

The Distribution of Ridgeless Least Squares Interpolators

Qiyang Han

*Department of Statistics,
Rutgers University,
Piscataway, NJ 08854, USA*

QH85@STAT.RUTGERS.EDU

Xiaocong Xu

*Data Sciences and Operations Department
University of Southern California
Los Angeles, CA 90089, USA*

XUXIAOCO@MARSHALL.USC.EDU

Editor: Mahdi Soltanolkotabi

Abstract

The Ridgeless minimum ℓ_2 -norm interpolator in overparametrized linear regression has attracted considerable attention in recent years in both machine learning and statistics communities. While it seems to defy conventional wisdom that overfitting leads to poor prediction, recent theoretical research on its ℓ_2 -type risks reveals that its norm minimizing property induces an ‘implicit regularization’ that helps prediction in spite of interpolation.

This paper takes a further step that aims at understanding its precise stochastic behavior as a statistical estimator. Specifically, we characterize the distribution of the Ridgeless interpolator in high dimensions, in terms of a Ridge estimator in an associated Gaussian sequence model with positive regularization, which provides a precise quantification of the prescribed implicit regularization in the most general distributional sense. Our distributional characterizations hold for general non-Gaussian random designs and extend uniformly to positively regularized Ridge estimators.

As a direct application, we obtain a complete characterization for a general class of weighted ℓ_q risks of the Ridge(less) estimators that are previously only known for $q = 2$ by random matrix methods. These weighted ℓ_q risks not only include the standard prediction and estimation errors, but also include the non-standard covariate shift settings. Our uniform characterizations further reveal a surprising feature of the commonly used generalized and k -fold cross-validation schemes: tuning the estimated ℓ_2 prediction risk by these methods alone lead to simultaneous optimal ℓ_2 in-sample, prediction and estimation risks, as well as the optimal length of debiased confidence intervals.

Keywords: comparison inequality, cross validation, minimum norm interpolator, random matrix theory, ridge regression, universality

1. Introduction

1.1 Overview

Consider the standard linear regression model

$$Y_i = X_i^\top \mu_0 + \xi_i, \quad 1 \leq i \leq m, \quad (1)$$

where we observe i.i.d. feature vectors $X_i \in \mathbb{R}^n$ and responses $Y_i \in \mathbb{R}$, and ξ_i 's are unobservable errors. For notational simplicity, we write $X = [X_1 \cdots X_m]^\top \in \mathbb{R}^{m \times n}$ as the design

matrix that collects all the feature vectors, and $Y = (Y_1, \dots, Y_m)^\top \in \mathbb{R}^m$ as the response vector. The feature vectors X_i 's are assumed to satisfy $\mathbb{E} X_1 = 0$ and $\text{Cov}(X_1) = \Sigma$, and the errors satisfy $\mathbb{E} \xi_1 = 0$ and $\text{Var}(\xi_1) = \sigma_\xi^2$.

Throughout this paper, we reserve m for the sample size, and n for the signal dimension. The aspect ratio ϕ , i.e., the number of samples per dimension, is then defined as $\phi \equiv m/n$. Accordingly, we refer to $\phi^{-1} > 1$ as the *overparametrized regime*, and $\phi^{-1} < 1$ as the *underparametrized regime*.

Within the linear model (1), the main object of interest is to recover/estimate the unknown signal $\mu_0 \in \mathbb{R}^n$. While a large class of regression techniques can be used for the purpose of signal recovery under various structural assumptions on μ_0 , here we will focus our attention on one widely used class of regression estimators, namely, the *Ridge estimator* (cf. Hoerl and Kennard (1970)) with regularization $\eta > 0$,

$$\hat{\mu}_\eta = \arg \min_{\mu \in \mathbb{R}^n} \left\{ \frac{1}{2n} \|Y - X\mu\|^2 + \frac{\eta}{2} \|\mu\|^2 \right\} = \frac{1}{n} \left(\frac{1}{n} X^\top X + \eta I_n \right)^{-1} X^\top Y, \quad (2)$$

and the *Ridgeless estimator* (also known as the *minimum-norm interpolator*),

$$\hat{\mu}_0 = \arg \min_{\mu \in \mathbb{R}^n} \{ \|\mu\|^2 : Y = X\mu \} = (X^\top X)^- X^\top Y, \quad (3)$$

which is almost surely (a.s.) well-defined in the overparametrized regime $\phi^{-1} > 1$. Here A^- is the Moore-Penrose pseudo-inverse of A . The notation $\hat{\mu}_0$ is justified since for $\phi^{-1} > 1$, $\hat{\mu}_\eta \rightarrow \hat{\mu}_0$ a.s. as $\eta \downarrow 0$.

From a conventional statistical point of view, the Ridgeless estimator seems far from an obviously good choice: As $\hat{\mu}_0$ perfectly interpolates the data, it is susceptible to high variability due to the widely recognized bias-variance tradeoff inherent in ‘optimal’ statistical estimators (James et al., 2021; Derumigny and Schmidt-Hieber, 2023). On the other hand, as the Ridgeless estimator $\hat{\mu}_0$ is the limit point of the gradient descent algorithm run on the squared loss in the overparametrized regime $\phi^{-1} > 1$, it provides a simple yet informative test case for understanding one major enigma of modern machine learning methods: these methods typically interpolate training data perfectly; still, they enjoy good generalization properties (Jacot et al., 2018; Du et al., 2019; Allen-Zhu et al., 2019; Belkin et al., 2019; Chizat et al., 2019; Zhang et al., 2021).

Inspired by this connection, recent years have witnessed a surge of interest in understanding the behavior of the Ridgeless estimator $\hat{\mu}_0$ and its closely related Ridge estimator $\hat{\mu}_\eta$, with an exclusive focus on their prediction risks, cf. Tulino et al. (2004); El Karoui (2013); Hsu et al. (2014); Dicker (2016); Dobriban and Wager (2018); El Karoui (2018); Advani et al. (2020); Belkin et al. (2020); Muthukumar et al. (2020); Wu and Xu (2020); Bartlett et al. (2020, 2021); Chang et al. (2021); Koehler et al. (2021); Richards et al. (2021); Hastie et al. (2022); Tsigler and Bartlett (2023); Cheng and Montanari (2024); Zhou et al. (2024). The readers are referred to Tsigler and Bartlett (2023, Sections 1.2 & 9) for a thorough review on the relation between various ℓ_2 risk results for the Ridge(less) estimator. A unique insight from these works is the existence of ‘implicit regularization’ within the Ridgeless interpolator $\hat{\mu}_0$, so that for certain scenarios of (Σ, μ_0) , the prediction risk of $\hat{\mu}_0$ could be small (i.e., benign overfitting) or even optimal (Kobak et al., 2020; Hastie et al., 2022; Tsigler and Bartlett, 2023).

Despite substantial progress in understanding the ℓ_2 risk behavior of the Ridgeless estimator $\widehat{\mu}_0$, our understanding of its stochastic behavior as a statistical estimator remains limited. This gap is particularly important if we aim to consider $\widehat{\mu}_0$ also as a ‘good’ estimator that can be applied in a broader context of statistical inference tasks, rather than merely viewing it as a theoretical proxy for modern interpolating learning algorithms.

The main goal of this paper is to advance our understanding of the precise stochastic behavior of the Ridge(less) estimator $\widehat{\mu}_\eta$. We achieve this by developing a high-dimensional *distributional* characterization in the so-called proportional regime where m and n is of the same order. This approach allows us to move beyond the exclusive focus in the existing literature on ℓ_2 -type risks of $\widehat{\mu}_\eta$. As will be clear, the distributional characterization of the Ridge(less) estimator $\widehat{\mu}_\eta$ not only provides a precise quantitative understanding of the ‘implicit regularization’ phenomenon for the Ridgeless interpolator $\widehat{\mu}_0$ in the most general distributional sense, but also unveils major new insights on the utility of the widely used cross-validation schemes in machine learning/statistics practice.

1.2 Distribution of Ridge(less) estimators

Before formally describing our high dimensional distributional characterization, it is insightful to consider the low dimensional regime $\phi^{-1} \ll 1$ where the sample size m far exceeds the signal dimension n . In this regime, with $(\bar{\eta}, \bar{\sigma}_\xi^2) \equiv (\eta, \sigma_\xi^2)/\phi$, using the closed form of (2) and the fact that $m^{-1}X^\top X \approx \Sigma$, we may safely regard $\widehat{\mu}_\eta \approx (\Sigma + \bar{\eta}I_n)^{-1}(\Sigma\mu_0 + m^{-1}X^\top \xi)$. Using central limit theorem for $m^{-1}X^\top \xi \stackrel{d}{\approx} \bar{\sigma}_\xi \cdot n^{-1/2}\Sigma^{1/2}g$ where $g \sim \mathcal{N}(0, I_n)$, we have

$$\widehat{\mu}_\eta \stackrel{d}{\approx} (\Sigma + \bar{\eta}I_n)^{-1}\Sigma^{1/2}(\Sigma^{1/2}\mu_0 + n^{-1/2} \cdot \bar{\sigma}_\xi g), \quad \phi^{-1} \ll 1. \quad (4)$$

A principled way to understand the above formula (4) is to consider an ‘effective regression problem’ in the *Gaussian sequence model*. Suppose for a given pair of (Σ, μ_0) and a noise level $\gamma > 0$, we observe

$$y_{(\Sigma, \mu_0)}^{\text{seq}}(\gamma) \equiv \Sigma^{1/2}\mu_0 + n^{-1/2} \cdot \gamma g, \quad g \sim \mathcal{N}(0, I_n). \quad (5)$$

The Ridge estimator $\widehat{\mu}_{(\Sigma, \mu_0)}^{\text{seq}}(\gamma; \tau)$ with regularization $\tau \geq 0$ in the Gaussian sequence model (5) is defined as

$$\begin{aligned} \widehat{\mu}_{(\Sigma, \mu_0)}^{\text{seq}}(\gamma; \tau) &\equiv \arg \min_{\mu \in \mathbb{R}^n} \left\{ \frac{1}{2} \|\Sigma^{1/2}\mu - y_{(\Sigma, \mu_0)}^{\text{seq}}(\gamma)\|^2 + \frac{\tau}{2} \|\mu\|^2 \right\} \\ &= (\Sigma + \tau I_n)^{-1}\Sigma^{1/2}(\Sigma^{1/2}\mu_0 + n^{-1/2} \cdot \gamma g). \end{aligned} \quad (6)$$

Here, the subscript (Σ, μ_0) emphasizes the dependence on the underlying Gaussian sequence model with covariance Σ and signal μ_0 . Comparing (4) and (6), it is clear that we may interpret (4) as $\widehat{\mu}_\eta \stackrel{d}{\approx} \widehat{\mu}_{(\Sigma, \mu_0)}^{\text{seq}}(\bar{\sigma}_\xi; \bar{\eta})$. In the proportional regime $\phi^{-1} \asymp 1$, the aforementioned interpretation still applies, but a crucial modification will be needed: the pair of the (scaled) original noise and regularization $(\bar{\sigma}_\xi, \bar{\eta})$ must be replaced by a pair of ‘*effective noise and regularization*’

$$(\gamma_{\eta,*}, \tau_{\eta,*}) \equiv \text{unique solution of the fixed point equation (13)} \quad (7)$$

when $\eta > 0$ and when $\eta = 0$ in the overparametrized regime (cf. Proposition 2).

More precisely, in the overparametrized regime $\phi^{-1} > 1$, under standard assumptions on (i) the design matrix $X = \Sigma^{1/2}Z$, where Z consists of independent mean 0, unit-variance and light-tailed entries, and (ii) the error vector ξ with light-tailed components, we show in Theorems 3 and 4 that the distribution $\widehat{\mu}_\eta$ can be characterized via $\widehat{\mu}_{(\Sigma, \mu_0)}^{\text{seq}}(\gamma_{\eta,*}; \tau_{\eta,*})$ in the following sense: for any 1-Lipschitz function $\mathbf{g} : \mathbb{R}^n \rightarrow \mathbb{R}$ and any $K > 0$, with high probability,

$$\sup_{\eta \in [0, K]} \left| \mathbf{g}(\widehat{\mu}_\eta) - \mathbb{E} \mathbf{g}(\widehat{\mu}_{(\Sigma, \mu_0)}^{\text{seq}}(\gamma_{\eta,*}; \tau_{\eta,*})) \right| \approx 0. \quad (8)$$

A particularly important technical aspect of (8) is that the distributional approximation (8) holds uniformly down to the interpolation regime $\eta = 0$ for $\phi^{-1} > 1$. This uniform guarantee will prove essential in the results ahead.

Interestingly, the distributional characterization (8) offers a principled approach to understand the ‘implicit regularization’ phenomenon for the Ridgeless interpolator $\widehat{\mu}_0$, through the lens of its distributionally equivalent Ridge estimator $\widehat{\mu}_{(\Sigma, \mu_0)}^{\text{seq}}(\gamma_{0,*}, \tau_{0,*})$ in the Gaussian sequence model (5). Specifically, the prescribed implicit regularization can be directly attributed to the quantity $\tau_{0,*} > 0$ that can be solved as the unique positive solution to the equation

$$\phi = \frac{1}{n} \text{tr} \left((\Sigma + \tau_{0,*} I_n)^{-1} \Sigma \right). \quad (9)$$

While this interpretation has been suggested in the context of ℓ_2 risks (Hastie et al., 2022; Cheng and Montanari, 2024) via a-posterior calculations, our theory (8) provides a formal justification for this equivalent understanding of the implicit regularization phenomenon for $\widehat{\mu}_0$ via $\widehat{\mu}_{(\Sigma, \mu_0)}^{\text{seq}}$, in the most precise and general distributional sense. The readers are referred to Section 1.5 for a more detailed comparison on the relation of our characterization of the implicit regularization via $\tau_{0,*}$ and a different line of interpretation in Bartlett et al. (2020, 2021); Tsigler and Bartlett (2023).

1.3 General ℓ_q -type risk formulae

As mentioned above, most prior works on the risk properties of Ridge(less) estimators $\widehat{\mu}_\eta$ have focused exclusively on ℓ_2 -type risks, leveraging random matrix theory (RMT) (Tulino et al., 2004; El Karoui, 2013; Dicker, 2016; Dobriban and Wager, 2018; El Karoui, 2018; Advani et al., 2020; Wu and Xu, 2020; Bartlett et al., 2021; Richards et al., 2021; Hastie et al., 2022; Cheng and Montanari, 2024). This RMT approach is viable due to a direct reduction of ℓ_2 -type risks of $\widehat{\mu}_\eta$ to the spectrum of X . In contrast, the more general ℓ_q risks depend not only on the spectrum but also on the structure of X ’s singular vectors in a highly nontrivial manner; therefore, the feasibility of a similar RMT-based analysis is in question.

Our uniform distributional theory in (8) is strong enough to characterize all ℓ_q risks of the Ridge(less) estimator. Specifically, for any $q \in [1, \infty)$ and a p.s.d. matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$, with high probability,

$$\frac{\|\mathbf{A}\widehat{\mu}_\eta - \mu_0\|_q}{n^{-1/2} \|\text{diag}(\Gamma_{\eta; (\Sigma, \|\mu_0\|)}^{\mathbf{A}})\|_{q/2}^{1/2} M_q} \approx 1, \quad (10)$$

where $M_q = \mathbb{E}^{1/q} |\mathcal{N}(0, 1)|^q$ and $\Gamma_{\eta; (\Sigma, \|\mu_0\|)}^{\mathbf{A}} = \mathbf{A}(\Sigma + \tau_{\eta, *}\mathbf{I}_n)^{-1}(\tilde{\gamma}_{\eta, *}^2(\|\mu_0\|)\Sigma + \tau_{\eta, *}^2\|\mu_0\|^2\mathbf{I}_n)(\Sigma + \tau_{\eta, *}\mathbf{I}_n)^{-1}\mathbf{A}$; see Theorem 6 for the precise definition of $\tilde{\gamma}_{\eta, *}^2(\|\mu_0\|)$ and the formal statement of the above result (10).

Beyond providing a precise characterization of all ℓ_q risks, the uniform nature of (8) also illuminates novel insights into certain global, qualitative behavior of the most commonly studied ℓ_2 risks for finite samples. To fix notation, we define

- (*prediction risk*) $R_{(\Sigma, \mu_0)}^{\text{pred}}(\eta) \equiv \|\Sigma^{1/2}(\hat{\mu}_\eta - \mu_0)\|^2$,
- (*estimation risk*) $R_{(\Sigma, \mu_0)}^{\text{est}}(\eta) \equiv \|\hat{\mu}_\eta - \mu_0\|^2$,
- (*in-sample risk*) $R_{(\Sigma, \mu_0)}^{\text{in}}(\eta) \equiv n^{-1}\|X(\hat{\mu}_\eta - \mu_0)\|^2$.

Using our uniform distributional characterization in (8), we show that for ‘most’ μ_0 ’s, the global optimum of $\eta \mapsto R_{(\Sigma, \mu_0)}^{\#}(\eta)$ for all $\# \in \{\text{pred, est, in}\}$ will be achieved approximately at the same point $\eta_* = \text{SNR}_{\mu_0}^{-1}$ with high probability¹ (cf. Theorem 54).

It must be stressed that, for different $\# \in \{\text{pred, est, in}\}$, the empirical risk curves $\eta \mapsto R_{(\Sigma, \mu_0)}^{\#}(\eta)$ concentrate on genuinely different deterministic counterparts $\eta \mapsto \bar{R}_{(\Sigma, \mu_0)}^{\#}(\eta)$ with different mathematical expressions (cf. Theorem 49). As such, there are no a priori reasons to expect that these risk curves share approximately the same global minimum. Remarkably, as a consequence of the approximate formulae for the deterministic risk curves $\eta \mapsto \bar{R}_{(\Sigma, \mu_0)}^{\#}(\eta)$ (cf. Theorem 8), we show that the curves $\eta \mapsto \bar{R}_{(\Sigma, \mu_0)}^{\#}(\eta)$ are qualitatively similar, in that they approximately behave locally like a quadratic function centered around $\eta_* = \text{SNR}_{\mu_0}^{-1}$ (cf. Proposition 9), at least for ‘most’ signal μ_0 ’s.

1.4 Cross-validation: optimality beyond prediction

The discussion in Section 1.3 naturally raises the question of how one can choose the optimal regularization in a data-driven manner. Here we study two widely used adaptive tuning methods, namely,

1. the generalized cross-validation scheme $\hat{\eta}^{\text{GCV}}$, and
2. the k -fold cross-validation scheme $\hat{\eta}^{\text{CV}}$.

The readers are referred to (22) and (24) for precise definitions and literature review of $\hat{\eta}^{\text{GCV}}, \hat{\eta}^{\text{CV}}$ in the context of Ridge regression.

By design, both methods $\hat{\eta}^{\text{GCV}}, \hat{\eta}^{\text{CV}}$ are intended to estimate the prediction risk, so it is natural to expect that they perform well for the task of prediction. Interestingly, the insight from Section 1.3 suggests a far broader utility of these adaptive tuning methods. Indeed, as all the empirical risk curves $\eta \mapsto R_{(\Sigma, \mu_0)}^{\#}(\eta)$ are approximately minimized at the same point $\eta_* = \text{SNR}_{\mu_0}^{-1}$, it is reasonable to conjecture that $\hat{\eta}^{\text{GCV}}, \hat{\eta}^{\text{CV}}$ could also yield optimal

1. Here $\text{SNR}_{\mu_0} = \|\mu_0\|^2/\sigma_\xi^2$ is the usual notion of signal-to-noise ratio; when $\mu_0 \neq 0$ and $\sigma_\xi^2 = 0$, we shall interpret $\text{SNR}_{\mu_0}^{-1} = 0$.

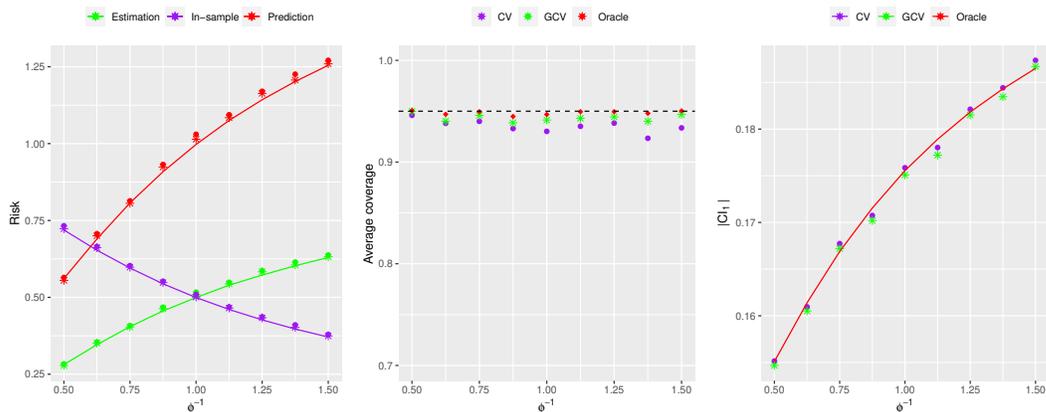


Figure 1: *Left panel*: Comparison between empirical risks and theoretical risks for $* = \text{GCV}$ and $\bullet = \text{CV}$ with $k = 5$. *Middle panel*: Averaged CI coverage $\mathcal{C}^{\text{dR}}(\hat{\eta}^\#)$ for $\# \in \{\text{GCV}, \text{CV}\}$ and the oracle $\mathcal{C}^{\text{dR}}(\eta_*)$. *Right panel*: CI length of $CI_1(\hat{\eta}^\#)$ for $\# \in \{\text{GCV}, \text{CV}\}$ and the oracle CI length. See Section 4 for the precise definitions.

performance for estimation and in-sample risks. We show in Theorems 10 and 12 that this is indeed the case: for ‘most’ signal μ_0 ’s and all $\# \in \{\text{pred}, \text{est}, \text{in}\}$, with high probability,

$$R_{(\Sigma, \mu_0)}^\#(\hat{\eta}^{\text{GCV}}), R_{(\Sigma, \mu_0)}^\#(\hat{\eta}^{\text{CV}}) \approx \min_{\eta \in [0, K]} R_{(\Sigma, \mu_0)}^\#(\eta). \quad (11)$$

A typical simulation for this phenomenon is reported in the left panel of Figure 1 above, where empirical risks tuned by $\hat{\eta}^{\text{GCV}}$ (in $*$) and $\hat{\eta}^{\text{CV}}$ (in \bullet) achieve optimal theoretical risks (in solid lines) for estimation and in-sample risks as well.

Even more surprisingly, the optimality of $\hat{\eta}^{\text{GCV}}, \hat{\eta}^{\text{CV}}$ extends to the much more challenging task of statistical inference. In fact, we show in Theorem 13 that within the so-called debiased Ridge scheme, these two adaptive tuning methods $\hat{\eta}^{\text{GCV}}, \hat{\eta}^{\text{CV}}$ yield an asymptotically valid construction of confidence intervals for the coordinates of μ_0 with the shortest possible length. This is numerically validated in the middle and right panels of Figure 1.

To the best of our knowledge, theoretical optimality properties for the cross-validation schemes beyond the realm of prediction accuracy has not been established in the literature, either for Ridge regression or for other regularized regression estimators.

On the other hand, in the related Lasso setting, some numerical evidence for the broader utility of cross-validation and other adaptive tuning methods is reported in Miolane and Montanari (2021, Figure 1). There it is shown that the SURE method, which is designed to tune in-sample risk, nearly matches the performance of k -fold cross-validation in prediction tasks, despite not being expected to perform well in prediction apriori. Our findings here in the context of Ridge regression can therefore be viewed as a first step toward understanding the broader potential of cross validation and other adaptive tuning schemes for a wider range of statistical inference problems.

1.5 Further literature

1.5.1 RELATION TO MEAN-FIELD ASYMPTOTICS

Our distributional theory (8) for the Ridge(less) estimator is closely related to a recent line of research that examines the mean-field behavior of statistical estimators in the proportional regime $m \asymp n$, see, e.g. Bayati and Montanari (2012); Donoho and Montanari (2016); Thrampoulidis et al. (2018); Sur and Candès (2019); Li and Wei (2021); Miolane and Montanari (2021); Liang and Sur (2022); Celentano et al. (2023); Han and Shen (2023) for an incomplete list and many more references can be found therein.

A key feature of this line of works is the use of a simplified ‘effective’ regression problem to understand the complicated behavior of the original statistical estimator. For instance, in the closely related Lasso setting, the ‘equivalence’ between the Lasso estimator $\widehat{\mu}_\eta^L$ in the linear model and a corresponding Lasso estimator in the sequence model $\widehat{\mu}_{\Sigma, \mu_0}^{\text{seq}, L}(\gamma_{\eta, *}, \tau_{\eta, *})$ has been established under Gaussian designs with positive regularization. This equivalence was first shown for ℓ_2 -type risks in Bayati and Montanari (2012), and later in the distributional sense akin to (8) in Miolane and Montanari (2021); Celentano et al. (2023). Such equivalence for Lasso is further extended to the interpolating regime in Li and Wei (2021) for the ℓ_2 risk under a standard Gaussian isotropic design. Our theory (8) here can thus be placed into a similar position as the progress made in Miolane and Montanari (2021); Celentano et al. (2023) over the Lasso risk characterization in Bayati and Montanari (2012), but now in the context of Ridge(less) estimator beyond a purely ℓ_2 risk as obtained in the references cited above.

While we have developed our distributional theory (8) primarily in the proportional regime $m \asymp n$, we conjecture that our theory (8) remains valid in the full nonparametric regime in which the ℓ_2 risk of the Ridge(less) estimator exceeds $\mathcal{O}(m^{-1/2})$. Some progress in this direction is made in Han (2023) in a related context of convex-constrained least squares estimator under a Gaussian design.

1.5.2 RELATION TO EXISTING INTERPRETATION OF ‘IMPLICIT REGULARIZATION’

A separate line of research (Bartlett et al., 2020, 2021; Tsigler and Bartlett, 2023) offers a different perspective on the implicit regularization phenomenon within the Ridgeless interpolator $\widehat{\mu}_0$. Specifically, by writing $X = [X_{\leq k}, X_{> k}]$ with ‘effective dimension’ k and expressing the Ridgeless interpolator as $\widehat{\mu}_0 = \bar{X}^\top (X_{\leq k} X_{\leq k}^\top + X_{> k} X_{> k}^\top)^{-1} Y$, this line of research identifies covariance structures Σ for which $X_{> k} X_{> k}^\top$ scales proportionally to the identity matrix (in a suitable sense). This implies that $X_{> k} X_{> k}^\top$ qualitatively plays the same role as if positive Ridge regularization were applied to the effective data $X_{\leq k}$. In particular, this line of theory suggests that the prediction risk of $\widehat{\mu}_0$ can indeed vanish (i.e., benign overfitting), provided that the eigen-decay of Σ is neither too fast nor too slow.

While this approach is insightful, it falls short of providing an *exact* understanding for the emergence of the implicit regularization phenomenon. This is so, as this approach seeks sufficient conditions for $X_{> k} X_{> k}^\top \asymp I$, and produces risk bounds for $\widehat{\mu}_0$ modulo unspecified multiplicative constants. In contrast, our characterization of the implicit regularization via (9) is exact up to the leading constant order, and is susceptible to be also exact in other

regimes as well; see Cheng and Montanari (2024) for some recent partial progress along this line.

Furthermore, both the approaches of Bartlett et al. (2020, 2021); Tsigler and Bartlett (2023) and Cheng and Montanari (2024) rely heavily on the closed form of the Ridgeless interpolator and thus do not generalize to more general interpolators. In contrast, a significant technical advantage of our characterization of implicit regularization via (9) lies in its natural connection to the mean-field theory for general regression estimators. This suggests a general paradigm to quantify the implicit regularization for a large class of interpolators via mean-field asymptotics. For instance, the minimum ℓ_1 -norm interpolator studied in Li and Wei (2021) demonstrates implicit regularization in the prediction risk that can be characterized via $\tau_{0,*}^L$ as the ‘Lasso’ version of (9). Our approach developed here for Ridgeless interpolator is expected to be useful for quantifying the implicit regularization phenomenon for a more general class of interpolators.

1.6 Organization

The rest of the paper is organized as follows. In Section 2, we present our main results on the distributional characterizations (8) of the Ridge(less) estimator $\hat{\mu}_\eta$. In Section 3, we provide a number of approximate ℓ_q risk formulae, and derive the optimal regularization for ℓ_2 risks. In Section 4, we give a formal validation for the two cross validation schemes mentioned above, both in terms of (11) and statistical inference via the debiased Ridge estimator. Due to the high technicalities involved in the proof of (8), a proof outline will be given in Section 5. All the proof details are then presented in Appendices A to G.

1.7 Notation

For any positive integer n , let $[n] = [1 : n]$ denote the set $\{1, \dots, n\}$. For $a, b \in \mathbb{R}$, $a \vee b \equiv \max\{a, b\}$ and $a \wedge b \equiv \min\{a, b\}$. For $a \in \mathbb{R}$, let $a_\pm \equiv (\pm a) \vee 0$. For $x \in \mathbb{R}^n$, let $\|x\|_p$ denote its p -norm ($0 \leq p \leq \infty$), and $B_{n;p}(R) \equiv \{x \in \mathbb{R}^n : \|x\|_p \leq R\}$. We simply write $\|x\| \equiv \|x\|_2$ and $B_n(R) \equiv B_{n;2}(R)$. For a matrix $M \in \mathbb{R}^{m \times n}$, let $\|M\|_{\text{op}}, \|M\|_F$ denote the spectral and Frobenius norm of M , respectively. I_n is reserved for an $n \times n$ identity matrix, written simply as I (in the proofs) if no confusion arises. For a square matrix $M \in \mathbb{R}^{n \times n}$, we let $\text{diag}(M) \equiv (M_{ii})_{i=1}^n \in \mathbb{R}^n$.

We use C_x to denote a generic constant that depends only on x , whose numeric value may change from line to line unless otherwise specified. $a \lesssim_x b$ and $a \gtrsim_x b$ mean $a \leq C_x b$ and $a \geq C_x b$, abbreviated as $a = \mathcal{O}_x(b)$, $a = \Omega_x(b)$ respectively; $a \asymp_x b$ means $a \lesssim_x b$ and $a \gtrsim_x b$, abbreviated as $a = \Theta_x(b)$. \mathcal{O} and \mathfrak{o} (resp. $\mathcal{O}_{\mathbf{P}}$ and $\mathfrak{o}_{\mathbf{P}}$) denote the usual big and small O notation (resp. in probability). For a random variable X , we use $\mathbb{P}^X, \mathbb{E}^X$ (resp. $\mathbb{P}^X, \mathbb{E}^X$) to indicate that the probability and expectation are taken with respect to X (resp. conditional on X).

For a measurable map $f : \mathbb{R}^n \rightarrow \mathbb{R}$, let $\|f\|_{\text{Lip}} \equiv \sup_{x \neq y} |f(x) - f(y)| / \|x - y\|$. f is called L -Lipschitz iff $\|f\|_{\text{Lip}} \leq L$. For a proper, closed convex function f defined on \mathbb{R}^n , its Moreau envelope $\mathbf{e}_f(\cdot; \tau)$ and proximal operator $\text{prox}_f(\cdot; \tau)$ for any $\tau > 0$ are defined by $\mathbf{e}_f(x; \tau) \equiv \min_{z \in \mathbb{R}^n} \left\{ \frac{1}{2\tau} \|x - z\|^2 + f(z) \right\}$ and $\text{prox}_f(x; \tau) \equiv \arg \min_{z \in \mathbb{R}^n} \left\{ \frac{1}{2\tau} \|x - z\|^2 + f(z) \right\}$.

Throughout this paper, for an invertible covariance matrix $\Sigma \in \mathbb{R}^{n \times n}$, we write $\mathcal{H}_\Sigma \equiv \text{tr}(\Sigma^{-1})/n$ as the harmonic mean of the eigenvalues of Σ .

2. Distribution of Ridge(less) estimators

2.1 Some definitions

For $K > 1$, let

$$\Xi_K \equiv [\mathbf{1}_{\phi^{-1} < 1+1/K} K^{-1}, K]. \quad (12)$$

This notation will be used throughout the paper for uniform-in- η statements. In particular, in the overparametrized regime $\phi^{-1} \geq 1 + 1/K$, we have $\Xi_K = [0, K]$.

Next, for $\gamma, \tau \geq 0$, recall $\hat{\mu}_{(\Sigma, \mu_0)}^{\text{seq}}(\gamma; \tau)$ in (6), and we define its associated estimation error $\text{err}_{(\Sigma, \mu_0)}(\gamma; \tau)$ and the degrees-of-freedom $\text{dof}_{(\Sigma, \mu_0)}(\gamma; \tau)$ as

$$\begin{cases} \text{err}_{(\Sigma, \mu_0)}(\gamma; \tau) \equiv \|\Sigma^{1/2}(\hat{\mu}_{(\Sigma, \mu_0)}^{\text{seq}}(\gamma; \tau) - \mu_0)\|^2, \\ \text{dof}_{(\Sigma, \mu_0)}(\gamma; \tau) \equiv \left\langle \frac{\gamma g}{\sqrt{n}}, \Sigma^{1/2}(\hat{\mu}_{(\Sigma, \mu_0)}^{\text{seq}}(\gamma; \tau) - \mu_0) \right\rangle. \end{cases}$$

The $\text{dof}_{(\Sigma, \mu_0)}(\gamma; \tau)$ defined above is naturally related to the usual notion of degrees-of-freedom (cf. Stein (1981); Efron (2004)) for $\hat{\mu}_{(\Sigma, \mu_0)}^{\text{seq}}(\gamma; \tau)$, in the sense that $\text{df}(\hat{\mu}_{(\Sigma, \mu_0)}^{\text{seq}}(\gamma; \tau)) \equiv \sum_{j=1}^n \frac{1}{\gamma^2/n} \text{Cov}((\Sigma^{1/2} \hat{\mu}_{(\Sigma, \mu_0)}^{\text{seq}})_j, y_{(\Sigma, \mu_0), j}^{\text{seq}}) = \frac{n}{\gamma^2} \mathbb{E} \text{dof}_{(\Sigma, \mu_0)}(\gamma; \tau)$.

2.2 Working assumptions

Assumption A. $X = Z\Sigma^{1/2}$, where (i) $Z \in \mathbb{R}^{m \times n}$ has independent, mean-zero, unit variance, uniformly sub-gaussian entries, and (ii) $\Sigma \in \mathbb{R}^{n \times n}$ is an invertible covariance matrix with eigenvalues $\lambda_1 \geq \dots \geq \lambda_n > 0$.

Here ‘uniform sub-gaussianity’ means $\sup_{i \in [m], j \in [n]} \|Z_{ij}\|_{\psi_2} \leq C$ for some universal $C > 0$, where ψ_2 is the Orlicz 2-norm (cf. van der Vaart and Wellner (1996, Section 2.2, pp. 95)).

We shall often write the Gaussian design as $Z = G$, where $G \in \mathbb{R}^{m \times n}$ consists of i.i.d. $\mathcal{N}(0, 1)$ entries.

Assumption B. $\xi = \sigma_\xi \cdot \xi_0$ for some ξ_0 with i.i.d. mean zero, unit variance and uniform sub-gaussian entries.

Remark 1. *The requirement on the noise level σ_ξ^2 will be specified in concrete results below. We assume sub-gaussian noise for simplicity, but our proofs use it only through the high probability events in Section A.2. With easy modifications, all results extend to more general noise distributions, including certain heavy-tailed or weakly dependent cases for which these events still hold.*

2.3 The fixed point equation

Fix $\eta \geq 0$. Consider the following fixed point equation in (γ, τ) :

$$\begin{cases} \phi \gamma^2 = \sigma_\xi^2 + \mathbb{E} \text{err}_{(\Sigma, \mu_0)}(\gamma; \tau), \\ \phi - \frac{\eta}{\tau} = \frac{1}{n} \text{tr}((\Sigma + \tau I_n)^{-1} \Sigma) = \frac{1}{\gamma^2} \mathbb{E} \text{dof}_{(\Sigma, \mu_0)}(\gamma; \tau). \end{cases} \quad (13)$$

Fixed point equations of the type described above have appeared in the general mean-field theory for high dimensional regularized least squares estimators (LSEs), see, e.g. Bayati

and Montanari (2012); Bu et al. (2021); Li and Wei (2021); Han (2023); Celentano et al. (2023) for a sample of this type of equations in the i.i.d. sampling setting, and Bao et al. (2025) in the i.n.i.d. sampling setting. A common theme of these works characterizes the behavior of the regularized LSE in the linear model—at various levels of generality—via a regularized LSE in the equivalent sequence model, whose ‘effective noise’ and ‘effective regularization’ are determined by the solution pair to the fixed point equation.

In the context of Ridge regression, the form of the fixed point equation (13) appeared in Bartlett et al. (2021); Cheng and Montanari (2024) for the purpose of characterizing ℓ_2 risks for the Ridge(less) estimator $\widehat{\mu}_\eta$. It is now well understood that for the purpose of distributional characterizations of $\widehat{\mu}_\eta$, further stability properties for the solution pair to the fixed point equation will be needed (Miolane and Montanari, 2021; Celentano et al., 2023; Han, 2023). We establish these properties for the solution to (13) in the following proposition. Recall Ξ_K from (12).

Proposition 2. *Recall $\mathcal{H}_\Sigma = \text{tr}(\Sigma^{-1})/n$. The following hold.*

1. *The fixed point equation (13) admits a unique solution $(\gamma_{\eta,*}, \tau_{\eta,*}) \in (0, \infty)^2$, for all $(m, n) \in \mathbb{N}^2$ when $\eta > 0$ and $m < n$ when $\eta = 0$.*
2. *Suppose $1/K \leq \phi^{-1} \leq K$ and $\|\Sigma\|_{\text{op}} \vee \mathcal{H}_\Sigma \leq K$ for some $K > 1$. Then there exists some $C = C(K) > 1$ such that uniformly in $\eta \in \Xi_K$,*

$$1/C \leq \tau_{\eta,*} \leq C, \quad 1/C \leq (-1)^{q+1} \partial_\eta^q \tau_{\eta,*} \leq C, \quad q \in \{1, 2\}.$$

If furthermore $1/K \leq \sigma_\xi^2 \leq K$ and $\|\mu_0\| \leq K$, then uniformly in $\eta \in \Xi_K$,

$$1/C \leq \gamma_{\eta,*} \leq C, \quad |\partial_\eta \gamma_{\eta,*}| \leq C.$$

3. *Suppose $1/K \leq \phi^{-1} \leq K$ and $\|\Sigma\|_{\text{op}} \vee \mathcal{H}_\Sigma \leq K$ for some $K > 1$. Then there exists some $C = C(K) > 1$ such that the following hold. For any $\varepsilon \in (0, 1/2]$, we may find some $\mathcal{U}_\varepsilon \subset B_n(1)$ with $\text{vol}(\mathcal{U}_\varepsilon)/\text{vol}(B_n(1)) \geq 1 - C\varepsilon^{-1}e^{-n\varepsilon^2/C}$,*

$$\sup_{\mu_0 \in \mathcal{U}_\varepsilon} \sup_{\eta \in \Xi_K} |\gamma_{\eta,*}^2 - \widetilde{\gamma}_{\eta,*}^2(\|\mu_0\|)| \leq \varepsilon,$$

where $\widetilde{\gamma}_{\eta,}^2(\|\mu_0\|) \equiv \sigma_\xi^2 \partial_\eta \tau_{\eta,*} + \|\mu_0\|^2 (\tau_{\eta,*} - \eta \partial_\eta \tau_{\eta,*}) > 0$. When $\Sigma = I_n$, we may take $\mathcal{U}_\varepsilon = B_n(1)$ and the above inequality holds with $\varepsilon = 0$.*

The above proposition combines parts of Propositions 23 and 52.

As an important qualitative consequence of (2), under the condition $\|\Sigma\|_{\text{op}} \vee \mathcal{H}_\Sigma \leq K$, the effective regularization $\eta \mapsto \tau_{\eta,*}$ is a strictly increasing and concave function of η . Moreover, in the overparametrized regime $\phi^{-1} > 1$, the quantity $\tau_{0,*}$ —also known as ‘implicit regularization’ in the literature (Bartlett et al., 2020, 2021; Hastie et al., 2022; Tsigler and Bartlett, 2023; Cheng and Montanari, 2024)—is strictly bounded away from zero.

The claim in (3) offers a useful approximate representation of the effective noise $\gamma_{\eta,*}^2$ in terms of the original noise σ_ξ^2 , the effective regularization $\tau_{\eta,*}$ and the signal energy $\|\mu_0\|$ without explicitly dependence of Σ . This representation will prove useful in understanding some qualitative aspects of the risk curves in Section 3.2 ahead.

2.4 Distribution of Ridge(less) estimators

In addition to $\widehat{\mu}_\eta$, we will also consider the distribution of the (scaled) residual \widehat{r}_η , defined by

$$\widehat{r}_\eta \equiv \frac{1}{\sqrt{n}}(Y - X\widehat{\mu}_\eta). \quad (14)$$

We define the ‘population’ version of \widehat{r}_η as

$$r_{\eta,*} \equiv \frac{\eta}{\phi\tau_{\eta,*}} \left(-\sqrt{\phi\gamma_{\eta,*}^2 - \sigma_\xi^2} \cdot \frac{h}{\sqrt{n}} + \frac{\xi}{\sqrt{n}} \right). \quad (15)$$

Here $h \sim \mathcal{N}(0, I_m)$ is independent of ξ .

We are now in a position to state our main results on the distributional results for the Ridge(less) estimator $\widehat{\mu}_\eta$ and the residual \widehat{r}_η .

First we work under the Gaussian design $Z = G$, and we write $\widehat{\mu}_\eta = \widehat{\mu}_{\eta;G}$, $\widehat{r}_\eta = \widehat{r}_{\eta;G}$. Recall $\mathcal{H}_\Sigma = \text{tr}(\Sigma^{-1})/n$ and Ξ_K from (12).

Theorem 3. *Suppose Assumption A holds with $Z = G$ and the following hold for some $K > 0$.*

- $1/K \leq \phi^{-1} \leq K$, $\|\Sigma\|_{\text{op}} \vee \mathcal{H}_\Sigma \leq K$.
- Assumption B holds with $\sigma_\xi^2 \in [1/K, K]$.

Then there exists some constant $C = C(K) > 0$ such that the following hold.

1. For any 1-Lipschitz function $\mathbf{g} : \mathbb{R}^n \rightarrow \mathbb{R}$ and $\varepsilon \in (0, 1/2]$,

$$\sup_{\mu_0 \in B_n(1)} \mathbb{P} \left(\sup_{\eta \in \Xi_K} |\mathbf{g}(\widehat{\mu}_{\eta;G}) - \mathbb{E} \mathbf{g}(\widehat{\mu}_{(\Sigma, \mu_0)}^{\text{seq}}(\gamma_{\eta,*}, \tau_{\eta,*}))| \geq \varepsilon \right) \leq Cne^{-n\varepsilon^4/C}.$$

Here we recall that $\widehat{\mu}_{\eta;G}$ is defined in (2) with $X = G\Sigma^{1/2}$, and that $\widehat{\mu}_{(\Sigma, \mu_0)}^{\text{seq}}$ with effective noise and regularization pair $(\gamma_{\eta,*}, \tau_{\eta,*})$ is defined in (6)-(7).

2. For any $\varepsilon \in (0, 1/2]$, $\xi \in \mathbb{R}^m$ satisfying $|\|\xi\|^2/m - \sigma_\xi^2| \leq \varepsilon^2/C$, and 1-Lipschitz function $\mathbf{h} : \mathbb{R}^m \rightarrow \mathbb{R}$ (which may depend on ξ),

$$\sup_{\mu_0 \in B_n(1)} \mathbb{P}^\xi \left(\sup_{\eta \in [1/K, K]} |\mathbf{h}(\widehat{r}_{\eta;G}) - \mathbb{E}^\xi \mathbf{h}(r_{\eta,*})| \geq \varepsilon \right) \leq Cne^{-n\varepsilon^4/C}.$$

The choice $\mu_0 \in B_n(1)$ is made merely for simplicity of presentation; it can be replaced by $\mu_0 \in B_n(R)$ with another constant C that depends further on R . The assumption $\mathcal{H}_\Sigma \lesssim 1$ is quite common in the literature of Ridge(less) regression; see, e.g., (Bartlett et al., 2021, Assumption 4.12) or a slight variant in (Montanari et al., 2025, Assumption 1). The major assumption in the above theorem is the Gaussianity on the design X . This may be lifted at the cost of a set of slightly stronger conditions.

Theorem 4. *Suppose Assumption A holds and the following hold for some $K > 0$.*

- $1/K \leq \phi^{-1} \leq K$, $\|\Sigma\|_{\text{op}} \vee \|\Sigma^{-1}\|_{\text{op}} \leq K$.
- Assumption B holds with $\sigma_\xi^2 \in [1/K, K]$.

Fix $\vartheta \in (0, 1/18)$. There exist some $C = C(K, \vartheta) > 0$ and two measurable sets $\mathcal{U}_\vartheta \subset B_n(1), \mathcal{E}_\vartheta \subset \mathbb{R}^m$ with $\min\{\text{vol}(\mathcal{U}_\vartheta)/\text{vol}(B_n(1)), \mathbb{P}(\xi \in \mathcal{E}_\vartheta)\} \geq 1 - Ce^{-n^{2\vartheta}/C}$, such that the following hold.

1. For any 1-Lipschitz function $\mathbf{g} : \mathbb{R}^n \rightarrow \mathbb{R}$, and $\varepsilon \in (0, 1/2]$,

$$\sup_{\mu_0 \in \mathcal{U}_\vartheta} \mathbb{P} \left(\sup_{\eta \in \Xi_K} |\mathbf{g}(\widehat{\mu}_\eta) - \mathbb{E} \mathbf{g}(\widehat{\mu}_{(\Sigma, \mu_0)}^{\text{seq}}(\gamma_{\eta, *}; \tau_{\eta, *}))| \geq \varepsilon \right) \leq C\varepsilon^{-13} n^{-1/6+3\vartheta}.$$

Here we recall that $\widehat{\mu}_\eta$ is defined in (2), and that $\widehat{\mu}_{(\Sigma, \mu_0)}^{\text{seq}}$ with effective noise and regularization pair $(\gamma_{\eta, *}, \tau_{\eta, *})$ is defined in (6)-(7).

2. For any $\varepsilon \in (0, 1/2]$, $\xi \in \mathcal{E}_\vartheta$ and 1-Lipschitz function $\mathbf{h} : \mathbb{R}^m \rightarrow \mathbb{R}$ (which may depend on ξ),

$$\sup_{\mu_0 \in \mathcal{U}_\vartheta} \mathbb{P}^\xi \left(\sup_{\eta \in [1/K, K]} |\mathbf{h}(\widehat{r}_\eta) - \mathbb{E}^\xi \mathbf{h}(r_{\eta, *})| \geq \varepsilon \right) \leq C\varepsilon^{-7} n^{-1/6+3\vartheta}.$$

Concrete forms of $\mathcal{U}_\vartheta, \mathcal{E}_\vartheta$ are specified in Proposition 40.

Remark 5. Compared with the Gaussian case (Theorem 3), which admits exponential tails via Gaussian concentration and a direct CGMT argument, the sub-Gaussian universality in Theorem 4 yields only polynomial rates. This stems from the quantitative comparison inequalities Han and Shen (2023) employed in the universality step. Extending these bounds to exponential decay would likely require methods beyond the comparison framework, which is beyond the scope of the present paper.

Theorems 3 and 4 are proved in Section C and Section D, respectively. Due to the high technicalities in the proof, a sketch is outlined in Section 5. These distributional results are the main input for all the applications developed in the subsequent sections. In particular, the flexibility in the choice of the test functions \mathbf{g} and \mathbf{h} allows us to obtain a variety of functionals of interest. By choosing the test functions appropriately, we derive (i) in Section 3 the ℓ_q -risk asymptotics for general $q \in [1, \infty)$, extending the classical ℓ_2 -risk formulas that are typically accessible via random matrix theory, and (ii) in Section 4 the optimality of cross-validation tuning rules.

We mention two particular important features on the theorems above:

1. The distributional characterizations for $\widehat{\mu}_\eta$ in both theorems above are uniformly valid down to the interpolation regime $\eta = 0$ for $\phi^{-1} > 1$. This uniform control will play a crucial role in our non-asymptotic analysis of cross-validation methods to be studied in Section 4 ahead.
2. The distribution of the residual \widehat{r}_η in (2) is formulated *conditional on the noise* ξ . A fundamental reason for adopting this formulation is that the distribution of \widehat{r}_η is *not* universal with respect to the law of ξ . In other words, one cannot simply assume Gaussianity of ξ in Theorem 3 in hope of proving universality of \widehat{r}_η in Theorem 4.

In the context of distributional characterizations for regularized regression estimators in the proportional regime, results in similar vein to Theorem 3 have been obtained in the closely related Lasso setting for isotropic $\Sigma = I_n$ in Miolane and Montanari (2021), and for general Σ in Celentano et al. (2023), both under Gaussian designs and with strictly non-vanishing regularization. A substantially simpler, isotropic ($\Sigma = I_n$) version of Theorem 4 is obtained in Han and Shen (2023) that holds pointwise in non-vanishing regularization level $\eta > 0$. As will be clear from the proof sketch in Section 5, in addition to the complications due to the implicit nature of the solution to the fixed point equation (13) for general Σ , the major difficulty in proving Theorems 3 and 4 rests in handling the singularity of the optimization problem (2) as $\eta \downarrow 0$.

3. General ℓ_q -type risk formulae

As a demonstration of the analytic power of Theorems 3 and 4, this section will be devoted to a detailed study for the ℓ_q -type risks for Ridgeless interpolators. We then conduct a more in-depth study of ℓ_2 risks, where techniques from RMT lead to a detailed characterization of the optimal regularization strategy for risk minimization.

3.1 Weighted ℓ_q risks and delocalization

We compute below the weighted ℓ_q risk $\|\mathbf{A}(\widehat{\mu}_\eta - \mu_0)\|_q$ for a well-behaved matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$ and $q \in [1, \infty)$. Recall Ξ_K from (12).

Theorem 6. *Suppose the same conditions in Theorem 4 hold for some $K > 0$. Fix $q \in [1, \infty)$ and a p.s.d. matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$ with $\|\mathbf{A}\|_{\text{op}} \vee \|\mathbf{A}^{-1}\|_{\text{op}} \leq K$. Then there exist constants $C > 1, \vartheta \in (0, 1/50)$ depending on K, q , and a measurable set $\mathcal{U}_\vartheta \subset B_n(1)$ with $\text{vol}(\mathcal{U}_\vartheta)/\text{vol}(B_n(1)) \geq 1 - Ce^{-n^\vartheta/C}$, such that*

$$\sup_{\mu_0 \in \mathcal{U}_\vartheta} \mathbb{P} \left(\sup_{\eta \in \Xi_K} \left| \frac{\|\mathbf{A}(\widehat{\mu}_\eta - \mu_0)\|_q}{\bar{R}_{(\Sigma, \mu_0); q}^{\mathbf{A}}(\eta)} - 1 \right| \geq n^{-\vartheta} \right) \leq Cn^{-1/7}.$$

Here $\bar{R}_{(\Sigma, \mu_0); q}^{\mathbf{A}}(\eta) \in \{ \mathbb{E} \|\mathbf{A}(\widehat{\mu}_{(\Sigma, \mu_0)}^{\text{seq}}(\gamma_{\eta, *}; \tau_{\eta, *}) - \mu_0)\|_q, n^{-1/2} \|\text{diag}(\Gamma_{\eta; (\Sigma, \|\mu_0\|)}^{\mathbf{A}})\|_{q/2}^{1/2} M_q \}$, where $M_q \equiv \mathbb{E}^{1/q} |\mathcal{N}(0, 1)|^q = 2^{1/2} \{\Gamma((q+1)/2)/\sqrt{\pi}\}^{1/q}$,

$$\Gamma_{\eta; (\Sigma, \|\mu_0\|)}^{\mathbf{A}} \equiv \mathbf{A}(\Sigma + \tau_{\eta, *} I_n)^{-1} \left(\tilde{\gamma}_{\eta, *}^2 (\|\mu_0\|) \Sigma + \tau_{\eta, *}^2 \|\mu_0\|^2 I_n \right) (\Sigma + \tau_{\eta, *} I_n)^{-1} \mathbf{A}, \quad (16)$$

and $\tilde{\gamma}_{\eta, *}^2(\|\mu_0\|)$ is defined in Proposition 2-(3).

The proof of the above theorem can be found in Section E. To the best of our knowledge, general weighted ℓ_q risks for the Ridge(less) estimator $\widehat{\mu}_\eta$ have not been available in the literature except for the special case $q = 2$, for which $\|\mathbf{A}(\widehat{\mu}_\eta - \mu_0)\|_2$ admits a closed-form expression in terms of the spectral statistics of X that facilitates direct applications of RMT techniques, cf. Tulino et al. (2004); El Karoui (2013); Dicker (2016); Dobriban and Wager (2018); El Karoui (2018); Advani et al. (2020); Wu and Xu (2020); Bartlett et al. (2021); Richards et al. (2021); Hastie et al. (2022); Cheng and Montanari (2024).

For $\mathbf{A} \neq \Sigma$, Theorem 6 above characterizes the out-of-distribution ℓ_q risk for the Ridge(less) estimators. This setting is naturally related to the covariate shift setting, where

ℓ_2 -type risks are studied in Patil et al. (2024); Tang et al. (2024) using random matrix methods in slightly different specific settings.

Let us remark that obtaining ℓ_q risks for $q \in [1, 2]$ via our Theorems 3 and 4 is relatively easy, as $x \mapsto \|x\|_q/n^{1/q-1/2}$ is 1-Lipschitz with respect to $\|\cdot\|$ for $q \in [1, 2]$. The stronger norm case $q \in (2, \infty)$ is significantly harder. In fact, we need additionally the following delocalization result for $\widehat{\mu}_\eta$.

Proposition 7. *Suppose the same conditions as in Theorem 6 hold for some $K > 0$. Fix $\vartheta \in (0, 1/2]$. Then there exist some constant $C = C(K, \vartheta) > 0$ and a measurable set $\mathcal{U}_\vartheta \subset B_n(1)$ with $\text{vol}(\mathcal{U}_\vartheta)/\text{vol}(B_n(1)) \geq 1 - Ce^{-n^{2\vartheta}/C}$, such that*

$$\sup_{\mu_0 \in \mathcal{U}_\vartheta} \mathbb{P} \left(\sup_{\eta \in \Xi_K} \|A(\widehat{\mu}_\eta - \mu_0)\|_\infty \geq Cn^{-1/2+\vartheta} \right) \leq Cn^{-100}.$$

The above proposition is a simplified version of Proposition 40, proved via the anisotropic local laws developed in Knowles and Yin (2017). In essence, delocalization allows us to apply Theorems 3 and 4 with a truncated version of the ℓ_q norm ($q > 2$) with a well-controlled Lipschitz constant with respect to ℓ_2 . Moreover, delocalization of $\widehat{\mu}_\eta$ also serves as a key technical ingredient in proving the universality Theorem 4; the readers are referred to Section 5 for a detailed account on the technical connection between delocalization and universality.

Convention on probability estimates:

1. When $Z = G$, $n^{-1/7}$ in Theorem 6 can be replaced by n^{-D} for any $D > 0$.
2. n^{-100} in Proposition 7 can be replaced by n^{-D} for any $D > 0$.

The cost will be a possibly enlarged constant $C > 0$ that depends further on D . This convention applies to other statements in the following sections in which the probability estimates $n^{-1/7}, n^{-100}$ appear.

3.2 ℓ_2 risk formulae and optimal regularization

In this subsection, we will study in some detail the behavior of various ℓ_2 risks associated with $\widehat{\mu}_\eta$. As will be clear below, a major analytic advantage of studying ℓ_2 risks is their close connection to techniques from RMT.

