$ where $\Pi_{[-1,1]}$ is the projection onto $[-1,1]$. By the continuity of the projection map and the bound $\sup_{\eta \in [-1,1]} |F_{\text{reg}}(\eta)| \leq 1$ (see Theorem 1), the above display yields

$$\hat{\eta}_G = \Pi_{[-1,1]}(F_{\text{reg}}(\hat{\eta}_H)) + o_{\mathbb{P}}(1) = F_{\text{reg}}(\hat{\eta}_H) + o_{\mathbb{P}}(1).$$

Next, let us show $\hat{\eta}_H = F_{\text{loss}}(\hat{\eta}_G) + o_{\mathbb{P}}(1)$ using the same argument. Using Lemma 14 with $\mathbf{F} = [\mathbf{h}/\alpha \mid \tilde{\mathbf{h}}/\tilde{\alpha}] \in \mathbb{R}^{p \times 2}$ and $\mathbf{z} = \mathbf{g}_i$, there exists a random vector $\mathbf{u}_i \in \mathbb{R}^2$ with conditional distribution

$$\mathbf{u}_i \mid \mathbf{z}, \boldsymbol{\theta}, \mathbf{G}_{-i}, I, \tilde{I} \stackrel{d}{=} \mathcal{N}(\mathbf{0}_2, \mathbf{I}_2)$$

and

$$\begin{aligned} \mathbb{E} \left[\left\| \begin{pmatrix} (\mathbf{g}_i^\top \mathbf{h} - \text{tr}[\mathbf{A}]\psi_i + \mathbf{h}^\top \mathbf{A} \mathbf{G}^\top \mathbf{D} \mathbf{e}_i) / \alpha \\ (\mathbf{g}_i^\top \tilde{\mathbf{h}} - \text{tr}[\tilde{\mathbf{A}}]\tilde{\psi}_i + \tilde{\mathbf{h}}^\top \tilde{\mathbf{A}} \mathbf{G}^\top \tilde{\mathbf{D}} \mathbf{e}_i) / \tilde{\alpha} \end{pmatrix} - \begin{pmatrix} \|\mathbf{h}\|^2 / \alpha^2 & \mathbf{h}^\top \tilde{\mathbf{h}} / \alpha \tilde{\alpha} \\ \tilde{\mathbf{h}}^\top \mathbf{h} / \alpha \tilde{\alpha} & \|\tilde{\mathbf{h}}\|^2 / \tilde{\alpha}^2 \end{pmatrix}^{1/2} \mathbf{u}_i \right\|^2 \right] \\ \leq C_6 \sum_{j=1}^p \mathbb{E} \frac{1}{\alpha^2} \left\| \frac{\partial \mathbf{h}}{\partial g_{ij}} \right\|^2 + \frac{1}{\tilde{\alpha}^2} \left\| \frac{\partial \tilde{\mathbf{h}}}{\partial g_{ij}} \right\|^2. \quad (32) \end{aligned}$$

Using $\|\mathbf{A}\|_{\text{op}} \leq (p\mu)^{-1}$ in (30) Lemma 10, noting $\mu^{-2} = o(n)$, it follows that

$$\begin{aligned} \mathbb{E} [\|\mathbf{h} \mathbf{A} \mathbf{G}^\top \mathbf{D}\|^2] &= \mathcal{O}_{\mathbb{P}}(n^{-1} \mu^{-2}) = o_{\mathbb{P}}(1) \\ \sum_{i \in [n]} \sum_{j \in [p]} \mathbb{E} \left[\left\| \frac{\partial \mathbf{h}}{\partial g_{ij}} \right\|^2 \right] &\leq 2 \mathbb{E} [\|\mathbf{A}\|_{\text{F}}^2 \|\boldsymbol{\psi}\|^2 + \|\mathbf{A} \mathbf{G}^\top \mathbf{D}\|_{\text{F}}^2 \|\mathbf{h}\|^2] = \mathcal{O}_{\mathbb{P}}(\mu^{-2}) = o_{\mathbb{P}}(n) \end{aligned}$$

so that summing over $i \in [n]$ the inequality (32), we get

$$\sum_{i \in [n]} \left\| \begin{pmatrix} (\mathbf{g}_i^\top \mathbf{h} - \text{tr}[\mathbf{A}]\psi_i) / \alpha \\ (\mathbf{g}_i^\top \tilde{\mathbf{h}} - \text{tr}[\tilde{\mathbf{A}}]\tilde{\psi}_i) / \tilde{\alpha} \end{pmatrix} - \begin{pmatrix} \|\mathbf{h}\|^2 / \alpha^2 & \mathbf{h}^\top \tilde{\mathbf{h}} / \alpha \tilde{\alpha} \\ \tilde{\mathbf{h}}^\top \mathbf{h} / \alpha \tilde{\alpha} & \|\tilde{\mathbf{h}}\|^2 / \tilde{\alpha}^2 \end{pmatrix}^{1/2} \mathbf{u}_i \right\|^2 = o_{\mathbb{P}}(n).$$

Appealing to the (1/2)-Hölder continuity for the matrix square root again, now using the two convergences $\|\mathbf{h}\|^2 \rightarrow \alpha^2$, $\|\tilde{\mathbf{h}}\|^2 \xrightarrow{P} \tilde{\alpha}^2$ and the concentration $\hat{\eta}_G = \mathbf{h}^\top \tilde{\mathbf{h}} / (\alpha \tilde{\alpha}) + o_{\mathbb{P}}(1)$, we get

$$\left\| \begin{pmatrix} \|\mathbf{h}\|^2 / \alpha^2 & \mathbf{h}^\top \tilde{\mathbf{h}} / \alpha \tilde{\alpha} \\ \tilde{\mathbf{h}}^\top \mathbf{h} / \alpha \tilde{\alpha} & \|\tilde{\mathbf{h}}\|^2 / \tilde{\alpha}^2 \end{pmatrix}^{1/2} - \begin{pmatrix} 1 & \hat{\eta}_G \\ \hat{\eta}_G & 1 \end{pmatrix}^{1/2} \right\|_{\text{op}} = o_{\mathbb{P}}(1).$$

Combined with the convergence $\text{tr}[\mathbf{A}] \xrightarrow{P} \kappa$, $\text{tr}[\tilde{\mathbf{A}}] \xrightarrow{P} \tilde{\kappa}$, noting that $\sum_{i \in [n]} \|\mathbf{u}_i\|^2$ and $\|\boldsymbol{\psi}\|^2 + \|\tilde{\boldsymbol{\psi}}\|^2$ are both $\mathcal{O}_{\mathbb{P}}(n)$, we are left with

$$\frac{1}{n} \sum_{i \in [n]} \left\| \begin{pmatrix} \mathbf{g}_i^\top \mathbf{h} - \kappa \psi_i \\ \mathbf{g}_i^\top \tilde{\mathbf{h}} - \tilde{\kappa} \tilde{\psi}_i \end{pmatrix} - \begin{pmatrix} \alpha \hat{u}_i \\ \tilde{\alpha} \hat{u}_i \end{pmatrix} \right\|^2 = o_{\mathbb{P}}(1) \quad \text{where} \quad \begin{pmatrix} \hat{u}_i \\ \hat{u}_i \end{pmatrix} = \begin{pmatrix} 1 & \hat{\eta}_G \\ \hat{\eta}_G & 1 \end{pmatrix}^{1/2} \mathbf{u}_i.$$

By the same argument that we used to bound Ξ_j and $\tilde{\Xi}_j$, using Lemma 10 and the fact that the marginal law of \hat{u}_i (and \tilde{u}_i) is $\mathcal{N}(0,1)$, we can show that the conditional expectation $\bar{\mathbb{E}}$ of the square of LHS is bounded from above by a constant C with high probability. Thus, the above approximation also holds in the conditional expectation $\bar{\mathbb{E}}$:

$$\frac{1}{n} \sum_{i \in [n]} \bar{\mathbb{E}} \left[\left\| \begin{pmatrix} \mathbf{g}_i^\top \mathbf{h} - \kappa \psi_i \\ \mathbf{g}_i^\top \tilde{\mathbf{h}} - \tilde{\kappa} \tilde{\psi}_i \end{pmatrix} - \begin{pmatrix} \alpha \hat{u}_i \\ \tilde{\alpha} \hat{u}_i \end{pmatrix} \right\|^2 \right] = o_{\mathbb{P}}(1).$$

Let us define $\Xi^i := \mathbf{g}_i^\top \mathbf{h} - \kappa \psi_i - \alpha \hat{u}_i$ and $\tilde{\Xi}^i := \mathbf{g}_i^\top \tilde{\mathbf{h}} - \tilde{\kappa} \tilde{\psi}_i - \tilde{\alpha} \tilde{u}_i$ so that the above display reads $\sum_{i \in [n]} \mathbb{E}[(\Xi^i)^2 + (\tilde{\Xi}^i)^2] = o(n)$. Since $\psi_i = \text{loss}'(z_i - \mathbf{g}_i^\top \mathbf{h})$ and $\tilde{\psi}_i = \widetilde{\text{loss}}'(z_i - \mathbf{g}_i^\top \tilde{\mathbf{h}})$ for all $i \in I \cap \tilde{I}$, the residuals can be written as

$$z_i - \mathbf{g}_i^\top \mathbf{h} = \text{prox}_{\text{loss}}(z_i - \alpha \hat{u}_i - \Xi^i; \kappa), \quad z_i - \mathbf{g}_i^\top \tilde{\mathbf{h}} = \text{prox}_{\widetilde{\text{loss}}}(z_i - \tilde{\alpha} \tilde{u}_i - \tilde{\Xi}^i; \tilde{\kappa})$$

for all $i \in I \cap \tilde{I}$. Since the map $x \mapsto \text{env}'_f(x; \tau) = f' \circ \text{prox}_f(x; \tau)$ is a composition of Lipschitz functions if f is convex and differentiable with Lipschitz derivative, the moment bound $\sum_{i \in [n]} \mathbb{E}[(\Xi^i)^2 + (\tilde{\Xi}^i)^2] = o(n)$ lets us approximate ψ_i and $\tilde{\psi}_i$ by $\text{env}'_{\text{loss}}$ and $\text{env}'_{\widetilde{\text{loss}}}$ as follows:

$$\frac{1}{n} \sum_{i \in I \cap \tilde{I}} \mathbb{E} \left\| \begin{pmatrix} \psi_i - \text{env}'_{\text{loss}}(z_i - \alpha \hat{u}_i; \kappa) \\ \tilde{\psi}_i - \text{env}'_{\widetilde{\text{loss}}}(z_i - \tilde{\alpha} \tilde{u}_i; \tilde{\kappa}) \end{pmatrix} \right\|^2 = o_{\mathbb{P}}(1) \quad \text{with} \quad \begin{pmatrix} \hat{u}_i \\ \tilde{u}_i \end{pmatrix} | \mathbf{z}, \boldsymbol{\theta}, \mathbf{G}_{-i}, I, \tilde{I} \stackrel{d}{=} \mathcal{N}(0_2, \begin{pmatrix} 1 & \hat{\eta}_G \\ \hat{\eta}_G & 1 \end{pmatrix})$$

Noting $\mathbb{E}[\|\boldsymbol{\psi}\|^2] = \mathcal{O}_{\mathbb{P}}(n)$ by Lemma 10, this lets us approximate $\mathbb{E}[\boldsymbol{\psi}^\top \tilde{\boldsymbol{\psi}}] / (p\beta\tilde{\beta})$ by the inner product of $(\text{env}'_{\text{loss}}, \text{env}'_{\widetilde{\text{loss}}})$:

$$\begin{aligned} \frac{\mathbb{E}[\boldsymbol{\psi}^\top \tilde{\boldsymbol{\psi}}]}{p\beta\tilde{\beta}} &= \frac{|I \cap \tilde{I}|}{p\beta\tilde{\beta}} \frac{1}{|I \cap \tilde{I}|} \mathbb{E}[\boldsymbol{\psi}^\top \tilde{\boldsymbol{\psi}}] \\ &= \frac{c\tilde{c}\tilde{\delta}}{\beta\tilde{\beta}} \frac{1}{|I \cap \tilde{I}|} \sum_{i \in I \cap \tilde{I}} \mathbb{E}[\text{env}'_{\text{loss}}(z_i - \alpha \hat{u}_i; \kappa) \cdot \text{env}'_{\widetilde{\text{loss}}}(z_i - \tilde{\alpha} \tilde{u}_i; \tilde{\kappa})] + o_{\mathbb{P}}(1). \end{aligned}$$

Since the marginal distribution of (u_i, \tilde{u}_i) given $(\boldsymbol{\theta}, \mathbf{z}, I, \tilde{I})$ is centered normal with unit variance and correlation $\hat{\eta}_G$, the above display reads

$$\frac{\mathbb{E}[\boldsymbol{\psi}^\top \tilde{\boldsymbol{\psi}}]}{p\beta\tilde{\beta}} = \frac{1}{|I \cap \tilde{I}|} \sum_{i \in I \cap \tilde{I}} F_{\text{loss}}(\hat{\eta}_G; z_i) + o_{\mathbb{P}}(1)$$

where $F_{\text{loss}}(\eta; z_i) : [-1, 1] \rightarrow \mathbb{R}$ is the function defined as:

$$F_{\text{loss}}(\eta; z_i) = \frac{c\tilde{c}\tilde{\delta}}{\beta\tilde{\beta}} \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \varphi(x) \varphi(y) \text{env}'_{\text{loss}}(z_i + \alpha x; \kappa) \text{env}'_{\widetilde{\text{loss}}}(z_i + \tilde{\alpha}(\eta x + \sqrt{1 - \eta^2} y); \tilde{\kappa}) dx dy.$$

By the same argument for $F_{\text{reg}}(\eta; \theta_j)$, the sequence $(F_{\text{loss}}(\eta; z_i))_{i \in [n]}$ are i.i.d. random variables with mean $\mathbb{E}[F_{\text{loss}}(\eta; z_i)] = F_{\text{loss}}(\eta)$ and the expectation of the absolute value $|F_{\text{loss}}(\eta; z_i)|$ is finite. Thus, by the weak law of large numbers, we have $\frac{1}{m} \sum_{i \in [m]} F_{\text{loss}}(\eta; z_i) \xrightarrow{P} \mathbb{E}[F_{\text{loss}}(\eta; z_1)] = F_{\text{loss}}(\eta)$ for any deterministic integer $m = m_n \rightarrow +\infty$. Then, we have that for any $\epsilon > 0$, denoting by \sum_K the sum over all possible value K taken by $I \cap \tilde{I}$,

$$\begin{aligned} &\mathbb{P}\left(\left| \frac{1}{|I \cap \tilde{I}|} \sum_{i \in I \cap \tilde{I}} F_{\text{loss}}(\eta; z_i) - F_{\text{loss}}(\eta) \right| > \epsilon\right) \\ &= \sum_K \mathbb{P}\left(\left| \frac{1}{|K|} \sum_{i \in K} F_{\text{loss}}(\eta; z_i) - F_{\text{loss}}(\eta) \right| > \epsilon, I \cap \tilde{I} = K\right) \quad (\text{by additivity of disjoint events}) \\ &= \sum_K \mathbb{P}\left(\left| \frac{1}{|K|} \sum_{i \in K} F_{\text{loss}}(\eta; z_i) - F_{\text{loss}}(\eta) \right| > \epsilon\right) \mathbb{P}(I \cap \tilde{I} = K) \quad (\text{by independence}) \\ &= \sum_{m \geq 0} \mathbb{P}\left(\left| \frac{1}{m} \sum_{i \in [m]} F_{\text{loss}}(\eta; z_i) - F_{\text{loss}}(\eta) \right| > \epsilon\right) \mathbb{P}(|I \cap \tilde{I}| = m) \quad (\text{since } (z_i)_{i \in K} \stackrel{d}{=} (z_i)_{i \in [m]} \text{ for } m = |K|). \end{aligned}$$

We now split the sum over $m \geq 0$ into two as follows:

$$\begin{aligned}
& \mathbb{P}\left(\left|\frac{1}{|I \cap \tilde{I}|} \sum_{i \in I \cap \tilde{I}} F_{\text{loss}}(\eta; z_i) - F_{\text{loss}}(\eta)\right| > \epsilon\right) \\
&= \left(\sum_{m \leq n\tilde{c}\tilde{c}/2} + \sum_{m > n\tilde{c}\tilde{c}/2}\right) \mathbb{P}\left(\left|\frac{1}{m} \sum_{i \in [M]} F_{\text{loss}}(\eta; z_i) - F_{\text{loss}}(\eta)\right| > \epsilon\right) \mathbb{P}\left(|I \cap \tilde{I}| = m\right) \\
&\leq \mathbb{P}(|I \cap \tilde{I}| \leq n\tilde{c}\tilde{c}/2) + \sup_{m > n\tilde{c}\tilde{c}/2} \mathbb{P}\left(\left|\frac{1}{m} \sum_{i \in [M]} F_{\text{loss}}(\eta; z_i) - F_{\text{loss}}(\eta)\right| > \epsilon\right) \mathbb{P}(|I \cap \tilde{I}| > n\tilde{c}\tilde{c}/2)
\end{aligned}$$

where $\mathbb{P}(|I \cap \tilde{I}| > n\tilde{c}\tilde{c}/2) \leq 1$ in the rightmost term. The second term converges to 0 by the weak law of large numbers, and the first by Chebyshev's inequality applied to the hypergeometric distribution.

