

# The generalization error of random features regression: Precise asymptotics and double descent curve

Song Mei\* and Andrea Montanari†

October 23, 2019

## Abstract

Deep learning methods operate in regimes that defy the traditional statistical mindset. Neural network architectures often contain more parameters than training samples, and are so rich that they can interpolate the observed labels, even if the latter are replaced by pure noise. Despite their huge complexity, the same architectures achieve small generalization error on real data.

This phenomenon has been rationalized in terms of a so-called ‘double descent’ curve. As the model complexity increases, the test error follows the usual U-shaped curve at the beginning, first decreasing and then peaking around the interpolation threshold (when the model achieves vanishing training error). However, it descends again as model complexity exceeds this threshold. The global minimum of the test error is found above the interpolation threshold, often in the extreme overparametrization regime in which the number of parameters is much larger than the number of samples. Far from being a peculiar property of deep neural networks, elements of this behavior have been demonstrated in much simpler settings, including linear regression with random covariates.

In this paper we consider the problem of learning an unknown function over the  $d$ -dimensional sphere  $\mathbb{S}^{d-1}$ , from  $n$  i.i.d. samples  $(\mathbf{x}_i, y_i) \in \mathbb{S}^{d-1} \times \mathbb{R}$ ,  $i \leq n$ . We perform ridge regression on  $N$  random features of the form  $\sigma(\mathbf{w}_a^\top \mathbf{x})$ ,  $a \leq N$ . This can be equivalently described as a two-layers neural network with random first-layer weights. We compute the precise asymptotics of the test error, in the limit  $N, n, d \rightarrow \infty$  with  $N/d$  and  $n/d$  fixed. This provides the first analytically tractable model that captures all the features of the double descent phenomenon without assuming ad hoc misspecification structures. In particular, above a critical value of the signal-to-noise ratio, minimum test error is achieved by extremely overparametrized interpolators, i.e., networks that have a number of parameters much larger than the sample size, and vanishing training error.

## Contents

<b>1</b>	<b>Introduction</b>	<b>2</b>
<b>2</b>	<b>Results and insights: An informal overview</b>	<b>6</b>
<b>3</b>	<b>Related literature</b>	<b>8</b>
<b>4</b>	<b>Main results</b>	<b>9</b>
4.1	Statement of main result . . . . .	11
4.2	Simplifying the asymptotic risk in special cases . . . . .	11
4.2.1	Ridgeless limit . . . . .	12
4.2.2	Highly overparametrized regime . . . . .	13
4.2.3	Large sample limit . . . . .	15
<b>5</b>	<b>Asymptotics of the training error</b>	<b>15</b>
5.1	Numerical illustrations . . . . .	16

---

\*Institute for Computational and Mathematical Engineering, Stanford University

†Department of Electrical Engineering and Department of Statistics, Stanford University

<b>6</b>	<b>An equivalent Gaussian covariates model</b>	<b>18</b>
<b>7</b>	<b>Proof of Theorem 2</b>	<b>19</b>
<b>A</b>	<b>Technical background and notations</b>	<b>27</b>
A.1	Notations . . . . .	27
A.2	Functional spaces over the sphere . . . . .	27
A.3	Gegenbauer polynomials . . . . .	28
A.4	Hermite polynomials . . . . .	29
<b>B</b>	<b>Proof of Proposition 7.1</b>	<b>29</b>
B.1	Proof of Lemma B.1 . . . . .	32
B.2	Proof of Lemma B.2 . . . . .	33
B.3	Proof of Lemma B.3 . . . . .	36
B.4	Proof of Lemma B.4 and B.5 . . . . .	38
B.5	Preliminary lemmas . . . . .	42
<b>C</b>	<b>Proof of Proposition 7.2</b>	<b>46</b>
C.1	Equivalence between Gaussian and sphere vectors . . . . .	46
C.2	Properties of the fixed point equations . . . . .	48
C.3	Key lemma: Stieltjes transforms are approximate fixed point . . . . .	49
C.4	Properties of Stieltjes transforms . . . . .	50
C.5	Leave-one-out argument: Proof of Lemma C.4 . . . . .	51
C.6	Proof of Proposition 7.2 . . . . .	59
<b>D</b>	<b>Proof of Proposition 7.3</b>	<b>61</b>
<b>E</b>	<b>Proof of Proposition 7.4</b>	<b>62</b>
E.1	Properties of the Stieltjes transforms and the log determinant . . . . .	62
E.2	Proof of Proposition 7.4 . . . . .	65
<b>F</b>	<b>Proof of Theorem 3, 4, and 5</b>	<b>66</b>
F.1	Proof of Theorem 3 . . . . .	66
F.2	Proof of Theorem 4 . . . . .	67
F.3	Proof of Theorem 5 . . . . .	67
<b>G</b>	<b>Proof of Proposition 4.1 and 4.2</b>	<b>68</b>
G.1	Proof of Proposition 4.1 . . . . .	68
G.2	Proof of Proposition 4.2 . . . . .	68
<b>H</b>	<b>Proof sketch for Theorem 6</b>	<b>71</b>

# 1 Introduction

Statistical lore recommends not to use models that have too many parameters since this will lead to ‘over-fitting’ and poor generalization. Indeed, a plot of the test error as a function of the model complexity often reveals a U-shaped curve. The test error first decreases because the model is less and less biased, but then increases because of a variance explosion [HTF09]. In particular, the interpolation threshold, i.e., the threshold in model complexity above which the training error vanishes (the model completely interpolates the data), corresponds to a large test error. It seems wise to keep the model complexity well below this threshold in order to obtain a small generalization error.

These classical prescriptions are in stark contrast with the current practice in deep learning. The number of parameters of modern neural networks can be much larger than the number of training samples, and the resulting models are often so complex that they can perfectly interpolate the data. Even more surprisingly, they can interpolate the data when the actual labels are replaced by pure noise [ZBH<sup>+</sup>16]. Despite such

a large complexity, these models have small test error and can outperform others trained in the classical underparametrized regime.

This behavior has been rationalized in terms of a so-called ‘double-descent’ curve [BMM18, BHMM18]. A plot of the test error as a function of the model complexity follows the traditional U-shaped curve until the interpolation threshold. However, after a peak at the interpolation threshold, the test error decreases, and attains a global minimum in the overparametrized regime. In fact, the minimum error often appears to be ‘at infinite complexity’: the more overparametrized is the model, the smaller is the error. It is conjectured that the good generalization behavior in this highly overparametrized regime is due to the implicit regularization induced by gradient descent learning: among all interpolating models, gradient descent selects the simplest one, in a suitable sense. An example of double descent curve is plotted in Fig. 1. The main contribution of this paper is to describe a natural, analytically tractable model leading to this generalization curve, and to derive precise formulae for the same curve, in a suitable asymptotic regime.

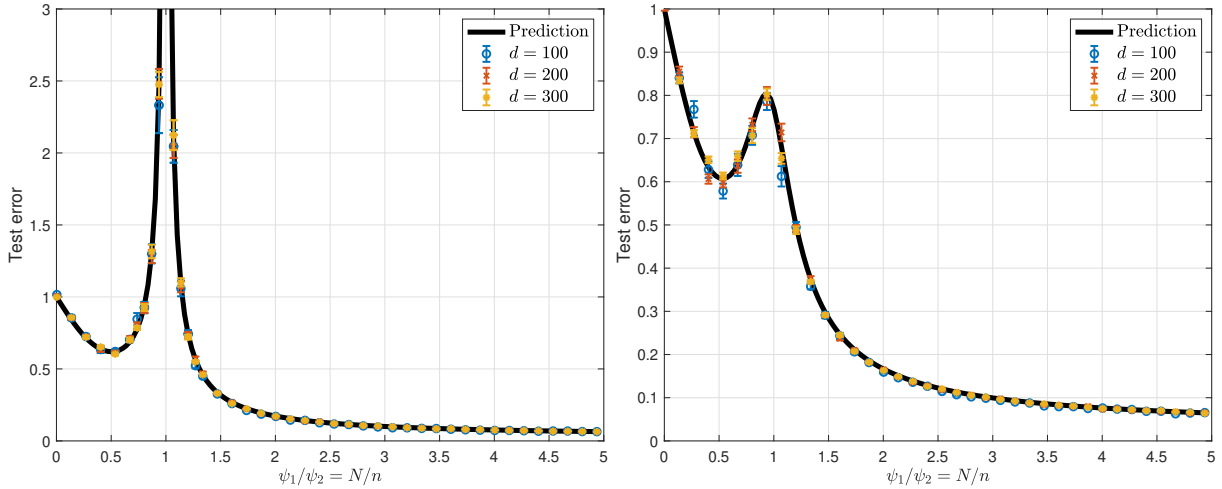


Figure 1: Random features ridge regression with ReLU activation ( $\sigma = \max\{x, 0\}$ ). Data are generated via  $y_i = \langle \beta_1, \mathbf{x}_i \rangle$  (zero noise) with  $\|\beta_1\|_2^2 = 1$ , and  $\psi_2 = n/d = 3$ . Left frame: regularization  $\lambda = 10^{-8}$  (we didn’t set  $\lambda = 0$  exactly for numerical stability). Right frame:  $\lambda = 10^{-3}$ . The continuous black line is our theoretical prediction, and the colored symbols are numerical results for several dimensions  $d$ . Symbols are averages over 20 instances and the error bars report the standard error of the means over these 20 instances.

The double-descent scenario is far from being specific to neural networks, and was instead demonstrated empirically in a variety of models including random forests and random features models [BHMM18]. Recently, several elements of this scenario were established analytically in simple least square regression, with certain probabilistic models for the random covariates [AS17, HMRT19, BHX19]. These papers consider a setting in which we are given i.i.d. samples  $(y_i, \mathbf{x}_i) \in \mathbb{R} \times \mathbb{R}^d$ ,  $i \leq n$ , where  $y_i$  is a response variable which depends on covariates  $\mathbf{x}_i$  via  $y_i = \langle \beta, \mathbf{x}_i \rangle + \varepsilon_i$ , with  $\mathbb{E}(\varepsilon_i) = 0$  and  $\mathbb{E}(\varepsilon_i^2) = \tau^2$ ; or in matrix notation,  $\mathbf{y} = \mathbf{X}\beta + \varepsilon$ . The authors study the test error of ‘ridgeless least square regression’  $\hat{\beta} = (\mathbf{X}^\top \mathbf{X})^\dagger \mathbf{X}^\top \mathbf{y}$ , and use random matrix theory to derive its precise asymptotics in the limit  $n, d \rightarrow \infty$  with  $d/n = \gamma$  fixed, when  $\mathbf{x}_i = \Sigma^{1/2} \mathbf{z}_i$  with  $\mathbf{z}_i$  a vector with i.i.d. entries.

Despite its simplicity, this random covariates model captures several features of the double descent scenario. In particular, the asymptotic generalization curve is U-shaped for  $\gamma < 1$ , diverging at the interpolation threshold  $\gamma = 1$ , and descends again after that threshold. The divergence at  $\gamma = 1$  is explained by an explosion in the variance, which is in turn related to a divergence of the condition number of the random matrix  $\mathbf{X}$ . At the same time, this simple model misses some interesting features that are observed in more complex settings: (i) In the Gaussian covariates model, the global minimum of the test error is achieved in the underparametrized regime  $\gamma < 1$ , unless ad-hoc misspecification structure is assumed; (ii) The number of parameters is tied to the covariates dimension  $d$  and hence the effects of overparametrization are not isolated from the effects of the ambient dimensions; (iii) Ridge regression, with some regularization  $\lambda > 0$ , is always found to outperform the ridgeless limit  $\lambda \rightarrow 0$ . Moreover, this linear model is not directly connected

to actual neural networks, which are highly nonlinear in the covariates  $\mathbf{x}_i$ .

In this paper, we study the random features model of Rahimi and Recht [RR08]. The random features model can be viewed either as a randomized approximation to kernel ridge regression, or as a two-layers neural networks with random first layer wights. We compute the precise asymptotics of the test error and show that it reproduces all the qualitative features of the double-descent scenario.

More precisely, we consider the problem of learning a function  $f_d \in L^2(\mathbb{S}^{d-1}(\sqrt{d}))$  on the  $d$ -dimensional sphere. (Here and below  $\mathbb{S}^{d-1}(r)$  denotes the sphere of radius  $r$  in  $d$  dimensions, and we set  $r = \sqrt{d}$  without loss of generality.) We are given i.i.d. data  $\{(\mathbf{x}_i, y_i)\}_{i \leq n} \sim_{iid} \mathbb{P}_{\mathbf{x}, y}$ , where  $\mathbf{x}_i \sim_{iid} \text{Unif}(\mathbb{S}^{d-1}(\sqrt{d}))$  and  $y_i = f_d(\mathbf{x}_i) + \varepsilon_i$ , with  $\varepsilon_i \sim_{iid} \mathbb{P}_\varepsilon$  independent of  $\mathbf{x}_i$ . The noise distribution satisfies  $\mathbb{E}_\varepsilon(\varepsilon_1) = 0$ ,  $\mathbb{E}_\varepsilon(\varepsilon_1^2) = \tau^2$ , and  $\mathbb{E}_\varepsilon(\varepsilon_1^4) < \infty$ . We fit these training data using the random features (RF) model, which is defined as the function class

$$\mathcal{F}_{\text{RF}}(\Theta) = \left\{ f(\mathbf{x}; \mathbf{a}, \Theta) \equiv \sum_{i=1}^N a_i \sigma(\langle \theta_i, \mathbf{x} \rangle / \sqrt{d}) : a_i \in \mathbb{R} \forall i \in [N] \right\}. \quad (1)$$

Here,  $\Theta \in \mathbb{R}^{N \times d}$  is a matrix whose  $i$ -th row is the vector  $\theta_i$ , which is chosen randomly, and independent of the data. In order to simplify some of the calculations below, we will assume the normalization  $\|\theta_i\|_2 = \sqrt{d}$ , which justifies the factor  $1/\sqrt{d}$  in the above expression, yielding  $\langle \theta_i, \mathbf{x}_j \rangle / \sqrt{d}$  of order one. As mentioned above, the functions in  $\mathcal{F}_{\text{RF}}(\Theta)$  are two-layers neural networks, except that the first layer is kept constant. A substantial literature draws connections between random features models, fully trained neural networks, and kernel methods. We refer to Section 3 for a summary of this line of work.

We learn the coefficients  $\mathbf{a} = (a_i)_{i \leq N}$  by performing ridge regression

$$\hat{\mathbf{a}}(\lambda) = \arg \min_{\mathbf{a} \in \mathbb{R}^N} \left\{ \frac{1}{n} \sum_{j=1}^n \left( y_j - \sum_{i=1}^N a_i \sigma(\langle \theta_i, \mathbf{x}_j \rangle / \sqrt{d}) \right)^2 + \frac{N\lambda}{d} \|\mathbf{a}\|_2^2 \right\}. \quad (2)$$

The choice of ridge penalty is motivated by the connection to kernel ridge regression, of which this method can be regarded as a finite-rank approximation. Further, the ridge regularization path is naturally connected to the path of gradient flow with respect to the mean square error  $\sum_{i \leq n} (y_i - f(\mathbf{x}_i; \mathbf{a}, \Theta))^2$ , starting at  $\mathbf{a} = 0$ . In particular, gradient flow converges to the ridgeless limit ( $\lambda \rightarrow 0$ ) of  $\hat{\mathbf{a}}(\lambda)$ , and there is a correspondence between positive  $\lambda$ , and early stopping in gradient descent [YRC07].

We are interested in the ‘prediction’ or ‘test’ error (which we will also call ‘generalization error,’ with a slight abuse of terminology), that is the mean square error on predicting  $f_d(\mathbf{x})$  for  $\mathbf{x} \sim \text{Unif}(\mathbb{S}^{d-1}(\sqrt{d}))$  a fresh sample independent of the training data  $\mathbf{X} = (\mathbf{x}_i)_{i \leq n}$ , noise  $\varepsilon = (\varepsilon_i)_{i \leq n}$ , and the random features  $\Theta = (\theta_a)_{a \leq N}$ :

$$R_{\text{RF}}(f_d, \mathbf{X}, \Theta, \lambda) = \mathbb{E}_{\mathbf{x}} \left[ \left( f_d(\mathbf{x}) - f(\mathbf{x}; \hat{\mathbf{a}}(\lambda), \Theta) \right)^2 \right]. \quad (3)$$

Notice that we do not take expectation with respect to the training data  $\mathbf{X}$ , the random features  $\Theta$  or the data noise  $\varepsilon$ . This is not very important, because we will show that  $R_{\text{RF}}(f_d, \mathbf{X}, \Theta, \lambda)$  concentrates around the expectation  $\bar{R}_{\text{RF}}(f_d, \lambda) \equiv \mathbb{E}_{\mathbf{X}, \Theta, \varepsilon} R_{\text{RF}}(f_d, \mathbf{X}, \Theta, \lambda)$ . We study the following setting

- The random features are uniformly distributed on a sphere:  $(\theta_i)_{i \leq N} \sim_{iid} \text{Unif}(\mathbb{S}^{d-1}(\sqrt{d}))$ .
- $N, n, d$  lie in a proportional asymptotics regime. Namely,  $N, n, d \rightarrow \infty$  with  $N/d \rightarrow \psi_1$ ,  $n/d \rightarrow \psi_2$  for some  $\psi_1, \psi_2 \in (0, \infty)$ .
- We consider two models for the regression function  $f_d$ : (1) A linear model:  $f_d(\mathbf{x}) = \beta_{d,0} + \langle \beta_{d,1}, \mathbf{x} \rangle$ , where  $\beta_{d,1} \in \mathbb{R}^d$  is arbitrary with  $\|\beta_{d,1}\|_2^2 = F_1^2$ ; (2) A nonlinear model:  $f_d(\mathbf{x}) = \beta_{d,0} + \langle \beta_{d,1}, \mathbf{x} \rangle + f_d^{\text{NL}}(\mathbf{x})$  where the nonlinear component  $f_d^{\text{NL}}(\mathbf{x})$  is a centered isotropic Gaussian process indexed by  $\mathbf{x} \in \mathbb{S}^{d-1}(\sqrt{d})$ . (Note that the linear model is a special case of the nonlinear one, but we prefer to keep the former distinct since it is purely deterministic.)

Within this setting, we are able to determine the precise asymptotics of the prediction error, as an explicit function of the dimension parameters  $\psi_1, \psi_2$ , the noise level  $\tau^2$ , the activation function  $\sigma$ , the regularization

parameter  $\lambda$ , and the power of linear and nonlinear components of  $f_d$ :  $F_1^2$  and  $F_\star^2 \equiv \lim_{d \rightarrow \infty} \mathbb{E}\{f_d^{\text{NL}}(\mathbf{x})^2\}$ . The resulting formulae are somewhat complicated, and we defer them to Section 4, limiting ourselves to give the general form of our result for the linear model.

**Theorem 1.** *(Linear truth, formulas omitted) Let  $\sigma : \mathbb{R} \rightarrow \mathbb{R}$  be weakly differentiable, with  $\sigma'$  be a weak derivative of  $\sigma$ . Assume  $|\sigma(u)|, |\sigma'(u)| \leq c_0 e^{c_1|u|}$  for some constants  $c_0, c_1 < \infty$ . Define the parameters  $\mu_0, \mu_1, \mu_\star, \zeta$ , and the signal-to-noise ratio  $\rho \in [0, \infty]$ , via*

$$\mu_0 = \mathbb{E}[\sigma(G)], \quad \mu_1 = \mathbb{E}[G\sigma(G)], \quad \mu_\star^2 = \mathbb{E}[\sigma(G)^2] - \mu_0^2 - \mu_1^2, \quad \zeta \equiv \mu_1^2/\mu_\star^2, \quad \rho \equiv F_1^2/\tau^2, \quad (4)$$

where expectation is taken with respect to  $G \sim \mathcal{N}(0, 1)$ . Assume  $\mu_0, \mu_1, \mu_\star \neq 0$ .

Then, for  $f_d$  linear, and in the setting described above, for any  $\lambda > 0$ , we have

$$R_{\text{RF}}(f_d, \mathbf{X}, \boldsymbol{\Theta}, \lambda) = (F_1^2 + \tau^2) \mathcal{R}(\rho, \zeta, \psi_1, \psi_2, \lambda/\mu_\star^2) + o_{d, \mathbb{P}}(1), \quad (5)$$

where  $\mathcal{R}(\rho, \zeta, \psi_1, \psi_2, \bar{\lambda})$  is explicitly given in Definition 1.

Section 4.1 also contains an analogous statement for the nonlinear model.

**Remark 1.** As usual, we can decompose the risk  $R_{\text{RF}}(f_d, \mathbf{X}, \boldsymbol{\Theta}, \lambda) = \|f_d - \hat{f}\|_{L^2}^2$  (where  $\hat{f}(\mathbf{x}) = f(\mathbf{x}; \hat{\mathbf{a}}(\lambda), \boldsymbol{\Theta})$ ) into a variance component  $\|\hat{f} - \mathbb{E}_\epsilon(\hat{f})\|_{L^2}^2$ , and a bias component  $\|f_d - \mathbb{E}_\epsilon(\hat{f})\|_{L^2}^2$ . The asymptotics of the variance component was computed already in [HMRT19, Section 7].

As it should be clear from the next sections, computing the full prediction error requires new technical ideas, and leads to new insights.

**Remark 2.** Theorem 1 and its generalizations stated below require  $\lambda > 0$  fixed as  $N, n, d \rightarrow \infty$ . We can then consider the ridgeless limit by taking  $\lambda \rightarrow 0$ . Let us stress that this does not necessarily yield the prediction risk of the min-norm least square estimator that is also given by the limit  $\hat{\mathbf{a}}(0+) \equiv \lim_{\lambda \rightarrow 0} \hat{\mathbf{a}}(\lambda)$  at  $N, n, d$  fixed. Denoting by  $\mathbf{Z} = \sigma(\mathbf{X}\boldsymbol{\Theta}^\top/\sqrt{d})/\sqrt{d}$  the design matrix, the latter is given by  $\hat{\mathbf{a}}(0+) = (\mathbf{Z}^\top \mathbf{Z})^\dagger \mathbf{Z}^\top \mathbf{y}/\sqrt{d}$ . While we conjecture that indeed this is the same as taking  $\lambda \rightarrow 0$  in the asymptotic expression of Theorem 1, establishing this rigorously would require proving that the limits  $\lambda \rightarrow 0$  and  $d \rightarrow \infty$  can be exchanged. We leave this to future work.

Figure 1 reports numerical results for learning a linear function  $f_d(\mathbf{x}) = \langle \boldsymbol{\beta}_1, \mathbf{x} \rangle$ ,  $\|\boldsymbol{\beta}_1\|_2^2 = 1$  with  $\mathbb{E}[\varepsilon^2] = 0$  using ReLU activation function  $\sigma(x) = \max\{x, 0\}$  and  $\psi_2 = n/d = 3$ . We use minimum  $\ell_2$ -norm least squares (the  $\lambda \rightarrow 0$  limit of Eq. (2), left figure) and regularized least squares with  $\lambda = 10^{-3}$  (right figure), and plot the prediction error as a function of the number of parameters per dimension  $\psi_1 = N/d$ . We compare the numerical results with the asymptotic formula  $\mathcal{R}(\infty, \zeta, \psi_1, \psi_2, \lambda/\mu_\star^2)$ . The agreement is excellent and displays all the key features of the double descent phenomenon, as discussed in the next section.

The proof of Theorem 1 builds on ideas from random matrix theory. A careful look at these arguments unveils an interesting phenomenon. While the random features  $\{\sigma(\langle \boldsymbol{\theta}_i, \mathbf{x} \rangle/\sqrt{d})\}_{i \leq d}$  are highly non-Gaussian, it is possible to construct a Gaussian covariates model with the same asymptotic prediction error as for the random features model. Apart from being mathematically interesting, this finding provides additional intuition for the behavior of random features models, and opens the way to some interesting future directions. In particular, [MRSY19] uses this Gaussian covariates proxy to analyze maximum margin classification using random features.

The rest of the paper is organized as follows:

- In Section 2 we summarize the main insights that can be extracted from the asymptotic theory, and illustrate them through plots.
- Section 3 provides a succinct overview of related work.
- Section 4 contains formal statements of our main results.
- Finally, in Section 7 we present the proof of the main theorem. The proofs of its supporting propositions are presented in the appendices.

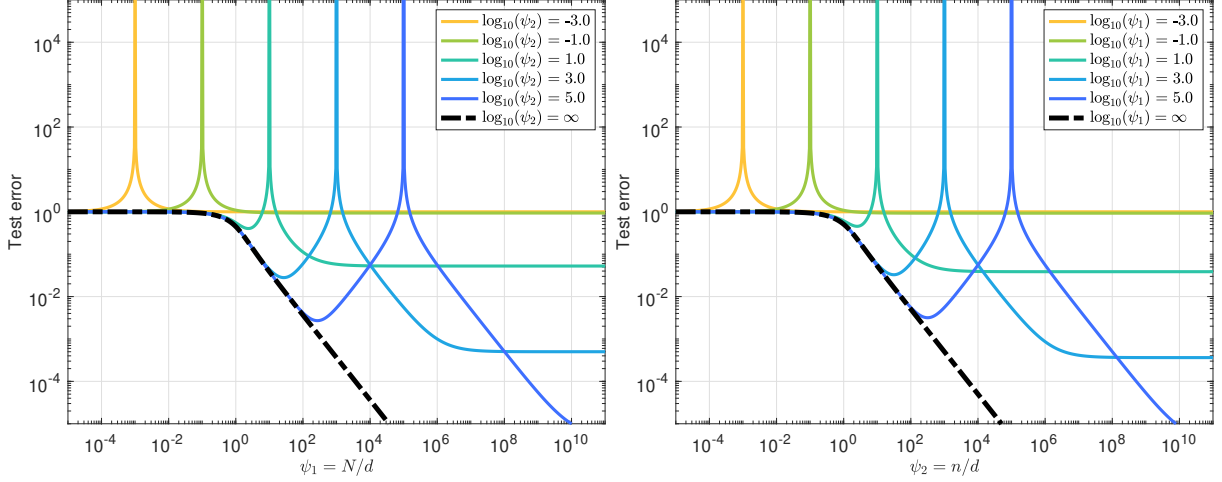


Figure 2: Analytical predictions for the test error of learning a linear function  $f_d(\mathbf{x}) = \langle \boldsymbol{\beta}_1, \mathbf{x} \rangle$  with  $\|\boldsymbol{\beta}_1\|_2^2 = 1$  using random features with ReLU activation function  $\sigma(x) = \max\{x, 0\}$ . Here we perform ridgeless regression ( $\lambda \rightarrow 0$ ). The signal-to-noise ratio is  $\|\boldsymbol{\beta}_1\|_2^2/\tau^2 \equiv \rho = 2$ . In the left figure, we plot the test error as a function of  $\psi_1 = N/d$ , and different curves correspond to different sample sizes ( $\psi_2 = n/d$ ). In the right figure, we plot the test error as a function of  $\psi_2 = n/d$ , and different curves correspond to different number of features ( $\psi_1 = N/d$ ).

## 2 Results and insights: An informal overview

Before explaining in detail our technical results –which we will do in Section 4– it is useful to pause and describe some consequences of the exact asymptotic formulae that we prove. Our focus here will be on insights that have a chance to hold more generally, beyond the specific setting studied here.

*Bias term also exhibits a singularity at the interpolation threshold.* A prominent feature of the double descent curve is the peak in test error at the interpolation threshold which, in the present case, is located at  $\psi_1 = \psi_2$ . In the linear regression model of [AS17, HMRT19, BHX19], this phenomenon is entirely explained by a peak (that diverges in the ridgeless limit  $\lambda \rightarrow 0$ ) in the variance of the estimator, while its bias is completely insensitive to this threshold.

In contrast, in the random features model studied here, both variance and bias have a peak at the interpolation threshold, diverging there when  $\lambda \rightarrow 0$ . This is apparent from Figure 1 which was obtained for  $\tau^2 = 0$ , and therefore in a setting in which the error is entirely due to bias. The fact that the double descent scenario persists in the noiseless limit is particularly important, especially in view of the fact that many machine learning tasks are usually considered nearly noiseless.

*Optimal prediction error is in the highly overparametrized regime.* Figure 2 (left) reports the predicted test error in the ridgeless limit  $\lambda \rightarrow 0$  (for a case with non-vanishing noise,  $\tau^2 > 0$ ) as a function of  $\psi_1 = N/d$ , for several values of  $\psi_2 = n/d$ . Figure 3 plots the predicted test error as a function of  $\psi_1$ , for fixed  $\psi_2$ , several values of  $\lambda > 0$ , and two values of the SNR. We repeatedly observe that: (i) For a fixed  $\lambda$ , the minimum of test error (over  $\psi_1$ ) is in the highly overparametrized regime  $\psi_1 \rightarrow \infty$ ; (ii) The global minimum (over  $\lambda$  and  $\psi_1$ ) of test error is achieved at a value of  $\lambda$  that depends on the SNR, but always at  $\psi_1 \rightarrow \infty$ ; (iii) In the ridgeless limit  $\lambda \rightarrow 0$ , the generalization curve is monotonically decreasing in  $\psi_1$  when  $\psi_1 > \psi_2$ .

To the best of our knowledge, this is the first natural and analytically tractable model which satisfies the following requirements: (1) Large overparametrization is necessary to achieve optimal prediction; (2) No special misspecification structure needs to be postulated.

*Non-vanishing regularization can hurt (at high SNR).* Figure 4 plots the predicted test error as a function of  $\lambda$ , for several values of  $\psi_1$ , with  $\psi_2$  fixed. The lower envelope of these curves is given by the curve at  $\psi_1 \rightarrow \infty$ , confirming that the optimal error is achieved in the highly overparametrized regime. However the dependence of this lower envelope on  $\lambda$  changes qualitatively, depending on the SNR. For small SNR, the



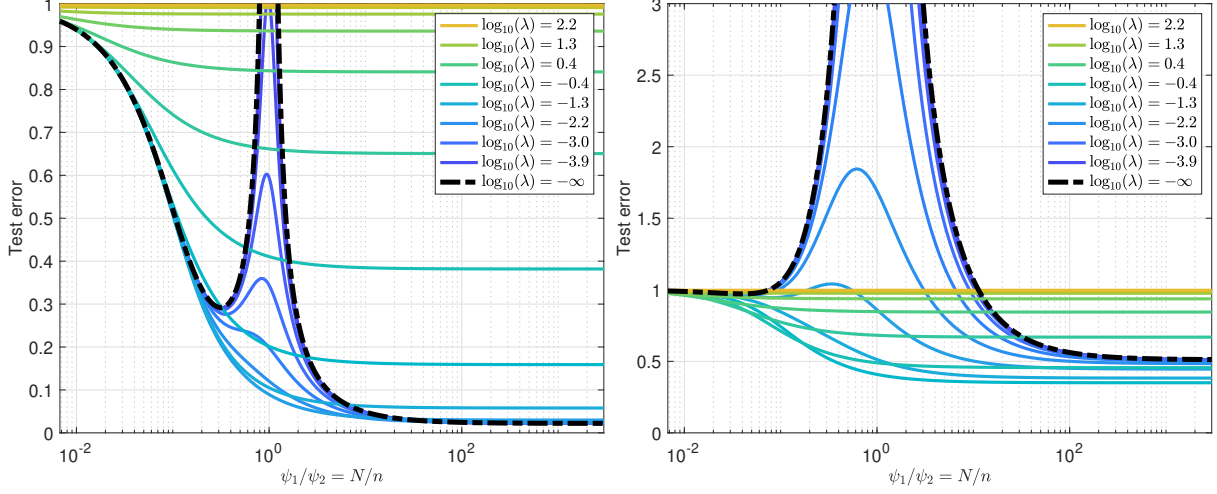


Figure 3: Analytical predictions for the test error of learning a linear function  $f_d(\mathbf{x}) = \langle \beta_1, \mathbf{x} \rangle$  with  $\|\beta_1\|_2^2 = 1$  using random features with ReLU activation function  $\sigma(x) = \max\{x, 0\}$ . The rescaled sample size is fixed to  $n/d \equiv \psi_2 = 10$ . Different curves are for different values of the regularization  $\lambda$ . On the left: high SNR  $\|\beta_1\|_2^2/\tau^2 \equiv \rho = 5$ . On the right: low SNR  $\rho = 1/5$ .

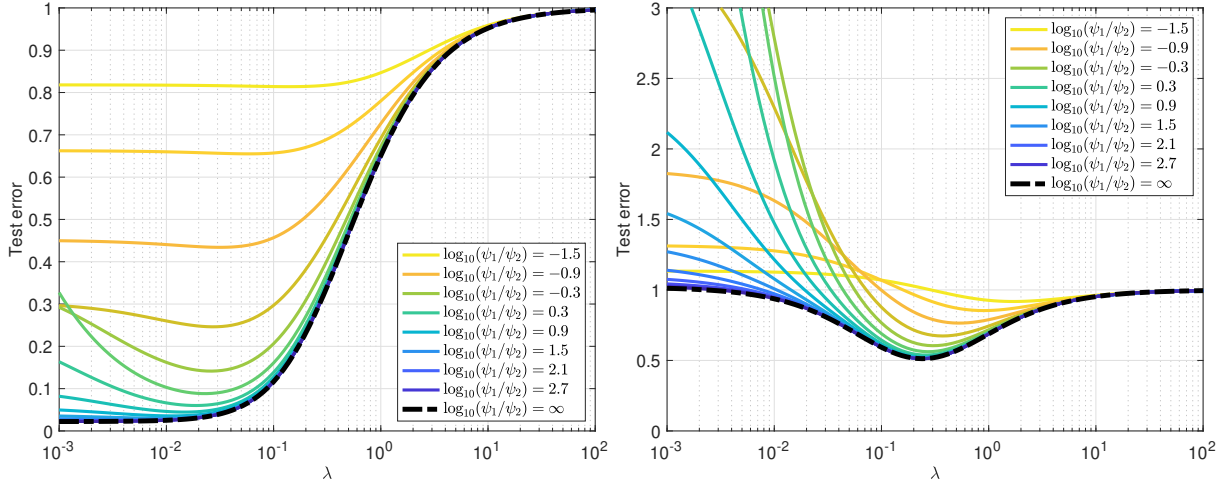


Figure 4: Analytical predictions for the test error of learning a linear function  $f_d(\mathbf{x}) = \langle \beta_1, \mathbf{x} \rangle$  with  $\|\beta_1\|_2^2 = 1$  using random features with ReLU activation function  $\sigma(x) = \max\{x, 0\}$ . The rescaled sample size is fixed to  $\psi_2 = n/d = 10$ . Different curves are for different values of the number of neurons  $\psi_1 = N/d$ . On the left: high SNR  $\|\beta_1\|_2^2/\tau^2 \equiv \rho = 5$ . On the right: low SNR  $\rho = 1/10$ .

global minimum is achieved as some  $\lambda > 0$ : regularization helps. However, for a large SNR the minimum error is achieved as  $\lambda \rightarrow 0$ . The optimal regularization is vanishingly small.

These two noise regime are separated by a phase transition at a critical SNR which we denote by  $\rho_\star$ . A characterization of this critical value is given in Section 4.2.2.

*Highly overparametrized interpolators are statistically optimal at high SNR.* This is a restatement of the last points. Notice that, in the overparametrized regime, the training error vanishes as  $\lambda \rightarrow 0$ , and the resulting model is a ‘near-interpolator’. (We cannot prove it is an exact interpolator because here we take  $\lambda \rightarrow 0$  after  $d \rightarrow \infty$ .) In the high-SNR regime of Figure 4, left frame, this strategy –namely, extreme overparametrization  $\psi_1 \rightarrow \infty$ , and interpolation limit  $\lambda \rightarrow 0$ – yields the globally minimum test error.

Following Remark 2, we expect the minimum- $\ell_2$  norm interpolator also to achieve asymptotically mini-

mum error.

### 3 Related literature

A recent stream of papers studied the generalization behavior of machine learning models in the interpolation regime. An incomplete list of references includes [BMM18, BRT18, LR18, BHMM18, RZ18]. The starting point of this line of work were the experimental results in [ZBH<sup>+</sup>16, BMM18], which showed that deep neural networks as well as kernel methods can generalize even if the prediction function interpolates all the data. It was proved that several machine learning models including kernel regression [BRT18] and kernel ridgeless regression [LR18] can generalize under certain conditions.

The double descent phenomenon, which is our focus in this paper, was first discussed in general terms in [BHMM18]. The same phenomenon was also empirically observed in [AS17, GJS<sup>+</sup>19]. The paper [KLS18] observes that the optimal amount of ridge regularization is sometimes vanishing, and provides an explanation in terms of noisy features. Analytical predictions confirming this scenario were obtained, within the linear regression model, in two concurrent papers [HMRT19, BHX19]. In particular, [HMRT19] derives the precise high-dimensional asymptotics of the prediction error, for a general model with correlated covariates. On the other hand, [BHX19] gives exact formula for any finite dimension, for a model with i.i.d. Gaussian covariates. The same papers also compute the double descent curve within other models, including over-specified linear model [HMRT19], and a Fourier series model [BHX19]. As mentioned in the introduction, the simple linear regression models of [HMRT19, BHX19] do not capture all the qualitative features of the double descent phenomenon in neural networks. In particular the observation that highly overparametrized models outperform other models trained in a more classical regime can only be recovered postulating specific misspecification models. The closest result to ours is the calculation of the variance term in the random features model in [HMRT19]. Bounds on the generalization error of overparametrized linear models were recently derived in [MVS19, BLLT19].

The random features model has been studied in considerable depth since the original work in [RR08]. A classical viewpoint suggests that  $\mathcal{F}_{\text{RF}}(\Theta)$  should be regarded as random approximation of the reproducing kernel Hilbert space  $\mathcal{F}_H$  defined by the kernel

$$H(\mathbf{x}, \mathbf{x}') = \mathbb{E}_{\boldsymbol{\theta} \sim \text{Unif}(\mathbb{S}^{d-1}(\sqrt{d}))} [\sigma(\langle \mathbf{x}, \boldsymbol{\theta} \rangle / \sqrt{d}) \sigma(\langle \mathbf{x}', \boldsymbol{\theta} \rangle / \sqrt{d})]. \quad (6)$$

Indeed  $\mathcal{F}_{\text{RF}}(\Theta)$  is an RKHS defined by the following finite-rank approximation of this kernel

$$H_N(\mathbf{x}, \mathbf{x}') = \frac{1}{N} \sum_{i=1}^N \sigma(\langle \mathbf{x}, \boldsymbol{\theta}_i \rangle / \sqrt{d}) \sigma(\langle \mathbf{x}', \boldsymbol{\theta}_i \rangle / \sqrt{d}). \quad (7)$$

The paper [RR08] showed the pointwise convergence of the empirical kernel  $H_N$  to  $H$ . Subsequent work [Bac17b] showed the convergence of the empirical kernel matrix to the population kernel in terms of operator norm and derived bound on the approximation error (see also [Bac13, AM15, RR17] for related work).

The setting in the present paper is quite different, since we take the limit of a large number neurons  $N \rightarrow \infty$ , together with large dimension  $d \rightarrow \infty$ . It is well-known that approximation using two-layers network suffers from the curse of dimensionality, in particular when first-layer weights are not trained [DHM89, Bac17a, VW18, GMMM19]. The recent paper [GMMM19] studies random features regression in a setting similar to ours, in the population limit  $n = \infty$ , but with  $N$  scaling as a general polynomial of  $d$ . It proves that, if  $N = O_d(d^{k+1-\delta})$  (for some  $\delta > 0$ ) then a random features model can only fit the projection of the true function  $f_d$  onto degree- $k$  polynomials. Here, we consider  $N, n = \Theta_d(d)$ , and therefore [GMMM19] implies that the prediction error of the random features model is lower bounded by the population risk achieved by the best linear predictor. The present results are of course much more precise (albeit limited to the proportional scaling) and indeed we observe that the nonlinear part of the function  $f_d$  effectively increases the noise level.

The relation between neural networks in the overparametrized regime and kernel methods has been studied in a number of recent papers. The connection between neural networks and random features models was pointed out originally in [Nea96, Wil97] and has attracted significant attention recently [HJ15, MRH<sup>+</sup>18, LBN<sup>+</sup>17, NXB<sup>+</sup>18, GAAR18]. The papers [DFS16, Dan17] showed that, for a certain initialization, gradient



descent training of overparametrized neural networks learns a function in an RKHS, which corresponds to the random features kernel. A recent line of work [JGH18, LL18, DZPS18, DLL<sup>+</sup>18, AZLS18, AZLL18, ADH<sup>+</sup>19, ZCZG18, OS19] studied the training dynamics of overparametrized neural networks under a second type of initialization, and showed that it learns a function in a different RKHS, which corresponds to the “neural tangent kernel”. A concurrent approach [MMN18, RVE18, CB18b, SS19, JMM19, Ngu19, RJBVE19, AOY19] studies the training dynamics of overparametrized neural networks under a third type of initialization, and showed that the dynamics of empirical distribution of weights follows Wasserstein gradient flow of a risk functional. The connection between neural tangent theory and Wasserstein gradient flow was studied in [CB18a, DL19, MMM19].

From a technical viewpoint, our analysis uses methods from random matrix theory. In particular, we use leave-one-out arguments to derive fixed point equations for the Stieltjes transform of certain spectral distributions. The general class of matrices we study are kernel inner product random matrices, namely matrices of the form  $\sigma(\mathbf{W}\mathbf{W}^\top/\sqrt{d})$ , where  $\mathbf{W}$  is a random matrix with i.i.d. entries, or similar. The paper [EK10] studied the spectrum of random kernel matrices when  $\sigma$  can be well approximated by a linear function and hence the spectrum converges to a scaled Marchenko-Pastur law. When  $\sigma$  cannot be approximated by a linear function, the spectrum of such matrices was studied in [CS13], and shown to converge to the free convolution of a Marchenko-Pastur law and a scaled semi-circular law. The extreme eigenvalue of the same random kernel matrix was studied in [FM19]. In the current paper, we need to consider an asymmetric kernel matrix  $\mathbf{Z} = \sigma(\mathbf{X}\mathbf{\Theta}^\top/\sqrt{d})/\sqrt{d}$ , whose asymptotic eigenvalue distribution was calculated in [PW17] (see also [LLC18] in the case when  $\mathbf{X}$  is deterministic).

The asymptotic spectral distribution is not sufficient to compute the asymptotic prediction error, which also depends on the eigenvectors of  $\mathbf{Z}$ . Our approach is to express the prediction error in terms of traces of products of  $\mathbf{Z}$  and other random matrices. We then express this traces as derivatives (with respect to certain auxiliary parameters) of the log-determinant of a certain block random matrix. We finally use the leave-one-out method to characterize the asymptotics of this log-determinant.

## 4 Main results

We begin by stating our assumptions and notations for the activation function  $\sigma$ . It is straightforward to check that these are satisfied by all commonly-used activations, including ReLU and sigmoid functions.

**Assumption 1.** *Let  $\sigma : \mathbb{R} \rightarrow \mathbb{R}$  be weakly differentiable, with weak derivative  $\sigma'$ . Assume  $|\sigma(u)|, |\sigma'(u)| \leq c_0 e^{c_1|u|}$  for some constants  $c_0, c_1 < \infty$ . Define*

$$\mu_0 \equiv \mathbb{E}\{\sigma(G)\}, \quad \mu_1 \equiv \mathbb{E}\{G\sigma(G)\}, \quad \mu_\star^2 \equiv \mathbb{E}\{\sigma(G)^2\} - \mu_0^2 - \mu_1^2, \quad (8)$$

where expectation is with respect to  $G \sim \mathcal{N}(0, 1)$ . Assuming  $0 < \mu_0^2, \mu_1^2, \mu_\star^2 < \infty$ , define  $\zeta$  by

$$\zeta \equiv \frac{\mu_1}{\mu_\star}. \quad (9)$$

We will consider sequences of parameters  $(N, n, d)$  that diverge proportionally to each other. When necessary, we can think such sequences to be indexed by  $d$ , with  $N = N(d)$ ,  $n = n(d)$  functions of  $d$ .

**Assumption 2.** *Defining  $\psi_{1,d} = N/d$  and  $\psi_{2,d} = n/d$ , we assume that the following limits exist in  $(0, \infty)$ :*

$$\lim_{d \rightarrow \infty} \psi_{1,d} = \psi_1, \quad \lim_{d \rightarrow \infty} \psi_{2,d} = \psi_2. \quad (10)$$

Our last assumption concerns the distribution of data  $(y, \mathbf{x})$ , and, in particular, the regression function  $f_d(\mathbf{x}) = \mathbb{E}[y|\mathbf{x}]$ . As stated in the introduction, we take  $f_d$  to be the sum of a deterministic linear component, and a nonlinear component that we assume to be random and isotropic.

**Assumption 3.** *We assume  $y_i = f_d(\mathbf{x}_i) + \varepsilon_i$ , where  $(\varepsilon_i)_{i \leq n} \sim_{iid} \mathbb{P}_\varepsilon$  independent of  $(\mathbf{x}_i)_{i \leq n}$ , with  $\mathbb{E}_\varepsilon(\varepsilon_1) = 0$ ,  $\mathbb{E}_\varepsilon(\varepsilon_1^2) = \tau^2$ ,  $\mathbb{E}_\varepsilon(\varepsilon_1^4) < \infty$ . Further*

$$f_d(\mathbf{x}) = \beta_{d,0} + \langle \beta_{d,1}, \mathbf{x} \rangle + f_d^{\text{NL}}(\mathbf{x}), \quad (11)$$

where  $\beta_{d,0} \in \mathbb{R}$  and  $\beta_{d,1} \in \mathbb{R}^d$  are deterministic with  $\lim_{d \rightarrow \infty} \beta_{d,0}^2 = F_0$ ,  $\lim_{d \rightarrow \infty} \|\beta_{d,1}\|_2^2 = F_1^2$ . The nonlinear component  $f_d^{\text{NL}}(\mathbf{x})$  is a centered Gaussian process indexed by  $\mathbf{x} \in \mathbb{S}^{d-1}(\sqrt{d})$ , with covariance

$$\mathbb{E}_{f_d^{\text{NL}}} \{f_d^{\text{NL}}(\mathbf{x}_1) f_d^{\text{NL}}(\mathbf{x}_2)\} = \Sigma_d(\langle \mathbf{x}_1, \mathbf{x}_2 \rangle / d) \quad (12)$$

satisfying  $\mathbb{E}_{\mathbf{x} \sim \text{Unif}(\mathbb{S}^{d-1}(\sqrt{d}))} \{\Sigma_d(x_1/\sqrt{d})\} = 0$ ,  $\mathbb{E}_{\mathbf{x} \sim \text{Unif}(\mathbb{S}^{d-1}(\sqrt{d}))} \{\Sigma_d(x_1/\sqrt{d})x_1\} = 0$ , and  $\lim_{d \rightarrow \infty} \Sigma_d(1) = F_\star^2$ . We define the signal-to-noise ratio parameter  $\rho$  by

$$\rho = \frac{F_1^2}{F_\star^2 + \tau^2}. \quad (13)$$

**Remark 3.** The last assumption covers, as a special case, deterministic linear functions  $f_d(\mathbf{x}) = \beta_{d,0} + \langle \beta_{d,1}, \mathbf{x} \rangle$ , but also a large class of random non-linear functions. As an example, let  $\mathbf{G} = (G_{ij})_{i,j \leq d}$ , where  $(G_{ij})_{i,j \leq d} \sim_{iid} \mathcal{N}(0, 1)$ , and consider the random quadratic function

$$f_d(x) = \beta_{d,0} + \langle \beta_{d,1}, \mathbf{x} \rangle + \frac{F_\star}{d} [\langle \mathbf{x}, \mathbf{G}\mathbf{x} \rangle - \text{Tr}(\mathbf{G})], \quad (14)$$

for some fixed  $F_\star \in \mathbb{R}$ . It is easy to check that this  $f_d$  satisfies Assumption 3, where the covariance function gives

$$\Sigma_d(\langle \mathbf{x}_1, \mathbf{x}_2 \rangle / d) = \frac{F_\star^2}{d^2} (\langle \mathbf{x}_1, \mathbf{x}_2 \rangle^2 - d).$$

Higher order polynomials can be constructed analogously (or using the expansion of  $f_d$  in spherical harmonics).

We also emphasize that the nonlinear part  $f_d^{\text{NL}}(\mathbf{x}_2)$ , although being random, is the same for all samples, and hence should not be confused with additive noise  $\varepsilon$ .

We finally introduce the formula for the asymptotic prediction error, denoted by  $\mathcal{R}(\rho, \zeta, \psi_1, \psi_2, \lambda)$  in Theorem 1.

**Definition 1** (Formula for the prediction error of random features regression). *Let the functions  $\nu_1, \nu_2 : \mathbb{C}_+ \rightarrow \mathbb{C}_+$  be uniquely defined by the following conditions: (i)  $\nu_1, \nu_2$  are analytic on  $\mathbb{C}_+$ ; (ii) For  $\Im(\xi) > 0$ ,  $\nu_1(\xi), \nu_2(\xi)$  satisfy the following equations*

$$\begin{aligned} \nu_1 &= \psi_1 \left( -\xi - \nu_2 - \frac{\zeta^2 \nu_2}{1 - \zeta^2 \nu_1 \nu_2} \right)^{-1}, \\ \nu_2 &= \psi_2 \left( -\xi - \nu_1 - \frac{\zeta^2 \nu_1}{1 - \zeta^2 \nu_1 \nu_2} \right)^{-1}; \end{aligned} \quad (15)$$

(iii)  $(\nu_1(\xi), \nu_2(\xi))$  is the unique solution of these equations with  $|\nu_1(\xi)| \leq \psi_1/\Im(\xi)$ ,  $|\nu_2(\xi)| \leq \psi_2/\Im(\xi)$  for  $\Im(\xi) > C$ , with  $C$  a sufficiently large constant.

Let

$$\chi \equiv \nu_1(\mathbf{i}(\psi_1 \psi_2 \bar{\lambda})^{1/2}) \cdot \nu_2(\mathbf{i}(\psi_1 \psi_2 \bar{\lambda})^{1/2}), \quad (16)$$

and

$$\begin{aligned} \mathcal{E}_0(\zeta, \psi_1, \psi_2, \bar{\lambda}) &\equiv -\chi^5 \zeta^6 + 3\chi^4 \zeta^4 + (\psi_1 \psi_2 - \psi_2 - \psi_1 + 1)\chi^3 \zeta^6 - 2\chi^3 \zeta^4 - 3\chi^3 \zeta^2 \\ &\quad + (\psi_1 + \psi_2 - 3\psi_1 \psi_2 + 1)\chi^2 \zeta^4 + 2\chi^2 \zeta^2 + \chi^2 + 3\psi_1 \psi_2 \chi \zeta^2 - \psi_1 \psi_2, \\ \mathcal{E}_1(\zeta, \psi_1, \psi_2, \bar{\lambda}) &\equiv \psi_2 \chi^3 \zeta^4 - \psi_2 \chi^2 \zeta^2 + \psi_1 \psi_2 \chi \zeta^2 - \psi_1 \psi_2, \\ \mathcal{E}_2(\zeta, \psi_1, \psi_2, \bar{\lambda}) &\equiv \chi^5 \zeta^6 - 3\chi^4 \zeta^4 + (\psi_1 - 1)\chi^3 \zeta^6 + 2\chi^3 \zeta^4 + 3\chi^3 \zeta^2 + (-\psi_1 - 1)\chi^2 \zeta^4 - 2\chi^2 \zeta^2 - \chi^2. \end{aligned} \quad (17)$$

We then define

$$\mathcal{B}(\zeta, \psi_1, \psi_2, \bar{\lambda}) \equiv \frac{\mathcal{E}_1(\zeta, \psi_1, \psi_2, \bar{\lambda})}{\mathcal{E}_0(\zeta, \psi_1, \psi_2, \bar{\lambda})}, \quad (18)$$

$$\mathcal{V}(\zeta, \psi_1, \psi_2, \bar{\lambda}) \equiv \frac{\mathcal{E}_2(\zeta, \psi_1, \psi_2, \bar{\lambda})}{\mathcal{E}_0(\zeta, \psi_1, \psi_2, \bar{\lambda})}, \quad (19)$$

$$\mathcal{R}(\rho, \zeta, \psi_1, \psi_2, \bar{\lambda}) \equiv \frac{\rho}{1+\rho} \mathcal{B}(\zeta, \psi_1, \psi_2, \bar{\lambda}) + \frac{1}{1+\rho} \mathcal{V}(\zeta, \psi_1, \psi_2, \bar{\lambda}). \quad (20)$$

The formula for the asymptotic risk can be easily evaluated numerically. In order to gain further insight, it can be simplified in some interesting special cases, as shown in Section 4.2.

#### 4.1 Statement of main result

We are now in position to state our main theorem, which generalizes Theorem 1 to the case in which  $f_d$  has a nonlinear component  $f_d^{\text{NL}}$ .

**Theorem 2.** *Let  $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_n)^\top \in \mathbb{R}^{n \times d}$  with  $(\mathbf{x}_i)_{i \in [n]} \sim_{iid} \text{Unif}(\mathbb{S}^{d-1}(\sqrt{d}))$  and  $\boldsymbol{\Theta} = (\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_N)^\top \in \mathbb{R}^{N \times d}$  with  $(\boldsymbol{\theta}_a)_{a \in [N]} \sim_{iid} \text{Unif}(\mathbb{S}^{d-1}(\sqrt{d}))$  independently. Let the activation function  $\sigma$  satisfy Assumption 1, and consider proportional asymptotics  $N/d \rightarrow \psi_1$ ,  $N/d \rightarrow \psi_2$ , as per Assumption 2. Finally, let the regression function  $\{f_d\}_{d \geq 1}$  and the response variables  $(y_i)_{i \in [n]}$  satisfy Assumption 3.*

*Then for any value of the regularization parameter  $\lambda > 0$ , the asymptotic prediction error of random features ridge regression satisfies*

$$\mathbb{E}_{\mathbf{X}, \boldsymbol{\Theta}, \boldsymbol{\varepsilon}, f_d^{\text{NL}}} \left| R_{\text{RF}}(f_d, \mathbf{X}, \boldsymbol{\Theta}, \lambda) - \left[ F_1^2 \mathcal{B}(\zeta, \psi_1, \psi_2, \lambda/\mu_\star^2) + (\tau^2 + F_\star^2) \mathcal{V}(\zeta, \psi_1, \psi_2, \lambda/\mu_\star^2) + F_\star^2 \right] \right| = o_d(1), \quad (21)$$

where  $\mathbb{E}_{\mathbf{X}, \boldsymbol{\Theta}, \boldsymbol{\varepsilon}, f_d^{\text{NL}}}$  denotes expectation with respect to data covariates  $\mathbf{X}$ , feature vectors  $\boldsymbol{\Theta}$ , data noise  $\boldsymbol{\varepsilon}$ , and  $f_d^{\text{NL}}$  the nonlinear part of the true regression function (as a Gaussian process), as per Assumption 3. The functions  $\mathcal{B}, \mathcal{V}$  are given in Definition 1.

**Remark 4.** If the regression function  $f_d(\mathbf{x})$  is linear (i.e.,  $f_d^{\text{NL}}(\mathbf{x}) = 0$ ), we recover Theorem 1, where  $\mathcal{R}$  is defined as per Eq. (20).

Numerical experiments suggest that Eq. (21) holds for any deterministic nonlinear functions  $f_d$  as well, and that the convergence in Eq. (21) is uniform over  $\lambda$  in compacts. We defer the study of these stronger properties to future work.

**Remark 5.** Note that the formula for a nonlinear truth, cf. Eq. (21), is almost identical to the one for a linear truth in Eq. (5). In fact, the only difference is that the prediction error increases by a term  $F_\star^2$ , and the noise level  $\tau^2$  is replaced by  $\tau^2 + F_\star^2$ .

Recall that the parameter  $F_\star^2$  is the variance of the nonlinear part  $\mathbb{E}_{f_d^{\text{NL}}}(f_d^{\text{NL}}(\mathbf{x})^2) \rightarrow F_\star^2$ . Hence, these changes can be interpreted by saying that random features regression (in  $N, n, d$  proportional regime) only estimates the linear component of  $f_d$  and the nonlinear component behaves similar to random noise. This finding is consistent with the results of [GMMM19] which imply, in particular,  $R_{\text{RF}}(f_d, \mathbf{X}, \boldsymbol{\Theta}, \lambda) \geq F_\star^2 + o_d(\mathbb{P}(1))$  for any  $n$  and for  $N = o_d(d^{2-\delta})$  for any  $\delta > 0$ .

Figure 5 illustrates the last remark. We report the simulated and predicted test error as a function of  $\psi_1/\psi_2 = N/n$ , for three different choices of the function  $f_d$  and noise level  $\tau^2$ . In all the settings, the total power of nonlinearity and noise is  $F_\star^2 + \tau^2 = 0.5$ , while the power of the linear component is  $F_1^2 = 1$ . The test errors in these three settings appear to be very close, as predicted by our theory.