Recall the notation $R_{(\Sigma, \mu_0)}^\#(\eta)$ defined in Section 1.3. Let their ‘theoretical’ versions be defined as follows:

- $\bar{R}_{(\Sigma, \mu_0)}^{\text{pred}}(\eta) \equiv \tau_{\eta, *}^2 \left\| (\Sigma + \tau_{\eta, *} I_n)^{-1} \Sigma^{1/2} \mu_0 \right\|^2 + \frac{\gamma_{\eta, *}^2}{n} \text{tr} \left(\Sigma^2 (\Sigma + \tau_{\eta, *} I_n)^{-2} \right).$
- $\bar{R}_{(\Sigma, \mu_0)}^{\text{est}}(\eta) \equiv \tau_{\eta, *}^2 \left\| (\Sigma + \tau_{\eta, *} I_n)^{-1} \mu_0 \right\|^2 + \frac{\gamma_{\eta, *}^2}{n} \text{tr} \left(\Sigma (\Sigma + \tau_{\eta, *} I_n)^{-2} \right).$
- $\bar{R}_{(\Sigma, \mu_0)}^{\text{in}}(\eta) \equiv \left(\frac{\eta \gamma_{\eta, *}}{\tau_{\eta, *}} \right)^2 + \phi \sigma_\xi^2 \cdot \left(1 - \frac{2\eta}{\phi \tau_{\eta, *}} \right).$

We also define the residual and its theoretical version as

- $R_{(\Sigma, \mu_0)}^{\text{res}}(\eta) \equiv n^{-1} \|Y - X \widehat{\mu}_\eta\|^2, \bar{R}_{(\Sigma, \mu_0)}^{\text{res}}(\eta) \equiv \left(\frac{\eta \gamma_{\eta, *}}{\tau_{\eta, *}} \right)^2.$

From Theorems 3 and 4, it is natural to expect that for $\# \in \{\text{pred, est, in, res}\}$,

$$\sup_{\eta \in \Xi^\#} |R_{(\Sigma, \mu_0)}^\#(\eta) - \bar{R}_{(\Sigma, \mu_0)}^\#(\eta)| \approx 0 \text{ with high probability.} \quad (17)$$

A rigorous statement of (17) is deferred to Theorem 49; its proof and the proofs for all other results in this section can be found in Section F.

Using the so-called Stieltjes transformation $\mathbf{m}(\cdot)$ in the RMT literature (defined formally via (41) in Section A.3 ahead), the following theorem provides an efficient RMT representation of $\bar{R}_{(\Sigma, \mu_0)}^\#(\eta)$ that holds for ‘most’ μ_0 ’s. Recall Ξ_K from (12).

Theorem 8. *Suppose $1/K \leq \phi^{-1} \leq K$, $\sigma_\xi^2 \in [0, K]$ and $\|\Sigma\|_{\text{op}} \vee \mathcal{H}_\Sigma \leq K$ for some $K > 0$. There exists some constant $C = C(K) > 0$ such that for any $\varepsilon \in (0, 1/2]$, we may find a measurable set $\mathcal{U}_\varepsilon \subset B_n(1)$ with $\text{vol}(\mathcal{U}_\varepsilon)/\text{vol}(B_n(1)) \geq 1 - C\varepsilon^{-1}e^{-n\varepsilon^2/C}$,*

$$\sup_{\mu_0 \in \mathcal{U}_\varepsilon} \sup_{\eta \in \Xi_K} |\bar{R}_{(\Sigma, \mu_0)}^\#(\eta) - \mathcal{R}_{(\Sigma, \mu_0)}^\#(\eta)| \leq \varepsilon, \quad \# \in \{\text{pred, est, in}\}. \quad (18)$$

Here with $\text{SNR}_{\mu_0} = \|\mu_0\|^2/\sigma_\xi^2$, $\mathbf{m}_\eta \equiv \mathbf{m}(-\eta/\phi)$ and $\mathbf{m}'_\eta \equiv \mathbf{m}'(-\eta/\phi)$,

- $\mathcal{R}_{(\Sigma, \mu_0)}^{\text{pred}}(\eta) \equiv \sigma_\xi^2 \cdot \left\{ \frac{1}{\mathbf{m}'_\eta} \left(\phi \cdot \text{SNR}_{\mu_0} \mathbf{m}_\eta - (\eta \cdot \text{SNR}_{\mu_0} - 1) \mathbf{m}'_\eta \right) - 1 \right\}$,
- $\mathcal{R}_{(\Sigma, \mu_0)}^{\text{est}}(\eta) \equiv \sigma_\xi^2 \cdot \left\{ \text{SNR}_{\mu_0} (1 - \phi) + \mathbf{m}_\eta + \frac{\eta}{\phi} (\eta \cdot \text{SNR}_{\mu_0} - 1) \mathbf{m}'_\eta \right\}$,
- $\mathcal{R}_{(\Sigma, \mu_0)}^{\text{in}}(\eta) \equiv \sigma_\xi^2 \cdot \frac{\eta^2}{\phi} \left(\phi \cdot \text{SNR}_{\mu_0} \mathbf{m}_\eta - (\eta \cdot \text{SNR}_{\mu_0} - 1) \mathbf{m}'_\eta \right) + \sigma_\xi^2 \cdot (\phi - 2\eta \mathbf{m}_\eta)$,
- $\mathcal{R}_{(\Sigma, \mu_0)}^{\text{res}}(\eta) \equiv \sigma_\xi^2 \cdot \frac{\eta^2}{\phi} \left(\phi \cdot \text{SNR}_{\mu_0} \mathbf{m}_\eta - (\eta \cdot \text{SNR}_{\mu_0} - 1) \mathbf{m}'_\eta \right)$.

When $\Sigma = I_n$, we may take $\mathcal{U}_\varepsilon = B_n(1)$ and (18) holds with $\varepsilon = 0$.

The RMT representation above yields the following crucial insight into the extremal behavior of the risk maps $\eta \mapsto \bar{R}_{(\Sigma, \mu_0)}^\#(\eta)$.

Proposition 9. *Suppose $1/K \leq \phi^{-1} \leq K$ and $\|\Sigma\|_{\text{op}} \vee \mathcal{H}_\Sigma \leq K$ for some $K > 0$. Then there exists some $C = C(K) > 0$ such that for all $\# \in \{\text{pred, est, in}\}$, the derivative formulae*

$$\partial_\eta \mathcal{R}_{(\Sigma, \mu_0)}^\#(\eta) = \sigma_\xi^2 \cdot \mathfrak{M}^\#(\eta) \cdot (\eta \cdot \text{SNR}_{\mu_0} - 1), \quad \eta \geq 0$$

hold for some measurable functions $\{\mathfrak{M}^\# : \mathbb{R}_{\geq 0} \rightarrow [1/C, C]\}$.

Consequently, for all $\# \in \{\text{pred, est, in}\}$, $\mathcal{R}_{(\Sigma, \mu_0)}^\#(\cdot)$ attains its global minimum at the same $\eta_* \equiv \text{SNR}_{\mu_0}^{-1} \in \Xi_K$, and $1/C \leq |\mathcal{R}_{(\Sigma, \mu_0)}^\#(\eta) - \mathcal{R}_{(\Sigma, \mu_0)}^\#(\eta_*)|/\{\|\mu_0\|^2(\eta - \eta_*)^2\} \leq C$.

A more general version of the above proposition with precise formulae for $\mathfrak{M}^\#$ can be found in Proposition 53. As the maps $\eta \mapsto \mathcal{R}_{(\Sigma, \mu_0)}^\#(\eta)$ are almost quadratic with the same global minimizer $\eta_* = \text{SNR}_{\mu_0}^{-1}$ for all $\# \in \{\text{pred, est, in}\}$, in view of (17) and Theorem 8, it is natural to expect that for ‘most’ signal μ_0 ’s and all $\# \in \{\text{pred, est, in}\}$,

$$R_{(\Sigma, \mu_0)}^\#(\eta_*) \approx \min_{\eta \in \Xi_L} R_{(\Sigma, \mu_0)}^\#(\eta) \text{ with high probability.} \quad (19)$$

A rigorous formulation of (19) is given in Theorem 54, which, along with its proof, is provided in Section F.4.

4. Cross-validation: optimality beyond prediction

This section is devoted to the validation of the broad optimality of two widely used cross-validation schemes beyond the prediction risk. Some consequences to statistical inference via debiased Ridge(less) estimators will also be discussed.

4.1 Estimation of effective noise and regularization

We shall first take a slight detour, by considering estimation of the effective regularization $\tau_{\eta,*}$ and the effective noise $\gamma_{\eta,*}$. We propose the following estimators:

$$\begin{cases} \widehat{\tau}_\eta \equiv \left\{ \frac{1}{m} \operatorname{tr} \left(\frac{1}{m} X X^\top + \frac{\eta}{\phi} I_m \right)^{-1} \right\}^{-1} = \left\{ \operatorname{tr} (X X^\top + \eta \cdot n I_m)^{-1} \right\}^{-1}, \\ \widehat{\gamma}_\eta \equiv \frac{\widehat{\tau}_\eta}{\sqrt{n}} \left(\eta^{-1} \|Y - X \widehat{\mu}_\eta\| \mathbf{1}_{\phi^{-1} < 1} + \|(X X^\top / n)^{-1} X \widehat{\mu}_\eta\| \mathbf{1}_{\phi^{-1} \geq 1} \right). \end{cases} \quad (20)$$

It can be easily shown that

$$\sup_{\eta \in \Xi_K} |\widehat{\tau}_\eta - \tau_{\eta,*}|, \quad \sup_{\eta \in \Xi_K} |\widehat{\gamma}_\eta - \gamma_{\eta,*}| \approx 0 \text{ with high probability.} \quad (21)$$

A rigorous statement of (21) is deferred to Theorem 57; its proof and the proofs for all other results in this section can be found in Section G. These estimators will not only be useful in their own rights, they will also play an important rule in understanding the generalized cross-validation scheme in the next subsection.

4.2 Validation of generalized cross-validation

Consider choosing η by minimizing the estimated effective noise $\widehat{\gamma}_\eta$ given in (20): for any $L > 0$,

$$\widehat{\eta}_L^{\text{GCV}} \in \arg \min_{\eta \in \Xi_L} \widehat{\gamma}_\eta. \quad (22)$$

Here we recall Ξ_K from (12). The subscript on L in $\widehat{\eta}_L^{\text{GCV}}$ will usually be suppressed for notational simplicity.

The proposal (22) is known in the literature as the *generalized cross validation* (Craven and Wahba, 1978/79; Golub et al., 1979), and is strongly tied to the so-called shortcut formula for leave-one-out cross validation that exists uniquely for Ridge regression, cf. (Hastie et al., 2022, Eqn. (46)). Here we take a different perspective on (22). From our developed theory, this tuning scheme is easily believed to “work” since

$$\widehat{\gamma}_\eta^2 \stackrel{\mathbb{P}}{\approx} \gamma_{\eta,*}^2 = \phi^{-1} (\sigma_\xi^2 + \bar{R}_{(\Sigma, \mu_0)}^{\text{pred}}(\eta)) \stackrel{\mathbb{P}}{\approx} \phi^{-1} (\sigma_\xi^2 + R_{(\Sigma, \mu_0)}^{\text{pred}}(\eta)). \quad (23)$$

So minimization of $\eta \mapsto \widehat{\gamma}_\eta$ is approximately the same as that of $\eta \mapsto R_{(\Sigma, \mu_0)}^{\text{pred}}(\eta)$, and therefore simultaneously of $\eta \mapsto R_{(\Sigma, \mu_0)}^\#(\eta)$ for $\# \in \{\text{est}, \text{in}\}$ as per (19). We make precise the foregoing heuristics in the following theorem.

Theorem 10. *Suppose Assumption A holds, and the following hold some $K > 0$.*

- $1/K \leq \phi^{-1} \leq K, \|\Sigma^{-1}\|_{\text{op}} \vee \|\Sigma\|_{\text{op}} \leq K$.

- Assumption B holds with either (i) $\sigma_\xi^2 \in [1/K, K]$ or (ii) $\sigma_\xi^2 \in [0, K]$ with $\phi^{-1} \geq 1 + 1/K$.

Fix $\delta \in (0, 1/2]$, $L \geq K/\delta^2$ and a small enough $\vartheta \in (0, 1/50)$. There exist a constant $C = C(K, L, \delta, \vartheta) > 0$ and a measurable set $\mathcal{U}_{\delta, \vartheta} \subset B_n(1) \setminus B_n(\delta)$ with $\text{vol}(\mathcal{U}_{\delta, \vartheta})/\text{vol}(B_n(1) \setminus B_n(\delta)) \geq 1 - Ce^{-n^\vartheta/C}$, such that for $\# \in \{\text{pred, est, in}\}$,

$$\sup_{\mu_0 \in \mathcal{U}_{\delta, \vartheta}} \mathbb{P} \left(R_{(\Sigma, \mu_0)}^\#(\widehat{\eta}_L^{\text{GCV}}) \geq \min_{\eta \in \Xi_L} R_{(\Sigma, \mu_0)}^\#(\eta) + n^{-\vartheta} \right) \leq Cn^{-1/7}.$$

Remark 11. Formally, the set $\mathcal{U}_{\delta, \vartheta}$ is defined as $\mathcal{U}_{\delta, \vartheta} \equiv \mathcal{U}_\vartheta \setminus B_n(\delta)$, where \mathcal{U}_ϑ is defined in Proposition 40. The cutoff $\|\mu_0\| \geq \delta$ excludes vanishing signals and ensures that $\eta_* = \text{SNR}_{\mu_0}^{-1}$ is uniformly bounded, so that $\eta_* \in \Xi_L$ whenever $L \geq K/\delta^2$.

Earlier low-dimension results for generalized cross validation in Ridge regression include Stone (1974, 1977); Craven and Wahba (1978/79); Li (1985, 1986, 1987); Dudoit and van der Laan (2005). In the proportional high-dimensional regime, Hastie et al. (2022); Patil et al. (2021) validate the optimality of $\widehat{\eta}^{\text{GCV}}$ with respect to the prediction risk $R_{(\Sigma, \mu_0)}^{\text{pred}}$ with increasing generality. In Theorem 10 above, we prove that the optimality of $\widehat{\eta}^{\text{GCV}}$ holds *simultaneously* for all the three indicated risks. To the best of our knowledge, such optimality of $\widehat{\eta}^{\text{GCV}}$ beyond the prediction risk has not been previously observed in the literature.

4.3 Validation of k -fold cross-validation

Next we consider the widely used k -fold cross-validation. We need some further notation:

- Let m_ℓ be the sample size of batch $\ell \in [k]$, so $\sum_{\ell \in [k]} m_\ell = m$. In the standard k -fold cross validation, we choose equal sized batch with $m_\ell = m/k$ (assumed to be integer without loss of generality).
- Let $X^{(\ell)} \in \mathbb{R}^{m_\ell \times n}$ (resp. $Y^{(\ell)} \in \mathbb{R}^{m_\ell}$) be the submatrix of X (resp. subvector of Y) that contains all rows corresponding to the training data in batch ℓ .
- In a similar fashion, let $X^{(-\ell)} \in \mathbb{R}^{(m-m_\ell) \times n}$ (resp. $Y^{(-\ell)} \in \mathbb{R}^{m-m_\ell}$) be the submatrix of X (resp. subvector of Y) that removes all rows corresponding to $X^{(\ell)}$ (resp. $Y^{(\ell)}$).

The k -fold cross-validation works as follows. For $\ell \in [k]$, let $\widehat{\mu}_\eta^{(\ell)} \equiv \arg \min_{\mu \in \mathbb{R}^n} \left\{ \frac{1}{2n} \|Y^{(-\ell)} - X^{(-\ell)}\mu\|^2 + \frac{\eta}{2} \|\mu\|^2 \right\}$ be the Ridge estimator over $(X^{(-\ell)}, Y^{(-\ell)})$ with regularization $\eta \geq 0$. We then pick the tuning parameter that minimizes the averaged test errors of $\widehat{\mu}_\eta^{(\ell)}$ over $(X^{(\ell)}, Y^{(\ell)})$: for any $L > 0$,

$$\widehat{\eta}_L^{\text{CV}} \in \arg \min_{\eta \in \Xi_L} \left\{ \frac{1}{k} \sum_{\ell \in [k]} \frac{1}{m_\ell} \|Y^{(\ell)} - X^{(\ell)}\widehat{\mu}_\eta^{(\ell)}\|^2 \right\} \equiv \arg \min_{\eta \in \Xi_L} R_{(\Sigma, \mu_0)}^{\text{CV}, k}(\eta). \quad (24)$$

We shall often omit the subscript L in $\widehat{\eta}_L^{\text{CV}}$.

Intuitively, due to the independence between $\widehat{\mu}_\eta^{(\ell)}$ and $(X^{(\ell)}, Y^{(\ell)})$, $R_{(\Sigma, \mu_0)}^{\text{CV}, k}(\eta)$ can be viewed as an estimator of the generalization error $R_{(\Sigma, \mu_0)}^{\text{pred}}(\eta) + \sigma_\xi^2$. So it is natural to expect

that $\hat{\eta}^{\text{CV}}$ approximately minimizes $\eta \mapsto R_{(\Sigma, \mu_0)}^{\text{pred}}(\eta)$. Based on the same heuristics as for $\hat{\eta}^{\text{GCV}}$ in (22), we may therefore expect that $\hat{\eta}^{\text{CV}}$ in (24) simultaneously provides optimal prediction, estimation and in-sample risks for ‘most’ signal μ_0 ’s. This is the content of the following theorem.

Theorem 12. *Suppose the same conditions as in Theorem 10 and $\max_{\ell \in [k]} m_\ell/n \leq 1/(2K)$ hold for some $K > 0$. Fix $\delta \in (0, 1/2]$, $L \geq K/\delta^2$ and a small enough $\vartheta \in (0, 1/50)$. Further assume $\min_{\ell \in [k]} m_\ell \geq \log^{2/\delta} m$. There exist a constant $C = C(K, L, \delta, \vartheta) > 0$ and a measurable set $\mathcal{U}_{\delta, \vartheta} \subset B_n(1) \setminus B_n(\delta)$ with $\text{vol}(\mathcal{U}_{\delta, \vartheta})/\text{vol}(B_n(1) \setminus B_n(\delta)) \geq 1 - Ce^{-n^\vartheta/C}$, such that for $\# \in \{\text{pred, est, in}\}$,*

$$\begin{aligned} \sup_{\mu_0 \in \mathcal{U}_{\delta, \vartheta}} \mathbb{P} \left(R_{(\Sigma, \mu_0)}^\#(\hat{\eta}_L^{\text{CV}}) \geq \min_{\eta \in \Xi_L} R_{(\Sigma, \mu_0)}^\#(\eta) + C \cdot \left\{ \frac{1}{k} \sum_{\ell \in [k]} \frac{1}{m_\ell^{(1-\delta)/2}} + \frac{1}{k} + n^{-\vartheta} \right\} \right) \\ \leq C(1 + \mathcal{L}_{\{m_\ell\}}) \cdot n^{-1/7}. \end{aligned}$$

Here $\mathcal{L}_{\{m_\ell\}} \equiv \sum_{\ell \in [k]} (m_\ell/m)^{-1}$.

Non-asymptotic results of this type for k -fold cross validation are previously obtained for $R_{(\Sigma, \mu_0)}^{\text{pred}}(\hat{\eta}^{\text{CV}})$ in the Lasso setting (Miolane and Montanari, 2021, Proposition 4.3) under isotropic $\Sigma = I_n$, where the range of the regularization must be strictly away from the interpolation regime. In contrast, our results above are valid down to $\eta = 0$ when $\phi^{-1} > 1$, and allow for general anisotropic Σ .

Interestingly, the error bound in the above theorem reflects the folklore tension between the bias and variance in the selection of k in the cross validation scheme (cf. (James et al., 2021, Chapter 5)):

- For a small number of k , $R_{(\Sigma, \mu_0)}^{\text{CV}, k}(\eta)$ is biased for estimating $R_{(\Sigma, \mu_0)}^{\text{pred}}(\eta)$; this corresponds to the term $\mathcal{O}(1/k)$ in the error bound, which is known to be of the optimal order in Ridge regression (cf. Liu and Dobriban (2019)).
- For a large number of k , $R_{(\Sigma, \mu_0)}^{\text{CV}, k}(\eta)$ has large fluctuations; this corresponds to the term $\mathcal{O}(k^{-1} \sum_{\ell \in [k]} m_\ell^{-(1-\delta)/2}) = \mathcal{O}((k/m)^{(1-\delta)/2})$ in the equal-sized case. By a central limit heuristic (cf. Kissel and Lei (2022); Austern and Zhou (2025)), we also expect this term to be of a near optimal order.

For the choice of k , it is instructive to consider the common equal-sized-folds case $m_\ell = m/k$. In this case, the error bound in Theorem 12 suggests that the optimal theoretical value of k is $k \sim m^{1/3}$ (when δ is small). In our numerical experiments, choosing $k = 5$ already yields cross-validation performance that is close to the theoretically optimal behavior (cf. Figure 1), whereas the theoretically prescribed optimal value of k offers limited practical gain.

4.4 Implications to statistical inference via $\hat{\mu}_\eta$

As Ridge(less) estimators $\hat{\mu}_\eta$ are in general biased, debiasing is necessary for statistical inference of μ_0 , cf. Bellec and Zhang (2023). Here the debiasing scheme for $\hat{\mu}_\eta$ can be

readily read off from the distributional characterizations in Theorems 3 and 4. Assuming known covariance Σ , let the debiased Ridge(less) estimator be defined as

$$\hat{\mu}_\eta^{\text{dR}} \equiv (\Sigma + \tau_{\eta,*} I) \Sigma^{-1} \hat{\mu}_\eta. \quad (25)$$

Note that $\tau_{\eta,*}$ and $\hat{\tau}_\eta$ is interchangeable in the above display due to known Σ . Using Theorems 3 and 4, we expect that $\hat{\mu}_\eta^{\text{dR}} \stackrel{d}{\approx} \mu_0 + \gamma_{\eta,*} \Sigma^{-1/2} g / \sqrt{n}$. This motivates the following confidence intervals for $\{\mu_{0,j}\}$:

$$\text{CI}_j(\eta) \equiv \left[\hat{\mu}_{\eta,j}^{\text{dR}} \pm \hat{\gamma}_\eta \cdot (\Sigma^{-1})_{jj}^{1/2} \cdot \frac{z_{\alpha/2}}{\sqrt{n}} \right], \quad j \in [n]. \quad (26)$$

Here z_α is the normal upper- α quantile defined via $\mathbb{P}(\mathcal{N}(0,1) > z_\alpha) = \alpha$. It is easy to see from the above definition that minimization of $\eta \mapsto \hat{\gamma}_\eta$ is equivalent to that of the CI length. As the former minimization procedure corresponds exactly to the proposal $\hat{\eta}^{\text{GCV}}$ in (22), we expect that $\{\text{CI}_j(\hat{\eta}^{\text{GCV}})\}$ provide the shortest (asymptotic) $(1 - \alpha)$ -CIs along the regularization path, and so do $\{\text{CI}_j(\hat{\eta}^{\text{CV}})\}$.

Below we give a rigorous statement on the above informal discussion. Let $\mathcal{C}^{\text{dR}}(\eta) \equiv n^{-1} \sum_{j=1}^n \mathbf{1}(\mu_{0,j} \in \text{CI}_j(\eta))$ denote the averaged coverage of $\{\text{CI}_j(\eta)\}$ for $\{\mu_{0,j}\}$. We have the following.

Theorem 13. *Suppose the same conditions as in Theorem 10 (resp. Theorem 12) for $\hat{\eta}^{\text{GCV}}$ (resp. $\hat{\eta}^{\text{CV}}$) hold for some $K > 0$. Fix $\alpha \in (0, 1/4]$, $\delta \in (0, 1/2]$, $L \geq K/\delta^2$ and a small enough $\vartheta \in (0, 1/50)$. There exist a constant $C = C(K, L, \delta, \vartheta) > 0$ and a measurable set $\mathcal{U}_{\delta,\vartheta} \subset B_n(1) \setminus B_n(\delta)$ with $\text{vol}(\mathcal{U}_{\delta,\vartheta})/\text{vol}(B_n(1) \setminus B_n(\delta)) \geq 1 - Ce^{-n^\vartheta/C}$, such that the CI length and the averaged coverage satisfy*

$$\begin{aligned} & \sup_{\mu_0 \in \mathcal{U}_{\delta,\vartheta}} \left\{ \mathbb{P} \left(\sqrt{n} z_{\alpha/2}^{-1} \cdot \max_{j \in [n]} |\text{CI}_j(\hat{\eta}_L^\#)| - \min_{\eta \in \Xi_L} |\text{CI}_j(\eta)| \geq C \mathcal{E}_n^\# \right) \right. \\ & \quad \left. \vee \mathbb{P} \left(|\mathcal{C}^{\text{dR}}(\hat{\eta}_L^\#) - (1 - \alpha)| \geq C (\mathcal{E}_n^\#)^{1/4} \right) \right\} \leq C \mathfrak{p}_n^\#. \end{aligned}$$

Here for $\# \in \{\text{GCV}, \text{CV}\}$, the quantities $\mathcal{E}_n^\#, \mathfrak{p}_n^\#$ are defined via

	$\mathcal{E}_n^\#$	$\mathfrak{p}_n^\#$
$\# = \text{GCV}$	$n^{-\vartheta}$	$n^{-1/7}$
$\# = \text{CV}$	$k^{-1} \sum_{\ell \in [k]} m_\ell^{-(1-\delta)/2} + k^{-1} + n^{-\vartheta}$	$(1 + \mathcal{L}_{\{m_\ell\}}) \cdot n^{-1/7}$

A somewhat non-standard special case of the above theorem is the noiseless setting $\sigma_\xi^2 = 0$ in the overparametrized regime $\phi^{-1} > 1$. In this case, exact recovery of μ_0 is impossible and our CI's above provide a precise scheme for partial recovery of μ_0 . As the effective noise $\phi \gamma_{\eta,*}^2(0) = \bar{R}_{(\Sigma, \mu_0)}^{\text{pred}}(\eta)$, Theorem 8 and Proposition 9 suggest that $\eta \mapsto \gamma_{\eta,*}^2(0)$ is approximately minimized at $\eta = 0$ for ‘most’ μ_0 's. This means that, in this noiseless case, the length of the adaptively tuned CIs is also approximately minimized at the interpolation regime for ‘most’ μ_0 's.

5. Proof outlines

5.1 Technical tools

The main technical tool we use for the proof of Theorem 3 is the following version of convex Gaussian min-max theorem, taken from (Miolane and Montanari, 2021, Corollary G.1).

Theorem 14 (Convex Gaussian Min-Max Theorem). *Suppose $D_u \in \mathbb{R}^{n_1+n_2}$, $D_v \in \mathbb{R}^{m_1+m_2}$ are compact sets, and $Q : D_u \times D_v \rightarrow \mathbb{R}$ is continuous. Let $G = (G_{ij})_{i \in [n_1], j \in [m_1]}$ with G_{ij} 's i.i.d. $\mathcal{N}(0, 1)$, and $g \sim \mathcal{N}(0, I_{n_1})$, $h \sim \mathcal{N}(0, I_{m_1})$ be independent Gaussian vectors. For $u \in \mathbb{R}^{n_1+n_2}$, $v \in \mathbb{R}^{m_1+m_2}$, write $u_1 \equiv u_{[n_1]} \in \mathbb{R}^{n_1}$, $v_1 \equiv v_{[m_1]} \in \mathbb{R}^{m_1}$. Define*

$$\begin{aligned}\Phi^p(G) &= \min_{u \in D_u} \max_{v \in D_v} \left(u_1^\top G v_1 + Q(u, v) \right), \\ \Phi^a(g, h) &= \min_{u \in D_u} \max_{v \in D_v} \left(\|v_1\| g^\top u_1 + \|u_1\| h^\top v_1 + Q(u, v) \right).\end{aligned}$$

Then the following hold.

1. For all $t \in \mathbb{R}$, $\mathbb{P}(\Phi^p(G) \leq t) \leq 2\mathbb{P}(\Phi^a(g, h) \leq t)$.
2. If $(u, v) \mapsto u_1^\top G v_1 + Q(u, v)$ satisfies the conditions of Sion's min-max theorem for the pair (D_u, D_v) a.s. (for instance, D_u, D_v are convex, and Q is convex-concave), then for any $t \in \mathbb{R}$, $\mathbb{P}(\Phi^p(G) \geq t) \leq 2\mathbb{P}(\Phi^a(g, h) \geq t)$.

Clearly, \geq (resp. \leq) in (1) (resp. (2)) can be replaced with $>$ (resp. $<$). In the proofs below, we shall assume without loss of generality that G, g, h are independent Gaussian matrix/vectors defined on the same probability space.

As mentioned above, the CGMT above has been utilized for deriving precise risk/distributional asymptotics for a number of canonical statistical estimators across various important models; we only refer the readers to Thrampoulidis et al. (2015, 2018); Salehi et al. (2019); Loureiro et al. (2021); Deng et al. (2022); Celentano et al. (2023); Han (2023); Liang and Sur (2022); Wang et al. (2022); Zhang et al. (2022); Montanari et al. (2025) for some selected references.

5.2 Reparametrization and further notation

Consider the reparametrization

$$w = \Sigma^{1/2}(\mu - \mu_0), \quad \hat{w}_{\eta; Z} \equiv \Sigma^{1/2}(\hat{\mu}_{\eta; Z} - \mu_0).$$

Then with

$$F(w) \equiv F_{(\Sigma, \mu_0)}(w) = \frac{1}{2} \|\mu_0 + \Sigma^{-1/2} w\|^2, \tag{27}$$

we have the following reparametrized version of $\hat{\mu}_{\eta; Z}$:

$$\hat{w}_{\eta; Z} = \begin{cases} \arg \min_{w \in \mathbb{R}^n} \{F(w) : Zw = \xi\}, & \eta = 0; \\ \arg \min_{w \in \mathbb{R}^n} \left\{ F(w) + \frac{1}{\eta} \cdot \frac{1}{2n} \|Zw - \xi\|^2 \right\}, & \eta > 0. \end{cases}$$

Next we give some further notation for cost functions. Let for $\eta \geq 0$,

$$\begin{aligned} h_{\eta;Z}(w, v) &\equiv \frac{1}{\sqrt{n}} \langle v, Zw - \xi \rangle + F(w) - \frac{\eta \|v\|^2}{2}, \\ \ell_{\eta}(w, v) &\equiv \frac{1}{\sqrt{n}} \left(-\|v\| \langle g, w \rangle + \|w\| \langle h, v \rangle - \langle v, \xi \rangle \right) + F(w) - \frac{\eta \|v\|^2}{2}, \end{aligned} \quad (28)$$

and for $L_v \in [0, \infty]$,

$$\begin{aligned} H_{\eta;Z}(w; L_v) &\equiv \max_{v \in B_n(L_v)} h_{\eta;Z}(w, v) \equiv \max_{v \in B_n(L_v)} \left\{ \frac{\langle v, Zw - \xi \rangle}{\sqrt{n}} + F(w) - \frac{\eta \|v\|^2}{2} \right\}, \\ L_{\eta}(w; L_v) &\equiv \max_{v \in B_n(L_v)} \ell_{\eta}(w, v) = \max_{\beta \in [0, L_v]} \left\{ \frac{\beta}{\sqrt{n}} \left(\|\|w\|h - \xi\| - \langle g, w \rangle \right) + F(w) - \frac{\eta \beta^2}{2} \right\}. \end{aligned} \quad (29)$$

We shall simply write $H_{\eta;Z}(\cdot) = H_{\eta;Z}(\cdot; \infty)$ and $L_{\eta}(\cdot) = L_{\eta}(\cdot; \infty)$. When $Z = G$, we sometimes write $h_{\eta;G} = h_{\eta}$ and $H_{\eta;G} = H_{\eta}$ for simplicity of notation.

Let the empirical noise σ_m^2 and its modified version be

$$\sigma_m^2 \equiv \frac{\|\xi\|^2}{\|h\|^2}, \quad \sigma_{\pm}^2(L_w) \equiv \left(\sigma_m^2 \pm 2L_w \frac{|\langle h, \xi \rangle|}{\|h\|^2} \right)_+. \quad (30)$$

Finally we define $D_{\eta, \pm}$ and its deterministic version \bar{D}_{η} as follows:

$$\begin{aligned} D_{\eta, \pm}(\beta, \gamma) &\equiv \frac{\beta}{2} \left(\gamma(\phi e_h^2 - e_g^2) + \frac{\sigma_{\pm}^2}{\gamma} \right) - \frac{\eta \beta^2}{2} + \mathbf{e}_F \left(\frac{\gamma}{\sqrt{n}} g; \frac{\gamma}{\beta} \right), \\ \bar{D}_{\eta}(\beta, \gamma) &\equiv \frac{\beta}{2} \left(\gamma(\phi - 1) + \frac{\sigma_{\xi}^2}{\gamma} \right) - \frac{\eta \beta^2}{2} + \mathbb{E} \mathbf{e}_F \left(\frac{\gamma}{\sqrt{n}} g; \frac{\gamma}{\beta} \right). \end{aligned} \quad (31)$$

Here recall \mathbf{e}_F is the Moreau envelope of F in (27). Note that $D_{\eta, \pm}$ depends on the choice of L_w , but for notational convenience we drop this dependence here.

5.3 Proof outline for Theorem 3 for $\eta = 0$

We shall outline below the main steps for the proof of Theorem 3 for $\eta = 0$ in the regime $\phi^{-1} > 1$ under a stronger condition $\|\Sigma^{-1}\|_{\text{op}} \lesssim 1$. The high level strategy of the proof shares conceptual similarities to Miolane and Montanari (2021); Celentano et al. (2023), but the details differ significantly.

(Step 1: Localization of the primal optimization). In this step, we show that for $L_w, L_v > 0$ such that $L_w \wedge L_v \gtrsim 1$, with high probability (w.h.p.),

$$\min_{w \in B_n(L_w)} H_0(w; L_v) = \min_{w \in \mathbb{R}^n} H_0(w). \quad (32)$$

A formal statement of the above localization can be found in Proposition 27. The key point here is that despite $\min_w H_0(w)$ optimizes a deterministic function with a random constraint, it can be efficiently rewritten (in a probabilistic sense) in a minimax form indexed by *compact sets* that facilitate the application of the convex Gaussian min-max Theorem 14.

(Step 2: Characterization of the Gordon cost optimum). In this step, we show that a suitably localized version of $\min_w L_0(w)$ concentrates around some *deterministic* quantity involving the function \bar{D}_0 in (31). In particular, we show in Theorem 28 that for $L_w, L_v \asymp 1$ chosen large enough, w.h.p.,

$$\min_{w \in B_n(L_w)} L_0(w; L_v) \approx \max_{\beta > 0} \min_{\gamma > 0} \bar{D}_0(\beta, \gamma). \quad (33)$$

The proof of (33) is fairly involved, as the minimax problem $\min_w L_0(w) = \min_w \max_v \ell_0(w, v)$ (and its suitably localized versions) cannot be computed exactly. We get around this technical issue by the following bracketing strategy:

- **(Step 2.1).** We show in Proposition 29 that for the prescribed choice of L_w, L_v , w.h.p., both

$$\max_{\beta > 0} \min_{\gamma > 0} D_{0,-}(\beta, \gamma) \leq \min_{w \in B_n(L_w)} L_0(w; L_v) \leq \max_{\beta > 0} \min_{\gamma > 0} D_{0,+}(\beta, \gamma),$$

and the localization

$$\max_{\beta > 0} \min_{\gamma > 0} D_{0,\pm}(\beta, \gamma) = \max_{1/C \leq \beta \leq C} \min_{1/C \leq \gamma \leq C} D_{0,\pm}(\beta, \gamma)$$

hold for some large $C > 0$.

- **(Step 2.2).** We show in Proposition 30 that for localized minimax problems, we may replace $D_{0,\pm}$ by \bar{D}_0 : w.h.p.,

$$\max_{1/C \leq \beta \leq C} \min_{1/C \leq \gamma \leq C} D_{0,\pm}(\beta, \gamma) \approx \max_{1/C \leq \beta \leq C} \min_{1/C \leq \gamma \leq C} \bar{D}_0(\beta, \gamma).$$

- **(Step 2.3).** We show in Proposition 31 that (de)localization holds for the (deterministic) max-min optimization problem with \bar{D}_0 :

$$\max_{\beta > 0} \min_{\gamma > 0} \bar{D}_0(\beta, \gamma) = \max_{1/C \leq \beta \leq C} \min_{1/C \leq \gamma \leq C} \bar{D}_0(\beta, \gamma).$$

Combining the above Steps 2.1-2.3 yields (33). An important step to prove the (de)localization claims above is to derive a priori estimates for the solutions of the fixed point equation (13) and its sample version, to be defined in (55). These estimates will be detailed in Section B.

(Step 3: Locating the global minimizer of the Gordon objective). In this step, we show that a suitably localized version of the Gordon objective $w \mapsto L_0(w)$ attains its global minimum approximately at $w_{0,*} \equiv \Sigma^{1/2}(\hat{\mu}_{(\Sigma, \mu_0)}^{\text{seq}}(\gamma_{0,*}; \tau_{0,*}) - \mu_0)$ in the following sense. For any $\varepsilon > 0$ and any $g : \mathbb{R}^n \rightarrow \mathbb{R}$ that is 1-Lipschitz with respect to $\|\cdot\|_{\Sigma^{-1}}$, let $D_{0;\varepsilon}(g) \equiv \{w \in \mathbb{R}^n : |g(w) - \mathbb{E}g(w_{0,*})| \geq \varepsilon\}$ be the ‘exceptional set’. We show in Theorem 32 that again for $L_w, L_v \asymp 1$ chosen large enough, w.h.p.,

$$\min_{w \in D_{0;\varepsilon}(g) \cap B_n(L_w)} L_0(w; L_v) \geq \max_{\beta > 0} \min_{\gamma > 0} \bar{D}_0(\beta, \gamma) + \Omega_\varepsilon(1). \quad (34)$$

The main challenge in proving (34) is partly attributed to the possible violation of strong convexity of the map $w \mapsto L_0(w; L_v)$, due to the necessity of working with non-Gaussian ξ 's. We will get around this technical issue in similar spirit to Step 2 by another bracketing strategy. In particular:

- (**Step 3.1**). In Lemma 33, we will use surrogate, strongly convex functions $L_{0,\pm}(\cdot; L_v)$, formally defined in (79), to provide a sufficiently tight bracket for $L_0(\cdot; L_v)$ over large enough compact sets.
- (**Step 3.2**). In Proposition 34, we show that the minimizers of $w \mapsto L_{0,\pm}(\cdot; L_v)$ can be computed exactly and are close enough to $w_{0,*}$.
- (**Step 3.3**). In Proposition 35, combined with the tight bracketing and certain apriori estimates, we then conclude that all minimizers of $w \mapsto L_0(\cdot; L_v)$ must be close to $w_{0,*}$.

With all the above steps, finally we prove (34) by (i) using the proximity of L_0 and its surrogate $L_{0,\pm}$ and (ii) exploiting the strong convexity of $L_{0,\pm}$.

(**Step 4: Putting pieces together and establishing uniform guarantees**). In this final step, we shall use the convex Gaussian min-max theorem to translate the estimates (33) in Step 2 and (34) in Step 3 to their counterparts with primal cost function H_0 . For the global cost optimum, with the help of the localization in (32), by choosing $L_w, L_v \asymp 1$, we have w.h.p.,

$$\min_{w \in \mathbb{R}^n} H_0(w) \stackrel{(32)}{=} \min_{w \in B_n(L_w)} H_0(w; L_v) \stackrel{\mathbb{P}}{\approx} \min_{w \in B_n(L_w)} L_0(w; L_v) \stackrel{(33)}{\approx} \max_{\beta > 0} \min_{\gamma > 0} \bar{D}_0(\beta, \gamma).$$

For the cost over the exceptional set, we have w.h.p.,

$$\begin{aligned} \min_{w \in D_{0;\varepsilon}(\mathbf{g}) \cap B_n(L_w)} H_0(w) &\geq \min_{w \in D_{0;\varepsilon}(\mathbf{g}) \cap B_n(L_w)} H_0(w; L_v) \\ &\stackrel{\mathbb{P}}{\geq} \min_{w \in D_{0;\varepsilon}(\mathbf{g}) \cap B_n(L_w)} L_0(w; L_v) \stackrel{(34)}{\geq} \max_{\beta > 0} \min_{\gamma > 0} \bar{D}_0(\beta, \gamma) + \Omega_\varepsilon(1). \end{aligned}$$

Combining the above two displays, we then conclude that w.h.p., $\hat{w}_0 \notin D_{0;\varepsilon}(\mathbf{g}) \cap B_n(L_w)$. Finally using apriori estimate on $\|\hat{w}_0\|$ we may conclude that w.h.p., $\hat{w}_0 \notin D_{0;\varepsilon}(\mathbf{g})$, i.e., $|\mathbf{g}(\hat{w}_0) - \mathbb{E} \mathbf{g}(w_{0,*})| \leq \varepsilon$.

The uniform guarantee in η is then proved by (i) extending the above arguments to include any positive $\eta > 0$, and (ii) establishing (high probability) Lipschitz continuity (w.r.t. $\|\cdot\|_{\Sigma^{-1}}$) of the maps $\eta \mapsto \hat{w}_\eta$ and $\eta \mapsto w_{\eta,*}$.

Details of the above outline are implemented in Section C.

5.4 Proof outline for Theorem 4 for $\eta = 0$

The main tool we will use to prove the universality Theorem 4 is the following set of comparison inequalities developed in Han and Shen (2023): Suppose Z matches the first two moments of G , and possesses enough high moments. Then for any measurable sets $\mathcal{S}_w \subset [-L_n/\sqrt{n}, L_n/\sqrt{n}]^n$, $\mathcal{S}_v \subset [-L_n/\sqrt{n}, L_n/\sqrt{n}]^m$, and any smooth test function $\mathbb{T} : \mathbb{R} \rightarrow \mathbb{R}$ (standardized with derivatives of order 1 in $\|\cdot\|_\infty$),

$$\begin{aligned} \left| \mathbb{E} \mathbb{T} \left(\min_{w \in \mathcal{S}_w} \max_{v \in \mathcal{S}_v} h_{\eta;Z}(w, v) \right) - \mathbb{E} \mathbb{T} \left(\min_{w \in \mathcal{S}_w} \max_{v \in \mathcal{S}_v} h_{\eta;G}(w, v) \right) \right| &\leq r_n(L_n), \\ \left| \mathbb{E} \mathbb{T} \left(\min_{w \in \mathcal{S}_w} H_{\eta;Z}(w) \right) - \mathbb{E} \mathbb{T} \left(\min_{w \in \mathcal{S}_w} H_{\eta;G}(w) \right) \right| &\leq r_n(L_n). \end{aligned} \quad (35)$$

Here $r_n(L_n) \rightarrow 0$ for $L_n = n^\vartheta$ with sufficiently small $\vartheta > 0$. The readers are referred to Theorems 38 and 39 for a precise statement of (35).

An important technical subtlety here is that while the first inequality in (35) holds down to $\eta = 0$, the second inequality does not. This is so because $\min_w H_{0;Z}(w)$, which minimizes a deterministic function under a random constraint due to the unbounded constraint in the maximization of v , is qualitatively different from $\min_w H_{\eta;Z}(w)$ for any $\eta > 0$.

Now we shall sketch how the comparison inequalities (35) lead to universality.

(Step 1: Universality of the global cost optimum). In this step, we shall use the first inequality in (35) to establish the universality of the global Gordon cost:

$$\min_{w \in \mathbb{R}^n} H_{0;Z}(w) = \min_{w \in \mathbb{R}^n} \max_{v \in \mathbb{R}^m} h_{0;Z}(w, v) \stackrel{\mathbb{P}}{\approx} \min_{w \in \mathbb{R}^n} \max_{v \in \mathbb{R}^m} h_{0;G}(w, v). \quad (36)$$

See Theorem 43 for a formal statement of (36).

The crux to establish (36) via the first inequality of (35) is to show that, the ranges of the minimum and the maximum of $\min_w \max_v h_{0;Z}(w, v)$ can be localized into an L_∞ ball of order close to $\mathcal{O}(1/\sqrt{n})$. This amounts to showing that the stationary points $(\widehat{w}_{0;Z}, \widehat{v}_{0;Z})$, where $\widehat{w}_{\eta;Z} = \Sigma^{1/2}(\widehat{\mu}_{\eta;Z} - \mu_0)$ and $\widehat{v}_{\eta;Z} = -n^{-1/2}(XX^\top/n + \eta I_m)^{-1}Y$ (cf. Eqn. (118)), are delocalized. We prove such delocalization properties in Proposition 40 for ‘most’ $\mu_0 \in B_n(1)$.

(Step 2: Universality of the cost over exceptional sets). In this step, we shall use the second inequality in (35) to establish the universality of the Gordon cost over exceptional sets $D_{0;\varepsilon}(\mathbf{g})$. In particular, we show in Theorem 44 that with $L_n = Cn^\vartheta$ for sufficiently small $\vartheta > 0$ and a large enough $C_0 > 0$, w.h.p.,

$$\min_{w \in D_{0;\varepsilon}(\mathbf{g}) \cap B_{(2,\infty)}(C_0, L_n/\sqrt{n})} H_{0;Z}(w) \geq \max_{\beta > 0} \min_{\gamma > 0} \overline{D}_0(\beta, \gamma) + \Omega_\varepsilon(1). \quad (37)$$

Here $B_{(2,\infty)}(C_0, L_n/\sqrt{n}) = B_n(C_0) \cap L_\infty(L_n/\sqrt{n})$. A technical difficulty to apply the second inequality of (35) rests in its singular behavior near the interpolation regime $\eta = 0$. Also, we note that for a general exceptional set $D_{0;\varepsilon}(\mathbf{g})$, the maximum over v in $\min_{w \in D_{0;\varepsilon}(\mathbf{g})} H_{0;Z}(w) = \min_{w \in D_{0;\varepsilon}(\mathbf{g})} \max_v h_{0;Z}(w, v)$ need not be delocalized, so the first inequality of (35) cannot be applied. This singularity issue will be resolved in two steps:

- **(Step 2.1).** First, we use the second inequality of (35) to show that, (37) is valid for a version with small enough $\eta > 0$:

$$\mathbb{P} \left(\min_{w \in D_{\eta;\varepsilon}(\mathbf{g}) \cap B_{(2,\infty)}(C_0, L_n/\sqrt{n})} H_{\eta;Z}(w) \geq \max_{\beta > 0} \min_{\gamma > 0} \overline{D}_\eta(\beta, \gamma) + \Omega_\varepsilon(1) \right) \geq 1 - c_\eta \cdot \mathfrak{o}(1).$$

See (130) for a precise statement. As expected, c_η blows up as $\eta \downarrow 0$.

- **(Step 2.2).** Next, by using the ‘stability’ of the set $D_{\eta;\varepsilon}(\mathbf{g})$ (cf. Lemma 45) and $\max_{\beta > 0} \min_{\gamma > 0} \overline{D}_\eta(\beta, \gamma)$ (cf. Eqn. (76)) with respect to η , for a small enough $\eta > 0$, we have the following series of inequalities:

$$\min_{w \in D_{0;\varepsilon}(\mathbf{g}) \cap B_{(2,\infty)}(C_0, \frac{L_n}{\sqrt{n}})} H_{0;Z}(w)$$

$$\begin{aligned}
 &\geq \min_{w \in D_{0;\varepsilon}(\mathbf{g}) \cap B_{(2,\infty)}(C_0, \frac{L_0}{\sqrt{n}})} H_{\eta;Z}(w) \quad (\text{by definition of } H_{\eta;Z}) \\
 &\geq \min_{w \in D_{\eta;\varepsilon_\eta}(\mathbf{g}) \cap B_{(2,\infty)}(C_0, \frac{L_0}{\sqrt{n}})} H_{\eta;Z}(w) \quad (\varepsilon_\eta \approx \varepsilon \text{ by Lemma 45}) \\
 &\stackrel{\mathbb{P}}{\geq} \max_{\beta > 0} \min_{\gamma > 0} \bar{D}_\eta(\beta, \gamma) + \Omega_\varepsilon(1) \quad (\text{by Step 2.1 above}) \\
 &\geq \max_{\beta > 0} \min_{\gamma > 0} \bar{D}_0(\beta, \gamma) - \mathcal{O}(\eta) + \Omega_\varepsilon(1) \quad (\text{by Eqn. (76)}).
 \end{aligned}$$

Now for a given $\varepsilon > 0$, we may choose $\eta > 0$ small enough so that the term $-\mathcal{O}(\eta)$ is absorbed into $\Omega_\varepsilon(1)$, and therefore concluding (37).

A complete proof of the above outline is detailed in Section D.

Acknowledgments

The research of Q. Han is partially supported by NSF grant DMS-2143468. Both authors would like to thank the referees for their helpful comments and suggestions that significantly improved the quality of the paper.

References

- Madhu S Advani, Andrew M Saxe, and Haim Sompolinsky. High-dimensional dynamics of generalization error in neural networks. *Neural Networks*, 132:428–446, 2020.
- Zeyuan Allen-Zhu, Yuanzhi Li, and Zhao Song. A convergence theory for deep learning via over-parameterization. In *International Conference on Machine Learning*, pages 242–252. PMLR, 2019.
- Morgane Austern and Wenda Zhou. Asymptotics of cross-validation. *Ann. Inst. Henri Poincaré Probab. Stat.*, 61(4), 2025.
- Zhidong Bai and Jack W. Silverstein. *Spectral analysis of large dimensional random matrices*. Springer Series in Statistics. Springer, New York, second edition, 2010.
- Zhigang Bao, Qiyang Han, and Xiacong Xu. A leave-one-out approach to approximate message passing. *Ann. Appl. Probab.*, 35(4):2716–2766, 2025.
- Peter L. Bartlett, Philip M. Long, Gábor Lugosi, and Alexander Tsigler. Benign overfitting in linear regression. *Proc. Natl. Acad. Sci. USA*, 117(48):30063–30070, 2020.
- Peter L. Bartlett, Andrea Montanari, and Alexander Rakhlin. Deep learning: a statistical viewpoint. *Acta Numer.*, 30:87–201, 2021.
- Mohsen Bayati and Andrea Montanari. The LASSO risk for Gaussian matrices. *IEEE Trans. Inform. Theory*, 58(4):1997–2017, 2012.
- Mikhail Belkin, Daniel Hsu, Siyuan Ma, and Soumik Mandal. Reconciling modern machine-learning practice and the classical bias-variance trade-off. *Proc. Natl. Acad. Sci. USA*, 116(32):15849–15854, 2019.

- Mikhail Belkin, Daniel Hsu, and Ji Xu. Two models of double descent for weak features. *SIAM J. Math. Data Sci.*, 2(4):1167–1180, 2020.
- Pierre C. Bellec and Cun-Hui Zhang. Debiasing convex regularized estimators and interval estimation in linear models. *Ann. Statist.*, 51(2):391–436, 2023.
- Stéphane Boucheron, Gábor Lugosi, and Pascal Massart. *Concentration inequalities: A nonasymptotic theory of independence*. Oxford University Press, Oxford, 2013.
- Zhiqi Bu, Jason M. Klusowski, Cynthia Rush, and Weijie J. Su. Algorithmic analysis and statistical estimation of SLOPE via approximate message passing. *IEEE Trans. Inform. Theory*, 67(1):506–537, 2021.
- Michael Celentano, Andrea Montanari, and Yuting Wei. The Lasso with general Gaussian designs with applications to hypothesis testing. *Ann. Statist.*, 51(5):2194–2220, 2023.
- Xiangyu Chang, Yingcong Li, Samet Oymak, and Christos Thrampoulidis. Provable benefits of overparameterization in model compression: From double descent to pruning neural networks. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pages 6974–6983, 2021.
- Chen Cheng and Andrea Montanari. Dimension free ridge regression. *Ann. Statist.*, 52(6):2879–2912, 2024.
- Lenaïc Chizat, Edouard Oyallon, and Francis Bach. On lazy training in differentiable programming. *Advances in Neural Information Processing Systems*, 32, 2019.
- Peter Craven and Grace Wahba. Smoothing noisy data with spline functions. Estimating the correct degree of smoothing by the method of generalized cross-validation. *Numer. Math.*, 31(4):377–403, 1978/79.
- Zeyu Deng, Abba Kammoun, and Christos Thrampoulidis. A model of double descent for high-dimensional binary linear classification. *Inf. Inference*, 11(2):435–495, 2022.
- Alexis Derumigny and Johannes Schmidt-Hieber. On lower bounds for the bias-variance trade-off. *Ann. Statist.*, 51(4):1510–1533, 2023.
- Lee H. Dicker. Ridge regression and asymptotic minimax estimation over spheres of growing dimension. *Bernoulli*, 22(1):1–37, 2016.
- Edgar Dobriban and Stefan Wager. High-dimensional asymptotics of prediction: ridge regression and classification. *Ann. Statist.*, 46(1):247–279, 2018.
- David Donoho and Andrea Montanari. High dimensional robust M-estimation: asymptotic variance via approximate message passing. *Probab. Theory Related Fields*, 166(3-4):935–969, 2016.
- Simon S. Du, Xiyu Zhai, Barnabas Póczos, and Aarti Singh. Gradient descent provably optimizes over-parameterized neural networks. In *International Conference on Learning Representations*, 2019.

- Sandrine Dudoit and Mark J. van der Laan. Asymptotics of cross-validated risk estimation in estimator selection and performance assessment. *Stat. Methodol.*, 2(2):131–154, 2005.
- Bradley Efron. The estimation of prediction error: covariance penalties and cross-validation. *J. Amer. Statist. Assoc.*, 99(467):619–642, 2004. With comments and a rejoinder by the author.
- Noureddine El Karoui. Asymptotic behavior of unregularized and ridge-regularized high-dimensional robust regression estimators: rigorous results. *arXiv preprint arXiv:1311.2445*, 2013.
- Noureddine El Karoui. On the impact of predictor geometry on the performance on high-dimensional ridge-regularized generalized robust regression estimators. *Probab. Theory Related Fields*, 170(1-2):95–175, 2018.
- Evarist Giné and Richard Nickl. *Mathematical foundations of infinite-dimensional statistical models*. Cambridge Series in Statistical and Probabilistic Mathematics, [40]. Cambridge University Press, New York, 2016.
- Gene H. Golub, Michael Heath, and Grace Wahba. Generalized cross-validation as a method for choosing a good ridge parameter. *Technometrics*, 21(2):215–223, 1979.
- Qiyang Han. Noisy linear inverse problems under convex constraints: Exact risk asymptotics in high dimensions. *Ann. Statist.*, 51(4):1611–1638, 2023.
- Qiyang Han and Yandi Shen. Universality of regularized regression estimators in high dimensions. *Ann. Statist.*, 51(4):1799–1823, 2023.
- Trevor Hastie, Andrea Montanari, Saharon Rosset, and Ryan J. Tibshirani. Surprises in high-dimensional ridgeless least squares interpolation. *Ann. Statist.*, 50(2):949–986, 2022.
- Arthur E. Hoerl and Robert W. Kennard. Ridge regression: Biased estimation for nonorthogonal problems. *Technometrics*, 12(1):55–67, 1970.
- Daniel Hsu, Sham M. Kakade, and Tong Zhang. Random design analysis of ridge regression. *Found. Comput. Math.*, 14(3):569–600, 2014.
- Arthur Jacot, Franck Gabriel, and Clément Hongler. Neural tangent kernel: Convergence and generalization in neural networks. *Advances in Neural Information Processing Systems*, 31, 2018.
- Gareth James, Daniela Witten, Trevor Hastie, and Robert Tibshirani. *An introduction to statistical learning—with applications in R*. Springer Texts in Statistics. Springer, New York, 2021. Second edition [of 3100153].
- Nicholas Kissel and Jing Lei. On high-dimensional gaussian comparisons for cross-validation. *arXiv preprint arXiv:2211.04958*, 2022.
- Antti Knowles and Jun Yin. Anisotropic local laws for random matrices. *Probab. Theory Related Fields*, 169(1-2):257–352, 2017.

- Dmitry Kobak, Jonathan Lomond, and Benoit Sanchez. The optimal ridge penalty for real-world high-dimensional data can be zero or negative due to the implicit ridge regularization. *J. Mach. Learn. Res.*, 21:Paper No. 169, 16, 2020.
- Frederic Koehler, Lijia Zhou, Danica J Sutherland, and Nathan Srebro. Uniform convergence of interpolators: Gaussian width, norm bounds and benign overfitting. *Advances in Neural Information Processing Systems*, 34:20657–20668, 2021.
- Ker-Chau Li. From Stein’s unbiased risk estimates to the method of generalized cross validation. *Ann. Statist.*, 13(4):1352–1377, 1985.
- Ker-Chau Li. Asymptotic optimality of C_L and generalized cross-validation in ridge regression with application to spline smoothing. *Ann. Statist.*, 14(3):1101–1112, 1986.
- Ker-Chau Li. Asymptotic optimality for C_p , C_L , cross-validation and generalized cross-validation: discrete index set. *Ann. Statist.*, 15(3):958–975, 1987.
- Yue Li and Yuting Wei. Minimum ℓ_1 -norm interpolators: Precise asymptotics and multiple descent. *arXiv preprint arXiv:2110.09502*, 2021.
- Tengyuan Liang and Pragya Sur. A precise high-dimensional asymptotic theory for boosting and minimum- ℓ_1 -norm interpolated classifiers. *Ann. Statist.*, 50(3):1669–1695, 2022.
- Sifan Liu and Edgar Dobriban. Ridge regression: Structure, cross-validation, and sketching. *arXiv preprint arXiv:1910.02373*, 2019.
- 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.
- Léo Miolane and Andrea Montanari. The distribution of the Lasso: uniform control over sparse balls and adaptive parameter tuning. *Ann. Statist.*, 49(4):2313–2335, 2021.
- Andrea Montanari, Feng Ruan, Youngtak Sohn, and Jun Yan. The generalization error of max-margin linear classifiers: benign overfitting and high dimensional asymptotics in the overparametrized regime. *Ann. Statist.*, 53(2):822–853, 2025.
- Vidya Muthukumar, Kailas Vodrahalli, Vignesh Subramanian, and Anant Sahai. Harmless interpolation of noisy data in regression. *IEEE Journal on Selected Areas in Information Theory*, 1(1):67–83, 2020.
- Pratik Patil, Yuting Wei, Alessandro Rinaldo, and Ryan Tibshirani. Uniform consistency of cross-validation estimators for high-dimensional ridge regression. In *International Conference on Artificial Intelligence and Statistics*, pages 3178–3186. PMLR, 2021.
- Pratik Patil, Jin-Hong Du, and Ryan J Tibshirani. Optimal ridge regularization for out-of-distribution prediction. *arXiv preprint arXiv:2404.01233*, 2024.

