

ONLINE SUPPLEMENTARY MATERIAL FOR  
**Multiple Descent in the Multiple Random Feature Model**

BY XURAN MENG, JIANFENG YAO AND YUAN CAO

The supplementary material is organized as follows:

- In Section I, we give the proof of Proposition A.2 which introduces the decomposition of the asymptotic excess risk.
- In Section II, we present the proof of Proposition A.4 by showing how the key terms in the risk decomposition are related to the logarithmic potential of the linear pencil matrix.
- In Section III, we establish basic properties of the fixed point equation (A.2) to justify the definition of  $\mathbf{m}(\xi; \mathbf{q}, \boldsymbol{\mu})$  below Definition A.5.
- In Section IV, we provide the proof of Proposition A.6, which extends the definition of  $m(\xi; \mathbf{q}, \boldsymbol{\mu})$  to  $\mathbb{C}_+$  and shows that it is the asymptotic limit of  $M_d(\xi; \mathbf{q}, \boldsymbol{\mu})$  as  $d \rightarrow \infty$ .
- In Section V, we give the proof of Proposition A.7 by relating the logarithmic potential  $G_d(\xi; \mathbf{q}, \boldsymbol{\mu})$  to  $g(\xi; \mathbf{q}, \boldsymbol{\mu})$  in (A.3).
- In Section VI, we prove Proposition A.8 to justify the definition of  $\nu(\xi)$  as the unique solution to the system (3.1).
- In Section VII, we display the proof of the lemmas and propositions given in Appendix C.

## I. Proof of Proposition A.2

Proposition A.2 gives a decomposition of the risk  $R_d(\mathbf{X}, \boldsymbol{\Theta}, \lambda, \beta_d, \varepsilon)$ . To prove this decomposition, we first introduce some additional notations and preliminary lemmas.

**Definition I.1.** *Define*

$$\mathbf{V}_0(F_0) = F_0 \mathbb{E}_{\mathbf{x}}[\boldsymbol{\sigma}(\mathbf{x})] \in \mathbb{R}^{N \times 1}, \quad \mathbf{V}(\beta_{1,d}) = \mathbb{E}_{\mathbf{x}}[\boldsymbol{\sigma}(\mathbf{x}) \mathbf{x}^\top \beta_{1,d}] \in \mathbb{R}^{N \times 1}, \quad \mathbf{U} = \mathbb{E}_{\mathbf{x}}[\boldsymbol{\sigma}(\mathbf{x}) \boldsymbol{\sigma}(\mathbf{x})^\top] \in \mathbb{R}^{N \times N},$$

where  $\mathbf{x}$  is a random vector uniformed distributed on the sphere  $\sqrt{d} \cdot \mathbb{S}^{d-1}$  and  $\boldsymbol{\sigma}(\mathbf{x})$  is defined in Definition A.1.  $\square$

Note that by the definition of  $\boldsymbol{\sigma}(\mathbf{x})$  in Definition A.1,  $\boldsymbol{\sigma}(\mathbf{x})$  also depends on the random feature parameter matrix  $\boldsymbol{\Theta}$ . Therefore,  $\mathbf{V}_0(F_0)$ ,  $\mathbf{V}(\beta_{1,d})$  and  $\mathbf{U}$  also depends on  $\boldsymbol{\Theta}$ . Now with these notations, and by the definition of  $\hat{\mathbf{a}}$  in (A.1), we can rewrite the risk as follows:

$$\begin{aligned} R_d(\mathbf{X}, \boldsymbol{\Theta}, \lambda, \beta_d, \varepsilon) &= \mathbb{E}_{\mathbf{x}}[\mathbf{x}^\top \beta_{1,d} + F_0 - \hat{\mathbf{a}}^\top \boldsymbol{\sigma}(\mathbf{x})]^2 \\ &= F_0^2 + F_{1,d}^2 - 2\mathbf{y}^\top \mathbf{Z} \mathbf{Y} [\mathbf{V}(\beta_{1,d}) + \mathbf{V}_0(F_0)] / \sqrt{d} + \mathbf{y}^\top [\mathbf{U}]_{\mathbf{Z}} \mathbf{y} / d. \end{aligned} \quad (\text{I.1})$$

Therefore, to prove Proposition A.2, it suffices to further decompose the terms  $\mathbf{U}$ ,  $\mathbf{V}(\beta_{1,d})$  and  $\mathbf{V}_0(F_0)$ . To handle these terms, we consider the Gegenbauer decomposition (Hua, 1963) of the nonlinear activation functions. For  $j = 1, 2$ , let  $\lambda_{d,k}(\sigma_j)$  be the coefficients of the Gegenbauer

decomposition of  $\sigma_j$ , i.e.,

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

where  $B(d,0) = 1$ ,  $B(d,k) = k^{-1}(2k+d-2) \binom{k+d-3}{k-1}$  with  $k \geq 1$ , and  $Q_k^{(d)}$ ,  $k \in \mathbb{N}$  are the Gegenbauer polynomials forms an orthogonal basis on  $L^2([-d,d], \tau_d)$ .  $\tau_d$  is the distribution of  $\langle \mathbf{x}_1, \mathbf{x}_2 \rangle$  where  $\mathbf{x}_1, \mathbf{x}_2 \sim \sqrt{d} \cdot \text{Unif}(\mathbb{S}^{d-1})$ . Then define

$$\mathbf{\Lambda}_{d,k} = \text{diag}(\lambda_{d,k}(\sigma_1) \mathbf{I}_{N_1}, \lambda_{d,k}(\sigma_2) \mathbf{I}_{N_2}), \quad k \in \mathbb{N} = \{0, 1, \dots\}. \quad (\text{I.2})$$

The following lemma decomposes the three terms in Definition I.1.

**Lemma I.2.** *With  $\mathbf{M}_1$  and  $\mathbf{M}_*$  in Definition A.1, and  $\mathbf{\Lambda}_{d,k}$  in equation (I.2), we have*

$$\begin{aligned} \mathbf{V}_0(F_0) &= F_0 \mathbf{\Lambda}_{d,0} \mathbf{1}_N, \\ \mathbf{V}(\beta_{1,d}) &= \mathbf{\Lambda}_{d,1} \mathbf{\Theta} \beta_{1,d} = \left( \frac{\mathbf{M}_1 + \mathbf{\Delta}'}{\sqrt{d}} \right) \mathbf{\Theta} \beta_{1,d}, \\ \mathbf{U} &= \mathbf{\Lambda}_{d,0} \mathbf{1}_N \mathbf{1}_N^\top \mathbf{\Lambda}_{d,0} + \mathbf{M}_1 \frac{\mathbf{\Theta} \mathbf{\Theta}^\top}{d} \mathbf{M}_1 + \mathbf{M}_* \mathbf{M}_* + \mathbf{\Delta}. \end{aligned}$$

where the remainder matrices  $\mathbf{\Delta}, \mathbf{\Delta}'$  satisfy  $\mathbb{E} \|\mathbf{\Delta}\|_{\text{op}}^2 \vee \mathbb{E} \|\mathbf{\Delta}'\|_{\text{op}}^2 = o_d(1)$ .

Lemma I.2 is proved in Section I.1. Plugging the decompositions in Lemma I.2 into (I.1) will then give a decomposition of the risk consisting of multiple terms. The next lemma establishes useful moment estimations for some of the terms in (I.1), which helps us get rid of the negligible terms in the decomposition.

**Lemma I.3.** *For any fixed  $k \in \mathbb{N} \setminus \{0\}$ , let  $\mathbf{\Gamma}_1 \in \mathbb{R}^{N \times N}$  and  $\mathbf{\Gamma}_2 \in \mathbb{R}^{n \times n}$  be symmetric random matrices with  $[\mathbb{E} \|\mathbf{\Gamma}_j\|_{\text{op}}^k]^{1/k} = O_d(1)$ ,  $j = 1, 2$ . Define*

$$\begin{aligned} \mathcal{B} &= \frac{1}{d} \mathbf{1}_n^\top [\mathbf{\Gamma}_1]_{\mathbf{Z}} \mathbf{1}_n, \\ \mathcal{C} &= 1 - \frac{2}{\sqrt{d}} \text{tr}(\mathbf{\Lambda}_{d,0} \mathbf{1}_N \mathbf{1}_n^\top \mathbf{Z} \mathbf{\Upsilon}) + \frac{1}{d} \mathbf{1}_n^\top [\mathbf{\Lambda}_{d,0} \mathbf{1}_N \mathbf{1}_N^\top \mathbf{\Lambda}_{d,0}]_{\mathbf{Z}} \mathbf{1}_n, \\ \mathcal{D} &= \frac{1}{d} \text{tr}([\mathbf{\Lambda}_{d,0} \mathbf{1}_N \mathbf{1}_N^\top \mathbf{\Lambda}_{d,0}]_{\mathbf{Z}} \mathbf{\Gamma}_2), \end{aligned}$$

where  $\mathbf{\Lambda}_{d,0}$  is defined in equation (I.2). Then if  $\sum_j \mu_{j,0}^2 > 0$ , for any fixed  $\lambda > 0$ , there exists a constant  $C > 0$  such that

$$(\mathbb{E} |\mathcal{B}|^k)^{1/k} \vee \mathbb{E} |\mathcal{C}| \vee (\mathbb{E} |\mathcal{D}|^k)^{1/k} = O_d(d^{-1} e^{C\sqrt{\log d}}) = o_d(1).$$

If  $\sum_j \mu_{j,0}^2 = 0$ , it still holds that  $(\mathbb{E} |\mathcal{D}|^k)^{1/k} = o_d(1)$ .

The proof of the lemma is given in Section I.2. To further decompose and calculate the risk, we also need to study the impact of fixed vector  $\beta_{1,d}$  on the risk. To do so, we aim to show that the risk only depends on  $F_{1,d}$  ( $= \|\beta_{1,d}\|_2$ ) due to rotation invariance of the learning problem. The result is given in the following lemma.

**Lemma I.4.** Suppose  $\tilde{\beta}_{1,d} \sim \text{Unif}(F_{1,d} \cdot \mathbb{S}^{d-1})$  is independent of  $(\mathbf{X}, \boldsymbol{\Theta}, \varepsilon)$ , and denote  $\tilde{\beta}_d = [F_0, \tilde{\beta}_{1,d}^\top]^\top$ . Then for any fixed  $\beta_{1,d}$ , under the assumptions of Proposition A.2, we have

$$\begin{aligned} \mathbb{E}_{\mathbf{X}, \boldsymbol{\Theta}, \varepsilon} |R_d(\mathbf{X}, \boldsymbol{\Theta}, \lambda, \beta_d, \varepsilon) - \bar{R}_d(\mathbf{X}, \boldsymbol{\Theta}, \lambda, F_{1,d}, \tau)| \\ = \mathbb{E}_{\mathbf{X}, \boldsymbol{\Theta}, \varepsilon, \tilde{\beta}_d} |R_d(\mathbf{X}, \boldsymbol{\Theta}, \lambda, \tilde{\beta}_d, \varepsilon) - \bar{R}_d(\mathbf{X}, \boldsymbol{\Theta}, \lambda, F_{1,d}, \tau)|, \\ \mathbb{E}_{\mathbf{X}, \boldsymbol{\Theta}} [\text{Var}_{\tilde{\beta}_d, \varepsilon}(R_d(\mathbf{X}, \boldsymbol{\Theta}, \lambda, \tilde{\beta}_d, \varepsilon))] = o_d(1). \end{aligned}$$

The proof of the lemma is given in Section I.3. Based on the above lemmas, we are ready to present the proof of Proposition A.2 as follows.

**Proof** [Proof of Proposition A.2] Let  $\tilde{\beta}_d = [F_0, \tilde{\beta}_{1,d}^\top]^\top$  with  $\tilde{\beta}_{1,d} \sim \text{Unif}(F_{1,d} \cdot \mathbb{S}^{d-1})$ . Then we have

$$\begin{aligned} \mathbb{E}_{\mathbf{X}, \boldsymbol{\Theta}, \varepsilon} |R_d(\mathbf{X}, \boldsymbol{\Theta}, \lambda, \beta_d, \varepsilon) - \bar{R}_d(\mathbf{X}, \boldsymbol{\Theta}, \lambda, F_{1,d}, \tau)| \\ = \mathbb{E}_{\mathbf{X}, \boldsymbol{\Theta}, \varepsilon, \tilde{\beta}_d} |R_d(\mathbf{X}, \boldsymbol{\Theta}, \lambda, \tilde{\beta}_d, \varepsilon) - \bar{R}_d(\mathbf{X}, \boldsymbol{\Theta}, \lambda, F_{1,d}, \tau)| \\ \leq \mathbb{E}_{\mathbf{X}, \boldsymbol{\Theta}, \varepsilon, \tilde{\beta}_d} |R_d(\mathbf{X}, \boldsymbol{\Theta}, \lambda, \tilde{\beta}_d, \varepsilon) - \mathbb{E}_{\varepsilon, \tilde{\beta}_d} R_d(\mathbf{X}, \boldsymbol{\Theta}, \lambda, \tilde{\beta}_d, \varepsilon)| \\ + \mathbb{E}_{\mathbf{X}, \boldsymbol{\Theta}} |\mathbb{E}_{\varepsilon, \tilde{\beta}_d} R_d(\mathbf{X}, \boldsymbol{\Theta}, \lambda, \tilde{\beta}_d, \varepsilon) - \bar{R}_d(\mathbf{X}, \boldsymbol{\Theta}, \lambda, F_{1,d}, \tau)| \\ \leq \mathbb{E}_{\mathbf{X}, \boldsymbol{\Theta}} \left[ \sqrt{\text{Var}_{\tilde{\beta}_d, \varepsilon}(R_d(\mathbf{X}, \boldsymbol{\Theta}, \lambda, \tilde{\beta}_d, \varepsilon))} \right] \\ + \mathbb{E}_{\mathbf{X}, \boldsymbol{\Theta}} |\mathbb{E}_{\varepsilon, \tilde{\beta}_d} R_d(\mathbf{X}, \boldsymbol{\Theta}, \lambda, \tilde{\beta}_d, \varepsilon) - \bar{R}_d(\mathbf{X}, \boldsymbol{\Theta}, \lambda, F_{1,d}, \tau)| \\ \leq \sqrt{\mathbb{E}_{\mathbf{X}, \boldsymbol{\Theta}} [\text{Var}_{\tilde{\beta}_d, \varepsilon}(R_d(\mathbf{X}, \boldsymbol{\Theta}, \lambda, \tilde{\beta}_d, \varepsilon))]} \\ + \mathbb{E}_{\mathbf{X}, \boldsymbol{\Theta}} |\mathbb{E}_{\varepsilon, \tilde{\beta}_d} R_d(\mathbf{X}, \boldsymbol{\Theta}, \lambda, \tilde{\beta}_d, \varepsilon) - \bar{R}_d(\mathbf{X}, \boldsymbol{\Theta}, \lambda, F_{1,d}, \tau)| \\ = o_d(1) + \mathbb{E}_{\mathbf{X}, \boldsymbol{\Theta}} |\mathbb{E}_{\varepsilon, \tilde{\beta}_d} R_d(\mathbf{X}, \boldsymbol{\Theta}, \lambda, \tilde{\beta}_d, \varepsilon) - \bar{R}_d(\mathbf{X}, \boldsymbol{\Theta}, \lambda, F_{1,d}, \tau)|, \end{aligned}$$

where the first equality follows by Lemma I.4, the first inequality follows by triangle inequality, the second and third inequalities are by Jensen's inequality, and the last equality follows by Lemma I.4 again. Therefore, to prove the proposition, it suffices to show that

$$\mathbb{E}_{\mathbf{X}, \boldsymbol{\Theta}} |\mathbb{E}_{\varepsilon, \tilde{\beta}_d} R_d(\mathbf{X}, \boldsymbol{\Theta}, \lambda, \tilde{\beta}_d, \varepsilon) - \bar{R}_d(\mathbf{X}, \boldsymbol{\Theta}, \lambda, F_{1,d}, \tau)| = o_d(1).$$

Similar to (I.1), we have

$$R_d(\mathbf{X}, \boldsymbol{\Theta}, \lambda, \tilde{\beta}_d, \varepsilon) = F_0^2 + F_{1,d}^2 - \frac{2\tilde{\mathbf{y}}^\top \mathbf{Z} \boldsymbol{\Upsilon} (\mathbf{V}(\tilde{\beta}_{1,d}) + \mathbf{V}_0(F_0))}{\sqrt{d}} + \frac{\tilde{\mathbf{y}}^\top [\mathbf{U}] \mathbf{Z} \tilde{\mathbf{y}}}{d}, \quad (\text{I.3})$$

where  $\tilde{\mathbf{y}} = \mathbf{1}_n F_0 + \mathbf{X} \tilde{\beta}_{1,d} + \varepsilon$ . From Lemma I.2, we further have

$$\mathbf{V}_0(F_0) F_0 = F_0^2 \boldsymbol{\Lambda}_{d,0} \mathbf{1}_N, \quad \mathbb{E}_{\tilde{\beta}_{1,d}} (\mathbf{V}(\tilde{\beta}_{1,d}) \tilde{\beta}_{1,d}^\top) = F_{1,d}^2 \left( \frac{\mathbf{M}_1 + \boldsymbol{\Delta}'}{\sqrt{d}} \right) \frac{\boldsymbol{\Theta}}{d}, \quad (\text{I.4})$$

and

$$\mathbf{U} = \boldsymbol{\Lambda}_{d,0} \mathbf{1}_N \mathbf{1}_N^\top \boldsymbol{\Lambda}_{d,0} + \mathbf{M}_1 \frac{\boldsymbol{\Theta} \boldsymbol{\Theta}^\top}{d} \mathbf{M}_1 + \mathbf{M}_* \mathbf{M}_* + \boldsymbol{\Delta}. \quad (\text{I.5})$$

By (I.3), (I.4), (I.5) and the definition of  $\bar{R}_d(\mathbf{X}, \boldsymbol{\Theta}, \lambda, F_{1,d}, \tau)$ , we obtain the following equation with direct calculation:

$$\begin{aligned} \mathbb{E}_{\tilde{\beta}_d, \varepsilon} R_d(\mathbf{X}, \boldsymbol{\Theta}, \lambda, \tilde{\beta}_d, \varepsilon) - \bar{R}_d(\mathbf{X}, \boldsymbol{\Theta}, \lambda, F_{1,d}, \tau) = \\ \underbrace{F_0^2 - \frac{2F_0^2}{\sqrt{d}} \text{tr}(\boldsymbol{\Lambda}_{d,0} \mathbf{1}_N \mathbf{1}_n^\top \mathbf{Z} \boldsymbol{\Upsilon}) + \frac{F_0^2}{d} \text{tr}([\boldsymbol{\Lambda}_{d,0} \mathbf{1}_N \mathbf{1}_N^\top \boldsymbol{\Lambda}_{d,0}]_{\mathbf{Z}} \mathbf{1}_n \mathbf{1}_n^\top)}_{I_1} \\ - \underbrace{\frac{2F_{1,d}^2}{d} \text{tr}\left(\boldsymbol{\Delta}' \frac{\boldsymbol{\Theta} \mathbf{X}^\top}{d} \mathbf{Z} \boldsymbol{\Upsilon}\right)}_{I_2} + \underbrace{\frac{F_0^2}{d} \text{tr}\left(\left[\mathbf{M}_1 \frac{\boldsymbol{\Theta} \boldsymbol{\Theta}^\top}{d} \mathbf{M}_1 + \mathbf{M}_* \mathbf{M}_*\right]_{\mathbf{Z}} \mathbf{1}_n \mathbf{1}_n^\top\right)}_{I_3} \\ + \underbrace{\frac{F_{1,d}^2}{d} \text{tr}\left([\boldsymbol{\Lambda}_{d,0} \mathbf{1}_N \mathbf{1}_N^\top \boldsymbol{\Lambda}_{d,0}]_{\mathbf{Z}} \frac{\mathbf{X} \mathbf{X}^\top}{d}\right)}_{I_4} + \underbrace{\frac{\tau^2}{d} \text{tr}([\boldsymbol{\Lambda}_{d,0} \mathbf{1}_N \mathbf{1}_N^\top \boldsymbol{\Lambda}_{d,0}]_{\mathbf{Z}})}_{I_5} \\ + \underbrace{\frac{F_{1,d}^2}{d} \text{tr}\left([\boldsymbol{\Delta}]_{\mathbf{Z}} \frac{\mathbf{X} \mathbf{X}^\top}{d}\right)}_{I_6} + \underbrace{\frac{F_{1,d}^2}{d} \text{tr}[\boldsymbol{\Delta}]_{\mathbf{Z}}}_{I_7} + \underbrace{\frac{F_0^2}{d} \text{tr}([\boldsymbol{\Delta}]_{\mathbf{Z}} \mathbf{1}_n \mathbf{1}_n^\top)}_{I_8}. \end{aligned}$$

We now show that all the terms  $I_1, \dots, I_8$  on the right hand side above are negligible terms. We note that by definition,  $\|\mathbf{Z} \boldsymbol{\Upsilon}\|_{\text{op}} = \|\mathbf{Z}(\mathbf{Z}^\top \mathbf{Z} + \lambda \mathbf{I})^{-1}\|_{\text{op}} \leq 1/(2\sqrt{\lambda})$  is deterministically bounded. Therefore we have

$$\mathbb{E}|I_2| \leq 2F_{1,d}^2 \cdot \mathbb{E} \left\| \left( \boldsymbol{\Delta}' \frac{\boldsymbol{\Theta} \mathbf{X}^\top}{d} \mathbf{Z} \boldsymbol{\Upsilon} \right) \right\|_{\text{op}} \leq O_d\left(\frac{1}{2\sqrt{\lambda}}\right) \cdot (\mathbb{E}\|\boldsymbol{\Delta}'\|_{\text{op}}^2)^{\frac{1}{2}} \cdot \left(\mathbb{E}\left\|\frac{\boldsymbol{\Theta} \mathbf{X}^\top}{d}\right\|_{\text{op}}^2\right)^{\frac{1}{2}} = o_d(1),$$

where the last equality follows by  $\mathbb{E}\|\boldsymbol{\Delta}'\|_{\text{op}}^2 = o_d(1)$  in Lemma I.2. Moreover, by definition, we have

$$\|[\boldsymbol{\Delta}]_{\mathbf{Z}}\|_{\text{op}} = \|\mathbf{Z} \boldsymbol{\Upsilon} \boldsymbol{\Delta} (\mathbf{Z} \boldsymbol{\Upsilon})^\top\|_{\text{op}} \leq \frac{1}{4\lambda} \|\boldsymbol{\Delta}\|_{\text{op}}.$$

Therefore, by Lemma I.2 that  $\mathbb{E}\|\boldsymbol{\Delta}\|_{\text{op}}^2 = o_d(1)$ , we have

$$\begin{aligned} \mathbb{E}|I_6| &\leq F_{1,d}^2 \cdot \mathbb{E} \left\| [\boldsymbol{\Delta}]_{\mathbf{Z}} \frac{\mathbf{X} \mathbf{X}^\top}{d} \right\|_{\text{op}} \leq F_{1,d}^2 \cdot \mathbb{E} \left[ \|[\boldsymbol{\Delta}]_{\mathbf{Z}}\|_{\text{op}} \cdot \left\| \frac{\mathbf{X} \mathbf{X}^\top}{d} \right\|_{\text{op}} \right] = O_d(\mathbb{E}\|\boldsymbol{\Delta}\|_{\text{op}}^2) = o_d(1), \\ \mathbb{E}|I_7| &\leq F_{1,d}^2 \cdot \mathbb{E}\|[\boldsymbol{\Delta}]_{\mathbf{Z}}\|_{\text{op}} \leq \frac{F_{1,d}^2}{4\lambda} \cdot \mathbb{E}\|\boldsymbol{\Delta}\|_{\text{op}} = o_d(1). \end{aligned}$$

For the remaining terms, we discuss them according to the value of  $F_0$ . When  $F_0 = 0$ , it is clear that  $I_1 = I_3 = 0$ . Note that under this situation, the condition  $\sum_j \mu_{j,0}^2 > 0$  in Lemma I.3 may not hold. If  $\sum_j \mu_{j,0}^2 = 0$ , from Lemma I.3, it still holds that

$$\mathbb{E}|I_4| = o_d(1), \quad \mathbb{E}|I_5| = o_d(1).$$

Therefore, when  $F_0 = 0$ , Proposition A.2 holds.

