

Dimension free ridge regression

Chen Cheng* Andrea Montanari†

October 16, 2022

Abstract

Random matrix theory has become a widely useful tool in high-dimensional statistics and theoretical machine learning. However, random matrix theory is largely focused on the proportional asymptotics in which the number of columns grows proportionally to the number of rows of the data matrix. This is not always the most natural setting in statistics where columns correspond to covariates and rows to samples.

With the objective to move beyond the proportional asymptotics, we revisit ridge regression (ℓ_2 -penalized least squares) on i.i.d. data (\mathbf{x}_i, y_i) , $i \leq n$, where \mathbf{x}_i is a feature vector and $y_i = \langle \boldsymbol{\beta}, \mathbf{x}_i \rangle + \varepsilon_i \in \mathbb{R}$ is a response. We allow the feature vector to be high-dimensional, or even infinite-dimensional, in which case it belongs to a separable Hilbert space, and assume either $\mathbf{z}_i := \boldsymbol{\Sigma}^{-1/2} \mathbf{x}_i$ to have i.i.d. entries, or to satisfy a certain convex concentration property.

Within this setting, we establish non-asymptotic bounds that approximate the bias and variance of ridge regression in terms of the bias and variance of an ‘equivalent’ sequence model (a regression model with diagonal design matrix). The approximation is up to multiplicative factors bounded by $(1 \pm \Delta)$ for some explicitly small Δ .

Previously, such an approximation result was known only in the proportional regime and only up to additive errors: in particular, it did not allow to characterize the behavior of the excess risk when this converges to 0. Our general theory recovers earlier results in the proportional regime (with better error rates). As a new application, we obtain a completely explicit and sharp characterization of ridge regression for Hilbert covariates with regularly varying spectrum. Finally, we analyze the overparametrized near-interpolation setting and obtain sharp ‘benign overfitting’ guarantees.

Contents

1	Introduction	2
2	Setting and characterization	4
3	Statement of main results	8
4	Applications	12
4.1	Proportional regime	12
4.2	Bounded varying spectrum	13
5	Numerical illustrations	16

*Department of Statistics, Stanford University

†Department of Electrical Engineering and Department of Statistics, Stanford University; School of Mathematics, Institute for Advanced Studies, Princeton

6	Proof of Theorem 1	17
7	Proof of Theorem 5	27
A	Proof of Proposition 2.2	36
B	Auxiliary lemmas	36
	B.1 Proof of Lemma 6.1	36
	B.2 Proof of Lemma 6.2	37
	B.3 Proof of Lemma 6.3	38
	B.4 Proof of Lemma 6.4	39
	B.5 Proof of Lemma 6.6	42
C	Proofs for Theorem 5	43
	C.1 Proof of Lemma 7.1	43
	C.2 Proof of Lemma 7.2	43
	C.3 Proof of Lemma 7.3	48
	C.4 Proof of Lemma 7.4	50
	C.5 Proof of Lemma 7.5	56
	C.6 Proof of Corollary 6.5	58
D	Proof of Theorem 2	60
E	Proof of Theorem 3	66
F	Proofs for proportional regime	68
	F.1 Proof of Proposition 4.1	68
	F.2 Proof of Proposition 4.2	69
G	Proofs for bounded varying spectrum regime	71
	G.1 Proof of Proposition 4.3	71
	G.2 Proof of Proposition 4.4	72
H	Proof of Theorem 4	74

1 Introduction

In regression modeling, we typically assume to be given data (\mathbf{x}_i, y_i) , $i \leq n$ that are i.i.d. samples from a common distribution \mathbb{P} , with \mathbf{x}_i a feature vector, and $y_i \in \mathbb{R}$ a scalar response. We would like to estimate a model $f : \mathbf{x} \mapsto f(\mathbf{x})$ to predict y_{new} from \mathbf{x}_{new} , where $(\mathbf{x}_{\text{new}}, y_{\text{new}}) \sim \mathbb{P}$ is a new sample from the same distribution. In this paper, we will focus on linear models whereby $f(\mathbf{x}) = \langle \hat{\boldsymbol{\beta}}, \mathbf{x} \rangle$, and use ridge regression for the estimator $\hat{\boldsymbol{\beta}}$. Denoting by \mathbf{X} the matrix with rows $\mathbf{x}_1, \dots, \mathbf{x}_n$, we have

$$\hat{\boldsymbol{\beta}}_\lambda := \arg \min_{\mathbf{b}} \{ \|\mathbf{y} - \mathbf{X}\mathbf{b}\|^2 + \lambda \|\mathbf{b}\|^2 \} \quad (1)$$

$$= (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^\top \mathbf{y}. \quad (2)$$

We will also be interested in the $\lambda \rightarrow 0+$ limit of this estimator which (in the overparametrized case) corresponds to the minimum norm interpolator of the data, and refer to it as ‘ridgeless regression.’ We will denote by $\beta := \arg \min_{\mathbf{b}} \mathbb{E}\{(y - \mathbf{b}^\top \mathbf{x})^2\}$ the population regressor.

Statistical theory studies this and similar estimators in three different regimes:

1. The classical low-dimensional setting in which $\mathbf{x}_i, \beta \in \mathbb{R}^d$ with d fixed and $n \rightarrow \infty$. In this regime, the empirical covariance $\mathbf{X}^\top \mathbf{X}/n$ converges to the population covariance $\Sigma := \mathbb{E}\{\mathbf{x}_1 \mathbf{x}_1^\top\}$ (provided the latter exists) and $\hat{\beta}_\lambda$ is asymptotically normal [VdV00].
2. The (by now) classical high-dimensional regime in which $\mathbf{x}_i, \beta \in \mathbb{R}^d$ with $d \gg n$ but: (i) the population covariance Σ is well conditioned, and (ii) the population regressor β is sparse. In this case it is advised to replace the ℓ_2 penalty $\|\mathbf{b}\|^2$ by a sparsity promoting penalty, e.g. $\|\mathbf{b}\|_1$ [Tib96, DET05]. In many ways, this regime is similar to the previous one, provided $n \gg s \log d$. While $\mathbf{X}^\top \mathbf{X}/n$ does not concentrate, its restrictions to subsets of $O(s)$ coordinates do [CT05].
3. The proportional regime in which $n \asymp d$. In this case $\mathbf{X}^\top \mathbf{X}/n$ does not concentrate, and $\hat{\beta}_\lambda$ is not consistent, and indeed consistent estimation is generally impossible. However, accurate characterizations of the ridge estimator and its risk can be derived using random matrix theory [Dic16, DW18, HMRT22]. Such characterizations answers the question of ε -consistency: for what sample size, and what data distributions does the ridge estimator achieves error $\mathbb{E}\{\|\hat{\beta}_\lambda - \beta\|^2\} \leq \varepsilon$? Similar characterizations hold for other estimators such as the Lasso [BM11, MM21, CMW20], robust M-estimators [BBEKY13, EKBB⁺13, EK18, DM16], and so on [BKM⁺19, TAH18, TPT21, CM22].

Despite the wealth of fascinating technical results in this area, this state of affairs leaves open many important questions.

First, it would be important have a unified theoretical framework that does not require the statistician to decide which asymptotics to use. For instance, in order to apply sharp asymptotics in the classical or proportional regimes, it is often assumed that a given pair (n, d) is in fact an element of a sequence $(n, d(n))$ with, respectively, either $d(n) \asymp 1$, or $d(n) \asymp n$.

In practice we are given a single pair, say $(n, d) = (1000, 50)$: should we interpret this as $d \asymp 1$, $d \asymp n$, or yet another regime that is not covered by current theory (e.g., $d \asymp n^{2/3}$)?

In fact, the distinction between three types of asymptotics outlined above is rather the consequence of the technical tools used to derive them, rather than a fundamental statistical phenomenon.

Second, the restriction $d = O(n)$ (or $s = O(n)$ in sparse regression) which is implied both by the proportional and by the classical asymptotics is artificial. While this condition might seem necessary for consistency at first sight (it might seem that at least d observations are required to estimate d parameters), as shown in [BLLT20, TB20] this is in fact not the case. Further, it is not even clear how to check in practice $d = O(n)$ for a given pair n, d .

Third, it would be important to remove the assumption of a well conditioned Σ , and derive precise asymptotics for general covariances. We would argue that the ill-conditioned case is most important in practice, since high-dimensional data have often low-dimensional structures.

Fourth, the proportional asymptotics is somewhat un-natural from a statistical viewpoint. Most statisticians are used to think of the data distribution is fixed (in particular, d is fixed), while we sample size n increases. In a standard proportional setting, one instead assumes $n, d \rightarrow \infty$ together with $n/d \rightarrow \delta$: the data distribution changes with the sample size.

Recent progress on several of these issues was achieved in the context of ridge regression. Among others, [HMRT22] derived a characterization for bias and variance in the proportional regime that

is *non-asymptotic*, i.e. holds up to an approximation error that is explicit and vanishes for large n, d . Using a different approach, [BLLT20, TB20] obtained bounds on bias and variance that hold for arbitrary (possibly infinite) dimension d , in terms of the decay of eigenvalue of Σ . These bounds allow to demonstrate ‘benign overfitting,’ i.e. choices of Σ, β (i.e. data distributions) such that minimum norm interpolator is consistent.

The results [HMRT22, BLLT20, TB20] have limitations. The characterization of the risk proved in [HMRT22] has sharp leading constants, but only holds for $C^{-1} \leq n/d \leq C$ with C a constant, and holds up to an additive error. However, this error terms can be larger than the actual excess risk when the latter vanishes. The bounds of [BLLT20, TB20], on the other hand, hold up to unspecified multiplicative constants. The proof techniques in these two sets of results are furthermore very different.

In this paper we attempt to provide a unified picture that covers these gaps, by extending the sharp characterization of ridge regression of [HMRT22] beyond the proportional regime. This will allow to recover the benign overfitting results of [BLLT20, TB20] (in several cases) with sharp constants. In doing so, we will extend random matrix theory analysis to cases with $d \gg n$ or $d = \infty$, without restrictions on the condition number of Σ . In the case $d = \infty$, the feature vectors \mathbf{x}_i are random elements in a separable Hilbert space, whose distribution is fixed (does not change with n), and whose covariance Σ is a trace class self-adjoint operator.

The rest of the paper is organized as follows. The next section describes the setting for our analysis, the main assumptions and the resulting asymptotic characterization. It also provides some intuition and connects our results to earlier work. Section 3 contains the formal statement of our general results, while Section 4 specializes our theorem to regimes of interest and develops tools to check its assumptions. Section 5 evaluates our characterization for certain choices of Σ, β , and compare the predictions with simulations. Finally, proof are presented in Sections 6 and 7, with most technical steps deferred to the appendices.

2 Setting and characterization

Ridge regression in Hilbert space. We consider the simple linear model

$$y_i = \mathbf{x}_i^\top \beta + \varepsilon_i, \quad (3)$$

where $\beta \in \mathbb{R}^d$ is the ground truth signal. The random features $\mathbf{x}_i \in \mathbb{R}^d$ and noise ε_i are independent, and the $(\mathbf{x}_i, \varepsilon_i)$ are i.i.d. samples with $1 \leq i \leq n$. We assume $\mathbf{x}_i, \varepsilon_i$ are mean zero with covariances $\text{Cov}(\mathbf{x}_i) = \Sigma$ and $\text{Var}(\varepsilon_i) = \tau^2$. Defining the data matrix

$$\mathbf{X} = \begin{bmatrix} - & \mathbf{x}_1^\top & - \\ - & \mathbf{x}_2^\top & - \\ & \vdots & \\ - & \mathbf{x}_n^\top & - \end{bmatrix} \in \mathbb{R}^{n \times d},$$

the response vector $\mathbf{y} = (y_1, \dots, y_n)^\top$ and the noise vector $\boldsymbol{\varepsilon} = (\varepsilon_1, \dots, \varepsilon_n)^\top$, we can write in matrix form

$$\mathbf{y} = \mathbf{X}\beta + \boldsymbol{\varepsilon}.$$

In this paper, we assume the dimension $d \in \mathbb{Z}_{\geq 0} \cup \{\infty\}$. When $d < \infty$, we are in the usual setup of linear model with finite dimensional features. In the case $d = \infty$, we assume that the

\mathbf{x}_i 's' are i.i.d. random vectors from a real, separable Hilbert space \mathcal{H} . We will use $\|\mathbf{x}\|$ to denote the norm and $\langle \mathbf{x}_1, \mathbf{x}_2 \rangle$ or $\mathbf{x}_1^\top \mathbf{x}_2$ to denote the scalar product in this space. Given a linear operator $\mathbf{A} : \mathcal{H} \rightarrow \mathcal{H}$, we denote by $\|\mathbf{A}\|$ the associated operator norm.

We will assume the covariance operator $\Sigma = \mathbb{E}[\mathbf{x}\mathbf{x}^\top]$ to be trace-class, namely

$$\text{Tr}(\Sigma) = \mathbb{E}\{\|\mathbf{x}_i\|^2\} < \infty,$$

and, without loss of generality, we also assume $\|\Sigma\| = 1$. Recall that, without loss of generality, one can always assume \mathcal{H} to be $\ell_2 := \{\mathbf{x} = (x_1, x_2, \dots) : \sum_{i=1}^{\infty} x_i^2 < \infty\}$ [Bré11].

For an estimator $\hat{\beta} = \hat{\beta}(\mathbf{X}, \mathbf{y})$ we define the excess risk as

$$\mathcal{R}_{\mathbf{X}}(\hat{\beta}; \beta) = \mathbb{E}_{\mathbf{x}_{\text{new}}, \mathbf{y}} \left[(\mathbf{x}_{\text{new}}^\top \hat{\beta} - \mathbf{x}_{\text{new}}^\top \beta)^2 \mid \mathbf{X} \right] = \mathbb{E}_{\mathbf{y}} \left[\|\hat{\beta} - \beta\|_{\Sigma}^2 \mid \mathbf{X} \right],$$

where \mathbf{x}_{new} is an independent copy of $\mathbf{x}_1, \dots, \mathbf{x}_n$ and $\|\mathbf{x}\|_{\Sigma}^2 := \mathbf{x}^\top \Sigma \mathbf{x}$. We will also refer to this as the ‘test error’ or the ‘generalization error’ (although the latter is actually given by the difference between $\mathcal{R}_{\mathbf{X}}$ and its empirical version.) Let us emphasize that in this definition, $\mathcal{R}_{\mathbf{X}}(\hat{\beta}; \beta)$ is a random quantity because it depends on the data \mathbf{X} : however, as we will prove, it concentrates around a non-random value.

The generalization error admits a variance-bias decomposition $\mathcal{R}_{\mathbf{X}}(\hat{\beta}; \beta) = \mathcal{V}_{\mathbf{X}}(\hat{\beta}; \beta) + \mathcal{B}_{\mathbf{X}}(\hat{\beta}; \beta)$, with

$$\mathcal{V}_{\mathbf{X}}(\hat{\beta}; \beta) = \text{Tr} \left(\Sigma \text{Cov}(\hat{\beta} \mid \mathbf{X}) \right), \quad \mathcal{B}_{\mathbf{X}}(\hat{\beta}; \beta) = \left\| \mathbb{E}_{\mathbf{y}}[\hat{\beta} \mid \mathbf{X}] - \beta \right\|_{\Sigma}^2.$$

For ridge regression, we can write explicit forms of variance and bias:

$$\mathcal{V}_{\mathbf{X}}(\lambda) = \tau^2 \text{Tr} \left(\Sigma \mathbf{X}^\top \mathbf{X} (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I})^{-2} \right), \quad (4a)$$

$$\mathcal{B}_{\mathbf{X}}(\lambda) = \lambda^2 \langle \beta, (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I})^{-1} \Sigma (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I})^{-1} \beta \rangle. \quad (4b)$$

Assumptions on the covariates distribution. We impose the following assumptions on the covariates \mathbf{x}_i throughout the paper.

Assumption 1. We assume $\mathbb{E}[\mathbf{x}_i] = \mathbf{0}$, $\Sigma := \mathbb{E}[\mathbf{x}_i \mathbf{x}_i^\top]$ is a trace class operator: $\text{Tr}(\Sigma) < \infty$ and (without loss of generality) $\|\Sigma\| = 1$. We denote its eigenvalues by $1 = \sigma_1 \geq \sigma_2 \geq \dots$ in non-increasing order. We assume $\|\beta\|_{\Sigma^{-1}} := \|\Sigma^{-1/2} \beta\| < \infty$.

We further assume $\mathbf{x}_i = \Sigma^{1/2} \mathbf{z}_i$ where the following hold.

I. There exist $\mathbf{d}_{\Sigma} := \mathbf{d}_{\Sigma}(n) \geq n$ such that, for all $1 \leq k \leq \min\{n, d\}$

$$\sum_{l=k}^{\mathbf{d}_{\Sigma}} \sigma_l \leq \mathbf{d}_{\Sigma} \sigma_k.$$

II. There exist $\mathbf{C}_{\mathbf{x}} > 0$, such that one of the following condition holds:

- (a) **Independent sub-Gaussian coordinates:** \mathbf{z}_i has independent but not necessarily identically distributed coordinates with uniformly bounded sub-Gaussian norm. Namely: each coordinate z_{ij} of \mathbf{z}_i satisfies $\mathbb{E}[z_{ij}] = 0$, $\text{Var}(z_{ij}) = 1$ and $\|z_{ij}\|_{\psi_2} := \sup_{p \geq 1} p^{-\frac{1}{2}} (\mathbb{E}[|z_{ij}|^p])^{\frac{1}{p}} \leq \mathbf{C}_{\mathbf{x}}$.
- (b) **Convex concentration:** allowing \mathbf{z}_i to have dependent coordinates, the following holds for any 1-Lipschitz convex function $\varphi : \mathbb{R}^d \rightarrow \mathbb{R}$, and for every $t > 0$

$$\mathbb{P}(|\varphi(\mathbf{z}_i) - \mathbb{E}\varphi(\mathbf{z}_i)| \geq t) \leq 2 \exp(-t^2/\mathbf{C}_{\mathbf{x}}^2).$$

The technical motivation for assumption II is to establish concentration of quadratic forms of \mathbf{z}_i , via Hanson-Wright inequality. We notice that the convex concentration property is implied by any of the following. (i) By Talagrand inequality, convex concentration holds for random vectors \mathbf{z}_i with independent bounded entries [BLM13, Theorem 7.12]. (ii) By Herbst’s argument, concentration of Lipschitz functions (and hence in particular convex concentration) holds for random vectors \mathbf{z}_i that satisfy a log-Sobolev inequality [BGL⁺14, Proposition 5.4.1]. (iii) Finally, as a special case of the last point, vectors \mathbf{z}_i with strongly log-concave probability density function satisfy this condition [BGL⁺14, Corollary 5.7.2].

The form of Hanson-Wright inequality that we will use is given below.

Lemma 2.1 (Hanson-Wright inequality [Ada15, RV13]). *Suppose $\mathbf{x} \in \mathbb{R}^d$ is a random copy of the features vector \mathbf{x}_i satisfying Assumption 1. Then there exists a universal constant $c_0 > 0$ such that, for any matrix $\mathbf{M} \in \mathbb{R}^{d \times d}$ with $\|\Sigma^{\frac{1}{2}} \mathbf{M} \Sigma^{\frac{1}{2}}\| < \infty$, we have*

$$\mathbb{P} \left(\left| \mathbf{x}^\top \mathbf{M} \mathbf{x} - \text{Tr}(\Sigma \mathbf{M}) \right| \geq t \right) \leq 2 \exp \left\{ -c_0 \min \left(\frac{t^2}{C_{\mathbf{x}}^4 \|\Sigma^{\frac{1}{2}} \mathbf{M} \Sigma^{\frac{1}{2}}\|_F^2}, \frac{t}{C_{\mathbf{x}}^2 \|\Sigma^{\frac{1}{2}} \mathbf{M} \Sigma^{\frac{1}{2}}\|} \right) \right\}.$$

Remark 2.1. The results of [Ada15, RV13] are stated for finite d . However, the inequality also holds for $d = \infty$ on the Hilbert space ℓ_2 by a standard approximation argument. Namely, one can project the vector \mathbf{x} on the span of the top k -eigenvectors of Σ , establish concentration, and take $k \rightarrow \infty$ at the end.

Effective variance and bias. An important observation of [HMRT22] is that variance $\mathcal{V}_{\mathbf{X}}$ and bias $\mathcal{B}_{\mathbf{X}}$ concentrate around some non-random quantities, that can be interpreted in terms of an ‘effective’ regression problem. While [HMRT22] proves such characterization in the proportional regime $n \asymp d$, here we will extend its validity and prove stronger guarantees.

Define the effective regularization λ_* as the unique non-negative solution of

$$n - \frac{\lambda}{\lambda_*} = \text{Tr}(\Sigma(\Sigma + \lambda_* \mathbf{I})^{-1}), \quad (5)$$

we that define the effective variance and bias as

$$\mathbf{V}_n(\lambda) := \frac{\tau^2 \text{Tr}(\Sigma^2(\Sigma + \lambda_* \mathbf{I})^{-2})}{n - \text{Tr}(\Sigma^2(\Sigma + \lambda_* \mathbf{I})^{-2})}, \quad (6)$$

$$\mathbf{B}_n(\lambda) := \frac{\lambda_*^2 \langle \boldsymbol{\beta}, (\Sigma + \lambda_* \mathbf{I})^{-2} \Sigma \boldsymbol{\beta} \rangle}{1 - n^{-1} \text{Tr}(\Sigma^2(\Sigma + \lambda_* \mathbf{I})^{-2})}, \quad (7)$$

$$\mathbf{R}_n(\lambda) := \mathbf{B}_n(\lambda) + \mathbf{V}_n(\lambda). \quad (8)$$

Our main result —stated in the next section— will establish dimension-free guarantees of the form

$$\mathcal{V}_{\mathbf{X}} = (1 + o_n(1)) \mathbf{V}_n \quad \mathcal{B}_{\mathbf{X}} = (1 + o_n(1)) \mathbf{B}_n \quad (9)$$

where the term These improve over earlier work in two important directions. First, they are *dimension free*, and in particular do not assume $n \asymp d$. Second, they provide *multiplicative approximations*, and hence retain their utility when the risk is small.

Bounds, interpretation, benign overfitting. Before stating our formal results relating $\mathcal{V}_{\mathbf{X}}$ to \mathbf{V}_n and $\mathcal{B}_{\mathbf{X}}$ to \mathbf{B}_n , it is useful to develop some intuition about the expressions (6), (7) and their immediate consequences. Note that, by Eq. (5), we necessarily have

$$\mathrm{Tr}(\boldsymbol{\Sigma}^2(\boldsymbol{\Sigma} + \lambda_{\star}\mathbf{I})^{-2}) < \mathrm{Tr}(\boldsymbol{\Sigma}(\boldsymbol{\Sigma} + \lambda_{\star}\mathbf{I})^{-1}) \leq n. \quad (10)$$

If we assume that inequality between the first and last term holds with a constant multiplicative factor, i.e. $\mathrm{Tr}(\boldsymbol{\Sigma}^2(\boldsymbol{\Sigma} + \lambda_{\star}\mathbf{I})^{-2}) \leq n(1 - c_{\star}^{-1})$ for some constant $c_{\star} \in (0, \infty)$, then we get

$$\mathbf{V}_n(\lambda) \leq \frac{c_{\star}\tau^2}{n} \mathrm{Tr}(\boldsymbol{\Sigma}^2(\boldsymbol{\Sigma} + \lambda_{\star}\mathbf{I})^{-2}), \quad (11)$$

$$\mathbf{B}_n(\lambda) \leq c_{\star}\lambda_{\star}^2 \langle \boldsymbol{\beta}, (\boldsymbol{\Sigma} + \lambda_{\star}\mathbf{I})^{-2} \boldsymbol{\Sigma} \boldsymbol{\beta} \rangle. \quad (12)$$

Comparing these bounds with the bias and variance of general ridge regression in Eqs. (4a), (4a), we observe that the right hand sides are (modulo the factor c_{\star}) the bias and variance of a modified ridge regression in which:

- The design matrix is non-random and given by $\boldsymbol{\Sigma}^{1/2}$ instead of \mathbf{X} .
- The regularization parameter is λ_{\star} instead of λ .
- The noise level is τ/\sqrt{n} instead of τ .

Even more explicit expressions can be obtained by writing the right-hand side of Eqs. (11), (12) in the basis that diagonalizes $\boldsymbol{\Sigma}$ as in the next proposition. A proof of this statement is in Appendix A.

Proposition 2.2. *Assume $\mathrm{Tr}(\boldsymbol{\Sigma}^2(\boldsymbol{\Sigma} + \lambda_{\star}\mathbf{I})^{-2}) \leq n(1 - c_{\star}^{-1})$, for $c_{\star} \in (1, \infty)$. Let $\boldsymbol{\Sigma} := \sum_{i \geq 1} \sigma_i \mathbf{v}_i \mathbf{v}_i^{\top}$ be the eigendecomposition of $\boldsymbol{\Sigma}$, and denote by $\boldsymbol{\beta}_{\leq k} := \sum_{i \leq k} \langle \boldsymbol{\beta}, \mathbf{v}_i \rangle \mathbf{v}_i$ the orthogonal projection of $\boldsymbol{\beta}$ onto the span of $\mathbf{v}_1, \dots, \mathbf{v}_k$, and by $\boldsymbol{\beta}_{> k} := \boldsymbol{\beta} - \boldsymbol{\beta}_{\leq k}$ its complement. Finally, let $k_{\star} := \max\{k : \sigma_k \geq \lambda_{\star}\}$, and define the tail effective rank parameters by*

$$r_q(k) := \sum_{\ell > k} \left(\frac{\sigma_{\ell}}{\sigma_{k+1}} \right)^q, \quad \bar{r}(k) := \frac{r_1(k)^2}{r_2(k)}. \quad (13)$$

Then, defining $b_k := \sigma_k/\sigma_{k+1}$, we have $2n \geq k_{\star} + r_1(k_{\star})/b_{k_{\star}}$ and

$$\mathbf{V}_n(\lambda) \leq c_{\star}\tau^2 \left(\frac{k_{\star}}{n} + \frac{r_2(k_{\star})}{n} \right) \leq c_{\star}\tau^2 \left(\frac{k_{\star}}{n} + \frac{4b_{k_{\star}}^2 n}{\bar{r}(k_{\star})} \right), \quad (14)$$

$$\mathbf{B}_n(\lambda) \leq c_{\star} \left(\sigma_{k_{\star}}^2 \|\boldsymbol{\beta}_{\leq k_{\star}}\|_{\boldsymbol{\Sigma}^{-1}}^2 + \|\boldsymbol{\beta}_{> k_{\star}}\|_{\boldsymbol{\Sigma}}^2 \right). \quad (15)$$

(We notice that if the singular values σ_k do not decay faster than exponentially, then b_k is of order one.) While these are only bounds on the theoretical characterization $\mathbf{B}_n(\lambda)$, $\mathbf{V}_n(\lambda)$ for bias and variance, our main results (Theorem 1 and Theorem 3) will allow to transfer them to the actual bias and variance $\mathcal{B}_{\mathbf{X}}(\lambda)$, $\mathcal{V}_{\mathbf{X}}(\lambda)$ (modulo additional error terms).

Remark 2.2. These bounds (more precisely, the bounds on $\mathcal{B}_{\mathbf{X}}(\lambda)$, $\mathcal{V}_{\mathbf{X}}(\lambda)$ that follow from these and Theorem 1) are closely related to the ones in [BLLT20, TB20], see in particular [TB20, Theorem 1]. It is worth pointing out two important differences. *First*, the bounds in Eqs. (14), (15) are somewhat more precise/explicit: there is no unspecified constant factor¹, no dependence on the condition number of $\sigma_1/\sigma_{k_{\star}}$, and no multiplicative factor depending on the probability. *Second*, Eqs. (14), (15) are only proved for the specific value of k_{\star} defined there.

¹The factor c_{\star} is explicit and, if useful, can be replaced by the original expression.

Remark 2.3. The bounds of Eqs. (14), (15) allow to characterize settings in which the excess test error (as predicted by our theory) vanishes. Indeed, for $V_n(\lambda)$ to vanish, it is sufficient that $k_\star/n \rightarrow 0$ and $\bar{r}(k_\star)/n \rightarrow \infty$. A simple sufficient condition for $B_n(\lambda) \rightarrow 0$ is that $\beta \in \text{span}(\mathbf{v}_1, \dots, \mathbf{v}_k)$ with $\sigma_k/\sigma_{k_\star} \rightarrow \infty$.

We will discuss special examples in Section 4, and show how our general results allow to derive more precise estimates of the risk in those cases.

Equivalent sequence model. The discussion above relies on the assumption $\text{Tr}(\Sigma^2(\Sigma + \lambda_\star \mathbf{I})^{-2}) \leq n(1 - c_\star^{-1})$, which implies the simple bounds (11), (12). However the interpretation in terms of a modified ridge regression problem holds for the exact formulas of Eqs. (6), (7). This interpretation was developed in the context of earlier work on the proportional asymptotics [DJM13], but it is useful to spell it out here for the present context.

In the modified model, we observe \mathbf{y}^s that is related to β according to

$$\mathbf{y}^s = \Sigma^{1/2} \beta + \frac{\omega}{\sqrt{n}} \mathbf{g}, \quad \mathbf{g} \sim \mathcal{N}(0, \mathbf{I}_d), \quad (16)$$

Without loss of generality, we can work in the basis in which Σ is diagonal, and therefore rewrite the above as $y_i^s = \sigma_i^{1/2} \beta_i + (\omega/\sqrt{n}) g_i$, which coincides with the definition of the classical sequence model [Tsy09].

We use ridge regression at regularization level λ_\star as defined in Eq. (5):

$$\hat{\beta}_\lambda^s := \text{argmin}_{\mathbf{b}} \{ \|\mathbf{y}^s - \Sigma^{1/2} \mathbf{b}\|^2 + \lambda_\star \|\mathbf{b}\|^2 \}. \quad (17)$$

Finally, choose the noise level ω to be the unique positive solution of

$$\omega^2 = \tau^2 + \mathbb{E}_{\mathbf{g}} \{ \|\hat{\beta}_\lambda^s - \beta\|_\Sigma^2 \}. \quad (18)$$

Then our theoretical prediction for the excess test error $R_n(\lambda)$ coincides with the excess test error of the sequence model:

$$R_n(\lambda) = \mathbb{E}_{\mathbf{g}} \{ \|\hat{\beta}_\lambda^s - \beta\|_\Sigma^2 \}. \quad (19)$$

Summarizing, the predicted test error for the original model is equal to the test error in the sequence model, albeit at a different value of the ridge regularization parameter and of the noise level. Needless to say, studying the sequence model is significantly simpler than the original model (3).

3 Statement of main results

Big-Oh notation. For two functions $f(\mathbf{x})$ and $g(\mathbf{x})$ (where \mathbf{x} can be a scalar or a vector), we write $f(\mathbf{x}) = \mathcal{O}_\alpha(g(\mathbf{x}))$ if there exists a constant C_α depending only on the value of α (also α can be either a scalar or a vector) such that $|f(\mathbf{x})| \leq C_\alpha |g(\mathbf{x})|$ for all \mathbf{x} . In particular, if the constant is universal we write $f(\mathbf{x}) = \mathcal{O}(g(\mathbf{x}))$. Similarly, we write $f(\mathbf{x}) = \Omega_\alpha(g(\mathbf{x}))$ if $|f(\mathbf{x})| \geq C_\alpha |g(\mathbf{x})|$ for all \mathbf{x} and some constant $C_\alpha > 0$. Finally, we write $f(\mathbf{x}) = \Theta_\alpha(g(\mathbf{x}))$ if we have both $f(\mathbf{x}) = \mathcal{O}_\alpha(g(\mathbf{x}))$ and $f(\mathbf{x}) = \Omega_\alpha(g(\mathbf{x}))$.

We will state three theorems: the first one for ridge regression with positive regularization $\lambda > 0$ (Theorem 1), and the other two for the ridgeless case $\lambda = 0+$ (Theorem 2 for the overparametrized regime, and Theorem 3 for the underparametrized one). Our approximation guarantees will depend on the pair Σ, β through the following three quantities (in the case $\lambda = 0+$, these quantities will be modified as described below):

1. The ratio between effective dimension and regularization parameter:

$$\chi_n(\lambda) := 1 + \frac{\sigma_{\lfloor \eta n \rfloor} \mathbf{d}_{\Sigma} \log^2(\mathbf{d}_{\Sigma})}{\lambda}. \quad (20)$$

Here η is a constant that only depends on $\mathbf{C}_{\mathbf{x}}$, and hence we will leave it implicit.

2. The ratio between regularization and effective regularization

$$\kappa := \min\left(\frac{\lambda}{n\lambda_{\star}}; 1 - \frac{\lambda}{n\lambda_{\star}}\right) > 0. \quad (21)$$

3. For a positive semi-definite operator \mathbf{Q} , define the modified population resolvent:

$$\mathcal{R}_0(\mu_0, \mu; \mathbf{Q}) := \text{Tr}\left(\Sigma^{\frac{1}{2}} \mathbf{Q} \Sigma^{\frac{1}{2}} (\mu_0 \mathbf{I} + \mu \Sigma)^{-1}\right). \quad (22)$$

Letting $\boldsymbol{\beta} = \Sigma^{1/2} \boldsymbol{\theta}$, $\|\boldsymbol{\theta}\| < \infty$, we consider the ratio

$$\rho(\lambda) := \frac{\mathcal{R}_0(\lambda_{\star}, 1; \boldsymbol{\theta} \boldsymbol{\theta}^{\top} / \|\boldsymbol{\theta}\|^2)}{\mathcal{R}_0(\lambda_{\star}, 1; \mathbf{I})} \in (0, 1]. \quad (23)$$

We next present our master theorem for ridge regression: its proof is postponed to Section 6.

Theorem 1 (Ridge regression). *Under Assumption 1, for any positive integers k and D , there exist constants $\eta = \eta(\mathbf{C}_{\mathbf{x}}) \in (0, 1/2)$ and $\mathbf{C} = \mathbf{C}(\mathbf{C}_{\mathbf{x}}, D) > 0$ such that the following hold. Define $\chi_n(\lambda), \kappa, \rho(\lambda)$ as above (with $\eta = \eta(\mathbf{C}_{\mathbf{x}})$ in Eq. (20)).*

If it holds that

$$\chi_n(\lambda)^3 \log^2 n \leq \mathbf{C} n \kappa^{4.5}, \quad n^{-2D+1} = \mathcal{O}\left(\sqrt{\frac{\kappa^3 \log^2 n}{\max\{n, \lambda\}}}\right),$$

then for all $n = \Omega_{k,D}(1)$, with probability $1 - \mathcal{O}_k(n^{-D+1})$ we have:

1. **Variance approximation.**

$$|\mathcal{V}_{\mathbf{X}}(\lambda) - \mathbf{V}_n(\lambda)| = \mathcal{O}_{k, \mathbf{C}_{\mathbf{x}}, D}\left(\frac{\chi_n(\lambda)^3 \log^2 n}{n^{1-\frac{1}{k}} \kappa^{9.5}}\right) \cdot \mathbf{V}_n(\lambda).$$

2. **Bias approximation.** *If we additionally have $\chi_n(\lambda)^3 \log^2 n \leq \mathbf{C} n \kappa^{4.5} \sqrt{\rho(\lambda)}$ and $\lambda k n^{-\frac{1}{k}} \leq n \kappa / 2$, for all $n = \Omega_{k,D}(1)$, we have*

$$|\mathcal{B}_{\mathbf{X}}(\lambda) - \mathbf{B}_n(\lambda)| = \mathcal{O}_{k, \mathbf{C}_{\mathbf{x}}, D}\left(\frac{\lambda_{\star}(\lambda)^{k+1}}{n \kappa^3} + \frac{\chi_n(\lambda)^3 \log^2 n}{\sqrt{\rho(\lambda)} n^{1-\frac{1}{k}} \kappa^{8.5}}\right) \cdot \mathbf{B}_n(\lambda).$$

Remark 3.1. The condition $\|\boldsymbol{\beta}\|_{\Sigma^{-1}} < \infty$ in Assumption 1 amounts to requiring that the coefficients of $\boldsymbol{\beta}$ in the basis of eigenvectors \mathbf{v}_i of Σ decay fast enough. Namely, it is equivalent to $\sum_i \langle \mathbf{v}_i, \boldsymbol{\beta} \rangle^2 / \sigma_i < \infty$. This condition appears to be a proof artifact and it would be interesting to relax it.

As mentioned above, the conditions on the isotropic random vectors \mathbf{z}_i in Assumption 1 are mainly imposed to be able to apply Hanson-Wright inequality (Lemma 2.1). It is an interesting research question to analyze ridge regression for covariates which do not satisfy this inequality.

We next consider the ridgeless limit for in the overparametrized case: recall that $\widehat{\beta}_\lambda$ coincides in this case with the minimum norm interpolator. In this case we need to modify the quantities defined above to measure the quality of our approximation. We begin by noting that Eq. (5) makes perfect sense in the case $\lambda = 0$ and we have $\lim_{\lambda \downarrow 0} \lambda_\star(\lambda) = \lambda_\star(0) > 0$. We then use the following definitions.

1. We replace $\chi_n(\lambda)$ of Eq. (20) by:

$$\chi'_n(\kappa) := 1 + \frac{\sigma_{\lfloor \eta n \rfloor} \mathbf{d}_\Sigma \log^2(\mathbf{d}_\Sigma)}{\kappa n \lambda_\star(0)}, \quad (24)$$

where κ will be introduced in the theorem statement.

2. The quantity $\rho(\lambda)$ defined in Eq. (23) has a well defined limit as $\lambda \downarrow 0$, given by

$$\rho(0) := \frac{\mathcal{R}_0(\lambda_\star(0), 1; \boldsymbol{\theta} \boldsymbol{\theta}^\top / \|\boldsymbol{\theta}\|^2)}{\mathcal{R}_0(\lambda_\star(0), 1; \mathbf{I})} \in (0, 1].$$

3. Finally we define

$$\mathbf{C}_\Sigma := 1 - \frac{1}{n} \text{Tr}(\Sigma^2 (\Sigma + \lambda_\star(0) \mathbf{I})^{-2}) \in (0, 1).$$

Before giving the statement, we introduce a piece of terminology. We say that A happens on the event E with probability at least $1 - \Delta$ if $\mathbb{P}(A^c \text{ and } E) \leq \Delta$ (and, as a consequence, $\mathbb{P}(A) \geq 1 - \Delta - \mathbb{P}(E^c)$).

Theorem 2 (Ridgeless regression in the overparameterized regime). *Suppose Assumption 1 holds with $n < d$. Further assume $\sigma_n > 0$, and let s_{\min} be the minimum nonzero eigenvalue of the sample covariance $\mathbf{X}^\top \mathbf{X} / n$. For any positive integers k and D , there exist constants $\eta = \eta(\mathbf{C}_\mathbf{x}) \in (0, 1/2)$ and $\mathbf{C}_1 = \mathbf{C}_1(\mathbf{C}_\mathbf{x}, D) > 0$, $\mathbf{C}_i = \mathbf{C}_i(k, \mathbf{C}_\mathbf{x}, D) > 0$, $i \in \{2, 3\}$, such that the following hold. Define $\chi'_n(\kappa)$, $\rho(0)$, \mathbf{C}_Σ as above.*

Let $\kappa > 0$ be such that the following hold

$$\kappa \leq \mathbf{C}_\Sigma^2 / 8, \quad \chi'_n(\kappa)^3 \log^2 n \leq \mathbf{C}_1 n \kappa^{4.5}, \quad n^{-2D+1} = \mathcal{O}\left(\sqrt{\frac{\kappa^3 \log^2 n}{\max\{n, \lambda\}}}\right).$$

Then, on the event $\{s_{\min} \geq 8\lambda_\star(0)\kappa\}$, the following hold with probability $1 - \mathcal{O}_k(n^{-D+1})$:

1. **Variance approximation.** *If in addition $\chi'_n(\kappa)^3 \log^2 n \leq \mathbf{C}_2 n^{1-\frac{1}{k}} \kappa^{9.5}$, then*

$$|\mathcal{V}_\mathbf{X}(0) - \mathbf{V}_n(0)| = \mathcal{O}_{k, \mathbf{C}_\mathbf{x}, D}\left(\kappa \cdot \left(\frac{\lambda_\star(0)}{s_{\min}} + \frac{1}{\mathbf{C}_\Sigma^2}\right) + \frac{\chi'_n(\kappa)^3 \log^2 n}{n^{1-\frac{1}{k}} \kappa^{9.5}}\right) \cdot \mathbf{V}_n(0).$$

2. **Bias approximation.** *If in addition $\chi'_n(\kappa)^3 \log^2 n \leq \mathbf{C}_1 n \kappa^{4.5} \sqrt{\rho(0)}$, $\lambda_\star(0) \kappa n^{-\frac{1}{k}} \leq 1/4$ and*

$$\frac{\lambda_\star(0)^{k+1}}{n \kappa^3} + \frac{\chi'_n(\kappa)^3 \log^2 n}{\sqrt{\rho(0)} n^{1-\frac{1}{k}} \kappa^{8.5}} \leq \mathbf{C}_3,$$

then, for $n = \Omega_{k, D}(1)$,

$$|\mathcal{B}_\mathbf{X}(0) - \mathbf{B}_n(0)|$$

$$\begin{aligned}
&= \mathcal{O}_{k, \mathbf{C}_x, D} \left(\frac{\kappa}{\mathbf{C}_\Sigma^2} + \frac{\lambda_\star(0)^{k+1}}{n\kappa^3} + \frac{\chi'_n(\kappa)^3 \log^2 n}{\sqrt{\rho(0)} n^{1-\frac{1}{k}} \kappa^{8.5}} \right) \cdot \mathbf{B}_n(0) \\
&\quad + \min \left\{ \mathcal{O} \left(\frac{\kappa \lambda_\star(0) \|\boldsymbol{\beta}\|^2}{s_{\min}} \right), \mathcal{O}_{\mathbf{C}_x, D}(\kappa^2 \lambda_\star(0)^2 \chi'_n(\kappa)^2) \|\boldsymbol{\theta}_{\leq n}\|^2 + \mathcal{O}_{\mathbf{C}_x, D}(\kappa \lambda_\star(0) \chi'_n(\kappa)) \|\boldsymbol{\beta}_{> n}\|^2 \right\}.
\end{aligned}$$

The proof of this theorem is presented in Appendix D.

Remark 3.2. Our approach to proving Theorem 2 consists in reducing the ridgeless case $\lambda = 0+$ to the case $\lambda > 0$, and appealing to Theorem 1. For instance, when controlling the variance, we will use triangular inequality

$$|\mathcal{V}_{\mathbf{X}}(0) - \mathbf{V}_n(0)| \leq |\mathcal{V}_{\mathbf{X}}(\lambda) - \mathbf{V}_n(\lambda)| + |\mathcal{V}_{\mathbf{X}}(0) - \mathcal{V}_{\mathbf{X}}(\lambda)| + |\mathbf{V}_n(0) - \mathbf{V}_n(\lambda)|.$$

We then use Theorem 1 to bound the first term by a quantity that diverges as $\lambda \downarrow 0$, and the main technical challenge is in bounding the other two terms by a quantity that vanishes faster than any polynomial as $\lambda \downarrow 0$.

Remark 3.3. In Theorem 2 we use the (random) minimum nonzero eigenvalue s_{\min} of the sample covariance $\mathbf{X}^\top \mathbf{X}/n$. To apply the theorem, we need to choose κ such that $\{s_{\min} \geq 8\lambda_\star(0)\kappa\}$ holds with high probability, and therefore we need a lower bound on s_{\min} that holds with high probability. In this paper, we will provide such lower bounds in two cases: (i) proportional regime and (ii) bounded varying spectrum, cf. Section 4. In proportional regime, $s_{\min} = \Omega_{|d/n-1|^{-1}}(1)$ by Bai-Yin law; and for bounded varying spectrum, $s_{\min} = \Omega(\sigma_n)$ (cf. Lemma G.1).

Beyond the two cases in the paper, it would be interesting to apply Theorem 2 with results lower bounding s_{\min} for other examples (e.g. for the kernel random matrices [HLCH19]).

In the underparameterized regime $d < n$, we have $\lim_{\lambda \downarrow 0} \lambda_\star(\lambda) = 0$ and therefore the previous bounds do not apply. In this case, we trivially have $\mathcal{B}_{\mathbf{X}}(0) = \mathbf{B}_n(0) = 0$. The proof for the variance approximation requires a different proof, which is presented in Appendix E.

Theorem 3 (Ridgeless regression in the underparameterized regime). *Suppose Assumption 1 holds with $n > d$, and further assume*

$$\mathbf{C}_\Sigma = \min \left(\frac{d}{n}, 1 - \frac{d}{n} \right) \in (0, 1).$$

For any positive integers k and D , there exist constants $\eta = \eta(\mathbf{C}_x) > 0$, $\mathbf{C}_1 = \mathbf{C}_1(\mathbf{C}_x, D) > 0$, $\mathbf{C}_2 = \mathbf{C}_2(k, \mathbf{C}_x, D) > 0$, such that the following hold.

If \mathbf{X} has rank d and s_{\min} is the minimum eigenvalue of the sample covariance $\mathbf{X}^\top \mathbf{X}/n$, then the following hold:

1. **Variance approximation.** Let ε be such that

$$\varepsilon \leq \mathbf{C}_\Sigma^2 \sigma_d / 4, \quad n^{-2D+1} = \mathcal{O} \left(\sqrt{\frac{\mathbf{C}_\Sigma^3 \log^2 n}{\max\{n, \lambda\}}} \right),$$

$$\chi_n(\varepsilon n)^3 \log^2 n \leq \mathbf{C}_1 n \mathbf{C}_\Sigma^{4.5}, \quad \chi_n(\varepsilon n)^3 \log^2 n \leq \mathbf{C}_2 n^{1-\frac{1}{k}} \mathbf{C}_\Sigma^{9.5}.$$

Then, on the event $\{s_{\min} \geq 2\varepsilon\}$, with probability $1 - \mathcal{O}_k(n^{-D+1})$:

$$|\mathcal{V}_{\mathbf{X}}(0) - \mathbf{V}_n(0)| = \mathcal{O}_{k, \mathbf{C}_x, D} \left(\varepsilon \cdot \left(\frac{1}{s_{\min}} + \frac{1}{\mathbf{C}_\Sigma^2 \sigma_d} \right) + \frac{\chi_n(\varepsilon n)^3 \log^2 n}{n^{1-\frac{1}{k}} \mathbf{C}_\Sigma^{9.5}} \right) \cdot \mathbf{V}_n(0).$$

2. **Bias approximation.** $\mathcal{B}_{\mathbf{X}}(0) = \mathbf{B}_n(0) = 0$ (this holds deterministically on the event $\text{rank}(\mathbf{X}) = d$).

4 Applications

4.1 Proportional regime

As a first application, we revisit the proportional regime that is defined by the following condition.

Assumption 2. *There exists a constant $M > 1$ such that $M^{-1} \leq d/n \leq M$ and $\sigma_d \geq M^{-1}$.*

This case is well studied and is not the main motivation of the present paper, but it is nevertheless important to compare our results to earlier work. We refer the reader to [Dic16, ASS20, DW18, WX20, RMR21] for background.

Among others, the results of [HMRT22] are more directly comparable to ours because they establish nonasymptotic bounds comparing variance and bias to the effective variance and bias of Eqs. (6) and (7), for both ridge and ridgeless regression. The proofs of [HMRT22] build on recent advances in random matrix theory, and in particular the anisotropic local law of [KY17].

Here we apply Theorems 1, 2 and 3 to the proportional regime. We note that, under assumption 2, the minimum eigenvalue of $\mathbf{X}^\top \mathbf{X}$ is, with high probability, of order n . In order for the ridge regularization to have a non-trivial effect, we need to choose $\lambda \asymp n$ as well, cf. (4a) and (4b). We will therefore assume λ/n bounded above and below (there is no loss of generality in using the same constant as in Eq. (2)). We will address the case $\lambda = 0+$ in a separate statement below.

Proposition 4.1. *Let Assumptions 1 and 2 hold, and further assume $\lambda/n \in [1/M, M]$. Then for any positive integers k and D , if $n = \Omega_{k,M,C_x,D}(1)$, with probability $1 - \mathcal{O}_k(n^{-D+1})$ we have*

$$\begin{aligned} |\mathcal{V}_{\mathbf{X}}(\lambda) - \mathbf{V}_n(\lambda)| &= \mathcal{O}_{k,M,C_x,D} \left(\frac{\log^8 n}{n^{1-\frac{1}{k}}} \right) \cdot \mathbf{V}_n(\lambda), \\ |\mathcal{B}_{\mathbf{X}}(\lambda) - \mathbf{B}_n(\lambda)| &= \mathcal{O}_{k,M,C_x,D} \left(\frac{\log^8 n}{n^{\frac{1}{2}-\frac{1}{k}}} \right) \cdot \mathbf{B}_n(\lambda). \end{aligned}$$

The proof of this result is presented in Appendix F.

We note that the rates $\mathcal{O}(n^{-1})$ and $\mathcal{O}(n^{-1/2})$ are optimal for variance and bias approximation—corresponding to fluctuations of the average law and local law for the resolvent [AEK⁺14, KY17]. The most direct comparison is with [HMRT22, Theorem 5]: let us point out two ways in which the present result improves over the earlier [HMRT22].

- In [HMRT22], the rate for variance approximation of ridge regression is $\mathcal{O}(n^{-1/2})$, while here we obtain the faster rate $\mathcal{O}(n^{-1})$.
- Error terms in [HMRT22] are additive, while Proposition 4.1 provides multiplicative error terms: the quality of approximation does not deteriorate in the interesting case in which bias and variance become small.

Note that [HMRT22] informally claimed that $n^{-1/2}$ is the optimal rate in the above estimates. While this is correct for the bias, for the variance Proposition 4.1 yields a faster rate. As related phenomenon arises for linear eigenvalue statistics of random matrices (i.e. statistics of the form $n^{-1} \sum_{i=1}^n \varphi(\lambda_i)$). While naively such statistics would have normal deviations of order $n^{-1/2}$, the actual deviations are of order n^{-1} because of eigenvalues correlations [LP09].

We finally consider the ridgeless case.

Proposition 4.2. *Let Assumptions 1 and 2 hold for $\mathbf{x}_i = \Sigma^{1/2} \mathbf{z}_i$, where \mathbf{z}_i has i.i.d. sub-Gaussian coordinates.*

1. **Overparameterized regime.** If additionally $d/n \geq 1 + M^{-1}$, then for all $n = \Omega_{M, \mathbf{C}_x, D}(1)$, with probability $1 - \mathcal{O}(n^{-D+1})$ we have

$$\begin{aligned} |\mathcal{V}_{\mathbf{X}}(0) - \mathbf{V}_n(0)| &= \mathcal{O}_{M, \mathbf{C}_x, D} \left(n^{-1/14} \right) \cdot \mathbf{V}_n(0), \\ |\mathcal{B}_{\mathbf{X}}(0) - \mathbf{B}_n(0)| &= \mathcal{O}_{M, \mathbf{C}_x, D} \left(n^{-1/28} \right) \cdot \mathbf{B}_n(0). \end{aligned}$$

2. **Underparameterized regime.** If additionally $d/n \leq 1 - M^{-1}$, then for all $n = \Omega_{M, \mathbf{C}_x, D}(1)$, with probability $1 - \mathcal{O}(n^{-D+1})$ we have

$$\begin{aligned} |\mathcal{V}_{\mathbf{X}}(0) - \mathbf{V}_n(0)| &= \mathcal{O}_{M, \mathbf{C}_x, D} \left(n^{-1/5} \right) \cdot \mathbf{V}_n(0), \\ \mathcal{B}_{\mathbf{X}}(0) &= \mathbf{B}_n(0) = 0. \end{aligned}$$

We do not expect the exponent $1/14$, $1/28$, $1/5$ in this statement to be tight. However, as in the positive λ case, also in this case the error is multiplicative.

4.2 Bounded varying spectrum

We next consider the highly overparametrized case $d \gg n$. Overparametrized ridge (or minimum norm) regression attracted significant attention recently because of the realization that many deep learning models are overparametrized and overfit the training data. This connection is reviewed in [BMR21, Bel21].

Here we consider covariate vectors \mathbf{x}_i taking values in a general Hilbert space with $d = \infty$, under Assumption 1 on the covariates distribution. This is most closely related to [BLLT20, TB20], and [KZSS21]. The last paper derives refined upper bounds using Gaussian width techniques, but is limited to the case of Gaussian covariates and, as for earlier results, is only accurate up to constant factors.

We will show that our general theory yields excess risk estimates that are accurate up to $1 + o_n(1)$ multiplicative errors. We impose the following condition on the spectrum of Σ .

Assumption 3 (Bounded varying spectrum). *There exists a monotone decreasing function $\psi : (0, 1] \rightarrow [1, \infty)$ with $\lim_{\delta \downarrow 0} \psi(\delta) = \infty$, such that $\sigma_{\lfloor \delta i \rfloor} / \sigma_i \leq \psi(\delta)$ for all $\delta \in (0, 1]$, $i \in \mathbb{N}$ and $\delta i \geq 1$.*

Recall that, by definition, for any $j \leq i$, $\sigma_j / \sigma_i \geq 1$. The bounded varying condition requires that, if i, j diverge proportionally, then the eigenvalue ratio σ_j / σ_i stays bounded. Note that this assumption is equivalent to $\sup_{i \geq 1} \sigma_{\lfloor \delta i \rfloor} / \sigma_i < \infty$ for every $\delta \in (0, 1]$, which is in turn equivalent to

$$\limsup_{i \rightarrow \infty} \frac{\sigma_{\lfloor \delta i \rfloor}}{\sigma_i} < \infty. \quad (25)$$

As special case, Assumption 3 holds if the sorted eigenvalues $(\sigma_1, \sigma_2, \dots)$ forms a so-called *regularly varying sequence*, namely for any $\delta \in (0, \infty)$,

$$\lim_{i \rightarrow \infty} \frac{\sigma_{\lfloor \delta i \rfloor}}{\sigma_i} = \psi(\delta),$$

where $\psi(\delta)$ is positive and finite for any δ . In other words, in the regularly varying case, the ratio σ_j / σ_i converges when i, j diverge proportionally.

A special case of regularly varying spectrum is given by Zipf's law whereby $\sigma_i = i^{-\alpha}$ for some $\alpha > 1$ (in this case $\psi(\delta) = \delta^{-\alpha}$). Regularly varying functions were characterized by Karamata

[Kar33] (for functions on the positive real line), and by Galambos and Seneta [GS73] (for the sequences, i.e. functions defined on the naturals). Namely all such sequences take the form

$$\sigma_i = i^{-\alpha} a_i \exp \left\{ \sum_{j=1}^i b_j/j \right\},$$

where a_i are arbitrary and converge to a positive limit as $i \rightarrow \infty$ and $b_i \rightarrow 0$.