- Dominic Richards, Jaouad Mourtada, and Lorenzo Rosasco. Asymptotics of ridge (less) regression under general source condition. In *International Conference on Artificial Intelligence and Statistics*, pages 3889–3897. PMLR, 2021.
- Mark Rudelson and Roman Vershynin. Smallest singular value of a random rectangular matrix. *Comm. Pure Appl. Math.*, 62(12):1707–1739, 2009.
- Fariborz Salehi, Ehsan Abbasi, and Babak Hassibi. The impact of regularization on high-dimensional logistic regression. *Advances in Neural Information Processing Systems*, 32, 2019.
- Charles M. Stein. Estimation of the mean of a multivariate normal distribution. *Ann. Statist.*, 9(6):1135–1151, 1981.
- M. Stone. Cross-validatory choice and assessment of statistical predictions. *J. Roy. Statist. Soc. Ser. B*, 36:111–147, 1974.
- M. Stone. Asymptotics for and against cross-validation. *Biometrika*, 64(1):29–35, 1977.
- Pragya Sur and Emmanuel J. Candès. A modern maximum-likelihood theory for high-dimensional logistic regression. *Proc. Natl. Acad. Sci.*, 116(29):14516–14525, 2019.
- Shange Tang, Jiayun Wu, Jianqing Fan, and Chi Jin. Benign overfitting in out-of-distribution generalization of linear models. *arXiv preprint arXiv:2412.14474*, 2024.
- Christos Thrampoulidis, Samet Oymak, and Babak Hassibi. Regularized linear regression: A precise analysis of the estimation error. In *Conference on Learning Theory*, pages 1683–1709. PMLR, 2015.
- Christos Thrampoulidis, Ehsan Abbasi, and Babak Hassibi. Precise error analysis of regularized M -estimators in high dimensions. *IEEE Trans. Inform. Theory*, 64(8):5592–5628, 2018.
- Alexander Tsigler and Peter L. Bartlett. Benign overfitting in ridge regression. *J. Mach. Learn. Res.*, 24:Paper No. [123], 76, 2023.
- Antonia M Tulino, Sergio Verdú, et al. Random matrix theory and wireless communications. *Foundations and Trends® in Communications and Information Theory*, 1(1):1–182, 2004.
- Aad van der Vaart and Jon A. Wellner. *Weak Convergence and Empirical Processes*. Springer Series in Statistics. Springer-Verlag, New York, 1996.
- Shuaiwen Wang, Haolei Weng, and Arian Maleki. Does SLOPE outperform bridge regression? *Inf. Inference*, 11(1):1–54, 2022.
- Denny Wu and Ji Xu. On the optimal weighted ℓ_2 regularization in overparameterized linear regression. *Advances in Neural Information Processing Systems*, 33:10112–10123, 2020.

Chiyuan Zhang, Samy Bengio, Moritz Hardt, Benjamin Recht, and Oriol Vinyals. Understanding deep learning (still) requires rethinking generalization. *Communications of the ACM*, 64(3):107–115, 2021.

Xianyang Zhang, Huijuan Zhou, and Hanxuan Ye. A modern theory for high-dimensional Cox regression models. *arXiv preprint arXiv:2204.01161*, 2022.

Lijia Zhou, Frederic Koehler, Danica J Sutherland, and Nathan Srebro. Optimistic rates: A unifying theory for interpolation learning and regularization in linear regression. *ACM/JMS Journal of Data Science*, 1(2):1–51, 2024.

Appendix A. Proof preliminaries

A.1 Some properties of \mathbf{e}_F and prox_F

We write $g_n \equiv g/\sqrt{n}$ in this subsection. First we give an explicit expression for $\mathbb{E} \text{err}_{(\Sigma, \mu_0)}(\gamma; \tau)$ and $\mathbb{E} \text{dof}_{(\Sigma, \mu_0)}(\gamma; \tau)$.

Lemma 15. *For any $(\gamma, \tau) \in (0, \infty)^2$,*

$$\begin{aligned} \mathbb{E} \text{err}_{(\Sigma, \mu_0)}(\gamma; \tau) &= \tau^2 \|(\Sigma + \tau I)^{-1} \Sigma^{1/2} \mu_0\|^2 + \gamma^2 \cdot n^{-1} \text{tr}(\Sigma^2(\Sigma + \tau I)^{-2}), \\ \mathbb{E} \text{dof}_{(\Sigma, \mu_0)}(\gamma; \tau) &= \gamma^2 \cdot n^{-1} \text{tr}(\Sigma(\Sigma + \tau I)^{-1}). \end{aligned}$$

Proof Using the closed-form of $\widehat{\mu}_{(\Sigma, \mu_0)}^{\text{seq}}$, we may compute

$$\Sigma^{1/2}(\widehat{\mu}_{(\Sigma, \mu_0)}^{\text{seq}}(\gamma; \tau) - \mu_0) = (\Sigma + \tau I)^{-1} \Sigma^{1/2}(-\tau \mu_0 + \gamma \Sigma^{1/2} g_n). \quad (38)$$

The claims follow from direct calculations. \blacksquare

Next we give explicit expression for $\text{prox}_F(\gamma g_n; \tau)$ and $\mathbf{e}_F(\gamma g_n; \tau)$.

Lemma 16. *It holds that*

$$\begin{aligned} \text{prox}_F(\gamma g_n; \tau) &= \Sigma^{1/2}(\widehat{\mu}_{(\Sigma, \mu_0)}^{\text{seq}}(\gamma; \tau) - \mu_0), \\ \mathbf{e}_F(\gamma g_n; \tau) &= \frac{1}{2\tau} \|\Sigma^{1/2} \widehat{\mu}_{(\Sigma, \mu_0)}^{\text{seq}}(\gamma; \tau) - y_{(\Sigma, \mu_0)}^{\text{seq}}(\gamma)\|^2 + \frac{1}{2} \|\widehat{\mu}_{(\Sigma, \mu_0)}^{\text{seq}}(\gamma; \tau)\|^2. \end{aligned}$$

Furthermore,

$$\begin{aligned} \mathbb{E} \mathbf{e}_F(\gamma g_n; \tau) &= \frac{1}{2\tau} (\mathbb{E} \text{err}_{(\Sigma, \mu_0)}(\gamma; \tau) - 2 \mathbb{E} \text{dof}_{(\Sigma, \mu_0)}(\gamma; \tau) + \gamma^2) \\ &\quad + \frac{1}{2} \left(\|(\Sigma + \tau I)^{-1} \Sigma \mu_0\|^2 + \gamma^2 \cdot \frac{1}{n} \text{tr}(\Sigma(\Sigma + \tau I)^{-2}) \right). \end{aligned}$$

Proof The two identities in the first display follows from the definition of F . For the second display, note that $\mathbb{E} \mathbf{e}_F(\gamma g_n; \tau)$ is equal to

$$\frac{1}{2\tau} (\mathbb{E} \text{err}_{(\Sigma, \mu_0)}(\gamma; \tau) - 2 \mathbb{E} \text{dof}_{(\Sigma, \mu_0)}(\gamma; \tau) + \gamma^2) + \frac{1}{2} \mathbb{E} \|\widehat{\mu}_{(\Sigma, \mu_0)}^{\text{seq}}(\gamma; \tau)\|^2.$$

Using $\mathbb{E} \|\widehat{\mu}_{(\Sigma, \mu_0)}^{\text{seq}}(\gamma; \tau)\|^2 = \|(\Sigma + \tau I)^{-1} \Sigma \mu_0\|^2 + \gamma^2 \cdot n^{-1} \text{tr}(\Sigma(\Sigma + \tau I)^{-2})$ to conclude. \blacksquare

The derivative formula below for \mathbf{e}_F will be useful.

Lemma 17. *It holds that*

$$\nabla_x \mathbf{e}_F(x; \tau) = \frac{1}{\tau} (x - \text{prox}_F(x; \tau)), \quad \partial_\tau \mathbf{e}_F(x; \tau) = -\frac{1}{2\tau^2} \|x - \text{prox}_F(x; \tau)\|^2.$$

Proof See e.g., (Thrampoulidis et al., 2018, Lemmas B.5 and D.1). \blacksquare

Finally we provide a concentration inequality for $\mathbf{e}_F(\gamma g_n; \tau)$.

Proposition 18. *There exists some universal constant $C > 0$ such that*

$$\mathbb{P}\left(\left|\mathbf{e}_F(\gamma g_n; \tau) - \mathbb{E} \mathbf{e}_F(\gamma g_n; \tau)\right| \geq C \left\{v \mathbb{E}^{1/2} \mathbf{e}_F(\gamma g_n; \tau) \sqrt{\frac{t}{n} + v^2 \cdot \frac{t}{n}}\right\}\right) \leq C e^{-t/C}$$

holds for any $t \geq 0$. Here $v^2 \equiv v^2(\gamma, \tau) \equiv \gamma^2(\tau \|(\Sigma + \tau I)^{-1}\|_{\text{op}}^2 + \|(\Sigma + \tau I)^{-1} \Sigma^{1/2}\|_{\text{op}}^2)$.

Proof Using that $\nabla_g \widehat{\mu}_{(\Sigma, \mu_0)}^{\text{seq}}(\gamma; \tau) = \frac{\gamma}{\sqrt{n}}(\Sigma + \tau I)^{-1} \Sigma^{1/2}$ and $\nabla_g y^{\text{seq}}(\gamma) = \frac{\gamma}{\sqrt{n}} I$,

$$\begin{aligned} \nabla_g \mathbf{e}_F(\gamma g_n; \tau) &= \frac{1}{\tau} \cdot \frac{\gamma}{\sqrt{n}} \left((\Sigma + \tau I)^{-1} \Sigma - I \right) \left(\Sigma^{1/2} \widehat{\mu}_{(\Sigma, \mu_0)}^{\text{seq}}(\gamma; \tau) - y_{(\Sigma, \mu_0)}^{\text{seq}}(\gamma) \right) \\ &\quad + \frac{\gamma}{\sqrt{n}} (\Sigma + \tau I)^{-1} \Sigma^{1/2} \nabla_g \widehat{\mu}_{(\Sigma, \mu_0)}^{\text{seq}}(\gamma; \tau). \end{aligned}$$

This means

$$\begin{aligned} \|\nabla_g \mathbf{e}_F(\gamma g_n; \tau)\|^2 &\leq 2\gamma^2 \cdot n^{-1} \left\{ \|(\Sigma + \tau I)^{-1}\|_{\text{op}}^2 \|\Sigma^{1/2} \widehat{\mu}_{(\Sigma, \mu_0)}^{\text{seq}}(\gamma; \tau) - y_{(\Sigma, \mu_0)}^{\text{seq}}(\gamma)\|^2 \right. \\ &\quad \left. + \|(\Sigma + \tau I)^{-1} \Sigma^{1/2}\|_{\text{op}}^2 \|\nabla_g \widehat{\mu}_{(\Sigma, \mu_0)}^{\text{seq}}(\gamma; \tau)\|^2 \right\} \\ &\leq 4\gamma^2 \cdot n^{-1} \left(\tau \|(\Sigma + \tau I)^{-1}\|_{\text{op}}^2 + \|(\Sigma + \tau I)^{-1} \Sigma^{1/2}\|_{\text{op}}^2 \right) \cdot \mathbf{e}_F(\gamma g_n; \tau). \quad (39) \end{aligned}$$

From here we may conclude by setting $H(g) \equiv \mathbf{e}_F(\gamma g_n; \tau)$ and $\Gamma^2 \equiv 4\gamma^2 n^{-1} (\tau \|(\Sigma + \tau I)^{-1}\|_{\text{op}}^2 + \|(\Sigma + \tau I)^{-1} \Sigma^{1/2}\|_{\text{op}}^2)$ in Proposition 59. \blacksquare

A.2 Some high probability events

Let

$$e_h^2 = \|h\|^2/m, \quad e_g^2 \equiv \|g\|^2/n. \quad (40)$$

For $M, \delta > 0$, consider the event

$$\begin{aligned} \mathcal{E}_0(M) &\equiv \{(\|G\|_{\text{op}}/\sqrt{n}) \vee [|(GG^\top/n)^{-1}|_{\text{op}} \mathbf{1}_{\eta=0}] \leq M\}, \\ \mathcal{E}_{1,0}(\delta) &\equiv \{|e_g^2 - 1| \vee |e_h^2 - 1| \vee |n^{-1/2} \langle \Sigma^{1/2} g, \mu_0 \rangle| \vee |n^{-1} \langle h, \xi \rangle| \leq \delta\}, \\ \mathcal{E}_{1,\xi}(\delta) &\equiv \{|\|\xi\|^2/m - \sigma_\xi^2| \leq \delta\}, \\ \mathcal{E}_1(\delta) &\equiv \mathcal{E}_{1,0}(\delta) \cap \mathcal{E}_{1,\xi}(\delta). \end{aligned}$$

Here in the definition of $\mathcal{E}_0(M)$, we interpret $\infty \cdot 0 = 0$. Typically we think of $M \asymp 1$ and $\delta \asymp 1/\sqrt{n}$.

Lemma 19. *Fix $\delta \in (0, 1/2)$ and $L_w > 0$. Then $\mathcal{E}_1(\delta) \subset \mathcal{E}_2(4(\sigma_\xi^2 + 1 + \phi^{-1} L_w)\delta, L_w)$, where $\mathcal{E}_2(\delta, L_w) \equiv \{|\sigma_\pm^2(L_w) - \sigma_\xi^2| \leq \delta\}$.*

Proof Using the definition of $\sigma_\pm^2(L_w)$ in (30), on $\mathcal{E}_1(\delta)$, we have

$$|\sigma_\pm^2(L_w) - \sigma_\xi^2| \leq \frac{\|\xi\|^2}{\|h\|^2} |e_h^2 - 1| + \left| \frac{\|\xi\|^2}{m} - \sigma_\xi^2 \right| + \frac{2L_w |\langle h, \xi \rangle|}{\|h\|^2} \leq 4(\sigma_\xi^2 + 1 + \phi^{-1} L_w)\delta.$$

The claim follows. \blacksquare

Lemma 20. *Suppose $1/K \leq \phi^{-1} - \mathbf{1}_{\eta=0} \leq K$. Then there exists some $C = C(K) > 0$ such that $\mathbb{P}(\mathcal{E}_0(C)) \geq 1 - Ce^{-n/C}$.*

Proof The claim for $\|G\|_{\text{op}}/\sqrt{n}$ follows from standard concentration estimates. The claim for $\|(GG^\top/n)^{-1}\|_{\text{op}}$ follows from, e.g., (Rudelson and Vershynin, 2009, Theorem 1.1). ■

Lemma 21. *Suppose $1/K \leq \phi^{-1} \leq K$, and $\|\mu_0\| \vee \|\Sigma\|_{\text{op}} \leq K$ for some $K > 0$, and Assumption B hold with $\sigma_\xi^2 > 0$. There exists some constant $C = C(K, \sigma_\xi) > 0$ such that for all $t \geq 0$, with $\delta(t, n) \equiv C(\sqrt{t/n} + t/n)$, for $\xi \in \mathcal{E}_{1, \xi}(\delta(t, n))$, we have $\mathbb{P}^\xi(\mathcal{E}_1(\delta(t, n))) \geq 1 - e^{-t}$.*

Proof The claim follows by standard concentration inequalities. ■

A.3 Some connections of the fixed point equation (13) to RMT

The second equation of (13) has a natural connection to RMT. To detail this connection, let $\widehat{\Sigma} \equiv \Sigma^{1/2}G^\top G\Sigma^{1/2}/m \in \mathbb{R}^{n \times n}$ and $\check{\Sigma} \equiv G\Sigma G^\top/m \in \mathbb{R}^{m \times m}$ be the sample covariance matrix and its dimension flipped, companion matrix. For $z \in \mathbb{C}^+ \equiv \{z \in \mathbb{C} : \Im z > 0\}$, let $\mathbf{m}_n(z) \equiv m^{-1} \text{tr}(\check{\Sigma} - zI_m)^{-1}$ and $\mathbf{m}(z)$ be the Stieltjes transforms of the empirical spectral distribution and the asymptotic eigenvalue density (cf. (Knowles and Yin, 2017, Definition 2.3)) of $\check{\Sigma}$, respectively. It is well-known that $\mathbf{m}(z)$ can be determined uniquely via the fixed point equation

$$z = -\frac{1}{\mathbf{m}(z)} + \frac{1}{\phi} \cdot \frac{1}{n} \text{tr} \left((I_n + \Sigma \mathbf{m}(z))^{-1} \Sigma \right). \quad (41)$$

See, e.g., (Knowles and Yin, 2017, Lemma 2.2) for more technical details and historical references. We also note that while the above equation is initially defined for $z \in \mathbb{C}^+$, it can be straightforwardly extended to the real axis provided that z lies outside the support of the asymptotic spectrum of $\check{\Sigma}$.

The following proposition provides a precise connection between the effective regularization $\tau_{\eta,*}$ defined via the second equation of (13), and the Stieltjes transform \mathbf{m} . This connection will prove important in some of the results ahead.

Proposition 22. *For any $\eta > 0$ and $\eta = 0$ with $\phi^{-1} > 1$,*

$$n^{-1} \text{tr}((\Sigma + \tau_{\eta,*} I_n)^{-1} \Sigma) = \phi - \eta \cdot \mathbf{m}(-\eta/\phi). \quad (42)$$

Proof By comparing (41) and the second equation of (13), we may identify the two equations by setting $\tau_{\eta,*} \equiv 1/\mathbf{m}(-z_\eta)$ with $z_\eta \equiv \eta/\phi$, as claimed. ■

While (42) appears somewhat purely algebraic, it actually admits a natural statistical interpretation. Suppose ξ is also Gaussian. We may then compute

$$\text{df}(\widehat{\mu}_\eta) = \sum_{j=1}^n \frac{\text{Cov}^X((X\widehat{\mu}_\eta)_j, Y_j)}{\sigma_\xi^2} = \text{tr}((\widehat{\Sigma} + z_\eta I_n)^{-1} \widehat{\Sigma}) = n(\phi - \eta \cdot \mathbf{m}_n(-z_\eta)). \quad (43)$$

Now comparing the above display with (42), we arrive at the following intriguing equivalence between the averaged law in RMT, and the proximity of $\widehat{\mu}_\eta$ and $\widehat{\mu}_{(\Sigma, \mu_0)}^{\text{seq}}(\gamma_{\eta,*}; \tau_{\eta,*})$ in terms of “degrees-of-freedom”:

$$\mathbf{m}_n(-z_\eta) \stackrel{\mathbb{P}}{\approx} \mathbf{m}(-z_\eta) \Leftrightarrow \text{df}(\widehat{\mu}_\eta)/n \stackrel{\mathbb{P}}{\approx} \text{df}(\widehat{\mu}_{(\Sigma, \mu_0)}^{\text{seq}}(\gamma_{\eta,*}; \tau_{\eta,*}))/n.$$

Appendix B. Properties of the fixed point equations

B.1 The fixed point equation (13)

Proposition 23. *The following hold.*

1. *The fixed point equation (13) admits a unique solution $(\gamma_{\eta,*}, \tau_{\eta,*}) \in (0, \infty)^2$, for all $(m, n) \in \mathbb{N}^2$ when $\eta > 0$ and $m < n$ when $\eta = 0$.*
2. *The following apriori bounds hold:*

$$\begin{aligned} \frac{1 - \phi + \sqrt{(1 - \phi)^2 + 4\mathcal{H}_\Sigma \eta}}{2\mathcal{H}_\Sigma} &\leq \tau_{\eta,*} \leq \inf_{k \in [0: \min\{m-1, n\}]} \left\{ \frac{\sum_{j>k} \lambda_j}{m-k} + \frac{n}{m-k} \cdot \eta \right\}, \\ \frac{\sigma_\xi^2}{\phi} &\leq \gamma_{\eta,*}^2 \leq \frac{\sigma_\xi^2 + \|\Sigma\|_{\text{op}} \|\mu_0\|^2}{\phi} \left(1 + \frac{\|\Sigma\|_{\text{op}}}{\tau_{\eta,*}} \right). \end{aligned}$$

3. *If $1/K \leq \phi^{-1} \leq K$ and $\|\Sigma\|_{\text{op}} \vee \mathcal{H}_\Sigma \leq K$ for some $K > 1$, then there exists some $C = C(K) > 1$ such that uniformly in $\eta \in \Xi_K$,*

$$1/C \leq \tau_{\eta,*} \leq C, \quad 1/C \leq (-1)^{q+1} \partial_\eta^q \tau_{\eta,*} \leq C, \quad q \in \{1, 2\}.$$

If furthermore $1/K \leq \sigma_\xi^2 \leq K$ and $\|\mu_0\| \leq K$, then uniformly in $\eta \in \Xi_K$,

$$1/C \leq \gamma_{\eta,*} \leq C, \quad |\partial_\eta \gamma_{\eta,*}| \leq C.$$

Proof We shall write $(\gamma_{\eta,*}, \tau_{\eta,*}) = (\gamma_*, \tau_*)$ for notational simplicity. All the constants in $\lesssim, \gtrsim, \asymp$ below may depend on K .

(1). First we prove the existence and uniqueness of τ_* . We rewrite the second equation of (13) as

$$\phi = \frac{1}{n} \text{tr}((\Sigma + \tau_* I)^{-1} \Sigma) + \frac{\eta}{\tau_*} = \frac{1}{n} \sum_{j=1}^n \frac{\lambda_j}{\lambda_j + \tau_*} + \frac{\eta}{\tau_*} \equiv \mathbf{f}(\tau_*). \quad (44)$$

Clearly $\mathbf{f}(\tau)$ is smooth, non-increasing, $\mathbf{f}(0) = 1 > \phi$ for $\eta = 0$ and $\mathbf{f}(0) = \infty$ for $\eta > 0$, and $\mathbf{f}(\infty) = 0$, so $\tau \mapsto \mathbf{f}(\tau) - \phi$ must admit a unique zero $\tau_* \in (0, \infty)$.

Next we prove the existence and uniqueness of γ_* . Using Lemma 15, the equation $\phi \gamma_*^2 = \sigma_\xi^2 + \mathbb{E} \text{err}_{(\Sigma, \mu_0)}(\gamma_*; \tau_*)$ reads

$$\phi = \frac{1}{\gamma_*^2} (\sigma_\xi^2 + \tau_*^2 \|(\Sigma + \tau_* I)^{-1} \Sigma^{1/2} \mu_0\|^2) + \frac{1}{n} \text{tr}((\Sigma + \tau_* I)^{-2} \Sigma^2). \quad (45)$$

As $n^{-1} \text{tr}((\Sigma + \tau_* I)^{-2} \Sigma^2) < n^{-1} \text{tr}((\Sigma + \tau_* I)^{-1} \Sigma) \leq \phi$ by (44) and the fact $\tau_* > 0$, the above equation admits a unique solution $\gamma_* \in (0, \infty)$, analytically given by

$$\gamma_*^2 = \frac{\sigma_\xi^2 + \tau_*^2 \|(\Sigma + \tau_* I)^{-1} \Sigma^{1/2} \mu_0\|^2}{\phi - \frac{1}{n} \text{tr}((\Sigma + \tau_* I)^{-2} \Sigma^2)} = \frac{\sigma_\xi^2 + \tau_*^2 \|(\Sigma + \tau_* I)^{-1} \Sigma^{1/2} \mu_0\|^2}{\frac{\eta}{\tau_*} + \frac{\tau_*}{n} \text{tr}((\Sigma + \tau_* I)^{-2} \Sigma)}. \quad (46)$$

(2). For the upper bound for τ_* , using the equation (44), we have

$$m = n\phi \leq k + \frac{1}{\tau_*} \sum_{j>k} \lambda_j + \frac{n\eta}{\tau_*}, \quad \forall k \in [0 : n], k \leq m - 1.$$

Solving for τ_* yields the desired upper bound. For the lower bound for τ_* , note that (44) leads to

$$\phi = 1 - \tau_* \cdot \frac{1}{n} \sum_{j=1}^n \frac{1}{\lambda_j + \tau_*} + \frac{\eta}{\tau_*} \geq 1 - \tau_* \mathcal{H}_\Sigma + \frac{\eta}{\tau_*},$$

or equivalently $\mathcal{H}_\Sigma \tau_*^2 + (\phi - 1) \tau_* - \eta \geq 0$. Solving this quadratic inequality yields the lower bound for τ_* .

On the other hand, the lower bound $\gamma_*^2 \geq \sigma_\xi^2 / \phi$ is trivial by (46). For the upper bound for γ_* , using that

$$\phi - \frac{1}{n} \text{tr}((\Sigma + \tau_* I)^{-2} \Sigma^2) \geq \phi - \frac{1}{n} \text{tr}((\Sigma + \tau_* I)^{-1} \Sigma) \cdot \max_{j \in [n]} \frac{\lambda_j}{\lambda_j + \tau_*} \geq \phi \cdot \frac{\tau_*}{\|\Sigma\|_{\text{op}} + \tau_*}, \quad (47)$$

and the first identity in (46), we have

$$\gamma_*^2 \leq \phi^{-1} (\sigma_\xi^2 + \|\Sigma\|_{\text{op}} \|\mu_0\|^2) (1 + \|\Sigma\|_{\text{op}} / \tau_*).$$

Collecting the bounds proves the claim.

(3). The claim on γ_*, τ_* is a simple consequence of (2). We shall prove the other claim on their derivatives. Viewing $\tau_* = \tau_*(\eta)$ and taking derivative with respect to η on both sides of (44) yield that, with $T_{-p,q}(\eta) \equiv n^{-1} \text{tr}((\Sigma + \tau_*(\eta) I)^{-p} \Sigma^q)$ for $p, q \in \mathbb{N}$,

$$0 = -T_{-2,1}(\eta) \cdot \tau'_*(\eta) + \frac{1}{\tau_*(\eta)} - \frac{\eta}{\tau_*^2(\eta)} \cdot \tau'_*(\eta).$$

Solving for $\tau'_*(\eta)$ yields that

$$\tau'_*(\eta) = \frac{\tau_*(\eta)}{\eta + \tau_*^2(\eta) \cdot T_{-2,1}(\eta)} \equiv \frac{\tau_*(\eta)}{G_0(\eta)}. \quad (48)$$

Further taking derivative with respect to η on both sides of the above display (48), we have

$$\begin{aligned} \tau_*''(\eta) &= \frac{1}{G_0^2(\eta)} (\tau'_*(\eta) G_0(\eta) - \tau_*(\eta) G_0'(\eta)) \\ &= \frac{1}{G_0^2(\eta)} \left\{ \tau_*(\eta) - \tau_*(\eta) \left(1 + 2\tau_*(\eta) \tau'_*(\eta) T_{-2,1}(\eta) - 2\tau_*^2(\eta) \tau'_*(\eta) T_{-3,1}(\eta) \right) \right\} \end{aligned}$$

$$= \frac{2\tau_*^2(\eta)\tau'_*(\eta)}{G_0^2(\eta)} \left(\tau_*(\eta)T_{-3,1}(\eta) - T_{-2,1}(\eta) \right) = -\frac{2\tau_*^2(\eta)\tau'_*(\eta)}{G_0^2(\eta)} T_{-3,2}(\eta). \quad (49)$$

Using the apriori estimate for $\tau_*(\eta)$ proved in (2), it follows that for $q \in \{1, 2\}$,

$$1 \lesssim \inf_{\eta \in \Xi_K} (-1)^{q+1} \tau_*^{(q)}(\eta) \leq \sup_{\eta \in \Xi_K} (-1)^{q+1} \tau_*^{(q)}(\eta) \lesssim 1. \quad (50)$$

For $\gamma'_*(\eta)$, let us define

$$G_1(\eta) \equiv \sigma_\xi^2 + \tau_*^2(\eta) \|(\Sigma + \tau_*(\eta)I)^{-1} \Sigma^{1/2} \mu_0\|^2, \quad G_2(\eta) \equiv \phi - n^{-1} \text{tr}((\Sigma + \tau_*(\eta)I)^{-2} \Sigma^2).$$

Then

$$\gamma'_*(\eta) = \frac{G'_1(\eta)G_2(\eta) - G_1(\eta)G'_2(\eta)}{2\gamma_*(\eta)G_2^2(\eta)}. \quad (51)$$

We shall now prove bounds for G_1, G'_1, G_2, G'_2 . First, using (47), we have

$$\sigma_\xi^2 \leq G_1(\eta) \leq \sigma_\xi^2 + \frac{\tau_*(\eta)}{2} \|\mu_0\|^2, \quad \phi \cdot \frac{\tau_*(\eta)}{\|\Sigma\|_{\text{op}} + \tau_*(\eta)} \leq G_2(\eta) \leq \phi.$$

In particular, uniformly in $\eta \in \Xi_K$,

$$G_1(\eta), G_2(\eta) \asymp 1. \quad (52)$$

The derivatives G'_1, G'_2 are

$$\begin{aligned} G'_1(\eta) &= 2\tau_*(\eta)\tau'_*(\eta) \|(\Sigma + \tau_*(\eta)I)^{-1} \Sigma^{1/2} \mu_0\|^2 \\ &\quad - 2\tau_*^2(\eta) \|(\Sigma + \tau_*(\eta)I)^{-3/2} \Sigma^{1/2} \mu_0\|^2 \cdot \tau'_*(\eta), \\ G'_2(\eta) &= 2 \cdot n^{-1} \text{tr}((\Sigma + \tau_*(\eta)I)^{-3} \Sigma^2) \cdot \tau'_*(\eta). \end{aligned}$$

Using the apriori estimates on $\tau_*(\eta)$ and (50), it now follows that

$$\sup_{\eta \in \Xi_K} \{|G'_1(\eta)| \vee |G'_2(\eta)|\} \lesssim 1. \quad (53)$$

Combining (51)-(53) and using apriori estimates on $\gamma_*(\eta)$, we arrive at

$$\sup_{\eta \in \Xi_K} |\gamma'_*(\eta)| \lesssim 1. \quad (54)$$

The claim follows by collecting (50) and (54). ■

B.2 Sample version of (13)

Let the sample version of (13) be defined by

$$\begin{cases} \phi e_h^2 \gamma^2 = \sigma_\pm^2(L_w) + \text{err}_{(\Sigma, \mu_0)}(\gamma; \tau), \\ (\phi e_h^2 - \frac{\eta}{\tau}) \cdot \gamma^2 = \text{dof}_{(\Sigma, \mu_0)}(\gamma; \tau). \end{cases} \quad (55)$$

Here recall that e_h^2 is defined in (40), and $\sigma_\pm^2(L_w)$ is defined in (30).

Proposition 24. $1/K \leq \phi^{-1}, \sigma_{\xi}^2 \leq K$ and $\|\mu_0\| \vee \|\Sigma\|_{\text{op}} \vee \mathcal{H}_{\Sigma} \leq K$ for some $K > 0$. There exist some $C, C_0 > 1$ depending on K , such that with $\delta \in (0, 1/C^{100})$, $1 \leq M \leq \sqrt{n}/C$ and $L_w \leq C$, on the event $\mathcal{E}_1(\delta) \cap \mathcal{E}_{\Delta, \Xi}(M)$, where

$$\mathcal{E}_{\Delta, \Xi}(M) \equiv \left\{ \max_{\ell=1,2} \sup_{\tau \geq 0} |\Delta_{\ell}(\tau)| \vee \max_{\ell=1,2} \sup_{\tau \geq 0} n^{-1/2} |\Xi_{\ell}(\tau) - \mathbb{E} \Xi_{\ell}(\tau)| \leq M \right\}$$

with $\Delta_{\ell}, \Xi_{\ell}$ defined in Lemmas 25 and 26 ahead, the following hold.

1. All solutions $(\gamma_{n,\eta,\pm}, \tau_{n,\eta,\pm})$ to the system of equations in (55) satisfy

$$1/C_0 \leq \tau_{n,\eta,\pm} \leq C_0, \quad 1/C_0 \leq \gamma_{n,\eta,\pm} \leq C_0$$

uniformly in $\eta \in \Xi_K$.

2. Moreover,

$$\sup_{\eta \in \Xi_K} \left\{ |\tau_{n,\eta,\pm} - \tau_{\eta,*}| \vee |\gamma_{n,\eta,\pm} - \gamma_{\eta,*}| \right\} \leq C_0 \cdot (M/\sqrt{n} + \delta).$$

We need two concentration lemmas before the proof of Proposition 24.

Lemma 25. Let $\Delta_{\ell}(\tau) \equiv -\ell \cdot \tau \langle (\Sigma + \tau I)^{-\ell} \Sigma^{\ell-1/2} \mu_0, g \rangle$ for $\ell = 1, 2$. Suppose that $\|\mu_0\| \vee \|\Sigma\|_{\text{op}} \vee \mathcal{H}_{\Sigma} \leq K$ for some $K > 0$. Then there exists some constant $C = C(K) > 1$ such that for $t \geq C \log(en)$,

$$\mathbb{P} \left(\max_{\ell=1,2} \sup_{\tau \geq 0} |\Delta_{\ell}(\tau)| \geq C\sqrt{t} \right) \leq e^{-t}.$$

Lemma 26. Let $\Xi_{\ell}(\tau) \equiv \|(\Sigma + \tau I)^{-\ell/2} \Sigma^{\ell/2} g\|^2$ for $\ell = 1, 2$. Suppose that $\|\Sigma\|_{\text{op}} \vee \mathcal{H}_{\Sigma} \leq K$ for some $K > 0$. Then there exists some constant $C = C(K) > 1$ such that for $t \geq C \log(en)$,

$$\mathbb{P} \left(\max_{\ell=1,2} \sup_{\tau \geq 0} |\Xi_{\ell}(\tau) - \mathbb{E} \Xi_{\ell}(\tau)| \geq C(\sqrt{nt} + t) \right) \leq e^{-t}.$$

The proofs of these lemmas are deferred to the next subsection.

Proof [Proof of Proposition 24] All the constants in $\lesssim, \gtrsim, \asymp$ and \mathcal{O} below may possibly depend on K . We often suppress the dependence of $\sigma_{\pm}^2(L_w)$ on L_w for simplicity.

(1). We shall write $(\gamma_{n,\eta,\pm}, \tau_{n,\eta,\pm})$ as (γ_n, τ_n) and $(\gamma_{\eta,*}, \tau_{\eta,*}) = (\gamma_*, \tau_*)$ for notational simplicity. Using (38), any solution (γ_n, τ_n) to the equations in (55) satisfies

$$\begin{cases} \phi e_h^2 - \frac{\eta}{\tau_n} + \frac{\Delta_1(\tau_n)}{\sqrt{n}\gamma_n} = \frac{1}{n} \text{tr} \left((\Sigma + \tau_n I)^{-1} \Sigma \right) + \frac{1}{n} (\text{id} - \mathbb{E}) \Xi_1(\tau_n), \\ \phi e_h^2 + \frac{\Delta_2(\tau_n)}{\sqrt{n}\gamma_n} = \frac{1}{\gamma_n^2} \left(\sigma_{\pm}^2 + \tau_n^2 \|(\Sigma + \tau_n I)^{-1} \Sigma^{1/2} \mu_0\|^2 \right) \\ \quad + \frac{1}{n} \text{tr} \left((\Sigma + \tau_n I)^{-2} \Sigma^2 \right) + \frac{1}{n} (\text{id} - \mathbb{E}) \Xi_2(\tau_n). \end{cases} \quad (56)$$

On the event $\mathcal{E}_1(\delta) \cap \mathcal{E}_{\Delta, \Xi}(M)$ with $\delta \in (0, 1/C^{100})$, $1 \leq M \leq \sqrt{n}/C$ and $L_w \leq C$, using Lemma 19, the second equation in (56) becomes

$$\phi + \mathcal{O}(M(1 \vee \gamma_n^{-1})/\sqrt{n} + \delta)$$

$$= \frac{1}{\gamma_n^2} \left(\sigma_{\pm}^2 + \tau_n^2 \|(\Sigma + \tau_n I)^{-1} \Sigma^{1/2} \mu_0\|^2 \right) + \frac{1}{n} \operatorname{tr} \left((\Sigma + \tau_n I)^{-2} \Sigma^2 \right) \gtrsim \frac{1}{\gamma_n^2}.$$

Rearranging terms we obtain the inequality

$$\frac{1}{\gamma_n^2} \lesssim 1 + \frac{M}{\sqrt{n}} + \frac{M}{\sqrt{n}\gamma_n} \Rightarrow \gamma_n \gtrsim \frac{1}{1 + M/\sqrt{n}} \gtrsim 1.$$

So with $\varepsilon_n \equiv \varepsilon_n(M, \delta) \equiv M/\sqrt{n} + \delta$, the equations in (56) reduce to

$$\begin{cases} \phi - \frac{\eta}{\tau_n} + \mathcal{O}(\varepsilon_n) = \frac{1}{n} \operatorname{tr} \left((\Sigma + \tau_n I)^{-1} \Sigma \right), \\ \phi + \mathcal{O}(\varepsilon_n) = \frac{1}{\gamma_n^2} \left(\sigma_{\xi}^2 + \tau_n^2 \|(\Sigma + \tau_n I)^{-1} \Sigma^{1/2} \mu_0\|^2 \right) + \frac{1}{n} \operatorname{tr} \left((\Sigma + \tau_n I)^{-2} \Sigma^2 \right). \end{cases} \quad (57)$$

The above equations match (13) up to the small perturbation $\mathcal{O}(\varepsilon_n) = \mathcal{O}(\varepsilon_n(M, \delta))$ that can be assimilated into the leading term ϕ with small enough $c_0 > 0$ such that $M \leq c_0 \sqrt{n}$. From here the existence (but not uniqueness) and apriori bounds for γ_n, τ_n can be established similarly to the proof of Proposition 23.

(2). Now we shall prove the claimed error bounds. By using (44) and the first equation of (57), we have

$$\frac{1}{n} \operatorname{tr} \left((\Sigma + \tau_n I)^{-1} \Sigma \right) + \frac{\eta}{\tau_n} = \frac{1}{n} \operatorname{tr} \left((\Sigma + \tau_* I)^{-1} \Sigma \right) + \frac{\eta}{\tau_*} + \mathcal{O}(\varepsilon_n).$$

Let $f(\tau) \equiv \frac{1}{n} \operatorname{tr} \left((\Sigma + \tau I)^{-1} \Sigma \right) + \frac{\eta}{\tau}$. Then it is easy to calculate $f'(\tau) = -\frac{1}{n} \operatorname{tr} \left((\Sigma + \tau I)^{-2} \Sigma \right) - \frac{\eta}{\tau^2} \leq 0$, and for any $C_0 > 1$,

$$\inf_{1/C_0 \leq \tau \leq C_0} |f'(\tau)| \geq \inf_{1/C_0 \leq \tau \leq C_0} \frac{1}{n} \operatorname{tr} \left((\Sigma + \tau I)^{-2} \Sigma \right) \geq (\|\Sigma\|_{\text{op}} + C_0)^{-2} \mathcal{H}_{\Sigma}^{-1}.$$

Now using the apriori estimates on τ_*, τ_n , we may conclude

$$\sup_{\eta \in \Xi_K} |\tau_n - \tau_*| \lesssim \varepsilon_n. \quad (58)$$

On the other hand, using (45) and the second equation of (57), we have

$$\begin{aligned} & \frac{1}{\gamma_n^2} \left(\sigma_{\xi}^2 + \tau_n^2 \|(\Sigma + \tau_n I)^{-1} \Sigma^{1/2} \mu_0\|^2 \right) + \frac{1}{n} \operatorname{tr} \left((\Sigma + \tau_n I)^{-2} \Sigma^2 \right) \\ &= \frac{1}{\gamma_*^2} \left(\sigma_{\xi}^2 + \tau_*^2 \|(\Sigma + \tau_* I)^{-1} \Sigma^{1/2} \mu_0\|^2 \right) + \frac{1}{n} \operatorname{tr} \left((\Sigma + \tau_* I)^{-2} \Sigma^2 \right) + \mathcal{O}(\varepsilon_n). \end{aligned} \quad (59)$$

Using the error bound in (58) and apriori estimates for τ_n, τ_* , and the fact that $\mathcal{H}_{\Sigma} \lesssim 1$, by an easy derivative estimate we have

- $\left| \frac{1}{n} \operatorname{tr} \left((\Sigma + \tau_n I)^{-2} \Sigma^2 \right) - \frac{1}{n} \operatorname{tr} \left((\Sigma + \tau_* I)^{-2} \Sigma^2 \right) \right| \lesssim \varepsilon_n$, and
- $\left| \tau_n^2 \|(\Sigma + \tau_n I)^{-1} \Sigma^{1/2} \mu_0\|^2 - \tau_*^2 \|(\Sigma + \tau_* I)^{-1} \Sigma^{1/2} \mu_0\|^2 \right| \lesssim \varepsilon_n$.

Now plugging these estimates into (59), with $\mathcal{C}_0 \equiv \sigma_\xi^2 + \tau_*^2 \|(\Sigma + \tau_* I)^{-1} \Sigma^{1/2} \mu_0\|^2$ satisfying $\mathcal{C}_0 \asymp 1$, we arrive at

$$\frac{\mathcal{C}_0 + \mathcal{O}(\varepsilon_n)}{\gamma_n^2} = \frac{\mathcal{C}_0}{\gamma_*^2} + \mathcal{O}(\varepsilon_n).$$

Using apriori estimates on γ_n, γ_* , we may then invert the above estimate into

$$\sup_{\eta \in \Xi_K} |\gamma_n - \gamma_*| \lesssim \varepsilon_n. \quad (60)$$

The claimed error bounds follow by combining (58) and (60). \blacksquare

B.3 Proofs of Lemmas 25 and 26

Proof [Proof of Lemma 25] We only handle the case $\ell = 1$. The case $\ell = 2$ is similar. Note that the assumption on μ_0 invariant over orthogonal transforms, so for notational simplicity we assume without loss of generality that Σ is diagonal. As $\sup_{\tau \geq Kn} |\Delta_1(\tau)| \leq \left| \sum_{j=1}^n \lambda_j^{1/2} \mu_{0,j} g_j \right| + C e_g \cdot n^{-1/2}$, a standard concentration for the first term shows for $t \geq 1$, with probability $1 - e^{-t}$,

$$\sup_{\tau \geq Kn} |\Delta_1(\tau)| \leq C_0 \sqrt{t}. \quad (61)$$

On the other hand, for $\varepsilon > 0$ to be chosen later, by taking an ε -net \mathcal{S}_ε of $[0, Kn]$, a union bound shows that with probability at least $1 - (Kn/\varepsilon + 1)e^{-t}$,

$$\begin{aligned} \sup_{\tau \in [0, Kn]} |\Delta_1(\tau)| &\leq \max_{\tau \in \mathcal{S}_\varepsilon} |\Delta_1(\tau)| + \sup_{\tau, \tau' \in [0, Kn]: |\tau - \tau'| \leq \varepsilon} |\Delta_1(\tau) - \Delta_1(\tau')| \\ &\leq C_1 \cdot \left(\sqrt{t} + \sqrt{n}(\sqrt{\log n} + \sqrt{t})\varepsilon \right). \end{aligned}$$

Here in the last inequality we used the simple estimate $\sup_{\tau \in [0, Kn]} |\partial_\tau \Delta_1(\tau)| \leq C\sqrt{n} \|\mu_0\| \|g\|_\infty$. Finally by choosing $\varepsilon \equiv \sqrt{t} / \{\sqrt{n}(\sqrt{\log n} + \sqrt{t})\}$, we conclude that for $t \geq C_2 \log(en)$, with probability $1 - e^{-t}$,

$$\sup_{\tau \in [0, Kn]} |\Delta_1(\tau)| \leq C_2 \sqrt{t}. \quad (62)$$

The claim follows by combining (61) and (62). \blacksquare

Proof [Proof of Lemma 26] We focus on the case $\ell = 1$ and will follow a similar idea used in the proof of Lemma 25 above. Similarly we assume Σ is diagonal without loss of generality. All the constants in $\lesssim, \gtrsim, \asymp$ below may depend on K .

First note by a standard concentration, for any $t \geq 1$, with probability at least $1 - e^{-t}$, $\sup_{\tau > Kn} |\Xi_1(\tau)| \lesssim e_g^2 \lesssim 1 + t/n$. Similarly we have $\sup_{\tau > Kn} \mathbb{E} |\Xi_1(\tau)| \lesssim 1$. This means for any $t \geq 1$, with probability at least $1 - e^{-t}$,

$$\sup_{\tau > Kn} \left(|\Xi_1(\tau)| \vee \mathbb{E} |\Xi_1(\tau)| \right) \lesssim 1 + t/n. \quad (63)$$

Next we handle the suprema over $[0, Kn]$ by discretization over an ε -net \mathcal{S}_ε . To this end, we shall establish a pointwise concentration. Note that $\|\nabla \Xi_1(\tau)\|^2 = 4\|(\Sigma + \tau I)^{-1} \Sigma g\|^2 \leq 4\Xi_1(\tau)$. An application of Proposition 59 then yields that, for each $\tau \geq 0$ and $t \geq 1$, with probability at least $1 - e^{-t}$,

$$|\Xi_1(\tau) - \mathbb{E} \Xi_1(\tau)| \leq C(\mathbb{E}^{1/2} \Xi_1(\tau) \cdot \sqrt{t} + t) \lesssim (\sqrt{nt} + t).$$

On the other hand, as $\sup_{\tau \in [0, Kn]} |\partial_\tau \Xi_1(\tau)| \lesssim n \|g\|_\infty^2$ and $\sup_{\tau \in [0, Kn]} |\partial_\tau \mathbb{E} \Xi_1(\tau)| \lesssim n \log n$, we deduce that with probability at least $1 - (Kn/\varepsilon + 1)e^{-t}$,

$$\begin{aligned} \sup_{\tau \in [0, Kn]} |\Xi_1(\tau) - \mathbb{E} \Xi_1(\tau)| &\leq \max_{\tau \in \mathcal{S}_\varepsilon} |\Xi_1(\tau) - \mathbb{E} \Xi_1(\tau)| + \sup_{\tau, \tau' \in [0, Kn]: |\tau - \tau'| \leq \varepsilon} |\Xi_1(\tau) - \Xi_1(\tau')| \\ &\quad + \sup_{\tau, \tau' \in [0, Kn]: |\tau - \tau'| \leq \varepsilon} |\mathbb{E} \Xi_1(\tau) - \mathbb{E} \Xi_1(\tau')| \\ &\lesssim \sqrt{nt} + t + n(\log n + t)\varepsilon. \end{aligned}$$

From here the claim follows by the same arguments used in the proof of Lemma 25 above. \blacksquare

Appendix C. Gaussian designs: Proof of Theorem 3

We assume without loss of generality that $\Sigma = \text{diag}(\lambda_1, \dots, \lambda_n)$, so $V = I$ unless otherwise specified. Recall $\mathcal{H}_\Sigma = \text{tr}(\Sigma^{-1})/n$.

C.1 Localization of the primal problem

Proposition 27. *Suppose $1/K \leq \phi^{-1} - \mathbf{1}_{\eta=0}, \sigma_\xi^2 \leq K$, and $\|\mu_0\| \vee \|\Sigma\|_{\text{op}} \leq K$ for some $K > 0$. Fix $M > 1, \delta \in (0, 1/2)$ and $\eta \geq 0$. On the event $\mathcal{E}_0(M) \cap \mathcal{E}_1(\delta)$, there exists some $C = C(K) > 0$ such that for any deterministic choice of (L_w, L_v) with*

$$L_w \wedge L_v \geq C\{1 + (\|\Sigma^{-1}\|_{\text{op}} M \mathbf{1}_{\phi^{-1} \geq 1+1/K} \wedge \eta^{-1}) \cdot M^2\},$$

we have $\min_{w \in B_n(L_w)} H_\eta(w; L_v) = \min_{w \in \mathbb{R}^n} H_\eta(w)$.

Proof Using the first-order optimality condition for the minimax problem

$$\min_{w \in \mathbb{R}^n} H_\eta(w) = \min_{w \in \mathbb{R}^n} \max_{v \in \mathbb{R}^m} \left\{ \frac{1}{\sqrt{n}} \langle v, Gw - \xi \rangle + F(w) - \frac{\eta \|v\|^2}{2} \right\}, \quad (64)$$

any saddle point (w_*, v_*) of (64) must satisfy $\nabla F(w_*) = -\frac{1}{\sqrt{n}} G^\top v_*$ and $\frac{1}{\sqrt{n}} (Gw_* - \xi) = \eta v_*$, or equivalently,

$$\begin{cases} w_* = -\Sigma^{1/2} \mu_0 + \frac{1}{n} \Sigma G^\top (\phi \check{\Sigma} + \eta I)^{-1} (G \Sigma^{1/2} \mu_0 + \xi), \\ v_* = -\frac{1}{\sqrt{n}} (\phi \check{\Sigma} + \eta I)^{-1} (G \Sigma^{1/2} \mu_0 + \xi). \end{cases}$$

Here recall $\check{\Sigma} = m^{-1} G \Sigma G^\top$. On the event $\mathcal{E}_0(M)$,

$$\|(\phi \check{\Sigma} + \eta I)^{-1}\|_{\text{op}} \lesssim_K \|\Sigma^{-1}\|_{\text{op}} M \mathbf{1}_{\phi^{-1} \geq 1+1/K} \wedge \eta^{-1}.$$

So on $\mathcal{E}_0(M) \cap \mathcal{E}_1(\delta)$,

$$\|w_*\| \vee \|v_*\| \lesssim_K 1 + (\|\Sigma^{-1}\|_{\text{op}} M \mathbf{1}_{\phi^{-1} \geq 1+1/K} \wedge \eta^{-1}) M^2.$$

This means that on the event $\mathcal{E}_0(M) \cap \mathcal{E}_1(\delta)$, for any L_w, L_v chosen as in the statement of the lemma,

$$\min_{w \in \mathbb{R}^n} H_\eta(w) = \min_{w \in B_n(L_w)} \max_{v \in B_m(L_v)} \left\{ \frac{1}{\sqrt{n}} \langle v, Gw - \xi \rangle + F(w) - \frac{\eta \|v\|^2}{2} \right\}.$$

The proof is complete by recalling the definition of $H_\eta(\cdot; L_v)$. ■

C.2 Characterization of the Gordon cost optimum

Theorem 28. *Suppose the following hold for some $K > 0$.*

- $1/K \leq \phi^{-1} \leq K$, $\|\mu_0\| \vee \|\Sigma\|_{\text{op}} \vee \mathcal{H}_\Sigma \leq K$.
- *Assumption B* with $\sigma_\xi^2 \in [1/K, K]$.

There exist some $C, C' > 1$ depending on K such that for any deterministic choice of $L_w, L_v \in [C, C^2]$, it holds for any $C' \log(en) \leq t \leq n/C'$, $\eta \in \Xi_K$ and $\xi \in \mathcal{E}_{1,\xi}(\sqrt{t/n})$,

$$\mathbb{P}^\xi \left(\left| \min_{w \in B_n(L_w)} L_\eta(w; L_v) - \max_{\beta > 0} \min_{\gamma > 0} \bar{D}_\eta(\beta, \gamma) \right| \geq \sqrt{t/n} \right) \leq C e^{-t/C}.$$

In the next subsection we will show that for large $L_v > 0$, the map $w \mapsto L_\eta(w; L_v)$ attains its global minimum in an ℓ_2 ball of constant order radius (under $\mathcal{H}_\Sigma \lesssim 1$) with high probability. This means that although the initial localization radius for the primal optimization may be highly suboptimal (which involves $\|\Sigma^{-1}\|_{\text{op}}$), the Gordon objective can be further localized into an ℓ_2 ball with constant order radius.

To prove Theorem 28, we shall first relate $\min_{w \in B_n(L_w)} L_\eta(w; L_v)$ to $\max_{\beta > 0} \min_{\gamma > 0} D_{\eta,\pm}(\beta, \gamma)$ and its localized versions.

Proposition 29. *Suppose $1/K \leq \phi^{-1}, \sigma_\xi^2 \leq K$, and $\|\mu_0\| \vee \|\Sigma\|_{\text{op}} \vee \mathcal{H}_\Sigma \leq K$ for some $K > 0$. There exists constant $C = C(K) > 1$ such that for any deterministic choice of $L_w, L_v \in [C, C^2]$, on the event $\mathcal{E}_1(\delta) \cap \mathcal{E}_{\Delta,\Xi}(M)$ (defined in Proposition 24) with $\delta \in (0, 1/C^{100})$ and $M \leq \sqrt{n}/C$, we have for any $\eta \in \Xi_K$,*

$$\max_{\beta > 0} \min_{\gamma > 0} D_{\eta,-}(\beta, \gamma) \leq \min_{w \in B_n(L_w)} L_\eta(w; L_v) \leq \max_{\beta > 0} \min_{\gamma > 0} D_{\eta,+}(\beta, \gamma),$$

and the following localization holds:

$$\max_{\beta > 0} \min_{\gamma > 0} D_{\eta,\pm}(\beta, \gamma) = \max_{1/C \leq \beta \leq C} \min_{1/C \leq \gamma \leq C} D_{\eta,\pm}(\beta, \gamma).$$

Proof We write $g_n \equiv g/\sqrt{n}$ in the proof.

(Step 1). Fix any $L_w, L_v > 0$. We may compute

$$\begin{aligned}
 & \min_{w \in B_n(L_w)} L_\eta(w; L_v) \\
 &= \min_{w \in B_n(L_w)} \max_{\beta \in [0, L_v]} \left\{ \frac{\beta}{\sqrt{n}} \left(\| \|w\|h - \xi \| - \langle g, w \rangle \right) + F(w) - \frac{\eta\beta^2}{2} \right\} \\
 &= \max_{\beta \in [0, L_v]} \min_{\gamma > 0} \left\{ \frac{\beta\gamma \|h\|^2}{2n} - \frac{\eta\beta^2}{2} + \min_{w \in B_n(L_w)} \left(\frac{\beta}{2\gamma} \frac{\| \|w\|h - \xi \|^2}{\|h\|^2} - \langle w, \beta g_n \rangle + F(w) \right) \right\}.
 \end{aligned} \tag{65}$$

Here in the last line we used Sion's min-max theorem to flip the order of minimum and maximum in $\min_{w \in B_n(L_w)} \max_{\beta \in [0, L_v]}$. The minimum over γ is achieved exactly at $\frac{\| \|w\|h - \xi \| / \|h\|}{\|h\| / \sqrt{n}}$, so when $\sigma_-^2 \neq 0$, using the simple inequality

$$\|w\|^2 + \sigma_-^2 \leq \| \|w\|h - \xi \|^2 / \|h\|^2 \leq \|w\|^2 + \sigma_+^2, \tag{66}$$

on the event $\mathcal{E}_1(\delta)$, we may further bound (65) as follows:

$$\begin{aligned}
 & \pm \min_{w \in B_n(L_w)} L_\eta(w; L_v) \leq \pm \max_{\beta \in [0, L_v]} \min_{\gamma > 0} \\
 & \left\{ \frac{\beta\gamma \|h\|^2}{2n} - \frac{\eta\beta^2}{2} + \min_{w \in B_n(L_w)} \left(\frac{\beta}{2\gamma} (\|w\|^2 + \sigma_\pm^2(L_w)) - \langle w, \beta g_n \rangle + F(w) \right) \right\}.
 \end{aligned} \tag{67}$$

We note that σ_\pm^2 depends on L_w , but this notational dependence will be dropped from now on for convenience.