This gives $\frac{1}{|I \cap \tilde{I}|} \sum_{i \in I \cap \tilde{I}} F_{\text{loss}}(\eta; z_i) \xrightarrow{\mathbb{P}} F_{\text{loss}}(\eta)$ pointwise for any $\eta \in [-1, 1]$. By the same argument we used for F_{reg} and $F_{\text{reg}}(\cdot; \theta_j)$, F_{loss} and $F_{\text{loss}}(\cdot; z_i)$ are non-decreasing and continuous. Thus, we can apply Lemma 15 with $f_n(\cdot) = \frac{1}{|I \cap \tilde{I}|} \sum_{i \in I \cap \tilde{I}} F_{\text{loss}}(\cdot; z_i)$, $f(\cdot) = F_{\text{loss}}(\cdot)$ and obtain the uniform convergence:

$$\sup_{\eta \in [-1, 1]} \left| \frac{1}{|I \cap \tilde{I}|} \sum_{i \in I \cap \tilde{I}} F_{\text{loss}}(\eta; z_i) - F_{\text{loss}}(\eta) \right| = o_{\mathbb{P}}(1).$$

Combined with $\frac{\mathbb{E}[\boldsymbol{\psi}^\top \tilde{\boldsymbol{\psi}}]}{p\beta\tilde{\beta}} = \frac{1}{|I \cap \tilde{I}|} \sum_{i \in I \cap \tilde{I}} F_{\text{loss}}(\hat{\eta}_G; z_i)$ and $\hat{\eta}_G \in [-1, 1]$, we are left with

$$\mathbb{E}[\boldsymbol{\psi}^\top \tilde{\boldsymbol{\psi}}]/(p\beta\tilde{\beta}) = F_{\text{loss}}(\hat{\eta}_G) + o_{\mathbb{P}}(1).$$

Recalling $\hat{\eta}_H = \Pi_{[-1, 1]}(\mathbb{E}[\boldsymbol{\psi}^\top \tilde{\boldsymbol{\psi}}]/(p\beta\tilde{\beta}))$ where $\Pi_{[-1, 1]}$ is the projection onto $[-1, 1]$, by the continuity of the projection map and $|F_{\text{loss}}(\eta)| \leq \sqrt{\tilde{c}\tilde{c}} \leq 1$ for all $\eta \in [-1, 1]$ (see Theorem 1), we finally obtain

$$\hat{\eta}_H = \Pi_{[-1, 1]}(F_{\text{loss}}(\hat{\eta}_G)) + o_{\mathbb{P}}(1) = F_{\text{loss}}(\hat{\eta}_G) + o_{\mathbb{P}}(1).$$

In summary, we have shown that $\hat{\eta}_G = F_{\text{reg}}(\hat{\eta}_H) + o_{\mathbb{P}}(1)$ and $\hat{\eta}_H = F_{\text{loss}}(\hat{\eta}_G) + o_{\mathbb{P}}(1)$. By the continuity of F_{reg} and F_{loss} , it holds that

$$\hat{\eta}_G = F_{\text{reg}} \circ F_{\text{loss}}(\hat{\eta}_G) + o_{\mathbb{P}}(1) \quad \text{and} \quad \hat{\eta}_H = F_{\text{loss}} \circ F_{\text{reg}}(\hat{\eta}_H) + o_{\mathbb{P}}(1).$$

Since $F_{\text{reg}} \circ F_{\text{loss}}$ and $F_{\text{loss}} \circ F_{\text{reg}}$ are $(c \wedge \tilde{c})$ -Lipschitz with $c \wedge \tilde{c} < 1$, we have

$$|\hat{\eta}_G - \eta_G| = |F_{\text{reg}} \circ F_{\text{loss}}(\hat{\eta}_G) + o_{\mathbb{P}}(1) - F_{\text{reg}} \circ F_{\text{loss}}(\eta_G)| \leq (c \wedge \tilde{c})|\hat{\eta}_G - \eta_G| + o_{\mathbb{P}}(1),$$

so $\hat{\eta}_G = \eta_G + o_{\mathbb{P}}(1)$, and similarly $\hat{\eta}_H = \eta_H + o_{\mathbb{P}}(1)$ for the fixed points η_G and η_H satisfying $\eta_G = F_{\text{reg}} \circ F_{\text{loss}}(\eta_G)$ and $\eta_H = F_{\text{loss}} \circ F_{\text{reg}}(\eta_H)$. Recalling $\hat{\eta}_G = \mathbf{h}^\top \tilde{\mathbf{h}}/(\alpha\tilde{\alpha}) + o_{\mathbb{P}}(1)$ and $\hat{\eta}_H = \boldsymbol{\psi}^\top \tilde{\boldsymbol{\psi}}/(p\beta\tilde{\beta}) + o_{\mathbb{P}}(1)$, this gives $\mathbf{h}^\top \tilde{\mathbf{h}}/(\alpha\tilde{\alpha}) \xrightarrow{\mathbb{P}} \eta_G$ and $\boldsymbol{\psi}^\top \tilde{\boldsymbol{\psi}}/(p\beta\tilde{\beta}) \xrightarrow{\mathbb{P}} \eta_H$, thereby completing the proof.

Finally, let us show the proximal approximation of estimators $(\sqrt{p}h_j, \sqrt{p}h_j)$ and residuals $(z_i - \mathbf{g}_i^\top \mathbf{h}, z_i - \mathbf{g}_i^\top \tilde{\mathbf{h}})$. Using the convergence $\boldsymbol{\psi}^\top \tilde{\boldsymbol{\psi}}/(p\beta\tilde{\beta}) \xrightarrow{\mathbb{P}} \eta_H$ that we have shown, applying the above

argument with $\widehat{\eta}_H$ replaced by η_H , we have

$$\frac{1}{p} \sum_{j \in [p]} \mathbb{E} \left[\underbrace{\left\| \begin{pmatrix} \sqrt{p} h_j \\ \sqrt{p} \widetilde{h}_j \end{pmatrix} - \begin{pmatrix} \text{prox}_{\text{reg}}(\theta_j + (\beta/\nu) \widehat{w}_j; \nu^{-1}) - \theta_j \\ \text{prox}_{\widetilde{\text{reg}}}(\theta_j + (\beta/\widetilde{\nu}) \widetilde{w}_j; \widetilde{\nu}^{-1}) - \theta_j \end{pmatrix} \right\|^2}_{\text{LHS}_j} \right] = o_{\mathbb{P}}(1).$$

Since $|X_n| = o_{\mathbb{P}}(1)$ is equivalent to $\mathbb{E}[1 \wedge |X_n|] = o(1)$ for any random variable X_n , the above display reads

$$\mathbb{E} \left[1 \wedge \frac{1}{p} \sum_{j \in [p]} \mathbb{E}[\text{LHS}_j] \right] = o(1)$$

where the expectation \mathbb{E} is with respect to $(\mathbf{z}, \boldsymbol{\theta}, I, \widetilde{I})$ (recall the definition $\mathbb{E}[\cdot] = \mathbb{E}[\cdot | \mathbf{z}, \boldsymbol{\theta}, I, \widetilde{I}]$). Note that the integrand is bounded from below as

$$1 \wedge \frac{1}{p} \sum_{j \in [p]} \mathbb{E}[\text{LHS}_j] \geq \frac{1}{p} \sum_{j \in [p]} \left(1 \wedge \mathbb{E}[\text{LHS}_j] \right) \geq \frac{1}{p} \sum_{j \in [p]} \mathbb{E}[1 \wedge \text{LHS}_j],$$

where the second inequality follows from the Jensen's inequality $\mathbb{E}[f(X)] \leq f(\mathbb{E}[X])$ applied with the concave function $f(x) = 1 \wedge x$ and $X = \text{LHS}_j$. Taking the expectation of the above display and using the tower property, we are left with

$$\frac{1}{p} \sum_{j \in [p]} \mathbb{E}[1 \wedge \text{LHS}_j] = o(1).$$

Since the marginal distribution of the integrand $(1 \wedge \text{LHS}_j)$ is the same for all $j \in [p]$ by symmetry, the LHS equals to $\mathbb{E}[1 \wedge \text{LHS}_{j'}]$ for any $j' \in [p]$. This gives $\max_{j \in [p]} \mathbb{E}[1 \wedge \text{LHS}_j] = o(1)$ and completes the proof of the joint approximation of estimators $(\sqrt{p} h_j, \sqrt{p} \widetilde{h}_j)$. The joint approximation of residual follows from the same argument so we omit the proof.

A.3 Proof of Theorem 4

We assume that reg and $\widetilde{\text{reg}}$ are μ -strongly convex for a fixed $\mu > 0$. Then [BS22, Theorem 5.1] implies that

$$\text{tr}[\mathbf{A}] \text{tr}[\mathbf{V}] - \text{df} = \mathcal{O}_{\mathbb{P}}(\sqrt{n}).$$

On the other hand, by the same argument in the proof of Lemma 12 with the diminishing ridge penalty μ_n replaced by the strongly convexity parameter $\mu > 0$, we have $\text{tr}[\mathbf{V}]/p \xrightarrow{\mathbb{P}} \nu > 0$ and $\text{tr}[\mathbf{A}] \xrightarrow{\mathbb{P}} \kappa$. Therefore, we get

$$\text{tr}[\mathbf{A}] - \text{df} / \text{tr}[\mathbf{V}] = \mathcal{O}_{\mathbb{P}}(n^{-1/2}) \quad \text{and} \quad \text{df} / \text{tr}[\mathbf{V}] = \mathcal{O}_{\mathbb{P}}(1).$$

Note in passing that the same things hold for $\widetilde{\text{df}}$ and $\text{tr}[\widetilde{\mathbf{V}}]$. By the Cauchy-Schwarz inequality, the error term due to the replacement of $(\text{df} / \text{tr}[\mathbf{V}], \widetilde{\text{df}} / \text{tr}[\widetilde{\mathbf{V}}])$ by $(\text{tr}[\mathbf{A}], \text{tr}[\widetilde{\mathbf{A}}])$ is

$$\begin{aligned} & \left| \sum_{i \in [n]} \left(r_i + \frac{\text{df}}{\text{tr}[\mathbf{V}]} \mathbf{1}_{\{i \in I\}} \psi_i \right)^\top \left(\widetilde{r}_i + \frac{\widetilde{\text{df}}}{\text{tr}[\widetilde{\mathbf{V}}]} \mathbf{1}_{\{i \in \widetilde{I}\}} \widetilde{\psi}_i \right) - \sum_{i \in [n]} \left(r_i + \text{tr}[\mathbf{A}] \mathbf{1}_{\{i \in I\}} \psi_i \right)^\top \left(\widetilde{r}_i + \text{tr}[\widetilde{\mathbf{A}}] \mathbf{1}_{\{i \in \widetilde{I}\}} \widetilde{\psi}_i \right) \right| \\ & \leq \sqrt{\sum_{i \in [n]} \left(r_i + \frac{\text{df}}{\text{tr}[\mathbf{V}]} \mathbf{1}_{\{i \in I\}} \psi_i \right)^2} \cdot \left| \frac{\widetilde{\text{df}}}{\text{tr}[\widetilde{\mathbf{V}}]} - \text{tr}[\widetilde{\mathbf{A}}] \right| \|\widetilde{\boldsymbol{\psi}}\| + \sqrt{\sum_{i \in [n]} \left(\widetilde{r}_i + \text{tr}[\widetilde{\mathbf{A}}] \mathbf{1}_{\{i \in \widetilde{I}\}} \widetilde{\psi}_i \right)^2} \cdot \left| \frac{\text{df}}{\text{tr}[\mathbf{V}]} - \text{tr}[\mathbf{A}] \right| \|\boldsymbol{\psi}\| \end{aligned}$$

$$= (\|\mathbf{z}\| + \mathcal{O}_{\mathbb{P}}(\sqrt{n})) \cdot \mathcal{O}_{\mathbb{P}}(1)$$

where the last equality follows from the following fact: $\|\mathbf{r}\|^2 \leq 2(\|\mathbf{z}\|^2 + \|\mathbf{G}\mathbf{h}\|^2)$, $\|\mathbf{G}\mathbf{h}\|^2 = \mathcal{O}_{\mathbb{P}}(n)$, $\|\boldsymbol{\psi}\|^2 = \mathcal{O}_{\mathbb{P}}(n)$, and $\text{df}/\text{tr}[\mathbf{V}] = \mathcal{O}_{\mathbb{P}}(1)$. Therefore, it suffices to show

$$n\mathbf{h}^{\top}\tilde{\mathbf{h}} + \|\mathbf{z}\|^2 - \sum_{i \in [n]} (z_i - \mathbf{g}_i^{\top}\mathbf{h} + \text{tr}[\mathbf{A}] \mathbf{1}_{\{i \in I\}} \psi_i)(z_i - \mathbf{g}_i^{\top}\tilde{\mathbf{h}} + \text{tr}[\tilde{\mathbf{A}}] \mathbf{1}_{\{i \in \tilde{I}\}} \tilde{\psi}_i) = \mathcal{O}_{\mathbb{P}}(1)\|\mathbf{z}\| + \mathcal{O}_{\mathbb{P}}(\sqrt{n}).$$

By simple algebra, the LHS can be decomposed into three terms $\xi_1 + \tilde{\xi}_1 + \xi_2$ with:

$$\begin{aligned} \xi_1 &= \sum_{i \in [n]} z_i(\mathbf{g}_i^{\top}\mathbf{h} - \text{tr}[\mathbf{A}] \mathbf{1}_{\{i \in I\}} \psi_i), & \tilde{\xi}_1 &= \sum_{i \in [n]} z_i(\mathbf{g}_i^{\top}\tilde{\mathbf{h}} - \text{tr}[\tilde{\mathbf{A}}] \mathbf{1}_{\{i \in \tilde{I}\}} \tilde{\psi}_i) \\ \xi_2 &= n\mathbf{h}^{\top}\tilde{\mathbf{h}} - \sum_{i \in [n]} (\mathbf{g}_i^{\top}\mathbf{h} - \text{tr}[\mathbf{A}] \mathbf{1}_{\{i \in I\}} \psi_i)(\mathbf{g}_i^{\top}\tilde{\mathbf{h}} - \text{tr}[\tilde{\mathbf{A}}] \mathbf{1}_{\{i \in \tilde{I}\}} \tilde{\psi}_i) \end{aligned}$$

For $(\xi_1, \tilde{\xi}_1)$, the derivative formula (35) and the argument in the proof of [BS22, Proposition 18] yield

$$\mathbb{E} \left[\frac{|\xi_1|}{\{\|\mathbf{h}\|^2 + p^{-1}\|\boldsymbol{\psi}\|^2\}^{1/2} \cdot \|\mathbf{z}\|} \right] + \mathbb{E} \left[\frac{|\tilde{\xi}_1|}{\{\|\tilde{\mathbf{h}}\|^2 + p^{-1}\|\tilde{\boldsymbol{\psi}}\|^2\}^{1/2} \cdot \|\mathbf{z}\|} \right] \leq C(\mu, \delta).$$

Since the denominators are $\mathcal{O}_{\mathbb{P}}(1)\|\mathbf{z}\|$, we have $\xi_1 + \tilde{\xi}_1 = \mathcal{O}_{\mathbb{P}}(1)\|\mathbf{z}\|$.

Below we bound ξ_2 using the following moment inequality, which is a variant of [Bel23, Theorem 7.2].

Lemma 16. *Let $\boldsymbol{\rho}, \tilde{\boldsymbol{\rho}} : \mathbb{R}^{K \times Q} \rightarrow \mathbb{R}^Q$ be locally Lipschitz functions and $\boldsymbol{\zeta}, \tilde{\boldsymbol{\zeta}} : \mathbb{R}^{K \times Q} \rightarrow \mathbb{R}^L$ be locally Lipschitz functions. If $(\mathbf{z}_k)_{k \in [K]}$ are i.i.d. $\mathcal{N}(\mathbf{0}_Q, \mathbf{I}_Q)$, we have*

$$\mathbb{E} \left| \frac{K\boldsymbol{\rho}^{\top}\tilde{\boldsymbol{\rho}} - \sum_{k \in [K]} (\mathbf{z}_k^{\top}\boldsymbol{\rho} - \sum_{q \in Q} \frac{\partial \rho_q}{z_{kq}})(\mathbf{z}_k^{\top}\tilde{\boldsymbol{\rho}} - \sum_{q \in Q} \frac{\partial \tilde{\rho}_q}{z_{kq}})}{\{\|\boldsymbol{\rho}\|^2 + \|\boldsymbol{\zeta}\|^2\}^{1/2} \{\|\tilde{\boldsymbol{\rho}}\|^2 + \|\tilde{\boldsymbol{\zeta}}\|^2\}^{1/2}} \right| \leq C_7(\sqrt{K}(1 + \sqrt{\mathbb{E}[\Xi + \tilde{\Xi}]) + \mathbb{E}[\Xi + \tilde{\Xi}])$$

where $\Xi := \frac{1}{\|\boldsymbol{\rho}\|^2 + \|\boldsymbol{\zeta}\|^2} \sum_{k \in [K]} \sum_{q \in [Q]} \left(\left\| \frac{\partial \boldsymbol{\rho}}{\partial g_{kq}} \right\|^2 + \left\| \frac{\partial \boldsymbol{\zeta}}{\partial g_{kq}} \right\|^2 \right)$.

(The proof of Lemma 16 is given in Appendix A.3.1.) By Lemma 16 with $(\boldsymbol{\rho}, \tilde{\boldsymbol{\rho}}) = (\mathbf{h}, \tilde{\mathbf{h}})$, $(\boldsymbol{\zeta}, \tilde{\boldsymbol{\zeta}}) = (\boldsymbol{\psi}/\sqrt{p}, \tilde{\boldsymbol{\psi}}/\sqrt{p})$, and $K = [n]$, we get

$$\mathbb{E} \left| \frac{n\mathbf{h}^{\top}\tilde{\mathbf{h}} - \sum_{i \in [n]} (\mathbf{g}_i^{\top}\mathbf{h} - \sum_{j \in [p]} \frac{\partial h_j}{\partial g_{ij}})(\mathbf{g}_i^{\top}\tilde{\mathbf{h}} - \sum_{j \in [p]} \frac{\partial \tilde{h}_j}{\partial g_{ij}})}{\{\|\mathbf{h}\|^2 + p^{-1}\|\boldsymbol{\psi}\|^2\}^{1/2} \{\|\tilde{\mathbf{h}}\|^2 + p^{-1}\|\tilde{\boldsymbol{\psi}}\|^2\}^{1/2}} \right| \leq C_8(\sqrt{n}(1 + \sqrt{\mathbb{E}[\Xi + \tilde{\Xi}]) + \mathbb{E}[\Xi + \tilde{\Xi}]).$$

where

$$\Xi = \frac{1}{\|\mathbf{h}\|^2 + p^{-1}\|\boldsymbol{\psi}\|^2} \sum_{i \in [n]} \sum_{j \in [p]} \left(\left\| \frac{\partial \mathbf{h}}{\partial g_{ij}} \right\|^2 + \frac{1}{p} \left\| \frac{\partial \boldsymbol{\psi}}{\partial g_{ij}} \right\|^2 \right) \quad (33)$$

Let us bound Ξ . By the derivative formula (29), we have

$$\sum_{ij} \|\partial_{ij}\mathbf{h}\|^2 \leq 2\|\mathbf{A}\|_{\text{F}}^2\|\boldsymbol{\psi}\|^2 + 2\|\mathbf{A}\mathbf{G}^{\top}\mathbf{D}\|_{\text{F}}^2\|\mathbf{h}\|^2, \quad \sum_{ij} \|\partial_{ij}\boldsymbol{\psi}\|^2 \leq 2\|\mathbf{D}\mathbf{G}\mathbf{A}\|_{\text{F}}^2\|\boldsymbol{\psi}\|^2 + 2\|\mathbf{V}\|_{\text{F}}^2\|\mathbf{h}\|^2$$

where $\mathbf{V} = \mathbf{D} - \mathbf{D}\mathbf{G}\mathbf{A}\mathbf{G}^\top\mathbf{D}$ and $\mathbf{D} = \text{diag}\{\text{loss}''(z_i - \mathbf{g}_i^\top\mathbf{h})\}$. By $\|\mathbf{A}\|_{\text{op}} \leq (p\mu)^{-1}$ and $\|\mathbf{D}\|_{\text{op}} \leq \|\text{loss}'\|_{\text{Lip}}$, Ξ is bounded from above by

$$\Xi \leq 2p\|\mathbf{A}\|_F^2 + 2\|\mathbf{A}\mathbf{G}^\top\mathbf{D}\|_F^2 + 2\|\mathbf{D}\mathbf{G}\mathbf{A}\|_F^2 + 2p^{-1}\|\mathbf{V}\|_F^2 \leq C(\delta, \|\text{loss}'\|_{\text{Lip}}) \cdot \|\mathbf{A}\|_{\text{op}}^2 \|\mathbf{G}\|_{\text{op}}^4 \quad (34)$$

By $\|\mathbf{A}\|_{\text{op}} \leq (p\mu)^{-1}$ in (30) and Lemma 10, we obtain $\bar{\mathbb{E}}[\Xi] \leq C_9$. This gives

$$\bar{\mathbb{E}} \left| \frac{n\mathbf{h}^\top\tilde{\mathbf{h}} - \sum_{i \in n} (\mathbf{g}_i^\top\mathbf{h} - \sum_{j \in [p]} \frac{\partial h_j}{\partial g_{ij}}) (\mathbf{g}_i^\top\tilde{\mathbf{h}} - \sum_{j \in [p]} \frac{\partial \tilde{h}_j}{\partial g_{ij}})}{\{\|\mathbf{h}\|^2 + p^{-1}\|\boldsymbol{\psi}\|^2\}^{1/2} \{\|\tilde{\mathbf{h}}\|^2 + p^{-1}\|\tilde{\boldsymbol{\psi}}\|^2\}^{1/2}} \right| \leq C_{10}\sqrt{n}$$