**Remark 6.** The terms  $\mathcal{B}$  and  $\mathcal{V}$  in Eq. (21) correspond to the limits of the bias and variance of the estimated function  $f(\mathbf{x}; \hat{\mathbf{a}}(\lambda), \boldsymbol{\Theta})$ , when the ground truth function  $f_d$  is linear. That is, for  $f_d$  to be a linear function, we have

$$\mathbb{E}_{\mathbf{x}} \left\{ [f_d(\mathbf{x}) - \mathbb{E}_{\boldsymbol{\varepsilon}} f(\mathbf{x}; \hat{\mathbf{a}}(\lambda), \boldsymbol{\Theta})]^2 \right\} = \mathcal{B}(\zeta, \psi_1, \psi_2, \lambda/\mu_\star^2) F_1^2 + o_d(\mathbb{P}(1)), \quad (22)$$

$$\mathbb{E}_{\mathbf{x}} \text{Var}_{\boldsymbol{\varepsilon}}(f(\mathbf{x}; \hat{\mathbf{a}}(\lambda), \boldsymbol{\Theta})) = \mathcal{V}(\zeta, \psi_1, \psi_2, \lambda/\mu_\star^2) \tau^2 + o_d(\mathbb{P}(1)). \quad (23)$$

#### 4.2 Simplifying the asymptotic risk in special cases

In order to gain further insight into the formula for the asymptotic risk  $\mathcal{R}(\rho, \zeta, \psi_1, \psi_2, \bar{\lambda})$ , we consider here three special cases that are particularly interesting:

1. The ridgeless limit  $\lambda \rightarrow 0+$ .
2. The highly overparametrized regime  $\psi_1 \rightarrow \infty$  (recall that  $\psi_1 = \lim_{d \rightarrow \infty} N/d$ ).

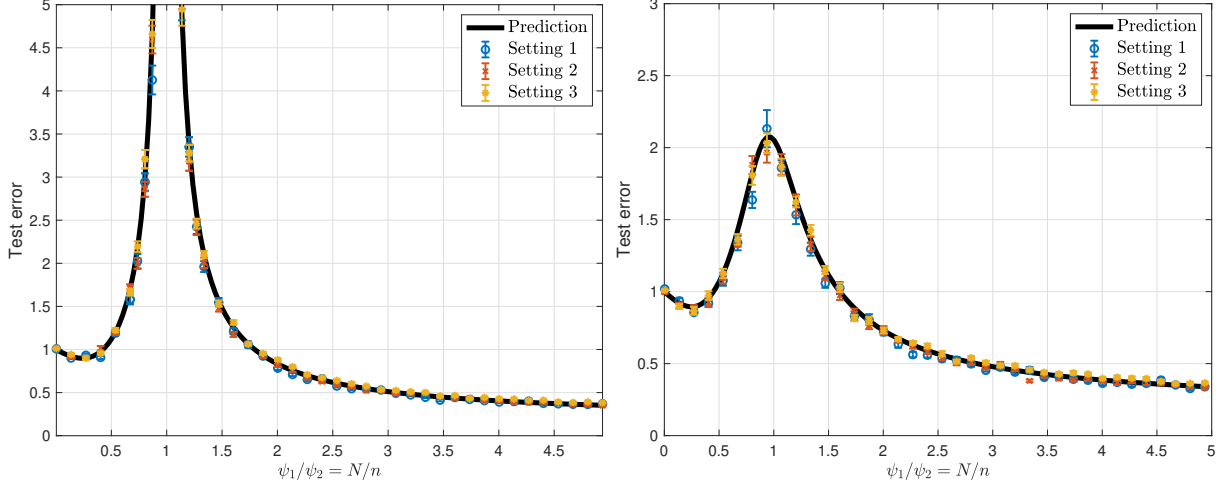


Figure 5: Random features regression with ReLU activation ( $\sigma = \max\{x, 0\}$ ). Data are generated according to one of three settings: (1)  $f_d(\mathbf{x}) = x_1$  and  $\mathbb{E}[\varepsilon^2] = 0.5$ ; (2)  $f_d(\mathbf{x}) = x_1 + (x_1^2 - 1)/2$  and  $\mathbb{E}[\varepsilon^2] = 0$ ; (3)  $f_d(\mathbf{x}) = x_1 + x_1 x_2 / \sqrt{2}$  and  $\mathbb{E}[\varepsilon^2] = 0$ . Within any of these settings, the total power of nonlinearity and noise is  $F_\star^2 + \tau^2 = 0.5$ , while the power of linear part is  $F_1^2 = 1$ . Left frame:  $\lambda = 10^{-8}$ . Right frame:  $\lambda = 10^{-3}$ . Here  $n = 300$ ,  $d = 100$ . The continuous black line is our theoretical prediction, and the colored symbols are numerical results. Symbols are averages over 20 instances and the error bars report the standard error of the means over these 20 instances.

3. The large sample limit  $\psi_2 \rightarrow \infty$  (recall that  $\psi_2 = \lim_{d \rightarrow \infty} n/d$ ).

Let us emphasize that these limits are taken *after* the limit  $N, n, d \rightarrow \infty$  with  $N/d \rightarrow \infty$  and  $n/d \rightarrow \infty$ . Hence, the correct interpretation of the highly overparametrized regime is not that the width  $N$  is infinite, but rather much larger than  $d$  (more precisely, larger than any constant times  $d$ ). Analogously, the large sample limit does not coincide with infinite sample size  $n$ , but instead sample size that is much larger than  $d$ .

#### 4.2.1 Ridgeless limit

The ridgeless limit  $\lambda \rightarrow 0+$  is important because it captures the asymptotic behavior the min-norm interpolation predictor (see also Remark 2.)

**Theorem 3.** *Under the assumptions of Theorem 2, set  $\psi \equiv \min\{\psi_1, \psi_2\}$  and define*

$$\chi \equiv -\frac{[(\psi\zeta^2 - \zeta^2 - 1)^2 + 4\zeta^2\psi]^{1/2} + (\psi\zeta^2 - \zeta^2 - 1)}{2\zeta^2}, \quad (24)$$

and

$$\begin{aligned} \mathcal{E}_{0,\text{rless}}(\zeta, \psi_1, \psi_2) &\equiv -\chi^5\zeta^6 + 3\chi^4\zeta^4 + (\psi_1\psi_2 - \psi_2 - \psi_1 + 1)\chi^3\zeta^6 - 2\chi^3\zeta^4 - 3\chi^3\zeta^2 \\ &\quad + (\psi_1 + \psi_2 - 3\psi_1\psi_2 + 1)\chi^2\zeta^4 + 2\chi^2\zeta^2 + \chi^2 + 3\psi_1\psi_2\chi\zeta^2 - \psi_1\psi_2, \\ \mathcal{E}_{1,\text{rless}}(\zeta, \psi_1, \psi_2) &\equiv \psi_2\chi^3\zeta^4 - \psi_2\chi^2\zeta^2 + \psi_1\psi_2\chi\zeta^2 - \psi_1\psi_2, \\ \mathcal{E}_{2,\text{rless}}(\zeta, \psi_1, \psi_2) &\equiv \chi^5\zeta^6 - 3\chi^4\zeta^4 + (\psi_1 - 1)\chi^3\zeta^6 + 2\chi^3\zeta^4 + 3\chi^3\zeta^2 + (-\psi_1 - 1)\chi^2\zeta^4 - 2\chi^2\zeta^2 - \chi^2, \end{aligned} \quad (25)$$

and

$$\mathcal{B}_{\text{rless}}(\zeta, \psi_1, \psi_2) \equiv \mathcal{E}_{1,\text{rless}}/\mathcal{E}_{0,\text{rless}}, \quad (26)$$

$$\mathcal{V}_{\text{rless}}(\zeta, \psi_1, \psi_2) \equiv \mathcal{E}_{2,\text{rless}}/\mathcal{E}_{0,\text{rless}}. \quad (27)$$

Then the asymptotic prediction error of random features ridgeless regression is given by

$$\lim_{\lambda \rightarrow 0} \lim_{d \rightarrow \infty} \mathbb{E}[R_{\text{RF}}(f_d, \mathbf{X}, \boldsymbol{\Theta}, \lambda)] = F_1^2 \mathcal{B}_{\text{rless}}(\zeta, \psi_1, \psi_2) + (\tau^2 + F_\star^2) \mathcal{V}_{\text{rless}}(\zeta, \psi_1, \psi_2) + F_\star^2. \quad (28)$$

The proof of this result can be found in Appendix F.

The next proposition establishes the main qualitative properties of the ridgeless limit.

**Proposition 4.1.** *Recall the bias and variance functions  $\mathcal{B}_{\text{rless}}$  and  $\mathcal{V}_{\text{rless}}$  defined in Eq. (26) and (27). Then, for any  $\zeta \in (0, \infty)$  and fixed  $\psi_2 \in (0, \infty)$ , we have*

1. *Small width limit  $\psi_1 \rightarrow 0$ :*

$$\lim_{\psi_1 \rightarrow 0} \mathcal{B}_{\text{rless}}(\zeta, \psi_1, \psi_2) = 1, \quad \lim_{\psi_1 \rightarrow 0} \mathcal{V}_{\text{rless}}(\zeta, \psi_1, \psi_2) = 0. \quad (29)$$

2. *Divergence at the interpolation threshold  $\psi_1 = \psi_2$ :*

$$\mathcal{B}_{\text{rless}}(\zeta, \psi_2, \psi_2) = \infty, \quad \mathcal{V}_{\text{rless}}(\zeta, \psi_2, \psi_2) = \infty. \quad (30)$$

3. *Large width limit  $\psi_1 \rightarrow \infty$  (here  $\chi$  is defined as per Eq. (24)):*

$$\lim_{\psi_1 \rightarrow \infty} \mathcal{B}_{\text{rless}}(\zeta, \psi_1, \psi_2) = (\psi_2 \chi \zeta^2 - \psi_2) / ((\psi_2 - 1) \chi^3 \zeta^6 + (1 - 3\psi_2) \chi^2 \zeta^4 + 3\psi_2 \chi \zeta^2 - \psi_2), \quad (31)$$

$$\lim_{\psi_1 \rightarrow \infty} \mathcal{V}_{\text{rless}}(\zeta, \psi_1, \psi_2) = (\chi^3 \zeta^6 - \chi^2 \zeta^4) / ((\psi_2 - 1) \chi^3 \zeta^6 + (1 - 3\psi_2) \chi^2 \zeta^4 + 3\psi_2 \chi \zeta^2 - \psi_2). \quad (32)$$

4. *Above the interpolation threshold (i.e. for  $\psi_1 \geq \psi_2$ ), the function  $\mathcal{B}_{\text{rless}}(\zeta, \psi_1, \psi_2)$  and  $\mathcal{V}_{\text{rless}}(\zeta, \psi_1, \psi_2)$  are strictly decreasing in the rescaled number of neurons  $\psi_1$ .*

The proof of this proposition is presented in Appendix G.1.

As anticipated, point 2 establishes an important difference with respect to the random covariates linear regression model of [AS17, HMRT19, BHX19]. While in those models the peak in prediction error is entirely due to a variance divergence, in the present setting both variance and bias diverge.

Another important difference is established in point 4: both bias and variance are monotonically decreasing above the interpolation threshold. This, again, contrasts with the behavior of simpler models, in which bias increases after the interpolation threshold, or after a somewhat larger point in the number of parameters per dimension (if misspecification is added).

This monotone decrease of the bias is crucial, and is at the origin of the observation that highly overparametrized models outperform underparametrized or moderately overparametrized ones. See Figure 6 for an illustration.

#### 4.2.2 Highly overparametrized regime

As the number of neurons  $N$  diverges (for fixed dimension  $d$ ), random features ridge regression is known to approach kernel ridge regression with respect to the kernel (6). It is therefore interesting what happens when  $N$  and  $d$  diverge together, but  $N$  is larger than any constant times  $d$ .

**Theorem 4.** *Under the assumptions of Theorem 2, define*

$$\omega \equiv - \frac{[(\psi_2 \zeta^2 - \zeta^2 - \bar{\lambda} \psi_2 - 1)^2 + 4\psi_2 \zeta^2 (\bar{\lambda} \psi_2 + 1)]^{1/2} + (\psi_2 \zeta^2 - \zeta^2 - \bar{\lambda} \psi_2 - 1)}{2(\bar{\lambda} \psi_2 + 1)}, \quad (33)$$

and

$$\mathcal{B}_{\text{wide}}(\zeta, \psi_2, \bar{\lambda}) = \frac{\psi_2 \omega - \psi_2}{(\psi_2 - 1) \omega^3 + (1 - 3\psi_2) \omega^2 + 3\psi_2 \omega - \psi_2}, \quad (34)$$

$$\mathcal{V}_{\text{wide}}(\zeta, \psi_2, \bar{\lambda}) = \frac{\omega^3 - \omega^2}{(\psi_2 - 1) \omega^3 + (1 - 3\psi_2) \omega^2 + 3\psi_2 \omega - \psi_2}. \quad (35)$$

Then the asymptotic prediction error of random features ridge regression, in the large width limit is given by

$$\lim_{\psi_1 \rightarrow \infty} \lim_{d \rightarrow \infty} \mathbb{E}[R_{\text{RF}}(f_d, \mathbf{X}, \boldsymbol{\Theta}, \lambda)] = F_1^2 \mathcal{B}_{\text{wide}}(\zeta, \psi_2, \lambda/\mu_\star^2) + (\tau^2 + F_\star^2) \mathcal{V}_{\text{wide}}(\zeta, \psi_2, \lambda/\mu_\star^2) + F_\star^2. \quad (36)$$

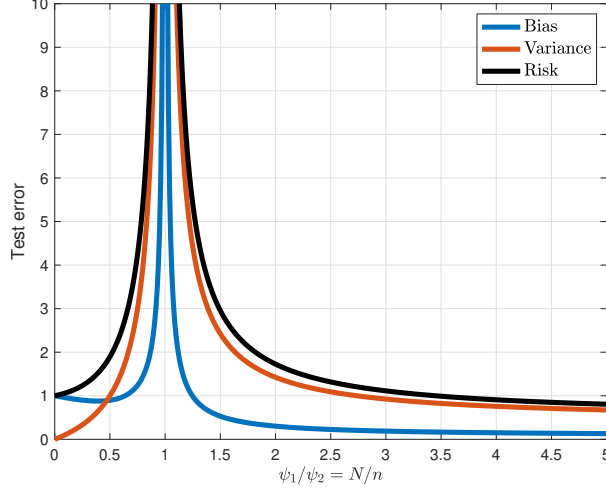


Figure 6: Analytical predictions of learning a linear function  $f_d(\mathbf{x}) = \langle \mathbf{x}, \boldsymbol{\beta}_1 \rangle$  with ReLU activation ( $\sigma = \max\{x, 0\}$ ) in the ridgeless limit ( $\lambda \rightarrow 0$ ). We take  $\|\boldsymbol{\beta}_1\|_2^2 = 1$  and  $\mathbb{E}[\varepsilon^2] = 1$ . We fix  $\psi_2 = 2$  and plot the bias, variance, and the test error as functions of  $\psi_1/\psi_2$ . Both the bias and the variance term diverge when  $\psi_1 = \psi_2$ , and decrease in  $\psi_1$  when  $\psi_1 > \psi_2$ .

The proof of this result can be found in Appendix F. Note that, as expected, the risk remains lower bounded by  $F_\star^2$ , even in the limit  $\psi_1 \rightarrow \infty$ . Naively, one could have expected to recover kernel ridge regression in this limit, and hence a method that can fit nonlinear functions. However, as shown in [GMMM19], random features methods can only learn linear functions for  $N = O_d(d^{2-\delta})$ .

As observed in Figures 2 to 4 (which have been obtained by applying Theorem 2), the minimum prediction error is often achieved by highly overparametrized networks  $\psi_1 \rightarrow \infty$ . It is natural to ask what is the effect of regularization on such networks. Somewhat surprisingly (and as anticipated in Section 2), we find that regularization does not always help. Namely, there exists a critical value  $\rho_\star$  of the signal-to-noise ratio, such that vanishing regularization is optimal for  $\rho > \rho_\star$ , and is not optimal for  $\rho < \rho_\star$ .

In order to state formally this result, we define the following quantities

$$\mathcal{R}_{\text{wide}}(\rho, \zeta, \psi_2, \bar{\lambda}) \equiv \frac{\rho}{1+\rho} \mathcal{B}_{\text{wide}}(\zeta, \psi_2, \bar{\lambda}) + \frac{1}{1+\rho} \mathcal{V}_{\text{wide}}(\zeta, \psi_2, \bar{\lambda}), \quad (37)$$

$$\omega_0(\zeta, \psi_2) \equiv - \frac{[(\psi_2 \zeta^2 - \zeta^2 - 1)^2 + 4\psi_2 \zeta^2]^{1/2} + (\psi_2 \zeta^2 - \zeta^2 - 1)}{2}, \quad (38)$$

$$\rho_\star(\zeta, \psi_2) \equiv \frac{\omega_0^2 - \omega_0}{(1 - \psi_2)\omega_0 + \psi_2}. \quad (39)$$

Notice in particular that  $\mathcal{R}_{\text{wide}}(\rho, \zeta, \psi_2, \lambda/\mu_\star^2)$  is the limiting value of the prediction error (right-hand side of (36)) up to an additive constant and an multiplicative constant.

**Proposition 4.2.** *Fix  $\zeta, \psi_2 \in (0, \infty)$  and  $\rho \in (0, \infty)$ . Then the function  $\bar{\lambda} \mapsto \mathcal{R}_{\text{wide}}(\rho, \zeta, \psi_2, \bar{\lambda})$  is either strictly increasing in  $\bar{\lambda}$ , or strictly decreasing first and then strictly increasing.*

Moreover, we have

$$\rho < \rho_\star(\zeta, \psi_2) \Rightarrow \arg \min_{\bar{\lambda} \geq 0} \mathcal{R}_{\text{wide}}(\rho, \zeta, \psi_2, \bar{\lambda}) = 0, \quad (40)$$

$$\rho > \rho_\star(\zeta, \psi_2) \Rightarrow \arg \min_{\bar{\lambda} \geq 0} \mathcal{R}_{\text{wide}}(\rho, \zeta, \psi_2, \bar{\lambda}) = \bar{\lambda}_\star(\zeta, \psi_2, \rho) > 0. \quad (41)$$

The proof of this proposition is presented in Appendix G.2, which also provides further information about this phase transition (and, in particular, an explicit expression for  $\bar{\lambda}_\star(\zeta, \psi_2, \rho)$ ).



### 4.2.3 Large sample limit

As the number of sample  $n$  goes to infinity, both training error (minus  $\tau^2$ ) and test error<sup>1</sup> converge to the approximation error using random features class to fit the true function  $f_d$ . It is therefore interesting what happens when  $n$  and  $d$  diverge together, but  $n$  is larger than any constant times  $d$ .

**Theorem 5.** *Under the assumptions of Theorem 2, define*

$$\omega \equiv -\frac{[(\psi_1\zeta^2 - \zeta^2 - \bar{\lambda}\psi_1 - 1)^2 + 4\psi_1\zeta^2(\bar{\lambda}\psi_1 + 1)]^{1/2} + (\psi_1\zeta^2 - \zeta^2 - \bar{\lambda}\psi_1 - 1)}{2(\bar{\lambda}\psi_1 + 1)}, \quad (42)$$

and

$$\mathcal{B}_{\text{lsamp}}(\zeta, \psi_1, \bar{\lambda}) = \frac{(\omega^3 - \omega^2)/\zeta^2 + \psi_1\omega - \psi_1}{(\psi_1 - 1)\omega^3 + (1 - 3\psi_1)\omega^2 + 3\psi_1\omega - \psi_1}. \quad (43)$$

Then the asymptotic prediction error of random features ridge regression, in the large width limit is given by

$$\lim_{\psi_2 \rightarrow \infty} \lim_{d \rightarrow \infty} \mathbb{E}[R_{\text{RF}}(f_d, \mathbf{X}, \boldsymbol{\Theta}, \lambda)] = F_1^2 \mathcal{B}_{\text{lsamp}}(\zeta, \psi_2, \lambda/\mu_\star^2) + F_\star^2. \quad (44)$$

The proof of this result can be found in Appendix F.

## 5 Asymptotics of the training error

Theorem 2 establishes the exact asymptotics of the test error in the random features model. However, the technical results obtained in the appendices allow us to characterize several other quantities of interest. Here we consider the behavior of the training error and of the norm of the parameters. We define the regularized training error by

$$L_{\text{RF}}(f_d, \mathbf{X}, \boldsymbol{\Theta}, \lambda) = \min_{\mathbf{a}} \left\{ \frac{1}{n} \sum_{i=1}^n \left( y_i - \sum_{j=1}^N a_j \sigma(\langle \boldsymbol{\theta}_j, \mathbf{x}_i \rangle / \sqrt{d}) \right)^2 + \frac{N\lambda}{d} \|\mathbf{a}\|_2^2 \right\}. \quad (45)$$

We also recall that  $\hat{\mathbf{a}}(\lambda)$  denotes the minimizer in the last expression, cf. Eq. (2). The next definition presents the asymptotic formulas for these quantities.

**Definition 2** (Asymptotic formula for training error of random features regression). *Let the functions  $\nu_1, \nu_2 : \mathbb{C}_+ \rightarrow \mathbb{C}_+$  be uniquely defined by the following conditions: (i)  $\nu_1, \nu_2$  are analytic on  $\mathbb{C}_+$ ; (ii) For  $\Im(\xi) > 0$ ,  $\nu_1(\xi), \nu_2(\xi)$  satisfy the following equations*

$$\begin{aligned} \nu_1 &= \psi_1 \left( -\xi - \nu_2 - \frac{\zeta^2 \nu_2}{1 - \zeta^2 \nu_1 \nu_2} \right)^{-1}, \\ \nu_2 &= \psi_2 \left( -\xi - \nu_1 - \frac{\zeta^2 \nu_1}{1 - \zeta^2 \nu_1 \nu_2} \right)^{-1}; \end{aligned} \quad (46)$$

(iii)  $(\nu_1(\xi), \nu_2(\xi))$  is the unique solution of these equations with  $|\nu_1(\xi)| \leq \psi_1/\Im(\xi)$ ,  $|\nu_2(\xi)| \leq \psi_2/\Im(\xi)$  for  $\Im(\xi) > C$ , with  $C$  a sufficiently large constant.

Let

$$\chi \equiv \nu_1(\mathbf{i}(\psi_1\psi_2\bar{\lambda})^{1/2}) \cdot \nu_2(\mathbf{i}(\psi_1\psi_2\bar{\lambda})^{1/2}), \quad (47)$$

<sup>1</sup>The difference between training error and test error is due to the fact that we define the former as  $\hat{\mathbb{E}}_n\{(y - \hat{f}(\mathbf{x}))^2\}$  and the latter as  $\mathbb{E}\{(f(\mathbf{x}) - \hat{f}(\mathbf{x}))^2\}$ .

and

$$\begin{aligned}
\mathcal{L} &= -i\nu_2(i(\psi_1\psi_2\bar{\lambda})^{1/2}) \cdot \left(\frac{\bar{\lambda}\psi_1}{\psi_2}\right)^{1/2} \cdot \left[\frac{\rho}{1+\rho} \cdot \frac{1}{1-\chi\zeta^2} + \frac{1}{1+\rho}\right], \\
\mathcal{A}_1 &= \frac{\rho}{1+\rho} \left[ -\chi^2(\chi\zeta^4 - \chi\zeta^2 + \psi_2\zeta^2 + \zeta^2 - \chi\psi_2\zeta^4 + 1) \right] \\
&\quad + \frac{1}{1+\rho} \left[ \chi^2(\chi\zeta^2 - 1)(\chi^2\zeta^4 - 2\chi\zeta^2 + \zeta^2 + 1) \right], \\
\mathcal{A}_0 &= -\chi^5\zeta^6 + 3\chi^4\zeta^4 + (\psi_1\psi_2 - \psi_2 - \psi_1 + 1)\chi^3\zeta^6 - 2\chi^3\zeta^4 - 3\chi^3\zeta^2 \\
&\quad + (\psi_1 + \psi_2 - 3\psi_1\psi_2 + 1)\chi^2\zeta^4 + 2\chi^2\zeta^2 + \chi^2 + 3\psi_1\psi_2\chi\zeta^2 - \psi_1\psi_2, \\
\mathcal{A} &= \mathcal{A}_1/\mathcal{A}_0.
\end{aligned} \tag{48}$$

We next state our asymptotic characterization of  $L_{\text{RF}}(f_d, \mathbf{X}, \boldsymbol{\Theta}, \lambda)$  and  $\|\hat{\mathbf{a}}(\lambda)\|_2^2$ .

**Theorem 6.** *Let  $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_n)^\top \in \mathbb{R}^{n \times d}$  with  $(\mathbf{x}_i)_{i \in [n]} \sim_{i.i.d.} \text{Unif}(\mathbb{S}^{d-1}(\sqrt{d}))$  and  $\boldsymbol{\Theta} = (\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_N)^\top \in \mathbb{R}^{N \times d}$  with  $(\boldsymbol{\theta}_a)_{a \in [N]} \sim_{i.i.d.} \text{Unif}(\mathbb{S}^{d-1}(\sqrt{d}))$  independently. Let the activation function  $\sigma$  satisfy Assumption 1, and consider proportional asymptotics  $N/d \rightarrow \psi_1$ ,  $N/d \rightarrow \psi_2$ , as per Assumption 2. Finally, let the regression function  $\{f_d\}_{d \geq 1}$  and the response variables  $(y_i)_{i \in [n]}$  satisfy Assumption 3.*

*Then for any value of the regularization parameter  $\lambda > 0$ , the asymptotic regularized training error and norm square of its minimizer satisfy*

$$\begin{aligned}
\mathbb{E}_{\mathbf{X}, \boldsymbol{\Theta}, \boldsymbol{\varepsilon}, f_d^{\text{NL}}} \left| L_{\text{RF}}(f_d, \mathbf{X}, \boldsymbol{\Theta}, \lambda) - (F_1^2 + F_\star^2 + \tau^2)\mathcal{L} \right| &= o_d(1), \\
\mathbb{E}_{\mathbf{X}, \boldsymbol{\Theta}, \boldsymbol{\varepsilon}, f_d^{\text{NL}}} \left| \mu_\star^2 \|\hat{\mathbf{a}}(\lambda)\|_2^2 - (F_1^2 + F_\star^2 + \tau^2)\mathcal{A} \right| &= o_d(1),
\end{aligned} \tag{49}$$

where  $\mathbb{E}_{\mathbf{X}, \boldsymbol{\Theta}, \boldsymbol{\varepsilon}, f_d^{\text{NL}}}$  denotes expectation with respect to data covariates  $\mathbf{X}$ , feature vectors  $\boldsymbol{\Theta}$ , data noise  $\boldsymbol{\varepsilon}$ , and  $f_d^{\text{NL}}$  the nonlinear part of the true regression function (as a Gaussian process), as per Assumption 3. The functions  $\mathcal{L}$  and  $\mathcal{A}$  are given in Definition 2.

## 5.1 Numerical illustrations

In this section, we illustrate Theorem 6 through numerical simulations. Figure 5.1 reports the theoretical prediction and numerical results for the regularized training error, the test error, and the norm of the coefficients  $\hat{\mathbf{a}}(\lambda)$ . We use a small non-zero value of the regularization parameter  $\lambda = 10^{-3}$ , fix the number of samples per dimension  $\psi_2 = n/d$ , and follow these quantities as a function of the overparameterization ratio  $\psi_1/\psi_2 = N/n$ .

As expected, the behavior of the training error strikingly different from the one of the test error. The training error is monotone decreasing in the overparameterization ratio  $N/n$ , and is close to zero in the overparameterized regime  $N/n > 1$  (it is not exactly vanishing because we use a small  $\lambda > 0$ ). In other words, the fitted model is nearly interpolating the data, and the peak in test error matches the interpolation threshold.

On the other hand, the penalty term  $\psi_1 \|\hat{\mathbf{a}}(\lambda)\|_2^2$  is non-monotone: it increases up to the interpolation threshold, then decreases for  $N/n > 1$ , and converges to a constant as  $\psi_1 \rightarrow \infty$ . If we take this as a proxy for the model complexity, the behavior of  $\psi_1 \|\hat{\mathbf{a}}(\lambda)\|_2^2$  provides useful intuition about descent of the generalization error. As the number of parameters increases beyond the interpolation threshold, the model complexity decreases instead of increasing.

We can confirm the intuition that the double descent of the test error is driven by the behavior of the model complexity  $\psi_1 \|\hat{\mathbf{a}}(\lambda)\|_2^2$ , by selecting  $\lambda$  in an optimal way. Following [HMRT19], we expect that the optimal regularization should produce a smaller value of  $\psi_1 \|\hat{\mathbf{a}}(\lambda)\|_2^2$ , and hence eliminate or reduce the double descent phenomenon. Indeed, this is illustrated in Figure 5.1 demonstrates the prediction of the regularized training error and the test error for two choices of  $\lambda$ :  $\lambda = 0$ , and an optimal  $\lambda$  such that the test error is minimized. When we choose an optimal  $\lambda$ , the test error becomes strictly decreasing as  $\psi_1 = N/d$  increases. We expect this is a generic phenomenon that also holds in other interesting models.

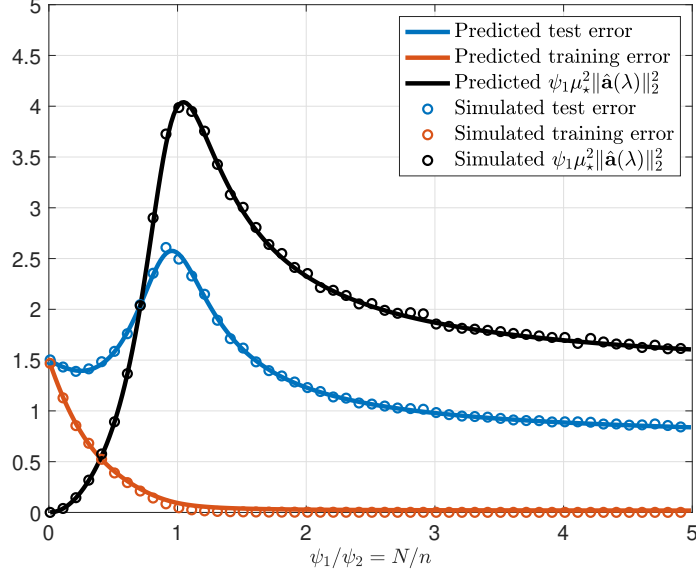


Figure 7: Analytical predictions and numerical simulations for the test error and regularized training error. Data are generated according to  $y_i = \langle \beta_1, \mathbf{x}_i \rangle + \varepsilon_i$  with  $\|\beta_1\|_2^2 = 1$  and  $\varepsilon_i \sim \mathcal{N}(0, \tau^2)$ ,  $\tau^2 = 0.5$ . We fit a random features model with ReLU activations ( $\sigma(x) = \max\{x, 0\}$ ) and ridge regularization parameter  $\lambda = 10^{-3}$ . In simulations we use  $d = 100$  and  $n = 300$ . We add  $\tau^2 = 0.5$  to the test error to make it comparable with training error. Symbols are averages over 20 instances.

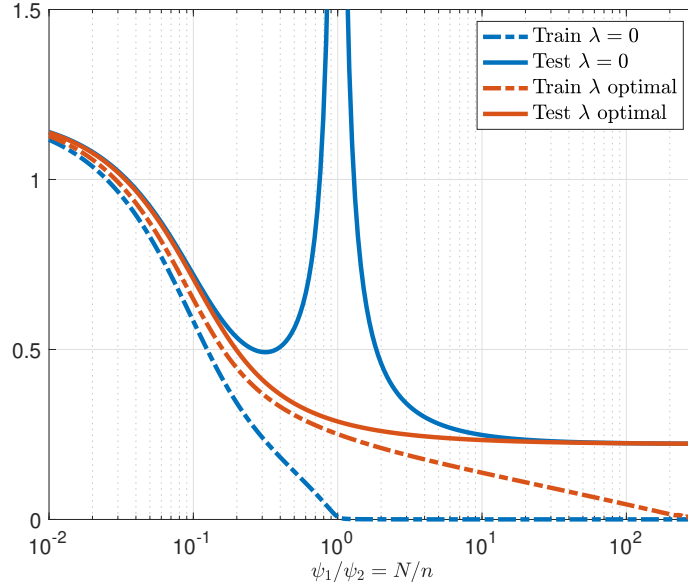


Figure 8: Analytical predictions and numerical simulations results for the test error and the regularized training error. Data are generated according to  $y_i = \langle \beta_1, \mathbf{x}_i \rangle + \varepsilon_i$  with  $\|\beta_1\|_2^2 = 1$  and  $\varepsilon_i \sim \mathcal{N}(0, \tau^2)$ ,  $\tau^2 = 0.2$ . We fit a random features model with ReLU activations ( $\sigma(x) = \max\{x, 0\}$ ). We fix  $\psi_2 = n/d = 10$ . We add  $\tau^2 = 0.2$  to the test error make it comparable with training error. In the optimal ridge setting, we choose  $\lambda$  for each value of  $\psi_1$  as to minimize the asymptotic test error.

## 6 An equivalent Gaussian covariates model

An exam of the proof of our main result, Theorem 2 reveals an interesting mathematical phenomenon. The random features model has the same asymptotic prediction error as a simpler model with Gaussian covariates and response that is linear in these covariates, provided we use a special covariance and signal structure.

The construction of the Gaussian covariates model proceeds as follows. Fix  $\beta_1 \in \mathbb{R}^d$ ,  $\|\beta_1\|_2^2 = F_1$  and  $\Theta = (\theta_1, \dots, \theta_N)^\top$  with  $(\theta_j)_{j \in [N]} \sim_{iid} \text{Unif}(\mathbb{S}^{d-1}(\sqrt{d}))$ . The joint distribution of  $(y, \mathbf{x}, \mathbf{u}) \in \mathbb{R} \times \mathbb{R}^d \times \mathbb{R}^N$  conditional on  $\Theta$  is defined by the following procedure:

1. Draw  $\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d)$ ,  $\varepsilon \sim \mathcal{N}(\mathbf{0}, \tau^2)$ , and  $\mathbf{w} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_N)$  independently, conditional on  $\Theta$ .
2. Let  $y = \langle \beta_1, \mathbf{x} \rangle + \varepsilon$ .
3. Let  $\mathbf{u} = (u_1, \dots, u_N)^\top$ ,  $u_j = \mu_0 + \mu_1 \langle \theta_j, \mathbf{x} \rangle / \sqrt{d} + \mu_\star w_j$ , for some  $0 < |\mu_0|, |\mu_1|, |\mu_\star| < \infty$ .

We will denote by  $\mathbb{P}_{y, \mathbf{x}, \mathbf{u} | \Theta}$  the probability distribution thus defined. As anticipated, this is a Gaussian covariates model. Indeed, the covariates vector  $\mathbf{u} \sim \mathcal{N}(\mathbf{0}, \Sigma)$  is Gaussian, with covariance  $\Sigma = \mu_0^2 \mathbf{1}\mathbf{1}^\top + \mu_1^2 \Theta \Theta^\top / d + \mu_\star^2 \mathbf{I}_N$ . Also  $(y, \mathbf{u})$  are jointly Gaussian and we can therefore write  $y = \langle \tilde{\beta}_1, \mathbf{u} \rangle + \tilde{\varepsilon}$ , for some new vector of coefficients  $\tilde{\beta}_1$ , and noise  $\tilde{\varepsilon}$  which is independent of  $\mathbf{u}$ .

Let  $\{(y_i, \mathbf{x}_i, \mathbf{u}_i)\}_{i \in [n]} | \Theta \sim_{iid} \mathbb{P}_{y, \mathbf{x}, \mathbf{u} | \Theta}$ . We learn a regression function  $\hat{f}(\mathbf{x}; \mathbf{a}, \Theta) = \langle \mathbf{u}, \mathbf{a} \rangle$ , by performing ridge regression

$$\hat{\mathbf{a}}(\lambda) = \arg \min_{\mathbf{a} \in \mathbb{R}^N} \left\{ \frac{1}{n} \sum_{i=1}^n (y_i - \langle \mathbf{u}_i, \mathbf{a} \rangle)^2 + \frac{N\lambda}{d} \|\mathbf{a}\|_2^2 \right\}. \quad (50)$$

The prediction error is defined by

$$R_{GC}(f_d, \mathbf{X}, \Theta, \lambda) = \mathbb{E}_{\mathbf{x}, \mathbf{z} | \Theta} [(f_d(\mathbf{x}) - \langle \mathbf{u}, \hat{\mathbf{a}}(\lambda) \rangle)^2]. \quad (51)$$

Remarkably, in the proportional asymptotics  $N, n, d \rightarrow \infty$  with  $N/d \rightarrow \psi_1, n/d \rightarrow \psi_2$ , the behavior of this model is the same as the one of the nonlinear random features model studied in the rest of the paper. In particular, the asymptotic prediction error  $\mathcal{R}$  is given by the same formula as in Definition 1.

**Theorem 7.** (Gaussian covariates prediction model) Define  $\zeta$  and the signal-to-noise ratio  $\rho \in [0, \infty]$  as

$$\zeta \equiv \mu_1^2 / \mu_\star^2, \quad \rho \equiv F_1^2 / \tau^2, \quad (52)$$

and assume  $\mu_0, \mu_1, \mu_\star \neq 0$ . Then, in the Gaussian covariates model described above, for any  $\lambda > 0$ , we have

$$R_{GC}(f_d, \mathbf{X}, \Theta, \lambda) = (F_1^2 + \tau^2) \mathcal{R}(\rho, \zeta, \psi_1, \psi_2, \lambda / \mu_\star^2) + o_{d, \mathbb{P}}(1), \quad (53)$$

where  $\mathcal{R}(\rho, \zeta, \psi_1, \psi_2, \bar{\lambda})$  is explicitly given in Definition 1.

The proof of Theorem 7 is almost the same as the one of Theorem 2 (with several simplifications, because of the greater amount of independence). To avoid repetitions, we will not present a proof here.

Figure 9 illustrates the content of Theorem 7 via numerical simulations. We report the simulated and predicted test error as a function of  $\psi_1 / \psi_2 = N / n$ . The theoretical prediction here is exactly the same as the one reported in Figure 5. However, numerical simulations were carried out with the Gaussian covariates model instead of random features. The agreement is excellent, as predicted by Theorem 7.

Why do the RF and GC models result in the same asymptotic prediction error? It is useful to provide a heuristic explanation of this interesting phenomenon. Consider an activation function  $\sigma : \mathbb{R} \rightarrow \mathbb{R}$ , with  $\mu_k = \mathbb{E}[\text{He}_k(G)\sigma(G)]$  and  $\mu_\star^2 = \mathbb{E}[\sigma^2(G)] - \mu_0^2 - \mu_1^2$  for  $G \sim \mathcal{N}(0, 1)$ . Define the nonlinear component of the activation function by  $\sigma^\perp(x) \equiv \sigma(x) - b_0 - b_1 x$ . Note that we have

$$\sigma(\langle \mathbf{x}_i, \theta_j \rangle / \sqrt{d}) = \mu_0 + \mu_1 \langle \mathbf{x}_i, \theta_j \rangle / \sqrt{d} + \mu_\star \tilde{w}_{ij}, \quad \tilde{w}_{ij} \equiv \frac{1}{\mu_\star} \sigma^\perp(\langle \mathbf{x}_i, \theta_j \rangle / \sqrt{d}),$$

$$u_j = \mu_0 + \mu_1 \langle \mathbf{x}_i, \theta_j \rangle / \sqrt{d} + \mu_\star w_{ij},$$

where  $(w_{ij})_{i \in [n], j \in [N]} \sim_{iid} \mathcal{N}(0, 1)$  independent of  $\mathbf{X}$  and  $\Theta$ . Note that the first two moments of  $\tilde{w}_{ij}$  match those of  $w_{ij}$ , i.e.  $\mathbb{E}_{\mathbf{x} | \Theta} \tilde{w}_{ij} = 0$ ,  $\mathbb{E}_{\mathbf{x} | \Theta} (\tilde{w}_{ij}^2) = 1$ . Further, for  $i \neq l$ ,  $\tilde{w}_{ij}$ ,  $\tilde{w}_{il}$  are nearly uncorrelated:  $\mathbb{E}_{\mathbf{x} | \Theta} \{\tilde{w}_{ij} \tilde{w}_{il}\} = O(\langle \theta_j, \theta_l \rangle / d) = O_{\mathbb{P}}(1/d)$ . It is therefore not unreasonable to imagine that they should behave as independents. The same intuition also appears in the analysis of the spectrum of kernel random matrices in [CS13, PW17].

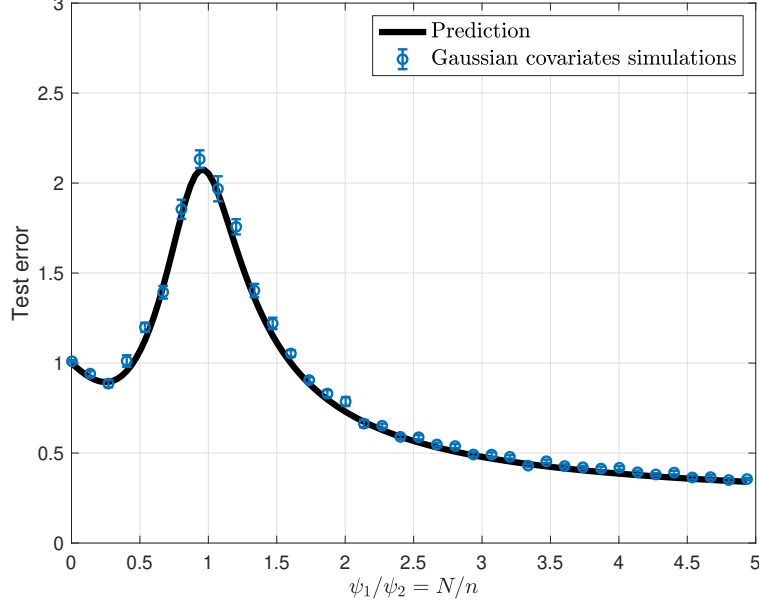


Figure 9: Predictions and numerical simulations for the test error of the Gaussian covariates model. We fit  $y_i = \langle \beta_1, \mathbf{x}_i \rangle + \varepsilon_i$  with  $\|\beta_1\|_2^2 = 1$  and  $\tau^2 = \mathbb{E}[\varepsilon_i^2] = 0.5$ , and parameters  $\mu_1 = 0.5$ ,  $\mu_\star = \sqrt{(\pi - 2)/(4\pi)}$ , and  $\lambda = 10^{-3}$ . This choice of parameters  $\mu_1$  and  $\mu_\star$  matches the corresponding parameters for ReLU activations. Here  $n = 300$ ,  $d = 100$ . The continuous black line is our theoretical prediction, and the colored symbols are numerical results. Symbols are averages over 20 instances and the error bars report the standard error of the means over 20 instances.

## 7 Proof of Theorem 2

As mentioned in the introduction, the proof of Theorem 2 relies on the following main steps. First we reduce the computation of the prediction error (in the high-dimensional limit  $N, n, d \rightarrow \infty$ ) to computing traces of products of certain kernel random matrices  $\mathbf{Q}, \mathbf{H}, \mathbf{Z}, \mathbf{Z}_1$  and their inverses (Step 1 below). Next we show that these traces can be obtained by taking derivatives of the log-determinant of a block-structured matrix  $\mathbf{A}$ , whose blocks are formed by  $\mathbf{Q}, \mathbf{H}, \mathbf{Z}, \mathbf{Z}_1$  (Step 3 below). Then we compute the the Stieltjes transform of  $\mathbf{A}$  and use it to characterize the asymptotics of the log-determinant (Step 2 below). Finally we simplify the formula of the limiting log-determinant and use it to derive the formula for the limiting risk function (Step 4 below).

This section presents a complete proof of Theorem 2, making use of technical propositions that formalize each of the steps above. The proofs of these propositions are deferred to the appendices.

### Step 1. Decompose the risk

The proposition below expresses the prediction risk in terms of  $\Psi_1, \Psi_2, \Psi_3$  which are traces of products of random matrices as defined in the proposition.

**Proposition 7.1** (Decomposition). *Let  $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_n)^\top \in \mathbb{R}^{n \times d}$  with  $(\mathbf{x}_i)_{i \in [n]} \sim_{iid} \text{Unif}(\mathbb{S}^{d-1}(\sqrt{d}))$ . Let  $\boldsymbol{\Theta} = (\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_n)^\top \in \mathbb{R}^{N \times d}$  with  $(\boldsymbol{\theta}_a)_{a \in [N]} \sim_{iid} \text{Unif}(\mathbb{S}^{d-1}(\sqrt{d}))$  to be independent of  $\mathbf{X}$ . Let  $\{f_d\}_{d \geq 1}$  and  $(y_i)_{i \in [n]}$  satisfy Assumption 3. Let activation function  $\sigma$  satisfy Assumption 1. Let  $N, n$ , and  $d$  satisfy Assumption 2. Then, for any  $\lambda > 0$ , we have*

$$\mathbb{E}_{\mathbf{X}, \boldsymbol{\Theta}, \varepsilon, f_d^{\text{NL}}} \left| R_{\text{RF}}(f_d, \mathbf{X}, \boldsymbol{\Theta}, \lambda) - \left[ F_1^2(1 - 2\Psi_1 + \Psi_2) + (F_\star^2 + \tau^2)\Psi_3 + F_\star^2 \right] \right| = o_d(1), \quad (54)$$

where

$$\begin{aligned}\Psi_1 &= \frac{1}{d} \text{Tr} \left[ \mathbf{Z}_1^\top \mathbf{Z} (\mathbf{Z}^\top \mathbf{Z} + \psi_1 \psi_2 \lambda \mathbf{I}_N)^{-1} \right], \\ \Psi_2 &= \frac{1}{d} \text{Tr} \left[ (\mathbf{Z}^\top \mathbf{Z} + \psi_1 \psi_2 \lambda \mathbf{I}_N)^{-1} (\mu_1^2 \mathbf{Q} + \mu_\star^2 \mathbf{I}_N) (\mathbf{Z}^\top \mathbf{Z} + \psi_1 \psi_2 \lambda \mathbf{I}_N)^{-1} \mathbf{Z}^\top \mathbf{H} \mathbf{Z} \right], \\ \Psi_3 &= \frac{1}{d} \text{Tr} \left[ (\mathbf{Z}^\top \mathbf{Z} + \psi_1 \psi_2 \lambda \mathbf{I}_N)^{-1} (\mu_1^2 \mathbf{Q} + \mu_\star^2 \mathbf{I}_N) (\mathbf{Z}^\top \mathbf{Z} + \psi_1 \psi_2 \lambda \mathbf{I}_N)^{-1} \mathbf{Z}^\top \mathbf{Z} \right],\end{aligned}\tag{55}$$

and

$$\begin{aligned}\mathbf{Q} &= \frac{1}{d} \boldsymbol{\Theta} \boldsymbol{\Theta}^\top, \\ \mathbf{H} &= \frac{1}{d} \mathbf{X} \mathbf{X}^\top, \\ \mathbf{Z} &= \frac{1}{\sqrt{d}} \sigma \left( \frac{1}{\sqrt{d}} \mathbf{X} \boldsymbol{\Theta}^\top \right), \\ \mathbf{Z}_1 &= \frac{\mu_1}{d} \mathbf{X} \boldsymbol{\Theta}^\top.\end{aligned}\tag{56}$$

**Step 2. The Stieltjes transform of a random block matrix.**

To calculate the quantities  $\Psi_1, \Psi_2, \Psi_3$ , first we study the Stieltjes transform of a block random matrix. Let  $\mathbf{Q}, \mathbf{H}, \mathbf{Z}, \mathbf{Z}_1$  be the matrices defined by Eq. (56). Define  $\mathbf{q} = (s_1, s_2, t_1, t_2, p) \in \mathbb{R}^5$ , and introduce a block matrix  $\mathbf{A} \in \mathbb{R}^{M \times M}$  with  $M = N + n$ , defined by

$$\mathbf{A} = \begin{bmatrix} s_1 \mathbf{I}_N + s_2 \mathbf{Q} & \mathbf{Z}^\top + p \mathbf{Z}_1^\top \\ \mathbf{Z} + p \mathbf{Z}_1 & t_1 \mathbf{I}_n + t_2 \mathbf{H} \end{bmatrix}.\tag{57}$$

We consider a set  $\mathcal{Q}$ , defined by

$$\mathcal{Q} = \{(s_1, s_2, t_1, t_2, p) : |s_2 t_2| \leq \mu_1^2 (1 + p)^2 / 2\}.\tag{58}$$

It is easy to see that  $\mathbf{0} \in \mathcal{Q}$ . We will restrict ourself to study the case  $\mathbf{q} \in \mathcal{Q}$ .

We consider sequences of matrices  $\mathbf{A}$  with  $n, N, d \rightarrow \infty$ . To be definite, we index elements of such sequences by the dimension  $d$ , and it is understood that  $\mathbf{A} = \mathbf{A}(d)$ ,  $n = n(d)$ ,  $N = N(d)$  depend on the dimension. We would like to calculate the asymptotic behavior of the Stieltjes transform

$$m_d(\xi; \mathbf{q}) = \mathbb{E}[M_d(\xi; \mathbf{q})],$$

where

$$M_d(\xi; \mathbf{q}) = \frac{1}{d} \text{Tr}[(\mathbf{A} - \xi \mathbf{I}_M)^{-1}].\tag{59}$$

We define the following function  $F(\cdot, \cdot; \xi; \mathbf{q}, \psi_1, \psi_2, \mu_1, \mu_\star) : \mathbb{C} \times \mathbb{C} \rightarrow \mathbb{C}$ :

$$F(m_1, m_2; \xi; \mathbf{q}, \psi_1, \psi_2, \mu_1, \mu_\star) \equiv \psi_1 \left( -\xi + s_1 - \mu_\star^2 m_2 + \frac{(1 + t_2 m_2) s_2 - \mu_1^2 (1 + p)^2 m_2}{(1 + s_2 m_1)(1 + t_2 m_2) - \mu_1^2 (1 + p)^2 m_1 m_2} \right)^{-1}.\tag{60}$$

**Proposition 7.2** (Stieltjes transform). *Let  $\sigma$  be an activation function satisfying Assumption 1. Consider the linear regime of Assumption 2. Consider a fixed  $\mathbf{q} \in \mathcal{Q}$ . Let  $m_1(\cdot; \mathbf{q}), m_2(\cdot; \mathbf{q}) : \mathbb{C}_+ \rightarrow \mathbb{C}_+$  be defined, for  $\Im(\xi) \geq C$  a sufficiently large constant, as the unique solution of the equations*

$$m_1 = F(m_1, m_2; \xi; \mathbf{q}, \psi_1, \psi_2, \mu_1, \mu_\star), \quad m_2 = F(m_2, m_1; \xi; \mathbf{q}, \psi_2, \psi_1, \mu_1, \mu_\star),\tag{61}$$

*subject to the condition  $|m_1| \leq \psi_1 / \Im(\xi)$ ,  $|m_2| \leq \psi_2 / \Im(\xi)$ . Extend this definition to  $\Im(\xi) > 0$  by requiring  $m_1, m_2$  to be analytic functions in  $\mathbb{C}_+$ . Define  $m(\xi; \mathbf{q}) = m_1(\xi; \mathbf{q}) + m_2(\xi; \mathbf{q})$ . Then for any  $\xi \in \mathbb{C}_+$  with  $\Im \xi > 0$ , we have*

$$\lim_{d \rightarrow \infty} \mathbb{E}[|M_d(\xi; \mathbf{q}) - m(\xi; \mathbf{q})|] = 0.\tag{62}$$

Further, for any compact set  $\Omega \subseteq \mathbb{C}_+$ , we have

$$\lim_{d \rightarrow \infty} \mathbb{E} \left[ \sup_{\xi \in \Omega} |M_d(\xi; \mathbf{q}) - m(\xi; \mathbf{q})| \right] = 0.\tag{63}$$



**Step 3. Compute  $\Psi_1, \Psi_2, \Psi_3$ .**

Recall the random matrix  $\mathbf{A} = \mathbf{A}(\mathbf{q})$  defined by Eq. (57). Let  $\text{Log}$  denote the complex logarithm with branch cut on the negative real axis. Let  $\{\lambda_i(\mathbf{A})\}_{i \in [M]}$  be the set of eigenvalues of  $\mathbf{A}$  in non-increasing order. For any  $\xi \in \mathbb{C}_+$ , we consider the quantity

$$G_d(\xi; \mathbf{q}) = \frac{1}{d} \sum_{i=1}^M \text{Log}(\lambda_i(\mathbf{A}(\mathbf{q})) - \xi).$$

Recall the definition of  $M_d(\xi; \mathbf{q})$  given in Eq. (59).