(Step 2). Consider the minimax optimization problem in (67):

$$\begin{aligned}
 & \max_{\beta > 0} \min_{\gamma > 0, w \in \mathbb{R}^n} \left\{ \frac{\beta\gamma \|h\|^2}{2n} - \frac{\eta\beta^2}{2} + \left(\frac{\beta}{2\gamma} (\|w\|^2 + \sigma_\pm^2) - \langle w, \beta g_n \rangle + F(w) \right) \right\} \\
 &= \max_{\beta > 0} \min_{\gamma > 0} \left\{ \frac{\beta}{2} \left(\gamma(\phi e_h^2 - e_g^2) + \frac{\sigma_\pm^2}{\gamma} \right) - \frac{\eta\beta^2}{2} + \mathbf{e}_F(\gamma g_n; \gamma/\beta) \right\}.
 \end{aligned} \tag{68}$$

Any saddle point $(\beta_{n,\eta,\pm}, \gamma_{n,\eta,\pm}, w_{n,\eta,\pm}) = (\beta_{n,\pm}, \gamma_{n,\pm}, w_{n,\pm})$ of the above program must satisfy the first-order optimality condition

$$\begin{cases} 0 = \frac{1}{2}(\gamma_{n,\pm}(\phi e_h^2 - e_g^2) + \frac{\sigma_\pm^2}{\gamma_{n,\pm}}) - \eta\beta_{n,\pm} + \partial_\beta \mathbf{e}_F(\gamma_{n,\pm} g_n; \gamma_{n,\pm}/\beta_{n,\pm}), \\ 0 = \frac{\beta_{n,\pm}}{2}((\phi e_h^2 - e_g^2) - \frac{\sigma_\pm^2}{\gamma_{n,\pm}}) + \partial_\gamma \mathbf{e}_F(\gamma_{n,\pm} g_n; \gamma_{n,\pm}/\beta_{n,\pm}), \\ w_{n,\pm} = \mathbf{prox}_F(\gamma_{n,\pm} g_n; \gamma_{n,\pm}/\beta_{n,\pm}). \end{cases} \tag{69}$$

Using the derivative formula in Lemma 17 and the form of \mathbf{prox}_F in Lemma 16, we may compute

$$\begin{cases} \partial_\beta \mathbf{e}_F(\gamma g_n; \gamma/\beta) = \frac{1}{2\gamma} \left(\mathbf{err}_{(\Sigma, \mu_0)}(\gamma; \gamma/\beta) - 2 \mathbf{dof}_{(\Sigma, \mu_0)}(\gamma; \gamma/\beta) + \gamma^2 e_g^2 \right), \\ \partial_\gamma \mathbf{e}_F(\gamma g_n; \gamma/\beta) = \frac{\beta}{2\gamma^2} \left(\gamma^2 e_g^2 - \mathbf{err}_{(\Sigma, \mu_0)}(\gamma; \gamma/\beta) \right). \end{cases} \tag{70}$$

$$\partial_\tau \mathbf{e}_F(\gamma g_n; \tau) = -\frac{1}{2\tau^2} \|\Sigma^{1/2} \hat{\mu} - y\|^2 + \frac{1}{\tau} \langle \Sigma^{1/2} \hat{\mu} - y, \Sigma^{1/2} \partial_\tau \hat{\mu} \rangle + \langle \hat{\mu}, \partial_\tau \hat{\mu} \rangle,$$

on the event $\mathcal{E}_1(\delta)$, we may estimate $|\partial_\gamma \mathbf{e}_F(\gamma g_n; \tau)| \vee |\partial_\tau \mathbf{e}_F(\gamma g_n; \tau)| \lesssim 1$. A similar estimate applies to the expectation versions, proving (74).

(Step 2). Next we show that for any $C_0 > 1$, there exists $C_1 > 0$ such that for $t \geq C_1 \log(en)$,

$$\mathbb{P} \left(\sup_{(\gamma, \tau) \in [C_0^{-1}, C_0]^2} |(\text{id} - \mathbb{E}) \mathbf{e}_F(\gamma g_n; \tau)| \geq C_1 (\sqrt{t/n} + t/n), \mathcal{E}_1(\delta) \right) \leq C_1 e^{-t/C_1}. \quad (75)$$

To prove the claim, we fix $\varepsilon > 0$ to be chosen later, and take an ε -net $\mathcal{S}(\varepsilon)$ for $[1/C_0, C_0]$. Then $|\mathcal{S}(\varepsilon)| \leq C_0/\varepsilon + 1$. So on the event $\mathcal{E}_1(\delta)$, using the estimate in (74) and a union bound via the pointwise concentration inequality in Proposition 18, for $t \geq 1$, with probability at least $1 - C\varepsilon^{-2}e^{-t/C}$,

$$\sup_{(\gamma, \tau) \in [C_0^{-1}, C_0]^2} |(\text{id} - \mathbb{E}) \mathbf{e}_F(\gamma g_n; \tau)| \lesssim \sup_{\gamma, \tau \in \mathcal{S}(\varepsilon)} |(\text{id} - \mathbb{E}) \mathbf{e}_F(\gamma g_n; \tau)| + C\varepsilon \lesssim \sqrt{\frac{t}{n}} + \frac{t}{n} + \varepsilon.$$

Here in the last inequality we used Lemma 16 to estimate $\sup_{(\gamma, \tau)} v^2(\gamma, \tau) \vee \sup_{(\gamma, \tau)} v^2(\gamma, \tau) \mathbb{E} \mathbf{e}_F(\gamma g_n; \tau) \lesssim 1$, where $v^2(\gamma, \tau)$ is defined in Proposition 18. The claim (75) follows by choosing $\varepsilon \equiv \sqrt{t/n} + t/n$ and some calculations.

(Step 3). By (75), for $t \geq C \log(en)$, on the event $\mathcal{E}_1(\delta)$, it holds with probability at least $1 - C_2 e^{-t/C_2}$ that

$$\begin{aligned} & \max_{1/C \leq \beta \leq C} \min_{1/C \leq \gamma \leq C} \mathbf{D}_{\eta, \pm}(\beta, \gamma) \\ &= \max_{1/C \leq \beta \leq C} \min_{1/C \leq \gamma \leq C} \left\{ \frac{\beta}{2} \left(\gamma(\phi e_h^2 - e_g^2) + \frac{\sigma_{\pm}^2}{\gamma} \right) - \frac{\eta \beta^2}{2} + \mathbf{e}_F(\gamma g_n; \gamma/\beta) \right\} \\ &= \max_{1/C \leq \beta \leq C} \min_{1/C \leq \gamma \leq C} \bar{\mathbf{D}}_\eta(\beta, \gamma) + \mathcal{O}(\sqrt{t/n} + t/n + \delta). \end{aligned}$$

The estimate in \mathcal{O} is uniform in $\eta \in \Xi_K$, so the claim follows. \blacksquare

Finally we delocalize the range constraints for β, γ in the deterministic minimax problem with $\bar{\mathbf{D}}_\eta$ in the above proposition.

Proposition 31. *Suppose $1/K \leq \phi^{-1}, \sigma_\xi^2 \leq K$, and $\|\mu_0\| \vee \|\Sigma\|_{\text{op}} \vee \mathcal{H}_\Sigma \leq K$ for some $K > 0$. There exists some $C = C(K) > 1$ such that for any $\eta \in \Xi_K$,*

$$\max_{\beta > 0} \min_{\gamma > 0} \bar{\mathbf{D}}_\eta(\beta, \gamma) = \max_{1/C \leq \beta \leq C} \min_{1/C \leq \gamma \leq C} \bar{\mathbf{D}}_\eta(\beta, \gamma).$$

Consequently,

$$\left| \max_{\beta > 0} \min_{\gamma > 0} \bar{\mathbf{D}}_\eta(\beta, \gamma) - \max_{\beta > 0} \min_{\gamma > 0} \bar{\mathbf{D}}_0(\beta, \gamma) \right| \leq C\eta. \quad (76)$$

Proof The proof is essentially a deterministic version of Step 2 in the proof of Proposition 29. We give some details below. We write $g_n \equiv g/\sqrt{n}$. First, using similar calculations as that of (70),

$$\begin{cases} \partial_\beta \mathbb{E} \mathbf{e}_F(\gamma g_n; \gamma/\beta) = \frac{1}{2\gamma} \left(\mathbb{E} \mathbf{err}_{(\Sigma, \mu_0)}(\gamma; \gamma/\beta) - 2 \mathbb{E} \mathbf{dof}_{(\Sigma, \mu_0)}(\gamma; \gamma/\beta) + \gamma^2 \right), \\ \partial_\gamma \mathbb{E} \mathbf{e}_F(\gamma g_n; \gamma/\beta) = \frac{\beta}{2\gamma^2} \left(\gamma^2 - \mathbb{E} \mathbf{err}_{(\Sigma, \mu_0)}(\gamma; \gamma/\beta) \right). \end{cases}$$

Then the first-order optimality condition for (β_*, γ_*) to be the saddle point of $\max_{\beta>0} \min_{\gamma>0} \bar{\mathbf{D}}_\eta(\beta, \gamma)$, i.e., a deterministic version of (71), is given by

$$\begin{cases} \phi \gamma_*^2 = \sigma_\xi^2 + \mathbb{E} \mathbf{err}_{(\Sigma, \mu_0)}(\gamma_*; \gamma_*/\beta_*), \\ \left(\phi - \frac{\eta}{\gamma_*/\beta_*} \right) \gamma_*^2 = \mathbb{E} \mathbf{dof}_{(\Sigma, \mu_0)}(\gamma_*; \gamma_*/\beta_*). \end{cases}$$

Finally using the apriori estimates in Proposition 23, we obtain a deterministic analogue of (72) in that $\gamma_* \asymp_K 1$, $\beta_* \asymp_K 1$. The claimed localization follows. The continuity follows by the definition of $\bar{\mathbf{D}}_\eta$ and the proven localization. \blacksquare

Proof [Proof of Theorem 28] By Propositions 29, 30 and 31, there exist $C, C' > 0$ such that for any $\delta \in (0, 1/C^{100})$, $M \leq \sqrt{n}/C$, $t \geq C' \log(en)$, $\xi \in \mathcal{E}_{1,\xi}(\delta)$ and $\eta \in \Xi_K$,

$$\begin{aligned} & \mathbb{P}^\xi \left[\left| \min_{w \in B_n(L_w)} L_\eta(w; L_v) - \max_{\beta>0} \min_{\gamma>0} \bar{\mathbf{D}}_\eta(\beta, \gamma) \right| \geq C(\sqrt{t/n} + t/n + \delta) \right] \\ & \leq C e^{-t/C} + \mathbb{P}^\xi \left(\mathcal{E}_{1,0}(\delta)^c \right) + \mathbb{P} \left(\mathcal{E}_{\Delta, \Xi}(M)^c \right). \end{aligned}$$

The claim now follows from the concentration estimates in Lemmas 21, 25 and 26, by choosing $M \equiv \sqrt{n}/C$ and $\delta \equiv C(\sqrt{t/n} + t/n)$, which is valid in the regime $t \leq n/C_0$ for large C_0 . \blacksquare

C.3 Locating the global minimizer of the Gordon objective

With $(\gamma_{\eta,*}, \tau_{\eta,*})$ denoting the unique solution to the system of equations (13), let

$$w_{\eta,*} \equiv \mathbf{prox}_F \left(\gamma_{\eta,*} g / \sqrt{n}; \tau_{\eta,*} \right) = \Sigma^{1/2} \left(\widehat{\mu}_{(\Sigma, \mu_0)}^{\text{seq}}(\gamma_{\eta,*}; \tau_{\eta,*}) - \mu_0 \right). \quad (77)$$

For any $\varepsilon > 0$, let the exceptional set be defined as

$$D_{\eta;\varepsilon}(\mathbf{g}) \equiv \{ w \in \mathbb{R}^n : |\mathbf{g}(w) - \mathbb{E} \mathbf{g}(w_{\eta,*})| \geq \varepsilon \}. \quad (78)$$

Theorem 32. *Suppose the following hold for some $K > 0$.*

- $1/K \leq \phi^{-1} \leq K$, $\|\mu_0\| \vee \|\Sigma\|_{\text{op}} \vee \mathcal{H}_\Sigma \leq K$.
- Assumption B holds with $\sigma_\xi^2 \in [1/K, K]$.

Fix any $\mathbf{g} : \mathbb{R}^n \rightarrow \mathbb{R}$ that is 1-Lipschitz with respect to $\|\cdot\|_{\Sigma^{-1}}$. There exist constants $C, C' > 10$ depending on K such that for $L_w, L_v \in [C, C^2]$, $C' \log(en) \leq t \leq n/C'$, $\xi \in \mathcal{E}_{1,\xi}(\sqrt{t/n})$ and $\eta \in \Xi_K$,

$$\mathbb{P}^\xi \left(\min_{w \in D_{\eta; C(t/n)^{1/4}}(\mathbf{g}) \cap B_n(L_w)} L_\eta(w; L_v) \leq \max_{\beta > 0} \min_{\gamma > 0} \bar{D}_\eta(\beta, \gamma) + \sqrt{t/n} \right) \leq C e^{-t/C}.$$

Roughly speaking, the above theorem will be proved by approximating L_η both from above and below by nicer strongly convex, surrogate functions whose minimizers can be directly located. Then we may relate the minimizer of L_η and those of the surrogate functions.

We first formally define these surrogate functions. For $L_w > 0, L_v > 0$, let

$$L_{\eta,\pm}(w; L_v) \equiv \max_{\beta \in [0, L_v]} \left\{ \frac{\beta}{\sqrt{n}} \left(\|h\| \sqrt{\|w\|^2 + \sigma_\pm^2(L_w)} - \langle g, w \rangle \right) - \frac{\eta\beta^2}{2} + F(w) \right\}. \quad (79)$$

Again we omit notational dependence of $L_{\eta,\pm}$ on L_w for simplicity.

The following lemma provides uniform (bracketing) approximation of L_η via $L_{\eta,\pm}$ on compact sets.

Lemma 33. *Fix $L_v > 0$. The following hold when $\sigma_-^2(L_w) \neq 0$.*

1. For any $w \in B_n(L_w)$, $L_{\eta,-}(w; L_v) \leq L_\eta(w; L_v) \leq L_{\eta,+}(w; L_v)$.
2. For any $L_w > 0$,

$$\sup_{w \in \mathbb{R}^n} |L_{\eta,+}(w; L_v) - L_{\eta,-}(w; L_v)| \leq \frac{4e_h}{\sigma_m} \cdot L_v L_w \frac{|\langle h, \xi \rangle|}{\|h\|^2}.$$

Proof The first claim (1) follows by the definition of $\sigma_\pm^2(L_w)$ in (30) and the simple inequality (66). For (2), note that

$$\begin{aligned} |L_{\eta,+}(w; L_v) - L_{\eta,-}(w; L_v)| &\leq L_v e_h \cdot \left| \sqrt{\|w\|^2 + \sigma_+^2(L_w)} - \sqrt{\|w\|^2 + \sigma_-^2(L_w)} \right| \\ &\leq L_v e_h \cdot \frac{|\sigma_+^2(L_w) - \sigma_-^2(L_w)|}{\sigma_+(L_w) + \sigma_-(L_w)} \leq \frac{4e_h}{\sigma_m} \cdot L_v L_w \frac{|\langle h, \xi \rangle|}{\|h\|^2}, \end{aligned}$$

as desired. ■

Next, we will study the properties of the global minimizers for $L_{\eta,\pm}$.

Proposition 34. *Suppose $1/K \leq \phi^{-1}, \sigma_\xi^2 \leq K$, and $\|\mu_0\| \vee \|\Sigma\|_{\text{op}} \vee \mathcal{H}_\Sigma \leq K$ for some $K > 0$. There exists some constant $C = \bar{C}(K) > 1$ such that for any deterministic choice of $L_w, L_v \in [C, C^2]$, on the event $\mathcal{E}_1(\delta) \cap \mathcal{E}_{\Delta, \Xi}(M)$ (defined in Proposition 24) with $\delta \in (0, 1/C^{100})$ and $M \leq \sqrt{n}/C$, for any $\eta \in \Xi_K$, the maps $w \mapsto L_{\eta,\pm}(w; L_v)$ attain its global minimum at $w_{n,\eta,\pm}$ with $\|w_{n,\eta,\pm}\|_{\Sigma^{-1}} \leq C$. Moreover, $\|w_{n,\eta,\pm} - w_{\eta,*}\|_{\Sigma^{-1}} \leq C(M/\sqrt{n} + \delta)^{1/2}$.*

Proof Note that the optimization problem

$$\begin{aligned}
 & \min_{w \in \mathbb{R}^n} L_{\eta, \pm}(w; L_v) \\
 &= \min_{w \in \mathbb{R}^n} \max_{\beta \in [0, L_v]} \left\{ \frac{\beta}{\sqrt{n}} \left(\|h\| \sqrt{\|w\|^2 + \sigma_{\pm}^2(L_w)} - \langle g, w \rangle \right) - \frac{\eta\beta^2}{2} + F(w) \right\} \\
 &\stackrel{(*)}{=} \max_{\beta \in [0, L_v]} \min_{w \in \mathbb{R}^n} \left\{ \frac{\beta}{\sqrt{n}} \left(\|h\| \sqrt{\|w\|^2 + \sigma_{\pm}^2(L_w)} - \langle g, w \rangle \right) - \frac{\eta\beta^2}{2} + F(w) \right\} \\
 &= \max_{\beta \in [0, L_v]} \min_{\gamma > 0, w \in \mathbb{R}^n} \left\{ \frac{\beta\gamma\|h\|^2}{2n} - \frac{\eta\beta^2}{2} + \left(\frac{\beta}{2\gamma} (\|w\|^2 + \sigma_{\pm}^2(L_w)) - \left\langle w, \frac{\beta}{\sqrt{n}}g \right\rangle + F(w) \right) \right\}.
 \end{aligned}$$

Here in (*) we used Sion's min-max theorem to exchange minimum and maximum, as the maximum is taken over a compact set. The difference of the above minimax problem compared to (68) rests in its range constraint on β . As proven in (72), all solutions $\beta_{n, \pm}$ to the unconstrained minimax problem (68) must satisfy $\beta_{n, \pm} \leq C$ on the event $\mathcal{E}_1(\delta) \cap \mathcal{E}_{\Delta, \Xi}(M)$. So on this event, for the choice $L_w, L_v \in [C, C^2]$ for some large $C > 0$, $\min_{w \in \mathbb{R}^n} L_{\eta, \pm}(w; L_v)$ exactly corresponds to (68), whose minimizers $w_{n, \pm}$ admit the apriori estimate (73) (with minor modifications that change $\|\cdot\|$ to the stronger estimate in $\|\cdot\|_{\Sigma^{-1}}$).

Next, for the error bound, using the last equation in (69) and the definition of $w_{\eta, *}$ in (77), along with the estimates in Proposition 24, we have

$$\begin{aligned}
 \|w_{n, \eta, \pm} - w_{\eta, *}\|_{\Sigma^{-1}}^2 &= \left\| \widehat{\mu}_{(\Sigma, \mu_0)}^{\text{seq}}(\gamma_{n, \eta, \pm}; \tau_{n, \eta, \pm}) - \widehat{\mu}_{(\Sigma, \mu_0)}^{\text{seq}}(\gamma_{\eta, *}; \tau_{\eta, *}) \right\|^2 \\
 &\lesssim_K |\gamma_{n, \eta, \pm} - \gamma_{\eta, *}| \vee |\tau_{n, \eta, \pm} - \tau_{\eta, *}| \leq C(M/\sqrt{n} + \delta),
 \end{aligned}$$

as desired. ■

Finally we shall relate back to the global minimizer of L_{η} . We note that the proposition below by itself is not formally used in the proof of Theorem 32, but will turn out to be useful in the proof of Theorem 3 ahead.

Proposition 35. *Suppose the conditions in Theorem 32 hold for some $K > 0$. There exist constants $C, C' > 1$ depending on K such that for $L_w, L_v \in [C, C^2]$, $C' \log(en) \leq t \leq n/C'$, $\xi \in \mathcal{E}_{1, \xi}(\sqrt{t/n})$ and $\eta \in \Xi_K$,*

$$\begin{aligned}
 & \mathbb{P}^{\xi} \left(\text{The map } w \mapsto L_{\eta}(w; L_v) \text{ attains its global minimum at } w_{n, \eta} \text{ with } \|w_{n, \eta}\|_{\Sigma^{-1}} \leq C, \right. \\
 & \left. \text{and } \|w_{n, \eta} - w_{\eta, *}\|_{\Sigma^{-1}} \leq C(t/n)^{1/4} \right) \geq 1 - Ce^{-t/C}.
 \end{aligned}$$

Proof Let us fix $\xi \in \mathcal{E}_{1, \xi}(\sqrt{t/n})$.

(Step 1). We first prove the apriori estimate for $\|w_{n, \eta}\|_{\Sigma^{-1}}$. To this end, for large enough $C_0, C'_0 > 0$ depending on K , we choose $L_w \equiv C_0, \delta \equiv 1/C_0^{100}$ and $M \equiv \delta\sqrt{n}$ in Proposition 34, it follows that

$$\begin{aligned}
 & \mathbb{P}^{\xi} \left(E_1 \equiv \left\{ \|w_{n, \eta, \pm}\|_{\Sigma^{-1}} \vee \|w_{n, \eta, \pm}\| \leq C_0/2, \right. \right. \\
 & \left. \left. L_{\eta, \pm}(w_{n, \eta, \pm}; L_v) = (68) \right\} \geq 1 - C_0 e^{-n/C_0}. \right. \quad (80)
 \end{aligned}$$

On the other hand, choosing $\delta \equiv \sqrt{t/n}$ with $C'_0 \log(en) \leq t \leq n/C'_0$ leads to

$$\mathbb{P}^\xi \left(E_2(t) \equiv \left\{ \|w_{n,\eta,\pm} - w_{\eta,*}\|_{\Sigma^{-1}} \leq C_0(t/n)^{1/4} \right\} \right) \geq 1 - C_0 e^{-t/C_0}. \quad (81)$$

On E_1 , we may characterize the value of $L_{\eta,\pm}(w_{n,\eta,\pm}; L_v)$ by applying Propositions 29-31: for $C'_0 \log(en) \leq t \leq n/C'_0$,

$$\mathbb{P}^\xi \left(E_3(t) \equiv \left\{ \left| L_{\eta,\pm}(w_{n,\eta,\pm}; L_v) - \max_{\beta>0} \min_{\gamma>0} \bar{D}_\eta(\beta, \gamma) \right| \leq C_0 \sqrt{t/n} \right\} \right) \geq 1 - C_0 e^{-t/C_0}. \quad (82)$$

Note by the strong convexity of $L_{\eta,\pm}(\cdot; L_v)$ with respect to $\|\cdot\|_{\Sigma^{-1}}$, we have

$$\inf_{w \in \mathbb{R}^n: \|w - w_{n,\eta,\pm}\|_{\Sigma^{-1}} \geq \sqrt{6C_0}(t/n)^{1/4}} L_{\eta,\pm}(w; L_v) - L_{\eta,\pm}(w_{n,\eta,\pm}; L_v) \geq 3C_0 \sqrt{t/n}.$$

This means on $E_3(t)$,

$$\begin{aligned} \inf_{w \in \mathbb{R}^n: \|w - w_{n,\eta,\pm}\|_{\Sigma^{-1}} \geq \sqrt{6C_0}(t/n)^{1/4}} L_{\eta,\pm}(w; L_v) &\geq \max_{\beta>0} \min_{\gamma>0} \bar{D}_\eta(\beta, \gamma) + 2C_0 \sqrt{t/n}, \\ L_{\eta,\pm}(w_{n,\eta,\pm}; L_v) &\leq \max_{\beta>0} \min_{\gamma>0} \bar{D}_\eta(\beta, \gamma) + C_0 \sqrt{t/n}. \end{aligned}$$

This in particular means on $E_1 \cap E_3(t)$,

$$\begin{aligned} w_{n,\eta,\pm} &\in \left\{ w \in \mathbb{R}^n : L_{\eta,\pm}(w; L_v) \leq \max_{\beta>0} \min_{\gamma>0} \bar{D}_\eta(\beta, \gamma) + C_0 \sqrt{t/n} \right\} \\ &\subset \left\{ w \in \mathbb{R}^n : \|w\|_{\Sigma^{-1}} \leq \sqrt{6C_0}(t/n)^{1/4} + C_0/2 \right\}. \end{aligned}$$

Consequently, by enlarging $C_0 > 0$ if necessary, using Lemma 33-(1), on $E_1 \cap E_3(t)$

$$\begin{aligned} &\left\{ w \in \mathbb{R}^n : L(w; L_v) \leq \max_{\beta>0} \min_{\gamma>0} \bar{D}_\eta(\beta, \gamma) + C_0 \sqrt{t/n} \right\} \\ &\subset \left\{ w \in \mathbb{R}^n : \|w\|_{\Sigma^{-1}} \leq 3C_0/5 \right\} \subset B_n(3C_0/4) \subsetneq B_n(C_0) = B_n(L_w). \end{aligned}$$

This implies, on $E_1 \cap E_3(t)$, we have $\|w_{n,\eta}\|_{\Sigma^{-1}} \vee \|w_{n,\eta}\| \leq 3C_0/4$, proving the apriori bound. **(Step 2)**. Next we establish the announced error bound. On the event $\mathcal{E}_{1,0}(\sqrt{t/n})$, by Lemma 33-(2),

$$\sup_{w \in B_n(C_0)} \left| L_\eta(w; L_v) - L_{\eta,\pm}(w; L_v) \right| \leq C_1 \sqrt{t/n}. \quad (83)$$

Consequently, on $E_1 \cap E_3(t) \cap \mathcal{E}_{1,0}(\sqrt{t/n})$,

$$\left| \min_{w \in \mathbb{R}^n} L_\eta(w; L_v) - \max_{\beta>0} \min_{\gamma>0} \bar{D}_\eta(\beta, \gamma) \right| \leq C_2 \sqrt{t/n}. \quad (84)$$

On this event, combining (83)-(84) with (82), and using again the strong convexity of $L_{\eta,+}(\cdot; L_v)$ respect to $\|\cdot\|_{\Sigma^{-1}}$, we have for $C_3 = 2\sqrt{(C_0 + C_1 + C_2)}$,

$$\inf_{w \in B_n(C_0): \|w - w_{n,\eta,+}\|_{\Sigma^{-1}} \geq C_3(t/n)^{1/4}} L_\eta(w; L_v) - \min_{w \in \mathbb{R}^n} L_\eta(w; L_v)$$

$$\begin{aligned}
 &\geq \inf_{w \in B_n(C_0): \|w - w_{n,\eta,+}\|_{\Sigma^{-1}} \geq C_3(t/n)^{1/4}} L_{\eta,+}(w; L_v) - \max_{\beta > 0} \min_{\gamma > 0} \bar{D}_\eta(\beta, \gamma) - (C_1 + C_2)\sqrt{t/n} \\
 &\geq \inf_{w \in B_n(C_0): \|w - w_{n,\eta,+}\|_{\Sigma^{-1}} \geq C_3(t/n)^{1/4}} L_{\eta,+}(w; L_v) - L_{\eta,+}(w_{n,+}; L_v) - (C_0 + C_1 + C_2)\sqrt{t/n} \\
 &\geq (C_3^2/2)\sqrt{t/n} - (C_0 + C_1 + C_2)\sqrt{t/n} = (C_0 + C_1 + C_2)\sqrt{t/n}.
 \end{aligned}$$

This means that $\|w_{n,\eta} - w_{n,\eta,+}\|_{\Sigma^{-1}} \leq C_3(t/n)^{1/4}$ on $E_1 \cap E_3(t) \cap \mathcal{E}_{1,0}(\sqrt{t/n})$. The claim follows by intersecting the prescribed event with $E_2(t)$ in (81) that controls the \mathbb{P}^ξ -probability of $\|w_{n,\eta,+} - w_{\eta,*}\|_{\Sigma^{-1}} \leq C_0(t/n)^{1/4}$. \blacksquare

Proof [Proof of Theorem 32] Fix $\xi \in \mathcal{E}_{1,\xi}(\sqrt{t/n})$, and $\varepsilon > 0$ to be chosen later on. First, as \mathbf{g} is Lipschitz with respect to $\|\cdot\|_{\Sigma^{-1}}$, by the Gaussian concentration inequality, there exists $C_0 = C_0(K) > 0$ such that for $t \geq 1$, on an event $E_0(t)$ with \mathbb{P}^ξ -probability at least $1 - e^{-t}$,

$$|\mathbf{g}(w_{\eta,*}) - \mathbb{E} \mathbf{g}(w_{\eta,*})| \leq C_0 \sqrt{t/n}.$$

Moreover, by Proposition 34 and Propositions 29-31, there exist some $C_1, C'_1 > 0$ depending on K such that for $C'_1 \log(en) \leq t \leq n/C'_1$, on an event $E_1(t)$ with \mathbb{P}^ξ -probability $1 - C_1 e^{-t/C_1}$, we have

1. $\|w_{n,\eta,-}\|_{\Sigma^{-1}} \vee \|w_{n,\eta,-}\| \leq C_1$, $\|w_{n,\eta,-} - w_{\eta,*}\|_{\Sigma^{-1}} \leq C_1(t/n)^{1/4}$, and
2. $|L_{\eta,-}(w_{n,\eta,-}; L_v) - \max_{\beta > 0} \min_{\gamma > 0} \bar{D}_\eta(\beta, \gamma)| \leq C_1 \sqrt{t/n}$.

Consequently, for $C'_1 \log(en) \leq t \leq n/C'_1$, on the event $E_0(t) \cap E_1(t)$, uniformly in $w \in D_{\eta;\varepsilon}(\mathbf{g}) \cap B_n(L_w)$,

$$\begin{aligned}
 \varepsilon &\leq |\mathbf{g}(w) - \mathbb{E} \mathbf{g}(w_{\eta,*})| \leq |\mathbf{g}(w) - \mathbf{g}(w_{\eta,*})| + |\mathbf{g}(w_{\eta,*}) - \mathbb{E} \mathbf{g}(w_{\eta,*})| \\
 &\leq \|w - w_{\eta,n,-}\|_{\Sigma^{-1}} + \|w_{\eta,n,-} - w_{\eta,*}\|_{\Sigma^{-1}} + C_0 \sqrt{t/n} \\
 &\leq \|w - w_{\eta,n,-}\|_{\Sigma^{-1}} + (C_0 + C_1)(t/n)^{1/4}.
 \end{aligned}$$

This implies that, for the prescribed range of t and on the event $E_0(t) \cap E_1(t)$,

$$\min_{w \in D_{\eta;\varepsilon}(\mathbf{g}) \cap B_n(L_w)} \|w - w_{\eta,n,-}\|_{\Sigma^{-1}} \geq (\varepsilon - (C_0 + C_1)(t/n)^{1/4})_+.$$

Using the strong convexity of $L_{\eta,-}(\cdot; L_v)$ with respect to $\|\cdot\|_{\Sigma^{-1}}$, we have for $C'_1 \log(en) \leq t \leq n/C'_1$, on the event $E_0(t) \cap E_1(t)$,

$$\begin{aligned}
 \min_{w \in D_{\eta;\varepsilon}(\mathbf{g}) \cap B_n(L_w)} L_\eta(w; L_v) &\geq \min_{w \in D_{\eta;\varepsilon}(\mathbf{g}) \cap B_n(L_w)} L_{\eta,-}(w; L_v) \\
 &\geq L_{\eta,-}(w_{\eta,n,-}; L_v) + \frac{1}{2}(\varepsilon - (C_0 + C_1)(t/n)^{1/4})_+^2 \\
 &\geq \max_{\beta > 0} \min_{\gamma > 0} \bar{D}_\eta(\beta, \gamma) + \frac{1}{2}(\varepsilon - (C_0 + C_1)(t/n)^{1/4})_+^2 - C_1 \sqrt{t/n}.
 \end{aligned}$$

Now we may choose $\varepsilon \equiv \varepsilon(t, n) \equiv (C_0 + C_1 + 2\sqrt{C_1})(t/n)^{1/4}$ to conclude by adjusting constants. \blacksquare

C.4 Proof of Theorem 3 for $\hat{\mu}_{\eta;G}$

Fix $\xi \in \mathcal{E}_{1,\xi}(\sqrt{t/n})$. All the constants in $\lesssim, \gtrsim, \asymp$ below may depend on K .

(Step 1). In this step, we will obtain an upper bound $\min_{w \in \mathbb{R}^n} H_\eta(w)$. By Proposition 27 and the concentration estimate in Lemma 20, there exists some $C_0 = C_0(K) > 0$ such that on an event E_0 with $\mathbb{P}^\xi(E_0) \geq 1 - C_0 e^{-n/C_0}$,

$$\min_{w \in \mathbb{R}^n} H_\eta(w) = \min_{w \in \mathbb{R}^n} H_\eta(w; L_0) = \min_{w \in B_n(L_0)} H_\eta(w) = \min_{w \in B_n(L_0)} H_\eta(w; L_0). \quad (85)$$

where

$$L_0 \equiv C_0 \left\{ 1 + \left(\|\Sigma^{-1}\|_{\text{op}} \mathbf{1}_{\phi^{-1} \geq 1+1/K} \wedge \eta^{-1} \right) \right\}. \quad (86)$$

Now we shall apply the convex(-side) Gaussian min-max theorem to obtain an upper bound for the right hand side of (85). Recall the definition of $h_\eta = h_{\eta;G}$ and ℓ_η in (28). Using Theorem 14-(2), for any $z \in \mathbb{R}$,

$$\begin{aligned} \mathbb{P}^\xi \left(\min_{w \in \mathbb{R}^n} H_\eta(w) \geq z \right) &\leq \mathbb{P}^\xi \left(\min_{w \in B_n(L_0)} H_\eta(w; L_0) \geq z \right) + \mathbb{P}^\xi(E_0^c) \\ &= \mathbb{P}^\xi \left(\min_{w \in B_n(L_0)} \max_{v \in B_m(L_0)} h_\eta(w, v) \geq z \right) + \mathbb{P}^\xi(E_0^c) \\ &\leq 2 \mathbb{P}^\xi \left(\min_{w \in B_n(L_0)} \max_{v \in B_m(L_0)} \ell_\eta(w, v) \geq z \right) + \mathbb{P}^\xi(E_0^c) \\ &= 2 \mathbb{P}^\xi \left(\min_{w \in B_n(L_0)} L_\eta(w; L_0) \geq z \right) + \mathbb{P}^\xi(E_0^c). \end{aligned} \quad (87)$$

By Proposition 35, there exist some $C_1, C'_1 > 0$ depending on K (which we assume without loss of generality $L_0 > C_1$ and C_1 exceeds the constants in Theorems 28 and 32), such that on an event E_1 with \mathbb{P}^ξ -probability at least $1 - C_1 e^{-n/C_1}$, the map $w \mapsto L_\eta(w; L_0)$ attains its global minimum in $B_n(C_1)$. We may now apply Theorem 28: with $z \equiv \bar{z}(t) = \max_{\beta > 0} \min_{\gamma > 0} \bar{D}_\eta(\beta, \gamma) + \sqrt{t/n}$, for $C'_1 \log(en) \leq t \leq n/C'_1$,

$$\begin{aligned} &\mathbb{P}^\xi \left(\min_{w \in B_n(L_0)} L_\eta(w; L_0) \geq \bar{z}(t) \right) \\ &\leq \mathbb{P}^\xi \left(\min_{w \in B_n(C_1)} L_\eta(w; L_0) \geq \bar{z}(t) \right) + \mathbb{P}^\xi(E_1^c) \leq C_1 e^{-t/C_1} + \mathbb{P}^\xi(E_1^c). \end{aligned} \quad (88)$$

Combining (87)-(88), by enlarging C_1 if necessary, for $C'_1 \log(en) \leq t \leq n/C'_1$, and $\eta \in \Xi_K$,

$$\mathbb{P}^\xi \left(\min_{w \in \mathbb{R}^n} H_\eta(w) \geq \max_{\beta > 0} \min_{\gamma > 0} \bar{D}_\eta(\beta, \gamma) + \sqrt{t/n} \right) \leq C_1 e^{-t/C_1}. \quad (89)$$

An entirely similar argument leads to a lower bound (which will be used later on):

$$\mathbb{P}^\xi \left(\min_{w \in \mathbb{R}^n} H_\eta(w) \leq \max_{\beta > 0} \min_{\gamma > 0} \bar{D}_\eta(\beta, \gamma) - \sqrt{t/n} \right) \leq C_1 e^{-t/C_1}. \quad (90)$$

(Step 2). In this step, we will obtain a lower bound on $\min_{w \in D_{\eta;\varepsilon}(\mathbf{g})} H_\eta(w)$ for the exceptional set $D_\varepsilon(\mathbf{g})$ defined in (78), with a suitable choice of ε . Let us take $C_2, C'_2 > 0$ to be

the constants in Theorem 32, and let $\varepsilon(t, n) \equiv C_2(t/n)^{1/4}$ for $C_2' \log(en) \leq t \leq n/C_2'$. To this end, using Theorem 14-(1) (that holds without convexity), for any $z \in \mathbb{R}$ and $L_v > 0$

$$\begin{aligned} \mathbb{P}^\xi \left(\min_{w \in B_n(L_0) \cap D_{\eta; \varepsilon}(\mathbf{g})} H_\eta(w) \leq z \right) &\leq \mathbb{P}^\xi \left(\min_{w \in B_n(L_0) \cap D_{\eta; \varepsilon}(\mathbf{g})} \max_{v \in B_m(L_v)} h_\eta(w, v) \leq z \right) \\ &\leq 2 \mathbb{P}^\xi \left(\min_{w \in B_n(L_0) \cap D_{\eta; \varepsilon}(\mathbf{g})} \max_{v \in B_m(L_v)} \ell_\eta(w, v) \leq z \right) \\ &= 2 \mathbb{P}^\xi \left(\min_{w \in B_n(L_0) \cap D_{\eta; \varepsilon}(\mathbf{g})} L_\eta(w; L_v) \leq z \right). \end{aligned}$$

By choosing $L_v \asymp 1$ of constant order but large enough, $\varepsilon \equiv \varepsilon(t, n)$ and $z \equiv \bar{z}(t) = \max_{\beta > 0} \min_{\gamma > 0} \bar{D}_\eta(\beta, \gamma) + 2\sqrt{t/n}$, we have for $C_2' \log(en) \leq t \leq n/C_2'$,

$$\mathbb{P}^\xi \left(\min_{w \in B_n(L_0) \cap D_{\eta; \varepsilon(t, n)}(\mathbf{g})} H_\eta(w) \leq \max_{\beta > 0} \min_{\gamma > 0} \bar{D}_\eta(\beta, \gamma) + 2\sqrt{\frac{t}{n}} \right) \leq 2C_2 e^{-t/C_2}. \quad (91)$$

(Step 3). Combining (91) and the localization in (85), there exist some $C_3, C_3' > 0$ depending on K such that for $C_3' \log(en) \leq t \leq n/C_3'$, on an event $E_3(t)$ with $\mathbb{P}^\xi(E_3(t)) \geq 1 - C_3 e^{-t/C_3}$,

$$\begin{aligned} \min_{w \in B_n(L_0) \cap D_{\eta; \varepsilon(t, n)}(\mathbf{g})} H_\eta(w) &\geq \max_{\beta > 0} \min_{\gamma > 0} \bar{D}_\eta(\beta, \gamma) + 2\sqrt{t/n} \\ &> \max_{\beta > 0} \min_{\gamma > 0} \bar{D}_\eta(\beta, \gamma) + \sqrt{t/n} \geq \min_{w \in \mathbb{R}^n} H_\eta(w) = \min_{w \in B_n(L_0)} H_\eta(w). \end{aligned}$$

So on $E_3(t)$, $\hat{w}_\eta \notin D_{\eta; \varepsilon(t, n)}(\mathbf{g}) \cap B_n(L_0)$, i.e., for $C_3' \log(en) \leq t \leq n/C_3'$,

$$\mathbb{P}^\xi \left(|\mathbf{g}(\hat{w}_\eta) - \mathbb{E} \mathbf{g}(w_{\eta, *})| \geq C_3(t/n)^{1/4} \right) \leq C_3 e^{-t/C_3}.$$

Using a change of variable and suitably adjusting the constant C_3 , for any 1-Lipschitz function $\mathbf{g}_0 : \mathbb{R}^n \rightarrow \mathbb{R}$, $\eta \in \Xi_K$ and $\varepsilon \in (0, 1/2]$,

$$\mathbb{P}^\xi \left(|\mathbf{g}_0(\hat{\mu}_\eta) - \mathbb{E} \mathbf{g}_0(\hat{\mu}_{(\Sigma, \mu_0)}^{\text{seq}}(\gamma_{\eta, *}; \tau_{\eta, *}))| \geq \varepsilon \right) \leq C_3 n e^{-n\varepsilon^4/C_3}.$$

(Step 4). In this step we shall establish uniform guarantees. We write $\hat{\mu}_{(\Sigma, \mu_0)}^{\text{seq}}(\gamma_{\eta, *}; \tau_{\eta, *}) = \hat{\mu}_{\eta; (\Sigma, \mu_0)}^{\text{seq}, *}$ in this part of the proof. First, in the case $\phi^{-1} \geq 1 + 1/K$, using $\hat{\mu}_\eta = n^{-1} X^\top (X X^\top / n + \eta I)^{-1} Y$, for $\eta_1, \eta_2 \in [0, K]$,

$$\begin{aligned} \|\hat{\mu}_{\eta_1} - \hat{\mu}_{\eta_2}\| &\lesssim n^{-1} \|G\|_{\text{op}} (\|G\|_{\text{op}} + \|\xi\|) \cdot \|(X X^\top / n + \eta_1 I)^{-1} - (X X^\top / n + \eta_2 I)^{-1}\|_{\text{op}} \\ &\lesssim \|\Sigma^{-1}\|_{\text{op}}^2 \cdot \left(1 + \frac{\|G\|_{\text{op}} + \|\xi\|}{\sqrt{n}} \right)^2 \cdot \|(G G^\top / n)^{-1}\|_{\text{op}}^2 \cdot |\eta_1 - \eta_2|. \end{aligned} \quad (92)$$

Here the last inequality follows by the fact that any p.s.d. matrix A , $\|(A + \eta_1 I)^{-1} - (A + \eta_2 I)^{-1}\|_{\text{op}} \leq \lambda_{\min}^{-2}(A) |\eta_1 - \eta_2|$. As $\|\Sigma^{-1}\|_{\text{op}} \lesssim n$ under $\mathcal{H}_\Sigma \leq K$, there exists $C_4 = C_4(K) > 0$ such that on an event E_4 with $\mathbb{P}^\xi(E_4) \geq 1 - C_4 e^{-n/C_4}$,

$$\|\hat{\mu}_{\eta_1} - \hat{\mu}_{\eta_2}\| \leq C_4 n^2 |\eta_1 - \eta_2|. \quad (93)$$

On the other hand, note that for $\eta_1, \eta_2 \in [0, K]$, using Proposition 23-(3),

$$\|\widehat{\mu}_{\eta_1;(\Sigma, \mu_0)}^{\text{seq},*} - \widehat{\mu}_{\eta_2;(\Sigma, \mu_0)}^{\text{seq},*}\| \lesssim (1 \vee e_g) \|\Sigma^{-1}\|_{\text{op}}^2 |\eta_1 - \eta_2|. \quad (94)$$

So we have

$$|\mathbb{E} \mathfrak{g}_0(\widehat{\mu}_{\eta_1;(\Sigma, \mu_0)}^{\text{seq},*}) - \mathbb{E} \mathfrak{g}_0(\widehat{\mu}_{\eta_2;(\Sigma, \mu_0)}^{\text{seq},*})| \leq C_4 n^2 |\eta_1 - \eta_2|. \quad (95)$$

Now by taking an $\varepsilon/(2C_4 n^2)$ -net Λ_ε of $[0, K]$ and a union bound,

$$\begin{aligned} & \mathbb{P}^\xi \left(\sup_{\eta \in [0, K]} |\mathfrak{g}_0(\widehat{\mu}_\eta) - \mathbb{E} \mathfrak{g}_0(\widehat{\mu}_{\eta;(\Sigma, \mu_0)}^{\text{seq},*})| \geq 2\varepsilon \right) \\ & \leq \mathbb{P}^\xi \left(\max_{\eta \in \Lambda_\varepsilon} |\mathfrak{g}_0(\widehat{\mu}_\eta) - \mathbb{E} \mathfrak{g}_0(\widehat{\mu}_{\eta;(\Sigma, \mu_0)}^{\text{seq},*})| \geq \varepsilon \right) + \mathbb{P}(E_4^c) \\ & \leq (1 + 2C_4 K n^2 / \varepsilon) \cdot C_3 n e^{-n\varepsilon^4/C_3} + C_4 e^{-n/C_4} \leq C \cdot \varepsilon^{-1} n^3 e^{-n\varepsilon^4/C}. \end{aligned} \quad (96)$$

By adjusting constants, we may replace n^3/ε by n . We then conclude by further taking expectation with respect to ξ , and noting that $\mathbb{P}(\xi \in \mathcal{E}_{1,\xi}(\sqrt{t/n})) \geq 1 - Ce^{-t/C}$.

Next, in the case $\phi^{-1} < 1 + 1/K$, we work with $\eta \in [1/K, K]$ and use the standard form of $\widehat{\mu}_\eta$ with $\widehat{\mu}_\eta = n^{-1}(X^\top X/n + \eta I)^{-1} X^\top Y$. As $\eta \geq 1/K$, the spectrum of the middle inverse matrix is bounded by $1/\eta \leq K$, so we may replicate the above calculations in (93) and (95) to reach a similar estimate as in (96). \square

C.5 Proof of Theorem 3 for $\widehat{r}_{\eta;G}$

Recall the cost function $h_\eta = h_{\eta;G}, \ell_\eta$ defined in (28). It is easy to see that

$$\widehat{v}_\eta \equiv \arg \max_{v \in \mathbb{R}^m} \min_{w \in \mathbb{R}^n} h_\eta(w, v) = \frac{1}{\sqrt{n}\eta} (G\widehat{w}_\eta - \xi) = -\frac{\widehat{r}_\eta}{\eta}. \quad (97)$$

We shall define the ‘population’ version of \widehat{v}_η as

$$v_{\eta,*} \equiv \frac{1}{\phi\tau_{\eta,*}} \left(\sqrt{\phi\gamma_{\eta,*}^2 - \sigma_\xi^2} \cdot \frac{h}{\sqrt{n}} - \frac{\xi}{\sqrt{n}} \right) \quad (98)$$

in the Gordon problem.

Proposition 36. *Suppose the following hold for some $K > 0$.*

- $1/K \leq \phi^{-1}, \eta \leq K, \|\mu_0\| \vee \|\Sigma\| \vee \mathcal{H}_\Sigma \leq K$.
- *Assumption B holds with $\sigma_\xi^2 \in [1/K, K]$.*

There exist constants $C, C' > 0$ depending on K such that for $C' \log(en) \leq t \leq n/C'$, $\eta \in [1/K, K]$ and $\xi \in \mathcal{E}_{1,\xi}(\sqrt{t/n})$,

$$\begin{aligned} & \mathbb{P}^\xi \left(\begin{aligned} & \text{The map } v \mapsto \ell_\eta(w_{\eta,*}; v) \text{ is } \eta\text{-strongly concave with unique maximizer } v_{\eta,n} \\ & \text{satisfying } \|v_{\eta,n}\| \leq C \text{ and } \|v_{\eta,n} - v_{\eta,*}\| \leq C\sqrt{t/n}. \\ & \text{Furthermore, } \left| \max_v \ell_\eta(w_{\eta,*}, v) - \max_{\beta>0} \min_{\gamma>0} \overline{D}_\eta(\beta, \gamma) \right| \leq C\sqrt{t/n}. \end{aligned} \right) \geq 1 - Ce^{-t/C}. \end{aligned}$$

We need the following before the proof of Proposition 36.

Lemma 37. *Suppose $1/K \leq \phi^{-1}, \sigma_\xi^2 \leq K$, and $\|\mu_0\| \vee \|\Sigma\|_{\text{op}} \vee \mathcal{H}_\Sigma \leq K$ for some $K > 0$. Recall $w_{\eta,*}$ defined in (77). Then there exist constants $C, C' > 0$ depending on K such that for $C' \log(en) \leq t \leq n/C'$, $\eta \in \Xi_K$ and $\xi \in \mathcal{E}_{1,\xi}(\sqrt{t/n})$,*

$$\mathbb{P}^\xi \left(\max \left\{ |(\text{id} - \mathbb{E})\langle g/\sqrt{n}, w_{\eta,*} \rangle|, |(\text{id} - \mathbb{E})\|w_{\eta,*}\|^2|, |(\text{id} - \mathbb{E})F(w_{\eta,*})|, \right. \right. \\ \left. \left. n^{-1}|(\text{id} - \mathbb{E})\|w_{\eta,*}\|h - \xi\|^2| \right\} \geq \sqrt{t/n} \right) \leq Ce^{-t/C}.$$

Proof All the constants in $\lesssim, \gtrsim, \asymp$ below may depend on K . Recall $w_{\eta,*} = (\Sigma + \tau_{\eta,*}I)^{-1}\Sigma^{1/2}(-\tau_{\eta,*}\mu_0 + \gamma_{\eta,*}\Sigma^{1/2}g/\sqrt{n})$. Under the assumed conditions, $\gamma_{\eta,*}, \tau_{\eta,*} \asymp 1$. We shall consider the four terms separately below.

For the first term, we have

$$n^{-1/2}|\langle g, w_{\eta,*} \rangle - \mathbb{E}\langle g, w_{\eta,*} \rangle| \leq \tau_{\eta,*} \cdot n^{-1/2}|\langle (\Sigma + \tau_{\eta,*}I)^{-1}\Sigma^{1/2}\mu_0, g \rangle| \\ + \gamma_{\eta,*} \cdot n^{-1}(\text{id} - \mathbb{E})\|(\Sigma + \tau_{\eta,*}I)^{-1/2}\Sigma^{1/2}g\|^2 \equiv A_{1,1} + A_{1,2}.$$

The concentration of the term $A_{1,1}$ can be handled using Gaussian tails and the fact that $\|(\Sigma + \tau_{\eta,*}I)^{-1}\Sigma^{1/2}\mu_0\|^2 \lesssim 1$. For the term $A_{1,2}$, with $H_1(g) \equiv \|(\Sigma + \tau_{\eta,*}I)^{-1/2}\Sigma^{1/2}g\|^2$, it is easy to evaluate $\|\nabla H_1(g)\|^2 = 4\|(\Sigma + \tau_{\eta,*}I)^{-1}\Sigma g\|^2 \leq 4H_1(g)$ and $\mathbb{E}H_1(g) \leq n$, so Proposition 59 applies to conclude the concentration of $A_{1,2}$.

For the second term, we may decompose

$$\|w_{\eta,*}\|^2 - \mathbb{E}\|w_{\eta,*}\|^2 \lesssim \tau_{\eta,*}\gamma_{\eta,*} \cdot n^{-1/2}|\langle (\Sigma + \tau_{\eta,*}I)^{-2}\Sigma^{3/2}\mu_0, g \rangle| \\ + \gamma_{\eta,*}^2 \cdot n^{-1}(\text{id} - \mathbb{E})\|(\Sigma + \tau_{\eta,*}I)^{-1}\Sigma g\|^2.$$

From here we may handle the concentration of the above two terms in a completely similar fashion to $A_{1,1}$ and $A_{1,2}$ above.

For the third term, recall that $\widehat{\mu}_{(\Sigma, \mu_0)}^{\text{seq}}(\gamma; \tau) = (\Sigma + \tau I)^{-1}\Sigma^{1/2}(\Sigma^{1/2}\mu_0 + \gamma g/\sqrt{n})$, so

$$|F(w_{\eta,*}) - \mathbb{E}F(w_{\eta,*})| = \frac{1}{2}|\|\widehat{\mu}_{(\Sigma, \mu_0)}^{\text{seq}}(\gamma_{\eta,*}; \tau_{\eta,*})\|^2 - \mathbb{E}\|\widehat{\mu}_{(\Sigma, \mu_0)}^{\text{seq}}(\gamma_{\eta,*}; \tau_{\eta,*})\|^2| \\ \lesssim \gamma_{\eta,*} \cdot n^{-1/2}|\langle (\Sigma + \tau_{\eta,*}I)^{-2}\Sigma^{3/2}\mu_0, g \rangle| + \gamma_{\eta,*}^2 \cdot n^{-1}(\text{id} - \mathbb{E})\|(\Sigma + \tau_{\eta,*}I)^{-1}\Sigma^{1/2}g\|^2.$$

The concentration properties of the two terms on the right hand side above can be handled similarly to the case for the second term.

For the last term, we have

$$n^{-1}|\|\|w_{\eta,*}\|h - \xi\|^2 - \mathbb{E}\|\|w_{\eta,*}\|h - \xi\|^2| \\ \lesssim n^{-1}|\|\|w_{\eta,*}\|^2\|h\|^2 - \mathbb{E}\|\|w_{\eta,*}\|^2\|h\|^2| + n^{-1}\|\|w_{\eta,*}\|\|\langle h, \xi \rangle|\| \equiv A_{4,1} + A_{4,2}.$$

On the other hand, on the event $\mathcal{E}_1(\sqrt{t/n})$,

$$A_{4,1} \lesssim (\|h\|^2/n)|\|w_{\eta,*}\|^2 - \mathbb{E}\|w_{\eta,*}\|^2| + n^{-1}\mathbb{E}\|w_{\eta,*}\|^2 \cdot |\|h\|^2 - m| \lesssim \sqrt{t/n},$$

and $A_{4,2} \lesssim (1 \vee e_g) \cdot n^{-1} |\langle h, \xi \rangle| \lesssim \sqrt{t/n}$. Combining the above estimates concludes the concentration claim for the last term. \blacksquare

Proof [Proof of Proposition 36] Fix $\xi \in \mathcal{E}_{1,\xi}(\sqrt{t/n})$. All the constants in $\lesssim, \gtrsim, \asymp$ below may depend on K .

(Step 1). In this step, we establish both the uniqueness and the apriori estimates for $v_{\eta,n}$. Using Lemma 37, we may choose a sufficiently large $C, C' > 0$ depending on K such that $C' \log(en) \leq t \leq n/C'$,

$$\mathbb{P}^{\xi} \left(E_0(t) \equiv \left\{ \max \left\{ |(\text{id} - \mathbb{E})\langle g/\sqrt{n}, w_{\eta,*} \rangle|, |(\text{id} - \mathbb{E})\|w_{\eta,*}\|^2|, |(\text{id} - \mathbb{E})F(w_{\eta,*})|, \right. \right. \right. \\ \left. \left. \left. n^{-1} |(\text{id} - \mathbb{E})\|w_{\eta,*}\|h - \xi|^2 \right\} \leq \sqrt{t/n} \right\} \right) \geq 1 - Ce^{-t/C}.$$

Therefore, on the event $E_0(t)$,

$$\langle g/\sqrt{n}, w_{\eta,*} \rangle \geq \mathbb{E}\langle g/\sqrt{n}, w_{\eta,*} \rangle - \sqrt{t/n} = \gamma_{\eta,*} \cdot n^{-1} \text{tr}((\Sigma + \tau_{\eta,*}I)^{-1}\Sigma) - \sqrt{t/n}.$$

Note that $n^{-1} \text{tr}((\Sigma + \tau_{\eta,*}I)^{-1}\Sigma) \gtrsim \mathcal{H}_{\Sigma}^{-1} \gtrsim 1$, by choosing sufficiently large C , we conclude $\langle g/\sqrt{n}, w_{\eta,*} \rangle > 0$ on the event $E_0(t)$. This implies that $v \mapsto \ell_{\eta}(w_{\eta,*}, v)$ is η -strongly concave with respect to $\|\cdot\|$, so $v_{\eta,n}$ exists uniquely on $E_0(t)$.

Next we derive apriori estimates. We claim that on $E_0(t)$, $v_{\eta,n} = \arg \max_{v \in \mathbb{R}^m} \ell_{\eta}(w_{\eta,*}, v)$ takes the following form:

$$v_{\eta,n} = \frac{1}{\sqrt{n\eta}} \left(1 - \frac{\langle g, w_{\eta,*} \rangle}{\|w_{\eta,*}\| \|h - \xi\|} \right)_+ \cdot (\|w_{\eta,*}\| h - \xi). \quad (99)$$

To see this, using the definition

$$\begin{aligned} v_{\eta,n} &= \arg \max_{v \in \mathbb{R}^m} \left\{ \frac{1}{\sqrt{n}} \left(-\|v\| \langle g, w_{\eta,*} \rangle + \|w_{\eta,*}\| \langle h, v \rangle - \langle v, \xi \rangle \right) - \frac{\eta \|v\|^2}{2} \right\} \\ &= \arg \max_{\alpha \geq 0} \left\{ \frac{\alpha}{\sqrt{n}} \left(-\langle g, w_{\eta,*} \rangle + \|w_{\eta,*}\| \|h - \xi\| \right) - \frac{\eta \alpha^2}{2} \right\} \cdot \frac{\|w_{\eta,*}\| h - \xi}{\|w_{\eta,*}\| \|h - \xi\|} \\ &= \frac{1}{\sqrt{n\eta}} \left(-\langle g, w_{\eta,*} \rangle + \|w_{\eta,*}\| \|h - \xi\| \right)_+ \cdot \frac{\|w_{\eta,*}\| h - \xi}{\|w_{\eta,*}\| \|h - \xi\|}. \end{aligned}$$

Some simple algebra leads to the expression in (99). The boundedness of $\|v_{\eta,n}\|$ then follows from the boundedness of $\|w_{\eta,*}\|$.