Since the denominator is $\mathcal{O}_{\mathbb{P}}(1)$, we have

$$n\mathbf{h}^\top\tilde{\mathbf{h}} - \sum_{i \in [n]} \left(\mathbf{g}_i^\top\mathbf{h} - \sum_{j \in [p]} \frac{\partial h_j}{\partial g_{ij}} \right) \left(\mathbf{g}_i^\top\tilde{\mathbf{h}} - \sum_{j \in [p]} \frac{\partial \tilde{h}_j}{\partial g_{ij}} \right) = \mathcal{O}_{\mathbb{P}}(\sqrt{n}).$$

By the derivative formula (29), noting that $\frac{\partial h_j}{\partial g_{ij}} = 0$ for all $i \notin I$ and $j \in [p]$, it holds that

$$\sum_{j \in [p]} \frac{\partial h_j}{\partial g_{ij}} = \begin{cases} 0 & i \notin I \\ \text{tr}[\mathbf{A}]\psi_i - \mathbf{h}^\top\mathbf{A}\mathbf{G}\mathbf{D}\mathbf{e}_i & i \in I \end{cases} \quad (35)$$

Combined with $\|\mathbf{h}^\top\mathbf{A}\mathbf{G}\mathbf{D}\|^2 = \mathcal{O}_{\mathbb{P}}(n^{-1})$, the Cauchy–Schwarz inequality leads to

$$\sum_{i \in [n]} \left(\mathbf{g}_i^\top\mathbf{h} - \sum_{j \in [p]} \frac{\partial h_j}{\partial g_{ij}} \right) \left(\mathbf{g}_i^\top\tilde{\mathbf{h}} - \sum_{j \in [p]} \frac{\partial \tilde{h}_j}{\partial g_{ij}} \right) = \sum_{i \in [n]} \left(\mathbf{g}_i^\top\mathbf{h} - \text{tr}[\mathbf{A}]\mathbf{1}_{\{i \in I\}}\psi_i \right) \left(\mathbf{g}_i^\top\tilde{\mathbf{h}} - \text{tr}[\tilde{\mathbf{A}}]\mathbf{1}_{\{i \in \tilde{I}\}}\tilde{\psi}_i \right) + \mathcal{O}_{\mathbb{P}}(1).$$

This gives $\xi_2 = \mathcal{O}_{\mathbb{P}}(\sqrt{n}) + \mathcal{O}_{\mathbb{P}}(1) = \mathcal{O}_{\mathbb{P}}(\sqrt{n})$ and completes the proof.

A.3.1 Proof of Lemma 16

Let us denote $\mathbf{f} = (\boldsymbol{\rho}^\top, \boldsymbol{\zeta}^\top)^\top \in \mathbb{R}^{Q+L}$ and $\tilde{\mathbf{f}} = (\tilde{\boldsymbol{\rho}}^\top, \tilde{\boldsymbol{\zeta}}^\top)^\top \in \mathbb{R}^{Q \times L}$. Let us write the product as

$$\frac{\boldsymbol{\rho}^\top\tilde{\boldsymbol{\rho}}}{\|\mathbf{f}\|\|\tilde{\mathbf{f}}\|} = \left\| \frac{1}{2} \left(\frac{\boldsymbol{\rho}}{\|\mathbf{f}\|} + \frac{\tilde{\boldsymbol{\rho}}}{\|\tilde{\mathbf{f}}\|} \right) \right\|^2 - \left\| \frac{1}{2} \left(\frac{\boldsymbol{\rho}}{\|\mathbf{f}\|} - \frac{\tilde{\boldsymbol{\rho}}}{\|\tilde{\mathbf{f}}\|} \right) \right\|^2$$

Note that $\frac{\boldsymbol{\rho}}{\|\mathbf{f}\|}$ and $\frac{\tilde{\boldsymbol{\rho}}}{\|\tilde{\mathbf{f}}\|}$ are bounded by 1 in the standard Euclid norm $\|\cdot\|$. Applying the χ -square type moment inequality [Bel23, Theorem 7.2] to these terms respectively, we observe that the following two terms are bounded from above by $\sqrt{K}\{1 + \mathbb{E}[\Xi + \tilde{\Xi}]\}^{1/2} + \mathbb{E}[\Xi + \tilde{\Xi}]$ up to some universal constant:

$$\begin{aligned} & \mathbb{E} \left| K \left\| \left(\frac{\boldsymbol{\rho}}{\|\mathbf{f}\|} + \frac{\tilde{\boldsymbol{\rho}}}{\|\tilde{\mathbf{f}}\|} \right) \right\|^2 - \sum_{k \in [K]} \left(\mathbf{z}_k^\top \left(\frac{\boldsymbol{\rho}}{\|\mathbf{f}\|} + \frac{\tilde{\boldsymbol{\rho}}}{\|\tilde{\mathbf{f}}\|} \right) - \sum_{q \in [Q]} \frac{\partial}{\partial z_{kq}} \left(\frac{\rho_q}{\|\mathbf{f}\|} + \frac{\tilde{\rho}_q}{\|\tilde{\mathbf{f}}\|} \right) \right)^2 \right| \\ & \mathbb{E} \left| K \left\| \left(\frac{\boldsymbol{\rho}}{\|\mathbf{f}\|} - \frac{\tilde{\boldsymbol{\rho}}}{\|\tilde{\mathbf{f}}\|} \right) \right\|^2 - \sum_{k \in [K]} \left(\mathbf{z}_k^\top \left(\frac{\boldsymbol{\rho}}{\|\mathbf{f}\|} - \frac{\tilde{\boldsymbol{\rho}}}{\|\tilde{\mathbf{f}}\|} \right) - \sum_{q \in [Q]} \frac{\partial}{\partial z_{kq}} \left(\frac{\rho_q}{\|\mathbf{f}\|} - \frac{\tilde{\rho}_q}{\|\tilde{\mathbf{f}}\|} \right) \right)^2 \right| \end{aligned}$$

Thus, by the triangle inequality, we are left with the bound of cross terms:

$$\mathbb{E} \left| K \frac{\boldsymbol{\rho}^\top\tilde{\boldsymbol{\rho}}}{\|\mathbf{f}\|\|\tilde{\mathbf{f}}\|} - \sum_{k \in [K]} b_k \tilde{b}_k \right| \lesssim \sqrt{K} \{1 + \mathbb{E}[\Xi + \tilde{\Xi}]\}^{1/2} + \mathbb{E}[\Xi + \tilde{\Xi}] \quad (36)$$

where $b_k = \mathbf{z}_k^\top \boldsymbol{\rho} / \|\mathbf{f}\| - \sum_{q \in [Q]} (\partial / \partial z_{kq})(\rho_q / \|\mathbf{f}\|)$. Expanding the derivative of the second term, b_k can be written as

$$b_k = \underbrace{\frac{\mathbf{z}_k^\top \boldsymbol{\rho}}{\|\mathbf{f}\|}}_{=: a_k} - \sum_{q \in [Q]} \rho_q \frac{\partial}{\partial z_{kq}} \frac{1}{\|\mathbf{f}\|} = a_k + \sum_{q \in [Q]} \rho_q \frac{\mathbf{f}^\top}{\|\mathbf{f}\|^3} \frac{\partial \mathbf{f}}{\partial z_{kq}}$$

Thus, by multiple applications of Cauchy–Schwarz inequality, using $\|\boldsymbol{\rho}\|^2 \leq \|\mathbf{f}\|^2$, we find that the error $\|\mathbf{b} - \mathbf{a}\|^2 = \sum_{k \in [K]} (b_k - a_k)^2$ is bounded from above by 2Ξ :

$$\|\mathbf{a} - \mathbf{b}\|^2 = \sum_{k \in [K]} \left(\sum_{q \in [Q]} \rho_q \frac{\mathbf{f}^\top}{\|\mathbf{f}\|^3} \frac{\partial \mathbf{f}}{\partial z_{kq}} \right)^2 \leq \sum_k \|\boldsymbol{\rho}\|^2 \sum_q \left(\frac{\mathbf{f}^\top}{\|\mathbf{f}\|^3} \frac{\partial \mathbf{f}}{\partial z_{kq}} \right)^2 \leq \frac{1}{\|\mathbf{f}\|^2} \sum_{k,q} \left\| \frac{\partial \mathbf{f}}{\partial z_{kq}} \right\|^2 \leq 2\Xi.$$

The same argument gives $\|\tilde{\mathbf{a}} - \tilde{\mathbf{b}}\|^2 \leq 2\tilde{\Xi}$.

Now we claim the following deterministic inequality for all $\mathbf{u}, \tilde{\mathbf{u}}, \mathbf{a}, \tilde{\mathbf{a}}, \mathbf{b}, \tilde{\mathbf{b}} \in \mathbb{R}^K$ with $\|\mathbf{u}\| \vee \|\tilde{\mathbf{u}}\| \leq 1$:

$$\begin{aligned} \left| |K\mathbf{u}^\top \tilde{\mathbf{u}} - \mathbf{a}^\top \tilde{\mathbf{a}}| - |K\mathbf{u}^\top \tilde{\mathbf{u}} - \mathbf{b}^\top \tilde{\mathbf{b}}| \right| &\leq (\|\mathbf{a} - \mathbf{b}\|^2 + \|\tilde{\mathbf{a}} - \tilde{\mathbf{b}}\|^2) + \sqrt{K}(\|\mathbf{a} - \mathbf{b}\| + \|\tilde{\mathbf{a}} - \tilde{\mathbf{b}}\|) \\ &\quad + 2^{-1}(|K\|\mathbf{u}\|^2 - \|\mathbf{b}\|^2| + |K\|\tilde{\mathbf{u}}\|^2 - \|\tilde{\mathbf{b}}\|^2|) \end{aligned} \quad (37)$$

We prove this inequality later. Applying this inequality with $\mathbf{u} = \boldsymbol{\rho} / \|\mathbf{f}\|$ and (\mathbf{a}, \mathbf{b}) defined above, using $\|\mathbf{a} - \mathbf{b}\|^2 \leq 2\Xi$, we get

$$\begin{aligned} |K\mathbf{u}^\top \tilde{\mathbf{u}} - \mathbf{a}^\top \tilde{\mathbf{a}}| &\leq 2(\Xi + \tilde{\Xi}) + \sqrt{2K}(\Xi^{1/2} + \tilde{\Xi}^{1/2}) \\ &\quad + 2^{-1}(|K\|\mathbf{u}\|^2 - \|\mathbf{b}\|^2| + |K\|\tilde{\mathbf{u}}\|^2 - \|\tilde{\mathbf{b}}\|^2|) \\ &\quad + |K\mathbf{u}^\top \tilde{\mathbf{u}} - \mathbf{b}^\top \tilde{\mathbf{b}}| \end{aligned}$$

Taking the expectation, using the moment bound (36), we are left with

$$\begin{aligned} \mathbb{E}|K\mathbf{u}^\top \tilde{\mathbf{u}} - \mathbf{a}^\top \tilde{\mathbf{a}}| &\lesssim \mathbb{E}[\Xi + \tilde{\Xi}] + \sqrt{K}\mathbb{E}[(\Xi^{1/2} + \tilde{\Xi}^{1/2})] \\ &\quad + \{1 + \mathbb{E}[\Xi]\}^{1/2} + \mathbb{E}[\Xi] + \{1 + \mathbb{E}[\tilde{\Xi}]\}^{1/2} + \mathbb{E}[\tilde{\Xi}] \\ &\quad + \{1 + \mathbb{E}[\Xi + \tilde{\Xi}]\}^{1/2} + \mathbb{E}[\Xi + \tilde{\Xi}] \end{aligned}$$

Using Jensen's inequality $\mathbb{E}[X^{1/2}] \leq \sqrt{\mathbb{E}[X]}$ for any non-negative random variable X and $\sqrt{a} + \sqrt{b} \leq \sqrt{2}\sqrt{a+b}$ for any non-negative scalars a, b , the RHS is bounded from above by $\sqrt{K}\{1 + \mathbb{E}[\Xi + \tilde{\Xi}]\}^{1/2} + \mathbb{E}[\Xi + \tilde{\Xi}]$ up to some universal constant. This finishes the proof.

Below we prove the deterministic inequality (37). By multiple applications of the triangle inequality and Cauchy–Schwarz inequality,

$$\begin{aligned} \left| |K\mathbf{u}^\top \tilde{\mathbf{u}} - \mathbf{a}^\top \tilde{\mathbf{a}}| - |K\mathbf{u}^\top \tilde{\mathbf{u}} - \mathbf{b}^\top \tilde{\mathbf{b}}| \right| &\leq |\mathbf{a}^\top \tilde{\mathbf{a}} - \mathbf{b}^\top \tilde{\mathbf{b}}| \\ &\leq |(\mathbf{a} - \mathbf{b})^\top (\tilde{\mathbf{a}} - \tilde{\mathbf{b}}) + (\mathbf{a} - \mathbf{b})^\top \tilde{\mathbf{b}} + \mathbf{b}^\top (\tilde{\mathbf{a}} - \tilde{\mathbf{b}})| \\ &\leq \|\mathbf{a} - \mathbf{b}\| \|\tilde{\mathbf{a}} - \tilde{\mathbf{b}}\| + \|\mathbf{a} - \mathbf{b}\| \|\tilde{\mathbf{b}}\| + \|\tilde{\mathbf{a}} - \tilde{\mathbf{b}}\| \|\mathbf{b}\| \\ &\leq \frac{\|\mathbf{a} - \mathbf{b}\|^2 + \|\tilde{\mathbf{a}} - \tilde{\mathbf{b}}\|^2}{2} + \|\mathbf{a} - \mathbf{b}\| \|\tilde{\mathbf{b}}\| + \|\tilde{\mathbf{a}} - \tilde{\mathbf{b}}\| \|\mathbf{b}\| \end{aligned}$$

Using $\|\mathbf{b}\| \leq \sqrt{\|\tilde{\mathbf{b}}\|^2 - K\|\tilde{\mathbf{u}}\|^2} + \sqrt{K}\|\tilde{\mathbf{u}}\|$ and $\|\mathbf{u}\| \leq 1$, $\|\mathbf{a} - \mathbf{b}\| \|\tilde{\mathbf{b}}\|$ can be bounded from above as

$$\|\mathbf{a} - \mathbf{b}\| \|\tilde{\mathbf{b}}\| \leq \|\mathbf{a} - \mathbf{b}\| \sqrt{\|\tilde{\mathbf{b}}\|^2 - K\|\tilde{\mathbf{u}}\|^2} + \sqrt{K}\|\mathbf{a} - \mathbf{b}\| \leq \frac{\|\mathbf{a} - \mathbf{b}\|^2}{2} + \frac{\|\tilde{\mathbf{b}}\|^2 - K\|\tilde{\mathbf{u}}\|^2}{2} + \sqrt{K}\|\mathbf{a} - \mathbf{b}\|.$$

The same argument gives $\|\tilde{\mathbf{a}} - \tilde{\mathbf{b}}\| \leq \frac{\|\tilde{\mathbf{a}} - \tilde{\mathbf{b}}\|^2}{2} + \frac{\|\mathbf{b}\|^{2-K}\|\mathbf{u}\|^2}{2} + \sqrt{K}\|\tilde{\mathbf{a}} - \tilde{\mathbf{b}}\|$. Putting them all together, we obtained the desired upper bound of $\|K\mathbf{u}^\top \tilde{\mathbf{u}} - \mathbf{a}^\top \tilde{\mathbf{a}}\| - \|K\mathbf{u}^\top \tilde{\mathbf{u}} - \mathbf{b}^\top \tilde{\mathbf{b}}\|$.

A.4 Proof of trace convergence

We assume without loss of generality that $I = [n]$. Throughout this section, we denote $\mathbf{h} = \hat{\mathbf{h}}_{\mu, K}$ and $\boldsymbol{\psi} = \boldsymbol{\psi}_{\mu, K}$ for simplicity.

Lemma 17. *For any $\mu \in (0, 1]$ with $\mu^{-1} = o(n^{1/4})$,*

$$\begin{aligned} \text{tr}[\mathbf{V}] &= \|\mathbf{h}\|^{-2} (\|\boldsymbol{\psi}\|^2 \text{tr}[\mathbf{A}] - \boldsymbol{\psi}^\top \mathbf{G} \mathbf{h}) + \mathcal{O}_{\mathbb{P}}(n^{1/2} \mu^{-1}) \\ \text{tr}[\mathbf{A}]^2 &= \frac{(\boldsymbol{\psi}^\top \mathbf{G} \mathbf{h})^2 + \|\mathbf{h}\|^2 (p \|\boldsymbol{\psi}\|^2 - \|\mathbf{G}^\top \boldsymbol{\psi}\|^2)}{\|\boldsymbol{\psi}\|^4} + \mathcal{O}_{\mathbb{P}}(n^{-1/2} \mu^{-2}) \end{aligned}$$

and $\mathbb{P}(\text{tr}[\mathbf{A}]^2 \leq C) \rightarrow 1$ for a constant $C > 0$ that only depend on (α, β, δ) .

Proof. First we show the stochastic representation of $\text{tr}[\mathbf{V}]$ by $(\mathbf{h}, \boldsymbol{\psi}, \text{tr}[\mathbf{A}])$.

Lemma 18 ([BS22]). *Let $\boldsymbol{\rho} : \mathbb{R}^{K \times Q} \rightarrow \mathbb{R}^Q$ and $\boldsymbol{\zeta} : \mathbb{R}^{K \times Q} \rightarrow \mathbb{R}^K$ be two locally Lipschitz functions with differentiable components. If $\mathbf{Z} \in \mathbb{R}^{K \times Q}$ has i.i.d. $\mathcal{N}(0, 1)$ entries, we have*

$$\mathbb{E} \left[\left(\frac{\boldsymbol{\zeta}^\top \mathbf{Z} \boldsymbol{\rho} - \sum_{kq} \frac{\partial(\zeta_k \rho_q)}{\partial g_{kq}}}{\|\boldsymbol{\zeta}\|^2 + \|\boldsymbol{\rho}\|^2} \right)^2 \right] \leq C_{11} (1 + \mathbb{E}[\Xi])$$

where $\Xi := \frac{1}{\|\boldsymbol{\rho}\|^2 + \|\boldsymbol{\zeta}\|^2} \sum_{k \in [K]} \sum_{q \in [Q]} \left(\left\| \frac{\partial \boldsymbol{\rho}}{\partial g_{kq}} \right\|^2 + \left\| \frac{\partial \boldsymbol{\zeta}}{\partial g_{kq}} \right\|^2 \right)$.