**Proposition 7.3.** *For  $\xi \in \mathbb{C}_+$  and  $\mathbf{q} \in \mathbb{R}^5$ , we have*

$$\frac{d}{d\xi} G_d(\xi; \mathbf{q}) = -\frac{1}{d} \sum_{i=1}^M (\lambda_i(\mathbf{A}) - \xi)^{-1} = -M_d(\xi; \mathbf{q}). \quad (64)$$

Moreover, for  $u \in \mathbb{R}$ , we have

$$\begin{aligned} \partial_p G_d(iu; \mathbf{0}) &= \frac{2}{d} \text{Tr} \left( (u^2 \mathbf{I}_N + \mathbf{Z}^\top \mathbf{Z})^{-1} \mathbf{Z}_1^\top \mathbf{Z} \right), \\ \partial_{s_1, t_1}^2 G_d(iu; \mathbf{0}) &= -\frac{1}{d} \text{Tr} \left( (u^2 \mathbf{I}_N + \mathbf{Z}^\top \mathbf{Z})^{-2} \mathbf{Z}^\top \mathbf{Z} \right), \\ \partial_{s_1, t_2}^2 G_d(iu; \mathbf{0}) &= -\frac{1}{d} \text{Tr} \left( (u^2 \mathbf{I}_N + \mathbf{Z}^\top \mathbf{Z})^{-2} \mathbf{Z}^\top \mathbf{H} \mathbf{Z} \right), \\ \partial_{s_2, t_1}^2 G_d(iu; \mathbf{0}) &= -\frac{1}{d} \text{Tr} \left( (u^2 \mathbf{I}_N + \mathbf{Z}^\top \mathbf{Z})^{-1} \mathbf{Q} (u^2 \mathbf{I}_N + \mathbf{Z}^\top \mathbf{Z})^{-1} \mathbf{Z}^\top \mathbf{Z} \right), \\ \partial_{s_2, t_2}^2 G_d(iu; \mathbf{0}) &= -\frac{1}{d} \text{Tr} \left( (u^2 \mathbf{I}_N + \mathbf{Z}^\top \mathbf{Z})^{-1} \mathbf{Q} (u^2 \mathbf{I}_N + \mathbf{Z}^\top \mathbf{Z})^{-1} \mathbf{Z}^\top \mathbf{H} \mathbf{Z} \right). \end{aligned} \quad (65)$$

**Proposition 7.4.** *Define*

$$\begin{aligned} \Xi(\xi, z_1, z_2; \mathbf{q}) &\equiv \log[(s_2 z_1 + 1)(t_2 z_2 + 1) - \mu_1^2(1+p)^2 z_1 z_2] - \mu_\star^2 z_1 z_2 \\ &\quad + s_1 z_1 + t_1 z_2 - \psi_1 \log(z_1/\psi_1) - \psi_2 \log(z_2/\psi_2) - \xi(z_1 + z_2) - \psi_1 - \psi_2. \end{aligned} \quad (66)$$

For  $\xi \in \mathbb{C}_+$  and  $\mathbf{q} \in \mathcal{Q}$  (c.f. Eq. (58)), let  $m_1(\xi; \mathbf{q}), m_2(\xi; \mathbf{q})$  be defined as the analytic continuation of solution of Eq. (61) as defined in Proposition 7.2. Define

$$g(\xi; \mathbf{q}) = \Xi(\xi, m_1(\xi; \mathbf{q}), m_2(\xi; \mathbf{q}); \mathbf{q}). \quad (67)$$

Consider proportional asymptotics  $N/d \rightarrow \psi_1$ ,  $N/d \rightarrow \psi_2$ , as per Assumption 2. Then for any fixed  $\xi \in \mathbb{C}_+$  and  $\mathbf{q} \in \mathcal{Q}$ , we have

$$\lim_{d \rightarrow \infty} \mathbb{E}[|G_d(\xi; \mathbf{q}) - g(\xi; \mathbf{q})|] = 0. \quad (68)$$

Moreover, for any fixed  $u \in \mathbb{R}_+$ , we have

$$\lim_{d \rightarrow \infty} \mathbb{E}[\|\partial_{\mathbf{q}} G_d(iu; \mathbf{0}) - \partial_{\mathbf{q}} g(iu; \mathbf{0})\|_2] = 0, \quad (69)$$

$$\lim_{d \rightarrow \infty} \mathbb{E}[\|\nabla_{\mathbf{q}}^2 G_d(iu; \mathbf{0}) - \nabla_{\mathbf{q}}^2 g(iu; \mathbf{0})\|_{\text{op}}] = 0. \quad (70)$$

By Eq. (55) and Eq. (65), we get

$$\begin{aligned} \Psi_1 &= \frac{1}{2} \partial_p G_d(\mathbf{i}(\psi_1 \psi_2 \lambda)^{1/2}; \mathbf{0}), \\ \Psi_2 &= -\mu_\star^2 \partial_{s_1, t_2} G_d(\mathbf{i}(\psi_1 \psi_2 \lambda)^{1/2}; \mathbf{0}) - \mu_1^2 \partial_{s_2, t_2} G_d(\mathbf{i}(\psi_1 \psi_2 \lambda)^{1/2}; \mathbf{0}), \\ \Psi_3 &= -\mu_\star^2 \partial_{s_1, t_1} G_d(\mathbf{i}(\psi_1 \psi_2 \lambda)^{1/2}; \mathbf{0}) - \mu_1^2 \partial_{s_2, t_1} G_d(\mathbf{i}(\psi_1 \psi_2 \lambda)^{1/2}; \mathbf{0}). \end{aligned}$$

Then by Eq. (69) and Eq. (70), we get

$$\begin{aligned}\mathbb{E}_{\mathbf{X}, \Theta} \left| \Psi_1 - \frac{1}{2} \partial_p g(\mathbf{i}(\psi_1 \psi_2 \lambda)^{1/2}; \mathbf{0}) \right| &= o_d(1), \\ \mathbb{E}_{\mathbf{X}, \Theta} \left| \Psi_2 + \left[ \mu_\star^2 \partial_{s_1, t_2} g(\mathbf{i}(\psi_1 \psi_2 \lambda)^{1/2}; \mathbf{0}) + \mu_1^2 \partial_{s_2, t_2} g(\mathbf{i}(\psi_1 \psi_2 \lambda)^{1/2}; \mathbf{0}) \right] \right| &= o_d(1), \\ \mathbb{E}_{\mathbf{X}, \Theta} \left| \Psi_3 + \left[ \mu_\star^2 \partial_{s_1, t_1} g(\mathbf{i}(\psi_1 \psi_2 \lambda)^{1/2}; \mathbf{0}) + \mu_1^2 \partial_{s_2, t_1} g(\mathbf{i}(\psi_1 \psi_2 \lambda)^{1/2}; \mathbf{0}) \right] \right| &= o_d(1).\end{aligned}$$

By Proposition 7.1, we get

$$\mathbb{E}_{\mathbf{X}, \Theta, \epsilon, f_d^{\text{NL}}} \left| R_{\text{RF}}(f_d, \mathbf{X}, \Theta, \lambda) - \overline{\mathcal{R}} \right| = o_d(1), \quad (71)$$

where

$$\overline{\mathcal{R}} = F_1^2 \mathcal{B} + (F_\star^2 + \tau^2) \mathcal{V} + F_\star^2. \quad (72)$$

$$\mathcal{B} = 1 - \partial_p g(\mathbf{i}(\psi_1 \psi_2 \lambda)^{1/2}; \mathbf{0}) - \mu_\star^2 \partial_{s_1, t_2} g(\mathbf{i}(\psi_1 \psi_2 \lambda)^{1/2}; \mathbf{0}) - \mu_1^2 \partial_{s_2, t_2} g(\mathbf{i}(\psi_1 \psi_2 \lambda)^{1/2}; \mathbf{0}), \quad (73)$$

$$\mathcal{V} = -\mu_\star^2 \partial_{s_1, t_1} g(\mathbf{i}(\psi_1 \psi_2 \lambda)^{1/2}; \mathbf{0}) - \mu_1^2 \partial_{s_2, t_1} g(\mathbf{i}(\psi_1 \psi_2 \lambda)^{1/2}; \mathbf{0}). \quad (74)$$

#### Step 4. Calculate explicitly $\mathcal{B}$ and $\mathcal{V}$ .

Next, we calculate derivatives of  $g(\xi; \mathbf{q})$  to give a more explicit expression for  $\mathcal{B}$  and  $\mathcal{V}$ .

**Lemma 7.1** (Formula for derivatives of  $g$ ). *For fixed  $\xi \in \mathbb{C}_+$  and  $\mathbf{q} \in \mathbb{R}^5$ , let  $m_1(\xi; \mathbf{q}), m_2(\xi; \mathbf{q})$  be defined as the analytic continuation of solution of Eq. (61) as defined in Proposition 7.2. Recall the definition of  $\Xi$  and  $g$  given in Eq. (66) and (67), i.e.,*

$$\begin{aligned}\Xi(\xi, z_1, z_2; \mathbf{q}) \equiv & \log[(s_2 z_1 + 1)(t_2 z_2 + 1) - \mu_1^2(1+p)^2 z_1 z_2] - \mu_\star^2 z_1 z_2 \\ & + s_1 z_1 + t_1 z_2 - \psi_1 \log(z_1/\psi_1) - \psi_2 \log(z_2/\psi_2) - \xi(z_1 + z_2) - \psi_1 - \psi_2.\end{aligned} \quad (75)$$

and

$$g(\xi; \mathbf{q}) = \Xi(\xi, m_1(\xi; \mathbf{q}), m_2(\xi; \mathbf{q}); \mathbf{q}). \quad (76)$$

Denoting

$$m_0 = m_0(\xi) \equiv m_1(\xi; \mathbf{0}) \cdot m_2(\xi; \mathbf{0}), \quad (77)$$

we have

$$\begin{aligned}\partial_p g(\xi; \mathbf{0}) &= 2m_0 \mu_1^2 / (m_0 \mu_1^2 - 1), \\ \partial_{s_1, t_1}^2 g(\xi; \mathbf{0}) &= [m_0^5 \mu_1^6 \mu_\star^2 - 3m_0^4 \mu_1^4 \mu_\star^2 + m_0^3 \mu_1^4 + 3m_0^3 \mu_1^2 \mu_\star^2 - m_0^2 \mu_1^2 - m_0^2 \mu_\star^2] / S, \\ \partial_{s_1, t_2}^2 g(\xi; \mathbf{0}) &= [(\psi_2 - 1)m_0^3 \mu_1^4 + m_0^3 \mu_1^2 \mu_\star^2 + (-\psi_2 - 1)m_0^2 \mu_1^2 - m_0^2 \mu_\star^2] / S, \\ \partial_{s_2, t_1}^2 g(\xi; \mathbf{0}) &= [(\psi_1 - 1)m_0^3 \mu_1^4 + m_0^3 \mu_1^2 \mu_\star^2 + (-\psi_1 - 1)m_0^2 \mu_1^2 - m_0^2 \mu_\star^2] / S, \\ \partial_{s_2, t_2}^2 g(\xi; \mathbf{0}) &= [-m_0^6 \mu_1^6 \mu_\star^4 + 2m_0^5 \mu_1^4 \mu_\star^4 + (\psi_1 \psi_2 - \psi_2 - \psi_1 + 1)m_0^4 \mu_1^6 \\ &\quad - m_0^4 \mu_1^4 \mu_\star^2 - m_0^4 \mu_1^2 \mu_\star^4 + (2 - 2\psi_1 \psi_2)m_0^3 \mu_1^4 \\ &\quad + (\psi_1 + \psi_2 + \psi_1 \psi_2 + 1)m_0^2 \mu_1^2 + m_0^2 \mu_\star^2] / [(m_0 \mu_1^2 - 1)S],\end{aligned} \quad (78)$$

where

$$\begin{aligned}S = & m_0^5 \mu_1^6 \mu_\star^4 - 3m_0^4 \mu_1^4 \mu_\star^4 + (\psi_1 + \psi_2 - \psi_1 \psi_2 - 1)m_0^3 \mu_1^6 \\ & + 2m_0^3 \mu_1^4 \mu_\star^2 + 3m_0^3 \mu_1^2 \mu_\star^4 + (3\psi_1 \psi_2 - \psi_2 - \psi_1 - 1)m_0^2 \mu_1^4 \\ & - 2m_0^2 \mu_1^2 \mu_\star^2 - m_0^2 \mu_\star^4 - 3\psi_1 \psi_2 m_0 \mu_1^2 + \psi_1 \psi_2.\end{aligned} \quad (79)$$

*Proof of Lemma 7.1.* For fixed  $\xi \in \mathbb{C}_+$  and  $\mathbf{q} \in \mathbb{R}^5$ , by the fixed point equation satisfied by  $m_1, m_2$  (c.f. Eq. (61)), we see that  $(m_1(\xi; \mathbf{q}), m_2(\xi; \mathbf{q}))$  is a stationary point of function  $\Xi(\xi, \cdot, \cdot; \mathbf{q})$ . Using the formula for implicit differentiation, we have

$$\begin{aligned}\partial_p g(\xi; \mathbf{q}) &= \partial_p \Xi(\xi, z_1, z_2; \mathbf{q})|_{(z_1, z_2) = (m_1(\xi; \mathbf{q}), m_2(\xi; \mathbf{q}))}, \\ \partial_{s_1, t_1}^2 g(\xi; \mathbf{q}) &= \mathbf{H}_{1,3} - \mathbf{H}_{1,[5,6]} \mathbf{H}_{[5,6],[5,6]}^{-1} \mathbf{H}_{[5,6],3}, \\ \partial_{s_1, t_2}^2 g(\xi; \mathbf{q}) &= \mathbf{H}_{1,4} - \mathbf{H}_{1,[5,6]} \mathbf{H}_{[5,6],[5,6]}^{-1} \mathbf{H}_{[5,6],4}, \\ \partial_{s_2, t_1}^2 g(\xi; \mathbf{q}) &= \mathbf{H}_{2,3} - \mathbf{H}_{2,[5,6]} \mathbf{H}_{[5,6],[5,6]}^{-1} \mathbf{H}_{[5,6],3}, \\ \partial_{s_2, t_2}^2 g(\xi; \mathbf{q}) &= \mathbf{H}_{2,4} - \mathbf{H}_{2,[5,6]} \mathbf{H}_{[5,6],[5,6]}^{-1} \mathbf{H}_{[5,6],4},\end{aligned}$$

where we have, for  $\mathbf{u} = (s_1, s_2, t_1, t_2, z_1, z_2)^\top$

$$\mathbf{H} = \nabla_{\mathbf{u}}^2 \Xi(\xi, z_1, z_2; \mathbf{q})|_{(z_1, z_2) = (m_1(\xi; \mathbf{q}), m_2(\xi; \mathbf{q}))}.$$

Basic algebra completes the proof.  $\square$

Define

$$\begin{aligned} \nu_1(\mathbf{i}\xi) &\equiv m_1(\mathbf{i}\xi\mu_\star; \mathbf{0}) \cdot \mu_\star, \\ \nu_2(\mathbf{i}\xi) &\equiv m_2(\mathbf{i}\xi\mu_\star; \mathbf{0}) \cdot \mu_\star. \end{aligned} \tag{80}$$

By the definition of analytic functions  $m_1$  and  $m_2$  (satisfying Eq. (61) and (60) with  $\mathbf{q} = \mathbf{0}$  as defined in Proposition 7.2), the definition of  $\nu_1$  and  $\nu_2$  in Eq. (80) above is equivalent to its definition in Definition 1 (as per Eq. (15)). Moreover, for  $\chi$  defined in Eq. (16) with  $\bar{\lambda} = \lambda/\mu_\star^2$  and  $m_0$  defined in Eq. (77), we have

$$\begin{aligned} \chi &= \nu_1(\mathbf{i}(\psi_1\psi_2\lambda/\mu_\star^2)^{1/2})\nu_2(\mathbf{i}(\psi_1\psi_2\lambda/\mu_\star^2)^{1/2}) \\ &= m_1(\mathbf{i}(\psi_1\psi_2\lambda)^{1/2}; \mathbf{0})m_2(\mathbf{i}(\psi_1\psi_2\lambda)^{1/2}; \mathbf{0}) \cdot \mu_\star^2 \\ &= m_0(\mathbf{i}(\psi_1\psi_2\lambda)^{1/2}) \cdot \mu_\star^2. \end{aligned} \tag{81}$$

Plugging in Eq. (78) and (79) into Eq. (73) and (74) and using Eq. (81), we can see that the expressions for  $\mathcal{B}$  and  $\mathcal{V}$  defined in Eq. (73) and (74) coincide with Eq. (18) and (19) where  $\mathcal{E}_0, \mathcal{E}_1, \mathcal{E}_2$  are provided in Eq. (17). Combining with Eq. (71) and (72) proves the theorem.

## Acknowledgements

This work was partially supported by grants NSF CCF-1714305, IIS-1741162, and ONR N00014-18-1-2729.

## References

- [ADH<sup>+</sup>19] Sanjeev Arora, Simon S Du, Wei Hu, Zhiyuan Li, and Ruosong Wang, *Fine-grained analysis of optimization and generalization for overparameterized two-layer neural networks*, arXiv:1901.08584 (2019).
- [AGZ09] Greg W. Anderson, Alice Guionnet, and Ofer Zeitouni, *An introduction to random matrices*, Cambridge University Press, 2009.
- [AM15] Ahmed Alaoui and Michael W Mahoney, *Fast randomized kernel ridge regression with statistical guarantees*, Advances in Neural Information Processing Systems, 2015, pp. 775–783.
- [AOY19] Dyego Araújo, Roberto I Oliveira, and Daniel Yukimura, *A mean-field limit for certain deep neural networks*, arXiv:1906.00193 (2019).
- [AS17] Madhu S Advani and Andrew M Saxe, *High-dimensional dynamics of generalization error in neural networks*, arXiv:1710.03667 (2017).
- [AZLL18] Zeyuan Allen-Zhu, Yuanzhi Li, and Yingyu Liang, *Learning and generalization in overparameterized neural networks, going beyond two layers*, arXiv:1811.04918 (2018).
- [AZLS18] Zeyuan Allen-Zhu, Yuanzhi Li, and Zhao Song, *A convergence theory for deep learning via over-parameterization*, arXiv:1811.03962 (2018).
- [Bac13] Francis Bach, *Sharp analysis of low-rank kernel matrix approximations*, Conference on Learning Theory, 2013, pp. 185–209.
- [Bac17a] ———, *Breaking the curse of dimensionality with convex neural networks*, The Journal of Machine Learning Research **18** (2017), no. 1, 629–681.

- [Bac17b] ———, *On the equivalence between kernel quadrature rules and random feature expansions*, The Journal of Machine Learning Research **18** (2017), no. 1, 714–751.
- [BHMM18] Mikhail Belkin, Daniel Hsu, Siyuan Ma, and Soumik Mandal, *Reconciling modern machine learning and the bias-variance trade-off*, arXiv preprint arXiv:1812.11118 (2018).
- [BHX19] Mikhail Belkin, Daniel Hsu, and Ji Xu, *Two models of double descent for weak features*, arXiv:1903.07571, 2019.
- [BLLT19] Peter L Bartlett, Philip M Long, Gábor Lugosi, and Alexander Tsigler, *Benign overfitting in linear regression*, arXiv:1906.11300 (2019).
- [BMM18] Mikhail Belkin, Siyuan Ma, and Soumik Mandal, *To understand deep learning we need to understand kernel learning*, arXiv:1802.01396, 2018.
- [BRT18] Mikhail Belkin, Alexander Rakhlin, and Alexandre B Tsybakov, *Does data interpolation contradict statistical optimality?*, arXiv: 1806.09471, 2018.
- [CB18a] Lenaïc Chizat and Francis Bach, *A note on lazy training in supervised differentiable programming*, arXiv:1812.07956 (2018).
- [CB18b] ———, *On the global convergence of gradient descent for over-parameterized models using optimal transport*, Advances in neural information processing systems, 2018, pp. 3036–3046.
- [Chi11] Theodore S Chihara, *An introduction to orthogonal polynomials*, Courier Corporation, 2011.
- [CS13] Xiuyuan Cheng and Amit Singer, *The spectrum of random inner-product kernel matrices*, Random Matrices: Theory and Applications **2** (2013), no. 04, 1350010.
- [Dan17] Amit Daniely, *Sgd learns the conjugate kernel class of the network*, Advances in Neural Information Processing Systems, 2017, pp. 2422–2430.
- [DFS16] Amit Daniely, Roy Frostig, and Yoram Singer, *Toward deeper understanding of neural networks: The power of initialization and a dual view on expressivity*, Advances In Neural Information Processing Systems, 2016, pp. 2253–2261.
- [DHM89] Ronald A DeVore, Ralph Howard, and Charles Micchelli, *Optimal nonlinear approximation*, Manuscripta mathematica **63** (1989), no. 4, 469–478.
- [DJ89] David L Donoho and Iain M Johnstone, *Projection-based approximation and a duality with kernel methods*, The Annals of Statistics (1989), 58–106.
- [DL19] Xialiang Dou and Tengyuan Liang, *Training neural networks as learning data-adaptive kernels: Provable representation and approximation benefits*, arXiv:1901.07114 (2019).
- [DLL<sup>+</sup>18] Simon S Du, Jason D Lee, Haochuan Li, Liwei Wang, and Xiyu Zhai, *Gradient descent finds global minima of deep neural networks*, arXiv:1811.03804 (2018).
- [DM16] Yash Deshpande and Andrea Montanari, *Sparse pca via covariance thresholding*, Journal of Machine Learning Research **17** (2016), 1–41.
- [DZPS18] Simon S Du, Xiyu Zhai, Barnabas Poczos, and Aarti Singh, *Gradient descent provably optimizes over-parameterized neural networks*, arXiv:1810.02054 (2018).
- [EF14] Costas Efthimiou and Christopher Frye, *Spherical harmonics in p dimensions*, World Scientific, 2014.
- [EK10] Noureddine El Karoui, *The spectrum of kernel random matrices*, The Annals of Statistics **38** (2010), no. 1, 1–50.
- [FM19] Zhou Fan and Andrea Montanari, *The spectral norm of random inner-product kernel matrices*, Probability Theory and Related Fields **173** (2019), no. 1-2, 27–85.

- [GAAR18] Adrià Garriga-Alonso, Laurence Aitchison, and Carl Edward Rasmussen, *Deep convolutional networks as shallow gaussian processes*, arXiv:1808.05587 (2018).
- [GJS<sup>+</sup>19] Mario Geiger, Arthur Jacot, Stefano Spigler, Franck Gabriel, Levent Sagun, Stéphane d’Ascoli, Giulio Biroli, Clément Hongler, and Matthieu Wyart, *Scaling description of generalization with number of parameters in deep learning*, arXiv:1901.01608 (2019).
- [GMMM19] Behrooz Ghorbani, Song Mei, Theodor Misiakiewicz, and Andrea Montanari, *Linearized two-layers neural networks in high dimension*, arXiv:1904.12191 (2019).
- [HJ15] Tamir Hazan and Tommi Jaakkola, *Steps toward deep kernel methods from infinite neural networks*, arXiv:1508.05133 (2015).
- [HMRT19] Trevor Hastie, Andrea Montanari, Saharon Rosset, and Ryan J Tibshirani, *Surprises in high-dimensional ridgeless least squares interpolation*, arXiv:1903.08560 (2019).
- [HTF09] Trevor Hastie, Robert Tibshirani, and Jerome Friedman, *The elements of statistical learning*, Springer, 2009.
- [JGH18] Arthur Jacot, Franck Gabriel, and Clément Hongler, *Neural tangent kernel: Convergence and generalization in neural networks*, Advances in neural information processing systems, 2018, pp. 8571–8580.
- [JMM19] Adel Javanmard, Marco Mondelli, and Andrea Montanari, *Analysis of a two-layer neural network via displacement convexity*, arXiv:1901.01375 (2019).
- [KLS18] Dmitry Kobak, Jonathan Lomond, and Benoit Sanchez, *Implicit ridge regularization provided by the minimum-norm least squares estimator when  $n \ll p$* , arXiv:1805.10939 (2018).
- [LBN<sup>+</sup>17] Jaehoon Lee, Yasaman Bahri, Roman Novak, Samuel S Schoenholz, Jeffrey Pennington, and Jascha Sohl-Dickstein, *Deep neural networks as gaussian processes*, arXiv:1711.00165 (2017).
- [LL18] Yuanzhi Li and Yingyu Liang, *Learning overparameterized neural networks via stochastic gradient descent on structured data*, Advances in Neural Information Processing Systems, 2018, pp. 8157–8166.
- [LLC18] Cosme Louart, Zhenyu Liao, and Romain Couillet, *A random matrix approach to neural networks*, The Annals of Applied Probability **28** (2018), no. 2, 1190–1248.
- [LR18] Tengyuan Liang and Alexander Rakhlin, *Just interpolate: Kernel” ridgeless” regression can generalize*, arXiv:1808.00387 (2018).
- [MMM19] Song Mei, Theodor Misiakiewicz, and Andrea Montanari, *Mean-field theory of two-layers neural networks: dimension-free bounds and kernel limit*, arXiv:1902.06015 (2019).
- [MMN18] Song Mei, Andrea Montanari, and Phan-Minh Nguyen, *A mean field view of the landscape of two-layer neural networks*, Proceedings of the National Academy of Sciences **115** (2018), no. 33, E7665–E7671.
- [MRH<sup>+</sup>18] Alexander G de G Matthews, Mark Rowland, Jiri Hron, Richard E Turner, and Zoubin Ghahramani, *Gaussian process behaviour in wide deep neural networks*, arXiv:1804.11271 (2018).
- [MRSY19] Andrea Montanari, Feng Ruan, Youngtak Sohn, and Jun Yan, *The prediction error of max-margin linear classifiers in high dimension*, In preparation, 2019.
- [MVS19] Vidya Muthukumar, Kailas Vodrahalli, and Anant Sahai, *Harmless interpolation of noisy data in regression*, arXiv:1903.09139 (2019).
- [Nea96] Radford M Neal, *Priors for infinite networks*, Bayesian Learning for Neural Networks, Springer, 1996, pp. 29–53.

- [Ngu19] Phan-Minh Nguyen, *Mean field limit of the learning dynamics of multilayer neural networks*, arXiv:1902.02880 (2019).
- [NXB<sup>+</sup>18] Roman Novak, Lechao Xiao, Yasaman Bahri, Jaehoon Lee, Greg Yang, Jiri Hron, Daniel A Abolafia, Jeffrey Pennington, and Jascha Sohl-Dickstein, *Bayesian deep convolutional networks with many channels are gaussian processes*.
- [OS19] Samet Oymak and Mahdi Soltanolkotabi, *Towards moderate overparameterization: global convergence guarantees for training shallow neural networks*, arXiv:1902.04674 (2019).
- [PW17] Jeffrey Pennington and Pratik Worah, *Nonlinear random matrix theory for deep learning*, Advances in Neural Information Processing Systems, 2017, pp. 2637–2646.
- [RJBVE19] Grant Rotskoff, Samy Jelassi, Joan Bruna, and Eric Vanden-Eijnden, *Neuron birth-death dynamics accelerates gradient descent and converges asymptotically*, International Conference on Machine Learning, 2019, pp. 5508–5517.
- [RR08] Ali Rahimi and Benjamin Recht, *Random features for large-scale kernel machines*, Advances in neural information processing systems, 2008, pp. 1177–1184.
- [RR17] Alessandro Rudi and Lorenzo Rosasco, *Generalization properties of learning with random features*, Advances in Neural Information Processing Systems, 2017, pp. 3215–3225.
- [RVE18] Grant M Rotskoff and Eric Vanden-Eijnden, *Neural networks as interacting particle systems: Asymptotic convexity of the loss landscape and universal scaling of the approximation error*, arXiv:1805.00915 (2018).
- [RZ18] Alexander Rakhlin and Xiyu Zhai, *Consistency of interpolation with laplace kernels is a high-dimensional phenomenon*, arXiv:1812.11167 (2018).
- [SS19] Justin Sirignano and Konstantinos Spiliopoulos, *Mean field analysis of neural networks: A central limit theorem*, Stochastic Processes and their Applications (2019).
- [Sze39] Szegő, Gabor, *Orthogonal polynomials*, vol. 23, American Mathematical Soc., 1939.
- [VW18] Santosh Vempala and John Wilmes, *Gradient descent for one-hidden-layer neural networks: Polynomial convergence and sq lower bounds*, arXiv:1805.02677 (2018).
- [Wil97] Christopher KI Williams, *Computing with infinite networks*, Advances in neural information processing systems, 1997, pp. 295–301.
- [YRC07] Yuan Yao, Lorenzo Rosasco, and Andrea Caponnetto, *On early stopping in gradient descent learning*, Constructive Approximation **26** (2007), no. 2, 289–315.
- [ZBH<sup>+</sup>16] Chiyuan Zhang, Samy Bengio, Moritz Hardt, Benjamin Recht, and Oriol Vinyals, *Understanding deep learning requires rethinking generalization*, arXiv:1611.03530 (2016).
- [ZCZG18] Difan Zou, Yuan Cao, Dongruo Zhou, and Quanquan Gu, *Stochastic gradient descent optimizes over-parameterized deep relu networks*, arXiv:1811.08888 (2018).



## A Technical background and notations

In this section we introduce some notations and technical background which will be useful for the proofs in the next sections. In particular, we will use decompositions in (hyper-)spherical harmonics on the  $\mathbb{S}^{d-1}(\sqrt{d})$  and in orthogonal polynomials on the real line. All of the properties listed below are classical: we will however prove a few facts that are slightly less standard. We refer the reader to [EF14, Sze39, Chi11] for further information on these topics. Expansions in spherical harmonics have been used in the past in the statistics literature, for instance in [DJ89, Bac17a].

### A.1 Notations

Let  $\mathbb{R}$  denote the set of real numbers,  $\mathbb{C}$  the set of complex numbers, and  $\mathbb{N} = \{0, 1, 2, \dots\}$  the set of natural numbers. For  $z \in \mathbb{C}$ , let  $\Re z$  and  $\Im z$  denote the real part and the imaginary part of  $z$  respectively. We denote by  $\mathbb{C}_+ = \{z \in \mathbb{C} : \Im z > 0\}$  the set of complex numbers with positive imaginary part. We denote by  $\mathbf{i} = \sqrt{-1}$  the imaginary unit. We denote by  $\mathbb{S}^{d-1}(r) = \{\mathbf{x} \in \mathbb{R}^d : \|\mathbf{x}\|_2 = r\}$  the set of  $d$ -dimensional vectors with radius  $r$ . For an integer  $k$ , let  $[k]$  denote the set  $\{1, 2, \dots, k\}$ .

Throughout the proofs, let  $O_d(\cdot)$  (respectively  $o_d(\cdot)$ ,  $\Omega_d(\cdot)$ ) denote the standard big-O (respectively little-o, big-Omega) notation, where the subscript  $d$  emphasizes the asymptotic variable. We denote by  $O_{d,\mathbb{P}}(\cdot)$  the big-O in probability notation:  $h_1(d) = O_{d,\mathbb{P}}(h_2(d))$  if for any  $\varepsilon > 0$ , there exists  $C_\varepsilon > 0$  and  $d_\varepsilon \in \mathbb{Z}_{>0}$ , such that

$$\mathbb{P}(|h_1(d)/h_2(d)| > C_\varepsilon) \leq \varepsilon, \quad \forall d \geq d_\varepsilon.$$

We denote by  $o_{d,\mathbb{P}}(\cdot)$  the little-o in probability notation:  $h_1(d) = o_{d,\mathbb{P}}(h_2(d))$ , if  $h_1(d)/h_2(d)$  converges to 0 in probability. We write  $h(d) = O_d(\text{Poly}(\log d))$ , if there exists a constant  $k$ , such that  $h(d) = O_d((\log d)^k)$ .

For a matrix  $\mathbf{A} \in \mathbb{R}^{n \times m}$ , we denote by  $\|\mathbf{A}\|_F = (\sum_{i \in [n], j \in [m]} A_{ij}^2)^{1/2}$  the Frobenius norm of  $\mathbf{A}$ ,  $\|\mathbf{A}\|_\star$  the nuclear norm of  $\mathbf{A}$ ,  $\|\mathbf{A}\|_{\text{op}}$  the operator norm of  $\mathbf{A}$ , and  $\|\mathbf{A}\|_{\max} = \max_{i \in [n], j \in [m]} |A_{ij}|$  the maximum norm of  $\mathbf{A}$ . For a matrix  $\mathbf{A} \in \mathbb{R}^{n \times n}$ , we denote by  $\text{Tr}(\mathbf{A}) = \sum_{i=1}^n A_{ii}$  the trace of  $\mathbf{A}$ . For two integers  $a$  and  $b$ , we denote by  $\text{Tr}_{[a,b]}(\mathbf{A}) = \sum_{i=a}^b A_{ii}$  the partial trace of  $\mathbf{A}$ . For two matrices  $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{n \times m}$ , let  $\mathbf{A} \odot \mathbf{B}$  denotes the element-wise product of  $\mathbf{A}$  and  $\mathbf{B}$ .

Let  $\mu_G$  denote the standard Gaussian measure. Let  $\gamma_d$  denote the uniform probability distribution on  $\mathbb{S}^{d-1}(\sqrt{d})$ . We denote by  $\mu_d$  the distribution of  $\langle \mathbf{x}_1, \mathbf{x}_2 \rangle / \sqrt{d}$  when  $\mathbf{x}_1, \mathbf{x}_2 \sim \mathbf{N}(\mathbf{0}, \mathbf{I}_d)$ ,  $\tau_d$  the distribution of  $\langle \mathbf{x}_1, \mathbf{x}_2 \rangle / \sqrt{d}$  when  $\mathbf{x}_1, \mathbf{x}_2 \sim \text{Unif}(\mathbb{S}^{d-1}(\sqrt{d}))$ , and  $\tilde{\tau}_d$  the distribution of  $\langle \mathbf{x}_1, \mathbf{x}_2 \rangle$  when  $\mathbf{x}_1, \mathbf{x}_2 \sim \text{Unif}(\mathbb{S}^{d-1}(\sqrt{d}))$ .

### A.2 Functional spaces over the sphere

For  $d \geq 1$ , we let  $\mathbb{S}^{d-1}(r) = \{\mathbf{x} \in \mathbb{R}^d : \|\mathbf{x}\|_2 = r\}$  denote the sphere with radius  $r$  in  $\mathbb{R}^d$ . We will mostly work with the sphere of radius  $\sqrt{d}$ ,  $\mathbb{S}^{d-1}(\sqrt{d})$  and will denote by  $\gamma_d$  the uniform probability measure on  $\mathbb{S}^{d-1}(\sqrt{d})$ . All functions in the following are assumed to be elements of  $L^2(\mathbb{S}^{d-1}(\sqrt{d}), \gamma_d)$ , with scalar product and norm denoted as  $\langle \cdot, \cdot \rangle_{L^2}$  and  $\|\cdot\|_{L^2}$ :

$$\langle f, g \rangle_{L^2} \equiv \int_{\mathbb{S}^{d-1}(\sqrt{d})} f(\mathbf{x}) g(\mathbf{x}) \gamma_d(d\mathbf{x}). \quad (82)$$

For  $\ell \in \mathbb{N}_{\geq 0}$ , let  $\tilde{V}_{d,\ell}$  be the space of homogeneous harmonic polynomials of degree  $\ell$  on  $\mathbb{R}^d$  (i.e. homogeneous polynomials  $q(\mathbf{x})$  satisfying  $\Delta q(\mathbf{x}) = 0$ ), and denote by  $V_{d,\ell}$  the linear space of functions obtained by restricting the polynomials in  $\tilde{V}_{d,\ell}$  to  $\mathbb{S}^{d-1}(\sqrt{d})$ . With these definitions, we have the following orthogonal decomposition

$$L^2(\mathbb{S}^{d-1}(\sqrt{d}), \gamma_d) = \bigoplus_{\ell=0}^{\infty} V_{d,\ell}. \quad (83)$$

The dimension of each subspace is given by

$$\dim(V_{d,\ell}) = B(d, \ell) = \frac{2\ell + d - 2}{\ell} \binom{\ell + d - 3}{\ell - 1}. \quad (84)$$

For each  $\ell \in \mathbb{N}_{\geq 0}$ , the spherical harmonics  $\{Y_{\ell j}^{(d)}\}_{1 \leq j \leq B(d, \ell)}$  form an orthonormal basis of  $V_{d, \ell}$ :

$$\langle Y_{ki}^{(d)}, Y_{sj}^{(d)} \rangle_{L^2} = \delta_{ij} \delta_{ks}.$$

Note that our convention is different from the more standard one, that defines the spherical harmonics as functions on  $\mathbb{S}^{d-1}(1)$ . It is immediate to pass from one convention to the other by a simple scaling. We will drop the superscript  $d$  and write  $Y_{\ell, j} = Y_{\ell, j}^{(d)}$  whenever clear from the context.

We denote by  $P_k$  the orthogonal projections to  $V_{d, k}$  in  $L^2(\mathbb{S}^{d-1}(\sqrt{d}), \gamma_d)$ . This can be written in terms of spherical harmonics as

$$P_k f(\mathbf{x}) \equiv \sum_{l=1}^{B(d, k)} \langle f, Y_{kl} \rangle_{L^2} Y_{kl}(\mathbf{x}). \quad (85)$$

Then for a function  $f \in L^2(\mathbb{S}^{d-1}(\sqrt{d}))$ , we have

$$f(\mathbf{x}) = \sum_{k=0}^{\infty} P_k f(\mathbf{x}) = \sum_{k=0}^{\infty} \sum_{l=1}^{B(d, k)} \langle f, Y_{kl} \rangle_{L^2} Y_{kl}(\mathbf{x}).$$

### A.3 Gegenbauer polynomials

The  $\ell$ -th Gegenbauer polynomial  $Q_{\ell}^{(d)}$  is a polynomial of degree  $\ell$ . Consistently with our convention for spherical harmonics, we view  $Q_{\ell}^{(d)}$  as a function  $Q_{\ell}^{(d)} : [-d, d] \rightarrow \mathbb{R}$ . The set  $\{Q_{\ell}^{(d)}\}_{\ell \geq 0}$  forms an orthogonal basis on  $L^2([-d, d], \tilde{\tau}_d)$  (where  $\tilde{\tau}_d$  is the distribution of  $\langle \mathbf{x}_1, \mathbf{x}_2 \rangle$  when  $\mathbf{x}_1, \mathbf{x}_2 \sim_{i.i.d.} \text{Unif}(\mathbb{S}^{d-1}(\sqrt{d}))$ ), satisfying the normalization condition:

$$\langle Q_k^{(d)}, Q_j^{(d)} \rangle_{L^2(\tilde{\tau}_d)} = \frac{1}{B(d, k)} \delta_{jk}. \quad (86)$$

In particular, these polynomials are normalized so that  $Q_{\ell}^{(d)}(d) = 1$ . As above, we will omit the superscript  $d$  when clear from the context (write it as  $Q_{\ell}$  for notation simplicity).

Gegenbauer polynomials are directly related to spherical harmonics as follows. Fix  $\mathbf{v} \in \mathbb{S}^{d-1}(\sqrt{d})$  and consider the subspace of  $V_{\ell}$  formed by all functions that are invariant under rotations in  $\mathbb{R}^d$  that keep  $\mathbf{v}$  unchanged. It is not hard to see that this subspace has dimension one, and coincides with the span of the function  $Q_{\ell}^{(d)}(\langle \mathbf{v}, \cdot \rangle)$ .

We will use the following properties of Gegenbauer polynomials

1. For  $\mathbf{x}, \mathbf{y} \in \mathbb{S}^{d-1}(\sqrt{d})$

$$\langle Q_j^{(d)}(\langle \mathbf{x}, \cdot \rangle), Q_k^{(d)}(\langle \mathbf{y}, \cdot \rangle) \rangle_{L^2(\mathbb{S}^{d-1}(\sqrt{d}), \gamma_d)} = \frac{1}{B(d, k)} \delta_{jk} Q_k^{(d)}(\langle \mathbf{x}, \mathbf{y} \rangle). \quad (87)$$

2. For  $\mathbf{x}, \mathbf{y} \in \mathbb{S}^{d-1}(\sqrt{d})$

$$Q_k^{(d)}(\langle \mathbf{x}, \mathbf{y} \rangle) = \frac{1}{B(d, k)} \sum_{i=1}^{B(d, k)} Y_{ki}^{(d)}(\mathbf{x}) Y_{ki}^{(d)}(\mathbf{y}). \quad (88)$$

Note in particular that property 2 implies that –up to a constant–  $Q_k^{(d)}(\langle \mathbf{x}, \mathbf{y} \rangle)$  is a representation of the projector onto the subspace of degree- $k$  spherical harmonics

$$(P_k f)(\mathbf{x}) = B(d, k) \int_{\mathbb{S}^{d-1}(\sqrt{d})} Q_k^{(d)}(\langle \mathbf{x}, \mathbf{y} \rangle) f(\mathbf{y}) \gamma_d(d\mathbf{y}). \quad (89)$$

For a function  $\sigma \in L^2([-\sqrt{d}, \sqrt{d}], \tau_d)$  (where  $\tau_d$  is the distribution of  $\langle \mathbf{x}_1, \mathbf{x}_2 \rangle / \sqrt{d}$  when  $\mathbf{x}_1, \mathbf{x}_2 \sim_{iid} \text{Unif}(\mathbb{S}^{d-1}(\sqrt{d}))$ ), denoting its spherical harmonics coefficients  $\lambda_{d,k}(\sigma)$  to be

$$\lambda_{d,k}(\sigma) = \int_{[-\sqrt{d}, \sqrt{d}]} \sigma(x) Q_k^{(d)}(\sqrt{d}x) \tau_d(x), \quad (90)$$

then we have the following equation holds in  $L^2([-\sqrt{d}, \sqrt{d}], \tau_d)$  sense

$$\sigma(x) = \sum_{k=0}^{\infty} \lambda_{d,k}(\sigma) B(d, k) Q_k^{(d)}(\sqrt{d}x).$$

## A.4 Hermite polynomials

The Hermite polynomials  $\{\text{He}_k\}_{k \geq 0}$  form an orthogonal basis of  $L^2(\mathbb{R}, \mu_G)$ , where  $\mu_G(dx) = e^{-x^2/2} dx / \sqrt{2\pi}$  is the standard Gaussian measure, and  $\text{He}_k$  has degree  $k$ . We will follow the classical normalization (here and below, expectation is with respect to  $G \sim \mathbf{N}(0, 1)$ ):

$$\mathbb{E}\{\text{He}_j(G) \text{He}_k(G)\} = k! \delta_{jk}. \quad (91)$$

As a consequence, for any function  $\sigma \in L^2(\mathbb{R}, \mu_G)$ , we have the decomposition

$$\sigma(x) = \sum_{k=1}^{\infty} \frac{\mu_k(\sigma)}{k!} \text{He}_k(x), \quad \mu_k(\sigma) \equiv \mathbb{E}\{\sigma(G) \text{He}_k(G)\}. \quad (92)$$

The Hermite polynomials can be obtained as high-dimensional limits of the Gegenbauer polynomials introduced in the previous section. Indeed, the Gegenbauer polynomials (up to a  $\sqrt{d}$  scaling in domain) are constructed by Gram-Schmidt orthogonalization of the monomials  $\{x^k\}_{k \geq 0}$  with respect to the measure  $\tau_d$ , while Hermite polynomial are obtained by Gram-Schmidt orthogonalization with respect to  $\mu_G$ . Since  $\tau_d \Rightarrow \mu_G$  (here  $\Rightarrow$  denotes weak convergence), it is immediate to show that, for any fixed integer  $k$ ,

$$\lim_{d \rightarrow \infty} \text{Coeff}\{Q_k^{(d)}(\sqrt{d}x) B(d, k)^{1/2}\} = \text{Coeff}\left\{\frac{1}{(k!)^{1/2}} \text{He}_k(x)\right\}. \quad (93)$$

Here and below, for  $P$  a polynomial,  $\text{Coeff}\{P(x)\}$  is the vector of the coefficients of  $P$ . As a consequence, for any fixed integer  $k$ , we have

$$\mu_k(\sigma) = \lim_{d \rightarrow \infty} \lambda_{d,k}(\sigma) (B(d, k) k!)^{1/2}, \quad (94)$$

where  $\mu_k(\sigma)$  and  $\lambda_{d,k}(\sigma)$  are given in Eq. (92) and (90).

## B Proof of Proposition 7.1

Throughout the proof of Proposition 7.1, we assume  $\psi_1 = \psi_{1,d} = N/d$  and  $\psi_2 = \psi_{2,d} = n/d$  for notation simplicity. The proof can be directly generalized to the case when  $\lim_{d \rightarrow \infty} N/d = \psi_1$  and  $\lim_{d \rightarrow \infty} n/d = \psi_2$ .

**Remark 7.** For any kernel function  $\Sigma_d$  satisfying Assumption 3, we can always find a sequence  $(F_{d,k}^2 \in \mathbb{R}_+)_{k \geq 2}$  satisfying: (1)  $\lim_{d \rightarrow \infty} \sum_{k \geq 2} F_{d,k}^2 = F_\star^2$ ; (2) Defining  $\beta_{d,k} \sim \mathbf{N}(\mathbf{0}, [F_{d,k}^2/B(d, k)] \mathbf{I}_{B(d, k)})$  independently for  $k \geq 2$ , and  $g_d^{\text{NL}}(\mathbf{x}) = \sum_{k \geq 2} \sum_{l \in [B(d, k)]} (\beta_{d,k})_l Y_{kl}^{(d)}(\mathbf{x})$ , then  $g_d^{\text{NL}}$  is a centered Gaussian process with covariance function  $\Sigma_d$ .

To prove this claim, we define the sequence  $(F_{d,k}^2)_{k \geq 2}$  to be the coefficients of Gegenbauer expansion of  $\Sigma_d$ :

$$\Sigma_d(x/\sqrt{d}) = \sum_{k=2}^{\infty} F_{d,k}^2 Q_k^{(d)}(\sqrt{d}x).$$

In the expansion, the zeroth and first order coefficients are 0, because, according to Assumption 3,

$$\mathbb{E}_{\mathbf{x} \sim \text{Unif}(\mathbb{S}^{d-1}(\sqrt{d}))}[\Sigma_d(x_1/\sqrt{d})] = \mathbb{E}_{\mathbf{x} \sim \text{Unif}(\mathbb{S}^{d-1}(\sqrt{d}))}[\Sigma_d(x_1/\sqrt{d})x_1] = 0.$$

To check point (1), we have  $\Sigma_d(1) = \sum_{k=2}^{\infty} F_{d,k}^2 Q_k^{(d)}(d) = \sum_{k=2}^{\infty} F_{d,k}^2$ , and by Assumption 3 we have  $\lim_{d \rightarrow \infty} \Sigma_d(1) = F_{\star}^2$ , so that (1) holds.

To check point (2), defining  $(\beta_{d,k})_{k \geq 2}$  and  $g_d^{\text{NL}}(\mathbf{x})$  accordingly, we have

$$\begin{aligned} \mathbb{E}_{\beta}[g_d^{\text{NL}}(\mathbf{x}_1)g_d^{\text{NL}}(\mathbf{x}_2)] &= \mathbb{E}_{\beta} \left[ \left( \sum_{k \geq 2} \sum_{l \in [B(d,k)]} (\beta_{d,k})_l Y_{kl}^{(d)}(\mathbf{x}_1) \right) \left( \sum_{k \geq 2} \sum_{l \in [B(d,k)]} (\beta_{d,k})_l Y_{kl}^{(d)}(\mathbf{x}_2) \right) \right] \\ &= \sum_{k \geq 2} F_{d,k}^2 Y_{kl}^{(d)}(\mathbf{x}_1) Y_{kl}^{(d)}(\mathbf{x}_2) / B(d,k) = \sum_{k \geq 2} F_{d,k}^2 Q_k^{(d)}(\langle \mathbf{x}_1, \mathbf{x}_2 \rangle) = \Sigma_d(\langle \mathbf{x}_1, \mathbf{x}_2 \rangle / d). \end{aligned}$$

**Remark 8.** Let us write the risk function  $R_{\text{RF}}(f_d, \mathbf{X}, \Theta, \lambda)$  as a function of  $\beta_1, \theta_1, \dots, \theta_n$  and de-emphasize its dependence on other variables, i.e.,

$$R_{\text{RF}}(f_d, \mathbf{X}, \Theta, \lambda) \equiv \tilde{R}(\beta_1, \theta_1, \dots, \theta_n).$$

Under Assumption 3, for any orthogonal matrix  $\mathcal{O} \in \mathbb{R}^{d \times d}$ , we have distribution equivalence

$$\tilde{R}(\beta_1, \theta_1, \dots, \theta_N) \stackrel{d}{=} \tilde{R}(\mathcal{O}\beta_1, \mathcal{O}\theta_1, \mathcal{O} \dots, \mathcal{O}\theta_N),$$

where the randomness is given by  $(\mathbf{X}, \Theta, \varepsilon, f_d^{\text{NL}})$ . Therefore, as long as we show Proposition 7.1 under the assumption that  $\beta_{d,1} \sim \text{Unif}(S^{d-1}(F_{d,1}))$  which is independent from all other random variables, then Proposition 7.1 immediately holds for any deterministic  $\beta_{d,1}$  with  $\|\beta_{d,1}\|_2^2 = F_{d,1}^2$ .

By Remark 7 and 8 above, in the following, we will prove Proposition 7.1 under Assumption 4 instead of Assumption 3.

**Assumption 4** (Reformulation of Assumption 3). *Let  $(F_{d,k}^2 \in \mathbb{R}_+)_{d \geq 1, k \geq 0}$  be an array of non-negative numbers. Let  $\beta_{d,0} = F_{d,0}$ ,  $\beta_{d,1} \sim \text{Unif}(S^{d-1}(F_{d,1}))$ , and  $\beta_{d,k} \sim \mathcal{N}(\mathbf{0}, [F_{d,k}^2/B(d,k)]\mathbf{I}_{B(d,k)})$  independently for  $k \geq 2$ . We assume the regression function to be*

$$f_d(\mathbf{x}) = \sum_{k \geq 0} \sum_{l \in [B(d,k)]} (\beta_{d,k})_l Y_{kl}^{(d)}(\mathbf{x}).$$

Assume  $y_i = f_d(\mathbf{x}_i) + \varepsilon_i$ , where  $\varepsilon_i \sim_{\text{iid}} \mathbb{P}_{\varepsilon}$ , with  $\mathbb{E}_{\varepsilon}(\varepsilon_1) = 0$ ,  $\mathbb{E}_{\varepsilon}(\varepsilon_1^2) = \tau^2$ , and  $\mathbb{E}_{\varepsilon}(\varepsilon_1^4) < \infty$ . Finally, assume

$$\begin{aligned} \lim_{d \rightarrow \infty} F_{d,0}^2 &= F_0^2 < \infty, \\ \lim_{d \rightarrow \infty} F_{d,1}^2 &= F_1^2 < \infty, \\ \lim_{d \rightarrow \infty} \sum_{k \geq 2} F_{d,k}^2 &= F_{\star}^2 < \infty. \end{aligned}$$

Proposition 7.1 is a direct consequence of the following three lemmas.

**Lemma B.1** (Decomposition). *Let  $\lambda_{d,k}(\sigma)$  be the Gegenbauer coefficients of function  $\sigma$ , i.e., we have*

$$\sigma(x) = \sum_{k=0}^{\infty} \lambda_{d,k}(\sigma) B(d,k) Q_k(\sqrt{d} \cdot x).$$

Under the assumptions of Proposition 7.1 (replacing Assumption 3 by Assumption 4), for any  $\lambda > 0$ , we have

$$\mathbb{E}_{\beta, \varepsilon}[R_{\text{RF}}(f_d, \mathbf{X}, \Theta, \lambda)] = \sum_{k=0}^{\infty} F_{d,k}^2 - 2 \sum_{k=0}^{\infty} F_{d,k}^2 S_{1k} + \sum_{k=0}^{\infty} F_{d,k}^2 S_{2k} + \tau^2 S_3, \quad (95)$$

where

$$\begin{aligned} S_{1k} &= \frac{1}{\sqrt{d}} \lambda_{d,k}(\sigma) \text{Tr} \left[ Q_k(\Theta \mathbf{X}^\top) \mathbf{Z} (\mathbf{Z}^\top \mathbf{Z} + \psi_1 \psi_2 \lambda \mathbf{I}_N)^{-1} \right], \\ S_{2k} &= \frac{1}{d} \text{Tr} \left[ (\mathbf{Z}^\top \mathbf{Z} + \psi_1 \psi_2 \lambda \mathbf{I}_N)^{-1} \mathbf{U} (\mathbf{Z}^\top \mathbf{Z} + \psi_1 \psi_2 \lambda \mathbf{I}_N)^{-1} \mathbf{Z}^\top Q_k(\mathbf{X} \mathbf{X}^\top) \mathbf{Z} \right], \\ S_3 &= \frac{1}{d} \text{Tr} \left[ (\mathbf{Z}^\top \mathbf{Z} + \psi_1 \psi_2 \lambda \mathbf{I}_N)^{-1} \mathbf{U} (\mathbf{Z}^\top \mathbf{Z} + \psi_1 \psi_2 \lambda \mathbf{I}_N)^{-1} \mathbf{Z}^\top \mathbf{Z} \right], \end{aligned} \quad (96)$$

and

$$\mathbf{U} = \sum_{k=0}^{\infty} \lambda_{d,k}(\sigma)^2 B(d, k) Q_k(\Theta \Theta^\top), \quad (97)$$

and  $\mathbf{Z}$  is given by Eq. (56).

**Lemma B.2.** Under the same definitions and assumptions of Proposition 7.1 and Lemma B.1, for any  $\lambda > 0$ , we have ( $\mathbb{E}$  is the expectation taken with respect to the randomness in  $\mathbf{X}$  and  $\Theta$ )

$$\begin{aligned} \mathbb{E}|1 - 2S_{10} + S_{20}| &= o_d(1), \\ \mathbb{E} \left[ \sup_{k \geq 2} |S_{1k}| \right] &= o_d(1), \\ \mathbb{E} \left[ \sup_{k \geq 2} |S_{2k} - S_3| \right] &= o_d(1), \\ \mathbb{E}|S_{11} - \Psi_1| &= o_d(1), \\ \mathbb{E}|S_{21} - \Psi_2| &= o_d(1), \\ \mathbb{E}|S_3 - \Psi_3| &= o_d(1), \end{aligned}$$

where  $S_{1k}, S_{2k}, S_3$  are given by Eq. (96), and  $\Psi_1, \Psi_2, \Psi_3$  are given by Eq. (55).

**Lemma B.3.** Under the assumptions of Proposition 7.1 (replacing Assumption 3 by Assumption 4), we have

$$\mathbb{E}_{\mathbf{X}, \Theta} \left[ \text{Var}_{\beta, \epsilon} \left( R_{\text{RF}}(f_d, \mathbf{X}, \Theta, \lambda) \middle| \mathbf{X}, \Theta \right)^{1/2} \right] = o_d(1). \quad (98)$$

We defer the proofs of these three lemmas to the following subsections.

By Lemma B.1, we get

$$\begin{aligned} \mathbb{E}_{\beta, \epsilon} [R_{\text{RF}}(f_d, \mathbf{X}, \Theta, \lambda)] &= \sum_{k=0}^{\infty} F_{d,k}^2 - 2 \sum_{k=0}^{\infty} F_{d,k}^2 S_{1k} + \sum_{k=0}^{\infty} F_{d,k}^2 S_{2k} + \tau^2 S_3 \\ &= F_{d,0}^2 (1 - 2S_{10} + S_{20}) + F_{d,1}^2 (1 - 2S_{11} + S_{21}) + \sum_{k=2}^{\infty} F_{d,k}^2 (1 - 2S_{1k} + S_{2k}) + \tau^2 S_3, \end{aligned}$$

By Lemma B.2 and Assumption 4, we get

$$\begin{aligned} &\mathbb{E}_{\mathbf{X}, \Theta} \left| \mathbb{E}_{\beta, \epsilon} [R_{\text{RF}}(f_d, \mathbf{X}, \Theta, \lambda)] - \left[ F_{d,1}^2 (1 - 2\Psi_1 + \Psi_2) + \left( \tau^2 + \sum_{k=2}^{\infty} F_{d,k}^2 \right) \Psi_3 + \sum_{k=2}^{\infty} F_{d,k}^2 \right] \right| \\ &\leq F_{d,0}^2 \cdot \mathbb{E}|1 - 2S_{10} + S_{20}| + F_{d,1}^2 \cdot \left[ \mathbb{E}|S_{11} - \Psi_1| + \mathbb{E}|S_{21} - \Psi_2| \right] \\ &\quad + \left( \sum_{k=2}^{\infty} F_{d,k}^2 \right) \cdot \sup_{k \geq 2} \left[ 2\mathbb{E}|S_{1k}| + \mathbb{E}|S_{2k} - \Psi_3| \right] + \tau^2 \mathbb{E}|S_3 - \Psi_3| \\ &= o_d(1). \end{aligned}$$

This proves the Eq. (54). Combining with Lemma B.3 (and  $\mathbb{E}[\Psi_1], \mathbb{E}[\Psi_2], \mathbb{E}[\Psi_3] = O_d(1)$ ) concludes the proof.

In the remaining part of this section, we prove Lemma B.1, B.2, and B.3. The proof of Lemma B.1 is relatively straightforward and is given in Section B.1. The proof of Lemma B.2 and B.3 is more complicated. We give their proof in Section B.2 and B.3. The proof of Lemma B.2 and B.3 depends on some other lemmas that is proved in Section B.4 and B.5.