It is easy to see that Assumption 3 holds beyond the case of regularly varying sequences. Consider for instance $\sigma_i = 3^{-s}$ for all $2^s \leq i < 2^{s+1}$, $s = 0, 1, \dots$.

Applying Theorems 1 and 2 to Σ with bounded varying spectrum, we obtain the following result, whose proofs are detailed in Appendix G.

Proposition 4.3. *Let Assumptions 1 and 3 hold. For any constants $M > 0$, $\gamma \in (0, 1/3)$, and positive integers k, D the following holds. If $\mathbf{d}_\Sigma \leq Mn^{1+\gamma}$ and $\lambda \in [n\lambda_*(0)/M, n\lambda_*(0)M]$, then for $n = \Omega_{k,M,\psi,\gamma,\mathbf{C}_x,D}(1)$, with probability $1 - \mathcal{O}_k(n^{-D+1})$*

$$|\mathcal{V}_\mathbf{X}(\lambda) - \mathbf{V}_n(\lambda)| = \mathcal{O}_{k,M,\psi,\mathbf{C}_x,D} \left(\frac{(\mathbf{d}_\Sigma/n)^3 \log^8 n}{n^{1-\frac{1}{k}}} \right) \cdot \mathbf{V}_n(\lambda).$$

If additionally $\mathbf{d}_\Sigma = \mathcal{O}_{M,\psi,\mathbf{C}_x}(n^{1+\gamma} (\rho(\lambda))^{1/6})$ $\lambda \in [n\lambda_*(0)/M, n\lambda_*(0)M]$ and $\lambda_*(0) = \mathcal{O}(1)$ (cf. Theorem 1 for the function ρ), with the same probability we have

$$|\mathcal{B}_\mathbf{X}(\lambda) - \mathbf{B}_n(\lambda)| = \mathcal{O}_{k,M,\psi,\mathbf{C}_x,D} \left(\frac{(\mathbf{d}_\Sigma/n)^3 \log^8 n}{\sqrt{\rho(\lambda)} n^{1-\frac{1}{k}}} \right) \cdot \mathbf{B}_n(\lambda).$$

Applying Theorem 2, we have the following conclusion for ridgeless regression.

Proposition 4.4. *Let Assumptions 1 and 3 hold. Suppose $\beta = \Sigma^{1/2}\theta$ with $\|\theta\| < \infty$. If we have $\lambda_*(0)/\sigma_n = \mathcal{O}(\log^{\mathcal{O}(1)} n)$ and $\mathbf{d}_\Sigma(n) = \mathcal{O}(n \log^{\mathcal{O}(1)} n)$, for any $n = \Omega_{\psi,\mathbf{C}_x,D}(1)$, it holds with probability $1 - \mathcal{O}(n^{-D+1})$ that*

$$|\mathcal{V}_\mathbf{X}(0) - \mathbf{V}_n(0)| = \mathcal{O}_{\psi,\mathbf{C}_x,D} \left(n^{-1/15} \right) \cdot \mathbf{V}_n(0).$$

Remark 4.1. The assumptions $\lambda_*(0)/\sigma_n = \mathcal{O}(\log^{\mathcal{O}(1)} n)$ and $\mathbf{d}_\Sigma = \mathcal{O}(n \log^{\mathcal{O}(1)} n)$ are primarily introduced to simplify the form of the statement. These two conditions can be relaxed to $\lambda_*(0)/\sigma_n = \mathcal{O}(n^{\bar{\gamma}})$ and $\mathbf{d}_\Sigma = \mathcal{O}(n^{1+\bar{\gamma}})$ for a sufficiently small $\bar{\gamma}$, but we do not pursue this generalization here.

Remark 4.2. It is possible to apply the upper/lower bounds on the bias of Theorem 2 to prove bounds on the bias in the setting of Proposition 4.4. However the resulting error term is larger than $(\sigma_n^2 \|\theta_{\leq n}\|^2 + \sigma_n \|\beta_{> n}\|^2)$, which is the size of the upper bound on $\mathbf{B}_n(0)$ in Proposition 2.2.

In order to illustrate the accuracy of our general framework, we apply Proposition 4.3 to derive sharp asymptotics for bias and variance in a number cases. In each of the case below, we scale the regularization parameter λ as $\lambda = \lambda_0(n) \cdot \nu$ for a certain explicit function $\lambda_0(n)$. The scaling $\lambda_0(n)$ is chosen so that the bias and variance retain a non-trivial dependence on ν for large n . We expect that the excess risk achieved by optimal regularization is also covered by this scaling (up to negligible corrections), but do not prove it formally here.

Theorem 4. *Let Assumption 1 hold. Then, for a fixed constant $\nu > 0$ and any positive integer D , the following events hold with probability $1 - \mathcal{O}(n^{-D})$ (the $o_n(1)$ errors may depend on D):*

1. **Regularly varying spectrum with $\alpha > 1$.** Assume $(\sigma_i)_{i \geq 1}$ is a regularly varying sequence with exponent $\alpha > 1$. As a consequence, $\sigma_i = i^{-\alpha} a_i \exp \left\{ \sum_{j=1}^i b_j/j \right\}$ with a_i converging to a positive limit and $b_i \rightarrow 0$. Define $\mathbf{c}_\star = \mathbf{c}_\star(\nu) > 0$ as the unique positive solution of

$$1 = \nu \mathbf{c}_\star^{-1} + \frac{\pi/\alpha}{\sin(\pi/\alpha)} \mathbf{c}_\star^{-1/\alpha}.$$

Then we have

$$\lambda_\star(\nu n^{1-\alpha}) = \mathbf{c}_\star \sigma_n (1 + o_n(1)), \quad (26)$$

$$\mathcal{V}_\mathbf{X}(\nu n^{1-\alpha}) = \frac{\tau^2 (1 - \nu \mathbf{c}_\star^{-1})(\alpha - 1)}{1 + \nu \mathbf{c}_\star^{-1}(\alpha - 1)} (1 + o_n(1)). \quad (27)$$

Let $F_\beta(x) = \sum_{k=1}^{\lfloor nx \rfloor} \langle \beta, \mathbf{v}_k \rangle^2$. If additionally β satisfies the following ‘‘polynomial-decay’’ property: for some $0 < \theta \leq 1$ that

$$\int_0^\infty x^\alpha dF_\beta(x) = \mathcal{O} \left(n^{1-\theta} \int_0^\infty x^\alpha (1 + \mathbf{c}_\star x^\alpha)^{-1} dF_\beta(x) \right),$$

we further have

$$\mathcal{B}_\mathbf{X}(\nu n^{1-\alpha}) = \frac{\sigma_n \mathbf{c}_\star^2 \alpha}{1 + \nu \mathbf{c}_\star^{-1}(\alpha - 1)} \int_0^\infty \frac{x^\alpha}{(1 + \mathbf{c}_\star x^\alpha)^2} dF_\beta(x) (1 + o_n(1)). \quad (28)$$

2. **Regularly varying spectrum with $\alpha = 1$.** Next consider the case $\sigma_i = i^{-1} a_i (1 + \log i)^{-\alpha'}$ for some $\alpha' > 1$ with a_i converging to a positive limit. Define $\mathbf{c}_\star = \mathbf{c}_\star(\nu) > 0$ as

$$\mathbf{c}_\star = \nu + \frac{1}{\alpha' - 1}.$$

We have

$$\lambda_\star(\nu \log^{1-\alpha'} n) = \mathbf{c}_\star \sigma_n \log n (1 + o_n(1)), \quad (29)$$

$$\mathcal{V}_\mathbf{X}(\nu \log^{1-\alpha'} n) = \frac{\tau^2}{\mathbf{c}_\star \log n} (1 + o_n(1)). \quad (30)$$

Let $F_\beta(x) = \sum_{k=1}^{\lfloor (n/\log n)x \rfloor} \langle \beta, \mathbf{v}_k \rangle^2$. If additionally β satisfies the following ‘‘rapid-decay’’ property: for some $0 < \theta \leq 1$ that

$$\int_0^\infty x dF_\beta(x) = \mathcal{O} \left(n^{1-\theta} \int_0^\infty x (1 + \mathbf{c}_\star x)^{-1} dF_\beta(x) \right).$$

then we further have

$$\mathcal{B}_\mathbf{X}(\nu \log^{1-\alpha'} n) = \mathbf{c}_\star^2 \sigma_n \log n \int_0^\infty \frac{x}{(1 + \mathbf{c}_\star x)^2} dF_\beta(x) (1 + o_n(1)). \quad (31)$$

3. **A non-regularly varying spectrum.** $\sigma_i = p^{-s}$ for all $q^s \leq i < q^{s+1}$, with $1 < q < p$ and $s = 0, 1, \dots$. Define s_\star such that $q^{s_\star} \leq n < q^{s_\star+1}$, and for positive integer r the following decreasing function in $t > 0$,

$$G_{p,q,r}(t) = \sum_{k=-\infty}^{\infty} \frac{q^k}{(1 + tp^k)^r}.$$

Let $\rho_\star = n/(q^{s_\star+1} - q^{s_\star}) \in [1/(q-1), q/(q-1))$. Then there exists a unique solution $\mathbf{c}_\star = \mathbf{c}_\star(\nu)$ to the following equation

$$1 = \nu \mathbf{c}_\star^{-1} + \rho_\star^{-1} \cdot G_{p,q,1}(\mathbf{c}_\star).$$

Then we have

$$\lambda_\star(\nu n p^{-s_\star}) = \mathbf{c}_\star \sigma_n (1 + o_n(1)), \quad (32)$$

$$\mathcal{V}_{\mathbf{X}}(\nu n p^{-s_\star}) = \frac{G_{p,q,2}(\mathbf{c}_\star) \tau^2}{\rho_\star - G_{p,q,2}(\mathbf{c}_\star)} (1 + o_n(1)). \quad (33)$$

Let $F_{\boldsymbol{\beta}}(x) = \sum_{k=1}^{\lceil x \rceil - 1} \langle \boldsymbol{\beta}, \mathbf{v}_k \rangle^2$. If additionally $\boldsymbol{\beta}$ satisfies the following ‘‘rapid-decay’’ property: for some $0 < \theta \leq 1$ that

$$\int_0^\infty p^{x-s_\star} dF_{\boldsymbol{\beta}}(x) = \mathcal{O} \left(n^{1-\theta} \int_0^\infty p^{x-s_\star} (1 + \mathbf{c}_\star p^{x-s_\star})^{-1} dF_{\boldsymbol{\beta}}(x) \right),$$

we further have

$$\mathcal{B}_{\mathbf{X}}(\nu n p^{-s_\star}) = \frac{\mathbf{c}_\star^2 \sigma_n}{1 - \rho_\star^{-1} G_{p,q,2}(\mathbf{c}_\star)} \int_0^\infty \frac{p^{x-s_\star}}{(1 + \mathbf{c}_\star p^{x-s_\star})^2} dF_{\boldsymbol{\beta}}(x) (1 + o_n(1)). \quad (34)$$

The proof of this theorem is presented in Appendix H.

Remark 4.3. In the case of a regularly varying spectrum with $\alpha > 1$, the bias vanishes with the sample size as $n^{-\alpha+o(1)}$ but the variance stays bounded away from zero as long as $\tau > 0$, cf. Eq. (27). In other words in this case overfitting is not benign and Theorem 4 quantifies precisely this claim.

On the other hand, in the case $\alpha = 1$, both bias and variance vanish for large n , and therefore we achieve benign overfitting. We must emphasize however that the variance decay is very slow, namely $\mathcal{V}_{\mathbf{X}}(\lambda) \asymp (\log n)^{-1}$, and hence the decay of the excess risk is at least as slow.

5 Numerical illustrations

In this section we evaluate numerically the theoretical prediction for variance and bias, cf. Eqs. (6), (7) and compare them with the results of numerical simulations with synthetic data. We carry out the simulations in the ridgeless limit $\lambda = 0+$ (corresponding to min-norm interpolation). This case is interesting because it is not covered by some of our theorems. Our numerical experiments suggest that the theoretical predictions of Eqs. (6), (7) hold in a broader domain of validity than the one that we are able to control rigorously.

We use Gaussian covariates \mathbf{x}_i . By rotational invariance, we can limit ourselves to diagonal covariance $\boldsymbol{\Sigma}$. We will consider two eigenvalue structures:

(I) **Regularly varying with $\alpha > 1$.** This is defined by $\sigma_i = i^{-\alpha}$ for all $i \geq 1$. This fits within the first case of Theorem 4.

(II) **Regularly varying with $\alpha = 1$.** This model is defined by $\sigma_i = i^{-1}(1 + \log i)^{-\alpha'}$, with $\alpha' > 1$. This fits within the second case of Theorem 4.

In all numerical experiments, we generate data according to the model (3) with a true parameters vector β concentrated on the top $d_0 = 100$ eigenvectors of Σ . More precisely, we will use $\beta = (1, 1, \dots, 1, 0, 0, \dots)$ where $\|\beta\|_0 = d_0 = 100$.

In Figure 1, we plot our theoretical predictions V_n, B_n, R_n for variance, bias and as a function of the sample size n , for the two models (I) and (II) defined above. We use $\lambda = 0+$. In each case, we consider several values of the exponents α, α' that control the decay of eigenvalues of Σ .

In Figure 2, we plot the same quantities at fixed sample size $n = 500$ and vary the regularization parameter λ . A few facts emerge from these figures:

- For both models, the bias of the minimum norm interpolator is a decreasing function of the sample size n , and appears to vanish as $n \rightarrow \infty$, see second row of Figure 1.
- In contrast, the variance exhibits a strikingly different behavior in the two covariance models, see first row of Figure 1. For model (I) (polynomial eigenvalue decay, with exponent $\alpha > 1$), the variance increases with n , and eventually stabilizes to a limit value. For model (II) (exponent $\alpha = 1$), the variance decreases with n , and appears to vanish, albeit very slowly, as $n \rightarrow \infty$.
- As a consequence of these points, the excess test error of minimum norm interpolation vanishes with sample size in model (II) but does not vanish in model (I). This behavior (and the one at previous points) is precisely quantified by Theorem 4 for $\lambda > 0$.
- Finally the dependence of bias and variance on λ is the expected one. As λ increases, bias increases but variance decreases. However, the balance between these two factors is non-trivial:
 - For the slowest eigenvalue decay (large α in model (I) or large α' in model (II)), the optimal λ is strictly positive.
 - On the other hand, for the fastest eigenvalue decay, the optimal λ vanishes. In these case interpolation is superior to ridge regression: we need to overfit to achieve the best test error.

The above discussion is based on evaluating the theoretical formulas for bias and variance, as given in Eqs. (6), (7). While our main result, Theorems 1, 2 guarantee that these formulas are accurate, it is important how accurate they are at small or moderate n , and whether random deviations modify the picture.

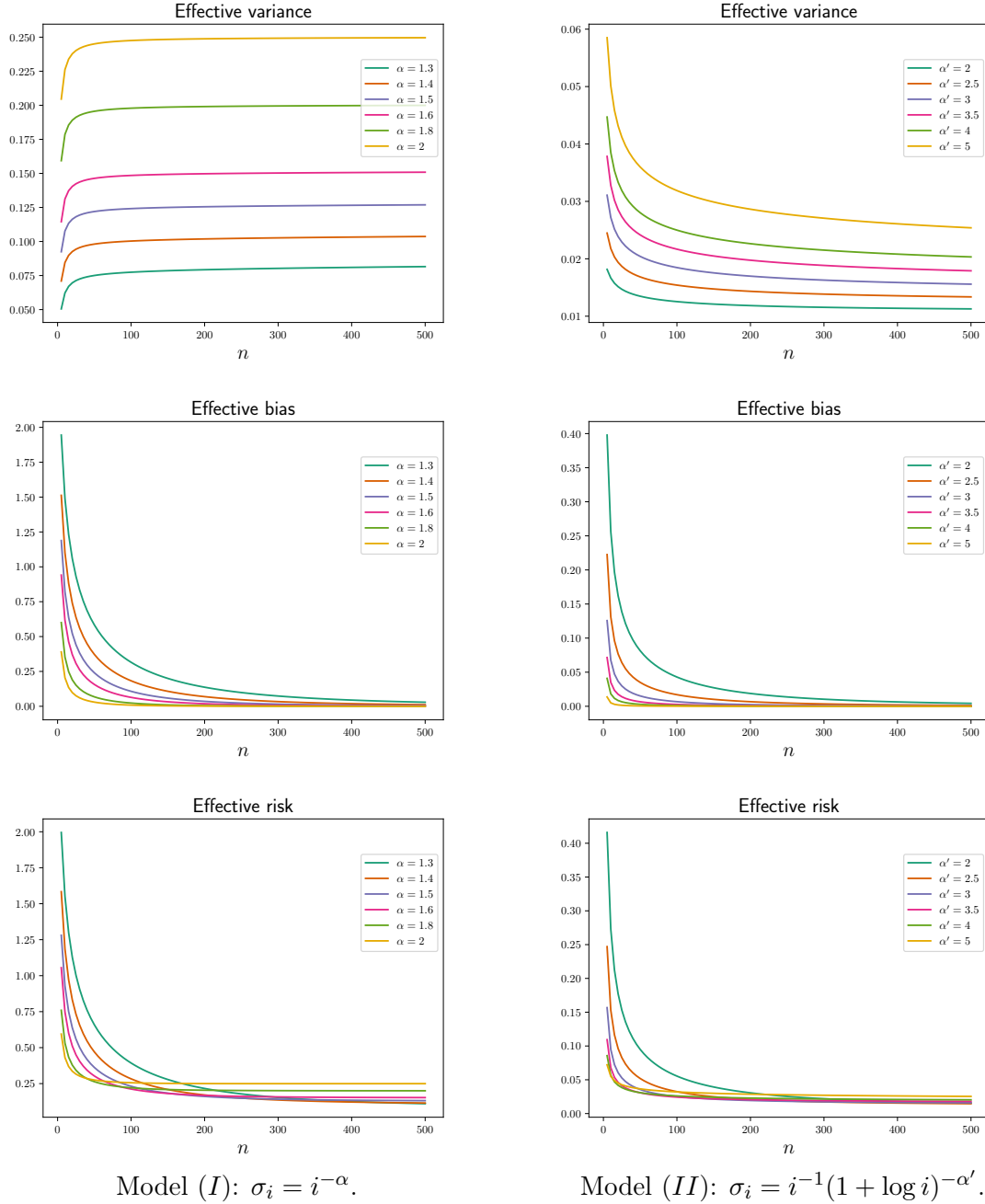
In Figure 3 we plot numerical simulations corroborating that $\mathcal{V}_{\mathbf{X}}, \mathcal{B}_{\mathbf{X}}$ do concentrate around V_n, B_n in models (I) and (II). As mentioned above, the predictions V_n, B_n appear to be accurate beyond what is guaranteed by Theorem 2, and the error appears to be a $(1 + o_n(1))$ multiplicative factor.

6 Proof of Theorem 1

Let $\mathcal{F}_k := \sigma(\mathbf{x}_1, \dots, \mathbf{x}_k)$ be the σ -field generated by the first k data points for $1 \leq k \leq n$, and \mathcal{F}_0 the trivial σ -field. We then have $\mathcal{V}_{\mathbf{X}}, \mathcal{B}_{\mathbf{X}} \in \mathcal{F}_n$ and $V_n, B_n \in \mathcal{F}_0$. Extending the previous notation of \mathcal{R}_0 in Eq. (22) to \mathcal{R}_k , we let

$$\mathcal{R}_k(\mu_0, \mu; \mathbf{Q}) = \text{Tr} \left(\Sigma^{\frac{1}{2}} \mathbf{Q} \Sigma^{\frac{1}{2}} (\mu_0 \mathbf{I} + \mu \Sigma + \mathbf{X}_k^{\top} \mathbf{X}_k)^{-1} \right), \quad \mathcal{F}_k(\mu_0, \mu; \mathbf{Q}) = \mu_0 \mathcal{R}_k(\mu_0, \mu; \mathbf{Q}), \quad (35)$$

where $\mu_0 > 0, \mu \geq 0$, \mathbf{Q} is a p.s.d. matrix with bounded spectral norm, and $\mathbf{X}_k = [\mathbf{x}_1, \dots, \mathbf{x}_k]^{\top} \in \mathbb{R}^{k \times d}$ is the partial data matrix comprising the first k rows of \mathbf{X} . By convention we set $\mathbf{X}_0^{\top} \mathbf{X}_0 := \mathbf{0}$



Model (I): $\sigma_i = i^{-\alpha}$.

Model (II): $\sigma_i = i^{-1}(1 + \log i)^{-\alpha'}$.

Figure 1: Effective variance, bias, and risk of minimum norm interpolation (a.k.a. ridgeless regression) for two covariance structures defined as models (I) and (II) (power law decay of the eigenvalues with exponents $\alpha > 1$ and $\alpha = 1$), as a function of the sample size n . In model (I) we let the noise level to be $\tau = 0.5$, and in (II) we take $\tau = 0.2$.

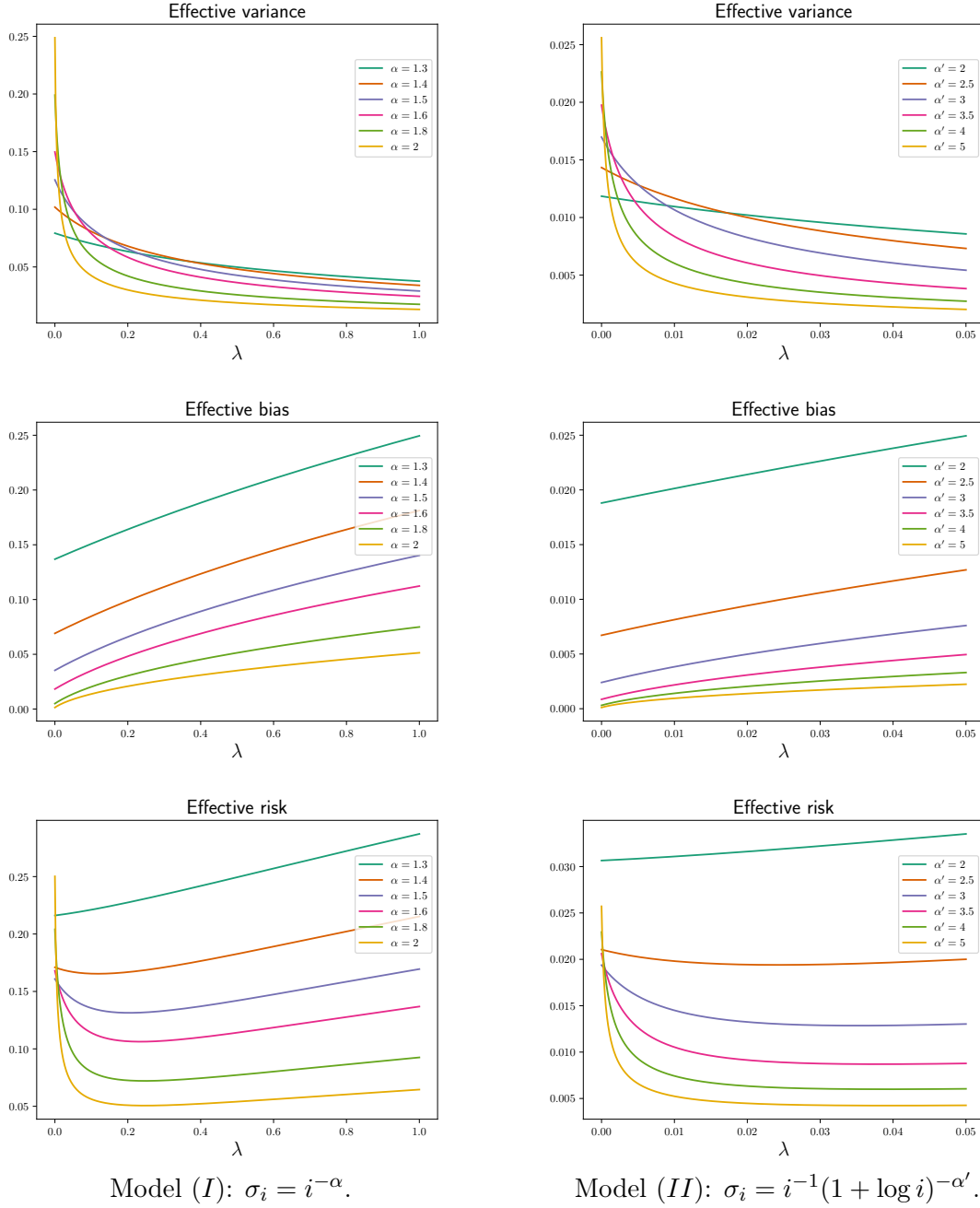


Figure 2: Effective variance, bias, and risk of minimum norm interpolation for two covariance structures defined as models (I) and (II). Here we fix $n = 500$ and vary the regularization parameter. In model (I) we let the noise size to be $\tau = 0.5$, and in (II) we take $\tau = 0.2$.

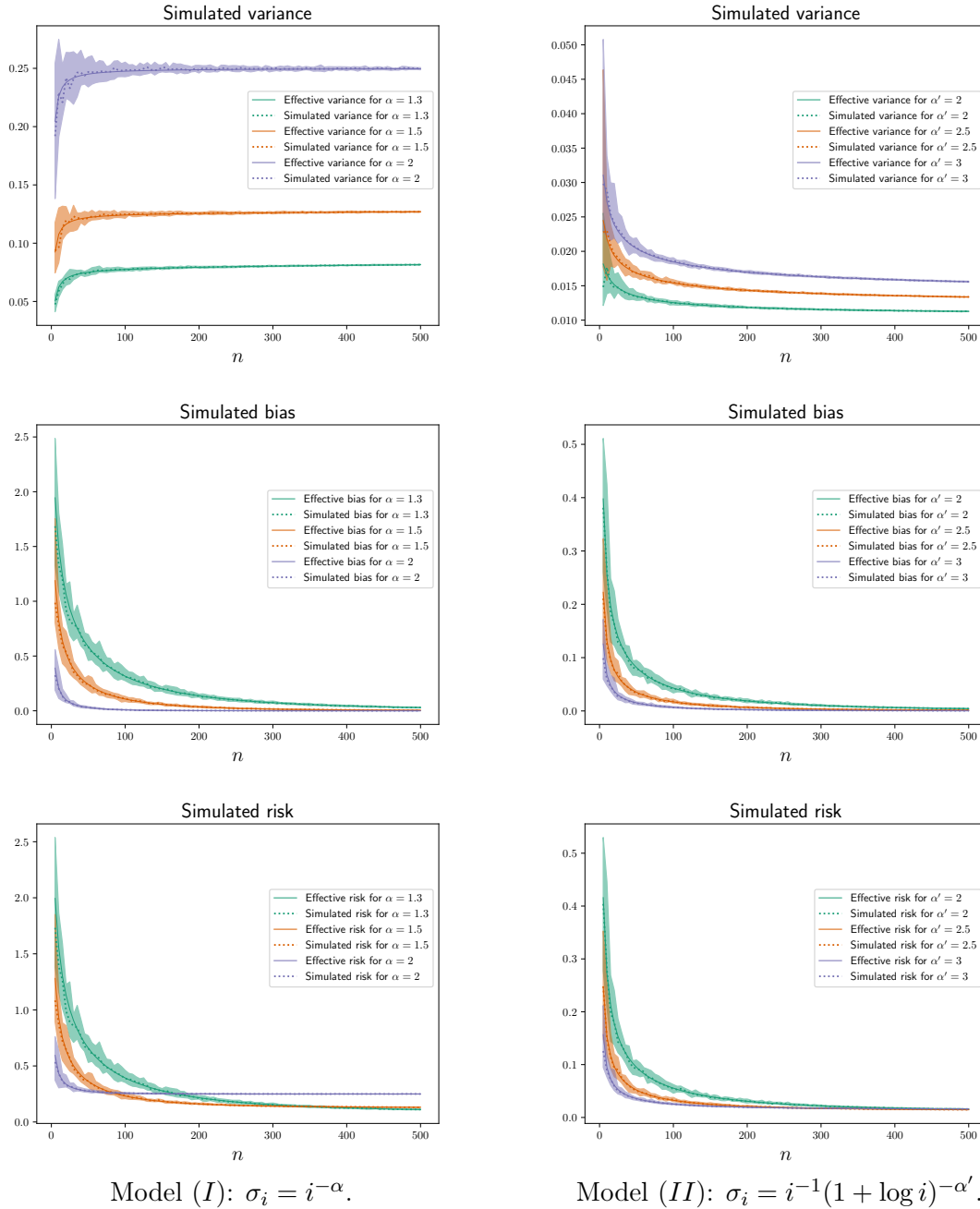


Figure 3: Simulation results for variance, bias, and risk of minimum norm interpolation for two covariance structures defined as models (I) and (II). In model (I) we let the noise size to be $\tau = 0.5$, and in (II) we take $\tau = 0.2$. For each n , we run 20 independent trials and take the median for the dotted lines. We also show the shaded areas between 10% and 90% quantiles.

when $k = 0$. An immediate consequence is that $\mathcal{R}_k, \mathcal{F}_k \in \mathcal{F}_k$. Define $\mu_\star := \mu_\star(\lambda, \mu)$ as the unique solution on of the following equation on (μ, ∞)

$$\mu_\star = \mu + \frac{n}{1 + \mathcal{R}_0(\lambda, \mu_\star; \mathbf{I})}. \quad (36)$$

For $\mu = 0$, this equation reduces to Eq. (5), via the change of variables $\mu_\star = \lambda/\lambda_\star$. For $\mu > 0$ existence and uniqueness follows by a similar argument to the case $\mu = 0$. Indeed, setting $\xi := (\mu_\star - \mu)^{-1}$, the equation is equivalent to $n\xi = 1 + \text{Tr}(\boldsymbol{\Sigma}(\mathbf{A} + \xi\boldsymbol{\Sigma}))$, where $\mathbf{A} := \lambda\mathbf{I} + \mu\boldsymbol{\Sigma}$. Existence and uniqueness follow since the left-hand side is monotone increasing and the right-hand side monotone decreasing in ξ .

In order to quantify the approximation errors $|\mathcal{V}_{\mathbf{X}} - \mathbf{V}_n|$ and $|\mathcal{B}_{\mathbf{X}} - \mathbf{B}_n|$, we will apply the following lemma (Lemma 6.1), which expresses the bias and variance $\mathcal{B}_{\mathbf{X}}, \mathcal{V}_{\mathbf{X}}$ in terms of derivatives of \mathcal{F}_n and \mathcal{F}_0 w.r.t. λ and μ .

Lemma 6.1. *For any $\lambda > 0, \mu \geq 0$, the quantity $\mu_\star > \mu$ is uniquely determined and we have*

$$\begin{aligned} \mathcal{V}_{\mathbf{X}}(\lambda) &= \tau^2 \cdot \left. \frac{\partial}{\partial \lambda} \mathcal{F}_n(\lambda, \mu; \mathbf{I}) \right|_{\mu=0}, & \mathcal{B}_{\mathbf{X}}(\lambda) &= -\lambda \cdot \left. \frac{\partial}{\partial \mu} \mathcal{F}_n(\lambda, \mu; \boldsymbol{\theta}\boldsymbol{\theta}^\top) \right|_{\mu=0}; \\ \mathbf{V}_n(\lambda) &= \tau^2 \cdot \left. \frac{\partial}{\partial \lambda} \mathcal{F}_0(\lambda, \mu_\star(\lambda, \mu); \mathbf{I}) \right|_{\mu=0}, & \mathbf{B}_n(\lambda) &= -\lambda \cdot \left. \frac{\partial}{\partial \mu} \mathcal{F}_0(\lambda, \mu_\star(\lambda, \mu); \boldsymbol{\theta}\boldsymbol{\theta}^\top) \right|_{\mu=0}, \end{aligned}$$

and $\mu_\star(\lambda, 0) = \lambda/\lambda_\star$.

Our proof strategy proceeds in four parts: (I) We show that —due to the regularity properties of \mathcal{F}_0 and \mathcal{F}_n — a bound on $|\mathcal{F}_0 - \mathcal{F}_n|$ implies a bound on the difference of their derivatives, and hence (via Lemma 6.1) on the error in approximating bias and variance; (II) We prove a bound on $|\mathcal{F}_0 - \mathcal{F}_n|$ interpolating between \mathcal{F}_0 and \mathcal{F}_n by adding one row at the time to \mathbf{X} ; (III), (IV) We apply these general bounds respectively to controlling variance and bias.

Recall that we defined $\boldsymbol{\theta} := \boldsymbol{\Sigma}^{-1/2}\boldsymbol{\beta}$, and assumed $\|\boldsymbol{\theta}\| < \infty$. By homogeneity, we can and will assume $\|\boldsymbol{\theta}\| = 1$ throughout the proof.

Part I: Reduction to function values approximation. The following lemma reduces controlling the difference of derivatives of \mathcal{F}_0 and \mathcal{F}_n to the less arduous task of bounding the difference in function values. Its proof is presented in Appendix B.2.

Lemma 6.2. *For any fixed $k \in \mathbb{N}, \delta \in \mathbb{R}_{\geq 0}$ and a $(k+1)$ -times continuously differentiable function $f(t)$ on $[0, k\delta]$, we have*

$$|f'(0)| \leq \mathcal{O}_k \left(\frac{\max_{0 \leq j \leq k} |f(j\delta)|}{\delta} + \sup_{t \in [0, k\delta]} |f^{(k+1)}(t)| \cdot \delta^k \right).$$

With the help of Lemmas 6.1 and 6.2, for any $\delta \in \mathbb{R}_{\geq 0}$, we can upper bound the variance and bias approximations by

$$\begin{aligned} |\mathcal{V}_{\mathbf{X}}(\lambda) - \mathbf{V}_n(\lambda)| &= \mathcal{O}_k \left(\tau^2 \cdot \max_{0 \leq j \leq k} \frac{|\mathcal{F}_n(\lambda + j\delta, 0; \mathbf{I}) - \mathcal{F}_0(\lambda + j\delta, \mu_\star(\lambda + j\delta, 0); \mathbf{I})|}{\delta} \right. \\ &\quad \left. + \tau^2 \delta^k \cdot \sup_{\lambda' \in [\lambda, \lambda + k\delta]} \left| \frac{\partial^{k+1}}{\partial \lambda'^{k+1}} \mathcal{F}_n(\lambda', 0; \mathbf{I}) - \frac{\partial^{k+1}}{\partial \lambda'^{k+1}} \mathcal{F}_0(\lambda', \mu_\star(\lambda', 0); \mathbf{I}) \right| \right), \quad (37) \end{aligned}$$

and

$$\begin{aligned}
|\mathcal{B}_{\mathbf{X}}(\lambda) - \mathbf{B}_n(\lambda)| &= \mathcal{O}_k \left(\lambda \cdot \max_{0 \leq j \leq k} \frac{|\mathcal{F}_n(\lambda, j\delta; \boldsymbol{\theta}\boldsymbol{\theta}^\top) - \mathcal{F}_0(\lambda, \mu_*(\lambda, j\delta); \boldsymbol{\theta}\boldsymbol{\theta}^\top)|}{\delta} \right. \\
&\quad \left. + \lambda \delta^k \cdot \sup_{\mu' \in [0, k\delta]} \left| \frac{\partial^{k+1}}{\partial \mu'^{k+1}} \mathcal{F}_n(\lambda, \mu'; \boldsymbol{\theta}\boldsymbol{\theta}^\top) - \frac{\partial^{k+1}}{\partial \mu'^{k+1}} \mathcal{F}_0(\lambda, \mu_*(\lambda, \mu'); \boldsymbol{\theta}\boldsymbol{\theta}^\top) \right| \right). \tag{38}
\end{aligned}$$

Before passing to bounding errors in function values, we provide upper bounds for higher order derivatives in Eqs. (37) and (38). Bounding the derivatives of \mathcal{F}_n is easier as we can easily write an explicit formula for the k -th derivative for any k . (The proof of this lemma is presented in Appendix B.3).

Lemma 6.3. *For any fixed $k \in \mathbb{N}$, we have for all $\lambda > 0$ and $\mu \geq 0$,*

$$\left| \frac{\partial^k}{\partial \lambda^k} \mathcal{F}_n(\lambda, 0; \mathbf{I}) \right| = \mathcal{O}_k \left(\frac{\mathcal{F}_n(\lambda, 0; \mathbf{I})}{\lambda^k} \right), \quad \left| \frac{\partial^k}{\partial \mu^k} \mathcal{F}_n(\lambda, \mu; \boldsymbol{\theta}\boldsymbol{\theta}^\top) \right| = \mathcal{O}_k \left(\frac{\mathcal{F}_n(\lambda, \mu; \boldsymbol{\theta}\boldsymbol{\theta}^\top)}{\lambda^k} \right).$$

Computing higher order derivatives of \mathcal{F}_0 is less straightforward because \mathcal{F}_0 depends on μ_* which itself depends implicitly depending on (λ, μ) . We postpone this proof to Appendix B.4.

Lemma 6.4. *Let Eq. (21) hold. Then, for any fixed $k \in \mathbb{N}$, we have for all $\lambda > 0$,*

$$\left| \frac{\partial^k}{\partial \lambda^k} \mathcal{F}_0(\lambda, \mu_*(\lambda, 0); \mathbf{I}) \right| = \mathcal{O}_k \left(\frac{\mathcal{F}_0(\lambda, \mu_*(\lambda, 0); \mathbf{I})}{\lambda^k \kappa^{2k}} \right),$$

and for all μ such that $0 \leq \mu \leq \mu_*(\lambda, \mu)/2$,

$$\left| \frac{\partial^k}{\partial \mu^k} \mathcal{F}_0(\lambda, \mu_*(\lambda, \mu); \boldsymbol{\theta}\boldsymbol{\theta}^\top) \right| = \mathcal{O}_k \left(\frac{\mathcal{F}_0(\lambda, \mu_*(\lambda, \mu); \boldsymbol{\theta}\boldsymbol{\theta}^\top)}{\mu_*(\lambda, \mu)^k \kappa^{2k}} \right).$$

Part II: Bounding errors in function values. We next proceed to bounding $|\mathcal{F}_n(\lambda, \mu; \mathbf{Q}) - \mathcal{F}_0(\lambda, \mu_*(\lambda, \mu); \mathbf{Q})|$ for a p.s.d. matrix \mathbf{Q} , which appears in Eqs. (37) and (38). Recall that $\mathcal{F}_i(\lambda, \mu; \mathbf{Q}) = \lambda \mathcal{R}_i(\lambda, \mu; \mathbf{Q})$.

The next theorem bounds $|\mathcal{R}_n(\lambda, \mu; \mathbf{Q}) - \mathcal{R}_0(\lambda, \mu_*(\lambda, \mu); \mathbf{Q})|$ and is the most important technical step in the proof of our main theorems. Its proof is outlined in Section 7, with several technical lemmas deferred to the appendices

Theorem 5. *Introduce the shorthand $\mathbf{R}_0(\mathbf{Q}) := \mathcal{R}_0(\lambda, \mu_*(\lambda, \mu); \mathbf{Q})$. Under Assumption 1, for any $\lambda > 0, \mu \geq 0$, p.s.d. matrix \mathbf{Q} with $\|\mathbf{Q}\| = 1$ and positive integer D , there exists constants $\eta = \eta(\mathbf{C}_{\mathbf{x}}) \in (0, 1/2)$, $\mathbf{C}_\alpha = \mathbf{C}_\alpha(\mathbf{C}_{\mathbf{x}}, D) > 0$, $\mathbf{C}_\beta = \mathbf{C}_\beta(\mathbf{C}_{\mathbf{x}}, D) > 0$ and $\mathbf{C}_\gamma = \mathbf{C}_\gamma(\mathbf{C}_{\mathbf{x}}, D)$ such that for*

$$\begin{aligned}
\gamma &:= \min \left\{ \frac{2}{n} \left(1 + \frac{\mathbf{C}_\gamma \mathbf{d}_{\Sigma} \sigma_{\lceil \eta n \rceil} \cdot \log n \log(\mathbf{d}_{\Sigma} n)}{\lambda} \right) + \frac{2}{\mu_*(\lambda, \mu)}, \frac{1}{\lambda} \right\}, \\
\alpha_1 &:= \mathbf{C}_\alpha \log n \cdot \sqrt{\gamma \mathbf{R}_0(\mathbf{I})}, \\
\alpha_2 &:= \mathbf{C}_\alpha \log n \cdot \sqrt{\gamma^3 \mathbf{R}_0(\mathbf{Q})}, \\
\beta_1 &:= \mathbf{C}_\beta \left(\sqrt{n \log n} \cdot \frac{\alpha_1 \gamma \mathbf{R}_0(\mathbf{Q}) + \alpha_2 (1 + \mathbf{R}_0(\mathbf{I}))}{1 + \mathbf{R}_0(\mathbf{I})^2} + n \cdot \left\{ \frac{\gamma^2 \mathbf{R}_0(\mathbf{Q}) + \alpha_1 \alpha_2}{1 + \mathbf{R}_0(\mathbf{I})^2} + \frac{\alpha_1^2 \gamma \mathbf{R}_0(\mathbf{Q})}{1 + \mathbf{R}_0(\mathbf{I})^3} \right\} + \frac{\gamma \mathbf{R}_0(\mathbf{Q})}{1 + \mathbf{R}_0(\mathbf{I})} \right),
\end{aligned}$$

$$\beta_2 := \frac{C_\beta n \beta_1}{1 + R_0(\mathbf{I})^2},$$

if $\alpha_1 \leq R_0(\mathbf{I})/8$, $\beta_1 \leq R_0(\mathbf{Q})/64$, $\gamma\beta_2(1 + R_0(\mathbf{I})) \leq 1/64$ and $n^{-D} = \mathcal{O}(\alpha_1/(1 + R_0(\mathbf{I})))$, for all $n = \Omega_D(1)$ with probability $1 - \mathcal{O}(n^{-D+1})$ we have

$$|\mathcal{R}_n(\lambda, \mu; \mathbf{Q}) - \mathcal{R}_0(\lambda, \mu_\star(\lambda, \mu); \mathbf{Q})| = \mathcal{O}(\gamma\beta_2(1 + \mathcal{R}_0(\lambda, \mu_\star(\lambda, \mu); \mathbf{I}))\mathcal{R}_0(\lambda, \mu_\star(\lambda, \mu); \mathbf{Q}) + \beta_1).$$

Let us emphasize that this theorem holds under weaker assumptions than Theorem 1, but the error bounds it provides are quite implicit. We can obtain more explicit bounds by imposing the assumptions of Theorem 1. We first define the generalized version of $\rho(\lambda)$ in Eq. (23) for any p.s.d. matrix \mathbf{Q} as

$$\rho(\lambda) := \frac{\mathcal{R}_0(\lambda_\star, 1; \mathbf{Q}/\|\mathbf{Q}\|)}{\mathcal{R}_0(\lambda_\star, 1; \mathbf{I})} \in (0, 1]. \quad (39)$$

The proof of this corollary is given in Appendix C.6.

Corollary 6.5. *Under Assumption 1, for any positive integers k , D and p.s.d. matrix \mathbf{Q} with $\|\mathbf{Q}\| = 1$, there exist constants $\eta = \eta(\mathbf{C}_\mathbf{x}) \in (0, 1/2)$ and $\mathbf{C} = \mathbf{C}(\mathbf{C}_\mathbf{x}, D) > 0$, such that the following hold. Define $\chi_n(\lambda)$, κ , $\rho(\lambda)$ as per Eqs. (20), (21), (39) (those quantities are defined for $\mu = 0$). If it holds that $\mu_\star(\lambda, \mu) \leq (1 - \kappa/2)^{-1}\mu_\star(\lambda, 0)$, and*

$$\chi_n(\lambda)^3 \log^2 n \leq Cn\kappa^{4.5} \sqrt{\rho(\lambda)}, \quad n^{-2D+1} = \mathcal{O}\left(\sqrt{\frac{\kappa^3 \log^2 n}{\max\{n, \lambda\}}}\right),$$

we then have for all $n = \Omega_D(1)$ with probability $1 - \mathcal{O}(n^{-D+1})$ that

$$|\mathcal{R}_n(\lambda, \mu; \mathbf{Q}) - \mathcal{R}_0(\lambda, \mu_\star(\lambda, \mu); \mathbf{Q})| \leq \mathcal{E}_n \cdot \mathcal{R}_0(\lambda, \mu_\star(\lambda, \mu); \mathbf{Q}),$$

where

$$\mathcal{E}_n = \mathcal{O}_{\mathbf{C}_\mathbf{x}, D}\left(\frac{\chi_n(\lambda)^3 \log^2 n}{n\kappa^{6.5}} \cdot \sqrt{\frac{R_0(\mathbf{I})}{R_0(\mathbf{Q})}}\right). \quad (40)$$

To further simplify the assumption $\mu_\star(\lambda, \mu) \leq (1 - \kappa/2)^{-1}\mu_\star(\lambda, 0)$ in Corollary 6.5, the next lemma will be helpful. We defer its proof to Appendix B.5.

Lemma 6.6. *For any fixed $\lambda > 0$, the function $\mu_\star(\lambda, \mu)$ is increasing in μ for all $\mu \geq 0$. Assuming Eq. (21), if $0 \leq \mu \leq n\kappa^3/2$, then*

$$\mu_\star(\lambda, \mu) \leq (1 - \kappa/2)^{-1}\mu_\star(\lambda, 0).$$

Part III: Approximation error for variance. We are now ready to combine our results in Part I and Part II to obtain approximation errors $|\mathcal{V}_\mathbf{X} - \mathbf{V}_n|$ and $|\mathcal{B}_\mathbf{X} - \mathbf{B}_n|$. For the variance, we want to take δ in Eq. (37) such that $k\delta \leq \lambda$. In this case, $[\lambda, \lambda + k\delta] \subset [\lambda, 2\lambda]$. Note that $\lambda_\star(\lambda)$ is an increasing function of λ . Further, by

$$n - \frac{\lambda}{\lambda_\star} = \text{Tr}(\boldsymbol{\Sigma}(\boldsymbol{\Sigma} + \lambda_\star \mathbf{I})^{-1}),$$

we know $\lambda \mapsto \mu_*(\lambda, 0) = \lambda/\lambda_*$ is an increasing function. Therefore $\lambda/\lambda_*(\lambda) \leq 2\lambda/\lambda_*(2\lambda)$, which implies $\lambda_*(\lambda) \leq \lambda_*(2\lambda) \leq 2\lambda_*(\lambda)$. For λ that satisfies Eq. (21), this guarantees that for any $\lambda' \in [\lambda, 2\lambda]$,

$$\frac{\lambda'}{n\lambda_*(\lambda')} \geq \frac{\lambda}{n\lambda_*(\lambda)} \geq \kappa,$$

and

$$\begin{aligned} 1 - \frac{\lambda'}{n\lambda_*(\lambda')} &= \frac{1}{n} \text{Tr}(\boldsymbol{\Sigma}(\boldsymbol{\Sigma} + \lambda_*(\lambda')\mathbf{I})^{-1}) \geq \frac{1}{n} \text{Tr}(\boldsymbol{\Sigma}(\boldsymbol{\Sigma} + 2\lambda_*(\lambda)\mathbf{I})^{-1}) \geq \frac{1}{2} \cdot \frac{1}{n} \text{Tr}(\boldsymbol{\Sigma}(\boldsymbol{\Sigma} + \lambda_*(\lambda)\mathbf{I})^{-1}) \\ &= \frac{1}{2} \left(1 - \frac{\lambda}{n\lambda_*(\lambda)}\right) \geq \kappa/2. \end{aligned}$$

Hence, for any $\lambda' \in [\lambda, 2\lambda]$, Eq. (21) still holds but with constant $\kappa' \geq \kappa/2$. Therefore, we can apply Corollary 6.5 for any $\lambda' \in [\lambda, 2\lambda]$ for $\mathbf{Q} = \mathbf{I}$ and $\mu = 0$, provided the following conditions hold

$$\chi_n(\lambda')^3 \log^2 n \leq Cn(\kappa/2)^{4.5} \sqrt{\rho(\lambda')} = Cn(\kappa/2)^{4.5},$$

where the last equality used the fact that $\rho(\lambda') = 1$ when $\mathbf{Q} = \mathbf{I}$. Finally, setting $C' := 2^{-4.5}C$ and using the fact that $\chi_n(\lambda')$ is decreasing in λ , it suffices to require

$$\chi_n(\lambda)^3 \log^2 n \leq C'n\kappa^{4.5},$$

which holds by the theorem's assumptions.

Hence, we can now apply Corollary 6.5 with $\mathbf{Q} = \mathbf{I}$, and it follows that with probability $1 - \mathcal{O}_k(n^{-D+1})$,

$$\begin{aligned} \max_{0 \leq j \leq k} \frac{|\mathcal{F}_n(\lambda + j\delta, 0; \mathbf{I}) - \mathcal{F}_0(\lambda + j\delta, \mu_*(\lambda + j\delta, 0); \mathbf{I})|}{\delta} &\leq \frac{2\lambda\mathcal{E}_n}{\delta} \cdot \max_{0 \leq j \leq k} \mathcal{R}_0(\lambda + j\delta, \mu_*(\lambda + j\delta, 0); \mathbf{I}) \\ &\leq \frac{2\mathcal{E}_n}{\delta} \mathcal{F}_0(\lambda, \mu_*(\lambda, 0); \mathbf{I}), \end{aligned}$$

where in the last inequality we use that $\mathcal{R}_0(\lambda, \mu_*(\lambda, 0); \mathbf{I}) = n/\mu_*(\lambda, 0) - 1$ is a decreasing function in λ as $\mu_*(\lambda, 0)$ is increasing in λ . Next by Lemmas 6.3 and 6.4 we obtain

$$\begin{aligned} &\sup_{\lambda' \in [\lambda, \lambda + k\delta]} \left| \frac{\partial^{k+1}}{\partial \lambda'^{k+1}} \mathcal{F}_n(\lambda', 0; \mathbf{I}) - \frac{\partial^{k+1}}{\partial \lambda'^{k+1}} \mathcal{F}_0(\lambda', \mu_*(\lambda', 0); \mathbf{I}) \right| \\ &= \mathcal{O}_k \left(\sup_{\lambda' \in [\lambda, 2\lambda]} \frac{\mathcal{F}_n(\lambda', 0; \mathbf{I}) + \mathcal{F}_0(\lambda', \mu_*(\lambda', 0); \mathbf{I})}{\lambda'^{k+1} \kappa^{2k+2}} \right) \\ &= \mathcal{O}_k \left(\sup_{\lambda' \in [\lambda, 2\lambda]} \frac{\lambda' \mathcal{R}_n(\lambda', 0; \mathbf{I}) + \lambda' \mathcal{R}_0(\lambda', \mu_*(\lambda', 0); \mathbf{I})}{\lambda^{k+1} \kappa^{2k+2}} \right) \\ &\stackrel{(i)}{=} \mathcal{O}_k \left(\frac{\lambda \mathcal{R}_n(\lambda, 0; \mathbf{I}) + \lambda \mathcal{R}_0(\lambda, \mu_*(\lambda, 0); \mathbf{I})}{\lambda^{k+1} \kappa^{2k+2}} \right) \\ &\stackrel{(ii)}{=} \mathcal{O}_k \left(\frac{(1 + \mathcal{E}_n) \mathcal{F}_0(\lambda, \mu_*(\lambda, 0); \mathbf{I})}{\lambda^{k+1} \kappa^{2k+2}} \right), \end{aligned}$$

where in (i) we use again that $\mathcal{R}_n(\lambda, \mu; \mathbf{I})$ and $\mathcal{R}_0(\lambda, \mu_*(\lambda, 0); \mathbf{I})$ are decreasing in λ and in (ii) we apply Corollary 6.5. Substituting the above displays into Eq. (37), we have

$$|\mathcal{V}_{\mathbf{X}}(\lambda) - \mathbf{V}_n(\lambda)| = \mathcal{O}_k \left(\left(\frac{\mathcal{E}_n}{\delta} + \frac{(1 + \mathcal{E}_n)\delta^k}{\lambda^{k+1} \kappa^{2k+2}} \right) \cdot \tau^2 \mathcal{F}_0(\lambda, \mu_*(\lambda, 0); \mathbf{I}) \right). \quad (41)$$

Finally, we use the fact that

$$\begin{aligned}
\mathcal{F}_0(\lambda, \mu_\star(\lambda, 0); \mathbf{I}) &= \lambda \text{Tr} \left(\boldsymbol{\Sigma} \left(\frac{\lambda}{\lambda_\star} \boldsymbol{\Sigma} + \lambda \mathbf{I} \right)^{-1} \right) = \lambda_\star \text{Tr} (\boldsymbol{\Sigma} (\boldsymbol{\Sigma} + \lambda_\star \mathbf{I})^{-1}) \\
&= \lambda \cdot \frac{\text{Tr} (\boldsymbol{\Sigma} (\boldsymbol{\Sigma} + \lambda_\star \mathbf{I})^{-1})}{n - \text{Tr} (\boldsymbol{\Sigma} (\boldsymbol{\Sigma} + \lambda_\star \mathbf{I})^{-1})} \leq \lambda \cdot \frac{\text{Tr} (\boldsymbol{\Sigma}^2 (\boldsymbol{\Sigma} + \lambda_\star \mathbf{I})^{-2})}{n - \text{Tr} (\boldsymbol{\Sigma} (\boldsymbol{\Sigma} + \lambda_\star \mathbf{I})^{-1})} \\
&\leq \lambda \cdot \frac{n - \text{Tr} (\boldsymbol{\Sigma}^2 (\boldsymbol{\Sigma} + \lambda_\star \mathbf{I})^{-2})}{n - \text{Tr} (\boldsymbol{\Sigma} (\boldsymbol{\Sigma} + \lambda_\star \mathbf{I})^{-1})} \cdot \frac{\text{Tr} (\boldsymbol{\Sigma}^2 (\boldsymbol{\Sigma} + \lambda_\star \mathbf{I})^{-2})}{n - \text{Tr} (\boldsymbol{\Sigma}^2 (\boldsymbol{\Sigma} + \lambda_\star \mathbf{I})^{-2})} \\
&= \lambda \cdot \frac{n - \text{Tr} (\boldsymbol{\Sigma}^2 (\boldsymbol{\Sigma} + \lambda_\star \mathbf{I})^{-2})}{n - \text{Tr} (\boldsymbol{\Sigma} (\boldsymbol{\Sigma} + \lambda_\star \mathbf{I})^{-1})} \cdot \frac{\partial}{\partial \lambda} \mathcal{F}_0(\lambda, \mu_\star(\lambda, \mu); \mathbf{I}) \Big|_{\mu=0},
\end{aligned}$$

and by Eq. (21),

$$\frac{n - \text{Tr} (\boldsymbol{\Sigma}^2 (\boldsymbol{\Sigma} + \lambda_\star \mathbf{I})^{-2})}{n - \text{Tr} (\boldsymbol{\Sigma} (\boldsymbol{\Sigma} + \lambda_\star \mathbf{I})^{-1})} \leq \frac{n}{n - \text{Tr} (\boldsymbol{\Sigma} (\boldsymbol{\Sigma} + \lambda_\star \mathbf{I})^{-1})} = \frac{n \lambda_\star}{\lambda} \leq \kappa^{-1}.$$

We therefore have, by Lemma 6.1, $\tau^2 \mathcal{F}_0(\lambda, \mu_\star(\lambda, 0); \mathbf{I}) \leq \lambda \kappa^{-1} \mathbf{V}_n(\lambda)$. Substituting in Eq. (41), we obtain

$$|\mathcal{V}_{\mathbf{X}}(\lambda) - \mathbf{V}_n(\lambda)| = \mathcal{O}_k \left(\frac{\lambda \mathcal{E}_n}{\delta \kappa} \cdot + \frac{(1 + \mathcal{E}_n) \delta^k}{\lambda^k \kappa^{2k+3}} \right) \cdot \mathbf{V}_n(\lambda).$$

By setting $\delta = \lambda \kappa^2 n^{-1/k}$, the condition $\delta \kappa \leq \lambda$ is satisfied for all $n = \Omega_k(1)$, which completes the proof for variance approximation with

$$|\mathcal{V}_{\mathbf{X}}(\lambda) - \mathbf{V}_n(\lambda)| = \mathcal{O}_k \left(\mathcal{E}_n \cdot n^{-\frac{1}{k}} \kappa^{-3} + n^{-1} \kappa^{-3} \right) \cdot \mathbf{V}_n(\lambda) = \mathcal{O}_{k, \mathbf{C}_{\mathbf{x}}, D} \left(\frac{\chi_n(\lambda)^3 \log^2 n}{n^{1-\frac{1}{k}} \kappa^{9.5}} \right) \cdot \mathbf{V}_n(\lambda),$$

where we use $\chi_n(\lambda) \geq 1$ in the final bound.