(Step 2). In this step, we establish the bound on $\|v_{\eta,n} - v_{\eta,*}\|$. The key observation is that we may rewrite $v_{\eta,*}$ defined via (98) into the following form

$$v_{\eta,*} = \frac{1}{\sqrt{n\eta}} \left(1 - \frac{\mathbb{E}\langle g, w_{\eta,*} \rangle}{\mathbb{E}^{1/2} \|w_{\eta,*}\| \cdot \|h - \xi\|} \right) \cdot (\mathbb{E}^{1/2} \|w_{\eta,*}\|^2 \cdot h - \xi). \quad (100)$$

This can be seen by observing

$$\begin{cases} \mathbb{E} \|w_{\eta,*}\|^2 = \mathbb{E} \text{err}_{(\Sigma, \mu_0)}(\gamma_{\eta,*}; \tau_{\eta,*}) = \phi \gamma_{\eta,*}^2 - \sigma_{\xi}^2, \\ \mathbb{E} \langle g, w_{\eta,*} \rangle = \frac{\sqrt{n}}{\gamma_{\eta,*}} \cdot \mathbb{E} \text{dof}_{(\Sigma, \mu_0)}(\gamma_{\eta,*}; \tau_{\eta,*}) = \sqrt{n} \gamma_{\eta,*} \cdot \left(\phi - \frac{\eta}{\tau_{\eta,*}} \right), \\ \mathbb{E}^{1/2} \|w_{\eta,*}\| \cdot \|h - \xi\|^2 = \sqrt{m} (\mathbb{E} \|w_{\eta,*}\|^2 + \sigma_{\xi}^2)^{1/2} = \sqrt{m} \phi \gamma_{\eta,*}, \end{cases} \quad (101)$$

and therefore $1 - \frac{\mathbb{E}\langle g, w_{\eta,*} \rangle}{\mathbb{E}^{1/2} \|\|w_{\eta,*}\| \cdot h - \xi\|^2} = \frac{\eta}{\phi\tau_{\eta,*}}$. Now with (99)-(100), we may use Lemma 37 to estimate

$$\begin{aligned} \|v_{\eta,n} - v_{\eta,*}\| &\leq \frac{1}{\sqrt{n\eta}} \left| \|w_{\eta,*}\| - \mathbb{E}^{1/2} \|w_{\eta,*}\|^2 \right| \cdot \|h\| \\ &\quad + \frac{1}{\sqrt{n\eta}} \left| \frac{\langle g, w_{\eta,*} \rangle}{\|w_{\eta,*}\| \|h - \xi\|} - \frac{\mathbb{E}\langle g, w_{\eta,*} \rangle}{\mathbb{E}^{1/2} \|\|w_{\eta,*}\| \cdot h - \xi\|^2} \right| \|\mathbb{E}^{1/2} \|w_{\eta,*}\|^2 \cdot h - \xi\| \\ &\equiv V_1 + V_2. \end{aligned} \tag{102}$$

We first handle the term V_1 . As $\mathbb{E}\|w_{\eta,*}\|^2 \geq \gamma_{\eta,*}^2 \text{tr}(\Sigma^2(\Sigma + \tau_{\eta,*})^{-2})/n \gtrsim 1$, on the event $E_0(t) \cap \mathcal{E}_{1,0}(\sqrt{t/n})$,

$$V_1 \lesssim \frac{\|h\|}{\sqrt{n}} \cdot \frac{|\|w_{\eta,*}\|^2 - \mathbb{E}\|w_{\eta,*}\|^2|}{\mathbb{E}^{1/2} \|w_{\eta,*}\|^2} \lesssim \sqrt{t/n}. \tag{103}$$

Next we handle V_2 . On the event $E_0(t) \cap \mathcal{E}_{1,0}(\sqrt{t/n})$,

$$\begin{aligned} V_2 &\lesssim \|\|w_{\eta,*}\|h - \xi\|^{-1} \cdot |\langle g, w_{\eta,*} \rangle - \mathbb{E}\langle g, w_{\eta,*} \rangle| \\ &\quad + \mathbb{E}\langle g, w_{\eta,*} \rangle \cdot \left| \|\|w_{\eta,*}\|h - \xi\|^{-1} - \mathbb{E}^{-1/2} \|\|w_{\eta,*}\| \cdot h - \xi\|^2 \right| \\ &\lesssim n^{-1/2} |\langle g, w_{\eta,*} \rangle - \mathbb{E}\langle g, w_{\eta,*} \rangle| + n^{-1/2} \left| \|\|w_{\eta,*}\|h - \xi\| - \mathbb{E}^{1/2} \|\|w_{\eta,*}\|h - \xi\|^2 \right| \\ &\lesssim n^{-1/2} |\langle g, w_{\eta,*} \rangle - \mathbb{E}\langle g, w_{\eta,*} \rangle| + n^{-1} \left| \|\|w_{\eta,*}\|h - \xi\|^2 - \mathbb{E}\|\|w_{\eta,*}\|h - \xi\|^2 \right| \\ &\lesssim \sqrt{t/n}. \end{aligned} \tag{104}$$

The desired estimate for $\|v_{\eta,n} - v_{\eta,*}\|$ follows from (102)-(104).

(Step 3). In this step, we prove the claimed bound on $|\max_v \ell_\eta(w_{\eta,*}, v) - \bar{D}_\eta(\beta_{\eta,*}, \gamma_{\eta,*})|$. First note that

$$\begin{aligned} &\max_{v \in \mathbb{R}^m} \ell_\eta(w_{\eta,*}, v) \\ &\equiv \max_{v \in \mathbb{R}^m} \left\{ \frac{1}{\sqrt{n}} \left(-\|v\| \langle g, w_{\eta,*} \rangle + \|w_{\eta,*}\| \langle h, v \rangle - \langle v, \xi \rangle \right) + F(w_{\eta,*}) - \frac{\eta\|v\|^2}{2} \right\} \\ &= \frac{1}{2n\eta} \left(\|\|w_{\eta,*}\|h - \xi\| - \langle g, w_{\eta,*} \rangle \right)_+^2 + F(w_{\eta,*}). \end{aligned} \tag{105}$$

On the other hand, with $\#_{\eta;(\Sigma, \mu_0)}^* \equiv \#_{(\Sigma, \mu_0)}(\gamma_{\eta,*}; \tau_{\eta,*})$, $\# \in \{\text{err}, \text{dof}\}$,

$$\mathbb{E} \mathbf{e}_F \left(\frac{\gamma_{\eta,*}}{\sqrt{n}} g; \frac{\gamma_{\eta,*}}{\beta_{\eta,*}} \right) = \frac{\beta_{\eta,*}}{2\gamma_{\eta,*}} \left(\mathbb{E} \text{err}_{\eta;(\Sigma, \mu_0)}^* - 2 \mathbb{E} \text{dof}_{\eta;(\Sigma, \mu_0)}^* + \gamma_{\eta,*}^2 \right) + \mathbb{E} F(w_{\eta,*}),$$

so we may rewrite $\max_{\beta>0} \min_{\gamma>0} \bar{D}_\eta(\beta, \gamma) = \bar{D}_\eta(\beta_{\eta,*}, \gamma_{\eta,*})$ as follows:

$$\begin{aligned} \bar{D}_\eta(\beta_{\eta,*}, \gamma_{\eta,*}) &= \frac{\beta_{\eta,*}}{2} \left(\gamma_{\eta,*}(\phi - 1) + \frac{\sigma_\xi^2}{\gamma_{\eta,*}} \right) - \frac{\eta\beta_{\eta,*}^2}{2} + \mathbb{E} \mathbf{e}_F \left(\frac{\gamma_{\eta,*}}{\sqrt{n}} g; \frac{\gamma_{\eta,*}}{\beta_{\eta,*}} \right) \\ &= \frac{\beta_{\eta,*}}{2\gamma_{\eta,*}} \left(\phi\gamma_{\eta,*}^2 + \sigma_\xi^2 + \mathbb{E} \text{err}_{\eta;(\Sigma, \mu_0)}^* - 2 \mathbb{E} \text{dof}_{\eta;(\Sigma, \mu_0)}^* \right) - \frac{\eta\beta_{\eta,*}^2}{2} + \mathbb{E} F(w_{\eta,*}) \end{aligned}$$

$$= \frac{\beta_{\eta,*}}{\gamma_{\eta,*}} \left(\phi \gamma_{\eta,*}^2 - \mathbb{E} \text{dof}_{\eta;(\Sigma, \mu_0)}^* \right) - \frac{\eta \beta_{\eta,*}^2}{2} + \mathbb{E} F(w_{\eta,*}).$$

Further using the second and third equations in (101), it now follows that

$$\begin{aligned} \max_{\beta > 0} \min_{\gamma > 0} \bar{D}_\eta(\beta, \gamma) &= \frac{\beta_{\eta,*}}{\sqrt{n}} \left(\mathbb{E}^{1/2} \left\| \|w_{\eta,*}\| h - \xi \right\|^2 - \mathbb{E} \langle g, w_{\eta,*} \rangle \right) - \frac{\eta \beta_{\eta,*}^2}{2} + \mathbb{E} F(w_{\eta,*}) \\ &= \frac{1}{2n\eta} \left(\mathbb{E}^{1/2} \left\| \|w_{\eta,*}\| h - \xi \right\|^2 - \mathbb{E} \langle g, w_{\eta,*} \rangle \right)^2 + \mathbb{E} F(w_{\eta,*}). \end{aligned} \quad (106)$$

Now combining (105) and (106), on the event $E_0(t) \cap \mathcal{E}_{1,0}(\sqrt{t/n})$, we may estimate

$$\begin{aligned} & \left| \max_{v \in \mathbb{R}^m} \ell_\eta(w_{\eta,*}, v) - \max_{\beta > 0} \min_{\gamma > 0} \bar{D}_\eta(\beta, \gamma) \right| \\ & \lesssim n^{-1/2} |\langle g, w_{\eta,*} \rangle - \mathbb{E} \langle g, w_{\eta,*} \rangle| + n^{-1/2} \left| \left\| \|w_{\eta,*}\| h - \xi \right\| - \mathbb{E}^{1/2} \left\| \|w_{\eta,*}\| h - \xi \right\|^2 \right| \\ & \quad + |F(w_{\eta,*}) - \mathbb{E} F(w_{\eta,*})| \lesssim \sqrt{t/n}, \end{aligned}$$

completing the proof. \blacksquare

Proof [Proof of Theorem 3 for \hat{r}_η] Fix $\xi \in \mathcal{E}_{1,\xi}(\sqrt{t/n})$. All the constants in $\lesssim, \gtrsim, \asymp$ below may depend on K . We sometimes write $\bar{\mathcal{D}}_\eta \equiv \max_{\beta > 0} \min_{\gamma > 0} \bar{D}_\eta(\beta, \gamma)$.

As $\hat{r}_\eta = -\eta \hat{v}_\eta$, we only need to study \hat{v}_η . Fix $\varepsilon > 0$, and any $\mathbf{h} : \mathbb{R}^m \rightarrow \mathbb{R}$, let

$$D_{\eta;\varepsilon}(\mathbf{h}) \equiv \{v \in \mathbb{R}^m : |\mathbf{h}(v) - \mathbb{E}^\xi \mathbf{h}(v_{\eta,*})| \geq \varepsilon\}.$$

(Step 1). In this step we establish the Gordon cost cap: there exist constants $C_1, C'_1 > 0$ depending on K such that for $C'_1 \log(en) \leq t \leq n/C'_1$,

$$\mathbb{P}^\xi \left(E_1(t)^c \equiv \left\{ \max_{v \in D_{\eta;C'_1(t/n)^{1/4}}(\mathbf{h})} \ell_\eta(w_{\eta,*}, v) \geq \bar{\mathcal{D}}_\eta - C_1^{-1} \sqrt{t/n} \right\} \right) \leq C_1 e^{-t/C_1}. \quad (107)$$

To this end, first note that by the Lipschitz property of \mathbf{h} , the Gaussian concentration and Proposition 36, there exist some $C_0, C'_0 > 0$ depending on K such that for $C'_0 \log(en) \leq t \leq n/C'_0$, on an event $E_{1,0}(t)$ with probability at least $1 - C_0 e^{-t/C_0}$, we have uniformly in $v \in D_{\eta;\varepsilon}(\mathbf{h})$,

$$\begin{aligned} \varepsilon &\leq |\mathbf{h}(v) - \mathbb{E}^\xi \mathbf{h}(v_{\eta,*})| \leq |\mathbf{h}(v) - \mathbf{h}(v_{\eta,*})| + |\mathbf{h}(v_{\eta,*}) - \mathbb{E}^\xi \mathbf{h}(v_{\eta,*})| \\ &\leq \|v - v_{\eta,n}\| + \|v_{\eta,*} - v_{\eta,n}\| + C \sqrt{t/n} \leq \|v - v_{\eta,n}\| + C_0 \sqrt{t/n}, \end{aligned}$$

and all the properties in Proposition 36 hold. In other word, on $E_{1,0}(t)$ with the prescribed range of t ,

$$\inf_{v \in D_{\eta;\varepsilon}(\mathbf{h})} \|v - v_{\eta,n}\| \geq (\varepsilon - C_0 \sqrt{t/n})_+.$$

Using the η -strong concavity of $v \mapsto \ell_\eta(w_{\eta,*}, v)$ on $E_{1,0}(t)$, we have

$$\max_{v \in D_{\eta;\varepsilon}(\mathbf{h})} \ell_\eta(w_{\eta,*}, v) \leq \ell_\eta(w_{\eta,*}, v_{\eta,n}) - \frac{\eta}{2} \inf_{v \in D_{\eta;\varepsilon}(\mathbf{h})} \|v - v_{\eta,n}\|^2$$

$$\leq \max_{\beta>0} \min_{\gamma>0} \bar{D}_\eta(\beta, \gamma) - \frac{\eta}{2} (\varepsilon - C_0 \sqrt{t/n})_+^2 + C_1 \sqrt{t/n}.$$

By choosing $\varepsilon \equiv \varepsilon_{\eta;v}(t, n) \equiv C_0 \sqrt{t/n} + 2\sqrt{C_1/\eta} \cdot (t/n)^{1/4}$, we have on $E_{1,0}(t)$,

$$\max_{v \in D_{\eta; \varepsilon_{\eta;v}(t, n)}(h)} \ell_\eta(w_{\eta,*}, v) \leq \max_{\beta>0} \min_{\gamma>0} \bar{D}_\eta(\beta, \gamma) - C_1 \sqrt{t/n}. \quad (108)$$

Adjusting constants proves the claim in (107).

(Step 2). In this step, we provide an upper bound for the original cost over exceptional set. More concretely, we will prove that there exist constants $C_2, C'_2 > 0$ depending on K such that for any $L_v > 0$, and $C'_2 \log(en) \leq t \leq n/C'_2$,

$$\begin{aligned} \mathbb{P}^\xi \left(E_2(t)^c \equiv \left\{ \max_{v \in D_{\eta; C_2(t/n)^{1/4}}(h) \cap B_m(L_v)} \min_{w \in \mathbb{R}^n} h_\eta(w, v) \right. \right. \\ \left. \left. \geq \bar{\mathcal{D}}_\eta - C_2^{-1} \sqrt{t/n} \right\} \right) \leq C_2 e^{-t/C_2}. \end{aligned} \quad (109)$$

To see this, first note by Proposition 35, there exists some $C_2 = C_2(K) > 0$ such that on an event $E_{2,0}$ with $\mathbb{P}^\xi(E_{2,0}) \geq 1 - C_2 e^{-n/C_2}$, $\|w_{\eta,*}\| \leq C_2$. So with $\bar{z}_{\eta;v}(t, n) \equiv \max_{\beta>0} \min_{\gamma>0} \bar{D}_\eta(\beta, \gamma) - C_1^{-1} \sqrt{t/n}$, for any $L_v > 0$, an application of Theorem 14-(1) yields that for $C'_1 \log(en) \leq t \leq n/C'_1$,

$$\begin{aligned} & \mathbb{P}^\xi \left(\max_{v \in D_{\eta; C_1(t/n)^{1/4}}(h) \cap B_m(L_v)} \min_{w \in \mathbb{R}^n} h_\eta(w, v) \geq \bar{z}_{\eta;v}(t, n) \right) \\ & \leq \mathbb{P}^\xi \left(\max_{v \in D_{\eta; C_1(t/n)^{1/4}}(h) \cap B_m(L_v)} \min_{w \in B_n(C_2)} h_\eta(w, v) \geq \bar{z}_{\eta;v}(t, n) \right) \\ & \leq 2 \mathbb{P}^\xi \left(\max_{v \in D_{\eta; C_1(t/n)^{1/4}}(h) \cap B_m(L_v)} \min_{w \in B_n(C_2)} \ell_\eta(w, v) \geq \bar{z}_{\eta;v}(t, n) \right) \\ & \leq 2 \mathbb{P}^\xi \left(\max_{v \in D_{\eta; C_1(t/n)^{1/4}}(h) \cap B_m(L_v)} \ell_\eta(w_{\eta,*}, v) \geq \bar{z}_{\eta;v}(t, n) \right) + 2 \mathbb{P}^\xi(E_{2,0}^c) \\ & \leq 2 \mathbb{P}^\xi \left(\max_{v \in D_{\eta; C_1(t/n)^{1/4}}(h)} \ell_\eta(w_{\eta,*}, v) \geq \bar{z}_{\eta;v}(t, n) \right) + 2 \mathbb{P}^\xi(E_{2,0}^c) \leq C e^{-t/C}, \end{aligned}$$

proving the claim (107) by possibly adjusting constants.

(Step 3). In this step, we recall a lower bound for the original cost optimum, essentially established in the Step 1 in the proof of Theorem 3. In particular, using (85), (86) and (90), there exist $C_3, C'_3, C''_3 > 0$ depending on K , such that for $C'_3 \log(en) \leq t \leq n/C'_3$,

$$\mathbb{P}^\xi \left(E_{3,0}(t)^c \equiv \left\{ \max_{v \in B_m(C'_3)} \min_{w \in \mathbb{R}^n} h_\eta(w, v) \leq \bar{\mathcal{D}}_\eta - C_3^{-1} \sqrt{t/n} \right\} \right) \leq C_3 e^{-t/C_3}, \quad (110)$$

and

$$\mathbb{P}^\xi \left(E_{3,1}^c \equiv \left\{ \max_{v \in B_m(C''_3)} \min_{w \in \mathbb{R}^n} h_\eta(w, v) = \max_{v \in \mathbb{R}^m} \min_{w \in \mathbb{R}^n} h_\eta(w, v) \right\} \right) \leq C_3 e^{-n/C_3}. \quad (111)$$

(Step 4). By choosing without loss of generality $C_3 > C_2$, on the event $E_2(t) \cap E_{3,0}(t) \cap E_{3,1}$, (109)-(111) yield that for any $C' \log(en) \leq t \leq n/C'$,

$$\begin{aligned} & \max_{v \in \mathbf{D}_{\eta; C_2(t/n)^{1/4}(\mathbf{h})} \cap B_m(C_3'')} \min_{w \in \mathbb{R}^n} h_\eta(w, v) \leq \max_{\beta > 0} \min_{\gamma > 0} \bar{\mathbf{D}}_\eta(\beta, \gamma) - C_2^{-1} \sqrt{t/n} \\ & < \max_{\beta > 0} \min_{\gamma > 0} \bar{\mathbf{D}}_\eta(\beta, \gamma) - C_3^{-1} \sqrt{t/n} \leq \max_{v \in B_m(C_3'')} \min_{w \in \mathbb{R}^n} h_\eta(w, v) = \max_{v \in \mathbb{R}^m} \min_{w \in \mathbb{R}^n} h_\eta(w, v). \end{aligned}$$

This means on the event $E_2(t) \cap E_{3,0}(t) \cap E_{3,1}$, $\hat{v}_\eta \notin \mathbf{D}_{\eta; C_2(t/n)^{1/4}(\mathbf{h})}$, i.e., there exist some $C_4, C_4' > 0$ depending on K such that for $C_4' \log(en) \leq t \leq n/C_4'$ and $1/K \leq \eta \leq K$,

$$\mathbb{P}^\xi \left(|\mathbf{h}(\hat{v}_\eta) - \mathbb{E}^\xi \mathbf{h}(v_{\eta,*})| \geq C_4(t/n)^{1/4} \right) \leq C_4 e^{-t/C_4}. \quad (112)$$

(Step 5). In this final step, we shall prove uniform version of the estimate (112). For $\eta_1, \eta_2 \in [1/K, K]$, using the definition of \hat{v}_η in (97),

$$\begin{aligned} & |\mathbf{h}(\hat{v}_{\eta_1}) - \mathbf{h}(\hat{v}_{\eta_2})| \leq \|\hat{v}_{\eta_1} - \hat{v}_{\eta_2}\| \\ & \leq n^{-1/2} \|\eta_1^{-1} G \hat{w}_{\eta_1} - \eta_2^{-1} G \hat{w}_{\eta_2}\| + (\|\xi\|/\sqrt{n}) \cdot |\eta_1^{-1} - \eta_2^{-1}| \\ & \leq \frac{\|G \hat{w}_{\eta_1}\| + \|\xi\|}{\sqrt{n}} \cdot |\eta_1^{-1} - \eta_2^{-1}| + \frac{1}{\sqrt{n} \eta_2} \|G(\hat{w}_{\eta_1} - \hat{w}_{\eta_2})\| \\ & \lesssim \left(1 + \|\hat{\mu}_{\eta_1}\| \frac{\|G\|_{\text{op}}}{\sqrt{n}}\right) \cdot |\eta_1 - \eta_2| + \frac{\|G\|_{\text{op}}}{\sqrt{n}} \cdot \|\hat{\mu}_{\eta_1} - \hat{\mu}_{\eta_2}\|. \end{aligned}$$

Using that $\|\hat{\mu}_\eta\| = \|n^{-1}(X^\top X/n + \eta I)^{-1} X^\top Y\| \leq \|X^\top Y\|/(n\eta) \lesssim (1 + \|G\|_{\text{op}}/\sqrt{n})^2$, we have

$$|\mathbf{h}(\hat{v}_{\eta_1}) - \mathbf{h}(\hat{v}_{\eta_2})| \lesssim (1 + \|G\|_{\text{op}}/\sqrt{n})^3 \cdot (|\eta_1 - \eta_2| \vee \|\hat{\mu}_{\eta_1} - \hat{\mu}_{\eta_2}\|). \quad (113)$$

In view of (93), there exists some $C_5 > 0$ depending on K , such that on an event $E_{5,1}$ with $\mathbb{P}^\xi(E_{5,1}) \geq 1 - C_5 e^{-n/C_5}$,

$$|\mathbf{h}(\hat{v}_{\eta_1}) - \mathbf{h}(\hat{v}_{\eta_2})| \leq C_5 n^2 |\eta_1 - \eta_2|. \quad (114)$$

On the other hand, using the definition of $v_{\eta,*}$ in (98), Proposition 23-(3) and the fact that $\phi \gamma_{\eta,*}^2 - \sigma_\xi^2 = \mathbb{E} \text{err}_{(\Sigma, \mu_0)}(\gamma_{\eta,*}; \tau_{\eta,*}) \geq \text{tr}((\Sigma + \tau_{\eta,*} I)^{-2} \Sigma^2) \gtrsim 1$, we have

$$\begin{aligned} & |\mathbb{E}^\xi \mathbf{h}(v_{\eta_1,*}) - \mathbb{E}^\xi \mathbf{h}(v_{\eta_2,*})| \leq \mathbb{E}^{1/2, \xi} \|v_{\eta_1,*} - v_{\eta_2,*}\|^2 \\ & \lesssim |\tau_{\eta_1,*}^{-1} \sqrt{\phi \gamma_{\eta_1,*}^2 - \sigma_\xi^2} - \tau_{\eta_2,*}^{-1} \sqrt{\phi \gamma_{\eta_2,*}^2 - \sigma_\xi^2}| + |\tau_{\eta_1,*}^{-1} - \tau_{\eta_2,*}^{-1}| \\ & \lesssim |\gamma_{\eta_1,*}^2 - \gamma_{\eta_2,*}^2| + |\tau_{\eta_1,*}^{-1} - \tau_{\eta_2,*}^{-1}| \leq C_5 |\eta_1 - \eta_2|. \end{aligned} \quad (115)$$

Now we may mimic the proof in (96) to conclude that, by possibly enlarging $C_5 > 0$, for any $\varepsilon \in (0, 1/2]$ and $\xi \in \mathcal{E}_{1, \xi}(\varepsilon^2/C_5)$,

$$\mathbb{P}^\xi \left(\sup_{\eta \in [1/K, K]} |\mathbf{h}(\hat{v}_\eta) - \mathbb{E}^\xi \mathbf{h}(v_{\eta,*})| \geq \varepsilon \right) \leq C_5 n e^{-n\varepsilon^4/C_5},$$

as desired. ■

Appendix D. Universality: Proof of Theorem 4

D.1 Comparison inequalities

For $f : \mathbb{R}^n \rightarrow \mathbb{R}$, let

$$\mathcal{H}_f(w, A) \equiv \frac{1}{2n} \|Aw - \xi\|^2 + f(w).$$

The following theorem is proved in (Han and Shen, 2023, Theorem 2.3).

Theorem 38. *Suppose $1/K \leq \phi^{-1} \leq K$ for some $K > 1$. Let $A_0, B_0 \in \mathbb{R}^{m \times n}$ be two random matrices with independent components, such that $\mathbb{E} A_{0;ij} = \mathbb{E} B_{0;ij} = 0$ and $\mathbb{E} A_{0;ij}^2 = \mathbb{E} B_{0;ij}^2$ for all $i \in [m], j \in [n]$. Further assume that*

$$M \equiv \max_{i \in [m], j \in [n]} (\mathbb{E} |A_{0;ij}|^6 + \mathbb{E} |B_{0;ij}|^6) < \infty.$$

Let $A \equiv A_0/\sqrt{n}$ and $B \equiv B_0/\sqrt{n}$. Then there exists some $C_0 = C_0(K, M) > 0$ such that the following hold: For any $\mathcal{S}_n \subset [-L_n, L_n]^n$ with $L_n \geq 1$, and any $\mathbb{T} \in C^3(\mathbb{R})$, we have

$$\left| \mathbb{E} \mathbb{T} \left(\min_{w \in \mathcal{S}_n} \mathcal{H}_f(w, A) \right) - \mathbb{E} \mathbb{T} \left(\min_{w \in \mathcal{S}_n} \mathcal{H}_f(w, B) \right) \right| \leq C_0 \cdot K_{\mathbb{T}} \cdot r_f(L_n).$$

Here $K_{\mathbb{T}} \equiv 1 + \max_{\ell \in [0:3]} \|\mathbb{T}^{(\ell)}\|_{\infty}$, and $r_f(L_n)$ is defined by

$$r_f(L_n) \equiv \inf_{\delta \in (0, n^{-5/2})} \left\{ \mathcal{N}_f(L_n, \delta) + \left(1 + \frac{1}{m} \sum_{i=1}^m \mathbb{E} |\xi_i|^3 \right)^{1/3} \cdot \frac{L_n^2 \log_+^{2/3}(L_n/\delta)}{n^{1/6}} \right\},$$

where $\mathcal{N}_f(L_n, \delta) \equiv \sup |f(w) - f(w')|$ with the supremum taken over all $w, w' \in [-L_n, L_n]^n$ such that $\|w - w'\|_{\infty} \leq \delta$. Consequently, for any $z \in \mathbb{R}, \varepsilon > 0$,

$$\mathbb{P} \left(\min_{w \in \mathcal{S}_n} \mathcal{H}_f(w, A) > z + 3\varepsilon \right) \leq \mathbb{P} \left(\min_{w \in \mathcal{S}_n} \mathcal{H}_f(w, B) > z + \varepsilon \right) + C_1 (1 \vee \varepsilon^{-3}) r_f(L_n).$$

Here $C_1 > 0$ is an absolute multiple of C_0 .

Let for $u \in \mathbb{R}^m, w \in \mathbb{R}^n, A \in \mathbb{R}^{m \times n}$ and a measurable function $Q : \mathbb{R}^m \times \mathbb{R}^n \rightarrow \mathbb{R}$

$$X(u, w; A) \equiv u^{\top} Aw + Q(u, w). \quad (116)$$

The following theorem is proved in (Han and Shen, 2023, Theorem 2.5).

Theorem 39. *Let $A, B \in \mathbb{R}^{m \times n}$ be two random matrices with independent entries and matching first two moments, i.e., $\mathbb{E} A_{ij}^{\ell} = \mathbb{E} B_{ij}^{\ell}$ for all $i \in [m], j \in [n], \ell = 1, 2$. There exists a universal constant $C_0 > 0$ such that the following hold. For any measurable subsets $\mathcal{S}_u \subset [-L_u, L_u]^m, \mathcal{S}_w \subset [-L_w, L_w]^n$ with $L_u, L_w \geq 1$, and any $\mathbb{T} \in C^3(\mathbb{R})$, we have*

$$\begin{aligned} & \left| \mathbb{E} \mathbb{T} \left(\max_{u \in \mathcal{S}_u} \min_{w \in \mathcal{S}_w} X(u, w; A) \right) - \mathbb{E} \mathbb{T} \left(\max_{u \in \mathcal{S}_u} \min_{w \in \mathcal{S}_w} X(u, w; B) \right) \right| \\ & \leq C_0 \cdot K_{\mathbb{T}} \cdot \inf_{\delta \in (0, 1)} \left\{ M_1 L \delta + \mathcal{N}_Q(L, \delta) + \log_+^{2/3}(L/\delta) \cdot (m+n)^{2/3} M_3^{1/3} L^2 \right\}. \end{aligned}$$

Here $K_{\mathbb{T}} \equiv 1 + \max_{\ell \in [0:3]} \|\mathbb{T}^{(\ell)}\|_{\infty}$, $L \equiv L_u + L_w$, $M_{\ell} \equiv \sum_{i \in [m], j \in [n]} (\mathbb{E} |A_{ij}|^{\ell} + \mathbb{E} |B_{ij}|^{\ell})$, and $\mathcal{N}_Q(L, \delta) \equiv \sup |Q(u, w) - Q(u', w')|$ with the supremum taken over all $u, u' \in [-L, L]^m, w, w' \in [-L, L]^n$ such that $\|u - u'\|_{\infty} \vee \|w - w'\|_{\infty} \leq \delta$. The conclusion continues to hold when max-min is flipped to min-max.

D.2 Delocalization

Recall that $\hat{\mu}_\eta$ defined in (3) can be rewritten as

$$\hat{\mu}_\eta = \arg \min_{\mu \in \mathbb{R}^n} \max_{v \in \mathbb{R}^m} \left\{ \frac{1}{2} \|\mu\|^2 + \frac{1}{\sqrt{n}} \langle v, X\mu - Y \rangle - \frac{\eta}{2} \|v\|^2 \right\}.$$

For any $\eta > 0$, we have the following closed form for $\hat{\mu}_\eta$:

$$\hat{\mu}_\eta = n^{-1} (X^\top X/n + \eta I_n)^{-1} X^\top Y, \quad \hat{v}_\eta = -(\sqrt{n}\eta)^{-1} (Y - X\hat{\mu}_\eta). \quad (117)$$

The above formula does not include the interpolating case $\eta = 0$ when $n > m$. To give an alternative expression, note that the first-order condition for the above minimax optimization is $\hat{\mu}_\eta = X^\top \hat{v}_\eta / \sqrt{n}$, $Y - X\hat{\mu}_\eta = -\sqrt{n}\eta \hat{v}_\eta$, or equivalently,

$$\hat{\mu}_\eta = n^{-1} X^\top (XX^\top/n + \eta I_m)^{-1} Y, \quad \hat{v}_\eta = -n^{-1/2} (XX^\top/n + \eta I_m)^{-1} Y. \quad (118)$$

The following proposition proves delocalization for $\hat{w}_\eta \equiv \Sigma^{1/2}(\hat{\mu}_\eta - \mu_0)$ and \hat{v}_η .

Proposition 40. *Suppose Assumption A holds and the following hold for some $K > 0$.*

- $1/K \leq \phi^{-1} \leq K$, $\|\Sigma^{-1}\|_{\text{op}} \vee \|\Sigma\|_{\text{op}} \leq K$.
- Assumption B holds with $\sigma_\xi^2 \in [1/K, K]$.

Fix $\vartheta \in (0, 1/2]$. Then there exist some constant $C = C(K, \vartheta) > 0$, two measurable sets $\mathcal{U}_\vartheta \subset B_n(1)$, $\mathcal{E}_\vartheta \subset \mathbb{R}^m$ with $\min\{\text{vol}(\mathcal{U}_\vartheta)/\text{vol}(B_n(1)), \mathbb{P}(\xi \in \mathcal{E}_\vartheta)\} \geq 1 - Ce^{-n^{2\vartheta}/C}$, such that

$$\sup_{\mu_0 \in \mathcal{U}_\vartheta, \xi \in \mathcal{E}_\vartheta} \mathbb{P}^\xi \left(\sup_{\eta \in \Xi_K} \left\{ \|\hat{w}_\eta\|_\infty \vee \|\hat{v}_\eta\|_\infty \right\} \geq Cn^{-1/2+\vartheta} \right) \leq Cn^{-100}.$$

The sets $\mathcal{U}_\vartheta, \mathcal{E}_\vartheta$ can be taken as

$$\begin{aligned} \mathcal{U}_\vartheta &\equiv \left\{ \mu_0 \in B_n(1) : \sup_{\eta \in \Xi_K} \left\| \Sigma^{1/2} \left(\mathbb{E} \hat{\mu}_{(\Sigma, \mu_0)}^{\text{seq}}(\gamma_{\eta, *}; \tau_{\eta, *}) - \mu_0 \right) \right\|_\infty \leq C_0 n^{-1/2+\vartheta} \right\}, \\ \mathcal{E}_\vartheta &\equiv \left\{ \xi \in \mathbb{R}^m : \|\xi\|_\infty \leq C_0 n^\vartheta, \left| \|\xi\|^2/m - \sigma_\xi^2 \right| \leq C_0 n^{-1/2+\vartheta} \right\} \end{aligned}$$

for some large enough $C_0 = C_0(K) > 0$.

Remark 41. Proposition 40 formalizes the delocalization required by our comparison argument: uniformly over $\eta \in \Xi_K$, $\hat{w}_\eta = \Sigma^{1/2}(\hat{\mu}_\eta - \mu_0)$ (and likewise \hat{v}_η) is small in ℓ_∞ . The set \mathcal{U}_ϑ encodes this coordinatewise control via the sequence model proxy, thereby ruling out highly localized signals. The event \mathcal{E}_ϑ imposes mild noise concentration in sup-norm and empirical variance; for i.i.d. sub-gaussian coordinates, this event holds with overwhelming probability.

Remark 42. Delocalization in the same sense of the above proposition holds for $\|\mathbf{P}\hat{\mu}_\eta + \mathbf{q}\|_\infty$ with any deterministic matrix $\mathbf{P} \in \mathbb{R}^{n \times n}$ and vector $\mathbf{q} \in \mathbb{R}^n$ satisfying $\|\mathbf{P}\|_{\text{op}} \vee \|\mathbf{q}\| \leq 1$, with a (slightly) different construction of \mathcal{U}_ϑ .

Proof [Proof of Proposition 40] All the constants in $\lesssim, \gtrsim, \asymp$ below may depend on K .
 (1). Let us consider delocalization for \widehat{w}_η . Using (118), for any $s \in [n]$,

$$\begin{aligned} \langle e_s, \widehat{w}_\eta \rangle &= n^{-1} \langle \Sigma^{1/2} e_s, X^\top (\phi \check{\Sigma} + \eta I_m)^{-1} X \mu_0 \rangle - \langle \Sigma^{1/2} e_s, \mu_0 \rangle \\ &\quad + n^{-1} \langle \Sigma^{1/2} e_s, X^\top (\phi \check{\Sigma} + \eta I_m)^{-1} \xi \rangle \equiv A_{1;s} + A_{2;s}. \end{aligned} \quad (119)$$

We first handle $A_{1;s}$. Let ρ be the asymptotic eigenvalue density of $\check{\Sigma} = X X^\top / m$ and fix $c > 0$. By (Knowles and Yin, 2017, Theorem 3.16-(i), Remark 3.17 and Lemma 4.4-(i)), for any small $\vartheta > 0$ and large $D > 0$,

$$\begin{aligned} \mathbb{P}^\xi \left(\left| m^{-1} \langle \Sigma^{1/2} e_s, X^\top (\check{\Sigma} - z I_m)^{-1} X \mu_0 \rangle \right. \right. \\ \left. \left. - \langle \Sigma^{1/2} e_s, \mathbf{m}(z) \Sigma (I_n + \mathbf{m}(z) \Sigma)^{-1} \mu_0 \rangle \right| \geq n^{-1/2+\vartheta} \sqrt{\Im \mathbf{m}(z) / \Im z} \right) \leq C n^{-D} \end{aligned}$$

holds for all $z \in [-1/c, 1/c] \times (0, 1/c]$. With $\kappa \equiv \kappa(z) \equiv \text{dist}(\Re z, \text{supp } \rho) \geq n^{-2/3+c}$, by further using the simple relation $\Im \mathbf{m}(z) / \Im z = \int \frac{\rho(dx)}{(\Re z - x)^2 + \Im^2 z} \leq \kappa^{-2}$, the error bound $n^{-1/2+\vartheta} \sqrt{\Im \mathbf{m}(z) / \Im z}$ in the above display can be replaced by $\kappa^{-1} n^{-1/2+\vartheta}$.

When $\phi^{-1} \geq 1 + 1/K$, according to (Bai and Silverstein, 2010, Theorem 6.3-(2)), $\text{supp } \rho \in (C_0^{-1}, C_0)$ for some constant $C_0 > 1$. Therefore, for $z \equiv z(b) \equiv -\eta/\phi + \sqrt{-1}b$ with a small enough $b > 0$ to be chosen later, it is easy to see that $\kappa \geq \kappa_0 \equiv (\eta/\phi) \vee C_0^{-1} \mathbf{1}_{\phi^{-1} \geq 1+1/K}$. Therefore, on an event $E_{1,0;s}(b)$ with $\mathbb{P}^\xi(E_{1,0;s}(b)) \geq 1 - C n^{-D}$,

$$\begin{aligned} \left| m^{-1} \langle \Sigma^{1/2} e_s, X^\top (\check{\Sigma} - z(0) I_m)^{-1} X \mu_0 \rangle \right. \\ \left. - \langle \Sigma^{1/2} e_s, \mathbf{m}(z(0)) \Sigma (I_n + \mathbf{m}(z(0)) \Sigma)^{-1} \mu_0 \rangle \right| \leq (I) + (II) + \kappa_0^{-1} n^{-1/2+\vartheta}, \end{aligned} \quad (120)$$

where

- (I) = $|m^{-1} \langle \Sigma^{1/2} e_s, X^\top (\check{\Sigma} - z(b) I_m)^{-1} X \mu_0 \rangle - m^{-1} \langle \Sigma^{1/2} e_s, X^\top (\check{\Sigma} - z(0) I_m)^{-1} X \mu_0 \rangle|$,
- (II) = $|\langle \Sigma^{1/2} e_s, \mathbf{m}(z(b)) \Sigma (I_n + \mathbf{m}(z(b)) \Sigma)^{-1} \mu_0 \rangle - \langle \Sigma^{1/2} e_s, \mathbf{m}(z(0)) \Sigma (I_n + \mathbf{m}(z(0)) \Sigma)^{-1} \mu_0 \rangle|$.

By a derivative calculation, it is easy to derive

$$(I) \lesssim (\|Z\|_{\text{op}} / \sqrt{n})^2 \cdot (\|(ZZ^\top / n)^{-1}\|_{\text{op}} \mathbf{1}_{\phi^{-1} \geq 1+1/K} \wedge \eta^{-1})^2 \cdot b.$$

Now by using the concentration result in (Rudelson and Vershynin, 2009, Theorem 1.1), on an event $E_{1,1;s}$ with $\mathbb{P}^\xi(E_{1,1;s}) \geq 1 - e^{-n/C}$, we have $(I) \leq Cb$.

For (II), using the boundedness of $\mathbf{m}(z(b))$ around 0 for $\phi^{-1} \geq 1 + 1/K$, we may estimate

$$\begin{aligned} (II) &\lesssim (\mathbf{1}_{\phi^{-1} \geq 1+1/K} \wedge \eta^{-1}) \cdot |\mathbf{m}(z(b)) - \mathbf{m}(z(0))| \\ &\leq (\mathbf{1}_{\phi^{-1} \geq 1+1/K} \wedge \eta^{-1}) \cdot \int_{C_0^{-1} \mathbf{1}_{\phi^{-1} \geq 1+1/K}}^{\infty} \frac{b}{|x - z(b)| |x - z(0)|} \rho(dx) \\ &\leq \left\{ C_0^2 \mathbf{1}_{\phi^{-1} \geq 1+1/K}^{-1} \wedge \eta^{-3} \right\} \cdot b. \end{aligned}$$

Combining the above estimates, for b chosen small enough, say, $b = n^{-100}$, on the event $E_{1,0;s}(n^{-100}) \cap E_{1,1;s}$,

$$|A_{1;s} - \langle \Sigma^{1/2} e_s, \mathbf{m}(-\eta/\phi) \Sigma (I_n + \mathbf{m}(-\eta/\phi) \Sigma)^{-1} \mu_0 - \mu_0 \rangle| \lesssim n^{-1/2+\vartheta}.$$

Using $\tau_{\eta,*}^{-1} = \mathbf{m}(-\eta/\phi)$ and the definition of $\widehat{\mu}_{(\Sigma,\mu_0)}^{\text{seq}}(\gamma_{\eta,*}; \tau_{\eta,*})$, recall $w_{\eta,*} = \Sigma^{1/2}(\widehat{\mu}_{(\Sigma,\mu_0)}^{\text{seq}}(\gamma_{\eta,*}; \tau_{\eta,*}) - \mu_0)$ defined in (77), we then have

$$\sup_{\mu_0 \in B_n(1)} \mathbb{P}^\xi \left(\max_{s \in [n]} |A_{1;s} - \langle e_s, \mathbb{E} w_{\eta,*} \rangle| \geq C n^{-1/2+\vartheta} \right) \leq C n^{-D}. \quad (121)$$

The term $A_{2;s}$ can be handled similarly, now reading off the (1, 2) element in (Knowles and Yin, 2017, Eqn. (3.10)), which shows that for any $\xi \in \mathbb{R}^m$,

$$\mathbb{P}^\xi \left(\max_{s \in [n]} |A_{2;s}| \geq C(\|\xi\|/\sqrt{m}) \cdot n^{-1/2+\vartheta} \right) \leq C n^{-D}. \quad (122)$$

Combining (119), (121) and (122), we have

$$\sup_{\mu_0 \in B_n(1), \xi \in \mathcal{E}_\vartheta} \mathbb{P}^\xi \left(\|\widehat{w}_\eta\|_\infty \geq \|\mathbb{E} w_{\eta,*}\|_\infty + C n^{-1/2+\vartheta} \right) \leq C n^{-D}. \quad (123)$$

Now we will construct $\mathcal{U}_\vartheta \subset B_n(1)$ with the desired volume estimate, and $\sup_{\mu_0 \in \mathcal{U}_\vartheta} \sup_{\eta \in \Xi_K} \|\mathbb{E} w_{\eta,*}\|_\infty \leq C n^{-1/2+\vartheta}$. To this end, we place a uniform prior on $\mu_0 \sim U_0 g_0 / \|g_0\|$, where $U_0 \sim \text{Unif}[0, 1]$ and $g_0 \sim \mathcal{N}(0, I_n)$ are independent of all other random variables. Then $\sup_{\eta \in \Xi_K} \|\mathbb{E} w_{\eta,*}\|_\infty \leq \sup_{\eta \in \Xi_K} \tau_{\eta,*} \|(\Sigma + \tau_{\eta,*} I_n)^{-1} \Sigma^{1/2} g_0\|_\infty / \|g_0\|$. Using Proposition 23-(3) and a standard Gaussian tail bound, $\mathbb{P}_{\mu_0}(\mathcal{U}_\vartheta \equiv \{\sup_{\eta \in \Xi_K} \|\mathbb{E} w_{\eta,*}\|_\infty \geq C_1 n^{-1/2+\vartheta}\}) \leq C e^{-n^{2\vartheta}/C}$. Moreover, $\mathbb{P}(\xi \notin \mathcal{E}_\vartheta) \leq e^{-n^{2\vartheta}/C}$. The pointwise-in- η delocalization claim on \widehat{w}_η follows. As $\eta \mapsto \|\widehat{w}_\eta\|_\infty$ is C -Lipschitz with exponentially high probability, the uniform version follows by a standard discretization and union bound argument.

(2). Let us consider delocalization for \widehat{v}_η . Using again (118), for any $t \in [m]$,

$$\begin{aligned} -\langle e_t, \widehat{v}_\eta \rangle &= n^{-1/2} \langle e_t, (\phi \check{\Sigma} + \eta I_m)^{-1} X \mu_0 \rangle + n^{-1/2} \langle e_t, (\phi \check{\Sigma} + \eta I_m)^{-1} \xi \rangle \\ &\equiv B_{1;t} + B_{2;t}. \end{aligned}$$

The term $B_{1;t}$ can be handled, by reading off the (2, 1) element in (Knowles and Yin, 2017, Eqn. (3.10)), which shows that

$$\sup_{\mu_0 \in B_n(1)} \mathbb{P}^\xi \left(\max_{t \in [m]} |B_{1;t}| \geq C n^{-1/2+\vartheta} \right) \leq C n^{-D}. \quad (124)$$

The term $B_{2;t}$ relies on the local law described by the (2, 2) element in (Knowles and Yin, 2017, Eqn. (3.10)): for any $\xi \in \mathbb{R}^m$,

$$\mathbb{P}^\xi \left(\max_{t \in [m]} |B_{2;t} - \phi^{-1} \mathbf{m}(-\eta/\phi) \xi_t| \geq C(\|\xi\|/\sqrt{m}) \cdot n^{-1/2+\vartheta} \right) \leq C n^{-D}. \quad (125)$$

Consequently, combining (124)-(125), we have

$$\sup_{\mu_0 \in B_n(1), \xi \in \mathcal{E}_\vartheta} \mathbb{P}^\xi \left(\|\widehat{v}_\eta\|_\infty \geq C n^{-1/2+\vartheta} \right) \leq C n^{-D}.$$

The claim follows. ■

D.3 Universality of the global cost optimum

Theorem 43. *Suppose Assumption A holds and the following hold for some $K > 0$.*

- $1/K \leq \phi^{-1} \leq K$, $\|\Sigma\|_{\text{op}} \vee \|\Sigma^{-1}\|_{\text{op}} \leq K$.
- Assumption B holds with $\sigma_\xi^2 \in [1/K, K]$.

Fix $\vartheta \in (0, 1/18)$. There exists some $C = C(K, \vartheta) > 0$ such that for $\rho_0 \leq 1/C$, $\eta \in \Xi_K$ and $\xi \in \mathcal{E}_\vartheta$,

$$\sup_{\mu_0 \in \mathcal{U}_\vartheta} \mathbb{P}^\xi \left(\left| \min_{w \in \mathbb{R}^n} H_{\eta; Z}(w) - \max_{\beta > 0} \min_{\gamma > 0} \bar{D}_\eta(\beta, \gamma) \right| \geq \rho_0 \right) \leq C \rho_0^{-3} \cdot n^{-1/6+3\vartheta}.$$

Here \mathcal{U}_ϑ is specified as in Proposition 40.

Proof Fix $\vartheta > 0$, $\mu_0 \in \mathcal{U}_\vartheta$ and $\xi \in \mathcal{E}_\vartheta$ as specified in Proposition 40. Let $L_n \equiv C_0 n^\vartheta$. By the same proposition, with \mathbb{P}^ξ -probability at least $1 - C_0 n^{-100}$,

$$\begin{aligned} \min_{w \in \mathbb{R}^n} H_{\eta; Z}(w) &= \min_{\|w\|_\infty \leq L_n/\sqrt{n}} \max_{\|v\|_\infty \leq L_n/\sqrt{n}} \left\{ \frac{1}{\sqrt{n}} \langle v, Zw \rangle - \frac{1}{\sqrt{n}} \langle v, \xi \rangle - \frac{\eta}{2} \|v\|^2 + F(w) \right\} \\ &= \min_{\|\tilde{w}\|_\infty \leq L_n} \max_{\|\tilde{v}\|_\infty \leq L_n} \left\{ \frac{1}{n^{3/2}} \langle \tilde{v}, Z\tilde{w} \rangle - \frac{1}{n} \langle \tilde{v}, \xi \rangle - \frac{\eta}{2n} \|\tilde{v}\|^2 + F(\tilde{w}/\sqrt{n}) \right\}, \end{aligned} \quad (126)$$

and

$$\min_{w \in \mathbb{R}^n} H_{\eta; G}(w) = \min_{\|\tilde{w}\|_\infty \leq L_n} \max_{\|\tilde{v}\|_\infty \leq L_n} \left\{ \frac{1}{n^{3/2}} \langle \tilde{v}, G\tilde{w} \rangle - \frac{1}{n} \langle \tilde{v}, \xi \rangle - \frac{\eta}{2n} \|\tilde{v}\|^2 + F(\tilde{w}/\sqrt{n}) \right\}. \quad (127)$$

By writing $Q(\tilde{v}, \tilde{w}) \equiv -\frac{1}{n} \langle \tilde{v}, \xi \rangle - \frac{\eta}{2n} \|\tilde{v}\|^2 + F(\tilde{w}/\sqrt{n})$, we have

$$\mathcal{N}_Q(L, \delta) \equiv \sup_{\substack{\|\tilde{v}\|_\infty \vee \|\tilde{v}'\|_\infty \leq L, \|\tilde{v} - \tilde{v}'\|_\infty \leq \delta, \\ \|\tilde{w}\|_\infty \vee \|\tilde{w}'\|_\infty \leq L, \|\tilde{w} - \tilde{w}'\|_\infty \leq \delta}} |Q(\tilde{v}, \tilde{w}) - Q(\tilde{v}', \tilde{w}')| \lesssim_K (1 \vee L) \delta \cdot \left(1 + \frac{\|\xi\|_1}{n} \right).$$

Now with $X_Q(\tilde{v}, \tilde{w}; Z) \equiv n^{-3/2} \langle \tilde{v}, Z\tilde{w} \rangle + Q(\tilde{v}, \tilde{w})$, for $\xi \in \mathcal{E}_\vartheta$, by applying Theorem 39, we have for any $\mathbb{T} \in C^3(\mathbb{R})$,

$$\begin{aligned} & \left| \mathbb{E}^\xi \mathbb{T} \left(\min_{\|\tilde{w}\|_\infty \leq L_n} \max_{\|\tilde{v}\|_\infty \leq L_n} X_Q(\tilde{v}, \tilde{w}; Z) \right) - \mathbb{E}^\xi \mathbb{T} \left(\min_{\|\tilde{w}\|_\infty \leq L_n} \max_{\|\tilde{v}\|_\infty \leq L_n} X_Q(\tilde{v}, \tilde{w}; G) \right) \right| \\ & \lesssim_K K_{\mathbb{T}} \cdot \inf_{\delta \in (0, 1)} \left\{ \sqrt{n} L_n \delta + L_n \delta + \log_+^{2/3}(L_n/\delta) \cdot n^{-1/6} L_n^2 \right\} \leq C_1 \cdot K_{\mathbb{T}} \cdot n^{-1/6+3\vartheta}. \end{aligned} \quad (128)$$

Replicating the last paragraph of proof of (Han and Shen, 2023, Theorem 2.3) (right above Section 4.3 therein), for any $z > 0, \rho_0 > 0$,

$$\begin{aligned} & \mathbb{P}^\xi \left(\min_{\|\tilde{w}\|_\infty \leq L_n} \max_{\|\tilde{v}\|_\infty \leq L_n} X_Q(\tilde{v}, \tilde{w}; Z) > z + 3\rho_0 \right) \\ & \leq \mathbb{P}^\xi \left(\min_{\|\tilde{w}\|_\infty \leq L_n} \max_{\|\tilde{v}\|_\infty \leq L_n} X_Q(\tilde{v}, \tilde{w}; G) > z + \rho_0 \right) + C \rho_0^{-3} n^{-1/6+3\vartheta}. \end{aligned}$$

Combined with (126)-(127), we have

$$\mathbb{P}^\xi \left(\min_{w \in \mathbb{R}^n} H_{\eta;Z}(w) > z + 3\rho_0 \right) \leq \mathbb{P}^\xi \left(\min_{w \in \mathbb{R}^n} H_{\eta;G}(w) > z + \rho_0 \right) + C_2 \rho_0^{-3} n^{-1/6+3\vartheta}.$$

In view of (89) (in Step 1 of the final proof of Theorem 3), for $\rho_0 \in (C_3 n^{-1/2+\vartheta}, 1/C_3)$, we take $z \equiv z_\eta \equiv \max_{\beta>0} \min_{\gamma>0} \bar{D}_\eta(\beta, \gamma)$ and $t \equiv \rho_0^2 n / C_3$ therein, so that for $\xi \in \mathcal{E}_\vartheta \subset \mathcal{E}_{1,\xi}(\rho_0/C_3^{1/2})$,

$$\mathbb{P}^\xi \left(\min_{w \in \mathbb{R}^n} H_{\eta;G}(w) > z_\eta + \rho_0 \right) \leq C_3 e^{-\rho_0^2 n / C_3}.$$

Combining the estimates, for $\xi \in \mathcal{E}_\vartheta$, $\rho_0 \in (C_3 n^{-1/2+\vartheta}, 1/C_3)$,

$$\mathbb{P}^\xi \left(\min_{w \in \mathbb{R}^n} H_{\eta;Z}(w) > z_\eta + 3\rho_0 \right) \leq C_4 \{ e^{-\rho_0^2 n / C_4} + \rho_0^{-3} n^{-1/6+3\vartheta} \}.$$

The first term above can be assimilated into the second one, and $\rho_0 \geq C_3 n^{-1/2+\vartheta}$ can be dropped. The lower bound follow similarly by utilizing (90). \blacksquare

D.4 Universality of the cost over exceptional sets

Theorem 44. *Suppose Assumption A holds and the following hold for some $K > 0$.*

- $1/K \leq \phi^{-1} \leq K$, $\|\Sigma\|_{\text{op}} \vee \|\Sigma^{-1}\|_{\text{op}} \leq K$.
- Assumption B with variance $\sigma_\xi^2 \in [1/K, K]$.

Fix $\vartheta \in (0, 1/18)$. Then there exists some $C = C(K, \vartheta) > 0$ such that for $\mathbf{g} : \mathbb{R}^n \rightarrow \mathbb{R}$ being 1-Lipschitz with respect to $\|\cdot\|_{\Sigma^{-1}}$, $\rho_0 \leq 1/C$, $\eta \in \Xi_K$ and $\xi \in \mathcal{E}_\vartheta$,

$$\sup_{\mu_0 \in \mathcal{U}_\vartheta} \mathbb{P}^\xi \left(\min_{w \in D_{\eta;C\rho_0^{1/2}}(\mathbf{g}) \cap B_{(2,\infty)}(C, \frac{L_n}{\sqrt{n}})} H_{\eta;Z}(w) \leq \max_{\beta>0} \min_{\gamma>0} \bar{D}_\eta(\beta, \gamma) + \rho_0 \right) \leq C \rho_0^{-6} \cdot n^{-1/6+3\vartheta}.$$

Here $B_{(2,\infty)}(C, L_n/\sqrt{n}) \equiv B_n(C) \cap L_\infty(L_n/\sqrt{n})$ with $L_n \equiv Cn^\vartheta$, and \mathcal{U}_ϑ is specified as in Proposition 40.

Proof Fix $\varepsilon, \vartheta > 0$, $\mu_0 \in \mathcal{U}_\vartheta$ and $\xi \in \mathcal{E}_\vartheta$ as specified in Proposition 40. We define a renormalized version of $D_{\varepsilon;\eta}(\mathbf{g})$ as

$$\tilde{D}_{\varepsilon;\eta}(\mathbf{g}) \equiv \{ \tilde{w} \in \mathbb{R}^n : |\mathbf{g}(\tilde{w}/\sqrt{n}) - \mathbb{E} \mathbf{g}(\tilde{w}_{\eta,*}/\sqrt{n})| \geq \varepsilon \},$$

where $\tilde{w}_{\eta,*} = \sqrt{n} w_{\eta,*}$.