## B.1 Proof of Lemma B.1

By the definition of prediction error given in Eq. (3), we have

$$\begin{aligned} R_{\text{RF}}(f_d, \mathbf{X}, \boldsymbol{\Theta}, \lambda) &= \mathbb{E}_{\mathbf{x}}[(f_d(\mathbf{x}) - \mathbf{y}^\top \mathbf{Z}(\mathbf{Z}^\top \mathbf{Z} + \psi_1 \psi_2 \lambda \mathbf{I}_N)^{-1} \boldsymbol{\sigma}(\mathbf{x})/\sqrt{d})^2] \\ &= \mathbb{E}_{\mathbf{x}}[f_d(\mathbf{x})^2] - 2\mathbf{y}^\top \mathbf{Z}(\mathbf{Z}^\top \mathbf{Z} + \psi_1 \psi_2 \lambda \mathbf{I}_N)^{-1} \mathbf{V}/\sqrt{d} \\ &\quad + \mathbf{y}^\top \mathbf{Z}(\mathbf{Z}^\top \mathbf{Z} + \psi_1 \psi_2 \lambda \mathbf{I}_N)^{-1} \mathbf{U}(\mathbf{Z}^\top \mathbf{Z} + \psi_1 \psi_2 \lambda \mathbf{I}_N)^{-1} \mathbf{Z}^\top \mathbf{y}/d, \end{aligned} \quad (99)$$

where

$$\begin{aligned} \boldsymbol{\sigma}(\mathbf{x}) &= (\sigma(\langle \boldsymbol{\theta}_1, \mathbf{x} \rangle/\sqrt{d}), \dots, \sigma(\langle \boldsymbol{\theta}_N, \mathbf{x} \rangle/\sqrt{d}))^\top \in \mathbb{R}^N, \\ \mathbf{y} &= (y_1, \dots, y_n)^\top = \mathbf{f} + \boldsymbol{\varepsilon} \in \mathbb{R}^n, \\ \mathbf{f} &= (f_d(\mathbf{x}_1), \dots, f_d(\mathbf{x}_n))^\top \in \mathbb{R}^n, \\ \boldsymbol{\varepsilon} &= (\varepsilon_1, \dots, \varepsilon_n)^\top \in \mathbb{R}^n, \end{aligned}$$

and  $\mathbf{V} = (V_1, \dots, V_N)^\top \in \mathbb{R}^N$ , and  $\mathbf{U} = (U_{ij})_{i,j \in [N]} \in \mathbb{R}^{N \times N}$ , with

$$\begin{aligned} V_i &= \mathbb{E}_{\mathbf{x}}[f_d(\mathbf{x})\sigma(\langle \boldsymbol{\theta}_i, \mathbf{x} \rangle/\sqrt{d})], \\ U_{ij} &= \mathbb{E}_{\mathbf{x}}[\sigma(\langle \boldsymbol{\theta}_i, \mathbf{x} \rangle/\sqrt{d})\sigma(\langle \boldsymbol{\theta}_j, \mathbf{x} \rangle/\sqrt{d})]. \end{aligned}$$

Taking expectation over  $\boldsymbol{\beta}$  and  $\boldsymbol{\varepsilon}$ , we get

$$\mathbb{E}_{\boldsymbol{\beta}, \boldsymbol{\varepsilon}}[R_{\text{RF}}(f_d, \mathbf{X}, \boldsymbol{\Theta}, \lambda)] = \sum_{k \geq 0} F_{d,k}^2 - 2T_1 + T_2 + T_3,$$

where

$$\begin{aligned} T_1 &= \mathbb{E}_{\boldsymbol{\beta}}[\mathbf{f}^\top \mathbf{Z}(\mathbf{Z}^\top \mathbf{Z} + \psi_1 \psi_2 \lambda \mathbf{I}_N)^{-1} \mathbf{V}]/\sqrt{d}, \\ T_2 &= \mathbb{E}_{\boldsymbol{\beta}}[\mathbf{f}^\top \mathbf{Z}(\mathbf{Z}^\top \mathbf{Z} + \psi_1 \psi_2 \lambda \mathbf{I}_N)^{-1} \mathbf{U}(\mathbf{Z}^\top \mathbf{Z} + \psi_1 \psi_2 \lambda \mathbf{I}_N)^{-1} \mathbf{Z}^\top \mathbf{f}]/d, \\ T_3 &= \mathbb{E}_{\boldsymbol{\varepsilon}}[\boldsymbol{\varepsilon}^\top \mathbf{Z}(\mathbf{Z}^\top \mathbf{Z} + \psi_1 \psi_2 \lambda \mathbf{I}_N)^{-1} \mathbf{U}(\mathbf{Z}^\top \mathbf{Z} + \psi_1 \psi_2 \lambda \mathbf{I}_N)^{-1} \mathbf{Z}^\top \boldsymbol{\varepsilon}]/d. \end{aligned}$$

**Term  $T_1$ .** Denote  $\mathbf{Y}_{k,\mathbf{x}} = (Y_{kl}(\mathbf{x}_i))_{i \in [n], l \in [B(d,k)]} \in \mathbb{R}^{n \times B(d,k)}$  and  $\mathbf{Y}_{k,\boldsymbol{\theta}} = (Y_{kl}(\boldsymbol{\theta}_a))_{a \in [N], l \in [B(d,k)]} \in \mathbb{R}^{N \times B(d,k)}$  where  $(Y_{kl})_{k \geq 0, l \in [B(d,k)]}$  is the set of spherical harmonics with domain  $\mathbb{S}^{d-1}(\sqrt{d})$  (c.f. Section A). Then we have

$$\mathbf{f} = \sum_{k=0}^{\infty} \mathbf{Y}_{k,\mathbf{x}} \boldsymbol{\beta}_k \in \mathbb{R}^n, \quad \mathbf{V} = \sum_{k=0}^{\infty} \lambda_{d,k}(\sigma) \mathbf{Y}_{k,\boldsymbol{\theta}} \boldsymbol{\beta}_k \in \mathbb{R}^N.$$

Therefore, by the assumption that  $\boldsymbol{\beta}_k \sim \mathcal{N}(\mathbf{0}, F_{d,k}^2/B(d,k))$  for  $k \geq 2$  and  $\boldsymbol{\beta}_1 \sim \text{Unif}(\mathbb{S}^{d-1}(F_{d,1}))$  independently, we have

$$\begin{aligned} T_1 &= \mathbb{E}_{\boldsymbol{\beta}} \left[ \left( \sum_{s=0}^{\infty} \boldsymbol{\beta}_s^\top \mathbf{Y}_{s,\mathbf{x}}^\top \right) \mathbf{Z}(\mathbf{Z}^\top \mathbf{Z} + \psi_1 \psi_2 \lambda \mathbf{I}_N)^{-1} \left( \sum_{t=0}^{\infty} \lambda_{d,t}(\sigma) \mathbf{Y}_{t,\boldsymbol{\theta}} \boldsymbol{\beta}_t \right) \right] / \sqrt{d} \\ &= \mathbb{E}_{\boldsymbol{\beta}} \left[ \text{Tr} \left( \left( \sum_{s,t=0}^{\infty} \lambda_{d,t}(\sigma) \mathbf{Y}_{t,\boldsymbol{\theta}} \boldsymbol{\beta}_t \boldsymbol{\beta}_s^\top \mathbf{Y}_{s,\mathbf{x}}^\top \right) \mathbf{Z}(\mathbf{Z}^\top \mathbf{Z} + \psi_1 \psi_2 \lambda \mathbf{I}_N)^{-1} \right) \right] / \sqrt{d} \\ &= \sum_{k=0}^{\infty} F_{d,k}^2 \lambda_{d,k}(\sigma) \cdot \text{Tr} \left( \left( \mathbf{Y}_{k,\boldsymbol{\theta}} \mathbf{Y}_{k,\mathbf{x}}^\top / B(d,k) \right) \mathbf{Z}(\mathbf{Z}^\top \mathbf{Z} + \psi_1 \psi_2 \lambda \mathbf{I}_N)^{-1} \right) / \sqrt{d} \\ &= \sum_{k=0}^{\infty} F_{d,k}^2 \lambda_{d,k}(\sigma) \cdot \text{Tr} \left[ Q_k(\boldsymbol{\Theta} \mathbf{X}^\top) \mathbf{Z}(\mathbf{Z}^\top \mathbf{Z} / d + \psi_1 \psi_2 \lambda \mathbf{I}_N)^{-1} \right] / \sqrt{d}. \end{aligned}$$



**Term  $T_2$ .** We have

$$\begin{aligned}
T_2 &= \mathbb{E}_\beta \left[ \left( \sum_{s=0}^{\infty} \beta_s^\top \mathbf{Y}_{s,\mathbf{x}}^\top \right) \mathbf{Z} (\mathbf{Z}^\top \mathbf{Z} + \psi_1 \psi_2 \lambda \mathbf{I}_N)^{-1} \mathbf{U} (\mathbf{Z}^\top \mathbf{Z} + \psi_1 \psi_2 \lambda \mathbf{I}_N)^{-1} \mathbf{Z}^\top \left( \sum_{t=0}^{\infty} \mathbf{Y}_{t,\mathbf{x}} \beta_t \right) \right] / d \\
&= \mathbb{E}_\beta \left[ \text{Tr} \left( \left( \sum_{t=0}^{\infty} \mathbf{Y}_{t,\mathbf{x}} \beta_t \right) \left( \sum_{s=0}^{\infty} \beta_s^\top \mathbf{Y}_{s,\mathbf{x}}^\top \right) \mathbf{Z} (\mathbf{Z}^\top \mathbf{Z} + \psi_1 \psi_2 \lambda \mathbf{I}_N)^{-1} \mathbf{U} (\mathbf{Z}^\top \mathbf{Z} + \psi_1 \psi_2 \lambda \mathbf{I}_N)^{-1} \mathbf{Z}^\top \right) \right] / d \\
&= \sum_{k=0}^{\infty} F_{d,k}^2 \text{Tr} \left( \left( \mathbf{Y}_{k,\mathbf{x}} \mathbf{Y}_{k,\mathbf{x}}^\top / B(d, k) \right) \mathbf{Z} (\mathbf{Z}^\top \mathbf{Z} + \psi_1 \psi_2 \lambda \mathbf{I}_N)^{-1} \mathbf{U} (\mathbf{Z}^\top \mathbf{Z} + \psi_1 \psi_2 \lambda \mathbf{I}_N)^{-1} \mathbf{Z}^\top \right) / d \\
&= \sum_{k=0}^{\infty} F_{d,k}^2 \cdot \text{Tr} \left[ (\mathbf{Z}^\top \mathbf{Z} + \psi_1 \psi_2 \lambda \mathbf{I}_N)^{-1} \mathbf{U} (\mathbf{Z}^\top \mathbf{Z} + \psi_1 \psi_2 \lambda \mathbf{I}_N)^{-1} \mathbf{Z}^\top Q_k (\mathbf{X} \mathbf{X}^\top) \mathbf{Z} \right] / d.
\end{aligned}$$

**Term  $T_3$ .** By the assumption that  $\varepsilon_i \sim_{iid} \mathbb{P}_\varepsilon$  with  $\mathbb{E}_\varepsilon(\varepsilon) = 0$  and  $\mathbb{E}_\varepsilon(\varepsilon_1^2) = \tau^2$ , we have

$$\begin{aligned}
T_3 &= \mathbb{E}_\varepsilon \left[ \text{Tr}(\varepsilon \varepsilon^\top (\mathbf{Z}^\top \mathbf{Z} + \psi_1 \psi_2 \lambda \mathbf{I}_N)^{-1} \mathbf{U} (\mathbf{Z}^\top \mathbf{Z} + \psi_1 \psi_2 \lambda \mathbf{I}_N)^{-1} \mathbf{Z}^\top \mathbf{Z}) \right] / d \\
&= \tau^2 \cdot \text{Tr} \left[ (\mathbf{Z}^\top \mathbf{Z} + \psi_1 \psi_2 \lambda \mathbf{I}_N)^{-1} \mathbf{U} (\mathbf{Z}^\top \mathbf{Z} + \psi_1 \psi_2 \lambda \mathbf{I}_N)^{-1} \mathbf{Z}^\top \mathbf{Z} \right] / d.
\end{aligned}$$

Combining these terms proves the lemma.

## B.2 Proof of Lemma B.2

We state two lemmas that are used to prove Lemma B.2 and Lemma B.3. Their proofs are given in Section B.4.

**Lemma B.4.** *Use the same definitions and assumptions as Proposition 7.1 and Lemma B.1. Define*

$$\begin{aligned}
A &= 1 - 2A_1 + A_2, \\
A_1 &= \frac{\lambda_{d,0}(\sigma)}{\sqrt{d}} \text{Tr} \left[ \mathbf{1}_N \mathbf{1}_n^\top \mathbf{Z} (\mathbf{Z}^\top \mathbf{Z} + \psi_1 \psi_2 \lambda \mathbf{I}_N)^{-1} \right], \\
A_2 &= \frac{\lambda_{d,0}(\sigma)^2}{d} \text{Tr} \left[ (\mathbf{Z}^\top \mathbf{Z} + \psi_1 \psi_2 \lambda \mathbf{I}_N)^{-1} \mathbf{1}_N \mathbf{1}_N^\top (\mathbf{Z}^\top \mathbf{Z} + \psi_1 \psi_2 \lambda \mathbf{I}_N)^{-1} \mathbf{Z}^\top \mathbf{1}_n \mathbf{1}_n^\top \mathbf{Z} \right].
\end{aligned}$$

Then for any  $\lambda > 0$ , we have

$$\mathbb{E}|A| = o_d(1).$$

**Lemma B.5.** *Use the same definitions and assumptions as Proposition 7.1 and Lemma B.1. Let  $(\mathbf{M}_\alpha)_{\alpha \in \mathcal{A}} \in \mathbb{R}^{n \times n}$  be a collection of symmetric matrices with  $\mathbb{E}[\sup_{\alpha \in \mathcal{A}} \|\mathbf{M}_\alpha\|_{\text{op}}^2]^{1/2} = O_d(1)$ . Define*

$$B_\alpha = \frac{\lambda_{d,0}(\sigma)^2}{d} \text{Tr} \left[ (\mathbf{Z}^\top \mathbf{Z} + \psi_1 \psi_2 \lambda \mathbf{I}_N)^{-1} \mathbf{1}_N \mathbf{1}_N^\top (\mathbf{Z}^\top \mathbf{Z} + \psi_1 \psi_2 \lambda \mathbf{I}_N)^{-1} \mathbf{Z}^\top \mathbf{M}_\alpha \mathbf{Z} \right]. \quad (100)$$

Then for any  $\lambda > 0$ , we have

$$\mathbb{E} \left[ \sup_{\alpha \in \mathcal{A}} |B_\alpha| \right] = o_d(1).$$

Since  $\lambda > 0$ , there exists a constant  $C < \infty$  depending on  $(\lambda, \psi_1, \psi_2)$  such that deterministically

$$\|\mathbf{Z} (\mathbf{Z}^\top \mathbf{Z} + \psi_1 \psi_2 \lambda \mathbf{I}_N)^{-1}\|_{\text{op}}, \|(\mathbf{Z}^\top \mathbf{Z} + \psi_1 \psi_2 \lambda \mathbf{I}_N)^{-1}\|_{\text{op}} \leq C.$$

By the property of Wishart matrices [AGZ09], we have (the definition of these matrices are given in Eq. (56))

$$\mathbb{E}[\|\mathbf{H}\|_{\text{op}}^2], \mathbb{E}[\|\mathbf{Q}\|_{\text{op}}^2], \mathbb{E}[\|\mathbf{Z}_1\|_{\text{op}}^2] = O_d(1).$$

These bounds are crucial in proving this lemma. In the following, we bound each terms one by one.

**Step 1. The term  $|1 - 2S_{10} + S_{20}|$ .** By Lemma B.10 given in section B.5, we can decompose

$$\mathbf{U} = \lambda_{d,0}^2 \mathbf{1}_N \mathbf{1}_N^\top + \mathbf{M},$$

with  $\mathbb{E}[\|\mathbf{M}\|_{\text{op}}^2] = O_d(1)$  (we have  $\mathbf{M} = \mu_1^2 \mathbf{Q} + \mu_*^2 (\mathbf{I}_N + \mathbf{\Delta})$  where  $\mathbb{E}[\|\mathbf{Q}\|_{\text{op}}^2] = O_d(1)$  and  $\mathbb{E}[\|\mathbf{I}_N + \mathbf{\Delta}\|_{\text{op}}^2] = O_d(1)$ ). Moreover, the terms  $S_{10}$  and  $S_{20}$  can be rewritten as

$$\begin{aligned} S_{10} &= \lambda_{d,0}(\sigma) \text{Tr}(\mathbf{1}_N \mathbf{1}_n^\top \mathbf{Z} (\mathbf{Z}^\top \mathbf{Z} + \psi_1 \psi_2 \lambda \mathbf{I}_N)^{-1}) / \sqrt{d}, \\ S_{20} &= \lambda_{d,0}(\sigma)^2 \text{Tr}((\mathbf{Z}^\top \mathbf{Z} + \psi_1 \psi_2 \lambda \mathbf{I}_N)^{-1} \mathbf{1}_N \mathbf{1}_N^\top (\mathbf{Z}^\top \mathbf{Z} + \psi_1 \psi_2 \lambda \mathbf{I}_N)^{-1} \mathbf{Z}^\top \mathbf{1}_n \mathbf{1}_n^\top \mathbf{Z}) / d \\ &\quad + \text{Tr}((\mathbf{Z}^\top \mathbf{Z} + \psi_1 \psi_2 \lambda \mathbf{I}_N)^{-1} \mathbf{M} (\mathbf{Z}^\top \mathbf{Z} + \psi_1 \psi_2 \lambda \mathbf{I}_N)^{-1} \mathbf{Z}^\top \mathbf{1}_n \mathbf{1}_n^\top \mathbf{Z}) / d \\ &= \lambda_{d,0}(\sigma)^2 \text{Tr}((\mathbf{Z}^\top \mathbf{Z} + \psi_1 \psi_2 \lambda \mathbf{I}_N)^{-1} \mathbf{1}_N \mathbf{1}_N^\top (\mathbf{Z}^\top \mathbf{Z} + \psi_1 \psi_2 \lambda \mathbf{I}_N)^{-1} \mathbf{Z}^\top \mathbf{1}_n \mathbf{1}_n^\top \mathbf{Z}) / d \\ &\quad + \text{Tr}((\mathbf{Z} \mathbf{Z}^\top + \psi_1 \psi_2 \lambda \mathbf{I}_n)^{-1} \mathbf{Z} \mathbf{M} \mathbf{Z}^\top (\mathbf{Z} \mathbf{Z}^\top + \psi_1 \psi_2 \lambda \mathbf{I}_n)^{-1} \mathbf{1}_n \mathbf{1}_n^\top) / d. \end{aligned}$$

The last equality holds because  $(\mathbf{Z}^\top \mathbf{Z} + \psi_1 \psi_2 \lambda \mathbf{I}_N)^{-1} \mathbf{Z}^\top = \mathbf{Z}^\top (\mathbf{Z} \mathbf{Z}^\top + \psi_1 \psi_2 \lambda \mathbf{I}_n)^{-1}$ . Define

$$A_1 = \lambda_{d,0}(\sigma) \text{Tr}(\mathbf{1}_N \mathbf{1}_n^\top \mathbf{Z} (\mathbf{Z}^\top \mathbf{Z} + \psi_1 \psi_2 \lambda \mathbf{I}_N)^{-1}) / \sqrt{d}, \quad (101)$$

$$A_2 = \lambda_{d,0}(\sigma)^2 \text{Tr}((\mathbf{Z}^\top \mathbf{Z} + \psi_1 \psi_2 \lambda \mathbf{I}_N)^{-1} \mathbf{1}_N \mathbf{1}_N^\top (\mathbf{Z}^\top \mathbf{Z} + \psi_1 \psi_2 \lambda \mathbf{I}_N)^{-1} \mathbf{Z}^\top \mathbf{1}_n \mathbf{1}_n^\top \mathbf{Z}) / d, \quad (102)$$

$$B = \text{Tr}((\mathbf{Z} \mathbf{Z}^\top + \psi_1 \psi_2 \lambda \mathbf{I}_n)^{-1} \mathbf{Z} \mathbf{M} \mathbf{Z}^\top (\mathbf{Z} \mathbf{Z}^\top + \psi_1 \psi_2 \lambda \mathbf{I}_n)^{-1} \mathbf{1}_n \mathbf{1}_n^\top) / d. \quad (103)$$

Then we have  $S_{10} = A_1$ ,  $S_{20} = A_2 + B$ , and

$$\mathbb{E}[|1 - 2S_{10} + S_{20}|] = \mathbb{E}[|1 - 2A_1 + A_2 + B|] \leq \mathbb{E}[|1 - 2A_1 + A_2|] + \mathbb{E}[|B|].$$

By Lemma B.4, we have

$$\mathbb{E}[|1 - 2A_1 + A_2|] = o_d(1).$$

Note  $\mathbb{E}[\|\mathbf{M}\|_{\text{op}}^2] = O_d(1)$ , by Lemma B.5 (when applying Lemma B.5, we change the role of  $N$  and  $n$ , and the role of  $\mathbf{\Theta}$  and  $\mathbf{X}$ ; this can be done because the role of  $\mathbf{\Theta}$  and  $\mathbf{X}$  is symmetric), we have

$$\mathbb{E}[|B|] = o_d(1).$$

This gives

$$\mathbb{E}[|1 - 2S_{10} + S_{20}|] = o_d(1).$$

**Step 2. The term  $\mathbb{E}[\sup_{k \geq 2} |S_{1k}|]$ .** Note that we have

$$\begin{aligned} \sup_{k \geq 2} |S_{1k}| &\leq \sup_{k \geq 2} \left[ |\sqrt{d} \lambda_{d,k}(\sigma)| \cdot \|Q_k(\mathbf{\Theta} \mathbf{X}^\top) \mathbf{Z} (\mathbf{Z}^\top \mathbf{Z} + \psi_1 \psi_2 \lambda \mathbf{I}_N)^{-1}\|_{\text{op}} \right] \\ &\leq \sup_{k \geq 2} \left[ |\sqrt{d} \lambda_{d,k}(\sigma)| \cdot \|Q_k(\mathbf{\Theta} \mathbf{X}^\top)\|_{\text{op}} \|\mathbf{Z} (\mathbf{Z}^\top \mathbf{Z} + \psi_1 \psi_2 \lambda \mathbf{I}_N)^{-1}\|_{\text{op}} \right] \\ &\leq \sup_{k \geq 2} \left[ C \cdot |\sqrt{d} \lambda_{d,k}(\sigma)| \cdot \|Q_k(\mathbf{\Theta} \mathbf{X}^\top)\|_{\text{op}} \right]. \end{aligned}$$

Further note  $\|\sigma\|_{L^2(\tau_d)}^2 = \sum_{k \geq 0} \lambda_{d,k}(\sigma)^2 B(d, k) = O_d(1)$ ,  $B(d, k) = \Theta(d^k)$ , and for fixed  $d$ ,  $B(d, k)$  is non-decreasing in  $k$  [GMMM19, Lemma 1]. Therefore

$$\sup_{k \geq 2} |\lambda_{d,k}(\sigma)| \leq \sup_{k \geq 2} \left[ \|\sigma\|_{L^2(\tau_d)} / \sqrt{B(d, k)} \right] = O_d(1/d).$$

Moreover, by Lemma B.9, and note that  $Q_k(\mathbf{\Theta} \mathbf{X}^\top)$  is a sub-matrix of  $Q_k(\mathbf{W} \mathbf{W}^\top)$  for  $\mathbf{W}^\top = [\mathbf{\Theta}^\top, \mathbf{X}^\top]$ , we have

$$\mathbb{E} \left[ \sup_{k \geq 2} \|Q_k(\mathbf{\Theta} \mathbf{X}^\top)\|_{\text{op}}^2 \right] = o_d(1).$$

This proves

$$\mathbb{E} \left[ \sup_{k \geq 2} |S_{1k}| \right] = o_d(1).$$

**Step 3. The term  $\mathbb{E}[\sup_{k \geq 2} |S_{2k} - S_3|]$ .** As discussed above, we can decompose  $\mathbf{U} = \lambda_{d,0}^2 \mathbf{1}_N \mathbf{1}_N^\top + \mathbf{M}$ , with  $\mathbb{E}[\|\mathbf{M}\|_{\text{op}}^2] = O_d(1)$ . Hence we can bound

$$\sup_{k \geq 2} |S_{2k} - S_3| \leq I_1 + I_2,$$

where

$$\begin{aligned} I_1 &= \sup_{k \geq 2} \left| \frac{\lambda_{d,0}^2}{d} \text{Tr} \left[ (\mathbf{Z}^\top \mathbf{Z} + \psi_1 \psi_2 \lambda \mathbf{I}_N)^{-1} \mathbf{1}_N \mathbf{1}_N^\top (\mathbf{Z}^\top \mathbf{Z} + \psi_1 \psi_2 \lambda \mathbf{I}_N)^{-1} \mathbf{Z}^\top (Q_k(\mathbf{X} \mathbf{X}^\top) - \mathbf{I}_N) \mathbf{Z} \right] \right|, \\ I_2 &= \sup_{k \geq 2} \left| \frac{\lambda_{d,0}^2}{d} \text{Tr} \left[ (\mathbf{Z}^\top \mathbf{Z} + \psi_1 \psi_2 \lambda \mathbf{I}_N)^{-1} \mathbf{M} (\mathbf{Z}^\top \mathbf{Z} + \psi_1 \psi_2 \lambda \mathbf{I}_N)^{-1} \mathbf{Z}^\top (Q_k(\mathbf{X} \mathbf{X}^\top) - \mathbf{I}_N) \mathbf{Z} \right] \right|. \end{aligned}$$

By Lemma B.9 we have  $\mathbb{E}[\sup_{k \geq 2} \|Q_k(\mathbf{X} \mathbf{X}^\top) - \mathbf{I}_N\|_{\text{op}}^2] = o_d(1)$ , and by Lemma B.5, we get

$$\mathbb{E}[|I_1|] = o_d(1).$$

Moreover, we have

$$\begin{aligned} \mathbb{E}[|I_2|] &\leq \mathbb{E} \left[ \sup_{k \geq 2} \left| \lambda_{d,0}^2 \text{Tr}((\mathbf{Z}^\top \mathbf{Z} + \psi_1 \psi_2 \lambda \mathbf{I}_N)^{-1} \mathbf{M} (\mathbf{Z}^\top \mathbf{Z} + \psi_1 \psi_2 \lambda \mathbf{I}_N)^{-1} \mathbf{Z}^\top (Q_k(\mathbf{X} \mathbf{X}^\top) - \mathbf{I}_N) \mathbf{Z}) / d \right| \right] \\ &\leq \mathbb{E} \left[ \sup_{k \geq 2} \lambda_{d,0}^2 \|\mathbf{Z}(\mathbf{Z}^\top \mathbf{Z} + \psi_1 \psi_2 \lambda \mathbf{I}_N)^{-1}\|_{\text{op}} \|\mathbf{M}\|_{\text{op}} \|(\mathbf{Z}^\top \mathbf{Z} + \psi_1 \psi_2 \lambda \mathbf{I}_N)^{-1} \mathbf{Z}^\top\|_{\text{op}} \|Q_k(\mathbf{X} \mathbf{X}^\top) - \mathbf{I}_N\|_{\text{op}} \right] \\ &\leq O_d(1) \cdot \mathbb{E}[\|\mathbf{M}\|_{\text{op}}^2]^{1/2} \cdot \mathbb{E} \left[ \sup_{k \geq 2} \|Q_k(\mathbf{X} \mathbf{X}^\top) - \mathbf{I}_N\|_{\text{op}}^2 \right]^{1/2} = o_d(1), \end{aligned}$$

and hence  $\mathbb{E}[\sup_{k \geq 2} |S_{2k} - S_3|] = o_d(1)$ .

**Step 4. The term  $\mathbb{E}|S_{11} - \Psi_1|$ .** This is direct by observing (see Eq. (94))

$$\lim_{d \rightarrow \infty} \sqrt{d} \lambda_{1,d}(\sigma) = \mu_1.$$

and (the definition of  $\mathbf{Z}_1$  is given in Eq. (56))

$$\mu_1 Q_1(\mathbf{X} \Theta^\top) = \mu_1 \mathbf{X} \Theta^\top / d = \mathbf{Z}_1,$$

**Step 5. The term  $\mathbb{E}|S_{21} - \Psi_2|$ .** Observing we have (see Eq. (56))

$$Q_1(\mathbf{X} \mathbf{X}^\top) = \mathbf{X} \mathbf{X}^\top / d = \mathbf{H},$$

By Lemma B.10 we have

$$\mathbf{U} = \lambda_{d,0}(\sigma)^2 \mathbf{1}_N \mathbf{1}_N^\top + \mu_1^2 \mathbf{Q} + \mu_\star^2 (\mathbf{I}_N + \Delta),$$

for  $\mathbb{E}[\|\Delta\|_{\text{op}}^2] = o_d(1)$ . Therefore

$$|S_{21} - \Psi_2| \leq I_3 + I_4,$$

where

$$\begin{aligned} I_3 &= \left| \frac{\lambda_{d,0}(\sigma)^2}{d} \text{Tr} \left[ (\mathbf{Z}^\top \mathbf{Z} + \psi_1 \psi_2 \lambda \mathbf{I}_N)^{-1} \mathbf{1}_N \mathbf{1}_N^\top (\mathbf{Z}^\top \mathbf{Z} + \psi_1 \psi_2 \lambda \mathbf{I}_N)^{-1} \mathbf{Z}^\top \mathbf{H} \mathbf{Z} \right] \right|, \\ I_4 &= \left| \frac{\mu_\star^2}{d} \text{Tr} \left[ (\mathbf{Z}^\top \mathbf{Z} + \psi_1 \psi_2 \lambda \mathbf{I}_N)^{-1} \Delta (\mathbf{Z}^\top \mathbf{Z} + \psi_1 \psi_2 \lambda \mathbf{I}_N)^{-1} \mathbf{Z}^\top \mathbf{H} \mathbf{Z} \right] \right|. \end{aligned}$$

By Lemma B.5 and noting that  $\mathbb{E}[\|\mathbf{H}\|_{\text{op}}^2] = O_d(1)$ , we have  $\mathbb{E}[|I_3|] = o_d(1)$ . Moreover, we have

$$\begin{aligned} \mathbb{E}[|I_4|] &\leq \mathbb{E} \left[ \mu_\star^2 \|\mathbf{Z}(\mathbf{Z}^\top \mathbf{Z} + \psi_1 \psi_2 \lambda \mathbf{I}_N)^{-1}\|_{\text{op}} \|\Delta\|_{\text{op}} \|(\mathbf{Z}^\top \mathbf{Z} + \psi_1 \psi_2 \lambda \mathbf{I}_N)^{-1} \mathbf{Z}^\top\|_{\text{op}} \|\mathbf{H}\|_{\text{op}} \right] \\ &\leq O_d(1) \cdot \mathbb{E}[\|\Delta\|_{\text{op}}^2] \cdot \mathbb{E}[\|\mathbf{H}\|_{\text{op}}^2] = o_d(1). \end{aligned}$$

**Step 6. The term  $\mathbb{E}|S_3 - \Psi_3|$ .** This term can be dealt with similarly to the term  $\mathbb{E}|S_{21} - \Psi_2|$ .

### B.3 Proof of Lemma B.3

Instead of assuming  $\beta_{d,1} \sim \text{Unif}(\mathbb{S}^{d-1}(F_{d,1}))$  as per Assumption 4, in the proof we will assume  $\beta_{d,1} \sim \mathcal{N}(\mathbf{0}, [F_{d,1}^2/d]\mathbf{I}_d)$ . Note for  $\beta_{d,1} \sim \mathcal{N}(\mathbf{0}, [F_{d,1}^2/d]\mathbf{I}_d)$ , we have  $F_{d,1}\beta_{d,1}/\|\beta_{d,1}\|_2 \sim \text{Unif}(\mathbb{S}^{d-1}(F_{d,1}))$ . Moreover, in high dimension,  $\|\beta_{d,1}\|_2$  concentrates tightly around  $F_{d,1}$ . Using these properties, it is not hard to translate the proof from Gaussian  $\beta_{d,1}$  to spherical  $\beta_{d,1}$ .

First we state a lemma that simplifies our calculations.

**Lemma B.6.** *Let  $\mathbf{A} \in \mathbb{R}^{n \times N}$  and  $\mathbf{B} \in \mathbb{R}^{n \times n}$ . Let  $\mathbf{g} = (g_1, \dots, g_n)^\top$  with  $g_i \sim_{iid} \mathbb{P}_g$ ,  $\mathbb{E}_g[g] = 0$ , and  $\mathbb{E}_g[g^2] = 1$ . Let  $\mathbf{h} = (h_1, \dots, h_N)^\top$  with  $h_i \sim_{iid} \mathbb{P}_h$ ,  $\mathbb{E}_h[h] = 0$ , and  $\mathbb{E}_h[h^2] = 1$ . Then we have*

$$\begin{aligned}\text{Var}(\mathbf{g}^\top \mathbf{A} \mathbf{h}) &= \|\mathbf{A}\|_F^2, \\ \text{Var}(\mathbf{g}^\top \mathbf{B} \mathbf{g}) &= \sum_{i=1}^n B_{ii}^2 (\mathbb{E}[g^4] - 3) + \|\mathbf{B}\|_F^2 + \text{Tr}(\mathbf{B}^2).\end{aligned}$$

*Proof of Lemma B.6.*

**Step 1. Term  $\mathbf{g}^\top \mathbf{A} \mathbf{h}$ .** Calculating the expectation, we have

$$\mathbb{E}[\mathbf{g}^\top \mathbf{A} \mathbf{h}] = 0.$$

Hence we have

$$\text{Var}(\mathbf{g}^\top \mathbf{A} \mathbf{h}) = \mathbb{E}[\mathbf{g}^\top \mathbf{A} \mathbf{h} \mathbf{h}^\top \mathbf{A}^\top \mathbf{g}] = \mathbb{E}[\text{Tr}(\mathbf{g} \mathbf{g}^\top \mathbf{A} \mathbf{h} \mathbf{h}^\top \mathbf{A}^\top)] = \text{Tr}(\mathbf{A} \mathbf{A}^\top) = \|\mathbf{A}\|_F^2.$$

**Step 2. Term  $\mathbf{g}^\top \mathbf{B} \mathbf{g}$ .** Calculating the expectation, we have

$$\mathbb{E}[\mathbf{g}^\top \mathbf{B} \mathbf{g}] = \mathbb{E}[\text{Tr}(\mathbf{B} \mathbf{g} \mathbf{g}^\top)] = \text{Tr}(\mathbf{B}).$$

Hence we have

$$\begin{aligned}\text{Var}(\mathbf{g}^\top \mathbf{B} \mathbf{g}) &= \left\{ \sum_{i_1, i_2, i_3, i_4} \mathbb{E}[g_{i_1} B_{i_1 i_2} g_{i_2} g_{i_3} B_{i_3 i_4} g_{i_4}] \right\} - \text{Tr}(\mathbf{B})^2 \\ &= \left\{ \left( \sum_{i_1=i_2=i_3=i_4} + \sum_{i_1=i_2 \neq i_3=i_4} + \sum_{i_1=i_3 \neq i_2=i_4} + \sum_{i_1=i_4 \neq i_2=i_3} \right) \mathbb{E}[g_{i_1} B_{i_1 i_2} g_{i_2} g_{i_3} B_{i_3 i_4} g_{i_4}] \right\} - \text{Tr}(\mathbf{B})^2 \\ &= \sum_{i=1}^n B_{ii}^2 \mathbb{E}[g^4] + \sum_{i \neq j} B_{ii} B_{jj} + \sum_{i \neq j} (B_{ij} B_{ij} + B_{ij} B_{ji}) - \text{Tr}(\mathbf{B})^2 \\ &= \sum_{i=1}^n B_{ii}^2 (\mathbb{E}[g^4] - 3) + \text{Tr}(\mathbf{B}^\top \mathbf{B}) + \text{Tr}(\mathbf{B}^2).\end{aligned}$$

This proves the lemma. □

We can rewrite the prediction risk to be

$$R_{\text{RF}}(f_d, \mathbf{X}, \boldsymbol{\Theta}, \lambda) = \sum_{k \geq 0} F_{d,k}^2 - 2\Gamma_1 + \Gamma_2 + \Gamma_3 - 2\Gamma_4 + 2\Gamma_5,$$

where

$$\begin{aligned}\Gamma_1 &= \mathbf{f}^\top \mathbf{Z} (\mathbf{Z}^\top \mathbf{Z} + \psi_1 \psi_2 \lambda \mathbf{I}_N)^{-1} \mathbf{V} / \sqrt{d}, \\ \Gamma_2 &= \mathbf{f}^\top \mathbf{Z} (\mathbf{Z}^\top \mathbf{Z} + \psi_1 \psi_2 \lambda \mathbf{I}_N)^{-1} \mathbf{U} (\mathbf{Z}^\top \mathbf{Z} + \psi_1 \psi_2 \lambda \mathbf{I}_N)^{-1} \mathbf{Z}^\top \mathbf{f} / d, \\ \Gamma_3 &= \boldsymbol{\varepsilon}^\top \mathbf{Z} (\mathbf{Z}^\top \mathbf{Z} + \psi_1 \psi_2 \lambda \mathbf{I}_N)^{-1} \mathbf{U} (\mathbf{Z}^\top \mathbf{Z} + \psi_1 \psi_2 \lambda \mathbf{I}_N)^{-1} \mathbf{Z}^\top \boldsymbol{\varepsilon} / d, \\ \Gamma_4 &= \boldsymbol{\varepsilon}^\top \mathbf{Z} (\mathbf{Z}^\top \mathbf{Z} + \psi_1 \psi_2 \lambda \mathbf{I}_N)^{-1} \mathbf{V} / \sqrt{d}, \\ \Gamma_5 &= \boldsymbol{\varepsilon}^\top \mathbf{Z} (\mathbf{Z}^\top \mathbf{Z} + \psi_1 \psi_2 \lambda \mathbf{I}_N)^{-1} \mathbf{U} (\mathbf{Z}^\top \mathbf{Z} + \psi_1 \psi_2 \lambda \mathbf{I}_N)^{-1} \mathbf{Z}^\top \mathbf{f} / d,\end{aligned}$$

and  $\mathbf{V} = (V_1, \dots, V_N)^\top \in \mathbb{R}^N$ , and  $\mathbf{U} = (U_{ij})_{ij \in [N]} \in \mathbb{R}^{N \times N}$ , with

$$\begin{aligned} V_i &= \mathbb{E}_{\mathbf{x}}[f_d(\mathbf{x})\sigma(\langle \boldsymbol{\theta}_i, \mathbf{x} \rangle / \sqrt{d})], \\ U_{ij} &= \mathbb{E}_{\mathbf{x}}[\sigma(\langle \boldsymbol{\theta}_i, \mathbf{x} \rangle / \sqrt{d})\sigma(\langle \boldsymbol{\theta}_j, \mathbf{x} \rangle / \sqrt{d})]. \end{aligned}$$

To show Eq. (98), we just need to show that, for  $k \in [5]$ , we have

$$\mathbb{E}_{\mathbf{X}, \boldsymbol{\Theta}}[\text{Var}_{\boldsymbol{\beta}, \boldsymbol{\epsilon}}(\Gamma_k)^{1/2}] = o_d(1).$$

In the following, we show the variance bound for  $\Gamma_1$ . The other terms can be dealt with similarly.

Denote  $\mathbf{Y}_{k, \mathbf{x}} = (Y_{kl}(\mathbf{x}_i))_{i \in [n], l \in [B(d, k)]} \in \mathbb{R}^{n \times B(d, k)}$  and  $\mathbf{Y}_{k, \boldsymbol{\theta}} = (Y_{kl}(\boldsymbol{\theta}_a))_{a \in [N], l \in [B(d, k)]} \in \mathbb{R}^{N \times B(d, k)}$  where  $(Y_{kl})_{k \geq 0, l \in [B(d, k)]}$  is the set of spherical harmonics with domain  $\mathbb{S}^{d-1}(\sqrt{d})$ . Denote  $\lambda_{d, k} = \lambda_{d, k}(\sigma)$  as the Gegenbauer coefficients of  $\sigma$ , i.e.,

$$\sigma(x) = \sum_{k=0}^{\infty} \lambda_{d, k}(\sigma) B(d, k) Q_k^{(d)}(\sqrt{d}x).$$

Then we have

$$\mathbf{f} = \sum_{k=0}^{\infty} \mathbf{Y}_{k, \mathbf{x}} \boldsymbol{\beta}_k, \quad \mathbf{V} = \sum_{k=0}^{\infty} \lambda_{d, k} \mathbf{Y}_{k, \boldsymbol{\theta}} \boldsymbol{\beta}_k. \quad (104)$$

Denote

$$\mathbf{R}_1 = \mathbf{Z}(\mathbf{Z}^\top \mathbf{Z} + \psi_1 \psi_2 \lambda \mathbf{I}_N)^{-1}.$$

Then

$$\Gamma_1 = \frac{1}{\sqrt{d}} \left( \sum_{k=0}^{\infty} \mathbf{Y}_{k, \mathbf{x}} \boldsymbol{\beta}_{d, k} \right)^\top \mathbf{R}_1 \left( \sum_{l=0}^{\infty} \lambda_{d, l} \mathbf{Y}_{l, \boldsymbol{\theta}} \boldsymbol{\beta}_{d, l} \right).$$

Calculating the variance of  $\Gamma_1$  with respect to  $\boldsymbol{\beta}_{d, k} \sim \mathbf{N}(\mathbf{0}, (F_{d, k}^2 / B(d, k)) \mathbf{I})$  for  $k \geq 1$  using Lemma B.6, we get

$$\begin{aligned} \text{Var}_{\boldsymbol{\beta}}(\Gamma_1) &= \frac{1}{d} \text{Var}_{\boldsymbol{\beta}} \left( \left( \sum_{k=0}^{\infty} \mathbf{Y}_{k, \mathbf{x}} \boldsymbol{\beta}_{d, k} \right)^\top \mathbf{R}_1 \left( \sum_{l=0}^{\infty} \lambda_{d, l} \mathbf{Y}_{l, \boldsymbol{\theta}} \boldsymbol{\beta}_{d, l} \right) \right) \\ &= \sum_{l \neq k} \frac{\lambda_{d, l}^2}{d} \text{Var}_{\boldsymbol{\beta}} \left( \boldsymbol{\beta}_{d, k}^\top \mathbf{Y}_{k, \mathbf{x}}^\top \mathbf{R}_1 \mathbf{Y}_{l, \boldsymbol{\theta}} \boldsymbol{\beta}_{d, l} \right) + \sum_{k \geq 1} \frac{\lambda_{d, k}^2}{d} \text{Var}_{\boldsymbol{\beta}} \left( \boldsymbol{\beta}_{d, k}^\top \mathbf{Y}_{k, \mathbf{x}}^\top \mathbf{R}_1 \mathbf{Y}_{k, \boldsymbol{\theta}} \boldsymbol{\beta}_{d, k} \right) \\ &= \sum_{l \neq k} F_{d, l}^2 F_{d, k}^2 \frac{\lambda_{d, l}^2}{d} \text{Tr} \left( \mathbf{R}_1^\top Q_k(\mathbf{X} \mathbf{X}^\top) \mathbf{R}_1 Q_l(\boldsymbol{\Theta} \boldsymbol{\Theta}^\top) \right) \\ &\quad + \sum_{k \geq 1} F_k^4 \frac{\lambda_{d, k}^2}{d} \left[ \text{Tr} \left( \mathbf{R}_1^\top Q_k(\mathbf{X} \mathbf{X}^\top) \mathbf{R}_1 Q_k(\boldsymbol{\Theta} \boldsymbol{\Theta}^\top) \right) + \text{Tr} \left( \mathbf{R}_1 Q_k(\boldsymbol{\Theta} \mathbf{X}^\top) \mathbf{R}_1 Q_k(\boldsymbol{\Theta} \mathbf{X}^\top) \right) \right]. \end{aligned}$$

We claim the following equality holds.

$$\sup_{k, l \geq 1} \mathbb{E}_{\mathbf{X}, \boldsymbol{\Theta}} \left| \frac{\lambda_{d, l}^2}{d} \text{Tr} \left( \mathbf{R}_1^\top Q_k(\mathbf{X} \mathbf{X}^\top) \mathbf{R}_1 Q_l(\boldsymbol{\Theta} \boldsymbol{\Theta}^\top) \right) \right| = o_d(1), \quad (105)$$

$$\sup_{k, l \geq 1} \mathbb{E}_{\mathbf{X}, \boldsymbol{\Theta}} \left| \frac{\lambda_{d, l}^2}{d} \text{Tr} \left( \mathbf{R}_1 Q_k(\boldsymbol{\Theta} \mathbf{X}^\top) \mathbf{R}_1 Q_l(\boldsymbol{\Theta} \mathbf{X}^\top) \right) \right| = o_d(1), \quad (106)$$

$$\sup_{k \geq 1} \mathbb{E}_{\mathbf{X}, \boldsymbol{\Theta}} \left| \frac{1}{d} \text{Tr} \left( \mathbf{R}_1^\top \mathbf{1}_n \mathbf{1}_n^\top \mathbf{R}_1 Q_k(\boldsymbol{\Theta} \boldsymbol{\Theta}^\top) \right) \right| = o_d(1), \quad (107)$$

$$\sup_{k \geq 1} \mathbb{E}_{\mathbf{X}, \boldsymbol{\Theta}} \left| \frac{1}{d} \text{Tr} \left( \mathbf{R}_1^\top Q_k(\mathbf{X} \mathbf{X}^\top) \mathbf{R}_1 \mathbf{1}_N \mathbf{1}_N^\top \right) \right| = o_d(1). \quad (108)$$

Assuming these equality holds, noticing  $\sup_{l \geq 0} \lambda_{d,l}^2 \leq \|\sigma\|_{L^2(\tau_d)}^2 = O_d(1)$ ,  $Q_0(\mathbf{X}\mathbf{X}^\top) = \mathbf{1}_n \mathbf{1}_n^\top$  and  $Q_0(\boldsymbol{\Theta}\boldsymbol{\Theta}^\top) = \mathbf{1}_N \mathbf{1}_N^\top$ , we have

$$\mathbb{E}_{\mathbf{X}, \boldsymbol{\Theta}}[\text{Var}_{\boldsymbol{\beta}}(\Gamma_1)] = \sum_{l \neq k} F_{d,l}^2 F_{d,k}^2 \cdot o_d(1) + \sum_{k \geq 1} F_k^4 \cdot o_d(1) = \left[ \sum_{k \geq 0} F_{d,k}^2 \right]^2 \cdot o_d(1) = o_d(1).$$

In the following, we prove claims Eq. (105) (106) (107) and (108).

**Show Eq. (105) and (106).** Note we have almost surely  $\|\mathbf{R}_1\|_{\text{op}} \leq C$  for some constant  $C$ . Hence we have

$$\begin{aligned} & \sup_{k, l \geq 1} \mathbb{E}_{\mathbf{X}, \boldsymbol{\Theta}} \left| \frac{\lambda_{d,l}^2}{d} \text{Tr} \left( \mathbf{R}_1^\top Q_k(\mathbf{X}\mathbf{X}^\top) \mathbf{R}_1 Q_l(\boldsymbol{\Theta}\boldsymbol{\Theta}^\top) \right) \right| \\ & \leq O_d(1) \cdot \left[ \sup_{l \geq 1} \lambda_{d,l}^2 \right] \cdot \sup_{k, l \geq 1} \left\{ \mathbb{E}[\|Q_k(\mathbf{X}\mathbf{X}^\top)\|_{\text{op}}^2]^{1/2} \mathbb{E}[\|Q_l(\boldsymbol{\Theta}\boldsymbol{\Theta}^\top)\|_{\text{op}}^2]^{1/2} \right\}. \end{aligned} \quad (109)$$

Note  $\|\sigma\|_{L^2(\tau_d)}^2 = \sum_{k \geq 0} \lambda_{d,k}^2 B(d, k) = O_d(1)$ ,  $B(d, k) = \Theta(d^k)$ , and for fixed  $d$ ,  $B(d, k)$  is non-decreasing in  $k$  [GMMM19, Lemma 1]. Therefore

$$\sup_{k \geq 1} |\lambda_{d,k}(\sigma)| \leq \sup_{k \geq 1} \left[ \|\sigma\|_{L^2(\tau_d)} / \sqrt{B(d, k)} \right] = O_d(1/\sqrt{d}).$$

Moreover, by Lemma B.9 and the operator norm bound for Wishart matrices [AGZ09], we have

$$\sup_{k \geq 1} \left\{ \mathbb{E}[\|Q_k(\mathbf{X}\mathbf{X}^\top)\|_{\text{op}}^2]^{1/2} \right\} = O_d(1). \quad (110)$$

Plugging these bound into Eq. (109), we get Eq. (105). The proof of Eq. (106) is the same as Eq. (105).

**Show Eq. (107) and (108).** Note we have

$$\begin{aligned} & \sup_{k \geq 1} \mathbb{E}_{\mathbf{X}, \boldsymbol{\Theta}} \left| \frac{1}{d} \text{Tr} \left( \mathbf{R}_1^\top \mathbf{1}_n \mathbf{1}_n^\top \mathbf{R}_1 Q_k(\boldsymbol{\Theta}\boldsymbol{\Theta}^\top) \right) \right| \\ & = \sup_{k \geq 1} \mathbb{E}_{\mathbf{X}, \boldsymbol{\Theta}} \left| \frac{1}{d} \text{Tr} \left( (\mathbf{Z}\mathbf{Z}^\top + \psi_1 \psi_2 \lambda \mathbf{I}_n)^{-1} \mathbf{1}_n \mathbf{1}_n^\top (\mathbf{Z}\mathbf{Z}^\top + \psi_1 \psi_2 \lambda \mathbf{I}_n)^{-1} \mathbf{Z} Q_k(\boldsymbol{\Theta}\boldsymbol{\Theta}^\top) \mathbf{Z}^\top \right) \right|. \end{aligned}$$

Note  $\mathbb{E}[\sup_{k \geq 1} \|Q_k(\boldsymbol{\Theta}\boldsymbol{\Theta}^\top)\|_{\text{op}}^2] = O_d(1)$  (Lemma B.9) and  $\lambda_{d,0}(\sigma) = \Theta_d(1)$  (by Assumption 1 and note that  $\mu_0(\sigma) = \lim_{d \rightarrow \infty} \lambda_{d,0}(\sigma)$  by Eq. (94)), by Lemma B.5 (when applying Lemma B.5, we change the role of  $N$  and  $n$ , and the role of  $\boldsymbol{\Theta}$  and  $\mathbf{X}$ ; this can be done because the role of  $\boldsymbol{\Theta}$  and  $\mathbf{X}$  is symmetric), we get

$$\sup_{k \geq 1} \mathbb{E}_{\mathbf{X}, \boldsymbol{\Theta}} \left| \frac{1}{d} \text{Tr} \left( \mathbf{R}_1^\top \mathbf{1}_n \mathbf{1}_n^\top \mathbf{R}_1 Q_k(\boldsymbol{\Theta}\boldsymbol{\Theta}^\top) \right) \right| = o_d(1).$$

This proves Eq. (107). The proof of Eq. (108) is the same as the proof of Eq. (107). This concludes the proof.

## B.4 Proof of Lemma B.4 and B.5

To prove Lemma B.4 and B.5, first we states a lemma that reformulate  $A_1, A_2$  and  $B_\alpha$  using Sherman-Morrison-Woodbury formula.

**Lemma B.7** (Simplifications using Sherman-Morrison-Woodbury formula). *Use the same definitions and assumptions as Proposition 7.1 and Lemma B.1. For  $\mathbf{M} \in \mathbb{R}^{N \times N}$ , define*

$$L_1 = \frac{1}{\sqrt{d}} \lambda_{d,0}(\sigma) \text{Tr} \left[ \mathbf{1}_N \mathbf{1}_n^\top \mathbf{Z} (\mathbf{Z}^\top \mathbf{Z} + \psi_1 \psi_2 \lambda \mathbf{I}_N)^{-1} \right], \quad (111)$$

$$L_2(\mathbf{M}) = \frac{1}{d} \text{Tr} \left[ (\mathbf{Z}^\top \mathbf{Z} + \psi_1 \psi_2 \lambda \mathbf{I}_N)^{-1} \mathbf{M} (\mathbf{Z}^\top \mathbf{Z} + \psi_1 \psi_2 \lambda \mathbf{I}_N)^{-1} \mathbf{Z}^\top \mathbf{1}_n \mathbf{1}_n^\top \mathbf{Z} \right]. \quad (112)$$

We then have

$$L_1 = 1 - \frac{K_{12} + 1}{K_{11}(1 - K_{22}) + (K_{12} + 1)^2}, \quad (113)$$

$$L_2(\mathbf{M}) = \psi_2 \frac{G_{11}(1 - K_{22})^2 + G_{22}(K_{12} + 1)^2 + 2G_{12}(K_{12} + 1)(1 - K_{22})}{(K_{11}(1 - K_{22}) + (K_{12} + 1)^2)^2}, \quad (114)$$

where

$$\begin{aligned} K_{11} &= \mathbf{T}_1^\top \mathbf{E}_0^{-1} \mathbf{T}_1, \\ K_{12} &= \mathbf{T}_1^\top \mathbf{E}_0^{-1} \mathbf{T}_2, \\ K_{22} &= \mathbf{T}_2^\top \mathbf{E}_0^{-1} \mathbf{T}_2, \\ G_{11} &= \mathbf{T}_1^\top \mathbf{E}_0^{-1} \mathbf{M} \mathbf{E}_0^{-1} \mathbf{T}_1, \\ G_{12} &= \mathbf{T}_1^\top \mathbf{E}_0^{-1} \mathbf{M} \mathbf{E}_0^{-1} \mathbf{T}_2, \\ G_{22} &= \mathbf{T}_2^\top \mathbf{E}_0^{-1} \mathbf{M} \mathbf{E}_0^{-1} \mathbf{T}_2, \end{aligned}$$

and

$$\begin{aligned} \varphi_d(x) &= \sigma(x) - \lambda_{d,0}(\sigma), \\ \mathbf{J} &= \frac{1}{\sqrt{d}} \varphi_d \left( \frac{1}{\sqrt{d}} \mathbf{X} \boldsymbol{\Theta}^\top \right), \\ \mathbf{E}_0 &= \mathbf{J}^\top \mathbf{J} + \psi_1 \psi_2 \lambda \mathbf{I}_N, \\ \mathbf{T}_1 &= \psi_2^{1/2} \lambda_{d,0}(\sigma) \mathbf{1}_N, \\ \mathbf{T}_2 &= \frac{1}{\sqrt{n}} \mathbf{J}^\top \mathbf{1}_n. \end{aligned}$$

*Proof of Lemma B.7.*

**Step 1. Term  $L_1$ .**

Note we have (denoting  $\lambda_{d,0} = \lambda_{d,0}(\sigma)$ )

$$\mathbf{Z} = \lambda_{d,0} \mathbf{1}_n \mathbf{1}_N^\top / \sqrt{d} + \mathbf{J}.$$

Hence we have (denoting  $\mathbf{T}_2 = \mathbf{J}^\top \mathbf{1}_n / \sqrt{n}$ )

$$\begin{aligned} L_1 &= \text{Tr} \left[ \lambda_{d,0} \mathbf{1}_N \mathbf{1}_n^\top (\lambda_{d,0} \mathbf{1}_n \mathbf{1}_N^\top / \sqrt{d} + \mathbf{J}) [(\lambda_{d,0} \mathbf{1}_n \mathbf{1}_N^\top / \sqrt{d} + \mathbf{J})^\top (\lambda_{d,0} \mathbf{1}_n \mathbf{1}_N^\top / \sqrt{d} + \mathbf{J}) + \psi_1 \psi_2 \lambda \mathbf{I}_N]^{-1} \right] / \sqrt{d} \\ &= \text{Tr} \left[ (\psi_2 \lambda_{d,0}^2 \mathbf{1}_N \mathbf{1}_N^\top + \psi_2^{1/2} \lambda_{d,0} \mathbf{1}_N \mathbf{T}_2^\top) [\psi_2 \lambda_{d,0}^2 \mathbf{1}_N \mathbf{1}_N^\top + \psi_2^{1/2} \lambda_{d,0} \mathbf{1}_N \mathbf{T}_2^\top + \psi_2^{1/2} \lambda_{d,0} \mathbf{T}_2 \mathbf{1}_N^\top + \mathbf{J}^\top \mathbf{J} + \psi_1 \psi_2 \lambda \mathbf{I}_N]^{-1} \right]. \end{aligned}$$

Define

$$\begin{aligned} \mathbf{E} &= \mathbf{Z}^\top \mathbf{Z} + \psi_1 \psi_2 \lambda \mathbf{I}_N = \mathbf{E}_0 + \mathbf{F}_1 \mathbf{F}_2^\top, \\ \mathbf{E}_0 &= \mathbf{J}^\top \mathbf{J} + \psi_1 \psi_2 \lambda \mathbf{I}_N, \\ \mathbf{F}_1 &= (\mathbf{T}_1, \mathbf{T}_1, \mathbf{T}_2), \\ \mathbf{F}_2 &= (\mathbf{T}_1, \mathbf{T}_2, \mathbf{T}_1), \\ \mathbf{T}_1 &= \psi_2^{1/2} \lambda_{d,0} \mathbf{1}_N, \\ \mathbf{T}_2 &= \mathbf{J}^\top \mathbf{1}_n / \sqrt{n}. \end{aligned}$$

By the Sherman-Morrison-Woodbury formula, we have

$$\mathbf{E}^{-1} = \mathbf{E}_0^{-1} - \mathbf{E}_0^{-1} \mathbf{F}_1 (\mathbf{I}_3 + \mathbf{F}_2^\top \mathbf{E}_0^{-1} \mathbf{F}_1)^{-1} \mathbf{F}_2^\top \mathbf{E}_0^{-1}.$$



Then we have

$$\begin{aligned}
L_1 &= \text{Tr} \left[ (\mathbf{T}_1 \mathbf{T}_1^\top + \mathbf{T}_1 \mathbf{T}_2^\top) (\mathbf{E}_0^{-1} - \mathbf{E}_0^{-1} \mathbf{F}_1 (\mathbf{I}_3 + \mathbf{F}_2^\top \mathbf{E}_0^{-1} \mathbf{F}_1)^{-1} \mathbf{F}_2^\top \mathbf{E}_0^{-1}) \right] \\
&= (\mathbf{T}_1^\top \mathbf{E}_0^{-1} \mathbf{T}_1 - \mathbf{T}_1^\top \mathbf{E}_0^{-1} \mathbf{F}_1 (\mathbf{I}_3 + \mathbf{F}_2^\top \mathbf{E}_0^{-1} \mathbf{F}_1)^{-1} \mathbf{F}_2^\top \mathbf{E}_0^{-1} \mathbf{T}_1) \\
&\quad + (\mathbf{T}_2^\top \mathbf{E}_0^{-1} \mathbf{T}_1 - \mathbf{T}_2^\top \mathbf{E}_0^{-1} \mathbf{F}_1 (\mathbf{I}_3 + \mathbf{F}_2^\top \mathbf{E}_0^{-1} \mathbf{F}_1)^{-1} \mathbf{F}_2^\top \mathbf{E}_0^{-1} \mathbf{T}_1) \\
&= (K_{11} - [K_{11}, K_{11}, K_{12}] (\mathbf{I}_3 + \mathbf{K})^{-1} [K_{11}, K_{12}, K_{11}]^\top) \\
&\quad + (K_{12} - [K_{12}, K_{12}, K_{22}] (\mathbf{I}_3 + \mathbf{K})^{-1} [K_{11}, K_{12}, K_{11}]^\top) \\
&= [K_{11}, K_{11}, K_{12}] (\mathbf{I}_3 + \mathbf{K})^{-1} [1, 0, 0]^\top \\
&\quad + [K_{12}, K_{12}, K_{22}] (\mathbf{I}_3 + \mathbf{K})^{-1} [1, 0, 0]^\top \\
&= (K_{12}^2 + K_{12} + K_{11} - K_{11} K_{22}) / (K_{12}^2 + 2K_{12} + K_{11} - K_{11} K_{22} + 1) \\
&= 1 - (K_{12} + 1) / [K_{11}(1 - K_{22}) + (K_{12} + 1)^2],
\end{aligned}$$

where

$$\begin{aligned}
K_{11} &= \mathbf{T}_1^\top \mathbf{E}_0^{-1} \mathbf{T}_1 = \psi_2 \lambda_{d,0}^2 \mathbf{1}_N^\top (\mathbf{J}^\top \mathbf{J} + \psi_1 \psi_2 \lambda \mathbf{I}_N)^{-1} \mathbf{1}_N, \\
K_{12} &= \mathbf{T}_1^\top \mathbf{E}_0^{-1} \mathbf{T}_2 = \lambda_{d,0} \mathbf{1}_N^\top (\mathbf{J}^\top \mathbf{J} + \psi_1 \psi_2 \lambda \mathbf{I}_N)^{-1} \mathbf{J}^\top \mathbf{1}_n / \sqrt{d}, \\
K_{22} &= \mathbf{T}_2^\top \mathbf{E}_0^{-1} \mathbf{T}_2 = \mathbf{1}_n^\top \mathbf{J} (\mathbf{J}^\top \mathbf{J} + \psi_1 \psi_2 \lambda \mathbf{I}_N)^{-1} \mathbf{J}^\top \mathbf{1}_n / n, \\
\mathbf{K} &= \begin{bmatrix} K_{11} & K_{11} & K_{12} \\ K_{12} & K_{12} & K_{22} \\ K_{11} & K_{11} & K_{12} \end{bmatrix}.
\end{aligned}$$

This prove Eq. (113).