When  $F_0 \neq 0$ ,  $\sum_j \mu_{j,0}^2 > 0$  holds from Assumption 3.4, the result for  $\mathcal{C}$  in Lemma I.3 gives the bound for  $I_1$ , the result for  $\mathcal{B}$  in Lemma I.3 gives the bounds on  $I_3$  and  $I_8$ , and the result for  $\mathcal{D}$  in

Lemma I.3 gives the bounds on  $I_4$  and  $I_5$ . Therefore we have

$$\mathbb{E}_{\mathbf{X}, \Theta} |\mathbb{E}_{\tilde{\beta}_d, \varepsilon} R_d(\mathbf{X}, \Theta, \lambda, \tilde{\beta}_d, \varepsilon) - \bar{R}_d(\mathbf{X}, \Theta, \lambda, F_{1,d}, \tau)| = o_d(1),$$

which proves Proposition A.2. ■

### I.1 Proof of Lemma I.2

The proof of Lemma I.2 is mainly based on the decomposition of the nonlinear activation function. We first present several classical lemmas about Gegenbauer polynomials and their relation to Hermite polynomials. The following lemma can be found in Mei and Montanari (2022) (see Lemma 9.4 and its proof in the reference).

**Lemma I.5.** *Let  $Q_k^{(d)}(\cdot)$ ,  $k \in \mathbb{N}$  be the Gegenbauer polynomials. The following properties hold:*

1. *For  $\mathbf{v}_1, \mathbf{v}_2 \in \sqrt{d} \cdot \mathbb{S}^{d-1}$ , suppose  $\mathbf{x} \sim \text{Unif}(\sqrt{d} \cdot \mathbb{S}^{d-1})$ , then for  $k, l \in \mathbb{N}$ ,*

$$\mathbb{E}_{\mathbf{x}} [Q_k^{(d)}(\mathbf{v}_1^\top \mathbf{x}) Q_l^{(d)}(\mathbf{x}^\top \mathbf{v}_2)] = \frac{\delta_{kl}}{B(d, k)} \cdot Q_k^{(d)}(\mathbf{v}_1^\top \mathbf{v}_2),$$

where  $\delta_{kl} = 1$  if  $k = l$  and  $\delta_{kl} = 0$  if  $k \neq l$ .

2. *For  $\Theta_1$  and  $\Theta_2$  defined in Section 2,  $Q_k^{(d)}(\cdot)$  the point wise function on matrices, the following equality holds:*

$$\begin{aligned} \mathbb{E} \left[ \sup_{k \geq 2} \|Q_k^{(d)}(\Theta_j \Theta_j^\top) - \mathbf{I}_{N_j}\|_{\text{op}}^2 \right] &= o_d(1), \quad j = 1, 2, \\ \mathbb{E} \left[ \sup_{k \geq 2} \|Q_k^{(d)}(\Theta_1 \Theta_2^\top)\|_{\text{op}}^2 \right] &= o_d(1). \end{aligned}$$

The next lemma gives the connection between the coefficients in Hermite polynomials  $H_k$  and the coefficients in Gegenbauer polynomials  $Q_k^{(d)}$ .

**Lemma I.6.** *Let  $Q_k^{(d)}(\cdot)$ ,  $H_k(\cdot)$ ,  $k \in \mathbb{N}$  be the Gegenbauer and Hermite polynomials respectively. For  $j = 1, 2$ , suppose that  $\sigma_j(x)$  has Gegenbauer decomposition*

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

and Hermite polynomial decomposition

$$\sigma_j(x) = \sum_{k=0}^{+\infty} \alpha_k(\sigma_j) / k! \cdot H_k(x).$$

Then for each  $k \in \mathbb{N}$ ,  $\lambda_{d,k}^2(\sigma_j) B(d, k) k! \rightarrow \alpha_k^2(\sigma_j)$  as  $d \rightarrow +\infty$ .

The proof of Lemma I.6 can be found in Appendix A.3 in Mei and Montanari (2022). Note that the orthogonality of the standard Hermite polynomials ( $H_1(x) = x$ ) implies that for  $G \sim N(0, 1)$ ,

$$\mathbb{E}[H_k(G) H_l(G)] = \delta_{kl} \cdot k!.$$

Based on this property, let  $\alpha_k(\sigma_j)$  be defined in Lemma I.6. Then for  $j = 1, 2$ , we have

$$\alpha_k(\sigma_j) = \mu_{j,k}, \quad k = 0, 1, \quad \mu_{j,*}^2 = \sum_{k \geq 2} \frac{\alpha_k^2(\sigma_j)}{k!},$$

where the constants  $\mu_{j,k}$  and  $\mu_{j,*}$  are defined in Definition 3.1. Therefore, by Lemma I.6, we further have

$$\sum_{k \geq 2} \lambda_{d,k}^2(\sigma_j) B(d, k) \rightarrow \mu_{j,*}^2. \quad (\text{I.6})$$

Recall that  $\boldsymbol{\sigma}(\mathbf{x}) = (\sigma_1(\mathbf{x}^\top \boldsymbol{\Theta}_1^\top / \sqrt{d}), \sigma_2(\mathbf{x}^\top \boldsymbol{\Theta}_2^\top / \sqrt{d}))^\top$ . Moreover, note that the zeroth order Gegenbauer polynomial  $Q_0^d(x) = 1$ . Therefore by Lemma I.5 and the Gegenbauer decomposition of  $\boldsymbol{\sigma}_j$  in Lemma I.6, we have

$$\mathbf{V}_0(F_0) = F_0 \mathbb{E}_{\mathbf{x}}[\boldsymbol{\sigma}(\mathbf{x}) \cdot Q_0^d(\mathbf{x}^\top \mathbf{1}_d)] = \frac{F_0}{B(d, 0)} \cdot \boldsymbol{\Lambda}_{d,0} \cdot Q_0^d(\boldsymbol{\Theta} \mathbf{1}_d) \cdot B(d, 0) = F_0 \boldsymbol{\Lambda}_{d,0} \mathbf{1}_N.$$

Here, the equality holds from the fact that  $Q_0^d(\mathbf{x}^\top \mathbf{1}_d) = 1$  and  $Q_0^d(\boldsymbol{\Theta} \mathbf{1}_d) = \mathbf{1}_N$ . Similarly,  $Q_1^d(x) = x/d$  holds. Again from Lemma I.5 and Lemma I.6, we have

$$\begin{aligned} \mathbf{V}(\boldsymbol{\beta}_{1,d}) &= \mathbb{E}_{\mathbf{x}} \boldsymbol{\sigma}(\mathbf{x}) \mathbf{x}^\top \boldsymbol{\beta}_{1,d} = d \cdot \mathbb{E}_{\mathbf{x}} \boldsymbol{\sigma}(\mathbf{x}) Q_1^d(\mathbf{x}^\top \boldsymbol{\beta}_{1,d}) = \frac{d}{B(d, 1)} \cdot \boldsymbol{\Lambda}_{d,1} \cdot Q_1^d(\boldsymbol{\Theta} \boldsymbol{\beta}_{1,d}) \cdot B(d, 1) \\ &= \boldsymbol{\Lambda}_{d,1} \boldsymbol{\Theta} \boldsymbol{\beta}_{1,d} = \left( \frac{\mathbf{M}_1 + \boldsymbol{\Delta}'}{\sqrt{d}} \right) \boldsymbol{\Theta} \boldsymbol{\beta}_{1,d}. \end{aligned}$$

Here,  $\boldsymbol{\Delta}' = \sqrt{d} \cdot \boldsymbol{\Lambda}_{d,1} - \mathbf{M}_1$ . From Lemma I.6, set  $k = 1$  and we have  $\sqrt{d} \lambda_{d,1}(\sigma_j) \rightarrow \mu_{j,1}$ . Thus  $\boldsymbol{\Delta}'$  satisfies  $\mathbb{E} \|\boldsymbol{\Delta}'\|_{\text{op}}^2 = o_d(1)$ . As for  $\mathbf{U} = \mathbb{E}_{\mathbf{x}}[\boldsymbol{\sigma}(\mathbf{x}) \boldsymbol{\sigma}(\mathbf{x})^\top]$ ,  $\mathbf{U}$  could be divided into the following block matrix:

$$\mathbf{U} = \begin{bmatrix} \mathbf{U}_{1,1} & \mathbf{U}_{1,2} \\ \mathbf{U}_{2,1} & \mathbf{U}_{2,2} \end{bmatrix},$$

where

$$\mathbf{U}_{i,j} = \mathbb{E}_{\mathbf{x}}[\sigma_i(\boldsymbol{\Theta}_i \mathbf{x} / \sqrt{d}) \sigma_j(\mathbf{x}^\top \boldsymbol{\Theta}_j^\top / \sqrt{d})], \quad i, j = 1, 2.$$

Now by Lemma I.5, we have

$$\mathbf{U}_{i,j} = \sum_{k=0}^{+\infty} \lambda_{d,k}(\sigma_i) \lambda_{d,k}(\sigma_j) B(d, k) Q_k^{(d)}(\boldsymbol{\Theta}_i \boldsymbol{\Theta}_j^\top), \quad i, j = 1, 2. \quad (\text{I.7})$$

Note that  $Q_0^{(d)}(x) = 1$ ,  $Q_1^{(d)}(x) = x/d$ , so that the first two terms in the decomposition (I.7) have a simple form. For  $k \geq 2$ , we approximate the terms using the approximation given in the second item of Lemma I.5. Consider first  $\mathbf{U}_{1,1}$ . We have

$$\mathbf{U}_{1,1} = \lambda_{d,0}^2(\sigma_1) \cdot \mathbf{1}_{N_1} \mathbf{1}_{N_1}^\top + \lambda_{d,1}^2(\sigma_1) \cdot B(d, 1) \cdot \frac{\boldsymbol{\Theta}_1 \boldsymbol{\Theta}_1^\top}{d} + \sum_{k=2}^{+\infty} \lambda_{d,k}^2(\sigma_1) \cdot B(d, k) \cdot Q_k^{(d)}(\boldsymbol{\Theta}_1 \boldsymbol{\Theta}_1^\top)$$

$$\begin{aligned}
 &= \lambda_{d,0}^2(\sigma_1) \cdot \mathbf{1}_{N_1} \mathbf{1}_{N_1}^\top + \lambda_{d,1}^2(\sigma_1) \cdot B(d, 1) \cdot \frac{\boldsymbol{\Theta}_1 \boldsymbol{\Theta}_1^\top}{d} + \sum_{k=2}^{+\infty} \lambda_{d,k}^2(\sigma_1) \cdot B(d, k) \cdot \mathbf{I}_{N_1} \\
 &\quad + \sum_{k=2}^{+\infty} \lambda_{d,k}^2(\sigma_1) \cdot B(d, k) \cdot [Q_k^{(d)}(\boldsymbol{\Theta}_1 \boldsymbol{\Theta}_1^\top) - \mathbf{I}_{N_1}],
 \end{aligned} \tag{I.8}$$

where we have used the fact that  $\sum_{k=2}^{+\infty} \lambda_{d,k}^2(\sigma_1) B(d, k) < +\infty$  for sufficiently large  $d$ , which is implied by (I.6). Moreover, by Lemma I.5, the convergence of this series also implies that

$$\mathbb{E} \left\| \sum_{k=2}^{+\infty} \lambda_{d,k}^2(\sigma_1) \cdot B(d, k) \cdot [Q_k^{(d)}(\boldsymbol{\Theta}_1 \boldsymbol{\Theta}_1^\top) - \mathbf{I}_{N_1}] \right\|_{\text{op}}^2 = o_d(1). \tag{I.9}$$

Therefore by (I.8) and (I.9), we have

$$\mathbb{E} \left\| \mathbf{U}_{1,1} - \lambda_{d,0}^2(\sigma_1) \cdot \mathbf{1}_{N_1} \mathbf{1}_{N_1}^\top - \lambda_{d,1}^2(\sigma_1) \cdot B(d, 1) \cdot \frac{\boldsymbol{\Theta}_1 \boldsymbol{\Theta}_1^\top}{d} - \sum_{k=2}^{+\infty} \lambda_{d,k}^2(\sigma_1) \cdot B(d, k) \cdot \mathbf{I}_{N_1} \right\|_{\text{op}}^2 = o_d(1). \tag{I.10}$$

Now by Lemma I.6 and equations (I.6), (I.10), we have

$$\mathbb{E} \left\| \mathbf{U}_{1,1} - \lambda_{d,0}^2(\sigma_1) \cdot \mathbf{1}_{N_1} \mathbf{1}_{N_1}^\top - \mu_{1,1}^2 \cdot \frac{\boldsymbol{\Theta}_1 \boldsymbol{\Theta}_1^\top}{d} - \mu_{1,*}^2 \cdot \mathbf{I}_{N_1} \right\|_{\text{op}}^2 = o_d(1).$$

This establishes the approximation for  $\mathbf{U}_{1,1}$ .

For the other sub-matrices  $\mathbf{U}_{1,2}$ ,  $\mathbf{U}_{2,1}$  and  $\mathbf{U}_{2,2}$ , the derivations are exactly the same, and we obtain the following results:

$$\begin{aligned}
 &\mathbb{E} \left\| \mathbf{U}_{1,2} - \lambda_{d,0}(\sigma_1) \lambda_{d,0}(\sigma_2) \mathbf{1}_{N_1} \mathbf{1}_{N_2}^\top - \mu_{1,1} \mu_{2,1} \cdot \frac{\boldsymbol{\Theta}_1 \boldsymbol{\Theta}_2^\top}{d} \right\|_{\text{op}}^2 = o_d(1), \\
 &\mathbb{E} \left\| \mathbf{U}_{2,1} - \lambda_{d,0}(\sigma_1) \lambda_{d,0}(\sigma_2) \mathbf{1}_{N_1} \mathbf{1}_{N_2}^\top - \mu_{1,1} \mu_{2,1} \cdot \frac{\boldsymbol{\Theta}_2 \boldsymbol{\Theta}_1^\top}{d} \right\|_{\text{op}}^2 = o_d(1), \\
 &\mathbb{E} \left\| \mathbf{U}_{2,2} - \lambda_{d,0}^2(\sigma_2) \mathbf{1}_{N_2} \mathbf{1}_{N_2}^\top - \mu_{2,1}^2 \frac{\boldsymbol{\Theta}_2 \boldsymbol{\Theta}_2^\top}{d} - \mu_{2,*}^2 \cdot \mathbf{I}_{N_2} \right\|_{\text{op}}^2 = o_d(1).
 \end{aligned}$$

Note that the collection of the approximations for the four blocks  $\mathbf{U}_{1,1}$ ,  $\mathbf{U}_{1,2}$ ,  $\mathbf{U}_{2,1}$  and  $\mathbf{U}_{2,2}$  gives the matrix

$$\boldsymbol{\Lambda}_{d,0} \mathbf{1}_N \mathbf{1}_N^\top \boldsymbol{\Lambda}_{d,0} + \mathbf{M}_1 \frac{\boldsymbol{\Theta} \boldsymbol{\Theta}^\top}{d} \mathbf{M}_1 + \mathbf{M}_* \mathbf{M}_*,$$

so finally we have

$$\mathbb{E} \left\| \mathbf{U} - \boldsymbol{\Lambda}_{d,0} \mathbf{1}_N \mathbf{1}_N^\top \boldsymbol{\Lambda}_{d,0} - \mathbf{M}_1 \frac{\boldsymbol{\Theta} \boldsymbol{\Theta}^\top}{d} \mathbf{M}_1 - \mathbf{M}_* \mathbf{M}_* \right\|_{\text{op}}^2 = o_d(1).$$

The proof of Lemma I.2 is complete.

## I.2 Proof of Lemma I.3

We first prove that  $(\mathbb{E}|\mathcal{D}|^k)^{1/k} = o_d(1)$  if  $\sum_j \mu_{j,0}^2 = 0$ . Note that the rank-1 matrix  $\mathbf{A}$  satisfies  $|\text{tr}(\mathbf{A})| = \|\mathbf{A}\|_{\text{op}}$ . Moreover,  $\sum_j \mu_{j,0}^2 = 0$  implies  $\|\Lambda_{d,0}\|_{\text{op}} = o_d(1)$ . We have  $(\mathbb{E}|\mathcal{D}|^k)^{1/k} = O_d(\|\Lambda_{d,0} \frac{1_N \mathbf{1}_N^\top}{d} \Lambda_{d,0}\|_{\text{op}}) \cdot (\mathbb{E}\|\Gamma_2\|_{\text{op}}^k)^{1/k} = o_d(1) \cdot O_d(1) = o_d(1)$ .

In the following proof of Lemma I.3, we have the condition  $\sum_j \mu_{j,0}^2 > 0$ . We separate the proof into two parts, estimating  $\mathcal{B}$  and  $\mathcal{C}$ , and  $\mathcal{D}$ , respectively.

### I.2.1 ESTIMATION FOR $\mathcal{B}$ AND $\mathcal{C}$

Let

$$L_1 = \frac{1}{\sqrt{d}} \text{tr}(\Lambda_{d,0} \mathbf{1}_N \mathbf{1}_n^\top \mathbf{Z} \Upsilon), \quad L_2(\Gamma) = \frac{1}{d} \text{tr}([\Gamma]_{\mathbf{Z}} \mathbf{1}_n \mathbf{1}_n^\top) = \frac{1}{d} \text{tr}(\mathbf{Z} \Upsilon \Gamma \Upsilon \mathbf{Z}^\top \mathbf{1}_n \mathbf{1}_n^\top),$$

where  $\Gamma \in \mathbb{R}^{N \times N}$  is a symmetric matrix. Then we have

$$\mathcal{B} = L_2(\Gamma), \quad \mathcal{C} = 1 - 2L_1 + L_2(\Lambda_{d,0} \mathbf{1}_N \mathbf{1}_N^\top \Lambda_{d,0}).$$

Define further the following terms:

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

where

$$\begin{aligned} \mathbf{J} &= \mathbf{Z} - \mathbf{1}_n \mathbf{1}_N^\top \Lambda_{d,0} / \sqrt{d}, & \mathbf{E}_0 &= \mathbf{J}^\top \mathbf{J} + \lambda \mathbf{I}_N, \\ \mathbf{T}_1 &= \psi_3^{1/2} \Lambda_{d,0} \mathbf{1}_N, & \mathbf{T}_2 &= \frac{1}{\sqrt{n}} \mathbf{J}^\top \mathbf{1}_n. \end{aligned}$$

We denote  $\psi_3 = n/d$  for notation simplification. The proof is organized in two steps:

1. Express  $\mathcal{B}$  and  $\mathcal{C}$  in function of  $K_{ij}$  and  $G_{ij}$ ,  $i, j \in \{1, 2\}$ .
2. Estimate the order of  $K_{ij}$  and  $G_{ij}$ , and show that  $\mathbb{E}|\mathcal{B}|$  and  $\mathbb{E}|\mathcal{C}|$  are both  $o_d(1)$ .