By Lemma 18 with $(\boldsymbol{\zeta}, \boldsymbol{\rho}) = (p^{-1/2} \boldsymbol{\psi}, \mathbf{h})$ and $\mathbf{Z} = \mathbf{G}$, we have

$$\mathbb{E} \left[\left(\frac{\frac{1}{\sqrt{p}} \boldsymbol{\psi}^\top \mathbf{G} \mathbf{h} - \frac{1}{\sqrt{p}} \sum_{ij} \frac{\partial(\psi_i h_j)}{\partial g_{ij}}}{\|\mathbf{h}\|^2 + p^{-1} \|\boldsymbol{\psi}\|^2} \right)^2 \right] \leq C_{12} (1 + \mathbb{E}[\Xi])$$

where Ξ is the same estimate as in (33). Using the upper estimate $\Xi \leq C_{13} \|\mathbf{A}\|_{\text{op}}^2 \|\mathbf{G}\|_{\text{op}}^4$ in (34), $\|\mathbf{A}\|_{\text{op}} \leq (p\mu)^{-1}$ and $\mathbb{E}[\|\mathbf{G}\|^4] \leq C(\delta)n^2$, we find that the RHS in the above display is bounded from away by $C(\delta)(1 + \mu^{-2})$. Since the denominator $\|\mathbf{h}\|^2 + p^{-1} \|\boldsymbol{\psi}\|^2$ is $\mathcal{O}_{\mathbb{P}}(1)$, we get

$$\boldsymbol{\psi}^\top \mathbf{G} \mathbf{h} - \sum_{ij} \frac{\partial \psi_i h_j}{\partial g_{ij}} = \sqrt{p} \cdot \mathcal{O}_{\mathbb{P}}(1) \cdot \mathcal{O}_{\mathbb{P}}(\mu^{-1}) = \mathcal{O}_{\mathbb{P}}(\sqrt{n} \mu^{-1}).$$

For the sum of derivative $\sum_{ij} \frac{\partial(\psi_i h_j)}{\partial g_{ij}}$, using $\|\mathbf{A}\|_{\text{op}} \leq (p\mu)^{-1}$ and $\|\mathbf{G}\|_{\text{op}} = \mathcal{O}_{\mathbb{P}}(\sqrt{n})$, we have

$$\begin{aligned} \sum_{ij} \frac{\partial(\psi_i h_j)}{\partial g_{ij}} &= \|\boldsymbol{\psi}\|^2 \text{tr}[\mathbf{A}] - \mathbf{h}^\top \mathbf{G}^\top \mathbf{D} \boldsymbol{\psi} - \boldsymbol{\psi}^\top \mathbf{D} \mathbf{G} \mathbf{A} \mathbf{h} - \|\mathbf{h}\|^2 \text{tr}[\mathbf{V}] \\ &= \|\boldsymbol{\psi}\|^2 \text{tr}[\mathbf{A}] + \mathcal{O}_{\mathbb{P}}(n^{1/2}) + \mathcal{O}_{\mathbb{P}}(\mu^{-1}) - \|\mathbf{h}\|^2 \text{tr}[\mathbf{V}] \end{aligned}$$

so that we are left with

$$\boldsymbol{\psi}^\top \mathbf{G} \mathbf{h} - \|\boldsymbol{\psi}\|^2 \text{tr}[\mathbf{A}] + \|\mathbf{h}\|^2 \text{tr}[\mathbf{V}] = \mathcal{O}_{\mathbb{P}}(\sqrt{n} \mu^{-1}).$$

Dividing by $\|\mathbf{h}\|^2$, noting $\|\mathbf{h}\|^{-1} = \mathcal{O}_{\mathbb{P}}(1)$ from $\|\mathbf{h}\|^2 \xrightarrow{\mathbb{P}} \alpha^2 > 0$, we obtain the representation of $\text{tr}[\mathbf{V}]$.

Next, we show the stochastic representation of $\text{tr}[\mathbf{A}^2]$ by $(\mathbf{G}, \mathbf{h}, \boldsymbol{\psi})$. By the stochastic representation of $\text{tr}[\mathbf{V}]$,

$$\begin{aligned} \left| \|\boldsymbol{\psi}\|^4 \text{tr}[\mathbf{A}]^2 - (\boldsymbol{\psi}^\top \mathbf{G}\mathbf{h} + \text{tr}[\mathbf{V}]\|\mathbf{h}\|^2)^2 \right| &\leq \mathcal{O}_{\mathbb{P}}(\sqrt{n}\mu^{-1}) \left| \|\boldsymbol{\psi}\|^2 \text{tr}[\mathbf{A}] + \boldsymbol{\psi}^\top \mathbf{G}\mathbf{h} + \text{tr}[\mathbf{V}]\|\mathbf{h}\|^2 \right| \\ &\leq \mathcal{O}_{\mathbb{P}}(\sqrt{n}\mu^{-1}) \mathcal{O}_{\mathbb{P}}(n\mu^{-1} + n + n) \\ &= \mathcal{O}_{\mathbb{P}}(n^{3/2}\mu^{-2}). \end{aligned}$$

By Lemma 16 with $\boldsymbol{\rho} = \tilde{\boldsymbol{\rho}} = \boldsymbol{\psi}/\sqrt{p}$ and $\boldsymbol{\zeta} = \tilde{\boldsymbol{\zeta}} = \mathbf{h}$, we have

$$\mathbb{E} \left[\frac{\left| p \|\frac{\boldsymbol{\psi}}{\sqrt{p}}\|^2 - \sum_{j \in [p]} \left(\frac{\boldsymbol{\psi}^\top \mathbf{G}\mathbf{e}_j}{\sqrt{p}} - \frac{1}{\sqrt{p}} \sum_{i \in [n]} \frac{\partial \psi_i}{\partial g_{ij}} \right)^2 \right|}{\|\mathbf{h}\|^2 + p^{-1}\|\boldsymbol{\psi}\|^2} \right] \leq C_{14}(\sqrt{n}(1 + \mathbb{E}[\Xi]^{1/2}) + \mathbb{E}[\Xi]),$$

where Ξ is the same estimate as in (33). Using $\mathbb{E}[\Xi] \leq C_{15}\mu^{-2}$ again, we get

$$\|\boldsymbol{\psi}\|^2 - \frac{1}{p} \sum_{j \in [p]} \left(\boldsymbol{\psi}^\top \mathbf{G}\mathbf{e}_j - \sum_{i \in [n]} \frac{\partial \psi_i}{\partial g_{ij}} \right)^2 = \mathcal{O}_{\mathbb{P}}(1) \cdot \mathcal{O}_{\mathbb{P}}(\sqrt{n}(1 + \mu^{-1}) + \mu^{-2}) = (n^{1/2}\mu^{-1})$$

Using $\sum_{i \in [n]} \frac{\partial \psi_i}{\partial g_{ij}} = -\boldsymbol{\psi}^\top \mathbf{D}\mathbf{G}\mathbf{A}\mathbf{e}_j - \text{tr}[\mathbf{V}]h_j$ by the derivative formula (29), it holds that

$$\begin{aligned} &\left\| \sum_{j \in [p]} \left(\boldsymbol{\psi}^\top \mathbf{G}\mathbf{e}_j - \sum_{i \in [n]} \frac{\partial \psi_i}{\partial g_{ij}} \right)^2 - \|\mathbf{G}^\top \boldsymbol{\psi} + \text{tr}[\mathbf{V}]\mathbf{h}\|^2 \right\| \\ &= \left\| \mathbf{G}^\top \boldsymbol{\psi} + \mathbf{A}^\top \mathbf{G}^\top \mathbf{D}\boldsymbol{\psi} + \text{tr}[\mathbf{V}]\mathbf{h} \right\|^2 - \|\mathbf{G}^\top \boldsymbol{\psi} + \text{tr}[\mathbf{V}]\mathbf{h}\|^2 \\ &\leq \|\mathbf{A}^\top \mathbf{G}^\top \mathbf{D}\boldsymbol{\psi}\| \left(\|\mathbf{A}^\top \mathbf{G}^\top \mathbf{D}\boldsymbol{\psi}\| + 2\|\mathbf{G}^\top \boldsymbol{\psi} + \text{tr}[\mathbf{V}]\mathbf{h}\| \right) \\ &= \mathcal{O}_{\mathbb{P}}(\mu^{-1})(\mathcal{O}_{\mathbb{P}}(\mu^{-1}) + \mathcal{O}_{\mathbb{P}}(n)) = \mathcal{O}_{\mathbb{P}}(\mu^{-1}n). \end{aligned}$$

Combining the above displays, we are left with

$$\|\boldsymbol{\psi}\|^2 - p^{-1}\|\mathbf{G}^\top \boldsymbol{\psi} + \text{tr}[\mathbf{V}]\mathbf{h}\|^2 = \mathcal{O}_{\mathbb{P}}(n^{1/2}\mu^{-1}) + \mathcal{O}_{\mathbb{P}}(\mu^{-1}n^{-1}) = \mathcal{O}_{\mathbb{P}}(n^{1/2}\mu^{-1}).$$

Therefore,

$$\begin{aligned} \|\boldsymbol{\psi}\|^4 \text{tr}[\mathbf{A}]^2 &= (\boldsymbol{\psi}^\top \mathbf{G}\mathbf{h} + \text{tr}[\mathbf{V}]\|\mathbf{h}\|^2)^2 + \mathcal{O}_{\mathbb{P}}(n^{3/2}\mu^{-2}) \\ &= (\boldsymbol{\psi}^\top \mathbf{G}\mathbf{h})^2 + \|\mathbf{h}\|^2(2\text{tr}[\mathbf{V}]\boldsymbol{\psi}^\top \mathbf{G}\mathbf{h} + \text{tr}[\mathbf{V}]^2\|\mathbf{h}\|^2) + \mathcal{O}_{\mathbb{P}}(n^{3/2}\mu^{-2}) \\ &= (\boldsymbol{\psi}^\top \mathbf{G}\mathbf{h})^2 + \|\mathbf{h}\|^2(\|\mathbf{G}^\top \boldsymbol{\psi} + \text{tr}[\mathbf{V}]\mathbf{h}\|^2 - \|\mathbf{G}^\top \boldsymbol{\psi}\|^2) + \mathcal{O}_{\mathbb{P}}(n^{3/2}\mu^{-2}) \\ &= (\boldsymbol{\psi}^\top \mathbf{G}\mathbf{h})^2 + \|\mathbf{h}\|^2(p\|\boldsymbol{\psi}\|^2 - \|\mathbf{G}^\top \boldsymbol{\psi}\|^2) + \mathcal{O}_{\mathbb{P}}(n^{3/2}\mu^{-1}) + \mathcal{O}_{\mathbb{P}}(n^{3/2}\mu^{-2}). \end{aligned}$$

Multiplying by $\|\boldsymbol{\psi}\|^{-4}$, which is $\mathcal{O}_{\mathbb{P}}(n^2)$ since $\|\boldsymbol{\psi}\|^2/p \xrightarrow{\mathbb{P}} \beta^2 > 0$, we obtain this representation of $\text{tr}[\mathbf{A}]^2$. The upper bound $\mathbb{P}(\text{tr}[\mathbf{A}]^2 \leq C_{16})$ follows from the stochastic representation and the convergences: $\|\boldsymbol{\psi}\|^2/p \xrightarrow{\mathbb{P}} \beta^2 > 0$, $\|\mathbf{h}\|^2 \xrightarrow{\mathbb{P}} \alpha^2 > 0$, $\|\mathbf{G}\|_{\text{op}}/\sqrt{n} \xrightarrow{\mathbb{P}} 1 + \sqrt{\delta} > 0$. \square

Lemma 19. *Letting $\bar{\mathbb{E}}[\cdot] = \mathbb{E}[\cdot | \mathbf{z}, \boldsymbol{\theta}]$ be the conditional expectation with respect to the design matrix \mathbf{G} , for any $\mu \in (0, 1]$ with $\mu^{-2} = o(n)$, we have*

$$\frac{(\boldsymbol{\psi}^\top \mathbf{G}\mathbf{h})^2 + \|\mathbf{h}\|^2(p\|\boldsymbol{\psi}\|^2 - \|\mathbf{G}^\top \boldsymbol{\psi}\|^2)}{\|\boldsymbol{\psi}\|^4} = \frac{(\bar{\mathbb{E}}[\boldsymbol{\psi}^\top \mathbf{G}\mathbf{h}])^2 + \bar{\mathbb{E}}[\|\mathbf{h}\|^2](p\bar{\mathbb{E}}[\|\boldsymbol{\psi}\|^2 - \|\mathbf{G}^\top \boldsymbol{\psi}\|^2])}{\bar{\mathbb{E}}[\|\boldsymbol{\psi}\|^2]^2} + \mathcal{O}_{\mathbb{P}}(n^{-1/2}\mu^{-1}).$$

Proof. By the Gaussian Poincaré inequality, we claim the following:

$$\begin{aligned}\bar{\mathbb{E}}[(\|\boldsymbol{\psi}\|^2 - \bar{\mathbb{E}}[\|\boldsymbol{\psi}\|^2])^2] &\leq C_{17}n\mu^{-2}, \\ \bar{\mathbb{E}}[(\|\mathbf{h}\|^2 - \bar{\mathbb{E}}[\|\mathbf{h}\|^2])^2] &\leq C_{18}n^{-1}\mu^{-2}, \\ \bar{\mathbb{E}}[(\boldsymbol{\psi}^\top \mathbf{G}\mathbf{h} - \bar{\mathbb{E}}[\boldsymbol{\psi}^\top \mathbf{G}\mathbf{h}])^2] &\leq C_{19}n\mu^{-2}, \\ \bar{\mathbb{E}}[(\|\mathbf{G}^\top \boldsymbol{\psi}\|^2 - \bar{\mathbb{E}}[\|\mathbf{G}^\top \boldsymbol{\psi}\|^2])^2] &\leq C_{20}n^3\mu^{-2}.\end{aligned}$$

First and second moment inequalities immediately follow from Lemma 13 with $\tilde{I} = I = [n]$. For the third moment inequality, the Gaussian Poincaré inequality gives the upper bound $\mathbb{E}[(\boldsymbol{\psi}^\top \mathbf{G}\mathbf{h} - \bar{\mathbb{E}}[\boldsymbol{\psi}^\top \mathbf{G}\mathbf{h}])^2] \leq \sum_{i,j} \mathbb{E}\left(\frac{\partial \boldsymbol{\psi}^\top \mathbf{G}\mathbf{h}}{\partial g_{ij}}\right)^2$, where

$$\begin{aligned}\text{RHS}_{ij} &:= \frac{\partial \boldsymbol{\psi}^\top \mathbf{G}\mathbf{h}}{\partial g_{ij}} = \left(-\mathbf{D}\mathbf{G}\mathbf{A}\mathbf{e}_j\psi_i - \mathbf{V}\mathbf{e}_i h_j\right)^\top \mathbf{G}\mathbf{h} + \psi_i h_j + \boldsymbol{\psi}^\top \mathbf{G}\mathbf{A}\left(\mathbf{e}_j\psi_i - \mathbf{G}^\top \mathbf{D}\mathbf{e}_i h_j\right) \\ &= \mathbf{e}_j^\top \mathbf{A}^\top \left(-\mathbf{G}^\top \mathbf{D}\mathbf{G}\mathbf{h} + \mathbf{G}^\top \boldsymbol{\psi}\right) \psi_i - \mathbf{e}_i^\top \left(\mathbf{V}^\top \mathbf{G}\mathbf{h} + \mathbf{D}\mathbf{G}\mathbf{A}^\top \mathbf{G}^\top \boldsymbol{\psi}\right) h_j + \psi_i h_j\end{aligned}$$

By $\mathbf{V} = \mathbf{D} - \mathbf{D}\mathbf{G}\mathbf{A}\mathbf{G}^\top \mathbf{D}$ and $\|\mathbf{A}\|_{\text{op}} \leq (p\mu)^{-1}$ in (30), we have $\sum_{i,j} \mathbb{E}[\text{RHS}_{ij}^2] \leq C_{21}n\mu^{-2}$. Next,

$$\mathbb{E}\left[\left(\|\mathbf{G}^\top \boldsymbol{\psi}\|^2 - \bar{\mathbb{E}}[\|\mathbf{G}^\top \boldsymbol{\psi}\|^2]\right)^2\right] \leq \sum_{i,j} \mathbb{E}\left(\frac{\partial \|\mathbf{G}^\top \boldsymbol{\psi}\|^2}{\partial g_{ij}}\right)^2 = 4 \sum_{i,j} \mathbb{E}(\boldsymbol{\psi}^\top \mathbf{G} \frac{\partial \mathbf{G}^\top \boldsymbol{\psi}}{\partial g_{ij}})^2$$

where

$$\begin{aligned}\boldsymbol{\psi}^\top \mathbf{G} \frac{\partial \mathbf{G}^\top \boldsymbol{\psi}}{\partial g_{ij}} &= \boldsymbol{\psi}^\top \mathbf{G}\mathbf{e}_j\psi_i + \boldsymbol{\psi}^\top \mathbf{G}\mathbf{G}^\top (-\mathbf{D}\mathbf{G}\mathbf{A}\mathbf{e}_j\psi_i - \mathbf{V}\mathbf{e}_i h_j) \\ &= \boldsymbol{\psi}^\top \mathbf{G}(\mathbf{I}_p - \mathbf{G}^\top \mathbf{D}\mathbf{G}\mathbf{A})\mathbf{e}_j\psi_i - \boldsymbol{\psi}^\top \mathbf{G}\mathbf{G}^\top \mathbf{V}\mathbf{e}_i h_j\end{aligned}$$

so $\sum_{i,j} \mathbb{E}[\text{RHS}_{ij}^2] \leq C_{22}n^3\mu^{-2}$. Thus, we get

$$\begin{aligned}\|\boldsymbol{\psi}\|^2 &= \bar{\mathbb{E}}[\|\boldsymbol{\psi}\|^2] + \mathcal{O}_{\mathbb{P}}(n^{1/2}\mu^{-1}) \\ \|\mathbf{h}\|^2 &= \bar{\mathbb{E}}[\|\mathbf{h}\|^2] + \mathcal{O}_{\mathbb{P}}(n^{-1/2}\mu^{-1}) \\ \boldsymbol{\psi}^\top \mathbf{G}\mathbf{h} &= \bar{\mathbb{E}}[\boldsymbol{\psi}^\top \mathbf{G}\mathbf{h}] + \mathcal{O}_{\mathbb{P}}(n^{1/2}\mu^{-1}) \\ \|\mathbf{G}^\top \boldsymbol{\psi}\|^2 &= \bar{\mathbb{E}}[\|\mathbf{G}^\top \boldsymbol{\psi}\|^2] + \mathcal{O}_{\mathbb{P}}(n^{3/2}\mu^{-1})\end{aligned}$$

Since $\mu^{-2} = o(n)$, the concentration of $\|\boldsymbol{\psi}\|^2$ on the conditional expectation $\bar{\mathbb{E}}[\|\boldsymbol{\psi}\|^2]$ and the convergence $\|\boldsymbol{\psi}\|^2/p \xrightarrow{\mathbb{P}} \beta^2$ yield $\bar{\mathbb{E}}[\|\boldsymbol{\psi}\|^2]/p \xrightarrow{\mathbb{P}} \beta^2 > 0$. This implies that $1/\|\boldsymbol{\psi}\|^2$ and $1/\bar{\mathbb{E}}[\|\boldsymbol{\psi}\|^2]$ are $\mathcal{O}_{\mathbb{P}}(n^{-1})$. Then, the error from replacing the denominator $\|\boldsymbol{\psi}\|^4$ by $\bar{\mathbb{E}}[\|\boldsymbol{\psi}\|^2]^2$ is estimated as

$$\begin{aligned}&(\|\boldsymbol{\psi}\|^{-4} - \bar{\mathbb{E}}[\|\boldsymbol{\psi}\|^2]^{-2})(\boldsymbol{\psi}^\top \mathbf{G}\mathbf{h})^2 + \|\mathbf{h}\|^2(p\|\boldsymbol{\psi}\|^2 - \|\mathbf{G}^\top \boldsymbol{\psi}\|^2) \\ &= \frac{(\|\boldsymbol{\psi}\|^2 - \bar{\mathbb{E}}[\|\boldsymbol{\psi}\|^2])(\|\boldsymbol{\psi}\|^2 + \bar{\mathbb{E}}[\|\boldsymbol{\psi}\|^2])}{\|\boldsymbol{\psi}\|^4 \bar{\mathbb{E}}[\|\boldsymbol{\psi}\|^2]^2} \cdot \mathcal{O}_{\mathbb{P}}(n^2) \\ &= \mathcal{O}_{\mathbb{P}}(n^{1/2}\mu^{-1})\mathcal{O}_{\mathbb{P}}(n \cdot n^{-4})\mathcal{O}_{\mathbb{P}}(n^2) = \mathcal{O}_{\mathbb{P}}(n^{-1/2}\mu^{-1}).\end{aligned}$$