Part IV: Approximation error for bias. Note that all the terms on the right-hand side of Eq. (38) are evaluated at the same value of λ . Hence, Eq. (21) applies to each of these terms. We claim that the assumptions of Corollary 6.5 apply to all of these terms, provided the following conditions hold

$$\mu_\star(\lambda, k\delta) \leq (1 - \kappa/2)^{-1} \mu_\star(\lambda, 0), \tag{42}$$

$$\chi_n(\lambda)^3 \log^2 n \leq C n \kappa^{4.5} \sqrt{\rho(\lambda)}. \tag{43}$$

Indeed, condition (42) implies $\mu_\star(\lambda, \mu) \leq (1 - \kappa/2)^{-1} \mu_\star(\lambda, 0)$ for all $\mu \in [0, k\delta]$ since $\mu \mapsto \mu_\star(\lambda, \mu)$ is monotone decreasing; finally, condition (43) is independent of μ .

then we can apply Lemmas 6.3 and 6.4 and invoke Corollary 6.5 with $\mathbf{Q} = \boldsymbol{\theta} \boldsymbol{\theta}^\top$. To be specific, by Corollary 6.5, we have with probability $1 - \mathcal{O}_k(n^{-D+1})$,

$$\begin{aligned}
\max_{0 \leq j \leq k} \frac{|\mathcal{F}_n(\lambda, j\delta; \boldsymbol{\theta} \boldsymbol{\theta}^\top) - \mathcal{F}_0(\lambda, \mu_\star(\lambda, j\delta); \boldsymbol{\theta} \boldsymbol{\theta}^\top)|}{\delta} &\leq \frac{\mathcal{E}_n}{\delta} \max_{j \in \{0, \dots, k\}} \mathcal{F}_0(\lambda, \mu_\star(\lambda, j\delta); \boldsymbol{\theta} \boldsymbol{\theta}^\top) \\
&\leq \frac{\mathcal{E}_n}{\delta} \mathcal{F}_0(\lambda, \mu_\star(\lambda, 0); \boldsymbol{\theta} \boldsymbol{\theta}^\top),
\end{aligned}$$

as $\mathcal{F}_0(\lambda, \mu_\star(\lambda, \mu); \boldsymbol{\theta}\boldsymbol{\theta}^\top)$ decreases with μ . By Lemmas 6.3 and 6.4 we obtain

$$\begin{aligned}
& \sup_{\mu' \in [0, k\delta]} \left| \frac{\partial^{k+1}}{\partial \mu'^{k+1}} \mathcal{F}_n(\lambda, \mu'; \boldsymbol{\theta}\boldsymbol{\theta}^\top) - \frac{\partial^{k+1}}{\partial \mu'^{k+1}} \mathcal{F}_0(\lambda, \mu_\star(\lambda, \mu'); \boldsymbol{\theta}\boldsymbol{\theta}^\top) \right| \\
&= \mathcal{O}_k \left(\sup_{\mu' \in [0, k\delta]} \frac{\mathcal{F}_n(\lambda, \mu'; \boldsymbol{\theta}\boldsymbol{\theta}^\top)}{\lambda^{k+1}} + \sup_{\mu' \in [0, k\delta]} \frac{\mathcal{F}_0(\lambda, \mu_\star(\lambda, \mu'); \boldsymbol{\theta}\boldsymbol{\theta}^\top)}{\mu_\star(\lambda, \mu')^{k+1} \kappa^{2k+2}} \right) \\
&\stackrel{(i)}{=} \mathcal{O}_k \left(\frac{\mathcal{F}_n(\lambda, 0; \boldsymbol{\theta}\boldsymbol{\theta}^\top)}{\lambda^{k+1}} + \frac{\mathcal{F}_0(\lambda, \mu_\star(\lambda, 0); \boldsymbol{\theta}\boldsymbol{\theta}^\top)}{\mu_\star(\lambda, 0)^{k+1} \kappa^{2k+2}} \right) \\
&\stackrel{(ii)}{=} \mathcal{O}_k \left(\frac{(1 + \mathcal{E}_n + \lambda_\star^{k+1} \kappa^{-2k-2}) \mathcal{F}_0(\lambda, \mu_\star(\lambda, 0); \boldsymbol{\theta}\boldsymbol{\theta}^\top)}{\lambda^{k+1}} \right),
\end{aligned}$$

where in the bound (i) we use the fact that $\mu_\star(\lambda, \mu)$ is increasing in μ (cf. Lemma 6.6) and $\mathcal{F}_k(\lambda, \mu; \mathbf{Q})$ is decreasing in μ when $\mu \geq 0$; in (ii) we use that $\mu_\star(\lambda, 0) = \lambda/\lambda_\star$. Combining the calculations above, we have from Eq. (38)

$$|\mathcal{B}_\mathbf{X}(\lambda) - \mathbf{B}_n(\lambda)| = \mathcal{O}_k \left(\left(\frac{\lambda \mathcal{E}_n}{\delta} + \frac{\delta^k (1 + \mathcal{E}_n + \lambda_\star^{k+1} \kappa^{-2k-2})}{\lambda^k} \right) \cdot \mathcal{F}_0(\lambda, \mu_\star(\lambda, 0); \boldsymbol{\theta}\boldsymbol{\theta}^\top) \right).$$

Then we make use of the following bound

$$\begin{aligned}
\mathcal{F}_0(\lambda, \mu_\star(\lambda, 0); \boldsymbol{\theta}\boldsymbol{\theta}^\top) &= \lambda \text{Tr} \left(\boldsymbol{\Sigma}^{\frac{1}{2}} \boldsymbol{\theta}\boldsymbol{\theta}^\top \boldsymbol{\Sigma}^{\frac{1}{2}} \left(\frac{\lambda}{\lambda_\star} \boldsymbol{\Sigma} + \lambda \mathbf{I} \right)^{-1} \right) = \lambda_\star \text{Tr} \left(\boldsymbol{\Sigma}^{\frac{1}{2}} \boldsymbol{\theta}\boldsymbol{\theta}^\top \boldsymbol{\Sigma}^{\frac{1}{2}} (\boldsymbol{\Sigma} + \lambda_\star \mathbf{I})^{-1} \right) \\
&= \lambda \cdot \frac{\boldsymbol{\theta}^\top (\boldsymbol{\Sigma} + \lambda_\star \mathbf{I})^{-1} \boldsymbol{\Sigma} \boldsymbol{\theta}}{n - \text{Tr}(\boldsymbol{\Sigma}(\boldsymbol{\Sigma} + \lambda_\star \mathbf{I})^{-1})} \leq \frac{\lambda}{n \lambda_\star} \cdot \frac{\lambda_\star^2 \boldsymbol{\theta}^\top (\boldsymbol{\Sigma} + \lambda_\star \mathbf{I})^{-2} \boldsymbol{\Sigma}^2 \boldsymbol{\theta}}{1 - n^{-1} \text{Tr}(\boldsymbol{\Sigma}(\boldsymbol{\Sigma} + \lambda_\star \mathbf{I})^{-1})} \\
&= \frac{\lambda}{n \lambda_\star} \cdot \frac{n - \text{Tr}(\boldsymbol{\Sigma}^2(\boldsymbol{\Sigma} + \lambda_\star \mathbf{I})^{-2})}{n - \text{Tr}(\boldsymbol{\Sigma}(\boldsymbol{\Sigma} + \lambda_\star \mathbf{I})^{-1})} \cdot \frac{\lambda_\star^2 \boldsymbol{\theta}^\top (\boldsymbol{\Sigma} + \lambda_\star \mathbf{I})^{-2} \boldsymbol{\Sigma}^2 \boldsymbol{\theta}}{1 - n^{-1} \text{Tr}(\boldsymbol{\Sigma}^2(\boldsymbol{\Sigma} + \lambda_\star \mathbf{I})^{-2})} \\
&= \frac{\lambda}{n \lambda_\star} \cdot \frac{n - \text{Tr}(\boldsymbol{\Sigma}^2(\boldsymbol{\Sigma} + \lambda_\star \mathbf{I})^{-2})}{n - \text{Tr}(\boldsymbol{\Sigma}(\boldsymbol{\Sigma} + \lambda_\star \mathbf{I})^{-1})} \cdot \mathbf{B}_n(\lambda),
\end{aligned}$$

where in the last line we used the definition of $\mathbf{B}_n(\lambda)$ in Eq. (7). By Eq. (21), we have

$$\frac{\lambda}{n \lambda_\star} \cdot \frac{n - \text{Tr}(\boldsymbol{\Sigma}^2(\boldsymbol{\Sigma} + \lambda_\star \mathbf{I})^{-2})}{n - \text{Tr}(\boldsymbol{\Sigma}(\boldsymbol{\Sigma} + \lambda_\star \mathbf{I})^{-1})} = 1 - \frac{1}{n} \text{Tr}(\boldsymbol{\Sigma}^2(\boldsymbol{\Sigma} + \lambda_\star \mathbf{I})^{-2}) \leq 1,$$

which reduces the approximation bound for bias to

$$|\mathcal{B}_\mathbf{X}(\lambda) - \mathbf{B}_n(\lambda)| = \mathcal{O}_k \left(\frac{\lambda \mathcal{E}_n}{\delta} + \frac{\delta^k (1 + \mathcal{E}_n + \lambda_\star^{k+1} \kappa^{-2k-2})}{\lambda^k} \right) \cdot \mathbf{B}_n(\lambda).$$

We again take $\delta = \lambda \kappa^2 n^{-\frac{1}{k}}$ and the bound becomes

$$|\mathcal{B}_\mathbf{X}(\lambda) - \mathbf{B}_n(\lambda)| = \mathcal{O}_{k, \mathbf{C}_\mathbf{x}, D} \left(\frac{\lambda_\star^{k+1}}{n \kappa^2} + \frac{\chi_n(\lambda)^3 \log^2 n}{n^{1-\frac{1}{k}} \kappa^{8.5}} \cdot \sqrt{\frac{\mathbf{R}_0(\mathbf{I})}{\mathbf{R}_0(\boldsymbol{\theta}\boldsymbol{\theta}^\top)}} \right) \cdot \mathbf{B}_n(\lambda).$$

This bounds hold under the conditions (42) to (43), which are implied by the following:

$$\mu_\star(\lambda, k \lambda \kappa^2 n^{-\frac{1}{k}}) \leq (1 - \kappa/2)^{-1} \mu_\star(\lambda, 0).$$

$$\chi_n(\lambda)^3 \log^2 n \leq C n \kappa^{4.5} \sqrt{\rho(\lambda)}.$$

For the first condition, we invoke Lemma 6.6 to obtain a sufficient requirement $\lambda k n^{-\frac{1}{k}} \leq n \kappa / 2$. For the last condition, it suffices to have $\chi_n(\lambda)^3 \log^2 n \leq C' n \kappa^{4.5} \sqrt{\rho(\lambda)}$ for the same C' defined in Part III.

7 Proof of Theorem 5

Part I: The iterative sequence. The proof is based on the following interpolating construction. We will construct a sequence of random variables $\mu_i \in \mathcal{F}_{i-1}$ for $i = 0, 1, \dots, n+1$ (where, by convention, $\mathcal{F}_{-1} = \mathcal{F}_0$ is the trivial σ -algebra) such that, defining

$$R_i(\mathbf{Q}) := \text{Tr} \left(\Sigma^{\frac{1}{2}} \mathbf{Q} \Sigma^{\frac{1}{2}} \left(\lambda \mathbf{I} + \mu_i \Sigma + \mathbf{X}_i^\top \mathbf{X}_i \right)^{-1} \right) = \mathcal{R}_i(\lambda, \mu_i; \mathbf{Q}), \quad (44)$$

we obtain that $R_i(\mathbf{Q})$ is approximately a martingale and, as a consequence, $R_0(\mathbf{Q}) \approx R_n(\mathbf{Q})$. We will further have $\mu_0 = \mu_\star(\lambda, \mu)$ and $\mu_n \approx \mu$, and therefore we obtain the desired claim $\mathcal{R}_0(\lambda, \mu_\star(\lambda, \mu); \mathbf{Q}) \approx \mathcal{R}_n(\lambda, \mu; \mathbf{Q})$.

Before formally defining the sequence $\{\mu_0, \dots, \mu_{n+1}\}$, we introduce some helpful notations. We first define the matrices $\mathbf{A}_i, \mathbf{B}_i \in \mathcal{F}_i$ for $0 \leq i \leq n$ as

$$\mathbf{A}_i := \Sigma^{\frac{1}{2}} \left(\lambda \mathbf{I} + \mu_i \Sigma + \mathbf{X}_i^\top \mathbf{X}_i \right)^{-1} \Sigma^{\frac{1}{2}}, \quad (45a)$$

$$\mathbf{B}_i := \Sigma^{\frac{1}{2}} \left(\lambda \mathbf{I} + \mu_{i+1} \Sigma + \mathbf{X}_i^\top \mathbf{X}_i \right)^{-1} \Sigma^{\frac{1}{2}}. \quad (45b)$$

Then we can write $R_i(\mathbf{Q}) = \text{Tr}(\mathbf{Q} \mathbf{A}_i)$. Similarly we define another sequence of functions by $S_i(\mathbf{Q}) := \text{Tr}(\mathbf{Q} \mathbf{B}_i)$.

Now we are ready to define the sequence $\mu_i \in \mathcal{F}_{i-1}$. We set the initial value $\mu_0 = \mu_\star(\lambda, \mu) \in \mathcal{F}_{-1}$ and thus $R_0(\mathbf{Q}) = \mathcal{R}_0(\lambda, \mu_\star(\lambda, \mu); \mathbf{Q})$. The sequence $(\mu_j)_{j \geq 1}$ is iteratively determined through the following equation

$$\mu_{i+1} = \mu_i - \frac{1}{1 + S_i(\mathbf{I})} = \mu_i - \frac{1}{1 + \text{Tr} \left(\Sigma^{\frac{1}{2}} \left(\lambda \mathbf{I} + \mu_{i+1} \Sigma + \mathbf{X}_i^\top \mathbf{X}_i \right)^{-1} \Sigma^{\frac{1}{2}} \right)} \quad \text{s.t. } \mathbf{B}_i \succ 0. \quad (46)$$

It is evident that if the solution μ_{i+1} exists and is unique (almost surely with respect to the random choice of μ_i), since $\mu_i \in \mathcal{F}_{i-1}$ and $\mathbf{X}_i \in \mathcal{F}_i$, it follows that $\mu_{i+1} \in \mathcal{F}_i$. The next lemma shows that the iteration via (46) is indeed well-defined. Its proof is in Appendix C.1.

Lemma 7.1. *There exists a unique strictly decreasing sequence $\mu_0 > \mu_1 > \mu_2 > \dots > \mu_n > \mu_{n+1}$ satisfying the update rule (46).*

Part II: Approximation to a martingale. We next explain what is the rationale for the iterative definition of Eq. (46), and how it will help us prove the theorem claim.

Since we want to upper bound $|R_n(\mathbf{Q}) - R_0(\mathbf{Q})|$, it makes sense to compute the difference $R_i(\mathbf{Q}) - R_{i-1}(\mathbf{Q})$,

$$\begin{aligned} R_i(\mathbf{Q}) - R_{i-1}(\mathbf{Q}) &= (R_i(\mathbf{Q}) - S_{i-1}(\mathbf{Q})) + (S_{i-1}(\mathbf{Q}) - R_{i-1}(\mathbf{Q})) \\ &= \underbrace{\text{Tr}(\mathbf{Q}(\mathbf{A}_i - \mathbf{B}_{i-1}))}_{\text{(I)}} + \underbrace{\text{Tr}(\mathbf{Q}(\mathbf{B}_{i-1} - \mathbf{A}_{i-1}))}_{\text{(II)}}. \end{aligned} \quad (47)$$

Using rgw definitions in Eqs. (45a) and (45b), we can further expand (I) by Sherman-Morrison formula

$$\begin{aligned}
\mathbf{A}_i - \mathbf{B}_{i-1} &= \Sigma^{\frac{1}{2}} \left\{ \left(\lambda \mathbf{I} + \mu_i \Sigma + \mathbf{X}_{i-1}^\top \mathbf{X}_{i-1} + \mathbf{x}_i \mathbf{x}_i^\top \right)^{-1} - \left(\lambda \mathbf{I} + \mu_i \Sigma + \mathbf{X}_{i-1}^\top \mathbf{X}_{i-1} \right)^{-1} \right\} \Sigma^{\frac{1}{2}} \\
&= - \frac{\Sigma^{\frac{1}{2}} \left(\lambda \mathbf{I} + \mu_i \Sigma + \mathbf{X}_{i-1}^\top \mathbf{X}_{i-1} \right)^{-1} \mathbf{x}_i \mathbf{x}_i^\top \left(\lambda \mathbf{I} + \mu_i \Sigma + \mathbf{X}_{i-1}^\top \mathbf{X}_{i-1} \right)^{-1} \Sigma^{\frac{1}{2}}}{1 + \mathbf{x}_i^\top \left(\lambda \mathbf{I} + \mu_i \Sigma + \mathbf{X}_{i-1}^\top \mathbf{X}_{i-1} \right)^{-1} \mathbf{x}_i} \\
&= - \frac{\Sigma^{\frac{1}{2}} \left(\lambda \mathbf{I} + \mu_i \Sigma + \mathbf{X}_{i-1}^\top \mathbf{X}_{i-1} \right)^{-1} \Sigma^{\frac{1}{2}} \mathbf{z}_i \mathbf{z}_i^\top \Sigma^{\frac{1}{2}} \left(\lambda \mathbf{I} + \mu_i \Sigma + \mathbf{X}_{i-1}^\top \mathbf{X}_{i-1} \right)^{-1} \Sigma^{\frac{1}{2}}}{1 + \mathbf{z}_i^\top \Sigma^{\frac{1}{2}} \left(\lambda \mathbf{I} + \mu_i \Sigma + \mathbf{X}_{i-1}^\top \mathbf{X}_{i-1} \right)^{-1} \Sigma^{\frac{1}{2}} \mathbf{z}_i} \\
&= - \frac{\mathbf{B}_{i-1} \mathbf{z}_i \mathbf{z}_i^\top \mathbf{B}_{i-1}}{1 + \mathbf{z}_i^\top \mathbf{B}_{i-1} \mathbf{z}_i},
\end{aligned}$$

and thus write

$$(I) = - \frac{\text{Tr}(\mathbf{Q} \mathbf{B}_{i-1} \mathbf{z}_i \mathbf{z}_i^\top \mathbf{B}_{i-1})}{1 + \mathbf{z}_i^\top \mathbf{B}_{i-1} \mathbf{z}_i}.$$

We can also compute (II) by noting that

$$\begin{aligned}
\mathbf{B}_{i-1} - \mathbf{A}_{i-1} &= \Sigma^{\frac{1}{2}} \left\{ \left(\lambda \mathbf{I} + \mu_i \Sigma + \mathbf{X}_{i-1}^\top \mathbf{X}_{i-1} \right)^{-1} - \left(\lambda \mathbf{I} + \mu_{i-1} \Sigma + \mathbf{X}_{i-1}^\top \mathbf{X}_{i-1} \right)^{-1} \right\} \Sigma^{\frac{1}{2}} \\
&= \Sigma^{\frac{1}{2}} \left\{ \left(\lambda \mathbf{I} + \mu_i \Sigma + \mathbf{X}_{i-1}^\top \mathbf{X}_{i-1} \right)^{-1} \cdot (\mu_{i-1} - \mu_i) \Sigma \cdot \left(\lambda \mathbf{I} + \mu_{i-1} \Sigma + \mathbf{X}_{i-1}^\top \mathbf{X}_{i-1} \right)^{-1} \right\} \Sigma^{\frac{1}{2}},
\end{aligned}$$

and therefore

$$(II) = (\mu_{i-1} - \mu_i) \cdot \text{Tr}(\mathbf{Q} \mathbf{B}_{i-1} \mathbf{A}_{i-1}).$$

It is now clear what is the motivation for defining μ_{i+1} as per Eq. (46). We hope to have (I)+(II) ≈ 0 . Under the approximation $\text{Tr}(\mathbf{Q} \mathbf{B}_{i-1} \mathbf{z}_i \mathbf{z}_i^\top \mathbf{B}_{i-1}) \approx \mathbb{E}[\text{Tr}(\mathbf{Q} \mathbf{B}_{i-1} \mathbf{z}_i \mathbf{z}_i^\top \mathbf{B}_{i-1}) \mid \mathcal{F}_{i-1}] = \text{Tr}(\mathbf{Q} \mathbf{B}_{i-1}^2) \approx \text{Tr}(\mathbf{Q} \mathbf{B}_{i-1} \mathbf{A}_{i-1})$, this is achieved when

$$\mu_i - \mu_{i-1} = - \frac{1}{1 + \mathbb{E}[\mathbf{z}_i^\top \mathbf{B}_{i-1} \mathbf{z}_i \mid \mathcal{F}_{i-1}]} = - \frac{1}{1 + \mathbf{S}_{i-1}(\mathbf{I})},$$

which recovers the iteration in Eq. (46).

Part III: Proof via stopping times. We next make the previous argument rigorous. For any scalars $\alpha_1, \alpha_2, \beta_1, \beta_2, \gamma > 0$ (in what follows, we'll use the notation $\Delta := (\alpha_1, \alpha_2, \beta_1, \beta_2, \gamma)$) we consider the events

$$E_i(\mathbf{Q}) := \left\{ \left| \mathbf{z}_i^\top \mathbf{B}_{i-1} \mathbf{z}_i - \mathbf{S}_{i-1}(\mathbf{I}) \right| \leq \alpha_1, \left| \mathbf{z}_i^\top \mathbf{B}_{i-1} \mathbf{Q} \mathbf{B}_{i-1} \mathbf{z}_i - \text{Tr}(\mathbf{Q} \mathbf{B}_{i-1}^2) \right| \leq \alpha_2, \|\mathbf{A}_i\| \leq \gamma \right\}, \quad (48a)$$

$$F_i(\mathbf{Q}) := \left\{ \max\{|\mathbf{R}_i(\mathbf{Q}) - \mathbf{R}_0(\mathbf{Q})|, |\mathbf{S}_i(\mathbf{Q}) - \mathbf{R}_0(\mathbf{Q})|\} \leq \beta_1, |\mu_{i+1} - \bar{\mu}_{i+1}| \leq \beta_2, \|\mathbf{B}_i\| \leq \gamma \right\}, \quad (48b)$$

where $\bar{\mu}_{i+1} = \mu \cdot (i+1)/n + \mu_\star(\lambda, \mu) \cdot (1 - (i+1)/n)$ is nonrandom. In particular we set $E_0(\mathbf{Q}) = \Omega$ so that $E_i(\mathbf{Q})$ and $F_i(\mathbf{Q})$ are well-defined for $0 \leq i \leq n$. It follows then $E_i(\mathbf{Q}), F_i(\mathbf{Q}) \in \mathcal{F}_i$. Next we can proceed to define two stopping times via

$$\{T_E(\mathbf{Q}) \geq k+1\} := \left(\bigcap_{i=0}^k E_i(\mathbf{Q}) \right) \cap \left(\bigcap_{i=0}^{k-1} F_i(\mathbf{Q}) \right), \quad (49a)$$

$$\{T_F(\mathbf{Q}) \geq k+1\} := \left(\bigcap_{i=0}^k E_i(\mathbf{Q}) \right) \cap \left(\bigcap_{i=0}^k F_i(\mathbf{Q}) \right), \quad (49b)$$

for $k = 0, 1, \dots, n$, with $T_E(\mathbf{Q}), T_F(\mathbf{Q}) \in \{0, 1, \dots, n+1\}$. One can easily check that $T_E(\mathbf{Q})$ and $T_F(\mathbf{Q})$ are indeed stopping times since the sets in the above displays are in \mathcal{F}_k , and another immediate consequence is that $T_E(\mathbf{Q}) \geq T_F(\mathbf{Q})$. These stopping times are helpful since the event $\{T_F(\mathbf{Q}) = n+1\}$ implies

$$\max_{0 \leq i \leq n} \{ |R_i(\mathbf{Q}) - R_0(\mathbf{Q})|, |S_i(\mathbf{Q}) - R_0(\mathbf{Q})| \} \leq \beta_1,$$

and thus if β_1 is much smaller than $R_0(\mathbf{Q})$, we can show $R_n(\mathbf{Q}) \approx R_0(\mathbf{Q})$ as desired. Therefore, we want to lower bound the probability for the event $\{T_F(\mathbf{Q}) = n+1\}$. We use the shorthand $p_{i,j}(T_1, T_2, \mathbf{Q}) := \mathbb{P}(T_1(\mathbf{Q}) \geq i, T_2(\mathbf{I}) \geq j)$ for $T_1, T_2 \in \{T_E, T_F\}$. By telescoping sum, we have

$$\begin{aligned} & \mathbb{P}(T_F(\mathbf{Q}) \geq 0, T_F(\mathbf{I}) \geq 0) - \mathbb{P}(T_F(\mathbf{Q}) = n+1, T_F(\mathbf{I}) = n+1) \\ &= p_{0,0}(T_F, T_F, \mathbf{Q}) - p_{n+1, n+1}(T_F, T_F, \mathbf{Q}) \\ &= \sum_{k=0}^n (p_{k,k}(T_F, T_F, \mathbf{Q}) - p_{k+1, k+1}(T_E, T_E, \mathbf{Q})) + \sum_{k=1}^{n+1} (p_{k,k}(T_E, T_E, \mathbf{Q}) - p_{k,k}(T_F, T_F, \mathbf{Q})) \\ &\leq \sum_{k=0}^n (p_{k,k}(T_F, T_F, \mathbf{Q}) - p_{k+1, k}(T_E, T_F, \mathbf{Q}) + p_{k,k}(T_F, T_F, \mathbf{Q}) - p_{k, k+1}(T_F, T_E, \mathbf{Q})) \\ &\quad + \sum_{k=1}^{n+1} (p_{k,k}(T_E, T_E, \mathbf{Q}) - p_{k,k}(T_F, T_E, \mathbf{Q}) + p_{k,k}(T_E, T_E, \mathbf{Q}) - p_{k,k}(T_E, T_F, \mathbf{Q})) \\ &\leq \sum_{k=0}^n (p_{k,k}(T_F, T_F, \mathbf{Q}) - p_{k+1, k}(T_E, T_F, \mathbf{Q}) + p_{k,k}(T_F, T_F, \mathbf{I}) - p_{k+1, k}(T_E, T_F, \mathbf{I})) \\ &\quad + \sum_{k=1}^{n+1} (p_{k,k}(T_E, T_E, \mathbf{Q}) - p_{k,k}(T_F, T_E, \mathbf{Q}) + p_{k,k}(T_E, T_E, \mathbf{I}) - p_{k,k}(T_F, T_E, \mathbf{I})), \end{aligned} \quad (50)$$

where in the last inequality we use $\mathbb{P}(A \cap B) - \mathbb{P}(A \cap B') \leq \mathbb{P}(B) - \mathbb{P}(B')$ for $B' \subset B$.

We are left with the task of bounding the two terms $p_{k,k}(T_F, T_F, \mathbf{Q}) - p_{k+1, k}(T_E, T_F, \mathbf{Q})$ and $p_{k,k}(T_E, T_E, \mathbf{Q}) - p_{k,k}(T_F, T_E, \mathbf{Q})$ for any p.s.d. \mathbf{Q} , and showing that they are small. Before doing this, we show that, by appropriately choosing γ , we have $\|\mathbf{A}_i\| \leq \gamma$ and $\|\mathbf{B}_i\| \leq \gamma$ with high probability. The proof of the next lemma is in Appendix C.2.

Lemma 7.2. *Under Assumption 1, for any positive integer D , there exists a fixed $\eta = \eta(\mathbf{C}_x) \in (0, 1/2)$, such that for all $n = \Omega_D(1)$, it holds with probability $1 - \mathcal{O}(n^{-D})$ that*

$$\left\| \Sigma^{\frac{1}{2}} \left(\lambda \mathbf{I} + \mathbf{X}^\top \mathbf{X} \right)^{-1} \Sigma^{\frac{1}{2}} \right\| \leq \frac{2}{n} \left(1 + \frac{\mathcal{O}_{\mathbf{C}_x, D} \left(d_\Sigma \sigma_{[\eta n]} \cdot \log n \log(d_\Sigma n) \right)}{\lambda} \right), \quad \text{for all } \lambda > 0;$$

additionally, under the same notations of Proposition 2.2, letting $\boldsymbol{\theta}_{\leq k} := \sum_{i \leq k} \langle \boldsymbol{\theta}, \mathbf{v}_i \rangle \mathbf{v}_i$ and $\boldsymbol{\theta}_{> k} := \boldsymbol{\theta} - \boldsymbol{\theta}_{\leq k}$, we have for all $\lambda > 0$,

$$\boldsymbol{\theta}^\top \Sigma^{\frac{1}{2}} \left(\lambda \mathbf{I} + \mathbf{X}^\top \mathbf{X} \right)^{-1} \Sigma^{\frac{1}{2}} \boldsymbol{\theta} \leq \frac{2}{n} \left(1 + \frac{\mathcal{O}_{\mathbf{C}_x, D} \left(d_\Sigma \sigma_{[\eta n]} \cdot \log n \log(d_\Sigma n) \right)}{\lambda} \right) \|\boldsymbol{\theta}_{\leq n}\|^2 + \frac{2 \|\boldsymbol{\beta}_{> n}\|^2}{\lambda}.$$

The next lemma—upper bounding the first term (I)—uses Hanson-Wright inequality to show concentration for events $E_i(\mathbf{Q})$ in (48a). A proof is in Appendix C.3.

Lemma 7.3. *Under Assumption 1, choose β_1, β_2 in Eq. (48b) so that $\beta_1 \leq R_0(\mathbf{Q})/4$ and $\beta_2 \leq \mu/2$. Then for any positive integer D , there exists constants $\eta = \eta(\mathbf{C}_x) \in (0, 1/2)$, $C_\alpha = C_\alpha(\mathbf{C}_x, D)$ and $C_\gamma = C_\gamma(\mathbf{C}_x, D)$ such that if we take*

$$\begin{aligned} \gamma &= \min \left\{ \frac{2}{n} \left(1 + \frac{C_\gamma \mathbf{d}_{\Sigma} \sigma_{\lfloor \eta n \rfloor} \cdot \log n \log(\mathbf{d}_{\Sigma} n)}{\lambda} \right) + \frac{2}{\mu_*(\lambda, \mu)}, \frac{1}{\lambda} \right\}, \\ \alpha_1 &= C_\alpha \log n \cdot \sqrt{\gamma R_0(\mathbf{I})}, \\ \alpha_2 &= C_\alpha \log n \cdot \sqrt{\gamma^3 R_0(\mathbf{Q})}, \end{aligned}$$

it holds for all $n = \Omega_D(1)$ that

$$p_{k,k}(T_F, T_F, \mathbf{Q}) - p_{k+1,k}(T_E, T_F, \mathbf{Q}) = \mathcal{O}(n^{-D}).$$

In addition, on the event $\{T_F(\mathbf{Q}) \geq k, T_F(\mathbf{I}) \geq k\} \in \mathcal{F}_{k-1}$ we have (using the shorthand $\mathbb{E}_{k-1}\{\cdot\} := \mathbb{E}\{\cdot | \mathcal{F}_{k-1}\}$)

$$\begin{aligned} &\mathbb{E}_{k-1} \left[\left| z_k^\top \mathbf{B}_{k-1} z_k - S_{k-1}(\mathbf{I}) \right| \mathbb{I} \left\{ \left| z_k^\top \mathbf{B}_{k-1} z_k - S_{k-1}(\mathbf{I}) \right| \geq \alpha_1 \right\} \right] \\ &= \mathcal{O}_{\mathbf{C}_x} \left(n^{-D} \cdot \sqrt{\gamma R_0(\mathbf{I})} \right) = \mathcal{O}_{\mathbf{C}_x, D} \left(n^{-D} \cdot \alpha_1 \right), \\ &\mathbb{E}_{k-1} \left[\left| z_k^\top \mathbf{B}_{k-1} \mathbf{Q} \mathbf{B}_{k-1} z_k - \text{Tr}(\mathbf{Q} \mathbf{B}_{k-1}^2) \right| \mathbb{I} \left\{ \left| z_k^\top \mathbf{B}_{k-1} \mathbf{Q} \mathbf{B}_{k-1} z_k - \text{Tr}(\mathbf{Q} \mathbf{B}_{k-1}^2) \right| \geq \alpha_2 \right\} \right] \\ &= \mathcal{O}_{\mathbf{C}_x} \left(n^{-D} \cdot \sqrt{\gamma^3 R_0(\mathbf{Q})} \right) = \mathcal{O}_{\mathbf{C}_x, D} \left(n^{-D} \cdot \alpha_2 \right). \end{aligned}$$

We then proceed to bound the term $p_{k,k}(T_E, T_E, \mathbf{Q}) - p_{k,k}(T_F, T_E, \mathbf{Q})$. The proof of the next lemma is in Appendix C.4.

Lemma 7.4. *Under Assumption 1, for any positive integer D , there exists a constant $C_\beta = C_\beta(\mathbf{C}_x, D) > 0$ such that the following holds. Consider $\alpha_1, \alpha_2, \gamma$ as defined in Lemma 7.3, and set β_1, β_2 by*

$$\begin{aligned} \beta_1 &= C_\beta \left(\sqrt{n \log n} \cdot \frac{\alpha_1 \gamma R_0(\mathbf{Q}) + \alpha_2 (1 + R_0(\mathbf{I}))}{1 + R_0(\mathbf{I})^2} + n \cdot \left\{ \frac{\gamma^2 R_0(\mathbf{Q}) + \alpha_1 \alpha_2}{1 + R_0(\mathbf{I})^2} + \frac{\alpha_1^2 \gamma R_0(\mathbf{Q})}{1 + R_0(\mathbf{I})^3} \right\} + \frac{\gamma R_0(\mathbf{Q})}{1 + R_0(\mathbf{I})} \right), \\ \beta_2 &= \frac{C_\beta n \beta_1}{1 + R_0(\mathbf{I})^2}. \end{aligned}$$

If $\alpha_1 \leq R_0(\mathbf{I})/4$, $\beta_1 \leq R_0(\mathbf{Q})/4$, $\beta_2 \leq \mu/2$ and $n^{-D} = \mathcal{O}(\alpha_1 / (1 + R_0(\mathbf{I})))$, then for all $1 \leq k \leq n+1$ and $n = \Omega_D(1)$,

$$p_{k,k}(T_E, T_E, \mathbf{Q}) - p_{k,k}(T_F, T_E, \mathbf{Q}) = \mathcal{O}(n^{-D}).$$

Applying Lemmas 7.3 and 7.4 to Eq. (50) (note that we can take $\mathbf{Q} = \mathbf{I}$), we have shown that $1 - \mathbb{P}(T_F(\mathbf{Q}) = n+1, T_F(\mathbf{I}) = n+1) = \mathbb{P}(T_F(\mathbf{Q}) \geq 0, T_F(\mathbf{I}) \geq 0) - \mathbb{P}(T_F(\mathbf{Q}) = n+1, T_F(\mathbf{I}) = n+1) = \mathcal{O}(n^{-D+1})$,

which implies by choosing the parameter Δ given by the above lemmas, with probability $1 - \mathcal{O}(n^{-D+1})$

$$|R_n(\mathbf{Q}) - R_0(\mathbf{Q})| \leq \beta_1, \quad |\mu_n - \bar{\mu}_n| \leq \beta_2.$$

Therefore, since $\bar{\mu}_n = \mu$, $\mu_0 = \mu_\star(\lambda, \mu)$, and recalling the definition of $R_k(\mathbf{Q})$, cf. Eq. (44), we have

$$\begin{aligned}
|\mathcal{R}_n(\lambda, \mu; \mathbf{Q}) - \mathcal{R}_0(\lambda, \mu_\star(\lambda, \mu); \mathbf{Q})| &\leq |\mathcal{R}_n(\lambda, \mu; \mathbf{Q}) - \mathcal{R}_n(\lambda, \mu_n; \mathbf{Q})| + |\mathcal{R}_n(\lambda, \mu_n; \mathbf{Q}) - \mathcal{R}_0(\lambda, \mu_\star(\lambda, \mu); \mathbf{Q})| \\
&= \left| (\mu_n - \mu) \cdot \text{Tr} \left(\mathbf{Q} \boldsymbol{\Sigma}^{\frac{1}{2}} (\lambda \mathbf{I} + \mu \boldsymbol{\Sigma} + \mathbf{X}^\top \mathbf{X})^{-1} \boldsymbol{\Sigma}^{\frac{1}{2}} \mathbf{A}_n \right) \right| + |R_n(\mathbf{Q}) - R_0(\mathbf{Q})| \\
&\stackrel{(i)}{\leq} \beta_2 \left\| \boldsymbol{\Sigma}^{\frac{1}{2}} (\lambda \mathbf{I} + \mu \boldsymbol{\Sigma} + \mathbf{X}^\top \mathbf{X})^{-1} \boldsymbol{\Sigma}^{\frac{1}{2}} \right\| |R_n(\mathbf{Q}) + |R_n(\mathbf{Q}) - R_0(\mathbf{Q})| \\
&\stackrel{(ii)}{\leq} \gamma \beta_2 R_n(\mathbf{Q}) + |R_n(\mathbf{Q}) - R_0(\mathbf{Q})| \leq \gamma \beta_2 R_0(\mathbf{Q}) + (1 + \gamma \beta_2) |R_n(\mathbf{Q}) - R_0(\mathbf{Q})| \\
&\leq \gamma \beta_2 R_0(\mathbf{Q}) + \beta_1 (1 + \gamma \beta_2) \stackrel{(ii)}{\leq} \frac{5}{4} \gamma \beta_2 R_0(\mathbf{Q}) + \beta_1,
\end{aligned}$$

where in (ii) we used Lemma 7.2; in (iii) we used the fact that $\beta_1 \leq R_0(\mathbf{Q})/4$ by assumption. We explain the inequality in (i) more carefully as it is less evident. Denoting by $\mathbf{B} = \boldsymbol{\Sigma}^{\frac{1}{2}} (\lambda \mathbf{I} + \mu \boldsymbol{\Sigma} + \mathbf{X}^\top \mathbf{X})^{-1} \boldsymbol{\Sigma}^{\frac{1}{2}}$, we first show \mathbf{B} and \mathbf{A}_n commute. Clearly commutativity holds if $\mu = \mu_n$, otherwise we have

$$\begin{aligned}
\mathbf{B} \mathbf{A}_n &= \boldsymbol{\Sigma}^{\frac{1}{2}} (\lambda \mathbf{I} + \mu \boldsymbol{\Sigma} + \mathbf{X}^\top \mathbf{X})^{-1} \boldsymbol{\Sigma}^{\frac{1}{2}} \cdot \boldsymbol{\Sigma}^{\frac{1}{2}} (\lambda \mathbf{I} + \mu_n \boldsymbol{\Sigma} + \mathbf{X}^\top \mathbf{X})^{-1} \boldsymbol{\Sigma}^{\frac{1}{2}} \\
&= (\mu_n - \mu)^{-1} \boldsymbol{\Sigma}^{\frac{1}{2}} \left\{ (\lambda \mathbf{I} + \mu \boldsymbol{\Sigma} + \mathbf{X}^\top \mathbf{X})^{-1} - (\lambda \mathbf{I} + \mu_n \boldsymbol{\Sigma} + \mathbf{X}^\top \mathbf{X})^{-1} \right\} \boldsymbol{\Sigma}^{\frac{1}{2}} \\
&= (\mu_n - \mu)^{-1} (\mathbf{B} - \mathbf{A}_n) = \mathbf{A}_n \mathbf{B}.
\end{aligned}$$

Noting that \mathbf{B} and \mathbf{A}_n are both p.s.d. compact self-adjoint operators in Hilbert space, commutativity implies they can be simultaneously orthogonally diagonalized and that $\mathbf{B}^{\frac{1}{2}}$ and $\mathbf{A}_n^{\frac{1}{2}}$ also commute. Consequently, combined with the fact that $\text{Tr}(\mathbf{A}_n \mathbf{C}) \leq \|\mathbf{A}_n\| \text{Tr}(\mathbf{C})$ for any p.s.d. matrix \mathbf{C} , we have (i) from

$$\begin{aligned}
\text{Tr} \left(\mathbf{Q} \boldsymbol{\Sigma}^{\frac{1}{2}} (\lambda \mathbf{I} + \mu \boldsymbol{\Sigma} + \mathbf{X}^\top \mathbf{X})^{-1} \boldsymbol{\Sigma}^{\frac{1}{2}} \mathbf{A}_n \right) &= \text{Tr}(\mathbf{Q} \mathbf{B} \mathbf{A}_n) = \text{Tr} \left(\mathbf{B} \cdot \mathbf{A}_n^{\frac{1}{2}} \mathbf{Q} \mathbf{A}_n^{\frac{1}{2}} \right) \leq \|\mathbf{B}\| \text{Tr}(\mathbf{Q} \mathbf{A}_n) \\
&= \left\| \boldsymbol{\Sigma}^{\frac{1}{2}} (\lambda \mathbf{I} + \mu \boldsymbol{\Sigma} + \mathbf{X}^\top \mathbf{X})^{-1} \boldsymbol{\Sigma}^{\frac{1}{2}} \right\| |R_n(\mathbf{Q})|.
\end{aligned}$$

We therefore proved the following. If $\beta_1 \leq R_0(\mathbf{Q})/4$ and $n^{-D} = \mathcal{O}(\alpha_1/(1 + R_0(\mathbf{I})))$, then

$$\beta_2 \leq \mu/2 \quad \Rightarrow \quad |\mathcal{R}_n(\lambda, \mu; \mathbf{Q}) - \mathcal{R}_0(\lambda, \mu_\star(\lambda, \mu); \mathbf{Q})| = \mathcal{O}(\gamma \beta_2 R_0(\mathbf{Q}) + \beta_1). \quad (51)$$

To remove the condition $\beta_2 \leq \mu/2$, we use the following estimate, proven in Appendix C.5.

Lemma 7.5. *Under Assumption 1, consider the parameter tuple $\Delta = (\alpha_1, \alpha_2, \beta_1, \beta_2, \gamma)$ defined in Lemmas 7.3 and 7.4. If $\alpha_1 \leq R_0(\mathbf{I})/8$, $\beta_1 \leq R_0(\mathbf{Q})/64$, $\gamma \beta_2 (1 + R_0(\mathbf{I})) \leq 1/64$, $n^{-D} = \mathcal{O}(\alpha_1/(1 + R_0(\mathbf{I})))$ and $\beta_2 > \mu/2$, then we have*

$$|\mathcal{R}_n(\lambda, \mu; \mathbf{Q}) - \mathcal{R}_0(\lambda, \mu_\star(\lambda, \mu); \mathbf{Q})| = \mathcal{O}(\gamma \beta_2 (1 + R_0(\mathbf{I})) R_0(\mathbf{Q}) + \beta_1) \quad (52)$$

Combining Eqs. (51) and (52), the proof is complete.

Acknowledgements

This work was supported by the NSF through award DMS-2031883, the Simons Foundation through Award 814639 for the Collaboration on the Theoretical Foundations of Deep Learning, the NSF grant

CCF-2006489, the ONR grant N00014-18-1-2729, and a grant from Eric and Wendy Schmidt at the Institute for Advanced Studies. C. Cheng is supported by the William R. Hewlett Stanford graduate fellowship.

Part of this work was carried out while A. Montanari was on partial leave from Stanford and a Chief Scientist at Ndata Inc dba Project N. The present research is unrelated to A. Montanari's activity while on leave.

References

- [Ada15] Radoslaw Adamczak, *A note on the Hanson-Wright inequality for random vectors with dependencies*, Electronic Communications in Probability **20** (2015), 1–13.
- [AEK⁺14] Bloemendal Alex, László Erdős, Antti Knowles, Horng-Tzer Yau, and Jun Yin, *Isotropic local laws for sample covariance and generalized wigner matrices*, Electronic Journal of Probability **19** (2014), 1–53.
- [Apo00] Tom M Apostol, *Calculating higher derivatives of inverses*, The American Mathematical Monthly **107** (2000), no. 8, 738–741.
- [ASS20] Madhu S Advani, Andrew M Saxe, and Haim Sompolinsky, *High-dimensional dynamics of generalization error in neural networks*, Neural Networks **132** (2020), 428–446.
- [BBEKY13] Derek Bean, Peter J Bickel, Noureddine El Karoui, and Bin Yu, *Optimal m -estimation in high-dimensional regression*, Proceedings of the National Academy of Sciences **110** (2013), no. 36, 14563–14568.
- [Bel21] Mikhail Belkin, *Fit without fear: remarkable mathematical phenomena of deep learning through the prism of interpolation*, Acta Numerica **30** (2021), 203–248.
- [BGL⁺14] Dominique Bakry, Ivan Gentil, Michel Ledoux, et al., *Analysis and geometry of markov diffusion operators*, vol. 103, Springer, 2014.
- [BKM⁺19] Jean Barbier, Florent Krzakala, Nicolas Macris, Léo Miolane, and Lenka Zdeborová, *Optimal errors and phase transitions in high-dimensional generalized linear models*, Proceedings of the National Academy of Sciences **116** (2019), no. 12, 5451–5460.
- [BLLT20] Peter L Bartlett, Philip M Long, Gábor Lugosi, and Alexander Tsigler, *Benign overfitting in linear regression*, Proceedings of the National Academy of Sciences **117** (2020), no. 48, 30063–30070.
- [BLM13] Stéphane Boucheron, Gábor Lugosi, and Pascal Massart, *Concentration inequalities: A nonasymptotic theory of independence*, Oxford University Press, 2013.
- [BM11] Mohsen Bayati and Andrea Montanari, *The lasso risk for gaussian matrices*, IEEE Transactions on Information Theory **58** (2011), no. 4, 1997–2017.
- [BMR21] Peter L Bartlett, Andrea Montanari, and Alexander Rakhlin, *Deep learning: a statistical viewpoint*, Acta numerica **30** (2021), 87–201.
- [Bré11] Haim Brézis, *Functional analysis, sobolev spaces and partial differential equations*, vol. 2, Springer, 2011.

- [BY08] Zhi-Dong Bai and Yong-Qua Yin, *Limit of the smallest eigenvalue of a large dimensional sample covariance matrix*, Advances In Statistics, World Scientific, 2008, pp. 108–127.
- [CM22] Michael Celentano and Andrea Montanari, *Fundamental barriers to high-dimensional regression with convex penalties*, The Annals of Statistics **50** (2022), no. 1, 170–196.
- [CMW20] Michael Celentano, Andrea Montanari, and Yuting Wei, *The lasso with general gaussian designs with applications to hypothesis testing*, arXiv:2007.13716 (2020).
- [CT05] Emmanuel J Candes and Terence Tao, *Decoding by linear programming*, IEEE transactions on information theory **51** (2005), no. 12, 4203–4215.
- [DET05] David L Donoho, Michael Elad, and Vladimir N Temlyakov, *Stable recovery of sparse overcomplete representations in the presence of noise*, IEEE Transactions on information theory **52** (2005), no. 1, 6–18.
- [Dic16] Lee H. Dicker, *Ridge regression and asymptotic minimax estimation over spheres of growing dimension*, Bernoulli **22** (2016), no. 1, 1–37.
- [DJM13] David L Donoho, Iain Johnstone, and Andrea Montanari, *Accurate prediction of phase transitions in compressed sensing via a connection to minimax denoising*, IEEE transactions on information theory **59** (2013), no. 6, 3396–3433.
- [DM16] David Donoho and Andrea Montanari, *High dimensional robust m -estimation: Asymptotic variance via approximate message passing*, Probability Theory and Related Fields **166** (2016), no. 3, 935–969.
- [DW18] Edgar Dobriban and Stefan Wager, *High-dimensional asymptotics of prediction: ridge regression and classification*, Annals of Statistics **46** (2018), no. 1, 247–279.
- [EK18] Noureddine El Karoui, *On the impact of predictor geometry on the performance on high-dimensional ridge-regularized generalized robust regression estimators*, Probability Theory and Related Fields **170** (2018), no. 1, 95–175.
- [EKBB⁺13] Noureddine El Karoui, Derek Bean, Peter J Bickel, Chinghay Lim, and Bin Yu, *On robust regression with high-dimensional predictors*, Proceedings of the National Academy of Sciences **110** (2013), no. 36, 14557–14562.
- [GS73] Janos Galambos and Eugene Seneta, *Regularly varying sequences*, Proceedings of the American Mathematical Society **41** (1973), no. 1, 110–116.
- [HLCH19] Lu-Jing Huang, Yin-Ting Liao, Lo-Bin Chang, and Chii-Ruey Hwang, *The smallest eigenvalues of random kernel matrices: Asymptotic results on the min kernel*, Statistics & Probability Letters **148** (2019), 23–29.
- [HMRT22] Trevor Hastie, Andrea Montanari, Saharon Rosset, and Ryan J Tibshirani, *Surprises in high-dimensional ridgeless least squares interpolation*, The Annals of Statistics **50** (2022), no. 2, 949–986.
- [Kar33] Jovan Karamata, *Sur un mode de croissance régulière. Théorèmes fondamentaux*, Bulletin de la Société Mathématique de France **61** (1933), 55–62.

- [KY17] Antti Knowles and Jun Yin, *Anisotropic local laws for random matrices*, Probability Theory and Related Fields **169** (2017), no. 1, 257–352.
- [KZSS21] 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** (2021), 20657–20668.
- [LP09] Anna Lytova and Leonid Pastur, *Central limit theorem for linear eigenvalue statistics of random matrices with independent entries*, The Annals of Probability **37** (2009), no. 5, 1778–1840.
- [MM21] Léo Miolane and Andrea Montanari, *The distribution of the lasso: Uniform control over sparse balls and adaptive parameter tuning*, The Annals of Statistics **49** (2021), no. 4, 2313–2335.
- [RMR21] Dominic Richards, Jaouad Mourtada, and Lorenzo Rosasco, *Asymptotics of ridge (less) regression under general source condition*, International Conference on Artificial Intelligence and Statistics, PMLR, 2021, pp. 3889–3897.
- [RV09] Mark Rudelson and Roman Vershynin, *Smallest singular value of a random rectangular matrix*, Communications on Pure and Applied Mathematics: A Journal Issued by the Courant Institute of Mathematical Sciences **62** (2009), no. 12, 1707–1739.
- [RV13] ———, *Hanson-Wright inequality and sub-Gaussian concentration*, Electronic Communications in Probability **18** (2013), 1–9.
- [T⁺15] Joel A Tropp et al., *An introduction to matrix concentration inequalities*, Foundations and Trends® in Machine Learning **8** (2015), no. 1-2, 1–230.
- [TAH18] Christos Thrampoulidis, Ehsan Abbasi, and Babak Hassibi, *Precise error analysis of regularized m -estimators in high dimensions*, IEEE Transactions on Information Theory **64** (2018), no. 8, 5592–5628.
- [TB20] Alexander Tsigler and Peter L Bartlett, *Benign overfitting in ridge regression*, arXiv:2009.14286 (2020).
- [Tib96] Robert Tibshirani, *Regression shrinkage and selection via the lasso*, Journal of the Royal Statistical Society: Series B (Methodological) **58** (1996), no. 1, 267–288.
- [TPT21] Hossein Taheri, Ramtin Pedarsani, and Christos Thrampoulidis, *Fundamental limits of ridge-regularized empirical risk minimization in high dimensions*, International Conference on Artificial Intelligence and Statistics, PMLR, 2021, pp. 2773–2781.
- [Tsy09] Alexandre B Tsybakov, *Introduction to nonparametric estimation*, Springer, 2009.
- [VdV00] Aad W Van der Vaart, *Asymptotic statistics*, vol. 3, Cambridge university press, 2000.
- [Ver12] R. Vershynin, *Introduction to the non-asymptotic analysis of random matrices*, Compressed Sensing: Theory and Applications (Y.C. Eldar and G. Kutyniok, eds.), Cambridge University Press, 2012, pp. 210–268.
- [WX20] Denny Wu and Ji Xu, *On the optimal weighted ℓ_2 regularization in overparameterized linear regression*, Advances in Neural Information Processing Systems **33** (2020), 10112–10123.

[Yas14] Pavel Yaskov, *Lower bounds on the smallest eigenvalue of a sample covariance matrix*, *Electronic Communications in Probability* **19** (2014), 1–10.