Denote  $\mathbf{F}_1 = [\mathbf{T}_1, \mathbf{T}_1, \mathbf{T}_2] \in \mathbb{R}^{N \times 3}$ ,  $\mathbf{F}_2 = [\mathbf{T}_1, \mathbf{T}_2, \mathbf{T}_1] \in \mathbb{R}^{N \times 3}$ , it is easy to see

$$\begin{aligned} \Upsilon &= (\mathbf{Z}^\top \mathbf{Z} + \lambda \mathbf{I}_N)^{-1} = ((\mathbf{J} + \mathbf{1}_n \mathbf{1}_N^\top \Lambda_{d,0})^\top (\mathbf{J} + \mathbf{1}_n \mathbf{1}_N^\top \Lambda_{d,0}) + \lambda \mathbf{I}_N)^{-1} \\ &= (\psi_3 \Lambda_{d,0} \mathbf{1}_N \mathbf{1}_N^\top + \psi_3^{1/2} \Lambda_{d,0} \mathbf{1}_N \mathbf{T}_2^\top + \psi_3^{1/2} \mathbf{T}_2 \mathbf{1}_N^\top \Lambda_{d,0} + \mathbf{J}^\top \mathbf{J} + \lambda \mathbf{I}_N)^{-1} \\ &= (\mathbf{E}_0 + \mathbf{F}_1 \mathbf{F}_2^\top)^{-1}. \end{aligned}$$

For  $L_1$ , replacing  $\mathbf{Z}$  by  $\mathbf{J} + \mathbf{1}_n \mathbf{1}_N^\top \Lambda_{d,0} / \sqrt{d}$ , we have

$$\begin{aligned} L_1 &= \text{tr}[(\psi_3 \Lambda_{d,0} \mathbf{1}_N \mathbf{1}_N^\top + \psi_3^{1/2} \Lambda_{d,0} \mathbf{1}_N \mathbf{T}_2^\top) \cdot (\mathbf{E}_0 + \mathbf{F}_1 \mathbf{F}_2^\top)^{-1}] \\ &= \text{tr}[(\mathbf{T}_1 \mathbf{T}_1^\top + \mathbf{T}_1 \mathbf{T}_2^\top) \cdot (\mathbf{E}_0 + \mathbf{F}_1 \mathbf{F}_2^\top)^{-1}]. \end{aligned} \tag{I.11}$$

By the Sherman-Morrison-Woodbury formula,

$$\Upsilon = (\mathbf{E}_0 + \mathbf{F}_1 \mathbf{F}_2^\top)^{-1} = \mathbf{E}_0^{-1} - \mathbf{E}_0^{-1} \mathbf{F}_1 (\mathbf{I}_3 + \mathbf{F}_2^\top \mathbf{E}_0^{-1} \mathbf{F}_1)^{-1} \mathbf{F}_2^\top \mathbf{E}_0^{-1}. \tag{I.12}$$



Plugging (I.12) into (I.11), we have

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

where

$$\mathbf{K} = \mathbf{F}_2^\top \mathbf{E}_0^{-1} \mathbf{F}_1 = \begin{bmatrix} K_{11} & K_{11} & K_{12} \\ K_{12} & K_{12} & K_{22} \\ K_{11} & K_{11} & K_{12} \end{bmatrix}.$$

Thus by simple calculation,

$$L_1 = 1 - \frac{K_{12} + 1}{K_{11}(1 - K_{22}) + (K_{12} + 1)^2}. \quad (\text{I.13})$$

As for  $L_2(\mathbf{\Gamma})$ , we have

$$\begin{aligned}
 \mathbf{Z}^\top \mathbf{1}_n \mathbf{1}_n^\top \mathbf{Z} / d &= (\mathbf{J} + \mathbf{1}_n \mathbf{1}_N^\top \mathbf{\Lambda}_{d,0} / \sqrt{d})^\top \mathbf{1}_n \mathbf{1}_n^\top (\mathbf{J} + \mathbf{1}_n \mathbf{1}_N^\top \mathbf{\Lambda}_{d,0} / \sqrt{d}) / d \\
 &= \psi_3 (\psi_3^{1/2} \mathbf{\Lambda}_{d,0} \mathbf{1}_N + \frac{1}{\sqrt{n}} \mathbf{J}^\top \mathbf{1}_n) (\psi_3^{1/2} \mathbf{\Lambda}_{d,0} \mathbf{1}_N + \frac{1}{\sqrt{n}} \mathbf{J}^\top \mathbf{1}_n)^\top \\
 &= \psi_3 (\mathbf{T}_1 + \mathbf{T}_2) (\mathbf{T}_1 + \mathbf{T}_2)^\top.
 \end{aligned}$$

Then after similar calculation by (I.12).

$$\begin{aligned}
 \mathcal{B} = L_2(\mathbf{\Gamma}) &= \frac{1}{d} \text{tr}(\mathbf{Z}^\top \mathbf{1}_n \mathbf{1}_n^\top \mathbf{Z} \mathbf{\Upsilon} \mathbf{\Gamma} \mathbf{\Upsilon}) = \text{tr}(\psi_3 (\mathbf{T}_1 + \mathbf{T}_2) (\mathbf{T}_1 + \mathbf{T}_2)^\top \mathbf{\Upsilon} \mathbf{\Gamma} \mathbf{\Upsilon}) \\
 &= \psi_3 (\mathbf{T}_1 + \mathbf{T}_2)^\top (\mathbf{E}_0 + \mathbf{F}_1 \mathbf{F}_2^\top)^{-1} \mathbf{\Gamma} (\mathbf{E}_0 + \mathbf{F}_1 \mathbf{F}_2^\top)^{-1} (\mathbf{T}_1 + \mathbf{T}_2) \\
 &= \psi_3 (\mathbf{T}_1 + \mathbf{T}_2)^\top (\mathbf{E}_0^{-1} - \mathbf{E}_0^{-1} \mathbf{F}_1 (\mathbf{I}_3 + \mathbf{F}_2^\top \mathbf{E}_0^{-1} \mathbf{F}_1)^{-1} \mathbf{F}_2^\top \mathbf{E}_0^{-1}) \\
 &\quad \cdot \mathbf{\Gamma} (\mathbf{E}_0^{-1} - \mathbf{E}_0^{-1} \mathbf{F}_1 (\mathbf{I}_3 + \mathbf{F}_2^\top \mathbf{E}_0^{-1} \mathbf{F}_1)^{-1} \mathbf{F}_2^\top \mathbf{E}_0^{-1}) (\mathbf{T}_1 + \mathbf{T}_2) \\
 &= \psi_3 \frac{G_{11}(1 - K_{22})^2 + G_{22}(K_{12} + 1)^2 + 2G_{12}(K_{12} + 1)(1 - K_{22})}{(K_{11}(1 - K_{22}) + (K_{12} + 1)^2)^2}. \quad (\text{I.14})
 \end{aligned}$$

When  $\mathbf{\Gamma} = \mathbf{\Lambda}_{d,0} \mathbf{1}_N \mathbf{1}_N^\top \mathbf{\Lambda}_{d,0}$ , the  $G_{11}$ ,  $G_{12}$  and  $G_{22}$  above can be given as

$$G_{11} = K_{11}^2 / \psi_3, \quad G_{12} = K_{11} K_{12} / \psi_3, \quad G_{22} = K_{12}^2 / \psi_3.$$

Then by (I.13), (I.14), we have

$$\mathcal{C} = 1 - 2L_1 + L_2(\mathbf{\Lambda}_{d,0} \mathbf{1}_N \mathbf{1}_N^\top \mathbf{\Lambda}_{d,0}) = \frac{(K_{12} + 1)^2}{(K_{11}(1 - K_{22}) + (K_{12} + 1)^2)^2}. \quad (\text{I.15})$$

We next estimate the order for  $K_{11}$ ,  $K_{12}$ ,  $K_{22}$ ,  $G_{11}$ ,  $G_{12}$  and  $G_{22}$  respectively. By the inequality  $\|[\mathbf{A} \ \mathbf{B}]\|_{\text{op}} \leq \|\mathbf{A}\|_{\text{op}} + \|\mathbf{B}\|_{\text{op}}$  for any matrices  $\mathbf{A}$  and  $\mathbf{B}$ , we have

$$\begin{aligned} \|\mathbf{J}\|_{\text{op}} &\leq \|\mathbf{Z}_1 - \lambda_{d,0}(\sigma_1)\mathbf{1}_n\mathbf{1}_{N_1}^\top/\sqrt{d}\|_{\text{op}} + \|\mathbf{Z}_2 - \lambda_{d,0}(\sigma_2)\mathbf{1}_n\mathbf{1}_{N_2}^\top/\sqrt{d}\|_{\text{op}} \\ &= O_{\mathbb{P}}(\exp(C\sqrt{\log d})), \end{aligned} \quad (\text{I.16})$$

where the last equality in (I.16) follows by Lemma C.5 in (Mei and Montanari, 2022). Moreover, for any fixed  $\lambda > 0$ , it also deterministically holds that

$$\|(\mathbf{J}^\top \mathbf{J} + \lambda \mathbf{I}_N)^{-1} \mathbf{J}^\top\|_{\text{op}} \leq 2/\sqrt{\lambda}, \quad \|(\mathbf{J}^\top \mathbf{J} + \lambda \mathbf{I}_N)^{-1}\|_{\text{op}} \leq 1/\lambda.$$

Now recall that

$$\begin{aligned} K_{11} &= \psi_3 \mathbf{1}_N^\top \mathbf{A}_{d,0} (\mathbf{J}^\top \mathbf{J} + \lambda \mathbf{I}_N)^{-1} \mathbf{A}_{d,0} \mathbf{1}_N, \\ K_{12} &= \mathbf{1}_N^\top \mathbf{A}_{d,0} (\mathbf{J}^\top \mathbf{J} + \lambda \mathbf{I}_N)^{-1} \mathbf{J}^\top \mathbf{1}_n / \sqrt{d}, \\ K_{22} &= \mathbf{1}_n^\top \mathbf{J} (\mathbf{J}^\top \mathbf{J} + \lambda \mathbf{I}_N)^{-1} \mathbf{J}^\top \mathbf{1}_n / n, \\ G_{11} &= \psi_3 \mathbf{1}_N^\top \mathbf{A}_{d,0} (\mathbf{J}^\top \mathbf{J} + \lambda \mathbf{I}_N)^{-1} \mathbf{\Gamma} (\mathbf{J}^\top \mathbf{J} + \lambda \mathbf{I}_N)^{-1} \mathbf{A}_{d,0} \mathbf{1}_N, \\ G_{12} &= \mathbf{1}_N^\top \mathbf{A}_{d,0} (\mathbf{J}^\top \mathbf{J} + \lambda \mathbf{I}_N)^{-1} \mathbf{\Gamma} (\mathbf{J}^\top \mathbf{J} + \lambda \mathbf{I}_N)^{-1} \mathbf{J}^\top \mathbf{1}_n / \sqrt{d}, \\ G_{22} &= \mathbf{1}_n^\top \mathbf{J} (\mathbf{J}^\top \mathbf{J} + \lambda \mathbf{I}_N)^{-1} \mathbf{\Gamma} (\mathbf{J}^\top \mathbf{J} + \lambda \mathbf{I}_N)^{-1} \mathbf{J}^\top \mathbf{1}_n / n. \end{aligned}$$

Therefore we deterministically have

$$|K_{12}| \leq \|(\mathbf{J}^\top \mathbf{J} + \lambda \mathbf{I}_N)^{-1} \mathbf{J}^\top\|_{\text{op}} \|\mathbf{1}_n \mathbf{1}_N^\top \mathbf{A}_{d,0} / \sqrt{d}\|_{\text{op}} = O_d(\sqrt{d/\lambda}). \quad (\text{I.17})$$

For  $K_{22}$ , by its definition, it is clear that  $K_{22} > 0$ . Moreover, we have

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

Therefore we have

$$0 < K_{22} \leq 1 - \frac{\lambda}{\|\mathbf{J}^\top \mathbf{J}\|_{\text{op}} + \lambda}. \quad (\text{I.18})$$

For  $K_{11}$ , the condition  $\mu_{1,0}^2 + \mu_{2,0}^2 > 0$  ensures that there exists  $j \in \{1, 2\}$  such that  $\mu_{j,0}^2 > 0$ . By Lemma I.6 (note that  $B(d, 0) = 1$ ), we have  $\lambda_{d,0}^2(\sigma_j) \rightarrow \mu_{j,0}^2$  as  $d \rightarrow +\infty$ . Therefore for large enough  $d$ , we have  $\lambda_{d,0}(\sigma_j) > \mu_{j,0}/2 > 0$ , and

$$\begin{aligned} K_{11} &\geq \psi_3 \mathbf{1}_N^\top \mathbf{A}_{d,0}^2 \mathbf{1}_N \lambda_{\min}((\mathbf{J}^\top \mathbf{J} + \lambda \mathbf{I}_N)^{-1}) \\ &\geq \psi_3 \cdot (\mu_{j,0}^2/4) \cdot N_j \cdot \lambda_{\min}((\mathbf{J}^\top \mathbf{J} + \lambda \mathbf{I}_N)^{-1}) \\ &= \frac{\Omega_d(d)}{\|\mathbf{J}^\top \mathbf{J}\|_{\text{op}} + \lambda}. \end{aligned} \quad (\text{I.19})$$

Plugging (I.17), (I.18), (I.19) into (I.14) then gives

$$|\mathcal{B}| = \frac{|G_{22}(1 + K_{12})^2 + G_{11}(1 - K_{22})^2 + 2G_{12}(1 + K_{12})(1 - K_{22})|}{[(1 + K_{12})^2 + K_{11} \cdot (1 - K_{22})]^2}$$

$$\begin{aligned}
 &\leq \frac{|G_{22}(1 + K_{12})^2 + G_{11}(1 - K_{22})^2 + 2G_{12}(1 + K_{12})(1 - K_{22})|}{[K_{11} \cdot (1 - K_{22})]^2} \\
 &\leq O_d(1) \cdot \frac{|G_{22}| \cdot d + |G_{12}| \cdot \sqrt{d} + |G_{11}|}{d^2/(\lambda + \|\mathbf{J}\mathbf{J}^\top\|_{\text{op}})^4},
 \end{aligned}$$

where we utilize the upper and lower bounds in (I.17), (I.18), (I.19) to obtain the last inequality. For  $G_{11}, G_{12}$  and  $G_{22}$ , we have

$$\begin{aligned}
 \mathbb{E}[|G_{11}|^k]^{1/k} &\leq \psi_3 \|(\mathbf{J}^\top \mathbf{J} + \lambda \mathbf{I}_N)^{-1}\|_{\text{op}} [\mathbb{E}\|\mathbf{\Gamma}\|_{\text{op}}^k]^{1/k} \|(\mathbf{J}^\top \mathbf{J} + \lambda \mathbf{I}_N)^{-1}\|_{\text{op}} \|\mathbf{\Lambda}_{d,0} \mathbf{1}_N \mathbf{1}_N^\top \mathbf{\Lambda}_{d,0}\|_{\text{op}} = O_d(d), \\
 \mathbb{E}[|G_{12}|^k]^{1/k} &\leq \|(\mathbf{J}^\top \mathbf{J} + \lambda \mathbf{I}_N)^{-1}\|_{\text{op}} [\mathbb{E}\|\mathbf{\Gamma}\|_{\text{op}}^k]^{1/k} \|(\mathbf{J}^\top \mathbf{J} + \lambda \mathbf{I}_N)^{-1} \mathbf{J}^\top\|_{\text{op}} \|\mathbf{1}_n \mathbf{1}_N^\top \mathbf{\Lambda}_{d,0} / \sqrt{d}\|_{\text{op}} = O_d(\sqrt{d}), \\
 \mathbb{E}[|G_{22}|^k]^{1/k} &\leq \|(\mathbf{J}^\top \mathbf{J} + \lambda \mathbf{I}_N)^{-1}\|_{\text{op}} [\mathbb{E}\|\mathbf{\Gamma}\|_{\text{op}}^k]^{1/k} \|(\mathbf{J}^\top \mathbf{J} + \lambda \mathbf{I}_N)^{-1} \mathbf{J}^\top \mathbf{J}\|_{\text{op}} \text{tr}(\mathbf{1}_n \mathbf{1}_n^\top / n) = O_d(1).
 \end{aligned}$$

Thus by the bounds above and the triangle inequality of the  $L_k$ -norm  $\mathbb{E}[|\cdot|^{1/k}]$ , we have

$$\begin{aligned}
 (\mathbb{E}|\mathcal{B}|^k)^{1/k} &\leq O_d(1) \cdot \frac{\mathbb{E}[|G_{22}|^k]^{1/k} \cdot d + \mathbb{E}[|G_{12}|^k]^{1/k} \cdot \sqrt{d} + \mathbb{E}[|G_{11}|^k]^{1/k}}{d^2/(\lambda + \|\mathbf{J}\mathbf{J}^\top\|_{\text{op}})^4} \\
 &= O_d(1) \cdot \frac{d}{d^2/(\lambda + \|\mathbf{J}\mathbf{J}^\top\|_{\text{op}})^4} = O_d\left(\frac{(\lambda + \|\mathbf{J}\mathbf{J}^\top\|_{\text{op}})^4}{d}\right) = O_d\left(\frac{\exp(C\sqrt{\log d})}{d}\right) \\
 &= o_d(1),
 \end{aligned}$$

and

$$\mathbb{E}|\mathcal{C}| = O_d\left(\frac{(\lambda + \|\mathbf{J}\mathbf{J}^\top\|_{\text{op}})^2}{d}\right) = O_d\left(\frac{\exp(C\sqrt{\log d})}{d}\right) = o_d(1).$$

This completes the proof.

### I.2.2 ESTIMATION FOR $\mathcal{D}$

The proof is similar to the calculations for  $\mathcal{B}$  and  $\mathcal{C}$  in the previous section. Also we use a set of similar notations as previously which may however have slightly different values. Let

$$\begin{aligned}
 \mathbf{J} &= \mathbf{Z} - \mathbf{1}_n \mathbf{1}_N^\top \mathbf{\Lambda}_{d,0} / \sqrt{d}, & \mathbf{E}_0 &= \mathbf{J}\mathbf{J}^\top + \lambda \mathbf{I}_n, \\
 \mathbf{T}_1 &= \mathbf{J} \mathbf{\Lambda}_{d,0} \mathbf{1}_N / \sqrt{d}, & \mathbf{T}_2 &= \mathbf{1}_n, \\
 K_{11} &= \mathbf{T}_1^\top \mathbf{E}_0^{-1} \mathbf{T}_1, & K_{12} &= \mathbf{T}_1^\top \mathbf{E}_0^{-1} \mathbf{T}_2, & K_{22} &= \mathbf{T}_2^\top \mathbf{E}_0^{-1} \mathbf{T}_2, \\
 G_{11} &= \mathbf{T}_1^\top \mathbf{E}_0^{-1} \mathbf{\Gamma} \mathbf{E}_0^{-1} \mathbf{T}_1, & G_{12} &= \mathbf{T}_1^\top \mathbf{E}_0^{-1} \mathbf{\Gamma} \mathbf{E}_0^{-1} \mathbf{T}_2, & G_{22} &= \mathbf{T}_2^\top \mathbf{E}_0^{-1} \mathbf{\Gamma} \mathbf{E}_0^{-1} \mathbf{T}_2,
 \end{aligned}$$

where  $\mathbf{\Gamma} \in \mathbb{R}^{n \times n}$  is a symmetric matrix. We express  $\mathcal{D}$  with the terms defined above. Recall that  $\mathbf{\Upsilon} = (\mathbf{Z}^\top \mathbf{Z} + \lambda \mathbf{I}_N)^{-1}$  and further define  $\mathbf{\Xi} = (\mathbf{Z}\mathbf{Z}^\top + \lambda \mathbf{I}_n)^{-1}$ . Clearly,  $\mathbf{Z}\mathbf{\Upsilon} = \mathbf{\Xi}\mathbf{Z}$ . Therefore we have

$$\mathcal{D} = \frac{1}{d} \text{tr}(\mathbf{Z}\mathbf{\Upsilon} \mathbf{\Lambda}_{d,0} \mathbf{1}_N \mathbf{1}_N^\top \mathbf{\Lambda}_{d,0} \mathbf{\Upsilon} \mathbf{Z}^\top \mathbf{\Gamma}) = \frac{1}{d} \text{tr}(\mathbf{\Xi} \mathbf{Z} \mathbf{\Lambda}_{d,0} \mathbf{1}_N \mathbf{1}_N^\top \mathbf{\Lambda}_{d,0} \mathbf{Z}^\top \mathbf{\Xi} \mathbf{\Gamma}). \quad (\text{I.20})$$

We proceed to calculate  $\mathbf{\Xi}$  and  $\mathbf{Z} \mathbf{\Lambda}_{d,0} \mathbf{1}_N \mathbf{1}_N^\top \mathbf{\Lambda}_{d,0} \mathbf{Z}^\top$ , respectively. Define  $c = \mathbf{1}_N^\top \mathbf{\Lambda}_{d,0}^2 \mathbf{1}_N / d = \Theta(1)$ ,  $\mathbf{F}_1 = [\mathbf{T}_1, \mathbf{T}_2, \mathbf{T}_2] \in \mathbb{R}^{n \times 3}$ ,  $\mathbf{F}_2 = [\mathbf{T}_2, \mathbf{T}_1, c\mathbf{T}_2] \in \mathbb{R}^{n \times 3}$ . Then we have

$$\mathbf{\Xi} = \left( (\mathbf{J} + \mathbf{1}_n \mathbf{1}_N^\top \mathbf{\Lambda}_{d,0} / \sqrt{d}) (\mathbf{J} + \mathbf{1}_n \mathbf{1}_N^\top \mathbf{\Lambda}_{d,0} / \sqrt{d})^\top + \lambda \mathbf{I}_n \right)^{-1}$$

$$= (\mathbf{E}_0 + \mathbf{F}_1 \mathbf{F}_2^\top)^{-1} = \mathbf{E}_0^{-1} - \mathbf{E}_0^{-1} \mathbf{F}_1 (\mathbf{I}_3 + \mathbf{F}_2^\top \mathbf{E}_0^{-1} \mathbf{F}_1)^{-1} \mathbf{F}_2^\top \mathbf{E}_0^{-1}, \quad (\text{I.21})$$

where the last equality follows from the Sherman-Morrison-Woodbury formula. Moreover, we have

$$\begin{aligned} \mathbf{Z} \mathbf{\Lambda}_{d,0} \mathbf{1}_N \mathbf{1}_N^\top \mathbf{\Lambda}_{d,0} \mathbf{Z}^\top &= \left( \mathbf{J} + \mathbf{1}_n \mathbf{1}_N^\top \mathbf{\Lambda}_{d,0} / \sqrt{d} \right) \mathbf{\Lambda}_{d,0} \mathbf{1}_N \mathbf{1}_N^\top \mathbf{\Lambda}_{d,0} \left( \mathbf{J} + \mathbf{1}_n \mathbf{1}_N^\top \mathbf{\Lambda}_{d,0} / \sqrt{d} \right)^\top \\ &= d (\mathbf{T}_1 \mathbf{T}_1^\top + c (\mathbf{T}_2 \mathbf{T}_1^\top + \mathbf{T}_1 \mathbf{T}_2^\top) + c^2 \mathbf{T}_2 \mathbf{T}_2^\top) \\ &= d \cdot (\mathbf{T}_1 + c \mathbf{T}_2) (\mathbf{T}_1 + c \mathbf{T}_2)^\top. \end{aligned} \quad (\text{I.22})$$

Plugging (I.21) and (I.22) into (I.20), we obtain

$$\begin{aligned} \mathcal{D} &= \text{tr} \left( (\mathbf{T}_1 + c \mathbf{T}_2)^\top (\mathbf{E}_0^{-1} - \mathbf{E}_0^{-1} \mathbf{F}_1 (\mathbf{I}_3 + \mathbf{F}_2^\top \mathbf{E}_0^{-1} \mathbf{F}_1)^{-1} \mathbf{F}_2^\top \mathbf{E}_0^{-1}) \right. \\ &\quad \left. \cdot \mathbf{\Gamma} (\mathbf{E}_0^{-1} - \mathbf{E}_0^{-1} \mathbf{F}_1 (\mathbf{I}_3 + \mathbf{F}_2^\top \mathbf{E}_0^{-1} \mathbf{F}_1)^{-1} \mathbf{F}_2^\top \mathbf{E}_0^{-1}) (\mathbf{T}_1 + c \mathbf{T}_2) \right). \end{aligned}$$

With similar calculation as in the proof of Lemma I.3, we obtain that

$$\mathcal{D} = \frac{G_{11}(1 + K_{12})^2 + G_{22}(c - K_{11})^2 + 2G_{12}(1 + K_{12})(c - K_{11})}{(1 + 2K_{12} + K_{12}^2 + cK_{22} - K_{11}K_{22})^2}. \quad (\text{I.23})$$

We then estimate the order for  $K_{11}$ ,  $K_{12}$ ,  $K_{22}$ ,  $G_{11}$ ,  $G_{12}$  and  $G_{22}$ , respectively. For  $K_{11}$ , apparently we have  $K_{11} > 0$ . Moreover,

$$\begin{aligned} c - K_{11} &= \frac{1}{d} \mathbf{1}_N^\top \mathbf{\Lambda}_{d,0} (\mathbf{I}_N - \mathbf{J}^\top (\mathbf{J} \mathbf{J}^\top + \lambda \mathbf{I}_n)^{-1} \mathbf{J}) \mathbf{\Lambda}_{d,0} \mathbf{1}_N \\ &\geq c (1 - \lambda_{\max}(\mathbf{J}^\top (\mathbf{J} \mathbf{J}^\top + \lambda \mathbf{I}_n)^{-1} \mathbf{J})) = \frac{c\lambda}{\lambda + \|\mathbf{J} \mathbf{J}^\top\|_{\text{op}}} > 0. \end{aligned}$$

Therefore we have

$$c \geq c - K_{11} \geq \frac{c\lambda}{\lambda + \|\mathbf{J} \mathbf{J}^\top\|_{\text{op}}} > 0. \quad (\text{I.24})$$