For the error from replacing the numerator with the conditional one, the error of replacing each term with the conditional one is given by

$$(\boldsymbol{\psi}^\top \mathbf{G}\mathbf{h})^2 = (\bar{\mathbb{E}}[\boldsymbol{\psi}^\top \mathbf{G}\mathbf{h}])^2 + \mathcal{O}_{\mathbb{P}}(n^{3/2}\mu^{-1})$$

$$\begin{aligned}\|\mathbf{h}\|^2\|\boldsymbol{\psi}\|^2 &= \bar{\mathbb{E}}[\|\mathbf{h}\|^2]\bar{\mathbb{E}}[\|\boldsymbol{\psi}\|^2] + \mathcal{O}_{\mathbb{P}}(n^{1/2}\mu^{-1}) \\ \|\mathbf{h}\|^2\|\mathbf{G}^\top\boldsymbol{\psi}\|^2 &= \bar{\mathbb{E}}[\|\mathbf{h}\|^2]\bar{\mathbb{E}}[\|\mathbf{G}^\top\boldsymbol{\psi}\|^2] + \mathcal{O}_{\mathbb{P}}(n^{3/2}\mu^{-1})\end{aligned}$$

Thus, noting that $\bar{\mathbb{E}}[\|\boldsymbol{\psi}\|^2]^{-1} = \mathcal{O}_{\mathbb{P}}(n^{-1})$, we get

$$\begin{aligned}\frac{(\boldsymbol{\psi}^\top\mathbf{G}\mathbf{h})^2 + \|\mathbf{h}\|^2(p\|\boldsymbol{\psi}\|^2 - \|\mathbf{G}^\top\boldsymbol{\psi}\|^2) - (\bar{\mathbb{E}}[\boldsymbol{\psi}^\top\mathbf{G}\mathbf{h}])^2 - \bar{\mathbb{E}}[\|\mathbf{h}\|^2](p\bar{\mathbb{E}}[\|\boldsymbol{\psi}\|^2] - \bar{\mathbb{E}}[\|\mathbf{G}^\top\boldsymbol{\psi}\|^2])}{\bar{\mathbb{E}}[\|\boldsymbol{\psi}\|^2]^2} \\ = \mathcal{O}_{\mathbb{P}}(n^{-1/2}\mu^{-1}).\end{aligned}$$

This finishes the proof. \square

Lemma 20. *Suppose $\mu \in (0, 1]$ and $\mu^{-1} = O(n^{1/8})$. Then, there exist a non-negative random variables $\hat{\kappa} \geq 0$, which is independent of \mathbf{G} , such that*

$$\mathbb{P}\left(|\operatorname{tr}[\mathbf{A}] - \hat{\kappa}| \leq 2n^{-1/16} \quad \text{and} \quad |\hat{\kappa}| \leq C\right) \rightarrow 1$$

where C is a constant depending on (α, β, δ) only.

Proof. By Lemma 17 and Lemma 19, there exists a random variable A_n , which is independent from \mathbf{G} , such that

$$\operatorname{tr}[\mathbf{A}]^2 = A_n + \mathcal{O}_{\mathbb{P}}(n^{-1/2}\mu^{-2}) + \mathcal{O}_{\mathbb{P}}(n^{-1/2}\mu^{-1}) = A_n + \mathcal{O}_{\mathbb{P}}(n^{-1/4}), \quad \text{and} \quad \mathbb{P}(|A_n| \leq C) \rightarrow 1$$

Noting $\operatorname{tr}[\mathbf{A}]^2 \geq 0$ and $\mathcal{O}_{\mathbb{P}}(n^{-1/4}) = o_{\mathbb{P}}(n^{-1/8})$, this implies that the event

$$\Omega := \{|\operatorname{tr}[\mathbf{A}]^2 - A_n| \leq n^{-1/8}\} \cap \{-n^{-1/8} \leq A_n \leq C\}$$

holds with high probability. Let us take $\hat{\kappa} := \sqrt{(A_n + 2n^{-1/8})_+}$. Note in passing that under the event Ω , we have $\hat{\kappa} = \sqrt{A_n + 2n^{-1/8}}$ and $\hat{\kappa} \leq \sqrt{C+1}$. By the non-negativeness $\operatorname{tr}[\mathbf{A}] \geq 0$ and the (1/2)-Hölder continuity of the square root $\mathbb{R}_{\geq 0} \ni x \mapsto \sqrt{x}$, under the event Ω , it holds that

$$|\operatorname{tr}[\mathbf{A}] - \hat{\kappa}| = |\sqrt{\operatorname{tr}[\mathbf{A}]^2} - \sqrt{A_n + 2n^{-1/8}}| \leq |\operatorname{tr}[\mathbf{A}]^2 - A_n - 2n^{-1/8}|^{1/2}$$

and the RHS is less than $|n^{-1/8} + 2n^{-1/8}|^{1/2} \leq 2n^{-1/16}$ by the triangle inequality. \square

A.4.1 Proof of of Lemma 12

Applying Lemma 14 with $M = 1$ and $(\mathbf{z}, \mathbf{F}(\mathbf{z})) = (\mathbf{g}_i, \mathbf{h})$ for each $i \in [n]$, using $\sum_{ij} \bar{\mathbb{E}}[\|\partial_{ij}\mathbf{h}\|^2] = \mathcal{O}_{\mathbb{P}}(\mu^{-2}) = o_{\mathbb{P}}(n)$ we get

$$\sum_{i \in [n]} (\mathbf{g}_i^\top \mathbf{h} - \operatorname{tr}[\mathbf{A}]\psi_i - \|\mathbf{h}\|\hat{u}_i)^2 = o_{\mathbb{P}}(n), \quad \hat{u}_i | \boldsymbol{\theta}, \mathbf{z}, \mathbf{G}_{-i} \stackrel{d}{=} \mathcal{N}(0, 1)$$

By Lemma 20, there exists a non-negative random variable $\hat{\kappa}$ independent from \mathbf{G} such that

$$\operatorname{tr}[\mathbf{A}] = \hat{\kappa} + o_{\mathbb{P}}(1), \quad \mathbb{P}(\hat{\kappa} \leq C) \rightarrow 1$$

for a positive constant C . Then, combined with $\|\mathbf{h}\| = \alpha + o_{\mathbb{P}}(1)$, we get

$$\frac{1}{2n} \sum_{i \in [n]} (\mathbf{g}_i^\top \mathbf{h} - \hat{\kappa}\psi_i - \alpha\hat{u}_i)^2 = o_{\mathbb{P}}(1) + (\hat{\kappa} - \operatorname{tr}[\mathbf{A}])^2 \frac{\|\boldsymbol{\psi}\|^2}{n} + (\|\mathbf{h}\| - \alpha)^2 \frac{\sum_{i \in [n]} \hat{u}_i^2}{n} = o_{\mathbb{P}}(1).$$

Furthermore, using the independence of $(\hat{\kappa}, \mathbf{G})$ and the upper bound $\hat{\kappa} \leq C$, by the same argument in the proof of Theorem 2, we can easily show that the conditional expectation $\bar{\mathbb{E}}[\cdot] = \mathbb{E}[\cdot | \mathbf{z}, \boldsymbol{\theta}]$ of the square of LHS is bounded by C with high probability for some constant C . This implies that the conditional expectation $\bar{\mathbb{E}}$ of LHS is also $o(1)$, i.e.,

$$\frac{1}{n} \sum_{i \in [n]} \bar{\mathbb{E}}[(\mathbf{g}_i^\top \mathbf{h} - \hat{\kappa} \psi_i - \alpha \hat{u}_i)^2] = o_{\mathbb{P}}(1).$$

Let us define $\Xi_i := \mathbf{g}_i^\top \mathbf{h} - \hat{\kappa} \psi_i - \alpha \hat{u}_i$ so that the above display reads $n^{-1} \sum_i \bar{\mathbb{E}}[\Xi_i^2] = o_{\mathbb{P}}(1)$. Noting $\psi_i = \text{loss}'(z_i - \mathbf{g}_i^\top \mathbf{h})$ for all $i \in [n]$, the residual can be written as

$$z_i - \mathbf{g}_i^\top \mathbf{h} = \text{prox}_{\text{loss}}(z_i - \alpha \hat{u}_i - \Xi_i; \hat{\kappa})$$

for all $i \in [n]$. Since $\text{prox}_f(\cdot)$ is 1-Lipschitz for any convex function, we have

$$\frac{1}{n} \sum_{i \in [n]} \bar{\mathbb{E}}\left[\left(z_i - \mathbf{g}_i^\top \mathbf{h} - \text{prox}_{\text{loss}}(z_i - \alpha \hat{u}_i; \hat{\kappa})\right)^2\right] \leq \frac{1}{n} \sum_{i \in [n]} \bar{\mathbb{E}}[\Xi_i^2] = o_{\mathbb{P}}(1).$$

Since loss' is Lipschitz, the above display lets us approximate $\psi_i = \text{loss}'(z_i - \mathbf{g}_i^\top \mathbf{h})$ by $\text{env}'_{\text{loss}}(z_i - \alpha \hat{u}_i; \hat{\kappa})$:

$$\frac{1}{n} \sum_{i \in [n]} \bar{\mathbb{E}}\left[\left(\psi_i - \text{env}'_{\text{loss}}(z_i - \alpha \hat{u}_i; \hat{\kappa})\right)^2\right] = o_{\mathbb{P}}(1).$$

Applying this approximation to the concentration $\beta^2 = \|\boldsymbol{\psi}\|^2/p + o_{\mathbb{P}}(1) = \bar{\mathbb{E}}[\|\boldsymbol{\psi}\|^2]/p + o_{\mathbb{P}}(1)$ by Lemma 13 with $I = \tilde{I} = [n]$, we get

$$\begin{aligned} \beta^2 &= \frac{n}{p} \frac{1}{n} \sum_{i \in [n]} \bar{\mathbb{E}}[\psi_i^2] + o_{\mathbb{P}}(1) = \delta \frac{1}{n} \sum_{i \in [n]} \bar{\mathbb{E}}\left[\text{env}'_{\text{loss}}(z_i - \alpha \hat{u}_i; \hat{\kappa})^2\right] + o_{\mathbb{P}}(1) \\ &= \delta \frac{1}{n} \sum_{i \in [n]} \int_{-\infty}^{\infty} \varphi(x) \text{env}'_{\text{loss}}(z_i + \alpha x; \hat{\kappa})^2 dx + o_{\mathbb{P}}(1). \end{aligned} \quad (38)$$

where $\varphi(x)$ is the pdf of $\mathcal{N}(0, 1)$. Let us define the functions $F, \hat{F} : [0, +\infty) \rightarrow \mathbb{R}$ by

$$F(\tau) := \beta^2 - \delta \mathbb{E}[\text{env}'_{\text{loss}}(Z + \alpha G; \tau)^2], \quad \hat{F}(\tau) := \beta^2 - \delta \frac{1}{n} \sum_{i \in [n]} \int_{-\infty}^{\infty} \varphi(x) \text{env}'_{\text{loss}}(z_i + \alpha x; \tau)^2 dx$$

so that $F(\kappa) = 0$ by (5b) in System 1a (with $c = 1$), while (38) reads $\hat{F}(\hat{\kappa}) = o_{\mathbb{P}}(1)$. Note in passing that the weak law of large numbers implies $\hat{F}(\tau) \xrightarrow{\mathbb{P}} F(\tau)$ pointwise, and F and \hat{F} are strictly increasing functions in τ since $-2^{-1} \text{env}'_{\text{loss}}(x; \tau)^2$ is the derivative of the convex function $\tau \mapsto \text{env}_{\text{loss}}(x; \tau)$, which is strictly convex under Assumption D-(4) (see [TAH18, Lemma 4.4]). Then, for any $\epsilon > 0$, we have

$$F(\kappa + \epsilon) > 0 = F(\kappa) > F(\kappa - \epsilon).$$

By the pointwise convergence $\hat{F}(\tau) \xrightarrow{\mathbb{P}} F(\tau)$, it holds that

$$\hat{F}(\kappa + \epsilon) > 2^{-1} F(\kappa + \epsilon) > 0 > 2^{-1} F(\kappa - \epsilon) > \hat{F}(\kappa - \epsilon).$$

with high probability. Then, combined with $\widehat{F}(\widehat{\kappa}) = o_{\mathbb{P}}(1)$, we have

$$\widehat{F}(\kappa + \epsilon) > 2^{-1}F(\kappa + \epsilon) > \widehat{F}(\widehat{\kappa}) > 2^{-1}F(\kappa - \epsilon) > \widehat{F}(\kappa - \epsilon).$$

with high probability. Since \widehat{F} is non-decreasing with probability 1, this gives $\mathbb{P}(|\widehat{\kappa} - \kappa| \leq \epsilon) \rightarrow 1$. Since we took $\epsilon > 0$ arbitrarily, we obtain $\widehat{\kappa} \xrightarrow{\mathbb{P}} \kappa$.

Going back to the place where we replaced $\text{tr}[\mathbf{A}]$ by $\widehat{\kappa}$, now replacing $\text{tr}[\mathbf{A}]$ by κ instead, we obtain

$$\frac{1}{n} \sum_{i \in [n]} (\mathbf{g}_i^\top \mathbf{h} - \kappa \psi_i - \alpha \widehat{u}_i)^2 = o_{\mathbb{P}}(1),$$

and this approximation also holds in $\bar{\mathbb{E}}$. Thus, we get the approximation of residual and ψ_i

$$\frac{1}{n} \sum_{i \in [n]} \bar{\mathbb{E}} \left[\left(z_i - \mathbf{g}_i^\top \mathbf{h} - \text{prox}_{\text{loss}}(z_i - \alpha \widehat{u}_i; \kappa) \right)^2 \right] = o_{\mathbb{P}}(1), \quad \frac{1}{n} \sum_{i \in [n]} \bar{\mathbb{E}} \left[\left(\psi_i - \text{env}'_{\text{loss}}(z_i - \alpha \widehat{u}_i; \kappa) \right)^2 \right] = o_{\mathbb{P}}(1)$$

Let us show $\text{tr}[\mathbf{V}]/p \xrightarrow{\mathbb{P}} \nu$. Using the concentration $\boldsymbol{\psi}^\top \mathbf{G} \mathbf{h} - \text{tr}[\mathbf{A}] \|\boldsymbol{\psi}\|^2 + \text{tr}[\mathbf{V}] \|\mathbf{h}\|^2 = o_{\mathbb{P}}(n)$ from Lemma 17 and $\|\mathbf{h}\|^2 \xrightarrow{\mathbb{P}} \alpha^2 > 0$, $\|\boldsymbol{\psi}\|^2/p \xrightarrow{\mathbb{P}} \beta^2$, $\text{tr}[\mathbf{A}] \xrightarrow{\mathbb{P}} \kappa$, we have

$$p^{-1} \text{tr}[\mathbf{V}] = \alpha^{-2} \kappa \beta^2 - \alpha^{-2} \boldsymbol{\psi}^\top \mathbf{G} \mathbf{h} / p + o_{\mathbb{P}}(1).$$

Recall that we have shown the concentration $\boldsymbol{\psi}^\top \mathbf{G} \mathbf{h} = \bar{\mathbb{E}}[\boldsymbol{\psi}^\top \mathbf{G} \mathbf{h}] + o_{\mathbb{P}}(n)$ in the proof of Lemma 19. Applying the proximal approximation of the residual $z_i - \mathbf{g}_i^\top \mathbf{h}$ and ψ_i to this, noting $\bar{\mathbb{E}}[\|\boldsymbol{\psi}\|^2] = \mathcal{O}_{\mathbb{P}}(n)$ and $\bar{\mathbb{E}}[\|\mathbf{G} \mathbf{h}\|^2] = \mathcal{O}_{\mathbb{P}}(n)$, we are left with

$$\begin{aligned} \frac{1}{n} \boldsymbol{\psi}^\top \mathbf{G} \mathbf{h} &= \frac{1}{n} \sum_{i \in [n]} \bar{\mathbb{E}}[\psi_i \cdot \mathbf{g}_i^\top \mathbf{h}] + o_{\mathbb{P}}(1) \\ &= \frac{1}{n} \sum_{i \in [n]} \bar{\mathbb{E}} \left[\text{env}'_{\text{loss}}(z_i - \alpha \widehat{u}_i; \kappa) (z_i - \text{prox}_{\text{loss}}(z_i - \alpha \widehat{u}_i; \kappa)) \right] + o_{\mathbb{P}}(1) \\ &= \frac{1}{n} \sum_{i \in [n]} \int_{-\infty}^{\infty} \varphi(x) \text{env}'_{\text{loss}}(z_i + \alpha x; \kappa) (z_i - \text{prox}_{\text{loss}}(z_i + \alpha x; \kappa)) dx + o_{\mathbb{P}}(1) \end{aligned}$$

so that the weak law of large numbers yields

$$\begin{aligned} \frac{1}{n} \boldsymbol{\psi}^\top \mathbf{G} \mathbf{h} &= \mathbb{E}[\text{env}'_{\text{loss}}(Z + \alpha G; \kappa) (Z - \text{prox}_{\text{loss}}(Z + \alpha G; \kappa))] + o_{\mathbb{P}}(1) \\ &= \mathbb{E}[\text{env}'_{\text{loss}}(Z + \alpha G; \kappa) (\alpha G + Z - \text{prox}_{\text{loss}}(Z + \alpha G; \kappa) - \alpha G)] + o_{\mathbb{P}}(1) \\ &= \kappa \mathbb{E}[\text{env}'_{\text{loss}}(\alpha G + Z; \kappa)^2] - \alpha \mathbb{E}[G \cdot \text{env}'_{\text{loss}}(\alpha G + Z; \kappa)] \\ &= \kappa \cdot \beta^2 / \delta - \alpha \cdot \nu \alpha / \delta. \end{aligned}$$

where we have used (5b) and (5d) in System 1a (with $c = 1$) for the last equation. Combined with $p^{-1} \text{tr}[\mathbf{V}] = \alpha^{-2} \kappa \beta^2 - \alpha^{-2} \boldsymbol{\psi}^\top \mathbf{G} \mathbf{h} / p + o_{\mathbb{P}}(1)$, this gives $p^{-1} \text{tr}[\mathbf{V}] = \nu + o_{\mathbb{P}}(1)$ and finishes the proof.