A Proof of Proposition 2.2

Since $\sigma_{k_*} \geq \lambda_* \geq \sigma_{k_*+1}$, we have

$$\begin{aligned} k_* + \frac{r_1(k_*)}{b_{k_*}} &= \sum_{l=1}^{k_*} \frac{\sigma_l}{\sigma_l} + \sum_{l=k_*+1}^d \frac{\sigma_l}{\sigma_{k_*}} \leq \sum_{l=1}^{k_*} \frac{\sigma_l + \lambda_*}{\sigma_l + \lambda_*} + \sum_{l=k_*+1}^d \frac{\sigma_l}{\lambda_*} \\ &\leq \sum_{l=1}^{k_*} \frac{2\sigma_l}{\sigma_l + \lambda_*} + \sum_{l=k_*+1}^d \frac{2\sigma_l}{\sigma_l + \lambda_*} = 2\text{Tr}(\boldsymbol{\Sigma}(\boldsymbol{\Sigma} + \lambda_*\mathbf{I})^{-1}) \leq 2n. \end{aligned}$$

Next we bound $V_n(\lambda)$. Recalling that $\text{Tr}(\boldsymbol{\Sigma}^2(\boldsymbol{\Sigma} + \lambda_*\mathbf{I})^{-2}) \leq n(1 - c_*^{-1})$, it then follows

$$\begin{aligned} V_n(\lambda) &= \frac{\tau^2 \text{Tr}(\boldsymbol{\Sigma}^2(\boldsymbol{\Sigma} + \lambda_*\mathbf{I})^{-2})}{n - \text{Tr}(\boldsymbol{\Sigma}^2(\boldsymbol{\Sigma} + \lambda_*\mathbf{I})^{-2})} \leq \frac{c_*\tau^2}{n} \cdot \left(\sum_{l=1}^{k_*} \frac{\sigma_l^2}{(\sigma_l + \lambda_*)^2} + \sum_{l=k_*+1}^d \frac{\sigma_l^2}{(\sigma_l + \lambda_*)^2} \right) \\ &\leq \frac{c_*\tau^2}{n} \cdot \left(k_* + \sum_{l=k_*+1}^d \frac{\sigma_l^2}{\lambda_*^2} \right) \leq c_*\tau^2 \left(\frac{k_*}{n} + \frac{r_2(k_*)}{n} \right) \stackrel{(i)}{\leq} c_*\tau^2 \left(\frac{k_*}{n} + \frac{4b_{k_*}^2 n}{\bar{r}(k_*)} \right), \end{aligned}$$

where in (i) we use the previous bound $r_1(k_*) \leq 2b_{k_*}n$. Finally, for the bias term, we have

$$\begin{aligned} B_n(\lambda) &= \frac{\lambda_*^2 \langle \boldsymbol{\beta}, (\boldsymbol{\Sigma} + \lambda_*\mathbf{I})^{-2} \boldsymbol{\Sigma} \boldsymbol{\beta} \rangle}{1 - n^{-1} \text{Tr}(\boldsymbol{\Sigma}^2(\boldsymbol{\Sigma} + \lambda_*\mathbf{I})^{-2})} \leq c_* \sum_{l=1}^d \frac{\lambda_*^2 \sigma_l}{(\sigma_l + \lambda_*)^2} \langle \boldsymbol{\beta}, \mathbf{v}_l \rangle^2 \\ &\leq c_* \left(\sum_{l=1}^{k_*} \lambda_*^2 \sigma_l^{-1} \langle \boldsymbol{\beta}, \mathbf{v}_l \rangle^2 + \sum_{l=k_*+1}^d \sigma_l \langle \boldsymbol{\beta}, \mathbf{v}_l \rangle^2 \right) \leq c_* \left(\sigma_{k_*}^2 \|\boldsymbol{\beta}_{\leq k_*}\|_{\boldsymbol{\Sigma}^{-1}}^2 + \|\boldsymbol{\beta}_{> k_*}\|_{\boldsymbol{\Sigma}}^2 \right). \end{aligned}$$

B Auxiliary lemmas

B.1 Proof of Lemma 6.1

The lemma follows by pure calculations.

Identities for $\mathcal{V}_{\mathbf{X}}(\lambda)$ and $\mathcal{B}_{\mathbf{X}}(\lambda)$. Substitute in Eq. (35), we have

$$\begin{aligned} \tau^2 \cdot \frac{\partial}{\partial \lambda} \mathcal{F}_n(\lambda, 0; \mathbf{I}) &= \tau^2 \cdot \frac{\partial}{\partial \lambda} \left(\lambda \text{Tr} \left(\boldsymbol{\Sigma}(\lambda \mathbf{I} + \mathbf{X}^\top \mathbf{X})^{-1} \right) \right) \\ &= \tau^2 \left\{ \text{Tr} \left(\boldsymbol{\Sigma}(\lambda \mathbf{I} + \mathbf{X}^\top \mathbf{X})^{-1} \right) - \lambda \text{Tr} \left(\boldsymbol{\Sigma}(\lambda \mathbf{I} + \mathbf{X}^\top \mathbf{X})^{-2} \right) \right\} \\ &= \tau^2 \text{Tr} \left(\boldsymbol{\Sigma} \mathbf{X}^\top \mathbf{X} (\lambda \mathbf{I} + \mathbf{X}^\top \mathbf{X})^{-2} \right) = \mathcal{V}_{\mathbf{X}}(\lambda), \end{aligned}$$

and similarly for the bias term

$$\begin{aligned} -\lambda \cdot \frac{\partial}{\partial \mu} \mathcal{F}_n(\lambda, 0; \boldsymbol{\theta} \boldsymbol{\theta}^\top) &= -\lambda \cdot \frac{\partial}{\partial \mu} \left(\lambda \text{Tr} \left(\boldsymbol{\Sigma}^{\frac{1}{2}} \boldsymbol{\theta} \boldsymbol{\theta}^\top \boldsymbol{\Sigma}^{\frac{1}{2}} (\lambda \mathbf{I} + \mu \boldsymbol{\Sigma} + \mathbf{X}^\top \mathbf{X})^{-1} \right) \right) \Big|_{\mu=0} \\ &= \lambda^2 \text{Tr} \left(\boldsymbol{\Sigma}^{\frac{1}{2}} \boldsymbol{\theta} \boldsymbol{\theta}^\top \boldsymbol{\Sigma}^{\frac{1}{2}} (\lambda \mathbf{I} + \mu \boldsymbol{\Sigma} + \mathbf{X}^\top \mathbf{X})^{-1} \boldsymbol{\Sigma} (\lambda \mathbf{I} + \mu \boldsymbol{\Sigma} + \mathbf{X}^\top \mathbf{X})^{-1} \right) \Big|_{\mu=0} \\ &= \lambda^2 \text{Tr} \left(\boldsymbol{\beta} \boldsymbol{\beta}^\top (\lambda \mathbf{I} + \mathbf{X}^\top \mathbf{X})^{-1} \boldsymbol{\Sigma} (\lambda \mathbf{I} + \mathbf{X}^\top \mathbf{X})^{-1} \right) = \mathcal{B}_{\mathbf{X}}(\lambda). \end{aligned}$$

Identities for $V_n(\lambda)$ and $B_n(\lambda)$. First we verify that $\mu_*(\lambda, 0) = \lambda/\lambda_*$. Set $\mu = 0$ in Eq. (36), we obtain

$$\mu_* = \frac{n}{1 + \mathcal{R}_0(\lambda, \mu_*; \mathbf{I})} = \frac{n}{1 + \text{Tr}(\boldsymbol{\Sigma}(\mu_*\boldsymbol{\Sigma} + \lambda\mathbf{I})^{-1})} = \frac{n}{1 + \mu_*^{-1}\text{Tr}(\boldsymbol{\Sigma}(\boldsymbol{\Sigma} + \frac{\lambda}{\mu_*}\mathbf{I})^{-1})},$$

and thus

$$n - \mu_* = \text{Tr}\left(\boldsymbol{\Sigma}\left(\boldsymbol{\Sigma} + \frac{\lambda}{\mu_*}\mathbf{I}\right)^{-1}\right),$$

which proves the claim comparing to Eq. (5). Further by (36), we can compute the derivatives

$$\frac{\partial}{\partial \lambda} \mu_*(\lambda, 0) = \frac{\text{Tr}(\boldsymbol{\Sigma}(\boldsymbol{\Sigma} + \lambda_*\mathbf{I})^{-2})}{n - \text{Tr}(\boldsymbol{\Sigma}^2(\boldsymbol{\Sigma} + \lambda_*\mathbf{I})^{-2})}, \quad (53a)$$

$$\frac{\partial}{\partial \mu} \mu_*(\lambda, 0) = \frac{n}{n - \text{Tr}(\boldsymbol{\Sigma}^2(\boldsymbol{\Sigma} + \lambda_*\mathbf{I})^{-2})}. \quad (53b)$$

We can then proceed to write

$$\begin{aligned} & \tau^2 \cdot \frac{\partial}{\partial \lambda} \mathcal{F}_0(\lambda, \mu_*(\lambda, 0); \mathbf{I}) \\ &= \tau^2 \cdot \frac{\partial}{\partial \lambda} (\lambda \text{Tr}(\boldsymbol{\Sigma}(\lambda\mathbf{I} + \mu_*\boldsymbol{\Sigma})^{-1})) \\ &= \tau^2 \left\{ \text{Tr}(\boldsymbol{\Sigma}(\lambda\mathbf{I} + \mu_*\boldsymbol{\Sigma})^{-1}) - \lambda \text{Tr}(\boldsymbol{\Sigma}(\lambda\mathbf{I} + \mu_*\boldsymbol{\Sigma})^{-2}) - \lambda \text{Tr}(\boldsymbol{\Sigma}^2(\lambda\mathbf{I} + \mu_*\boldsymbol{\Sigma})^{-2}) \cdot \frac{\partial}{\partial \lambda} \mu_*(\lambda, 0) \right\} \\ &= \frac{\tau^2}{\mu_*} \text{Tr}(\boldsymbol{\Sigma}^2(\boldsymbol{\Sigma} + \lambda_*\mathbf{I})^{-2}) \cdot \left(1 - \frac{\lambda_* \text{Tr}(\boldsymbol{\Sigma}(\boldsymbol{\Sigma} + \lambda_*\mathbf{I})^{-2})}{n - \text{Tr}(\boldsymbol{\Sigma}^2(\boldsymbol{\Sigma} + \lambda_*\mathbf{I})^{-2})} \right) \\ &= \frac{\tau^2}{\mu_*} \text{Tr}(\boldsymbol{\Sigma}^2(\boldsymbol{\Sigma} + \lambda_*\mathbf{I})^{-2}) \cdot \frac{n - \text{Tr}(\boldsymbol{\Sigma}(\boldsymbol{\Sigma} + \lambda_*\mathbf{I})^{-1})}{n - \text{Tr}(\boldsymbol{\Sigma}^2(\boldsymbol{\Sigma} + \lambda_*\mathbf{I})^{-2})} \\ &\stackrel{(i)}{=} \frac{\tau^2 \text{Tr}(\boldsymbol{\Sigma}^2(\boldsymbol{\Sigma} + \lambda_*\mathbf{I})^{-2})}{n - \text{Tr}(\boldsymbol{\Sigma}^2(\boldsymbol{\Sigma} + \lambda_*\mathbf{I})^{-2})} = V_n(\lambda), \end{aligned}$$

where in (i) we use Eq. (5) which implies $\mu_* = n - \text{Tr}(\boldsymbol{\Sigma}(\boldsymbol{\Sigma} + \lambda_*\mathbf{I})^{-1})$. For the bias we can compute

$$\begin{aligned} & -\lambda \cdot \frac{\partial}{\partial \mu} \mathcal{F}_0(\lambda, \mu_*(\lambda, 0); \boldsymbol{\theta}\boldsymbol{\theta}^\top) \\ &= \lambda^2 \text{Tr}\left(\boldsymbol{\Sigma}^{\frac{1}{2}}\boldsymbol{\theta}\boldsymbol{\theta}^\top\boldsymbol{\Sigma}^{\frac{1}{2}}(\lambda\mathbf{I} + \mu_*\boldsymbol{\Sigma})^{-1}\boldsymbol{\Sigma}(\lambda\mathbf{I} + \mu_*\boldsymbol{\Sigma})^{-1}\right) \cdot \frac{n}{n - \text{Tr}(\boldsymbol{\Sigma}^2(\boldsymbol{\Sigma} + \lambda_*\mathbf{I})^{-2})} \\ &= \frac{\lambda_*^2 \boldsymbol{\beta}^\top (\boldsymbol{\Sigma} + \lambda_*\mathbf{I})^{-2} \boldsymbol{\Sigma} \boldsymbol{\beta}}{1 - n^{-1} \text{Tr}(\boldsymbol{\Sigma}^2(\boldsymbol{\Sigma} + \lambda_*\mathbf{I})^{-2})} = B_n(\lambda). \end{aligned}$$

The proof is complete.

B.2 Proof of Lemma 6.2

The lemma is an analogue of [HMRT22, Lemma. 5], which requires two-sided differentiability around 0 and makes use of higher order central difference operators from numerical analysis. Here we apply

a more straightforward argument. For any $0 \leq j \leq k$, by Taylor expansion with Lagrange remainder, we can write

$$f(j\delta) = \sum_{l=0}^k j^l \cdot \frac{\delta^l}{j!} f^{(l)}(0) + j^{k+1} \cdot \frac{\delta^{k+1}}{(k+1)!} f^{(k+1)}(t_j),$$

for some $t_j \in [0, j\delta]$. We can write the $k+1$ equations in matrix form,

$$\underbrace{\begin{bmatrix} 1 & 0 & 0 & \cdots & 0 \\ 1 & 1 & 1 & \cdots & 1 \\ 1 & 2 & 4 & \cdots & 2^k \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & k & k^2 & \cdots & k^k \end{bmatrix}}_{:=\mathbf{V}_k} \begin{bmatrix} f(0) \\ f'(0)\delta \\ f''(0)\delta^2/2 \\ \vdots \\ f^{(k)}(0)\delta^k/k! \end{bmatrix} + \frac{\delta^{k+1}}{(k+1)!} \begin{bmatrix} 0 \\ f^{(k+1)}(t_1) \\ 2^{k+1}f^{(k+1)}(t_2) \\ \vdots \\ k^{k+1}f^{(k+1)}(t_k) \end{bmatrix} = \begin{bmatrix} f(0) \\ f(\delta) \\ f(2\delta) \\ \vdots \\ f(k\delta) \end{bmatrix}.$$

The Vandermonde matrix \mathbf{V}_k is invertible, and therefore we can write

$$\begin{bmatrix} f(0) \\ f'(0)\delta \\ f''(0)\delta^2/2 \\ \vdots \\ f^{(k)}(0)\delta^k/k! \end{bmatrix} = \mathbf{V}_k^{-1} \begin{bmatrix} f(0) \\ f(\delta) \\ f(2\delta) \\ \vdots \\ f(k\delta) \end{bmatrix} - \frac{\delta^{k+1}}{(k+1)!} \mathbf{V}_k^{-1} \begin{bmatrix} 0 \\ f^{(k+1)}(t_1) \\ 2^{k+1}f^{(k+1)}(t_2) \\ \vdots \\ k^{k+1}f^{(k+1)}(t_k) \end{bmatrix}.$$

Denote by $\|\mathbf{M}\|_\infty$ the ℓ_∞ -induced operator norm, we thus have

$$\begin{aligned} |f'(0)\delta| &\leq \left\| \begin{bmatrix} f(0) \\ f'(0)\delta \\ f''(0)\delta^2/2 \\ \vdots \\ f^{(k)}(0)\delta^k/k! \end{bmatrix} \right\|_\infty \leq \|\mathbf{V}_k^{-1}\|_\infty \cdot \left(\left\| \begin{bmatrix} f(0) \\ f(\delta) \\ f(2\delta) \\ \vdots \\ f(k\delta) \end{bmatrix} \right\|_\infty + \frac{\delta^{k+1}}{(k+1)!} \left\| \begin{bmatrix} 0 \\ f^{(k+1)}(t_1) \\ 2^{k+1}f^{(k+1)}(t_2) \\ \vdots \\ k^{k+1}f^{(k+1)}(t_k) \end{bmatrix} \right\|_\infty \right) \\ &= \mathcal{O}_k \left(\max_{0 \leq j \leq k} |f(j\delta)| + \sup_{t \in [0, k\delta]} |f^{(k+1)}(t)| \cdot \delta^{k+1} \right). \end{aligned}$$

Dividing δ from both sides completes the proof.

B.3 Proof of Lemma 6.3

Part I: Derivative w.r.t. λ . By Lemma 6.1,

$$\frac{\partial}{\partial \lambda} \mathcal{F}_n(\lambda, 0; \mathbf{I}) = \mathcal{V}_{\mathbf{X}}(\lambda)/\tau^2 = \text{Tr} \left(\boldsymbol{\Sigma} \mathbf{X}^\top \mathbf{X} (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I})^{-2} \right),$$

we can easily write out derivatives with respect to λ up to any order $k \geq 1$ as

$$\frac{\partial^k}{\partial \lambda^k} \mathcal{F}_n(\lambda, 0; \mathbf{I}) = \mathcal{O}_k \left(\text{Tr} \left(\boldsymbol{\Sigma} \mathbf{X}^\top \mathbf{X} (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I})^{-1-k} \right) \right),$$

and therefore

$$\begin{aligned} \left| \frac{\partial^k}{\partial \lambda^k} \mathcal{F}_n(\lambda, 0; \mathbf{I}) \right| &\leq \mathcal{O}_k \left(\frac{1}{\lambda^{k-1}} \text{Tr} \left(\boldsymbol{\Sigma} \mathbf{X}^\top \mathbf{X} (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I})^{-2} \right) \right) \leq \mathcal{O}_k \left(\frac{1}{\lambda^{k-1}} \text{Tr} \left(\boldsymbol{\Sigma} (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I})^{-1} \right) \right) \\ &= \mathcal{O}_k \left(\frac{\mathcal{F}_n(\lambda, 0; \mathbf{I})}{\lambda^k} \right). \end{aligned}$$

Part II: Derivative w.r.t. μ . We can directly compute that

$$\begin{aligned}
\left| \frac{\partial^k}{\partial \mu^k} \mathcal{F}_n(\lambda, \mu; \boldsymbol{\theta}\boldsymbol{\theta}^\top) \right| &= \left| \frac{\partial^k}{\partial \mu^k} \cdot \lambda \text{Tr} \left(\boldsymbol{\Sigma}^{\frac{1}{2}} \boldsymbol{\theta}\boldsymbol{\theta}^\top \boldsymbol{\Sigma}^{\frac{1}{2}} (\lambda \mathbf{I} + \mu \boldsymbol{\Sigma} + \mathbf{X}^\top \mathbf{X})^{-1} \right) \right| \\
&= \mathcal{O}_k \left(\lambda \text{Tr} \left(\boldsymbol{\Sigma}^{\frac{1}{2}} \boldsymbol{\theta}\boldsymbol{\theta}^\top \boldsymbol{\Sigma}^{\frac{1}{2}} \left((\lambda \mathbf{I} + \mu \boldsymbol{\Sigma} + \mathbf{X}^\top \mathbf{X})^{-1} \boldsymbol{\Sigma} \right)^k (\lambda \mathbf{I} + \mu \boldsymbol{\Sigma} + \mathbf{X}^\top \mathbf{X})^{-1} \right) \right) \\
&= \mathcal{O}_k \left(\lambda \text{Tr} \left(\boldsymbol{\theta}\boldsymbol{\theta}^\top \left(\boldsymbol{\Sigma}^{\frac{1}{2}} (\lambda \mathbf{I} + \mu \boldsymbol{\Sigma} + \mathbf{X}^\top \mathbf{X})^{-1} \boldsymbol{\Sigma}^{\frac{1}{2}} \right)^{k+1} \right) \right) \\
&\stackrel{(i)}{=} \mathcal{O}_k \left(\lambda^{1-k} \text{Tr} \left(\boldsymbol{\theta}\boldsymbol{\theta}^\top \boldsymbol{\Sigma}^{\frac{1}{2}} (\lambda \mathbf{I} + \mu \boldsymbol{\Sigma} + \mathbf{X}^\top \mathbf{X})^{-1} \boldsymbol{\Sigma}^{\frac{1}{2}} \right) \right) = \mathcal{O}_k \left(\frac{\mathcal{F}_n(\lambda, \mu; \boldsymbol{\theta}\boldsymbol{\theta}^\top)}{\lambda^k} \right),
\end{aligned}$$

where in (i) we use $\|\boldsymbol{\Sigma}\| = 1$.

B.4 Proof of Lemma 6.4

Part I: Derivative w.r.t. λ . Note that

$$\mathcal{F}_0(\lambda, \mu_\star(\lambda, 0); \mathbf{I}) = \lambda \text{Tr} \left(\boldsymbol{\Sigma} (\lambda \mathbf{I} + \mu_\star(\lambda, 0) \boldsymbol{\Sigma})^{-1} \right) = \lambda_\star \text{Tr} \left(\boldsymbol{\Sigma} (\boldsymbol{\Sigma} + \lambda_\star \mathbf{I})^{-1} \right).$$

Combining with the fixed-point equation (5) that determines λ_\star , we further get

$$\mathcal{F}_0(\lambda, \mu_\star(\lambda, 0); \mathbf{I}) = n\lambda_\star - \lambda.$$

Therefore, for all $k \geq 1$.

$$\frac{\partial}{\partial \lambda^k} \mathcal{F}_0(\lambda, \mu_\star(\lambda, 0); \mathbf{I}) = n \cdot \frac{\partial^k \lambda_\star}{\partial \lambda^k} - \mathbb{I}\{k = 1\}, \quad (54)$$

and it boils down to controlling higher order derivatives of λ_\star w.r.t. λ . Of course, we need to first show that we can actually write $\lambda_\star = \lambda_\star(\lambda)$ locally by implicit function theorem. Since

$$\lambda = \lambda_\star \cdot (n - \text{Tr}(\boldsymbol{\Sigma}(\boldsymbol{\Sigma} + \lambda_\star \mathbf{I})^{-1}))$$

which is clearly a increasing function of λ_\star on the right hand side, and thus $\partial \lambda / \partial \lambda_\star > 0$ and the implicit function theorem applies. To calculate the higher order derivative of the inverse function, we apply the formula for higher order derivatives of inverse function [Apo00]

$$\frac{\partial^k \lambda_\star}{\partial \lambda^k} = \left| \frac{\partial \lambda}{\partial \lambda_\star} \right|^{1-2k} \cdot \sum_{\substack{m_1+m_2+\dots+m_p=k-1 \\ m_1+2m_2+\dots+pm_p=2k-2}} \mathcal{O}_k \left(\prod_{l=1}^p \left(\frac{\partial^l \lambda}{\partial \lambda_\star^l} \right)^{m_l} \right). \quad (55)$$

To further upper bound the above display, we need a lower bound for the derivative $\partial \lambda / \partial \lambda_\star$ and upper bounds for higher order derivatives $\partial^l \lambda / \partial \lambda_\star^l$. Using the Leibniz rule, we can compute that

$$\begin{aligned}
\frac{\partial^l \lambda}{\partial \lambda_\star^l} &= \sum_{r=0}^l \binom{l}{r} \frac{\partial^r \lambda_\star}{\partial \lambda_\star^r} \cdot \frac{\partial^{l-r} (n - \text{Tr}(\boldsymbol{\Sigma}(\boldsymbol{\Sigma} + \lambda_\star \mathbf{I})^{-1}))}{\partial \lambda_\star^{l-r}} \\
&= \lambda_\star \cdot \frac{\partial^l (n - \text{Tr}(\boldsymbol{\Sigma}(\boldsymbol{\Sigma} + \lambda_\star \mathbf{I})^{-1}))}{\partial \lambda_\star^l} + l \cdot \frac{\partial^{l-1} (n - \text{Tr}(\boldsymbol{\Sigma}(\boldsymbol{\Sigma} + \lambda_\star \mathbf{I})^{-1}))}{\partial \lambda_\star^{l-1}}.
\end{aligned}$$

For $l = 1$, since

$$\begin{aligned}\frac{\partial \lambda}{\partial \lambda_\star} &= \lambda_\star \cdot \text{Tr}(\boldsymbol{\Sigma}(\boldsymbol{\Sigma} + \lambda_\star \mathbf{I})^{-2}) + n - \text{Tr}(\boldsymbol{\Sigma}(\boldsymbol{\Sigma} + \lambda_\star \mathbf{I})^{-1}) = n - \text{Tr}(\boldsymbol{\Sigma}^2(\boldsymbol{\Sigma} + \lambda_\star \mathbf{I})^{-2}) \\ &\geq n - \text{Tr}(\boldsymbol{\Sigma}(\boldsymbol{\Sigma} + \lambda_\star \mathbf{I})^{-1}) = \frac{\lambda}{\lambda_\star} \geq n\kappa,\end{aligned}$$

we have $n\kappa \leq \partial \lambda / \partial \lambda_\star \leq n$. When $l \geq 2$, we get

$$\begin{aligned}\frac{\partial^l \lambda}{\partial \lambda_\star^l} &= (-1)^{l-1} l! \cdot \lambda_\star \text{Tr}(\boldsymbol{\Sigma}(\boldsymbol{\Sigma} + \lambda_\star \mathbf{I})^{-l-1}) + (-1)^{l-2} l! \cdot \text{Tr}(\boldsymbol{\Sigma}(\boldsymbol{\Sigma} + \lambda_\star \mathbf{I})^{-l}) = (-1)^{l-2} l! \cdot \text{Tr}(\boldsymbol{\Sigma}^2(\boldsymbol{\Sigma} + \lambda_\star \mathbf{I})^{-l-1}) \\ &= \mathcal{O}_l\left(\left\|(\boldsymbol{\Sigma} + \lambda_\star \mathbf{I})^{-l+1}\right\| \cdot \text{Tr}(\boldsymbol{\Sigma}^2(\boldsymbol{\Sigma} + \lambda_\star \mathbf{I})^{-2})\right) = \mathcal{O}_l\left(\frac{n}{\lambda_\star^{l-1}}\right).\end{aligned}$$

Substituting the above displays into Eq. (55) yields

$$\frac{\partial^k \lambda_\star}{\partial \lambda^k} = \left(\frac{1}{(n\kappa)^{2k-1}}\right) \cdot \sum_{\substack{m_1+m_2+\dots+m_p=k-1 \\ m_1+2m_2+\dots+pm_p=2k-2}} \mathcal{O}_k\left(\prod_{l=1}^p \mathcal{O}_l\left(\frac{n^{m_l}}{\lambda_\star^{lm_l-m_l}}\right)\right) = \mathcal{O}_k\left(\frac{1}{n^k \lambda_\star^{k-1} \cdot \kappa^{2k-1}}\right)$$

Taken collectively with Eq. (54) and $\mathcal{F}_0(\lambda, \mu_\star(\lambda, 0); \mathbf{I}) = \lambda_\star \text{Tr}(\boldsymbol{\Sigma}(\boldsymbol{\Sigma} + \lambda_\star \mathbf{I})^{-1}) \geq \kappa n \lambda_\star$, we obtain for all $k \geq 2$,

$$\left|\frac{\partial^k}{\partial \lambda^k} \mathcal{F}_0(\lambda, \mu_\star(\lambda, 0); \mathbf{I})\right| = n \left|\frac{\partial^k \lambda_\star}{\partial \lambda^k}\right| = \mathcal{O}_k\left(\frac{\mathcal{F}_0(\lambda, \mu_\star(\lambda, 0); \mathbf{I})}{n^k \lambda_\star^k \kappa^{2k}}\right) = \mathcal{O}_k\left(\frac{\mathcal{F}_0(\lambda, \mu_\star(\lambda, 0); \mathbf{I})}{\lambda^k \kappa^{2k}}\right),$$

where we use Assumption (21) again for the final bound. This is also valid for $k = 1$ as

$$\begin{aligned}\left|\frac{\partial}{\partial \lambda} \mathcal{F}_0(\lambda, \mu_\star(\lambda, 0); \mathbf{I})\right| &= n \left|\frac{\partial \lambda_\star}{\partial \lambda}\right| + 1 = \mathcal{O}\left(\frac{\mathcal{F}_0(\lambda, \mu_\star(\lambda, 0); \mathbf{I})}{\lambda \kappa^2}\right) + 1 = \mathcal{O}\left(\frac{n\lambda_\star - \lambda}{\lambda \kappa^2}\right) + 1 \\ &= \mathcal{O}\left(\frac{n\lambda_\star - \lambda}{\lambda \kappa^2}\right) = \mathcal{O}\left(\frac{\mathcal{F}_0(\lambda, \mu_\star(\lambda, 0); \mathbf{I})}{\lambda \kappa^2}\right),\end{aligned}$$

where we use $\mathcal{F}_0(\lambda, \mu_\star(\lambda, 0); \mathbf{I})/\lambda = n\lambda_\star/\lambda - 1$ and

$$\frac{n\lambda_\star - \lambda}{\lambda \kappa^2} \geq ((1 - \kappa)^{-1} - 1) \kappa^{-2} \geq \kappa^{-1} \geq 1.$$

Part II: Derivative w.r.t. μ . Now we fix λ and allow μ be take nonzero values. We will also use the shorthand $\mu_\star = \mu_\star(\lambda, \mu)$. Similar to the previous part, we apply Faà di Bruno's formula to \mathcal{F}_0 and bound

$$\left|\frac{\partial^k}{\partial \mu^k} \mathcal{F}_0(\lambda, \mu_\star; \boldsymbol{\theta} \boldsymbol{\theta}^\top)\right| = \sum_{m_1+2m_2+\dots+pm_p=k} \mathcal{O}_k\left(\frac{\partial^{m_1+\dots+m_p}}{\partial \mu_\star^{m_1+\dots+m_p}} \mathcal{F}_0(\lambda, \mu_\star; \boldsymbol{\theta} \boldsymbol{\theta}^\top) \cdot \prod_{l=1}^p \left(\frac{\partial^l \mu_\star}{\partial \mu^l}\right)^{m_l}\right). \quad (56)$$

For any $1 \leq l \leq k-1$, we have

$$\left|\frac{\partial^l}{\partial \mu_\star^l} \mathcal{F}_0(\lambda, \mu_\star; \boldsymbol{\theta} \boldsymbol{\theta}^\top)\right| = \mathcal{O}_l\left(\lambda \text{Tr}\left(\boldsymbol{\Sigma}^{\frac{1}{2}} \boldsymbol{\theta} \boldsymbol{\theta}^\top \boldsymbol{\Sigma}^{\frac{1}{2}} \cdot \boldsymbol{\Sigma}^l (\lambda \mathbf{I} + \mu_\star \boldsymbol{\Sigma})^{-1-l}\right)\right) = \mathcal{O}_l\left(\frac{\mathcal{F}_0(\lambda, \mu_\star; \boldsymbol{\theta} \boldsymbol{\theta}^\top)}{\mu_\star^l}\right).$$

To bound higher order derivatives $\partial^l \mu_\star / \partial \mu^l$, we apply again the formula for higher order derivatives of inverse function. Of course, this would first require showing the existence of inverse function

by implicit function theorem, which will be evident as we will provide a lower bound for $|\partial\mu/\partial\mu_\star|$ below. By [Apo00], we have for all $1 \leq l \leq k-1$,

$$\left| \frac{\partial^l \mu_\star}{\partial \mu^l} \right| = \left| \frac{\partial \mu}{\partial \mu_\star} \right|^{1-2l} \cdot \sum_{\substack{m_1+m_2+\dots+m_p=l-1 \\ m_1+2m_2+\dots+pm_p=2l-2}} \mathcal{O}_l \left(\prod_{r=1}^p \left(\frac{\partial^r \mu}{\partial \mu_\star^r} \right)^{m_r} \right). \quad (57)$$

This is a more manageable formula as we can explicitly write μ as a function of μ_\star

$$\mu = \mu_\star - \frac{n}{1 + \mathcal{R}_0(\lambda, \mu_\star; \mathbf{I})} = \mu_\star - \frac{n}{1 + \text{Tr}(\boldsymbol{\Sigma}(\lambda \mathbf{I} + \mu_\star \boldsymbol{\Sigma})^{-1})}.$$

We can compute the first order derivative as

$$\frac{\partial \mu}{\partial \mu_\star} = 1 - \frac{n \text{Tr}(\boldsymbol{\Sigma}^2(\lambda \mathbf{I} + \mu_\star \boldsymbol{\Sigma})^{-2})}{(1 + \text{Tr}(\boldsymbol{\Sigma}(\lambda \mathbf{I} + \mu_\star \boldsymbol{\Sigma})^{-1}))^2} = 1 - \frac{(\mu_\star - \mu) \cdot \text{Tr}(\boldsymbol{\Sigma}^2(\lambda \mathbf{I} + \mu_\star \boldsymbol{\Sigma})^{-2})}{1 + \text{Tr}(\boldsymbol{\Sigma}(\lambda \mathbf{I} + \mu_\star \boldsymbol{\Sigma})^{-1})},$$

which, together with $0 \leq \mu \leq \mu_\star/2$, implies a lower bound

$$\begin{aligned} \frac{\partial \mu}{\partial \mu_\star} &\geq 1 - \frac{\mu_\star \text{Tr}(\boldsymbol{\Sigma}^2(\lambda \mathbf{I} + \mu_\star \boldsymbol{\Sigma})^{-2})}{1 + \text{Tr}(\boldsymbol{\Sigma}(\lambda \mathbf{I} + \mu_\star \boldsymbol{\Sigma})^{-1})} = \frac{1 + \lambda \text{Tr}(\boldsymbol{\Sigma}(\lambda \mathbf{I} + \mu_\star \boldsymbol{\Sigma})^{-2})}{1 + \text{Tr}(\boldsymbol{\Sigma}(\lambda \mathbf{I} + \mu_\star \boldsymbol{\Sigma})^{-1})} \\ &\geq \frac{1}{1 + \text{Tr}(\boldsymbol{\Sigma}(\lambda \mathbf{I} + \mu_\star \boldsymbol{\Sigma})^{-1})} = \frac{\mu_\star - \mu}{n} \geq \frac{\mu_\star}{2n}. \end{aligned}$$

To further bound higher order derivatives, we again appeal to Faà di Bruno's formula. Use the shorthand $\mathcal{R}_0 = \mathcal{R}_0(\lambda, \mu_\star; \mathbf{I})$, we have for all $r \geq 1$,

$$\left| \frac{\partial^r \mu}{\partial \mu_\star^r} \right| = \sum_{m_1+2m_2+\dots+pm_p=r} \mathcal{O}_r \left(\frac{\partial^{m_1+\dots+m_p}}{\partial \mathcal{R}_0^{m_1+\dots+m_p}} \frac{n}{1 + \mathcal{R}_0} \cdot \prod_{s=1}^p \left(\frac{\partial^s \mathcal{R}_0}{\partial \mu_\star^s} \right)^{m_s} \right).$$

Making use of the following two bounds,

$$\begin{aligned} \frac{\partial^s}{\partial \mathcal{R}_0^s} \frac{n}{1 + \mathcal{R}_0} &= \mathcal{O}_s \left(\frac{n}{(1 + \mathcal{R}_0)^{s+1}} \right) = \mathcal{O}_s \left(\frac{\mu_\star}{(1 + \mathcal{R}_0)^s} \right), \\ \frac{\partial^s \mathcal{R}_0}{\partial \mu_\star^s} &= \mathcal{O}_s (\text{Tr}(\boldsymbol{\Sigma}^{s+1}(\lambda \mathbf{I} + \mu_\star \boldsymbol{\Sigma})^{-s-1})) = \mathcal{O}_s \left(\frac{\mathcal{R}_0}{\mu_\star^s} \right), \end{aligned}$$

we can further obtain

$$\left| \frac{\partial^r \mu}{\partial \mu_\star^r} \right| = \sum_{m_1+2m_2+\dots+pm_p=r} \mathcal{O}_r \left(\frac{\mu_\star}{(1 + \mathcal{R}_0)^{m_1+\dots+m_p}} \cdot \prod_{s=1}^p \frac{\mathcal{R}_0^{m_s}}{\mu_\star^{sm_s}} \right) = \mathcal{O}_r \left(\frac{1}{\mu_\star^{r-1}} \right).$$

Taking the above displays into Eq. (57) and use the condition $\mu_\star/n \geq \kappa$, we have

$$\left| \frac{\partial^l \mu_\star}{\partial \mu^l} \right| = \mathcal{O}_l \left(\frac{1}{\kappa^{2l-1}} \right) \cdot \sum_{\substack{m_1+m_2+\dots+m_p=l-1 \\ m_1+2m_2+\dots+pm_p=2l-2}} \mathcal{O}_l \left(\prod_{r=1}^p \mathcal{O}_r \left(\frac{1}{\mu_\star^{rm_r-m_r}} \right) \right) = \mathcal{O}_l \left(\frac{1}{\mu_\star^{l-1} \kappa^{2l-1}} \right).$$

Finally, taking the above display back into Eq. (56) yields

$$\begin{aligned} \left| \frac{\partial^k}{\partial \mu^k} \mathcal{F}_0(\lambda, \mu_\star; \boldsymbol{\theta} \boldsymbol{\theta}^\top) \right| &= \sum_{m_1+2m_2+\dots+pm_p=k} \mathcal{O}_k \left(\frac{\mathcal{F}_0(\lambda, \mu_\star; \boldsymbol{\theta} \boldsymbol{\theta}^\top)}{\mu_\star^{m_1+\dots+m_p}} \cdot \prod_{l=1}^p \mathcal{O}_l \left(\frac{1}{\mu_\star^{lm_l-m_l} \kappa^{2lm_l-m_l}} \right) \right) \\ &= \mathcal{O}_k \left(\frac{\mathcal{F}_0(\lambda, \mu_\star; \boldsymbol{\theta} \boldsymbol{\theta}^\top)}{\mu_\star^k \kappa^{2k}} \right). \end{aligned}$$

B.5 Proof of Lemma 6.6

First we show $\mu_*(\lambda, \mu)$ is increasing in μ when $\mu \geq 0$. To this end, we consider the function

$$f(t) = t - \frac{n}{1 + \mathcal{R}_0(\lambda, t; \mathbf{I})}.$$

By Eq. (36), we have $f(\mu_*(\lambda, \mu)) = \mu$ for all $\mu \geq 0$. Further, we prove $f(t)$ is increasing in $[\mu_*(\lambda, 0), \infty)$. We write

$$f'(t) = 1 - \frac{n \text{Tr}(\boldsymbol{\Sigma}^2(\lambda \mathbf{I} + t \boldsymbol{\Sigma})^{-2})}{(1 + \text{Tr}(\boldsymbol{\Sigma}(\lambda \mathbf{I} + t \boldsymbol{\Sigma})^{-1}))^2} \stackrel{(i)}{=} 1 - \frac{(t - f(t)) \cdot \text{Tr}(\boldsymbol{\Sigma}^2(\lambda \mathbf{I} + t \boldsymbol{\Sigma})^{-2})}{1 + \text{Tr}(\boldsymbol{\Sigma}(\lambda \mathbf{I} + t \boldsymbol{\Sigma})^{-1})},$$

where in (i) we use that $t - f(t) = n/(1 + \mathcal{R}_0(\lambda, t; \mathbf{I}))$. Define

$$g(t) := \frac{\text{Tr}(\boldsymbol{\Sigma}^2(\lambda \mathbf{I} + t \boldsymbol{\Sigma})^{-2})}{1 + \text{Tr}(\boldsymbol{\Sigma}(\lambda \mathbf{I} + t \boldsymbol{\Sigma})^{-1})},$$

we have

$$f'(t) - f(t)g(t) = 1 - \frac{t \text{Tr}(\boldsymbol{\Sigma}^2(\lambda \mathbf{I} + t \boldsymbol{\Sigma})^{-2})}{1 + \text{Tr}(\boldsymbol{\Sigma}(\lambda \mathbf{I} + t \boldsymbol{\Sigma})^{-1})} = \frac{1 + \lambda \text{Tr}(\boldsymbol{\Sigma}(\lambda \mathbf{I} + t \boldsymbol{\Sigma})^{-2})}{1 + \text{Tr}(\boldsymbol{\Sigma}(\lambda \mathbf{I} + t \boldsymbol{\Sigma})^{-1})} > 0,$$

and therefore $e^{-g(t)}f(t)$ is increasing. As $f(\mu_*(\lambda, 0)) = 0$ (cf. Eq. (36)), we must have $f(t) \geq 0$ for all $t \geq \mu_*(\lambda, 0)$. Substituting back into the above display with $g(t) \geq 0$ yields

$$f'(t) \geq f'(t) - f(t)g(t) > 0.$$

We then proceed to show a sufficient condition for $\mu_*(\lambda, \mu) \leq (1 - \kappa/2)^{-1}\mu_*(\lambda, 0)$ is $0 \leq \mu \leq n\kappa^3/2$ under Assumption (21). Provided with monotonicity of $f(t)$, the desired condition $\mu_*(\lambda, \mu) \leq (1 - \kappa/2)^{-1}\mu_*(\lambda, 0)$ is essentially equivalent to $\mu = f(\mu_*(\lambda, \mu)) \leq f((1 - \kappa/2)^{-1}\mu_*(\lambda, 0))$. Together with $\mu_*(\lambda, 0) = n/(1 + \mathcal{R}_0(\lambda, \mu_*(\lambda, 0); \mathbf{I}))$, we obtain a lower bound for the right hand side

$$\begin{aligned} & f((1 - \kappa/2)^{-1}\mu_*(\lambda, 0)) \\ &= (1 - \kappa/2)^{-1}\mu_*(\lambda, 0) - \frac{n}{1 + \text{Tr}(\boldsymbol{\Sigma}(\lambda \mathbf{I} + (1 - \kappa/2)^{-1}\mu_*(\lambda, 0)\boldsymbol{\Sigma})^{-1})} \\ &= \frac{n(1 - \kappa/2)^{-1}}{1 + \text{Tr}(\boldsymbol{\Sigma}(\lambda \mathbf{I} + \mu_*(\lambda, 0)\boldsymbol{\Sigma})^{-1})} - \frac{n(1 - \kappa/2)^{-1}}{(1 - \kappa/2)^{-1} + \text{Tr}(\boldsymbol{\Sigma}((1 - \kappa/2)\lambda \mathbf{I} + \mu_*(\lambda, 0)\boldsymbol{\Sigma})^{-1})} \\ &\stackrel{(i)}{\geq} \frac{n(1 - \kappa/2)^{-1}}{1 + \text{Tr}(\boldsymbol{\Sigma}(\lambda \mathbf{I} + \mu_*(\lambda, 0)\boldsymbol{\Sigma})^{-1})} - \frac{n(1 - \kappa/2)^{-1}}{(1 - \kappa/2)^{-1} + \text{Tr}(\boldsymbol{\Sigma}(\lambda \mathbf{I} + \mu_*(\lambda, 0)\boldsymbol{\Sigma})^{-1})} \\ &= \frac{n((1 - \kappa/2)^{-1} - 1)}{(1 + \text{Tr}(\boldsymbol{\Sigma}(\lambda \mathbf{I} + \mu_*(\lambda, 0)\boldsymbol{\Sigma})^{-1})) \cdot (1 + (1 - \kappa/2)\text{Tr}(\boldsymbol{\Sigma}(\lambda \mathbf{I} + \mu_*(\lambda, 0)\boldsymbol{\Sigma})^{-1}))} \\ &\stackrel{(ii)}{\geq} \frac{n((1 - \kappa/2)^{-1} - 1)}{(1 + \text{Tr}(\boldsymbol{\Sigma}(\lambda \mathbf{I} + \mu_*(\lambda, 0)\boldsymbol{\Sigma})^{-1}))^2} \\ &= \frac{((1 - \kappa/2)^{-1} - 1)\mu_*(\lambda, 0)^2}{n}, \end{aligned}$$

where in (i) and (ii) we use two times the trivial bound $1 - \kappa/2 \leq 1$. By Assumption (21),

$$\frac{\mu_*(\lambda, 0)}{n} = \frac{\lambda}{n\lambda_*} \geq \kappa,$$

we know

$$f((1 - \kappa/2)^{-1}\mu_*(\lambda, 0)) \geq n \cdot \kappa^2 ((1 - \kappa/2)^{-1} - 1) \geq n \cdot \kappa^3/2,$$

and thus a sufficient condition for $\mu_*(\lambda, \mu) \leq (1 - \kappa/2)^{-1}\mu_*(\lambda, 0)$ is $\mu \leq n\kappa^3/2$.

C Proofs for Theorem 5

C.1 Proof of Lemma 7.1

We define $\bar{\mu}_i := \inf \left\{ \mu \mid \Sigma^{\frac{1}{2}} (\lambda \mathbf{I} + \mu \Sigma + \mathbf{X}_i^{\top} \mathbf{X}_i)^{-1} \Sigma^{\frac{1}{2}} \succ 0 \right\}$. Note that, by construction $\bar{\mu}_{i+1} \leq \bar{\mu}_i$. Let $\varphi \in \mathbb{R}^d$ be the leading normalized eigenvector of Σ . If $\mu \leq -(\lambda + \|\mathbf{X}_i \varphi\|^2) / \|\Sigma\|$, it follows that

$$\varphi^{\top} \left(\lambda \mathbf{I} + \mu \Sigma + \mathbf{X}_i^{\top} \mathbf{X}_i \right) \varphi = \lambda + \mu \|\Sigma\| + \|\mathbf{X}_i \varphi\|^2 \leq 0,$$

which implies $\bar{\mu}_i \geq -(\lambda + \|\mathbf{X}_i \varphi\|^2) / \|\Sigma\| > -\infty$. The update rule is equivalent to solving the equation

$$\mu_{i+1} + \frac{1}{1 + \text{Tr} \left(\Sigma (\lambda \mathbf{I} + \mu_{i+1} \Sigma + \mathbf{X}_i^{\top} \mathbf{X}_i)^{-1} \right)} = \mu_i, \quad \mu_{i+1} \in (\bar{\mu}_i, \infty).$$

For all $t \in (\bar{\mu}_i, \infty)$, let

$$f(t) = t + \frac{1}{1 + \text{Tr} \left(\Sigma (\lambda \mathbf{I} + t \Sigma + \mathbf{X}_i^{\top} \mathbf{X}_i)^{-1} \right)}.$$

In this given domain, $\Sigma^{\frac{1}{2}} (\lambda \mathbf{I} + t \Sigma + \mathbf{X}_i^{\top} \mathbf{X}_i)^{-1} \Sigma^{\frac{1}{2}} \succ 0$ and thus $\text{Tr} \left(\Sigma (\lambda \mathbf{I} + t \Sigma + \mathbf{X}_i^{\top} \mathbf{X}_i)^{-1} \right)$ is decreasing in t (this can be seen by computing its derivative with respect to t), which further implies $f(t)$ is strictly increasing in t . Since

$$\lim_{t \downarrow \bar{\mu}_i} \text{Tr} \left(\Sigma (\lambda \mathbf{I} + t \Sigma + \mathbf{X}_i^{\top} \mathbf{X}_i)^{-1} \right) = \infty,$$

and we have

$$\lim_{t \downarrow \bar{\mu}_i} f(t) = \bar{\mu}_i < \mu_i, \quad f(\mu_i) > \mu_i.$$

(The first inequality follows since $\mu_i \in (\bar{\mu}_{i-1}, \infty)$ and $\bar{\mu}_i \leq \bar{\mu}_{i-1}$.) Thus, there must be a unique $\mu_{i+1} \in (\bar{\mu}_i, \mu_i)$ that solves $f(\mu_{i+1}) = \mu_i$, proving the lemma.

C.2 Proof of Lemma 7.2

Without loss of generality, we can always assume $d \geq n$ or simply $d = \infty$ by embedding \mathbb{R}^d into the Hilbert space ℓ_2 since we always have

$$\sum_{l=k}^d \sigma_l \leq d_{\Sigma} \sigma_k,$$

when $\sigma_k = 0$. We write the spectral decomposition of Σ as

$$\Sigma = \sum_{i=1}^d \sigma_i \mathbf{v}_i \mathbf{v}_i^{\top},$$

with $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_n \geq \dots$ where $\{\mathbf{v}_i\}$ form an orthogonal basis of eigenvectors. For any $k \leq n$, define the projection operators

$$\mathbf{P}_k := \sum_{i=1}^k \mathbf{v}_i \mathbf{v}_i^\top, \quad \mathbf{P}_k^\perp := \mathbf{I} - \mathbf{P}_k = \sum_{i=k+1}^d \mathbf{v}_i \mathbf{v}_i^\top,$$

and we write

$$\mathbf{X} = \mathbf{X} \mathbf{P}_k + \mathbf{X} \mathbf{P}_k^\perp := \mathbf{U}_k + \mathbf{W}_k.$$

Part I: Decomposing into the top and lower eigenspaces. By writing $\mathbf{X} = \mathbf{U}_k + \mathbf{W}_k$, we can have the following inequality:

Lemma C.1. *For any $1 \leq k \leq n-1$,*

$$\lambda \mathbf{I} + \mathbf{X}^\top \mathbf{X} \succeq \frac{\lambda}{2} \mathbf{I} + \left(1 + \frac{2 \|\mathbf{W}_k^\top \mathbf{W}_k\|}{\lambda}\right)^{-1} \mathbf{U}_k^\top \mathbf{U}_k.$$

Proof. Note that

$$\begin{aligned} \lambda \mathbf{I} + \mathbf{X}^\top \mathbf{X} &= \lambda \mathbf{I} + \mathbf{U}_k^\top \mathbf{U}_k + \mathbf{U}_k^\top \mathbf{W}_k + \mathbf{W}_k^\top \mathbf{U}_k + \mathbf{W}_k^\top \mathbf{W}_k \\ &\succeq \frac{\lambda}{2} \mathbf{I} + \mathbf{U}_k^\top \mathbf{U}_k + \mathbf{U}_k^\top \mathbf{W}_k + \mathbf{W}_k^\top \mathbf{U}_k + \left(1 + \frac{\lambda}{2 \|\mathbf{W}_k^\top \mathbf{W}_k\|}\right) \mathbf{W}_k^\top \mathbf{W}_k \\ &= \frac{\lambda}{2} \mathbf{I} + \left(1 + \frac{2 \|\mathbf{W}_k^\top \mathbf{W}_k\|}{\lambda}\right)^{-1} \mathbf{U}_k^\top \mathbf{U}_k + \mathbf{C}_k, \end{aligned}$$

where

$$\mathbf{C}_k = \left(1 + \frac{\lambda}{2 \|\mathbf{W}_k^\top \mathbf{W}_k\|}\right)^{-1} \mathbf{U}_k^\top \mathbf{U}_k + \mathbf{U}_k^\top \mathbf{W}_k + \mathbf{W}_k^\top \mathbf{U}_k + \left(1 + \frac{\lambda}{2 \|\mathbf{W}_k^\top \mathbf{W}_k\|}\right) \mathbf{W}_k^\top \mathbf{W}_k = \mathbf{D}_k^\top \mathbf{D}_k \succeq 0,$$

with

$$\mathbf{D}_k = \left(1 + \frac{\lambda}{2 \|\mathbf{W}_k^\top \mathbf{W}_k\|}\right)^{-\frac{1}{2}} \mathbf{U}_k + \left(1 + \frac{\lambda}{2 \|\mathbf{W}_k^\top \mathbf{W}_k\|}\right)^{\frac{1}{2}} \mathbf{W}_k.$$

□

To apply the above lemma, we need to further provide an upper bound on $\|\mathbf{W}_k^\top \mathbf{W}_k\|$, which we summarize as the following result.