**Step 2. Term  $L_2(\mathbf{M})$ .** We have

$$\begin{aligned}
\mathbf{Z}^\top \mathbf{1}_n \mathbf{1}_n^\top \mathbf{Z} / d &= (\lambda_{d,0} \mathbf{1}_n \mathbf{1}_N^\top / \sqrt{d} + \mathbf{J})^\top \mathbf{1}_n \mathbf{1}_n^\top (\lambda_{d,0} \mathbf{1}_n \mathbf{1}_N^\top / \sqrt{d} + \mathbf{J}) / d \\
&= \psi_2^2 \lambda_{d,0}^2 \mathbf{1}_N \mathbf{1}_N^\top + \psi_2 \mathbf{T}_2 \cdot \sqrt{\psi_2} \lambda_{d,0} \mathbf{1}_N^\top + \psi_2 \sqrt{\psi_2} \lambda_{d,0} \mathbf{1}_N \mathbf{T}_2^\top + \psi_2 \mathbf{T}_2 \mathbf{T}_2^\top = \psi_2 (\mathbf{T}_1 + \mathbf{T}_2) (\mathbf{T}_1 + \mathbf{T}_2)^\top.
\end{aligned}$$

As a result, we have

$$\begin{aligned}
L_2(\mathbf{M}) &= \psi_2 \cdot (\mathbf{T}_1 + \mathbf{T}_2)^\top \mathbf{E}^{-1} \mathbf{M} \mathbf{E}^{-1} (\mathbf{T}_1 + \mathbf{T}_2) \\
&= \psi_2 \cdot (\mathbf{T}_1 + \mathbf{T}_2)^\top (\mathbf{I}_N - \mathbf{E}_0^{-1} \mathbf{F}_1 (\mathbf{I}_3 + \mathbf{F}_2^\top \mathbf{E}_0^{-1} \mathbf{F}_1)^{-1} \mathbf{F}_2^\top) \\
&\quad \cdot (\mathbf{E}_0^{-1} \mathbf{M} \mathbf{E}_0^{-1}) (\mathbf{I}_N - \mathbf{F}_2 (\mathbf{I}_3 + \mathbf{F}_1^\top \mathbf{E}_0^{-1} \mathbf{F}_2)^{-1} \mathbf{F}_1^\top \mathbf{E}_0^{-1}) (\mathbf{T}_1 + \mathbf{T}_2).
\end{aligned}$$

Simplifying this formula using simple algebra proves Eq. (114).  $\square$

*Proof of Lemma B.4.*

**Step 1. Term  $A_1$ .** By Lemma B.7, we get

$$A_1 = 1 - (K_{12} + 1) / (K_{11}(1 - K_{22}) + (K_{12} + 1)^2), \quad (115)$$

where

$$\begin{aligned}
K_{11} &= \mathbf{T}_1^\top \mathbf{E}_0^{-1} \mathbf{T}_1 = \psi_2 \lambda_{d,0}^2 \mathbf{1}_N^\top (\mathbf{J}^\top \mathbf{J} + \psi_1 \psi_2 \lambda \mathbf{I}_N)^{-1} \mathbf{1}_N, \\
K_{12} &= \mathbf{T}_1^\top \mathbf{E}_0^{-1} \mathbf{T}_2 = \lambda_{d,0} \mathbf{1}_N^\top (\mathbf{J}^\top \mathbf{J} + \psi_1 \psi_2 \lambda \mathbf{I}_N)^{-1} \mathbf{J}^\top \mathbf{1}_n / \sqrt{d}, \\
K_{22} &= \mathbf{T}_2^\top \mathbf{E}_0^{-1} \mathbf{T}_2 = \mathbf{1}_n^\top \mathbf{J} (\mathbf{J}^\top \mathbf{J} + \psi_1 \psi_2 \lambda \mathbf{I}_N)^{-1} \mathbf{J}^\top \mathbf{1}_n / n.
\end{aligned}$$

**Step 2. Term  $A_2$ .**

Note that we have

$$A_2 = \text{Tr}((\mathbf{Z}^\top \mathbf{Z} + \psi_1 \psi_2 \lambda \mathbf{I}_N)^{-1} \mathbf{U}_0 (\mathbf{Z}^\top \mathbf{Z} + \psi_1 \psi_2 \lambda \mathbf{I}_N)^{-1} \mathbf{Z}^\top \mathbf{1}_n \mathbf{1}_n^\top \mathbf{Z}) / d,$$

where

$$\mathbf{U}_0 = \lambda_{d,0}(\sigma)^2 \mathbf{1}_N \mathbf{1}_N^\top = \mathbf{T}_1 \mathbf{T}_1^\top / \psi_2. \quad (116)$$

By Lemma B.7, we have

$$A_2 = \psi_2 [G_{11}(1 - K_{22})^2 + G_{22}(K_{12} + 1)^2 + 2G_{12}(K_{12} + 1)(1 - K_{22})] / (K_{11}(1 - K_{22}) + (K_{12} + 1)^2)^2, \quad (117)$$

where

$$\begin{aligned} G_{11} &= \mathbf{T}_1^\top \mathbf{E}_0^{-1} \mathbf{U}_0 \mathbf{E}_0^{-1} \mathbf{T}_1 = K_{11}^2 / \psi_2, \\ G_{12} &= \mathbf{T}_1^\top \mathbf{E}_0^{-1} \mathbf{U}_0 \mathbf{E}_0^{-1} \mathbf{T}_2 = K_{11} K_{12} / \psi_2, \\ G_{22} &= \mathbf{T}_2^\top \mathbf{E}_0^{-1} \mathbf{U}_0 \mathbf{E}_0^{-1} \mathbf{T}_2 = K_{12}^2 / \psi_2. \end{aligned}$$

We can simplify  $S_{20}$  in Eq. (117) further, and get

$$A_2 = (K_{11}(1 - K_{22}) + K_{12}^2 + K_{12})^2 / (K_{11}(1 - K_{22}) + (K_{12} + 1)^2)^2. \quad (118)$$

### Step 3. Combining $A_1$ and $A_2$

By Eq. (115) and (118), we have

$$A = 1 - 2A_1 + A_2 = (K_{12} + 1)^2 / (K_{11}(1 - K_{22}) + (K_{12} + 1)^2)^2 \geq 0.$$

For term  $K_{12}$ , we have

$$|K_{12}| \leq \lambda_{d,0} \|(\mathbf{J}^\top \mathbf{J} + \psi_1 \psi_2 \lambda \mathbf{I}_N)^{-1} \mathbf{J}^\top\|_{\text{op}} \|\mathbf{1}_n \mathbf{1}_N^\top / \sqrt{d}\|_{\text{op}} = O_d(\sqrt{d}).$$

For term  $K_{11}$ , we have

$$K_{11} \geq \psi_2 \lambda_{d,0}^2 N \lambda_{\min}((\mathbf{J}^\top \mathbf{J} + \psi_1 \psi_2 \lambda \mathbf{I}_N)^{-1}) = \Omega_d(d) / (\|\mathbf{J}^\top \mathbf{J}\|_{\text{op}} + \psi_1 \psi_2 \lambda).$$

For term  $K_{22}$ , we have

$$\begin{aligned} 1 \geq 1 - K_{22} &= \mathbf{1}_n^\top (\mathbf{I}_n - \mathbf{J}(\mathbf{J}^\top \mathbf{J} + \psi_1 \psi_2 \lambda \mathbf{I}_N)^{-1} \mathbf{J}^\top) \mathbf{1}_n / n \geq 1 - \lambda_{\max}(\mathbf{J}(\mathbf{J}^\top \mathbf{J} + \psi_1 \psi_2 \lambda \mathbf{I}_N)^{-1} \mathbf{J}^\top) \\ &\geq \psi_1 \psi_2 \lambda / (\psi_1 \psi_2 \lambda + \|\mathbf{J}^\top \mathbf{J}\|_{\text{op}}) > 0. \end{aligned}$$

As a result, we have

$$1 / (K_{11}(1 - K_{22}) + (K_{12} + 1)^2)^2 = O_d(d^{-2}) \cdot (1 + \|\mathbf{J}\|_{\text{op}}^8),$$

and hence

$$A = O_d(1/d) \cdot (1 + \|\mathbf{J}\|_{\text{op}}^8)$$

Lemma B.13 in Section B.5 provides an upper bound on the operator norm of  $\|\mathbf{J}\|_{\text{op}}$ , which gives  $\|\mathbf{J}\|_{\text{op}} = O_{d,\mathbb{P}}(\exp\{C(\log d)^{1/2}\})$  (note  $\mathbf{J}$  can be regarded as a sub-matrix of  $\mathbf{K}$  in Lemma B.13, so that  $\|\mathbf{J}\|_{\text{op}} \leq \|\mathbf{K}\|_{\text{op}}$ ). Using this bound, we get

$$A = o_{d,\mathbb{P}}(1).$$

It is easy to see that  $0 \leq A \leq 1$ . Hence the high probability bound translates to an expectation bound. This proves the lemma.  $\square$

*Proof of Lemma B.5.* For notation simplicity, we prove this lemma under the case when  $\mathcal{A} = \{\alpha\}$  which is a singleton. We denote  $B = B_\alpha$ . The proof can be directly generalized to the case for arbitrary set  $\mathcal{A}$ .

By Lemma B.7 (when applying Lemma B.7, we change the role of  $N$  and  $n$ , and the role of  $\Theta$  and  $\mathbf{X}$ ; this can be done because the role of  $\Theta$  and  $\mathbf{X}$  is symmetric), we have

$$B = \psi_2 \frac{G_{11}(1 - K_{22})^2 + G_{22}(K_{12} + 1)^2 + 2G_{12}(K_{12} + 1)(1 - K_{22})}{(K_{11}(1 - K_{22}) + (K_{12} + 1)^2)^2}, \quad (119)$$

where

$$\begin{aligned} K_{11} &= \mathbf{T}_1^\top \mathbf{E}_0^{-1} \mathbf{T}_1 = \psi_2 \lambda_{d,0}(\sigma)^2 \mathbf{1}_N^\top (\mathbf{J}^\top \mathbf{J} + \psi_1 \psi_2 \lambda \mathbf{I}_N)^{-1} \mathbf{1}_N, \\ K_{12} &= \mathbf{T}_1^\top \mathbf{E}_0^{-1} \mathbf{T}_2 = \lambda_{d,0}(\sigma) \mathbf{1}_N^\top (\mathbf{J}^\top \mathbf{J} + \psi_1 \psi_2 \lambda \mathbf{I}_N)^{-1} \mathbf{J}^\top \mathbf{1}_n / \sqrt{d}, \\ K_{22} &= \mathbf{T}_2^\top \mathbf{E}_0^{-1} \mathbf{T}_2 = \mathbf{1}_n^\top \mathbf{J}(\mathbf{J}^\top \mathbf{J} + \psi_1 \psi_2 \lambda \mathbf{I}_N)^{-1} \mathbf{J}^\top \mathbf{1}_n / n, \\ G_{11} &= \mathbf{T}_1^\top \mathbf{E}_0^{-1} \mathbf{M} \mathbf{E}_0^{-1} \mathbf{T}_1 = \psi_2 \lambda_{d,0}(\sigma)^2 \mathbf{1}_N^\top (\mathbf{J}^\top \mathbf{J} + \psi_1 \psi_2 \lambda \mathbf{I}_N)^{-1} \mathbf{M}(\mathbf{J}^\top \mathbf{J} + \psi_1 \psi_2 \lambda \mathbf{I}_N)^{-1} \mathbf{1}_N, \\ G_{12} &= \mathbf{T}_1^\top \mathbf{E}_0^{-1} \mathbf{M} \mathbf{E}_0^{-1} \mathbf{T}_2 = \lambda_{d,0}(\sigma) \mathbf{1}_N^\top (\mathbf{J}^\top \mathbf{J} + \psi_1 \psi_2 \lambda \mathbf{I}_N)^{-1} \mathbf{M}(\mathbf{J}^\top \mathbf{J} + \psi_1 \psi_2 \lambda \mathbf{I}_N)^{-1} \mathbf{J}^\top \mathbf{1}_n / \sqrt{d}, \\ G_{22} &= \mathbf{T}_2^\top \mathbf{E}_0^{-1} \mathbf{M} \mathbf{E}_0^{-1} \mathbf{T}_2 = \mathbf{1}_n^\top \mathbf{J}(\mathbf{J}^\top \mathbf{J} + \psi_1 \psi_2 \lambda \mathbf{I}_N)^{-1} \mathbf{M}(\mathbf{J}^\top \mathbf{J} + \psi_1 \psi_2 \lambda \mathbf{I}_N)^{-1} \mathbf{J}^\top \mathbf{1}_n / n. \end{aligned}$$

Note we have shown in the proof of Lemma B.4 that

$$\begin{aligned} K_{11} &= \Omega_d(d)/(\psi_1\psi_2\lambda + \|\mathbf{J}\|_{\text{op}}^2), \\ K_{12} &= O_d(\sqrt{d}), \\ 1 &\geq 1 - K_{22} \geq \psi_1\psi_2\lambda/(\psi_1\psi_2\lambda + \|\mathbf{J}\|_{\text{op}}^2), \\ 1/(K_{11}(1 - K_{22}) + (K_{12} + 1)^2)^2 &= O_d(d^{-2}) \cdot (1 \vee \|\mathbf{J}\|_{\text{op}}^8). \end{aligned}$$

Lemma B.13 provides an upper bound on the operator norm of  $\|\mathbf{J}\|_{\text{op}}$ , which gives  $\|\mathbf{J}\|_{\text{op}} = O_{d,\mathbb{P}}(\exp\{C(\log d)^{1/2}\})$ . Using this bound, we get for any  $\varepsilon > 0$

$$\begin{aligned} (1 - K_{22})^2/(K_{11}(1 - K_{22}) + (K_{12} + 1)^2)^2 &= O_{d,\mathbb{P}}(d^{-2+\varepsilon}), \\ (K_{12} + 1)^2/(K_{11}(1 - K_{22}) + (K_{12} + 1)^2)^2 &= O_{d,\mathbb{P}}(d^{-1+\varepsilon}), \\ |(K_{12} + 1)(1 - K_{22})|/(K_{11}(1 - K_{22}) + (K_{12} + 1)^2)^2 &= O_{d,\mathbb{P}}(d^{-3/2+\varepsilon}). \end{aligned}$$

Since all the quantities above are deterministically bounded by a constant, these high probability bounds translate to expectation bounds.

Moreover, we have

$$\begin{aligned} \mathbb{E}[G_{11}^2]^{1/2} &\leq \psi_2\lambda_{d,0}(\sigma)^2(\psi_1\psi_2\lambda)^{-2}\mathbb{E}[\|\mathbf{M}\|_{\text{op}}^2]^{1/2}\|\mathbf{1}_N\mathbf{1}_N^\top\|_{\text{op}} = O_d(d), \\ \mathbb{E}[G_{22}^2]^{1/2} &\leq O_d(1) \cdot \mathbb{E}[\|\mathbf{M}\|_{\text{op}}^2]^{1/2}\|\mathbf{1}_n\mathbf{1}_n^\top/n\|_{\text{op}} = O_d(1), \\ \mathbb{E}[G_{12}^2]^{1/2} &\leq O_d(1) \cdot \lambda_{d,0}(\sigma)\mathbb{E}[\|\mathbf{M}\|_{\text{op}}^2]^{1/2}\|\mathbf{1}_n\mathbf{1}_N^\top/\sqrt{d}\|_{\text{op}} = O_d(d^{1/2}). \end{aligned}$$

Plugging in the above bounds into Equation (119), we have

$$\mathbb{E}[|B|] = o_d(1).$$

This proves the lemma.  $\square$

## B.5 Preliminary lemmas

We denote by  $\mu_d$  the probability law of  $\langle \mathbf{x}_1, \mathbf{x}_2 \rangle / \sqrt{d}$  when  $\mathbf{x}_1, \mathbf{x}_2 \sim_{iid} \mathbf{N}(\mathbf{0}, \mathbf{I}_d)$ . Note that  $\mu_d$  is symmetric, and  $\int x^2 \mu_d(dx) = 1$ . By the central limit theorem,  $\mu_d$  converges weakly to  $\mu_G$  as  $d \rightarrow \infty$ , where  $\mu_G$  is the standard Gaussian measure. In fact, we have the following stronger convergence result.

**Lemma B.8.** *For any  $\lambda \in [-\sqrt{d}/2, \sqrt{d}/2]$ , we have*

$$\int e^{\lambda x} \mu_d(dx) \leq e^{\lambda^2}. \quad (120)$$

Further, let  $f : \mathbb{R} \rightarrow \mathbb{R}$  be a continuous function such that  $|f(x)| \leq c_0 \exp(c_1|x|)$  for some constants  $c_0, c_1 < \infty$ . Then

$$\lim_{d \rightarrow \infty} \int f(x) \mu_d(dx) = \int f(x) \mu_G(dx). \quad (121)$$

*Proof of Lemma B.8.* In order to prove Eq. (120), we note that the left hand side is given by

$$\begin{aligned} \mathbb{E}\{e^{\lambda \langle \mathbf{x}_1, \mathbf{x}_2 \rangle / \sqrt{d}}\} &= \frac{1}{(2\pi)^d} \int \exp\left\{-\frac{1}{2}\|\mathbf{x}_1\|_2^2 - \frac{1}{2}\|\mathbf{x}_2\|_2^2 + \frac{\lambda}{\sqrt{d}}\langle \mathbf{x}_1, \mathbf{x}_2 \rangle\right\} d\mathbf{x}_1 d\mathbf{x}_2 \\ &= \left[\det\begin{pmatrix} 1 & -\lambda/\sqrt{d} \\ -\lambda/\sqrt{d} & 1 \end{pmatrix}\right]^{-d/2} = \left(1 - \frac{\lambda^2}{d}\right)^{-d/2} \\ &\leq e^{\lambda^2}, \end{aligned}$$

where the last inequality holds for  $|\lambda| \leq \sqrt{d}/2$  using the fact that  $(1 - x)^{-1} \leq e^{2x}$  for  $x \in [0, 1/4]$ .

In order to prove (121), let  $X_d \sim \mu_d$ , and  $G \sim \mathbf{N}(0, 1)$ . Since  $\mu_d$  converges weakly to  $\mathbf{N}(0, 1)$ , we can construct such random variables so that  $X_d \rightarrow G$  almost surely. Hence  $f(X_d) \rightarrow f(G)$  almost surely. However  $|f(X_d)| \leq c_0 \exp(c_1|X_d|)$  which is a uniformly integrable family by the previous point, implying  $\mathbb{E}f(X_d) \rightarrow \mathbb{E}f(G)$  as claimed.  $\square$

The next lemma is a reformulation of Proposition 3 in [GMMM19]. We present it in a stronger form, but it can be easily derived from the proof of Proposition 3 in [GMMM19]. This lemma was first proved in [EK10] in the Gaussian case.

**Lemma B.9.** *Let  $\Theta = (\theta_1, \dots, \theta_N)^\top \in \mathbb{R}^{N \times d}$  with  $(\theta_a)_{a \in [N]} \sim_{iid} \text{Unif}(\mathbb{S}^{d-1}(\sqrt{d}))$ . Assume  $1/c \leq N/d \leq c$  for some constant  $c \in (0, \infty)$ . Then*

$$\mathbb{E} \left[ \sup_{k \geq 2} \|Q_k(\Theta \Theta^\top) - \mathbf{I}_N\|_{\text{op}}^2 \right] = o_d(1).$$

The next lemma can be easily derived from Lemma B.9. Again, this lemma was first proved in [EK10] in the Gaussian case.

**Lemma B.10.** *Let  $\Theta = (\theta_1, \dots, \theta_N)^\top \in \mathbb{R}^{N \times d}$  with  $(\theta_a)_{a \in [N]} \sim_{iid} \text{Unif}(\mathbb{S}^{d-1}(\sqrt{d}))$ . Let activation function  $\sigma$  satisfies Assumption 1. Assume  $1/c \leq N/d \leq c$  for some constant  $c \in (0, \infty)$ . Denote*

$$U = \left( \mathbb{E}_{\mathbf{x} \sim \text{Unif}(\mathbb{S}^{d-1}(\sqrt{d}))} [\sigma(\langle \theta_a, \mathbf{x} \rangle / \sqrt{d}) \sigma(\langle \theta_b, \mathbf{x} \rangle / \sqrt{d})] \right)_{a,b \in [N]} \in \mathbb{R}^{N \times N}.$$

Then we can rewrite the matrix  $U$  to be

$$U = \lambda_{d,0}(\sigma)^2 \mathbf{1}_N \mathbf{1}_N^\top + \mu_1^2 Q + \mu_\star^2 (\mathbf{I}_N + \Delta),$$

with  $Q = \Theta \Theta^\top / d$  and  $\mathbb{E}[\|\Delta\|_{\text{op}}^2] = o_d(1)$ .

The next several lemmas establish general bounds on the operator norm of random kernel matrices which is of independent interest.

**Lemma B.11.** *Let  $\sigma : \mathbb{R} \rightarrow \mathbb{R}$  be an activation function satisfying Assumption 1, i.e.,  $|\sigma(u)|, |\sigma'(u)| \leq c_0 e^{c_1|u|}$  for some constants  $c_0, c_1 \in (0, \infty)$ . Let  $(\bar{\mathbf{z}}_i)_{i \in [M]} \sim_{iid} \mathbf{N}(\mathbf{0}, \mathbf{I}_d)$ . Assume  $0 < 1/c_2 \leq M/d \leq c_2 < \infty$  for some constant  $c_2 \in (0, \infty)$ . Consider the random matrix  $\bar{\mathbf{R}} \in \mathbb{R}^{M \times M}$  defined by*

$$\bar{R}_{ij} = \mathbf{1}_{i \neq j} \cdot \sigma(\langle \bar{\mathbf{z}}_i, \bar{\mathbf{z}}_j \rangle / \sqrt{d}) / \sqrt{d}. \quad (122)$$

Then there exists a constant  $C$  depending uniquely on  $c_0, c_1, c_2$ , and a sequence of numbers  $(\bar{\eta}_d)_{d \geq 1}$  with  $|\bar{\eta}_d| \leq C \exp\{C(\log d)^{1/2}\}$ , such that

$$\|\bar{\mathbf{R}} - \bar{\eta}_d \mathbf{1}_M \mathbf{1}_M^\top / \sqrt{d}\|_{\text{op}} = O_{d,\mathbb{P}}(\exp\{C(\log d)^{1/2}\}). \quad (123)$$

*Proof of Lemma B.11.* By Lemma B.8 and Markov inequality, we have, for any  $i \neq j$  and all  $0 \leq t \leq \sqrt{d}$ ,

$$\mathbb{P}\left(\langle \bar{\mathbf{z}}_i, \bar{\mathbf{z}}_j \rangle / \sqrt{d} \geq t\right) \leq e^{-t^2/4}. \quad (124)$$

Hence

$$\begin{aligned} & \mathbb{P}\left(\max_{1 \leq i < j \leq M} \left| \frac{1}{\sqrt{d}} \langle \bar{\mathbf{z}}_i, \bar{\mathbf{z}}_j \rangle \right| \geq 16\sqrt{\log M}\right) \\ & \leq \frac{M^2}{2} \max_{1 \leq i < j \leq M} \mathbb{P}\left(\left| \frac{1}{\sqrt{d}} \langle \bar{\mathbf{z}}_i, \bar{\mathbf{z}}_j \rangle \right| \geq 16\sqrt{\log M}\right) \leq M^2 \exp\{-4(\log M)\} \leq \frac{1}{M^2}. \end{aligned} \quad (125)$$

We define  $\tilde{\sigma} : \mathbb{R} \rightarrow \mathbb{R}$  as follows: for  $|u| \leq \bar{x} \equiv 16\sqrt{\log d}$ , define  $\tilde{\sigma}(u) \equiv \sigma(u)e^{-c_1|u|}/c_0$ ; for  $u > \bar{x}$ , define  $\tilde{\sigma}(u) = \tilde{\sigma}(\bar{x})$ ; for  $u < -\bar{x}$ , define  $\tilde{\sigma}(u) = \tilde{\sigma}(-\bar{x})$ . Then  $\tilde{\sigma}$  is a 1-bounded-Lipschitz function on  $\mathbb{R}$ . Define  $\tilde{\eta}_d = \mathbb{E}_{\bar{\mathbf{x}}, \bar{\mathbf{y}} \sim \mathbf{N}(\mathbf{0}, \mathbf{I}_d)} [\tilde{\sigma}(\langle \bar{\mathbf{x}}, \bar{\mathbf{y}} \rangle / \sqrt{d})]$  and  $\bar{\eta}_d = \tilde{\eta}_d c_0 e^{c_1 \bar{x}}$ . Since we have  $|\tilde{\eta}_d| \leq \max_u |\tilde{\sigma}(u)| \leq 1$ , we have

$$|\bar{\eta}_d| = O_d(\exp\{C(\log d)^{1/2}\}). \quad (126)$$

Moreover, we define  $\bar{\mathbf{K}}, \tilde{\mathbf{K}} \in \mathbb{R}^{M \times M}$  by

$$\begin{aligned} \tilde{K}_{ij} &= \mathbf{1}_{i \neq j} \cdot (\tilde{\sigma}(\langle \bar{\mathbf{z}}_i, \bar{\mathbf{z}}_j \rangle / \sqrt{d}) - \tilde{\eta}_d) / \sqrt{d}, \\ \bar{K}_{ij} &= \mathbf{1}_{i \neq j} \cdot (\sigma(\langle \bar{\mathbf{z}}_i, \bar{\mathbf{z}}_j \rangle / \sqrt{d}) - \bar{\eta}_d) / \sqrt{d}. \end{aligned} \quad (127)$$

By [DM16, Lemma 20], there exists a constant  $C$  such that

$$\mathbb{P}(\|\tilde{\mathbf{K}}\|_{\text{op}} \geq C) \leq Ce^{-d/C}.$$

Note that [DM16, Lemma 20] considers one specific choice of  $\tilde{\sigma}$ , but the proof applies unchanged to any 1-Lipschitz function with zero expectation under the measure  $\mu_d$ , where  $\mu_d$  is the distribution of  $\langle \bar{\mathbf{x}}, \bar{\mathbf{y}} \rangle / \sqrt{d}$  for  $\bar{\mathbf{x}}, \bar{\mathbf{y}} \sim \mathbf{N}(\mathbf{0}, \mathbf{I}_d)$ .

Defining the event  $\mathcal{G} \equiv \{|\langle \bar{\mathbf{z}}_i, \bar{\mathbf{z}}_j \rangle / \sqrt{d}| \leq 16\sqrt{\log d}, \forall 1 \leq i < j \leq M\}$ , we have

$$\mathbb{P}(\|\bar{\mathbf{K}}\|_{\text{op}} \geq C c_0 e^{c_1 |\bar{x}|}) \leq \mathbb{P}(\|\bar{\mathbf{K}}\|_{\text{op}} \geq C c_0 e^{c_1 |\bar{x}|}; \mathcal{G}) + \mathbb{P}(\mathcal{G}^c) \leq \mathbb{P}(\|\tilde{\mathbf{K}}\|_{\text{op}} \geq C) + \frac{1}{M^2} = o_d(1). \quad (128)$$

By Eq. (122) and (127), we have

$$\bar{\mathbf{R}} = \bar{\mathbf{K}} - \bar{\eta}_d \mathbf{I}_M / \sqrt{d} + \bar{\eta}_d \mathbf{1}_M \mathbf{1}_M^\top / \sqrt{d}.$$

By Eq. (128) and (126), we have

$$\|\bar{\mathbf{R}} - \bar{\eta}_d \mathbf{1}_M \mathbf{1}_M^\top / \sqrt{d}\|_{\text{op}} = \|\bar{\mathbf{K}} - \bar{\eta}_d \mathbf{I}_M / \sqrt{d}\|_{\text{op}} \leq \|\bar{\mathbf{K}}\|_{\text{op}} + \bar{\eta}_d / \sqrt{d} = O_{d,\mathbb{P}}(\exp\{C(\log d)^{1/2}\}).$$

This completes the proof.  $\square$

**Lemma B.12.** *Let  $\sigma : \mathbb{R} \rightarrow \mathbb{R}$  be an activation function satisfying Assumption 1, i.e.,  $|\sigma(u)|, |\sigma'(u)| \leq c_0 e^{c_1 |u|}$  for some constants  $c_0, c_1 \in (0, \infty)$ . Let  $(\bar{\mathbf{z}}_i)_{i \in [M]} \sim_{i.i.d} \mathbf{N}(\mathbf{0}, \mathbf{I}_d)$ . Assume  $0 < 1/c_2 \leq M/d \leq c_2 < \infty$  for some constant  $c_2 \in (0, \infty)$ . Define  $\mathbf{z}_i = \sqrt{d} \cdot \bar{\mathbf{z}}_i / \|\bar{\mathbf{z}}_i\|_2$ . Consider two random matrices  $\mathbf{R}, \bar{\mathbf{R}} \in \mathbb{R}^{M \times M}$  defined by*

$$\begin{aligned} \bar{R}_{ij} &= \mathbf{1}_{i \neq j} \cdot \sigma(\langle \bar{\mathbf{z}}_i, \bar{\mathbf{z}}_j \rangle / \sqrt{d}) / \sqrt{d}, \\ R_{ij} &= \mathbf{1}_{i \neq j} \cdot \sigma(\langle \mathbf{z}_i, \mathbf{z}_j \rangle / \sqrt{d}) / \sqrt{d}. \end{aligned}$$

*Then there exists a constant  $C$  depending uniquely on  $c_0, c_1, c_2$ , such that*

$$\|\bar{\mathbf{R}} - \mathbf{R}\|_{\text{op}} = O_{d,\mathbb{P}}(\exp\{C(\log d)^{1/2}\}).$$

*Proof of Lemma B.12.* In the proof of this lemma, we assume  $\sigma$  has continuous derivatives. In the case when  $\sigma$  is only weak differentiable, the proof is the same, except that we need to replace the mean value theorem to its integral form.

Define  $r_i = \sqrt{d} / \|\bar{\mathbf{z}}_i\|_2$ , and

$$\tilde{R}_{ij} = \mathbf{1}_{i \neq j} \cdot \sigma(r_i \langle \bar{\mathbf{z}}_i, \bar{\mathbf{z}}_j \rangle / \sqrt{d}) / \sqrt{d}.$$

By the concentration of  $\chi$ -squared distribution, it is easy to see that

$$\max_{i \in [M]} |r_i - 1| = O_{d,\mathbb{P}}((\log d)^{1/2} / d^{1/2}).$$

Moreover, we have (for  $\zeta_i$  between  $r_i$  and 1)

$$|\bar{R}_{ij} - \tilde{R}_{ij}| \leq |\sigma'(\zeta_i \langle \bar{\mathbf{z}}_i, \bar{\mathbf{z}}_j \rangle / \sqrt{d})| \cdot |\langle \bar{\mathbf{z}}_i, \bar{\mathbf{z}}_j \rangle / \sqrt{d}| \cdot |r_i - 1| / \sqrt{d}.$$

By Eq. (125), we have

$$\begin{aligned} \max_{i \neq j \in [M]} |\langle \bar{\mathbf{z}}_i, \bar{\mathbf{z}}_j \rangle / \sqrt{d}| &= O_{d,\mathbb{P}}((\log d)^{1/2}), \\ \max_{i \neq j \in [M]} |\zeta_i \langle \bar{\mathbf{z}}_i, \bar{\mathbf{z}}_j \rangle / \sqrt{d}| &= O_{d,\mathbb{P}}((\log d)^{1/2}). \end{aligned}$$

Moreover by the assumption that  $|\sigma'(u)| \leq c_0 e^{c_1 |u|}$ , we have

$$\max_{i \neq j \in [M]} |\sigma'(\zeta_i \langle \bar{\mathbf{z}}_i, \bar{\mathbf{z}}_j \rangle / \sqrt{d})| \cdot |\langle \bar{\mathbf{z}}_i, \bar{\mathbf{z}}_j \rangle / \sqrt{d}| = O_{d,\mathbb{P}}(\exp\{C(\log d)^{1/2}\}).$$

This gives

$$\max_{i \neq j \in [M]} |\bar{R}_{ij} - \tilde{R}_{ij}| = O_{d,\mathbb{P}}(\exp\{C(\log d)^{1/2}\}/d).$$

Using similar argument, we can show that

$$\max_{i \neq j \in [M]} |R_{ij} - \tilde{R}_{ij}| = O_{d,\mathbb{P}}(\exp\{C(\log d)^{1/2}\}/d),$$

which gives

$$\max_{i \neq j \in [M]} |R_{ij} - \bar{R}_{ij}| = O_{d,\mathbb{P}}(\exp\{C(\log d)^{1/2}\}/d).$$

This gives

$$\|\mathbf{R} - \bar{\mathbf{R}}\|_{\text{op}} \leq \|\mathbf{R} - \bar{\mathbf{R}}\|_F \leq d \cdot \max_{i \neq j \in [M]} |R_{ij} - \bar{R}_{ij}| = O_{d,\mathbb{P}}(\exp\{C(\log d)^{1/2}\}).$$

This proves the lemma.  $\square$

**Lemma B.13.** *Let  $\sigma : \mathbb{R} \rightarrow \mathbb{R}$  be an activation function satisfying Assumption 1, i.e.,  $|\sigma(u)|, |\sigma'(u)| \leq c_0 e^{c_1|u|}$  for some constants  $c_0, c_1 \in (0, \infty)$ . Let  $(\mathbf{z}_i)_{i \in [M]} \sim_{iid} \text{Unif}(\mathbb{S}^{d-1}(\sqrt{d}))$ . Assume  $0 < 1/c_2 \leq M/d \leq c_2 < \infty$  for some constant  $c_2 \in (0, \infty)$ . Define  $\lambda_{d,0} = \mathbb{E}_{\mathbf{z}_1, \mathbf{z}_2 \sim \text{Unif}(\mathbb{S}^{d-1}(\sqrt{d}))}[\sigma(\langle \mathbf{z}_1, \mathbf{z}_2 \rangle / \sqrt{d})]$ , and  $\varphi_d(u) = \sigma(u) - \lambda_{d,0}$ . Consider the random matrix  $\mathbf{K} \in \mathbb{R}^{M \times M}$  with*

$$K_{ij} = \mathbf{1}_{i \neq j} \cdot \frac{1}{\sqrt{d}} \varphi_d\left(\frac{1}{\sqrt{d}} \langle \mathbf{z}_i, \mathbf{z}_j \rangle\right).$$

Then there exists a constant  $C$  depending uniquely on  $c_0, c_1, c_2$ , such that

$$\|\mathbf{K}\|_{\text{op}} \leq O_{d,\mathbb{P}}(\exp\{C(\log d)^{1/2}\}).$$

*Proof of Lemma B.13.* We construct  $(\mathbf{z}_i)_{i \in [M]}$  by normalizing a collection of independent Gaussian random vectors. Let  $(\bar{\mathbf{z}}_i)_{i \in [M]} \sim_{iid} \mathbf{N}(\mathbf{0}, \mathbf{I}_d)$  and denote  $\mathbf{z}_i = \sqrt{d} \cdot \bar{\mathbf{z}}_i / \|\bar{\mathbf{z}}_i\|_2$  for  $i \in [M]$ . Then we have  $(\mathbf{z}_i)_{i \in [M]} \sim_{iid} \text{Unif}(\mathbb{S}^{d-1}(\sqrt{d}))$ .

Consider two random matrices  $\bar{\mathbf{R}}, \mathbf{R} \in \mathbb{R}^{M \times M}$  defined by

$$\begin{aligned} \bar{R}_{ij} &= \mathbf{1}_{i \neq j} \cdot \sigma(\langle \bar{\mathbf{z}}_i, \bar{\mathbf{z}}_j \rangle / \sqrt{d}) / \sqrt{d}, \\ R_{ij} &= \mathbf{1}_{i \neq j} \cdot \sigma(\langle \mathbf{z}_i, \mathbf{z}_j \rangle / \sqrt{d}) / \sqrt{d}. \end{aligned}$$

By Lemma B.11, there exists a sequence  $(\bar{\eta}_d)_{d \geq 0}$  with  $|\bar{\eta}_d| \leq C \exp\{C(\log d)^{1/2}\}$ , such that

$$\|\bar{\mathbf{R}} - \bar{\eta}_d \mathbf{1}_M \mathbf{1}_M^\top / \sqrt{d}\|_{\text{op}} = O_{d,\mathbb{P}}(\exp\{C(\log d)^{1/2}\}).$$

Moreover, by Lemma B.12, we have,

$$\|\bar{\mathbf{R}} - \mathbf{R}\|_{\text{op}} \leq O_{d,\mathbb{P}}(\exp\{C(\log d)^{1/2}\}),$$

which gives,

$$\|\mathbf{R} - \bar{\eta}_d \mathbf{1}_M \mathbf{1}_M^\top / \sqrt{d}\|_{\text{op}} = O_{d,\mathbb{P}}(\exp\{C(\log d)^{1/2}\}).$$

Note we have

$$\mathbf{R} = \mathbf{K} + \lambda_{d,0} \mathbf{1}_M \mathbf{1}_M^\top / \sqrt{d} - \lambda_{d,0} \mathbf{I}_M / \sqrt{d}.$$

Moreover, note that  $\lim_{d \rightarrow \infty} \lambda_{d,0} = \mathbb{E}_{G \sim \mathbf{N}(0,1)}[\sigma(G)]$  so that  $\sup_d |\lambda_{d,0}| \leq C$ . Therefore, denoting  $\kappa_d = \lambda_{d,0} - \bar{\eta}_d$ , we have

$$\begin{aligned} \|\mathbf{K} + \kappa_d \mathbf{1}_M \mathbf{1}_M^\top / \sqrt{d}\|_{\text{op}} &= \|\mathbf{R} - \bar{\eta}_d \mathbf{1}_M \mathbf{1}_M^\top / \sqrt{d} + \lambda_{d,0} \mathbf{I}_M / \sqrt{d}\|_{\text{op}} \\ &\leq \|\mathbf{R} - \bar{\eta}_d \mathbf{1}_M \mathbf{1}_M^\top / \sqrt{d}\|_{\text{op}} + \lambda_{d,0} / \sqrt{d} = O_{d,\mathbb{P}}(\exp\{C(\log d)^{1/2}\}). \end{aligned} \tag{129}$$

Notice that

$$|\mathbf{1}_M^\top \mathbf{K} \mathbf{1}_M / M| \leq \frac{C}{M^{3/2}} \left| \sum_{i \neq j} \varphi_d(\langle \mathbf{z}_i, \mathbf{z}_j \rangle / \sqrt{d}) \right| \leq \frac{C}{M} \sum_{i=1}^M \left| \sum_{j: j \neq i} \varphi_d(\langle \mathbf{z}_i, \mathbf{z}_j \rangle / \sqrt{d}) / \sqrt{M} \right| \equiv \frac{C}{M} \sum_{i=1}^M |V_i|,$$

where

$$V_i = \frac{1}{\sqrt{M}} \sum_{j: j \neq i} \varphi_d(\langle \mathbf{z}_i, \mathbf{z}_j \rangle / \sqrt{d}).$$

Note  $\mathbb{E}[\varphi_d(\langle \mathbf{z}_i, \mathbf{z}_j \rangle / \sqrt{d})] = 0$  for  $i \neq j$  so that  $\mathbb{E}[\varphi_d(\langle \mathbf{z}_i, \mathbf{z}_{j_1} \rangle / \sqrt{d}) \varphi_d(\langle \mathbf{z}_i, \mathbf{z}_{j_2} \rangle / \sqrt{d})] = 0$  for  $i, j_1, j_2$  distinct. Calculating the second moment, we have

$$\sup_{i \in [M]} \mathbb{E}[V_i^2] = \sup_{i \in [M]} \mathbb{E} \left[ \left( \sum_{j: j \neq i} \varphi_d(\langle \mathbf{z}_i, \mathbf{z}_j \rangle / \sqrt{d}) / \sqrt{M} \right)^2 \right] = \sup_{i \in [M]} \frac{1}{M} \sum_{j: j \neq i} \mathbb{E}[\varphi_d(\langle \mathbf{z}_i, \mathbf{z}_j \rangle / \sqrt{d})^2] = O_d(1).$$

Therefore, we have

$$\mathbb{E}[(\mathbf{1}_M^\top \mathbf{K} \mathbf{1}_M / M)^2] \leq \frac{C^2}{M^2} \sum_{i,j=1}^M \mathbb{E}[|V_i| \cdot |V_j|] \leq \frac{C^2}{M^2} \sum_{i,j=1}^M \mathbb{E}[(V_i^2 + V_j^2)/2] \leq C^2 \sup_{i \in [M]} \mathbb{E}[V_i^2] = O_d(1).$$

This gives

$$|\mathbf{1}_M^\top \mathbf{K} \mathbf{1}_M / M| = O_{d,\mathbb{P}}(1).$$

Combining this equation with Eq. (129), we get

$$\begin{aligned} & \|\kappa_d \mathbf{1}_M \mathbf{1}_M^\top / \sqrt{d}\|_{\text{op}} = |\langle \mathbf{1}_M, (\kappa_d \mathbf{1}_M \mathbf{1}_M^\top / \sqrt{d}) \mathbf{1}_M \rangle / M| \\ & \leq |\langle \mathbf{1}_M, (\mathbf{K} + \kappa_d \mathbf{1}_M \mathbf{1}_M^\top / \sqrt{d}) \mathbf{1}_M \rangle / M| + |\mathbf{1}_M^\top \mathbf{K} \mathbf{1}_M / M| \\ & \leq \|\mathbf{K} + \kappa_d \mathbf{1}_M \mathbf{1}_M^\top / \sqrt{d}\|_{\text{op}} + |\mathbf{1}_M^\top \mathbf{K} \mathbf{1}_M / M| = O_{d,\mathbb{P}}(\exp\{C(\log d)^{1/2}\}), \end{aligned}$$

and hence

$$\|\mathbf{K}\|_{\text{op}} \leq \|\mathbf{K} + \kappa_d \mathbf{1}_M \mathbf{1}_M^\top / \sqrt{d}\|_{\text{op}} + \|\kappa_d \mathbf{1}_M \mathbf{1}_M^\top / \sqrt{d}\|_{\text{op}} = O_{d,\mathbb{P}}(\exp\{C(\log d)^{1/2}\}).$$

This proves the lemma.  $\square$

## C Proof of Proposition 7.2

This section is organized as follows. We prove Proposition 7.2 in Section C.6. We collect the elements to prove Proposition 7.2 in Section C.1, C.2, C.3, C.4, and C.5. In Section C.1, we show that the Stieltjes transform is stable when replacing the distribution of  $\mathbf{x}_i, \boldsymbol{\theta}_a$  from uniform distribution on the sphere to Gaussian distribution. In Section C.2, we give some properties for the fixed point equation Eq. (61) defined in the statement of Proposition 7.2. In Section C.3 we states the key lemma (Lemma C.4): Stieltjes transform approximately satisfies the fixed point equation, when  $\mathbf{x}_i, \boldsymbol{\theta}_a \sim_{iid} \mathbf{N}(\mathbf{0}, \mathbf{I}_d)$  and  $\varphi$  is a polynomial with  $\mathbb{E}_{G \sim \mathbf{N}(0,1)}[\varphi(G)] = 0$ . In Section C.4 we give some properties of Stieltjes transform used to prove Lemma C.4. In Section C.5, we prove Lemma C.4 using leave-one-out argument.

### C.1 Equivalence between Gaussian and sphere vectors

Let  $(\bar{\boldsymbol{\theta}}_a)_{a \in [N]} \sim_{iid} \mathbf{N}(0, \mathbf{I}_d)$ ,  $(\bar{\mathbf{x}}_i)_{i \in [n]} \sim_{iid} \mathbf{N}(0, \mathbf{I}_d)$ . We denote by  $\bar{\boldsymbol{\Theta}} \in \mathbb{R}^{N \times d}$  the matrix whose  $a$ -th row is given by  $\bar{\boldsymbol{\theta}}_a$ , and by  $\bar{\mathbf{X}} \in \mathbb{R}^{n \times d}$  the matrix whose  $i$ -th row is given by  $\bar{\mathbf{x}}_i$ . We denote by  $\boldsymbol{\Theta} \in \mathbb{R}^{N \times d}$  the matrix whose  $a$ -th row is given by  $\boldsymbol{\theta}_a = \sqrt{d} \cdot \bar{\boldsymbol{\theta}}_a / \|\bar{\boldsymbol{\theta}}_a\|_2$ , and by  $\mathbf{X} \in \mathbb{R}^{n \times d}$  the matrix whose  $i$ -th row is given by  $\mathbf{x}_i = \sqrt{d} \cdot \bar{\mathbf{x}}_i / \|\bar{\mathbf{x}}_i\|_2$ . Then we have  $(\mathbf{x}_i)_{i \in [n]} \sim_{iid} \text{Unif}(\mathbb{S}^{d-1}(\sqrt{d}))$  and  $(\boldsymbol{\theta}_a)_{a \in [N]} \sim_{iid} \text{Unif}(\mathbb{S}^{d-1}(\sqrt{d}))$  independently.

We consider activation functions  $\sigma, \varphi : \mathbb{R} \rightarrow \mathbb{R}$  with  $\varphi(x) = \sigma(x) - \mathbb{E}_{G \sim \mathcal{N}(0,1)}[\sigma(G)]$ . We define the following matrices (where  $\mu_1$  is the first Hermite coefficients of  $\sigma$ )

$$\bar{\mathbf{J}} \equiv \frac{1}{\sqrt{d}} \varphi\left(\frac{1}{\sqrt{d}} \bar{\mathbf{X}} \bar{\boldsymbol{\Theta}}^\top\right), \quad \mathbf{Z} \equiv \frac{1}{\sqrt{d}} \sigma\left(\frac{1}{\sqrt{d}} \mathbf{X} \boldsymbol{\Theta}^\top\right), \quad (130)$$

$$\bar{\mathbf{J}}_1 \equiv \frac{\mu_1}{d} \bar{\mathbf{X}} \bar{\boldsymbol{\Theta}}^\top, \quad \mathbf{Z}_1 \equiv \frac{\mu_1}{d} \mathbf{X} \boldsymbol{\Theta}^\top, \quad (131)$$

$$\bar{\mathbf{Q}} \equiv \frac{1}{d} \bar{\boldsymbol{\Theta}} \bar{\boldsymbol{\Theta}}^\top, \quad \mathbf{Q} \equiv \frac{1}{d} \boldsymbol{\Theta} \boldsymbol{\Theta}^\top, \quad (132)$$

$$\bar{\mathbf{H}} \equiv \frac{1}{d} \bar{\mathbf{X}} \bar{\mathbf{X}}^\top, \quad \mathbf{H} \equiv \frac{1}{d} \mathbf{X} \mathbf{X}^\top, \quad (133)$$

as well as the block matrix  $\bar{\mathbf{A}}, \mathbf{A} \in \mathbb{R}^{M \times M}$ ,  $M = N + n$ , defined by

$$\bar{\mathbf{A}} = \begin{bmatrix} s_1 \mathbf{I}_N + s_2 \bar{\mathbf{Q}} & \bar{\mathbf{J}}^\top + p \bar{\mathbf{J}}_1^\top \\ \bar{\mathbf{J}} + p \bar{\mathbf{J}}_1 & t_1 \mathbf{I}_n + t_2 \bar{\mathbf{H}} \end{bmatrix}, \quad \mathbf{A} = \begin{bmatrix} s_1 \mathbf{I}_N + s_2 \mathbf{Q} & \mathbf{Z}^\top + p \mathbf{Z}_1^\top \\ \mathbf{Z} + p \mathbf{Z}_1 & t_1 \mathbf{I}_n + t_2 \mathbf{H} \end{bmatrix}, \quad (134)$$

and the Stieltjes transforms  $\bar{M}_d(\boldsymbol{\xi}; \mathbf{q})$  and  $M_d(\boldsymbol{\xi}; \mathbf{q})$ , defined by

$$\bar{M}_d(\boldsymbol{\xi}; \mathbf{q}) = \frac{1}{d} \text{Tr}[(\bar{\mathbf{A}} - \boldsymbol{\xi} \mathbf{I}_M)^{-1}], \quad M_d(\boldsymbol{\xi}; \mathbf{q}) = \frac{1}{d} \text{Tr}[(\mathbf{A} - \boldsymbol{\xi} \mathbf{I}_M)^{-1}]. \quad (135)$$

The readers could keep in mind: a quantity with an overline corresponds to the case when features and data are Gaussian, while a quantity without overline usually corresponds to the case when features and data are on the sphere.

**Lemma C.1.** *Let  $\sigma$  be a fixed polynomial. Let  $\varphi(x) = \sigma(x) - \mathbb{E}_{G \sim \mathcal{N}(0,1)}[\sigma(G)]$ . Consider the linear regime of Assumption 2. For any fixed  $\mathbf{q} \in \mathcal{Q}$  and for any  $\xi_0 > 0$ , we have*

$$\mathbb{E} \left[ \sup_{\Im \xi \geq \xi_0} |\bar{M}_d(\boldsymbol{\xi}; \mathbf{q}) - M_d(\boldsymbol{\xi}; \mathbf{q})| \right] = o_d(1).$$

*Proof of Lemma C.1.*

**Step 1. Show that the resolvent is stable with respect to nuclear norm perturbation.**

We define

$$\Delta(\mathbf{A}, \bar{\mathbf{A}}, \boldsymbol{\xi}) = M_d(\boldsymbol{\xi}; \mathbf{q}) - \bar{M}_d(\boldsymbol{\xi}; \mathbf{q}).$$

Then we have deterministically

$$|\Delta(\mathbf{A}, \bar{\mathbf{A}}, \boldsymbol{\xi})| \leq |M_d(\boldsymbol{\xi}; \mathbf{q})| + |\bar{M}_d(\boldsymbol{\xi}; \mathbf{q})| \leq 4(\psi_1 + \psi_2)/\Im \boldsymbol{\xi}.$$

Moreover, we have

$$\begin{aligned} |\Delta(\mathbf{A}, \bar{\mathbf{A}}, \boldsymbol{\xi})| &= |\text{Tr}((\mathbf{A} - \boldsymbol{\xi} \mathbf{I})^{-1} (\mathbf{A} - \bar{\mathbf{A}}) (\bar{\mathbf{A}} - \boldsymbol{\xi} \mathbf{I})^{-1})|/d \\ &\leq \|(\mathbf{A} - \boldsymbol{\xi} \mathbf{I})^{-1} (\bar{\mathbf{A}} - \boldsymbol{\xi} \mathbf{I})^{-1}\|_{\text{op}} \|\mathbf{A} - \bar{\mathbf{A}}\|_*/d \\ &\leq \|\mathbf{A} - \bar{\mathbf{A}}\|_*/(d(\Im \boldsymbol{\xi})^2). \end{aligned}$$

Therefore, if we can show  $\|\mathbf{A} - \bar{\mathbf{A}}\|_*/d = o_{d, \mathbb{P}}(1)$ , then  $\mathbb{E}[\sup_{\Im \boldsymbol{\xi} \geq \xi_0} |\Delta(\mathbf{A}, \bar{\mathbf{A}}, \boldsymbol{\xi})|] = o_d(1)$ .