(Step 1). Let $L_n \equiv C_0 n^\vartheta$. For any $z \in \mathbb{R}$ and $\rho_0 > 0$, with $Z_n \equiv Z/\sqrt{n}$,

$$\begin{aligned} & \mathbb{P}^\xi \left(\min_{w \in D_{\eta;\varepsilon}(\mathbf{g}) \cap B_{(2,\infty)}(C_0, \frac{L_n}{\sqrt{n}})} H_{\eta;Z}(w) \leq z + \rho_0 \right) \\ &= \mathbb{P}^\xi \left(\min_{w \in D_{\eta;\varepsilon}(\mathbf{g}) \cap B_{(2,\infty)}(C_0, \frac{L_n}{\sqrt{n}})} \left\{ F(w) + \frac{1}{2n\eta} \|Zw - \xi\|^2 \right\} \leq z + \rho_0 \right) \end{aligned} \tag{129}$$

$$= \mathbb{P}^\xi \left(\min_{\tilde{w} \in \tilde{D}_{\eta; \varepsilon}(\mathfrak{g}) \cap B_{(2, \infty)}(\sqrt{n}C_0, L_n)} \left\{ \eta F(\tilde{w}/\sqrt{n}) + \frac{1}{2n} \|Z_n \tilde{w} - \xi\|^2 \right\} \leq \eta(z + \rho_0) \right).$$

Now we may apply Theorem 38. To do so, let us write $f(\tilde{w}) \equiv \eta F(\tilde{w}/\sqrt{n})$ to match the notation. Then a simple calculation leads to

$$\mathcal{M}(L, \delta) \equiv \sup_{\|\tilde{w}\|_\infty \vee \|\tilde{w}'\|_\infty \leq L, \|\tilde{w} - \tilde{w}'\|_\infty \leq \delta} |f(\tilde{w}) - f(\tilde{w}')| \lesssim_K (1 \vee L)\delta,$$

Consequently, an application of Theorem 38 leads to

$$\begin{aligned} & \text{RHS of (129)} - C_1 (1 \vee (\eta\rho_0)^{-3}) L_n^2 n^{-1/6} \log^{2/3}(L_n n) \\ & \leq \mathbb{P}^\xi \left(\min_{\tilde{w} \in \tilde{D}_{\eta; \varepsilon}(\mathfrak{g}) \cap B_{(2, \infty)}(\sqrt{n}C_0, L_n)} \left\{ \eta F(\tilde{w}/\sqrt{n}) + \frac{1}{2n} \|G_n \tilde{w} - \xi\|^2 \right\} \leq \eta(z + 3\rho_0) \right) \\ & \leq \mathbb{P}^\xi \left(\min_{w \in D_{\eta; \varepsilon}(\mathfrak{g}) \cap B_{(2, \infty)}(C_0, \frac{L_n}{\sqrt{n}})} H_{\eta; G}(w) \leq z + 3\rho_0 \right) \\ & \leq \mathbb{P}^\xi \left(\min_{w \in D_{\eta; \varepsilon}(\mathfrak{g}) \cap B_n(C_0)} H_{\eta; G}(w) \leq z + 3\rho_0 \right). \end{aligned}$$

Here in the last inequality we simply drop the L_∞ constraint. Now for $C_2 n^{-1/2+\vartheta} \leq \rho_0 \leq 1/C_2$, by choosing $z \equiv z_\eta \equiv \max_{\beta>0} \min_{\gamma>0} \bar{D}_\eta(\beta, \gamma)$ and $t \equiv 2\rho_0^2 n/C_3$ in Theorem 32, where C_3 is the constant therein, we have

$$\begin{aligned} & \mathbb{P}^\xi \left(\min_{w \in D_{\eta; C_4 \rho_0^{1/2}}(\mathfrak{g}) \cap B_{(2, \infty)}(C, \frac{L_n}{\sqrt{n}})} H_{\eta; Z}(w) \leq \max_{\beta>0} \min_{\gamma>0} \bar{D}_\eta(\beta, \gamma) + \rho_0 \right) \\ & \leq C \left\{ e^{-\rho_0^2 n/C_3} + (\eta\rho_0)^{-3} \cdot n^{-1/6+3\vartheta} \right\} \leq C_4 \cdot (\eta\rho_0)^{-3} \cdot n^{-1/6+3\vartheta}. \end{aligned} \quad (130)$$

The constraints $\rho_0 \geq C_2 n^{-1/2+\vartheta}$ can be removed by enlarging C_4 if necessary.

(Step 2). In this step we shall trade the dependence of the above bound with respect to $\eta > 0$ with a possible worsened dependence on ρ_0 , primarily in the regime $\phi^{-1} \geq 1 + 1/K$. Fix $\eta_0 \in \Xi_K$. Let $\eta > 0$ be chosen later and $\eta_1 \equiv \eta_0 + \eta$. Without loss of generality we assume $\eta_0, \eta_1 \in \Xi_K$, so by (76) in Proposition 31, $|z_{\eta_1} - z_{\eta_0}| \leq C_5 \eta$. By enlarging C_5 if necessary we assume that C_5 exceeds the constant in Lemma 45. Using Lemma 45, for $\varepsilon = 2C_4 \rho_0^{1/2}$, with the choice $\eta = C_4 \rho_0 / C_5 \leq C_4 \rho_0^{1/2} / C_5$ (we assume without loss of generality $\rho_0 \leq 1$),

$$\begin{aligned} & \mathbb{P}^\xi \left(\min_{w \in D_{\eta_0; \varepsilon}(\mathfrak{g}) \cap B_{(2, \infty)}(C_0, \frac{L_n}{\sqrt{n}})} H_{\eta_0; Z}(w) \leq z_{\eta_0} + \rho_0 \right) \\ & \leq \mathbb{P}^\xi \left(\min_{w \in D_{\eta_0; \varepsilon}(\mathfrak{g}) \cap B_{(2, \infty)}(C_0, \frac{L_n}{\sqrt{n}})} H_{\eta_1; Z}(w) \leq z_{\eta_0} + \rho_0 \right) \quad (\text{since } H_{\eta_1; Z} \leq H_{\eta_0; Z}) \\ & \leq \mathbb{P}^\xi \left(\min_{w \in D_{\eta_1; (\varepsilon - C_5 \eta)_+}(\mathfrak{g}) \cap B_{(2, \infty)}(C_0, \frac{L_n}{\sqrt{n}})} H_{\eta_1; Z}(w) \leq z_{\eta_0} + \rho_0 \right) \quad (\text{by Lemma 45}) \\ & \leq \mathbb{P}^\xi \left(\min_{w \in D_{\eta_1; (\varepsilon - C_5 \eta)_+}(\mathfrak{g}) \cap B_{(2, \infty)}(C_0, \frac{L_n}{\sqrt{n}})} H_{\eta_1; Z}(w) \leq z_{\eta_1} + C_5 \eta + \rho_0 \right) \end{aligned}$$

$$\begin{aligned}
 &\leq \mathbb{P}^\xi \left(\min_{w \in D_{\eta; C_0 \rho_0^{1/2}}(\mathbf{g}) \cap B_{(2, \infty)}(C_0, \frac{L_n}{\sqrt{n}})} H_{\eta; Z}(w) \leq \max_{\beta > 0} \min_{\gamma > 0} \bar{D}_\eta(\beta, \gamma) + C\rho_0 \right) \\
 &\leq C \cdot (\eta_0 + \rho_0)^{-3} \rho_0^{-3} \cdot n^{-1/6+3\vartheta} \leq C \cdot \rho_0^{-6} n^{-1/6+3\vartheta}.
 \end{aligned}$$

The proof is complete by adjusting constants. \blacksquare

Lemma 45. *Suppose $\|\mu_0\| \vee \|\Sigma\|_{\text{op}} \vee \|\Sigma^{-1}\|_{\text{op}} \leq K$. Let $\mathbf{g} : \mathbb{R}^n \rightarrow \mathbb{R}$ be 1-Lipschitz with respect to $\|\cdot\|_{\Sigma^{-1}}$. Then there exists some constant $C = C(K) > 0$ such that for any $\varepsilon > 0, \eta_0, \eta_1 \in \Xi_K$ with $\eta_1 \geq \eta_0$, we have $D_{\eta_0; \varepsilon}(\mathbf{g}) \subset D_{\eta_1; (\varepsilon - C(\eta_1 - \eta_0))_+}(\mathbf{g})$.*

Proof Using the definition of $w_{\eta, *}$ in (77), we have

$$\begin{aligned}
 &|\mathbb{E} \mathbf{g}(w_{\eta_1, *}) - \mathbb{E} \mathbf{g}(w_{\eta_0, *})| \leq \mathbb{E} \|w_{\eta_1, *} - w_{\eta_0, *}\|_{\Sigma^{-1}} \\
 &= \mathbb{E} \left\| \hat{\mu}_{(\Sigma, \mu_0)}^{\text{seq}}(\gamma_{\eta_1, *}; \tau_{\eta_1, *}) - \hat{\mu}_{(\Sigma, \mu_0)}^{\text{seq}}(\gamma_{\eta_0, *}; \tau_{\eta_0, *}) \right\| \leq C \cdot (\eta_1 - \eta_0).
 \end{aligned}$$

Here the last inequality follows from the calculations in (94). So for any $w \in D_{\eta_0; \varepsilon}(\mathbf{g})$,

$$\varepsilon \leq |\mathbf{g}(w) - \mathbb{E} \mathbf{g}(w_{\eta_0, *})| \leq |\mathbf{g}(w) - \mathbb{E} \mathbf{g}(w_{\eta_1, *})| + C(\eta_1 - \eta_0).$$

This means $w \in D_{\eta_1; (\varepsilon - C(\eta_1 - \eta_0))_+}(\mathbf{g})$, as desired. \blacksquare

D.5 Proof of the universality Theorem 4 for $\hat{\mu}_{\eta; Z}$

Fix $\vartheta > 0, \mu_0 \in \mathcal{U}_\vartheta$ and $\xi \in \mathcal{E}_\vartheta$. Let $L_n \equiv C_0 n^\vartheta$, and $E_0 \equiv \{\hat{w}_{n; Z} \in B_{(2, \infty)}(C_0, L_n/\sqrt{n}) = B_n(C_0) \cap L_\infty(L_n/\sqrt{n})\}$. We assume that C_0 exceeds the constants in Proposition 40 and Theorem 44. By Proposition 40 and a simple ℓ_2 estimate, $\mathbb{P}^\xi(E_0^c) \leq C_0 n^{-100}$. We further let $z_\eta \equiv \max_{\beta > 0} \min_{\gamma > 0} \bar{D}_\eta(\beta, \gamma)$ for $\eta \geq 0$.

Let $\mathbf{g} : \mathbb{R}^n \rightarrow \mathbb{R}$ be 1-Lipschitz with respect to $\|\cdot\|_{\Sigma^{-1}}$. Then for $\rho_0 \leq 1/C_0$ and $\eta \in \Xi_K$, we have

$$\begin{aligned}
 &\mathbb{P}^\xi (\hat{w}_{\eta; Z} \in D_{\eta; C_0 \rho_0^{1/2}}(\mathbf{g})) \\
 &\leq \mathbb{P}^\xi (\hat{w}_{\eta; Z} \in D_{\eta; C_0 \rho_0^{1/2}}(\mathbf{g}) \cap B_{(2, \infty)}(C_0, L_n/\sqrt{n})) + \mathbb{P}^\xi(E_0^c) \\
 &\leq \mathbb{P}^\xi \left(\min_{w \in B_{(2, \infty)}(C_0, \frac{L_n}{\sqrt{n}})} H_{\eta; Z}(w) \geq z_\eta + \rho_0 \right) \\
 &\quad + \mathbb{P}^\xi \left(\min_{w \in D_{\eta; C_0 \rho_0^{1/2}}(\mathbf{g}) \cap B_{(2, \infty)}(C_0, \frac{L_n}{\sqrt{n}})} H_{\eta; Z}(w) \leq z_\eta + 2\rho_0 \right) + C_0 n^{-100}.
 \end{aligned}$$

Here in the last inequality we used the simple fact that

$$\begin{aligned}
 &\left\{ \min_{w \in B_{(2, \infty)}(C_0, \frac{L_n}{\sqrt{n}})} H_{\eta; Z}(w) < z_\eta + \rho_0 \right\} \cap \left\{ \min_{w \in D_{\eta; C_0 \rho_0^{1/2}}(\mathbf{g}) \cap B_{(2, \infty)}(C_0, \frac{L_n}{\sqrt{n}})} H_{\eta; Z}(w) > z_\eta + 2\rho_0 \right\} \\
 &\subset \left\{ \hat{w}_{\eta; Z} \notin D_{\eta; C_0 \rho_0^{1/2}}(\mathbf{g}) \cap B_{(2, \infty)}(C_0, L_n/\sqrt{n}) \right\}.
 \end{aligned}$$

Invoking Theorems 43 and 44, by enlarging C_0 if necessary, we have for $\rho_0 \leq 1/C_0$ and $\eta \in \Xi_K$,

$$\mathbb{P}^\xi \left(|\mathbf{g}(\widehat{w}_{\eta;Z}) - \mathbb{E} \mathbf{g}(w_{\eta,*})| \geq \rho_0^{1/2} \right) \leq C_0 \cdot \rho_0^{-6} n^{-1/6+3\vartheta},$$

or equivalently, for $\mathbf{g}_0 : \mathbb{R}^n \rightarrow \mathbb{R}$ being 1-Lipschitz with respect to $\|\cdot\|$,

$$\mathbb{P}^\xi \left(\left| \mathbf{g}_0(\widehat{\mu}_{\eta;Z}) - \mathbb{E} \mathbf{g}_0(\widehat{\mu}_{(\Sigma,\mu_0)}^{\text{seq}}(\gamma_{\eta,*}; \tau_{\eta,*})) \right| \geq \rho_0 \right) \leq C_0 \cdot \rho_0^{-12} n^{-1/6+3\vartheta}.$$

Now we may follow Step 4 in the proof of Theorem 3 to strengthen the above statement to a uniform one in η ; we only sketch the differences below. Using (92) with G therein replaced by Z , and the assumption $\|\Sigma^{-1}\|_{\text{op}} \leq K$, we arrive at a modified form of (93): on an event E_1 with $\mathbb{P}^\xi(E_1) \geq 1 - C_1 e^{-n/C_1}$, for any $\eta_1, \eta_2 \in \Xi_K$,

$$\|\widehat{\mu}_{\eta_1;Z} - \widehat{\mu}_{\eta_2;Z}\| \leq C_1 |\eta_1 - \eta_2|. \quad (131)$$

Using (94) with $\|\Sigma^{-1}\|_{\text{op}} \leq K$, we arrive at a modified form of (95): for any $\eta_1, \eta_2 \in \Xi_K$,

$$\left| \mathbb{E} \mathbf{g}_0(\widehat{\mu}_{(\Sigma,\mu_0)}^{\text{seq}}(\gamma_{\eta_1,*}; \tau_{\eta_1,*})) - \mathbb{E} \mathbf{g}_0(\widehat{\mu}_{(\Sigma,\mu_0)}^{\text{seq}}(\gamma_{\eta_2,*}; \tau_{\eta_2,*})) \right| \leq C_1 |\eta_1 - \eta_2|.$$

Now using a standard discretization and a union bound, we have

$$\mathbb{P}^\xi \left(\sup_{\eta \in \Xi_K} \left| \mathbf{g}_0(\widehat{\mu}_{\eta;Z}) - \mathbb{E} \mathbf{g}_0(\widehat{\mu}_{(\Sigma,\mu_0)}^{\text{seq}}(\gamma_{\eta,*}; \tau_{\eta,*})) \right| \geq \rho_0 \right) \leq C_2 \cdot \rho_0^{-13} n^{-1/6+3\vartheta}.$$

The proof is complete by taking expectation with respect to ξ and note that $\mathbb{P}(\xi \in \mathcal{E}_\vartheta) \geq 1 - C e^{-n^{2\vartheta}/C}$ as in Proposition 40. \square

D.6 Proof of the universality Theorem 4 for $\widehat{r}_{\eta;Z}$

Proposition 46. *Suppose Assumption A holds and the following hold for some $K > 0$.*

- $1/K \leq \phi^{-1}, \eta \leq K, \|\Sigma\|_{\text{op}} \vee \|\Sigma^{-1}\|_{\text{op}} \leq K$.
- *Assumption B holds with $\sigma_\xi^2 \in [1/K, K]$.*

Fix $\vartheta \in (0, 1/18)$. Then there exists some $C = C(K, \vartheta) > 0$ such that for any 1-Lipschitz function $\mathbf{h} : \mathbb{R}^m \rightarrow \mathbb{R}$, $\rho_0 \leq 1/C$, $\eta \in \Xi_K$ and $\xi \in \mathcal{E}_\vartheta$,

$$\sup_{\mu_0 \in \mathcal{U}_\vartheta} \mathbb{P}^\xi \left(\max_{\substack{v \in D_{\eta;C\rho_0} \\ \eta; C\rho_0}} \min_{\substack{w \in \mathbb{R}^n \\ L_\infty(\frac{L_n}{\sqrt{n}})}} \mathbf{h}_{\eta;Z}(w, v) \geq \max_{\beta > 0} \min_{\gamma > 0} \overline{D}_\eta(\beta, \gamma) - \rho_0 \right) \leq C \rho_0^{-3} n^{-1/6+3\vartheta}.$$

Here $L_n \equiv Cn^\vartheta$, and \mathcal{U}_ϑ is specified as in Proposition 40.

Proof Fix $\varepsilon, \vartheta > 0$, $\mu_0 \in \mathcal{U}_\vartheta$ and $\xi \in \mathcal{E}_\vartheta$ as specified in Proposition 40. We define a renormalized version of $D_{\varepsilon;\eta}(\mathbf{h})$ as

$$\widetilde{D}_{\varepsilon;\eta}(\mathbf{h}) \equiv \left\{ \widetilde{r} \in \mathbb{R}^m : |\mathbf{h}(\widetilde{r}/\sqrt{n}) - \mathbb{E}^\xi \mathbf{h}(\widetilde{r}_{\eta,*}/\sqrt{n})| \geq \varepsilon \right\},$$

where $\tilde{r}_{\eta,*} = \sqrt{n}r_{\eta,*}$. Let $L_n = C_0n^\vartheta$ and $Q(\tilde{v}, \tilde{w})$ be defined as in the proof of Theorem 43. Then we have,

$$\begin{aligned}
 & \max_{v \in D_{\eta;\varepsilon}(\mathbf{h}) \cap L_\infty(L_n/\sqrt{n})} \min_{w \in L_\infty(L_n/\sqrt{n})} h_{\eta;Z}(w, v) \\
 &= \max_{\tilde{v} \in \tilde{D}_{\eta;\varepsilon}(\mathbf{h}) \cap L_\infty(L_n)} \min_{\tilde{w} \in L_\infty(L_n)} h_{\eta;Z}(\tilde{w}/\sqrt{n}, \tilde{v}/\sqrt{n}) \\
 &= \max_{\tilde{v} \in \tilde{D}_{\eta;\varepsilon}(\mathbf{h}) \cap L_\infty(L_n)} \min_{\tilde{w} \in L_\infty(L_n)} \left\{ \frac{1}{n^{3/2}} \langle \tilde{v}, Z\tilde{w} \rangle - \frac{1}{n} \langle \tilde{v}, \xi \rangle + F(\tilde{w}/\sqrt{n}) - \frac{\eta \|\tilde{v}\|^2}{2n} \right\} \\
 &= \max_{\tilde{v} \in \tilde{D}_{\eta;\varepsilon}(\mathbf{h}) \cap L_\infty(L_n)} \min_{\tilde{w} \in L_\infty(L_n)} \left\{ \frac{1}{n^{3/2}} \langle \tilde{v}, Z\tilde{w} \rangle + Q(\tilde{v}, \tilde{w}) \right\}.
 \end{aligned}$$

Using the comparison inequality in Theorem 39 and a similar calculation as in (128), with $s_n(\rho_0) \equiv \rho_0^{-3}n^{-1/6+3\vartheta}$,

$$\begin{aligned}
 & \mathbb{P}^\xi \left(\max_{v \in D_{\eta;\varepsilon}(\mathbf{h}) \cap L_\infty(L_n/\sqrt{n})} \min_{w \in L_\infty(L_n/\sqrt{n})} h_{\eta;Z}(w, v) \geq z - \rho_0 \right) \\
 &= \mathbb{P}^\xi \left(\max_{\tilde{v} \in \tilde{D}_{\eta;\varepsilon}(\mathbf{h}) \cap L_\infty(L_n)} \min_{\tilde{w} \in L_\infty(L_n)} \left\{ \frac{1}{n^{3/2}} \langle \tilde{v}, Z\tilde{w} \rangle + Q(\tilde{v}, \tilde{w}) \right\} \geq z - \rho_0 \right) \\
 &\leq \mathbb{P}^\xi \left(\max_{\tilde{v} \in \tilde{D}_{\eta;\varepsilon}(\mathbf{h}) \cap L_\infty(L_n)} \min_{\tilde{w} \in L_\infty(L_n)} \left\{ \frac{1}{n^{3/2}} \langle \tilde{v}, G\tilde{w} \rangle + Q(\tilde{v}, \tilde{w}) \right\} \geq z - 3\rho_0 \right) + C s_n(\rho_0) \\
 &= \mathbb{P}^\xi \left(\max_{v \in D_{\eta;\varepsilon}(\mathbf{h}) \cap L_\infty(L_n/\sqrt{n})} \min_{w \in L_\infty(L_n/\sqrt{n})} h_{\eta;G}(w, v) \geq z - 3\rho_0 \right) + C_1 s_n(\rho_0).
 \end{aligned}$$

Using the convex Gaussian min-max theorem (cf. Theorem 14),

$$\begin{aligned}
 & \mathbb{P}^\xi \left(\max_{v \in D_{\eta;\varepsilon}(\mathbf{h}) \cap L_\infty(L_n/\sqrt{n})} \min_{w \in L_\infty(L_n/\sqrt{n})} h_{\eta;Z}(w, v) \geq z - \rho_0 \right) \\
 &\leq 2 \mathbb{P} \left(\max_{v \in D_{\eta;\varepsilon}(\mathbf{h}) \cap L_\infty(L_n/\sqrt{n})} \min_{w \in L_\infty(L_n/\sqrt{n})} \ell_\eta(w, v) \geq z - 3\rho_0 \right) + C_1 s_n(\rho_0). \quad (132)
 \end{aligned}$$

On the other hand, using the definition of $w_{\eta,*}$ in (77), and the fact that for any $\mu_0 \in \mathcal{U}_\vartheta$, $\|\mathbb{E} w_{\eta,*}\|_\infty \leq L_n/\sqrt{n}$, we have $\mathbb{P}(\|w_{\eta,*}\|_\infty \geq L_n/\sqrt{n}) \leq C e^{-n^{2\vartheta}/C}$. Combined with (132), we have

$$\begin{aligned}
 & \mathbb{P}^\xi \left(\max_{v \in D_{\eta;\varepsilon}(\mathbf{h}) \cap L_\infty(L_n/\sqrt{n})} \min_{w \in L_\infty(L_n/\sqrt{n})} h_{\eta;Z}(w, v) \geq z - \rho_0 \right) \\
 &\leq 2 \mathbb{P} \left(\max_{v \in D_{\eta;\varepsilon}(\mathbf{h}) \cap L_\infty(L_n/\sqrt{n})} \ell_\eta(w_{\eta,*}, v) \geq z - 3\rho_0 \right) + C_2 s_n(\rho_0).
 \end{aligned}$$

In view of (108), now by choosing $z \equiv z_\eta \equiv \max_{\beta>0} \min_{\gamma>0} \bar{D}_\eta(\beta, \gamma)$ and $\varepsilon \equiv C_3 \rho_0^{1/2}$, for $\rho_0 \geq C_4 n^{-1/2+\vartheta}$, $\xi \in \mathcal{E}_\vartheta \subset \mathcal{E}_{1,\xi}(\rho_0/C)$, it follows that

$$\begin{aligned}
 & \mathbb{P}^\xi \left(\max_{v \in D_{\eta;C_3\rho_0^{1/2}}(\mathbf{h}) \cap L_\infty(L_n/\sqrt{n})} \min_{w \in \mathbb{R}^n} h_{\eta;Z}(w, v) \geq z_\eta - \rho_0 \right) \\
 &\leq \mathbb{P}^\xi \left(\max_{v \in D_{\eta;C_3\rho_0^{1/2}}(\mathbf{h}) \cap L_\infty(L_n/\sqrt{n})} \min_{w \in L_\infty(L_n/\sqrt{n})} h_{\eta;Z}(w, v) \geq z_\eta - \rho_0 \right) \leq C_4 s_n(\rho_0).
 \end{aligned}$$

The claim follows by adjusting constants. \blacksquare

Proof [Proof of Theorem 4 for $\widehat{r}_{\eta;Z}$] Fix $\vartheta > 0$, $\mu_0 \in \mathcal{U}_\vartheta$ and $\xi \in \mathcal{E}_\vartheta$ as specified in Proposition 40. We continue writing $z_\eta \equiv \max_{\beta>0} \min_{\gamma>0} \overline{D}_\eta(\beta, \gamma)$ in the proof. Using the delocalization results in Proposition 40, on an event E_0 with $\mathbb{P}^\xi(E_0) \geq 1 - C_0 n^{-100}$, we have $\|\widehat{w}_{\eta;Z}\|_\infty \vee \|\widehat{r}_{\eta;Z}\|_\infty \leq L_n/\sqrt{n}$ with $L_n = C_0 n^\vartheta$. Using Theorem 43, for $\rho_0 \leq 1/C$, and $\eta \in \Xi_K$, by possibly adjusting $C_0 > 0$,

$$\begin{aligned} & \mathbb{P}^\xi \left(\max_{v \in L_\infty(L_n/\sqrt{n})} \min_{w \in \mathbb{R}^m} h_{\eta;Z}(w, v) \leq z_\eta - \rho_0/2 \right) \\ & \leq \mathbb{P}^\xi \left(\min_{w \in \mathbb{R}^m} H_{\eta;Z}(w) \leq z_\eta - \rho_0/2 \right) + \mathbb{P}^\xi(E_0^c) \leq C_0 \rho_0^{-3} \cdot n^{-1/6+3\vartheta}. \end{aligned} \quad (133)$$

Let us take $C_1 > 0$ to be the constant in Proposition 46. By noting that

$$\begin{aligned} & \left\{ \max_{v \in L_\infty(L_n/\sqrt{n})} \min_{w \in \mathbb{R}^m} h_{\eta;Z}(w, v) > z_\eta - \rho_0/2 \right\} \\ & \cap \left\{ \max_{v \in D_{\eta;C_1\rho_0^{1/2}}(\mathbf{h}) \cap L_\infty(L_n/\sqrt{n})} \min_{w \in \mathbb{R}^n} h_{\eta;Z}(w, v) < z_\eta - \rho_0 \right\} \\ & \subset \left\{ \widehat{v}_{\eta;Z} \notin D_{\eta;C_1\rho_0^{1/2}}(\mathbf{h}) \cap L_\infty(L_n/\sqrt{n}) \right\}, \end{aligned}$$

it follows from (133) and Proposition 46 that

$$\begin{aligned} & \mathbb{P}^\xi \left(\widehat{v}_{\eta;Z} \in D_{\eta;C_1\rho_0^{1/2}}(\mathbf{h}) \right) \\ & \leq \mathbb{P}^\xi \left(\widehat{v}_{\eta;Z} \in D_{\eta;C_1\rho_0^{1/2}}(\mathbf{h}) \cap L_\infty(L_n/\sqrt{n}) \right) + \mathbb{P}^\xi \left(\widehat{v}_{\eta;Z} \notin L_\infty(L_n/\sqrt{n}) \right) \\ & \leq \mathbb{P}^\xi \left(\max_{v \in L_\infty(L_n/\sqrt{n})} \min_{w \in \mathbb{R}^m} h_{\eta;Z}(w, v) \leq z_\eta - \rho_0/2 \right) \\ & \quad + \mathbb{P}^\xi \left(\max_{v \in D_{\eta;C_1\rho_0^{1/2}}(\mathbf{h}) \cap L_\infty(L_n/\sqrt{n})} \min_{w \in \mathbb{R}^n} h_{\eta;Z}(w, v) \geq z_\eta - \rho_0 \right) + \mathbb{P}^\xi(E_0^c) \\ & \leq C \rho_0^{-3} \cdot n^{-1/6+3\vartheta}. \end{aligned}$$

Finally we only need to extend the above display to a uniform control over $\eta \in [1/K, K]$ by continuity arguments similar to Step 5 of the proof of Theorem 3 for $\widehat{r}_{\eta;G}$. By (113) (where G therein is replaced by Z) and (131), on an event E_1 with $\mathbb{P}^\xi(E_1) \geq 1 - C e^{-n/C}$, for any $\eta_1, \eta_2 \in [1/K, K]$,

$$\|\widehat{r}_{\eta_1;Z} - \widehat{r}_{\eta_2;Z}\| \leq C |\eta_1 - \eta_2|.$$

On the other hand, (115) remains valid, so we may proceed with an ε -net argument over $[1/K, K]$ to conclude. \blacksquare

Appendix E. Proof of Theorem 6

To keep notation simple, we work with $\mathbf{A} = I_n$ and write $\Gamma_{\eta;(\Sigma, \|\mu_0\|)}^{I_n} = \Gamma_{\eta;(\Sigma, \|\mu_0\|)}$. The general case follows from minor modifications.

Lemma 47. *Suppose the conditions in Theorem 6 hold for some $K > 0$. Fix $q \in [1, \infty)$. There exists some constant $c = c(K, q) > 0$ such that $n^{\frac{1}{2} - \frac{1}{q}} \mathbb{E} \|\widehat{\mu}_{(\Sigma, \mu_0)}^{\text{seq}}(\gamma_{\eta, *}; \tau_{\eta, *}) - \mu_0\|_q \geq c$ uniformly in $\eta \in \Xi_K$.*

Proof We may write $\mathbb{E} \|\widehat{\mu}_{(\Sigma, \mu_0)}^{\text{seq}}(\gamma_{\eta, *}; \tau_{\eta, *}) - \mu_0\|_q = \mathbb{E} \left(\sum_{j=1}^n |a_j + b_j g_j|^q \right)^{1/q}$ for some $a_j, b_j \in \mathbb{R}$ with $b_j \asymp 1$, and $g_j \sim \mathcal{N}(0, 1/n)$ not necessarily independent of each other. So for some $c_j \in \mathbb{R}$,

$$\mathbb{E} \|\widehat{\mu}_{(\Sigma, \mu_0)}^{\text{seq}}(\gamma_{\eta, *}; \tau_{\eta, *}) - \mu_0\|_q \gtrsim \mathbb{E} \left(\sum_{j=1}^n |c_j + g_j|^q \right)^{1/q}.$$

If $\sum_{j=1}^n |c_j|^q \geq C_0 \sum_{j=1}^n \mathbb{E} |g_j|^q$ for a large enough $C_0 > 0$, the lower bound follows trivially. Otherwise, with $Z \equiv \sum_{j=1}^n |c_j + g_j|^q$, we have $\mathbb{E} Z \geq \sum_{j=1}^n \inf_{c \in \mathbb{R}} \mathbb{E} |c + g_j|^q \gtrsim n^{1-q/2}$ and $\mathbb{E} Z^2 \lesssim \mathbb{E} \left(\sum_{j=1}^n (|g_j|^q + \mathbb{E} |g_j|^q) \right)^2 \lesssim (n^{1-q/2})^2$, so by Paley-Zygmund inequality, $\mathbb{P}(Z \geq \mathbb{E} Z/2) \geq (\mathbb{E} Z)^2 / (4 \mathbb{E} Z^2) \geq c_0$ for some $c_0 > 0$. In other words, on an event E_0 with $\mathbb{P}(E_0) \geq c_0$, $Z \geq c_0 n^{1-q/2}$. Using the above display, this means that $\mathbb{E} \|\widehat{\mu}_{(\Sigma, \mu_0)}^{\text{seq}}(\gamma_{\eta, *}; \tau_{\eta, *}) - \mu_0\|_q \gtrsim \mathbb{E} Z^{1/q} \geq \mathbb{E} Z^{1/q} \mathbf{1}_{E_0} \gtrsim n^{1/q-1/2}$. \blacksquare

Lemma 48. *Suppose the conditions in Theorem 6 hold for some $K > 0$. Fix $q \in [1, \infty)$. Then there exist constants $C = C(K, q) > 1$, $\vartheta = \vartheta(q) \in (0, 1/50)$, and a measurable set $\mathcal{U}_\vartheta \subset B_n(1)$ with $\text{vol}(\mathcal{U}_\vartheta) / \text{vol}(B_n(1)) \geq 1 - C e^{-n^\vartheta/C}$, such that*

$$\sup_{\mu_0 \in \mathcal{U}_\vartheta} n^{\frac{1}{2} - \frac{1}{q}} \left| \mathbb{E} \|\widehat{\mu}_{(\Sigma, \mu_0)}^{\text{seq}}(\gamma_{\eta, *}; \tau_{\eta, *}) - \mu_0\|_q - n^{-\frac{1}{2}} \|\text{diag}(\Gamma_{\eta; (\Sigma, \|\mu_0\|)})\|_{q/2}^{1/2} M_q \right| \leq C n^{-\vartheta}.$$

Here $M_q = \mathbb{E}^{1/q} |\mathcal{N}(0, 1)|^q$.

Proof We write $\tau_{\eta, *} = \tau_\eta$, $\gamma_{\eta, *} = \gamma_\eta$ for notational simplicity in the proof. All the constants in $\lesssim, \gtrsim, \asymp$ below may depend on K, q . Recall the general fact $\|x\|_q \leq n^{-\frac{1}{2} + \frac{1}{q\wedge 2}} \|x\|_{q\wedge 2}^{\frac{2}{q\wedge 2}} \|x\|_\infty^{1 - \frac{2}{q\wedge 2}}$ for $x \in \mathbb{R}^n$ and $q \in (0, \infty)$.

By Proposition 52 below, for any $\vartheta \in (0, 1/2)$, there exists some $\mathcal{U}_\vartheta \subset B_n(1)$ with $\text{vol}(\mathcal{U}_\vartheta) / \text{vol}(B_n(1)) \geq 1 - C e^{-n^{1-2\vartheta}/C}$, such that $\sup_{\mu_0 \in \mathcal{U}_\vartheta} \sup_{\eta \in \Xi_K} |\gamma_\eta^2 - \tilde{\gamma}_\eta^2(\|\mu_0\|)| \leq n^{-\vartheta}$. Consequently, uniformly in $\mu_0 \in \mathcal{U}_\vartheta$ and $\eta \in \Xi_K$,

$$\begin{aligned} & \left| \mathbb{E} \|\widehat{\mu}_{(\Sigma, \mu_0)}^{\text{seq}}(\gamma_\eta; \tau_\eta) - \mu_0\|_q - \mathbb{E} \|\widehat{\mu}_{(\Sigma, \mu_0)}^{\text{seq}}(\tilde{\gamma}_\eta(\|\mu_0\|); \tau_\eta) - \mu_0\|_q \right| \\ & \lesssim |\gamma_\eta - \tilde{\gamma}_\eta(\|\mu_0\|)| \cdot n^{-1/2} \mathbb{E} \|(\Sigma + \tau_\eta I)^{-1} \Sigma^{1/2} g\|_q \\ & \leq n^{-\frac{1}{2} - \vartheta} \cdot n^{-\frac{1}{2} + \frac{1}{q\wedge 2}} \cdot \mathbb{E} \left\{ \|(\Sigma + \tau_\eta I)^{-1} \Sigma^{1/2} g\|_{q\wedge 2}^{\frac{2}{q\wedge 2}} \cdot \|(\Sigma + \tau_\eta I)^{-1} \Sigma^{1/2} g\|_\infty^{1 - \frac{2}{q\wedge 2}} \right\} \\ & \lesssim n^{-\frac{1}{2} + \frac{1}{q} - \vartheta} (\log n)^{\frac{1}{2} - \frac{1}{q\wedge 2}}. \end{aligned} \tag{134}$$

For $g' \in \mathbb{R}^n$, let $\mathbf{f}(g') \equiv n^{-1/2} \|(\Sigma + \tau_\eta I)^{-1} g'\|_q$, and

$$\mathbf{F}_{\|\mu_0\|}(g') \equiv n^{-1/2} \|(\Sigma + \tau_\eta I)^{-1} (-\tau_\eta \|\mu_0\| g' + \tilde{\gamma}_\eta(\|\mu_0\|) \Sigma^{1/2} g)\|_q,$$

$$\mathbf{F}_{\|\mu_0\|,0}(g') \equiv \left\| (\Sigma + \tau_\eta I)^{-1} \left(-\tau_\eta \|\mu_0\| \frac{g'}{\|g'\|} + \tilde{\gamma}_\eta(\|\mu_0\|) \Sigma^{1/2} \frac{g}{\sqrt{n}} \right) \right\|_q.$$

Then for $g'_1, g'_2 \in \mathbb{R}^n$,

$$|\mathbf{F}_{\|\mu_0\|}(g'_1) - \mathbf{F}_{\|\mu_0\|}(g'_2)| \vee |\mathbf{f}(g'_1) - \mathbf{f}(g'_2)| \lesssim n^{-1+\frac{1}{q\lambda^2}} \|g'_1 - g'_2\|.$$

By Gaussian concentration inequality, for any $\vartheta \in (0, 1/2)$, we may find some $\mathcal{G}_{\vartheta, \|\mu_0\|} \subset \mathbb{R}^n$ with $\mathbb{P}(g_0 \in \mathcal{G}_{\vartheta, \|\mu_0\|}) \geq 1 - Ce^{-n^{2\vartheta}/C}$, $g_0 \sim \mathcal{N}(0, I_n)$, such that uniformly in $g' \in \mathcal{G}_{\vartheta, \|\mu_0\|}$,

$$\begin{aligned} \max \left\{ \left| \|g'\| - \sqrt{n} \right|, n^{1-\frac{1}{q\lambda^2}} |\mathbf{F}_{\|\mu_0\|}(g') - \mathbb{E}_{g_0} \mathbf{F}_{\|\mu_0\|}(g_0)|, \right. \\ \left. n^{1-\frac{1}{q\lambda^2}} |\mathbf{f}(g') - \mathbb{E} \mathbf{f}(g_0)| \right\} \leq n^\vartheta. \end{aligned} \quad (135)$$

As $\mathbb{E} \mathbf{f}(g_0) = n^{-1/2} \mathbb{E} \|(\Sigma + \tau_\eta I)g\|_q \lesssim n^{-\frac{1}{2}+\frac{1}{q}} (\log n)^{\frac{1}{2}-\frac{1}{q\sqrt{2}}}$, for ϑ small enough, uniformly in $g' \in \mathcal{G}_{\vartheta, \|\mu_0\|}$,

$$\begin{aligned} |\mathbf{F}_{\|\mu_0\|}(g') - \mathbf{F}_{\|\mu_0\|,0}(g')| &\lesssim |\mathbf{f}(g')| \cdot |1 - \sqrt{n}/\|g'\|| \\ &\lesssim (\mathbb{E} \mathbf{f}(g_0) + n^{-1+\frac{1}{q\lambda^2}+\vartheta}) \cdot n^{-\frac{1}{2}+\vartheta} \lesssim n^{-1+\frac{1}{q}+\vartheta} (\log n)^{\frac{1}{2}-\frac{1}{q\sqrt{2}}}. \end{aligned} \quad (136)$$

Combining (135)-(136), for ϑ small enough,

$$\sup_{g' \in \mathcal{G}_{\vartheta, \|\mu_0\|}} n^{1-\frac{1}{q\lambda^2}} |\mathbf{F}_{\|\mu_0\|,0}(g') - \mathbb{E}_{g_0} \mathbf{F}_{\|\mu_0\|}(g_0)| \lesssim n^\vartheta. \quad (137)$$

Now let $\partial \mathcal{G}_{\vartheta, \|\mu_0\|} \equiv \{g'/\|g'\| : g' \in \mathcal{G}_{\vartheta, \|\mu_0\|}\} \subset \partial B_n(1)$. Using that $\{g_0 \in \mathcal{G}_{\vartheta, \|\mu_0\|}\} \subset \{g_0/\|g_0\| \in \partial \mathcal{G}_{\vartheta, \|\mu_0\|}\}$, we have $\mathbb{P}(g_0/\|g_0\| \in \partial \mathcal{G}_{\vartheta, \|\mu_0\|}) \geq \mathbb{P}(g_0 \in \mathcal{G}_{\vartheta, \|\mu_0\|}) \geq 1 - Ce^{-n^{2\vartheta}/C}$. So with

$$\mathcal{V}_\vartheta \equiv \{\mu_0 = U_0 g' : U_0 \in [0, 1], g' \in \partial \mathcal{G}_{\vartheta, U_0}\} \subset B_n(1),$$

we have $\mathbb{P}(\text{Unif}(B_n(1)) \in \mathcal{V}_\vartheta) = \mathbb{E}_{U_0} \mathbb{P}_{g_0}(g_0/\|g_0\| \in \partial \mathcal{G}_{\vartheta, U_0}) \geq 1 - Ce^{-n^{2\vartheta}/C}$. In other words, for this constructed set \mathcal{V}_ϑ , we have the desired volume estimate $\text{vol}(\mathcal{V}_\vartheta)/\text{vol}(B_n(1)) \geq 1 - Ce^{-n^{2\vartheta}/C}$, and by (137),

$$n^{1-\frac{1}{q\lambda^2}} \sup_{\mu_0 \in \mathcal{V}_\vartheta} \left| \mathbb{E} \|\hat{\mu}_{(\Sigma, \mu_0)}^{\text{seq}}(\tilde{\gamma}_\eta(\|\mu_0\|); \tau_\eta) - \mu_0\|_q - \mathbb{E}_{g_0} \mathbf{F}_{\|\mu_0\|}(g_0) \right| \lesssim n^\vartheta. \quad (138)$$

On the other hand, using the definition of $\Gamma_{\eta;(\Sigma, \|\mu_0\|)}$ in (16), we may compute

$$\mathbb{E}_{g_0} \mathbf{F}_{\|\mu_0\|}(g_0) = \mathbb{E} \left\| \Gamma_{\eta;(\Sigma, \|\mu_0\|)}^{1/2} g/\sqrt{n} \right\|_q. \quad (139)$$

Combining (134), (138) and (139), for ϑ chosen small enough,

$$\begin{aligned} \sup_{\mu_0 \in \mathcal{U}_\vartheta \cap \mathcal{V}_\vartheta} n^{1/2-1/q} \left| \mathbb{E} \|\hat{\mu}_{(\Sigma, \mu_0)}^{\text{seq}}(\gamma_\eta; \tau_\eta) - \mu_0\|_q - \mathbb{E} \left\| \Gamma_{\eta;(\Sigma, \|\mu_0\|)}^{1/2} g/\sqrt{n} \right\|_q \right| \\ \lesssim n^{\frac{1}{2}-\frac{1}{q}} \cdot \left(n^{-\frac{1}{2}+\frac{1}{q}-\vartheta} (\log n)^{\frac{1}{2}-\frac{1}{q\sqrt{2}}} + n^{-1+\frac{1}{q\lambda^2}+\vartheta} \right) \end{aligned}$$

$$= n^{-\vartheta} (\log n)^{\frac{1}{2} - \frac{1}{q\sqrt{2}}} + n^{-\frac{1}{2} - \frac{1}{q} + \frac{1}{q\sqrt{2}} + \vartheta} \lesssim n^{-\vartheta/2}.$$

The claim follows from Lemma 60. \blacksquare

Proof [Proof of Theorem 6] We write $\widehat{\mu}_{(\Sigma, \mu_0)}^{\text{seq}}(\gamma_{\eta, *}; \tau_{\eta, *}) = \widehat{\mu}_{\eta; (\Sigma, \mu_0)}^{\text{seq}, *}$ in the proof.

First we consider $1 \leq q \leq 2$. This is the easy case, as $\mathbf{g}_q(x) \equiv \|x - \mu_0\|_q / n^{1/q-1/2}$ is 1-Lipschitz with respect to $\|\cdot\|$. So applying Theorems 3 and 4 verifies the existence of some small $\vartheta > 0$ such that for some $\mathcal{U}_\vartheta \subset B_n(1)$ with $\text{vol}(\mathcal{U}_\vartheta) / \text{vol}(B_n(1)) \geq 1 - Ce^{-n^\vartheta/C}$,

$$\sup_{\mu_0 \in \mathcal{U}_\vartheta} \mathbb{P} \left(\sup_{\eta \in \Xi_K} n^{\frac{1}{2} - \frac{1}{q}} \|\widehat{\mu}_\eta - \mu_0\|_q - \mathbb{E} \|\widehat{\mu}_{\eta; (\Sigma, \mu_0)}^{\text{seq}, *}\|_q \geq n^{-\vartheta} \right) \leq Cn^{-1/7}.$$

The ratio formulation follows from Lemmas 47 and 48 by further intersecting \mathcal{U}_ϑ and the set therein.

Next we consider $q \in (2, \infty)$. Let $L_n \equiv n^{\vartheta_1}$ for some ϑ_1 to be chosen later. Using Proposition 40 and its proofs below (123), for $\vartheta_1 > 0$ chosen small enough, we may find some $\mathcal{U}_{\vartheta_1} \subset B_n(1)$ with the desired volume estimate, such that $\sup_{\mu_0 \in \mathcal{U}_{\vartheta_1}} \sup_{\eta \in \Xi_K} \|\mathbb{E} \widehat{\mu}_{\eta; (\Sigma, \mu_0)}^{\text{seq}, *}\|_\infty - \mu_0\|_\infty \leq L_n / \sqrt{n}$, and

$$\sup_{\mu_0 \in \mathcal{U}_{\vartheta_1}} \mathbb{P} \left(\sup_{\eta \in \Xi_K} \left\{ \|\widehat{\mu}_\eta - \mu_0\|_\infty \vee \|\widehat{\mu}_{\eta; (\Sigma, \mu_0)}^{\text{seq}, *}\|_\infty \right\} \geq \frac{L_n}{\sqrt{n}} \right) \leq Cn^{-2D}, \quad (140)$$

where we choose $D > 0$ sufficiently large. Recall for $x \in \mathbb{R}^n$ and $q > 2$, $\|x\|_q \leq \|x\|^{2/q} \|x\|_\infty^{1-2/q}$. This motivates the choice

$$\mathbf{g}_q(x) \equiv \left[\left(\frac{L_n}{\sqrt{n}} \right)^{\frac{2}{q}-1} \left\| (x - \mu_0) \wedge \left(\frac{L_n}{\sqrt{n}} \right)^{\frac{2}{q}-1} \vee \left\{ - \left(\frac{L_n}{\sqrt{n}} \right)^{\frac{2}{q}-1} \right\} \right\|_q \right]^{\frac{q}{2}},$$

which verifies that \mathbf{g}_q is 1-Lipschitz with respect to $\|\cdot\|$. Using (140),

$$\inf_{\mu_0 \in \mathcal{U}_{\vartheta_1}} \mathbb{P} \left(\mathbf{g}_q(\widehat{\mu}_\eta) = n^{(1-\frac{q}{2})\vartheta_1} \cdot \left\{ n^{\frac{1}{2} - \frac{1}{q}} \|\widehat{\mu}_\eta - \mu_0\|_q \right\}^{\frac{q}{2}}, \forall \eta \in \Xi_K \right) \geq 1 - Cn^{-D}, \quad (141)$$

and with $E_{\mu_0} \equiv \left\{ \sup_{\eta \in \Xi_K} \|\widehat{\mu}_{\eta; (\Sigma, \mu_0)}^{\text{seq}, *}\|_\infty \leq L_n / \sqrt{n} \right\}$,

$$\begin{aligned} & \sup_{\mu_0 \in \mathcal{U}_{\vartheta_1}} \sup_{\eta \in \Xi_K} \left| \mathbb{E} \mathbf{g}_q(\widehat{\mu}_{\eta; (\Sigma, \mu_0)}^{\text{seq}, *}) - n^{(1-\frac{q}{2})\vartheta_1} \mathbb{E} \left\{ n^{\frac{1}{2} - \frac{1}{q}} \|\widehat{\mu}_{\eta; (\Sigma, \mu_0)}^{\text{seq}, *}\|_q \right\}^{\frac{q}{2}} \right| \\ &= n^{(1-\frac{q}{2})\vartheta_1} \sup_{\mu_0 \in \mathcal{U}_{\vartheta_1}} \sup_{\eta \in \Xi_K} \mathbb{E} \left\{ n^{\frac{1}{2} - \frac{1}{q}} \|\widehat{\mu}_{\eta; (\Sigma, \mu_0)}^{\text{seq}, *}\|_q \right\}^{\frac{q}{2}} \mathbf{1}_{E_{\mu_0}^c} \\ & \quad + \sup_{\mu_0 \in \mathcal{U}_{\vartheta_1}} \sup_{\eta \in \Xi_K} \mathbb{E} \mathbf{g}_q(\widehat{\mu}_{\eta; (\Sigma, \mu_0)}^{\text{seq}, *}) \mathbf{1}_{E_{\mu_0}^c} \lesssim n^{-D}. \end{aligned} \quad (142)$$

As the map $g \mapsto \|\widehat{\mu}_{\eta; (\Sigma, \mu_0)}^{\text{seq}, *}\|_q = \|(\Sigma + \tau_{\eta, *} I)^{-1} (-\tau_{\eta, *} \mu_0 + \gamma_{\eta, *} \Sigma^{1/2} g / \sqrt{n})\|_q$ is $Cn^{-1/2}$ -Lipschitz with respect to $\|\cdot\|$, Gaussian concentration yields

$$\mathbb{P} \left(n^{1/2} \|\widehat{\mu}_{\eta; (\Sigma, \mu_0)}^{\text{seq}, *}\|_q - \mathbb{E} \|\widehat{\mu}_{\eta; (\Sigma, \mu_0)}^{\text{seq}, *}\|_q \geq n^{\vartheta_1} \right) \leq Cn^{-2D}.$$

Using the Lipschitz property of the maps, we may strengthen the above inequality to a uniform control over $\eta \in \Xi_K$. This means uniformly in $\mu_0 \in \mathcal{U}_{\vartheta_1}, \eta \in \Xi_K$,

$$\left| \mathbb{E} \left\{ n^{\frac{1}{2}-\frac{1}{q}} \|\widehat{\mu}_{\eta;(\Sigma, \mu_0)}^{\text{seq},*} - \mu_0\|_q \right\}^{\frac{q}{2}} - \left\{ n^{\frac{1}{2}-\frac{1}{q}} \mathbb{E} \|\widehat{\mu}_{\eta;(\Sigma, \mu_0)}^{\text{seq},*} - \mu_0\|_q \right\}^{\frac{q}{2}} \right| \lesssim n^{-\frac{1}{q}+\vartheta_1}. \quad (143)$$

Combining (142)-(143), we have uniformly in $\mu_0 \in \mathcal{U}_{\vartheta_1}, \eta \in \Xi_K$,

$$\left| \mathbb{E} \mathfrak{g}_q(\widehat{\mu}_{\eta;(\Sigma, \mu_0)}^{\text{seq},*}) - n^{(1-\frac{q}{2})\vartheta_1} \left\{ n^{\frac{1}{2}-\frac{1}{q}} \mathbb{E} \|\widehat{\mu}_{\eta;(\Sigma, \mu_0)}^{\text{seq},*} - \mu_0\|_q \right\}^{q/2} \right| \leq C n^{(2-\frac{q}{2})\vartheta_1-\frac{1}{q}}. \quad (144)$$

Combining (141) and (144) proves the existence of some small ϑ_2 and some $\mathcal{U}_{\vartheta_2} \subset B_n(1)$ with the desired volume estimate, such that

$$\sup_{\mu_0 \in \mathcal{U}_{\vartheta_2}} \mathbb{P} \left(n^{\frac{1}{2}-\frac{1}{q}} \sup_{\eta \in \Xi_K} \left| \|\widehat{\mu}_\eta - \mu_0\|_q - \mathbb{E} \|\widehat{\mu}_{\eta;(\Sigma, \mu_0)}^{\text{seq},*} - \mu_0\|_q \right| \geq n^{-\vartheta_2} \right) \leq C n^{-1/7}.$$

The ratio formulation follows again from Lemmas 47 and 48. ■

Appendix F. Proofs for Section 3.2

F.1 A rigorous version of (17) and its proof

The following theorem follows easily from Theorems 3 and 4.

Theorem 49. *Suppose Assumption A holds and the following hold for some $K > 0$.*

- $1/K \leq \phi^{-1} \leq K$, $\|\Sigma^{-1}\|_{\text{op}} \vee \|\Sigma\|_{\text{op}} \leq K$.
- Assumption B holds with $\sigma_\xi^2 \in [1/K, K]$.

Fix a small enough $\vartheta \in (0, 1/50)$. Then there exist a constant $C = C(K, \vartheta) > 1$, and a measurable set $\mathcal{U}_\vartheta \subset B_n(1)$ with $\text{vol}(\mathcal{U}_\vartheta)/\text{vol}(B_n(1)) \geq 1 - C e^{-n^\vartheta/C}$, such that for any $\varepsilon \in (0, 1/2]$, and $\# \in \{\text{pred}, \text{est}, \text{in}, \text{res}\}$,

$$\sup_{\mu_0 \in \mathcal{U}_\vartheta} \mathbb{P} \left(\sup_{\eta \in \Xi^\#} |R_{(\Sigma, \mu_0)}^\#(\eta) - \bar{R}_{(\Sigma, \mu_0)}^\#(\eta)| \geq \varepsilon \right) \leq C \cdot \begin{cases} n e^{-n\varepsilon^4/C}, & Z = G; \\ \varepsilon^{-c_0} n^{-1/6.5}, & \text{otherwise.} \end{cases}$$

Here $\Xi^\# = \Xi_K$ for $\# \in \{\text{pred}, \text{est}\}$ and $\Xi^\# = [1/K, K]$ for $\# \in \{\text{in}, \text{res}\}$, and $c_0 > 0$ is universal. Moreover, when $Z = G$, the supremum in the above display extends to $\mu_0 \in B_n(1)$, and the constant $C > 0$ does not depend on ϑ .

Remark 50. 1. For $\# \in \{\text{in}, \text{res}\}$, we may take $\Xi^\# = \Xi_K$ at the cost of an worsened probability estimate $C(n e^{-n\varepsilon^{c_0}/C} + \varepsilon^{-c_0} n^{-1/6.5} \mathbf{1}_{Z \neq G})$, cf. Lemma 55.

2. The closest non-asymptotic results on exact risk characterizations related to our Theorem 49, appear to be those presented in (i) (Hastie et al., 2022, Theorems 2 and 5), which proved non-asymptotic additive approximations $R_{(\Sigma, \mu_0)}^{\text{pred}}(\eta) = \bar{R}_{(\Sigma, \mu_0)}^{\text{pred}}(\eta) + \mathbf{o}_{\mathbf{P}}(1)$,

and (ii) (Cheng and Montanari, 2024, Theorems 1 and 2), which provided substantially refined, multiplicative approximations $R_{(\Sigma, \mu_0)}^{\text{pred}}(\eta)/\bar{R}_{(\Sigma, \mu_0)}^{\text{pred}}(\eta) = 1 + \mathfrak{o}_{\mathbf{P}}(1)$ that hold beyond the proportional regime. Both works Hastie et al. (2022); Cheng and Montanari (2024) leverage the closed form of the Ridge(less) estimator $\hat{\mu}_\eta$ to analyze the bias and variance terms in $R_{(\Sigma, \mu_0)}^{\text{pred}}(\eta)$, by means of calculus for the resolvent of the sample covariance. Their analysis works under $\eta \gg n^{-c_0}$ for some suitable $c_0 > 0$. For the case $\# = \text{pred}$, Theorem 49 above complements the results in Hastie et al. (2022); Cheng and Montanari (2024) by providing uniform control in η when $\phi^{-1} > 1$ (under a set of different conditions).

Proof [Proof of Theorem 49] For ϑ chosen small enough, we fix $\mu_0 \in \mathcal{U}_\vartheta$, where \mathcal{U}_ϑ is specified in Theorem 4. We omit the subscripts in $R_{(\Sigma, \mu_0)}^\#(\eta) = R^\#(\eta)$, $\bar{R}_{(\Sigma, \mu_0)}^\#(\eta) = \bar{R}^\#(\eta)$, and write $\hat{\mu}_{(\Sigma, \mu_0)}^{\text{seq}}(\gamma_{\eta, *}; \tau_{\eta, *}) = \hat{\mu}_{\eta; (\Sigma, \mu_0)}^{\text{seq}, *}$ in the proof. All the constants in $\lesssim, \gtrsim, \asymp$ and \mathcal{O} below may possibly depend on K .