Lemma C.2. *Let Assumption 1 holds, we have for any $1 \leq k \leq n-1$ with probability $1 - \mathcal{O}(n^{-D})$ that,*

$$\|\mathbf{W}_k \mathbf{W}_k^\top\| = \mathcal{O}_{\mathbf{C}_x, D}(\mathbf{d}_\Sigma \sigma_k \cdot \log n \log(\mathbf{d}_\Sigma n)).$$

Proof. Let $\zeta_i = \mathbf{P}_k^\perp \mathbf{x}_i \mathbf{x}_i^\top \mathbf{P}_k^\perp \in \mathbb{R}^{d \times d}$, we can write

$$\mathbf{S}_k := \mathbf{W}_k^\top \mathbf{W}_k = \sum_{i=1}^n \mathbf{P}_k^\perp \mathbf{x}_i \mathbf{x}_i^\top \mathbf{P}_k^\perp = \sum_{i=1}^n \zeta_i.$$

Since $\|\zeta_i\| = \|\mathbf{P}_k^\perp \mathbf{x}_i\|^2$, we can apply Hanson-Wright inequality (cf. Lemma 2.1) and conclude that

$$\begin{aligned} \mathbb{P}\left(\left|\|\zeta_i\| - \text{Tr}\left(\mathbf{P}_k^\perp \boldsymbol{\Sigma}\right)\right| \geq t\right) &= \mathbb{P}\left(\left|\mathbf{x}_i^\top \mathbf{P}_k^\perp \mathbf{x}_i - \text{Tr}\left(\mathbf{P}_k^\perp \boldsymbol{\Sigma}\right)\right| \geq t\right) \\ &\leq 2 \exp\left\{-\Omega\left(\min\left\{\frac{t^2}{C_{\mathbf{x}}^4 \left\|\boldsymbol{\Sigma}^{\frac{1}{2}} \mathbf{P}_k^\perp \boldsymbol{\Sigma}^{\frac{1}{2}}\right\|_F^2}, \frac{t}{C_{\mathbf{x}}^2 \left\|\boldsymbol{\Sigma}^{\frac{1}{2}} \mathbf{P}_k^\perp \boldsymbol{\Sigma}^{\frac{1}{2}}\right\|}\right\}\right)\right\}. \end{aligned}$$

For $t = \Theta_{C_{\mathbf{x}}, D}(\|\boldsymbol{\Sigma}^{\frac{1}{2}} \mathbf{P}_k^\perp \boldsymbol{\Sigma}^{\frac{1}{2}}\|_F \log n)$ we have with probability $1 - \mathcal{O}(n^{-D})$ that for all $i = 1, 2, \dots, n$

$$\|\zeta_i\| \leq \text{Tr}\left(\mathbf{P}_k^\perp \boldsymbol{\Sigma}\right) + \Theta_{C_{\mathbf{x}}, D}\left(\left\|\boldsymbol{\Sigma}^{\frac{1}{2}} \mathbf{P}_k^\perp \boldsymbol{\Sigma}^{\frac{1}{2}}\right\|_F \log n\right) = \mathcal{O}_{C_{\mathbf{x}}, D}\left(\text{Tr}\left(\mathbf{P}_k^\perp \boldsymbol{\Sigma}\right) \log n\right),$$

where the last inequality follows from $\|\boldsymbol{\Sigma}\| = 1$ and

$$\left\|\boldsymbol{\Sigma}^{\frac{1}{2}} \mathbf{P}_k^\perp \boldsymbol{\Sigma}^{\frac{1}{2}}\right\|_F = \sqrt{\text{Tr}\left(\mathbf{P}_k^\perp \boldsymbol{\Sigma} \mathbf{P}_k^\perp \boldsymbol{\Sigma}\right)} \leq \text{Tr}\left(\mathbf{P}_k^\perp \boldsymbol{\Sigma}\right).$$

In the next step, we will adopt a standard truncation argument and apply a matrix concentration inequality. By setting $L_k := \Theta_{C_{\mathbf{x}}, D}(\text{Tr}\left(\mathbf{P}_k^\perp \boldsymbol{\Sigma}\right) \log n)$, $\tilde{\zeta}_i := \zeta_i \mathbb{I}\{\|\zeta_i\| \leq L_k\}$ and considering

$$\tilde{\mathbf{S}}_k := \sum_{i=1}^k \tilde{\zeta}_i = \sum_{i=1}^k \zeta_i \mathbb{I}\{\|\zeta_i\| \leq L_k\},$$

we have $\tilde{\mathbf{S}}_k = \mathbf{S}_k$ with probability $1 - \mathcal{O}(n^{-D})$. In order to bound $\|\tilde{\mathbf{S}}_k\|$, we will use matrix Bernstein inequality. Since we know $\|\tilde{\zeta}_i\| \leq L_k$ by construction, we only need to upper bound the matrix variance. The $\tilde{\zeta}_i$'s are independent symmetric random matrices and therefore we have

$$\text{Var}(\tilde{\mathbf{S}}_k) \preceq \sum_{i=1}^n \mathbb{E}[\tilde{\zeta}_i^2] \stackrel{(i)}{\preceq} \sum_{i=1}^n L_k \mathbb{E}[\tilde{\zeta}_i] \stackrel{(ii)}{\preceq} \sum_{i=1}^n L_k \mathbb{E}[\zeta_i] \stackrel{(iii)}{\preceq} n L_k \cdot \mathbf{P}_k^\perp \boldsymbol{\Sigma} \mathbf{P}_k^\perp =: \mathbf{V}_k,$$

where in (i) we use $\|\tilde{\zeta}_i\| \leq L_k$, in (ii) we apply $\tilde{\zeta}_i \preceq \zeta_i$ and lastly in (iii) we use $\mathbb{E}[\zeta_i] = \mathbb{E}[\mathbf{P}_k^\perp \mathbf{x}_i \mathbf{x}_i^\top \mathbf{P}_k^\perp] = \mathbf{P}_k^\perp \boldsymbol{\Sigma} \mathbf{P}_k^\perp$. It then follows that $\|\mathbf{V}_k\| \leq n L_k \|\mathbf{P}_k^\perp \boldsymbol{\Sigma} \mathbf{P}_k^\perp\| = n \sigma_{k+1} L_k \leq n \sigma_k L_k =: v_k$. Combine with the bound on the intrinsic dimension under Assumption 1,

$$\text{intdim}(\mathbf{V}_k) = \frac{\text{Tr}(\mathbf{V}_k)}{\|\mathbf{V}_k\|} = \frac{\sum_{l=k+1}^{\infty} \sigma_l}{\sigma_{k+1}} \leq \mathbf{d}_{\boldsymbol{\Sigma}},$$

we can thus deduce from the Bernstein inequality with intrinsic dimension [T⁺15, Theorem 7.3.1] that for $t \geq \sqrt{v_k} + L_k/3$

$$\mathbb{P}(\|\tilde{\mathbf{S}}_k - \mathbb{E}[\tilde{\mathbf{S}}_k]\| \geq t) \leq 4 \mathbf{d}_{\boldsymbol{\Sigma}} \cdot \exp\left(\frac{-t^2/2}{v_k + L_k t/3}\right).$$

Finally, by further bounding the mean

$$\|\mathbb{E}[\tilde{\mathbf{S}}_k]\| \leq \|\mathbb{E}[\tilde{\zeta}_i]\| \leq n \cdot \|\mathbb{E}[\zeta_i]\| = n \cdot \left\|\mathbf{P}_k^\perp \boldsymbol{\Sigma} \mathbf{P}_k^\perp\right\| \leq n \sigma_k,$$

we can obtain with probability $1 - \mathcal{O}(n^{-D})$,

$$\begin{aligned}\|\tilde{\mathbf{S}}_k\| &= \mathcal{O}_D((\sqrt{v_k} + L_k) \log(\mathbf{d}_\Sigma n)) + \|\mathbb{E}[\tilde{\mathbf{S}}_k]\| \\ &\stackrel{(i)}{=} \mathcal{O}_{\mathbf{C}_x, D} \left(\left\{ \sqrt{n\sigma_k \cdot \text{Tr}(\mathbf{P}_k^\perp \Sigma)} \log n + \text{Tr}(\mathbf{P}_k^\perp \Sigma) \log n \right\} \cdot \log(\mathbf{d}_\Sigma n) + n\sigma_k \right) \\ &\stackrel{(ii)}{=} \mathcal{O}_{\mathbf{C}_x, D} \left(\left\{ \sqrt{n\sigma_k \cdot \mathbf{d}_\Sigma \sigma_k \log n} + \mathbf{d}_\Sigma \sigma_k \log n \right\} \cdot \log(\mathbf{d}_\Sigma n) + n\sigma_k \right)\end{aligned}$$

where in (i) we make use of $v_k = n\sigma_k L_k$, and apply Assumption 1 for the spectrum in (ii). Next by the fact that $\mathbf{d}_\Sigma \geq n$, we can further write

$$\|\tilde{\mathbf{S}}_k\| = \mathcal{O}_{\mathbf{C}_x, D}(\mathbf{d}_\Sigma \sigma_k \cdot \log n \log(\mathbf{d}_\Sigma n)).$$

The proof is complete as $\mathbf{W}_k \mathbf{W}_k^\top = \mathbf{S}_k = \tilde{\mathbf{S}}_k$ holds with probability $1 - \mathcal{O}(n^{-D})$. \square

To bound the norm of $\Sigma^{\frac{1}{2}} (\lambda \mathbf{I} + \mathbf{X}^\top \mathbf{X})^{-1} \Sigma^{\frac{1}{2}}$, we apply Lemmas C.1 and C.2 and obtain

$$\begin{aligned}\lambda \mathbf{I} + \mathbf{X}^\top \mathbf{X} &\succeq \frac{\lambda}{2} \mathbf{I} + \left(1 + \frac{2 \|\mathbf{W}_k^\top \mathbf{W}_k\|}{\lambda} \right)^{-1} \mathbf{U}_k^\top \mathbf{U}_k \\ &\succeq \frac{\lambda}{2} \mathbf{I} + \left(1 + \frac{\mathcal{O}_{\mathbf{C}_x, D}(\mathbf{d}_\Sigma \sigma_k \cdot \log n \log(\mathbf{d}_\Sigma n))}{\lambda} \right)^{-1} \mathbf{U}_k^\top \mathbf{U}_k.\end{aligned}$$

Therefore by block matrix inverse, we can further get

$$\begin{aligned}&\Sigma^{\frac{1}{2}} (\lambda \mathbf{I} + \mathbf{X}^\top \mathbf{X})^{-1} \Sigma^{\frac{1}{2}} \\ &\preceq \Sigma^{\frac{1}{2}} \left(\frac{\lambda}{2} \mathbf{I} + \left(1 + \frac{\mathcal{O}_{\mathbf{C}_x, D}(\mathbf{d}_\Sigma \sigma_k \cdot \log n \log(\mathbf{d}_\Sigma n))}{\lambda} \right)^{-1} \mathbf{U}_k^\top \mathbf{U}_k \right)^{-1} \Sigma^{\frac{1}{2}} \\ &= \Sigma^{\frac{1}{2}} \left(\frac{\lambda}{2} \mathbf{I} + \left(1 + \frac{\mathcal{O}_{\mathbf{C}_x, D}(\mathbf{d}_\Sigma \sigma_k \cdot \log n \log(\mathbf{d}_\Sigma n))}{\lambda} \right)^{-1} \mathbf{P}_k \Sigma^{\frac{1}{2}} \mathbf{Z}^\top \mathbf{Z} \Sigma^{\frac{1}{2}} \mathbf{P}_k \right)^{-1} \Sigma^{\frac{1}{2}} \\ &\preceq \left(1 + \frac{\mathcal{O}_{\mathbf{C}_x, D}(\mathbf{d}_\Sigma \sigma_k \cdot \log n \log(\mathbf{d}_\Sigma n))}{\lambda} \right) \left(\mathbf{P}_k \mathbf{Z}^\top \mathbf{Z} \mathbf{P}_k \right)^\dagger + \frac{2 \mathbf{P}_k^\perp \Sigma \mathbf{P}_k^\perp}{\lambda},\end{aligned}\tag{58}$$

where $\mathbf{X} = \mathbf{Z} \Sigma^{\frac{1}{2}}$. Define the matrix $\mathbf{V}_k = [\mathbf{v}_1 \ \cdots \ \mathbf{v}_k] \in \mathbb{R}^{d \times k}$, we can then write $\mathbf{P}_k = \mathbf{V}_k \mathbf{V}_k^\top$. Thus by exploiting the block matrix structure, it follows that

$$\left\| \Sigma^{\frac{1}{2}} (\lambda \mathbf{I} + \mathbf{X}^\top \mathbf{X})^{-1} \Sigma^{\frac{1}{2}} \right\| \leq \left(1 + \frac{\mathcal{O}_{\mathbf{C}_x, D}(\mathbf{d}_\Sigma \sigma_k \cdot \log n \log(\mathbf{d}_\Sigma n))}{\lambda} \right) \lambda_{\min} \left(\mathbf{V}_k^\top \mathbf{Z}^\top \mathbf{Z} \mathbf{V}_k \right)^{-1} + \frac{2\sigma_k}{\lambda}.\tag{59}$$

Substituting $\boldsymbol{\theta} = \Sigma^{-1/2} \boldsymbol{\beta}$ into Eq. (58), we also obtain

$$\begin{aligned}&\boldsymbol{\theta}^\top \Sigma^{\frac{1}{2}} (\lambda \mathbf{I} + \mathbf{X}^\top \mathbf{X})^{-1} \Sigma^{\frac{1}{2}} \boldsymbol{\theta} \\ &\leq \left(1 + \frac{\mathcal{O}_{\mathbf{C}_x, D}(\mathbf{d}_\Sigma \sigma_k \cdot \log n \log(\mathbf{d}_\Sigma n))}{\lambda} \right) \lambda_{\min} \left(\mathbf{V}_k^\top \mathbf{Z}^\top \mathbf{Z} \mathbf{V}_k \right)^{-1} \|\boldsymbol{\theta}_{\leq k}\|^2 + \frac{2 \|\boldsymbol{\beta}_{> k}\|^2}{\lambda}.\end{aligned}\tag{60}$$

Part II: Lower bounding the smallest eigenvalue $\lambda_{\min}(\mathbf{V}_k^\top \mathbf{Z}^\top \mathbf{Z} \mathbf{V}_k)$. The last step is then to provide a lower bound for the smallest eigenvalue of $\mathbf{V}_k^\top \mathbf{Z}^\top \mathbf{Z} \mathbf{V}_k$. Consider $\tilde{\mathbf{Z}} = [\tilde{\mathbf{z}}_1 \ \cdots \ \tilde{\mathbf{z}}_n]^\top \in \mathbb{R}^{n \times k}$ with

$$\tilde{\mathbf{z}}_i = \mathbf{V}_k^\top \mathbf{z}_i = \begin{bmatrix} \langle \mathbf{z}, \mathbf{v}_1 \rangle \\ \vdots \\ \langle \mathbf{z}, \mathbf{v}_{k-1} \rangle \end{bmatrix}.$$

We therefore need to lower bound $\lambda_{\min}(\tilde{\mathbf{Z}}^\top \tilde{\mathbf{Z}})$ where $\tilde{\mathbf{Z}}$ has i.i.d. rows $\tilde{\mathbf{z}}_i$ in \mathbb{R}^k . An immediate consequence is that $\mathbb{E}[\tilde{\mathbf{z}}_i] = \mathbf{0}$ and $\text{Var}(\tilde{\mathbf{z}}_i) = \mathbf{I}_k$. Moreover, for any unit vector $\boldsymbol{\varphi} \in \mathbb{R}^k$, we can apply Hanson-Wright (cf. Lemma 2.1) and deduce that for any $t \geq 0$,

$$\begin{aligned} \mathbb{P}(|\langle \tilde{\mathbf{z}}_i, \boldsymbol{\varphi} \rangle^2 - 1| \geq t) &= \mathbb{P}\left(\left| \mathbf{z}_i^\top \mathbf{V}_k \boldsymbol{\varphi} \boldsymbol{\varphi}^\top \mathbf{V}_k^\top \mathbf{z}_i - \text{Tr}\left(\mathbf{V}_k \boldsymbol{\varphi} \boldsymbol{\varphi}^\top \mathbf{V}_k^\top\right) \right| \geq t\right) \\ &\leq 2 \exp\left\{-\Omega\left(\min\left\{\frac{t^2}{\mathbf{C}_x^4 \|\mathbf{V}_k \boldsymbol{\varphi} \boldsymbol{\varphi}^\top \mathbf{V}_k^\top\|_F^2}, \frac{t}{\mathbf{C}_x^2 \|\mathbf{V}_k \boldsymbol{\varphi} \boldsymbol{\varphi}^\top \mathbf{V}_k^\top\|}\right\}\right)\right\} \\ &= 2 \exp(-\Omega_{\mathbf{C}_x}(\min\{t^2, t\})), \end{aligned} \quad (61)$$

where we use the fact that $\|\mathbf{V}_k \boldsymbol{\varphi} \boldsymbol{\varphi}^\top \mathbf{V}_k^\top\|_F = \|\mathbf{V}_k \boldsymbol{\varphi} \boldsymbol{\varphi}^\top \mathbf{V}_k^\top\| \leq 1$. Thus we can bound the fourth moment of $\langle \tilde{\mathbf{z}}_i, \boldsymbol{\varphi} \rangle$ by

$$\begin{aligned} \mathbb{E}[\langle \tilde{\mathbf{z}}_i, \boldsymbol{\varphi} \rangle^4] &= \int_0^\infty 2t \mathbb{P}(\langle \tilde{\mathbf{z}}_i, \boldsymbol{\varphi} \rangle^2 \geq t) dt \\ &\leq 1 + \int_0^\infty 2(t+1) \mathbb{P}(\langle \tilde{\mathbf{z}}_i, \boldsymbol{\varphi} \rangle^2 \geq t+1) dt \\ &\leq 1 + 4 \int_0^\infty (t+1) \exp(-\Omega_{\mathbf{C}_x}(\min\{t^2, t\})) dt = \mathcal{O}_{\mathbf{C}_x}(1). \end{aligned}$$

Clearly the above bound holds uniformly for all $\boldsymbol{\varphi} \in \mathbb{S}^{k-1}$ from the unit sphere in \mathbb{R}^k . Since the upper bound $\mathcal{O}_{\mathbf{C}_x}(1)$ does not depend on k , we can appeal to [Yas14, Theorem 2.2] and obtain that with probability $1 - \mathcal{O}(n^{-D})$

$$\lambda_{\min}(n^{-1} \tilde{\mathbf{Z}}^\top \tilde{\mathbf{Z}}) \geq 1 - \mathcal{O}_{\mathbf{C}_x}\left(\sqrt{\frac{k}{n}}\right) - \mathcal{O}_D\left(\sqrt{\frac{\log n}{n}}\right).$$

Therefore, if we choose $k = \lfloor \eta n \rfloor$ for some fixed η such that $\mathcal{O}_{\mathbf{C}_x}(\sqrt{\eta}) \leq 1/4$, it holds for $n = \Omega_D(1)$ that

$$\lambda_{\min}(n^{-1} \tilde{\mathbf{Z}}^\top \tilde{\mathbf{Z}}) \geq 1 - \frac{1}{4} - \frac{1}{4} = \frac{1}{2},$$

and we therefore conclude the proof by taking $k = \lfloor \eta n \rfloor$ as above and substituting into Eq. (59)

$$\begin{aligned} \left\| \boldsymbol{\Sigma}^{\frac{1}{2}} \left(\lambda \mathbf{I} + \mathbf{X}^\top \mathbf{X} \right)^{-1} \boldsymbol{\Sigma}^{\frac{1}{2}} \right\| &\leq \frac{2}{n} \left(1 + \frac{\mathcal{O}_{\mathbf{C}_x, D}(\mathbf{d}_{\boldsymbol{\Sigma}} \sigma_k \cdot \log n \log(\mathbf{d}_{\boldsymbol{\Sigma}} n))}{\lambda} \right) + \frac{2\sigma_k}{\lambda} \\ &= \frac{2}{n} \left(1 + \frac{\mathcal{O}_{\mathbf{C}_x, D}(\mathbf{d}_{\boldsymbol{\Sigma}} \sigma_k \cdot \log n \log(\mathbf{d}_{\boldsymbol{\Sigma}} n))}{\lambda} \right), \end{aligned}$$

where in the last line we use the fact that $\mathbf{d}_{\boldsymbol{\Sigma}} \geq n$ and therefore $\sigma_k = \mathcal{O}(\mathbf{d}_{\boldsymbol{\Sigma}} \sigma_k \cdot \log n \log(\mathbf{d}_{\boldsymbol{\Sigma}} n)/n)$.

Similarly for Eq. (60), we have

$$\begin{aligned} \boldsymbol{\theta}^\top \boldsymbol{\Sigma}^{\frac{1}{2}} \left(\lambda \mathbf{I} + \mathbf{X}^\top \mathbf{X} \right)^{-1} \boldsymbol{\Sigma}^{\frac{1}{2}} \boldsymbol{\theta} &\leq \frac{2}{n} \left(1 + \frac{\mathcal{O}_{\mathbf{C}_x, D}(\mathbf{d}_\Sigma \sigma_k \cdot \log n \log(\mathbf{d}_\Sigma n))}{\lambda} \right) \|\boldsymbol{\theta}_{\leq k}\|^2 + \frac{2 \|\boldsymbol{\beta}_{>k}\|^2}{\lambda} \\ &\leq \frac{2}{n} \left(1 + \frac{\mathcal{O}_{\mathbf{C}_x, D}(\mathbf{d}_\Sigma \sigma_k \cdot \log n \log(\mathbf{d}_\Sigma n))}{\lambda} \right) \|\boldsymbol{\theta}_{\leq n}\|^2 + \frac{2 \|\boldsymbol{\beta}_{>n}\|^2}{\lambda}, \end{aligned}$$

where in the last line we use the fact that for all $k+1 \leq i \leq n$,

$$\langle \boldsymbol{\beta}, \mathbf{v}_i \rangle^2 = \sigma_i \langle \boldsymbol{\theta}, \mathbf{v}_i \rangle^2 \leq \sigma_k \langle \boldsymbol{\theta}, \mathbf{v}_i \rangle^2 = \mathcal{O}(\mathbf{d}_\Sigma \sigma_k \cdot \log n \log(\mathbf{d}_\Sigma n) / n) \langle \boldsymbol{\theta}, \mathbf{v}_i \rangle^2.$$

The proof is complete.

C.3 Proof of Lemma 7.3

We apply Hanson-Wright inequality in Lemma 2.1 and get

$$\mathbb{P} \left(\left| \mathbf{z}_k^\top \mathbf{B}_{k-1} \mathbf{z}_k - S_{k-1}(\mathbf{I}) \right| \geq t \mid \mathbf{B}_{k-1} \right) \leq 2 \exp \left\{ -\Omega \left(\min \left\{ \frac{t^2}{\mathbf{C}_x^4 \|\mathbf{B}_{k-1}\|_F^2}, \frac{t}{\mathbf{C}_x^2 \|\mathbf{B}_{k-1}\|} \right\} \right) \right\}$$

and

$$\begin{aligned} &\mathbb{P} \left(\left| \mathbf{z}_k^\top \mathbf{B}_{k-1} \mathbf{Q} \mathbf{B}_{k-1} \mathbf{z}_k - \text{Tr}(\mathbf{Q} \mathbf{B}_{k-1}^2) \right| \geq t \mid \mathbf{B}_{k-1} \right) \\ &\leq 2 \exp \left\{ -\Omega \left(\min \left\{ \frac{t^2}{\mathbf{C}_x^4 \|\mathbf{B}_{k-1} \mathbf{Q} \mathbf{B}_{k-1}\|_F^2}, \frac{t}{\mathbf{C}_x^2 \|\mathbf{B}_{k-1} \mathbf{Q} \mathbf{B}_{k-1}\|} \right\} \right) \right\}. \end{aligned}$$

In particular, on the event $\{T_F(\mathbf{Q}) \geq k, T_F(\mathbf{I}) \geq k\}$, we have

$$|S_{k-1}(\mathbf{I}) - R_0(\mathbf{I})| \leq \beta_1 \leq \frac{1}{4} R_0(\mathbf{I}), \quad |S_{k-1}(\mathbf{Q}) - R_0(\mathbf{Q})| \leq \beta_1 \leq \frac{1}{4} R_0(\mathbf{Q}), \quad \|\mathbf{B}_{k-1}\| \leq \gamma,$$

which further implies that

$$\|\mathbf{B}_{k-1}\| \leq \|\mathbf{B}_{k-1}\|_F = \sqrt{\text{Tr}(\mathbf{B}_{k-1}^2)} \leq \sqrt{\|\mathbf{B}_{k-1}\| \cdot S_{k-1}(\mathbf{I})} = \mathcal{O}(\sqrt{\gamma R_0(\mathbf{I})}),$$

and

$$\begin{aligned} \|\mathbf{B}_{k-1} \mathbf{Q} \mathbf{B}_{k-1}\| &\leq \|\mathbf{B}_{k-1} \mathbf{Q} \mathbf{B}_{k-1}\|_F = \sqrt{\text{Tr}(\mathbf{B}_{k-1} \mathbf{Q} \mathbf{B}_{k-1}^2 \mathbf{Q} \mathbf{B}_{k-1})} \\ &\leq \sqrt{\left\| \mathbf{Q}^{\frac{1}{2}} \mathbf{B}_{k-1}^2 \mathbf{Q}^{\frac{1}{2}} \right\| \cdot \text{Tr}(\mathbf{B}_{k-1} \mathbf{Q} \mathbf{B}_{k-1})} \\ &\leq \sqrt{\|\mathbf{B}_{k-1}\|^2 \cdot \text{Tr}(\mathbf{Q}^{\frac{1}{2}} \mathbf{B}_{k-1}^2 \mathbf{Q}^{\frac{1}{2}})} \\ &\stackrel{(i)}{\leq} \sqrt{\|\mathbf{B}_{k-1}\|^3 \cdot S_{k-1}(\mathbf{Q})} = \mathcal{O}(\sqrt{\gamma^3 R_0(\mathbf{Q})}), \end{aligned}$$

Substituting the above bounds into the Hanson-Wright inequalities, we have conditioning on $H_k := \{T_F(\mathbf{Q}) \geq k, T_F(\mathbf{I}) \geq k\}$ for some constant $\mathbf{C} = \mathbf{C}(\mathbf{C}_x, D)$ that

$$\exp \left\{ -\Omega \left(\frac{\mathbf{C} \log n \cdot \sqrt{\gamma R_0(\mathbf{I})}}{\mathbf{C}_x^2 \|\mathbf{B}_{k-1}\|} \right) \right\} = \mathcal{O}(n^{-D}),$$

$$\exp \left\{ -\Omega \left(\frac{\mathbb{C} \log n \cdot \sqrt{\gamma^3 \mathbf{R}_0(\mathbf{Q})}}{\mathbb{C}_{\mathbf{x}}^2 \|\mathbf{B}_{k-1} \mathbf{Q} \mathbf{B}_{k-1}\|} \right) \right\} = \mathcal{O}(n^{-D}),$$

and therefore it holds with probability $1 - \mathcal{O}(n^{-D})$ that

$$\begin{aligned} \left| \mathbf{z}_k^\top \mathbf{B}_{k-1} \mathbf{z}_k - \mathbf{S}_{k-1}(\mathbf{I}) \right| &\leq \mathbb{C} \log n \cdot \sqrt{\gamma \mathbf{R}_0(\mathbf{I})} =: \alpha_1, \\ \left| \mathbf{z}_k^\top \mathbf{B}_{k-1} \mathbf{Q} \mathbf{B}_{k-1} \mathbf{z}_k - \text{Tr}(\mathbf{Q} \mathbf{B}_{k-1}^2) \right| &\leq \mathbb{C} \log n \cdot \sqrt{\gamma^3 \mathbf{R}_0(\mathbf{Q})} =: \alpha_2. \end{aligned}$$

The Hanson-Wright inequalities also give the following upper bounds on the expectations conditioning on the tail event when $\|\mathbf{B}_{k-1}\| \leq \gamma$. In particular, we would have

$$\begin{aligned} &\mathbb{E}_{k-1} \left[\left| \mathbf{z}_k^\top \mathbf{B}_{k-1} \mathbf{z}_k - \mathbf{S}_{k-1}(\mathbf{I}) \right| \mathbb{I} \left\{ \left| \mathbf{z}_k^\top \mathbf{B}_{k-1} \mathbf{z}_k - \mathbf{S}_{k-1}(\mathbf{I}) \right| \geq \alpha_1 \right\} \right] \mathbb{I}(H_k) \\ &= \int_{\mathbb{C} \log n \cdot \sqrt{\gamma \mathbf{R}_0(\mathbf{I})}}^{\infty} \mathbb{P} \left(\left| \mathbf{z}_k^\top \mathbf{B}_{k-1} \mathbf{z}_k - \mathbf{S}_{k-1}(\mathbf{I}) \right| \geq t \mid \mathbf{B}_{k-1} \right) \mathbb{I}(H_k) dt \\ &\leq \int_{\mathbb{C} \log n \cdot \sqrt{\gamma \mathbf{R}_0(\mathbf{I})}}^{\infty} 2 \exp \left\{ -\Omega \left(\min \left\{ \frac{t^2}{\mathbb{C}_{\mathbf{x}}^4 \|\mathbf{B}_{k-1}\|_F^2}, \frac{t}{\mathbb{C}_{\mathbf{x}}^2 \|\mathbf{B}_{k-1}\|} \right\} \right) \right\} \mathbb{I}(H_k) dt \\ &\leq \int_{\mathbb{C} \log n \cdot \sqrt{\gamma \mathbf{R}_0(\mathbf{I})}}^{\infty} 2 \exp \left\{ -\Omega \left(\frac{t}{\mathbb{C}_{\mathbf{x}}^2 \|\mathbf{B}_{k-1}\|_F} \right) \right\} \mathbb{I}(H_k) dt \\ &\leq \mathcal{O}(\mathbb{C}_{\mathbf{x}}^2 \|\mathbf{B}_{k-1}\|_F) \cdot \mathcal{O}(n^{-D}) = \mathcal{O}_{\mathbb{C}_{\mathbf{x}}} \left(n^{-D} \cdot \sqrt{\gamma \mathbf{R}_0(\mathbf{I})} \right). \end{aligned}$$

Similarly it also holds that

$$\begin{aligned} &\mathbb{E}_{k-1} \left[\left| \mathbf{z}_k^\top \mathbf{B}_{k-1} \mathbf{Q} \mathbf{B}_{k-1} \mathbf{z}_k - \text{Tr}(\mathbf{Q} \mathbf{B}_{k-1}^2) \right| \mathbb{I} \left\{ \left| \mathbf{z}_k^\top \mathbf{B}_{k-1} \mathbf{Q} \mathbf{B}_{k-1} \mathbf{z}_k - \text{Tr}(\mathbf{Q} \mathbf{B}_{k-1}^2) \right| \geq \alpha_2 \right\} \right] \mathbb{I}(H_k) \\ &= \mathcal{O}_{\mathbb{C}_{\mathbf{x}}} \left(n^{-D} \cdot \sqrt{\gamma^3 \mathbf{R}_0(\mathbf{Q})} \right). \end{aligned}$$

To finish the proof, we now only need to show $\|\mathbf{A}_k\| \leq \gamma$ holds with probability $1 - \mathcal{O}(n^{-D})$. We provide upper bounds for small $k \leq n/2$ and large $k > n/2$ separately. Under the assumption $\beta_2 \leq \mu/2$, we make use of the fact that $H_k \subset F_{k-1}(\mathbf{Q})$ which enables us to derive

$$\left| \mu_k - \mu_{\star}(\lambda, \mu) + \frac{k}{1 + \mathbf{R}_0(\mathbf{I})} \right| \leq \beta_2 \leq \frac{\mu}{2}, \quad (62)$$

and thus for all $k \leq n/2$,

$$\mu_k \geq \mu_{\star}(\lambda, \mu) - \frac{k}{1 + \mathbf{R}_0(\mathbf{I})} - \frac{\mu}{2} = \frac{\mu}{2} + \frac{n-k}{1 + \mathbf{R}_0(\mathbf{I})} \geq \frac{\mu_{\star}(\lambda, \mu)}{2},$$

which in particular implies for $k \leq n/2$ that

$$\|\mathbf{A}_k\| \leq \left\| \boldsymbol{\Sigma}^{\frac{1}{2}} (\lambda \mathbf{I} + \mu_k \boldsymbol{\Sigma})^{-1} \boldsymbol{\Sigma}^{\frac{1}{2}} \right\| \leq \frac{2}{\mu_{\star}(\lambda, \mu)}.$$

On the other hand, if $k > n/2$, we can still deduce from Eq. (62) that $\mu_k \geq \mu/2 > 0$ and thus

$$\mathbf{A}_k = \boldsymbol{\Sigma}^{\frac{1}{2}} \left(\lambda \mathbf{I} + \mu_k \boldsymbol{\Sigma} + \mathbf{X}_k^\top \mathbf{X}_k \right)^{-1} \boldsymbol{\Sigma}^{\frac{1}{2}} \leq \boldsymbol{\Sigma}^{\frac{1}{2}} \left(\lambda \mathbf{I} + \mathbf{X}_k^\top \mathbf{X}_k \right)^{-1} \boldsymbol{\Sigma}^{\frac{1}{2}}.$$

Applying Lemma 7.2, we obtain with probability $1 - \mathcal{O}(n^{-D})$ for all $k > n/2$,

$$\|\mathbf{A}_k\| \leq \left\| \boldsymbol{\Sigma}^{\frac{1}{2}} \left(\lambda \mathbf{I} + \mathbf{X}_{\lceil n/2 \rceil}^{\top} \mathbf{X}_{\lceil n/2 \rceil} \right)^{-1} \boldsymbol{\Sigma}^{\frac{1}{2}} \right\| = \frac{2}{n} \left(1 + \frac{\mathcal{O}_{\mathbf{C}_x, D}(\mathbf{d}_{\boldsymbol{\Sigma}} \sigma_{\lceil \eta n \rceil} \cdot \log n \log(\mathbf{d}_{\boldsymbol{\Sigma}} n))}{\lambda} \right),$$

for some $\eta = \eta(\mathbf{C}_x)$. Combine with the trivial bound $\|\mathbf{A}_k\| \leq 1/\lambda$, we conclude that $\|\mathbf{A}_k\| \leq \gamma$ with probability $1 - \mathcal{O}(n^{-D})$, provided we take

$$\gamma = \min \left\{ \frac{2}{n} \left(1 + \frac{\mathcal{O}_{\mathbf{C}_x, D}(\mathbf{d}_{\boldsymbol{\Sigma}} \sigma_{\lceil \eta n \rceil} \cdot \log n \log(\mathbf{d}_{\boldsymbol{\Sigma}} n))}{\lambda} \right) + \frac{2}{\mu_{\star}(\lambda, \mu)}, \frac{1}{\lambda} \right\}.$$

The proof is completed by noting that

$$p_{k,k}(T_F, T_F, \mathbf{Q}) - p_{k+1,k}(T_E, T_F, \mathbf{Q}) = \mathbb{P}(E_k(\mathbf{Q})^c; H_k) = \mathcal{O}(n^{-D}).$$

C.4 Proof of Lemma 7.4

We begin by noticing that

$$p_{k,k}(T_E, T_E, \mathbf{Q}) - p_{k,k}(T_F, T_E, \mathbf{Q}) = \mathbb{P}(T_E(\mathbf{Q}) \geq k, T_E(\mathbf{I}) \geq k; F_{k-1}^c(\mathbf{Q})). \quad (63)$$

We therefore need to control $|\mathbf{R}_{k-1}(\mathbf{Q}) - \mathbf{R}_0(\mathbf{Q})|$ and $|\mathbf{S}_{k-1}(\mathbf{Q}) - \mathbf{R}_0(\mathbf{Q})|$, as well as $|\mu_k - \bar{\mu}_k|$ and $\|\mathbf{B}_{k-1}\|$.

Part I: Decomposing into martingale part and bias part. Recall the calculations for Eq. (47), we have

$$\mathbf{R}_i(\mathbf{Q}) - \mathbf{R}_{i-1}(\mathbf{Q}) = -\frac{\text{Tr}(\mathbf{Q} \mathbf{B}_{i-1} \mathbf{z}_i \mathbf{z}_i^{\top} \mathbf{B}_{i-1})}{1 + \mathbf{z}_i^{\top} \mathbf{B}_{i-1} \mathbf{z}_i} + \frac{\text{Tr}(\mathbf{Q} \mathbf{B}_{i-1} \mathbf{A}_{i-1})}{1 + \mathbf{S}_{i-1}(\mathbf{I})}$$

Define the stopping time

$$\bar{T} = T_E(\mathbf{Q}) \wedge T_E(\mathbf{I}),$$

and on the event $\{T_E(\mathbf{Q}) \geq k, T_E(\mathbf{I}) \geq k\} = \{\bar{T} \geq k\}$, it holds

$$\mathbf{R}_{k-1}(\mathbf{Q}) - \mathbf{R}_0(\mathbf{Q}) = (\mathbf{R}_{k-1}(\mathbf{Q}) - \mathbf{R}_0(\mathbf{Q})) \mathbb{I}\{\bar{T} \geq k\} = \sum_{i=1}^{k-1} (\mathbf{R}_i(\mathbf{Q}) - \mathbf{R}_{i-1}(\mathbf{Q})) \mathbb{I}\{\bar{T} \geq i+1\}.$$

For each of the summand, we can decompose it into two parts—the martingale difference part $D_i(\mathbf{Q}, \bar{T})$ and a bias part $B_i(\mathbf{Q}, \bar{T})$ —to be specific, we can write

$$(\mathbf{R}_i(\mathbf{Q}) - \mathbf{R}_{i-1}(\mathbf{Q})) \mathbb{I}\{\bar{T} \geq i+1\} = D_i(\mathbf{Q}, \bar{T}) + B_i(\mathbf{Q}, \bar{T}),$$

where by setting $G_i := F_i(\mathbf{Q}) \cap F_i(\mathbf{I}) \in \mathcal{F}_i$ and $S_i := \{\bar{T} = i\} \cap G_{i-1} \in \mathcal{F}_i$, the explicit forms of D_i and B_i are (recall that $\mathbb{E}_i(\cdot) := \mathbb{E}(\cdot | \mathcal{F}_i)$):

$$\begin{aligned} D_i(\mathbf{Q}, \bar{T}) &:= -\frac{\text{Tr}(\mathbf{Q} \mathbf{B}_{i-1} \mathbf{z}_i \mathbf{z}_i^{\top} \mathbf{B}_{i-1})}{1 + \mathbf{z}_i^{\top} \mathbf{B}_{i-1} \mathbf{z}_i} \mathbb{I}\{\bar{T} \geq i+1\} - \frac{\text{Tr}(\mathbf{Q} \mathbf{B}_{i-1}^2)}{1 + \mathbf{S}_{i-1}(\mathbf{I})} \mathbb{I}(S_i) \\ &+ \mathbb{E}_{i-1} \left[\frac{\text{Tr}(\mathbf{Q} \mathbf{B}_{i-1} \mathbf{z}_i \mathbf{z}_i^{\top} \mathbf{B}_{i-1})}{1 + \mathbf{z}_i^{\top} \mathbf{B}_{i-1} \mathbf{z}_i} \mathbb{I}\{\bar{T} \geq i+1\} + \frac{\text{Tr}(\mathbf{Q} \mathbf{B}_{i-1}^2)}{1 + \mathbf{S}_{i-1}(\mathbf{I})} \mathbb{I}(S_i) \right], \end{aligned} \quad (64)$$

and

$$B_i(\mathbf{Q}, \bar{T}) := \frac{\text{Tr}(\mathbf{Q}\mathbf{B}_{i-1}\mathbf{A}_{i-1})}{1 + \mathbf{S}_{i-1}(\mathbf{I})} \mathbb{I}\{\bar{T} \geq i + 1\} + \frac{\text{Tr}(\mathbf{Q}\mathbf{B}_{i-1}^2)}{1 + \mathbf{S}_{i-1}(\mathbf{I})} \mathbb{I}(S_i) \\ - \mathbb{E}_{i-1} \left[\frac{\text{Tr}(\mathbf{Q}\mathbf{B}_{i-1}\mathbf{z}_i\mathbf{z}_i^\top\mathbf{B}_{i-1})}{1 + \mathbf{z}_i^\top\mathbf{B}_{i-1}\mathbf{z}_i} \mathbb{I}\{\bar{T} \geq i + 1\} + \frac{\text{Tr}(\mathbf{Q}\mathbf{B}_{i-1}^2)}{1 + \mathbf{S}_{i-1}(\mathbf{I})} \mathbb{I}(S_i) \right]. \quad (65)$$

Since \bar{T} is a stopping time, one can easily have that $D_i(\mathbf{Q}, \bar{T})$ is a martingale difference sequence for $i = 0, 1, \dots, n$. (We note in passing that the above decomposition is similar but does not coincide with the standard Doob decomposition. In particular $B_i(\mathbf{Q}, \bar{T})$ is not measurable on \mathcal{F}_{i-1} . We find the present decomposition more convenient.)

Part II: Controlling the martingale part. We will show $D_i(\mathbf{Q}, \bar{T})$ is bounded and thus by concentration inequality for bounded martingale differences, we can obtain an upper bound for the sum of the $D_i(\mathbf{Q}, \bar{T})$'s. To this end, we use the fact that if for some $m_{i-1} \in \mathcal{F}_{i-1}$

$$|D_i(\mathbf{Q}, \bar{T}) - m_{i-1}| \leq M,$$

then $|D_i(\mathbf{Q}, \bar{T})| = |D_i(\mathbf{Q}, \bar{T}) - \mathbb{E}[D_i(\mathbf{Q}, \bar{T}) | \mathcal{F}_{i-1}]| \leq 2M$. Substitute the following $m_{i-1} \in \mathcal{F}_{i-1}$

$$m_{i-1} = -\frac{\text{Tr}(\mathbf{Q}\mathbf{B}_{i-1}^2)}{1 + \mathbf{S}_{i-1}(\mathbf{I})} \mathbb{I}\{\bar{T} \geq i\} \mathbb{I}(G_{i-1}) + \mathbb{E}_{i-1} \left[\frac{\text{Tr}(\mathbf{Q}\mathbf{B}_{i-1}\mathbf{z}_i\mathbf{z}_i^\top\mathbf{B}_{i-1})}{1 + \mathbf{z}_i^\top\mathbf{B}_{i-1}\mathbf{z}_i} \mathbb{I}\{\bar{T} \geq i + 1\} + \frac{\text{Tr}(\mathbf{Q}\mathbf{B}_{i-1}^2)}{1 + \mathbf{S}_{i-1}(\mathbf{I})} \mathbb{I}(S_i) \right]$$

into the previous display, we have

$$|D_i(\mathbf{Q}, \bar{T}) - m_{i-1}| \\ \stackrel{(i)}{\leq} \left| \frac{\text{Tr}(\mathbf{Q}\mathbf{B}_{i-1}\mathbf{z}_i\mathbf{z}_i^\top\mathbf{B}_{i-1})}{1 + \mathbf{z}_i^\top\mathbf{B}_{i-1}\mathbf{z}_i} - \frac{\text{Tr}(\mathbf{Q}\mathbf{B}_{i-1}^2)}{1 + \mathbf{S}_{i-1}(\mathbf{I})} \right| \mathbb{I}\{\bar{T} \geq i + 1\} \\ \leq \left| \frac{\mathbf{z}_i^\top\mathbf{B}_{i-1}\mathbf{Q}\mathbf{B}_{i-1}\mathbf{z}_i - \text{Tr}(\mathbf{Q}\mathbf{B}_{i-1}^2)}{1 + \mathbf{z}_i^\top\mathbf{B}_{i-1}\mathbf{z}_i} \right| \mathbb{I}\{\bar{T} \geq i + 1\} + \left| \frac{\text{Tr}(\mathbf{Q}\mathbf{B}_{i-1}^2) \cdot (\mathbf{z}_i^\top\mathbf{B}_{i-1}\mathbf{z}_i - \mathbf{S}_{i-1}(\mathbf{I}))}{(1 + \mathbf{z}_i^\top\mathbf{B}_{i-1}\mathbf{z}_i)(1 + \mathbf{S}_{i-1}(\mathbf{I}))} \right| \mathbb{I}\{\bar{T} \geq i + 1\},$$

where in (i) we use $\{\bar{T} \geq i + 1\} \subset F_{i-1}(\mathbf{Q}) \cap F_{i-1}(\mathbf{I}) = G_{i-1}$, and therefore

$$\mathbb{I}\{\bar{T} \geq i\} \mathbb{I}(G_{i-1}) = \mathbb{I}\{\bar{T} = i\} \mathbb{I}(G_{i-1}) + \mathbb{I}\{\bar{T} \geq i + 1\} \mathbb{I}(G_{i-1}) = \mathbb{I}(S_i) + \mathbb{I}\{\bar{T} \geq i + 1\}.$$

Recalling our assumptions for α_1 , we observe that on the event $\{\bar{T} \geq i + 1\} \subset E_i(\mathbf{Q})$,

$$\left| \mathbf{z}_i^\top\mathbf{B}_{i-1}\mathbf{z}_i - \mathbf{S}_{i-1}(\mathbf{I}) \right| \leq \alpha_1 \leq \frac{1}{4}\mathbf{R}_0(\mathbf{I}), \quad \left| \mathbf{z}_i^\top\mathbf{B}_{i-1}\mathbf{Q}\mathbf{B}_{i-1}\mathbf{z}_i - \text{Tr}(\mathbf{Q}\mathbf{B}_{i-1}^2) \right| \leq \alpha_2, \quad (66)$$

and on the event $\{\bar{T} \geq i + 1\} \subset F_{i-1}(\mathbf{I})$ and $\{\bar{T} \geq i + 1\} \subset F_{i-1}(\mathbf{Q})$ by assumptions on β_1 ,

$$|\mathbf{S}_{i-1}(\mathbf{I}) - \mathbf{R}_0(\mathbf{I})| \leq \beta_1 \leq \frac{1}{4}\mathbf{R}_0(\mathbf{I}), \quad |\mathbf{S}_{i-1}(\mathbf{Q}) - \mathbf{R}_0(\mathbf{Q})| \leq \beta_1 \leq \frac{1}{4}\mathbf{R}_0(\mathbf{Q}), \quad (67)$$

and finally on the event $\{\bar{T} \geq i + 1\} \subset G_{i-1}$ it holds

$$\text{Tr}(\mathbf{Q}\mathbf{B}_{i-1}^2) \leq \|\mathbf{B}_{i-1}\| \cdot \text{Tr}(\mathbf{Q}\mathbf{B}_{i-1}) \leq \gamma\mathbf{S}_{i-1}(\mathbf{Q}) = \mathcal{O}(\gamma\mathbf{R}_0(\mathbf{Q})). \quad (68)$$

Putting together bounds in Eqs. (66), (67), (68) and making use of the fact that $\{\bar{T} \geq i+1\} \subset E_i(\mathbf{Q}) \cap F_{i-1}(\mathbf{I}) \cap F_{i-1}(\mathbf{Q})$ yield

$$|D_i(\mathbf{Q}, \bar{T}) - m_{i-1}| \leq \left| \frac{\alpha_2}{1 + \frac{1}{2}\mathbf{R}_0(\mathbf{I})} \right| + \left| \frac{\mathcal{O}(\gamma\mathbf{R}_0(\mathbf{Q})) \cdot \alpha_1}{(1 + \frac{1}{2}\mathbf{R}_0(\mathbf{I})) (1 + \frac{3}{4}\mathbf{R}_0(\mathbf{I}))} \right| = \mathcal{O} \left(\frac{\alpha_1\gamma\mathbf{R}_0(\mathbf{Q}) + \alpha_2(1 + \mathbf{R}_0(\mathbf{I}))}{1 + \mathbf{R}_0(\mathbf{I})^2} \right).$$

Then we can apply Azuma-Hoeffding inequality and obtain

$$\max_{k \leq n} \left| \sum_{i=1}^k D_i(\mathbf{Q}, \bar{T}) \right| = \mathcal{O}_D \left(\sqrt{n \log n} \cdot \frac{\alpha_1\gamma\mathbf{R}_0(\mathbf{Q}) + \alpha_2(1 + \mathbf{R}_0(\mathbf{I}))}{1 + \mathbf{R}_0(\mathbf{I})^2} \right),$$

with probability $1 - \mathcal{O}(n^{-D})$.

Part III: Controlling the bias part. Now we proceed to bound the bias part $|B_i(\mathbf{Q}, \bar{T})|$ in Eq. (65). We can write an upper bound

$$\begin{aligned} |B_i(\mathbf{Q}, \bar{T})| &\leq \underbrace{\left| \frac{\text{Tr}(\mathbf{Q}\mathbf{B}_{i-1}\mathbf{A}_{i-1})}{1 + \mathbf{S}_{i-1}(\mathbf{I})} - \frac{\text{Tr}(\mathbf{Q}\mathbf{B}_{i-1}^2)}{1 + \mathbf{S}_{i-1}(\mathbf{I})} \right|}_{\text{(I)}} \mathbb{I}\{\bar{T} \geq i+1\} \\ &+ \underbrace{\left| \frac{\text{Tr}(\mathbf{Q}\mathbf{B}_{i-1}^2)}{1 + \mathbf{S}_{i-1}(\mathbf{I})} \mathbb{I}\{\bar{T} \geq i\} \mathbb{I}(G_{i-1}) - \mathbb{E}_{i-1} \left[\frac{\text{Tr}(\mathbf{Q}\mathbf{B}_{i-1}\mathbf{z}_i\mathbf{z}_i^\top\mathbf{B}_{i-1})}{1 + \mathbf{z}_i^\top\mathbf{B}_{i-1}\mathbf{z}_i} \mathbb{I}\{\bar{T} \geq i+1\} + \frac{\text{Tr}(\mathbf{Q}\mathbf{B}_{i-1}^2)}{1 + \mathbf{S}_{i-1}(\mathbf{I})} \mathbb{I}(S_i) \right] \right|}_{\text{(II)}}. \end{aligned}$$

Using the fact that

$$\mathbf{B}_{i-1} - \mathbf{A}_{i-1} = (\mu_{i-1} - \mu_i)\mathbf{B}_{i-1}\mathbf{A}_{i-1} = \frac{\mathbf{B}_{i-1}\mathbf{A}_{i-1}}{1 + \mathbf{S}_i(\mathbf{I})},$$

we have

$$\text{(I)} \leq \left| \frac{\text{Tr}(\mathbf{Q}\mathbf{B}_{i-1}(\mathbf{A}_{i-1} - \mathbf{B}_{i-1}))}{1 + \mathbf{S}_{i-1}(\mathbf{I})} \right| \mathbb{I}\{\bar{T} \geq i+1\} = \left| \frac{\text{Tr}(\mathbf{Q}\mathbf{B}_{i-1}^2\mathbf{A}_{i-1})}{(1 + \mathbf{S}_{i-1}(\mathbf{I}))^2} \right| \mathbb{I}\{\bar{T} \geq i+1\}.$$

Upper bounding the term $\text{Tr}(\mathbf{Q}\mathbf{B}_{i-1}^2\mathbf{A}_{i-1})$ requires more careful treatment. Note that \mathbf{B}_{i-1} and \mathbf{A}_{i-1} commute, as follows from the observation that

$$\begin{aligned} &\mathbf{B}_{i-1}\mathbf{A}_{i-1} \\ &= (1 + \mathbf{S}_i(\mathbf{I})) \cdot \Sigma^{\frac{1}{2}} \left\{ \left(\lambda\mathbf{I} + \mu_i\Sigma + \mathbf{X}_{i-1}^\top\mathbf{X}_{i-1} \right)^{-1} \cdot (\mu_{i-1} - \mu_i)\Sigma \cdot \left(\lambda\mathbf{I} + \mu_{i-1}\Sigma + \mathbf{X}_{i-1}^\top\mathbf{X}_{i-1} \right)^{-1} \right\} \Sigma^{\frac{1}{2}} \\ &= (1 + \mathbf{S}_i(\mathbf{I})) (\mathbf{B}_{i-1} - \mathbf{A}_{i-1}) \\ &= (1 + \mathbf{S}_i(\mathbf{I})) \cdot \Sigma^{\frac{1}{2}} \left\{ \left(\lambda\mathbf{I} + \mu_{i-1}\Sigma + \mathbf{X}_{i-1}^\top\mathbf{X}_{i-1} \right)^{-1} \cdot (\mu_{i-1} - \mu_i)\Sigma \cdot \left(\lambda\mathbf{I} + \mu_i\Sigma + \mathbf{X}_{i-1}^\top\mathbf{X}_{i-1} \right)^{-1} \right\} \Sigma^{\frac{1}{2}} \\ &= \mathbf{A}_{i-1}\mathbf{B}_{i-1}. \end{aligned}$$

Since \mathbf{A}_{i-1} and \mathbf{B}_{i-1} are both p.s.d. compact self-adjoint operators in Hilbert space, commutativity implies they can be simultaneously orthogonally diagonalized, which further implies that $\mathbf{A}_{i-1}^{\frac{1}{2}}$ and $\mathbf{B}_{i-1}^{\frac{1}{2}}$ also commute. Therefore

$$\text{Tr}(\mathbf{Q}\mathbf{B}_{i-1}^2\mathbf{A}_{i-1}) = \text{Tr} \left(\mathbf{Q}^{\frac{1}{2}}\mathbf{A}_{i-1}^{\frac{1}{2}}\mathbf{B}_{i-1}^2\mathbf{A}_{i-1}^{\frac{1}{2}}\mathbf{Q}^{\frac{1}{2}} \right) \leq \|\mathbf{B}_{i-1}\|^2 \cdot \text{Tr} \left(\mathbf{Q}^{\frac{1}{2}}\mathbf{A}_{i-1}\mathbf{Q}^{\frac{1}{2}} \right) = \|\mathbf{B}_{i-1}\|^2 \cdot \mathbf{R}_{i-1}(\mathbf{Q}),$$

and thus

$$(I) \leq \frac{\|\mathbf{B}_{i-1}\|^2 \cdot \mathbf{R}_{i-1}(\mathbf{Q})}{(1 + \mathbf{S}_{i-1}(\mathbf{I}))^2} \mathbb{I}\{\bar{T} \geq i+1\} = \mathcal{O}\left(\frac{\gamma^2 \mathbf{R}_0(\mathbf{Q})}{1 + \mathbf{R}_0(\mathbf{I})^2}\right), \quad (69)$$

where in the last inequality we use Eq. (67) on the event $\{\bar{T} \geq i+1\} \subset F_{i-1}(\mathbf{I})$, while $\{\bar{T} \geq i+1\} \subset F_{i-1}(\mathbf{Q})$ also implies

$$|\mathbf{R}_{i-1}(\mathbf{Q}) - \mathbf{R}_0(\mathbf{Q})| \leq \beta_1 \leq \frac{1}{4} \mathbf{R}_0(\mathbf{Q}).$$

Next to bound (II), we make use of the fact that

$$\frac{\text{Tr}(\mathbf{Q}\mathbf{B}_{i-1}^2)}{1 + \mathbf{S}_{i-1}(\mathbf{I})} \mathbb{I}\{\bar{T} \geq i\} \mathbb{I}(G_{i-1}) = \mathbb{E}_{i-1} \left[\frac{\text{Tr}(\mathbf{Q}\mathbf{B}_{i-1}\mathbf{z}_i\mathbf{z}_i^\top\mathbf{B}_{i-1})}{1 + \mathbf{S}_{i-1}(\mathbf{I})} \mathbb{I}\{\bar{T} \geq i\} \mathbb{I}(G_{i-1}) \right],$$

and therefore

$$\begin{aligned} (II) &= \left| \mathbb{E}_{i-1} \left[\frac{\text{Tr}(\mathbf{Q}\mathbf{B}_{i-1}\mathbf{z}_i\mathbf{z}_i^\top\mathbf{B}_{i-1})}{1 + \mathbf{S}_{i-1}(\mathbf{I})} \mathbb{I}\{\bar{T} \geq i\} \mathbb{I}(G_{i-1}) - \frac{\text{Tr}(\mathbf{Q}\mathbf{B}_{i-1}\mathbf{z}_i\mathbf{z}_i^\top\mathbf{B}_{i-1})}{1 + \mathbf{z}_i^\top\mathbf{B}_{i-1}\mathbf{z}_i} \mathbb{I}\{\bar{T} \geq i+1\} - \frac{\text{Tr}(\mathbf{Q}\mathbf{B}_{i-1}^2)}{1 + \mathbf{S}_{i-1}(\mathbf{I})} \mathbb{I}(S_i) \right] \right| \\ &= \left| \mathbb{E}_{i-1} \left[\frac{\mathbf{z}_i^\top\mathbf{B}_{i-1}\mathbf{Q}\mathbf{B}_{i-1}\mathbf{z}_i \cdot (\mathbf{z}_i^\top\mathbf{B}_{i-1}\mathbf{z}_i - \mathbf{S}_{i-1}(\mathbf{I}))}{(1 + \mathbf{S}_{i-1}(\mathbf{I})) (1 + \mathbf{z}_i^\top\mathbf{B}_{i-1}\mathbf{z}_i)} \mathbb{I}\{\bar{T} \geq i+1\} + \frac{\mathbf{z}_i^\top\mathbf{B}_{i-1}\mathbf{Q}\mathbf{B}_{i-1}\mathbf{z}_i - \text{Tr}(\mathbf{Q}\mathbf{B}_{i-1}^2)}{1 + \mathbf{S}_{i-1}(\mathbf{I})} \mathbb{I}(S_i) \right] \right| \\ &\leq \underbrace{\left| \mathbb{E}_{i-1} \left[\left(\frac{\mathbf{z}_i^\top\mathbf{B}_{i-1}\mathbf{Q}\mathbf{B}_{i-1}\mathbf{z}_i \cdot (\mathbf{z}_i^\top\mathbf{B}_{i-1}\mathbf{z}_i - \mathbf{S}_{i-1}(\mathbf{I}))}{(1 + \mathbf{S}_{i-1}(\mathbf{I})) (1 + \mathbf{z}_i^\top\mathbf{B}_{i-1}\mathbf{z}_i)} - \frac{\text{Tr}(\mathbf{Q}\mathbf{B}_{i-1}^2) \cdot (\mathbf{z}_i^\top\mathbf{B}_{i-1}\mathbf{z}_i - \mathbf{S}_{i-1}(\mathbf{I}))}{(1 + \mathbf{S}_{i-1}(\mathbf{I}))^2} \right) \mathbb{I}\{\bar{T} \geq i+1\} \right] \right|}_{(III)} \\ &\quad + \underbrace{\left| \mathbb{E}_{i-1} \left[\left(\frac{\text{Tr}(\mathbf{Q}\mathbf{B}_{i-1}^2) \cdot (\mathbf{z}_i^\top\mathbf{B}_{i-1}\mathbf{z}_i - \mathbf{S}_{i-1}(\mathbf{I}))}{(1 + \mathbf{S}_{i-1}(\mathbf{I}))^2} + \frac{\mathbf{z}_i^\top\mathbf{B}_{i-1}\mathbf{Q}\mathbf{B}_{i-1}\mathbf{z}_i - \text{Tr}(\mathbf{Q}\mathbf{B}_{i-1}^2)}{1 + \mathbf{S}_{i-1}(\mathbf{I})} \right) \mathbb{I}(S_i) \right] \right|}_{(IV)}, \end{aligned} \quad (70)$$

where in the last inequality we use that

$$\begin{aligned} &\mathbb{E}_{i-1} \left[\frac{\text{Tr}(\mathbf{Q}\mathbf{B}_{i-1}^2) \cdot (\mathbf{z}_i^\top\mathbf{B}_{i-1}\mathbf{z}_i - \mathbf{S}_{i-1}(\mathbf{I}))}{(1 + \mathbf{S}_{i-1}(\mathbf{I}))^2} \mathbb{I}\{\bar{T} \geq i\} \cap G_{i-1} \right] \\ &= \frac{\text{Tr}(\mathbf{Q}\mathbf{B}_{i-1}^2)}{(1 + \mathbf{S}_{i-1}(\mathbf{I}))^2} \mathbb{I}(\{\bar{T} \geq i\} \cap G_{i-1}) \cdot \mathbb{E}_{i-1} \left[\mathbf{z}_i^\top\mathbf{B}_{i-1}\mathbf{z}_i - \mathbf{S}_{i-1}(\mathbf{I}) \right] = 0. \end{aligned}$$