A.5 Convergence of error vector norm squared under Assumption D

In this section we verify that Assumption C and Assumption D-(1)-(3) are sufficient for [TAH18, Theorem 4.1] to hold. Comparing our assumptions and the condition assumed in the theorem, it suffices to show that the conditions (10) and (12) in [TAH18] can be omitted.

Indeed, the authors used the condition (10) to show that any $\hat{\mathbf{u}} \in \partial \text{loss}(\mathbf{z} - \mathbf{G}\hat{\mathbf{h}}^B)/\sqrt{p}$ belongs to a compact set with high probability, where $\hat{\mathbf{h}}^B$ is a ‘‘bounded’’ estimator. More precisely, $\hat{\mathbf{h}}^B$ is a solution to the constrained optimization problem $\min_{\|\mathbf{h}\| \leq K_\alpha} \text{obj}(\mathbf{h})$ for a positive constant $K_\alpha > 0$ where $\text{obj}(\mathbf{h}) = \sum_{i \in [n]} \text{loss}(z_i - \mathbf{g}_i^\top \mathbf{h}) + \sum_{j \in [p]} \text{reg}(\sqrt{p}h_j + \theta_j)$ is the objective function for the original unconstrained M-estimator. Since loss is differentiable with Lipschitz derivative loss' by Assumption C, using the same argument in Lemma 10, the norm of $\hat{\mathbf{u}}$ is bounded as $\|\hat{\mathbf{u}}\| \leq (\|\text{loss}'(\mathbf{z})\| + \|\text{loss}'\|_{\text{Lip}}\|\mathbf{G}\|_{\text{op}}K_\alpha)/\sqrt{p}$. Since \mathbf{G} has i.i.d. $\mathcal{N}(0, 1)$ entries while $\text{loss}'(z_i)$ has a finite second moment by Assumption D-(1), this gives $\|\hat{\mathbf{u}}\|^2 \leq C$ with probability approaching to 1 for a positive constant C .

Next, we claim that the condition (12) in [TAH18] is not necessary. Noting that the condition (12) is used for [TAH18, Assumption 2-(b)], it suffices to show that the assumption 2-(b) is satisfied given our assumptions. Here we restate the assumption 2-(b) for convenience:

$$\forall \tau > 0, \quad \lim_{c \rightarrow +\infty} c^2/(2\tau) - \mathbb{E}[\text{env}_{\text{reg}}(cH + \Theta) - \text{reg}(\Theta)] = +\infty.$$

By the condition $\mathbb{P}(\Theta \neq 0) > 0$ in Assumption D-(3), either $\mathbb{P}(\Theta > 0) > 0$ or $\mathbb{P}(\Theta < 0) > 0$ holds. Let us consider the case $\mathbb{P}(\Theta > 0) > 0$. Define a measurable function $u(H, \Theta)$ of (H, Θ) as follows:

$$u(H, \Theta) := -\min(\Theta, 1)I\{\Theta > 0 \text{ and } H < 0\}.$$

Note that $u(H, \Theta)$ is always bounded as $|u(H, \Theta)| \leq 1$. By the definition of Moreau envelope, i.e., $\text{env}_{\text{reg}}(cH + \Theta) := \text{argmin}_{p \in \mathbb{R}} (cH + \Theta - p)^2/(2\tau) + \text{reg}(p)$, taking the point $p = \Theta + u(H, \Theta)$, we get the lower estimate as follows:

$$\begin{aligned} \frac{c^2}{2\tau} - \mathbb{E}[\text{env}_{\text{reg}}(cH + \Theta) - \text{reg}(\Theta)] &\geq \frac{c^2}{2\tau} - \mathbb{E}\left[\frac{(cH - u(H, \Theta))^2}{2\tau} + \text{reg}(\Theta + u(H, \Theta)) - \text{reg}(\Theta)\right] \\ &= \frac{\mathbb{E}[uH]}{\tau} \cdot c - \frac{\mathbb{E}[u^2]}{2\tau} - \mathbb{E}[\text{reg}(\Theta + u) - \text{reg}(\Theta)]. \end{aligned}$$

Thus it suffices to show that $\mathbb{E}[uH] > 0$ and $\mathbb{E}[\text{reg}(\Theta + u) - \text{reg}(\Theta)]$ is finite for the RHS to diverge as $c \rightarrow +\infty$. By the definition of $u = u(\Theta, H)$, the expectation $\mathbb{E}[uH]$ can be written as

$$\mathbb{E}[uH] = \mathbb{E}[-H \min(\Theta, 1)I\{\Theta > 0, H < 0\}] = \mathbb{E}[|H| \min(\Theta, 1)I\{\Theta > 0, H < 0\}].$$

Here $|H| \min(\Theta, 1)$ is always strictly positive under the event $\{\Theta > 0, H > 0\}$, and this event has positive probability $\mathbb{P}(\Theta > 0, H > 0) = \mathbb{P}(\Theta > 0)\mathbb{P}(H > 0) = \mathbb{P}(\Theta > 0) \cdot 2^{-1}$ by the assumption $\mathbb{P}(\Theta > 0) > 0$ and the independence of Θ and $H \sim \mathcal{N}(0, 1)$. This means that $\mathbb{E}[uH]$ is strictly positive. Regarding $\mathbb{E}[\text{reg}(\Theta + u) - \text{reg}(\Theta)]$, we have

$$\mathbb{E}[\text{reg}(\Theta + u) - \text{reg}(\Theta)] = \mathbb{E}\left[\{\text{reg}(\Theta - \min(\Theta, 1)) - \text{reg}(\Theta)\}I\{\Theta > 0, H < 0\}\right].$$

Note that $0 < \Theta - \min(\Theta, 1) < \Theta$ for all $\Theta > 0$. Then, by the convexity of reg and the condition $\text{reg}(0) = \min_x \text{reg}(x)$ in Assumption C, under the event $I\{\Theta > 0, H < 0\}$ it holds that

$$0 > \text{reg}(\Theta - \min(\Theta, 1)) - \text{reg}(\Theta) > -d_\Theta \min(\Theta, 1) > -d_\Theta \quad \text{for all } d_\Theta \in \partial \text{reg}(\Theta)$$

By Assumption D-(1), d_Θ has a finite second moment for any choice of sub-derivative d_Θ . Therefore, we have $0 > \mathbb{E}[\text{reg}(\Theta + u) - \text{reg}(\Theta)] > -\mathbb{E}[d_\Theta I\{\Theta > 0, H < 0\}] > -\infty$ so that $\mathbb{E}[\text{reg}(\Theta + u) - \text{reg}(\Theta)]$ is finite.

In the other case $\mathbb{P}(\Theta < 0) > 0$, we may take $u(H, \Theta) := \min(-\Theta, 1)I\{\Theta < 0, H > 0\}$. Then the same argument leads to $\mathbb{E}[uH] > 0$ and $|\mathbb{E}[\text{reg}(\Theta + u) - \text{reg}(\Theta)]| < +\infty$.

A.6 Convergence of loss gradient norm squared under Assumption D

Let $\boldsymbol{\psi} = \text{loss}'(\mathbf{z} - \mathbf{G}\mathbf{h})$ and loss^* be the conjugate of loss . By the same argument in [TAH18], restricting the range of $\boldsymbol{\psi}$ to a compact set so that the strong duality holds, we observe that $\boldsymbol{\psi}$ is a solution to the following min-max problem with probability approaching to 1:

$$\max_{\boldsymbol{\psi} \in \mathbb{R}^n} \min_{\mathbf{h} \in \mathbb{R}^p} \boldsymbol{\psi}^\top (\mathbf{z} - \mathbf{G}\mathbf{h}) - \text{loss}^*(\boldsymbol{\psi}) + \text{reg}(\sqrt{p}\mathbf{h} + \boldsymbol{\theta}) = \max_{\boldsymbol{\psi}} \boldsymbol{\psi}^\top \mathbf{z} - \text{loss}^*(\boldsymbol{\psi}) + \boldsymbol{\psi}^\top \mathbf{G}\boldsymbol{\theta}/\sqrt{p} - \text{reg}^*(\mathbf{G}^\top \boldsymbol{\psi}/\sqrt{p})$$

If we write $\hat{\mathbf{u}} = \boldsymbol{\psi}/\sqrt{p} \in \mathbb{R}^n$ then $\hat{\mathbf{u}}$ is the M-estimator of the form:

$$\hat{\mathbf{u}} \in \underset{\mathbf{u} \in \mathbb{R}^n}{\text{argmin}} \mathbf{F}(-\mathbf{G}^\top \mathbf{u}) + \mathbf{L}(\sqrt{n}\mathbf{u}), \quad \text{where} \quad \begin{aligned} \mathbf{F} : \mathbb{R}^p &\mapsto \mathbb{R}, & \mathbf{v} &\mapsto \text{reg}^*(-\mathbf{v}) + \boldsymbol{\theta}^\top \mathbf{v} \\ \mathbf{L} : \mathbb{R}^n &\mapsto \mathbb{R}, & \mathbf{u} &\mapsto \text{loss}^*\left(\sqrt{\frac{p}{n}}\mathbf{u}\right) - \sqrt{\frac{p}{n}}\mathbf{z}^\top \mathbf{u} \end{aligned}$$

For any $\mathbf{H} \sim \mathcal{N}(\mathbf{0}_p, \mathbf{I}_p)$, $\mathbf{G} \sim \mathcal{N}(\mathbf{0}_n, \mathbf{I}_n)$, for all $c \in \mathbb{R}$ and $\tau > 0$ it holds that

$$\begin{aligned} \frac{1}{p} (\text{env}_{\mathbf{F}}(c\mathbf{H}; \tau) - \mathcal{F}(\mathbf{0})) &\xrightarrow{p} \frac{c^2}{2\tau} - \mathbb{E}[\text{env}_{\text{reg}}(\frac{c}{\tau}H + \Theta_0; \frac{1}{\tau}) - \text{reg}(0)] \\ \frac{1}{n} (\text{env}_{\mathbf{L}}(c\mathbf{G}; \tau) - \mathcal{L}(\mathbf{0})) &\xrightarrow{p} \frac{c^2}{2\tau} - \mathbb{E}[\text{env}_{\text{loss}}(\sqrt{\delta}\frac{c}{\tau}G + Z; \frac{\delta}{\tau}) - \text{loss}(0)] \end{aligned}$$

Thus, letting $F(c, \tau) := \mathbb{E}[\text{env}_{\text{reg}}(cH + \Theta; \tau)]$ and $L(c, \tau) := \mathbb{E}[\text{env}_{\text{loss}}(cG + \Theta; \tau)]$, [TAH18] implies that $\|\hat{\mathbf{u}}\|^2 \xrightarrow{p} \tilde{\alpha}^2$ where $\tilde{\alpha}$ is the minimizer of the min-max optimization problem

$$\inf_{\tilde{\alpha}, \tilde{\tau}_g} \sup_{\tilde{\beta}, \tilde{\tau}_h} \frac{\tilde{\beta}\tilde{\tau}_g}{2} + \delta^{-1} \left(\frac{\tilde{\alpha}^2 \tilde{\beta}}{2\tilde{\tau}_g} - F\left(\frac{\tilde{\alpha}\tilde{\beta}}{\tilde{\tau}_g}, \frac{\tilde{\beta}}{\tilde{\tau}_g}\right) \right) - \frac{\tilde{\alpha}\tilde{\tau}_h}{2} - \frac{\tilde{\alpha}\tilde{\beta}^2}{2\tilde{\tau}_h} + \left(\frac{\tilde{\alpha}\tilde{\beta}^2}{2\tilde{\tau}_h} - L(\sqrt{\delta}\tilde{\beta}, \delta\frac{\tilde{\tau}_h}{\tilde{\alpha}}) \right)$$

Thus, by the change of variables $(\tilde{\alpha}, \tilde{\beta}, \tilde{\tau}_g, \tilde{\tau}_h) \mapsto (\beta, \alpha/\sqrt{\delta}, \tau_h/\sqrt{\delta}, \tau_g/\delta)$ and multiplying the potential by $-\delta$, we are left with $\|\hat{\mathbf{u}}\|^2 \xrightarrow{p} \beta_*^2$ where β_* is the minimizer of

$$\sup_{\beta\tau_h} \inf_{\alpha, \tau_g} -\frac{\alpha\tau_h}{2} - \frac{\beta^2\alpha}{2\tau_g} + F\left(\frac{\beta\alpha}{\tau_h}, \frac{\alpha}{\tau_h}\right) + \frac{\beta\tau_g}{2} + \delta L\left(\alpha, \frac{\tau_g}{\beta}\right),$$

which is the same potential in [TAH18].

A.7 Explicit expressions for degrees of freedom

Table 5: Explicit formulae for the degrees of freedom and residual degrees of freedom of the estimator $\widehat{\boldsymbol{\theta}}_I$ defined using (3). Here the loss functions $\text{loss}(r)$ are: 1) square loss: $r^2/2$, 2) Huber loss: $r^2/2$ for $|r| \leq 1$ and $|r| - 1/2$ for $|r| > 1$; and the regularization functions $\text{reg}(b)$ are: 1) ridge penalty: $\frac{\lambda_1}{2}b^2$, 2) lasso penalty: $\lambda_2|b|$, and 3) elastic net penalty: $\frac{\lambda_1}{2}b^2 + \lambda_2|b|$. And other quantities are: 1) $\widehat{S}_I = \{j \in [p]: \widehat{\boldsymbol{\theta}}_I(j) \neq 0\}$ is the set of active variables of $\widehat{\boldsymbol{\theta}}_I$ and $\mathbf{X}_{\widehat{S}_I}$ is the submatrix of \mathbf{X}_I made of columns indexed in \widehat{S}_I , 2) $\widehat{T}_I = \{i \in I: \text{loss}''(y_i - \mathbf{x}_i^\top \widehat{\boldsymbol{\theta}}_I) > 0\}$ is the set of detected inliers (active observations), and 3) the matrix $\mathbf{D}_I = \text{diag}(\text{loss}''(\mathbf{y}_I - \mathbf{X}_I \widehat{\boldsymbol{\theta}}_I))$.

Loss (loss)	Regularizer (reg)	Degrees of freedom (df _I)	Residual degrees of freedom (tr[V _I])
Square	Ridge	$\text{tr}[(\mathbf{X}_I^\top \mathbf{X}_I + \lambda_1 \mathbf{I})^{-1} \mathbf{X}_I^\top \mathbf{X}_I]$	$ I - \text{df}_I$
Square	Lasso	$ \widehat{S}_I $	$ I - \text{df}_I$
Square	Elastic net	$\text{tr}[(\mathbf{X}_{\widehat{S}_I}^\top \mathbf{X}_{\widehat{S}_I} + \lambda_1 \mathbf{I})^{-1} \mathbf{X}_{\widehat{S}_I}^\top \mathbf{X}_{\widehat{S}_I}]$	$ I - \text{df}_I$
Huber	Ridge	$\text{tr}[(\mathbf{X}_I^\top \mathbf{D}_I \mathbf{X}_I + \lambda_1 \mathbf{I})^{-1} \mathbf{X}_I^\top \mathbf{D}_I \mathbf{X}_I]$	$ \widehat{T}_I - \text{df}_I$
Huber	Lasso	$ \widehat{S}_I $	$ \widehat{T}_I - \text{df}_I$
Huber	Elastic net	$\text{tr}[(\mathbf{X}_{\widehat{S}_I}^\top \mathbf{D}_I \mathbf{X}_{\widehat{S}_I} + \lambda_1 \mathbf{I})^{-1} \mathbf{X}_{\widehat{S}_I}^\top \mathbf{D}_I \mathbf{X}_{\widehat{S}_I}]$	$ \widehat{T}_I - \text{df}_I$
Convex	Convex	$\text{tr}[(\partial/\partial \mathbf{y}_I) \mathbf{X}_I \widehat{\boldsymbol{\theta}}_I]$	$\text{tr}[(\partial/\partial \mathbf{y}_I) \text{loss}'(\mathbf{y}_I - \mathbf{X}_I \widehat{\boldsymbol{\theta}}_I)]$

B Proofs in Section 4

B.1 Proof of Proposition 7

Let us fix $\psi = (c\delta)^{-1}$ and take the derivative of $\mathcal{R}_M = M^{-1}\alpha^2 + (1 - M^{-1})\alpha^2\eta_G$ with respect to $\phi = \delta^{-1}$. Note that $(\alpha, \beta, \kappa, \nu)$ are all fixed. Below, we derive the partial derivative of η_G with respect to ϕ . Using ψ and ϕ , we observe that η_G is the unique solution to the fixed point equation

$$\eta_G = F_{\text{reg}} \circ F_{\text{loss}}(\eta_G; \phi).$$

Here, F_{reg} does not depend on ϕ but F_{loss} depends on ϕ as:

$$F_{\text{loss}}(\eta_G; \phi) = \frac{\phi}{\psi^2 \beta^2} \cdot \mathbb{E}[\text{env}'_{\text{loss}}(\alpha G + Z; \kappa) \cdot \text{env}'_{\text{loss}}(\alpha \widetilde{G} + Z; \kappa)],$$

Since the map $\eta_G \mapsto F_{\text{reg}} \circ F_{\text{loss}}(\eta_G; \phi)$ is differentiable and c -Lipschitz with $c < 1$ (see Theorem 1), the implicit function theorem implies that η_G is differentiable with respect to ϕ and the derivative satisfies:

$$\frac{\partial \eta_G}{\partial \phi} = (F_{\text{reg}} \circ F_{\text{loss}})'(\eta_G) \frac{\partial \eta_G}{\partial \phi} + F'_{\text{reg}}(\eta_H) \frac{F_{\text{loss}}(\eta_G)}{\phi}.$$

Rearranging the above display, we get

$$\frac{\partial \eta_G}{\partial \phi} = \frac{F'_{\text{reg}}(\eta_H)}{\phi} \frac{F_{\text{loss}}(\eta_G)}{1 - (F_{\text{reg}} \circ F_{\text{loss}})'(\eta_G)} = \frac{F'_{\text{reg}}(\eta_H)}{\phi} \frac{\eta_H}{1 - (F_{\text{reg}} \circ F_{\text{loss}})'(\eta_G)} \geq 0.$$

where the last inequality follows from the fact that F_{reg} is non-decreasing and η_H is non-negative (see Theorem 1 for a homogeneous case). Combined with $\partial_\psi \mathcal{R}_M = (1 - M^{-1})\alpha^2 \frac{\partial \eta_G}{\partial \phi}$, we observe that \mathcal{R}_M is non-decreasing in ϕ . Therefore, for any two $\phi_1 \leq \phi_2$, we have that

$$\mathcal{R}_M(\phi_2, \psi) \geq \mathcal{R}_M(\phi_1, \psi) \geq \inf_{\psi \geq \phi_2} \mathcal{R}_M(\phi_1, \psi) \geq \inf_{\psi \geq \phi_1} \mathcal{R}_M(\phi_1, \psi) \quad \text{for all all } \psi \geq \phi_2.$$

This gives $\inf_{\psi \geq \phi_2} \mathcal{R}_M(\phi_2, \psi) \geq \inf_{\psi \geq \phi_1} \mathcal{R}_M(\phi_1, \psi)$ for any $\phi_1 \leq \phi_2$, which means the map $\phi \mapsto \inf_{\psi \geq \phi} \mathcal{R}_M(\phi, \psi)$ is non-decreasing.