(1). Consider the case $\# = \text{pred}$. We omit the superscript pred as well. Using Theorem 4-(1) with $\mathbf{g}(x) = \|\Sigma^{1/2}(x - \mu_0)\|$, on an event E_0 with $\mathbb{P}(E_0^c) \leq C_0 \varepsilon^{-c_0} n^{-1/6.5}$,

$$\sup_{\eta \in \Xi_K} \left| \sqrt{R(\eta)} - \mathbb{E} \|\Sigma^{1/2}(\hat{\mu}_{\eta; (\Sigma, \mu_0)}^{\text{seq}, *} - \mu_0)\| \right| \leq \varepsilon.$$

By Gaussian-Poincaré inequality, $0 \leq \bar{R}(\eta) - (\mathbb{E} \|\Sigma^{1/2}(\hat{\mu}_{\eta; (\Sigma, \mu_0)}^{\text{seq}, *} - \mu_0)\|)^2 = \text{Var}(\|\Sigma^{1/2}(\hat{\mu}_{\eta; (\Sigma, \mu_0)}^{\text{seq}, *} - \mu_0)\|) \lesssim n^{-1}$. As $\bar{R}(\eta) \asymp 1$ uniformly in $\eta \in \Xi_K$, on E_0 ,

$$\sup_{\eta \in \Xi_K} |R^{1/2}(\eta) - \bar{R}^{1/2}(\eta)| \leq \varepsilon + C'_0 n^{-1}. \quad (145)$$

On the other hand, using both the standard form $\hat{\mu}_\eta = n^{-1}(X^\top X/n + \eta I_n)^{-1} X^\top Y$ and the alternative form $\hat{\mu}_\eta = n^{-1} X^\top (X X^\top/n + \eta I_m)^{-1} Y$, we have

$$\sup_{\eta \in \Xi_K} \|\hat{\mu}_\eta\| \lesssim \left(\|(ZZ^\top/n)^{-1}\|_{\text{op}} \mathbf{1}_{\phi^{-1} \geq 1+1/K}^{-1} \wedge 1 \right) \cdot \left(1 + \frac{\|Z\|_{\text{op}} + \|\xi\|}{\sqrt{n}} \right)^2. \quad (146)$$

Consequently, on an event E_1 with $\mathbb{P}(E_1^c) \leq C_1 e^{-n/C_1}$,

$$\sup_{\eta \in \Xi_K} \|\hat{\mu}_\eta\| \leq C_1. \quad (147)$$

Finally, using (145) and (147), on $E_0 \cap E_1$,

$$\sup_{\eta \in \Xi_K} |R(\eta) - \bar{R}(\eta)| \lesssim \sup_{\eta \in \Xi_K} |R^{1/2}(\eta) - \bar{R}^{1/2}(\eta)| \left(1 + \sup_{\eta \in \Xi_K} \|\hat{\mu}_\eta\| \right) \lesssim \varepsilon + n^{-1}.$$

The claim follows. The case $\# = \text{est}$ follows from minor modifications so will be omitted.

(2). Consider the case $\# = \text{res}$. We omit the superscript res as well. Further fix $\xi \in \mathcal{E}_\vartheta$ as specified in Theorem 4 (the concrete form of \mathcal{E}_ϑ is given in Proposition 40). Using the same Theorem 4-(2) with $\mathbf{h}(x) = \|x\|$,

$$\mathbb{P}^\xi \left(\sup_{\eta \in [1/K, K]} \left| \|\hat{r}_\eta\| - \mathbb{E}^\xi \|r_{\eta, *}\| \right| \geq \varepsilon \right) \leq C \varepsilon^{-c_0} \cdot n^{-1/6.5}.$$

By Gaussian-Poincaré inequality, $0 \leq \mathbb{E}^\xi \|r_{\eta,*}\|^2 - (\mathbb{E}^\xi \|r_{\eta,*}\|)^2 = \text{Var}^\xi(\|r_{\eta,*}\|) \lesssim 1/n$. Combined with the fact that $\mathbb{E}^\xi \|r_{\eta,*}\|^2 = (\eta\gamma_{\eta,*}/\tau_{\eta,*})^2 + \mathcal{O}(\|\xi\|^2/m - \sigma_\xi^2)$, for $\eta \in [1/K, K]$, using the stability estimate in Proposition 23-(3),

$$|\mathbb{E}^\xi \|r_{\eta,*}\| - \eta\gamma_{\eta,*}/\tau_{\eta,*}| \lesssim |(\mathbb{E}^\xi \|r_{\eta,*}\|)^2 - (\eta\gamma_{\eta,*}/\tau_{\eta,*})^2| \lesssim n^{-1/2+\vartheta}.$$

So for $\varepsilon \in (Cn^{-1/2+\vartheta}, 1/C]$,

$$\mathbb{P}^\xi \left(\sup_{\eta \in [1/K, K]} \left| \|\widehat{r}_\eta\| - \eta\gamma_{\eta,*}/\tau_{\eta,*} \right| \geq \varepsilon \right) \leq C\varepsilon^{-c_0} \cdot n^{-1/6.5}.$$

Now taking expectation over ξ , for the same range of ε ,

$$\mathbb{P} \left(\sup_{\eta \in [1/K, K]} \left| \|\widehat{r}_\eta\| - \eta\gamma_{\eta,*}/\tau_{\eta,*} \right| \geq \varepsilon \right) \leq C\varepsilon^{-c_0} \cdot n^{-1/6.5}. \quad (148)$$

On the other hand, using (146),

$$\sup_{\eta \in [1/K, K]} \|\widehat{r}_\eta\| \lesssim \left(\|(ZZ^\top/n)^{-1}\|_{\text{op}} \mathbf{1}_{\phi^{-1} \geq 1+1/K} \wedge 1 \right) \cdot \left(1 + \frac{\|Z\|_{\text{op}} + \|\xi\|}{\sqrt{n}} \right)^3.$$

Consequently, on an event E_3 with $\mathbb{P}(E_3^c) \leq C_3 e^{-n/C_3}$, $\sup_{\eta \in [1/K, K]} \|\widehat{r}_\eta\| \leq C_3$, and therefore

$$\sup_{\eta \in [1/K, K]} \left| \|\widehat{r}_\eta\|^2 - (\eta\gamma_{\eta,*}/\tau_{\eta,*})^2 \right| \leq C_3 \cdot \sup_{\eta \in [1/K, K]} \left| \|\widehat{r}_\eta\| - \eta\gamma_{\eta,*}/\tau_{\eta,*} \right|.$$

The claim follows. The case $\# = \text{in}$ proceeds similarly, but with the function now taken as $\mathbf{h}(x) = \|x - \xi/\sqrt{n}\|$, and the claim follows by computing that

$$\begin{aligned} \mathbb{E}^\xi \|r_{\eta,*} - \xi/\sqrt{n}\|^2 &= \phi \cdot \left\{ \left(\frac{\eta}{\phi\tau_{\eta,*}} \right)^2 (\phi\gamma_{\eta,*}^2 - \sigma_\xi^2) + \frac{\|\xi\|^2}{m} \cdot \left(\frac{\eta}{\phi\tau_{\eta,*}} - 1 \right)^2 \right\} \\ &= \left(\frac{\eta\gamma_{\eta,*}}{\tau_{\eta,*}} \right)^2 + \phi\sigma_\xi^2 \cdot \left[\left(\frac{\eta}{\phi\tau_{\eta,*}} - 1 \right)^2 - \left(\frac{\eta}{\phi\tau_{\eta,*}} \right)^2 \right] + \mathcal{O}(\|\xi\|^2/m - \sigma_\xi^2). \end{aligned}$$

The proof is complete. \blacksquare

F.2 Proof of Theorem 8

Lemma 51. *Suppose $1/K \leq \phi^{-1} \leq K$, and $\|\Sigma\|_{\text{op}} \vee \mathcal{H}_\Sigma \leq K$ for some $K > 0$. Then with $g \sim \mathcal{N}(0, I_n)$, there exists some $C = C(K) > 0$ such that for $\varepsilon \in (0, 1)$, and $q \in \{0, 1/2\}$,*

$$\mathbb{P} \left(\sup_{\eta \in \Xi_K} \left| \|(\Sigma + \tau_{\eta,*}I)^{-1}\Sigma^q g/\|g\|\|^2 - n^{-1} \text{tr}((\Sigma + \tau_{\eta,*}I)^{-2}\Sigma^{2q}) \right| > \varepsilon \right) \leq C\varepsilon^{-1} e^{-n\varepsilon^2/C}.$$

Proof We only prove the case $q = 1/2$. All the constants in $\lesssim, \gtrsim, \asymp$ below may depend on K . We write $A_\eta \equiv (\Sigma + \tau_{\eta,*}I)^{-2}\Sigma$ for notational simplicity. Note that

$$\begin{aligned} & \left| \|(\Sigma + \tau_{\eta,*}I)^{-1}\Sigma^{1/2}g/\|g\|\|^2 - n^{-1} \text{tr}((\Sigma + \tau_{\eta,*}I)^{-2}\Sigma) \right| \\ &= n^{-1} \left| e_g^{-2} \|A_\eta^{1/2}g\|^2 - \mathbb{E}\|A_\eta^{1/2}g\|^2 \right| \\ &\lesssim e_g^{-2} \cdot n^{-1} \left| \|A_\eta^{1/2}g\|^2 - \mathbb{E}\|A_\eta^{1/2}g\|^2 \right| + |e_g^{-2} - 1|. \end{aligned}$$

Here in the last inequality we used $\mathbb{E}\|A_\eta^{1/2}g\|^2 \lesssim n$. As

- $\|A_\eta\|_F^2 = \text{tr}((\Sigma + \tau_{\eta,*}I)^{-4}\Sigma^2) \lesssim n(1 \wedge \tau_{\eta,*})^{-4} \asymp n$, and
- $\|A_\eta\|_F^2 \gtrsim \text{tr}(\Sigma^2) \cdot (1 \vee \tau_{\eta,*})^{-4} \gtrsim n$,

we have uniformly in $\eta \in [0, K]$, $\|A_\eta\|_F \asymp \sqrt{n}$. It is easy to see that $\|A_\eta\|_{\text{op}} \asymp 1$. So by Hanson-Wright inequality, there exists some constant $C_1 = C_1(K)$ such that for $\varepsilon \in (0, 1)$,

$$\begin{aligned} & \mathbb{P}\left(\left|\|(\Sigma + \tau_{\eta,*}I)^{-1}\Sigma^{1/2}g/\|g\|\|^2 - n^{-1}\text{tr}((\Sigma + \tau_{\eta,*}I)^{-2}\Sigma)\right| > \varepsilon\right) \\ & \leq \mathbb{P}\left(|n^{-1}(\|A_\eta^{1/2}g\|^2 - \mathbb{E}\|A_\eta^{1/2}g\|^2)| > \varepsilon/4\right) + \mathbb{P}(|e_g^{-2} - 1| > \varepsilon/2) + \mathbb{P}(e_g^2 \leq 1/2) \\ & \leq C_1 e^{-n\varepsilon^2/C_1}. \end{aligned}$$

On the other hand, for any $\eta_1, \eta_2 \in \Xi_K$, using Proposition 23-(3),

$$\begin{aligned} & \left|\|(\Sigma + \tau_{\eta_1,*}I)^{-1}\Sigma^{1/2}g/\|g\|\|^2 - \|(\Sigma + \tau_{\eta_2,*}I)^{-1}\Sigma^{1/2}g/\|g\|\|^2\right| \lesssim |\eta_1 - \eta_2|, \\ & n^{-1}\left|\text{tr}((\Sigma + \tau_{\eta_1,*}I)^{-2}\Sigma) - \text{tr}((\Sigma + \tau_{\eta_2,*}I)^{-2}\Sigma)\right| \lesssim |\eta_1 - \eta_2|, \end{aligned}$$

so we may conclude by a standard discretization and union bound argument. \blacksquare

Proposition 52. *The following hold with $\mathbf{m}_\eta \equiv \mathbf{m}(-\eta/\phi)$, $\mathbf{m}'_\eta \equiv \mathbf{m}'(-\eta/\phi)$.*

1. $\tau_{\eta,*} = 1/\mathbf{m}_\eta$ and $\partial_\eta \tau_{\eta,*} = \mathbf{m}'_\eta/(\phi \mathbf{m}_\eta^2)$.
2. It holds that

$$\begin{aligned} \frac{1}{n} \text{tr}((\Sigma + \tau_{\eta,*}I)^{-2}\Sigma) &= \frac{\phi \mathbf{m}_\eta^2}{\mathbf{m}'_\eta} (\mathbf{m}_\eta - (\eta/\phi) \mathbf{m}'_\eta), \\ \frac{1}{n} \text{tr}((\Sigma + \tau_{\eta,*}I)^{-2}) &= \frac{\phi \mathbf{m}_\eta^2}{\mathbf{m}'_\eta} ((\phi^{-1} - 1) \mathbf{m}'_\eta + 2(\eta/\phi) \cdot \mathbf{m}_\eta \mathbf{m}'_\eta - \mathbf{m}_\eta^2). \end{aligned}$$

3. Suppose $1/K \leq \phi^{-1} \leq K$, and $\|\Sigma\|_{\text{op}} \vee \mathcal{H}_\Sigma \leq K$ for some $K > 0$. There exists some constant $C = C(K) > 0$ such that the following hold. For any $\varepsilon \in (0, 1/2]$, for some $\mathcal{U}_\varepsilon \subset B_n(1)$ with $\text{vol}(\mathcal{U}_\varepsilon)/\text{vol}(B_n(1)) \geq 1 - C\varepsilon^{-1}e^{-n\varepsilon^2/C}$,

$$\sup_{\mu_0 \in \mathcal{U}_\varepsilon} \sup_{\eta \in \Xi_K} \left| \gamma_{\eta,*}^2 - \frac{\sigma_\xi^2 \mathbf{m}'_\eta + \|\mu_0\|^2 (\phi \mathbf{m}_\eta - \eta \mathbf{m}'_\eta)}{\phi \mathbf{m}_\eta^2} \right| \leq \varepsilon.$$

When $\Sigma = I_n$, we may take $\mathcal{U}_\varepsilon = B_n(1)$ and the above inequality holds with $\varepsilon = 0$.

Proof (1) follows from definition so we focus on (2)-(3).

(2). Differentiating both sides of (42) with respect to η yields that

$$-n^{-1} \text{tr}((\Sigma + \tau_{\eta,*}I)^{-2}\Sigma) \cdot \partial_\eta \tau_{\eta,*} = -(\mathbf{m}_\eta - (\eta/\phi) \mathbf{m}'_\eta).$$

Now using $\partial_\eta \tau_{\eta,*} = \mathbf{m}'_\eta/(\phi \mathbf{m}_\eta^2)$ to obtain the formula for $n^{-1} \text{tr}((\Sigma + \tau_{\eta,*}I)^{-2}\Sigma)$.

Next, using that $\phi - \frac{\eta}{\tau_{\eta,*}} = n^{-1} \operatorname{tr}((\Sigma + \tau_{\eta,*}I)^{-1}\Sigma) = 1 - \tau_{\eta,*} \cdot n^{-1} \operatorname{tr}((\Sigma + \tau_{\eta,*}I)^{-1})$, we may solve

$$n^{-1} \operatorname{tr}((\Sigma + \tau_{\eta,*}I)^{-1}) = \mathbf{m}_\eta(1 - \phi + \eta \cdot \mathbf{m}_\eta).$$

Differentiating with respect to η on both sides of the above display, we obtain

$$\begin{aligned} -n^{-1} \operatorname{tr}((\Sigma + \tau_{\eta,*}I)^{-2}) \cdot \partial_\eta \tau_{\eta,*} &= -\phi^{-1} \mathbf{m}'_\eta (1 - \phi + \eta \cdot \mathbf{m}_\eta) + \mathbf{m}_\eta \cdot (\mathbf{m}_\eta - (\eta/\phi) \mathbf{m}'_\eta) \\ &= -(\phi^{-1} - 1) \mathbf{m}'_\eta - 2(\eta/\phi) \cdot \mathbf{m}_\eta \mathbf{m}'_\eta + \mathbf{m}_\eta^2, \end{aligned}$$

proving the second identity.

(3). Let $\mu_0 \equiv U_0 g_0 / \|g_0\|$, where $U_0 \sim \operatorname{Unif}[0, 1]$ and $g_0 \sim \mathcal{N}(0, I_n)$ are independent variables. Then μ_0 is uniformly distributed on $B_n(1)$. For some $\varepsilon > 0$ to be chosen later, let

$$\mathcal{G}_\varepsilon \equiv \left\{ g \in \mathbb{R}^n : \sup_{\eta \in \Xi_K} \left| \left\| (\Sigma + \tau_{\eta,*}I)^{-1} \Sigma^{1/2} \frac{g}{\|g\|} \right\|^2 - \frac{1}{n} \operatorname{tr}((\Sigma + \tau_{\eta,*}I)^{-2} \Sigma) \right| \leq \varepsilon \right\}. \quad (149)$$

Let $\mathcal{U}_\varepsilon \equiv \{Ug/\|g\| : U \in [0, 1], g \in \mathcal{G}_\varepsilon\} \subset B_n(1)$. Using Lemma 51, there exists some constant $C_0 = C_0(K) > 0$ such that $\operatorname{vol}(\mathcal{U}_\varepsilon)/\operatorname{vol}(B_n(1)) = \mathbb{P}_{\mu_0}(\mu_0 \in \mathcal{U}_\varepsilon) \geq 1 - C_0 \varepsilon^{-1} e^{-n\varepsilon^2/C_0}$, and moreover,

$$\sup_{\mu_0 \in \mathcal{U}_\varepsilon} \sup_{\eta \in \Xi_K} \left| \left\| (\Sigma + \tau_{\eta,*}I)^{-1} \Sigma^{1/2} \mu_0 \right\|^2 - \|\mu_0\|^2 \cdot n^{-1} \operatorname{tr}((\Sigma + \tau_{\eta,*}I)^{-2} \Sigma) \right| \leq \varepsilon.$$

Note that when $\Sigma = I_n$, the above estimate holds for all $\mu_0 \in B_n(1)$ with $\varepsilon = 0$.

Combining the above display with the formula (46) for $\gamma_{\eta,*}^2$, and the fact that the denominator therein is of order 1 (depending on K), we have

$$\sup_{\mu_0 \in \mathcal{U}_\varepsilon} \sup_{\eta \in \Xi_K} \left| \gamma_{\eta,*}^2 - \frac{\sigma_\xi^2 + \|\mu_0\|^2 \tau_{\eta,*}^2 \cdot \frac{1}{n} \operatorname{tr}((\Sigma + \tau_{\eta,*}I)^{-2} \Sigma)}{\frac{\eta}{\tau_{\eta,*}} + \tau_{\eta,*} \cdot \frac{1}{n} \operatorname{tr}((\Sigma + \tau_{\eta,*}I)^{-2} \Sigma)} \right| \leq C_1 \varepsilon.$$

Now using (2), the second term in the above display equals to

$$\frac{\sigma_\xi^2 + \|\mu_0\|^2 \cdot \frac{\phi}{\mathbf{m}'_\eta} (\mathbf{m}_\eta - \frac{\eta}{\phi} \mathbf{m}'_\eta)}{\eta \mathbf{m}_\eta + \frac{\phi \mathbf{m}_\eta}{\mathbf{m}'_\eta} (\mathbf{m}_\eta - \frac{\eta}{\phi} \mathbf{m}'_\eta)} = \frac{\sigma_\xi^2 \mathbf{m}'_\eta + \|\mu_0\|^2 (\phi \mathbf{m}_\eta - \eta \mathbf{m}'_\eta)}{\phi \mathbf{m}_\eta^2}.$$

The claim follows by adjusting constants. ■

Proof [Proof of Theorem 8] As $\bar{R}_{(\Sigma, \mu_0)}^{\operatorname{pred}}(\eta) = \phi \gamma_{\eta,*}^2 - \sigma_\xi^2$, directly invoking Proposition 52-(3) yields the claim for $\bar{R}_{(\Sigma, \mu_0)}^{\operatorname{pred}}(\eta)$.

Next we handle $\bar{R}_{(\Sigma, \mu_0)}^{\operatorname{est}}(\eta)$. Note that

$$\bar{R}_{(\Sigma, \mu_0)}^{\operatorname{est}}(\eta) = \tau_{\eta,*}^2 \left\| (\Sigma + \tau_{\eta,*}I)^{-1} \mu_0 \right\|^2 + \gamma_{\eta,*}^2 \cdot n^{-1} \operatorname{tr}((\Sigma + \tau_{\eta,*}I)^{-2} \Sigma). \quad (150)$$

Using a similar construction as in the proof of Proposition 52 via the help of Lemma 51, this time with $q = 0$ therein, we may find some $\mathcal{U}_\varepsilon \subset B_n(1)$ with the desired volume estimate, such that both Proposition 52-(3) and

$$\sup_{\mu_0 \in \mathcal{U}_\varepsilon} \sup_{\eta \in \Xi_K} \left| \left\| (\Sigma + \tau_{\eta,*}I)^{-1} \mu_0 \right\|^2 - \|\mu_0\|^2 \cdot n^{-1} \operatorname{tr}((\Sigma + \tau_{\eta,*}I)^{-2}) \right| \leq \varepsilon \quad (151)$$

hold. Combining (150)-(151), we may set

$$\begin{aligned}\mathcal{R}_{(\Sigma, \mu_0)}^{\text{est}}(\eta) &\equiv \tau_{\eta, *}^2 \|\mu_0\|^2 \cdot n^{-1} \text{tr} \left((\Sigma + \tau_{\eta, *} I)^{-2} \right) \\ &\quad + (\phi \mathbf{m}_\eta^2)^{-1} \left(\sigma_\xi^2 \mathbf{m}'_\eta + \|\mu_0\|^2 (\phi \mathbf{m}_\eta - \eta \mathbf{m}'_\eta) \right) \cdot n^{-1} \text{tr} \left((\Sigma + \tau_{\eta, *} I)^{-2} \Sigma \right) \\ &\equiv R_{2,1} + R_{2,2}.\end{aligned}$$

By Proposition 52-(2), we may compute $R_{2,1}, R_{2,2}$ separately:

$$\begin{aligned}R_{2,1} &= \|\mu_0\|^2 \cdot \frac{\phi}{\mathbf{m}'_\eta} \cdot \left((\phi^{-1} - 1) \mathbf{m}'_\eta + 2(\eta/\phi) \cdot \mathbf{m}_\eta \mathbf{m}'_\eta - \mathbf{m}_\eta^2 \right) \\ &= \|\mu_0\|^2 (1 - \phi) + \left\{ 2\|\mu_0\|^2 \eta \mathbf{m}_\eta - \|\mu_0\|^2 \phi \cdot \frac{\mathbf{m}_\eta^2}{\mathbf{m}'_\eta} \right\}, \\ R_{2,2} &= \frac{1}{\mathbf{m}'_\eta} \left(\sigma_\xi^2 \mathbf{m}'_\eta + \|\mu_0\|^2 (\phi \mathbf{m}_\eta - \eta \mathbf{m}'_\eta) \right) \cdot (\mathbf{m}_\eta - (\eta/\phi) \mathbf{m}'_\eta) \\ &= \sigma_\xi^2 (\mathbf{m}_\eta - (\eta/\phi) \mathbf{m}'_\eta) + \phi^{-1} \|\mu_0\|^2 \eta^2 \mathbf{m}'_\eta - \left\{ 2\|\mu_0\|^2 \eta \mathbf{m}_\eta - \|\mu_0\|^2 \phi \cdot \frac{\mathbf{m}_\eta^2}{\mathbf{m}'_\eta} \right\}.\end{aligned}$$

Consequently,

$$\begin{aligned}\mathcal{R}_{(\Sigma, \mu_0)}^{\text{est}}(\eta) &= \|\mu_0\|^2 (1 - \phi) + \sigma_\xi^2 (\mathbf{m}_\eta - (\eta/\phi) \mathbf{m}'_\eta) + \phi^{-1} \|\mu_0\|^2 \eta^2 \mathbf{m}'_\eta \\ &= \sigma_\xi^2 \cdot \left\{ \text{SNR}_{\mu_0} (1 - \phi) + \mathbf{m}_\eta + (\eta/\phi) (\eta \cdot \text{SNR}_{\mu_0} - 1) \mathbf{m}'_\eta \right\}.\end{aligned}$$

The claims for $\mathcal{R}_{(\Sigma, \mu_0)}^{\text{in}}(\eta)$ and $\mathcal{R}_{(\Sigma, \mu_0)}^{\text{res}}(\eta)$ follow from Proposition 52-(3). \blacksquare

F.3 Proof of Proposition 9

We will prove the following version of Proposition 9, where $\mathfrak{M}^\#$ is represented via $\tau_{\eta, *}$ instead of \mathbf{m} . In the proof below, we will also verify the representation of $\mathfrak{M}^\#$ via \mathbf{m} as stated in Proposition 9.

Proposition 53. *Recall $\text{SNR}_{\mu_0} = \|\mu_0\|^2 / \sigma_\xi^2$. Then for $\# \in \{\text{pred}, \text{est}, \text{in}\}$,*

$$\partial_\eta \mathcal{R}_{(\Sigma, \mu_0)}^\#(\eta) = \sigma_\xi^2 \cdot \mathfrak{M}^\#(\eta) \cdot (\eta \cdot \text{SNR}_{\mu_0} - 1).$$

Here with $T_{-p, q}(\eta) \equiv n^{-1} \text{tr} \left((\Sigma + \tau_*(\eta) I)^{-p} \Sigma^q \right)$ for $p, q \in \mathbb{N}$,

$$\mathfrak{M}^\#(\eta) \equiv \begin{cases} \phi(-\tau_*''(\eta)), & \# = \text{pred}; \\ 2(\tau_*'(\eta))^2 (T_{-3,1}(\eta) + \tau_*'(\eta) T_{-2,1}(\eta) T_{-3,2}(\eta)), & \# = \text{est}; \\ \frac{2(\tau_*'(\eta))^2}{\tau_*^2(\eta)} \left(\eta^2 \tau_*'(\eta) T_{-3,2}(\eta) + \tau_*^3(\eta) T_{-2,1}^2(\eta) \right), & \# = \text{in}. \end{cases}$$

Suppose further $1/K \leq \phi^{-1} \leq K$ and $\|\Sigma\|_{\text{op}} \vee \mathcal{H}_\Sigma \leq K$ for some $K > 0$. Then there exists some $C = C(K) > 0$ such that uniformly in $\eta \in \Xi_K$ and for all $\# \in \{\text{pred}, \text{est}, \text{in}\}$,

1. $1/C \leq \mathfrak{M}^\#(\eta) \leq C$, and
2. if $\eta_* \equiv \text{SNR}_{\mu_0}^{-1} \in \Xi_K$, then $1/C \leq |\mathcal{R}_{(\Sigma, \mu_0)}^\#(\eta) - \mathcal{R}_{(\Sigma, \mu_0)}^\#(\eta_*)| / \|\mu_0\|^2 (\eta - \eta_*)^2 \leq C$.

Proof In the proof we write $\tau_{\eta,*} = \tau_\eta$. Recall the notation $\mathbf{m}_\eta = \mathbf{m}(-\eta/\phi)$, $\mathbf{m}'_\eta = \mathbf{m}'(-\eta/\phi)$, and we naturally write $\mathbf{m}''_\eta \equiv \mathbf{m}''(-\eta/\phi)$. By differentiating with respect to η for both sides of $\mathbf{m}_\eta = 1/\tau_\eta$, with some calculations we have

$$\mathbf{m}_\eta = \tau_\eta^{-1}, \quad \mathbf{m}'_\eta = \phi \tau'_\eta / \tau_\eta^2, \quad \mathbf{m}''_\eta = -\phi^2 (\tau''_\eta \tau_\eta - 2(\tau'_\eta)^2) / \tau_\eta^3. \quad (152)$$

Using ρ , we may also write $\mathbf{m}_\eta^{(q)} = q! \int \frac{\rho(dx)}{(x + \eta/\phi)^{q+1}}$ for $q \in \mathbb{N}$. Here by convention $0! = 1$.

(1). Using the formula for $\mathcal{R}_{(\Sigma, \mu_0)}^{\text{pred}}$,

$$\begin{aligned} \partial_\eta \mathcal{R}_{(\Sigma, \mu_0)}^{\text{pred}}(\eta) &= \sigma_\xi^2 \cdot \partial_\eta \left\{ \mathbf{m}_\eta^{-2} (\phi \cdot \text{SNR}_{\mu_0} \mathbf{m}_\eta - (\eta \cdot \text{SNR}_{\mu_0} - 1) \mathbf{m}'_\eta) \right\} \\ &= \phi^{-1} \sigma_\xi^2 \cdot \mathbf{m}_\eta^{-3} (\mathbf{m}_\eta \mathbf{m}''_\eta - 2(\mathbf{m}'_\eta)^2) \cdot (\eta \cdot \text{SNR}_{\mu_0} - 1). \end{aligned}$$

Some calculations show that

$$\mathbf{m}_\eta \mathbf{m}''_\eta - 2(\mathbf{m}'_\eta)^2 = \tau_\eta^{-3} \phi^2 (-\tau''_\eta) = 2 \left\{ \int \frac{\rho(dx)}{(x + \eta/\phi)} \int \frac{\rho(dx)}{(x + \eta/\phi)^3} - \left(\int \frac{\rho(dx)}{(x + \eta/\phi)^2} \right)^2 \right\},$$

so the identity follows.

(2). Using the formula for $\mathcal{R}_{(\Sigma, \mu_0)}^{\text{est}}$,

$$\begin{aligned} \partial_\eta \mathcal{R}_{(\Sigma, \mu_0)}^{\text{est}}(\eta) &= \sigma_\xi^2 \cdot \partial_\eta (\text{SNR}_{\mu_0} (1 - \phi) + \mathbf{m}_\eta + (\eta/\phi)(\eta \cdot \text{SNR}_{\mu_0} - 1) \mathbf{m}'_\eta) \\ &= \phi^{-1} \sigma_\xi^2 \cdot (2\mathbf{m}'_\eta - (\eta/\phi) \mathbf{m}''_\eta) \cdot (\eta \cdot \text{SNR}_{\mu_0} - 1). \end{aligned} \quad (153)$$

To compute the second term in the above display, recall the identity for τ'_η, τ''_η in (48)-(49). Also recall $G_0(\eta) = \eta + \tau_\eta^2 T_{-2,1}(\eta) = \tau_\eta / \tau'_\eta$ defined in (48). Then

$$\begin{aligned} 2\mathbf{m}'_\eta - \frac{\eta}{\phi} \mathbf{m}''_\eta &= \frac{\phi}{\tau_\eta} \left\{ \frac{2\tau'_\eta}{\tau_\eta} \left(1 - \frac{\eta \tau'_\eta}{\tau_\eta} \right) + \frac{\eta \tau''_\eta}{\tau_\eta} \right\} = \frac{\phi}{\tau_\eta} \left\{ \frac{2\tau'_\eta}{\tau_\eta} \frac{\tau_\eta^2 T_{-2,1}(\eta)}{G_0(\eta)} - \frac{2\eta \tau_\eta \tau'_\eta}{G_0^2(\eta)} T_{-3,2}(\eta) \right\} \\ &= \frac{2\phi \tau'_\eta}{G_0(\eta)} \left(T_{-2,1}(\eta) - \frac{\eta}{G_0(\eta)} T_{-3,2}(\eta) \right) \stackrel{(*)}{=} \frac{2\phi \tau'_\eta}{G_0(\eta)} \left\{ \tau_\eta T_{-3,1}(\eta) + \left(1 - \frac{\eta}{G_0(\eta)} \right) T_{-3,2}(\eta) \right\} \\ &= 2\phi (\tau'_\eta)^2 (T_{-3,1}(\eta) + \tau'_\eta T_{-2,1}(\eta) T_{-3,2}(\eta)). \end{aligned}$$

Here in (*) we used $T_{-2,1}(\eta) - T_{-3,2}(\eta) = \tau_\eta T_{-3,1}(\eta)$. The claimed identity follows by combining the above display and (153). Using ρ , we may write

$$2\mathbf{m}'_\eta - (\eta/\phi) \mathbf{m}''_\eta = 2 \int \frac{x}{(x + \eta/\phi)^3} \rho(dx).$$

(3). Using the formula for $\mathcal{R}_{(\Sigma, \mu_0)}^{\text{in}}$,

$$\partial_\eta \mathcal{R}_{(\Sigma, \mu_0)}^{\text{in}}(\eta) = \sigma_\xi^2 \cdot \partial_\eta \left\{ \phi^{-1} \eta^2 (\phi \cdot \text{SNR}_{\mu_0} \mathbf{m}_\eta - (\eta \cdot \text{SNR}_{\mu_0} - 1) \mathbf{m}'_\eta) + (\phi - 2\eta \mathbf{m}_\eta) \right\}$$

$$= \sigma_\xi^2 \cdot (2\mathbf{m}_\eta - 4(\eta/\phi)\mathbf{m}'_\eta + \phi^{-2}\eta^2\mathbf{m}''_\eta) \cdot (\eta \cdot \text{SNR}_{\mu_0} - 1). \quad (154)$$

The second term in the above display requires some non-trivial calculations:

$$\begin{aligned} 2\mathbf{m}_\eta - \frac{4\eta}{\phi}\mathbf{m}'_\eta + \frac{\eta^2}{\phi^2}\mathbf{m}''_\eta &= \frac{1}{\tau_\eta} \left\{ 2 - 4\eta \frac{\tau'_\eta}{\tau_\eta} - \eta^2 \frac{\tau''_\eta}{\tau_\eta} + 2\eta^2 \left(\frac{\tau'_\eta}{\tau_\eta} \right)^2 \right\} \\ &= \frac{1}{\tau_\eta} \left\{ 2 - \frac{4\eta}{G_0(\eta)} + \eta^2 \frac{2\tau_\eta \tau'_\eta T_{-3,2}(\eta)}{G_0^2(\eta)} + \frac{2\eta^2}{G_0^2(\eta)} \right\} \\ &= \frac{2}{\tau_\eta G_0^2(\eta)} \{ G_0^2(\eta) - 2\eta G_0(\eta) + \eta^2 \tau_\eta \tau'_\eta T_{-3,2}(\eta) + \eta^2 \}. \end{aligned}$$

Expanding the $G_0(\eta)$ terms in the bracket using $G_0(\eta) = \eta + \tau_\eta^2 T_{-2,1}(\eta)$, with some calculations we arrive at

$$2\mathbf{m}_\eta - \frac{4\eta}{\phi}\mathbf{m}'_\eta + \frac{\eta^2}{\phi^2}\mathbf{m}''_\eta = \frac{2}{G_0^2(\eta)} \left(\eta^2 \tau'_\eta T_{-3,2}(\eta) + \tau_\eta^3 T_{-2,1}^2(\eta) \right).$$

The claimed identity follows by combining the above display and (154). Using ρ ,

$$2\mathbf{m}_\eta - \frac{4\eta}{\phi}\mathbf{m}'_\eta + \frac{\eta^2}{\phi^2}\mathbf{m}''_\eta = 2 \int \frac{x^2}{(x + \eta/\phi)^3} \rho(dx).$$

Finally, the claimed first two-sided bound on $\mathfrak{M}^\#$ follows from Proposition 23, and the second bound follows by using the fundamental theorem of calculus. \blacksquare

F.4 A rigorous version of (19) and its proof

The theorem below presents a rigorous formulation of (19).

Theorem 54. *Suppose Assumptions A-B hold, and $\|\Sigma^{-1}\|_{\text{op}} \vee \|\Sigma\|_{\text{op}} \leq K$ for some $K > 0$. Fix a small enough $\vartheta \in (0, 1/50)$. The following hold for all $\# \in \{\text{pred}, \text{est}, \text{in}\}$.*

1. (**Noisy case**). *Suppose $1/K \leq \phi^{-1} \leq K$ and $1/K \leq \sigma_\xi^2 \leq K$. Fix $\delta \in (0, 1/2]$ and $L \geq K/\delta^2$. There exist a constant $C = C(K, L, \delta, \vartheta) > 0$ and a measurable set $\mathcal{U}_{\delta, \vartheta} \subset B_n(1) \setminus B_n(\delta)$ with $\text{vol}(\mathcal{U}_{\delta, \vartheta})/\text{vol}(B_n(1) \setminus B_n(\delta)) \geq 1 - Ce^{-n^\vartheta/C}$, such that*

$$\sup_{\mu_0 \in \mathcal{U}_{\delta, \vartheta}} \mathbb{P} \left(\inf_{\eta' \in \Xi_L: |\eta' - \text{SNR}_{\mu_0}^{-1}| \geq \delta} |R_{(\Sigma, \mu_0)}^\#(\eta') - \min_{\eta \in \Xi_L} R_{(\Sigma, \mu_0)}^\#(\eta)| < \frac{1}{C} \right) \leq Cn^{-1/7}.$$

2. (**Noiseless case**). *Suppose $1 + 1/K \leq \phi^{-1} \leq K$ and $\sigma_\xi^2 = 0$. There exist a constant $C = C(K, \vartheta) > 0$ and a measurable set $\mathcal{U}_\vartheta \subset B_n(1)$ with $\text{vol}(\mathcal{U}_\vartheta)/\text{vol}(B_n(1)) \geq 1 - Ce^{-n^\vartheta/C}$, such that*

$$\sup_{\mu_0 \in \mathcal{U}_\vartheta} \mathbb{P} \left(R_{(\Sigma, \mu_0)}^\#(0) \geq \min_{\eta \in [0, K]} R_{(\Sigma, \mu_0)}^\#(\eta) + n^{-\vartheta} \right) \leq Cn^{-1/7}.$$

We need a few lemmas to prove Theorem 54.

The following lemma gives a technical extension of Theorem 49 for $\# \in \{\text{pred}, \text{est}, \text{in}\}$ under $\sigma_\xi^2 \approx 0$ when $\phi^{-1} > 1$. For $\# = \text{in}$, the extension also allows uniform control over $\eta \approx 0$ under both the above small variance scenario with $\phi^{-1} > 1$, and under the original conditions.

Lemma 55. *Suppose Assumption A holds and the following hold for some $K > 0$.*

- $1 + 1/K \leq \phi^{-1} \leq K$, $\|\Sigma^{-1}\|_{\text{op}} \vee \|\Sigma\|_{\text{op}} \leq K$.
- Assumption B with $\sigma_\xi^2 \in [0, K]$.

Fix a small enough $\vartheta \in (0, 1/50)$. Then there exist a constant $C = C(K, \vartheta) > 1$, and a measurable set $\mathcal{U}_\vartheta \subset B_n(1)$ with $\text{vol}(\mathcal{U}_\vartheta)/\text{vol}(B_n(1)) \geq 1 - Ce^{-n^\vartheta/C}$, such that for any $\varepsilon \in (0, 1/2]$, and $\# \in \{\text{pred}, \text{est}, \text{in}, \text{res}\}$,

$$\sup_{\mu_0 \in \mathcal{U}_\vartheta} \mathbb{P} \left(\sup_{\eta \in \Xi_K} |R_{(\Sigma, \mu_0)}^\#(\eta, \sigma_\xi) - \bar{R}_{(\Sigma, \mu_0)}^\#(\eta, \sigma_\xi)| \geq \varepsilon \right) \leq C \cdot \begin{cases} ne^{-n\varepsilon^{c_0}/C}, & Z = G; \\ \varepsilon^{-c_0} n^{-1/6.5}, & \text{otherwise.} \end{cases}$$

Proof All the constants in $\lesssim, \gtrsim, \asymp$ below may possibly depend on K .

(**Part 1**). We shall first extend the claim of Theorem 49 for $\# = \text{pred}$ to $\sigma_\xi^2 \in [0, K]$ in the case $\phi^{-1} \geq 1 + 1/K$. Note that uniformly in $\eta \in [0, K]$, for $\sigma_\xi, \sigma'_\xi \in [0, K]$,

$$\|\hat{\mu}_\eta(\sigma_\xi) - \hat{\mu}_\eta(\sigma'_\xi)\| \lesssim |\sigma_\xi - \sigma'_\xi| \cdot n^{-1} \|Z\|_{\text{op}} \|\xi_0\| \cdot \|(ZZ^\top/n)^{-1}\|_{\text{op}}. \quad (155)$$

Using the estimate (146), uniformly in $\eta \in [0, K]$, for all $\sigma_\xi, \sigma'_\xi \in [0, K]$,

$$\begin{aligned} |R_{(\Sigma, \mu_0)}^{\text{pred}}(\eta, \sigma_\xi) - R_{(\Sigma, \mu_0)}^{\text{pred}}(\eta, \sigma'_\xi)| &\lesssim \|\hat{\mu}_\eta(\sigma_\xi) - \hat{\mu}_\eta(\sigma'_\xi)\| \cdot (\|\hat{\mu}_\eta(\sigma_\xi)\| + \|\hat{\mu}_\eta(\sigma'_\xi)\| + \|\mu_0\|) \\ &\lesssim |\sigma_\xi - \sigma'_\xi| \cdot \|(ZZ^\top/n)^{-1}\|_{\text{op}}^2 \cdot \left(1 + \frac{\|Z\|_{\text{op}} + \|\xi_0\|}{\sqrt{n}}\right)^4. \end{aligned}$$

So on an event E_1 with $\mathbb{P}(E_1) \geq 1 - C_1 e^{-n/C_1}$, for $\sigma_\xi, \sigma'_\xi \in [0, K]$,

$$\sup_{\eta \in [0, K]} |R_{(\Sigma, \mu_0)}^{\text{pred}}(\eta, \sigma_\xi) - R_{(\Sigma, \mu_0)}^{\text{pred}}(\eta, \sigma'_\xi)| \leq C_1 \cdot |\sigma_\xi - \sigma'_\xi|.$$

On the other hand, using Lemma 56-(2),

$$\sup_{\eta \in [0, K]} |\bar{R}_{(\Sigma, \mu_0)}^{\text{pred}}(\eta, \sigma_\xi) - \bar{R}_{(\Sigma, \mu_0)}^{\text{pred}}(\eta, \sigma'_\xi)| \leq C_1 \cdot |\sigma_\xi - \sigma'_\xi|.$$

Using the above two displays, for any $\varepsilon > 0$, by choosing $\sigma'_\xi \equiv \varepsilon/(2C_1)$, we have for any $\sigma_\xi \leq \sigma'_\xi$,

$$\begin{aligned} &\mathbb{P} \left(\sup_{\eta \in [0, K]} |R_{(\Sigma, \mu_0)}^{\text{pred}}(\eta, \sigma_\xi) - \bar{R}_{(\Sigma, \mu_0)}^{\text{pred}}(\eta, \sigma_\xi)| \geq 2\varepsilon \right) \\ &\leq \mathbb{P} \left(\sup_{\eta \in [0, K]} |R_{(\Sigma, \mu_0)}^{\text{pred}}(\eta, \sigma'_\xi) - \bar{R}_{(\Sigma, \mu_0)}^{\text{pred}}(\eta, \sigma'_\xi)| \geq \varepsilon \right) + C_1 e^{-n/C_1}. \end{aligned} \quad (156)$$

The first term on the right hand side of the above display can be handled by the proven claim in Theorem 49, upon noting that (i) the constant C therein depends on K polynomially, and here we choose K to be larger than $2C_1/\varepsilon$; (ii) $(n/\varepsilon)^{C'} e^{-n\varepsilon^{C'}} \wedge 1 \leq ne^{-n\varepsilon^{C''}}$ holds for C'' chosen much larger than C' .

The extension of the claim of Theorem 49 for $\# = \text{est}$ to $\sigma_\xi^2 \in [0, K]$ follows a similar proof with minor modifications, so we omit the details.

(Part 2). Next we consider the case $\# = \text{in}$. We need to extend the corresponding claim of Theorem 49 to both $\sigma_\xi^2 \in [0, K]$ and $\eta \in [0, K]$.

We first verify the (high probability) Lipschitz continuity of the maps $\sigma_\xi \mapsto R_{(\Sigma, \mu_0)}^{\text{in}}(\eta, \sigma_\xi), \bar{R}_{(\Sigma, \mu_0)}^{\text{in}}(\eta, \sigma_\xi)$. Note that uniformly in $\eta \in [0, K]$, by virtue of (155), for any $\sigma_\xi, \sigma'_\xi \in [0, K]$,

$$\begin{aligned} & |R_{(\Sigma, \mu_0)}^{\text{in}}(\eta, \sigma_\xi) - R_{(\Sigma, \mu_0)}^{\text{in}}(\eta, \sigma'_\xi)| \\ & \lesssim \left(1 + \frac{\|Z\|_{\text{op}}}{\sqrt{n}}\right)^2 \cdot \|\widehat{\mu}_\eta(\sigma_\xi) - \widehat{\mu}_\eta(\sigma'_\xi)\| \cdot (\|\widehat{\mu}_\eta(\sigma_\xi)\| + \|\widehat{\mu}_\eta(\sigma'_\xi)\| + \|\mu_0\|) \\ & \lesssim |\sigma_\xi - \sigma'_\xi| \cdot \|(ZZ^\top/n)^{-1}\|_{\text{op}}^2 \cdot \left(1 + \frac{\|Z\|_{\text{op}} + \|\xi_0\|}{\sqrt{n}}\right)^6. \end{aligned}$$

This verifies the high probability Lipschitz property of $\sigma_\xi \mapsto R_{(\Sigma, \mu_0)}^{\text{in}}(\eta, \sigma_\xi)$. The Lipschitz property of $\sigma_\xi \mapsto \bar{R}_{(\Sigma, \mu_0)}^{\text{in}}(\eta, \sigma_\xi)$ is easily verified. From here we may use a similar argument to (156) to conclude the extension of the claim of Theorem 49 for $\# = \text{in}$ to $\sigma_\xi^2 \in [0, K]$.

Finally we verify the (high probability) Lipschitz continuity of the maps $\eta \mapsto R_{(\Sigma, \mu_0)}^{\text{in}}(\eta, \sigma_\xi), \bar{R}_{(\Sigma, \mu_0)}^{\text{in}}(\eta, \sigma_\xi)$. Using the estimates (92) (with G replaced by Z) and (146), uniformly in $\sigma_\xi \in [0, K]$ and $\eta_1, \eta_2 \in [0, K]$,

$$\begin{aligned} & |R_{(\Sigma, \mu_0)}^{\text{in}}(\eta_1, \sigma_\xi) - R_{(\Sigma, \mu_0)}^{\text{in}}(\eta_2, \sigma_\xi)| \\ & \lesssim \left(1 + \frac{\|Z\|_{\text{op}}}{\sqrt{n}}\right)^2 \cdot \|\widehat{\mu}_{\eta_1}(\sigma_\xi) - \widehat{\mu}_{\eta_2}(\sigma_\xi)\| \cdot (\|\widehat{\mu}_{\eta_1}(\sigma_\xi)\| + \|\widehat{\mu}_{\eta_2}(\sigma_\xi)\| + \|\mu_0\|) \\ & \lesssim \left(1 + \frac{\|Z\|_{\text{op}} + \|\xi_0\|}{\sqrt{n}}\right)^6 \cdot \|(ZZ^\top/n)^{-1}\|_{\text{op}}^3 \cdot |\eta_1 - \eta_2|. \end{aligned}$$

The Lipschitz property of $\sigma_\xi \mapsto \bar{R}_{(\Sigma, \mu_0)}^{\text{in}}(\eta, \sigma_\xi)$ is again easily verified. Again from here we may argue similarly to (156) to extend the claim of Theorem 49 for $\# = \text{in}$ to $\eta \in [0, K]$. The case for $\# = \text{res}$ is similar so we omit repetitive details. \blacksquare

Lemma 56. *Suppose $\phi^{-1} > 1$. The following hold.*

1. *The system of equations*

$$\begin{cases} \phi\gamma^2 = \mathbb{E} \text{err}_{(\Sigma, \mu_0)}(\gamma; \tau), \\ \phi - \frac{\eta}{\tau} = \gamma^{-2} \mathbb{E} \text{dof}_{(\Sigma, \mu_0)}(\gamma; \tau) = \frac{1}{n} \text{tr}((\Sigma + \tau I)^{-1} \Sigma) \end{cases}$$

admit a unique solution $(\gamma_{\eta,}(0), \tau_{\eta,*}(0)) \in [0, \infty) \times (0, \infty)$.*

2. It holds that $\tau_{\eta,*}(0) = \tau_{\eta,*}(\sigma_\xi)$. If furthermore $1 + 1/K \leq \phi^{-1} \leq K$ and $\|\Sigma\|_{\text{op}} \vee \mathcal{H}_\Sigma \leq K$ for some $K > 0$, then there exists some $C = C(K) > 0$ such that $|\gamma_{\eta,*}^2(\sigma_\xi) - \gamma_{\eta,*}^2(0)| \leq C\sigma_\xi^2$.

Proof The claim (1) follows verbatim from the proof of Proposition 23-(1) by setting $\sigma_\xi^2 = 0$ therein. The claim (2) follows by using the formula (46). \blacksquare

Proof [Proof of Theorem 54] Let $\mathcal{U}_\vartheta \subset B_n(1)$ be as specified in Theorem 49 or 4. In view of its explicit form given in Proposition 40, with $\mathcal{U}_{\delta,\vartheta} \equiv \mathcal{U}_\vartheta \cap (B_n(1) \setminus B_n(\delta))$, the volume estimates $\min \{ \text{vol}(\mathcal{U}_\vartheta)/\text{vol}(B_n(1)), \text{vol}(\mathcal{U}_{\delta,\vartheta})/\text{vol}(B_n(1) \setminus B_n(\delta)) \} \geq 1 - Ce^{-n^\vartheta/C}$ hold.

On the other hand, using the construction around (149), we may find some $\mathcal{V}_\varepsilon \subset B_n(1), \mathcal{V}_{\varepsilon,\delta} \subset B_n(1) \setminus B_n(\delta)$ (for the latter, we take $U_0 \sim \text{Unif}(\delta, 1)$ therein) with $\min \{ \text{vol}(\mathcal{V}_\varepsilon)/\text{vol}(B_n(1)), \text{vol}(\mathcal{V}_{\varepsilon,\delta})/\text{vol}(B_n(1) \setminus B_n(\delta)) \} \geq 1 - C\varepsilon^{-1}e^{-n\varepsilon^2/C}$, such that for $\# \in \{\text{pred, est, in}\}$,

$$\sup_{\mu_0 \in \{\mathcal{V}_\varepsilon, \mathcal{V}_{\varepsilon,\delta}\}} \sup_{\eta \in \Xi_L} |\bar{R}_{(\Sigma, \mu_0)}^\#(\eta) - \mathcal{R}_{(\Sigma, \mu_0)}^\#(\eta)| \leq \varepsilon. \quad (157)$$

Now let

$$\mathcal{W}_{\varepsilon,\vartheta} \equiv \mathcal{U}_\vartheta \cap \mathcal{V}_\varepsilon, \quad \mathcal{W}_{\varepsilon,\delta,\vartheta} \equiv \mathcal{U}_{\delta,\vartheta} \cap \mathcal{V}_{\varepsilon,\delta}. \quad (158)$$

Then we have the volume estimates $\min \{ \text{vol}(\mathcal{W}_{\varepsilon,\vartheta})/\text{vol}(B_n(1)), \text{vol}(\mathcal{W}_{\varepsilon,\delta,\vartheta})/\text{vol}(B_n(1) \setminus B_n(\delta)) \} \geq 1 - C\varepsilon^{-1}e^{-n\varepsilon^2/C} - Ce^{-n^\vartheta/C}$.

Moreover, by Proposition 53, provided $\eta_* \equiv \sigma_\xi^2/\|\mu_0\|^2 = \text{SNR}_{\mu_0}^{-1} \in \Xi_L$,

$$\|\mu_0\|^2/C_0 \leq \frac{|\mathcal{R}_{(\Sigma, \mu_0)}^\#(\eta) - \mathcal{R}_{(\Sigma, \mu_0)}^\#(\eta_*)|}{(\eta - \eta_*)^2} \leq C_0\|\mu_0\|^2 \quad (159)$$

holds uniformly in $\eta \in \Xi_L$ for some $C_0 > 0$.

(**Noisy case** $\sigma_\xi^2 \in [1/K, K]$). Fix $\mu_0 \in \mathcal{W}_{\varepsilon,\delta,\vartheta}$. Under the assumed conditions, $\eta_* \in \Xi_L$. So using the estimates (157) and (159), for any $\eta' \geq \eta_*$,

$$\bar{R}_{(\Sigma, \mu_0)}^\#(\eta') - \inf_{\eta \in \Xi_L} \bar{R}_{(\Sigma, \mu_0)}^\#(\eta) \geq \mathcal{R}_{(\Sigma, \mu_0)}^\#(\eta') - \inf_{\eta \in \Xi_L} \mathcal{R}_{(\Sigma, \mu_0)}^\#(\eta) - 2\varepsilon \geq \frac{\delta^2(\eta' - \eta_*)^2}{C_0} - 2\varepsilon.$$

Combined with a similar inequality for $\eta' \leq \eta_*$, we conclude that for any $\mu_0 \in \mathcal{W}_{\varepsilon,\delta,\vartheta}$ and $\eta' \in \Xi_L$,

$$|\bar{R}_{(\Sigma, \mu_0)}^\#(\eta') - \inf_{\eta \in \Xi_L} \bar{R}_{(\Sigma, \mu_0)}^\#(\eta)| \geq \frac{\delta^2(\eta' - \eta_*)^2}{C_0} - 2\varepsilon.$$

Now for $|\eta' - \eta_*| \geq \Delta$, choosing $\varepsilon \equiv \varepsilon_0 \equiv \delta^2\Delta^2/(4C_0)$, we have

$$\inf_{\mu_0 \in \mathcal{W}_{\varepsilon,\delta,\vartheta}} \inf_{\eta' \in \Xi_L: |\eta' - \eta_*| \geq \Delta} |\bar{R}_{(\Sigma, \mu_0)}^\#(\eta') - \inf_{\eta \in \Xi_L} \bar{R}_{(\Sigma, \mu_0)}^\#(\eta)| \geq \frac{\delta^2\Delta^2}{2C_0}.$$

From here the claim follows from Theorem 49.

(**Noiseless case** $\sigma_\xi^2 = 0$). In this case, (159) implies that the map $\eta \mapsto \mathcal{R}_{(\Sigma, \mu_0)}^\#(\eta)$ attains global minimum at $\eta = 0$. So together with (157), it implies that uniformly in $\mu_0 \in \mathcal{W}_{\varepsilon, \vartheta}$,

$$\left| \min_{\eta \in [0, K]} \bar{R}_{(\Sigma, \mu_0)}^\#(\eta) - \bar{R}_{(\Sigma, \mu_0)}^\#(0) \right| \leq \varepsilon.$$

From here the claim follows from Lemma 55 that holds for $\sigma_\xi = 0$. \blacksquare

Appendix G. Proofs for Section 4

G.1 A rigorous version of (21) and its proof

Theorem 57. *Suppose Assumption A holds, and $1/K \leq \phi^{-1} \leq K$, $\|\Sigma^{-1}\|_{\text{op}} \vee \|\Sigma\|_{\text{op}} \leq K$ hold for some $K > 0$.*

1. *For any small $\varepsilon > 0$, there exists some $C_1 = C_1(K, \varepsilon) > 0$ such that*

$$\mathbb{P} \left(\sup_{\eta \in \Xi_K} |\hat{\tau}_\eta - \tau_{\eta, *}| \geq n^{-1/2+\varepsilon} \right) \leq C_1 n^{-100}.$$

2. *Suppose further Assumption B holds with either (i) $\sigma_\xi^2 \in [1/K, K]$ or (ii) $\sigma_\xi^2 \in [0, K]$ with $1 + 1/K \leq \phi^{-1} \leq K$. Fix a small enough constant $\vartheta \in (0, 1/50)$. Then there exist a constant $C_2 = C_2(K, \vartheta) > 1$, and a measurable set $\mathcal{U}_\vartheta \subset B_n(1)$ with $\text{vol}(\mathcal{U}_\vartheta)/\text{vol}(B_n(1)) \geq 1 - Ce^{-n^\vartheta/C}$, such that*

$$\sup_{\mu_0 \in \mathcal{U}_\vartheta} \mathbb{P} \left(\sup_{\eta \in \Xi_K} |\hat{\gamma}_\eta - \gamma_{\eta, *}| \geq n^{-\vartheta} \right) \leq C_2 n^{-1/7}.$$

Proof [Proof of Theorem 57 for $\hat{\tau}_\eta$] All the constants in $\lesssim, \gtrsim, \asymp$ may depend on K .