To control (III), we note that

$$\begin{aligned} &\frac{\mathbf{z}_i^\top\mathbf{B}_{i-1}\mathbf{Q}\mathbf{B}_{i-1}\mathbf{z}_i \cdot (\mathbf{z}_i^\top\mathbf{B}_{i-1}\mathbf{z}_i - \mathbf{S}_{i-1}(\mathbf{I}))}{(1 + \mathbf{S}_{i-1}(\mathbf{I})) (1 + \mathbf{z}_i^\top\mathbf{B}_{i-1}\mathbf{z}_i)} - \frac{\text{Tr}(\mathbf{Q}\mathbf{B}_{i-1}^2) \cdot (\mathbf{z}_i^\top\mathbf{B}_{i-1}\mathbf{z}_i - \mathbf{S}_{i-1}(\mathbf{I}))}{(1 + \mathbf{S}_{i-1}(\mathbf{I}))^2} \\ &= \frac{\{\mathbf{z}_i^\top\mathbf{B}_{i-1}\mathbf{Q}\mathbf{B}_{i-1}\mathbf{z}_i \cdot (1 + \mathbf{S}_{i-1}(\mathbf{I})) - \text{Tr}(\mathbf{Q}\mathbf{B}_{i-1}^2) \cdot (1 + \mathbf{z}_i^\top\mathbf{B}_{i-1}\mathbf{z}_i)\} \cdot (\mathbf{z}_i^\top\mathbf{B}_{i-1}\mathbf{z}_i - \mathbf{S}_{i-1}(\mathbf{I}))}{(1 + \mathbf{S}_{i-1}(\mathbf{I}))^2 (1 + \mathbf{z}_i^\top\mathbf{B}_{i-1}\mathbf{z}_i)} \end{aligned}$$

$$\begin{aligned}
&= \frac{1}{(1 + \mathbf{S}_{i-1}(\mathbf{I}))^2 (1 + \mathbf{z}_i^\top \mathbf{B}_{i-1} \mathbf{z}_i)} \cdot \left\{ \left(\mathbf{z}_i^\top \mathbf{B}_{i-1} \mathbf{Q} \mathbf{B}_{i-1} \mathbf{z}_i - \text{Tr}(\mathbf{Q} \mathbf{B}_{i-1}^2) \right) \cdot (1 + \mathbf{S}_{i-1}(\mathbf{I})) \right. \\
&\quad \left. - \text{Tr}(\mathbf{Q} \mathbf{B}_{i-1}^2) \cdot \left(\mathbf{z}_i^\top \mathbf{B}_{i-1} \mathbf{z}_i - \mathbf{S}_{i-1}(\mathbf{I}) \right) \right\} \cdot \left(\mathbf{z}_i^\top \mathbf{B}_{i-1} \mathbf{z}_i - \mathbf{S}_{i-1}(\mathbf{I}) \right).
\end{aligned}$$

We again make use of the bounds in Eqs. (66), (67) and (68) on the event $\{\bar{T} \geq i+1\}$, which implies

$$\begin{aligned}
\text{(III)} &\leq \left| \mathbb{E}_{i-1} \left[\frac{(\alpha_2 \cdot (1 + \mathbf{S}_{i-1}(\mathbf{I})) + \text{Tr}(\mathbf{Q} \mathbf{B}_{i-1}^2) \cdot \alpha_1) \cdot \alpha_1}{(1 + \mathbf{S}_{i-1}(\mathbf{I}))^2 (1 + \mathbf{z}_i^\top \mathbf{B}_{i-1} \mathbf{z}_i)} \mathbb{I}\{\bar{T} \geq i+1\} \right] \right| \\
&= \mathcal{O} \left(\frac{\alpha_1 \alpha_2 (1 + \mathbf{R}_0(\mathbf{I})) + \alpha_1^2 \gamma \mathbf{R}_0(\mathbf{Q})}{1 + \mathbf{R}_0(\mathbf{I})^3} \right). \tag{71}
\end{aligned}$$

Finally for term (IV) in Eq. (70), we can control it by

$$\begin{aligned}
\text{(IV)} &\leq \frac{\mathcal{O}(\gamma \mathbf{R}_0(\mathbf{Q}))}{1 + \mathbf{R}_0(\mathbf{I})^2} \cdot \mathbb{E}_{i-1} \left[\left| \mathbf{z}_i^\top \mathbf{B}_{i-1} \mathbf{z}_i - \mathbf{S}_{i-1}(\mathbf{I}) \right| \mathbb{I}(S_i) \right] \\
&\quad + \frac{1}{1 + \mathbf{R}_0(\mathbf{I})} \cdot \mathbb{E}_{i-1} \left[\left| \mathbf{z}_i^\top \mathbf{B}_{i-1} \mathbf{Q} \mathbf{B}_{i-1} \mathbf{z}_i - \text{Tr}(\mathbf{Q} \mathbf{B}_{i-1}^2) \right| \mathbb{I}(S_i) \right].
\end{aligned}$$

Recall $S_i = \{\bar{T} = i\} \cap F_{i-1}(\mathbf{Q}) \cap F_{i-1}(\mathbf{I})$, which implies $\{T_F(\mathbf{Q}) \geq i, T_F(\mathbf{I}) \geq i\}$ holds but at least one of $E_i(\mathbf{Q})$ and $E_i(\mathbf{I})$ doesn't hold. This allows us to invoke Lemma 7.3 and conclude that $\mathbb{P}(S_i | \mathcal{F}_{i-1}) = \mathcal{O}(n^{-D})$. Moreover, we can further deduce from Lemma 7.3 that

$$\begin{aligned}
&\mathbb{E}_{i-1} \left[\left| \mathbf{z}_i^\top \mathbf{B}_{i-1} \mathbf{z}_i - \mathbf{S}_{i-1}(\mathbf{I}) \right| \mathbb{I}(S_i) \right] \\
&\leq \alpha_1 \mathbb{P}(S_i | \mathcal{F}_{i-1}) + \mathbb{E}_{i-1} \left[\left| \mathbf{z}_i^\top \mathbf{B}_{i-1} \mathbf{z}_i - \mathbf{S}_{i-1}(\mathbf{I}) \right| \mathbb{I} \left\{ \left| \mathbf{z}_i^\top \mathbf{B}_{i-1} \mathbf{z}_i - \mathbf{S}_{i-1}(\mathbf{I}) \right| \geq \alpha_1 \right\} \right] \mathbb{I}\{T_F(\mathbf{Q}) \geq i, T_F(\mathbf{I}) \geq i\} \\
&= \mathcal{O}_{\mathbf{C}_{\mathbf{x}}, D} (n^{-D} \cdot \alpha_1).
\end{aligned}$$

Similarly, we also have

$$\mathbb{E}_{i-1} \left[\left| \mathbf{z}_i^\top \mathbf{B}_{i-1} \mathbf{Q} \mathbf{B}_{i-1} \mathbf{z}_i - \text{Tr}(\mathbf{Q} \mathbf{B}_{i-1}^2) \right| \mathbb{I}(S_i) \right] = \mathcal{O}_{\mathbf{C}_{\mathbf{x}}, D} (n^{-D} \cdot \alpha_2).$$

Combining the above displays, we obtain that

$$\text{(IV)} = \mathcal{O}_{\mathbf{C}_{\mathbf{x}}, D} \left(\frac{n^{-D} \alpha_1 \gamma \mathbf{R}_0(\mathbf{Q})}{1 + \mathbf{R}_0(\mathbf{I})^2} + \frac{n^{-D} \alpha_2}{1 + \mathbf{R}_0(\mathbf{I})} \right)$$

Now applying the assumption that

$$n^{-D} = \mathcal{O} \left(\frac{\alpha_1}{1 + \mathbf{R}_0(\mathbf{I})} \right),$$

we obtain

$$\text{(IV)} = \mathcal{O}_{\mathbf{C}_{\mathbf{x}}, D} \left(\frac{\alpha_1^2 \gamma \mathbf{R}_0(\mathbf{Q})}{1 + \mathbf{R}_0(\mathbf{I})^3} + \frac{\alpha_1 \alpha_2}{1 + \mathbf{R}_0(\mathbf{I})^2} \right) = \mathcal{O}_{\mathbf{C}_{\mathbf{x}}, D} \left(\frac{\alpha_1 \alpha_2 (1 + \mathbf{R}_0(\mathbf{I})) + \alpha_1^2 \gamma \mathbf{R}_0(\mathbf{Q})}{1 + \mathbf{R}_0(\mathbf{I})^3} \right). \tag{72}$$

Substitute Eqs. (71) and (72) into Eq. (70) we have

$$\text{(II)} = \mathcal{O}_{\mathbf{C}_{\mathbf{x}}, D} \left(\frac{\alpha_1 \alpha_2 (1 + \mathbf{R}_0(\mathbf{I})) + \alpha_1^2 \gamma \mathbf{R}_0(\mathbf{Q})}{1 + \mathbf{R}_0(\mathbf{I})^3} \right),$$

and together with Eq. (69) we obtain

$$|B_i(\mathbf{Q}, \bar{T})| = \mathcal{O}_{\mathbf{C}_{\mathbf{x}}, D} \left(\frac{\gamma^2 \mathbf{R}_0(\mathbf{Q}) + \alpha_1 \alpha_2}{1 + \mathbf{R}_0(\mathbf{I})^2} + \frac{\alpha_1^2 \gamma \mathbf{R}_0(\mathbf{Q})}{1 + \mathbf{R}_0(\mathbf{I})^3} \right).$$

Part IV: Combining the results. Hence, by combining results in part III and IV, we have with probability $1 - \mathcal{O}(n^{-D})$ that

$$\begin{aligned} & |(\mathbf{R}_{k-1}(\mathbf{Q}) - \mathbf{R}_0(\mathbf{Q}))\mathbb{I}\{\bar{T} \geq k\}| \leq \left| \sum_{i=1}^{k-1} D_i(\mathbf{Q}, \bar{T}) \right| + \left| \sum_{i=1}^{k-1} B_i(\mathbf{Q}, \bar{T}) \right| \\ & \leq \mathcal{O}_D \left(\sqrt{n \log n} \cdot \frac{\alpha_1 \gamma \mathbf{R}_0(\mathbf{Q}) + \alpha_2(1 + \mathbf{R}_0(\mathbf{I}))}{1 + \mathbf{R}_0(\mathbf{I})^2} \right) + n \cdot \mathcal{O}_{\mathbf{C}_x, D} \left(\frac{\gamma^2 \mathbf{R}_0(\mathbf{Q}) + \alpha_1 \alpha_2}{1 + \mathbf{R}_0(\mathbf{I})^2} + \frac{\alpha_1^2 \gamma \mathbf{R}_0(\mathbf{Q})}{1 + \mathbf{R}_0(\mathbf{I})^3} \right) \\ & = \mathcal{O}_{\mathbf{C}_x, D} \left(\sqrt{n \log n} \cdot \frac{\alpha_1 \gamma \mathbf{R}_0(\mathbf{Q}) + \alpha_2(1 + \mathbf{R}_0(\mathbf{I}))}{1 + \mathbf{R}_0(\mathbf{I})^2} + n \cdot \left\{ \frac{\gamma^2 \mathbf{R}_0(\mathbf{Q}) + \alpha_1 \alpha_2}{1 + \mathbf{R}_0(\mathbf{I})^2} + \frac{\alpha_1^2 \gamma \mathbf{R}_0(\mathbf{Q})}{1 + \mathbf{R}_0(\mathbf{I})^3} \right\} \right). \end{aligned}$$

We can first see $\|\mathbf{B}_{k-1}\| \leq \gamma$ holds with probability $1 - \mathcal{O}(n^{-D})$, which follows via exactly the same argument as in Appendix C.3 for $\|\mathbf{A}_k\| \leq \gamma$ by invoking Lemma 7.2. Moreover, we have on $\{\bar{T} \geq k\}$ that

$$\begin{aligned} |\mathbf{S}_{k-1}(\mathbf{Q}) - \mathbf{R}_{k-1}(\mathbf{Q})| &= |\text{Tr}(\mathbf{Q}(\mathbf{B}_{k-1} - \mathbf{A}_{k-1}))| = \frac{\text{Tr}(\mathbf{Q}\mathbf{B}_{k-1}\mathbf{A}_{k-1})}{1 + \mathbf{S}_{k-1}(\mathbf{I})} \\ &\leq \frac{\|\mathbf{B}_{k-1}\| \mathbf{R}_{k-1}(\mathbf{Q})}{1 + \mathbf{S}_{k-1}(\mathbf{I})} \stackrel{(i)}{\leq} \frac{\gamma \mathbf{R}_{k-1}(\mathbf{Q})}{1 + \mathbf{R}_{k-1}(\mathbf{I})} = \mathcal{O} \left(\frac{\gamma \mathbf{R}_0(\mathbf{Q})}{1 + \mathbf{R}_0(\mathbf{I})} \right), \end{aligned}$$

where in (i) we apply $\mu_{k-1} \geq \mu_k$ which indicates $\mathbf{R}_{k-1}(\mathbf{I}) \leq \mathbf{S}_{k-1}(\mathbf{I})$. Therefore, by setting a constant $\mathbf{C}_\beta := \mathbf{C}_\beta(\mathbf{C}_x, D)$ large enough and take

$$\begin{aligned} \beta_1 &= \mathbf{C}_\beta \left(\sqrt{n \log n} \cdot \frac{\alpha_1 \gamma \mathbf{R}_0(\mathbf{Q}) + \alpha_2(1 + \mathbf{R}_0(\mathbf{I}))}{1 + \mathbf{R}_0(\mathbf{I})^2} + n \cdot \left\{ \frac{\gamma^2 \mathbf{R}_0(\mathbf{Q}) + \alpha_1 \alpha_2}{1 + \mathbf{R}_0(\mathbf{I})^2} + \frac{\alpha_1^2 \gamma \mathbf{R}_0(\mathbf{Q})}{1 + \mathbf{R}_0(\mathbf{I})^3} \right\} + \frac{\gamma \mathbf{R}_0(\mathbf{Q})}{1 + \mathbf{R}_0(\mathbf{I})} \right), \\ \beta_2 &= \frac{\mathbf{C}_\beta n \beta_1}{1 + \mathbf{R}_0(\mathbf{I})^2}, \end{aligned}$$

if this satisfies the assumption $\beta_1 \leq \mathbf{R}_0(\mathbf{I})/4$, we can conclude that with probability $1 - \mathcal{O}(n^{-D})$,

$$|(\mathbf{R}_{k-1}(\mathbf{Q}) - \mathbf{R}_0(\mathbf{Q}))\mathbb{I}\{\bar{T} \geq k\}| \leq \beta_1, \quad |(\mathbf{S}_{k-1}(\mathbf{Q}) - \mathbf{R}_0(\mathbf{Q}))\mathbb{I}\{\bar{T} \geq k\}| \leq \beta_1.$$

Further, since

$$\begin{aligned} |(\mu_k - \bar{\mu}_k)\mathbb{I}\{\bar{T} \geq k\}| &= \left| \sum_{i=0}^{k-1} \left(\frac{1}{1 + \mathbf{S}_i(\mathbf{I})} - \frac{1}{1 + \mathbf{R}_0(\mathbf{I})} \right) \mathbb{I}\{\bar{T} \geq k\} \right| \\ &\leq \sum_{i=0}^{k-1} \frac{|\mathbf{S}_i(\mathbf{I}) - \mathbf{R}_0(\mathbf{I})|}{(1 + \mathbf{S}_i(\mathbf{I}))(1 + \mathbf{R}_0(\mathbf{I}))} \mathbb{I}\{\bar{T} \geq k\} = \mathcal{O} \left(\frac{n \beta_1}{1 + \mathbf{R}_0(\mathbf{I})^2} \right), \end{aligned}$$

taking \mathbf{C}_β large will guarantee $|(\mu_k - \bar{\mu}_k)\mathbb{I}\{\bar{T} \geq k\}| \leq \beta_2$. Combining the above displays, we see on the event $\{\bar{T} \geq k\}$, it holds with probability $1 - \mathcal{O}(n^{-D})$ that

$$\max\{|\mathbf{R}_{k-1}(\mathbf{Q}) - \mathbf{R}_0(\mathbf{Q})|, |\mathbf{S}_{k-1}(\mathbf{Q}) - \mathbf{R}_0(\mathbf{Q})|\} \leq \beta_1, \quad |\mu_k - \bar{\mu}_k| \leq \beta_2, \quad \|\mathbf{B}_{k-1}\| \leq \gamma,$$

which is exactly the event $F_{k-1}(\mathbf{Q})$ (cf. Eq. (48b)). Substituting into Eq. (63) completes the proof.

C.5 Proof of Lemma 7.5

As $\beta_2 > \mu/2$, we cannot directly apply Lemmas 7.3 and 7.4. We will instead use a perturbation argument, reducing ourselves to the case $\beta_2 \leq \mu/2$. We will define a second sequence μ'_i following the recursion Eq. (46) but with a different initialization $\mu'_0 = \mu_\star(\lambda, \mu')$ with $\mu' := 64\beta_2 > \mu$. We use the notations

$$\begin{aligned} \mathbf{A}'_i &:= \boldsymbol{\Sigma}^{\frac{1}{2}} \left(\lambda \mathbf{I} + \mu'_i \boldsymbol{\Sigma} + \mathbf{X}_i^\top \mathbf{X}_i \right)^{-1} \boldsymbol{\Sigma}^{\frac{1}{2}}, \\ \mathbf{B}'_i &:= \boldsymbol{\Sigma}^{\frac{1}{2}} \left(\lambda \mathbf{I} + \mu'_{i+1} \boldsymbol{\Sigma} + \mathbf{X}_i^\top \mathbf{X}_i \right)^{-1} \boldsymbol{\Sigma}^{\frac{1}{2}}, \end{aligned}$$

and also denote by $\mathbf{R}'_i(\mathbf{Q}) = \mathcal{R}_i(\lambda, \mu'_i; \mathbf{Q}) := \text{Tr}(\mathbf{Q} \mathbf{A}'_i)$. For this second iteration, we define a parameter tuple $\Delta' = (\alpha'_1, \alpha'_2, \beta'_1, \beta'_2, \gamma)$ defined in Lemmas 7.3 and 7.4 as

$$\gamma' = \min \left\{ \frac{2}{n} \left(1 + \frac{\mathbf{C}_\gamma \mathbf{d}_{\boldsymbol{\Sigma}} \sigma_{[\eta n]} \cdot \log n \log(\mathbf{d}_{\boldsymbol{\Sigma}} n)}{\lambda} \right) + \frac{2}{\mu_\star(\lambda, \mu')}, \frac{1}{\lambda} \right\}, \quad (73a)$$

$$\alpha'_1 = \mathbf{C}_\alpha \log n \cdot \sqrt{\gamma' \mathbf{R}'_0(\mathbf{I})}, \quad (73b)$$

$$\alpha'_2 = \mathbf{C}_\alpha \log n \cdot \sqrt{\gamma'^3 \mathbf{R}'_0(\mathbf{Q})}, \quad (73c)$$

$$\begin{aligned} \beta'_1 &= \mathbf{C}_\beta \left(\sqrt{n \log n} \cdot \frac{\alpha'_1 \gamma' \mathbf{R}'_0(\mathbf{Q}) + \alpha'_2 (1 + \mathbf{R}'_0(\mathbf{I}))}{1 + \mathbf{R}'_0(\mathbf{I})^2} + n \cdot \left\{ \frac{\gamma'^2 \mathbf{R}'_0(\mathbf{Q}) + \alpha'_1 \alpha'_2}{1 + \mathbf{R}'_0(\mathbf{I})^2} + \frac{\alpha'^2_1 \gamma' \mathbf{R}'_0(\mathbf{Q})}{1 + \mathbf{R}'_0(\mathbf{I})^3} \right\} \right. \\ &\quad \left. + \frac{\gamma' \mathbf{R}'_0(\mathbf{Q})}{1 + \mathbf{R}'_0(\mathbf{I})} \right), \end{aligned} \quad (73d)$$

$$\beta'_2 = \frac{\mathbf{C}_\beta n \beta'_1}{1 + \mathbf{R}'_0(\mathbf{I})^2}. \quad (73e)$$

We want to show $\alpha'_1 \leq \mathbf{R}'_0(\mathbf{I})/4$, $\beta'_1 \leq \mathbf{R}'_0(\mathbf{Q})/4$, $\beta'_2 \leq \mu'/2$ and $n^{-D} = \mathcal{O}(\alpha'_1 / (1 + \mathbf{R}'_0(\mathbf{I})))$ so that Lemmas 7.3 and 7.4 are valid for λ, μ' and Δ' . To prove this claim, we need the following result bounding the perturbation of μ_\star .

Lemma C.3. *For any $\lambda > 0$ and $\mu \geq 0$,*

$$0 \leq \frac{\partial \mu_\star(\lambda, \mu)}{\partial \mu} \leq 1 + \mathcal{R}_0(\lambda, \mu_\star(\lambda, \mu); \mathbf{I}).$$

Proof. Taking derivatives w.r.t. μ on both sides of

$$\mu = \mu_\star - \frac{n}{1 + \text{Tr}(\boldsymbol{\Sigma}(\lambda \mathbf{I} + \mu_\star \boldsymbol{\Sigma})^{-1})},$$

we have

$$1 = \frac{\partial \mu_\star}{\partial \mu} \cdot \left(1 - \frac{(\mu_\star - \mu) \text{Tr}(\boldsymbol{\Sigma}^2(\lambda \mathbf{I} + \mu_\star \boldsymbol{\Sigma})^{-2})}{1 + \text{Tr}(\boldsymbol{\Sigma}(\lambda \mathbf{I} + \mu_\star \boldsymbol{\Sigma})^{-1})} \right).$$

Further

$$\begin{aligned} 1 - \frac{(\mu_\star - \mu) \text{Tr}(\boldsymbol{\Sigma}^2(\lambda \mathbf{I} + \mu_\star \boldsymbol{\Sigma})^{-2})}{1 + \text{Tr}(\boldsymbol{\Sigma}(\lambda \mathbf{I} + \mu_\star \boldsymbol{\Sigma})^{-1})} &\geq 1 - \frac{\mu_\star \text{Tr}(\boldsymbol{\Sigma}^2(\lambda \mathbf{I} + \mu_\star \boldsymbol{\Sigma})^{-2})}{1 + \text{Tr}(\boldsymbol{\Sigma}(\lambda \mathbf{I} + \mu_\star \boldsymbol{\Sigma})^{-1})} = \frac{1 + \lambda \text{Tr}(\boldsymbol{\Sigma}(\lambda \mathbf{I} + \mu_\star \boldsymbol{\Sigma})^{-2})}{1 + \text{Tr}(\boldsymbol{\Sigma}(\lambda \mathbf{I} + \mu_\star \boldsymbol{\Sigma})^{-1})} \\ &\geq \frac{1}{1 + \mathcal{R}_0(\lambda, \mu_\star(\lambda, \mu); \mathbf{I})}, \end{aligned}$$

which gives the desired bounds $0 \leq \partial \mu_\star / \partial \mu \leq 1 + \mathcal{R}_0(\lambda, \mu_\star(\lambda, \mu); \mathbf{I})$. \square

Recalling that $\mu_*(\lambda, \mu)$ is increasing in μ and therefore $\mathcal{R}_0(\lambda, \mu_*(\lambda, \mu); \mathbf{I})$ is decreasing in μ , as a direct consequence of Lemma C.3 we have

$$\begin{aligned} 0 \leq \mu_*(\lambda, \mu') - \mu_*(\lambda, \mu) &= \int_{\mu}^{\mu'} \frac{\partial \mu_*}{\partial \nu}(\lambda, \nu) d\nu \\ &\leq \int_{\mu}^{\mu'} (1 + \mathcal{R}_0(\lambda, \mu_*(\lambda, \nu); \mathbf{I})) d\nu \leq (\mu' - \mu)(1 + \mathcal{R}_0(\lambda, \mu_*(\lambda, \mu); \mathbf{I})) \leq \mu'(1 + \mathbf{R}_0(\mathbf{I})), \end{aligned} \quad (74)$$

and further

$$0 \leq \mathbf{R}_0(\mathbf{Q}) - \mathbf{R}'_0(\mathbf{Q}) = (\mu_*(\lambda, \mu') - \mu_*(\lambda, \mu)) \text{Tr}(\mathbf{Q} \mathbf{A}_0 \mathbf{A}'_0) \leq \gamma \mu'(1 + \mathbf{R}_0(\mathbf{I})) \mathbf{R}'_0(\mathbf{Q}),$$

where in the last inequality we used $\|\mathbf{A}_0\| \leq \min\{1/\mu_*(\lambda, \mu), 1/\lambda\} \leq \gamma$. Substituting $\mu' = 64\beta_2$ and using the condition $\gamma\beta_2(1 + \mathbf{R}_0(\mathbf{I})) \leq 1/64$, it then follows that

$$\frac{1}{2} \mathbf{R}_0(\mathbf{Q}) \leq \mathbf{R}'_0(\mathbf{Q}) \leq \mathbf{R}_0(\mathbf{Q}), \quad \forall \text{ p.s.d. } \mathbf{Q}.$$

Using the last inequalities in Eqs. (73a) to (73c), it follows immediately that

$$\gamma' \leq \gamma, \quad \alpha'_1 \leq \alpha_1, \quad \alpha'_2 \leq \alpha_2.$$

We then first see that $\alpha_1 \leq \mathbf{R}_0(\mathbf{I})/8$ implies $\alpha'_1 \leq \alpha_1 \leq \mathbf{R}'_0(\mathbf{I})/4$. For β'_1 and β'_2 , using $1 + \mathbf{R}'_0(\mathbf{I})^k \geq 2^{-k}(1 + \mathbf{R}_0(\mathbf{I}))$ with $k = 1, 2, 3$, we can deduce from Eqs. (73d) and (73e) that

$$\beta'_1 \leq 8\beta_1, \quad \beta'_2 \leq \frac{4C_\beta n \beta'_1}{1 + \mathbf{R}_0(\mathbf{I})^2} \leq 32\beta_2.$$

The last inequality verifies $\beta'_2 \leq 32\beta_2 = \mu'/2$. The condition $\beta_1 \leq \mathbf{R}_0(\mathbf{Q})/64$ implies $\beta'_1 \leq 8\beta_1 \leq \mathbf{R}_0(\mathbf{Q})/8 \leq \mathbf{R}'_0(\mathbf{Q})/4$.

Finally we need to show $n^{-D} = \mathcal{O}(\alpha'_1/(1 + \mathbf{R}'_0(\mathbf{I})))$. From Eq. (74), we can obtain that

$$\mu_*(\lambda, \mu') \leq \mu_*(\lambda, \mu) + 64\beta_2(1 + \mathbf{R}_0(\mathbf{I})) \leq \mu_*(\lambda, \mu) + \frac{1}{\gamma}.$$

Recalling that $\gamma \leq 2/\mu_*(\lambda, \mu)$, we then know $\mu_*(\lambda, \mu') \leq 3/\gamma$ and thus

$$\gamma' \geq \min \left\{ \frac{2}{\mu_*(\lambda, \mu')}, \frac{1}{\lambda} \right\} \geq \frac{2}{3}\gamma.$$

Together with $\mathbf{R}_0(\mathbf{I}) = \Theta(\mathbf{R}'_0(\mathbf{I}))$, we then show $\alpha_1 = \mathcal{O}(\alpha'_1)$ and further that $n^{-D} = \mathcal{O}(\alpha_1/(1 + \mathbf{R}_0(\mathbf{I}))) = \mathcal{O}(\alpha'_1/(1 + \mathbf{R}'_0(\mathbf{I})))$. Hence, we can apply Lemmas 7.3 and 7.4, and by Eq. (51)

$$|\mathcal{R}_n(\lambda, \mu'; \mathbf{Q}) - \mathcal{R}_0(\lambda, \mu_*(\lambda, \mu'); \mathbf{Q})| = \mathcal{O}(\gamma' \beta'_2 \mathbf{R}'_0(\mathbf{Q}) + \beta'_1) = \mathcal{O}(\gamma \beta_2 \mathbf{R}_0(\mathbf{Q}) + \beta_1). \quad (75)$$

In order to finish the perturbation argument, we bound

$$\begin{aligned} |\mathcal{R}_n(\lambda, \mu; \mathbf{Q}) - \mathcal{R}_n(\lambda, \mu'; \mathbf{Q})| &\leq \left| (\mu' - \mu) \cdot \text{Tr} \left(\mathbf{Q} \Sigma^{\frac{1}{2}} (\lambda \mathbf{I} + \mu \Sigma + \mathbf{X}^\top \mathbf{X})^{-1} \Sigma (\lambda \mathbf{I} + \mu' \Sigma + \mathbf{X}^\top \mathbf{X})^{-1} \right) \right| \\ &\leq \mu' \left\| \Sigma^{\frac{1}{2}} (\lambda \mathbf{I} + \mathbf{X}^\top \mathbf{X})^{-1} \Sigma \right\| |\mathcal{R}_n(\lambda, \mu'; \mathbf{Q})| \\ &\stackrel{(i)}{=} \mathcal{O} \left(\gamma \beta_2 \left(\mathcal{R}_0(\lambda, \mu_*(\lambda, \mu'); \mathbf{Q}) + |\mathcal{R}_n(\lambda, \mu'; \mathbf{Q}) - \mathcal{R}_0(\lambda, \mu_*(\lambda, \mu'); \mathbf{Q})| \right) \right) \end{aligned}$$

$$\begin{aligned}
&= \mathcal{O}(\gamma\beta_2(\mathbf{R}_0(\mathbf{Q}) + \gamma\beta_2\mathbf{R}_0(\mathbf{Q}) + \beta_1)) \\
&= \mathcal{O}(\gamma\beta_2\mathbf{R}_0(\mathbf{Q})), \tag{76}
\end{aligned}$$

where in (i) we apply Lemma 7.2 and in the last line we use $\beta_1 = \mathcal{O}(\mathbf{R}_0(\mathbf{Q}))$ and $\gamma\beta_2 = \mathcal{O}((1 + \mathbf{R}_0(\mathbf{I}))^{-1}) = \mathcal{O}(1)$. Similarly, invoke Lemma C.3 and we have

$$\begin{aligned}
|\mathcal{R}_0(\lambda, \mu_*(\lambda, \mu'); \mathbf{Q}) - \mathcal{R}_0(\lambda, \mu_*(\lambda, \mu); \mathbf{Q})| &\leq (\mu_*(\lambda, \mu) - \mu_*(\lambda, \mu'))\gamma\mathcal{R}_0(\lambda, \mu_*(\lambda, \mu); \mathbf{Q}) \\
&\leq \mu'(1 + \mathbf{R}_0(\mathbf{I}))\gamma\mathbf{R}_0(\mathbf{Q}) \\
&= \mathcal{O}(\gamma\beta_2(1 + \mathbf{R}_0(\mathbf{I}))\mathbf{R}_0(\mathbf{Q})). \tag{77}
\end{aligned}$$

By triangular inequality, we deduce from Eqs. (75), (76) and (77) that

$$\begin{aligned}
&|\mathcal{R}_n(\lambda, \mu; \mathbf{Q}) - \mathcal{R}_0(\lambda, \mu_*(\lambda, \mu); \mathbf{Q})| \\
&= \mathcal{O}(\gamma\beta_2\mathbf{R}_0(\mathbf{Q}) + \beta_1) + \mathcal{O}(\gamma\beta_2\mathbf{R}_0(\mathbf{Q})) + \mathcal{O}(\gamma\beta_2(1 + \mathbf{R}_0(\mathbf{I}))\mathbf{R}_0(\mathbf{Q})) \\
&= \mathcal{O}(\gamma\beta_2(1 + \mathbf{R}_0(\mathbf{I}))\mathbf{R}_0(\mathbf{Q}) + \beta_1).
\end{aligned}$$

C.6 Proof of Corollary 6.5

We first derive upper bounds for the parameter $\Delta = (\alpha_1, \alpha_2, \beta_1, \beta_2, \gamma)$ in Theorem 5. Since $\mu_*(\lambda, 0) \leq \mu_*(\lambda, \mu) \leq (1 - \kappa/2)^{-1}\mu_*(\lambda, 0)$, we have

$$\frac{n}{1 + \mathbf{R}_0(\mathbf{I})} \leq \mu_*(\lambda, \mu) \leq \frac{1}{1 - \kappa/2}\mu_*(\lambda, 0) = \frac{1}{1 - \kappa/2} \cdot \frac{n}{1 + \mathcal{R}_0(\lambda, \mu_*(\lambda, 0); \mathbf{I})} \leq \frac{1}{1 - \kappa/2} \cdot \frac{n}{1 + \mathbf{R}_0(\mathbf{I})}.$$

On the other hand, by Eq. (21) and the fact that $\mu_*(\lambda, 0) = \lambda/\lambda_*$, we know

$$\frac{1}{1 + \kappa^{-1}} \leq \kappa \leq \frac{\lambda}{n\lambda_*} \leq \frac{1}{1 + \mathbf{R}_0(\mathbf{I})} \leq \frac{1}{1 - \kappa/2} \frac{\lambda}{n\lambda_*} = \frac{1 - \kappa}{1 - \kappa/2} \leq \frac{1}{1 + \kappa/2},$$

which implies $\kappa/2 \leq \mathbf{R}_0(\mathbf{I}) \leq \kappa^{-1}$. Generalizing the definition of Eq. (23) to $\mu > 0$ and arbitrary \mathbf{Q} , we let $\rho := \mathbf{R}_0(\mathbf{Q})/\mathbf{R}_0(\mathbf{I}) \in (0, 1]$.

Upper bound for γ . First we notice that $\mathbf{d}_\Sigma \geq n$ by Assumption 1 and therefore

$$\log n \log(\mathbf{d}_\Sigma n) \leq \log(\mathbf{d}_\Sigma) \cdot \log(\mathbf{d}_\Sigma^2) = \mathcal{O}(\log^2(\mathbf{d}_\Sigma)),$$

which yields

$$\gamma = \mathcal{O}\left(\frac{2}{n} \left\{ 1 + \frac{\mathcal{O}_{\mathbf{C}_x, D}(\sigma_{\lfloor \eta n \rfloor} \mathbf{d}_\Sigma \log^2(\mathbf{d}_\Sigma))}{\lambda} \right\} + \frac{2}{\mu_*(\lambda, \mu)}\right).$$

Since $\mu_*(\lambda, \mu) \geq \mu_*(\lambda, 0) \geq n\kappa$, we can write

$$\gamma = \frac{1}{n\kappa} \cdot \mathcal{O}_{\mathbf{C}_x, D} \left(1 + \frac{\sigma_{\lfloor \eta n \rfloor} \mathbf{d}_\Sigma \log^2(\mathbf{d}_\Sigma)}{\lambda} \right) = \mathcal{O}_{\mathbf{C}_x, D} \left(\frac{\chi_n(\lambda)}{n\kappa} \right). \tag{78}$$

Upper bounds for α_1 and α_2 . We know that $\mathbf{R}_0(\mathbf{I}) \leq \kappa^{-1}$ and $\mathbf{R}_0(\mathbf{Q}) = \rho\mathbf{R}_0(\mathbf{I})$. As a consequence, we have

$$\alpha_1 = \mathcal{O}_{\mathbf{C}_x, D} \left(\log n \cdot \sqrt{\frac{\gamma}{\kappa}} \right), \quad \alpha_2 = \mathcal{O}_{\mathbf{C}_x, D} \left(\log n \cdot \gamma \sqrt{\frac{\gamma\rho}{\kappa}} \right). \tag{79}$$

Upper bounds for β_1 and β_2 . Using the bounds in the previous displays, we can write

$$\begin{aligned}
\beta_1 &= C_\beta \left(\sqrt{n \log n} \cdot \frac{\alpha_1 \gamma R_0(\mathbf{Q}) + \alpha_2 (1 + R_0(\mathbf{I}))}{1 + R_0(\mathbf{I})^2} + n \cdot \left\{ \frac{\gamma^2 R_0(\mathbf{Q}) + \alpha_1 \alpha_2}{1 + R_0(\mathbf{I})^2} + \frac{\alpha_1^2 \gamma R_0(\mathbf{Q})}{1 + R_0(\mathbf{I})^3} \right\} + \frac{\gamma R_0(\mathbf{Q})}{1 + R_0(\mathbf{I})} \right) \\
&= \mathcal{O}_{\mathbf{C}_{\mathbf{x}}, D} \left(\sqrt{n \log n} \cdot (\alpha_1 \gamma \rho + \alpha_2) + n \cdot \{ (\gamma^2 \rho + \alpha_1 \alpha_2) + \alpha_1^2 \gamma \rho \} + \gamma \rho \right) \\
&= \mathcal{O}_{\mathbf{C}_{\mathbf{x}}, D} \left(\sqrt{n \log n} \cdot \mathcal{O}_{\mathbf{C}_{\mathbf{x}}, D} \left(\log n \cdot \gamma \sqrt{\frac{\gamma \rho}{\kappa}} \right) + n \cdot \mathcal{O}_{\mathbf{C}_{\mathbf{x}}, D} \left(\log^2 n \cdot \gamma^2 \sqrt{\frac{\rho}{\kappa}} \right) + \gamma \rho \right) \\
&= \mathcal{O}_{\mathbf{C}_{\mathbf{x}}, D} \left(\sqrt{n (\log n)^3 \gamma^3} + n (\log n)^2 \gamma^2 + \gamma \sqrt{\kappa} \right) \cdot \sqrt{\frac{\rho}{\kappa}}.
\end{aligned}$$

Substituting in Eq. (78), we can further bound

$$\beta_1 = \mathcal{O}_{\mathbf{C}_{\mathbf{x}}, D} \left(\frac{\sqrt{\rho} \chi_n(\lambda)^2 \log^2 n}{n \kappa^{2.5}} \right). \quad (80)$$

As for β_2 , we simply bound it by

$$\beta_2 = \mathcal{O}_{\mathbf{C}_{\mathbf{x}}, D} (n \beta_1) = \mathcal{O}_{\mathbf{C}_{\mathbf{x}}, D} \left(\frac{\sqrt{\rho} \chi_n(\lambda)^2 \log^2 n}{\kappa^{2.5}} \right). \quad (81)$$

Upper bound for resolvent approximation. Recall the approximation bound we have in Theorem 5,

$$\begin{aligned}
|\mathcal{R}_n(\lambda, \mu; \mathbf{Q}) - \mathcal{R}_0(\lambda, \mu_\star(\lambda, \mu); \mathbf{Q})| &= \mathcal{O}(\gamma \beta_2 (1 + \mathcal{R}_0(\lambda, \mu_\star(\lambda, \mu); \mathbf{I})) \mathcal{R}_0(\lambda, \mu_\star(\lambda, \mu); \mathbf{Q}) + \beta_1) \\
&= \mathcal{O}(\gamma \beta_2 \kappa^{-1} \mathcal{R}_0(\lambda, \mu_\star(\lambda, \mu); \mathbf{Q}) + \beta_1),
\end{aligned}$$

because $\mathcal{R}_0(\lambda, \mu_\star(\lambda, \mu); \mathbf{I}) = R_0(\mathbf{I}) \leq \kappa^{-1}$. And thus

$$\begin{aligned}
|\mathcal{R}_n(\lambda, \mu; \mathbf{Q}) - \mathcal{R}_0(\lambda, \mu_\star(\lambda, \mu); \mathbf{Q})| &= \mathcal{O}(\gamma \beta_2 \kappa^{-1} \cdot \rho \kappa^{-1} + \beta_1) \\
&= \mathcal{O}_{\mathbf{C}_{\mathbf{x}}, D} \left(\frac{\sqrt{\rho^3} \chi_n(\lambda)^3 \log^2 n}{n \kappa^{5.5}} \right) + \mathcal{O}_{\mathbf{C}_{\mathbf{x}}, D} \left(\frac{\sqrt{\rho} \chi_n(\lambda)^2 \log^2 n}{n \kappa^{2.5}} \right) \\
&= \mathcal{O}_{\mathbf{C}_{\mathbf{x}}, D} \left(\frac{\sqrt{\rho} \chi_n(\lambda)^3 \log^2 n}{n \kappa^{5.5}} \right).
\end{aligned}$$

As $R_0(\mathbf{Q}) = \rho R_0(\mathbf{I}) \geq \rho \kappa / 2$, we can also write

$$|\mathcal{R}_n(\lambda, \mu; \mathbf{Q}) - \mathcal{R}_0(\lambda, \mu_\star(\lambda, \mu); \mathbf{Q})| = \mathcal{O}_{\mathbf{C}_{\mathbf{x}}, D} \left(\frac{\chi_n(\lambda)^3 \log^2 n}{n \sqrt{\rho} \cdot \kappa^{6.5}} \right) \mathcal{R}_0(\lambda, \mu_\star(\lambda, \mu); \mathbf{Q}).$$

Simplifying the conditions. Finally we conclude the proof by simplifying the conditions $\alpha_1 \leq R_0(\mathbf{I})/8$, $\beta_1 \leq R_0(\mathbf{Q})/64$, $\gamma \beta_2 (1 + R_0(\mathbf{I})) \leq 1/64$ and $n^{-D} = \mathcal{O}(\alpha_1 / (1 + R_0(\mathbf{I})))$. As $R_0(\mathbf{I}) \geq \kappa/2$, by Eq. (79) it is sufficient to have the first condition once

$$\frac{\chi_n(\lambda) \log^2 n}{\kappa^4} \leq Cn,$$

for some sufficiently small constant $C = C(\mathbf{C}_x, D)$. Recall that $R_0(\mathbf{Q}) \geq \rho\kappa/2$. Therefore, by Eq. (80), the second requirement can be deduced from

$$\frac{\chi_n(\lambda)^2 \log^2 n}{\kappa^{3.5}} \leq C' n \sqrt{\rho},$$

for some sufficiently small constant $C' = C'(\mathbf{C}_x, D)$. By Eqs. (78) and (81), we can derive $\gamma\beta_2(1 + R_0(\mathbf{I})) \leq 1/64$ from

$$\frac{\chi_n(\lambda)^3 \log^2 n}{\kappa^{4.5}} \leq C'' n / \sqrt{\rho},$$

for some constant $C'' = C''(\mathbf{C}_x, D)$. For the last condition, we need a lower bound for $\alpha_1 = C_\alpha \log n \cdot \sqrt{\gamma R_0(\mathbf{I})}$. As $\mu_\star(\lambda, \mu) \leq (1 - \kappa/2)^{-1} \mu_\star(\lambda, 0) \leq (1 - \kappa/2)^{-1} n \leq 2n$, it follows that

$$\begin{aligned} \gamma &= \min \left\{ \frac{2}{n} \left(1 + \frac{C_\gamma d_\Sigma \sigma_{\lfloor \eta n \rfloor} \cdot \log n \log(d_\Sigma n)}{\lambda} \right) + \frac{2}{\mu_\star(\lambda, \mu)}, \frac{1}{\lambda} \right\} \\ &= \Omega \left(\min \left\{ \frac{1}{\mu_\star(\lambda, \mu)}, \frac{1}{\lambda} \right\} \right) = \Omega \left(\min \left\{ \frac{1}{n}, \frac{1}{\lambda} \right\} \right). \end{aligned}$$

With $\kappa/2 \leq R_0(\mathbf{I}) \leq \kappa^{-1}$, we obtain

$$\frac{\alpha_1}{1 + R_0(\mathbf{I})} = \Omega \left(\sqrt{\frac{\kappa^3 \log^2 n}{\max\{n, \lambda\}}} \right).$$

It is then sufficient to have

$$n^{-D} = \mathcal{O} \left(\sqrt{\frac{\kappa^3 \log^2 n}{\max\{n, \lambda\}}} \right).$$

D Proof of Theorem 2

To apply triangle inequalities

$$\begin{aligned} |\mathcal{V}_\mathbf{X}(0) - \mathbf{V}_n(0)| &\leq |\mathcal{V}_\mathbf{X}(\lambda) - \mathbf{V}_n(\lambda)| + |\mathcal{V}_\mathbf{X}(0) - \mathcal{V}_\mathbf{X}(\lambda)| + |\mathbf{V}_n(0) - \mathbf{V}_n(\lambda)|, \\ |\mathcal{B}_\mathbf{X}(0) - \mathbf{B}_n(0)| &\leq |\mathcal{B}_\mathbf{X}(\lambda) - \mathbf{B}_n(\lambda)| + |\mathcal{B}_\mathbf{X}(0) - \mathcal{B}_\mathbf{X}(\lambda)| + |\mathbf{B}_n(0) - \mathbf{B}_n(\lambda)|, \end{aligned} \quad (82)$$

we define $\lambda = \kappa n \lambda_\star$ and bound each term separately. By homogeneity, we will assume $\|\boldsymbol{\theta}\| = 1$ throughout the proof.

Part I: Bounding $|\mathcal{V}_\mathbf{X}(0) - \mathcal{V}_\mathbf{X}(\lambda)|$ and $|\mathcal{B}_\mathbf{X}(0) - \mathcal{B}_\mathbf{X}(\lambda)|$. Assume $\mathbf{X}^\top \mathbf{X}$ has rank r with eigendecomposition $\mathbf{X}^\top \mathbf{X} = \mathbf{U} \mathbf{D} \mathbf{U}^\top$ where $\mathbf{U} \in \mathbb{R}^{d \times r}$ has orthonormal columns and \mathbf{D} is a diagonal matrix with entries $s_1 \geq \dots \geq s_r > 0$. Note that $s_r = n s_{\min}$.

For the variance term, by the elementary inequality $|1/x - x/(x + \lambda)^2| \leq 2\lambda/x^2$ for all $x, \lambda > 0$, we have by Eq. (4a)

$$\begin{aligned} |\mathcal{V}_\mathbf{X}(0) - \mathcal{V}_\mathbf{X}(\lambda)| &= \left| \tau^2 \text{Tr} \left(\boldsymbol{\Sigma} \mathbf{U} (\mathbf{D}^{-1} - \mathbf{D}(\mathbf{D} + \lambda \mathbf{I})^{-2}) \mathbf{U}^\top \right) \right| \\ &\leq \tau^2 \text{Tr} \left(\frac{2\lambda}{s_r} \cdot \boldsymbol{\Sigma} \mathbf{U} \mathbf{D}^{-1} \mathbf{U}^\top \right) = \frac{2\kappa \lambda_\star(\lambda)}{s_{\min}} \cdot \mathcal{V}_\mathbf{X}(0), \end{aligned} \quad (83)$$

where in the last equality we use $\lambda = \kappa n \lambda_\star(\lambda)$ and $s_r = n s_{\min}$. The next lemma bounds the difference between $\lambda_\star(0)$ and $\lambda_\star(\lambda)$.

Lemma D.1. Under the assumptions of Theorem 2, for λ such that $\lambda = \kappa n \lambda_*(\lambda)$ it holds that

$$\lambda_*(0) \leq \lambda_*(\lambda) \leq \left(1 + \frac{2\kappa}{C_\Sigma}\right) \lambda_*(0) \leq 2\lambda_*(0).$$

Proof. Since

$$\lambda = \lambda_* \cdot (n - \text{Tr}(\Sigma(\Sigma + \lambda_* \mathbf{I})^{-1})),$$

we can compute that

$$\begin{aligned} \frac{\partial \lambda}{\partial \lambda_*} &= \lambda_* \cdot \text{Tr}(\Sigma(\Sigma + \lambda_* \mathbf{I})^{-2}) + n - \text{Tr}(\Sigma(\Sigma + \lambda_* \mathbf{I})^{-1}) = n - \text{Tr}(\Sigma^2(\Sigma + \lambda_*(\lambda) \mathbf{I})^{-2}) \\ &\geq n - \text{Tr}(\Sigma^2(\Sigma + \lambda_*(0) \mathbf{I})^{-2}) \geq C_\Sigma n, \end{aligned}$$

and thus

$$\lambda_*(\lambda) = \lambda_*(0) + \int_0^\lambda \frac{\partial \lambda_*(\lambda)}{\partial \lambda} d\lambda \leq \lambda_*(0) + \frac{\lambda}{C_\Sigma n} = \lambda_*(0) + \frac{\lambda}{C_\Sigma n \lambda_*(\lambda)} \cdot \lambda_*(\lambda) = \lambda_*(0) + \frac{\kappa}{C_\Sigma} \cdot \lambda_*(\lambda).$$

Rearranging terms, using $\kappa \leq C_\Sigma^2/8 \leq C_\Sigma/2$ and the fact that $(1-x)^{-1} \leq 1+2x$ for $0 \leq x \leq 1/2$ conclude the proof. \square

Returning to the bound of the variance term, we can thus further derive the upper bound

$$|\mathcal{V}_{\mathbf{X}}(0) - \mathcal{V}_{\mathbf{X}}(\lambda)| \leq \frac{4\kappa\lambda_*(0)}{s_{\min}} \cdot \mathcal{V}_{\mathbf{X}}(0).$$

Using the fact that $\kappa \leq s_{\min}/(8\lambda_*(0))$, we further have

$$|\mathcal{V}_{\mathbf{X}}(0) - \mathcal{V}_{\mathbf{X}}(\lambda)| \leq \left(1 - \frac{4\kappa\lambda_*(0)}{s_{\min}}\right)^{-1} \cdot \frac{4\kappa\lambda_*(0)}{s_{\min}} \cdot \mathcal{V}_{\mathbf{X}}(\lambda) \leq \frac{8\kappa\lambda_*(0)}{s_{\min}} \cdot \mathcal{V}_{\mathbf{X}}(\lambda). \quad (84)$$

Now we look at the bias term. From Eq. (4b), we first have

$$\begin{aligned} \mathcal{B}_{\mathbf{X}}(0) &= \lim_{\lambda \downarrow 0} \lambda^2 \text{Tr} \left(\beta \beta^\top (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I})^{-1} \Sigma (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I})^{-1} \right) \\ &= \lim_{\lambda \downarrow 0} \left\| (\mathbf{X}^\top \mathbf{X} / \lambda + \mathbf{I})^{-1} \beta \right\|_\Sigma^2 \\ &= \lim_{\lambda \downarrow 0} \left\| \left(\mathbf{I} + \mathbf{U} ((\mathbf{D}/\lambda + \mathbf{I})^{-1} - \mathbf{I}) \mathbf{U}^\top \right) \beta \right\|_\Sigma^2 = \left\| (\mathbf{I} - \mathbf{U} \mathbf{U}^\top) \beta \right\|_\Sigma^2. \end{aligned}$$

By triangle inequality, it thus follows

$$\begin{aligned} \left| \mathcal{B}_{\mathbf{X}}(0)^{\frac{1}{2}} - \mathcal{B}_{\mathbf{X}}(\lambda)^{\frac{1}{2}} \right| &\leq \left\| \mathbf{U} (\mathbf{D}/\lambda + \mathbf{I})^{-1} \mathbf{U}^\top \beta \right\|_\Sigma \leq \frac{\lambda}{\lambda + s_r} \|\beta\| = \frac{\kappa \lambda_*(\lambda)}{\kappa \lambda_*(\lambda) + s_{\min}} \|\beta\| \\ &\leq \frac{2\kappa \lambda_*(0)}{2\kappa \lambda_*(0) + s_{\min}} \|\beta\|, \end{aligned}$$

where in the last line we invoke Lemma D.1 and use $\lambda_*(\lambda) \leq 2\lambda_*(0)$.