B.2 Derivation of System 2

We will first reformulate Systems 1a and 1b in a slightly different set of parameters. Since the purpose of this reparameterization is to match with existing work, we will also consider regularizer reg with an explicit regularization level λ . The mapping with respect to the parameters in Systems 1a and 1b is as follows: $a = \frac{\lambda}{\beta}$, $\tau = \sqrt{c\delta} \frac{\beta}{\nu}$. Under this parameterization, note that $\frac{\beta}{\nu} = \frac{\tau}{\sqrt{c\delta}}$ and $\frac{\lambda}{\nu} = \frac{a\tau}{\sqrt{c\delta}}$.

Remark 9 (Scaling differences in design). It is worth remarking that in the literature on risk characterization of regularized M-estimator under proportional asymptotics, the scaling of λ can be slightly different, up to a factor of $\sqrt{c\delta}$. One of the reasons for the differences is how the design matrix \mathbf{X} is scaled. We assume that the entries of $\mathbf{X} \in \mathbb{R}^{n \times p}$ each have variance $1/p$ and thus each row has a unit average norm squared. It is also common to assume that the entries of \mathbf{X} each have variance $1/n$ and thus each column has a unit average norm squared. This brings in a factor of $\sqrt{\delta}$. For the subsampled design $\mathbf{X}_I \in \mathbb{R}^{k \times p}$, we get an additional factor of \sqrt{c} . Consequently, some expressions may appear different up to this scaling.

Derivation of System 2. The derivation is straightforward. We will use some simple relationships between proximal operators and Moreau envelopes. Recall from (2) that $\text{env}'_f(x; \tau) = \frac{1}{\tau}(x - \text{prox}_f(x; \tau))$ so that the derivative identity $\text{env}''_f(x; \tau) = \frac{1}{\tau}(1 - \text{prox}'_f(x; \tau))$ holds for almost every x by the non-expansiveness of the proximal operator.

(1) Case of $m = \ell$. The original system of equations is given by System 1a. For regularizers of the form λreg , from (5a), we have

$$\begin{aligned} \alpha^2 &= \mathbb{E} \left[\left(\frac{\lambda}{\nu} \cdot \text{env}'_{\text{reg}} \left(\frac{\beta}{\nu} H + \Theta; \frac{\lambda}{\nu} \right) - \frac{\beta}{\nu} H \right)^2 \right] \\ &= \mathbb{E} \left[\left(\text{prox}_{\text{reg}} \left(\frac{\beta}{\nu} H + \Theta; \frac{\lambda}{\nu} \right) - \Theta \right)^2 \right] \\ &= \mathbb{E} \left[\left(\text{prox}_{\text{reg}} \left(\frac{\tau}{\sqrt{c\delta}} H + \Theta; \frac{a\tau}{\sqrt{c\delta}} \right) - \Theta \right)^2 \right] \end{aligned} \quad (39)$$

where the variables τ and a are defined as:

$$\tau = \sqrt{c\delta} \frac{\beta}{\nu} \stackrel{(5b), (5d)}{=} \frac{\alpha \sqrt{\mathbb{E}[\text{env}'_{\text{loss}}(\alpha G + Z; \kappa)^2]}}{\mathbb{E}[\text{env}'_{\text{loss}}(\alpha G + Z; \kappa) \cdot G]}, \quad (40)$$

$$a = \frac{\lambda}{\beta} \stackrel{(5b)}{=} \frac{\lambda}{\sqrt{c\delta} \sqrt{\mathbb{E}[\text{env}'_{\text{loss}}(\alpha G + Z; \kappa)^2]}}. \quad (41)$$

For squared loss $\text{loss}(x) = x^2/2$, since $\text{env}'_{\text{loss}}(x; \tau) = x/(1 + \tau)$ from Table 6, squaring the final expression in (40) yields

$$\tau^2 = \alpha^2 + \sigma^2. \quad (42)$$

Thus, τ^2 and α^2 are the limiting total and excess risks, respectively. Combining (39) and (42) gives the first desired equation (18a).

Similarly, (5c) after dividing by τ yields

$$\frac{\kappa\beta}{\tau} = \frac{\beta}{\nu\tau} - \frac{\lambda}{\nu\tau} \mathbb{E} \left[\text{env}'_{\text{reg}} \left(\frac{\tau}{\sqrt{c\delta}} H + \Theta; \frac{a\tau}{\sqrt{c\delta}} \right) \cdot H \right].$$

Multiplying both sides by $\sqrt{c\delta}a\tau$ and noting that $a\beta = \lambda$, $\frac{\beta}{\nu} = \frac{\tau}{\sqrt{c\delta}}$, $\frac{\lambda}{\nu} = \frac{a\tau}{\sqrt{c\delta}}$ then gives

$$\sqrt{c\delta}\kappa\lambda = a\tau - \frac{(a\tau)^2}{\tau} \mathbb{E}[\text{env}'_{\text{reg}}(\frac{\tau}{\sqrt{c\delta}}H + \Theta; \frac{a\tau}{\sqrt{c\delta}}) \cdot H]. \quad (43)$$

From Stein's lemma, we have

$$\mathbb{E}[\text{env}'_{\text{reg}}(\tau H + \Theta; \kappa) \cdot H] = \tau \mathbb{E}[\text{env}''_{\text{reg}}(\tau H + \Theta; \kappa)], \quad (44)$$

and so (43) reduces to

$$\begin{aligned} \sqrt{c\delta}\kappa\lambda &\stackrel{(44)}{=} a\tau - \frac{(a\tau)^2}{\sqrt{c\delta}} \mathbb{E}[\text{env}''_{\text{reg}}(\frac{\tau}{\sqrt{c\delta}}H + \Theta; \frac{a\tau}{\sqrt{c\delta}})] \\ &= a\tau - a\tau \mathbb{E}[1 - \text{prox}'_{\text{reg}}(\frac{\tau}{\sqrt{c\delta}}H + \Theta; \frac{a\tau}{\sqrt{c\delta}})] \\ &= a\tau \mathbb{E}[\text{prox}'_{\text{reg}}(\frac{\tau}{\sqrt{c\delta}}H + \Theta; \frac{a\tau}{\sqrt{c\delta}})]. \end{aligned} \quad (45)$$

Using similar manipulations as above, (40) reduces to

$$\tau = \frac{\kappa \sqrt{\mathbb{E}[\text{env}'_{\text{loss}}(\alpha G + Z; \kappa)^2]}}{\mathbb{E}[1 - \text{prox}'_{\text{loss}}(\alpha G + Z; \kappa)]}. \quad (46)$$

Combining (45) and (46), we get

$$0 = \kappa\lambda \left(1 - \frac{a\tau}{\sqrt{c\delta}\kappa\lambda} \mathbb{E}[\text{prox}'_{\text{reg}}(\frac{\tau}{\sqrt{c\delta}}H + \Theta; \frac{a\tau}{\sqrt{c\delta}})]\right). \quad (47)$$

From (41), under squared loss, we also have

$$1 + \kappa = \frac{\sqrt{c\delta}a\tau}{\lambda}. \quad (48)$$

Combining (47) with (48) then leads to the second desired equation (18b).

(2) Case of $m \neq \ell$. Step 2. (η_G, η_H) is the solution to the following 2-scalar fix-point equations:

$$\eta_G = \frac{\mathbb{E}[(\frac{\lambda}{\nu} \cdot \text{env}'_{\text{reg}}(\frac{\beta}{\nu}H + \Theta; \frac{\lambda}{\nu}) - \frac{\beta}{\nu}H) \cdot (\frac{\lambda}{\nu} \cdot \text{env}'_{\text{reg}}(\frac{\beta}{\nu}\tilde{H} + \Theta; \frac{\lambda}{\nu}) - \frac{\beta}{\nu}\tilde{H})]}{\mathbb{E}[(\frac{\lambda}{\nu} \cdot \text{env}'_{\text{reg}}(\frac{\beta}{\nu}H + \Theta; \frac{\lambda}{\nu}) - \frac{\beta}{\nu}H)^2]}$$

$$\eta_H = c \cdot \frac{\mathbb{E}[\text{env}'_{\text{loss}}(\alpha G + Z; \kappa) \cdot \text{env}'_{\text{loss}}(\alpha \tilde{G} + Z; \kappa)]}{\mathbb{E}[\text{env}'_{\text{loss}}(\alpha G + Z; \kappa)^2]},$$

where $\begin{pmatrix} G \\ \tilde{G} \end{pmatrix} \sim N\left(\begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} 1 & \eta_G \\ \eta_G & 1 \end{pmatrix}\right)$ and $\begin{pmatrix} H \\ \tilde{H} \end{pmatrix} \sim N\left(\begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} 1 & \eta_H \\ \eta_H & 1 \end{pmatrix}\right)$. Similarly, by variable substitution, we can rewrite the equations as:

$$\eta_G \alpha^2 = \mathbb{E}[(\text{prox}_{\text{reg}}(\frac{\tau}{\sqrt{c\delta}}H + \Theta; \frac{a\tau}{\sqrt{c\delta}}) - \Theta)(\text{prox}_{\text{reg}}(\frac{\tau}{\sqrt{c\delta}}\tilde{H} + \Theta; \frac{a\tau}{\sqrt{c\delta}}) - \Theta)] \quad (49a)$$

$$\eta_H = c \cdot \frac{\mathbb{E}[\text{env}'_{\text{loss}}(\alpha G + Z; \kappa) \cdot \text{env}'_{\text{loss}}(\alpha \tilde{G} + Z; \kappa)]}{\mathbb{E}[\text{env}'_{\text{loss}}(\alpha G + Z; \kappa)^2]}. \quad (49b)$$

Since $\text{env}'_{\text{loss}}(x; \tau) = x/(1 + \tau)$ for squared loss, we obtain the third desired equation (19a). \square

B.3 Proximal operators and Moreau envelopes

Table 6: Proximal operators, Moreau envelopes, and their derivatives for ridge (first row) and lasso (second row) regularizers and Huber (third row) loss considered in Sections 4.2.1 and 4.2.2.

$f(x)$	$\text{prox}_f(x; \tau)$	$\text{prox}'_f(x; \tau)$	$\text{env}_f(x; \tau)$	$\text{env}'_f(x; \tau)$
$\frac{1}{2}x^2$	$\frac{x}{1+\tau}$	$\frac{1}{1+\tau}$	$\frac{1}{2} \frac{x^2}{1+\tau}$	$\frac{x}{1+\tau}$
$ x $	$(x - \tau)_+ \text{sign}(x)$	$\mathbf{1}\{ x \geq \tau\}$	$\begin{cases} \frac{1}{2\tau}x^2 & x < \tau \\ x - \frac{1}{2}\tau & x \geq \tau \end{cases}$	$\min\left\{\frac{ x }{\tau}, 1\right\} \text{sign}(x)$
$\begin{cases} \frac{x^2}{2} & x \leq 1 \\ x - \frac{1}{2} & x > 1 \end{cases}$	$\begin{cases} \frac{x}{1+\tau} & x \leq 1 + \tau \\ x - \tau \text{sign}(x) & x > 1 + \tau \end{cases}$	$\begin{cases} \frac{1}{1+\tau} & x \leq 1 + \tau \\ 1 & x > 1 + \tau \end{cases}$	$\begin{cases} \frac{1}{2} \frac{x^2}{1+\tau} & x \leq 1 + \tau \\ x - \tau - \frac{1}{2} & x > 1 + \tau \end{cases}$	$\left(\frac{ x }{1+\tau} - 1\right)_+ \text{sign}(x)$

C Proofs in Section 5

C.1 Proof of Equation (26)

Applying Stein's lemma to (23b) and Cauchy–Schwarz inequality, we have

$$\begin{aligned}
1 &= \frac{1}{c\delta} \mathbb{E} \left[\text{prox}'_{\text{reg}} \left(\Theta + \frac{\tau}{\sqrt{c\delta}} H; \frac{a\tau}{\sqrt{c\delta}} \right) \right] \\
&= \frac{1}{\sqrt{c\delta}} \frac{1}{\tau} \mathbb{E} \left[H \cdot \left(\text{prox}_{\text{reg}} \left(\Theta + \frac{\tau}{\sqrt{c\delta}} H; \frac{a\tau}{\sqrt{c\delta}} \right) - \Theta \right) \right] \\
&\leq \frac{1}{\sqrt{c\delta}\tau} \mathbb{E} \left[\left(\text{prox}_{\text{reg}} \left(\Theta + \frac{\tau}{\sqrt{c\delta}} H; \frac{a\tau}{\sqrt{c\delta}} \right) - \Theta \right)^2 \right]^{1/2} \\
&= \frac{1}{\sqrt{c\delta}\tau} \sqrt{\tau^2 - \sigma^2} \quad (\text{by (23a)}).
\end{aligned}$$

Taking the square of both sides and rearranging the resulting τ^2 term, we get the lower estimate $\tau^2 \geq \sigma^2/(1 - c\delta)$ as desired.

C.2 Derivation of System 4 as a limit of System 2 as $\lambda \rightarrow 0^+$

Fix $\text{reg}(x) = |x|^q$ for $q \in \{1, 2\}$. Let $(a_\lambda, \tau_\lambda, \xi_\lambda)$ be the solution to System 2 for any $\lambda > 0$ and let (a_*, τ_*, ξ_*) be the solution to System 4. Note that previous papers showed that $(a_\lambda, \tau_\lambda)$ converges to the solution (a_*, τ_*) as $\lambda \rightarrow (0)^+$; $q = 1$ case is given by the proof of Lemma A.1 in [LW21] while the $q = 2$ case immediately follows from the explicit formulae of $(a_\lambda, \tau_\lambda)$ and (a_*, τ_*) . Below we claim the continuity of ξ , i.e., $\xi_\lambda \rightarrow \xi_*$. Denoting $\eta_\lambda = c\xi_\lambda^2/\tau_\lambda^2$ and $\eta_* = c\xi_*^2/\tau_*^2$, what we want to show is $\eta_\lambda \rightarrow \eta_*$.

Observe that the systems for ξ_λ and ξ_* read to $\eta_\lambda = F(\eta_\lambda; \tau_\lambda, a_\lambda)$ and $\eta_* = F(\eta_*; \tau_*, a_*)$ where

$$F(\eta; \tau, a) := \frac{c}{\tau^2} \left\{ \mathbb{E} \left[\left(\text{prox}_{\text{reg}} \left(\Theta + \frac{\tau}{\sqrt{c\delta}} H; \frac{a\tau}{\sqrt{c\delta}} \right) - \Theta \right) \cdot \left(\text{prox}_{\text{reg}} \left(\Theta + \frac{\tau}{\sqrt{c\delta}} \tilde{H}; \frac{a\tau}{\sqrt{c\delta}} \right) - \Theta \right) \right] + \sigma^2 \right\}$$

with (H, \tilde{H}) being the mean zero jointly normals such that $\mathbb{E}[H^2] = \mathbb{E}[\tilde{H}^2] = 1$ and $\mathbb{E}[H\tilde{H}] = \eta$. Note that the map $\eta \mapsto F(\eta; \tau_\lambda, a_\lambda)$ is c -Lipschitz over $[-1, 1]$ by the same argument in Appendix A.1, while the map $(\tau, a) \mapsto F(\eta; \tau, a)$ is continuous over $(0, \infty)^2$ for any $\eta \in [-1, 1]$ by the moment assumption in Assumption D-1 and the dominated convergence theorem. Then, $|\eta_\lambda - \eta_*|$ is bounded from above as

$$|\eta_\lambda - \eta_*| = |F(\eta_\lambda; \tau_\lambda, a_\lambda) - F(\eta_*, \tau_*, a_*)|$$

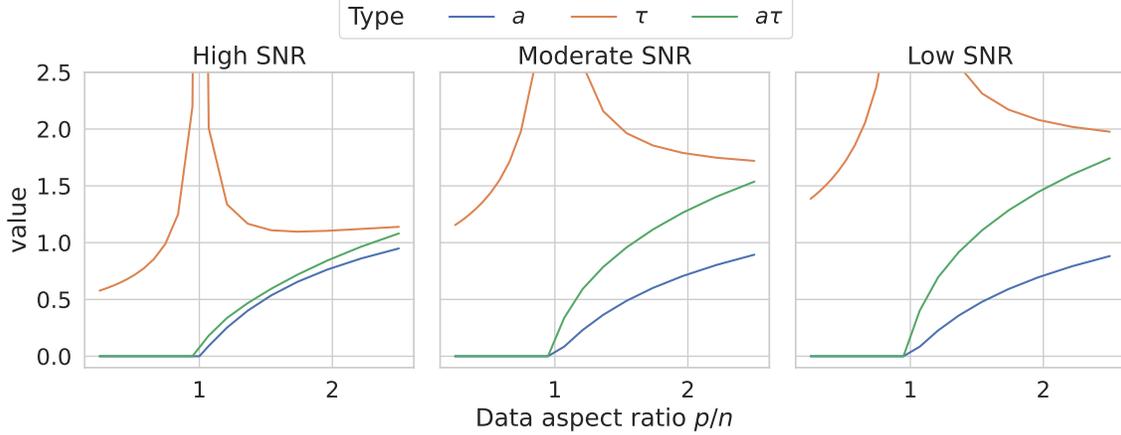


Figure 8: Fixed-point quantities for lassoless at different data aspect ratios p/n . The data model is given by (50) with signal strength $\rho = 0.5$ and sparsity levels $s = 0.2$ at different noise levels σ . *Left*: High SNR $\sigma = 0.5$. *Middle*: Modest SNR $\sigma = 1$. *Right*: Low SNR $\sigma = 1.2$.

$$\begin{aligned} &\leq |F(\eta_\lambda; \tau_\lambda, a_\lambda) - F(\eta_*, \tau_\lambda, a_\lambda)| + |F(\eta_*, \tau_\lambda, a_\lambda) - F(\eta_*, \tau_*, a_*)| \\ &\leq c|\eta_\lambda - \eta_*| + |F(\eta_*, \tau_\lambda, a_\lambda) - F(\eta_*, \tau_*, a_*)| \end{aligned}$$

so that $|\eta_\lambda - \eta_*| \leq (1 - c)^{-1}|F(\eta_*, \tau_\lambda, a_\lambda) - F(\eta_*, \tau_*, a_*)|$ holds. This upper bounds converges to 0 as $\lambda \rightarrow 0$ by the continuity of $F(\eta_*, \tau, a)$ in (τ, a) and the convergence $(\tau_\lambda, a_\lambda) \rightarrow (\tau_*, a_*)$.