Similarly, for  $K_{12}$  and  $K_{22}$  we have

$$|K_{12}| \leq \|(\mathbf{J} \mathbf{J}^\top + \lambda \mathbf{I}_n)^{-1} \mathbf{J}^\top\|_{\text{op}} \|\mathbf{1}_n \mathbf{1}_N^\top \mathbf{\Lambda}_{d,0} / \sqrt{d}\|_{\text{op}} = O_d(\sqrt{d/\lambda}), \quad (\text{I.25})$$

$$K_{22} \geq n \lambda_{\min}((\mathbf{J} \mathbf{J}^\top + \lambda \mathbf{I}_n)^{-1}) = \Omega(d) / (\|\mathbf{J} \mathbf{J}^\top\|_{\text{op}} + \lambda). \quad (\text{I.26})$$

Plugging (I.24), (I.25), (I.26) into (I.23) then gives

$$\begin{aligned} |\mathcal{D}| &= \frac{|G_{11}(1 + K_{12})^2 + G_{22}(c - K_{11})^2 + 2G_{12}(1 + K_{12})(c - K_{11})|}{[(1 + K_{12})^2 + K_{22} \cdot (c - K_{11})]^2} \\ &\leq \frac{|G_{11}(1 + K_{12})^2 + G_{22}(c - K_{11})^2 + 2G_{12}(1 + K_{12})(c - K_{11})|}{[K_{22} \cdot (c - K_{11})]^2} \\ &\leq O_d(1) \cdot \frac{|G_{11}| \cdot d + |G_{12}| \cdot \sqrt{d} + |G_{22}| \cdot c^2}{d^2 / (\lambda + \|\mathbf{J} \mathbf{J}^\top\|_{\text{op}})^4}, \end{aligned}$$

where we utilize the upper and lower bounds in (I.24), (I.25), (I.26) to obtain the last inequality. Now recall that

$$\begin{aligned} G_{11} &= \mathbf{1}_N^\top \mathbf{A}_{d,0} \mathbf{J}^\top (\mathbf{J} \mathbf{J}^\top + \lambda \mathbf{I}_n)^{-1} \mathbf{\Gamma} (\mathbf{J} \mathbf{J}^\top + \lambda \mathbf{I}_n)^{-1} \mathbf{J} \mathbf{A}_{d,0} \mathbf{1}_N / d, \\ G_{12} &= \mathbf{1}_N^\top \mathbf{A}_{d,0} \mathbf{J}^\top (\mathbf{J} \mathbf{J}^\top + \lambda \mathbf{I}_n)^{-1} \mathbf{\Gamma} (\mathbf{J} \mathbf{J}^\top + \lambda \mathbf{I}_n)^{-1} \mathbf{1}_n / \sqrt{d}, \\ G_{22} &= \mathbf{1}_n (\mathbf{J} \mathbf{J}^\top + \lambda \mathbf{I}_n)^{-1} \mathbf{\Gamma} (\mathbf{J} \mathbf{J}^\top + \lambda \mathbf{I}_n)^{-1} \mathbf{1}_n. \end{aligned}$$

Therefore we have

$$\mathbb{E}[|G_{11}|^k]^{1/k} \leq c \cdot O_d(\lambda^{-1}) \cdot [\mathbb{E}\|\mathbf{\Gamma}\|_{\text{op}}^k]^{1/k} = O_d(1), \quad (\text{I.27})$$

$$\mathbb{E}[|G_{12}|^k]^{1/k} \leq O_d(\lambda^{-3/2}) \cdot [\mathbb{E}\|\mathbf{\Gamma}\|_{\text{op}}^k]^{1/k} \cdot \|\mathbf{1}_n \mathbf{1}_N^\top \mathbf{A}_{d,0} / \sqrt{d}\|_{\text{op}} = O_d(\sqrt{d}), \quad (\text{I.28})$$

$$\mathbb{E}[|G_{22}|^k]^{1/k} \leq O_d(\lambda^{-2}) \cdot [\mathbb{E}\|\mathbf{\Gamma}\|_{\text{op}}^k]^{1/k} \cdot \|\mathbf{1}_n \mathbf{1}_n^\top\|_{\text{op}} = O_d(d). \quad (\text{I.29})$$

By the triangle inequality of the  $L_k$ -norm  $\mathbb{E}[|\cdot|^k]^{1/k}$ , we have

$$\begin{aligned} (\mathbb{E}|\mathcal{D}|^k)^{1/k} &\leq O_d(1) \cdot \frac{\mathbb{E}[|G_{11}|^k]^{1/k} \cdot d + \mathbb{E}[|G_{12}|^k]^{1/k} \cdot \sqrt{d} + \mathbb{E}[|G_{22}|^k]^{1/k} \cdot c^2}{d^2 / (\lambda + \|\mathbf{J} \mathbf{J}^\top\|_{\text{op}})^4} \\ &= O_d(1) \cdot \frac{d}{d^2 / (\lambda + \|\mathbf{J} \mathbf{J}^\top\|_{\text{op}})^4} \\ &= O_d\left(\frac{(\lambda + \|\mathbf{J} \mathbf{J}^\top\|_{\text{op}})^4}{d}\right) = O_d\left(\frac{\exp(C\sqrt{\log d})}{d}\right) = o_d(1), \end{aligned}$$

where the first equality follows by (I.27), (I.28) and (I.29). This completes the proof.

### I.3 Proof of Lemma I.4

The first result follows by the rotation invariance of the learning problem. For any  $\beta_d = [F_0, \beta_{1,d}^\top]^\top$  and  $\tilde{\beta}_d = [F_0, \tilde{\beta}_{1,d}^\top]^\top$  with  $\beta_{1,d}, \tilde{\beta}_{1,d} \in F_{1,d} \cdot \mathbb{S}^{d-1}$ , there exists an orthogonal matrix  $\mathbf{P}$  such that  $\mathbf{P}\beta_{1,d} = \tilde{\beta}_{1,d}$ . Then by definition, we have

$$R_d(\mathbf{X}\mathbf{P}, \mathbf{\Theta}\mathbf{P}, \lambda, \beta_d, \varepsilon) = R_d(\mathbf{X}, \mathbf{\Theta}, \lambda, \tilde{\beta}_d, \varepsilon).$$

Moreover, it is easy to check that

$$\bar{R}_d(\mathbf{X}\mathbf{P}, \mathbf{\Theta}\mathbf{P}, \lambda, F_{1,d}, \tau) = \bar{R}_d(\mathbf{X}, \mathbf{\Theta}, \lambda, F_{1,d}, \tau).$$

Since  $(\mathbf{X}\mathbf{P}, \mathbf{\Theta}\mathbf{P}) \stackrel{\text{d}}{=} (\mathbf{X}, \mathbf{\Theta})$ , we see that conditional to  $\beta_d, \tilde{\beta}_d$ , we have

$$R_d(\mathbf{X}, \mathbf{\Theta}, \lambda, \beta_d, \varepsilon) - \bar{R}_d(\mathbf{X}, \mathbf{\Theta}, \lambda, F_{1,d}, \tau) \stackrel{\text{d}}{=} R_d(\mathbf{X}, \mathbf{\Theta}, \lambda, \tilde{\beta}_d, \varepsilon) - \bar{R}_d(\mathbf{X}, \mathbf{\Theta}, \lambda, F_{1,d}, \tau).$$

This implies the first result in Lemma I.4.

If we assume that  $\tilde{\beta}_{1,d} \sim \mathcal{N}(\mathbf{0}, [F_{1,d}^2/d] \mathbf{I}_d)$ , then  $F_{1,d} \cdot \tilde{\beta}_{1,d} / \|\tilde{\beta}_{1,d}\|_2 \sim F_{1,d} \cdot \text{Unif}(\mathbb{S}^{d-1})$ . The proof of the second result in Lemma I.4 from Gaussian  $\tilde{\beta}_{1,d}$  to spherical  $\tilde{\beta}_{1,d}$  differs by the factor  $\|\tilde{\beta}_{1,d}\|_2 / F_{1,d}$ . Note that in high dimensions, the norm of Gaussian  $\tilde{\beta}_{1,d}$  ( $\|\tilde{\beta}_{1,d}\|_2$ ) is tightly concentrated on  $F_{1,d}$ . Therefore, it is not hard to translate the proof from Gaussian version to spherical version.

Based on the analysis above, without loss of generality we could assume  $\tilde{\beta}_{1,d} \sim N(\mathbf{0}, [F_{1,d}^2/d]\mathbf{I}_d)$  in the following of the proof. The lemma below helps us further handle the quadratic form of the variance which appears later.

**Lemma I.7.** *Let  $\mathbf{A} \in \mathbb{R}^{d \times d}$ , and define the random vector  $\mathbf{h} \sim N(0, (F_{1,d}^2/d)\mathbf{I}_d)$ . Then we have*

$$\text{Var}_{\mathbf{h}}(\mathbf{h}^\top \mathbf{A} \mathbf{h}) = \frac{F_{1,d}^4}{d^2} (\|\mathbf{A}\|_F^2 + \text{tr}(\mathbf{A}^2)).$$

The proof of Lemma I.7 is given at the end of this section. With this lemma, we are well-prepared to prove the second result in Lemma I.4. Recall the definitions

$$\boldsymbol{\sigma}(\mathbf{x}) = (\sigma_1(\mathbf{x}^\top \boldsymbol{\Theta}_1^\top / \sqrt{d}), \sigma_2(\mathbf{x}^\top \boldsymbol{\Theta}_2^\top / \sqrt{d}))^\top \in \mathbb{R}^N, \quad \boldsymbol{\Upsilon} = (\mathbf{Z}^\top \mathbf{Z} + \lambda \mathbf{I}_N)^{-1},$$

and  $\tilde{\mathbf{V}} = \mathbb{E}_{\mathbf{x}} \boldsymbol{\sigma}(\mathbf{x})(\mathbf{x}^\top \tilde{\beta}_{1,d} + F_0)$ ,  $\mathbf{U} = \mathbb{E}_{\mathbf{x}} \boldsymbol{\sigma}(\mathbf{x}) \boldsymbol{\sigma}(\mathbf{x})^\top$ . By the definition of the risk  $R_d(\mathbf{X}, \boldsymbol{\Theta}, \lambda, \tilde{\beta}_{1,d}, \boldsymbol{\varepsilon})$ , we have

$$\begin{aligned} R_d(\mathbf{X}, \boldsymbol{\Theta}, \lambda, \tilde{\beta}_{1,d}, \boldsymbol{\varepsilon}) &= \mathbb{E}_{\mathbf{x}} (\mathbf{x}^\top \tilde{\beta}_d + F_0 - \hat{\mathbf{a}}(\lambda)^\top \boldsymbol{\sigma}(\mathbf{x}))^2 \\ &= F_0^2 + F_{1,d}^2 - 2\Gamma_1 + \Gamma_2 + \Gamma_3 - 2\Gamma_4 + 2\Gamma_5, \end{aligned} \quad (\text{I.30})$$

where

$$\begin{aligned} \mathbf{f} &= \mathbf{X} \tilde{\beta}_{1,d} + \mathbf{1}_n F_0, & \Gamma_1 &= \mathbf{f}^\top \mathbf{Z} \boldsymbol{\Upsilon} \tilde{\mathbf{V}} / \sqrt{d}, & \Gamma_2 &= \mathbf{f}^\top [\mathbf{U}]_{\mathbf{Z}} \mathbf{f} / d, \\ \Gamma_2 &= \boldsymbol{\varepsilon}^\top [\mathbf{U}]_{\mathbf{Z}} \boldsymbol{\varepsilon} / d, & \Gamma_4 &= \boldsymbol{\varepsilon}^\top \mathbf{Z} \boldsymbol{\Upsilon} \tilde{\mathbf{V}} / \sqrt{d}, & \Gamma_5 &= \boldsymbol{\varepsilon}^\top [\mathbf{U}]_{\mathbf{Z}} \mathbf{f} / d. \end{aligned}$$

Note that the terms  $F_0^2$  and  $F_{1,d}^2$  in (I.30) are constants, and therefore do not contribute to the variance of  $R_d(\mathbf{X}, \boldsymbol{\Theta}, \lambda, \tilde{\beta}_{1,d}, \boldsymbol{\varepsilon})$ . In the following, we aim to show that  $\mathbb{E}_{\mathbf{X}, \boldsymbol{\Theta}} [\text{Var}_{\tilde{\beta}_{1,d}, \boldsymbol{\varepsilon}}(\Gamma_k)] = o_d(1)$  for  $k \in [5]$ . Consider first the variance of  $\Gamma_1$ . From Lemma I.2, we have  $\tilde{\mathbf{V}} = \boldsymbol{\Lambda}_{d,1} \boldsymbol{\Theta} \tilde{\beta}_{1,d} + \boldsymbol{\Lambda}_{d,0} \mathbf{1}_N F_0$ . Then

$$\begin{aligned} \text{Var}_{\tilde{\beta}_{1,d}}(\Gamma_1) &= \text{Var}_{\tilde{\beta}_{1,d}} \left( (\mathbf{X} \tilde{\beta}_{1,d} + \mathbf{1}_n F_0)^\top \mathbf{Z} \boldsymbol{\Upsilon} (\boldsymbol{\Lambda}_{d,1} \boldsymbol{\Theta} \tilde{\beta}_{1,d} + \boldsymbol{\Lambda}_{d,0} \mathbf{1}_N F_0) / \sqrt{d} \right) \\ &= \frac{1}{d} \text{Var}_{\tilde{\beta}_{1,d}} \left( \tilde{\beta}_{1,d}^\top \mathbf{X}^T \mathbf{Z} \boldsymbol{\Upsilon} \boldsymbol{\Lambda}_{d,1} \boldsymbol{\Theta} \tilde{\beta}_{1,d} + \tilde{\beta}_{1,d}^\top \mathbf{X} \mathbf{Z} \boldsymbol{\Upsilon} \boldsymbol{\Lambda}_{d,0} \mathbf{1}_N F_0 \right. \\ &\quad \left. + F_0 \mathbf{1}_n^\top \mathbf{Z} \boldsymbol{\Upsilon} \boldsymbol{\Lambda}_{d,1} \boldsymbol{\Theta} \tilde{\beta}_{1,d} + F_0 \mathbf{1}_n^\top \mathbf{Z} \boldsymbol{\Upsilon} \boldsymbol{\Lambda}_{d,0} \mathbf{1}_N F_0 \right) \\ &\leq \frac{4}{d} \text{Var}_{\tilde{\beta}_{1,d}} \left( \tilde{\beta}_{1,d}^\top \mathbf{X}^T \mathbf{Z} \boldsymbol{\Upsilon} \boldsymbol{\Lambda}_{d,1} \boldsymbol{\Theta} \tilde{\beta}_{1,d} \right) + \frac{4}{d} \text{Var}_{\tilde{\beta}_{1,d}} \left( \tilde{\beta}_{1,d}^\top \mathbf{X} \mathbf{Z} \boldsymbol{\Upsilon} \boldsymbol{\Lambda}_{d,0} \mathbf{1}_N F_0 \right) \\ &\quad + \frac{4}{d} \text{Var}_{\tilde{\beta}_{1,d}} \left( F_0 \mathbf{1}_n^\top \mathbf{Z} \boldsymbol{\Upsilon} \boldsymbol{\Lambda}_{d,1} \boldsymbol{\Theta} \tilde{\beta}_{1,d} \right) + 0 \\ &\leq 4F_{1,d}^4 \cdot \underbrace{\frac{1}{d^3} \left( \|\mathbf{X}^T \mathbf{Z} \boldsymbol{\Upsilon} \boldsymbol{\Lambda}_{d,1} \boldsymbol{\Theta}\|_F^2 + \text{tr}(\mathbf{X}^T \mathbf{Z} \boldsymbol{\Upsilon} \boldsymbol{\Lambda}_{d,1} \boldsymbol{\Theta} \mathbf{X}^T \mathbf{Z} \boldsymbol{\Upsilon} \boldsymbol{\Lambda}_{d,1} \boldsymbol{\Theta}) \right)}_{I_1} \\ &\quad + 4F_{1,d}^2 F_0^2 \cdot \underbrace{\left( \frac{1}{d} \text{tr} \left( \frac{\mathbf{X} \mathbf{X}^\top}{d} [\boldsymbol{\Lambda}_{d,0} \mathbf{1}_N \mathbf{1}_N^\top \boldsymbol{\Lambda}_{d,0}]_{\mathbf{Z}} \right) + \frac{1}{d} \text{tr} \left( \mathbf{1}_n \mathbf{1}_n^\top \left[ \boldsymbol{\Lambda}_{d,1} \frac{\boldsymbol{\Theta} \boldsymbol{\Theta}^\top}{d} \boldsymbol{\Lambda}_{d,1} \right]_{\mathbf{Z}} \right) \right)}_{I_2}. \end{aligned}$$

The first inequality holds from  $\text{Var}(a+b) \leq 2\text{Var}(a) + 2\text{Var}(b)$ ,  $I_1$  comes from Lemma I.7 and  $I_2$  comes from  $\text{Var}(a) \leq \mathbb{E}a^2$ . Note that  $\|\boldsymbol{\Lambda}_{d,1}\|_{\text{op}} = O_d(1/\sqrt{d})$  and  $\|\mathbf{Z} \boldsymbol{\Upsilon}\|_{\text{op}} \leq 1/(2\sqrt{\lambda})$ , we conclude

that

$$\begin{aligned} \mathbb{E}_{\mathbf{X}, \Theta} |I_1| &\leq \left| \frac{1}{d^3} \mathbb{E}_{\mathbf{X}, \Theta} \text{tr}(\mathbf{X}^T \mathbf{Z} \Upsilon \Lambda_{d,1} \Theta \Theta^\top \Lambda_{d,1} \Upsilon \mathbf{Z}^\top \mathbf{X}) \right| + \left| \frac{1}{d^3} \mathbb{E}_{\mathbf{X}, \Theta} \text{tr}(\mathbf{X}^T \mathbf{Z} \Upsilon \Lambda_{d,1} \Theta \mathbf{X}^T \mathbf{Z} \Upsilon \Lambda_{d,1} \Theta) \right| \\ &\leq \frac{1}{4\lambda} \mathbb{E}_{\mathbf{X}, \Theta} \left\| \Lambda_{d,1} \frac{\Theta \Theta^\top}{d} \Lambda_{d,1} \right\|_{\text{op}} \cdot \left\| \frac{\mathbf{X} \mathbf{X}^\top}{d} \right\|_{\text{op}} + \frac{1}{4\lambda} \mathbb{E}_{\mathbf{X}, \Theta} \left\| \frac{\Lambda_{d,1} \Theta \mathbf{X}^\top}{d} \right\|_{\text{op}}^2 = o_d(1). \end{aligned}$$

Furthermore from Lemma [I.3](#), we have

$$\mathbb{E}_{\mathbf{X}, \Theta} |I_2| = \mathbb{E}_{\mathbf{X}, \Theta} \left| \frac{1}{d} \text{tr} \left( \frac{\mathbf{X} \mathbf{X}^\top}{d} [\Lambda_{d,0} \mathbf{1}_N \mathbf{1}_N^\top \Lambda_{d,0}]_{\mathbf{Z}} \right) + \frac{1}{d} \text{tr} \left( \mathbf{1}_n \mathbf{1}_n^\top \left[ \Lambda_{d,1} \frac{\Theta \Theta^\top}{d} \Lambda_{d,1} \right]_{\mathbf{Z}} \right) \right| = o_d(1).$$

Thus we obtain  $\mathbb{E}_{\mathbf{X}, \Theta} (\text{Var}_{\beta_d}(\Gamma_1)) = o_d(1)$ . Similarly, we have for  $\Gamma_2$ ,

$$\begin{aligned} \text{Var}_{\tilde{\beta}_{1,d}}(\Gamma_2) &= \text{Var}_{\tilde{\beta}_{1,d}} \left( (\mathbf{X} \tilde{\beta}_{1,d} + \mathbf{1}_n F_0)^\top [\mathbf{U}]_{\mathbf{Z}} (\mathbf{X} \tilde{\beta}_{1,d} + \mathbf{1}_n F_0) / d \right) \\ &= \frac{1}{d^2} \text{Var}_{\tilde{\beta}_{1,d}} \left( \tilde{\beta}_{1,d}^\top \mathbf{X}^T [\mathbf{U}]_{\mathbf{Z}} \mathbf{X} \tilde{\beta}_{1,d} + \tilde{\beta}_{1,d}^\top \mathbf{X}^\top [\mathbf{U}]_{\mathbf{Z}} \mathbf{1}_n F_0 \right. \\ &\quad \left. + F_0 \mathbf{1}_n^\top [\mathbf{U}]_{\mathbf{Z}} \mathbf{X} \tilde{\beta}_{1,d} + F_0 \mathbf{1}_n^\top [\mathbf{U}]_{\mathbf{Z}} \mathbf{1}_n F_0 \right) \\ &\leq \frac{4}{d^2} \text{Var}_{\tilde{\beta}_{1,d}} \left( \tilde{\beta}_{1,d}^\top \mathbf{X}^T [\mathbf{U}]_{\mathbf{Z}} \mathbf{X} \tilde{\beta}_{1,d} \right) + \frac{8}{d^2} \text{Var}_{\tilde{\beta}_{1,d}} \left( \tilde{\beta}_{1,d}^\top \mathbf{X}^\top [\mathbf{U}]_{\mathbf{Z}} \mathbf{1}_n F_0 \right) \\ &\leq 8F_{1,d}^4 \cdot \underbrace{\frac{1}{d^4} \text{tr} \left( \mathbf{X}^T [\mathbf{U}]_{\mathbf{Z}} \mathbf{X} \mathbf{X}^T [\mathbf{U}]_{\mathbf{Z}} \mathbf{X} \right)}_{I_3} + 8F_{1,d}^2 F_0^2 \cdot \underbrace{\frac{1}{d^2} \text{tr} \left( [\mathbf{U}]_{\mathbf{Z}} \mathbf{1}_n \mathbf{1}_n^\top [\mathbf{U}]_{\mathbf{Z}} \frac{\mathbf{X} \mathbf{X}^\top}{d} \right)}_{I_4}. \end{aligned}$$