Additionally with $\mathcal{B}_{\mathbf{X}}(0) \leq \|\beta\|^2$ and

$$\mathcal{B}_{\mathbf{X}}(\lambda) = \lambda^2 \text{Tr} \left(\beta \beta^\top (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I})^{-1} \Sigma (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I})^{-1} \right) \leq \text{Tr} \left(\beta \beta^\top (\mathbf{X}^\top \mathbf{X} / \lambda + \mathbf{I})^{-2} \right) \leq \text{Tr} \left(\beta \beta^\top \right) = \|\beta\|^2,$$

we conclude that

$$|\mathcal{B}_{\mathbf{X}}(0) - \mathcal{B}_{\mathbf{X}}(\lambda)| = \left| \mathcal{B}_{\mathbf{X}}(0)^{\frac{1}{2}} + \mathcal{B}_{\mathbf{X}}(\lambda)^{\frac{1}{2}} \right| \cdot \left| \mathcal{B}_{\mathbf{X}}(0)^{\frac{1}{2}} - \mathcal{B}_{\mathbf{X}}(\lambda)^{\frac{1}{2}} \right| \leq \frac{4\kappa\lambda_{\star}(0) \|\boldsymbol{\beta}\|^2}{2\kappa\lambda_{\star}(0) + s_{\min}} \leq \frac{4\kappa\lambda_{\star}(0) \|\boldsymbol{\beta}\|^2}{s_{\min}}. \quad (85)$$

We obtain an alternative upper bound for $|\mathcal{B}_{\mathbf{X}}(0) - \mathcal{B}_{\mathbf{X}}(\lambda)|$ in the following way. Note that

$$\begin{aligned} \left| \mathcal{B}_{\mathbf{X}}(0)^{\frac{1}{2}} - \mathcal{B}_{\mathbf{X}}(\lambda)^{\frac{1}{2}} \right| &\leq \left\| \mathbf{U}(\mathbf{D}/\lambda + \mathbf{I})^{-1} \mathbf{U}^{\top} \boldsymbol{\beta} \right\|_{\boldsymbol{\Sigma}} = \lambda \left\| \boldsymbol{\Sigma}^{1/2} (\mathbf{X}^{\top} \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^{\top} (\mathbf{X} \mathbf{X}^{\top})^{-1} \mathbf{X} \boldsymbol{\Sigma}^{1/2} \boldsymbol{\theta} \right\| \\ &= \lambda \sqrt{\boldsymbol{\theta}^{\top} \boldsymbol{\Sigma}^{1/2} \mathbf{X}^{\top} (\mathbf{X} \mathbf{X}^{\top})^{-1} \mathbf{X} (\mathbf{X}^{\top} \mathbf{X} + \lambda \mathbf{I})^{-1} \boldsymbol{\Sigma} (\mathbf{X}^{\top} \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^{\top} (\mathbf{X} \mathbf{X}^{\top})^{-1} \mathbf{X} \boldsymbol{\Sigma}^{1/2} \boldsymbol{\theta}} \\ &\leq \lambda \sqrt{\left\| (\mathbf{X}^{\top} \mathbf{X} + \lambda \mathbf{I})^{-\frac{1}{2}} \boldsymbol{\Sigma} (\mathbf{X}^{\top} \mathbf{X} + \lambda \mathbf{I})^{-\frac{1}{2}} \right\| \cdot \sqrt{\boldsymbol{\theta}^{\top} \boldsymbol{\Sigma}^{1/2} (\mathbf{X}^{\top} \mathbf{X} + \lambda \mathbf{I})^{-1} \boldsymbol{\Sigma}^{1/2} \boldsymbol{\theta}}} \\ &= \lambda \sqrt{\left\| \boldsymbol{\Sigma}^{\frac{1}{2}} (\mathbf{X}^{\top} \mathbf{X} + \lambda \mathbf{I})^{-1} \boldsymbol{\Sigma}^{\frac{1}{2}} \right\| \cdot \sqrt{\boldsymbol{\theta}^{\top} \boldsymbol{\Sigma}^{1/2} (\mathbf{X}^{\top} \mathbf{X} + \lambda \mathbf{I})^{-1} \boldsymbol{\Sigma}^{1/2} \boldsymbol{\theta}}}. \end{aligned}$$

We next apply Lemma 7.2, which implies that with probability $1 - \mathcal{O}(n^{-D})$

$$\begin{aligned} \left| \mathcal{B}_{\mathbf{X}}(0)^{\frac{1}{2}} - \mathcal{B}_{\mathbf{X}}(\lambda)^{\frac{1}{2}} \right| &\leq \sqrt{\frac{1}{n} \mathcal{O}_{\mathbf{C}_{\mathbf{x}}, D}(\lambda \chi'_n(\kappa))} \cdot \sqrt{\frac{1}{n} \mathcal{O}_{\mathbf{C}_{\mathbf{x}}, D}(\lambda \chi'_n(\kappa)) \|\boldsymbol{\theta}_{\leq n}\|^2 + 2 \|\boldsymbol{\beta}_{> n}\|^2} \\ &= \sqrt{\mathcal{O}_{\mathbf{C}_{\mathbf{x}}, D}(\kappa^2 \lambda_{\star}(0)^2 \chi'_n(\kappa)^2) \|\boldsymbol{\theta}_{\leq n}\|^2 + \mathcal{O}_{\mathbf{C}_{\mathbf{x}}, D}(\kappa \lambda_{\star}(0) \chi'_n(\kappa)) \|\boldsymbol{\beta}_{> n}\|^2}. \end{aligned}$$

Using the same argument, we can also bound

$$\begin{aligned} \mathcal{B}_{\mathbf{X}}(\lambda) &= \lambda^2 \text{Tr} \left(\boldsymbol{\beta} \boldsymbol{\beta}^{\top} (\mathbf{X}^{\top} \mathbf{X} + \lambda \mathbf{I})^{-1} \boldsymbol{\Sigma} (\mathbf{X}^{\top} \mathbf{X} + \lambda \mathbf{I})^{-1} \right) \\ &\leq \lambda^2 \left\| \boldsymbol{\Sigma}^{1/2} (\mathbf{X}^{\top} \mathbf{X} + \lambda \mathbf{I})^{-1} \boldsymbol{\Sigma}^{1/2} \right\| \cdot \boldsymbol{\theta}^{\top} \boldsymbol{\Sigma}^{1/2} (\mathbf{X}^{\top} \mathbf{X} + \lambda \mathbf{I})^{-1} \boldsymbol{\Sigma}^{1/2} \boldsymbol{\theta} \\ &= \mathcal{O}_{\mathbf{C}_{\mathbf{x}}, D}(\kappa^2 \lambda_{\star}(0)^2 \chi'_n(\kappa)^2) \|\boldsymbol{\theta}_{\leq n}\|^2 + \mathcal{O}_{\mathbf{C}_{\mathbf{x}}, D}(\kappa \lambda_{\star}(0) \chi'_n(\kappa)) \|\boldsymbol{\beta}_{> n}\|^2, \end{aligned}$$

and we can therefore conclude that

$$\begin{aligned} |\mathcal{B}_{\mathbf{X}}(0) - \mathcal{B}_{\mathbf{X}}(\lambda)| &= \left| \mathcal{B}_{\mathbf{X}}(0)^{\frac{1}{2}} + \mathcal{B}_{\mathbf{X}}(\lambda)^{\frac{1}{2}} \right| \cdot \left| \mathcal{B}_{\mathbf{X}}(0)^{\frac{1}{2}} - \mathcal{B}_{\mathbf{X}}(\lambda)^{\frac{1}{2}} \right| \\ &\leq \left(\left| \mathcal{B}_{\mathbf{X}}(0)^{\frac{1}{2}} - \mathcal{B}_{\mathbf{X}}(\lambda)^{\frac{1}{2}} \right| + 2 \mathcal{B}_{\mathbf{X}}(\lambda)^{\frac{1}{2}} \right) \cdot \left| \mathcal{B}_{\mathbf{X}}(0)^{\frac{1}{2}} - \mathcal{B}_{\mathbf{X}}(\lambda)^{\frac{1}{2}} \right| \\ &= \mathcal{O}_{\mathbf{C}_{\mathbf{x}}, D}(\kappa^2 \lambda_{\star}(0)^2 \chi'_n(\kappa)^2) \|\boldsymbol{\theta}_{\leq n}\|^2 + \mathcal{O}_{\mathbf{C}_{\mathbf{x}}, D}(\kappa \lambda_{\star}(0) \chi'_n(\kappa)) \|\boldsymbol{\beta}_{> n}\|^2. \quad (86) \end{aligned}$$

Combining Eqs. (85) and (86), we finally have

$$\begin{aligned} &|\mathcal{B}_{\mathbf{X}}(0) - \mathcal{B}_{\mathbf{X}}(\lambda)| \\ &= \min \left\{ \mathcal{O} \left(\frac{\kappa \lambda_{\star}(0) \|\boldsymbol{\beta}\|^2}{s_{\min}} \right), \mathcal{O}_{\mathbf{C}_{\mathbf{x}}, D}(\kappa^2 \lambda_{\star}(0)^2 \chi'_n(\kappa)^2) \|\boldsymbol{\theta}_{\leq n}\|^2 + \mathcal{O}_{\mathbf{C}_{\mathbf{x}}, D}(\kappa \lambda_{\star}(0) \chi'_n(\kappa)) \|\boldsymbol{\beta}_{> n}\|^2 \right\}. \quad (87) \end{aligned}$$

Part II: Bounding $|\mathbf{V}_n(0) - \mathbf{V}_n(\lambda)|$ and $|\mathbf{B}_n(0) - \mathbf{B}_n(\lambda)|$. Note that

$$0 \geq \frac{\partial \text{Tr} (\boldsymbol{\Sigma}^2 (\boldsymbol{\Sigma} + \lambda_{\star} \mathbf{I})^{-2})}{\partial \lambda_{\star}} = -2 \text{Tr} (\boldsymbol{\Sigma}^2 (\boldsymbol{\Sigma} + \lambda_{\star} \mathbf{I})^{-3}) \geq -\frac{2}{\lambda_{\star}(0)} \text{Tr} (\boldsymbol{\Sigma}^2 (\boldsymbol{\Sigma} + \lambda_{\star}(0) \mathbf{I})^{-2}),$$

we can apply Lemma D.1 and obtain

$$\begin{aligned}
\mathrm{Tr} \left(\boldsymbol{\Sigma}^2 (\boldsymbol{\Sigma} + \lambda_*(0) \mathbf{I})^{-2} \right) &\geq \mathrm{Tr} \left(\boldsymbol{\Sigma}^2 (\boldsymbol{\Sigma} + \lambda_*(\lambda) \mathbf{I})^{-2} \right) \\
&\geq \mathrm{Tr} \left(\boldsymbol{\Sigma}^2 (\boldsymbol{\Sigma} + \lambda_*(0) \mathbf{I})^{-2} \right) - \frac{2(\lambda_*(\lambda) - \lambda_*(0))}{\lambda_*(0)} \mathrm{Tr} \left(\boldsymbol{\Sigma}^2 (\boldsymbol{\Sigma} + \lambda_*(0) \mathbf{I})^{-2} \right) \\
&\geq \left(1 - \frac{4\kappa}{\mathbf{C}_\Sigma} \right) \cdot \mathrm{Tr} \left(\boldsymbol{\Sigma}^2 (\boldsymbol{\Sigma} + \lambda_*(0) \mathbf{I})^{-2} \right).
\end{aligned}$$

We then have

$$\begin{aligned}
\mathbf{V}_n(0) \geq \mathbf{V}_n(\lambda) &= \frac{n - \mathrm{Tr} \left(\boldsymbol{\Sigma}^2 (\boldsymbol{\Sigma} + \lambda_*(0) \mathbf{I})^{-2} \right)}{n - \mathrm{Tr} \left(\boldsymbol{\Sigma}^2 (\boldsymbol{\Sigma} + \lambda_*(\lambda) \mathbf{I})^{-2} \right)} \cdot \frac{\mathrm{Tr} \left(\boldsymbol{\Sigma}^2 (\boldsymbol{\Sigma} + \lambda_*(\lambda) \mathbf{I})^{-2} \right)}{\mathrm{Tr} \left(\boldsymbol{\Sigma}^2 (\boldsymbol{\Sigma} + \lambda_*(0) \mathbf{I})^{-2} \right)} \cdot \mathbf{V}_n(0) \\
&\stackrel{(i)}{\geq} \frac{\mathbf{C}_\Sigma n}{\mathbf{C}_\Sigma n + \frac{4\kappa}{\mathbf{C}_\Sigma} \cdot \mathrm{Tr} \left(\boldsymbol{\Sigma}^2 (\boldsymbol{\Sigma} + \lambda_*(0) \mathbf{I})^{-2} \right)} \cdot \left(1 - \frac{4\kappa}{\mathbf{C}_\Sigma} \right) \cdot \mathbf{V}_n(0) \\
&\stackrel{(ii)}{\geq} \frac{\mathbf{C}_\Sigma^2}{\mathbf{C}_\Sigma^2 + 4\kappa} \cdot \left(1 - \frac{4\kappa}{\mathbf{C}_\Sigma} \right) \cdot \mathbf{V}_n(0),
\end{aligned}$$

where we use $n - \mathrm{Tr} \left(\boldsymbol{\Sigma}^2 (\boldsymbol{\Sigma} + \lambda_*(0) \mathbf{I})^{-2} \right) \geq \mathbf{C}_\Sigma n$ in (i) and $n \geq \mathrm{Tr} \left(\boldsymbol{\Sigma}^2 (\boldsymbol{\Sigma} + \lambda_*(0) \mathbf{I})^{-2} \right)$ in (ii). By the elementary inequality $1 - (1-a)(1-b) \leq a+b$ for all $0 \leq a, b \leq 1$, we can thus derive that

$$\begin{aligned}
|\mathbf{V}_n(0) - \mathbf{V}_n(\lambda)| &\leq \left\{ 1 - \left(1 - \frac{4\kappa}{\mathbf{C}_\Sigma^2 + 4\kappa} \right) \cdot \left(1 - \frac{4\kappa}{\mathbf{C}_\Sigma} \right) \right\} \cdot \mathbf{V}_n(0) \\
&\leq \left(\frac{4\kappa}{\mathbf{C}_\Sigma^2 + 4\kappa} + \frac{4\kappa}{\mathbf{C}_\Sigma} \right) \cdot \mathbf{V}_n(0) \leq \frac{8\kappa}{\mathbf{C}_\Sigma^2} \cdot \mathbf{V}_n(0).
\end{aligned} \tag{88}$$

For the bias term, we first similarly derive

$$\boldsymbol{\beta}^\top (\boldsymbol{\Sigma} + \lambda_*(0) \mathbf{I})^{-2} \boldsymbol{\Sigma} \boldsymbol{\beta} \geq \boldsymbol{\beta}^\top (\boldsymbol{\Sigma} + \lambda_*(\lambda) \mathbf{I})^{-2} \boldsymbol{\Sigma} \boldsymbol{\beta} \geq \left(1 - \frac{4\kappa}{\mathbf{C}_\Sigma} \right) \cdot \boldsymbol{\beta}^\top (\boldsymbol{\Sigma} + \lambda_*(0) \mathbf{I})^{-2} \boldsymbol{\Sigma} \boldsymbol{\beta}.$$

Note that

$$\begin{aligned}
\left| \frac{\mathbf{B}_n(\lambda)}{\mathbf{B}_n(0)} - 1 \right| &= \left| \frac{n - \mathrm{Tr} \left(\boldsymbol{\Sigma}^2 (\boldsymbol{\Sigma} + \lambda_*(0) \mathbf{I})^{-2} \right)}{n - \mathrm{Tr} \left(\boldsymbol{\Sigma}^2 (\boldsymbol{\Sigma} + \lambda_*(\lambda) \mathbf{I})^{-2} \right)} \cdot \frac{\boldsymbol{\beta}^\top (\boldsymbol{\Sigma} + \lambda_*(\lambda) \mathbf{I})^{-2} \boldsymbol{\Sigma} \boldsymbol{\beta}}{\boldsymbol{\beta}^\top (\boldsymbol{\Sigma} + \lambda_*(0) \mathbf{I})^{-2} \boldsymbol{\Sigma} \boldsymbol{\beta}} \cdot \frac{\lambda_*(\lambda)^2}{\lambda_*(0)^2} - 1 \right| \\
&\leq \max \left\{ 1 - \frac{n - \mathrm{Tr} \left(\boldsymbol{\Sigma}^2 (\boldsymbol{\Sigma} + \lambda_*(0) \mathbf{I})^{-2} \right)}{n - \mathrm{Tr} \left(\boldsymbol{\Sigma}^2 (\boldsymbol{\Sigma} + \lambda_*(\lambda) \mathbf{I})^{-2} \right)} \cdot \frac{\boldsymbol{\beta}^\top (\boldsymbol{\Sigma} + \lambda_*(\lambda) \mathbf{I})^{-2} \boldsymbol{\Sigma} \boldsymbol{\beta}}{\boldsymbol{\beta}^\top (\boldsymbol{\Sigma} + \lambda_*(0) \mathbf{I})^{-2} \boldsymbol{\Sigma} \boldsymbol{\beta}}; \frac{\lambda_*(\lambda)^2}{\lambda_*(0)^2} - 1 \right\},
\end{aligned}$$

From the previous calculations for the variance term, we know

$$1 - \frac{n - \mathrm{Tr} \left(\boldsymbol{\Sigma}^2 (\boldsymbol{\Sigma} + \lambda_*(0) \mathbf{I})^{-2} \right)}{n - \mathrm{Tr} \left(\boldsymbol{\Sigma}^2 (\boldsymbol{\Sigma} + \lambda_*(\lambda) \mathbf{I})^{-2} \right)} \cdot \frac{\boldsymbol{\beta}^\top (\boldsymbol{\Sigma} + \lambda_*(\lambda) \mathbf{I})^{-2} \boldsymbol{\Sigma} \boldsymbol{\beta}}{\boldsymbol{\beta}^\top (\boldsymbol{\Sigma} + \lambda_*(0) \mathbf{I})^{-2} \boldsymbol{\Sigma} \boldsymbol{\beta}} \leq \frac{8\kappa}{\mathbf{C}_\Sigma^2},$$

and by Lemma D.1 we have

$$\frac{\lambda_*(\lambda)^2}{\lambda_*(0)^2} - 1 \leq \left(1 + \frac{2\kappa}{\mathbf{C}_\Sigma} \right)^2 - 1 \leq \frac{4\kappa}{\mathbf{C}_\Sigma} + \frac{2\kappa}{\mathbf{C}_\Sigma} \cdot \frac{2\kappa}{\mathbf{C}_\Sigma} \leq \frac{6\kappa}{\mathbf{C}_\Sigma}.$$

In the last inequality, recall $\kappa \leq \mathbf{C}_\Sigma^2/8 \leq \mathbf{C}_\Sigma/2$. Putting together, we have error of the bias term bounded by

$$|\mathbf{B}_n(0) - \mathbf{B}_n(\lambda)| \leq \frac{8\kappa}{\mathbf{C}_\Sigma^2} \cdot \mathbf{B}_n(0). \tag{89}$$

Part III: Variance approximation when $\lambda = 0$. Recalling that $\lambda = \kappa n \lambda_*(\lambda)$, we want to invoke Theorem 1 to bound $|\mathcal{V}_{\mathbf{X}}(\lambda) - \mathbf{V}_n(\lambda)|$. Note that by Lemma D.1 it holds $\lambda_*(\lambda) = \Theta(\lambda_*(0))$ and thus

$$\begin{aligned}\chi_n(\lambda) &= 1 + \frac{\sigma_{\lfloor \eta n \rfloor} \mathbf{d}_{\Sigma} \log^2(\mathbf{d}_{\Sigma})}{\lambda} = 1 + \frac{\sigma_{\lfloor \eta n \rfloor} \mathbf{d}_{\Sigma} \log^2(\mathbf{d}_{\Sigma})}{\kappa n \lambda_*(\lambda)} = \Theta \left(1 + \frac{\sigma_{\lfloor \eta n \rfloor} \mathbf{d}_{\Sigma} \log^2(\mathbf{d}_{\Sigma})}{\kappa n \lambda_*(0)} \right) \\ &= \Theta(\chi'_n(\kappa)).\end{aligned}$$

Hence the conditions hold for Theorem 1 by taking $\mathbf{C}_1 = \Theta(\mathbf{C})$, and we have for some constant $\mathbf{C}' := \mathbf{C}'(k, \mathbf{C}_{\mathbf{x}}, D) > 0$,

$$|\mathcal{V}_{\mathbf{X}}(\lambda) - \mathbf{V}_n(\lambda)| \leq \mathbf{C}' \cdot \frac{\chi'_n(\kappa)^3 \log^2 n}{n^{1-\frac{1}{k}} \kappa^{9.5}} \cdot \mathbf{V}_n(\lambda).$$

Substituting the above display and Eqs. (84), (88) into Eq. (82) yields

$$\begin{aligned}|\mathcal{V}_{\mathbf{X}}(0) - \mathbf{V}_n(0)| &\leq \mathbf{C}' \cdot \frac{\chi'_n(\kappa)^3 \log^2 n}{n^{1-\frac{1}{k}} \kappa^{9.5}} \cdot \mathbf{V}_n(\lambda) + \frac{8\kappa \lambda_*(0)}{s_{\min}} \cdot \mathcal{V}_{\mathbf{X}}(\lambda) + \frac{8\kappa}{\mathbf{C}_{\Sigma}^2} \cdot \mathbf{V}_n(0) \\ &\leq \mathbf{C}' \cdot \frac{\chi'_n(\kappa)^3 \log^2 n}{n^{1-\frac{1}{k}} \kappa^{9.5}} \cdot \mathbf{V}_n(\lambda) + \frac{8\kappa \lambda_*(0)}{s_{\min}} \cdot \left(1 + \mathbf{C}' \cdot \frac{\chi'_n(\kappa)^3 \log^2 n}{n^{1-\frac{1}{k}} \kappa^{9.5}} \right) \mathbf{V}_n(\lambda) + \frac{8\kappa}{\mathbf{C}_{\Sigma}^2} \cdot \mathbf{V}_n(0) \\ &= \left\{ \left(1 + \frac{8\kappa \lambda_*(0)}{s_{\min}} \right) \left(1 + \mathbf{C}' \cdot \frac{\chi'_n(\kappa)^3 \log^2 n}{n^{1-\frac{1}{k}} \kappa^{9.5}} \right) - 1 \right\} \cdot \mathbf{V}_n(\lambda) + \frac{8\kappa}{\mathbf{C}_{\Sigma}^2} \cdot \mathbf{V}_n(0) \\ &\leq \left\{ \left(1 + \frac{8\kappa \lambda_*(0)}{s_{\min}} \right) \left(1 + \mathbf{C}' \cdot \frac{\chi'_n(\kappa)^3 \log^2 n}{n^{1-\frac{1}{k}} \kappa^{9.5}} \right) - 1 \right\} \cdot \left(1 + \frac{8\kappa}{\mathbf{C}_{\Sigma}^2} \right) \mathbf{V}_n(0) + \frac{8\kappa}{\mathbf{C}_{\Sigma}^2} \cdot \mathbf{V}_n(0) \\ &\leq \left\{ \left(1 + \frac{8\kappa \lambda_*(0)}{s_{\min}} \right) \left(1 + \mathbf{C}' \cdot \frac{\chi'_n(\kappa)^3 \log^2 n}{n^{1-\frac{1}{k}} \kappa^{9.5}} \right) \left(1 + \frac{8\kappa}{\mathbf{C}_{\Sigma}^2} \right) - 1 \right\} \cdot \mathbf{V}_n(0) \\ &\leq \left(\exp \left(\frac{8\kappa \lambda_*(0)}{s_{\min}} + \frac{8\kappa}{\mathbf{C}_{\Sigma}^2} + \mathbf{C}' \cdot \frac{\chi'_n(\kappa)^3 \log^2 n}{n^{1-\frac{1}{k}} \kappa^{9.5}} \right) - 1 \right) \cdot \mathbf{V}_n(0).\end{aligned}$$

Since $\kappa \leq s_{\min}/(8\lambda_*(0))$ and $\kappa \leq \mathbf{C}_{\Sigma}^2/8$, if we additionally assume

$$\frac{\chi'_n(\kappa)^3 \log^2 n}{n^{1-\frac{1}{k}} \kappa^{9.5}} \leq \frac{1}{\mathbf{C}'},$$

we can then conclude that

$$\exp \left(\frac{8\kappa \lambda_*(0)}{s_{\min}} + \frac{8\kappa}{\mathbf{C}_{\Sigma}^2} + \mathbf{C}' \cdot \frac{\chi'_n(\kappa)^3 \log^2 n}{n^{1-\frac{1}{k}} \kappa^{9.5}} \right) - 1 = \mathcal{O} \left(\frac{8\kappa \lambda_*(0)}{s_{\min}} + \frac{8\kappa}{\mathbf{C}_{\Sigma}^2} + \mathbf{C}' \cdot \frac{\chi'_n(\kappa)^3 \log^2 n}{n^{1-\frac{1}{k}} \kappa^{9.5}} \right),$$

and

$$|\mathcal{V}_{\mathbf{X}}(0) - \mathbf{V}_n(0)| = \mathcal{O}_{k, \mathbf{C}_{\mathbf{x}}, D} \left(\kappa \cdot \left(\frac{\lambda_*(0)}{s_{\min}} + \frac{1}{\mathbf{C}_{\Sigma}^2} \right) + \frac{\chi'_n(\kappa)^3 \log^2 n}{n^{1-\frac{1}{k}} \kappa^{9.5}} \right) \cdot \mathbf{V}_n(0).$$

We meet this assumption by setting $\mathbf{C}_2 = 1/\mathbf{C}'$.

Part IV: Bias approximation when $\lambda = 0$. To apply Theorem 1 when $\lambda = \kappa n \lambda_*(\lambda)$, we first note that the condition $\lambda \kappa n^{-\frac{1}{k}} \leq n \kappa / 2$ is equivalent to $\lambda_*(\lambda) \kappa n^{-\frac{1}{k}} \leq 1/2$, and by Lemma D.1 it suffices to have $\lambda_*(0) \kappa n^{-\frac{1}{k}} \leq 1/4$, which holds by assumption. Since we know $\chi'_n(\kappa) = \Theta(\chi_n(\lambda))$ from the previous part of the proof, we only need to additionally verify that $\lambda_*(0) = \Theta(\lambda_*(\lambda))$ and $\rho(0) = \Theta(\rho(\lambda))$. The first relation is a direct consequence of Lemma D.1, and for the second claim we observe that

$$\rho(\lambda) = \frac{\mathcal{R}_0(\lambda_*(\lambda), 1; \boldsymbol{\theta} \boldsymbol{\theta}^\top)}{\mathcal{R}_0(\lambda_*(\lambda), 1; \mathbf{I})} = \frac{\text{Tr} \left(\boldsymbol{\Sigma}^{\frac{1}{2}} \boldsymbol{\theta} \boldsymbol{\theta}^\top \boldsymbol{\Sigma}^{\frac{1}{2}} (\boldsymbol{\Sigma} + \lambda_*(\lambda) \mathbf{I})^{-1} \right)}{\text{Tr} \left(\boldsymbol{\Sigma} (\boldsymbol{\Sigma} + \lambda_*(\lambda) \mathbf{I})^{-1} \right)}.$$

As for any p.s.d. \mathbf{Q} ,

$$0 \geq \frac{\partial \mathcal{R}_0(\lambda_*, 1; \mathbf{Q})}{\partial \lambda_*} = -\text{Tr} \left(\boldsymbol{\Sigma}^{\frac{1}{2}} \mathbf{Q} \boldsymbol{\Sigma}^{\frac{1}{2}} (\boldsymbol{\Sigma} + \lambda_* \mathbf{I})^{-2} \right) \geq -\frac{1}{\lambda_*(0)} \text{Tr} \left(\boldsymbol{\Sigma}^{\frac{1}{2}} \mathbf{Q} \boldsymbol{\Sigma}^{\frac{1}{2}} (\boldsymbol{\Sigma} + \lambda_*(0) \mathbf{I})^{-1} \right),$$

we have

$$\mathcal{R}_0(\lambda_*(0), 1; \mathbf{Q}) \geq \mathcal{R}_0(\lambda_*(\lambda), 1; \mathbf{Q}) \geq \mathcal{R}_0(\lambda_*(0), 1; \mathbf{Q}) - \frac{\lambda_*(\lambda) - \lambda_*(0)}{\lambda_*(0)} \cdot \mathcal{R}_0(\lambda_*(0), 1; \mathbf{Q}).$$

Therefore by Lemma D.1 and $\kappa \leq \mathbf{C}_\Sigma^2/8$, we can obtain

$$\left| \frac{\mathcal{R}_0(\lambda_*(\lambda), 1; \mathbf{Q})}{\mathcal{R}_0(\lambda_*(0), 1; \mathbf{Q})} - 1 \right| \leq \left| \frac{\lambda_*(\lambda)}{\lambda_*(0)} - 1 \right| \leq \frac{2\kappa}{\mathbf{C}_\Sigma} \leq \frac{\mathbf{C}_\Sigma}{4} \leq \frac{1}{4},$$

which implies $\mathcal{R}_0(\lambda_*(\lambda), 1; \mathbf{Q})/\mathcal{R}_0(\lambda_*(0), 1; \mathbf{Q}) = \Theta(1)$ and therefore $\rho(0) = \Theta(\rho(\lambda))$. Now we are able to invoke Theorem 2, yielding for some constant $C' := C'(k, \mathbf{C}_x, D) > 0$,

$$|\mathcal{B}_X(\lambda) - \mathbf{B}_n(\lambda)| \leq C' \cdot \left(\frac{\lambda_*(0)^{k+1}}{n \kappa^3} + \frac{\chi'_n(\kappa)^3 \log^2 n}{\sqrt{\rho(0)} n^{1-\frac{1}{k}} \kappa^{8.5}} \right) \cdot \mathbf{B}_n(\lambda).$$

Now we can substitute the above bound and Eq. (89) into Eq. (82),

$$\begin{aligned} & |\mathcal{B}_X(0) - \mathbf{B}_n(0)| \\ & \leq C' \cdot \left(\frac{\lambda_*(0)^{k+1}}{n \kappa^3} + \frac{\chi'_n(\kappa)^3 \log^2 n}{\sqrt{\rho(0)} n^{1-\frac{1}{k}} \kappa^{8.5}} \right) \cdot \mathbf{B}_n(\lambda) + |\mathcal{B}_X(0) - \mathcal{B}_X(\lambda)| + \frac{8\kappa}{\mathbf{C}_\Sigma^2} \cdot \mathbf{B}_n(0) \\ & \leq \left\{ \left(1 + C' \cdot \left(\frac{\lambda_*(0)^{k+1}}{n \kappa^3} + \frac{\chi'_n(\kappa)^3 \log^2 n}{\sqrt{\rho(0)} n^{1-\frac{1}{k}} \kappa^{8.5}} \right) \right) \left(1 + \frac{8\kappa}{\mathbf{C}_\Sigma^2} \right) - 1 \right\} \cdot \mathbf{B}_n(0) + |\mathcal{B}_X(0) - \mathcal{B}_X(\lambda)| \\ & \leq \left(\exp \left(C' \cdot \left(\frac{\lambda_*(0)^{k+1}}{n \kappa^3} + \frac{\chi'_n(\kappa)^3 \log^2 n}{\sqrt{\rho(0)} n^{1-\frac{1}{k}} \kappa^{8.5}} \right) + \frac{8\kappa}{\mathbf{C}_\Sigma^2} \right) - 1 \right) \cdot \mathbf{B}_n(0) + |\mathcal{B}_X(0) - \mathcal{B}_X(\lambda)|. \end{aligned}$$

Similar to previous calculations for the variance approximation, setting $\mathbf{C}_3 = 1/C'$ and thus

$$\frac{\lambda_*(0)^{k+1}}{n \kappa^3} + \frac{\chi'_n(\kappa)^3 \log^2 n}{\sqrt{\rho(0)} n^{1-\frac{1}{k}} \kappa^{8.5}} \leq \frac{1}{C'}.$$

Substituting in Eq. (87), it then holds that

$$|\mathcal{B}_X(0) - \mathbf{B}_n(0)|$$

$$\begin{aligned}
&= \mathcal{O}_{k, \mathbf{C}_\Sigma, D} \left(\frac{\kappa}{\mathbf{C}_\Sigma^2} + \frac{\lambda_\star(0)^{k+1}}{n\kappa^3} + \frac{\chi'_n(\kappa)^3 \log^2 n}{\sqrt{\rho(0)} n^{1-\frac{1}{k}} \kappa^{8.5}} \right) \cdot \mathbf{B}_n(0) + |\mathcal{B}_\mathbf{X}(0) - \mathcal{B}_\mathbf{X}(\lambda)| \\
&= \mathcal{O}_{k, \mathbf{C}_\Sigma, D} \left(\frac{\kappa}{\mathbf{C}_\Sigma^2} + \frac{\lambda_\star(0)^{k+1}}{n\kappa^3} + \frac{\chi'_n(\kappa)^3 \log^2 n}{\sqrt{\rho(0)} n^{1-\frac{1}{k}} \kappa^{8.5}} \right) \cdot \mathbf{B}_n(0) \\
&\quad + \min \left\{ \mathcal{O} \left(\frac{\kappa \lambda_\star(0) \|\boldsymbol{\beta}\|^2}{s_{\min}} \right), \mathcal{O}_{\mathbf{C}_\Sigma, D} (\kappa^2 \lambda_\star(0)^2 \chi'_n(\kappa)^2) \|\boldsymbol{\theta}_{\leq n}\|^2 + \mathcal{O}_{\mathbf{C}_\Sigma, D} (\kappa \lambda_\star(0) \chi'_n(\kappa)) \|\boldsymbol{\beta}_{> n}\|^2 \right\}.
\end{aligned}$$

E Proof of Theorem 3

We follow the same proof strategy in Appendix D for the overparameterized regime, taking $\lambda = \varepsilon n$.

Part I: Bounding $|\mathcal{V}_\mathbf{X}(0) - \mathcal{V}_\mathbf{X}(\lambda)|$. As we assume $\mathbf{X}^\top \mathbf{X}$ has rank d , we can write its eigendecomposition $\mathbf{X}^\top \mathbf{X} = \mathbf{U} \mathbf{D} \mathbf{U}^\top$ with $\mathbf{U} \in \mathbb{R}^{d \times d}$ an orthogonal matrix and \mathbf{D} is a diagonal matrix with entries $s_1 \geq \dots \geq s_d > 0$. In this case, $s_d = n s_{\min}$. Substitute $\lambda = \varepsilon n$ into Eq. (83) instead of $\lambda = \kappa n \lambda_\star(\lambda)$, we have

$$|\mathcal{V}_\mathbf{X}(0) - \mathcal{V}_\mathbf{X}(\lambda)| \leq \frac{2\varepsilon}{s_{\min}} \cdot \mathcal{V}_\mathbf{X}(0). \quad (90)$$

Part II: Bounding $|\mathbf{V}_n(0) - \mathbf{V}_n(\lambda)|$. Similar to the overparameterized case, we can control the growth of $\lambda_\star(\lambda)$ by

Lemma E.1. *Under the assumptions of Theorem 3, for λ such that $\lambda = \varepsilon n$ it holds that*

$$0 = \lambda_\star(0) \leq \lambda_\star(\lambda) \leq \frac{\varepsilon}{\mathbf{C}_\Sigma}.$$

Proof. By the proof of Lemma D.1, we have

$$\frac{\partial \lambda}{\partial \lambda_\star} = n - \text{Tr}(\boldsymbol{\Sigma}^2 (\boldsymbol{\Sigma} + \lambda_\star(\lambda) \mathbf{I})^{-2}) \geq n - d \geq \mathbf{C}_\Sigma n,$$

and thus

$$0 \leq \lambda_\star(0) \leq \lambda_\star(\lambda) = \int_0^\lambda \frac{\partial \lambda_\star(\nu)}{\partial \lambda} d\nu \leq \frac{\lambda}{\mathbf{C}_\Sigma n} = \frac{\varepsilon}{\mathbf{C}_\Sigma}.$$

□

In this case, note that

$$0 \geq \frac{\partial \text{Tr}(\boldsymbol{\Sigma}^2 (\boldsymbol{\Sigma} + \lambda_\star \mathbf{I})^{-2})}{\partial \lambda_\star} = -2 \text{Tr}(\boldsymbol{\Sigma}^2 (\boldsymbol{\Sigma} + \lambda_\star \mathbf{I})^{-3}) \geq -\frac{2}{\sigma_d} \text{Tr}(\boldsymbol{\Sigma}^2 (\boldsymbol{\Sigma} + \lambda_\star(0) \mathbf{I})^{-2}) = -\frac{2d}{\sigma_d},$$

we can apply Lemma E.1 and obtain for $\lambda_\star(0) = 0$,

$$d = \text{Tr}(\boldsymbol{\Sigma}^2 (\boldsymbol{\Sigma} + \lambda_\star(0) \mathbf{I})^{-2}) \geq \text{Tr}(\boldsymbol{\Sigma}^2 (\boldsymbol{\Sigma} + \lambda_\star(\lambda) \mathbf{I})^{-2}) \geq \left(1 - \frac{2\varepsilon}{\mathbf{C}_\Sigma \sigma_d}\right) \cdot d.$$

We then have

$$\mathbf{V}_n(0) \geq \mathbf{V}_n(\lambda) = \frac{n - \text{Tr}(\boldsymbol{\Sigma}^2 (\boldsymbol{\Sigma} + \lambda_\star(0) \mathbf{I})^{-2})}{n - \text{Tr}(\boldsymbol{\Sigma}^2 (\boldsymbol{\Sigma} + \lambda_\star(\lambda) \mathbf{I})^{-2})} \cdot \frac{\text{Tr}(\boldsymbol{\Sigma}^2 (\boldsymbol{\Sigma} + \lambda_\star(\lambda) \mathbf{I})^{-2})}{\text{Tr}(\boldsymbol{\Sigma}^2 (\boldsymbol{\Sigma} + \lambda_\star(0) \mathbf{I})^{-2})} \cdot \mathbf{V}_n(0)$$

$$\begin{aligned}
&= \frac{n-d}{n - \text{Tr}(\boldsymbol{\Sigma}^2(\boldsymbol{\Sigma} + \lambda_*(\lambda)\mathbf{I})^{-2})} \cdot \frac{\text{Tr}(\boldsymbol{\Sigma}^2(\boldsymbol{\Sigma} + \lambda_*(\lambda)\mathbf{I})^{-2})}{d} \cdot \mathbf{V}_n(0) \\
&\geq \frac{\mathbf{C}_\Sigma n}{\mathbf{C}_\Sigma n + \frac{2\varepsilon}{\mathbf{C}_\Sigma \sigma_d} \cdot d} \cdot \left(1 - \frac{2\varepsilon}{\mathbf{C}_\Sigma \sigma_d}\right) \cdot \mathbf{V}_n(0) \\
&\geq \frac{\mathbf{C}_\Sigma^2 \sigma_d}{\mathbf{C}_\Sigma^2 \sigma_d + 2\varepsilon} \cdot \left(1 - \frac{2\varepsilon}{\mathbf{C}_\Sigma \sigma_d}\right) \cdot \mathbf{V}_n(0),
\end{aligned}$$

where in the last line we use $n \geq d$. Again by the elementary inequality $1 - (1-a)(1-b) \leq a+b$ for all $0 \leq a, b \leq 1$,

$$\begin{aligned}
|\mathbf{V}_n(0) - \mathbf{V}_n(\lambda)| &\leq \left\{1 - \left(1 - \frac{2\varepsilon}{\mathbf{C}_\Sigma^2 \sigma_d + 2\varepsilon}\right) \cdot \left(1 - \frac{2\varepsilon}{\mathbf{C}_\Sigma \sigma_d}\right)\right\} \cdot \mathbf{V}_n(0) \\
&\leq \left(\frac{2\varepsilon}{\mathbf{C}_\Sigma^2 \sigma_d + 2\varepsilon} + \frac{2\varepsilon}{\mathbf{C}_\Sigma \sigma_d}\right) \cdot \mathbf{V}_n(0) \leq \frac{4\varepsilon}{\mathbf{C}_\Sigma^2 \sigma_d} \cdot \mathbf{V}_n(0). \tag{91}
\end{aligned}$$

Part III: Variance approximation. Taking $\lambda = \varepsilon n$, we want to invoke Theorem 1 to bound $|\mathcal{V}_\mathbf{X}(\lambda) - \mathbf{V}_n(\lambda)|$. Using Lemma E.1, we know

$$\mathbf{C}_\Sigma \leq 1 - \frac{d}{n} \leq 1 - \frac{1}{n} \text{Tr}(\boldsymbol{\Sigma}(\boldsymbol{\Sigma} + \lambda_*(\lambda)\mathbf{I})^{-1}) \leq 1 - \frac{d}{n} \cdot \frac{\sigma_d}{\sigma_d + \varepsilon/\mathbf{C}_\Sigma} \leq 1 - \frac{\mathbf{C}_\Sigma^2 \sigma_d}{\mathbf{C}_\Sigma^2 \sigma_d + \varepsilon}.$$

Since by assumption $\varepsilon \leq \mathbf{C}_\Sigma^2 \sigma_d/4 \leq \mathbf{C}_\Sigma \sigma_d$, Eq. (21) holds with $\kappa = \mathbf{C}_\Sigma/2$, because

$$\mathbf{C}_\Sigma \leq \frac{\lambda}{n\lambda_*(\lambda)} = 1 - \frac{1}{n} \text{Tr}(\boldsymbol{\Sigma}(\boldsymbol{\Sigma} + \lambda_*(\lambda)\mathbf{I})^{-1}) \leq 1 - \frac{\mathbf{C}_\Sigma}{2}.$$

Thus by Theorem 1, we have for some constant $C' := C'(k, \mathbf{C}_x, D) > 0$,

$$|\mathcal{V}_\mathbf{X}(\lambda) - \mathbf{V}_n(\lambda)| \leq C' \cdot \frac{\chi_n(\varepsilon n)^3 \log^2 n}{n^{1-\frac{1}{k}} \mathbf{C}_\Sigma^{9.5}} \cdot \mathbf{V}_n(\lambda).$$

Combining the above display with Eqs. (90), (91) yields

$$\begin{aligned}
|\mathcal{V}_\mathbf{X}(0) - \mathbf{V}_n(0)| &\leq C' \cdot \frac{\chi_n(\varepsilon n)^3 \log^2 n}{n^{1-\frac{1}{k}} \mathbf{C}_\Sigma^{9.5}} \cdot \mathbf{V}_n(\lambda) + \frac{2\varepsilon}{s_{\min}} \cdot \mathcal{V}_\mathbf{X}(\lambda) + \frac{4\varepsilon}{\mathbf{C}_\Sigma^2 \sigma_d} \cdot \mathbf{V}_n(0) \\
&\leq \left\{ \left(1 + \frac{2\varepsilon}{s_{\min}}\right) \left(1 + C' \cdot \frac{\chi_n(\varepsilon n)^3 \log^2 n}{n^{1-\frac{1}{k}} \mathbf{C}_\Sigma^{9.5}}\right) - 1 \right\} \cdot \mathbf{V}_n(\lambda) + \frac{4\varepsilon}{\mathbf{C}_\Sigma^2 \sigma_d} \cdot \mathbf{V}_n(0) \\
&\leq \left\{ \left(1 + \frac{2\varepsilon}{s_{\min}}\right) \left(1 + C' \cdot \frac{\chi_n(\varepsilon n)^3 \log^2 n}{n^{1-\frac{1}{k}} \mathbf{C}_\Sigma^{9.5}}\right) \left(1 + \frac{4\varepsilon}{\mathbf{C}_\Sigma^2 \sigma_d}\right) - 1 \right\} \cdot \mathbf{V}_n(0) \\
&\leq \left(\exp\left(\frac{2\varepsilon}{s_{\min}} + \frac{4\varepsilon}{\mathbf{C}_\Sigma^2 \sigma_d} + C' \cdot \frac{\chi_n(\varepsilon n)^3 \log^2 n}{n^{1-\frac{1}{k}} \mathbf{C}_\Sigma^{9.5}}\right) - 1 \right) \cdot \mathbf{V}_n(0).
\end{aligned}$$

Since $\varepsilon \leq s_{\min}/2$ and $\varepsilon \leq \mathbf{C}_\Sigma^2 \sigma_d/4$, if we additionally assume

$$\frac{\chi_n(\varepsilon n)^3 \log^2 n}{n^{1-\frac{1}{k}} \mathbf{C}_\Sigma^{9.5}} \leq \frac{1}{C'},$$

we can then conclude that

$$\exp\left(\frac{2\varepsilon}{s_{\min}} + \frac{4\varepsilon}{\mathbf{C}_\Sigma^2 \sigma_d} + C' \cdot \frac{\chi_n(\varepsilon n)^3 \log^2 n}{n^{1-\frac{1}{k}} \mathbf{C}_\Sigma^{9.5}}\right) - 1 = \mathcal{O}\left(\frac{\varepsilon}{s_{\min}} + \frac{\varepsilon}{\mathbf{C}_\Sigma^2 \sigma_d} + C' \cdot \frac{\chi_n(\varepsilon n)^3 \log^2 n}{n^{1-\frac{1}{k}} \mathbf{C}_\Sigma^{9.5}}\right),$$

and the proof is complete with $\mathbf{C}_2 = 1/C'$.

F Proofs for proportional regime

F.1 Proof of Proposition 4.1

To apply Theorem 1, we first provide upper bounds for $\mathbf{d}_\Sigma(n)$ and κ implying that Assumptions 1 and Eq. (21) hold. Throughout we use the shorthand $\lambda_p = \lambda/n \in [1/M, M]$.

Lemma F.1. *Under Assumption 2 and $\lambda = n\lambda_p$, Assumptions 1 and (21) hold for*

$$\begin{aligned}\mathbf{d}_\Sigma(n) &= \mathcal{O}_M(n), \\ \kappa &= \Omega_M(1).\end{aligned}$$

For such \mathbf{d}_Σ and κ , $\chi_n(\lambda) = \mathcal{O}_{\lambda_p, M}(\log^2 n)$.

Proof. By Assumption 2 we know $d \leq Mn$ and therefore for any $1 \leq k \leq \min\{n, d\}$,

$$\sum_{l=k}^d \sigma_l \leq d\sigma_k \leq Mn\sigma_k =: \mathbf{d}_\Sigma \sigma_k.$$

Using $\lambda = n\lambda_p$ into Eq. (5), we have

$$1 - \frac{\lambda_p}{\lambda_\star} = \frac{1}{n} \text{Tr}(\Sigma(\Sigma + \lambda_\star \mathbf{I})^{-1}).$$

which implies

$$1 - \frac{\lambda_p}{\lambda_\star} \leq \frac{d}{n} \cdot \frac{1}{1 + \lambda_\star} \leq \frac{M}{\lambda_\star}.$$

This implies $\lambda_\star \leq \lambda_p + M$ and therefore

$$1 - \frac{\lambda}{n\lambda_\star} = 1 - \frac{\lambda_p}{\lambda_\star} \leq 1 - \frac{\lambda_p}{\lambda_p + M}.$$

On the other hand,

$$1 - \frac{\lambda}{n\lambda_\star} = 1 - \frac{\lambda_p}{\lambda_\star} \geq \frac{d}{n} \cdot \frac{\sigma_d}{\sigma_d + \lambda_\star} \geq \frac{1}{M + M^2\lambda_\star} \geq \frac{1}{M + M^2\lambda_p + M^3},$$

and therefore we have

$$\kappa := \min \left\{ \frac{\lambda_p}{\lambda_p + M}, \frac{1}{M + M^2\lambda_p + M^3} \right\} = \Omega_M(1).$$

Finally, we can bound $\chi_n(\lambda)$ as $\mathbf{d}_\Sigma = \mathcal{O}_M(n)$, and thus

$$\chi_n(\lambda) = 1 + \frac{\sigma_{\lfloor m \rfloor} \mathbf{d}_\Sigma \log^2(\mathbf{d}_\Sigma)}{n\lambda_p} = \mathcal{O}_M \left(1 + \frac{\log^2 n}{\lambda_p} \right) = \mathcal{O}_M(\log^2 n).$$

□

For any unit vector $\mathbf{u} \in \mathbb{R}^d$, since

$$\begin{aligned}
\mathcal{F}_0(\lambda, \mu_\star(\lambda, 0); \mathbf{u}\mathbf{u}^\top) &= \lambda \text{Tr} \left(\mathbf{u}\mathbf{u}^\top \boldsymbol{\Sigma} (\lambda \mathbf{I} + \mu_\star(\lambda, 0) \boldsymbol{\Sigma})^{-1} \right) \\
&\geq \frac{\lambda}{\lambda M + \mu_\star(\lambda, 0)} \text{Tr} \left(\mathbf{u}\mathbf{u}^\top \right) = \frac{\lambda}{d(\lambda M + \mu_\star(\lambda, 0))} \text{Tr}(\mathbf{I}) \\
&\geq \frac{\lambda + \mu_\star(\lambda, 0)}{d(\lambda M + \mu_\star(\lambda, 0))} \cdot \lambda \text{Tr} \left(\boldsymbol{\Sigma} (\lambda \mathbf{I} + \mu_\star(\lambda, 0) \boldsymbol{\Sigma})^{-1} \right) \\
&\geq \frac{1}{dM} \mathcal{F}_0(\lambda, \mu_\star(\lambda, 0); \mathbf{I}) \geq \frac{1}{nM^2} \mathcal{F}_0(\lambda, \mu_\star(\lambda, 0); \mathbf{I}), \tag{92}
\end{aligned}$$

we have $\rho(\lambda) = \Omega_M(n^{-1})$. Together with Lemma F.1, since $n = \Omega_{M, \mathbf{C}_x, D}(1)$, the following conditions in Theorem 1 hold

$$\chi_n(\lambda)^3 \log^2 n \leq Cn\kappa^{4.5} \min \left\{ 1, \sqrt{\rho(\lambda)} \right\}, \quad n^{-2D+1} = \mathcal{O} \left(\sqrt{\frac{\kappa^3 \log^2 n}{\max \{n, \lambda\}}} \right).$$

Additionally, $\lambda k n^{-\frac{1}{k}} \leq n\kappa/2$ is equivalent to $\lambda_p k n^{-\frac{1}{k}} \leq \kappa/2$, which holds for $n = \Omega_{k, M}(1)$. Finally, by using $\lambda_\star(\lambda) \leq \lambda_p + M = \mathcal{O}_M(1)$, as shown above, we can conclude from Theorem 1 and Lemma F.1 that, for $n = \Omega_{k, M, \mathbf{C}_x, D}(1)$, with probability $1 - \mathcal{O}_k(n^{-D+1})$,

$$\begin{aligned}
|\mathcal{V}_{\mathbf{X}}(\lambda) - \mathbf{V}_n(\lambda)| &= \mathcal{O}_{k, \mathbf{C}_x, D} \left(\frac{\chi_n(\lambda)^3 \log^2 n}{n^{1-\frac{1}{k}} \kappa^{9.5}} \right) \cdot \mathbf{V}_n(\lambda) = \mathcal{O}_{k, M, \mathbf{C}_x, D} \left(\frac{\log^8 n}{n^{1-\frac{1}{k}}} \right) \cdot \mathbf{V}_n(\lambda), \\
|\mathcal{B}_{\mathbf{X}}(\lambda) - \mathbf{B}_n(\lambda)| &= \mathcal{O}_{k, \mathbf{C}_x, D} \left(\frac{\lambda_\star(\lambda)^{k+1}}{n\kappa^3} + \frac{\chi_n(\lambda)^3 \log^2 n}{\sqrt{\rho(\lambda)} n^{1-\frac{1}{k}} \kappa^{8.5}} \right) \cdot \mathbf{B}_n(\lambda) = \mathcal{O}_{k, M, \mathbf{C}_x, D} \left(\frac{\log^8 n}{n^{\frac{1}{2}-\frac{1}{k}}} \right) \cdot \mathbf{B}_n(\lambda).
\end{aligned}$$

The proof is complete.