**Step 2. Show that  $\|\mathbf{A} - \bar{\mathbf{A}}\|_*/d = o_{d, \mathbb{P}}(1)$ .**

Denote  $\mathbf{Z}_0 = \mathbb{E}_{G \sim \mathcal{N}(0,1)}[\sigma(G)] \mathbf{1}_n \mathbf{1}_N^\top / \sqrt{d}$  and  $\mathbf{Z}_* = \varphi(\mathbf{X} \boldsymbol{\Theta}^\top / \sqrt{d}) / \sqrt{d}$ . Then we have  $\mathbf{Z} = \mathbf{Z}_0 + \mathbf{Z}_*$ , and

$$\begin{aligned} \mathbf{A} - \bar{\mathbf{A}} &= s_2 \begin{bmatrix} \mathbf{Q} - \bar{\mathbf{Q}} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} + t_2 \begin{bmatrix} \mathbf{0} - \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{H} - \bar{\mathbf{H}} \end{bmatrix} + p \begin{bmatrix} \mathbf{0} & \mathbf{Z}_1^\top - \bar{\mathbf{J}}_1^\top \\ \mathbf{Z}_1 - \bar{\mathbf{J}}_1 & \mathbf{0} \end{bmatrix} \\ &\quad + \begin{bmatrix} \mathbf{0} & \mathbf{Z}_*^\top - \bar{\mathbf{J}}^\top \\ \mathbf{Z}_* - \bar{\mathbf{J}} & \mathbf{0} \end{bmatrix} + \begin{bmatrix} \mathbf{0} & \mathbf{Z}_0^\top \\ \mathbf{Z}_0 & \mathbf{0} \end{bmatrix}. \end{aligned}$$



Since  $\mathbf{q} = (s_1, s_2, t_1, t_2, p)$  is fixed, we have

$$\begin{aligned} \frac{1}{d} \|\mathbf{A} - \bar{\mathbf{A}}\|_* \leq & C \left[ \frac{1}{\sqrt{d}} \|\bar{\mathbf{Q}} - \mathbf{Q}\|_F + \frac{1}{\sqrt{d}} \|\bar{\mathbf{H}} - \mathbf{H}\|_F + \frac{1}{\sqrt{d}} \|\bar{\mathbf{J}}_1 - \mathbf{Z}_1\|_F + \frac{1}{\sqrt{d}} \|\bar{\mathbf{J}} - \mathbf{Z}_*\|_F \right. \\ & \left. + \frac{1}{d} \left\| \begin{bmatrix} \mathbf{0} & \mathbf{Z}_0^\top \\ \mathbf{Z}_0 & \mathbf{0} \end{bmatrix} \right\|_* \right]. \end{aligned}$$

The nuclear norm of the term involving  $\mathbf{Z}_0$  can be easily bounded by

$$\begin{aligned} \frac{1}{d} \left\| \begin{bmatrix} \mathbf{0} & \mathbf{Z}_0^\top \\ \mathbf{Z}_0 & \mathbf{0} \end{bmatrix} \right\|_* &= \frac{1}{d^{3/2}} |\mathbb{E}_{G \sim \mathcal{N}(0,1)}[\sigma(G)]| \cdot \left\| \begin{bmatrix} \mathbf{0} & \mathbf{1}_N \mathbf{1}_n^\top \\ \mathbf{1}_n \mathbf{1}_N^\top & \mathbf{0} \end{bmatrix} \right\|_* \\ &\leq \frac{\sqrt{2}}{d^{3/2}} |\mathbb{E}_{G \sim \mathcal{N}(0,1)}[\sigma(G)]| \cdot \left\| \begin{bmatrix} \mathbf{0} & \mathbf{1}_N \mathbf{1}_n^\top \\ \mathbf{1}_n \mathbf{1}_N^\top & \mathbf{0} \end{bmatrix} \right\|_F = o_d(1). \end{aligned}$$

For term  $\bar{\mathbf{H}} - \mathbf{H}$ , denoting  $\mathbf{D}_x = \text{diag}(\sqrt{d}/\|\bar{\mathbf{x}}_1\|_2, \dots, \sqrt{d}/\|\bar{\mathbf{x}}_n\|_2)$ , we have

$$\|\bar{\mathbf{H}} - \mathbf{H}\|_F / \sqrt{d} \leq \|\bar{\mathbf{H}} - \mathbf{H}\|_{\text{op}} \leq \|\mathbf{I}_n - \mathbf{D}_x\|_{\text{op}} \|\bar{\mathbf{H}}\|_{\text{op}} (1 + \|\mathbf{D}_x\|_{\text{op}}) = o_{d,\mathbb{P}}(1),$$

where we used the fact that  $\|\mathbf{D}_x - \mathbf{I}_n\|_{\text{op}} = o_{d,\mathbb{P}}(1)$  and  $\|\mathbf{H}\|_{\text{op}} = O_{d,\mathbb{P}}(1)$ . Similar argument shows that

$$\|\bar{\mathbf{Q}} - \mathbf{Q}\|_F / \sqrt{d} = o_{d,\mathbb{P}}(1), \quad \|\bar{\mathbf{J}}_1 - \mathbf{Z}_1\|_F / \sqrt{d} = o_{d,\mathbb{P}}(1).$$

**Step 3. Bound for  $\|\bar{\mathbf{J}} - \mathbf{Z}_*\|_F / \sqrt{d}$ .**

Define  $\bar{\mathbf{Z}}_* = \varphi(\mathbf{D}_x \bar{\mathbf{X}} \bar{\boldsymbol{\Theta}}^\top / \sqrt{d}) / \sqrt{d}$ . Define  $r_i = \sqrt{d} / \|\bar{\mathbf{x}}_i\|_2$ . We have (for  $\zeta_{ia}$  between  $r_i$  and 1)

$$\begin{aligned} \bar{\mathbf{Z}}_* - \bar{\mathbf{J}} &= \left( \varphi(r_i \langle \bar{\mathbf{x}}_i, \bar{\boldsymbol{\theta}}_a \rangle / \sqrt{d}) / \sqrt{d} - \varphi(\langle \bar{\mathbf{x}}_i, \bar{\boldsymbol{\theta}}_a \rangle / \sqrt{d}) / \sqrt{d} \right)_{i \in [n], a \in [N]} \\ &= \left( (r_i - 1) (\langle \bar{\mathbf{x}}_i, \bar{\boldsymbol{\theta}}_a \rangle / \sqrt{d}) \varphi'(\zeta_{ia} \langle \bar{\mathbf{x}}_i, \bar{\boldsymbol{\theta}}_a \rangle / \sqrt{d}) / \sqrt{d} \right)_{i \in [n], a \in [N]} \\ &= (\mathbf{D}_x - \mathbf{I}_n) \bar{\varphi}(\boldsymbol{\Xi} \odot (\bar{\mathbf{X}} \bar{\boldsymbol{\Theta}}^\top / \sqrt{d})) / \sqrt{d}, \end{aligned}$$

where  $\boldsymbol{\Xi} = (\zeta_{ia})_{i \in [n], a \in [N]}$ , and  $\bar{\varphi}(x) = x\varphi'(x)$  (so  $\bar{\varphi}$  is a polynomial). It is easy to see that

$$\|\mathbf{D}_x - \mathbf{I}_n\|_{\text{op}} = \max_i |r_i - 1| = O_{d,\mathbb{P}}(\sqrt{\log d} / \sqrt{d}), \quad \|\boldsymbol{\Xi}\|_{\max} = O_{d,\mathbb{P}}(1), \quad \|\bar{\mathbf{X}} \bar{\boldsymbol{\Theta}}^\top / \sqrt{d}\|_{\max} = O_{d,\mathbb{P}}(\sqrt{\log d}).$$

Therefore, we have

$$\begin{aligned} \|\bar{\mathbf{Z}}_* - \bar{\mathbf{J}}\|_F / \sqrt{d} &= \|(\mathbf{D}_x - \mathbf{I}_n) \bar{\varphi}(\boldsymbol{\Xi} \odot (\bar{\mathbf{X}} \bar{\boldsymbol{\Theta}}^\top / \sqrt{d}))\|_F / d \\ &\leq \|\mathbf{D}_x - \mathbf{I}_n\|_{\text{op}} \|\bar{\varphi}(\boldsymbol{\Xi} \odot (\bar{\mathbf{X}} \bar{\boldsymbol{\Theta}}^\top / \sqrt{d}))\|_F / d \\ &\leq C(\varphi) \cdot \|\mathbf{D}_x - \mathbf{I}_n\|_{\text{op}} (1 + \|\boldsymbol{\Xi}\|_{\max} \|\bar{\mathbf{X}} \bar{\boldsymbol{\Theta}}^\top / \sqrt{d}\|_{\max})^{\deg(\varphi)} = O_{d,\mathbb{P}}((\log d)^{\deg(\varphi)+1} / \sqrt{d}) = o_{d,\mathbb{P}}(1). \end{aligned}$$

This proves the lemma.  $\square$

## C.2 Properties of the fixed point equations

In this section we establish some useful properties of the fixed point characterization (61), where  $\mathbf{F}$  is defined via Eq. (60). For the sake of simplicity, we will write  $\mathbf{m} = (m_1, m_2)$  and introduce the function  $\mathbf{F}(\cdot; \xi; \mathbf{q}, \psi_1, \psi_2, \mu_1, \mu_*) : \mathbb{C} \times \mathbb{C} \rightarrow \mathbb{C} \times \mathbb{C}$  via

$$\mathbf{F}(\mathbf{m}; \xi; \mathbf{q}, \psi_1, \psi_2, \mu_1, \mu_*) = \begin{bmatrix} \mathbf{F}(m_1, m_2; \xi; \mathbf{q}, \psi_1, \psi_2, \mu_1, \mu_*) \\ \mathbf{F}(m_2, m_1; \xi; \mathbf{q}, \psi_2, \psi_1, \mu_1, \mu_*) \end{bmatrix}. \quad (136)$$

Since  $\mathbf{q}, \psi_1, \psi_2, \mu_1, \mu_*$  are mostly fixed through what follows, we will drop them from the argument of  $\mathbf{F}$  unless necessary. In these notations, Eq. (61) reads

$$\mathbf{m} = \mathbf{F}(\mathbf{m}; \xi). \quad (137)$$

**Lemma C.2.** *If  $\xi \in \mathbb{C}_+$ , the  $F(\cdot; \xi)$  maps  $\mathbb{C}_+ \times \mathbb{C}_+$  to  $\mathbb{C}_+$ .*

*Proof.* Assume  $\mu_1 \neq 0$  (the proof goes along the same line for  $\mu_1 = 0$ , with some simplifications). We rewrite Eq. (60) as

$$F(m_1, m_2; \xi) = \frac{\psi_1}{-\xi + s_1 + H(m_1, m_2)}, \quad (138)$$

$$H(m_1, m_2) = -\mu_\star^2 m_1 + \frac{1}{m_1 + \frac{1+t_2 m_2}{1+(t_2 s_2 - \mu_1^2(1+p)^2)m_2}}. \quad (139)$$

Consider  $\Im(m_1), \Im(m_2) > 0$ . Note that  $z \mapsto (1 + t_2 z)/(s_2 + (t_2 s_2 - \mu_1^2(1+p)^2)z)$  maps  $\mathbb{C}_+ \rightarrow \mathbb{C}_+$ . Hence  $\Im[(m_1 + (1 + t_2 m_2)/(1 + (t_2 s_2 - \mu_1^2(1+p)^2)m_2))] > 0$ , whence the fraction in Eq. (139) has negative imaginary part, and therefore  $\Im(H) \leq 0$  (note that  $\mu_\star^2 > 0$ ). From Eq. (138), we get  $\Im(F) > 0$  as claimed.  $\square$

**Lemma C.3.** *Let  $\mathbb{D}(r) = \{z : |z| < r\}$  be the disk of radius  $r$  in the complex plane. Then, there exists  $r_0 > 0$  such that, for any  $r, \delta > 0$  there exists  $\xi_0 = \xi_0(r, \delta, \mathbf{q}, \psi_1, \psi_2, \mu_1, \mu_\star) > 0$  such that, if  $\Im(\xi) \geq \xi_0$ , then  $\mathbf{F}$  maps  $\mathbb{D}(r_0) \times \mathbb{D}(r_0)$  into  $\mathbb{D}(r) \times \mathbb{D}(r)$  and further is Lipschitz continuous, with Lipschitz constant at most  $\delta$  on that domain.*

*In particular, if  $\Im(\xi) \geq \xi_0 > 0$ , then Eq. (61) admits a unique solution with  $|m_1|, |m_2| < r_0$ . For  $\Im(\xi) > 0$ , this solution further satisfies  $\Im(m_1), \Im(m_2) > 0$  and  $|m_1| \leq \psi_1/\Im(\xi)$ ,  $|m_2| \leq \psi_2/\Im(\xi)$ .*

*Finally, for  $\Im(\xi) > \xi_0$ , the solution  $\mathbf{m}(\xi; \mathbf{q}, \psi_1, \psi_2, \mu_\star, \mu_1)$  is continuously differentiable in  $\mathbf{q}, \psi_1, \psi_2, \mu_\star, \mu_1$ .*

*Proof of Lemma C.3.* Consider the definition (60), which we rewrite as in Eq. (138). It is easy to see that, for  $r_0 = r_0(\mathbf{q})$  small enough,  $|H(\mathbf{m})| \leq 2$  for  $\mathbf{m} \in \mathbb{D}(r_0) \times \mathbb{D}(r_0)$ . Therefore  $|F(\mathbf{m}; \xi)| \leq \psi_1/(\Im(\xi) - 2) < r$  provided  $\xi_0 > 2 + (\psi_1/r)$ . By eventually enlarging  $\xi_0$ , we get the same bound for  $\|\mathbf{F}(\mathbf{m}; \xi)\|_\infty$ . The existence of a unique fixed point follows by Banach fixed point theorem.

In order to prove the Lipschitz continuity of  $\mathbf{F}$  in this domain, notice that  $\mathbf{F}$  is differentiable and

$$\nabla_{\mathbf{m}} F(\mathbf{m}; \xi) = \frac{\psi_1}{(-\xi + s_1 + H(\mathbf{m}))^2} \nabla_{\mathbf{m}} H(\mathbf{m}). \quad (140)$$

By eventually reducing  $r_0$ , we can ensure  $\|\nabla_{\mathbf{m}} H(\mathbf{m})\|_2 \leq C(\mathbf{q})$  for all  $\mathbf{m} \in \mathbb{D}(r_0) \times \mathbb{D}(r_0)$ , whence in the same domain  $\|\nabla_{\mathbf{m}} F(\mathbf{m}; \xi)\|_2 \leq C\psi_1/(\Im(\xi) - 2)^2$  which implies the desired claim.

Finally assume  $\Im(\xi) > 0$ . Since by Lemma C.2  $\mathbf{F}$  maps  $\mathbb{C}_+^2$  to  $\mathbb{C}_+^2$  we can repeat the argument above with  $\mathbb{D}(r_0)$  replaced by  $\mathbb{D}_+(r_0) \equiv \mathbb{D}(r_0) \cap \mathbb{C}_+$ , and  $\mathbb{D}(r)$  replaced by  $\mathbb{D}_+(r)$ . We therefore conclude that the fixed  $\mathbf{m}(\xi; s, t)$  must have  $\Im(m_1) > 0$ ,  $\Im(m_2) > 0$ . Hence, as shown in the proof of Lemma C.2,  $\Im(H(\mathbf{m})) \leq 0$ , and therefore

$$|m_1| \leq \frac{\psi_1}{|\Im(-\xi + H(\mathbf{m}))|} \leq \frac{\psi_1}{\Im(\xi)}. \quad (141)$$

The same argument implies the bound  $|m_2| \leq \psi_2/\Im(\xi)$  as well.

Differentiability in the parameters follows immediately from the implicit using the fact that  $\mathbf{F}(\mathbf{m}; \xi; \mathbf{q}, \psi_1, \psi_2, \mu_\star, \mu_1)$  (with an abuse of notation, we added the dependence on) is differentiable with derivatives bounded by  $2\delta$  with respect to the parameters in  $\mathbb{D}(r) \times \mathbb{D}(r) \times \mathcal{N}$ , with  $\mathcal{N}$  a neighborhood of  $(\mathbf{q}^*, \psi_1^*, \psi_2^*, \mu_\star^*, \mu_1^*)$ , provided  $\Im(\xi) > \xi_0(r, \delta, \mathbf{q}^*, \psi_1^*, \psi_2^*, \mu_\star^*, \mu_1^*)$ .  $\square$

### C.3 Key lemma: Stieltjes transforms are approximate fixed point

Recall that  $(\bar{\theta}_a)_{a \in [N]} \sim_{iid} \mathbf{N}(0, \mathbf{I}_d)$ ,  $(\bar{\mathbf{x}}_i)_{i \in [n]} \sim_{iid} \mathbf{N}(0, \mathbf{I}_d)$ . We denote by  $\bar{\Theta} \in \mathbb{R}^{N \times d}$  the matrix whose  $a$ -th row is given by  $\bar{\theta}_a$ , and by  $\bar{\mathbf{X}} \in \mathbb{R}^{n \times d}$  the matrix whose  $i$ -th row is given by  $\bar{\mathbf{x}}_i$ . We consider a polynomial activation functions  $\varphi : \mathbb{R} \rightarrow \mathbb{R}$ . Denote  $\mu_k = \mathbb{E}[\varphi(G)\text{He}_k(G)]$  and  $\mu_\star^2 = \sum_{k \geq 2} \mu_k^2/k!$ . We define the following

matrices

$$\bar{\mathbf{J}} \equiv \frac{1}{\sqrt{d}} \varphi \left( \frac{1}{\sqrt{d}} \bar{\mathbf{X}} \bar{\boldsymbol{\Theta}}^\top \right), \quad (142)$$

$$\bar{\mathbf{J}}_1 \equiv \frac{\mu_1}{d} \bar{\mathbf{X}} \bar{\boldsymbol{\Theta}}^\top, \quad (143)$$

$$\bar{\mathbf{Q}} \equiv \frac{1}{d} \bar{\boldsymbol{\Theta}} \bar{\boldsymbol{\Theta}}^\top, \quad (144)$$

$$\bar{\mathbf{H}} \equiv \frac{1}{d} \bar{\mathbf{X}} \bar{\mathbf{X}}^\top, \quad (145)$$

as well as the block matrix  $\bar{\mathbf{A}} \in \mathbb{R}^{M \times M}$ ,  $M = N + n$ , defined by

$$\bar{\mathbf{A}} = \begin{bmatrix} s_1 \mathbf{I}_N + s_2 \bar{\mathbf{Q}} & \bar{\mathbf{J}}^\top + p \bar{\mathbf{J}}_1^\top \\ \bar{\mathbf{J}} + p \bar{\mathbf{J}}_1 & t_1 \mathbf{I}_n + t_2 \bar{\mathbf{H}} \end{bmatrix}. \quad (146)$$

In what follows, we will write  $\mathbf{q} = (s_1, s_2, t_1, t_2, p)$ .

We would like to calculate the asymptotic behavior of the following partial Stieltjes transforms

$$\begin{aligned} \bar{m}_{1,d}(\xi; \mathbf{q}) &= \frac{N}{d} \mathbb{E} \{ (\bar{\mathbf{A}} - \xi \mathbf{I}_M)^{-1}_{11} \} = \mathbb{E} [\bar{M}_{1,d}(\xi; \mathbf{q})], \\ \bar{m}_{2,d}(\xi; \mathbf{q}) &= \frac{n}{d} \mathbb{E} \{ (\bar{\mathbf{A}} - \xi \mathbf{I}_M)^{-1}_{N+1, N+1} \} = \mathbb{E} [\bar{M}_{2,d}(\xi; \mathbf{q})], \end{aligned} \quad (147)$$

where

$$\begin{aligned} \bar{M}_{1,d}(\xi; \mathbf{q}) &= \frac{1}{d} \text{Tr}_{[1, N]} [(\bar{\mathbf{A}} - \xi \mathbf{I}_M)^{-1}], \\ \bar{M}_{2,d}(\xi; \mathbf{q}) &= \frac{1}{d} \text{Tr}_{[N+1, N+n]} [(\bar{\mathbf{A}} - \xi \mathbf{I}_M)^{-1}]. \end{aligned} \quad (148)$$

Here, the partial trace notation  $\text{Tr}_{[a, b]}$  is defined as follows: for a matrix  $\mathbf{K} \in \mathbb{C}^{M \times M}$  and  $1 \leq a \leq b \leq M$ , define

$$\text{Tr}_{[a, b]}(\mathbf{K}) = \sum_{i=a}^b K_{ii}.$$

We denote by  $\psi_{1,d} = N/d$ ,  $\psi_{2,d} = n/d$  the aspect ratios at finite  $d$ . By Assumption 2,  $\psi_{i,d} \rightarrow \psi_i \in (0, \infty)$  for  $i \in \{1, 2\}$ . The crucial step consists in showing that the expected Stieltjes transforms  $\bar{m}_{1,d}, \bar{m}_{2,d}$  are approximate solutions of the fixed point equations (61).

**Lemma C.4.** *Assume that  $\varphi$  is a fixed polynomial with  $\mathbb{E}[\varphi(G)] = 0$  and  $\mu_1 \equiv \mathbb{E}[\varphi(G)G] \neq 0$ . Consider the linear regime Assumption 2. Then for any  $\mathbf{q} \in \mathcal{Q}$  and for any  $\xi_0 > 0$ , there exists  $C = C(\xi_0, \mathbf{q}, \psi_1, \psi_2, \varphi)$  which is uniformly bounded when  $(\mathbf{q}, \psi_1, \psi_2)$  is in a compact set, and a function  $\text{err}(d)$ , such that for all  $\xi \in \mathbb{C}_+$  with  $\Im(\xi) > \xi_0$ , we have*

$$\left| \bar{m}_{1,d} - F(\bar{m}_{1,d}, \bar{m}_{2,d}; \xi; \mathbf{q}, \psi_{1,d}, \psi_{2,d}, \mu_1, \mu_\star) \right| \leq C \cdot \text{err}(d), \quad (149)$$

$$\left| \bar{m}_{2,d} - F(\bar{m}_{2,d}, \bar{m}_{1,d}; \xi; \mathbf{q}, \psi_{2,d}, \psi_{1,d}, \mu_1, \mu_\star) \right| \leq C \cdot \text{err}(d), \quad (150)$$

with  $\lim_{d \rightarrow \infty} \text{err}(d) \rightarrow 0$ .

## C.4 Properties of Stieltjes transforms

The functions  $\xi \mapsto \bar{m}_{i,d}(\xi; \mathbf{q})/\psi_{i,d}$ ,  $i \in \{1, 2\}$ , can be shown to be Stieltjes transforms of certain probability measures on the reals line  $\mathbb{R}$  [HMRT19]. As such, they enjoy several useful properties (see, e.g., [AGZ09]). The next three lemmas are standard, and already stated in [HMRT19]. We reproduce them here without proof for the reader's convenience: although the present definition of the matrix  $\bar{\mathbf{A}}$  is slightly more general, the proofs are unchanged.

**Lemma C.5** (Lemma 7 in [HMRT19]). *The functions  $\xi \mapsto \overline{m}_{1,d}(\xi)$ ,  $\xi \mapsto \overline{m}_{2,d}(\xi)$  have the following properties:*

- (a)  $\xi \in \mathbb{C}_+$ , then  $|\overline{m}_{i,d}| \leq \psi_i/\Im(\xi)$  for  $i \in \{1, 2\}$ .
- (b)  $\overline{m}_{1,d}, \overline{m}_{2,d}$  are analytic on  $\mathbb{C}_+$  and map  $\mathbb{C}_+$  into  $\mathbb{C}_+$ .
- (c) Let  $\Omega \subseteq \mathbb{C}_+$  be a set with an accumulation point. If  $\overline{m}_{a,d}(\xi) \rightarrow m_a(\xi)$  for all  $\xi \in \Omega$ , then  $m_a(\xi)$  has a unique analytic continuation to  $\mathbb{C}_+$  and  $\overline{m}_{a,d}(\xi) \rightarrow m_a(\xi)$  for all  $\xi \in \mathbb{C}_+$ . Moreover, the convergence is uniform over compact sets  $\Omega \subseteq \mathbb{C}_+$ .

**Lemma C.6** (Lemma 8 in [HMRT19]). *Let  $\mathbf{W} \in \mathbb{R}^{M \times M}$  be a symmetric matrix, and denote by  $\mathbf{w}_i$  its  $i$ -th column, with the  $i$ -th entry set to 0. Let  $\mathbf{W}^{(i)} \equiv \mathbf{W} - \mathbf{w}_i \mathbf{e}_i^\top - \mathbf{e}_i \mathbf{w}_i^\top$ , where  $\mathbf{e}_i$  is the  $i$ -th element of the canonical basis (in other words,  $\mathbf{W}^{(i)}$  is obtained from  $\mathbf{W}$  by zeroing all elements in the  $i$ -th row and column except on the diagonal). Finally, let  $\xi \in \mathbb{C}_+$  with  $\Im(\xi) \geq \xi_0 > 0$ . Then for any subset  $S \subseteq [M]$ , we have*

$$\left| \text{Tr}_S[(\mathbf{W} - \xi \mathbf{I}_M)^{-1}] - \text{Tr}_S[(\mathbf{W}^{(i)} - \xi \mathbf{I}_M)^{-1}] \right| \leq \frac{3}{\xi_0}. \quad (151)$$

The next lemma establishes the concentration of Stieltjes transforms to its mean, whose proof is the same as the proof of Lemma 9 in [HMRT19].

**Lemma C.7** (Concentration). *Let  $\Im(\xi) \geq \xi_0 > 0$  and consider the partial Stieltjes transforms  $M_{i,d}(\xi; \mathbf{q})$  and  $\overline{M}_{i,d}(\xi; \mathbf{q})$  as per Eq. (135). Then there exists  $c_0 = c_0(\xi_0)$  such that, for  $i \in \{1, 2\}$ ,*

$$\mathbb{P}(|\overline{M}_{i,d}(\xi; \mathbf{q}) - \mathbb{E}\overline{M}_{i,d}(\xi; \mathbf{q})| \geq u) \leq 2e^{-c_0 d u^2}, \quad (152)$$

$$\mathbb{P}(|M_{i,d}(\xi; \mathbf{q}) - \mathbb{E}M_{i,d}(\xi; \mathbf{q})| \geq u) \leq 2e^{-c_0 d u^2}. \quad (153)$$

*In particular, if  $\Im(\xi) > 0$ , then  $|\overline{M}_{i,d}(\xi; \mathbf{q}) - \mathbb{E}\overline{M}_{i,d}(\xi; \mathbf{q})| \rightarrow 0$ ,  $|M_{i,d}(\xi; \mathbf{q}) - \mathbb{E}M_{i,d}(\xi; \mathbf{q})| \rightarrow 0$  almost surely and in  $L^1$ .*

**Lemma C.8** (Lemma 5 in [GMMM19]). *Assume  $\sigma$  is an activation function with  $\sigma(u)^2 \leq c_0 \exp(c_1 u^2/2)$  for some constants  $c_0 > 0$  and  $c_1 < 1$  (this is implied by Assumption 1). Then*

$$(a) \mathbb{E}_{G \sim \mathcal{N}(0,1)}[\sigma(G)^2] < \infty.$$

$$(b) \text{ Let } \|\mathbf{w}\|_2 = 1. \text{ Then there exists } d_0 = d_0(c_1) \text{ such that, for } \mathbf{x} \sim \text{Unif}(\mathbb{S}^{d-1}(\sqrt{d})),$$

$$\sup_{d \geq d_0} \mathbb{E}_{\mathbf{x}}[\sigma(\langle \mathbf{w}, \mathbf{x} \rangle)^2] < \infty. \quad (154)$$

$$(c) \text{ Let } \|\mathbf{w}\|_2 = 1. \text{ Then there exists a coupling of } G \sim \mathcal{N}(0,1) \text{ and } \mathbf{x} \sim \text{Unif}(\mathbb{S}^{d-1}(\sqrt{d})) \text{ such that}$$

$$\lim_{d \rightarrow \infty} \mathbb{E}_{\mathbf{x}, G}[(\sigma(\langle \mathbf{w}, \mathbf{x} \rangle) - \sigma(G))^2] = 0. \quad (155)$$

## C.5 Leave-one-out argument: Proof of Lemma C.4

Throughout the proof, we write  $F(d) = O_d(G(d))$  if there exists a constant  $C = C(\xi_0, \mathbf{q}, \psi_1, \psi_2, \varphi)$  which is uniformly bounded when  $(\xi_0, \mathbf{q}, \psi_1, \psi_2)$  is in a compact set, such that  $|F(d)| \leq C \cdot |G(d)|$ . We write  $F(d) = o_d(G(d))$  if for any  $\varepsilon > 0$ , there exists a constant  $C = C(\varepsilon, \xi_0, \mathbf{q}, \psi_1, \psi_2, \varphi)$  which is uniformly bounded when  $(\xi_0, \mathbf{q}, \psi_1, \psi_2)$  is in a compact set, such that  $|F(d)| \leq \varepsilon \cdot |G(d)|$  for any  $d \geq C$ . We use  $C$  to denote generically such a constant, that can change from line to line.

We write  $F(d) = O_{d,\mathbb{P}}(G(d))$  if for any  $\delta > 0$ , there exists constant  $K = K(\delta, \xi_0, \mathbf{q}, \psi_1, \psi_2, \varphi)$ ,  $d_0 = d_0(\delta, \xi_0, \mathbf{q}, \psi_1, \psi_2, \varphi)$  which are uniformly bounded when  $(\xi_0, \mathbf{q}, \psi_1, \psi_2)$  is in a compact set, such that  $\mathbb{P}(|F(d)| > K|G(d)|) \leq \delta$  for any  $d \geq d_0$ . We write  $F(d) = o_{d,\mathbb{P}}(G(d))$  if for any  $\varepsilon, \delta > 0$ , there exists constant  $d_0 = d_0(\varepsilon, \delta, \xi_0, \mathbf{q}, \psi_1, \psi_2, \varphi)$  which are uniformly bounded when  $(\xi_0, \mathbf{q}, \psi_1, \psi_2)$  is in a compact set, such that  $\mathbb{P}(|F(d)| > \varepsilon|G(d)|) \leq \delta$  for any  $d \geq d_0$ .

We will assume  $p = 0$  throughout the proof. For  $p \neq 0$ , the lemma holds by viewing  $\bar{\mathbf{J}} + p\bar{\mathbf{J}}_1 = \varphi_\star(\mathbf{X}\boldsymbol{\Theta}^\top/\sqrt{d})/\sqrt{d}$  as a new kernel inner product matrix, with  $\varphi_\star(x) = \varphi(x) + p\mu_1x$ .

**Step 1. Calculate the Schur complement and define some notations.**

Let  $\bar{\mathbf{A}}_{\cdot,N} \in \mathbb{R}^{M-1}$  be the  $N$ -th column of  $\bar{\mathbf{A}}$ , with the  $N$ -th entry removed. We further denote by  $\bar{\mathbf{B}} \in \mathbb{R}^{(M-1) \times (M-1)}$  be the the matrix obtained from  $\bar{\mathbf{A}}$  by removing the  $N$ -th column and  $N$ -th row. Applying Schur complement formula with respect to element  $(N, N)$ , we get

$$\bar{m}_{1,d} = \psi_{1,d} \mathbb{E} \left\{ \left( -\xi + s_1 + s_2 \|\bar{\boldsymbol{\theta}}_N\|_2^2/d - \bar{\mathbf{A}}_{\cdot,N}^\top (\bar{\mathbf{B}} - \xi \mathbf{I}_{M-1})^{-1} \bar{\mathbf{A}}_{\cdot,N} \right)^{-1} \right\}. \quad (156)$$

We decompose the vectors  $\bar{\boldsymbol{\theta}}_a, \bar{\mathbf{x}}_i$  in the components along  $\bar{\boldsymbol{\theta}}_N$  and the orthogonal component:

$$\bar{\boldsymbol{\theta}}_a = \eta_a \frac{\bar{\boldsymbol{\theta}}_N}{\|\bar{\boldsymbol{\theta}}_N\|_2} + \tilde{\boldsymbol{\theta}}_a, \quad \langle \bar{\boldsymbol{\theta}}_N, \tilde{\boldsymbol{\theta}}_a \rangle = 0, \quad a \in [N-1], \quad (157)$$

$$\bar{\mathbf{x}}_i = u_i \frac{\bar{\boldsymbol{\theta}}_N}{\|\bar{\boldsymbol{\theta}}_N\|} + \tilde{\mathbf{x}}_i, \quad \langle \bar{\boldsymbol{\theta}}_N, \tilde{\mathbf{x}}_i \rangle = 0, \quad i \in [n]. \quad (158)$$

Note that  $\{\eta_a\}_{a \in [N-1]}, \{u_i\}_{i \in [n]} \sim iid \mathbf{N}(0, 1)$  are independent of all the other random variables, and  $\{\tilde{\boldsymbol{\theta}}_a\}_{a \in [N-1]}, \{\tilde{\mathbf{x}}_i\}_{i \in [n]}$  are conditionally independent given  $\bar{\boldsymbol{\theta}}_N$ , with  $\tilde{\boldsymbol{\theta}}_a, \tilde{\mathbf{x}}_i \sim \mathbf{N}(\mathbf{0}, \mathbf{P}_\perp)$ , where  $\mathbf{P}_\perp$  is the projector orthogonal to  $\bar{\boldsymbol{\theta}}_N$ .

With this decomposition we have

$$\bar{Q}_{a,b} = \frac{1}{d} \left( \eta_a \eta_b + \langle \tilde{\boldsymbol{\theta}}_a, \tilde{\boldsymbol{\theta}}_b \rangle \right), \quad a, b \in [N-1], \quad (159)$$

$$\bar{J}_{i,a} = \frac{1}{\sqrt{d}} \varphi \left( \frac{1}{\sqrt{d}} \langle \tilde{\mathbf{x}}_i, \tilde{\boldsymbol{\theta}}_a \rangle + \frac{1}{\sqrt{d}} u_i \eta_a \right), \quad a \in [N-1], i \in [n], \quad (160)$$

$$\bar{H}_{ij} = \frac{1}{d} \left( u_i u_j + \langle \tilde{\mathbf{x}}_i, \tilde{\mathbf{x}}_j \rangle \right), \quad i, j \in [n]. \quad (161)$$

Further we have  $\bar{\mathbf{A}}_{\cdot,N} = (\bar{A}_{1,N}, \dots, \bar{A}_{M-1,N})^\top \in \mathbb{R}^{M-1}$  with

$$\bar{A}_{i,N} = \begin{cases} \frac{1}{d} s_2 \eta_i \|\bar{\boldsymbol{\theta}}_N\|_2 & \text{if } i \leq N-1, \\ \frac{1}{\sqrt{d}} \varphi \left( \frac{1}{\sqrt{d}} u_i \|\bar{\boldsymbol{\theta}}_N\|_2 \right) & \text{if } i \geq N. \end{cases} \quad (162)$$

We next write  $\bar{\mathbf{B}}$  as the sum of three terms:

$$\bar{\mathbf{B}} = \tilde{\mathbf{B}} + \boldsymbol{\Delta} + \mathbf{E}_0, \quad (163)$$

where

$$\tilde{\mathbf{B}} = \begin{bmatrix} s_1 \mathbf{I}_N + s_2 \tilde{\mathbf{Q}} & \tilde{\mathbf{J}}^\top \\ \tilde{\mathbf{J}} & t_1 \mathbf{I}_n + t_2 \tilde{\mathbf{H}} \end{bmatrix}, \quad \boldsymbol{\Delta} = \begin{bmatrix} \frac{s_2}{d} \boldsymbol{\eta} \boldsymbol{\eta}^\top & \frac{\mu_1}{d} \boldsymbol{\eta} \mathbf{u}^\top \\ \frac{\mu_1}{d} \mathbf{u} \boldsymbol{\eta}^\top & \frac{\mu_1^2}{d} \mathbf{u} \mathbf{u}^\top \end{bmatrix}, \quad \mathbf{E}_0 = \begin{bmatrix} \mathbf{0} & \mathbf{E}_1^\top \\ \mathbf{E}_1 & \mathbf{0} \end{bmatrix}, \quad (164)$$

where

$$\tilde{Q}_{a,b} = \frac{1}{d} \langle \tilde{\boldsymbol{\theta}}_a, \tilde{\boldsymbol{\theta}}_b \rangle, \quad a, b \in [N-1], \quad (165)$$

$$\tilde{J}_{i,a} = \frac{1}{\sqrt{d}} \varphi \left( \frac{1}{\sqrt{d}} \langle \tilde{\mathbf{x}}_i, \tilde{\boldsymbol{\theta}}_a \rangle \right), \quad a \in [N-1], i \in [n], \quad (166)$$

$$\tilde{H}_{ij} = \frac{1}{d} \langle \tilde{\mathbf{x}}_i, \tilde{\mathbf{x}}_j \rangle, \quad i, j \in [n]. \quad (167)$$

Further, for  $i \in [n], a \in [N-1]$ ,

$$E_{1,ia} = \frac{1}{\sqrt{d}} \left[ \varphi \left( \frac{1}{\sqrt{d}} \langle \tilde{\mathbf{x}}_i, \tilde{\boldsymbol{\theta}}_a \rangle + \frac{1}{\sqrt{d}} u_i \eta_a \right) - \varphi \left( \frac{1}{\sqrt{d}} \langle \tilde{\mathbf{x}}_i, \tilde{\boldsymbol{\theta}}_a \rangle \right) - \frac{\mu_1}{\sqrt{d}} u_i \eta_a \right] \quad (168)$$

$$= \frac{1}{\sqrt{d}} \left[ \varphi_\perp \left( \frac{1}{\sqrt{d}} \langle \tilde{\mathbf{x}}_i, \tilde{\boldsymbol{\theta}}_a \rangle + \frac{1}{\sqrt{d}} u_i \eta_a \right) - \varphi_\perp \left( \frac{1}{\sqrt{d}} \langle \tilde{\mathbf{x}}_i, \tilde{\boldsymbol{\theta}}_a \rangle \right) \right], \quad (169)$$

where  $\varphi_\perp(x) \equiv \varphi(x) - \mu_1 x$ .

**Step 2. Perturbation bound for the Schur complement.**

Denote

$$\omega_1 = \left( -\xi + s_1 + s_2 \|\bar{\boldsymbol{\theta}}_N\|_2^2/d - \bar{\mathbf{A}}_{\cdot,N}^\top (\bar{\mathbf{B}} - \xi \mathbf{I}_{M-1})^{-1} \bar{\mathbf{A}}_{\cdot,N} \right)^{-1}, \quad (170)$$

$$\omega_2 = \left( -\xi + s_1 + s_2 - \bar{\mathbf{A}}_{\cdot,N}^\top (\tilde{\mathbf{B}} + \boldsymbol{\Delta} - \xi \mathbf{I}_{M-1})^{-1} \bar{\mathbf{A}}_{\cdot,N} \right)^{-1}. \quad (171)$$

Note we have  $\bar{m}_{1,d} = \psi_{1,d} \mathbb{E}[\omega_1]$ . Combining Lemma C.9, C.10, and C.11 below, we have

$$|\omega_1 - \omega_2| \leq O_d(1) \cdot \left| \|\bar{\boldsymbol{\theta}}_N\|_2^2/d - 1 \right| + O_d(1) \cdot \|\bar{\mathbf{A}}_{\cdot,N}\|_2^2 \cdot \|\mathbf{E}_1\|_{\text{op}} = o_d(1).$$

Moreover, by Lemma C.9,  $|\omega_1 - \omega_2|$  is deterministically bounded by  $2/\xi_0$ . This gives

$$|\bar{m}_{1,d} - \psi_{1,d} \mathbb{E}[\omega_2]| \leq \psi_{1,d} \mathbb{E}[|\omega_1 - \omega_2|] = o_d(1). \quad (172)$$

**Lemma C.9.** *Using the definitions of  $\omega_1$  and  $\omega_2$  as in Eq. (170) and (171), for  $\Im \xi \geq \xi_0$ , we have*

$$|\omega_1 - \omega_2| \leq \left[ s_2 \|\bar{\boldsymbol{\theta}}_N\|_2^2/d - 1/\xi_0^2 + 2\|\bar{\mathbf{A}}_{\cdot,N}\|_2^2 \|\mathbf{E}_1\|_{\text{op}}/\xi_0^4 \right] \wedge [2/\xi_0].$$

*Proof of Lemma C.9.* Note that

$$\Im(-\omega_1^{-1}) \geq \Im \xi + \Im(\bar{\mathbf{A}}_{\cdot,N}^\top (\bar{\mathbf{B}} - \xi \mathbf{I}_{M-1})^{-1} \bar{\mathbf{A}}_{\cdot,N}) \geq \Im \xi > \xi_0.$$

Hence we have  $|\omega_1| \leq 1/\xi_0$ , and, using a similar argument,  $|\omega_2| \leq 1/\xi_0$ . Hence we get the bound  $|\omega_1 - \omega_2| \leq 2/\xi_0$ .

Denote

$$\omega_{1.5} = \left( -\xi + s_1 + s_2 - \bar{\mathbf{A}}_{\cdot,N}^\top (\bar{\mathbf{B}} - \xi \mathbf{I}_{M-1})^{-1} \bar{\mathbf{A}}_{\cdot,N} \right)^{-1},$$

we get

$$|\omega_1 - \omega_{1.5}| = s_2 \left| \omega_1 (\|\bar{\boldsymbol{\theta}}_N\|_2^2/d - 1) \omega_{1.5} \right| \leq s_2 \left| \|\bar{\boldsymbol{\theta}}_N\|_2^2/d - 1 \right| / \xi_0^2.$$

Moreover, we have

$$\begin{aligned} |\omega_{1.5} - \omega_2| &= |\omega_{1.5} \omega_2 \bar{\mathbf{A}}_{\cdot,N}^\top [(\tilde{\mathbf{B}} + \boldsymbol{\Delta} - \xi \mathbf{I}_{M-1})^{-1} - (\bar{\mathbf{B}} + \boldsymbol{\Delta} + \mathbf{E}_0 - \xi \mathbf{I}_{M-1})^{-1}] \bar{\mathbf{A}}_{\cdot,N}| \\ &= |\omega_{1.5} \omega_2 \bar{\mathbf{A}}_{\cdot,N}^\top (\tilde{\mathbf{B}} + \boldsymbol{\Delta} - \xi \mathbf{I}_{M-1})^{-1} \mathbf{E}_0 (\tilde{\mathbf{B}} + \boldsymbol{\Delta} + \mathbf{E}_0 - \xi \mathbf{I}_{M-1})^{-1} \bar{\mathbf{A}}_{\cdot,N}| \\ &\leq (1/\xi_0^2) \cdot \|\bar{\mathbf{A}}_{\cdot,N}\|_2^2 (1/\xi_0^2) \|\mathbf{E}_0\|_{\text{op}} \leq 2\|\mathbf{E}_1\|_{\text{op}} \|\bar{\mathbf{A}}_{\cdot,N}\|_2^2 / \xi_0^4. \end{aligned}$$

This proves the lemma.  $\square$

**Lemma C.10.** *Under the assumptions of Lemma C.4, we have*

$$\|\mathbf{E}_1\|_{\text{op}} = O_{d,\mathbb{P}}(\text{Poly}(\log d)/d^{1/2}). \quad (173)$$

*Proof of Lemma C.10.* Define  $\mathbf{z}_i = \tilde{\boldsymbol{\theta}}_i$  for  $i \in [N-1]$ ,  $\mathbf{z}_i = \tilde{\mathbf{x}}_{i-N+1}$  for  $N \leq i \leq M-1$ , and  $\zeta_i = \eta_i$  for  $i \in [N-1]$ ,  $\zeta_i = u_{i-N+1}$  for  $N \leq i \leq M-1$ . Consider the symmetric matrix  $\mathbf{E} \in \mathbb{R}^{(M-1) \times (M-1)}$  with  $E_{ii} = 0$ , and

$$E_{ij} = \frac{1}{\sqrt{d}} \left[ \varphi_\perp \left( \frac{1}{\sqrt{d}} \langle \mathbf{z}_i, \mathbf{z}_j \rangle + \frac{1}{\sqrt{d}} \zeta_i \zeta_j \right) - \varphi_\perp \left( \frac{1}{\sqrt{d}} \langle \mathbf{z}_i, \mathbf{z}_j \rangle \right) \right]. \quad (174)$$

Since  $\mathbf{E}_1$  is a sub-matrix of  $\mathbf{E}$ , we have  $\|\mathbf{E}_1\|_{\text{op}} \leq \|\mathbf{E}\|_{\text{op}}$ . By the intermediate value theorem

$$\mathbf{E} = \frac{1}{\sqrt{d}} \boldsymbol{\Xi} \mathbf{F}_1 \boldsymbol{\Xi} + \frac{1}{2d} \boldsymbol{\Xi}^2 \mathbf{F}_2 \boldsymbol{\Xi}^2, \quad (175)$$

$$\boldsymbol{\Xi} \equiv \text{diag}(\zeta_1, \dots, \zeta_{M-1}), \quad (176)$$

$$F_{1,ij} \equiv \frac{1}{\sqrt{d}} \varphi'_\perp \left( \frac{1}{\sqrt{d}} \langle \mathbf{z}_i, \mathbf{z}_j \rangle \right) \mathbf{1}_{i \neq j}, \quad (177)$$

$$F_{2,ij} \equiv \frac{1}{\sqrt{d}} \varphi''_\perp(\tilde{z}_{ij}) \mathbf{1}_{i \neq j}, \quad \tilde{z}_{ij} \in \left[ \frac{1}{\sqrt{d}} \langle \mathbf{z}_i, \mathbf{z}_j \rangle, \frac{1}{\sqrt{d}} \langle \mathbf{z}_i, \mathbf{z}_j \rangle + \frac{1}{\sqrt{d}} \zeta_i \zeta_j \right]. \quad (178)$$

Hence we get

$$\|\mathbf{E}\|_{\text{op}} \leq (\|\mathbf{F}_1\|_{\text{op}}/\sqrt{d})\|\Xi\|_{\text{op}}^2 + (\|\mathbf{F}_2\|_{\text{op}}/d)\|\Xi\|_{\text{op}}^4.$$

Note that  $\varphi''_{\perp}(x) = \varphi''(x)$  is a polynomial with some fixed degree  $\bar{k}$ . Therefore we have

$$\begin{aligned} \mathbb{E}\{\|\mathbf{F}_2\|_F^2\} &= [M(M-1)/d] \cdot \mathbb{E}[\varphi''_{\perp}(\tilde{z}_{12})^2] \leq O_d(d) \cdot \mathbb{E}[(1 + |\tilde{z}_{12}|)^{2\bar{k}}] \\ &\leq O_d(d) \cdot \left\{ \mathbb{E}[(1 + |\langle \mathbf{z}_i, \mathbf{z}_j \rangle|/\sqrt{d})^{2\bar{k}}] + \mathbb{E}[(1 + |\langle \mathbf{z}_i, \mathbf{z}_j \rangle + \zeta_i \zeta_j/\sqrt{d}|)^{2\bar{k}}] \right\} = O_d(d). \end{aligned}$$

Moreover, by the fact that  $\varphi'_{\perp}$  is a polynomial with  $\mathbb{E}[\varphi'_{\perp}(G)] = 0$ , and by Theorem 1.7 in [FM19], we have  $\|\mathbf{F}_1\|_{\text{op}} = O_{d,\mathbb{P}}(1)$ . By the concentration bound for  $\chi$ -squared random variable, we get  $\|\Xi\|_{\text{op}} = O_{d,\mathbb{P}}(\sqrt{\log d})$ . Therefore, we have

$$\|\mathbf{E}\|_{\text{op}} \leq O_{d,\mathbb{P}}(d^{-1/2})O_{d,\mathbb{P}}(\text{Poly}(\log d)) + O_{d,\mathbb{P}}(d^{-1/2})O_{d,\mathbb{P}}(\text{Poly}(\log d)) = O_{d,\mathbb{P}}(\text{Poly}(\log d)/d^{-1/2}).$$

This proves the lemma.  $\square$

**Lemma C.11.** *Under the assumptions of Lemma C.4, we have*

$$\|\bar{\mathbf{A}}_{\cdot,N}\|_2 = O_{d,\mathbb{P}}(1). \quad (179)$$

*Proof of Lemma C.11.* Recall the definition of  $\bar{\mathbf{A}}_{\cdot,N}$  as in Eq. (162). Denote  $\mathbf{a}_1 = s_2 \boldsymbol{\eta} \|\bar{\boldsymbol{\theta}}_N\|_2/d \in \mathbb{R}^{N-1}$ , and  $\mathbf{a}_2 = \varphi(\mathbf{u} \|\bar{\boldsymbol{\theta}}_N\|_2/\sqrt{d})/\sqrt{d} \in \mathbb{R}^n$ , where  $\boldsymbol{\eta} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_{N-1})$ , and  $\mathbf{u} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_n)$ . Then  $\bar{\mathbf{A}}_{\cdot,N} = (\mathbf{a}_1; \mathbf{a}_2) \in \mathbb{R}^{n+N-1}$ .

For  $\mathbf{a}_1$ , note we have  $\|\mathbf{a}_1\|_2 = |s_2| \cdot \|\boldsymbol{\eta}\|_2 \|\bar{\boldsymbol{\theta}}_N\|_2/d$  where  $\boldsymbol{\eta} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_{N-1})$  and  $\bar{\boldsymbol{\theta}}_N \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d)$  are independent. Hence we have

$$\mathbb{E}[\|\mathbf{a}_1\|_2^2] = s_2^2 \mathbb{E}[\|\boldsymbol{\eta}\|_2^2 \|\bar{\boldsymbol{\theta}}_N\|_2^2]/d^2 = O_d(1).$$

For  $\mathbf{a}_2$ , note  $\varphi$  is a polynomial with some fixed degree  $\bar{k}$ , hence we have

$$\mathbb{E}[\|\mathbf{a}_2\|_2^2] = \mathbb{E}[\varphi(u_i \|\bar{\boldsymbol{\theta}}_N\|_2/\sqrt{d})^2] = O_d(1).$$

This proves the lemma.  $\square$

### Step 3. Simplification using Sherman-Morrison-Woodbury formula.

Notice that  $\boldsymbol{\Delta}$  is a matrix with rank at most two. Indeed

$$\boldsymbol{\Delta} = \mathbf{U} \mathbf{M} \mathbf{U}^T \in \mathbb{R}^{(M-1) \times (M-1)}, \quad \mathbf{U} = \frac{1}{\sqrt{d}} \begin{bmatrix} \boldsymbol{\eta} & \mathbf{0} \\ \mathbf{0} & \mathbf{u} \end{bmatrix} \in \mathbb{R}^{(M-1) \times 2}, \quad \mathbf{M} = \begin{bmatrix} s_2 & \mu_1 \\ \mu_1 & t_2 \end{bmatrix} \in \mathbb{R}^{2 \times 2}. \quad (180)$$

Since we assumed  $\mathbf{q} \in \mathcal{Q}$  so that  $|s_2 t_2| \leq \mu_1^2/2$ , the matrix  $\mathbf{M}$  is invertible with  $\|\mathbf{M}^{-1}\|_{\text{op}} \leq C$ .

Recall the definition of  $\omega_2$  in Eq. (171). By the Sherman-Morrison-Woodbury formula, we get

$$\omega_2 = \left( -\xi + s_1 + s_2 - v_1 + \mathbf{v}_2^T (\mathbf{M}^{-1} + \mathbf{V}_3)^{-1} \mathbf{v}_2 \right)^{-1}, \quad (181)$$

where

$$v_1 = \bar{\mathbf{A}}_{\cdot,N}^T (\tilde{\mathbf{B}} - \xi \mathbf{I}_{M-1})^{-1} \bar{\mathbf{A}}_{\cdot,N}, \quad v_2 = \mathbf{U}^T (\tilde{\mathbf{B}} - \xi \mathbf{I}_{M-1})^{-1} \bar{\mathbf{A}}_{\cdot,N}, \quad \mathbf{V}_3 = \mathbf{U}^T (\tilde{\mathbf{B}} - \xi \mathbf{I}_{M-1})^{-1} \mathbf{U}. \quad (182)$$

We define

$$\bar{v}_1 = s_2^2 \bar{m}_{1,d} + (\mu_1^2 + \mu_*^2) \bar{m}_{2,d}, \quad \bar{\mathbf{v}}_2 = \begin{bmatrix} s_2 \bar{m}_{1,d} \\ \mu_1 \bar{m}_{2,d} \end{bmatrix}, \quad \bar{\mathbf{V}}_3 = \begin{bmatrix} \bar{m}_{1,d} & 0 \\ 0 & \bar{m}_{2,d} \end{bmatrix}, \quad (183)$$

and

$$\omega_3 = \left( -\xi + s_1 + s_2 - \bar{v}_1 + \bar{\mathbf{v}}_2^T (\mathbf{M}^{-1} + \bar{\mathbf{V}}_3)^{-1} \bar{\mathbf{v}}_2 \right)^{-1}. \quad (184)$$

By auxiliary Lemmas C.12, C.13, and C.14 below, we get

$$\mathbb{E}[|\omega_2 - \omega_3|] = o_d(1),$$

Combining with Eq. (172) we get

$$|\bar{m}_{1,d} - \psi_{1,d}\omega_3| = o_d(1).$$

Elementary algebra simplifying Eq. (184) gives  $\psi_{1,d}\omega_3 = \mathbf{F}(\bar{m}_{1,d}, \bar{m}_{2,d}; \xi; \mathbf{q}, \psi_{1,d}, \psi_{2,d}, \mu_1, \mu_*)$ . This proves Eq. (149) in Lemma C.4. Eq. (150) follows by the same argument (exchanging  $N$  and  $n$ ). In the rest of this section, we prove auxiliary Lemmas C.12, C.13, and C.14.

**Lemma C.12.** *Using the formula of  $\omega_2$  and  $\omega_3$  as in Eq. (181) and (184), for  $\Im\xi \geq \xi_0$ , we have*

$$|\omega_2 - \omega_3| \leq O_d(1) \cdot \left\{ \left[ |v_1 - \bar{v}_1| + \|\bar{\mathbf{v}}_2\|_2^2 \|(\mathbf{M}^{-1} + \mathbf{V}_3)^{-1}\|_{\text{op}} \|(\mathbf{M}^{-1} + \bar{\mathbf{V}}_3)^{-1}\|_{\text{op}} \|\mathbf{V}_3 - \bar{\mathbf{V}}_3\|_{\text{op}} \right. \right. \\ \left. \left. + (\|\mathbf{v}_2\|_2 + \|\bar{\mathbf{v}}_2\|_2) \|(\mathbf{M}^{-1} + \mathbf{V}_3)^{-1}\|_{\text{op}} \|\mathbf{v}_2 - \bar{\mathbf{v}}_2\|_2 \right] \wedge 1 \right\}.$$

*Proof of Lemma C.12.* Denote

$$\omega_{2.5} = \left( -\xi + s_1 + s_2 - \bar{v}_1 + \mathbf{v}_2^\top (\mathbf{M}^{-1} + \mathbf{V}_3)^{-1} \mathbf{v}_2 \right)^{-1}$$

We have

$$|\omega_2 - \omega_{2.5}| = |\omega_2(v_1 - \bar{v}_1)\omega_{2.5}| \leq |v_1 - \bar{v}_1|/\xi_0^2.$$

Moreover, we have

$$|\omega_{2.5} - \omega_3| \leq (1/\xi_0^2) |\mathbf{v}_2^\top (\mathbf{M}^{-1} + \mathbf{V}_3)^{-1} \mathbf{v}_2 - \bar{\mathbf{v}}_2^\top (\mathbf{M}^{-1} + \bar{\mathbf{V}}_3)^{-1} \bar{\mathbf{v}}_2| \\ \leq (1/\xi_0^2) \left\{ (\|\mathbf{v}_2\|_2 + \|\bar{\mathbf{v}}_2\|_2) \|(\mathbf{M}^{-1} + \mathbf{V}_3)^{-1}\|_{\text{op}} \|\mathbf{v}_2 - \bar{\mathbf{v}}_2\|_2 \right. \\ \left. + \|\bar{\mathbf{v}}_2\|_2^2 \|(\mathbf{M}^{-1} + \mathbf{V}_3)^{-1}\|_{\text{op}} \|(\mathbf{M}^{-1} + \bar{\mathbf{V}}_3)^{-1}\|_{\text{op}} \|\mathbf{V}_3 - \bar{\mathbf{V}}_3\|_{\text{op}} \right\}.$$

Combining with  $|\omega_2 - \omega_3| \leq |\omega_2| + |\omega_3| \leq 2/\xi_0$  proves the lemma.  $\square$

**Lemma C.13.** *Under the assumptions of Lemma C.4, we have (following the notations of Eq. (182) and (183))*

$$\|\bar{\mathbf{v}}_2\|_2 = O_d(1), \tag{185}$$

$$|v_1 - \bar{v}_1| = o_{d,\mathbb{P}}(1), \tag{186}$$

$$\|\mathbf{v}_2 - \bar{\mathbf{v}}_2\|_2 = o_{d,\mathbb{P}}(1), \tag{187}$$

$$\|\mathbf{V}_3 - \bar{\mathbf{V}}_3\|_{\text{op}} = o_{d,\mathbb{P}}(1). \tag{188}$$

*Proof of Lemma C.13.* The first bound is because (see Lemma C.5 for the boundedness of  $\bar{m}_{1,d}$  and  $\bar{m}_{2,d}$ )

$$\|\bar{\mathbf{v}}_2\|_2 \leq |s_2| \cdot |\bar{m}_{1,d}| + |\mu_1| \cdot |\bar{m}_{2,d}| \leq (\psi_1 + \psi_2)(|s_2| + |\mu_1|)/\xi_0 = O_d(1).$$

In the following, we limit ourselves to proving Eq. (186), since Eq. (187) and (188) follow by similar arguments.