The first inequality holds from  $\text{Var}(a+b) \leq 2\text{Var}(a) + 2\text{Var}(b)$ ,  $I_3$  comes from Lemma [I.7](#) and the symmetric of  $\mathbf{X}^T [\mathbf{U}]_{\mathbf{Z}} \mathbf{X}$ , and  $I_4$  comes from  $\text{Var}(a) \leq \mathbb{E}a^2$ . Define  $\Gamma_{\mathbf{U}} = \mathbf{U} - \Lambda_{d,0} \mathbf{1}_N \mathbf{1}_N^\top \Lambda_{d,0}$ , from Lemma [I.2](#),  $\mathbb{E} \|\Gamma_{\mathbf{U}}\|_{\text{op}}^2 = O_d(1)$ . By replacing  $\mathbf{U}$  by  $\Lambda_{d,0} \mathbf{1}_N \mathbf{1}_N^\top \Lambda_{d,0} + \Gamma_{\mathbf{U}}$  in the terms  $I_3$  and  $I_4$ , we obtain the following equalities:

$$\begin{aligned} I_3 &= \frac{1}{d^2} \text{tr} \left( [\Lambda_{d,0} \mathbf{1}_N \mathbf{1}_N^\top \Lambda_{d,0} + \Gamma_{\mathbf{U}}]_{\mathbf{Z}} \frac{\mathbf{X} \mathbf{X}^\top}{d} [\Lambda_{d,0} \mathbf{1}_N \mathbf{1}_N^\top \Lambda_{d,0} + \Gamma_{\mathbf{U}}]_{\mathbf{Z}} \frac{\mathbf{X} \mathbf{X}^\top}{d} \right) \\ &= \frac{1}{d^2} \text{tr} \left( \underbrace{[\Lambda_{d,0} \mathbf{1}_N \mathbf{1}_N^\top \Lambda_{d,0}]_{\mathbf{Z}} \frac{\mathbf{X} \mathbf{X}^\top}{d} [\Lambda_{d,0} \mathbf{1}_N \mathbf{1}_N^\top \Lambda_{d,0}]_{\mathbf{Z}} \frac{\mathbf{X} \mathbf{X}^\top}{d}}_{J_1} \right) \\ &\quad + \underbrace{\frac{2}{d^2} \text{tr} \left( [\Lambda_{d,0} \mathbf{1}_N \mathbf{1}_N^\top \Lambda_{d,0}]_{\mathbf{Z}} \frac{\mathbf{X} \mathbf{X}^\top}{d} [\Gamma_{\mathbf{U}}]_{\mathbf{Z}} \frac{\mathbf{X} \mathbf{X}^\top}{d} \right)}_{J_2} + \underbrace{\frac{1}{d^2} \text{tr} \left( [\Gamma_{\mathbf{U}}]_{\mathbf{Z}} \frac{\mathbf{X} \mathbf{X}^\top}{d} [\Gamma_{\mathbf{U}}]_{\mathbf{Z}} \frac{\mathbf{X} \mathbf{X}^\top}{d} \right)}_{J_3}, \\ I_4 &= \frac{1}{d^2} \text{tr} \left( [\Lambda_{d,0} \mathbf{1}_N \mathbf{1}_N^\top \Lambda_{d,0} + \Gamma_{\mathbf{U}}]_{\mathbf{Z}} \mathbf{1}_n \mathbf{1}_n^\top [\Lambda_{d,0} \mathbf{1}_N \mathbf{1}_N^\top \Lambda_{d,0} + \Gamma_{\mathbf{U}}]_{\mathbf{Z}} \frac{\mathbf{X} \mathbf{X}^\top}{d} \right) \\ &= \frac{1}{d^2} \text{tr} \left( \underbrace{[\Lambda_{d,0} \mathbf{1}_N \mathbf{1}_N^\top \Lambda_{d,0}]_{\mathbf{Z}} \mathbf{1}_n \mathbf{1}_n^\top [\Lambda_{d,0} \mathbf{1}_N \mathbf{1}_N^\top \Lambda_{d,0}]_{\mathbf{Z}} \frac{\mathbf{X} \mathbf{X}^\top}{d}}_{K_1} \right) \end{aligned}$$

$$+ \underbrace{\frac{2}{d^2} \text{tr} \left( [\mathbf{\Lambda}_{d,0} \mathbf{1}_N \mathbf{1}_N^\top \mathbf{\Lambda}_{d,0}]_{\mathbf{Z}} \mathbf{1}_n \mathbf{1}_n^\top [\mathbf{\Gamma}_{\mathbf{U}}]_{\mathbf{Z}} \frac{\mathbf{X} \mathbf{X}^\top}{d} \right)}_{K_2} + \underbrace{\frac{1}{d^2} \text{tr} \left( [\mathbf{\Gamma}_{\mathbf{U}}]_{\mathbf{Z}} \mathbf{1}_n \mathbf{1}_n^\top [\mathbf{\Gamma}_{\mathbf{U}}]_{\mathbf{Z}} \frac{\mathbf{X} \mathbf{X}^\top}{d} \right)}_{K_3}.$$

We investigate the terms  $K_i$ ,  $i = 1, 2, 3$ . The investigation of terms  $J_i$ ,  $i = 1, 2, 3$  are quite similar, we omit the proof for  $J_i$  for brevity. Consider first the term  $K_2$ . Due to  $\mathbb{E} \|\mathbf{\Gamma}_{\mathbf{U}}\|_{\text{op}}^2 = O_d(1)$ , it is true that

$$\begin{aligned} \left( \mathbb{E} \left\| \frac{\mathbf{1}_n \mathbf{1}_n^\top}{d} [\mathbf{\Gamma}_{\mathbf{U}}]_{\mathbf{Z}} \frac{\mathbf{X} \mathbf{X}^\top}{d} \right\|_{\text{op}}^2 \right)^{1/2} &= O_d(1) \cdot \left( \mathbb{E} \left\| \mathbf{\Gamma}_{\mathbf{U}} \frac{\mathbf{X} \mathbf{X}^\top}{d} \right\|_{\text{op}}^2 \right)^{1/2} \\ &= O_d(1) \cdot \left( \mathbb{E} \|\mathbf{\Gamma}_{\mathbf{U}}\|_{\text{op}}^2 \right)^{1/2} \cdot \left( \mathbb{E} \left\| \frac{\mathbf{X} \mathbf{X}^\top}{d} \right\|_{\text{op}}^2 \right)^{1/2} = O_d(1). \end{aligned}$$

The second equality comes from the independence of  $\mathbf{\Gamma}_{\mathbf{U}}$  and  $\mathbf{X}$ . Note that for any rank 1 matrix  $\mathbf{A}$ ,  $|\text{tr} \mathbf{A}| = \|\mathbf{A}\|_{\text{op}}$ , the term  $K_2$  has the property

$$\begin{aligned} \mathbb{E}_{\mathbf{X}, \mathbf{\Theta}} |K_2| &= \mathbb{E}_{\mathbf{X}, \mathbf{\Theta}} \left| \frac{2}{d^2} \text{tr} \left( [\mathbf{\Lambda}_{d,0} \mathbf{1}_N \mathbf{1}_N^\top \mathbf{\Lambda}_{d,0}]_{\mathbf{Z}} \mathbf{1}_n \mathbf{1}_n^\top [\mathbf{\Gamma}_{\mathbf{U}}]_{\mathbf{Z}} \frac{\mathbf{X} \mathbf{X}^\top}{d} \right) \right| \\ &\leq \mathbb{E}_{\mathbf{X}, \mathbf{\Theta}} \left( \frac{2}{d} \left\| [\mathbf{\Lambda}_{d,0} \mathbf{1}_N \mathbf{1}_N^\top \mathbf{\Lambda}_{d,0}]_{\mathbf{Z}} \right\|_{\text{op}} \cdot \left\| \frac{\mathbf{1}_n \mathbf{1}_n^\top}{d} [\mathbf{\Gamma}_{\mathbf{U}}]_{\mathbf{Z}} \frac{\mathbf{X} \mathbf{X}^\top}{d} \right\|_{\text{op}} \right) \\ &\leq \frac{2}{d} \left( \mathbb{E}_{\mathbf{X}, \mathbf{\Theta}} \left\| [\mathbf{\Lambda}_{d,0} \mathbf{1}_N \mathbf{1}_N^\top \mathbf{\Lambda}_{d,0}]_{\mathbf{Z}} \right\|_{\text{op}}^2 \cdot \mathbb{E}_{\mathbf{X}, \mathbf{\Theta}} \left\| \frac{\mathbf{1}_n \mathbf{1}_n^\top}{d} [\mathbf{\Gamma}_{\mathbf{U}}]_{\mathbf{Z}} \frac{\mathbf{X} \mathbf{X}^\top}{d} \right\|_{\text{op}}^2 \right)^{1/2} \\ &= o_d(1) \cdot O_d(1) = o_d(1). \end{aligned}$$

The equality comes from the estimation of  $\mathcal{D}$  in Lemma I.3. For the term  $K_1$ , it is true from the estimation of  $\mathcal{D}$  in Lemma I.3 that

$$\begin{aligned} \frac{1}{d} \left( \mathbb{E} \left\| \frac{\mathbf{1}_n \mathbf{1}_n^\top}{d} [\mathbf{\Lambda}_{d,0} \mathbf{1}_N \mathbf{1}_N^\top \mathbf{\Lambda}_{d,0}]_{\mathbf{Z}} \frac{\mathbf{X} \mathbf{X}^\top}{d} \right\|_{\text{op}}^2 \right)^{1/2} \\ \leq \frac{1}{d} \left\| \frac{\mathbf{1}_n \mathbf{1}_n^\top}{d} \right\|_{\text{op}} \left( \mathbb{E} \left\| [\mathbf{\Lambda}_{d,0} \mathbf{1}_N \mathbf{1}_N^\top \mathbf{\Lambda}_{d,0}]_{\mathbf{Z}} \frac{\mathbf{X} \mathbf{X}^\top}{d} \right\|_{\text{op}}^2 \right)^{1/2} = O_d(1) \cdot o_d(1) = o_d(1). \end{aligned}$$

By repeating the arguments used previously for the term  $K_2$  but for the consideration of  $K_1$ , we have

$$\mathbb{E}_{\mathbf{X}, \mathbf{\Theta}} |K_1| = \frac{1}{d^2} \mathbb{E}_{\mathbf{X}, \mathbf{\Theta}} \left| \text{tr} \left( [\mathbf{\Lambda}_{d,0} \mathbf{1}_N \mathbf{1}_N^\top \mathbf{\Lambda}_{d,0}]_{\mathbf{Z}} \mathbf{1}_n \mathbf{1}_n^\top [\mathbf{\Lambda}_{d,0} \mathbf{1}_N \mathbf{1}_N^\top \mathbf{\Lambda}_{d,0}]_{\mathbf{Z}} \frac{\mathbf{X} \mathbf{X}^\top}{d} \right) \right| = o_d(1).$$

For the term  $K_3$ , similarly we have

$$\begin{aligned} \mathbb{E}_{\mathbf{X}, \mathbf{\Theta}} |K_3| &= \mathbb{E}_{\mathbf{X}, \mathbf{\Theta}} \left| \frac{1}{d^2} \text{tr} \left( [\mathbf{\Gamma}_{\mathbf{U}}]_{\mathbf{Z}} \mathbf{1}_n \mathbf{1}_n^\top [\mathbf{\Gamma}_{\mathbf{U}}]_{\mathbf{Z}} \frac{\mathbf{X} \mathbf{X}^\top}{d} \right) \right| \\ &\leq \frac{1}{d^2} \mathbb{E}_{\mathbf{X}, \mathbf{\Theta}} \left( \left\| [\mathbf{\Gamma}_{\mathbf{U}}]_{\mathbf{Z}} \mathbf{1}_n \mathbf{1}_n^\top \right\|_{\text{op}} \cdot \left\| [\mathbf{\Gamma}_{\mathbf{U}}]_{\mathbf{Z}} \frac{\mathbf{X} \mathbf{X}^\top}{d} \right\|_{\text{op}} \right) \end{aligned}$$



$$\begin{aligned}
 &\leq \frac{1}{d^2} \left( \mathbb{E}_{\mathbf{X}, \Theta} \left\| [\Gamma \mathbf{U}]_{\mathbf{Z}} \mathbf{1}_n \mathbf{1}_n^\top \right\|_{\text{op}}^2 \cdot \mathbb{E}_{\mathbf{X}, \Theta} \left\| [\Gamma \mathbf{U}]_{\mathbf{Z}} \frac{\mathbf{X} \mathbf{X}^\top}{d} \right\|_{\text{op}}^2 \right)^{1/2} \\
 &\leq \frac{1}{d^2} (\mathbb{E}_{\mathbf{X}, \Theta} \left\| [\Gamma \mathbf{U}]_{\mathbf{Z}} \mathbf{1}_n \mathbf{1}_n^\top \right\|_{\text{op}}^2)^{1/2} \cdot O_d(1) \\
 &= o_d(1) \cdot O_d(1) = o_d(1).
 \end{aligned}$$

Now we conclude that  $K_1$ ,  $K_2$  and  $K_3$  are all small terms under the expectation over  $\mathbf{X}$  and  $\Theta$ , we immediately get that

$$\mathbb{E}_{\mathbf{X}, \Theta} |I_4| = o_d(1).$$

Similarly we get  $\mathbb{E}_{\mathbf{X}, \Theta} |I_3| = o_d(1)$ , thus we conclude that  $\mathbb{E}_{\mathbf{X}, \Theta} \text{Var}_{\tilde{\beta}_{1,d}}(\Gamma_2) = o_d(1)$ . We omit the other terms for brevity. The proof of Lemma I.4 is complete.

### Proof of Lemma I.7

We have

$$\mathbb{E}[\mathbf{h}^\top \mathbf{A} \mathbf{h}] = \mathbb{E} \text{tr}(\mathbf{A} \mathbf{h} \mathbf{h}^\top) = \frac{F_{1,d}^2}{d} \text{tr}(\mathbf{A}).$$

Hence we have

$$\begin{aligned}
 \text{Var}(\mathbf{h}^\top \mathbf{A} \mathbf{h}) &= \sum_{i_1, i_2, i_3, i_4} \mathbb{E} \left[ \mathbf{h}_{i_1} \mathbf{A}_{i_1, i_2} \mathbf{h}_{i_2} \mathbf{h}_{i_3} \mathbf{A}_{i_3, i_4} \mathbf{h}_{i_4} \right] - \frac{F_{1,d}^4}{d^2} \text{tr}(\mathbf{A})^2 \\
 &= \left\{ \left( \sum_{\substack{i_1=i_2, i_3=i_4 \\ i_1=i_3}} + \sum_{\substack{i_1=i_2, i_3=i_4 \\ i_1 \neq i_3}} + \sum_{\substack{i_1=i_3, i_2=i_4 \\ i_1 \neq i_2}} + \sum_{\substack{i_1=i_4, i_2=i_3 \\ i_1 \neq i_2}} \right) \mathbb{E} \left[ \mathbf{h}_{i_1} \mathbf{A}_{i_1, i_2} \mathbf{h}_{i_2} \mathbf{h}_{i_3} \mathbf{A}_{i_3, i_4} \mathbf{h}_{i_4} \right] \right\} \\
 &\quad - \frac{F_{1,d}^4}{d^2} \text{tr}(\mathbf{A})^2 \\
 &= \frac{F_{1,d}^4}{d^2} \left( \sum_{i=1}^d \mathbf{A}_{i,i}^2 \cdot \frac{d^2}{F_{1,d}^4} (\mathbb{E} h_i^4) + \sum_{i \neq j} \mathbf{A}_{i,i} \mathbf{A}_{j,j} + \sum_{i \neq j} (\mathbf{A}_{i,j} \mathbf{A}_{i,j} + \mathbf{A}_{i,j} \mathbf{A}_{j,i}) - \text{tr}(\mathbf{A})^2 \right) \\
 &= \frac{F_{1,d}^4}{d^2} \left( \sum_{i=1}^d 3 \mathbf{A}_{i,i}^2 + \sum_{i \neq j} (\mathbf{A}_{i,j} \mathbf{A}_{i,j} + \mathbf{A}_{i,j} \mathbf{A}_{j,i}) + \sum_{i \neq j} \mathbf{A}_{i,i} \mathbf{A}_{j,j} - \text{tr}(\mathbf{A})^2 \right) \\
 &= \frac{F_{1,d}^4}{d^2} \left( \sum_{i,j} (\mathbf{A}_{i,j} \mathbf{A}_{i,j} + \mathbf{A}_{i,j} \mathbf{A}_{j,i}) + \sum_{i,j} \mathbf{A}_{i,i} \mathbf{A}_{j,j} - \text{tr}(\mathbf{A})^2 \right) = \frac{F_{1,d}^4}{d^2} (\|\mathbf{A}\|_F^2 + \text{tr}(\mathbf{A}^2)).
 \end{aligned}$$

This proves Lemma I.7.

## II. Proof of Proposition A.4

We first recall the following matrix differential rules:

$$\frac{\partial \det(\mathbf{Y})}{\partial x} = \det(\mathbf{Y}) \cdot \text{tr} \left( \mathbf{Y}^{-1} \cdot \frac{\partial \mathbf{Y}}{\partial x} \right), \quad \frac{\partial \mathbf{Y}^{-1}}{\partial x} = -\mathbf{Y}^{-1} \frac{\partial \mathbf{Y}}{\partial x} \mathbf{Y}^{-1}. \quad (\text{II.1})$$

Let  $q_i, q_j$  be the elements in the vector  $\mathbf{q}$ . Now the matrix  $\mathbf{A} = \mathbf{A}(\mathbf{q}, \boldsymbol{\mu})$  (see Definition A.3) is linear in  $\mathbf{q}$ , thus  $\frac{\partial^2 \mathbf{A}}{\partial q_i \partial q_j} = 0$ . Therefore by the definition  $G_d(\xi) = \frac{1}{d} \log \det(\mathbf{A} - \xi \mathbf{I})$  and the matrix derivative rules in (II.1), we have

$$\frac{\partial G_d}{\partial q_i} = \frac{1}{d} \text{tr} \left( (\mathbf{A} - \xi \mathbf{I})^{-1} \frac{\partial \mathbf{A}}{\partial q_i} \right), \quad (\text{II.2})$$

$$\frac{\partial^2 G_d}{\partial q_i \partial q_j} = \frac{1}{d} \text{tr} \left( \frac{\partial (\mathbf{A} - \xi \mathbf{I})^{-1}}{\partial q_j} \frac{\partial \mathbf{A}}{\partial q_i} \right) = -\frac{1}{d} \text{tr} \left( (\mathbf{A} - \xi \mathbf{I})^{-1} \frac{\partial \mathbf{A}}{\partial q_j} (\mathbf{A} - \xi \mathbf{I})^{-1} \frac{\partial \mathbf{A}}{\partial q_i} \right). \quad (\text{II.3})$$

By the Schur complement formula, we further have

$$(\mathbf{A}(\mathbf{0}, \boldsymbol{\mu}) - \xi \mathbf{I}_P)^{-1} = \begin{bmatrix} -\xi \mathbf{I}_N & \mathbf{Z}^\top \\ \mathbf{Z} & -\xi \mathbf{I}_n \end{bmatrix}^{-1} = \begin{bmatrix} * & \mathbf{Z}^\top (\mathbf{Z} \mathbf{Z}^\top - \xi^2 \mathbf{I}_n)^{-1} \\ (\mathbf{Z} \mathbf{Z}^\top - \xi^2 \mathbf{I}_n)^{-1} \mathbf{Z} & * \end{bmatrix},$$

where we use  $*$  to hide the irrelevant blocks in the matrix inverse. Moreover, by definition, it holds that

$$\begin{aligned} \frac{\partial \mathbf{A}(\mathbf{q}, \boldsymbol{\mu})}{\partial q_1} &= \begin{bmatrix} \mathbf{0} & \frac{1}{d} \mathbf{M}_1 \boldsymbol{\Theta} \mathbf{X}^\top \\ \frac{1}{d} \mathbf{X} \boldsymbol{\Theta}^\top \mathbf{M}_1 & \mathbf{0} \end{bmatrix}, \quad \frac{\partial \mathbf{A}(\mathbf{q}, \boldsymbol{\mu})}{\partial q_4} = \begin{bmatrix} \mathbf{M}_1 \frac{\boldsymbol{\Theta} \boldsymbol{\Theta}^\top}{d} \mathbf{M}_1 & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix}, \\ \frac{\partial \mathbf{A}(\mathbf{q}, \boldsymbol{\mu})}{\partial q_2} &= \begin{bmatrix} \mathbf{M}_*^2 & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix}, \quad \frac{\partial \mathbf{A}(\mathbf{q}, \boldsymbol{\mu})}{\partial q_3} = \begin{bmatrix} \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{I}_n \end{bmatrix}, \quad \frac{\partial \mathbf{A}(\mathbf{q}, \boldsymbol{\mu})}{\partial q_5} = \begin{bmatrix} \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \frac{\mathbf{X} \mathbf{X}^\top}{d} \end{bmatrix}. \end{aligned}$$

By plugging these derivatives into (II.2) and (II.3), we continue the calculation with  $\xi = \xi^*$ . Note that we have the identity  $(\mathbf{Z} \mathbf{Z}^\top - (\xi^*)^2 \mathbf{I}_n)^{-1} \mathbf{Z} = \mathbf{Z} (\mathbf{Z}^\top \mathbf{Z} + \lambda \mathbf{I}_N)^{-1} = \mathbf{Z} \boldsymbol{\Upsilon}$ . Then by (II.2), we have

$$\begin{aligned} \left. \frac{\partial G_d(\xi^*; \mathbf{q}, \boldsymbol{\mu})}{\partial q_1} \right|_{\mathbf{q}=\mathbf{0}} &= \frac{1}{d} \text{tr} \left( \begin{bmatrix} * & \mathbf{Z}^\top (\mathbf{Z} \mathbf{Z}^\top - (\xi^*)^2 \mathbf{I}_n)^{-1} \\ (\mathbf{Z} \mathbf{Z}^\top - (\xi^*)^2 \mathbf{I}_n)^{-1} \mathbf{Z} & * \end{bmatrix} \begin{bmatrix} \mathbf{0} & \frac{1}{d} \mathbf{M}_1 \boldsymbol{\Theta} \mathbf{X}^\top \\ \frac{1}{d} \mathbf{X} \boldsymbol{\Theta}^\top \mathbf{M}_1 & \mathbf{0} \end{bmatrix} \right) \\ &= \frac{1}{d} \text{tr} \left( \begin{bmatrix} * & \boldsymbol{\Upsilon} \mathbf{Z}^\top \\ \mathbf{Z} \boldsymbol{\Upsilon} & * \end{bmatrix} \begin{bmatrix} \mathbf{0} & \frac{1}{d} \mathbf{M}_1 \boldsymbol{\Theta} \mathbf{X}^\top \\ \frac{1}{d} \mathbf{X} \boldsymbol{\Theta}^\top \mathbf{M}_1 & \mathbf{0} \end{bmatrix} \right) = \frac{2}{d} \text{tr} \mathbf{M}_1 \frac{\boldsymbol{\Theta} \mathbf{X}^\top}{d} \mathbf{Z} \boldsymbol{\Upsilon}, \end{aligned}$$