F.2 Proof of Proposition 4.2

Overparameterized regime. When $d/n \geq 1 + M^{-1}$, by

$$n = \text{Tr} \left(\boldsymbol{\Sigma} (\boldsymbol{\Sigma} + \lambda_\star(0) \mathbf{I})^{-1} \right) \geq \frac{d}{1 + M\lambda_\star(0)},$$

we can deduce that $\lambda_\star(0) \geq M^{-1} \cdot (d/n - 1) \geq M^{-2}$. Hence,

$$n - \text{Tr} \left(\boldsymbol{\Sigma}^2 (\boldsymbol{\Sigma} + \lambda_\star(0) \mathbf{I})^{-2} \right) \geq n - \frac{1}{1 + \lambda_\star(0)} \cdot \text{Tr} \left(\boldsymbol{\Sigma} (\boldsymbol{\Sigma} + \lambda_\star(0) \mathbf{I})^{-1} \right) \geq \frac{\lambda_\star(0)}{1 + \lambda_\star(0)} \cdot n \geq \frac{1}{M^2 + 1} \cdot n,$$

and therefore, in Theorem 2 we can take $\mathbf{C}_{\boldsymbol{\Sigma}} \geq 1/(M^2 + 1) = \Theta_M(1)$. By Eq. (92) we know $\rho(0) = \Omega_M(n^{-1})$. By [BY08, RV09], we know when $n = \Omega_{\mathbf{C}_x, M, D}(1)$, with probability $1 - \mathcal{O}(n^{-D+1})$ we have $s_{\min} = \Omega_{M, \mathbf{C}_x, D}(1)$. Substituting $\lambda_\star(0) = \Omega_M(1)$ and $\mathbf{d}_{\boldsymbol{\Sigma}}(n) = \mathcal{O}_M(n)$ (c.f. Lemma F.1) into $\chi'_n(\kappa)$, we get for $\kappa = \mathcal{O}(1)$,

$$\chi'_n(\kappa) = \mathcal{O}_M \left(\frac{\log^2 n}{\kappa} \right).$$

Thus, by taking $\kappa = n^{-1/14}$, the conditions below hold for $n = \Omega_{k, M, \mathbf{C}_x, D}(1)$ given $k \geq 15$,

$$\kappa \leq \min \left\{ s_{\min}/(8\lambda_\star(0)), \mathbf{C}_{\boldsymbol{\Sigma}}^2/8 \right\}, \quad n^{-2D+1} = \mathcal{O} \left(\sqrt{\frac{\kappa^3 \log^2 n}{\max \{n, \lambda\}}} \right),$$

$$\chi'_n(\kappa)^3 \log^2 n \leq C_2 n^{1-\frac{1}{k}} \kappa^{9.5},$$

and by taking $\kappa = n^{-1/28}$, the following additional conditions hold when $n = \Omega_{k,M,C_{\mathbf{x}},D}(1)$ given $k \geq 29$,

$$\chi'_n(\kappa)^3 \log^2 n \leq C_1 n \kappa^{4.5} \min \left\{ 1, \sqrt{\rho(0)} \right\}, \quad \frac{\lambda_*(0)^{k+1}}{n \kappa^3} + \frac{\chi'_n(\kappa)^3 \log^2 n}{\sqrt{\rho(0)} n^{1-\frac{1}{k}} \kappa^{8.5}} \leq C_3.$$

We can then invoke Theorem 2 by taking $\kappa = n^{-1/14}$ for variance approximation and $n^{-1/28}$ for bias approximation. Therefore, we can conclude that for $k \geq 29$,

$$\begin{aligned} |\mathcal{V}_{\mathbf{X}}(0) - \mathbf{V}_n(0)| &= \mathcal{O}_{k,M,C_{\mathbf{x}},D} \left(n^{-1/14} + \frac{\log^8 n}{n^{1.5/14-\frac{1}{k}}} \right) \cdot \mathbf{V}_n(0), \\ |\mathcal{B}_{\mathbf{X}}(0) - \mathbf{B}_n(0)| &= \mathcal{O}_{k,M,C_{\mathbf{x}},D} \left(n^{-1/28} + n^{-25/28} + \frac{\log^8 n}{n^{2.5/28-\frac{1}{k}}} \right) \cdot \mathbf{B}_n(0) + \mathcal{O} \left(\|\boldsymbol{\beta}\|^2 n^{-1/28} \right). \end{aligned}$$

Use again $\lambda_*(0) = \Omega_M(1)$ and $\mathbf{C}_{\boldsymbol{\Sigma}} = \Omega_M(1)$, we know

$$\mathbf{B}_n(0) = \frac{\boldsymbol{\beta}^\top (\boldsymbol{\Sigma}/\lambda_*(0) + \mathbf{I})^{-2} \boldsymbol{\Sigma} \boldsymbol{\beta}}{1 - n^{-1} \text{Tr}(\boldsymbol{\Sigma}^2 (\boldsymbol{\Sigma} + \lambda_*(0) \mathbf{I})^{-2})} = \Omega_M(\|\boldsymbol{\beta}\|^2).$$

We conclude the proof by fixing $k \geq 57$, and thus

$$\begin{aligned} |\mathcal{V}_{\mathbf{X}}(0) - \mathbf{V}_n(0)| &= \mathcal{O}_{M,C_{\mathbf{x}},D} \left(n^{-1/14} \right) \cdot \mathbf{V}_n(0), \\ |\mathcal{B}_{\mathbf{X}}(0) - \mathbf{B}_n(0)| &= \mathcal{O}_{M,C_{\mathbf{x}},D} \left(n^{-1/28} \right) \cdot \mathbf{B}_n(0). \end{aligned}$$

Underparameterized regime. Suppose $M^{-1} \leq d/n \leq 1 - M^{-1}$, we can invoke Theorem 3 with $\mathbf{C}_{\boldsymbol{\Sigma}} = M^{-1}$. By [BY08] we have $s_{\min} = \Omega_{M,C_{\mathbf{x}},D}(1)$. Also as we can take $\mathbf{d}_{\boldsymbol{\Sigma}}(n) = n$ in this case, we have

$$\chi_n(\varepsilon n) \leq 1 + \frac{n \log^2 n}{\varepsilon n} = \mathcal{O} \left(\frac{\log^2 n}{\kappa} \right),$$

and therefore the conditions below hold for $n = \Omega_{k,M,C_{\mathbf{x}},D}(1)$ by taking $\varepsilon = n^{-1/4}$ when $k \geq 5$,

$$\varepsilon \leq \min \left\{ s_{\min}/2, \mathbf{C}_{\boldsymbol{\Sigma}}^2 \sigma_d/4 \right\}, \quad \chi_n(\varepsilon n)^3 \log^2 n \leq C_1 n \mathbf{C}_{\boldsymbol{\Sigma}}^{4.5}, \quad n^{-2D+1} = \mathcal{O} \left(\sqrt{\frac{\mathbf{C}_{\boldsymbol{\Sigma}}^3 \log^2 n}{\max \{n, \lambda\}}} \right),$$

$$\chi_n(\varepsilon n)^3 \log^2 n \leq C_2 n^{1-\frac{1}{k}} \mathbf{C}_{\boldsymbol{\Sigma}}^{9.5}.$$

We thus have

$$|\mathcal{V}_{\mathbf{X}}(0) - \mathbf{V}_n(0)| = \mathcal{O}_{k,C_{\mathbf{x}},D} \left(n^{-1/4} \cdot \left(\frac{1}{s_{\min}} + \frac{1}{\mathbf{C}_{\boldsymbol{\Sigma}}^2 \sigma_d} \right) + \frac{\log^8 n}{n^{\frac{1}{4}-\frac{1}{k}} \mathbf{C}_{\boldsymbol{\Sigma}}^{9.5}} \right) \cdot \mathbf{V}_n(0).$$

By fixing $k > 20$, we know for all $n = \Omega_{M,C_{\mathbf{x}},D}(1)$,

$$|\mathcal{V}_{\mathbf{X}}(0) - \mathbf{V}_n(0)| = \mathcal{O}_{M,C_{\mathbf{x}},D} \left(\frac{1}{n^{\frac{1}{5}}} \right) \cdot \mathbf{V}_n(0).$$

G Proofs for bounded varying spectrum regime

G.1 Proof of Proposition 4.3

Throughout this proof, we will use the shorthand $\lambda_{\text{bv}} := \lambda/(n\lambda_*(0)) \in [1/M, M]$. We begin by lower bounding the constant κ of Eq. (21). Since

$$\begin{aligned} \text{Tr}(\mathbf{\Sigma}(\mathbf{\Sigma} + \lambda_*(0)\mathbf{I})^{-1}) &= n, \\ \text{Tr}(\mathbf{\Sigma}(\mathbf{\Sigma} + \sigma_{2n}\mathbf{I})^{-1}) &\geq \sum_{i=1}^{2n} \frac{\sigma_i}{\sigma_i + \sigma_{2n}} \geq n, \end{aligned} \tag{93}$$

we know that $\lambda_*(0) \geq \sigma_{2n}$ and therefore $\psi(\delta)\lambda_*(0) \geq \psi(\delta)\sigma_{2n} \geq \sigma_{\lfloor 2\delta n \rfloor}$ for any $\delta \in (0, 1]$. We then have

$$\begin{aligned} \text{Tr}(\mathbf{\Sigma}(\mathbf{\Sigma} + \psi(\delta)\lambda_*(0)\mathbf{I})^{-1}) &\leq 2\delta n + \sum_{i=\lfloor 2\delta n \rfloor+1}^{\infty} \frac{\sigma_i}{\sigma_i + \psi(\delta)\lambda_*(0)} \\ &\leq 2\delta n + \frac{\sigma_{\lfloor 2\delta n \rfloor} + \lambda_*(0)}{\sigma_{\lfloor 2\delta n \rfloor} + \psi(\delta)\lambda_*(0)} \cdot \text{Tr}(\mathbf{\Sigma}(\mathbf{\Sigma} + \lambda_*(0)\mathbf{I})^{-1}) \\ &\leq \left(\frac{1}{2} + 2\delta + \frac{1}{\psi(\delta)} \right) \cdot n, \end{aligned}$$

where in the last inequality we use $\psi(\delta)\lambda_*(0) \geq \sigma_{\lfloor 2\delta n \rfloor}$. Further

$$n - \frac{\lambda}{\psi(\delta)\lambda_*(0)} = \left(1 - \frac{\lambda_{\text{bv}}}{\psi(\delta)} \right) \cdot n,$$

and therefore, using the previous inequality, we conclude the following. If $\delta > 0$ is such that

$$1 - \frac{\lambda_{\text{bv}}}{\psi(\delta)} \geq \frac{1}{2} + 2\delta + \frac{1}{\psi(\delta)},$$

then $\lambda_*(\lambda) \leq \psi(\delta)\lambda_*(0)$. Let $\delta_0 = \delta_0(M, \psi)$ be defined follows

$$\delta_0 := \sup \left\{ \delta \in (0, 1/2) : 2\delta + \frac{1+M}{\psi(\delta)} \leq \frac{1}{2} \right\}.$$

Then $\lambda_*(0) \leq \lambda_*(\lambda) \leq \psi(\delta_0)\lambda_*(0)$. Hence

$$\frac{\lambda}{n\lambda_*(\lambda)} = \frac{\lambda_*(0)\lambda_{\text{bv}}}{\lambda_*(\lambda)} \geq \frac{\lambda_{\text{bv}}}{\psi(\delta_0)},$$

and

$$1 - \frac{\lambda}{n\lambda_*(\lambda)} = \frac{1}{n} \text{Tr}(\mathbf{\Sigma}(\mathbf{\Sigma} + \lambda_*(\lambda)\mathbf{I})^{-1}) \geq \frac{1}{n} \text{Tr}(\mathbf{\Sigma}(\mathbf{\Sigma} + \psi(\delta_0)\lambda_*(0)\mathbf{I})^{-1}) \geq \frac{1}{\psi(\delta_0)}.$$

Putting together the above displays, we conclude Eq (21) holds for $\kappa \geq \min\{\lambda_{\text{bv}}, 1\}/\psi(\delta_0) = \Omega_{M,\psi}(1)$ when $n = \Omega_{M,\psi}(1)$ (recall that we need $2\delta_0 n \geq 1$).

To verify the conditions of Theorem 1, we first assume $\mathbf{d}_{\mathbf{\Sigma}} = \mathcal{O}(n^{1+\gamma})$ for $\gamma \in [0, 1/3)$,

$$\chi_n(n\lambda_*(0)\lambda_{\text{bv}}) = 1 + \frac{\sigma_{\lfloor \eta n \rfloor} \mathbf{d}_{\mathbf{\Sigma}} \log^2(\mathbf{d}_{\mathbf{\Sigma}})}{n\lambda_*(0)\lambda_{\text{bv}}} \leq 1 + \frac{\sigma_{\lfloor \eta n \rfloor} \mathbf{d}_{\mathbf{\Sigma}} \log^2(\mathbf{d}_{\mathbf{\Sigma}})}{n\sigma_{2n}\lambda_{\text{bv}}} \leq 1 + \frac{\psi(\eta/4)\mathbf{d}_{\mathbf{\Sigma}} \log^2(\mathbf{d}_{\mathbf{\Sigma}})}{n\lambda_{\text{bv}}}$$

$$= \mathcal{O}_{M,\psi,\mathbf{C}_x} \left(\frac{\mathbf{d}_\Sigma \log^2 n}{n} \right) = \mathcal{O}_{M,\psi,\mathbf{C}_x} (n^\gamma \log^2 n),$$

and with $\kappa = \Omega_{M,\psi}(1)$, the conditions

$$\chi_n(n\lambda_\star(0)\lambda_{\mathbf{b}_v})^3 \log^2 n \leq Cn\kappa^{4.5}, \quad n^{-2D+1} = \mathcal{O} \left(\sqrt{\frac{\kappa^3 \log^2 n}{\max\{n, \lambda\}}} \right),$$

hold if $n = \Omega_{M,\psi,\gamma,\mathbf{C}_x,D}(1)$. We then can apply Theorem 1 to approximate the variance. Given any positive integer k , if $n = \Omega_{k,M,\psi,\gamma,\mathbf{C}_x,D}(1)$, it holds with probability $1 - \mathcal{O}_k(n^{-D+1})$ that

$$|\mathcal{V}_X(\lambda) - \mathbf{V}_n(\lambda)| = \mathcal{O}_{k,M,\psi,\mathbf{C}_x,D} \left(\frac{(\mathbf{d}_\Sigma/n)^3 \log^8 n}{n^{1-\frac{1}{k}}} \right) \cdot \mathbf{V}_n(\lambda).$$

If additionally $\mathbf{d}_\Sigma = \mathcal{O}_{M,\psi,\mathbf{C}_x}(n^{1+\gamma}\lambda^{1/6})$, we have

$$\chi_n(\lambda)^3 \log^2 n \leq Cn\kappa^{4.5} \sqrt{\rho(\lambda)},$$

when $n = \Omega_{M,\psi,\gamma,\mathbf{C}_x,D}(1)$. The condition $\lambda k n^{-\frac{1}{k}} \leq n\kappa/2$ is equivalent to $\lambda_\star(0)\lambda_{\mathbf{b}_v} k n^{-1/k} \leq \kappa/2$, which holds when $n = \Omega_{k,M,\psi}(1)$ since we have assumed $\lambda_\star(0) = \mathcal{O}(1)$. Therefore, we can appeal to the bias approximation result in Theorem 1, yielding

$$|\mathcal{B}_X(\lambda) - \mathbf{B}_n(\lambda)| = \mathcal{O}_{k,M,\psi,\mathbf{C}_x,D} \left(\frac{(\mathbf{d}_\Sigma/n)^3 \log^8 n}{\sqrt{\rho(\lambda)} n^{1-\frac{1}{k}}} \right) \cdot \mathbf{B}_n(\lambda).$$

G.2 Proof of Proposition 4.4

We provide the following bounds for the quantities in Theorem 2.

Lemma G.1. *Under the same Assumptions of Theorem 2, we can take*

$$\mathbf{C}_\Sigma = \Omega_\psi(1),$$

when $n = \Omega(1)$. For $\kappa = \mathcal{O}(1)$, we have

$$\chi'_n(\kappa) = \mathcal{O}_\psi \left(\frac{\log^{2+\mathcal{O}(1)} n}{\kappa} \right).$$

In addition, $s_{\min} = \Omega_{\psi,\mathbf{C}_x}(\sigma_n)$ with probability $1 - \mathcal{O}(n^{-D+1})$.

Proof. Since

$$\begin{aligned} n - \text{Tr} \left(\Sigma^2 (\Sigma + \lambda_\star(0)\mathbf{I})^{-2} \right) &\geq n - \sum_{i=1}^{\lfloor n/2 \rfloor} \frac{\sigma_i}{\sigma_i + \lambda_\star(0)} - \frac{\sigma_{\lfloor n/2 \rfloor}}{\sigma_{\lfloor n/2 \rfloor} + \lambda_\star(0)} \sum_{i=\lfloor n/2 \rfloor + 1}^{\infty} \frac{\sigma_i}{\sigma_i + \lambda_\star(0)} \\ &= \frac{\lambda_\star(0)}{\sigma_{\lfloor n/2 \rfloor} + \lambda_\star(0)} \sum_{i=\lfloor n/2 \rfloor + 1}^{\infty} \frac{\sigma_i}{\sigma_i + \lambda_\star(0)}, \end{aligned}$$

where in the last line we use $\text{Tr}(\mathbf{\Sigma}(\mathbf{\Sigma} + \lambda_\star(0)\mathbf{I})^{-1}) = n$. Since

$$\sum_{i=\lfloor n/2 \rfloor + 1}^{\infty} \frac{\sigma_i}{\sigma_i + \lambda_\star(0)} = n - \sum_{i=1}^{\lfloor n/2 \rfloor} \frac{\sigma_i}{\sigma_i + \lambda_\star(0)} \geq n - \frac{n}{2} = \frac{n}{2},$$

and $\lambda_\star(0) \geq \sigma_{2n}$ from Eq. (93), we know

$$n - \text{Tr}(\mathbf{\Sigma}^2(\mathbf{\Sigma} + \lambda_\star(0)\mathbf{I})^{-2}) \geq \frac{\sigma_{2n}}{\sigma_{\lfloor n/2 \rfloor} + \sigma_{2n}} \cdot \frac{n}{2} \geq \frac{1}{2\psi(1/4) + 2} \cdot n.$$

We can hence take $\mathbf{C}_\Sigma := (2\psi(1/4) + 2)^{-1} = \Omega_\psi(1)$.

Substituting $\lambda_\star(0) \geq \sigma_{2n}$ and $\mathbf{d}_\Sigma = \mathcal{O}(n^{1+\gamma})$ for some $1 \leq \gamma < 1/2$ into Eq. (24), we have for $\kappa = \mathcal{O}(1)$,

$$\chi'_n(\kappa) = \mathcal{O}\left(\frac{\sigma_{\lfloor \eta n \rfloor} \mathbf{d}_\Sigma \log^2 n}{\kappa n \sigma_{2n}}\right) = \mathcal{O}_\psi\left(\frac{\log^{2+\mathcal{O}(1)} n}{\kappa}\right).$$

Finally, to get a lower bound for s_{\min} , we use Cauchy interlacing theorem which implies

$$\begin{aligned} s_{\min} &\geq \lambda_n\left(\frac{\mathbf{X}^\top \mathbf{X}}{n}\right) = \lambda_n\left(\frac{\mathbf{\Sigma}^{\frac{1}{2}} \mathbf{Z}^\top \mathbf{Z} \mathbf{\Sigma}^{\frac{1}{2}}}{n}\right) \geq \lambda_n\left(\frac{\mathbf{P}_k \mathbf{\Sigma}^{\frac{1}{2}} \mathbf{Z}^\top \mathbf{Z} \mathbf{\Sigma}^{\frac{1}{2}} \mathbf{P}_k}{n}\right) \\ &= \lambda_n\left(\frac{\mathbf{Z} \mathbf{P}_k \mathbf{\Sigma} \mathbf{P}_k \mathbf{Z}^\top}{n}\right). \end{aligned}$$

where \mathbf{P}_k is the projection to the space spanned by the top k eigenvectors. Let $k \geq n$, we further have

$$s_{\min} \geq \sigma_k \cdot \lambda_n\left(\frac{\mathbf{Z} \mathbf{V}_k \mathbf{V}_k^\top \mathbf{Z}^\top}{n}\right),$$

where $\mathbf{P}_k = \mathbf{V}_k \mathbf{V}_k^\top$ and $\mathbf{V}_k = [\mathbf{v}_1 \ \cdots \ \mathbf{v}_k] \in \mathbb{R}^{d \times k}$ with \mathbf{v}_i being the i -th eigenvector of $\mathbf{\Sigma}$. Since $\mathbf{Z} \mathbf{V}_k$ is a $n \times k$ random matrix with i.i.d. isotropic and sub-Gaussian rows. When $k \geq n$, by [Ver12, Thm. 5.58, generalized version in Sec. 5.7], we have

$$\lambda_n\left(\frac{\mathbf{Z} \mathbf{V}_k \mathbf{V}_k^\top \mathbf{Z}^\top}{n}\right) \geq \left((1 - \zeta) \sqrt{\frac{k}{n}} - \mathcal{O}_{\mathbf{C}_x}(1) - \mathcal{O}_{\mathbf{C}_x, D}\left(\sqrt{\frac{\log n}{n}}\right) \right)^2,$$

with probability at least $1 - \mathcal{O}(n^{-D+1})$, where ζ is the random variable

$$\zeta := \max_{1 \leq i \leq n} \left| \frac{\|\mathbf{V}_k^\top \mathbf{z}_i\|^2}{k} - 1 \right|.$$

By Hanson-Wright in Lemma 2.1 and similar to the argument in Eq. (61), we have

$$\mathbb{P}\left(\left| \frac{\|\mathbf{V}_k^\top \mathbf{z}_i\|^2}{k} - 1 \right| \geq t\right) = 2 \exp(-\Omega_{\mathbf{C}_x}(k \cdot \min\{t^2, t\})).$$

Given the above sharp concentration of ζ , we can therefore conclude by taking $k = \lfloor \mathbf{C}(\mathbf{C}_x)n \rfloor$ for some $\mathbf{C} > 0$, and $n = \Omega_{\mathbf{C}_x, D}(1)$, we have with probability $1 - \mathcal{O}(n^{-D+1})$ that

$$\lambda_n\left(\frac{\mathbf{Z} \mathbf{V}_k \mathbf{V}_k^\top \mathbf{Z}^\top}{n}\right) \geq 1,$$

and therefore $s_{\min} \geq \sigma_k \geq \sigma_n / \psi(\mathbf{C})$ by Assumption 3. \square

By Lemma G.1 and the assumption $\lambda_\star(0)/\sigma_n = \mathcal{O}(\log^{\mathcal{O}(1)} n)$, we know by taking $\kappa = n^{-1/14}$, the conditions below hold for $n = \Omega_{k,\psi,\mathbf{C}_\mathbf{x},D}(1)$ whenever $k \geq 15$,

$$\kappa \leq \min \left\{ s_{\min}/(8\lambda_\star(0)), \mathbf{C}_\Sigma^2/8 \right\}, \quad \chi'_n(\kappa)^3 \log^2 n \leq \mathbf{C}_1 n \kappa^{4.5}, \quad n^{-2D+1} = \mathcal{O} \left(\sqrt{\frac{\kappa^3 \log^2 n}{\max \{n, \lambda\}}} \right),$$

$$\chi'_n(\kappa)^3 \log^2 n \leq \mathbf{C}_2 n^{1-\frac{1}{k}} \kappa^{9.5}.$$

Therefore by the variance approximation in Theorem 2, it holds

$$|\mathcal{V}_\mathbf{X}(0) - \mathbf{V}_n(0)| = \mathcal{O}_{k,\psi,\mathbf{C}_\mathbf{x},D} \left(n^{-1/14} \log^{\mathcal{O}(1)}(n) + \frac{\log^{8+\mathcal{O}(1)} n}{n^{1.5/14-\frac{1}{k}}} \right) \cdot \mathbf{V}_n(0).$$

Fixing $k \geq 29$, we have with probability $1 - \mathcal{O}(n^{-D+1})$, we know as $n = \Omega_{\psi,\mathbf{C}_\mathbf{x},D}(1)$,

$$|\mathcal{V}_\mathbf{X}(0) - \mathbf{V}_n(0)| = \mathcal{O}_{\psi,\mathbf{C}_\mathbf{x},D} \left(n^{-1/15} \right) \cdot \mathbf{V}_n(0).$$

For the bias approximation with the assumption $\rho(0) = \Omega(n^{-2+\gamma})$, the following additional conditions hold by taking $\kappa = n^{-\gamma/28}$ when $n = \Omega_{k,\psi,\mathbf{C}_\mathbf{x},D}(1)$ given $k \geq 29/\gamma$,

$$\chi'_n(\kappa)^3 \log^2 n \leq \mathbf{C}_1 n \kappa^{4.5} \min \left\{ 1, \sqrt{\rho(0)} \right\}, \quad \frac{\lambda_\star(0)^{k+1}}{n \kappa^3} + \frac{\chi'_n(\kappa)^3 \log^2 n}{\sqrt{\rho(0)} n^{1-\frac{1}{k}} \kappa^{8.5}} \leq \mathbf{C}_3.$$

In verifying the second condition above, we use $\lambda_\star(0) = \mathcal{O}(\sigma_n \log^{\mathcal{O}(1)} n) = \mathcal{O}(\sigma_n \log^{\mathcal{O}(1)} n)$. By Lemma G.1, we also have

$$\kappa \lambda_\star(0) \chi'_n(\kappa) = \mathcal{O}_\psi(\sigma_n \log^{2+\mathcal{O}(1)}).$$

We can then write out the bias approximation result applying Theorem 2

$$|\mathcal{B}_\mathbf{X}(0) - \mathbf{B}_n(0)| = \mathcal{O}_{k,\psi,\mathbf{C}_\mathbf{x},D} \left(n^{-\gamma/28} + n^{-25\gamma/28} \log^{\mathcal{O}(k)} n + \frac{\log^{8+\mathcal{O}(1)} n}{n^{2.5\gamma/28-\frac{1}{k}}} \right) \cdot \mathbf{B}_n(0)$$

$$+ \mathcal{O}_{\psi,\mathbf{C}_\mathbf{x},D} \left(\sigma_n^2 \log^{4+\mathcal{O}(1)} \|\boldsymbol{\theta}_{\leq n}\|^2 + \sigma_n \log^{2+\mathcal{O}(1)} \|\boldsymbol{\beta}_{> n}\|^2 \right).$$

Fixing $k \geq 57/\gamma$, we conclude the proof with

$$|\mathcal{B}_\mathbf{X}(0) - \mathbf{B}_n(0)| = \mathcal{O}_{\psi,\mathbf{C}_\mathbf{x},D} \left(n^{-\gamma/29} \right) \cdot \mathbf{B}_n(0) + \log^{\mathcal{O}(1)} \cdot \mathcal{O}_{\psi,\mathbf{C}_\mathbf{x},D} \left(\sigma_n^2 \|\boldsymbol{\theta}_{\leq n}\|^2 + \sigma_n \|\boldsymbol{\beta}_{> n}\|^2 \right).$$

H Proof of Theorem 4

Define the following increasing function in t ,

$$f_n(t; \lambda) = 1 - \frac{\lambda}{n \cdot t \sigma_n} - \frac{1}{n} \text{Tr} \left(\boldsymbol{\Sigma} (\boldsymbol{\Sigma} + t \sigma_n \mathbf{I})^{-1} \right).$$

Case I: regularly varying spectrum when $\alpha > 1$. In the first case, we set $\lambda = \nu n \sigma_n$. For any $t > 0$, we can compute that

$$f_n(t; \lambda) = 1 - \frac{\nu}{t} - \frac{1}{n} \text{Tr} \left(\mathbf{\Sigma} (\mathbf{\Sigma} + t \sigma_n \mathbf{I})^{-1} \right).$$

We will first show $\mathbf{d}_{\mathbf{\Sigma}}(n) = \mathcal{O}_{\mathbf{\Sigma}}(n)$ and $\lambda_{\star} = \Theta_{\nu}(\lambda/n)$, and then we can invoke Proposition 4.3 for variance approximation. For simplicity, we will suppress the dependence on sequences $\{a_i\}$ and $\{b_i\}$ in the big-O and big- Ω notations. For instance, we will just write for all $n = \Omega_{\alpha}(1)$, $|b_n| \leq \alpha$.

We first upper bound $\mathbf{d}_{\mathbf{\Sigma}}$. Note that

$$\sum_{l=k}^d \sigma_l = \sum_{l=k}^{\infty} l^{-\alpha} a_l \exp \left\{ \sum_{j=1}^l b_j/j \right\} = \sigma_k \cdot \sum_{l=k}^{\infty} \left(\frac{l}{k} \right)^{-\alpha} \cdot \frac{a_l}{a_k} \cdot \exp \left\{ \sum_{j=k+1}^l b_j/j \right\}.$$

As a_l converges to a positive limit, we have $a_l/a_k = \mathcal{O}(1)$. For $k = \Omega_{\alpha}(1)$ such that $|b_l| \leq \alpha/2$ for all $l \geq k$, we can further derive that

$$\begin{aligned} \sum_{l=k}^d \sigma_l &\leq \sigma_k \cdot \mathcal{O}(1) \cdot \sum_{l=k}^{\infty} \left(\frac{l}{k} \right)^{-\alpha} \cdot \exp \left\{ \frac{\alpha}{2} \sum_{j=k+1}^l j^{-1} \right\} \\ &\leq \sigma_k \cdot \mathcal{O}(1) \cdot \sum_{l=k}^{\infty} \left(\frac{l}{k} \right)^{-\alpha} \cdot \exp \left\{ \frac{\alpha}{2} \int_k^l \frac{dt}{t} \right\} \\ &= \sigma_k \cdot \mathcal{O}(1) \cdot \sum_{l=k}^{\infty} \left(\frac{l}{k} \right)^{\alpha/2-\alpha} = k \sigma_k \cdot \mathcal{O}(1) \cdot \sum_{l=k}^{\infty} \frac{1}{k} \left(\frac{l}{k} \right)^{-\alpha/2} \\ &\leq k \sigma_k \cdot \mathcal{O}(1) \cdot \left(\frac{1}{k} + \int_1^{\infty} t^{-\alpha/2} dt \right) \\ &= \mathcal{O}_{\alpha}(k \sigma_k). \end{aligned}$$

This implies for all $n = \Omega_{\alpha}(1)$, we can take $\mathbf{d}_{\mathbf{\Sigma}}(n) = \mathcal{O}_{\alpha}(n)$. Next we show $\lambda_{\star} = \Theta_{\nu}(\lambda/n)$. Note that

$$\begin{aligned} &\lim_{n \rightarrow \infty} \frac{1}{n} \text{Tr} \left(\mathbf{\Sigma} (\mathbf{\Sigma} + t \sigma_n \mathbf{I})^{-1} \right) \\ &= \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{l=1}^{\infty} \frac{\sigma_l}{\sigma_l + t \sigma_n} = \lim_{M \rightarrow \infty} \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{M^{-1}n \leq l \leq Mn} \frac{\sigma_l}{\sigma_l + t \sigma_n} \\ &= \lim_{M \rightarrow \infty} \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{M^{-1}n \leq l \leq Mn} \frac{1}{1 + t(l/n)^{\alpha}} = \int_0^{\infty} \frac{1}{1 + tx^{\alpha}} dx = t^{-1/\alpha} \cdot \frac{1}{\alpha} \int_0^{\infty} \frac{u^{1/\alpha-1}}{1+u} du \\ &= t^{-1/\alpha} \cdot \frac{\text{Beta}(1/\alpha, 1-1/\alpha)}{\alpha} = t^{-1/\alpha} \cdot \frac{\Gamma(1/\alpha)\Gamma(1-1/\alpha)}{\alpha\Gamma(1)} \\ &\stackrel{(i)}{=} t^{-1/\alpha} \frac{\pi/\alpha}{\sin(\pi/\alpha)}, \end{aligned} \tag{94}$$

where in (i) we use the reflection formula for Γ function. Recall that we define $\mathbf{c}_{\star} = \mathbf{c}_{\star}(\nu)$ as the unique solution of

$$1 = \nu \mathbf{c}_{\star}^{-1} + \frac{\pi/\alpha}{\sin(\pi/\alpha)} \mathbf{c}_{\star}^{-1/\alpha},$$

it then follows from the above displays that

$$\begin{aligned}\lim_{n \rightarrow \infty} f_n(\mathbf{c}_*; \lambda) &= 1 - \nu \mathbf{c}_*^{-1} - \lim_{n \rightarrow \infty} \frac{1}{n} \text{Tr} \left(\boldsymbol{\Sigma} (\boldsymbol{\Sigma} + \mathbf{c}_* \sigma_n \mathbf{I})^{-1} \right) \\ &= 1 - \nu \mathbf{c}_*^{-1} - \frac{\pi/\alpha}{\sin(\pi/\alpha)} \mathbf{c}_*^{-1/\alpha} = 0.\end{aligned}$$

By the definition of λ_* in (5), we can write $f_n(\lambda_*/\sigma_n; \lambda) = 0$. Combining with the above limit, we can then conclude that

$$\lambda_* = \mathbf{c}_* \sigma_n (1 + o_n(1)).$$

Substituting into Eq. (23), we further have

$$\begin{aligned}\rho(\lambda) &= \frac{\boldsymbol{\theta}^\top \boldsymbol{\Sigma}^{\frac{1}{2}} (\boldsymbol{\Sigma} + \lambda_* \mathbf{I})^{-1} \boldsymbol{\Sigma}^{\frac{1}{2}} \boldsymbol{\theta}}{\|\boldsymbol{\theta}\|^2 \text{Tr}(\boldsymbol{\Sigma} (\boldsymbol{\Sigma} + \lambda_* \mathbf{I})^{-1})} = \frac{\boldsymbol{\beta}^\top (\boldsymbol{\Sigma} + \lambda_* \mathbf{I})^{-1} \boldsymbol{\beta}}{\|\boldsymbol{\beta}\|_{\boldsymbol{\Sigma}^{-1}}^2 (n - \lambda/\lambda_*)} \\ &= \frac{\sum_{l=1}^{\infty} (\sigma_l + \lambda_*)^{-1} \langle \boldsymbol{\beta}, \mathbf{v}_l \rangle^2}{n(1 - \nu \mathbf{c}_*^{-1}) \sum_{l=1}^{\infty} \sigma_l^{-1} \langle \boldsymbol{\beta}, \mathbf{v}_l \rangle^2} \cdot (1 + o_n(1)) \\ &= \frac{\sum_{l=1}^{\infty} \sigma_n (\sigma_l + \lambda_*)^{-1} \langle \boldsymbol{\beta}, \mathbf{v}_l \rangle^2}{n(1 - \nu \mathbf{c}_*^{-1}) \sum_{l=1}^{\infty} \sigma_n \sigma_l^{-1} \langle \boldsymbol{\beta}, \mathbf{v}_l \rangle^2} \cdot (1 + o_n(1)) \\ &= \frac{\sum_{l=1}^{\infty} (l/n)^\alpha (1 + \mathbf{c}_* (l/n)^\alpha)^{-1} \langle \boldsymbol{\beta}, \mathbf{v}_l \rangle^2}{n(1 - \nu \mathbf{c}_*^{-1}) \sum_{l=1}^{\infty} (l/n)^\alpha \langle \boldsymbol{\beta}, \mathbf{v}_l \rangle^2} \cdot (1 + o_n(1)) \\ &= \frac{\int_0^\infty x^\alpha (1 + \mathbf{c}_* x^\alpha)^{-1} dF_{\boldsymbol{\beta}}(x)}{n(1 - \nu \mathbf{c}_*^{-1}) \int_0^\infty x^\alpha dF_{\boldsymbol{\beta}}(x)} \cdot (1 + o_n(1)).\end{aligned}$$

Therefore, under the additional condition for some $0 < \theta \leq 1$ that

$$\int_0^\infty x^\alpha dF_{\boldsymbol{\beta}}(x) = \mathcal{O} \left(n^{1-\theta} \int_0^\infty x^\alpha (1 + \mathbf{c}_* x^\alpha)^{-1} dF_{\boldsymbol{\beta}}(x) \right),$$

we have $\rho(\lambda) = \Omega(n^{-2+\theta})$. By choosing $\gamma = (1 - \theta)/3$, we can invoke Proposition 4.3. Choosing a sufficiently large k yields $\mathcal{V}_{\mathbf{X}}(\lambda) = \mathbf{V}_n(\lambda)(1 + o_n(1))$ and $\mathcal{B}_{\mathbf{X}}(\lambda) = \mathbf{B}_n(\lambda)(1 + o_n(1))$.

In the next step, we derive explicit asymptotic formulas for \mathbf{V}_n and \mathbf{B}_n . Similar to the previous calculations in Eq. (94), we can compute that

$$\begin{aligned}&\lim_{n \rightarrow \infty} \frac{1}{n} \text{Tr} \left(\boldsymbol{\Sigma}^2 (\boldsymbol{\Sigma} + t \sigma_n \mathbf{I})^{-2} \right) \\ &= \int_0^\infty \frac{1}{(1 + tx^\alpha)^2} dx = t^{-1/\alpha} \cdot \frac{1}{\alpha} \int_0^\infty \frac{u^{1/\alpha-1}}{(1+u)^2} du \\ &= t^{-1/\alpha} \cdot \frac{\text{Beta}(1/\alpha, 2 - 1/\alpha)}{\alpha} = t^{-1/\alpha} \cdot \frac{\Gamma(1/\alpha) \Gamma(2 - 1/\alpha)}{\alpha \Gamma(2)} = t^{-1/\alpha} \cdot \frac{\Gamma(1/\alpha) \Gamma(1 - 1/\alpha)}{\alpha} \cdot \left(1 - \frac{1}{\alpha} \right) \\ &= t^{-1/\alpha} \frac{\pi/\alpha}{\sin(\pi/\alpha)} \cdot \left(1 - \frac{1}{\alpha} \right).\end{aligned}$$

and further $n^{-1} \text{Tr}(\boldsymbol{\Sigma}^2 (\boldsymbol{\Sigma} + \lambda_* \mathbf{I})^{-2}) \rightarrow (1 - \nu \mathbf{c}_*^{-1})(1 - \alpha^{-1})$. This then gives the variance

$$\mathbf{V}_n(\lambda) = \frac{\tau^2 n^{-1} \text{Tr}(\boldsymbol{\Sigma}^2 (\boldsymbol{\Sigma} + \lambda_* \mathbf{I})^{-2})}{1 - n^{-1} \text{Tr}(\boldsymbol{\Sigma}^2 (\boldsymbol{\Sigma} + \lambda_* \mathbf{I})^{-2})} = \frac{\tau^2 (1 - \nu \mathbf{c}_*^{-1})(\alpha - 1)}{1 + \nu \mathbf{c}_*^{-1}(\alpha - 1)} (1 + o_n(1)).$$

For the bias term, we can similarly write

$$\begin{aligned}
\lambda_\star^2 \langle \boldsymbol{\beta}, (\boldsymbol{\Sigma} + \lambda_\star \mathbf{I})^{-2} \boldsymbol{\Sigma} \boldsymbol{\beta} \rangle &= (1 + o_n(1)) \cdot \sigma_n \mathbf{c}_\star^2 \sum_{l=1}^{\infty} \frac{\sigma_l \sigma_n}{(\sigma_l + \lambda_\star)^2} \langle \boldsymbol{\beta}, \mathbf{v}_l \rangle^2 \\
&= (1 + o_n(1)) \cdot \sigma_n \mathbf{c}_\star^2 \sum_{l=1}^{\infty} \frac{(l/n)^\alpha}{(1 + \mathbf{c}_\star (l/n)^\alpha)^2} \langle \boldsymbol{\beta}, \mathbf{v}_l \rangle^2 \\
&= \sigma_n \mathbf{c}_\star^2 \int_0^\infty \frac{x^\alpha}{(1 + \mathbf{c}_\star x^\alpha)^2} dF_\beta(x) (1 + o_n(1)).
\end{aligned}$$

Together with $\text{Tr}(\boldsymbol{\Sigma}^2(\boldsymbol{\Sigma} + \lambda_\star \mathbf{I})^{-2}) = n(1 - \nu \mathbf{c}_\star^{-1})(1 - \alpha^{-1})(1 + o_n(1))$, we conclude the proof for this case.

Case II: regularly varying spectrum when $\alpha = 1$. Setting $\lambda = \nu n \sigma_n \log n$. For any $t > 0$, we can compute that

$$f_n(t; \lambda) = 1 - \frac{\nu \log n}{t} - \frac{1}{n} \text{Tr}(\boldsymbol{\Sigma}(\boldsymbol{\Sigma} + t \sigma_n \mathbf{I})^{-1}).$$

We first verify Assumption 1 holds. With $a_i = \mathcal{O}(1)$ bounded and $\alpha' > 1$, we indeed have that $\text{Tr}(\boldsymbol{\Sigma}) < \infty$ as

$$\text{Tr}(\boldsymbol{\Sigma}) = \sum_{l=1}^{\infty} \frac{a_l}{l(1 + \log l)^{\alpha'}} = \mathcal{O}\left(1 + \int_1^\infty \frac{M}{t(1 + \log t)^{\alpha'}} dt\right) = \mathcal{O}\left(-\frac{(1 + \log t)^{1-\alpha'}}{\alpha' - 1} \Big|_{t=1}^\infty\right) = \mathcal{O}\left(\frac{1}{\alpha' - 1}\right).$$

Since the sequence $\{a_i\}$ converge to a positive limit, we have for $k = \Omega(1)$

$$\begin{aligned}
\sum_{l=k}^d \sigma_l &= \Theta\left(\sum_{l=k}^{\infty} \frac{1}{l(1 + \log l)^{\alpha'}}\right) = \Theta\left(\sigma_k + \int_k^\infty \frac{1}{t(1 + \log t)^{\alpha'}} dt\right) = \Theta\left(\sigma_k + \frac{(1 + \log k)^{1-\alpha'}}{\alpha' - 1}\right) \\
&= \Theta_{\alpha'}(k \log k \sigma_k),
\end{aligned}$$

and therefore, we can take $\mathbf{d}_\Sigma(n) = \Theta_{\alpha'}(n \log n)$. We proceed to compute λ_\star . Taking any $t > 0$,

$$\begin{aligned}
&\lim_{n \rightarrow \infty} \frac{1}{n} \text{Tr}(\boldsymbol{\Sigma}(\boldsymbol{\Sigma} + t \sigma_n \log n \mathbf{I})^{-1}) \\
&= \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{l=1}^{\infty} \frac{\sigma_l}{\sigma_l + t \sigma_n \log n} = \lim_{M \rightarrow \infty} \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{l \geq M^{-1}n} \frac{\sigma_l}{\sigma_l + t \sigma_n \log n} \\
&= \lim_{M \rightarrow \infty} \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{l \geq M^{-1}n} \frac{\sigma_l}{t \sigma_n \log n} = \lim_{M \rightarrow \infty} \lim_{n \rightarrow \infty} \frac{(1 + \log n - \log M)^{1-\alpha'}}{t(\alpha' - 1)(\log n)^{1-\alpha'}} \\
&= \frac{1}{t(\alpha' - 1)}.
\end{aligned}$$

Recalling that \mathbf{c}_\star solves

$$1 = \nu \mathbf{c}_\star^{-1} + (\alpha' - 1)^{-1} \mathbf{c}_\star^{-1},$$

we then have

$$\lim_{n \rightarrow \infty} f_n(\mathbf{c}_\star \log n; \lambda) = 1 - \nu \mathbf{c}_\star^{-1} - \lim_{n \rightarrow \infty} \frac{1}{n} \text{Tr}(\boldsymbol{\Sigma}(\boldsymbol{\Sigma} + \mathbf{c}_\star \sigma_n \log n \mathbf{I})^{-1})$$

$$= 1 - \nu \mathbf{c}_\star^{-1} - (\alpha' - 1)^{-1} \mathbf{c}_\star^{-1} = 0,$$

and consequently by Eq. (5),

$$\lambda_\star = \mathbf{c}_\star \sigma_n \log n (1 + o_n(1)).$$

Taking the above display into Eq. (23), we get

$$\begin{aligned} \rho(\lambda) &= \frac{\boldsymbol{\theta}^\top \boldsymbol{\Sigma}^{\frac{1}{2}} (\boldsymbol{\Sigma} + \lambda_\star \mathbf{I})^{-1} \boldsymbol{\Sigma}^{\frac{1}{2}} \boldsymbol{\theta}}{\|\boldsymbol{\theta}\|^2 \text{Tr}(\boldsymbol{\Sigma}(\boldsymbol{\Sigma} + \lambda_\star \mathbf{I})^{-1})} = \frac{\sum_{l=1}^{\infty} (\sigma_l + \lambda_\star)^{-1} \langle \boldsymbol{\beta}, \mathbf{v}_l \rangle^2}{n(1 - \nu \mathbf{c}_\star^{-1}) \sum_{l=1}^{\infty} \sigma_l^{-1} \langle \boldsymbol{\beta}, \mathbf{v}_l \rangle^2} \cdot (1 + o_n(1)) \\ &= \frac{\sum_{l=1}^{\infty} \sigma_n (\sigma_l + \mathbf{c}_\star \sigma_n \log n)^{-1} \langle \boldsymbol{\beta}, \mathbf{v}_l \rangle^2}{n(1 - \nu \mathbf{c}_\star^{-1}) \sum_{l=1}^{\infty} \sigma_n \sigma_l^{-1} \langle \boldsymbol{\beta}, \mathbf{v}_l \rangle^2} \cdot (1 + o_n(1)) \\ &= \frac{\sum_{l=1}^{\infty} l/n \cdot (1 + l \cdot (\mathbf{c}_\star \log n/n))^{-1} \langle \boldsymbol{\beta}, \mathbf{v}_l \rangle^2}{n(1 - \nu \mathbf{c}_\star^{-1}) \sum_{l=1}^{\infty} l/n \cdot \langle \boldsymbol{\beta}, \mathbf{v}_l \rangle^2} \cdot (1 + o_n(1)) \\ &= \frac{\sum_{l=1}^{\infty} l \cdot (\log n/n) \cdot (1 + l \cdot (\mathbf{c}_\star \log n/n))^{-1} \langle \boldsymbol{\beta}, \mathbf{v}_l \rangle^2}{n(1 - \nu \mathbf{c}_\star^{-1}) \sum_{l=1}^{\infty} l \cdot (\log n/n) \cdot \langle \boldsymbol{\beta}, \mathbf{v}_l \rangle^2} \cdot (1 + o_n(1)) \\ &= \frac{\int_0^{\infty} x (1 + \mathbf{c}_\star x)^{-1} dF_{\boldsymbol{\beta}}(x)}{n(1 - \nu \mathbf{c}_\star^{-1}) \int_0^{\infty} x dF_{\boldsymbol{\beta}}(x)} \cdot (1 + o_n(1)), \end{aligned}$$

where in the last line we use $F_{\boldsymbol{\beta}}(x) = \sum_{k=1}^{\lfloor (n/\log n)x \rfloor} \langle \boldsymbol{\beta}, \mathbf{v}_k \rangle^2$. Thus we can have $\rho(\lambda) = \Omega(n^{-2+\theta})$ provided the condition

$$\int_0^{\infty} x dF_{\boldsymbol{\beta}}(x) = \mathcal{O}\left(n^{1-\theta} \int_0^{\infty} x (1 + \mathbf{c}_\star x)^{-1} dF_{\boldsymbol{\beta}}(x)\right).$$

Setting $\gamma = (1 - \theta)/3$, we can invoke Proposition 4.3 and obtain $\mathcal{V}_{\mathbf{X}}(\lambda) = \mathbf{V}_n(\lambda)(1 + o_n(1))$, $\mathcal{B}_{\mathbf{X}}(\lambda) = \mathbf{B}_n(\lambda)(1 + o_n(1))$.

For the variance $\mathbf{V}_n(\lambda)$, we note

$$\begin{aligned} &\lim_{n \rightarrow \infty} \frac{\log n}{n} \text{Tr}(\boldsymbol{\Sigma}^2 (\boldsymbol{\Sigma} + t \sigma_n \log n \mathbf{I})^{-2}) \\ &= \lim_{n \rightarrow \infty} \frac{\log n}{n} \sum_{l=1}^{\infty} \frac{\sigma_l^2}{(\sigma_l + t \sigma_n \log n)^2} = \lim_{M \rightarrow \infty} \lim_{n \rightarrow \infty} \frac{\log n}{n} \sum_{M^{-1}n/\log n \leq l \leq Mn/\log n} \frac{\sigma_l^2}{(\sigma_l + t \sigma_n \log n)^2} \\ &= \lim_{M \rightarrow \infty} \lim_{n \rightarrow \infty} \frac{\log n}{n} \sum_{M^{-1}n/\log n \leq l \leq Mn/\log n} \frac{1}{(1 + tl \log n/n)^2} = \int_0^{\infty} \frac{1}{(1 + tx)^2} dx = \frac{1}{t}. \end{aligned}$$

Substituting in λ_\star , we thus have $n^{-1} \text{Tr}(\boldsymbol{\Sigma}^2 (\boldsymbol{\Sigma} + \lambda_\star \mathbf{I})^{-2}) = (1 + o_n(1))/(\mathbf{c}_\star \log n)$, which further implies that

$$\mathbf{V}_n(0) = \frac{\tau^2 n^{-1} \text{Tr}(\boldsymbol{\Sigma}^2 (\boldsymbol{\Sigma} + \lambda_\star(0) \mathbf{I})^{-2})}{1 - n^{-1} \text{Tr}(\boldsymbol{\Sigma}^2 (\boldsymbol{\Sigma} + \lambda_\star(0) \mathbf{I})^{-2})} = \frac{\tau^2}{\mathbf{c}_\star \log n} (1 + o_n(1)).$$

Finally for the bias, we have

$$\lambda_\star^2 \langle \boldsymbol{\beta}, (\boldsymbol{\Sigma} + \lambda_\star \mathbf{I})^{-2} \boldsymbol{\Sigma} \boldsymbol{\beta} \rangle = (1 + o_n(1)) \cdot \mathbf{c}_\star^2 \sigma_n \log n \sum_{l=1}^{\infty} \frac{\sigma_l \sigma_n \log n}{(\sigma_l + \lambda_\star)^2} \langle \boldsymbol{\beta}, \mathbf{v}_l \rangle^2$$

$$\begin{aligned}
&= (1 + o_n(1)) \cdot \mathbf{c}_\star^2 \sigma_n \log n \sum_{l=1}^{\infty} \frac{(l \log n/n)}{(1 + \mathbf{c}_\star l \log n/n)^2} \langle \boldsymbol{\beta}, \mathbf{v}_l \rangle^2 \\
&= \mathbf{c}_\star^2 \sigma_n \log n \int_0^{\infty} \frac{x}{(1 + \mathbf{c}_\star x)^2} dF_{\boldsymbol{\beta}}(x) (1 + o_n(1)).
\end{aligned}$$

Combining with $n^{-1} \text{Tr}(\boldsymbol{\Sigma}^2(\boldsymbol{\Sigma} + \lambda_\star \mathbf{I})^{-2}) = (1 + o_n(1))/(\mathbf{c}_\star \log n)$, it holds that

$$\mathbf{B}_n(\lambda) = \mathbf{c}_\star^2 \sigma_n \log n \int_0^{\infty} \frac{x}{(1 + \mathbf{c}_\star x)^2} dF_{\boldsymbol{\beta}}(x) (1 + o_n(1)).$$

Case III: a non-regularly varying spectrum. Take $\lambda = \nu n \sigma_n$. For any $t > 0$, we can compute that

$$f_n(t; \lambda) = 1 - \frac{\nu}{t} - \frac{1}{n} \text{Tr}(\boldsymbol{\Sigma}(\boldsymbol{\Sigma} + t \sigma_n \mathbf{I})^{-1}).$$

If $\sigma_k = p^{-s}$, we can easily have $\sigma_l \leq p^{-r-s}$ if $q^r k \leq l < q^{r+1} k$. This immediately yields

$$\sum_{l=k}^d \sigma_l \leq \sum_{r=0}^{\infty} (q^{r+1} k - q^r k) \cdot p^{-r-s} \leq k \sigma_k \sum_{r=0}^{\infty} \frac{q^{r+1}}{p^r} = \mathcal{O}_{p,q}(k \sigma_k),$$

as $q < p$ and the geometric sum converges. We can thus take $\mathbf{d}_{\boldsymbol{\Sigma}}(n) = \mathcal{O}_{p,q}(n)$. For λ_\star , using that

$$\text{Tr}(\boldsymbol{\Sigma}(\boldsymbol{\Sigma} + t \sigma_n \mathbf{I})^{-1}) = \sum_{l=0}^{\infty} \frac{(q^{l+1} - q^l) p^{-l}}{p^{-l} + t p^{-s_\star}} = (q^{s_\star+1} - q^{s_\star}) \cdot \sum_{l=0}^{\infty} \frac{q^{l-s_\star}}{1 + t p^{l-s_\star}}.$$

Since $s_\star \rightarrow \infty$ as n tends to infinity, we have

$$\text{Tr}(\boldsymbol{\Sigma}(\boldsymbol{\Sigma} + t \sigma_n \mathbf{I})^{-1}) = n \rho_\star^{-1} \cdot G_{p,q,1}(t) (1 + o_n(1)).$$

Hence

$$\lim_{n \rightarrow \infty} f_n(t; \lambda) = 1 - \nu t^{-1} - \rho_\star^{-1} \cdot G_{p,q,1}(t).$$

While the right hand side is increasing in t ranging in $(-\infty, 1)$. There exists a unique $\mathbf{c}_\star = \mathbf{c}_\star(\nu)$ solving

$$\lim_{n \rightarrow \infty} f_n(\mathbf{c}_\star; \lambda) = 0,$$

and substituting into Eq. (5) yields

$$\lambda_\star = \mathbf{c}_\star \sigma_n (1 + o_n(1)).$$

Next we compute $\rho(\lambda)$ from Eq. (23),

$$\begin{aligned}
\rho(\lambda) &= \frac{\boldsymbol{\theta}^\top \boldsymbol{\Sigma}^{\frac{1}{2}} (\boldsymbol{\Sigma} + \lambda_\star \mathbf{I})^{-1} \boldsymbol{\Sigma}^{\frac{1}{2}} \boldsymbol{\theta}}{\|\boldsymbol{\theta}\|^2 \text{Tr}(\boldsymbol{\Sigma}(\boldsymbol{\Sigma} + \lambda_\star \mathbf{I})^{-1})} = \frac{\sum_{l=1}^{\infty} (\sigma_l + \lambda_\star)^{-1} \langle \boldsymbol{\beta}, \mathbf{v}_l \rangle^2}{n(1 - \nu \mathbf{c}_\star^{-1}) \sum_{l=1}^{\infty} \sigma_l^{-1} \langle \boldsymbol{\beta}, \mathbf{v}_l \rangle^2} \cdot (1 + o_n(1)) \\
&= \frac{\sum_{l=1}^{\infty} \sigma_n (\sigma_l + \lambda_\star)^{-1} \langle \boldsymbol{\beta}, \mathbf{v}_l \rangle^2}{n(1 - \nu \mathbf{c}_\star^{-1}) \sum_{l=1}^{\infty} \sigma_n \sigma_l^{-1} \langle \boldsymbol{\beta}, \mathbf{v}_l \rangle^2} \cdot (1 + o_n(1))
\end{aligned}$$

$$\begin{aligned}
&= \frac{\sum_{s=0}^{\infty} p^{s-s_*} / (1 + \mathbf{c}_* p^{s-s_*}) \sum_{l=q^s}^{q^{s+1}-1} \langle \boldsymbol{\beta}, \mathbf{v}_l \rangle^2}{n(1 - \nu \mathbf{c}_*^{-1}) \sum_{s=0}^{\infty} p^{s-s_*} \sum_{l=q^s}^{q^{s+1}-1} \langle \boldsymbol{\beta}, \mathbf{v}_l \rangle^2} \cdot (1 + o_n(1)) \\
&= \frac{\int_0^{\infty} p^{x-s_*} / (1 + \mathbf{c}_* p^{x-s_*}) dF_{\boldsymbol{\beta}}(x)}{n(1 - \nu \mathbf{c}_*^{-1}) \int_0^{\infty} p^{x-s_*} dF_{\boldsymbol{\beta}}(x)} \cdot (1 + o_n(1)).
\end{aligned}$$

Given the ‘‘rapid-decay’’ property

$$\int_0^{\infty} p^{x-s_*} dF_{\boldsymbol{\beta}}(x) = \mathcal{O} \left(n^{1-\theta} \int_0^{\infty} p^{x-s_*} (1 + \mathbf{c}_* p^{x-s_*})^{-1} dF_{\boldsymbol{\beta}}(x) \right),$$

we have $\rho(\lambda) = \Omega(n^{-2+\theta})$ and Proposition 4.3 holds with $\gamma = (1 - \theta)/3$, implying that $\mathcal{V}_{\mathbf{X}}(\lambda) = \mathbf{V}_n(\lambda)(1 + o_n(1))$ and $\mathcal{B}_{\mathbf{X}}(\lambda) = \mathbf{B}_n(\lambda)(1 + o_n(1))$.

To compute the effective variance $\mathbf{V}_n(\lambda)$, we first note that

$$\text{Tr}(\boldsymbol{\Sigma}^2(\boldsymbol{\Sigma} + t\sigma_n \mathbf{I})^{-2}) = (q^{s_*+1} - q^{s_*}) \cdot \sum_{l=0}^{\infty} \frac{q^{l-s_*}}{(1 + tp^{l-s_*})^2} = n\rho_*^{-1} \cdot G_{p,q,2}(t)(1 + o_n(1)).$$

Thus

$$\mathbf{V}_n(\lambda) = \frac{\tau^2 \text{Tr}(\boldsymbol{\Sigma}^2(\boldsymbol{\Sigma} + \lambda_* \mathbf{I})^{-2})}{n - \text{Tr}(\boldsymbol{\Sigma}^2(\boldsymbol{\Sigma} + \lambda_* \mathbf{I})^{-2})} = \frac{G_{p,q,2}(\mathbf{c}_*) \tau^2}{\rho_* - G_{p,q,2}(\mathbf{c}_*)} (1 + o_n(1)).$$

For the bias term, we have

$$\begin{aligned}
\lambda_*^2 \langle \boldsymbol{\beta}, (\boldsymbol{\Sigma} + \lambda_* \mathbf{I})^{-2} \boldsymbol{\Sigma} \boldsymbol{\beta} \rangle &= (1 + o_n(1)) \mathbf{c}_*^2 \sigma_n \sum_{l=0}^{\infty} \frac{\sigma_l \sigma_n}{(\sigma_l + \mathbf{c}_* \sigma_n)^2} \langle \boldsymbol{\beta}, \mathbf{v}_l \rangle^2 \\
&= (1 + o_n(1)) \mathbf{c}_*^2 \sigma_n \sum_{s=0}^{\infty} \left\{ \frac{p^{s-s_*}}{(1 + \mathbf{c}_* p^{s-s_*})^2} \cdot \sum_{l=q^s}^{q^{s+1}-1} \langle \boldsymbol{\beta}, \mathbf{v}_l \rangle^2 \right\} \\
&= \mathbf{c}_*^2 \sigma_n \int_0^{\infty} \frac{p^{x-s_*}}{(1 + \mathbf{c}_* p^{x-s_*})^2} dF_{\boldsymbol{\beta}}(x) (1 + o_n(1)).
\end{aligned}$$

We conclude the proof for $\mathbf{B}_n(\lambda)$ by substituting in $n^{-1} \text{Tr}(\boldsymbol{\Sigma}^2(\boldsymbol{\Sigma} + \lambda_* \mathbf{I})^{-2}) = (1 + o_n(1)) \rho_*^{-1} G_{p,q,2}(\mathbf{c}_*)$.