C.3 Additional details for Remark 8

Expanding the product term in the fixed point equation (24a) by

$$\text{prox}\left(\Theta + \frac{\tau}{\sqrt{c\delta}}H; \frac{a\tau}{\sqrt{c\delta}}\right) - \Theta = \frac{\tau}{\sqrt{c\delta}}H - \frac{a\tau}{\sqrt{c\delta}}\text{env}'_{\text{reg}}\left(\Theta + \frac{\tau}{\sqrt{c\delta}}H; \frac{a\tau}{\sqrt{c\delta}}\right),$$

we have

$$\begin{aligned} \xi^2 - \sigma^2 &= \frac{\tau^2}{c\delta}\mathbb{E}[H\tilde{H}] - 2\frac{a\tau^2}{c\delta}\mathbb{E}\left[\tilde{H}\text{env}'_{\text{reg}}\left(\Theta + \frac{\tau}{\sqrt{c\delta}}H; \frac{a\tau}{\sqrt{c\delta}}\right)\right] \\ &\quad + \frac{(a\tau)^2}{c\delta}\mathbb{E}\left[\text{env}'_{\text{reg}}\left(\Theta + \frac{\tau}{\sqrt{c\delta}}H; \frac{a\tau}{\sqrt{c\delta}}\right) \cdot \text{env}'_{\text{reg}}\left(\Theta + \frac{\tau}{\sqrt{c\delta}}\tilde{H}; \frac{a\tau}{\sqrt{c\delta}}\right)\right], \end{aligned}$$

where $\frac{\tau^2}{c\delta}\mathbb{E}[H\tilde{H}] = \frac{\tau^2}{c\delta}\eta_H = \frac{\xi^2}{\delta}$ by the definition $\eta_H = c\xi^2/\tau^2$. For the second term, realizing $\tilde{H} = \eta_H H + \sqrt{1 - \eta_H^2}\bar{H}$ for a standard normal $\bar{H} \sim \mathcal{N}(0, 1)$ independent of (H, \tilde{H}, Θ) , noting that $\text{env}'_{\text{reg}}(\Theta + \frac{\tau}{\sqrt{c\delta}}H)$ is bounded in second moment by Assumption D-(1), we have

$$\begin{aligned} \mathbb{E}\left[\tilde{H} \cdot \text{env}'_{\text{reg}}\left(\Theta + \frac{\tau}{\sqrt{c\delta}}H; \frac{a\tau}{\sqrt{c\delta}}\right)\right] &= \eta_H \mathbb{E}\left[H \cdot \text{env}'_{\text{reg}}\left(\Theta + \frac{\tau}{\sqrt{c\delta}}H; \frac{a\tau}{\sqrt{c\delta}}\right)\right] \\ &= \eta_H \frac{\tau}{\sqrt{c\delta}} \mathbb{E}\left[\text{env}''_{\text{reg}}\left(\Theta + \frac{\tau}{\sqrt{c\delta}}H; \frac{a\tau}{\sqrt{c\delta}}\right)\right] && \text{(by Stein's lemma)} \\ &= \eta_H \frac{\tau}{\sqrt{c\delta}} \frac{\sqrt{c\delta}}{a\tau} \mathbb{E}\left[1 - \text{prox}'_{\text{reg}}\left(\Theta + \frac{\tau}{\sqrt{c\delta}}H; \frac{a\tau}{\sqrt{c\delta}}\right)\right] \\ &= \frac{\eta_H}{a}(1 - c\delta) && \text{(by (23b)).} \end{aligned}$$

Combined with $\eta_H = c\xi^2/\tau^2$, we have

$$\frac{a\tau^2}{c\delta}\mathbb{E}\left[\tilde{H}\text{env}'_{\text{reg}}\left(\Theta + \frac{\tau}{\sqrt{c\delta}}H\right)\right] = \frac{\xi^2}{\delta}(1 - c\delta).$$

Therefore, we get

$$\xi^2 \left(1 - \frac{1}{\delta} + 2 \frac{(1 - c\delta)}{\delta} \right) = \sigma^2 + \frac{(a\tau)^2}{c\delta} \mathbb{E} \left[\text{env}'_{\text{reg}} \left(\Theta + \frac{\tau}{\sqrt{c\delta}} H; \frac{a\tau}{\sqrt{c\delta}} \right) \cdot \text{env}'_{\text{reg}} \left(\Theta + \frac{\tau}{\sqrt{c\delta}} \tilde{H}; \frac{a\tau}{\sqrt{c\delta}} \right) \right].$$

This means that $\xi^2 \rightarrow \frac{\delta}{\delta-1} \sigma^2$ holds if and only if the rightmost term converges to 0 as $c \rightarrow (\delta^{-1})^-$. This is true if reg is Lipschitz and $\lim_{c \rightarrow (\delta^{-1})^-} a\tau = 0$, as $|\text{env}'_{\text{reg}}(x; \tau)| \leq \|\text{reg}\|_{\text{Lip}}$ for all $x \in \mathbb{R}$ and $\tau > 0$. However, we are not able to provably establish that $a\tau \rightarrow 0$ as $c \rightarrow \delta^{-1}$. For the lasso regularizer, we observe in Figure 8 that $a\tau \rightarrow 0$ appears to hold as $c \rightarrow \delta^{-1}$.

D Additional numerical simulations

D.1 Minimum ℓ_1 -norm interpolator

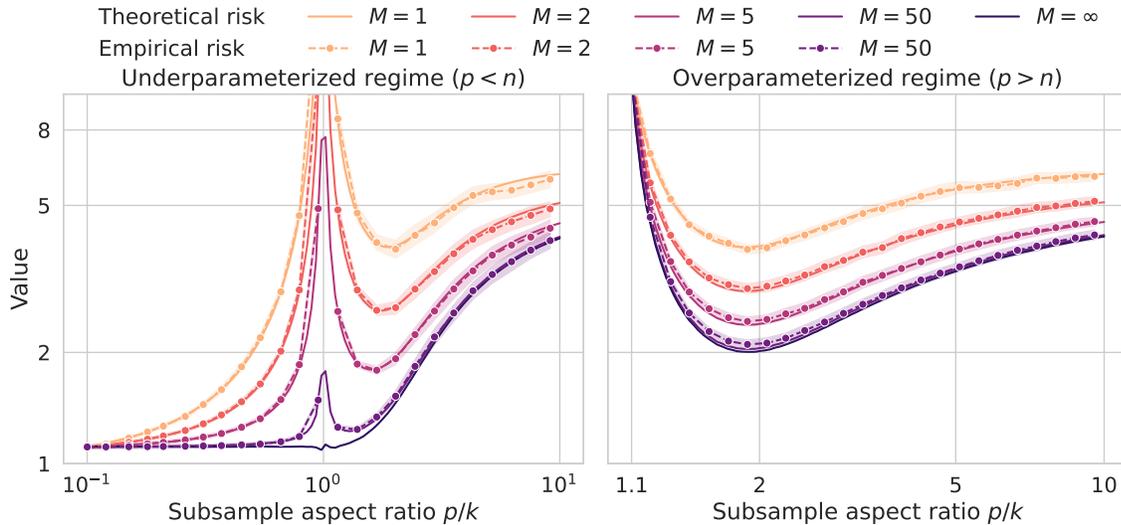


Figure 9: The prediction risk for lasso ensemble at different subsample aspect ratios p/k with regularization parameter $\lambda = 0.001$ and varying ensemble size M . The solid lines represent the theoretical risks, the dashed lines represent the empirical risks averaged over 50 simulations, and the shaded regions represent the standard errors. The data model is given by (50) with signal strength $\rho = 2$, noise level $\sigma = 1$, and support proportion $s = 0.1$. *Left*: underparameterized regime when $p/n = 0.1$ and $n = 2000$. *Right*: overparameterized regime when $p/n = 1.1$ and $n = 500$.

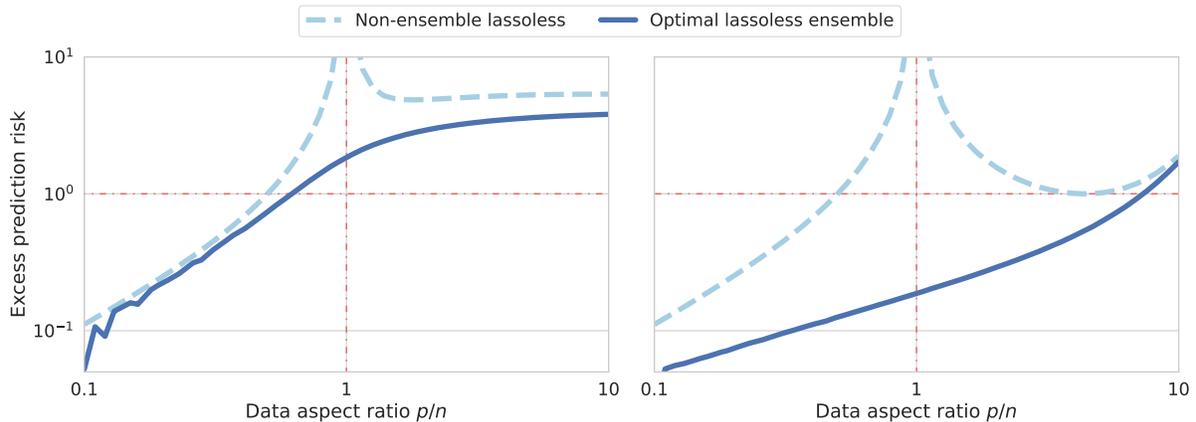


Figure 10: **Optimal subsample risk of the lassoless ensemble is monotonic in the data aspect ratio.** Excess risk of the lasso and optimal lasso ensemble, at different data aspect ratios p/n ranging from 0.1 to 10. The data model is given by (50) with signal strength $\rho = 2$, noise level $\sigma = 1$, data aspect ratio $p/n = 0.1$, feature size $p = 500$, and varying support proportion s . *Left*: dense regime with $s = 0.9$. *Right*: sparse regime with $s = 0.01$.

D.2 Optimal subbagging versus optimal regularization

D.2.1 Square loss, lasso regularizer

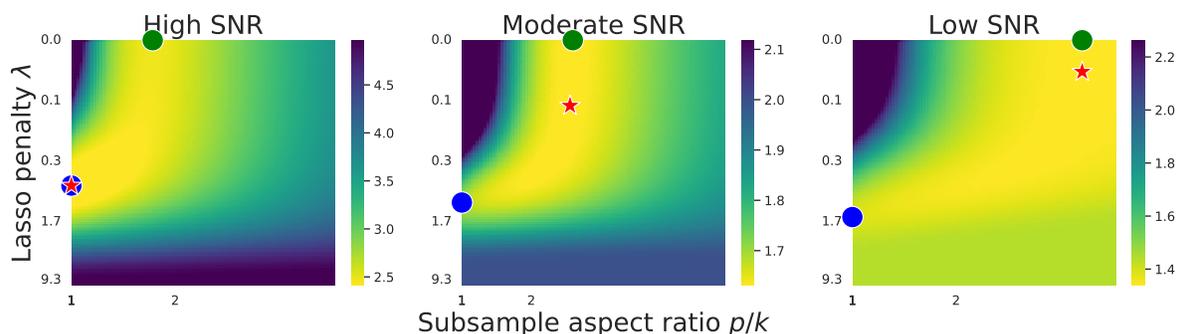


Figure 11: Heatmaps of theoretical prediction risk in λ and p/k of full lasso ensemble in the overparameterized regime ($p/n = 1.1$) and sparse regime ($s = 0.2$). The data model is given by (50) with support proportion $s = 0.2$ and noise level $\sigma = 1$ at different signal levels ρ . *Left*: high SNR $\rho = 2$ (optimal lasso is better than optimal subsample lassoless). *Middle*: moderate SNR $\rho = 1$ (optimal subsample lasso is better than optimal lasso and optimal subsample lassoless). *Right*: low SNR $\rho = 0.67$ (optimal subsample lassoless is better than optimal lasso).

D.2.2 Huber loss, unregularized

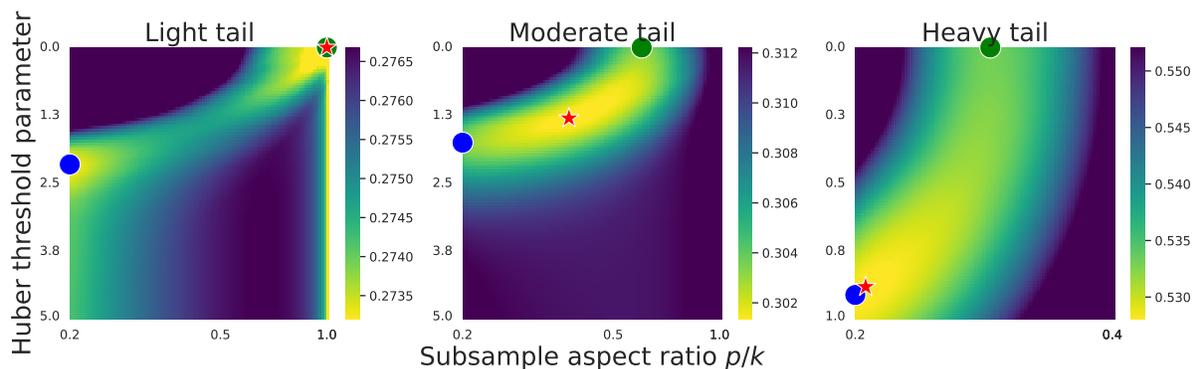


Figure 12: Heatmaps of theoretical prediction risk in Huber threshold parameter and subsample ratio k/n of full Huber ensemble in the underparameterized regime ($p/n = 0.2$). *Left*: noise follows Student's t distribution t_{20} (optimal subsample Huberless is better than optimal Huber). *Middle*: noise follows Student's t distribution t_{10} (optimal subsample Huber is better than both optimal subsample Huberless and optimal Huber). *Right*: noise follows Student's t distribution t_2 (optimal Huber is better than optimal subsample Huberless).

D.3 Huber regression with ℓ_1 -regularization

D.3.1 Varying threshold parameter

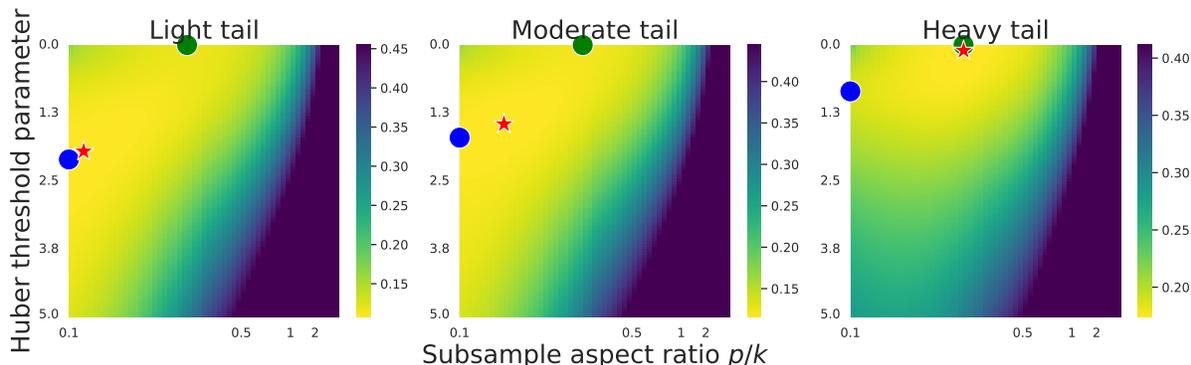


Figure 13: Heatmaps of theoretical prediction risk in Huber threshold parameter and subsample ratio k/n of full ℓ_1 -regularized Huber ensemble with lasso regularization level 0.5 in the underparameterized regime ($p/n = 0.1$). *Left*: noise follows Student's t distribution t_{20} . *Middle*: noise follows Student's t distribution t_{10} . *Right*: noise follows Student's t distribution t_2 .

D.3.2 Varying regularization level

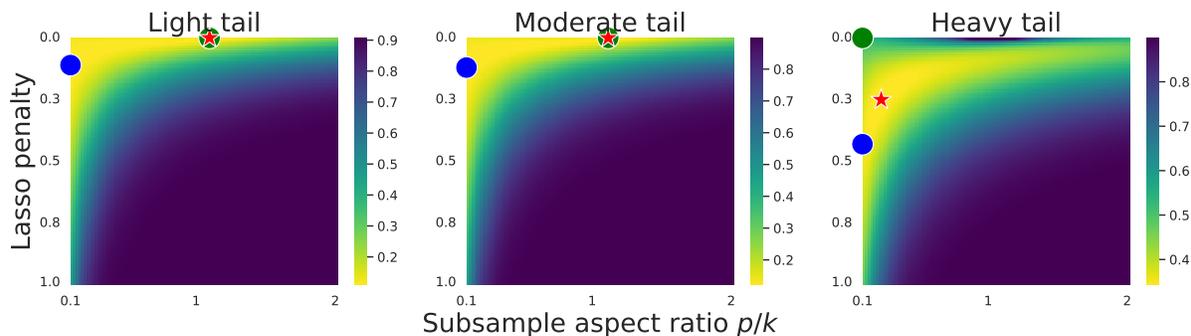


Figure 14: Heatmaps of theoretical prediction risk in lasso regularization level and subsample ratio k/n of full ℓ_1 -regularized Huber ensemble with Huber threshold parameter 10 in the underparameterized regime ($p/n = 0.1$). *Left*: noise follows Student's t distribution t_{20} . *Middle*: noise follows Student's t distribution t_{10} . *Right*: noise follows Student's t distribution t_2 .

D.4 Details of numerical experiments

D.4.1 Data model for lasso experiments

We consider a linear model whose signal variables are generated from two-point distribution:

$$y = \mathbf{x}^\top \boldsymbol{\theta} + z, \quad \mathbf{x} \sim \mathcal{N}(\mathbf{0}, p^{-1} \mathbf{I}_p), \quad z \sim \mathcal{N}(0, \sigma^2), \quad \boldsymbol{\theta} \stackrel{i.i.d.}{\sim} s \mathcal{P}_{\rho/\sqrt{s}} + (1-s) \mathcal{P}_0, \quad (50)$$

where \mathcal{P}_c denotes the Dirac measure at point $c \in \mathbb{R}$, and $\rho > 0$ is some given quantity that determines the signal magnitude. Here, $s \in (0, 1)$ is the support proportion. The signal-to-noise-ratio (SNR) under the above model obeys

$$\text{SNR} = \frac{\mathbb{E}[(\mathbf{x}^\top \boldsymbol{\theta})^2]}{\sigma^2} = \frac{s \cdot (\rho^2/s)}{\sigma^2} = \frac{\rho^2}{\sigma^2}.$$

D.4.2 Data model for ℓ_1 -regularized Huber experiments

We consider the ℓ_1 -regularized Huber regression

$$\widehat{\boldsymbol{\theta}}_I \in \operatorname{argmin}_{\boldsymbol{\theta} \in \mathbb{R}^p} \sum_{i \in I} \text{Huber}(y_i - \mathbf{x}_i^\top \boldsymbol{\theta}; \rho) + \lambda \|\boldsymbol{\theta}\|_1$$

where $\lambda \geq 0$ is a regularization parameter and $\text{Huber}(\cdot; \rho)$ is the Huber loss with threshold parameter $\rho > 0$, which is defined as follows [Hub64]:

$$\text{For all } x \in \mathbb{R}, \quad \text{Huber}(x; \rho) := \operatorname{env}_{|\cdot|}(x; \rho) = \begin{cases} \frac{1}{2\rho} x^2 & \text{if } |x| \leq \rho \\ |x| - \frac{1}{2}\rho & \text{if } |x| > \rho. \end{cases} \quad (51)$$

Observe that $\text{Huber}(\cdot; \rho)$ behaves like the squared loss for large ρ and like the absolute loss for small ρ . More precisely, for all $x \in \mathbb{R}$ it holds that

$$\lim_{\rho \rightarrow 0^+} \text{Huber}(x; \rho) = |x| \quad \text{and} \quad \lim_{\rho \rightarrow +\infty} \rho \cdot \text{Huber}(x; \rho) = x^2/2.$$

We consider the linear model $y = \mathbf{x}^\top \boldsymbol{\theta} + z$ where the marginal distribution of the signal is set to the mixture $\epsilon \mathcal{N}(0, 1) + (1 - \epsilon) \mathcal{P}_0$ (recall that \mathcal{P}_0 denotes the Dirac measure at 0) with support proportion $\epsilon = 0.1$, while the noise z follows Student's t-distribution t_d for some degree of freedom $d \geq 2$, which will be specified for each numerical simulation.