Similarly, by (II.3), we have

$$\begin{aligned} -\left. \frac{\partial^2 G_d(\xi^*; \mathbf{q}, \boldsymbol{\mu})}{\partial q_4 \partial q_5} \right|_{\mathbf{q}=\mathbf{0}} &= \text{tr} \left( \begin{bmatrix} * & \boldsymbol{\Upsilon} \mathbf{Z}^\top \\ \mathbf{Z} \boldsymbol{\Upsilon} & * \end{bmatrix} \begin{bmatrix} \mathbf{M}_1 \frac{\boldsymbol{\Theta} \boldsymbol{\Theta}^\top}{d} \mathbf{M}_1 & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \begin{bmatrix} * & \boldsymbol{\Upsilon} \mathbf{Z}^\top \\ \mathbf{Z} \boldsymbol{\Upsilon} & * \end{bmatrix} \begin{bmatrix} \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \frac{\mathbf{X} \mathbf{X}^\top}{d} \end{bmatrix} \right) \\ &= \frac{1}{d} \text{tr} \mathbf{Z} \boldsymbol{\Upsilon} \mathbf{M}_1 \frac{\boldsymbol{\Theta} \boldsymbol{\Theta}^\top}{d} \mathbf{M}_1 \boldsymbol{\Upsilon} \mathbf{Z}^\top \frac{\mathbf{X} \mathbf{X}^\top}{d}, \\ -\left. \frac{\partial^2 G_d(\xi^*; \mathbf{q}, \boldsymbol{\mu})}{\partial q_2 \partial q_5} \right|_{\mathbf{q}=\mathbf{0}} &= \text{tr} \left( \begin{bmatrix} * & \boldsymbol{\Upsilon} \mathbf{Z}^\top \\ \mathbf{Z} \boldsymbol{\Upsilon} & * \end{bmatrix} \begin{bmatrix} \mathbf{M}_* \mathbf{M}_* & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \begin{bmatrix} * & \boldsymbol{\Upsilon} \mathbf{Z}^\top \\ \mathbf{Z} \boldsymbol{\Upsilon} & * \end{bmatrix} \begin{bmatrix} \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \frac{\mathbf{X} \mathbf{X}^\top}{d} \end{bmatrix} \right) \\ &= \frac{1}{d} \text{tr} \mathbf{Z} \boldsymbol{\Upsilon} \mathbf{M}_* \mathbf{M}_* \boldsymbol{\Upsilon} \mathbf{Z}^\top \frac{\mathbf{X} \mathbf{X}^\top}{d}, \\ -\left. \frac{\partial^2 G_d(\xi^*; \mathbf{q}, \boldsymbol{\mu})}{\partial q_3 \partial q_4} \right|_{\mathbf{q}=\mathbf{0}} &= \text{tr} \left( \begin{bmatrix} * & \boldsymbol{\Upsilon} \mathbf{Z}^\top \\ \mathbf{Z} \boldsymbol{\Upsilon} & * \end{bmatrix} \begin{bmatrix} \mathbf{M}_1 \frac{\boldsymbol{\Theta} \boldsymbol{\Theta}^\top}{d} \mathbf{M}_1 & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \begin{bmatrix} * & \boldsymbol{\Upsilon} \mathbf{Z}^\top \\ \mathbf{Z} \boldsymbol{\Upsilon} & * \end{bmatrix} \begin{bmatrix} \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{I}_n \end{bmatrix} \right) \\ &= \frac{1}{d} \text{tr} \mathbf{Z} \boldsymbol{\Upsilon} \mathbf{M}_1 \frac{\boldsymbol{\Theta} \boldsymbol{\Theta}^\top}{d} \mathbf{M}_1 \boldsymbol{\Upsilon} \mathbf{Z}^\top, \end{aligned}$$

$$\begin{aligned}
 -\frac{\partial^2 G_d(\xi^*; \mathbf{q}, \boldsymbol{\mu})}{\partial q_2 \partial q_3} \Big|_{\mathbf{q}=\mathbf{0}} &= \text{tr} \left( \begin{bmatrix} * & \boldsymbol{\Upsilon} \mathbf{Z}^\top \\ \mathbf{Z} \boldsymbol{\Upsilon} & * \end{bmatrix} \begin{bmatrix} \mathbf{M}_* \mathbf{M}_* & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \begin{bmatrix} * & \boldsymbol{\Upsilon} \mathbf{Z}^\top \\ \mathbf{Z} \boldsymbol{\Upsilon} & * \end{bmatrix} \begin{bmatrix} \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{I}_n \end{bmatrix} \right) \\
 &= \frac{1}{d} \text{tr} \mathbf{Z} \boldsymbol{\Upsilon} \mathbf{M}_* \mathbf{M}_* \boldsymbol{\Upsilon} \mathbf{Z}^\top.
 \end{aligned}$$

The above equations complete the proof of Proposition A.4.

### III. Properties of the fixed point equation

In this section, we justify the definition of  $\mathbf{m}(\xi; \mathbf{q}, \boldsymbol{\mu})$  below Definition A.5, by proving that there exists a constant  $\xi_0 > 0$ , such that the fixed point equation (A.2) has a unique solution defined on  $\{\xi : \Im(\xi) > \xi_0\}$  satisfying  $|m_j(\xi)| \leq 2\psi_j/\xi_0$  for  $j = 1, 2, 3$ . The result is given in the following lemma.

**Lemma III.1.** *Let  $\mathbf{F}(\mathbf{m}; \xi, \mathbf{q}, \boldsymbol{\mu})$ ,  $\mathbf{q} \in \mathcal{Q}$  be defined in Definition A.5, and  $\mathbb{D}(r) = \{z : |z| < r\}$  be the disk of radius  $r$  in the complex plane. There exists  $\xi_0 > 0$  such that, for any  $\xi \in \mathbb{C}_+$  with  $\Im(\xi) > \xi_0$ ,  $\mathbf{F}(\mathbf{m}; \xi, \mathbf{q}, \boldsymbol{\mu})$  is  $1/2$ -Lipschitz continuous with respect to the  $\ell_2$  norm, and the map  $\mathbf{m} \mapsto \mathbf{F}(\mathbf{m}; \xi, \mathbf{q}, \boldsymbol{\mu})$  admits a unique fixed point in  $\mathbb{D}(2\psi_1/\xi_0) \times \mathbb{D}(2\psi_2/\xi_0) \times \mathbb{D}(2\psi_3/\xi_0)$ .*

Lemma III.1 demonstrates that our definition of  $\mathbf{m}$  in Subsection A.3 as the unique fixed point of  $\mathbf{F}$  is valid.

**Proof** [Proof of Lemma III.1] We prove the existence and uniqueness of the solution by the Banach fixed point theorem when  $\Im(\xi) \geq \xi_0$  for some sufficiently large  $\xi_0$ . To do so, we want to show that

1.  $\mathbf{F}(\cdot; \mathbf{q}, \boldsymbol{\mu})$  maps domain  $\mathbb{D}(2\psi_1/\xi_0) \times \mathbb{D}(2\psi_2/\xi_0) \times \mathbb{D}(2\psi_3/\xi_0)$  into itself.
2.  $\mathbf{F}(\cdot; \mathbf{q}, \boldsymbol{\mu})$  is Lipschitz continuous with a Lipschitz constant smaller than 1.

For  $\mathbf{F}_1(\cdot; \mathbf{q}, \boldsymbol{\mu})$ , by Definition A.5, we have

$$\mathbf{F}_1(\mathbf{m}; \xi, \mathbf{q}, \boldsymbol{\mu}) = \frac{\psi_1}{-\xi + q_2 \mu_{1,*}^2 + H_1(\mathbf{m}; \mathbf{q}, \boldsymbol{\mu})},$$

where

$$H_1(\mathbf{m}; \mathbf{q}, \boldsymbol{\mu}) = -\mu_{1,*}^2 m_3 + \frac{1}{m_1 + \frac{-\mu_{2,1}^2(1+q_1)^2 m_2 m_3 + (1+\mu_{2,1}^2 m_2 q_4)(1+m_3 q_5)}{\mu_{1,1}^2 q_4(1+m_3 q_5) - \mu_{1,1}^2(1+q_1)^2 m_3}}. \quad (\text{III.1})$$

Note that  $q_4, q_5 \leq (1+q_1)/2$ . Thus for small enough  $r_0$ , we have for any  $\mathbf{m} \in \mathbb{D}(r_0)^3$

$$|H_1(\mathbf{m}; \mathbf{q}, \boldsymbol{\mu})| \leq 2 + 2|q_4| \mu_{1,1}^2. \quad (\text{III.2})$$

Now as long as  $\xi_0 \geq 4 + 4|q_4| \mu_{1,1}^2$ , it is clear that for  $\xi$  with  $\Im(\xi) \geq \xi_0$  we have

$$\Im(\xi) \geq \xi_0/2 + \xi_0/2 \geq \xi_0/2 + 2 + 2|q_4| \mu_{1,1}^2 \geq \xi_0/2 + |H_1(\mathbf{m}; \mathbf{q}, \boldsymbol{\mu})|. \quad (\text{III.3})$$

Therefore,

$$|\mathbf{F}_1(\mathbf{m}; \xi, \mathbf{q}, \boldsymbol{\mu})| \leq \frac{\psi_1}{|\Im(\xi - q_2 \mu_{1,*}^2 - H_1(\mathbf{m}; \mathbf{q}, \boldsymbol{\mu}))|}$$

$$\leq \frac{\psi_1}{\Im(\xi) - |H_1(\mathbf{m}; \mathbf{q}, \boldsymbol{\mu})|} \leq \frac{2\psi_1}{\xi_0},$$

where the last inequality follows from (III.3).

Similarly, for  $F_2$  and  $F_3$  we show that  $|F_2(\mathbf{m}; \xi, \mathbf{q}, \boldsymbol{\mu})| \leq 2\psi_2/\xi_0$  provided  $\Im(\xi) \geq \xi_0 \geq 4 + 4|q_4|\mu_{2,1}^2$ , and  $|F_3(\mathbf{m}; \xi, \mathbf{q}, \boldsymbol{\mu})| \leq 2\psi_3/\xi_0$  provided  $\Im(\xi) \geq \xi_0 \geq 4 + 4|q_5|$ . Therefore if  $\xi_0$  satisfies  $2\max\{\psi_1, \psi_2, \psi_3\}/\xi_0 \leq r_0$  and  $\xi_0 \geq 4 + 4\max\{|q_4|\mu_{1,1}^2, |q_4|\mu_{2,1}^2, |q_5|\}$ ,  $\mathbf{F}$  maps domain  $\mathbb{D}(2\psi_1/\xi_0) \times \mathbb{D}(2\psi_2/\xi_0) \times \mathbb{D}(2\psi_3/\xi_0)$  into itself.

As for the Lipschitz continuity of  $\mathbf{F}(\cdot; \mathbf{q}, \boldsymbol{\mu})$ , note that

$$\nabla_{\mathbf{m}} F_1(\mathbf{m}; \xi, \mathbf{q}, \boldsymbol{\mu}) = -\frac{\psi_1}{(-\xi + q_2\mu_{1,*}^2 + H_1(\mathbf{m}; \mathbf{q}, \boldsymbol{\mu}))^2} \cdot \nabla_{\mathbf{m}} H_1(\mathbf{m}; \mathbf{q}, \boldsymbol{\mu}).$$

With the same calculation as above, it is easy to see that when  $\xi_0$  is sufficiently large,  $\|\nabla_{\mathbf{m}} H_1(\mathbf{m}; \mathbf{q}, \boldsymbol{\mu})\|_2 \leq C(\mathbf{q}, \boldsymbol{\mu})$  for all  $\mathbf{m} \in \mathbb{D}(2\psi_1/\xi_0) \times \mathbb{D}(2\psi_2/\xi_0) \times \mathbb{D}(2\psi_3/\xi_0)$ , where  $C(\mathbf{q}, \boldsymbol{\mu})$  is a constant that only depends on  $\mathbf{q}$  and  $\boldsymbol{\mu}$ . Thus for such  $\xi_0$  and  $\xi$  with  $\Im(\xi) \geq \xi_0$ ,

$$\|\nabla_{\mathbf{m}} F_1(\mathbf{m}; \xi, \mathbf{q}, \boldsymbol{\mu})\|_2 \leq \frac{C(\mathbf{q}, \boldsymbol{\mu}) \cdot \psi_1}{\Im(\xi) - |H_1(\mathbf{m}; \mathbf{q}, \boldsymbol{\mu})|} \leq \frac{4C(\mathbf{q}, \boldsymbol{\mu}) \cdot \psi_1}{\xi_0} \leq \frac{1}{4},$$

where we again utilize (III.3). We can apply the same argument for  $F_2$  and  $F_3$ , and conclude that  $\mathbf{F}$  is  $1/2$ -Lipschitz on  $\mathbf{m} \in \mathbb{D}(2\psi_1/\xi_0) \times \mathbb{D}(2\psi_2/\xi_0) \times \mathbb{D}(2\psi_3/\xi_0)$ . Therefore by Banach fixed point theorem, there exists a unique fixed point of  $\mathbf{F}$ . Thus the solution of the implicit equations defined in Definition A.5 exists and is unique.  $\blacksquare$

## IV. Proof of Proposition A.6

The proof of Proposition A.6 is split into several sections. In Sections IV.1 and IV.2, we give some useful preliminary results. In Section IV.3, we show that  $\mathbf{m}(\xi; \mathbf{q}, \boldsymbol{\mu})$  is analytic on  $\{\xi : \Im(\xi) \geq \xi_0\}$ , and then prove the first and second conclusions of Proposition A.6. In Section IV.4, we prove the point convergence of  $M_d(\xi; \mathbf{q}, \boldsymbol{\mu})$  to  $\mathbf{m}(\xi; \mathbf{q}, \boldsymbol{\mu})$  under the additional assumption that  $\sigma_j(x)$ ,  $j = 1, 2$  are polynomials. In Section IV.5, we extend this point convergence result to general activation functions satisfying Assumption 3.2. In Section IV.6, we conclude the proof by showing the uniform convergence of  $M_d(\xi; \mathbf{q}, \boldsymbol{\mu})$  to  $\mathbf{m}(\xi; \mathbf{q}, \boldsymbol{\mu})$  on compact sets.

### IV.1 Equivalence between Gaussian and spherical versions

The first step in the proof of Proposition A.6 is to relate the Stieltjes transform  $M_d(\xi; \mathbf{q}, \boldsymbol{\mu})$  to the Stieltjes transform corresponding to Gaussian data and Gaussian random features.

**Definition IV.1.** Let  $(\bar{\boldsymbol{\theta}}_a)_{a \in [N]}$  be i.i.d. standard Gaussian random vectors distributed as  $N(\mathbf{0}, \mathbf{I}_d)$ , and  $\bar{\boldsymbol{\Theta}} \in \mathbb{R}^{N \times d}$  be the matrix whose  $a^{\text{th}}$  row is given by  $\bar{\boldsymbol{\theta}}_a$ . Similarly, we denote  $(\bar{\mathbf{x}}_i)_{i \in [n]} \sim_{\text{iid}} N(\mathbf{0}, \mathbf{I}_d)$ , and let  $\bar{\mathbf{X}} \in \mathbb{R}^{n \times d}$  be the matrix whose  $i^{\text{th}}$  row is  $\bar{\mathbf{x}}_i$ .  $\square$

Given these definitions, our original data inputs and random feature parameters which are distributed uniformly on the sphere  $\sqrt{d} \cdot \mathbb{S}^{d-1}$  can be represented as

$$\mathbf{x}_i = \sqrt{d} \cdot \frac{\bar{\mathbf{x}}_i}{\|\bar{\mathbf{x}}_i\|_2} \sim \text{Unif}(\sqrt{d} \cdot \mathbb{S}^{d-1}), \text{ and } \boldsymbol{\theta}_a = \sqrt{d} \cdot \frac{\bar{\boldsymbol{\theta}}_a}{\|\bar{\boldsymbol{\theta}}_a\|_2} \sim \text{Unif}(\sqrt{d} \cdot \mathbb{S}^{d-1}) \quad (\text{IV.1})$$

for all  $i \in [n]$  and  $a \in [N]$ . We can now consider the ‘‘Gaussian version’’ of the learning problem, where the data inputs are  $(\bar{\mathbf{x}}_i)_{i \in [n]}$ , and the double random feature model uses random parameters  $(\bar{\boldsymbol{\theta}}_a)_{a \in [N]}$  and activation functions

$$\phi_j(x) = \sigma_j(x) - \mathbb{E}_{G \sim \mathcal{N}(0,1)}[\sigma_j(G)], \quad j = 1, 2. \quad (\text{IV.2})$$

For this version of the learning problem, we can similarly construct the linear pencil matrix  $\bar{\mathbf{A}}(\mathbf{q}, \boldsymbol{\mu})$ , which is the counterpart of the linear pencil matrix  $\mathbf{A}(\mathbf{q}, \boldsymbol{\mu})$  defined in Definition A.3.

**Definition IV.2.** *The linear pencil matrix  $\bar{\mathbf{A}}(\mathbf{q}, \boldsymbol{\mu}) \in \mathbb{R}^{P \times P}$  ( $P = N + n$ ) is defined as*

$$\bar{\mathbf{A}}(\mathbf{q}, \boldsymbol{\mu}) = \begin{bmatrix} q_2 \mu_{1,*}^2 \mathbf{I}_{N_1} + q_4 \mu_{1,1}^2 \frac{\bar{\boldsymbol{\Theta}}_1 \bar{\boldsymbol{\Theta}}_1^\top}{d} & q_4 \mu_{1,1} \mu_{2,1} \frac{\bar{\boldsymbol{\Theta}}_1 \bar{\boldsymbol{\Theta}}_2^\top}{d} & \mathbf{J}_1^\top + q_1 \tilde{\mathbf{J}}_1^\top \\ q_4 \mu_{1,1} \mu_{2,1} \frac{\bar{\boldsymbol{\Theta}}_2 \bar{\boldsymbol{\Theta}}_1^\top}{d} & q_2 \mu_{2,*}^2 \mathbf{I}_{N_2} + q_4 \mu_{2,1}^2 \frac{\bar{\boldsymbol{\Theta}}_2 \bar{\boldsymbol{\Theta}}_2^\top}{d} & \mathbf{J}_2^\top + q_1 \tilde{\mathbf{J}}_2^\top \\ \mathbf{J}_1 + q_1 \tilde{\mathbf{J}}_1 & \mathbf{J}_2 + q_1 \tilde{\mathbf{J}}_2 & q_3 \mathbf{I}_n + q_5 \frac{\bar{\mathbf{X}} \bar{\mathbf{X}}^\top}{d} \end{bmatrix},$$

where  $\mathbf{J}_j = \phi_j(\bar{\mathbf{X}} \bar{\boldsymbol{\Theta}}_j^\top / \sqrt{d}) / \sqrt{d}$ ,  $\tilde{\mathbf{J}}_j = \frac{\mu_{j,1}}{d} \bar{\mathbf{X}} \bar{\boldsymbol{\Theta}}_j^\top$ ,  $j = 1, 2$ .  $\square$

We also define  $\bar{M}_d(\xi; \mathbf{q}, \boldsymbol{\mu}) = \frac{1}{d} \text{tr}[(\bar{\mathbf{A}}(\mathbf{q}, \boldsymbol{\mu}) - \xi \mathbf{I}_P)^{-1}]$  as the counterpart of the Stieltjes transform  $M_d(\xi; \mathbf{q}, \boldsymbol{\mu})$ . The following lemma establishes the equivalence between the two versions of the Stieltjes transforms.

**Lemma IV.3.** *Suppose that  $\sigma_j(x)$ ,  $j = 1, 2$ , are polynomials. Then for any fixed  $\mathbf{q}$  and  $\xi \in \mathbb{C}_+$ , we have*

$$\mathbb{E}|\bar{M}_d(\xi; \mathbf{q}, \boldsymbol{\mu}) - M_d(\xi; \mathbf{q}, \boldsymbol{\mu})| = o_d(1).$$

**Proof** [Proof of Lemma IV.3] Define

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

and write  $M_d(\xi; \mathbf{q}, \boldsymbol{\mu})$  and  $\bar{M}_d(\xi; \mathbf{q}, \boldsymbol{\mu})$  as  $M_d(\xi)$  and  $\bar{M}_d(\xi)$  to simplify the notation. Then by definition we have

$$\begin{aligned} |\Delta(\mathbf{A}, \bar{\mathbf{A}}, \xi)| &= |\text{tr}[(\mathbf{A} - \xi \mathbf{I})^{-1}(\mathbf{A} - \bar{\mathbf{A}})(\bar{\mathbf{A}} - \xi \mathbf{I})^{-1}]|/d \\ &\leq \|(\mathbf{A} - \xi \mathbf{I})^{-1}(\bar{\mathbf{A}} - \xi \mathbf{I})^{-1}\|_{\text{op}} \|\mathbf{A} - \bar{\mathbf{A}}\|_* / d \\ &\leq \|\mathbf{A} - \bar{\mathbf{A}}\|_* \cdot \frac{1}{d} \cdot \frac{1}{(\Im(\xi))^2}, \end{aligned} \quad (\text{IV.3})$$

where  $\|\cdot\|_*$  is the nuclear norm, the first inequality follows from the fact that  $\text{tr}(\mathbf{U}\mathbf{V}) \leq \|\mathbf{U}\|_{\text{op}} \|\mathbf{V}\|_*$  for all  $\mathbf{U} \in \mathbb{C}^{N \times N}$  and Hermite  $\mathbf{V} \in \mathbb{C}^{N \times N}$ , and the second inequality follows from the fact that  $\mathbf{A}$  and  $\bar{\mathbf{A}}$  are real matrices. Because

$$|M_d(\xi)| = \frac{1}{d} \left| \text{tr}(\mathbf{A} - \xi \mathbf{I})^{-1} \right| \leq \frac{P}{d} \|\mathbf{A} - \xi \mathbf{I}\|_{\text{op}} \leq P/(d \cdot \Im(\xi)),$$

$$|\overline{M}_d(\xi)| = \frac{1}{d} \left| \text{tr}(\overline{\mathbf{A}} - \xi \mathbf{I})^{-1} \right| \leq \frac{P}{d} \|\overline{\mathbf{A}} - \xi \mathbf{I}\|_{\text{op}} \leq P/(d \cdot \Im(\xi)),$$

$|\Delta(\mathbf{A}, \overline{\mathbf{A}}, \xi)|$  is deterministically upper bounded:

$$|\Delta(\mathbf{A}, \overline{\mathbf{A}}, \xi)| \leq |M_d(\xi)| + |\overline{M}_d(\xi)| \leq 2P/(d \cdot \Im(\xi)). \quad (\text{IV.4})$$

Therefore, if we can prove  $\|\mathbf{A} - \overline{\mathbf{A}}\|_*/d = o_{\mathbb{P}}(1)$ , then according to (IV.3) and (IV.4), we can conclude that  $\mathbb{E}|\Delta(\mathbf{A}, \overline{\mathbf{A}}, \xi)| = o_d(1)$  by the dominated convergence theorem. To this end, we first recall the notations in Definitions A.1 and A.3 that for  $j = 1, 2$ ,

$$\mathbf{Z}_j = \sigma_j \left( \mathbf{X} \boldsymbol{\Theta}_j^\top / \sqrt{d} \right) / \sqrt{d} \in \mathbb{R}^{n \times N_j}, \quad \tilde{\mathbf{Z}}_j = \frac{\mu_{j,1}}{d} \mathbf{X} \boldsymbol{\Theta}_j^\top.$$