Let κ_0 be defined in the same way as in the proof of Proposition 40. Using a similar local law and continuity argument as in the proof of that proposition, on an event E_0 with $\mathbb{P}(E_0) \geq 1 - Cn^{-D}$,

$$\sup_{\eta \in \Xi_K} |m^{-1} \text{tr}(\check{\Sigma} + (\eta/\phi)I)^{-1} - \mathbf{m}(-\eta/\phi)| \lesssim \kappa_0^{-1} n^{-1/2+\varepsilon}.$$

So on $E_0 \cap \mathcal{E}(C_1)$, where $\mathcal{E}(C_1) \equiv \{\|Z\|_{\text{op}}/\sqrt{n} \leq C_1\}$ with $\mathbb{P}(\mathcal{E}(C_1)) \geq 1 - Ce^{-n/C}$, uniformly in $\eta \in \Xi_K$,

$$|\hat{\tau}_\eta - \tau_{\eta, *}| \leq \frac{|\frac{1}{m} \text{tr}(\check{\Sigma} + \frac{\eta}{\phi}I)^{-1} - \mathbf{m}(-\frac{\eta}{\phi})|}{\frac{1}{m} \text{tr}(\check{\Sigma} + \frac{\eta}{\phi}I)^{-1} \cdot \mathbf{m}(-\frac{\eta}{\phi})} \lesssim \left\{ C_1^2 \mathbf{1}_{\phi^{-1} \geq 1+1/K} \wedge \eta^{-1} \right\} \cdot \kappa_0^{-1} n^{-1/2+\varepsilon}.$$

Here in the last inequality, we use the following estimate for $\mathbf{m}(z)$: As \mathbf{m} is the Stieltjes transform of ρ (cf. (Knowles and Yin, 2017, Lemma 2.2)), $\mathbf{m}(z) \geq 0$ for $z \leq 0$, and

$$\frac{1}{\mathbf{m}(z)} = (-z) + \frac{1}{m} \text{tr} \left((I + \Sigma \mathbf{m}(z))^{-1} \Sigma \right) \lesssim 1 + |z|.$$

The claim follows. \blacksquare

Proof [Proof of Theorem 57 for $\hat{\gamma}_\eta$] All the constants in $\lesssim, \gtrsim, \asymp$ may depend on K .

Using Theorem 49, the stability of $\tau_{\eta,*}$ in Proposition 23, and the proven fact in (1) on $\hat{\tau}_\eta$, it holds for $\varepsilon \in (0, 1/2]$ that

$$\mathbb{P}\left(\sup_{\eta \in [1/K, K]} |\eta^{-1} \hat{\tau}_\eta \|\hat{r}_\eta(\sigma_\xi)\| - \gamma_{\eta,*}(\sigma_\xi)| \geq \varepsilon\right) \leq C_1 \varepsilon^{-c_0} n^{-1/6.5}. \quad (160)$$

Next we consider extension to $\eta \in [0, K]$ in the regime $\phi^{-1} \geq 1 + 1/K$. By KKT condition, we have $n^{-1} X^\top (Y - X \hat{\mu}_\eta) = \eta \hat{\mu}_\eta$, so a.s. $\hat{r}_\eta/\eta = (Y - X \mu_\eta)/(\sqrt{n}\eta) = \sqrt{n}(X X^\top)^{-1} X \hat{\mu}_\eta$ for any $\eta > 0$. So we only need to verify the high probability Lipschitz continuity for $\eta \mapsto \sqrt{n} \hat{\tau}_\eta (X X^\top)^{-1} X \hat{\mu}_\eta$: for any $\eta_1, \eta_2 \in [0, K]$, using the estimate (92) (with G replaced by Z) we obtain, for some universal $c_0 > 1$,

$$\begin{aligned} & \left| \sqrt{n} \hat{\tau}_{\eta_1} \|(X X^\top)^{-1} X \hat{\mu}_{\eta_1}\| - \sqrt{n} \hat{\tau}_{\eta_2} \|(X X^\top)^{-1} X \hat{\mu}_{\eta_2}\| \right| \\ & \lesssim \|(Z Z^\top/n)^{-1}\|_{\text{op}} \cdot (\|Z\|_{\text{op}}/\sqrt{n}) \cdot \left(|\hat{\tau}_{\eta_1} - \hat{\tau}_{\eta_2}| \cdot \|\hat{\mu}_{\eta_1}\| + |\hat{\tau}_{\eta_2}| \cdot \|\hat{\mu}_{\eta_1} - \hat{\mu}_{\eta_2}\| \right) \\ & \lesssim \left(1 + \frac{\|Z\|_{\text{op}} + \|\xi_0\|}{\sqrt{n}} + \|(Z Z^\top/n)^{-1}\|_{\text{op}} \right)^{c_0} \cdot |\eta_2 - \eta_1|. \end{aligned}$$

Finally we consider extension to $\sigma_\xi^2 \in [0, K]$ in the same regime $\phi^{-1} \geq 1 + 1/K$ by verifying a similar high probability uniform-in- η Lipschitz continuity property for $\sigma_\xi \mapsto \sqrt{n}(X X^\top)^{-1} X \hat{\mu}_\eta(\sigma_\xi)$: for any $\sigma_\xi, \sigma'_\xi \in [0, K]$, using the estimate (155),

$$\begin{aligned} & \sup_{\eta \in [0, K]} \left| \sqrt{n} \hat{\tau}_\eta \|(X X^\top)^{-1} X \hat{\mu}_\eta(\sigma_\xi)\| - \sqrt{n} \hat{\tau}_\eta \|(X X^\top)^{-1} X \hat{\mu}_\eta(\sigma'_\xi)\| \right| \\ & \lesssim \|(Z Z^\top/n)^{-1}\|_{\text{op}}^2 \cdot (\|Z\|_{\text{op}}/\sqrt{n}) \cdot \sup_{\eta \in [0, K]} \|\hat{\mu}_\eta(\sigma_\xi) - \hat{\mu}_\eta(\sigma'_\xi)\| \\ & \lesssim \left(1 + \frac{\|Z\|_{\text{op}} + \|\xi_0\|}{\sqrt{n}} + \|(Z Z^\top/n)^{-1}\|_{\text{op}} \right)^{c_0} \cdot |\sigma_\xi - \sigma'_\xi|. \end{aligned}$$

The claimed bound follows. \blacksquare

G.2 Proof of Theorem 10

Recall we have $\gamma_{\eta,*}^2 = \phi^{-1}(\sigma_\xi^2 + \bar{R}_{(\Sigma, \mu_0)}^{\text{pred}}(\eta))$. For both the case $\sigma_\xi^2 \in [1/K, K]$ and $\sigma_\xi^2 \in [0, K]$ with $\phi^{-1} \geq 1 + 1/K$, we take $\mathcal{W}_{\varepsilon, \delta, \vartheta} \subset B_n(1) \setminus B_n(\delta)$ as constructed in (158) in the proof of Theorem 54, with $\varepsilon \equiv \varepsilon_n \equiv n^{-\vartheta}$. Fix $\mu_0 \in \mathcal{W}_{\varepsilon, \delta, \vartheta}$, then $\eta_* = \text{SNR}_{\mu_0}^{-1} \in \Xi_L$. Using Theorems 8 and 57, on an event E_0 with $\mathbb{P}(E_0^c) \leq C n^{-1/7}$,

$$\sup_{\eta \in \Xi_L} \left| \hat{\gamma}_\eta^2 - \phi^{-1}(\sigma_\xi^2 + \mathcal{R}_{(\Sigma, \mu_0)}^{\text{pred}}(\eta)) \right| \leq \varepsilon. \quad (161)$$

This in particular implies that on E_0 , both the following inequalities hold:

$$\phi \hat{\gamma}_{\hat{\eta}^{\text{GCV}}}^2 - \sigma_\xi^2 - \phi \varepsilon \leq \mathcal{R}_{(\Sigma, \mu_0)}^{\text{pred}}(\hat{\eta}^{\text{GCV}}) \leq \phi \hat{\gamma}_{\hat{\eta}^{\text{GCV}}}^2 - \sigma_\xi^2 + \phi \varepsilon,$$

$$\phi \min_{\eta \in \Xi_L} \hat{\gamma}_\eta^2 - \sigma_\xi^2 - \phi\varepsilon \leq \min_{\eta \in \Xi_L} \mathcal{R}_{(\Sigma, \mu_0)}^{\text{pred}}(\eta) \leq \phi \min_{\eta \in \Xi_L} \hat{\gamma}_\eta^2 - \sigma_\xi^2 + \phi\varepsilon. \quad (162)$$

Using the definition of $\hat{\eta}^{\text{GCV}}$ which gives $\hat{\gamma}_{\hat{\eta}^{\text{GCV}}}^2 = \min_{\eta \in \Xi_L} \hat{\gamma}_\eta^2$, the above two displays can be used to relate $\mathcal{R}_{(\Sigma, \mu_0)}^{\text{pred}}(\hat{\eta}^{\text{GCV}})$ and $\min_{\eta \in \Xi_L} \mathcal{R}_{(\Sigma, \mu_0)}^{\text{pred}}(\eta)$: on the event E_0 ,

$$|\mathcal{R}_{(\Sigma, \mu_0)}^{\text{pred}}(\hat{\eta}^{\text{GCV}}) - \min_{\eta \in \Xi_L} \mathcal{R}_{(\Sigma, \mu_0)}^{\text{pred}}(\eta)| \leq 2\phi\varepsilon. \quad (163)$$

As $\eta_* \in \Xi_L$, $\min_{\eta \in \Xi_L} \mathcal{R}_{(\Sigma, \mu_0)}^\#(\eta) = \mathcal{R}_{(\Sigma, \mu_0)}^\#(\eta_*)$ for $\# \in \{\text{pred}, \text{est}, \text{in}\}$. Consequently, by the second inequality in Proposition 53, we have on the event E_0 ,

$$|\hat{\eta}^{\text{GCV}} - \eta_*| \leq \frac{C}{\|\mu_0\|} |\mathcal{R}_{(\Sigma, \mu_0)}^{\text{pred}}(\hat{\eta}^{\text{GCV}}) - \mathcal{R}_{(\Sigma, \mu_0)}^{\text{pred}}(\eta_*)|^{1/2} \leq C_1 \varepsilon^{1/2}. \quad (164)$$

This means on E_0 , for both $\# \in \{\text{est}, \text{in}\}$,

$$|\mathcal{R}_{(\Sigma, \mu_0)}^\#(\hat{\eta}^{\text{GCV}}) - \min_{\eta \in \Xi_L} \mathcal{R}_{(\Sigma, \mu_0)}^\#(\eta)| \leq C_2 \varepsilon.$$

We may conclude from here by virtues of Theorems 49 and 8, together with Lemma 55. \square

G.3 Proof of Theorem 12

Lemma 58. *Consider the following version of (13) with sample size $m - m_\ell$:*

$$\begin{cases} \frac{m - m_\ell}{n} \cdot \gamma^2 = \sigma_\xi^2 + \mathbb{E} \text{err}_{(\Sigma, \mu_0)}(\gamma; \tau), \\ \left(\frac{m - m_\ell}{n} - \frac{\eta}{\tau}\right) \cdot \gamma^2 = \mathbb{E} \text{dof}_{(\Sigma, \mu_0)}(\gamma; \tau). \end{cases} \quad (165)$$

1. *The fixed point equation (165) admits a unique solution $(\gamma_{\eta, *}^{(\ell)}, \tau_{\eta, *}^{(\ell)}) \in (0, \infty)^2$, for all $(m, n) \in \mathbb{N}^2$ when $\eta > 0$ and $m < n$ when $\eta = 0$.*
2. *Further suppose $1/K \leq \phi^{-1}, \sigma_\xi^2 \leq K$, $m_\ell/n \leq 1/(2K)$ and $\|\Sigma^{-1}\|_{\text{op}} \vee \|\Sigma\|_{\text{op}} \leq K$ for some $K > 10$. Then there exists some $C = C(K) > 1$ such that uniformly in $\eta \in \Xi_K$, $1/C \leq \gamma_{\eta, *}^{(\ell)}, \tau_{\eta, *}^{(\ell)} \leq C$. Moreover,*

$$|\gamma_{\eta, *}^{(\ell)} - \gamma_{\eta, *}| \vee |\tau_{\eta, *}^{(\ell)} - \tau_{\eta, *}| \leq \frac{C m_\ell}{n}.$$

Proof All the constants in $\lesssim, \gtrsim, \asymp$ below may depend on K . We only need to prove (2). The method of proof is similar to that of Proposition 23-(3). Instead of considering (165), we shall consider the system of equations

$$\begin{cases} \phi - \alpha = \frac{1}{\gamma^2} (\sigma_\xi^2 + \tau^2 \|(\Sigma + \tau I)^{-1} \Sigma^{1/2} \mu_0\|^2) + \frac{1}{n} \text{tr}((\Sigma + \tau I)^{-2} \Sigma^2), \\ \phi - \alpha = \frac{1}{n} \text{tr}((\Sigma + \tau I)^{-1} \Sigma) + \frac{\eta}{\tau}, \end{cases} \quad (166)$$

indexed by $\alpha \geq 0$. For $\alpha \in [0, 1/(2K)]$, the solution $(\gamma_{\eta, *}(\alpha), \tau_{\eta, *}(\alpha))$ exists uniquely for $\eta > 0$ and also for $\eta = 0$ if additionally $m < n$. Moreover, using the apriori estimate in Proposition 23-(2), we have uniformly in $\eta \in \Xi_K$ and $\alpha \in [0, 1/(2K)]$, $\gamma_{\eta, *}(\alpha), \tau_{\eta, *}(\alpha) \asymp 1$.

Now differentiating on both sides of the second equation in (166) with respect to α , we obtain

$$1 = \left(n^{-1} \operatorname{tr} \left((\Sigma + \tau_{\eta,*}(\alpha)I)^{-2} \Sigma \right) + \eta \tau_{\eta,*}^{-2}(\alpha) \right) \cdot \tau'_{\eta,*}(\alpha).$$

This means uniformly in $\eta \in \Xi_K$ and $\alpha \in [0, 1/(2K)]$, $\tau'_{\eta,*}(\alpha) \asymp 1$. Next, using the first equation in (166), we obtain

$$\gamma_{\eta,*}^2(\alpha) = \frac{\sigma_\xi^2 + \tau_{\eta,*}^2(\alpha) \|(\Sigma + \tau_{\eta,*}(\alpha)I)^{-1} \Sigma^{1/2} \mu_0\|^2}{\phi - \alpha - \frac{1}{n} \operatorname{tr} \left((\Sigma + \tau_{\eta,*}(\alpha)I)^{-2} \Sigma^2 \right)} \equiv \frac{G_{1,\eta}(\alpha)}{G_{2,\eta}(\alpha)}.$$

Using similar calculations as in (52)-(53), we have uniformly in $\eta \in \Xi_K$ and $\alpha \in [0, 1/(2K)]$, $G_{1,\eta}(\alpha), G_{2,\eta}(\alpha) \asymp 1$, and $|G'_{1,\eta}(\alpha)| \vee |G'_{2,\eta}(\alpha)| \lesssim 1$. This concludes the claim. \blacksquare

Proof [Proof of Theorem 12] All the constants in $\lesssim, \gtrsim, \asymp$ below may depend on K, L .

As $\|Y^{(\ell)} - X^{(\ell)} \widehat{\mu}_\eta^{(\ell)}\|^2 = \|Z^{(\ell)} \Sigma^{1/2} (\mu_0 - \widehat{\mu}_\eta^{(\ell)}) + \xi^{(\ell)}\|^2$ and $\widehat{\mu}_\eta^{(\ell)}$ is independent of $(Z^{(\ell)}, \xi^{(\ell)})$, by using Lemma 61 first conditionally on $(Z^{(-\ell)}, \xi^{(-\ell)})$ and then further taking expectation over $(Z^{(-\ell)}, \xi^{(-\ell)})$, we have for $0 < \varrho \leq 1$,

$$\begin{aligned} \mathbb{P} \left(E_{0,\ell}^c(\eta) \equiv \left\{ |m_\ell^{-1} \|Y^{(\ell)} - X^{(\ell)} \widehat{\mu}_\eta^{(\ell)}\|^2 - (\|\Sigma^{1/2} (\widehat{\mu}_\eta^{(\ell)} - \mu_0)\|^2 + \sigma_\xi^2) \right\} \right. \\ \left. \geq C_0 (\sigma_\xi^2 \vee \|\Sigma^{1/2} (\widehat{\mu}_\eta^{(\ell)} - \mu_0)\|^2) m_\ell^{-(1-\varrho)/2} \right\} \leq C_0 e^{-m_\ell^\varrho / C_0}. \end{aligned}$$

Here $C_0 > 0$ is a universal constant. Using similar arguments as in (147) (by noting that the normalization in $\widehat{\mu}_\eta^{(\ell)}$ is still n), there exists some constant $C_1 > 0$ such that for any $\ell \in [k]$, on an event $E_{1,\ell}$ with $\mathbb{P}(E_{1,\ell}^c) \leq C_1 e^{-m_\ell / C_1}$, $\sup_{\eta \in \Xi_L} \|\widehat{\mu}_\eta^{(\ell)}\| \leq C_1$. This means that for any $\eta \in \Xi_L$, on the event $\cap_{\ell \in [k]} (E_{0,\ell}(\eta) \cap E_{1,\ell})$,

$$\max_{\ell \in [k]} m_\ell^{(1-\varrho)/2} \cdot |m_\ell^{-1} \|Y^{(\ell)} - X^{(\ell)} \widehat{\mu}_\eta^{(\ell)}\|^2 - (\|\Sigma^{1/2} (\widehat{\mu}_\eta^{(\ell)} - \mu_0)\|^2 + \sigma_\xi^2)| \leq C'_1. \quad (167)$$

On the other hand, using Theorem 49, we may find some $\mathcal{U}_{\vartheta;\ell} \subset B_n(1)$ with $\operatorname{vol}(\mathcal{U}_{\vartheta;\ell}) / \operatorname{vol}(B_n(1)) \geq 1 - C_2 e^{-n^\vartheta / C_2}$, such that for $\varepsilon \in (0, 1/2]$, on an event $E_{2,\ell}(\varepsilon)$ with $\mathbb{P}(E_{2,\ell}^c(\varepsilon)) \leq C_2 (n e^{-n\varepsilon^4 / C_2} + \varepsilon^{-c_0} n^{-1/6.5} \mathbf{1}_{Z \neq G})$, for $\mu_0 \in \mathcal{U}_{\vartheta;\ell}$,

$$\sup_{\eta \in \Xi_L} \left| \|\Sigma^{1/2} (\widehat{\mu}_\eta^{(\ell)} - \mu_0)\|^2 - \left\{ \frac{m - m_\ell}{n} (\gamma_{\eta,*}^{(\ell)})^2 - \sigma_\xi^2 \right\} \right| \leq \varepsilon.$$

Here $\gamma_{\eta,*}^{(\ell)}$ is taken from Lemma 58, and we extend the definition to $\ell = 0$ with $\widehat{\mu}_\eta^{(0)} \equiv \widehat{\mu}_\eta$ and $\gamma_{\eta,*}^{(0)} \equiv \gamma_{\eta,*}$. Using the statement (2) of the same Lemma 58, on the event $E_{2,\ell}(\varepsilon)$, we then have

$$\sup_{\eta \in \Xi_L} \left| \|\Sigma^{1/2} (\widehat{\mu}_\eta^{(\ell)} - \mu_0)\|^2 - \left\{ \phi \gamma_{\eta,*}^2 - \sigma_\xi^2 \right\} \right| \leq \varepsilon + \frac{C_2 m_\ell}{n}.$$

Replacing $\phi\gamma_{\eta,*}^2 - \sigma_\xi^2$ by $R_{(\Sigma,\mu_0)}^{\text{pred}}(\eta) = \|\Sigma^{1/2}(\widehat{\mu}_\eta - \mu_0)\|^2$ yields that, on $\cap_{\ell \in [0:k]} E_{2,\ell}(\varepsilon)$,

$$\sup_{\eta \in \Xi_L} \left| \|\Sigma^{1/2}(\widehat{\mu}_\eta^{(\ell)} - \mu_0)\|^2 - R_{(\Sigma,\mu_0)}^{\text{pred}}(\eta) \right| \leq 2\varepsilon + \frac{C_2 m_\ell}{n}. \quad (168)$$

Combining (167)-(168), for $\mu_0 \in \mathcal{U}_\vartheta \equiv \cap_{\ell \in [0:k]} \mathcal{U}_{\vartheta;\ell}$, $\varepsilon \in (0, 1/2]$ and $\eta \in \Xi_L$,

$$\begin{aligned} & \mathbb{P} \left(\left| R_{(\Sigma,\mu_0)}^{\text{CV},k}(\eta) - (R_{(\Sigma,\mu_0)}^{\text{pred}}(\eta) + \sigma_\xi^2) \right| \geq C'_2 \cdot \left\{ \frac{1}{k} \sum_{\ell \in [k]} \frac{1}{m_\ell^{(1-\varrho)/2}} + \frac{1}{k} + \varepsilon \right\} \right) \\ & \leq C'_2 \cdot \begin{cases} \sum_{\ell \in [k]} e^{-m_\ell^\varrho/C_0} + kn e^{-n\varepsilon^4/C_2}, & Z = G; \\ \sum_{\ell \in [k]} e^{-m_\ell^\varrho/C_0} + \varepsilon^{-c_0} \cdot kn^{-1/6.5}, & \text{otherwise.} \end{cases} \end{aligned} \quad (169)$$

Now we strengthen the estimate (169) into a uniform version. It is easy to verify that on an event $E_{3,\ell}$ with $\mathbb{P}(E_{3,\ell}^c) \leq C_3 e^{-m_\ell/C_3}$, $\|Z^{(\ell)}\|_{\text{op}} \leq C_3(\sqrt{m_\ell} + \sqrt{n})$, $\|\xi^{(\ell)}\| \leq C_3\sqrt{m_\ell}$, and for $\eta_1, \eta_2 \in \Xi_L$, $\|\widehat{\mu}_{\eta_1}^{(\ell)} - \widehat{\mu}_{\eta_2}^{(\ell)}\| \leq C_3|\eta_1 - \eta_2|$. So on $\cap_{\ell \in [k]} (E_{1,\ell} \cap E_{3,\ell})$, for $\eta_1, \eta_2 \in \Xi_L$,

$$\begin{aligned} |R_{(\Sigma,\mu_0)}^{\text{CV},k}(\eta_1) - R_{(\Sigma,\mu_0)}^{\text{CV},k}(\eta_2)| & \lesssim \frac{1}{k} \sum_{\ell \in [k]} \frac{1}{m_\ell} \left\| \|Z^{(\ell)}\|_{\text{op}} \|\widehat{\mu}_{\eta_1}^{(\ell)} - \widehat{\mu}_{\eta_2}^{(\ell)}\| \cdot (\|Z^{(\ell)}\|_{\text{op}} + \|\xi^{(\ell)}\|) \right\| \\ & \lesssim \frac{1}{k} \sum_{\ell \in [k]} \frac{m_\ell + n}{m_\ell} \cdot |\eta_1 - \eta_2| \leq C'_3 \cdot \frac{1}{k} \sum_{\ell \in [k]} \frac{n}{m_\ell} \cdot |\eta_1 - \eta_2|, \end{aligned}$$

and

$$|R_{(\Sigma,\mu_0)}^{\text{pred}}(\eta_1) - R_{(\Sigma,\mu_0)}^{\text{pred}}(\eta_2)| \leq C'_3 |\eta_1 - \eta_2|.$$

From here, using (i) (169) along with a discretization and union bound that strengthens (169) to a uniform control, and (ii) Theorem 49 which replaces $R_{(\Sigma,\mu_0)}^{\text{pred}}(\eta)$ by $\bar{R}_{(\Sigma,\mu_0)}^{\text{pred}}(\eta)$, we obtain for $\mu_0 \in \mathcal{U}_\vartheta$ and $\varepsilon \in (0, 1/2]$,

$$\begin{aligned} & \mathbb{P} \left(\sup_{\eta \in \Xi_L} \left| R_{(\Sigma,\mu_0)}^{\text{CV},k}(\eta) - (\bar{R}_{(\Sigma,\mu_0)}^{\text{pred}}(\eta) + \sigma_\xi^2) \right| \geq C''_3 \cdot \left\{ \frac{1}{k} \sum_{\ell \in [k]} \frac{1}{m_\ell^{(1-\varrho)/2}} + \frac{1}{k} + \varepsilon \right\} \equiv \varepsilon_{\{m_\ell\}} \right) \\ & \leq \mathfrak{p}_0 \equiv \frac{C''_3}{\varepsilon k} \sum_{\ell \in [k]} \frac{n}{m_\ell} \cdot \begin{cases} \sum_{\ell \in [k]} e^{-m_\ell^\varrho/C_0} + kn e^{-n\varepsilon^4/C_2}, & Z = G; \\ \sum_{\ell \in [k]} e^{-m_\ell^\varrho/C_0} + \varepsilon^{-c_0} \cdot kn^{-1/6.5}, & \text{otherwise.} \end{cases} \end{aligned}$$

Now with the same $\mathcal{W}_{\varepsilon,\delta,\vartheta}$ as in (158) using $\varepsilon \equiv \varepsilon_n \equiv n^{-\vartheta}$, for any $\mu_0 \in \mathcal{U}_\vartheta \cap \mathcal{W}_{\varepsilon,\delta,\vartheta}$, we may further replace $\bar{R}_{(\Sigma,\mu_0)}^{\text{pred}}(\eta)$ by $\mathcal{R}_{(\Sigma,\mu_0)}^{\text{pred}}(\eta)$ in the above display (with a possibly slightly larger $\varepsilon_{\{m_\ell\}}$, but for notational simplicity we abuse this notation). In summary, for any $\mu_0 \in \mathcal{U}_\vartheta \cap \mathcal{W}_{\varepsilon,\delta,\vartheta}$, on an event E_4 with $\mathbb{P}(E_4^c) \leq \mathfrak{p}_0$,

$$\sup_{\eta \in \Xi_L} \left| R_{(\Sigma,\mu_0)}^{\text{CV},k}(\eta) - (\mathcal{R}_{(\Sigma,\mu_0)}^{\text{pred}}(\eta) + \sigma_\xi^2) \right| \leq \varepsilon_{\{m_\ell\}}.$$

From here, using similar arguments as in (162)-(163), on the event E_4 ,

$$\left| \mathcal{R}_{(\Sigma,\mu_0)}^{\text{pred}}(\widehat{\eta}^{\text{CV}}) - \min_{\eta \in \Xi_L} \mathcal{R}_{(\Sigma,\mu_0)}^{\text{pred}}(\eta) \right| \leq 2\varepsilon_{\{m_\ell\}}.$$

Similar to (164), on the event E_4 , we have

$$|\hat{\eta}^{\text{CV}} - \eta_*| \leq C_4 \cdot \varepsilon_{\{m_\ell\}}^{1/2}. \quad (170)$$

From here we may argue along the same lines as those following (164) in the proof of Theorem 10 to conclude with probability estimated at \mathbf{p}_0 , by further noting that $\mathcal{U}_\vartheta \cap \mathcal{W}_{\varepsilon, \delta, \vartheta}$ satisfies the desired volume estimate. Under the further condition $\min_{\ell \in [k]} m_\ell \geq \log^{2/\delta} m$, by taking $\varrho = \delta$, \mathbf{p}_0 simplifies as indicated in the statement of the theorem for n large. \blacksquare

G.4 Proof of Theorem 13

We only prove the case for $\# = \text{GCV}$; the other case is similar. All constants in $\lesssim, \gtrsim, \asymp$ and \mathcal{O} may possibly depend on K, L . Let $\mathcal{W}_{\varepsilon, \delta, \vartheta} \subset B_n(1) \setminus B_n(\delta)$ be as constructed in (158) with $\varepsilon \equiv \varepsilon_n \equiv n^{-\vartheta}$.

(1). We first prove the statement for the length of the CI. Note that $|\text{CI}_j(\eta)| = 2\hat{\gamma}_\eta(\Sigma^{-1})_{jj}^{1/2} z_{\alpha/2} / \sqrt{n}$. By Theorem 57-(2), on an event E_0 with the probability indicated therein,

$$\max_{j \in [n]} \sup_{\eta \in \Xi_L} \left| |\text{CI}_j(\eta)| - 2\gamma_{\eta, *}(\Sigma^{-1})_{jj}^{1/2} \frac{z_{\alpha/2}}{\sqrt{n}} \right| \leq \frac{2\|\Sigma^{-1}\|_{\text{op}}^{1/2} z_{\alpha/2}}{\sqrt{n}} \sup_{\eta \in \Xi_L} |\hat{\gamma}_\eta - \gamma_{\eta, *}| \lesssim \frac{z_{\alpha/2}}{\sqrt{n}} \cdot \varepsilon.$$

Consequently, on the event E_0 , for any $\mu_0 \in \mathcal{W}_{\varepsilon, \delta, \vartheta}$,

$$\begin{aligned} & \sqrt{n} z_{\alpha/2}^{-1} \cdot \max_{j \in [n]} \left| |\text{CI}_j(\hat{\eta}^{\text{GCV}})| - \min_{\eta \in \Xi_L} |\text{CI}_j(\eta)| \right| \\ & \lesssim |\gamma_{\hat{\eta}^{\text{GCV}}, *}| - \min_{\eta \in \Xi_L} |\gamma_{\eta, *}| + \varepsilon \\ & \lesssim |\gamma_{\hat{\eta}^{\text{GCV}}, *}^2| - \min_{\eta \in \Xi_L} |\gamma_{\eta, *}^2| + \varepsilon \quad (\text{using Proposition 23-(3)}) \\ & \asymp \left| \bar{R}_{(\Sigma, \mu_0)}^{\text{pred}}(\hat{\eta}^{\text{GCV}}) - \min_{\eta \in \Xi_L} \bar{R}_{(\Sigma, \mu_0)}^{\text{pred}}(\eta) \right| + \varepsilon \quad (\text{using definition of } \gamma_{\eta, *}^2) \\ & \lesssim \left| \mathcal{R}_{(\Sigma, \mu_0)}^{\text{pred}}(\hat{\eta}^{\text{GCV}}) - \min_{\eta \in \Xi_L} \mathcal{R}_{(\Sigma, \mu_0)}^{\text{pred}}(\eta) \right| + \varepsilon \quad (\text{using Theorem 8}). \end{aligned}$$

As in the proof of Theorem 10, for $\sigma_\xi^2 \leq K$, $\eta_* = \text{SNR}_{\mu_0}^{-1} \in \Xi_L$, so by using Proposition 53-(2), on the event E_0 , for any $\mu_0 \in \mathcal{W}_{\varepsilon, \delta, \vartheta}$,

$$\begin{aligned} & \sqrt{n} z_{\alpha/2}^{-1} \cdot \max_{j \in [n]} \left| |\text{CI}_j(\hat{\eta}^{\text{GCV}})| - \min_{\eta \in \Xi_L} |\text{CI}_j(\eta)| \right| \\ & \lesssim \left| \mathcal{R}_{(\Sigma, \mu_0)}^{\text{pred}}(\hat{\eta}^{\text{GCV}}) - \mathcal{R}_{(\Sigma, \mu_0)}^{\text{pred}}(\eta_*) \right| + \varepsilon \lesssim |\hat{\eta}^{\text{GCV}} - \eta_*|^2 + \varepsilon. \end{aligned}$$

The above reasoning also proves that on the same event E_0 , for any $\mu_0 \in \mathcal{W}_{\varepsilon, \delta, \vartheta}$,

$$|\gamma_{\hat{\eta}^{\text{GCV}}, *} - \gamma_{\eta_*, *}| \lesssim |\hat{\eta}^{\text{GCV}} - \eta_*|^2 + \varepsilon.$$

From here, in view of (164), by adjusting constants, on an event E_1 with $\mathbb{P}(E_1^c) \leq C_1 n^{-1/7}$, it holds that

$$\sqrt{n} z_{\alpha/2}^{-1} \cdot \max_{j \in [n]} \left| |\text{CI}_j(\hat{\eta}^{\text{GCV}})| - \min_{\eta \in \Xi_L} |\text{CI}_j(\eta)| \right| \vee |\gamma_{\hat{\eta}^{\text{GCV}}, *} - \gamma_{\eta_*, *}| \vee |\hat{\eta}^{\text{GCV}} - \eta_*|^2 \leq \varepsilon. \quad (171)$$

This proves the claim for the length of the CI.

(2). Next we prove the statement for the coverage. We note that a similar Lipschitz continuity argument as in the proof of Lemma 55 shows that for any 1-Lipschitz $\mathbf{g} : \mathbb{R}^n \rightarrow \mathbb{R}$, on an event $E_2(\mathbf{g})$ with $\mathbb{P}(E_2(\mathbf{g})^c) \leq Cn^{-1/7}$,

$$\sup_{\eta \in \Xi_L} |\mathbf{g}(\widehat{\mu}_\eta^{\text{dR}}) - \mathbb{E} \mathbf{g}(\mu_0 + \gamma_{\eta_*} \Sigma^{-1/2} g / \sqrt{n})| \leq \varepsilon. \quad (172)$$

On the other hand, using the Lipschitz continuity of $\eta \mapsto \tau_{\eta_*}$ in Proposition 23-(3),

$$\begin{aligned} |\mathbf{g}(\widehat{\mu}_{\widehat{\eta}^{\text{GCV}}}^{\text{dR}}) - \mathbf{g}(\widehat{\mu}_{\eta_*}^{\text{dR}})| &\leq \|\widehat{\mu}_{\widehat{\eta}^{\text{GCV}}}^{\text{dR}} - \widehat{\mu}_{\eta_*}^{\text{dR}}\| \lesssim |\tau_{\widehat{\eta}^{\text{GCV}}} - \tau_{\eta_*}| \sup_{\eta \in \Xi_L} \|\widehat{\mu}_\eta\| + \|\widehat{\mu}_{\widehat{\eta}^{\text{GCV}}} - \widehat{\mu}_{\eta_*}\| \\ &\lesssim \|\widehat{\eta}^{\text{GCV}} - \eta_*\| \sup_{\eta \in \Xi_L} \|\widehat{\mu}_\eta\| + \|\widehat{\mu}_{\widehat{\eta}^{\text{GCV}}} - \widehat{\mu}_{\eta_*}\|. \end{aligned}$$

So by enlarging C_1 if necessary, we may assume without loss of generality that on $E_1 \cap E_2(\mathbf{g})$,

$$|\mathbf{g}(\widehat{\mu}_{\widehat{\eta}^{\text{GCV}}}^{\text{dR}}) - \mathbf{g}(\widehat{\mu}_{\eta_*}^{\text{dR}})| \leq C_1 \varepsilon^{1/2}. \quad (173)$$

Now we shall make a good choice of \mathbf{g} in (171). Let $\Delta \in (0, 1)$ and $\mathbf{g}_{0,\Delta} : \mathbb{R} \rightarrow [0, 1]$ be a function such that $\mathbf{g}_{0,\Delta} = 1$ on $[-1, 1]$, $\mathbf{g}_{0,\Delta} = 0$ on $\mathbb{R} \setminus (-1 - \Delta, 1 + \Delta)$, and linearly interpolated in $(-1 - \Delta, -1) \cup (1, 1 + \Delta)$. Let

$$\mathbf{g}(u) \equiv \frac{\Delta}{n} \sum_{j=1}^n \mathbf{g}_{0,\Delta} \left(\frac{u_j - \mu_{0,j}}{(\gamma_{\eta_*} + \varepsilon)(\Sigma^{-1})_{jj}^{1/2} z_{\alpha/2} / \sqrt{n}} \right). \quad (174)$$

It is easy to verify the Lipschitz property of \mathbf{g} : for any $u_1, u_2 \in \mathbb{R}^n$, $|\mathbf{g}(u_1) - \mathbf{g}(u_2)| \lesssim n^{-1/2} \Delta \|\mathbf{g}_{0,\Delta}\|_{\text{Lip}} \sum_{j=1}^n |u_{1,j} - u_{2,j}| \lesssim \|u_1 - u_2\|$. Consequently, we may apply (172) with \mathbf{g} defined in (174) to obtain that on the event $E_1 \cap E_2(\mathbf{g})$,

$$\begin{aligned} \mathcal{E}^{\text{dR}}(\widehat{\eta}^{\text{GCV}}) &= \frac{1}{n} \sum_{j=1}^n \mathbf{1} \left(\widehat{\mu}_{\widehat{\eta}^{\text{GCV},j}}^{\text{dR}} \in \left[\mu_{0,j} \pm \widehat{\eta}^{\text{GCV}} (\Sigma^{-1})_{jj}^{1/2} \frac{z_{\alpha/2}}{\sqrt{n}} \right] \right) \\ &\leq \frac{1}{n} \sum_{j=1}^n \mathbf{1} \left(\widehat{\mu}_{\widehat{\eta}^{\text{GCV},j}}^{\text{dR}} \in \left[\mu_{0,j} \pm (\gamma_{\eta_*} + \varepsilon) (\Sigma^{-1})_{jj}^{1/2} \frac{z_{\alpha/2}}{\sqrt{n}} \right] \right) \\ &\leq \Delta^{-1} \cdot \mathbf{g}(\widehat{\mu}_{\widehat{\eta}^{\text{GCV}}}^{\text{dR}}) \quad (\text{using } \mathbf{1}_{[-1,1]} \leq \mathbf{g}_{0,\Delta}) \\ &\leq \Delta^{-1} \cdot \mathbf{g}(\widehat{\mu}_{\eta_*}^{\text{dR}}) + \mathcal{O}(\varepsilon^{1/2}/\Delta) \quad (\text{by (173)}) \\ &\leq \Delta^{-1} \cdot \mathbb{E} \mathbf{g}(\mu_0 + \gamma_{\eta_*} \Sigma^{-1/2} g / \sqrt{n}) + \mathcal{O}(\varepsilon^{1/2}/\Delta). \end{aligned} \quad (175)$$

Now using $\mathbf{g}_{0,\Delta} \leq \mathbf{1}_{[-1-\Delta, 1+\Delta]}$ and the anti-concentration of the standard normal random variable, we may compute

$$\begin{aligned} \Delta^{-1} \cdot \mathbb{E} \mathbf{g}(\mu_0 + \gamma_{\eta_*} \Sigma^{-1/2} g / \sqrt{n}) &= \mathbb{E} \mathbf{g}_{0,\Delta} \left(\frac{\gamma_{\eta_*}}{\gamma_{\eta_*} + \varepsilon} \cdot \frac{g}{z_{\alpha/2}} \right) \\ &\leq \mathbb{P} \left(\mathcal{N}(0, 1) \in \left[\pm z_{\alpha/2} \cdot (1 + \varepsilon/\gamma_{\eta_*}) \cdot (1 + \Delta) \right] \right) \leq 1 - \alpha + \mathcal{O}(\varepsilon + \Delta). \end{aligned} \quad (176)$$

Combining the above two displays (175)-(176), on the event $E_1 \cap E_2(\mathbf{g})$,

$$\mathcal{E}^{\text{dR}}(\widehat{\eta}^{\text{GCV}}) \leq 1 - \alpha + \mathcal{O}(\varepsilon + \Delta + \varepsilon^{1/2}/\Delta).$$

Finally choosing $\Delta = \varepsilon^{1/4}$ to conclude the upper control. The lower control can be proved similarly so we omit the details. \square

Appendix H. Auxiliary results

Proposition 59. *Let $H : \mathbb{R}^n \rightarrow \mathbb{R}_{\geq 0}$ be a non-negative, differentiable function. Suppose there exists some deterministic $\Gamma > 0$ such that $\|\nabla H(g)\|^2 \leq \Gamma^2 H(g)$ almost surely for $g \sim \mathcal{N}(0, I_n)$. Then there exists some universal constant $C > 0$ such that for all $t \geq 0$,*

$$\mathbb{P}\left(|H(g) - \mathbb{E}H(g)|/C \geq \Gamma \mathbb{E}^{1/2} H(g) \cdot \sqrt{t} + \Gamma^2 \cdot t\right) \leq Ce^{-t/C}.$$

Proof The method of proof via the Gaussian log-Sobolev inequality and the Herbst's argument is well known. We give some details for the convenience of the reader. Let $Z \equiv H(g) - \mathbb{E}H(g)$ be the centered version of H , and $G(g) \equiv \lambda Z = \lambda(H(g) - \mathbb{E}H(g))$. Then $\|\nabla G(g)\|^2 = \lambda^2 \|\nabla H(g)\|^2 \leq \lambda^2 \Gamma^2 \cdot H(g) = \lambda^2 \Gamma^2 \cdot (Z + \mathbb{E}H(g))$. By the Gaussian log-Sobolev inequality (see e.g., (Boucheron et al., 2013, Theorem 5.4), or (Giné and Nickl, 2016, Theorem 2.5.6)),

$$\text{Ent}(e^{\lambda Z}) = \mathbb{E}[\lambda Z e^{\lambda Z}] - \mathbb{E} e^{\lambda Z} \log \mathbb{E} e^{\lambda Z} \leq \frac{1}{2} \mathbb{E} [\lambda^2 \Gamma^2 (Z + \mathbb{E}H(g)) e^{\lambda Z}].$$

With $m_Z(\lambda) \equiv \mathbb{E} e^{\lambda Z}$ denoting the moment generation function of Z , the above inequality is equivalent to

$$\lambda m'_Z(\lambda) - m_Z(\lambda) \log m_Z(\lambda) \leq \frac{\Gamma^2 \lambda^2}{2} \left(m'_Z(\lambda) + \mathbb{E}H(g) \cdot m_Z(\lambda) \right).$$

Now dividing $\lambda^2 m_\lambda(Z)$ on both sides of the above display, we have $(\log m_Z(\lambda)/\lambda)' \leq \frac{\Gamma^2}{2} (\log m_Z(\lambda) + \lambda \mathbb{E}H(g))'$. Integrating both sides with the condition $\lim_{\lambda \downarrow 0} (\log m_Z(\lambda)/\lambda) = 0$ and $\log m_Z(\lambda) = 0$, we arrive at $\log m_Z(\lambda) \leq \frac{\Gamma^2}{2} (\lambda \log m_Z(\lambda) + \lambda^2 \mathbb{E}H(g))$. Solving for $\log m_Z(\lambda)$ and using the standard method to convert to tail bound yield the claimed inequality. \blacksquare

Lemma 60. *Let $\Sigma \in \mathbb{R}^{n \times n}$ be an invertible covariance matrix with $\|\Sigma\|_{\text{op}} \vee \|\Sigma^{-1}\|_{\text{op}} \leq K$ for some $K > 0$. Then for any $q \in [1, \infty)$, there exists some $C = C(K, q) > 0$ such that*

$$\left| \frac{\mathbb{E}\|\mathcal{N}(0, \Sigma)\|_q}{\|\text{diag}(\Sigma)\|_{q/2}^{1/2} M_q} - 1 \right| \leq C n^{-\frac{1}{q\sqrt{2}}} \sqrt{\log n}.$$

where $M_q \equiv \mathbb{E}^{1/q} |\mathcal{N}(0, 1)|^q = 2^{1/2} \{\Gamma((q+1)/2)/\sqrt{\pi}\}^{1/q}$.

Proof Let $g \sim \mathcal{N}(0, I_n)$. We first prove that for some $C_0 > 1$,

$$n^{\frac{1}{q\sqrt{2}}}/C_0 \leq \mathbb{E}\|\Sigma^{1/2}g\|_q \leq C_0 n^{\frac{1}{q}}. \quad (177)$$

The upper bound in the above display is trivial. For the lower bound, using $\|x\| \leq n^{\frac{1}{2} - \frac{1}{q\sqrt{2}}} \|x\|_q$, we find $\mathbb{E}\|\Sigma^{1/2}g\|_q \geq n^{-\frac{1}{2} + \frac{1}{q\sqrt{2}}} \mathbb{E}\|\Sigma^{1/2}g\| \gtrsim n^{\frac{1}{q\sqrt{2}}}$. This proves (177).

As $\|x\|_q \leq n^{-\frac{1}{2} + \frac{1}{q\sqrt{2}}} \|x\|$, the map $g \mapsto \|\Sigma^{1/2}g\|_q$ is $\|\Sigma\|_{\text{op}}^{1/2} n^{-\frac{1}{2} + \frac{1}{q\sqrt{2}}}$ -Lipschitz with respect to $\|\cdot\|$. So by Gaussian concentration, for any $t \geq 0$,

$$\mathbb{P}\left(E(t)^c \equiv \left\{ n^{\frac{1}{2} - \frac{1}{q\sqrt{2}}} \left| \|\Sigma^{1/2}g\|_q - \mathbb{E}\|\Sigma^{1/2}g\|_q \right| \geq C\sqrt{t} \right\}\right) \leq Ce^{-t/C}.$$

Consequently, using the above concentration and (177),

$$\begin{aligned} \mathbb{E}\|\Sigma^{1/2}g\|_q^q &\leq \mathbb{E}\|\Sigma^{1/2}g\|_q^q \mathbf{1}_{E(t)} + \mathbb{E}^{1/2}\|\Sigma^{1/2}g\|_q^{2q} \cdot \mathbb{P}^{1/2}(E(t)^c) \\ &\leq (\mathbb{E}\|\Sigma^{1/2}g\|_q + C\sqrt{t})^q + C \cdot n^{1/q} \mathbb{P}^{1/2}(E(t)^c) \\ &\leq (\mathbb{E}\|\Sigma^{1/2}g\|_q)^q \cdot \left\{ (1 + Cn^{-\frac{1}{q\sqrt{2}}}\sqrt{t})^q + C \cdot n^{\frac{1}{q} - \frac{1}{q\sqrt{2}}} \mathbb{P}^{1/2}(E(t)^c) \right\}. \end{aligned}$$

By choosing $t = C_1 \log n$ for some sufficiently large $C_1 > 0$, we have

$$\frac{\mathbb{E}\|\mathcal{N}(0, \Sigma)\|_q}{\|\text{diag}(\Sigma)\|_{q/2}^{1/2} M_q} = \frac{\mathbb{E}\|\Sigma^{1/2}g\|_q}{\mathbb{E}^{1/q}\|\Sigma^{1/2}g\|_q^q} \geq (1 - Cn^{-\frac{1}{q\sqrt{2}}}\sqrt{\log n})_+.$$

The upper bound follows similarly. ■

Lemma 61. *Let $Z \in \mathbb{R}^{m \times n}$ be a random matrix with independent, mean-zero, unit variance, uniformly sub-gaussian components. Suppose the coordinates of ξ are i.i.d. mean zero and uniformly subgaussian with variance $\sigma_\xi^2 > 0$, and are independent of Z . Then there exists some universal constant $C > 0$ such that for any $b \in \mathbb{R}^n$ and $0 < \varrho \leq 1$, with probability at least $1 - Ce^{-m^\varrho/C}$,*

$$|m^{-1}\|Zb + \xi\|^2 - (\|b\|^2 + \sigma_\xi^2)| \leq C \cdot (\sigma_\xi^2 \vee \|b\|^2) \cdot m^{-(1-\varrho)/2}.$$

Proof Let $Z_1, \dots, Z_m \in \mathbb{R}^n$ be the rows of Z . Then

$$\frac{1}{m}\|Zb + \xi\|^2 = \|b\|^2 \frac{1}{m} \sum_{i=1}^m \left\langle Z_i, \frac{b}{\|b\|} \right\rangle^2 + \frac{2\sigma_\xi\|b\|}{m} \sum_{i=1}^n \frac{\xi_i}{\sigma_\xi} \left\langle Z_i, \frac{b}{\|b\|} \right\rangle + \sigma_\xi^2 \frac{\|\xi/\sigma_\xi\|^2}{m}.$$

Using standard concentration estimates, with probability at least $1 - Ce^{-m^\varrho/C}$,

- $\left| \|b\|^2 \frac{1}{m} \sum_{i=1}^m \left\langle Z_i, \frac{b}{\|b\|} \right\rangle^2 - \|b\|^2 \right| \leq C\|b\|^2 \cdot m^{-(1-\varrho)/2},$
- $\left| \frac{2\sigma_\xi\|b\|}{m} \sum_{i=1}^n \frac{\xi_i}{\sigma_\xi} \left\langle Z_i, \frac{b}{\|b\|} \right\rangle \right| \leq C\sigma_\xi\|b\| \cdot m^{-(1-\varrho)/2},$
- $\left| \sigma_\xi^2 \frac{\|\xi/\sigma_\xi\|^2}{m} - \sigma_\xi^2 \right| \leq C\sigma_\xi^2 \cdot m^{-(1-\varrho)/2}.$

Collecting the bounds to conclude. ■

Appendix I. Simulation details for Figure 1 and additional simulations

I.1 Common numerical settings

We set $\Sigma = 1.99 \cdot I_n + 0.01 \cdot \mathbf{1}_n \mathbf{1}_n^\top$, with $\mathbf{1}_n$ representing an n -dimensional all one vector. The random design matrix Z and the error ξ are both generated by t -distribution with 10 degrees of freedom, scaled by $\sqrt{0.8}$. This scaling choice ensures that Z_{ij} and ξ_i have mean zero and variance one. The concrete choice of the signal dimension n , the sample size m , and μ_0 will be specified later.

I.2 Simulation details for Figure 1

We investigate the efficacy of two cross validation schemes in Section 4, namely $\hat{\eta}^{\text{GCV}}$ in (22) and $\hat{\eta}^{\text{CV}}$ in (24). We keep the sample size fixed at $m = 500$, and allow the signal dimension n to vary so that the aspect ratio $\phi = m/n$ ranges from $[0.5, 1.5]$. To facilitate the tuning process, we employ 31 equidistant η 's within the range of $[0, 1.5]$. Moreover, the k -fold cross validation scheme $\hat{\eta}^{\text{CV}}$ is carried out with the default choice $k = 5$.

To empirically verify Theorem 10 and 12, we report in the left panel of Figure 1 the empirical risks $R_{(\Sigma, \mu_0)}^{\#}(\hat{\eta}^{\text{GCV}}), R_{(\Sigma, \mu_0)}^{\#}(\hat{\eta}^{\text{CV}})$ for all $\# \in \{\text{pred, est, in}\}$. All the empirical risk curves are found to concentrate around their theoretical optimal counterparts $\mathcal{R}_{(\Sigma, \mu_0)}^{\#}(\eta_*)$. We note again that as $\hat{\eta}^{\text{GCV}}$ and $\hat{\eta}^{\text{CV}}$ are designed to tune the prediction risk, it is not surprising that $R_{(\Sigma, \mu_0)}^{\text{pred}}(\hat{\eta}^{\text{GCV}}), R_{(\Sigma, \mu_0)}^{\text{pred}}(\hat{\eta}^{\text{CV}})$ concentrate around $\mathcal{R}_{(\Sigma, \mu_0)}^{\text{pred}}(\eta_*)$. The major surprise appears to be that $\hat{\eta}^{\text{GCV}}$ and $\hat{\eta}^{\text{CV}}$ also provide optimal tuning for estimation and in-sample risks, both theoretically validated in our Theorems 10 and 12 and empirically confirmed here.

To empirically verify Theorem 13, we report in the middle and right panels of Figure 1 the averaged coverage and length for the 95%-debiased Ridge CI's with cross-validation, namely $\{\text{CI}_j(\hat{\eta}^{\#})\}$ for $\# \in \{\text{GCV, CV}\}$, and with oracle tuning $\eta_* = \text{SNR}_{\mu_0}^{-1}$. For the middle panel, we observe that adaptive tuning via $\hat{\eta}^{\text{GCV}}$ and $\hat{\eta}^{\text{CV}}$ both provide approximate nominal coverage for a moderate sample size m and signal dimension n . For the right panel, as the lengths of $\{\text{CI}_j(\hat{\eta}^{\#})\}$ are solely determined by $\hat{\gamma}_{\hat{\eta}^{\#}}$, we report here only the length of $\text{CI}_1(\hat{\eta}^{\#})$. We observe that the CI length for both $\text{CI}_1(\hat{\eta}^{\text{GCV}}), \text{CI}_1(\hat{\eta}^{\text{CV}})$ are also in excellent agreement to the oracle length across different aspect ratios.

I.3 Validation of (19)

We next verify the optimal oracle regularization rule in (19) (see Theorem 54 for a rigorous formulation) by simulation. We use $m = 100, n = 200$, and a unit vector μ_0 chosen randomly (and then fixed) from the sphere $\partial B_n(1)$. For this setting, we plot both the theoretical risk curve $\eta \mapsto \bar{R}_{(\Sigma, \mu_0)}^{\#}(\eta)$ and the empirical risk curve $\eta \mapsto R_{(\Sigma, \mu_0)}^{\#}(\eta)$ for all $\# \in \{\text{pred, est, in}\}$. The left panel of Figure 2 reports the noisy case with noise level $\sigma_{\xi}^2 = 1$ and $\text{SNR}_{\mu_0}^{-1} = 1$, while the middle panel reports the noiseless case $\sigma_{\xi}^2 = 0$ with $\text{SNR}_{\mu_0}^{-1} = 0$. These plots show excellent agreement with (19) in that the global minimum of both the theoretical and empirical risk curves is attained roughly at $\eta_* = \text{SNR}_{\mu_0}^{-1}$.

In order to demonstrate the validity of the above phenomenon for ‘most’ μ_0 's, as claimed in Theorem 54, we uniformly generate 500 different μ_0 's over $\partial B_n(1)$. For each μ_0 , we discretize $\eta \in [0, 1.5]$ into 160 grid points and select the empirical optimal value $\eta^{\#}$ by minimizing the empirical prediction, estimation, and in-sample risks. The difference between the empirical optimal $\eta^{\#}$ and the theoretically optimal tuning η_* is depicted in the right panel of Figure 2 through a boxplot of $\eta^{\#} - \eta_*$. It is easily seen that, for all three risks, these differences are highly concentrated around 0.

We finally explain how the theoretical risk curves $\eta \mapsto \bar{R}_{(\Sigma, \mu_0)}^{\#}(\eta)$ are computed in practice. For each fixed η , we solve the fixed-point system (13) for $(\gamma_{\eta, *}, \tau_{\eta, *})$ as follows. The second equation in (13) involves only the scalar variable τ ; under our assumptions, its

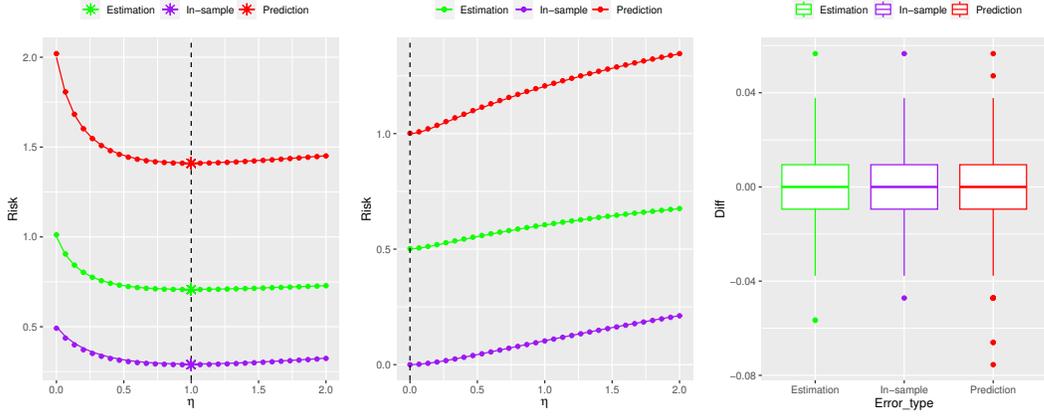


Figure 2: Validation of (19) (see also Theorem 54 for a rigorous formulation). The theoretical risks $\bar{R}_{(\Sigma, \mu_0)}^\#(\eta)$ are computed by solving (13), and the empirical risks $R_{(\Sigma, \mu_0)}^\#(\eta)$ are computed via Monte Carlo simulation over 200 repetitions. *Left panel*: noisy case with minimal empirical risks attained at $\eta_* = \text{SNR}_{\mu_0}^{-1} = 1$ (marked with *). *Middle panel*: noiseless case with all risks minimized at the interpolation regime $\eta_* = \text{SNR}_{\mu_0}^{-1} = 0$. *Right panel*: differences between the global minimizer of the empirical risk curves and the oracle η_* are concentrated around 0 over 500 different μ_0 's.

right-hand side is monotone in τ , so the solution $\tau_{\eta, *}$ is unique. We therefore solve this one-dimensional fixed-point equation for $\tau_{\eta, *}$ by a standard bisection method on a prescribed interval, up to a given numerical tolerance. Once $\tau_{\eta, *}$ is obtained, we plug it into the first equation in (13) to compute $\gamma_{\eta, *}$. The resulting pair $(\gamma_{\eta, *}, \tau_{\eta, *})$ is then substituted into the closed-form expressions for $\bar{R}_{(\Sigma, \mu_0)}^\#(\eta)$.