Let  $\mathbf{R} \equiv (\tilde{\mathbf{B}} - \xi \mathbf{I}_{M-1})^{-1}$ . Then we have  $\|\mathbf{R}\|_{\text{op}} \leq 1/\xi_0$ . Define  $\mathbf{a}, \mathbf{h}$  as

$$\mathbf{a} = \bar{\mathbf{A}}_{\cdot,N} = \left[ \frac{1}{d} s_2 \boldsymbol{\eta}^\top \|\bar{\boldsymbol{\theta}}_N\|_2, \frac{1}{\sqrt{d}} \varphi\left(\frac{1}{\sqrt{d}} \mathbf{u}^\top \|\bar{\boldsymbol{\theta}}_N\|_2\right) \right]^\top, \\ \mathbf{h} = \left[ \frac{1}{\sqrt{d}} s_2 \boldsymbol{\eta}^\top, \frac{1}{\sqrt{d}} \varphi(\mathbf{u}^\top) \right]^\top.$$

Then by the definition of  $v_1$  in Eq. (182), we have  $v_1 = \mathbf{a}^\top \mathbf{R} \mathbf{a}$ . Note we have

$$\|\mathbf{h} - \mathbf{a}\|_2 \leq (s_2 \|\boldsymbol{\eta}\|_2 + \|\varphi'(\mathbf{u} \odot \boldsymbol{\xi})\|_2) \cdot \|\bar{\boldsymbol{\theta}}_N\|_2 / \sqrt{d} - 1 / \sqrt{d},$$

for some  $\boldsymbol{\xi} = (\xi_1, \dots, \xi_n)^\top$  with  $\xi_i$  between  $\|\bar{\boldsymbol{\theta}}_N\|_2 / \sqrt{d}$  and 1. Since  $\|\bar{\boldsymbol{\theta}}_N\|_2 / \sqrt{d} - 1 = O_{d,\mathbb{P}}(\sqrt{\log d} / \sqrt{d})$ ,  $\|\boldsymbol{\eta}\|_2 = O_{d,\mathbb{P}}(\sqrt{d})$ , and  $\|\varphi'(\mathbf{u} \cdot \boldsymbol{\xi})\|_2 = O_{d,\mathbb{P}}(\text{Poly}(\log d) \cdot \sqrt{d})$ , we have

$$\|\mathbf{h} - \mathbf{a}\|_2 = o_{d,\mathbb{P}}(1).$$



By Lemma C.11 we have  $\|\mathbf{a}\|_2 = O_{d,\mathbb{P}}(1)$  and hence  $\|\mathbf{h}\|_2 = O_{d,\mathbb{P}}(1)$ . Combining all these bounds, we have

$$|v_1 - \mathbf{h}^\top \mathbf{R} \mathbf{h}| = |\mathbf{a}^\top \mathbf{R} \mathbf{a} - \mathbf{h}^\top \mathbf{R} \mathbf{h}| \leq (\|\mathbf{a}\|_2 + \|\mathbf{h}\|_2) \|\mathbf{h} - \mathbf{a}\|_2 \|\mathbf{R}\|_{\text{op}} = o_{d,\mathbb{P}}(1). \quad (189)$$

Denote by  $\mathbf{D}$  the covariance matrix of  $\mathbf{h}$ . Since  $\mathbf{h}$  has independent elements,  $\mathbf{D}$  a diagonal matrix with  $\max_i D_{ii} = \max_i \text{Var}(h_i) \leq C/d$ . Since  $\mathbb{E}[\mathbf{h}] = \mathbf{0}$ , we have

$$\mathbb{E}\{\mathbf{h}^\top \mathbf{R} \mathbf{h} | \mathbf{R}\} = \text{Tr}(\mathbf{D} \mathbf{R}). \quad (190)$$

We next compute  $\text{Var}(\mathbf{h}^\top \mathbf{R} \mathbf{h} | \mathbf{R})$ , (for a complex matrix, denote  $\mathbf{R}^\top$  to be the transpose of  $\mathbf{R}$ , and  $\mathbf{R}^*$  to be the conjugate transpose of  $\mathbf{R}$ )

$$\begin{aligned} \text{Var}(\mathbf{h}^\top \mathbf{R} \mathbf{h} | \mathbf{R}) &= \left\{ \sum_{i_1, i_2, i_3, i_4} \mathbb{E}[h_{i_1} R_{i_1 i_2} h_{i_2} h_{i_3} R_{i_3 i_4}^* h_{i_4}] \right\} - |\text{Tr}(\mathbf{D} \mathbf{R})|^2 \\ &= \left\{ \left( \sum_{i_1=i_2=i_3=i_4} + \sum_{i_1=i_2 \neq i_3=i_4} + \sum_{i_1=i_3 \neq i_2=i_4} + \sum_{i_1=i_4 \neq i_2=i_3} \right) \mathbb{E}[h_{i_1} R_{i_1 i_2} h_{i_2} h_{i_3} R_{i_3 i_4}^* h_{i_4}] \right\} - |\text{Tr}(\mathbf{D} \mathbf{R})|^2 \\ &= \sum_{i=1}^{M-1} |R_{ii}|^2 \mathbb{E}[h_i^4] + \sum_{i \neq j} D_{ii} R_{ii} D_{jj} R_{jj}^* + \sum_{i \neq j} D_{ii} D_{jj} (R_{ij} R_{ij}^* + R_{ij} R_{ji}^*) - |\text{Tr}(\mathbf{D} \mathbf{R})|^2 \\ &= \sum_{i=1}^{M-1} |R_{ii}|^2 (\mathbb{E}[h_i^4] - 3\mathbb{E}[h_i^2]^2) + \text{Tr}(\mathbf{D} \mathbf{R}^\top \mathbf{D} \mathbf{R}^*) + \text{Tr}(\mathbf{D} \mathbf{R} \mathbf{D} \mathbf{R}^*). \end{aligned}$$

Note that we have  $\max_i |\mathbb{E}[h_i^4] - 3\mathbb{E}[h_i^2]^2| = O_d(1/d^2)$ , so that

$$\sum_{i=1}^{M-1} |R_{ii}|^2 (\mathbb{E}[h_i^4] - 3\mathbb{E}[h_i^2]^2) \leq O_d(1/d^2) \cdot \|\mathbf{R}\|_F^2 \leq O_d(1/d) \|\mathbf{R}\|_{\text{op}}^2 = O_d(1/d).$$

Moreover, we have

$$|\text{Tr}(\mathbf{D} \mathbf{R}^\top \mathbf{D} \mathbf{R}^*) + \text{Tr}(\mathbf{D} \mathbf{R} \mathbf{D} \mathbf{R}^*)| \leq \|\mathbf{D} \mathbf{R}\|_F^2 + \|\mathbf{D} \mathbf{R}\|_F \|\mathbf{D} \mathbf{R}^*\|_F \leq 2\|\mathbf{D}\|_{\text{op}}^2 \|\mathbf{R}\|_F^2 = O_d(1/d),$$

which gives

$$\text{Var}(\mathbf{h}^\top \mathbf{R} \mathbf{h} | \mathbf{R}) = O_d(1/d),$$

and therefore

$$|\mathbf{h}^\top \mathbf{R} \mathbf{h} - \text{Tr}(\mathbf{D} \mathbf{R})| = O_{d,\mathbb{P}}(d^{-1/2}). \quad (191)$$

Combining Eq. (191) and (189), we obtain

$$|v_1 - \text{Tr}(\mathbf{D} \mathbf{R})| \leq |\mathbf{a}^\top \mathbf{R} \mathbf{a} - \mathbf{h}^\top \mathbf{R} \mathbf{h}| + |\mathbf{h}^\top \mathbf{R} \mathbf{h} - \text{Tr}(\mathbf{D} \mathbf{R})| = o_{d,\mathbb{P}}(1). \quad (192)$$

Finally, notice that

$$\text{Tr}(\mathbf{D} \mathbf{R}) = s_2^2 \text{Tr}_{[1, N-1]}((\tilde{\mathbf{B}} - \xi \mathbf{I}_{M-1})^{-1})/d + (\mu_1^2 + \mu_*^2) \text{Tr}_{[N, M-1]}((\tilde{\mathbf{B}} - \xi \mathbf{I}_{M-1})^{-1})/d.$$

By Lemmas C.6, C.7, and Lemma C.15 (which will be stated and proved later), we have

$$\begin{aligned} |\text{Tr}_{[1, N-1]}((\tilde{\mathbf{B}} - \xi \mathbf{I}_{M-1})^{-1})/d - \bar{m}_{1,d}| &= o_{d,\mathbb{P}}(1), \\ |\text{Tr}_{[N, M-1]}((\tilde{\mathbf{B}} - \xi \mathbf{I}_{M-1})^{-1})/d - \bar{m}_{2,d}| &= o_{d,\mathbb{P}}(1), \end{aligned}$$

so that

$$|\text{Tr}(\mathbf{D} \mathbf{R}) - \bar{v}_1| = o_{d,\mathbb{P}}(1).$$

Combining with Eq. (192) proves Eq. (186).  $\square$

The following lemma is the analog of Lemma B.7 and B.8 in [CS13].

**Lemma C.14.** Under the assumptions of Lemma C.4, we have (using the definitions in Eq. (180), (182) and (183))

$$\|(\mathbf{M}^{-1} + \mathbf{V}_3)^{-1}\|_{\text{op}} = O_{d,\mathbb{P}}(1), \quad (193)$$

$$\|(\mathbf{M}^{-1} + \bar{\mathbf{V}}_3)^{-1}\|_{\text{op}} = O_d(1). \quad (194)$$

*Proof of Lemma C.14.*

**Step 1. Bounding**  $\|(\mathbf{M}^{-1} + \mathbf{V}_3)^{-1}\|_{\text{op}}$ . By Sherman-Morrison-Woodbury formula, we have

$$\begin{aligned} (\mathbf{M}^{-1} + \mathbf{V}_3)^{-1} &= (\mathbf{M}^{-1} + \mathbf{U}^\top (\tilde{\mathbf{B}} - \xi \mathbf{I}_{M-1})^{-1} \mathbf{U})^{-1} \\ &= \mathbf{M} - \mathbf{M} \mathbf{U}^\top (\tilde{\mathbf{B}} - \xi \mathbf{I}_{M-1} + \mathbf{U} \mathbf{M} \mathbf{U}^\top)^{-1} \mathbf{U} \mathbf{M}. \end{aligned}$$

Note we have  $\|\mathbf{M}\|_{\text{op}} = O_d(1)$ , and  $\|(\tilde{\mathbf{B}} - \xi \mathbf{I}_{M-1} + \mathbf{U} \mathbf{M} \mathbf{U}^\top)^{-1}\|_{\text{op}} \leq 1/\xi_0 = O_d(1)$ . Therefore, by the concentration of  $\|\boldsymbol{\eta}\|_2/\sqrt{d}$  and  $\|\mathbf{u}\|_2/\sqrt{d}$ , we have

$$(\mathbf{M}^{-1} + \mathbf{V}_3)^{-1} = O_d(1) \cdot (1 + \|\mathbf{U}\|_{\text{op}}^2) = O_d(1)(1 + \|\boldsymbol{\eta}\|_2/\sqrt{d} + \|\mathbf{u}\|_2/\sqrt{d}) = O_{d,\mathbb{P}}(1).$$

**Step 2. Bounding**  $\|(\mathbf{M}^{-1} + \bar{\mathbf{V}}_3)^{-1}\|_{\text{op}}$ . Define  $\mathbf{G} = \mathbf{M}^{1/2} \mathbf{V}_3 \mathbf{M}^{1/2}$  and  $\bar{\mathbf{G}} = \mathbf{M}^{1/2} \bar{\mathbf{V}}_3 \mathbf{M}^{1/2}$ . By Lemma C.13, we have

$$\|\mathbf{G} - \bar{\mathbf{G}}\|_{\text{op}} = o_{d,\mathbb{P}}(1). \quad (195)$$

By the bound  $\|(\mathbf{M}^{-1} + \mathbf{V}_3)^{-1}\|_{\text{op}} = O_{d,\mathbb{P}}(1)$ , we get

$$\|(\mathbf{I}_2 + \mathbf{G})^{-1}\|_{\text{op}} = \|\mathbf{M}^{-1/2}(\mathbf{M}^{-1} + \mathbf{V}_3)^{-1}\mathbf{M}^{-1/2}\|_{\text{op}} \leq \|(\mathbf{M}^{-1} + \mathbf{V}_3)^{-1}\| \cdot \|\mathbf{M}^{-1/2}\|_{\text{op}}^2 = O_{d,\mathbb{P}}(1). \quad (196)$$

Note we have

$$(\mathbf{I}_2 + \bar{\mathbf{G}})^{-1} - (\mathbf{I}_2 + \mathbf{G})^{-1} = (\mathbf{I}_2 + \bar{\mathbf{G}})^{-1}(\mathbf{G} - \bar{\mathbf{G}})(\mathbf{I}_2 + \mathbf{G})^{-1},$$

so that

$$(\mathbf{I}_2 + \bar{\mathbf{G}})^{-1} = \{\mathbf{I}_2 - (\mathbf{G} - \bar{\mathbf{G}})(\mathbf{I}_2 + \mathbf{G})^{-1}\}(\mathbf{I}_2 + \mathbf{G})^{-1}.$$

Combining with Eq. (195) and (196), we get

$$\|(\mathbf{I}_2 + \bar{\mathbf{G}})^{-1}\|_{\text{op}} \leq \|\mathbf{I}_2 - (\mathbf{G} - \bar{\mathbf{G}})(\mathbf{I}_2 + \mathbf{G})^{-1}\|_{\text{op}} \|(\mathbf{I}_2 + \mathbf{G})^{-1}\|_{\text{op}} = O_{d,\mathbb{P}}(1) = O_d(1).$$

The last equality holds because  $\|(\mathbf{I}_2 + \bar{\mathbf{G}})^{-1}\|_{\text{op}}$  is deterministic. Hence we have

$$\|(\mathbf{M}^{-1} + \bar{\mathbf{V}}_3)^{-1}\|_{\text{op}} = \|\mathbf{M}^{1/2}(\mathbf{I}_2 + \bar{\mathbf{G}})^{-1}\mathbf{M}^{1/2}\|_{\text{op}} \leq \|(\mathbf{I}_2 + \bar{\mathbf{G}})^{-1}\|_{\text{op}} \|\mathbf{M}^{1/2}\|_{\text{op}}^2 = O_d(1).$$

This proves the lemma.  $\square$

**Lemma C.15.** Follow the assumptions of Lemma C.4. Let  $\bar{\mathbf{X}} = (\bar{\mathbf{x}}_1, \dots, \bar{\mathbf{x}}_n)^\top \in \mathbb{R}^{n \times d}$  with  $(\bar{\mathbf{x}}_i)_{i \in [n]} \sim_{iid} \mathbf{N}(\mathbf{0}, \mathbf{I}_d)$ , and  $\bar{\boldsymbol{\Theta}} = (\bar{\boldsymbol{\theta}}_1, \dots, \bar{\boldsymbol{\theta}}_N)^\top \in \mathbb{R}^{N \times d}$  with  $(\bar{\boldsymbol{\theta}}_a)_{a \in [N]} \sim_{iid} \mathbf{N}(\mathbf{0}, \mathbf{I}_d)$ . Let  $\tilde{\mathbf{X}} = (\tilde{\mathbf{x}}_1, \dots, \tilde{\mathbf{x}}_n)^\top \in \mathbb{R}^{n \times (d-1)}$  with  $(\tilde{\mathbf{x}}_i)_{i \in [n]} \sim_{iid} \mathbf{N}(\mathbf{0}, \mathbf{I}_{d-1})$ , and  $\tilde{\boldsymbol{\Theta}} = (\tilde{\boldsymbol{\theta}}_1, \dots, \tilde{\boldsymbol{\theta}}_N)^\top \in \mathbb{R}^{N \times (d-1)}$  with  $(\tilde{\boldsymbol{\theta}}_a)_{a \in [N]} \sim_{iid} \mathbf{N}(\mathbf{0}, \mathbf{I}_{d-1})$ . Denote

$$\bar{\mathbf{J}} \equiv \frac{1}{\sqrt{d}} \varphi\left(\frac{1}{\sqrt{d}} \bar{\mathbf{X}} \bar{\boldsymbol{\Theta}}^\top\right), \quad \tilde{\mathbf{J}} \equiv \frac{1}{\sqrt{d}} \varphi\left(\frac{1}{\sqrt{d}} \tilde{\mathbf{X}} \tilde{\boldsymbol{\Theta}}^\top\right), \quad (197)$$

$$\bar{\mathbf{Q}} \equiv \frac{1}{d} \bar{\boldsymbol{\Theta}} \bar{\boldsymbol{\Theta}}^\top, \quad \tilde{\mathbf{Q}} \equiv \frac{1}{d} \tilde{\boldsymbol{\Theta}} \tilde{\boldsymbol{\Theta}}^\top, \quad (198)$$

$$\bar{\mathbf{H}} \equiv \frac{1}{d} \bar{\mathbf{X}} \bar{\mathbf{X}}^\top, \quad \tilde{\mathbf{H}} \equiv \frac{1}{d} \tilde{\mathbf{X}} \tilde{\mathbf{X}}^\top, \quad (199)$$

as well as the block matrix  $\bar{\mathbf{A}}, \tilde{\mathbf{A}} \in \mathbb{R}^{M \times M}$ ,  $M = N + n$ , defined by

$$\bar{\mathbf{A}} = \begin{bmatrix} s_1 \mathbf{I}_N + s_2 \bar{\mathbf{Q}} & \bar{\mathbf{J}}^\top \\ \bar{\mathbf{J}} & t_1 \mathbf{I}_n + t_2 \bar{\mathbf{H}} \end{bmatrix} \quad \tilde{\mathbf{A}} = \begin{bmatrix} s_1 \mathbf{I}_N + s_2 \tilde{\mathbf{Q}} & \tilde{\mathbf{J}}^\top \\ \tilde{\mathbf{J}} & t_1 \mathbf{I}_n + t_2 \tilde{\mathbf{H}} \end{bmatrix}. \quad (200)$$

Then for any  $\xi \in \mathbb{C}_+$  with  $\Im \xi \geq \xi_0 > 0$ , we have

$$\frac{1}{d} \mathbb{E} \left| \text{Tr}_{[1,N]}[(\bar{\mathbf{A}} - \xi \mathbf{I}_M)^{-1}] - \text{Tr}_{[1,N]}[(\tilde{\mathbf{A}} - \xi \mathbf{I}_M)^{-1}] \right| = o_d(1), \quad (201)$$

$$\frac{1}{d} \mathbb{E} \left| \text{Tr}_{[N+1,M]}[(\bar{\mathbf{A}} - \xi \mathbf{I}_M)^{-1}] - \text{Tr}_{[N+1,M]}[(\tilde{\mathbf{A}} - \xi \mathbf{I}_M)^{-1}] \right| = o_d(1). \quad (202)$$

*Proof of Lemma C.15.*

**Step 1. The Schur complement.**

We denote  $\bar{\mathbf{A}}_{ij}$  and  $\tilde{\mathbf{A}}_{ij}$  for  $i, j \in [2]$  to be

$$\bar{\mathbf{A}} = \begin{bmatrix} \bar{\mathbf{A}}_{11} & \bar{\mathbf{A}}_{12} \\ \bar{\mathbf{A}}_{21} & \bar{\mathbf{A}}_{22} \end{bmatrix} = \begin{bmatrix} s_1 \mathbf{I}_N + s_2 \bar{\mathbf{Q}} & \bar{\mathbf{J}}^\top \\ \bar{\mathbf{J}} & t_1 \mathbf{I}_n + t_2 \bar{\mathbf{H}} \end{bmatrix}, \quad \tilde{\mathbf{A}} = \begin{bmatrix} \tilde{\mathbf{A}}_{11} & \tilde{\mathbf{A}}_{12} \\ \tilde{\mathbf{A}}_{21} & \tilde{\mathbf{A}}_{22} \end{bmatrix} = \begin{bmatrix} s_1 \mathbf{I}_N + s_2 \tilde{\mathbf{Q}} & \tilde{\mathbf{J}}^\top \\ \tilde{\mathbf{J}} & t_1 \mathbf{I}_n + t_2 \tilde{\mathbf{H}} \end{bmatrix}.$$

Define

$$\bar{\omega} = \frac{1}{d} \text{Tr}_{[1,N]}[(\bar{\mathbf{A}} - \xi \mathbf{I}_M)^{-1}], \quad \tilde{\omega} = \frac{1}{d} \text{Tr}_{[1,N]}[(\tilde{\mathbf{A}} - \xi \mathbf{I}_M)^{-1}],$$

and

$$\begin{aligned} \bar{\Omega} &= \left( \bar{\mathbf{A}}_{11} - \xi \mathbf{I}_N - \bar{\mathbf{A}}_{12} (\bar{\mathbf{A}}_{22} - \xi \mathbf{I}_n)^{-1} \bar{\mathbf{A}}_{21} \right)^{-1}, \\ \tilde{\Omega} &= \left( \tilde{\mathbf{A}}_{11} - \xi \mathbf{I}_N - \tilde{\mathbf{A}}_{12} (\tilde{\mathbf{A}}_{22} - \xi \mathbf{I}_n)^{-1} \tilde{\mathbf{A}}_{21} \right)^{-1}. \end{aligned}$$

Then we have

$$\bar{\omega} = \frac{1}{d} \text{Tr}(\bar{\Omega}), \quad \tilde{\omega} = \frac{1}{d} \text{Tr}(\tilde{\Omega}).$$

Define

$$\begin{aligned} \Omega_1 &= \left( \tilde{\mathbf{A}}_{11} - \xi \mathbf{I}_N - \bar{\mathbf{A}}_{12} (\bar{\mathbf{A}}_{22} - \xi \mathbf{I}_n)^{-1} \bar{\mathbf{A}}_{21} \right)^{-1}, \\ \Omega_2 &= \left( \tilde{\mathbf{A}}_{11} - \xi \mathbf{I}_N - \tilde{\mathbf{A}}_{12} (\bar{\mathbf{A}}_{22} - \xi \mathbf{I}_n)^{-1} \bar{\mathbf{A}}_{21} \right)^{-1}, \\ \Omega_3 &= \left( \tilde{\mathbf{A}}_{11} - \xi \mathbf{I}_N - \tilde{\mathbf{A}}_{12} (\bar{\mathbf{A}}_{22} - \xi \mathbf{I}_n)^{-1} \tilde{\mathbf{A}}_{21} \right)^{-1}, \end{aligned}$$

Then it's easy to see that  $\|\bar{\Omega}\|_{\text{op}}, \|\Omega_1\|_{\text{op}}, \|\Omega_2\|_{\text{op}}, \|\Omega_3\|_{\text{op}}, \|\tilde{\Omega}\|_{\text{op}} \leq 1/\xi_0$ .

Calculating their difference, we have

$$\begin{aligned} \left| \frac{1}{d} \text{Tr}(\bar{\Omega}) - \frac{1}{d} \text{Tr}(\Omega_1) \right| &= \left| \frac{1}{d} \text{Tr}(\bar{\Omega}(\tilde{\mathbf{A}}_{11} - \bar{\mathbf{A}}_{11})\tilde{\Omega}) \right| \leq O_d(1) \cdot \frac{1}{d} \|\tilde{\mathbf{A}}_{11} - \bar{\mathbf{A}}_{11}\|_*, \\ \left| \frac{1}{d} \text{Tr}(\Omega_1) - \frac{1}{d} \text{Tr}(\Omega_2) \right| &\leq O_d(1) \cdot \frac{1}{d} \|(\tilde{\mathbf{A}}_{12} - \bar{\mathbf{A}}_{12})(\bar{\mathbf{A}}_{22} - \xi \mathbf{I}_n)^{-1} \bar{\mathbf{A}}_{21}\|_*, \\ \left| \frac{1}{d} \text{Tr}(\Omega_2) - \frac{1}{d} \text{Tr}(\Omega_3) \right| &\leq O_d(1) \cdot \frac{1}{d} \|(\tilde{\mathbf{A}}_{12} - \bar{\mathbf{A}}_{12})(\bar{\mathbf{A}}_{22} - \xi \mathbf{I}_n)^{-1} \tilde{\mathbf{A}}_{21}\|_*, \\ \left| \frac{1}{d} \text{Tr}(\Omega_3) - \frac{1}{d} \text{Tr}(\tilde{\Omega}) \right| &\leq O_d(1) \cdot \frac{1}{d} \|\tilde{\mathbf{A}}_{12}(\bar{\mathbf{A}}_{22} - \xi \mathbf{I}_n)^{-1}(\bar{\mathbf{A}}_{22} - \tilde{\mathbf{A}}_{22})(\tilde{\mathbf{A}}_{22} - \xi \mathbf{I}_n)^{-1} \tilde{\mathbf{A}}_{21}\|_*. \end{aligned}$$

**Step 2. Bounding the differences.**

First, we have

$$\bar{\mathbf{A}}_{11} - \tilde{\mathbf{A}}_{11} = s_2(\bar{\mathbf{Q}} - \tilde{\mathbf{Q}}) = s_2(\bar{\theta}_{id}\bar{\theta}_{jd}/d)_{i,j \in [N]} = s_2 \boldsymbol{\eta} \boldsymbol{\eta}^\top / d,$$

where  $\boldsymbol{\eta} = (\bar{\theta}_{1d}, \dots, \bar{\theta}_{Nd})^\top \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_N)$ . This gives

$$\|\bar{\mathbf{A}}_{11} - \tilde{\mathbf{A}}_{11}\|_*/d = s_2 \|\boldsymbol{\eta}\|_2^2/d^2 = o_{d,\mathbb{P}}(1),$$

and therefore

$$\left| \frac{1}{d} \text{Tr}(\bar{\Omega}) - \frac{1}{d} \text{Tr}(\Omega_1) \right| = o_{d,\mathbb{P}}(1).$$

By Theorem 1.7 in [FM19], and by the fact that  $\varphi$  is a polynomial with  $\mathbb{E}[\varphi(G)] = 0$ , we have

$$\|\bar{\mathbf{A}}_{12}\|_{\text{op}} = \|\bar{\mathbf{J}}\|_{\text{op}} = O_{d,\mathbb{P}}(1), \quad \|\tilde{\mathbf{A}}_{12}\|_{\text{op}} = O_{d,\mathbb{P}}(1).$$

It is also easy to see that

$$\|(\bar{\mathbf{A}}_{22} - \xi \mathbf{I}_n)^{-1}\|_{\text{op}}, \|(\tilde{\mathbf{A}}_{22} - \xi \mathbf{I}_n)^{-1}\|_{\text{op}} \leq 1/\xi_0 = O_d(1).$$

Moreover, we have

$$\bar{\mathbf{A}}_{22} - \tilde{\mathbf{A}}_{22} = t_2(\bar{\mathbf{H}} - \tilde{\mathbf{H}}) = t_2(\bar{x}_{id}\bar{x}_{jd}/d)_{i,j \in [n]} = t_2 \mathbf{u} \mathbf{u}^\top / d,$$

where  $\mathbf{u} = (\bar{x}_{1d}, \dots, \bar{x}_{nd})^\top \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_n)$ . This gives

$$\begin{aligned} & \|\tilde{\mathbf{A}}_{12}(\bar{\mathbf{A}}_{22} - \xi \mathbf{I}_n)^{-1}(\bar{\mathbf{A}}_{22} - \tilde{\mathbf{A}}_{22})(\tilde{\mathbf{A}}_{22} - \xi \mathbf{I}_n)^{-1}\tilde{\mathbf{A}}_{21}\|_*/d \\ & \leq t_2 \|\tilde{\mathbf{A}}_{12}(\bar{\mathbf{A}}_{22} - \xi \mathbf{I}_n)^{-1}\mathbf{u}\|_2 \|\tilde{\mathbf{A}}_{12}(\tilde{\mathbf{A}}_{22} - \xi \mathbf{I}_n)^{-1}\mathbf{u}\|_2 / d^2 \\ & \leq t_2 \|\tilde{\mathbf{A}}_{12}\|_{\text{op}}^2 \|(\bar{\mathbf{A}}_{22} - \xi \mathbf{I}_n)^{-1}\|_{\text{op}}^2 \|\mathbf{u}\|_2^2 / d^2 \\ & = O_{d,\mathbb{P}}(1) \cdot \|\mathbf{u}\|_2^2 / d^2 = o_{d,\mathbb{P}}(1), \end{aligned}$$

and therefore

$$\left| \frac{1}{d} \text{Tr}(\mathbf{\Omega}_3) - \frac{1}{d} \text{Tr}(\tilde{\mathbf{\Omega}}) \right| = o_{d,\mathbb{P}}(1).$$

By Lemma C.10, defining

$$\mathbf{E} = \bar{\mathbf{A}}_{12} - \tilde{\mathbf{A}}_{12} - \mu_1 \mathbf{u} \boldsymbol{\eta}^\top / d,$$

we have  $\|\mathbf{E}\|_{\text{op}} = O_d(\text{Poly}(\log d)/\sqrt{d})$ . Therefore, we get

$$\begin{aligned} & \|(\tilde{\mathbf{A}}_{12} - \bar{\mathbf{A}}_{12})(\bar{\mathbf{A}}_{22} - \xi \mathbf{I}_n)^{-1}\bar{\mathbf{A}}_{21}\|_*/d \\ & \leq \|(\mu_1 \mathbf{u} \boldsymbol{\eta}^\top / d)(\bar{\mathbf{A}}_{22} - \xi \mathbf{I}_n)^{-1}\bar{\mathbf{A}}_{21}\|_*/d + \|\mathbf{E}(\bar{\mathbf{A}}_{22} - \xi \mathbf{I}_n)^{-1}\bar{\mathbf{A}}_{21}\|_*/d \\ & \leq \mu_1 \|\boldsymbol{\eta}^\top (\bar{\mathbf{A}}_{22} - \xi \mathbf{I}_n)^{-1}\bar{\mathbf{A}}_{21}\|_2 \|\mathbf{u}\|_2 / d^2 + \|\mathbf{E}(\bar{\mathbf{A}}_{22} - \xi \mathbf{I}_n)^{-1}\bar{\mathbf{A}}_{21}\|_{\text{op}} \\ & \leq \mu_1 \|\boldsymbol{\eta}\|_2 \|(\bar{\mathbf{A}}_{22} - \xi \mathbf{I}_n)^{-1}\|_{\text{op}} \|\bar{\mathbf{A}}_{21}\|_{\text{op}} \|\mathbf{u}\|_2 / d^2 + \|\mathbf{E}\|_{\text{op}} \|(\bar{\mathbf{A}}_{22} - \xi \mathbf{I}_n)^{-1}\|_{\text{op}} \|\bar{\mathbf{A}}_{21}\|_{\text{op}} \\ & = o_{d,\mathbb{P}}(1), \end{aligned}$$

and therefore

$$\left| \frac{1}{d} \text{Tr}(\mathbf{\Omega}_1) - \frac{1}{d} \text{Tr}(\mathbf{\Omega}_2) \right|, \left| \frac{1}{d} \text{Tr}(\mathbf{\Omega}_2) - \frac{1}{d} \text{Tr}(\mathbf{\Omega}_3) \right| = o_{d,\mathbb{P}}(1).$$

Combining all these bounds establishes Eq. (201). Finally, Eq. (202) can be shown using the same argument.  $\square$

## C.6 Proof of Proposition 7.2

### Step 1. Polynomial activation function $\sigma$ .

First we consider the case when  $\sigma$  is a fixed polynomial with  $\mathbb{E}_{G \sim \mathcal{N}(0,1)}[\sigma(G)G] \neq 0$ . Let  $\varphi(u) = \sigma(u) - \mathbb{E}[\sigma(G)]$ , and let  $\bar{\mathbf{m}}_d \equiv (\bar{m}_{1,d}, \bar{m}_{2,d})$  (whose definition is given by Eq. (147) and (148)), and recall that  $\psi_{1,d} \rightarrow \psi_1$  and  $\psi_{2,d} \rightarrow \psi_2$  as  $d \rightarrow \infty$ . By Lemma C.4, together with the continuity of  $\mathbf{F}$  with respect to  $\psi_1, \psi_2$ , cf. Lemma C.3, we have, for any  $\xi_0 > 0$ , there exists  $C = C(\xi_0, \mathbf{q}, \psi_1, \psi_2, \varphi)$  and  $\text{err}(d) \rightarrow 0$  such that for all  $\xi \in \mathbb{C}_+$  with  $\Im \xi \geq \xi_0$ ,

$$\|\bar{\mathbf{m}}_d - \mathbf{F}(\bar{\mathbf{m}}_d; \xi)\|_2 \leq C \cdot \text{err}(d), \quad (203)$$

Let  $\mathbf{m}$  be the solution of the fixed point equation (61) defined in the statement of the Proposition 7.2. This solution is well defined by Lemma C.3. By the Lipschitz property of  $\mathbf{F}$  in Lemma C.3, for Lipschitz constant  $\delta = 1/2$  (there exists a larger  $\xi_0$  depending on  $\delta$ ), such that for  $\Im \xi \geq \xi_0$  for some large  $\xi_0$

$$\|\bar{\mathbf{m}}_d - \mathbf{m}\|_2 \leq \|\mathbf{F}(\bar{\mathbf{m}}_d; \xi) - \mathbf{F}(\mathbf{m}; \xi)\|_2 + C \cdot \text{err}(d) \leq \frac{1}{2} \|\bar{\mathbf{m}}_d - \mathbf{m}\|_2 + C \cdot \text{err}(d), \quad (204)$$

which yields for some large  $\xi_0$

$$\sup_{\Im \xi \geq \xi_0} \|\bar{\mathbf{m}}_d(\xi) - \mathbf{m}(\xi)\|_2 = o_d(1).$$

By the property of Stieltjes transform as in Lemma C.5 (c), we have

$$\|\bar{\mathbf{m}}_d(\xi) - \mathbf{m}(\xi)\|_2 = o_d(1), \quad \forall \xi \in \mathbb{C}_+.$$

By the concentration result of Lemma C.7, for  $\bar{M}_d(\xi) = d^{-1}\text{Tr}[(\bar{\mathbf{A}} - \xi \mathbf{I}_M)^{-1}]$ , we also have

$$\mathbb{E}|\bar{M}_d(\xi) - m(\xi)| = o_d(1), \quad \forall \xi \in \mathbb{C}_+. \quad (205)$$

Then we use Lemma C.1 to transfer this property from  $\bar{M}_d$  to  $M_d$ . Recall the definition of resolvent  $M_d(\xi; \mathbf{q})$  in sphere case in Eq. (59). Combining Lemma C.1 with Eq. (205), we have

$$\mathbb{E}|M_d(\xi) - m(\xi)| = o_d(1), \quad \forall \xi \in \mathbb{C}_+. \quad (206)$$

## Step 2. General activation function $\sigma$ satisfying Assumption 1.

Next consider the case of a general function  $\sigma$  as in the theorem statement satisfying Assumption 1. Fix  $\varepsilon > 0$  and let  $\tilde{\sigma}$  is a polynomial be such that  $\|\sigma - \tilde{\sigma}\|_{L^2(\tau_d)} \leq \varepsilon$ , where  $\tau_d$  is the marginal distribution of  $\langle \mathbf{x}, \boldsymbol{\theta} \rangle / \sqrt{d}$  for  $\mathbf{x}, \boldsymbol{\theta} \sim_{iid} \text{Unif}(\mathbb{S}^{d-1}(\sqrt{d}))$ . In order to construct such a polynomial, consider the expansion of  $\sigma$  in the orthogonal basis of Hermite polynomials

$$\sigma(x) = \sum_{k=0}^{\infty} \frac{\mu_k}{k!} \text{He}_k(x). \quad (207)$$

Since this series converges in  $L^2(\mu_G)$ , we can choose  $\bar{k} < \infty$  such that, letting  $\tilde{\sigma}(x) = \sum_{k=0}^{\bar{k}} (\mu_k/k!) \text{He}_k(x)$ , we have  $\|\sigma - \tilde{\sigma}\|_{L^2(\mu_G)}^2 \leq \varepsilon/2$ . By Lemma C.8 (cf. Eq. (155)) we therefore have  $\|\sigma - \tilde{\sigma}\|_{L^2(\tau_d)}^2 \leq \varepsilon$  for all  $d$  large enough.

Write  $\mu_k(\tilde{\sigma}) = \mathbb{E}[\tilde{\sigma}(G)\text{He}_k(G)]$  and  $\mu_*(\tilde{\sigma})^2 = \sum_{k=2}^{\bar{k}} \mu_k^2/k!$ . Notice that, by construction we have  $\mu_0(\tilde{\sigma}) = \mu_0(\sigma)$ ,  $\mu_1(\tilde{\sigma}) = \mu_1(\sigma)$  and  $|\mu_*(\tilde{\sigma})^2 - \mu_*(\sigma)^2| \leq \varepsilon$ . Let  $\tilde{m}_{1,d}, \tilde{m}_{2,d}$  be the Stieltjes transforms associated to activation  $\tilde{\sigma}$ , and  $\tilde{m}_1, \tilde{m}_2$  be the solution of the corresponding fixed point equation (61) (with  $\mu_* = \mu_*(\tilde{\sigma})$  and  $\mu_1 = \mu_1(\tilde{\sigma})$ ), and  $\tilde{m} = \tilde{m}_1 + \tilde{m}_2$ . Denoting by  $\tilde{\mathbf{A}}$  the matrix obtained by replacing the  $\sigma$  in  $\mathbf{A}$  to be  $\tilde{\sigma}$ , and  $\tilde{M}_d(\xi) = (1/d)\text{Tr}[(\tilde{\mathbf{A}} - \xi \mathbf{I})^{-1}]$ . Step 1 of this proof implies

$$\mathbb{E}|\tilde{M}_d(\xi) - \tilde{m}(\xi)| = o_d(1), \quad \forall \xi \in \mathbb{C}_+. \quad (208)$$

Further, by continuity of the solution of the fixed point equation with respect to  $\mu_*, \mu_1$  when  $\Im \xi \geq \xi_0$  for some large  $\xi_0$  (as stated in Lemma C.3), we have for  $\Im \xi \geq \xi_0$ ,

$$|\tilde{m}(\xi) - m(\xi)| \leq C(\xi, \mathbf{q})\varepsilon. \quad (209)$$

Eq. (209) also holds for any  $\xi \in \mathbb{C}_+$ , by the property of Stieltjes transform as in Lemma C.5 (c).

Moreover, we have (for  $C$  independent of  $d, \sigma, \tilde{\sigma}$  and  $\varepsilon$ , but depend on  $\xi$  and  $\mathbf{q}$ )

$$\begin{aligned} \mathbb{E}\left[\left|M_d(\xi) - \tilde{M}_d(\xi)\right|\right] &\leq \frac{1}{d}\mathbb{E}\left[\left|\text{Tr}[(\mathbf{A} - \xi \mathbf{I})^{-1}(\tilde{\mathbf{A}} - \mathbf{A})(\tilde{\mathbf{A}} - \xi \mathbf{I})^{-1}]\right|\right] \\ &\leq \frac{1}{d}\mathbb{E}\left[\|(\tilde{\mathbf{A}} - \xi \mathbf{I})^{-1}(\mathbf{A} - \xi \mathbf{I})^{-1}\|_{\text{op}}\|\tilde{\mathbf{A}} - \mathbf{A}\|_*\right] \\ &\leq [1/(\xi_0^2 d)] \cdot \mathbb{E}[\|\tilde{\mathbf{A}} - \mathbf{A}\|_*] \leq [1/(\xi_0^2 \sqrt{d})] \cdot \mathbb{E}\{\|\tilde{\mathbf{A}} - \mathbf{A}\|_F^2\}^{1/2} \leq C(\xi, \mathbf{q}) \cdot \|\sigma - \tilde{\sigma}\|_{L^2(\tau_d)}. \end{aligned}$$

Therefore

$$\limsup_{d \rightarrow \infty} \mathbb{E}[|M_d(\xi) - \tilde{M}_d(\xi)|] \leq C(\xi, \mathbf{q})\varepsilon, \quad \forall \xi \in \mathbb{C}_+. \quad (210)$$

Combining Eq. (208), (209), and (210), we obtain

$$\limsup_{d \rightarrow \infty} \mathbb{E}|M_d(\xi) - m(\xi)| \leq C(\xi, \mathbf{q})\varepsilon, \quad \forall \xi \in \mathbb{C}_+.$$

Taking  $\varepsilon \rightarrow 0$  proves Eq. (62).

**Step 3. Uniform convergence in compact sets (Eq. (63)).**

Note  $m_d(\xi; \mathbf{q}) = \mathbb{E}[M_d(\xi; \mathbf{q})]$  is an analytic function on  $\mathbb{C}_+$ . By Lemma C.5 (c), for any compact set  $\Omega \subseteq \mathbb{C}_+$ , we have

$$\lim_{d \rightarrow \infty} \left[ \sup_{\xi \in \Omega} |\mathbb{E}[M_d(\xi; \mathbf{q})] - m(\xi; \mathbf{q})| \right] = 0. \quad (211)$$

In the following, we show the concentration of  $M_d(\xi; \mathbf{q})$  around its expectation uniformly in the compact set  $\Omega \subset \mathbb{C}_+$ . Define  $L = \sup_{\xi \in \Omega} (1/\Im \xi^2)$ . Since  $\Omega \subset \mathbb{C}_+$  is a compact set, we have  $L < \infty$ , and  $M_d(\xi; \mathbf{q})$  (as a function of  $\xi$ ) is  $L$ -Lipschitz on  $\Omega$ . Moreover, for any  $\varepsilon > 0$ , there exists a finite set  $\mathcal{N}(\varepsilon, \Omega) \subseteq \mathbb{C}_+$  which is an  $\varepsilon/L$ -covering of  $\Omega$ . That is, for any  $\xi \in \Omega$ , there exists  $\xi_* \in \mathcal{N}(\varepsilon, \Omega)$  such that  $|\xi - \xi_*| \leq \varepsilon/L$ . Since  $M_d(\xi; \mathbf{q})$  (as a function of  $\xi$ ) is  $L$ -Lipschitz on  $\Omega$ , we have

$$\sup_{\xi \in \Omega} \inf_{\xi_* \in \mathcal{N}(\varepsilon, \Omega)} |M_d(\xi; \mathbf{q}) - M_d(\xi_*; \mathbf{q})| \leq \varepsilon. \quad (212)$$

By the concentration of  $M_d(\xi_*; \mathbf{q})$  to its expectation as per Lemma C.7, we have

$$|M_d(\xi_*; \mathbf{q}) - \mathbb{E}[M_d(\xi_*; \mathbf{q})]| = o_{d, \mathbb{P}}(1),$$

and since  $\mathcal{N}(\varepsilon, \Omega)$  is a finite set, we have

$$\sup_{\xi_* \in \mathcal{N}(\varepsilon, \Omega)} |M_d(\xi_*; \mathbf{q}) - \mathbb{E}[M_d(\xi_*; \mathbf{q})]| = o_{d, \mathbb{P}}(1). \quad (213)$$

Combining Eq. (212) and (213), we have

$$\limsup_{d \rightarrow \infty} \sup_{\xi \in \Omega} |M_d(\xi; \mathbf{q}) - M_d(\xi; \mathbf{q})| \leq \varepsilon.$$

Letting  $\varepsilon \rightarrow 0$  proves Eq. (63).

## D Proof of Proposition 7.3

We can see Eq. (64) is trivial. In the following, we prove Eq. (65).

For any fixed  $\mathbf{q} \in \mathbb{R}^5$ ,  $\xi \in \mathbb{C}_+$  and a fixed instance  $\mathbf{A}(\mathbf{q})$ , the determinant can be represented as  $\det(\mathbf{A}(\mathbf{q}) - \xi \mathbf{I}_M) = r(\mathbf{q}, \xi) \exp(i\theta(\mathbf{q}, \xi))$  for  $\theta(\mathbf{q}, \xi) \in (-\pi, \pi]$ . Without loss of generality, we assume for this fixed  $\mathbf{q}$  and  $\xi$ , we have  $\theta(\mathbf{q}, \xi) \neq \pi$ , and then  $\text{Log}(\det(\mathbf{A}(\mathbf{q}) - \xi \mathbf{I}_M)) = \log r(\mathbf{q}, \xi) + i\theta(\mathbf{q}, \xi)$  (when  $\theta(\mathbf{q}, \xi) = \pi$ , we use another definition of Log notation, and the proof is the same). For this  $\mathbf{q}$ ,  $\xi$ , and  $\mathbf{A}(\mathbf{q})$ , there exists some integer  $k = k(\mathbf{q}, \xi) \in \mathbb{N}$ , such that

$$\sum_{i=1}^M \text{Log}(\lambda_i(\mathbf{A}(\mathbf{q})) - \xi) = \text{Log} \det(\mathbf{A}(\mathbf{q}) - \xi \mathbf{I}_M) + 2\pi i k(\mathbf{q}, \xi).$$

Moreover, the set of eigenvalues of  $\mathbf{A}(\mathbf{q}) - \xi \mathbf{I}_M$  and  $\det(\mathbf{A}(\mathbf{q}) - \xi \mathbf{I}_M)$  are continuous with respect to  $\mathbf{q}$ . Therefore, for any perturbation  $\Delta \mathbf{q}$  with  $\|\Delta \mathbf{q}\|_2 \leq \varepsilon$  and  $\varepsilon$  small enough, we have  $k(\mathbf{q} + \Delta \mathbf{q}, \xi) = k(\mathbf{q}, \xi)$ . As a result, we have

$$\partial_{q_i} \left[ \sum_{i=1}^M \text{Log}(\lambda_i(\mathbf{A}(\mathbf{q})) - \xi) \right] = \partial_{q_i} \text{Log} \left[ \det(\mathbf{A}(\mathbf{q}) - \xi \mathbf{I}_M) \right] = \text{Tr} \left[ (\mathbf{A}(\mathbf{q}) - \xi \mathbf{I}_M)^{-1} \partial_{q_i} \mathbf{A}(\mathbf{q}) \right].$$

Moreover,  $\mathbf{A}(\mathbf{q})$  (defined as in Eq. (57)) is a linear matrix function of  $\mathbf{q}$ , which gives  $\partial_{q_i, q_j} \mathbf{A}(\mathbf{q}) = \mathbf{0}$ . Hence we have

$$\begin{aligned} & \partial_{q_i, q_j}^2 \left[ \sum_{i=1}^M \text{Log}(\lambda_i(\mathbf{A}(\mathbf{q})) - \xi) \right] \\ &= \partial_{q_i, q_j}^2 \text{Log} \left[ \det(\mathbf{A}(\mathbf{q}) - \xi \mathbf{I}_M) \right] \\ &= \partial_{q_j} \text{Tr} \left[ (\mathbf{A}(\mathbf{q}) - \xi \mathbf{I}_M)^{-1} \partial_{q_i} \mathbf{A}(\mathbf{q}) \right] \\ &= \text{Tr} \left[ (\mathbf{A}(\mathbf{q}) - \xi \mathbf{I}_M)^{-1} \partial_{q_j} \mathbf{A}(\mathbf{q}) (\mathbf{A}(\mathbf{q}) - \xi \mathbf{I}_M)^{-1} \partial_{q_i} \mathbf{A}(\mathbf{q}) \right]. \end{aligned}$$

Note

$$\begin{aligned}\partial_{s_1}\mathbf{A}(\mathbf{0}) &= \begin{bmatrix} \mathbf{I}_N & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix}, & \partial_{s_2}\mathbf{A}(\mathbf{0}) &= \begin{bmatrix} \mathbf{Q} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix}, \\ \partial_{t_1}\mathbf{A}(\mathbf{0}) &= \begin{bmatrix} \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{I}_n \end{bmatrix}, & \partial_{t_2}\mathbf{A}(\mathbf{0}) &= \begin{bmatrix} \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{H} \end{bmatrix}, & \partial_p\mathbf{A}(\mathbf{0}) &= \begin{bmatrix} \mathbf{0} & \mathbf{Z}_1^\top \\ \mathbf{Z}_1 & \mathbf{0} \end{bmatrix},\end{aligned}$$

and using the formula for block matrix inversion, we have

$$(\mathbf{A}(\mathbf{0}) - iu\mathbf{I}_M)^{-1} = \begin{bmatrix} (-iu\mathbf{I}_N - i\mathbf{Z}^\top\mathbf{Z}/u)^{-1} & (u^2\mathbf{I}_N + \mathbf{Z}^\top\mathbf{Z})^{-1}\mathbf{Z}^\top \\ \mathbf{Z}(u^2\mathbf{I}_N + \mathbf{Z}^\top\mathbf{Z})^{-1} & (-iu\mathbf{I}_n - i\mathbf{Z}\mathbf{Z}^\top/u)^{-1} \end{bmatrix}.$$

With simple algebra, we can show the proposition holds.

## E Proof of Proposition 7.4

### E.1 Properties of the Stieltjes transforms and the log determinant

**Lemma E.1.** For  $\xi \in \mathbb{C}_+$  and  $\mathbf{q} \in \mathcal{Q}$  (c.f. Eq. (58)), let  $m_1(\xi; \mathbf{q}), m_2(\xi; \mathbf{q})$  be defined as the analytic continuation of solution of Eq. (61) as defined in Proposition 7.2. Denote  $\xi = \xi_r + iK$  for some fixed  $\xi_r \in \mathbb{R}$ . Then we have

$$\lim_{K \rightarrow \infty} |m_1(\xi; \mathbf{q})\xi + \psi_1| = 0, \quad \lim_{K \rightarrow \infty} |m_2(\xi; \mathbf{q})\xi + \psi_2| = 0.$$

*Proof of Lemma E.1.* Define  $\bar{m}_1 = -\psi_1/\xi$ , and  $\bar{m}_2 = -\psi_2/\xi$ ,  $\bar{\mathbf{m}} = (\bar{m}_1, \bar{m}_2)^\top$ , and  $\mathbf{m} = (m_1, m_2)^\top$ . Let  $\mathbf{F}$  be defined as in Eq. (60),  $\mathbf{F}$  be defined as in Eq. (136), and  $\mathbf{H}$  defined as in Eq. (139). By simple calculus we can see that

$$\lim_{K \rightarrow \infty} \mathbf{H}(\bar{\mathbf{m}}) = 1,$$

so that

$$\xi[\bar{\mathbf{m}} - \mathbf{F}(\bar{\mathbf{m}}; \xi)] = [\psi_1, \psi_2]^\top \cdot \frac{s_1 + \mathbf{H}(\bar{\mathbf{m}})}{\xi - s_1 - \mathbf{H}(\bar{\mathbf{m}})} \rightarrow 0, \quad \text{as } K \rightarrow \infty.$$

Moreover, by Lemma C.3, for any  $r > 0$ , there exists sufficiently large  $\xi_0$ , so that for any  $\Im \xi = K \geq \xi_0$ ,  $\mathbf{F}(\mathbf{m}; \xi)$  is  $1/2$ -Lipschitz on domain  $\mathbf{m} \in \mathbb{D}(r) \times \mathbb{D}(r)$ . Therefore, for  $\Im \xi = K \geq \xi_0$ , we have (note we have  $\mathbf{m} = \mathbf{F}(\mathbf{m}; \xi)$ )

$$\begin{aligned}\|\bar{\mathbf{m}} - \mathbf{m}\|_2 &= \|\mathbf{F}(\bar{\mathbf{m}}; \xi) - \mathbf{F}(\mathbf{m}; \xi) + \bar{\mathbf{m}} - \mathbf{F}(\bar{\mathbf{m}}; \xi)\|_2 \\ &\leq \|\mathbf{F}(\bar{\mathbf{m}}; \xi) - \mathbf{F}(\mathbf{m}; \xi)\|_2 + \|\bar{\mathbf{m}} - \mathbf{F}(\bar{\mathbf{m}}; \xi)\|_2 \\ &\leq \|\bar{\mathbf{m}} - \mathbf{m}\|_2/2 + \|\bar{\mathbf{m}} - \mathbf{F}(\bar{\mathbf{m}}; \xi)\|_2,\end{aligned}$$

so that

$$\xi\|\bar{\mathbf{m}} - \mathbf{m}\|_2 \leq 2\xi\|\bar{\mathbf{m}} - \mathbf{F}(\bar{\mathbf{m}}; \xi)\|_2 \rightarrow 0, \quad \text{as } K \rightarrow \infty.$$

This proves the lemma. □

**Lemma E.2.** Follow the notations and settings of Proposition 7.4. For any fixed  $\mathbf{q}$ , we have

$$\lim_{K \rightarrow \infty} \sup_{d \geq 1} \mathbb{E}|G_d(iK; \mathbf{q}) - (\psi_1 + \psi_2)\text{Log}(-iK)| = 0, \quad (214)$$

$$\lim_{K \rightarrow \infty} |g(iK; \mathbf{q}) - (\psi_1 + \psi_2)\text{Log}(-iK)| = 0. \quad (215)$$

*Proof of Lemma E.2.*

**Step 1. Asymptotics of  $G_d(\mathbf{i}K; \mathbf{q})$ .** First we look at the real part. We have

$$\begin{aligned}
& \left| \Re \left[ \frac{1}{M} \sum_{i=1}^M \text{Log}(\lambda_i(\mathbf{A}) - \mathbf{i}K) - \text{Log}(-\mathbf{i}K) \right] \right| \\
&= \left| \Re \left[ \frac{1}{M} \sum_{i=1}^M \text{Log}(1 + \mathbf{i}\lambda_i(\mathbf{A})/K) \right] \right| = \frac{1}{2M} \sum_{i=1}^M \log(1 + \lambda_i(\mathbf{A})^2/K^2) \\
&\leq \frac{1}{2MK^2} \sum_{i=1}^M \lambda_i(\mathbf{A})^2 = \frac{\|\mathbf{A}\|_F^2}{2MK^2}.
\end{aligned}$$

For the imaginary part, we have

$$\begin{aligned}
& \left| \Im \left[ \frac{1}{M} \sum_{i=1}^M \text{Log}(\lambda_i(\mathbf{A}) - \mathbf{i}K) - \text{Log}(-\mathbf{i}K) \right] \right| \\
&= \left| \Im \left[ \frac{1}{M} \sum_{i=1}^M \text{Log}(1 + \mathbf{i}\lambda_i(\mathbf{A})/K) \right] \right| = \left| \frac{1}{M} \sum_{i=1}^M \arctan(\lambda_i(\mathbf{A})/K) \right| \\
&\leq \frac{1}{MK} \sum_{i=1}^M |\lambda_i(\mathbf{A})| \leq \frac{\|\mathbf{A}\|_F}{M^{1/2}K}.
\end{aligned}$$

As a result, we have

$$\begin{aligned}
& \mathbb{E} \left| \frac{1}{M} \sum_{i=1}^M \text{Log}(\lambda_i(\mathbf{A}) - \mathbf{i}K) - \text{Log}(-\mathbf{i}K) \right| \\
&\leq \mathbb{E} \left| \Re \left[ \frac{1}{M} \sum_{i=1}^M \text{Log}(\lambda_i(\mathbf{A}) - \mathbf{i}K) - \text{Log}(-\mathbf{i}K) \right] \right| + \mathbb{E} \left| \Im \left[ \frac{1}{M} \sum_{i=1}^M \text{Log}(\lambda_i(\mathbf{A}) - \mathbf{i}K) - \text{Log}(-\mathbf{i}K) \right] \right| \\
&\leq \frac{\mathbb{E}[\|\mathbf{A}\|_F^2]}{2MK^2} + \frac{\mathbb{E}[\|\mathbf{A}\|_F^2]^{1/2}}{M^{1/2}K}.
\end{aligned}$$

Note that

$$\frac{1}{M} \mathbb{E}[\|\mathbf{A}\|_F^2] \leq \frac{1}{M} \left( \mathbb{E}\|s_1 \mathbf{I}_N + s_2 \mathbf{Q}\|_F^2 + \mathbb{E}\|t_1 \mathbf{I}_n + t_2 \mathbf{H}\|_F^2 + 2\mathbb{E}[\|\mathbf{Z} + p\mathbf{Z}_1\|_F^2] \right) = O_d(1).$$

This proves Eq. (214).

**Step 2. Asymptotics of  $g(\mathbf{i}K; \mathbf{q})$ .**

Note that we have formula

$$g(\mathbf{i}K; \mathbf{q}) = \Xi(\mathbf{i}K, m_1(\mathbf{i}K; \mathbf{q}), m_2(\mathbf{i}K; \mathbf{q}); \mathbf{q}),$$

where the formula of  $\Xi$  is given by Eq. (66). Define

$$\begin{aligned}
\Xi_1(z_1, z_2; \mathbf{q}) &\equiv \log[(s_2 z_1 + 1)(t_2 z_2 + 1) - \mu_1^2(1+p)^2 z_1 z_2] - \mu_*^2 z_1 z_2 + s_1 z_1 + t_1 z_2, \\
\Xi_2(\xi, z_1, z_2) &\equiv -\psi_1 \log(z_1/\psi_1) - \psi_2 \log(z_2/\psi_2) - \xi(z_1 + z_2) - \psi_1 - \psi_2.
\end{aligned} \tag{216}$$

Then we have  $\Xi(\xi, z_1, z_2; \mathbf{q}) = \Xi_1(z_1, z_2; \mathbf{q}) + \Xi_2(\xi, z_1, z_2)$ . It is easy to see that for any fixed  $\mathbf{q}$ , we have

$$\lim_{z_1, z_2 \rightarrow 0} \Xi_1(z_1, z_2, \mathbf{q}) = 0.$$

Moreover, we have

$$\begin{aligned}
& |\Xi_2(\mathbf{i}K, m_1(\mathbf{i}K), m_2(\mathbf{i}K)) - \Xi_2(\mathbf{i}K, \mathbf{i}\psi_1/K, \mathbf{i}\psi_2/K)| \\
&\leq |\psi_1| \log(-\mathbf{i}K m_1(\mathbf{i}K)/\psi_1) + |\psi_2| \log(-\mathbf{i}K m_2(\mathbf{i}K)/\psi_2) + |\mathbf{i}K m_1(\mathbf{i}K) + \psi_1| + |\mathbf{i}K m_2(\mathbf{i}K) + \psi_2|.
\end{aligned}$$



By Lemma E.1

$$\lim_{K \rightarrow \infty} |\mathbf{i}K m_1(\mathbf{i}K) + \psi_1| = \lim_{K \rightarrow \infty} |\mathbf{i}K m_2(\mathbf{i}K) + \psi_2| = 0.$$

Hence

$$\begin{aligned} \lim_{K \rightarrow \infty} |\Xi_2(\mathbf{i}K, m_1(\mathbf{i}K), m_2(\mathbf{i}K)) - \Xi_2(\mathbf{i}K, \mathbf{i}\psi_1/K, \mathbf{i}\psi_2/K)| &= 0, \\ \lim_{K \rightarrow \infty} \Xi_1(m_1(\mathbf{i}K), m_2(\mathbf{i}K), \mathbf{q}) &= 0. \end{aligned}$$

Noting that we have  $\Xi_2(\mathbf{i}K, \mathbf{i}\psi_1/K, \mathbf{i}\psi_2/K) = (\psi_1 + \psi_2)\text{Log}(-\mathbf{i}K)$ . This proves the lemma.  $\square$

**Lemma E.3.** *Follow the notations and settings of Proposition 7.4. For fixed  $u \in \mathbb{R}_+$ , we have*

$$\begin{aligned} \limsup_{d \rightarrow \infty} \sup_{\mathbf{q} \in \mathbb{R}^5} \mathbb{E} \|\nabla_{\mathbf{q}} G_d(\mathbf{i}u; \mathbf{q}) - \nabla_{\mathbf{q}} g(\mathbf{i}u; \mathbf{q})\|_2 &< \infty, \\ \limsup_{d \rightarrow \infty} \sup_{\mathbf{q} \in \mathbb{R}^5} \mathbb{E} \|\nabla_{\mathbf{q}}^2 G_d(\mathbf{i}u; \mathbf{q}) - \nabla_{\mathbf{q}}^2 g(\mathbf{i}u; \mathbf{q})\|_{\text{op}} &< \infty, \\ \limsup_{d \rightarrow \infty} \sup_{\mathbf{q} \in \mathbb{R}^5} \mathbb{E} \|\nabla_{\mathbf{q}}^3 G_d(\mathbf{i}u; \mathbf{q}) - \nabla_{\mathbf{q}}^3 g(\mathbf{i}u; \mathbf{q})\|_{\text{op}} &< \infty. \end{aligned}$$

*Proof of Lemma E.3.* Define  $\mathbf{q} = (s_1, s_2, t_1, t_2, p) = (q_1, q_2, q_3, q_4, q_5)$ , and

$$\mathbf{S}_1 = \begin{bmatrix} \mathbf{I}_N & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix}, \quad \mathbf{S}_2 = \begin{bmatrix} \mathbf{Q} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix}, \quad \mathbf{S}_3 = \begin{bmatrix} \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{I}_n \end{bmatrix}, \quad \mathbf{S}_4 = \begin{bmatrix} \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{H} \end{bmatrix}, \quad \mathbf{S}_5 = \begin{bmatrix} \mathbf{0} & \mathbf{Z}_1^T \\ \mathbf{Z}_1 & \mathbf{0} \end{bmatrix}.$$

Then by the bound on the operator norm of Wishart matrix [AGZ09], for any fixed  $k \in \mathbb{N}$ , we have

$$\limsup_{d \rightarrow \infty} \sup_{i \in [5]} \mathbb{E}[\|\mathbf{S}_i\|_{\text{op}}^{2k}] < \infty.$$

Moreover, define  $\mathbf{R} = (\mathbf{A} - \mathbf{i}u\mathbf{I}_M)^{-1}$ . Then we have almost surely  $\sup_{\mathbf{q}} \|\mathbf{R}\|_{\text{op}} \leq 1/u$ .