We also remind readers that  $\mathbf{J}_j = \phi_j(\overline{\mathbf{X}} \overline{\boldsymbol{\Theta}}_j^\top / \sqrt{d}) / \sqrt{d}$ ,  $\tilde{\mathbf{J}}_j = \frac{\mu_{j,1}}{d} \overline{\mathbf{X}} \overline{\boldsymbol{\Theta}}_j^\top$  are the ‘‘Gaussian version’’ counterparts of  $\mathbf{Z}_j$  and  $\tilde{\mathbf{Z}}_j$  respectively. We further denote  $\mathbf{Z}_{j,0} = \mu_{j,0} \mathbf{1}_n \mathbf{1}_{N_j} / \sqrt{d}$  and let  $\mathbf{Z}_{j,\star} = \mathbf{Z}_j - \mathbf{Z}_{j,0}$  for  $j = 1, 2$ . Then by the definition of the functions  $\phi_1, \phi_2$ , clearly we have  $\mathbf{Z}_{j,\star} = \phi_j(\mathbf{X} \boldsymbol{\Theta}_j^\top / \sqrt{d}) / \sqrt{d}$  for  $j = 1, 2$ . With these notations, we can rewrite  $\mathbf{A} - \overline{\mathbf{A}}$  as follows:

$$\begin{aligned} \mathbf{A} - \overline{\mathbf{A}} = & q_5 \begin{bmatrix} \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \frac{\mathbf{X} \mathbf{X}^\top - \overline{\mathbf{X}} \overline{\mathbf{X}}^\top}{d} \end{bmatrix} + q_4 \begin{bmatrix} \frac{\mathbf{M}_1 \boldsymbol{\Theta} \boldsymbol{\Theta}^\top \mathbf{M}_1 - \mathbf{M}_1 \overline{\boldsymbol{\Theta}} \overline{\boldsymbol{\Theta}}^\top \mathbf{M}_1}{d} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \\ & + q_1 \begin{bmatrix} \mathbf{0} & [\tilde{\mathbf{Z}}_1, \tilde{\mathbf{Z}}_2]^\top - [\tilde{\mathbf{J}}_1, \tilde{\mathbf{J}}_2]^\top \\ [\tilde{\mathbf{Z}}_1, \tilde{\mathbf{Z}}_2]^\top - [\tilde{\mathbf{J}}_1, \tilde{\mathbf{J}}_2] & \mathbf{0} \end{bmatrix} + \begin{bmatrix} \mathbf{0} & [\mathbf{Z}_{1,0}, \mathbf{Z}_{2,0}]^\top \\ [\mathbf{Z}_{1,0}, \mathbf{Z}_{2,0}] & \mathbf{0} \end{bmatrix} \\ & + \begin{bmatrix} \mathbf{0} & [\mathbf{Z}_{1,\star}, \mathbf{Z}_{2,\star}]^\top - [\mathbf{J}_1, \mathbf{J}_2]^\top \\ [\mathbf{Z}_{1,\star}, \mathbf{Z}_{2,\star}] - [\mathbf{J}_1, \mathbf{J}_2] & \mathbf{0} \end{bmatrix}. \end{aligned}$$

Then by the triangle inequality and Cauchy-Schwarz inequality, we have

$$\frac{\|\mathbf{A} - \overline{\mathbf{A}}\|_*}{d} = O_{\mathbb{P}}(I_1 + I_2 + I_3 + I_4 + I_5),$$

where

$$\begin{aligned} I_1 &= \frac{1}{\sqrt{d}} \left\| \begin{bmatrix} \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \frac{\mathbf{X} \mathbf{X}^\top - \overline{\mathbf{X}} \overline{\mathbf{X}}^\top}{d} \end{bmatrix} \right\|_F, \\ I_2 &= \frac{1}{\sqrt{d}} \left\| \begin{bmatrix} \frac{\mathbf{M}_1 \boldsymbol{\Theta} \boldsymbol{\Theta}^\top \mathbf{M}_1 - \mathbf{M}_1 \overline{\boldsymbol{\Theta}} \overline{\boldsymbol{\Theta}}^\top \mathbf{M}_1}{d} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \right\|_F, \\ I_3 &= \frac{1}{\sqrt{d}} \left\| \begin{bmatrix} \mathbf{0} & [\tilde{\mathbf{Z}}_1, \tilde{\mathbf{Z}}_2]^\top - [\tilde{\mathbf{J}}_1, \tilde{\mathbf{J}}_2]^\top \\ [\tilde{\mathbf{Z}}_1, \tilde{\mathbf{Z}}_2]^\top - [\tilde{\mathbf{J}}_1, \tilde{\mathbf{J}}_2] & \mathbf{0} \end{bmatrix} \right\|_F, \\ I_4 &= \frac{1}{d} \left\| \begin{bmatrix} \mathbf{0} & [\mathbf{Z}_{1,0}, \mathbf{Z}_{2,0}]^\top \\ [\mathbf{Z}_{1,0}, \mathbf{Z}_{2,0}] & \mathbf{0} \end{bmatrix} \right\|_*, \\ I_5 &= \frac{1}{\sqrt{d}} \left\| \begin{bmatrix} \mathbf{0} & [\mathbf{Z}_{1,\star}, \mathbf{Z}_{2,\star}]^\top - [\mathbf{J}_1, \mathbf{J}_2]^\top \\ [\mathbf{Z}_{1,\star}, \mathbf{Z}_{2,\star}] - [\mathbf{J}_1, \mathbf{J}_2] & \mathbf{0} \end{bmatrix} \right\|_F. \end{aligned}$$

In the following, we bound the terms  $I_1, \dots, I_5$  separately. For  $I_1$ , let  $\mathbf{D}_\mathbf{x} = \text{diag}(\sqrt{d}/\|\bar{\mathbf{x}}_1\|_2, \dots, \sqrt{d}/\|\bar{\mathbf{x}}_n\|_2)$ . Then we have  $\mathbf{X} = \mathbf{D}_\mathbf{x} \bar{\mathbf{X}}$  by (IV.1), and

$$\begin{aligned}
 I_1 &= \frac{1}{\sqrt{d}} \left\| \frac{\mathbf{X} \mathbf{X}^\top - \bar{\mathbf{X}} \bar{\mathbf{X}}^\top}{d} \right\|_F \leq \left\| \frac{\mathbf{X} \mathbf{X}^\top - \bar{\mathbf{X}} \bar{\mathbf{X}}^\top}{d} \right\|_{\text{op}} = \left\| \frac{\mathbf{D}_\mathbf{x} \bar{\mathbf{X}} \bar{\mathbf{X}}^\top \mathbf{D}_\mathbf{x} - \bar{\mathbf{X}} \bar{\mathbf{X}}^\top}{d} \right\|_{\text{op}} \\
 &= \left\| \frac{(\mathbf{D}_\mathbf{x} - \mathbf{I}_n) \bar{\mathbf{X}} \bar{\mathbf{X}}^\top (\mathbf{D}_\mathbf{x} + \mathbf{I}_n) + \bar{\mathbf{X}} \bar{\mathbf{X}}^\top \mathbf{D}_\mathbf{x} - \mathbf{D}_\mathbf{x} \bar{\mathbf{X}} \bar{\mathbf{X}}^\top}{d} \right\|_{\text{op}} \\
 &\leq \|\mathbf{D}_\mathbf{x} - \mathbf{I}_n\|_{\text{op}} \cdot \left\| \frac{\bar{\mathbf{X}} \bar{\mathbf{X}}^\top}{d} \right\|_{\text{op}} \cdot (1 + \|\mathbf{D}_\mathbf{x}\|_{\text{op}}) + \left\| \frac{\bar{\mathbf{X}} \bar{\mathbf{X}}^\top \mathbf{D}_\mathbf{x} - \mathbf{D}_\mathbf{x} \bar{\mathbf{X}} \bar{\mathbf{X}}^\top}{d} \right\|_{\text{op}} \\
 &= \left\| \frac{\bar{\mathbf{X}} \bar{\mathbf{X}}^\top \mathbf{D}_\mathbf{x} - \mathbf{D}_\mathbf{x} \bar{\mathbf{X}} \bar{\mathbf{X}}^\top}{d} \right\|_{\text{op}} + o_{\mathbb{P}}(1),
 \end{aligned} \tag{IV.5}$$

where the first inequality holds since the average of the  $d$  squared eigenvalues of  $(\mathbf{X} \mathbf{X}^\top - \bar{\mathbf{X}} \bar{\mathbf{X}}^\top)/d$  is bounded by the largest one of them, and the last equality follows from  $\|\mathbf{D}_\mathbf{x} - \mathbf{I}_n\|_{\text{op}} = O_{\mathbb{P}}\left(\sqrt{\frac{\log d}{d}}\right)$  and  $\|\mathbf{D}_\mathbf{x}\|_{\text{op}} = O_{\mathbb{P}}(1)$ , which are direct consequences of the definition of  $\mathbf{D}_\mathbf{x}$ . We further let  $\tilde{\mathbf{D}}_\mathbf{x}$  be the matrix whose elements  $(\tilde{\mathbf{D}}_\mathbf{x})_{ij}$  satisfy  $(\tilde{\mathbf{D}}_\mathbf{x})_{ij} = (\mathbf{D}_\mathbf{x})_{jj} - (\mathbf{D}_\mathbf{x})_{ii}$  for  $i, j \in [n]$ . Then we have  $\|\tilde{\mathbf{D}}_\mathbf{x}\|_{\text{max}} = o_{\mathbb{P}}(1)$ , and

$$\left\| \frac{\bar{\mathbf{X}} \bar{\mathbf{X}}^\top \mathbf{D}_\mathbf{x} - \mathbf{D}_\mathbf{x} \bar{\mathbf{X}} \bar{\mathbf{X}}^\top}{d} \right\|_{\text{op}} = \left\| \tilde{\mathbf{D}}_\mathbf{x} \odot \frac{\bar{\mathbf{X}} \bar{\mathbf{X}}^\top}{d} \right\|_{\text{op}} \leq \|\tilde{\mathbf{D}}_\mathbf{x}\|_{\text{max}} \cdot \left\| \frac{\bar{\mathbf{X}} \bar{\mathbf{X}}^\top}{d} \right\|_{\text{op}} = o_{\mathbb{P}}(1). \tag{IV.6}$$

Plugging (IV.6) into (IV.5) completes the proof of  $I_1 = o_{\mathbb{P}}(1)$ . Similarly, it can be shown that  $I_2$  and  $I_3$  are both  $o_{\mathbb{P}}(1)$ . For  $I_4$ , by the definition that  $\mathbf{Z}_{j,0} = \mu_{j,0} \mathbf{1}_n \mathbf{1}_{N_j} / \sqrt{d}$ ,  $j = 1, 2$ , it is clear that  $\mathbf{Z}_{j,0}$  is rank-one and  $\|\mathbf{Z}_{j,0}\|_{\text{op}} = O_d(\sqrt{d})$ . Therefore we have

$$I_4 = \frac{1}{d} \left\| \begin{bmatrix} \mathbf{0} & [\mathbf{Z}_{1,0}, \mathbf{Z}_{2,0}]^\top \\ [\mathbf{Z}_{1,0}, \mathbf{Z}_{2,0}] & \mathbf{0} \end{bmatrix} \right\|_{\star} = o_d(1).$$

Finally, to prove  $I_5 = o_d(1)$ , it clearly suffices to show that

$$\frac{1}{\sqrt{d}} \|\mathbf{Z}_{j,\star} - \mathbf{J}_j\|_F = o_{\mathbb{P}}(1), \quad j = 1, 2.$$

Define  $\bar{\mathbf{Z}}_{j,\star} = \phi_j(\mathbf{X} \bar{\boldsymbol{\Theta}}_j^\top / \sqrt{d}) / \sqrt{d} = \phi_j(\mathbf{D}_\mathbf{x} \bar{\mathbf{X}} \bar{\boldsymbol{\Theta}}_j^\top / \sqrt{d}) / \sqrt{d}$ ,  $j \in [N]$  and  $r_i = \sqrt{d} / \|\bar{\mathbf{x}}_i\|_2$ ,  $i \in [n]$ . By the mean value theorem, for  $j = 1, 2$ ,  $a \in [N_j]$  and  $i \in [n]$ , there exists  $\zeta_{ia}$  between  $r_i$  and 1, such that

$$\begin{aligned}
 \bar{\mathbf{Z}}_{j,\star} - \mathbf{J}_j &= [\phi_j(r_i \langle \bar{\mathbf{x}}_i, \bar{\boldsymbol{\theta}}_a \rangle / \sqrt{d}) / \sqrt{d} - \phi_j(\langle \bar{\mathbf{x}}_i, \bar{\boldsymbol{\theta}}_a \rangle / \sqrt{d}) / \sqrt{d}]_{i \in [n], a \in [N_j]} \\
 &= [(r_i - 1) \langle \bar{\mathbf{x}}_i, \bar{\boldsymbol{\theta}}_a \rangle / \sqrt{d}] \phi'_j(\zeta_{ij} \langle \bar{\mathbf{x}}_i, \bar{\boldsymbol{\theta}}_a \rangle / \sqrt{d}) / \sqrt{d}]_{i \in [n], a \in [N_j]} \\
 &= (\mathbf{D}_\mathbf{x} - \mathbf{I}_n) \bar{\phi}_j(\boldsymbol{\zeta} \odot (\bar{\mathbf{X}} \bar{\boldsymbol{\Theta}}_j^\top / \sqrt{d})) / \sqrt{d},
 \end{aligned}$$

where  $\boldsymbol{\zeta} = (\zeta_{ij})_{i \in [n], a \in [N_j]}$  and  $\bar{\phi}_j(x) = x \phi'_j(x)$ . By Bernstein-type concentration inequalities (Ver-shynin, 2010), we have

$$\|\mathbf{D}_\mathbf{x} - \mathbf{I}_n\|_{\text{op}} = O_{\mathbb{P}}\left(\sqrt{\frac{\log d}{d}}\right), \quad \|\boldsymbol{\zeta}\|_{\text{max}} = O_{\mathbb{P}}(1), \quad \|\bar{\mathbf{X}} \bar{\boldsymbol{\Theta}}_j^\top / \sqrt{d}\|_{\text{max}} = O_{\mathbb{P}}(\sqrt{\log d}).$$

Moreover, note that we currently assume that the activation functions  $\sigma_j$ ,  $j = 1, 2$  are fixed polynomials, which implies that  $\phi_j$  are also fixed polynomials. Therefore, there exists a constant  $M_0 \in \mathbb{N}$  such that

$$\|\bar{\mathbf{Z}}_{j,*} - \mathbf{J}_j\|_F / \sqrt{d} \leq \|\mathbf{D}_x - \mathbf{I}_n\|_{\text{op}} \|\bar{\phi}_j(\zeta \odot (\bar{\mathbf{X}} \bar{\boldsymbol{\Theta}}_j^\top / \sqrt{d}))\|_F / d = O_{\mathbb{P}}((\log d)^{M_0} / \sqrt{d}) = o_{\mathbb{P}}(1).$$

With exactly the same argument, it can be shown that  $\|\bar{\mathbf{Z}}_{j,*} - \mathbf{Z}_{j,*}\|_F / \sqrt{d} = o_{\mathbb{P}}(1)$  (recall  $\mathbf{Z}_{j,*} = \phi_j(\mathbf{X} \boldsymbol{\Theta}^\top / \sqrt{d}) / \sqrt{d}$ ). Therefore we have

$$\frac{1}{\sqrt{d}} \|\mathbf{Z}_{j,*} - \mathbf{J}_j\|_F \leq \frac{1}{\sqrt{d}} \|\bar{\mathbf{Z}}_{j,*} - \mathbf{Z}_{j,*}\|_F + \frac{1}{\sqrt{d}} \|\bar{\mathbf{Z}}_{j,*} - \mathbf{J}_j\|_F = o_{\mathbb{P}}(1)$$

for  $j = 1, 2$ . Finally  $I_5 = o_{\mathbb{P}}(1)$  and the proof of Lemma IV.3 is complete.  $\blacksquare$

## IV.2 Calculation of the resolvent equations

Lemma IV.3 and its proof show the readers that the Stieltjes transforms of the empirical eigenvalue distributions of  $\mathbf{A}$  and  $\bar{\mathbf{A}}$  share the same asymptotics. Based on this result, we can equivalently consider the ‘‘Gaussian version’’ counterpart of the learning problem. Therefore, throughout Section IV.2, we directly consider the matrices  $\bar{\mathbf{X}}$  and  $\bar{\boldsymbol{\Theta}}$ , whose elements are independently generated from standard normal  $N(0, 1)$ . In addition, the activation functions for the two types of random features are  $\phi_j(x) = \sigma_j(x) - \mu_{j,0}$ ,  $j = 1, 2$ , and the linear pencil matrix is  $\bar{\mathbf{A}}(\mathbf{q}, \boldsymbol{\mu})$  is given in Definition IV.2.

Moreover, for  $j = 1, 2$ , let  $\Phi_j(x) = \phi_j(x) + q_1 \mu_{j,1} x$ , and it is easy to see  $\mathbf{J}_j^\top + q_1 \tilde{\mathbf{J}}_j^\top = \Phi_j(\bar{\mathbf{X}} \bar{\boldsymbol{\Theta}}_j^\top / \sqrt{d}) / \sqrt{d}$ . We further denote  $\phi_{j,0} \triangleq \mathbb{E}_{G \sim N(0,1)} \{\Phi_j(G)\}$ ,  $\phi_{j,1} \triangleq \mathbb{E}_{G \sim N(0,1)} \{G \Phi_j(G)\}$ ,  $\phi_{j,*} \triangleq \mathbb{E}_{G \sim N(0,1)} \{\Phi_j(G)^2\} - \phi_{j,0}^2 - \phi_{j,1}^2$ . By these definitions, it is easy to see that  $\phi_{j,0} = 0$ ,  $\phi_{j,1}^2 = \mu_{j,1}^2 (1 + q_1)^2$ ,  $\phi_{j,*} = \mu_{j,*}^2$ . Importantly, the property that  $\mathbb{E}_{G \sim N(0,1)} \{\Phi_j(G)\} = \phi_{j,0} = 0$  enables the application of the following lemma, which is summarized from Section 4.3, Step 2 in Cheng and Singer (2013).

**Lemma IV.4.** *Suppose that  $\Phi$  is a polynomial satisfying  $\mathbb{E}_{G \sim N(0,1)} \{\Phi(G)\} = 0$ ,  $\mathbb{E}_{G \sim N(0,1)} \{G \Phi(G)\} = 0$  and  $\bar{\mathbf{x}}_i, \bar{\boldsymbol{\theta}}_a \in \mathbb{R}^d$ ,  $i \in [n]$ ,  $a \in [N]$  are standard Gaussian vectors. Define matrix  $\mathbf{E} \in \mathbb{R}^{n \times N}$  elementwisely as*

$$(\mathbf{E}_j)_{i,a} = \frac{1}{\sqrt{d}} \left[ \Phi \left( \frac{1}{\sqrt{d}} \langle \bar{\mathbf{x}}_i, \bar{\boldsymbol{\theta}}_a \rangle \right) - \Phi \left( \frac{1}{\sqrt{d}} \langle (\bar{\mathbf{x}}_i)_{[1:d-1]}, (\bar{\boldsymbol{\theta}}_a)_{[1:d-1]} \rangle \right) \right]$$

for  $i \in [n]$ ,  $a \in [N]$ . Then  $\|\mathbf{E}\|_{\text{op}} = o_{\mathbb{P}}(1)$ .

Lemma IV.4 formally shows the intuitive result that under the setting where  $d, n$  grows proportionally, removing one entry in the random vectors does not change the asymptotic limit of polynomials. This enables us to apply the standard leave-one-out argument in random matrix theory.



Our goal in this part of the proof is to calculate the resolvent equations of the Stieltjes transforms corresponding to the pencil matrix  $\overline{\mathbf{A}}(\mathbf{q}, \boldsymbol{\mu})$ . To do so, we define the following terms:

$$\begin{aligned}\overline{m}_{1,d}(\xi; \mathbf{q}, \boldsymbol{\mu}) &= \mathbb{E}[\overline{M}_{1,d}(\xi; \mathbf{q}, \boldsymbol{\mu})], & \overline{M}_{1,d}(\xi; \mathbf{q}, \boldsymbol{\mu}) &= \frac{1}{d} \text{tr}_{[1:N_1]}[(\overline{\mathbf{A}}(\mathbf{q}, \boldsymbol{\mu}) - \xi \mathbf{I}_P)^{-1}], \\ \overline{m}_{2,d}(\xi; \mathbf{q}, \boldsymbol{\mu}) &= \mathbb{E}[\overline{M}_{2,d}(\xi; \mathbf{q}, \boldsymbol{\mu})], & \overline{M}_{2,d}(\xi; \mathbf{q}, \boldsymbol{\mu}) &= \frac{1}{d} \text{tr}_{[N_1+1:N]}[(\overline{\mathbf{A}}(\mathbf{q}, \boldsymbol{\mu}) - \xi \mathbf{I}_P)^{-1}], \\ \overline{m}_{3,d}(\xi; \mathbf{q}, \boldsymbol{\mu}) &= \mathbb{E}[\overline{M}_{3,d}(\xi; \mathbf{q}, \boldsymbol{\mu})], & \overline{M}_{3,d}(\xi; \mathbf{q}, \boldsymbol{\mu}) &= \frac{1}{d} \text{tr}_{[N+1:P]}[(\overline{\mathbf{A}}(\mathbf{q}, \boldsymbol{\mu}) - \xi \mathbf{I}_P)^{-1}].\end{aligned}$$

Standard argument in random matrix theory then gives us the concentration result

$$\mathbb{E}|\overline{M}_{i,d}(\xi; \mathbf{q}, \boldsymbol{\mu}) - \overline{m}_{i,d}(\xi; \mathbf{q}, \boldsymbol{\mu})| = o_d(1)$$

for any fixed  $\xi \in \mathbb{C}_+$ . Therefore, denoting  $\overline{m}_d(\xi) = \sum_{i=1}^3 \overline{m}_{i,d}(\xi)$  and  $\overline{M}_d(\xi) = \sum_{i=1}^3 \overline{M}_{i,d}(\xi)$ , (we drop the argument  $\mathbf{q}, \boldsymbol{\mu}$  for simplicity), we have

$$\mathbb{E}|\overline{M}_d(\xi) - \overline{m}_d(\xi)| = o_d(1) \tag{IV.7}$$

for any fixed  $\xi \in \mathbb{C}_+$ . A proof of this concentration can be found in [Hastie et al. \(2022\)](#); [Mei and Montanari \(2022\)](#). Based on (IV.7), to study  $\overline{M}_d(\xi; \mathbf{q}, \boldsymbol{\mu})$ , which is the Stieltjes transform of the empirical eigenvalue distribution of  $\overline{\mathbf{A}}(\mathbf{q}, \boldsymbol{\mu})$ , it suffices to derive the resolvent equations for  $\overline{m}_d(\xi; \mathbf{q}, \boldsymbol{\mu})$ . This is done in the following lemma.

**Lemma IV.5.** *Let  $\overline{\mathbf{m}}_d(\xi) = [\overline{m}_{1,d}(\xi), \overline{m}_{2,d}(\xi), \overline{m}_{3,d}(\xi)]^\top$ . Then for any fixed  $\xi \in \mathbb{C}_+$ , the following property holds:*

$$\|\overline{\mathbf{m}}_d(\xi) - \mathbf{F}(\overline{\mathbf{m}}_d(\xi))\|_2 = o_d(1).$$

**Proof** [Proof of Lemma IV.5] Since  $\overline{\mathbf{m}}_d(\xi), \mathbf{F}(\overline{\mathbf{m}}_d(\xi)) \in \mathbb{C}^3$ , Lemma IV.5 essentially contains three results showing that the first, second, and third elements of  $\overline{\mathbf{m}}_d(\xi) - \mathbf{F}(\overline{\mathbf{m}}_d(\xi))$  are all asymptotically zero. Since the proofs of the three results are almost the same, we mainly focus on the proof of the first result. The proof consists of three steps. The first step is to use the Schur complement formula to calculate  $\overline{m}_{1,d}$ . The second step is to simplify the formula of  $\overline{m}_{1,d}$ . The third step is to give the recursive equations of  $\overline{m}_{1,d}$  based on the result of step 2.