Therefore

$$\begin{aligned} \sup_{\mathbf{q}} \mathbb{E} |\partial_{q_i} G_d(\mathbf{i}u; \mathbf{q})| &= \sup_{\mathbf{q}} \frac{1}{d} \mathbb{E} |\text{Tr}(\mathbf{R} \mathbf{S}_i)| \leq \sup_{\mathbf{q}} \frac{1}{u} \mathbb{E} [\|\mathbf{S}_i\|_{\text{op}}] = O_d(1), \\ \sup_{\mathbf{q}} \mathbb{E} |\partial_{q_i, q_j}^2 G_d(\mathbf{i}u; \mathbf{q})| &= \sup_{\mathbf{q}} \frac{1}{d} \mathbb{E} |\text{Tr}(\mathbf{R} \mathbf{S}_i \mathbf{R} \mathbf{S}_j)| \leq \sup_{\mathbf{q}} \frac{1}{u^2} (\mathbb{E} [\|\mathbf{S}_i\|_{\text{op}}^2] \mathbb{E} [\|\mathbf{S}_j\|_{\text{op}}^2])^{1/2} = O_d(1), \\ \sup_{\mathbf{q}} \mathbb{E} |\partial_{q_i, q_j, q_l}^3 G_d(\mathbf{i}u; \mathbf{q})| &= \sup_{\mathbf{q}} \frac{1}{d} \left[ \mathbb{E} |\text{Tr}(\mathbf{R} \mathbf{S}_i \mathbf{R} \mathbf{S}_j \mathbf{R} \mathbf{S}_l)| + \mathbb{E} |\text{Tr}(\mathbf{R} \mathbf{S}_i \mathbf{R} \mathbf{S}_l \mathbf{R} \mathbf{S}_j)| \right] \\ &\leq 2 \sup_{\mathbf{q}} \frac{1}{u^3} \left[ \mathbb{E} [\|\mathbf{S}_i\|_{\text{op}}^4] \mathbb{E} [\|\mathbf{S}_j\|_{\text{op}}^4] \mathbb{E} [\|\mathbf{S}_l\|_{\text{op}}^4] \right]^{1/4} = O_d(1). \end{aligned}$$

Similarly we can show that for fixed  $u > 0$ , we have  $\sup_{\mathbf{q} \in \mathbb{R}^5} \|\nabla_{\mathbf{q}}^j g(\mathbf{i}u; \mathbf{q})\| < \infty$  for  $j = 1, 2, 3$ . The lemma holds by that

$$\limsup_{d \rightarrow \infty} \sup_{\mathbf{q} \in \mathbb{R}^5} \mathbb{E} \|\nabla_{\mathbf{q}}^j G_d(\mathbf{i}u; \mathbf{q}) - \nabla_{\mathbf{q}}^j g(\mathbf{i}u; \mathbf{q})\| \leq \limsup_{d \rightarrow \infty} \sup_{\mathbf{q} \in \mathbb{R}^5} \left[ \mathbb{E} \|\nabla_{\mathbf{q}}^j G_d(\mathbf{i}u; \mathbf{q})\| + \|\nabla_{\mathbf{q}}^j g(\mathbf{i}u; \mathbf{q})\| \right] < \infty$$

for  $j = 1, 2, 3$ .  $\square$

**Lemma E.4.** *Let  $f \in C^2([a, b])$ . Then we have*

$$\sup_{x \in [a, b]} |f'(x)| \leq \left| \frac{f(a) - f(b)}{a - b} \right| + \frac{1}{2} \sup_{x \in [a, b]} |f''(x)| \cdot |a - b|.$$

*Proof.* Let  $x_0 \in [a, b]$ . Performing Taylor expansion of  $f(a)$  and  $f(b)$  at  $x_0$ , we have (for  $c_1 \in [a, x_0]$  and  $c_2 \in [x_0, b]$ )

$$\begin{aligned} f(a) &= f(x_0) + f'(x_0)(a - x_0) + f''(c_1)(a - x_0)^2/2, \\ f(b) &= f(x_0) + f'(x_0)(b - x_0) + f''(c_2)(b - x_0)^2/2. \end{aligned}$$

Then we have

$$|f'(x_0)| = \left| \frac{f(a) - f(b)}{a - b} + \frac{f''(c_1)(a - x_0)^2 - f''(c_2)(b - x_0)^2}{2(a - b)} \right| \leq \left| \frac{f(a) - f(b)}{a - b} \right| + \frac{1}{2} \sup_{x \in [a, b]} |f''(x)| \cdot |a - b|.$$

This proves the lemma.  $\square$

## E.2 Proof of Proposition 7.4

By the expression of  $\Xi$  in Eq. (66), we have

$$\begin{aligned}\partial_{z_1}\Xi(\xi, z_1, z_2; \mathbf{q}) &= \frac{s_2(t_2 z_2 + 1) - \mu_1^2(1+p)^2 z_2}{(s_2 z_1 + 1)(t_2 z_2 + 1) - \mu_1^2(1+p)^2 z_1 z_2} - \mu_\star^2 z_2 + s_1 - \psi_1/z_1 - \xi, \\ \partial_{z_2}\Xi(\xi, z_1, z_2; \mathbf{q}) &= \frac{t_2(s_2 z_1 + 1) - \mu_1^2(1+p)^2 z_1}{(s_2 z_1 + 1)(t_2 z_2 + 1) - \mu_1^2(1+p)^2 z_1 z_2} - \mu_\star^2 z_1 + s_2 - \psi_2/z_2 - \xi.\end{aligned}$$

By fixed point equation (61) with  $F$  defined in (60), we obtain so that

$$\nabla_{(z_1, z_2)}\Xi(\xi, z_1, z_2; \mathbf{q})|_{(z_1, z_2)=(m_1(\xi; \mathbf{q}), m_2(\xi; \mathbf{q}))} = \mathbf{0}.$$

As a result, by the definition of  $g$  given in Eq. (67), and by formula for implicit differentiation, we have

$$\frac{d}{d\xi}g(\xi; \mathbf{q}) = -m(\xi; \mathbf{q}).$$

Hence, for any  $\xi \in \mathbb{C}_+$  and  $K \in \mathbb{R}$  and compact continuous path  $\phi(\xi, iK)$

$$g(\xi; \mathbf{q}) - g(iK; \mathbf{q}) = \int_{\phi(\xi, iK)} m(\eta; \mathbf{q}) d\eta. \quad (217)$$

By Proposition 7.3, for any  $\xi \in \mathbb{C}_+$  and  $K \in \mathbb{R}$ , we have

$$G_d(\xi; \mathbf{q}) - G_d(iK; \mathbf{q}) = \int_{\phi(\xi, iK)} M_d(\eta; \mathbf{q}) d\eta. \quad (218)$$

Combining Eq. (218) with Eq. (217), we get

$$\mathbb{E}[|G_d(\xi; \mathbf{q}) - g(\xi; \mathbf{q})|] \leq \int_{\phi(\xi, iK)} \mathbb{E}|M_d(\eta; \mathbf{q}) - m(\eta; \mathbf{q})| d\eta + \mathbb{E}|G_d(iK; \mathbf{q}) - g(iK; \mathbf{q})|. \quad (219)$$

By Proposition 7.2, we have

$$\lim_{d \rightarrow \infty} \int_{\phi(\xi, iK)} \mathbb{E}|M_d(\eta; \mathbf{q}) - m(\eta; \mathbf{q})| d\eta = 0. \quad (220)$$

By Lemma E.2, we have

$$\lim_{K \rightarrow \infty} \sup_{d \geq d_0} \mathbb{E}|G_d(iK; \mathbf{q}) - g(iK; \mathbf{q})| = 0. \quad (221)$$

Combining Eq. (219), (220) and (221), we get Eq. (68).

For fixed  $\xi \in \mathbb{C}_+$ , define  $f_d(\mathbf{q}) = G_d(\xi, \mathbf{q}) - g(\xi; \mathbf{q})$ . By Lemma E.4, there exists some generic constant  $C$ , such that

$$\sup_{\mathbf{q} \in \mathbb{B}(\mathbf{0}, \varepsilon)} \|\nabla f_d(\mathbf{q})\|_2 \leq C \left[ \varepsilon^{-1} \sup_{\mathbf{q} \in \mathbb{B}(\mathbf{0}, \varepsilon)} |f_d(\mathbf{q})| + \varepsilon \sup_{\mathbf{q} \in \mathbb{B}(\mathbf{0}, \varepsilon)} \|\nabla^2 f_d(\mathbf{q})\|_{\text{op}} \right].$$

By Eq. (68) and Lemma E.3, we have

$$\limsup_{d \rightarrow \infty} \mathbb{E} \left[ \sup_{\mathbf{q} \in \mathbb{B}(\mathbf{0}, \varepsilon)} \|\nabla f_d(\mathbf{q})\|_2 \right] \leq C' \varepsilon.$$

Sending  $\varepsilon \rightarrow 0$  gives Eq. (69). Using similar argument we get Eq. (70).

## F Proof of Theorem 3, 4, and 5

### F.1 Proof of Theorem 3

To prove this theorem, we just need to show that

$$\begin{aligned}\lim_{\bar{\lambda} \rightarrow 0} \mathcal{B}(\zeta, \psi_1, \psi_2, \bar{\lambda}) &= \mathcal{B}_{\text{rless}}(\zeta, \psi_1, \psi_2), \\ \lim_{\bar{\lambda} \rightarrow 0} \mathcal{V}(\zeta, \psi_1, \psi_2, \bar{\lambda}) &= \mathcal{V}_{\text{rless}}(\zeta, \psi_1, \psi_2).\end{aligned}$$

More specifically, we just need to show that, the formula for  $\chi$  defined in Eq. (16) as  $\bar{\lambda} \rightarrow 0$  coincides with the formula for  $\chi$  defined in Eq. (24). By the relationship of  $\chi$  and  $m_1 m_2$  as per Eq. (81), we just need to show the lemma below.

**Lemma F.1.** *Let Assumption 1 and 2 hold. For fixed  $\xi \in \mathbb{C}_+$ , let  $m_1(\xi)$  and  $m_2(\xi)$  be defined by*

$$\begin{aligned}m_1(\xi) &= \lim_{d \rightarrow \infty} \frac{1}{d} \mathbb{E} \{ \text{Tr}_{[1, N]} [([\mathbf{0}, \mathbf{Z}^\top; \mathbf{Z}, \mathbf{0}] - \xi \mathbf{I}_M)^{-1}] \}, \\ m_2(\xi) &= \lim_{d \rightarrow \infty} \frac{1}{d} \mathbb{E} \{ \text{Tr}_{[N+1, M]} [([\mathbf{0}, \mathbf{Z}^\top; \mathbf{Z}, \mathbf{0}] - \xi \mathbf{I}_M)^{-1}] \}.\end{aligned}\tag{222}$$

By Proposition 7.2 this is equivalently saying  $m_1(\xi), m_2(\xi)$  is the analytic continuation of solution of Eq. (61) as defined in Proposition 7.2, when  $\mathbf{q} = \mathbf{0}$ . Defining  $\psi = \min(\psi_1, \psi_2)$ , we have

$$\lim_{u \rightarrow 0} [m_1(iu)m_2(iu)] = -\frac{[(\psi\zeta^2 - \zeta^2 - 1)^2 + 4\zeta^2\psi]^{1/2} + (\psi\zeta^2 - \zeta^2 - 1)}{2\mu_\star^2\zeta^2}.\tag{223}$$

*Proof of Lemma F.1.* In the following, we consider the case  $\psi_2 > \psi_1$ . The proof for the case  $\psi_2 < \psi_1$  is the same, and the case  $\psi_1 = \psi_2$  is simpler. By Proposition 7.2,  $m_1 = m_1(iu)$  and  $m_2 = m_2(iu)$  must satisfy Eq. (61) for  $\xi = iu$  and  $\mathbf{q} = \mathbf{0}$ . A reformulation for Eq. (61) for  $\mathbf{q} = \mathbf{0}$  yields

$$\frac{-\mu_1^2 m_1 m_2}{1 - \mu_1^2 m_1 m_2} - \mu_\star^2 m_1 m_2 - \psi_1 - iu \cdot m_1 = 0,\tag{224}$$

$$\frac{-\mu_1^2 m_1 m_2}{1 - \mu_1^2 m_1 m_2} - \mu_\star^2 m_1 m_2 - \psi_2 - iu \cdot m_2 = 0.\tag{225}$$

Defining  $m_0(iu) = m_1(iu)m_2(iu)$ . Then  $m_0$  must satisfy the following equation

$$-u^2 m_0 = \left( \frac{-\mu_1^2 m_0}{1 - \mu_1^2 m_0} - \mu_\star^2 m_0 - \psi_1 \right) \left( \frac{-\mu_1^2 m_0}{1 - \mu_1^2 m_0} - \mu_\star^2 m_0 - \psi_2 \right).$$

Note we must have  $|m_0(iu)| \leq |m_1(iu)| \cdot |m_2(iu)| \leq \psi_1 \psi_2 / u^2$ , and hence  $|u^2 m_0| = O_u(1)$  (as  $u \rightarrow 0$ ). This implies that

$$\frac{-\mu_1^2 m_0}{1 - \mu_1^2 m_0} - \mu_\star^2 m_0 = O_u(1),$$

and hence  $m_0 = O_u(1)$ . Taking the difference between Eq. (224) and (225), we get

$$m_2 - m_1 = -(\psi_2 - \psi_1)/(iu).\tag{226}$$

This implies one of  $m_1$  and  $m_2$  should be of order  $1/u$  and the other one should be of order  $u$ , as  $u \rightarrow 0$ .

By definition of  $m_1$  and  $m_2$  in Eq. (222), we have

$$\begin{aligned}m_1(iu) &= iu \lim_{d \rightarrow \infty} \frac{1}{d} \mathbb{E} \{ \text{Tr}[(\mathbf{Z}^\top \mathbf{Z} + u^2 \mathbf{I}_N)^{-1}] \}, \\ m_2(iu) &= iu \lim_{d \rightarrow \infty} \frac{1}{d} \mathbb{E} \{ \text{Tr}[(\mathbf{Z} \mathbf{Z}^\top + u^2 \mathbf{I}_N)^{-1}] \}.\end{aligned}$$

When  $n > N$ ,  $(\mathbf{Z}\mathbf{Z}^\top + u^2\mathbf{I}_N)$  has  $(n - N)$  number of eigenvalues that are  $u^2$ , and therefore we must have  $m_2(iu) = O_u(1/u)$ . Hence  $m_1(iu) = O_u(u)$ . Moreover, when  $u > 0$ ,  $m_1(iu)$  and  $m_2(iu)$  are purely imaginary and  $\Im m_1(iu), \Im m_2(iu) > 0$ . This implies that  $m_0(iu)$  must be a real number which is non-positive.

By Eq. (224) and  $\lim_{u \rightarrow 0} iu \cdot m_1(iu) = 0$ , all the accumulation points of  $m_1(iu)m_2(iu)$  as  $u \rightarrow 0$  should satisfy the quadratic equation

$$\frac{-\mu_1^2 m_\star}{1 - \mu_1^2 m_\star} - \mu_\star^2 m_\star - \psi_1 = 0.$$

Note that the above equation has only one non-positive solution, and  $m_0(iu)$  for any  $u > 0$  must be non-positive. Therefore  $\lim_{u \rightarrow 0} m_1(iu)m_2(iu)$  must exist and be the non-positive solution of the above quadratic equation. The right hand side of Eq. (223) gives the non-positive solution of the above quadratic equation.  $\square$

## F.2 Proof of Theorem 4

To prove this theorem, we just need to show that

$$\begin{aligned} \lim_{\psi_1 \rightarrow \infty} \mathcal{B}(\zeta, \psi_1, \psi_2, \bar{\lambda}) &= \mathcal{B}_{\text{wide}}(\zeta, \psi_2, \bar{\lambda}), \\ \lim_{\psi_1 \rightarrow \infty} \mathcal{V}(\zeta, \psi_1, \psi_2, \bar{\lambda}) &= \mathcal{V}_{\text{wide}}(\zeta, \psi_2, \bar{\lambda}). \end{aligned}$$

This follows by simple calculus and a lemma below.

**Lemma F.2.** *Let Assumption 1 and 2 hold. For fixed  $\xi \in \mathbb{C}_+$ , let  $m_1(\xi; \psi_1)$  and  $m_2(\xi; \psi_1)$  satisfies*

$$\begin{aligned} m_1(\xi; \psi_1) &= \lim_{d \rightarrow \infty, N/d \rightarrow \psi_1} \frac{1}{d} \mathbb{E}\{\text{Tr}_{[1, N]}[(\mathbf{0}, \mathbf{Z}^\top; \mathbf{Z}, \mathbf{0}] - \xi \mathbf{I}_M)^{-1}\}, \\ m_2(\xi; \psi_1) &= \lim_{d \rightarrow \infty, N/d \rightarrow \psi_1} \frac{1}{d} \mathbb{E}\{\text{Tr}_{[N+1, M]}[(\mathbf{0}, \mathbf{Z}^\top; \mathbf{Z}, \mathbf{0}] - \xi \mathbf{I}_M)^{-1}\}. \end{aligned}$$

By Proposition 7.2 this is equivalently saying  $m_1(\xi; \psi_1), m_2(\xi; \psi_1)$  is the analytic continuation of solution of Eq. (61) as defined in Proposition 7.2, when  $\mathbf{q} = \mathbf{0}$ . Then we have

$$\begin{aligned} &\lim_{\psi_1 \rightarrow \infty} [m_1(i(\psi_1 \psi_2 \mu_\star^2 \bar{\lambda})^{1/2}; \psi_1) m_2(i(\psi_1 \psi_2 \mu_\star^2 \bar{\lambda})^{1/2}; \psi_1)] \\ &= - \frac{[(\psi_2 \zeta^2 - \zeta^2 - (\bar{\lambda} \psi_2 + 1))^2 + 4 \zeta^2 \psi_2 (\bar{\lambda} \psi_2 + 1)]^{1/2} + (\psi_2 \zeta^2 - \zeta^2 - (\bar{\lambda} \psi_2 + 1))}{2 \mu_\star^2 \zeta^2 (\bar{\lambda} \psi_2 + 1)}. \end{aligned}$$

The proof of this lemma is similar to the proof of Lemma F.1.

## F.3 Proof of Theorem 5

To prove this theorem, we just need to show that

$$\begin{aligned} \lim_{\psi_2 \rightarrow \infty} \mathcal{B}(\zeta, \psi_1, \psi_2, \bar{\lambda}) &= \mathcal{B}_{\text{lsamp}}(\zeta, \psi_1, \bar{\lambda}), \\ \lim_{\psi_2 \rightarrow \infty} \mathcal{V}(\zeta, \psi_1, \psi_2, \bar{\lambda}) &= 0. \end{aligned}$$

This follows by simple calculus and a lemma below (this lemma is symmetric to Lemma F.2).

**Lemma F.3.** *Let Assumption 1 and 2 hold. For fixed  $\xi \in \mathbb{C}_+$ , let  $m_1(\xi; \psi_2)$  and  $m_2(\xi; \psi_2)$  satisfies*

$$\begin{aligned} m_1(\xi; \psi_2) &= \lim_{d \rightarrow \infty, n/d \rightarrow \psi_2} \frac{1}{d} \mathbb{E}\{\text{Tr}_{[1, N]}[(\mathbf{0}, \mathbf{Z}^\top; \mathbf{Z}, \mathbf{0}] - \xi \mathbf{I}_M)^{-1}\}, \\ m_2(\xi; \psi_2) &= \lim_{d \rightarrow \infty, n/d \rightarrow \psi_2} \frac{1}{d} \mathbb{E}\{\text{Tr}_{[N+1, M]}[(\mathbf{0}, \mathbf{Z}^\top; \mathbf{Z}, \mathbf{0}] - \xi \mathbf{I}_M)^{-1}\}. \end{aligned}$$

By Proposition 7.2 this is equivalently saying  $m_1(\xi; \psi_2), m_2(\xi; \psi_2)$  is the analytic continuation of solution of Eq. (61) as defined in Proposition 7.2, when  $\mathbf{q} = \mathbf{0}$ . Then we have

$$\begin{aligned} & \lim_{\psi_2 \rightarrow \infty} [m_1(i(\psi_1 \psi_2 \mu_\star^2 \bar{\lambda})^{1/2}; \psi_1) m_2(i(\psi_1 \psi_2 \mu_\star^2 \bar{\lambda})^{1/2}; \psi_1)] \\ &= - \frac{[(\psi_1 \zeta^2 - \zeta^2 - (\bar{\lambda} \psi_1 + 1))^2 + 4\zeta^2 \psi_1 (\bar{\lambda} \psi_1 + 1)]^{1/2} + (\psi_1 \zeta^2 - \zeta^2 - (\bar{\lambda} \psi_1 + 1))}{2\mu_\star^2 \zeta^2 (\bar{\lambda} \psi_1 + 1)}. \end{aligned}$$

The proof of this lemma is similar to the proof of Lemma F.1.

## G Proof of Proposition 4.1 and 4.2

### G.1 Proof of Proposition 4.1

**Proof of Point (1).** When  $\psi_1 \rightarrow 0$ , we have  $\chi = O(\psi_1)$ , so that  $\mathcal{E}_{1,\text{rless}} = -\psi_1 \psi_2 + O(\psi_1^2)$ ,  $\mathcal{E}_{2,\text{rless}} = O(\psi_1^2)$  and  $\mathcal{E}_{0,\text{rless}} = -\psi_1 \psi_2 + O(\psi_1^2)$ . This proves Point (1).

**Proof of Point (2).** When  $\psi_1 = \psi_2$ , substituting the expression for  $\chi$  into  $\mathcal{E}_{0,\text{rless}}$ , we can see that  $\mathcal{E}_{0,\text{rless}}(\zeta, \psi_2, \psi_2) = 0$ . We also see that  $\mathcal{E}_{1,\text{rless}}(\zeta, \psi_2, \psi_2) \neq 0$  and  $\mathcal{E}_{2,\text{rless}}(\zeta, \psi_2, \psi_2) \neq 0$ . This proves Point (2).

**Proof of Point (3).** When  $\psi_1 > \psi_2$ , we have

$$\begin{aligned} \lim_{\psi_1 \rightarrow \infty} \mathcal{E}_{0,\text{rless}}(\zeta, \psi_1, \psi_2)/\psi_1 &= (\psi_2 - 1)\chi^3 \zeta^6 + (1 - 3\psi_2)\chi^2 \zeta^4 + 3\psi_2 \chi \zeta^2 - \psi_2, \\ \lim_{\psi_1 \rightarrow \infty} \mathcal{E}_{1,\text{rless}}(\zeta, \psi_1, \psi_2)/\psi_1 &= \psi_2 \chi \zeta^2 - \psi_2, \\ \lim_{\psi_1 \rightarrow \infty} \mathcal{E}_{2,\text{rless}}(\zeta, \psi_1, \psi_2)/\psi_1 &= \chi^3 \zeta^6 - \chi^2 \zeta^4, \end{aligned}$$

This proves Point (3).

**Proof of Point (4).** For  $\psi_1 > \psi_2$ , taking derivative of  $\mathcal{B}_{\text{rless}}$  and  $\mathcal{V}_{\text{rless}}$  with respect to  $\psi_1$ , we have

$$\begin{aligned} \partial_{\psi_1} \mathcal{B}_{\text{rless}}(\zeta, \psi_1, \psi_2) &= (\partial_{\psi_1} \mathcal{E}_{1,\text{rless}} \cdot \mathcal{E}_{0,\text{rless}} - \partial_{\psi_1} \mathcal{E}_{0,\text{rless}} \cdot \mathcal{E}_{1,\text{rless}})/\mathcal{E}_{0,\text{rless}}^2, \\ \partial_{\psi_1} \mathcal{V}_{\text{rless}}(\zeta, \psi_1, \psi_2) &= (\partial_{\psi_1} \mathcal{E}_{2,\text{rless}} \cdot \mathcal{E}_{0,\text{rless}} - \partial_{\psi_1} \mathcal{E}_{0,\text{rless}} \cdot \mathcal{E}_{2,\text{rless}})/\mathcal{E}_{0,\text{rless}}^2. \end{aligned}$$

It is easy to check that when  $\psi_1 > \psi_2$ , the functions  $\partial_{\psi_1} \mathcal{E}_{1,\text{rless}} \cdot \mathcal{E}_{0,\text{rless}} - \partial_{\psi_1} \mathcal{E}_{0,\text{rless}} \cdot \mathcal{E}_{1,\text{rless}}$  and  $\partial_{\psi_1} \mathcal{E}_{2,\text{rless}} \cdot \mathcal{E}_{0,\text{rless}} - \partial_{\psi_1} \mathcal{E}_{0,\text{rless}} \cdot \mathcal{E}_{2,\text{rless}}$  are functions of  $\zeta$  and  $\psi_2$ , and are independent of  $\psi_1$  (note when  $\psi_1 > \psi_2$ ,  $\chi$  is a function of  $\psi_2$  and doesn't depend on  $\psi_1$ ). Therefore,  $\mathcal{B}_{\text{rless}}(\zeta, \cdot, \psi_2)$  and  $\mathcal{V}_{\text{rless}}(\zeta, \cdot, \psi_2)$  as functions of  $\psi_1$  must be strictly increasing, strictly decreasing, or staying constant on the interval  $\psi_1 \in (\psi_2, \infty)$ . However, we know  $\mathcal{B}_{\text{rless}}(\zeta, \psi_2, \psi_2) = \mathcal{V}_{\text{rless}}(\zeta, \psi_2, \psi_2) = \infty$ , and  $\mathcal{B}_{\text{rless}}(\zeta, \infty, \psi_2)$  and  $\mathcal{V}_{\text{rless}}(\zeta, \infty, \psi_2)$  are finite. Therefore, we must have that  $\mathcal{B}_{\text{rless}}$  and  $\mathcal{V}_{\text{rless}}$  are strictly decreasing on  $\psi_1 \in (\psi_2, \infty)$ .

### G.2 Proof of Proposition 4.2

In Proposition G.1 given by the following, we give a more precise description of the behavior of  $\mathcal{R}_{\text{wide}}$ , which is stronger than Proposition 4.2.

**Proposition G.1.** Denote

$$\begin{aligned}
\overline{\mathcal{R}}_{\text{wide}}(u, \rho, \psi_2) &= \frac{\psi_2 \rho + u^2}{(1 + \rho)(\psi_2 - 2u\psi_2 + u^2\psi_2 - u^2)}, \\
\omega(\bar{\lambda}, \zeta, \psi_2) &= - \frac{[(\psi_2 \zeta^2 - \zeta^2 - \bar{\lambda}\psi_2 - 1)^2 + 4\psi_2 \zeta^2 (\bar{\lambda}\psi_2 + 1)]^{1/2} + (\psi_2 \zeta^2 - \zeta^2 - \bar{\lambda}\psi_2 - 1)}{2(\bar{\lambda}\psi_2 + 1)}, \\
\omega_0(\zeta, \psi_2) &= - \frac{[(\psi_2 \zeta^2 - \zeta^2 - 1)^2 + 4\psi_2 \zeta^2]^{1/2} + (\psi_2 \zeta^2 - \zeta^2 - 1)}{2}, \\
\omega_1(\rho, \psi_2) &= - \frac{(\psi_2 \rho - \rho - 1) + [(\psi_2 \rho - \rho - 1)^2 + 4\psi_2 \rho]^{1/2}}{2}, \\
\rho_\star(\zeta, \psi_2) &= \frac{\omega_0^2 - \omega_1}{(1 - \psi_2)\omega_0 + \psi_2}, \\
\zeta_\star^2(\rho, \psi_2) &= \frac{\omega_1^2 - \omega_1}{\omega_1 - \psi_2\omega_1 + \psi_2}, \\
\bar{\lambda}_\star(\zeta, \psi_2, \rho) &= \frac{\zeta^2\psi_2 - \zeta^2\omega_1\psi_2 + \zeta^2\omega_1 + \omega_1 - \omega_1^2}{(\omega_1^2 - \omega_1)\psi_2}.
\end{aligned}$$

Fix  $\zeta, \psi_2 \in (0, \infty)$  and  $\rho \in (0, \infty)$ . Then the function  $\bar{\lambda} \mapsto \mathcal{R}_{\text{wide}}(\rho, \zeta, \psi_2, \bar{\lambda})$  is either strictly increasing in  $\bar{\lambda}$ , or strictly decreasing first and then strictly increasing.

Moreover, For any  $\rho < \rho_\star(\zeta, \psi_2)$ , we have

$$\begin{aligned}
\arg \min_{\bar{\lambda} \geq 0} \mathcal{R}_{\text{wide}}(\rho, \zeta, \bar{\lambda}, \psi_2) &= 0, \\
\min_{\bar{\lambda} \geq 0} \mathcal{R}_{\text{wide}}(\rho, \zeta, \bar{\lambda}, \psi_2) &= \overline{\mathcal{R}}_{\text{wide}}(\omega_0(\zeta, \psi_2), \rho, \psi_2).
\end{aligned}$$

For any  $\rho \geq \rho_\star(\zeta, \psi_2)$ , we have

$$\begin{aligned}
\arg \min_{\bar{\lambda} \geq 0} \mathcal{R}_{\text{wide}}(\rho, \zeta, \bar{\lambda}, \psi_2) &= \bar{\lambda}_\star(\zeta, \psi_2, \rho), \\
\min_{\bar{\lambda} \geq 0} \mathcal{R}_{\text{wide}}(\rho, \zeta, \bar{\lambda}, \psi_2) &= \overline{\mathcal{R}}_{\text{wide}}(\omega_1(\rho, \psi_2), \rho, \psi_2).
\end{aligned}$$

Minimizing over  $\bar{\lambda}$  and  $\zeta$ , we have

$$\min_{\zeta, \bar{\lambda} \geq 0} \mathcal{R}_{\text{wide}}(\rho, \zeta, \bar{\lambda}, \psi_2) = \overline{\mathcal{R}}_{\text{wide}}(\omega_1(\rho, \psi_2), \rho, \psi_2).$$

The minimizer is achieved for any  $\zeta^2 \geq \zeta_\star^2(\rho, \psi_2)$ , and  $\bar{\lambda} = \bar{\lambda}_\star(\zeta, \psi_2, \rho)$ .

In the following, we prove Proposition G.1. It is easy to see that

$$\mathcal{R}_{\text{wide}}(\rho, \zeta, \bar{\lambda}, \psi_2) = \overline{\mathcal{R}}_{\text{wide}}(\omega(\bar{\lambda}, \zeta, \psi_2), \rho, \psi_2).$$

Hence we study the properties of  $\overline{\mathcal{R}}_{\text{wide}}$  first.

**Step 1. Properties of the function  $\overline{\mathcal{R}}_{\text{wide}}$ .** Calculating the derivative of  $\overline{\mathcal{R}}_{\text{wide}}$  with respect to  $u$ , we have

$$\partial_u \overline{\mathcal{R}}_{\text{wide}}(u, \rho, \psi_2) = -2\psi_2[u^2 + (\psi_2\rho - \rho - 1)u - \psi_2\rho] / [(1 + \rho)(\psi_2 - 2u\psi_2 + u^2\psi_2 - u^2)^2].$$

Note the equation

$$u^2 + (\psi_2\rho - \rho - 1)u - \psi_2\rho = 0$$

has one negative and one positive solution, and  $\omega_1$  is the negative solution of the above equation. Therefore, when  $u \leq \omega_1$ ,  $\overline{\mathcal{R}}_{\text{wide}}$  will be strictly decreasing in  $u$ ; when  $0 \geq u \geq \omega_1$ ,  $\overline{\mathcal{R}}_{\text{wide}}$  will be strictly increasing in  $u$ . Therefore, we have

$$\arg \min_{u \in (-\infty, 0]} \overline{\mathcal{R}}_{\text{wide}}(u, \rho, \psi_2) = \omega_1(\rho, \psi_2).$$

**Step 2. Properties of the function  $\mathcal{R}_{\text{wide}}$ .** For fixed  $(\zeta, \rho, \psi_2)$ , we look at the minimizer over  $\bar{\lambda}$  of the function  $\mathcal{R}_{\text{wide}}(\rho, \zeta, \bar{\lambda}, \psi_2) = \mathcal{R}_{\text{wide}}(\omega(\bar{\lambda}, \zeta, \psi_2), \rho, \psi_2)$ . The minimum  $\min_{\bar{\lambda} \geq 0} \mathcal{R}_{\text{wide}}(\rho, \zeta, \bar{\lambda}, \psi_2)$  could be different from the minimum  $\min_{u \in (-\infty, 0]} \mathcal{R}_{\text{wide}}(u, \rho, \psi_2)$ , since  $\arg \min_{u \in (-\infty, 0]} \mathcal{R}_{\text{wide}}(u, \rho, \psi_2) = \omega_1(\rho, \psi_2)$  may not be achievable by  $\omega(\bar{\lambda}, \zeta, \psi_2)$  when  $\bar{\lambda} \geq 0$ .

One observation is that  $\omega(\cdot, \psi_2, \zeta)$  as a function of  $\bar{\lambda}$  is always negative and increasing.

**Lemma G.1.** *Let*

$$\omega(\bar{\lambda}, \psi_2, \zeta) = -\frac{[(\psi_2 \zeta^2 - \zeta^2 - \bar{\lambda} \psi_2 - 1)^2 + 4\psi_2 \zeta^2 (\bar{\lambda} \psi_2 + 1)]^{1/2} + (\psi_2 \zeta^2 - \zeta^2 - \bar{\lambda} \psi_2 - 1)}{2(\bar{\lambda} \psi_2 + 1)}.$$

*Then for any  $\psi_2 \in (0, \infty)$ ,  $\zeta \in (0, \infty)$  and  $\bar{\lambda} > 0$ , we have*

$$\begin{aligned} \omega(\bar{\lambda}, \psi_2, \zeta) &< 0, \\ \partial_{\bar{\lambda}} \omega(\bar{\lambda}, \psi_2, \zeta) &> 0. \end{aligned}$$

Let us for now admit this lemma holds. When  $\rho$  is such that  $\omega_1 > \omega_0$  (i.e.  $\rho < \rho_*(\zeta, \psi_2)$ ), then we can choose  $\bar{\lambda} = \bar{\lambda}_*(\zeta, \psi_2, \rho) > 0$  such that  $\omega(\bar{\lambda}, \zeta, \psi_2) = \omega(\bar{\lambda}_*, \zeta, \psi_2) = \omega_1(\rho, \psi_2)$ , and then  $\mathcal{R}_{\text{wide}}(\rho, \zeta, \bar{\lambda}_*(\zeta, \psi_2, \rho), \psi_2) = \mathcal{R}_{\text{wide}}(\omega_1(\rho, \psi_2), \rho, \psi_2)$  gives the minimum of  $\mathcal{R}_{\text{wide}}$  optimizing over  $\bar{\lambda} \in [0, \infty)$ . When  $\rho$  is such that  $\omega_1 < \omega_0$  (i.e.  $\rho > \rho_*(\zeta, \psi_2)$ ), there is not a  $\bar{\lambda}$  such that  $\omega(\bar{\lambda}, \zeta, \psi_2) = \omega_1(\rho, \psi_2)$  holds. Therefore, the best we can do is to take  $\bar{\lambda} = 0$ , and then  $\mathcal{R}_{\text{wide}}(\rho, \zeta, 0, \psi_2) = \mathcal{R}_{\text{wide}}(\omega_0(\rho, \psi_2), \rho, \psi_2)$  gives the minimum of  $\mathcal{R}_{\text{wide}}$  optimizing over  $\bar{\lambda} \in [0, \infty)$ .

Finally, when we minimize  $\mathcal{R}_{\text{wide}}(\rho, \zeta, \bar{\lambda}, \psi_2)$  jointly over  $\zeta$  and  $\bar{\lambda}$ , note that as long as  $\zeta^2 \geq \zeta_*^2$ , we can choose  $\bar{\lambda} = \bar{\lambda}_*(\zeta, \psi_2, \rho) > 0$  such that  $\omega(\bar{\lambda}, \zeta, \psi_2) = \omega(\bar{\lambda}_*, \zeta, \psi_2) = \omega_1(\rho, \psi_2)$ , and then  $\mathcal{R}_{\text{wide}}(\rho, \zeta, \bar{\lambda}_*(\zeta, \psi_2, \rho), \psi_2) = \mathcal{R}_{\text{wide}}(\omega_1(\rho, \psi_2), \rho, \psi_2)$  gives the minimum of  $\mathcal{R}_{\text{wide}}$  optimizing over  $\bar{\lambda} \in [0, \infty)$  and  $\zeta \in (0, \infty)$ . This proves the proposition. In the following, we prove Lemma G.1.

*Proof of Lemma G.1.* It is easy to see that  $\omega(\bar{\lambda}, \psi_2, \zeta) < 0$ . In the following, we show  $\partial_{\bar{\lambda}} \omega(\bar{\lambda}, \psi_2, \zeta) > 0$ .

**Step 1. When  $\psi_2 \geq 1$ .** We have

$$\omega = -\frac{[(\psi_2 \zeta^2 - \zeta^2 - \bar{\lambda} \psi_2 - 1)^2 + 4\psi_2 \zeta^2 (\bar{\lambda} \psi_2 + 1)]^{1/2} + (\psi_2 \zeta^2 - \zeta^2 - \bar{\lambda} \psi_2 - 1)}{2(\bar{\lambda} \psi_2 + 1)}.$$

Then we have

$$\partial_{\bar{\lambda}} \omega = \frac{(\psi_2 - 1)[(\bar{\lambda} \psi_2 - \psi_2 \zeta^2 + \zeta^2 + 1)^2 + 4\psi_2 \zeta^2 (\bar{\lambda} \psi_2 + 1)]^{1/2} + (\bar{\lambda} \psi_2^2 + \bar{\lambda} \psi_2 + (\psi_2 - 1)^2 \zeta^2 + \psi_2 + 1)}{2\psi_2^2 \bar{\lambda} [\bar{\lambda}^2 \psi_2^2 (\bar{\lambda} \psi_2 - \psi_2 \zeta^2 + \zeta^2 + 1)^2 + 4\bar{\lambda}^2 \psi_2^3 \zeta^2 (\bar{\lambda} \psi_2 + 1)]^{1/2} (\bar{\lambda} \psi_2 + 1)^2}$$

It is easy to see that, when  $\bar{\lambda} > 0$  and  $\psi_2 > 1$ , both the denominator and numerator is positive, so that  $\partial_{\bar{\lambda}} \omega > 0$ .

**Step 2. When  $\psi_2 < 1$ .** Note  $\omega$  is the negative solution of the quadratic equation

$$(\bar{\lambda} \psi_2 + 1)\omega^2 + (\psi_2 \zeta^2 - \zeta^2 - \bar{\lambda} \psi_2 - 1)\omega - \psi_2 \zeta^2 = 0.$$

Differentiating the quadratic equation with respect to  $\bar{\lambda}$ , we have

$$\psi_2 \omega^2 + 2(\bar{\lambda} \psi_2 + 1)\omega \partial_{\bar{\lambda}} \omega - \psi_2 \omega + (\psi_2 \zeta^2 - \zeta^2 - \bar{\lambda} \psi_2 - 1)\partial_{\bar{\lambda}} \omega = 0,$$

which gives

$$\partial_{\bar{\lambda}} \omega = (\psi_2 \omega - \psi_2 \omega^2) / [2(\bar{\lambda} \psi_2 + 1)\omega + \psi_2 \zeta^2 - \zeta^2 - \bar{\lambda} \psi_2 - 1] = (\psi_2 \omega - \psi_2 \omega^2) / [(\bar{\lambda} \psi_2 + 1)(2\omega - 1) + (\psi_2 - 1)\zeta^2].$$

We can see that, since  $\omega < 0$ , when  $\psi_2 < 1$ , both the denominator and numerator is negative. This proves  $\partial_{\bar{\lambda}} \omega > 0$  when  $\psi_2 < 1$ .  $\square$

## H Proof sketch for Theorem 6

In this section, we sketch the calculations of Theorem 6. We assume  $\psi_{1,d} \equiv N/d = \psi_1$  and  $\psi_{2,d} \equiv n/d = \psi_2$  are constants independent of  $d$ .

### Step 1. The expectation of regularized training error.

By Eq. (45), the regularized training error of random features regression gives

$$\begin{aligned} L_{\text{RF}}(f_d, \mathbf{X}, \boldsymbol{\Theta}, \lambda) &= \min_{\mathbf{a}} \left[ \frac{1}{n} \sum_{i=1}^n \left( y_i - \sum_{j=1}^N a_j \sigma(\langle \boldsymbol{\theta}_j, \mathbf{x}_i \rangle / \sqrt{d}) \right)^2 + \lambda \psi_1 \|\mathbf{a}\|_2^2 \right] \\ &= \min_{\mathbf{a}} \left[ \frac{1}{n} \|\mathbf{y} - \sqrt{d} \mathbf{Z} \mathbf{a}\|^2 + \lambda \psi_1 \|\mathbf{a}\|_2^2 \right] \\ &= \frac{1}{n} \|\mathbf{y} - \mathbf{Z}(\mathbf{Z}^\top \mathbf{Z} + \psi_1 \psi_2 \lambda \mathbf{I}_N)^{-1} \mathbf{Z}^\top \mathbf{y}\|^2 + \lambda \psi_1 \|(\mathbf{Z}^\top \mathbf{Z} + \psi_1 \psi_2 \lambda \mathbf{I}_N)^{-1} \mathbf{Z}^\top \mathbf{y}\|_2^2 / d \\ &= \frac{1}{n} \left[ \|\mathbf{y}\|_2^2 - \mathbf{y}^\top \mathbf{Z}(\mathbf{Z}^\top \mathbf{Z} + \psi_1 \psi_2 \lambda \mathbf{I}_N)^{-1} \mathbf{Z}^\top \mathbf{y} \right]. \end{aligned}$$

Its expectation with respect to  $f_d$  (that satisfies Assumption 4) and  $\boldsymbol{\varepsilon}$  gives

$$\begin{aligned} &\mathbb{E}_{\boldsymbol{\beta}, \boldsymbol{\varepsilon}}[L_{\text{RF}}(f_d, \mathbf{X}, \boldsymbol{\Theta}, \lambda)] \\ &= \frac{1}{n} \left[ \mathbb{E}_{\boldsymbol{\beta}, \boldsymbol{\varepsilon}}[\|\mathbf{y}\|_2^2] - \mathbb{E}_{\boldsymbol{\beta}, \boldsymbol{\varepsilon}}[\mathbf{y}^\top \mathbf{Z}(\mathbf{Z}^\top \mathbf{Z} + \psi_1 \psi_2 \lambda \mathbf{I}_N)^{-1} \mathbf{Z}^\top \mathbf{y}] \right] \\ &= \mathbb{E}_{\boldsymbol{\beta}}[\|f_d\|_{L^2}^2] + \tau^2 - \frac{1}{n} \mathbb{E}_{\boldsymbol{\beta}} \left[ \mathbf{f}^\top \mathbf{Z}(\mathbf{Z}^\top \mathbf{Z} + \psi_1 \psi_2 \lambda \mathbf{I}_N)^{-1} \mathbf{Z}^\top \mathbf{f} \right] - \frac{1}{n} \mathbb{E}_{\boldsymbol{\varepsilon}} \left[ \boldsymbol{\varepsilon}^\top \mathbf{Z}(\mathbf{Z}^\top \mathbf{Z} + \psi_1 \psi_2 \lambda \mathbf{I}_N)^{-1} \mathbf{Z}^\top \boldsymbol{\varepsilon} \right] \\ &= \mathbb{E}_{\boldsymbol{\beta}}[\|f_d\|_{L^2}^2] + \tau^2 - \frac{1}{n} \mathbb{E}_{\boldsymbol{\beta}} \left[ \left( \sum_{k=0}^{\infty} \mathbf{Y}_{\mathbf{x},k} \beta_k \right)^\top \mathbf{Z}(\mathbf{Z}^\top \mathbf{Z} + \psi_1 \psi_2 \lambda \mathbf{I}_N)^{-1} \mathbf{Z}^\top \left( \sum_{k=0}^{\infty} \mathbf{Y}_{\mathbf{x},k} \beta_k \right) \right] \\ &\quad - \frac{\tau^2}{n} \text{Tr} \left( (\mathbf{Z}^\top \mathbf{Z} + \psi_1 \psi_2 \lambda \mathbf{I}_N)^{-1} \mathbf{Z}^\top \mathbf{Z} \right) \\ &= \sum_{k=0}^{\infty} F_k^2 + \tau^2 - \frac{1}{n} \sum_{k=0}^{\infty} F_k^2 \text{Tr} \left( (\mathbf{Z}^\top \mathbf{Z} + \psi_1 \psi_2 \lambda \mathbf{I}_N)^{-1} \mathbf{Z}^\top Q_k(\mathbf{X} \mathbf{X}^\top) \mathbf{Z} \right) \\ &\quad - \frac{\tau^2}{n} \text{Tr} \left( (\mathbf{Z}^\top \mathbf{Z} + \psi_1 \psi_2 \lambda \mathbf{I}_N)^{-1} \mathbf{Z}^\top \mathbf{Z} \right) \end{aligned}$$

It can be shown that the coefficients before  $F_0^2$  is asymptotically vanishing, and by Lemma B.9, we have  $\mathbb{E}[\sup_{k \geq 2} \|\mathbf{Q}_k(\mathbf{X} \mathbf{X}^\top) - \mathbf{I}_n\|_{\text{op}}^2] = o_d(1)$ . Hence we get

$$\begin{aligned} \mathbb{E}_{\boldsymbol{\beta}, \boldsymbol{\varepsilon}}[L_{\text{RF}}(f_d, \mathbf{X}, \boldsymbol{\Theta}, \lambda)] &= F_1^2 \left\{ 1 - \frac{1}{n} \text{Tr} \left( (\mathbf{Z}^\top \mathbf{Z} + \psi_1 \psi_2 \lambda \mathbf{I}_N)^{-1} \mathbf{Z}^\top \mathbf{H} \mathbf{Z} \right) \right\} \\ &\quad + (F_\star^2 + \tau^2) \cdot \left\{ 1 - \frac{1}{n} \text{Tr} \left( (\mathbf{Z}^\top \mathbf{Z} + \psi_1 \psi_2 \lambda \mathbf{I}_N)^{-1} \mathbf{Z}^\top \mathbf{Z} \right) \right\} + o_{d, \mathbb{P}}(1). \end{aligned}$$

Using the fact that

$$(\mathbf{Z}^\top \mathbf{Z} + \psi_1 \psi_2 \lambda \mathbf{I}_N)^{-1} \mathbf{Z}^\top = \mathbf{Z}^\top (\mathbf{Z} \mathbf{Z}^\top + \psi_1 \psi_2 \lambda \mathbf{I}_n)^{-1},$$

we have

$$\begin{aligned} &\mathbb{E}_{\boldsymbol{\beta}, \boldsymbol{\varepsilon}}[L_{\text{RF}}(f_d, \mathbf{X}, \boldsymbol{\Theta}, \lambda)] \\ &= F_1^2 \cdot \frac{\psi_1 \lambda}{d} \text{Tr} \left( (\mathbf{Z} \mathbf{Z}^\top + \psi_1 \psi_2 \lambda \mathbf{I}_n)^{-1} \mathbf{H} \right) + (F_\star^2 + \tau^2) \cdot \frac{\psi_1 \lambda}{d} \text{Tr} \left( (\mathbf{Z} \mathbf{Z}^\top + \psi_1 \psi_2 \lambda \mathbf{I}_n)^{-1} \right) + o_{d, \mathbb{P}}(1). \end{aligned}$$

### Step 2. The norm square of minimizers.

We have

$$\begin{aligned} \|\mathbf{a}\|_2^2 &= \|\mathbf{y}^\top \mathbf{Z}(\mathbf{Z}^\top \mathbf{Z} + \psi_1 \psi_2 \lambda \mathbf{I}_N)^{-1}\|_2^2 / d \\ &= \mathbf{y}^\top \mathbf{Z}(\mathbf{Z}^\top \mathbf{Z} + \psi_1 \psi_2 \lambda \mathbf{I}_N)^{-2} \mathbf{Z}^\top \mathbf{y} / d, \end{aligned}$$



so that

$$\begin{aligned}
\mathbb{E}_{\beta, \epsilon}[\|\mathbf{a}\|_2^2] &= \mathbb{E}_{\beta}[\mathbf{f}^\top \mathbf{Z}(\mathbf{Z}^\top \mathbf{Z} + \psi_1 \psi_2 \lambda \mathbf{I}_N)^{-2} \mathbf{Z}^\top \mathbf{f}]/d + \psi_1 \mathbb{E}_{\epsilon}[\epsilon^\top \mathbf{Z}(\mathbf{Z}^\top \mathbf{Z} + \psi_1 \psi_2 \lambda \mathbf{I}_N)^{-2} \mathbf{Z}^\top \epsilon]/d \\
&= \mathbb{E}_{\beta} \left[ \left( \sum_{k=0}^{\infty} \mathbf{Y}_{\mathbf{x}, k} \beta_k \right)^\top \mathbf{Z}(\mathbf{Z}^\top \mathbf{Z} + \psi_1 \psi_2 \lambda \mathbf{I}_N)^{-2} \mathbf{Z}^\top \left( \sum_{k=0}^{\infty} \mathbf{Y}_{\mathbf{x}, k} \beta_k \right) \right] / d \\
&\quad + \tau^2 \text{Tr} \left( (\mathbf{Z}^\top \mathbf{Z} + \psi_1 \psi_2 \lambda \mathbf{I}_N)^{-2} \mathbf{Z}^\top \mathbf{Z} \right) / d \\
&= \sum_{k=0}^{\infty} F_k^2 \cdot \text{Tr} \left( \mathbf{Z}^\top \mathbf{Z} + \psi_1 \psi_2 \lambda \mathbf{I}_N \right)^{-2} \mathbf{Z}^\top Q_k (\mathbf{X} \mathbf{X}^\top \mathbf{Z}) / d \\
&\quad + \tau^2 \text{Tr} \left( (\mathbf{Z}^\top \mathbf{Z} + \psi_1 \psi_2 \lambda \mathbf{I}_N)^{-2} \mathbf{Z}^\top \mathbf{Z} \right) / d \\
&= F_1^2 \text{Tr} \left( (\mathbf{Z}^\top \mathbf{Z} + \psi_1 \psi_2 \lambda \mathbf{I}_N)^{-2} \mathbf{Z}^\top \mathbf{H} \mathbf{Z} \right) / d + (F_\star^2 + \tau^2) \cdot \text{Tr} \left( (\mathbf{Z}^\top \mathbf{Z} + \psi_1 \psi_2 \lambda \mathbf{I}_N)^{-2} \mathbf{Z}^\top \mathbf{Z} \right) / d + o_{d, \mathbb{P}}(1).
\end{aligned}$$

**Step 3. The derivatives of the log determinant.**

Define  $\mathbf{q} = (s_1, s_2, t_1, t_2, p) \in \mathbb{R}^5$ , and introduce a block matrix  $\mathbf{A} \in \mathbb{R}^{M \times M}$  with  $M = N + n$ , defined by

$$\mathbf{A} = \begin{bmatrix} s_1 \mathbf{I}_N + s_2 \mathbf{Q} & \mathbf{Z}^\top + p \mathbf{Z}_1^\top \\ \mathbf{Z} + p \mathbf{Z}_1 & t_1 \mathbf{I}_n + t_2 \mathbf{H} \end{bmatrix}. \quad (227)$$

For any  $\xi \in \mathbb{C}_+$ , we consider the quantity

$$G_d(\xi; \mathbf{q}) = \frac{1}{d} \sum_{i=1}^M \log(\lambda_i(\mathbf{A}(\mathbf{q})) - \xi).$$

With simple algebra, we can show that

$$\begin{aligned}
\partial_{t_1} G_d(iu; \mathbf{0}) &= \frac{i u}{d} \text{Tr} \left( (u^2 \mathbf{I}_n + \mathbf{Z} \mathbf{Z}^\top)^{-1} \right), \\
\partial_{t_2} G_d(iu; \mathbf{0}) &= \frac{i u}{d} \text{Tr} \left( (u^2 \mathbf{I}_n + \mathbf{Z} \mathbf{Z}^\top)^{-1} \mathbf{H} \right), \\
\partial_{s_1, t_1}^2 G_d(iu; \mathbf{0}) &= -\frac{1}{d} \text{Tr} \left( (u^2 \mathbf{I}_N + \mathbf{Z}^\top \mathbf{Z})^{-2} \mathbf{Z}^\top \mathbf{Z} \right), \\
\partial_{s_1, t_2}^2 G_d(iu; \mathbf{0}) &= -\frac{1}{d} \text{Tr} \left( (u^2 \mathbf{I}_N + \mathbf{Z}^\top \mathbf{Z})^{-2} \mathbf{Z}^\top \mathbf{H} \mathbf{Z} \right).
\end{aligned} \quad (228)$$

Hence, we have

$$\begin{aligned}
\mathbb{E}[L_{\text{RF}}(f_d, \mathbf{X}, \boldsymbol{\Theta}, \lambda)] &= -F_1^2 \cdot i \left( \frac{\psi_1 \lambda}{\psi_2} \right)^{1/2} \partial_{t_2} \mathbb{E}[G_d(i(\lambda \psi_1 \psi_2)^{1/2}; \mathbf{0})] \\
&\quad - (F_\star^2 + \tau^2) \cdot i \left( \frac{\psi_1 \lambda}{\psi_2} \right)^{1/2} \partial_{t_1} \mathbb{E}[G_d(i(\lambda \psi_1 \psi_2)^{1/2}; \mathbf{0})] + o_d(1),
\end{aligned}$$

and

$$\mathbb{E}[\|\mathbf{a}\|_2^2] = -F_1^2 \partial_{s_1, t_2}^2 \mathbb{E}[G_d(i(\lambda \psi_1 \psi_2)^{1/2}; \mathbf{0})] - (F_\star^2 + \tau^2) \cdot \partial_{s_1, t_1}^2 \mathbb{E}[G_d(i(\lambda \psi_1 \psi_2)^{1/2}; \mathbf{0})] + o_d(1).$$

By Lemma E.3, we get

$$\begin{aligned}
\mathbb{E}[L_{\text{RF}}(f_d, \mathbf{X}, \boldsymbol{\Theta}, \lambda)] &= -F_1^2 \cdot i \left( \frac{\psi_1 \lambda}{\psi_2} \right)^{1/2} \partial_{t_2} g(i(\lambda \psi_1 \psi_2)^{1/2}; \mathbf{0}) \\
&\quad - (F_\star^2 + \tau^2) \cdot i \left( \frac{\psi_1 \lambda}{\psi_2} \right)^{1/2} \partial_{t_1} g(i(\lambda \psi_1 \psi_2)^{1/2}; \mathbf{0}) + o_d(1),
\end{aligned}$$

and

$$\mathbb{E}_{\beta, \epsilon}[\|\mathbf{a}\|_2^2] = -F_1^2 \partial_{s_1, t_2}^2 g(i(\lambda \psi_1 \psi_2)^{1/2}; \mathbf{0}) - (F_\star^2 + \tau^2) \cdot \partial_{s_1, t_1}^2 g(i(\lambda \psi_1 \psi_2)^{1/2}; \mathbf{0}) + o_{d, \mathbb{P}}(1),$$

where  $g$  is given in Eq. (67). The derivatives of  $g$  can be obtained by differentiating Eq. (66) and using Daskin's theorem. The theorem then follows by simple calculus.