**Step 1.** We first use a leave-one-out argument to calculate  $\overline{m}_{1,d}$ . Let  $\overline{\mathbf{A}}_{\cdot, N_1} \in \mathbb{R}^{P-1}$  be the  $N_1^{\text{th}}$  column of  $\overline{\mathbf{A}}$ , with the  $N_1^{\text{th}}$  entry removed. We further denote by  $\overline{\mathbf{B}} \in \mathbb{R}^{(P-1) \times (P-1)}$  the sub-matrix of  $\overline{\mathbf{A}}$  obtained by removing the  $N_1^{\text{th}}$  row and  $N_1^{\text{th}}$  column in  $\overline{\mathbf{A}}$ . We can then treat  $\overline{\mathbf{A}}$  as a  $2 \times 2$  block matrix formed by  $\overline{\mathbf{A}}_{\cdot, N_1}$ ,  $\overline{\mathbf{A}}_{\cdot, N_1}^\top$ ,  $\overline{\mathbf{B}}$ , and  $\overline{\mathbf{A}}_{N_1, N_1} = q_2 \mu_{1,*}^2 + q_4 \mu_{1,1}^2 \|\overline{\boldsymbol{\theta}}_{N_1}\|_2^2 / d$ . Then by the Schur complement formula, we get

$$\overline{m}_{1,d} = \psi_1 \mathbb{E} \left( -\xi + q_2 \mu_{1,*}^2 + q_4 \mu_{1,1}^2 \|\overline{\boldsymbol{\theta}}_{N_1}\|_2^2 / d - \overline{\mathbf{A}}_{\cdot, N_1}^\top (\overline{\mathbf{B}} - \xi \mathbf{I}_{P-1})^{-1} \overline{\mathbf{A}}_{\cdot, N_1} \right)^{-1}. \tag{IV.8}$$

We decompose the vectors  $\bar{\boldsymbol{\theta}}_a$ ,  $a \in [N]$  and  $\bar{\mathbf{x}}_i$ ,  $i \in [n]$  into components along the direction of  $\bar{\boldsymbol{\theta}}_{N_1}$  and other orthogonal directions:

$$\begin{aligned}\bar{\boldsymbol{\theta}}_a &= \eta_a \frac{\bar{\boldsymbol{\theta}}_{N_1}}{\|\bar{\boldsymbol{\theta}}_{N_1}\|} + \tilde{\boldsymbol{\theta}}_a, \quad \langle \bar{\boldsymbol{\theta}}_{N_1}, \tilde{\boldsymbol{\theta}}_a \rangle = 0, \quad a \in [N] \setminus \{N_1\}, \\ \bar{\mathbf{x}}_i &= u_i \frac{\bar{\boldsymbol{\theta}}_{N_1}}{\|\bar{\boldsymbol{\theta}}_{N_1}\|} + \tilde{\mathbf{x}}_i, \quad \langle \bar{\boldsymbol{\theta}}_{N_1}, \tilde{\mathbf{x}}_i \rangle = 0, \quad i \in [n].\end{aligned}\tag{IV.9}$$

Note that for any  $a \in [N] \setminus \{N_1\}$  and  $i \in [n]$ ,  $\eta_a$ ,  $u_i$  are standard Gaussian and are independent of  $\tilde{\boldsymbol{\theta}}_a$  and  $\tilde{\mathbf{x}}_i$ . Moreover,  $\tilde{\boldsymbol{\theta}}_a$  and  $\tilde{\mathbf{x}}_i$  are conditionally independent on each other given  $\bar{\boldsymbol{\theta}}_{N_1}$ , with  $\tilde{\boldsymbol{\theta}}_a, \tilde{\mathbf{x}}_i \sim N(0, P_\perp)$ , where  $P_\perp$  is the projector orthogonal to  $\bar{\boldsymbol{\theta}}_{N_1}$ . We can then use the coefficients  $\eta_a$ ,  $a \in [N] \setminus \{N_1\}$  and  $u_i$ ,  $i \in [n]$  to represent the entries of  $\bar{\mathbf{A}}_{\cdot, N_1}$ . We have  $\bar{\mathbf{A}}_{\cdot, N_1} = [\bar{\mathbf{A}}_{1, N_1}, \dots, \bar{\mathbf{A}}_{P-1, N_1}]^\top \in \mathbb{R}^{P-1}$  with

$$\bar{\mathbf{A}}_{i, N_1} = \begin{cases} \frac{q_4 \mu_{1,1}^2 \eta_i}{d} \|\bar{\boldsymbol{\theta}}_{N_1}\|_2, & \text{if } i \in [1, N_1 - 1], \\ \frac{q_4 \mu_{1,1} \mu_{2,1} \eta_{i+1}}{d} \|\bar{\boldsymbol{\theta}}_{N_1}\|_2, & \text{if } i \in [N_1, N - 1], \\ \frac{1}{\sqrt{d}} \Phi_1\left(\frac{1}{\sqrt{d}} u_{i-N+1} \|\bar{\boldsymbol{\theta}}_{N_1}\|_2\right), & \text{if } i \geq N. \end{cases}\tag{IV.10}$$

To calculate the resolvent equations, we need to further represent the matrix  $\bar{\mathbf{B}}$  in (IV.8) with  $\eta_a$ ,  $\tilde{\boldsymbol{\theta}}_a$ ,  $u_i$ , and  $\tilde{\mathbf{x}}_i$  for  $a \in [N] \setminus \{N_1\}$  and  $i \in [n]$ . Below we first list some additional notations for easier reference. We write  $\boldsymbol{\eta}_1 = [\eta_1, \dots, \eta_{N_1-1}] \in \mathbb{R}^{N_1-1}$ ,  $\boldsymbol{\eta}_2 = [\eta_{N_1+1}, \dots, \eta_N] \in \mathbb{R}^{N_2}$ ,  $\boldsymbol{\eta} = [\boldsymbol{\eta}_1^\top, \boldsymbol{\eta}_2^\top]^\top \in \mathbb{R}^{N-1}$ ,  $\mathbf{u} = [u_1, \dots, u_n]^\top \in \mathbb{R}^n$ ,  $\tilde{\boldsymbol{\theta}}_1 = [\tilde{\boldsymbol{\theta}}_1, \dots, \tilde{\boldsymbol{\theta}}_{N_1-1}]^\top$ ,  $\tilde{\boldsymbol{\theta}}_2 = [\tilde{\boldsymbol{\theta}}_{N_1+1}, \dots, \tilde{\boldsymbol{\theta}}_N]^\top$ ,  $\tilde{\boldsymbol{\theta}} = \begin{bmatrix} \tilde{\boldsymbol{\theta}}_1 \\ \tilde{\boldsymbol{\theta}}_2 \end{bmatrix} \in \mathbb{R}^{(N-1) \times d}$ ,  $\tilde{\mathbf{M}}_1 = \begin{bmatrix} \mu_{1,1} \mathbf{I}_{N_1-1} & \\ & \mu_{2,1} \mathbf{I}_{N_2} \end{bmatrix}$  and  $\tilde{\mathbf{M}}_* = \begin{bmatrix} \mu_{1,*} \mathbf{I}_{N_1-1} & \\ & \mu_{2,*} \mathbf{I}_{N_2} \end{bmatrix}$ . Now with (IV.9) and the notations above, we can decompose  $\bar{\mathbf{B}}_{[1:N-1], [1:N-1]}$  as follows:

$$\bar{\mathbf{B}}_{[1:N-1], [1:N-1]} = q_2 \tilde{\mathbf{M}}_* \tilde{\mathbf{M}}_* + \frac{q_4}{d} \tilde{\mathbf{M}}_1 \tilde{\boldsymbol{\theta}} \tilde{\boldsymbol{\theta}}^\top \tilde{\mathbf{M}}_1 + \frac{q_4}{d} \tilde{\mathbf{M}}_1 \boldsymbol{\eta} \boldsymbol{\eta}^\top \tilde{\mathbf{M}}_1.\tag{IV.11}$$

Moreover, for  $i, j \in [n]$  we define

$$(\tilde{\mathbf{H}})_{ij} = \frac{1}{d} \langle \tilde{\mathbf{x}}_i, \tilde{\mathbf{x}}_j \rangle.$$

Then we can decompose  $\bar{\mathbf{B}}_{[N:P-1], [N:P-1]}$  into

$$\bar{\mathbf{B}}_{[N:P-1], [N:P-1]} = q_3 \mathbf{I}_n + q_5 \tilde{\mathbf{H}} + \frac{q_5}{d} \mathbf{u} \mathbf{u}^\top.\tag{IV.12}$$

For  $\bar{\mathbf{B}}_{[N:P-1], [1:N-1]}$ , by definition we see that the elements in  $\bar{\mathbf{B}}_{[N:P-1], [1:N-1]}$  are  $(\mathbf{Z})_{i,a}$  for  $a \in [N] \setminus \{N_1\}$  and  $i \in [n]$ . Therefore, we have

$$\begin{aligned}(\mathbf{Z})_{i,a} &= \frac{1}{\sqrt{d}} \Phi_j\left(\frac{1}{\sqrt{d}} \langle \bar{\mathbf{x}}_i, \bar{\boldsymbol{\theta}}_a \rangle\right) = \frac{1}{\sqrt{d}} \Phi_j\left(\frac{1}{\sqrt{d}} \langle \tilde{\mathbf{x}}_i, \tilde{\boldsymbol{\theta}}_a \rangle + \frac{1}{d} u_i \eta_a\right) \\ &= \frac{1}{\sqrt{d}} \Phi_j\left(\frac{1}{\sqrt{d}} \langle \tilde{\mathbf{x}}_i, \tilde{\boldsymbol{\theta}}_a \rangle\right) + \frac{\phi_{j,1}}{d} u_i \eta_a + \frac{1}{\sqrt{d}} \left[ \Phi_{j,\perp}\left(\frac{1}{\sqrt{d}} \langle \tilde{\mathbf{x}}_i, \tilde{\boldsymbol{\theta}}_a \rangle + \frac{1}{\sqrt{d}} u_i \eta_a\right) - \Phi_{j,\perp}\left(\frac{1}{\sqrt{d}} \langle \tilde{\mathbf{x}}_i, \tilde{\boldsymbol{\theta}}_a \rangle\right) \right],\end{aligned}$$

where  $\Phi_{j,\perp}(x) = \Phi_j(x) - \phi_{j,1}x$ ,  $j = 1$  when  $a \leq N_1 - 1$  and  $j = 2$  when  $a \geq N_1 + 1$ . By the symmetry of  $\bar{\mathbf{B}}$ , we can then decompose  $\bar{\mathbf{B}}_{[N:P-1],[1:N-1]}$  and  $\bar{\mathbf{B}}_{[1:N-1],[N:P-1]}^\top$  into

$$\bar{\mathbf{B}}_{[N:P-1],[1:N-1]} = \bar{\mathbf{B}}_{[1:N-1],[N:P-1]}^\top = \tilde{\mathbf{Z}} + \frac{1}{d}\mathbf{u}\boldsymbol{\eta}\mathbf{M}_\phi + [\mathbf{E}_1, \mathbf{E}_2], \quad (\text{IV.13})$$

where we define

$$\begin{aligned} \tilde{\mathbf{Z}} &= [\tilde{\mathbf{Z}}_1, \tilde{\mathbf{Z}}_2], \quad (\tilde{\mathbf{Z}}_1)_{i,a} = \frac{1}{\sqrt{d}}\Phi_1\left(\frac{1}{\sqrt{d}}\langle \tilde{\mathbf{x}}_i, \tilde{\boldsymbol{\theta}}_a \rangle\right), \quad (\tilde{\mathbf{Z}}_2)_{i,a} = \frac{1}{\sqrt{d}}\Phi_2\left(\frac{1}{\sqrt{d}}\langle \tilde{\mathbf{x}}_i, \tilde{\boldsymbol{\theta}}_a \rangle\right), \\ \mathbf{M}_\phi &= \begin{bmatrix} \phi_{1,1}\mathbf{I}_{N_1-1} & \\ & \phi_{2,1}\mathbf{I}_{N_2} \end{bmatrix}, \quad (\mathbf{E}_j)_{i,a} = \frac{1}{\sqrt{d}}\left[\Phi_{j,\perp}\left(\frac{1}{\sqrt{d}}\langle \tilde{\mathbf{x}}_i, \tilde{\boldsymbol{\theta}}_a \rangle\right) + \frac{1}{\sqrt{d}}u_i\eta_a\right] - \Phi_{j,\perp}\left(\frac{1}{\sqrt{d}}\langle \tilde{\mathbf{x}}_i, \tilde{\boldsymbol{\theta}}_a \rangle\right) \end{aligned}$$

for  $a \in [N] \setminus \{N_1\}$ ,  $i \in [n]$  and  $j \in \{1, 2\}$ . Combining (IV.11), (IV.12) and (IV.13), we decompose  $\bar{\mathbf{B}}$  into

$$\bar{\mathbf{B}} = \tilde{\mathbf{B}} + \boldsymbol{\Delta} + \mathbf{E} \in \mathbb{R}^{(P-1) \times (P-1)}, \quad (\text{IV.14})$$

where

$$\begin{aligned} \tilde{\mathbf{B}} &= \begin{bmatrix} q_2\tilde{\mathbf{M}}_*\tilde{\mathbf{M}}_* + \frac{q_4}{d}\tilde{\mathbf{M}}_1\tilde{\boldsymbol{\Theta}}\tilde{\boldsymbol{\Theta}}^\top\tilde{\mathbf{M}}_1 & \tilde{\mathbf{Z}}^\top \\ \tilde{\mathbf{Z}} & q_3\mathbf{I}_n + q_5\tilde{\mathbf{H}} \end{bmatrix} \\ &= \begin{bmatrix} q_2\mu_{1,*}^2\mathbf{I}_{N_1-1} + \frac{q_4\mu_{1,1}^2}{d}\tilde{\boldsymbol{\Theta}}_1\tilde{\boldsymbol{\Theta}}_1^\top & \frac{q_4\mu_{1,1}\mu_{2,1}}{d}\tilde{\boldsymbol{\Theta}}_1\tilde{\boldsymbol{\Theta}}_2^\top & \tilde{\mathbf{Z}}_1^\top \\ \frac{q_4\mu_{1,1}\mu_{2,1}}{d}\tilde{\boldsymbol{\Theta}}_2\tilde{\boldsymbol{\Theta}}_1^\top & q_2\mu_{2,*}^2\mathbf{I}_{N_2} + \frac{q_4\mu_{2,1}^2}{d}\tilde{\boldsymbol{\Theta}}_2\tilde{\boldsymbol{\Theta}}_2^\top & \tilde{\mathbf{Z}}_2^\top \\ \tilde{\mathbf{Z}}_1 & \tilde{\mathbf{Z}}_2 & q_3\mathbf{I}_n + q_5\tilde{\mathbf{H}} \end{bmatrix}, \\ \boldsymbol{\Delta} &= \begin{bmatrix} \frac{q_4}{d}\tilde{\mathbf{M}}_1\boldsymbol{\eta}\boldsymbol{\eta}^\top\tilde{\mathbf{M}}_1 & \frac{1}{d}\mathbf{M}_\phi\boldsymbol{\eta}\mathbf{u}^\top \\ \frac{1}{d}\mathbf{u}\boldsymbol{\eta}\mathbf{M}_\phi & \frac{q_5}{d}\mathbf{u}\mathbf{u}^\top \end{bmatrix} = \begin{bmatrix} \frac{q_4\mu_{1,1}^2}{d}\boldsymbol{\eta}_1\boldsymbol{\eta}_1^\top & \frac{q_4\mu_{1,1}\mu_{2,1}}{d}\boldsymbol{\eta}_1\boldsymbol{\eta}_2^\top & \frac{\phi_{1,1}}{d}\boldsymbol{\eta}_1\mathbf{u}^\top \\ \frac{q_4\mu_{1,1}\mu_{2,1}}{d}\boldsymbol{\eta}_2\boldsymbol{\eta}_1^\top & \frac{q_4\mu_{2,1}^2}{d}\boldsymbol{\eta}_2\boldsymbol{\eta}_2^\top & \frac{\phi_{2,1}}{d}\boldsymbol{\eta}_2\mathbf{u}^\top \\ \frac{\phi_{1,1}}{d}\mathbf{u}\boldsymbol{\eta}_1^\top & \frac{\phi_{2,1}}{d}\mathbf{u}\boldsymbol{\eta}_2^\top & \frac{q_5}{d}\mathbf{u}\mathbf{u}^\top \end{bmatrix}, \quad \mathbf{E} = \begin{bmatrix} \mathbf{0} & \mathbf{0} & \mathbf{E}_1^\top \\ \mathbf{0} & \mathbf{0} & \mathbf{E}_2^\top \\ \mathbf{E}_1 & \mathbf{E}_2 & \mathbf{0} \end{bmatrix}. \end{aligned}$$

Clearly, by the definition of  $\tilde{\mathbf{B}}$ , the Stieltjes transform corresponding to  $\tilde{\mathbf{B}}$  shares the same asymptotics as the Stieltjes transform corresponding to  $\bar{\mathbf{A}}$ .

**Step 2.** According to our analysis in **Step 1**, we can then calculate  $\bar{m}_{1,d}$  by (IV.8), in which the terms  $\bar{\mathbf{A}}_{\cdot,N_1}$  and  $\bar{\mathbf{B}}$  have the decompositions (IV.10) and (IV.14) respectively. In this step, we aim to further simplify the calculation by getting rid of the terms  $\|\bar{\boldsymbol{\theta}}_{N_1}\|_2^2/d$  in (IV.8) and  $\mathbf{E}$  in (IV.14). Define

$$\begin{aligned} w_0 &= \left(-\xi + q_2\mu_{1,*}^2 + q_4\mu_{1,1}^2 - \bar{\mathbf{A}}_{\cdot,N_1}^\top(\bar{\mathbf{B}} - \xi\mathbf{I}_{P-1})^{-1}\bar{\mathbf{A}}_{\cdot,N_1}\right)^{-1}, \\ w_1 &= \left(-\xi + q_2\mu_{1,*}^2 + q_4\mu_{1,1}^2\|\bar{\boldsymbol{\theta}}_{N_1}\|_2^2/d - \bar{\mathbf{A}}_{\cdot,N_1}^\top(\bar{\mathbf{B}} - \xi\mathbf{I}_{P-1})^{-1}\bar{\mathbf{A}}_{\cdot,N_1}\right)^{-1}, \\ w_2 &= \left(-\xi + q_2\mu_{1,*}^2 + q_4\mu_{1,1}^2 - \bar{\mathbf{A}}_{\cdot,N_1}^\top(\tilde{\mathbf{B}} + \boldsymbol{\Delta} - \xi\mathbf{I}_{P-1})^{-1}\bar{\mathbf{A}}_{\cdot,N_1}\right)^{-1}. \end{aligned}$$

Then by (IV.8), we have  $\bar{m}_{1,d} = \psi_1\mathbb{E}w_1$ . We now give an upper bound of  $|w_1 - w_2|$ . Recall that we consider a fixed  $\xi \in \mathbb{C}_+$ . Since  $\bar{\mathbf{B}}$  is a real symmetric matrix, by diagonalizing  $\bar{\mathbf{B}}$ , it is easy to see that  $\Im(\bar{\mathbf{A}}_{\cdot,N_1}^\top(\bar{\mathbf{B}} - \xi\mathbf{I}_{P-1})^{-1}\bar{\mathbf{A}}_{\cdot,N_1}) \geq 0$ . Therefore, we deterministically have

$$\Im(-w_1^{-1}) = \Im(\xi) + \Im(\bar{\mathbf{A}}_{\cdot,N_1}^\top(\bar{\mathbf{B}} - \xi\mathbf{I}_{P-1})^{-1}\bar{\mathbf{A}}_{\cdot,N_1}) \geq \Im(\xi).$$

Thus we have  $|w_1| \leq 1/\Im(\xi)$ . Using a similar argument, we have  $\max\{|w_0|, |w_1|, |w_2|\} \leq 1/\Im(\xi)$ , which indicates that  $|w_1 - w_2| \leq 2/\Im(\xi)$ . Moreover, we have

$$\begin{aligned} |w_1 - w_2| &\leq |w_1 - w_0| + |w_0 - w_2| \\ &\leq q_4 \mu_{1,1}^2 |w_1 (\bar{\boldsymbol{\theta}}_{N_1} \|_2^2 / d - 1) w_0| + |w_1 w_2 \bar{\mathbf{A}}_{\cdot, N_1}^\top ((\bar{\mathbf{B}} - \xi \mathbf{I}_{P-1})^{-1} - (\tilde{\mathbf{B}} + \boldsymbol{\Delta} - \xi \mathbf{I}_{P-1})^{-1}) \bar{\mathbf{A}}_{\cdot, N_1}| \\ &\leq q_4 \mu_{1,1}^2 \|\bar{\boldsymbol{\theta}}_{N_1}\|_2^2 / d - 1 / \Im^2(\xi) + 2 \|\bar{\mathbf{A}}_{\cdot, N_1}\|_2^2 \|\mathbf{E}\|_{\text{op}} / \Im^4(\xi) \end{aligned}$$

By Lemma IV.4, we have  $\|\mathbf{E}_1\|_{\text{op}} = o_{\mathbb{P}}(1)$ ,  $\|\mathbf{E}_2\|_{\text{op}} = o_{\mathbb{P}}(1)$ . It is also easy to see that  $\|\bar{\mathbf{A}}_{\cdot, N_1}\|_2^2 = O_{\mathbb{P}}(1)$  and  $\|\bar{\boldsymbol{\theta}}_{N_1}\|_2^2 / d - 1 = o_{\mathbb{P}}(1)$ . Therefore we have

$$|w_1 - w_2| = o_{\mathbb{P}}(1).$$

Combining with the fact that  $|w_1 - w_2|$  is deterministically bounded by  $2/\Im(\xi)$ , by the dominated convergence theorem, we have

$$\mathbb{E}|w_1 - w_2| = o_d(1).$$

Therefore  $\bar{m}_{1,d} = \psi_1 \mathbb{E} w_2 + o_d(1)$ , and the derivation of the resolvent equations reduces to the calculation of  $\mathbb{E} w_2$ .

**Step 3.** We calculate  $\mathbb{E} w_2$  to get the resolvent equations. For simplicity, we give some notations which will be used later. Let

$$\mathbf{v} = \bar{\mathbf{A}}_{\cdot, N_1}, \quad \mathbf{v}_i = \bar{\mathbf{A}}_{i, N_1} = \begin{cases} \frac{q_4 \mu_{1,1}^2 \eta_i}{d} \|\bar{\boldsymbol{\theta}}_{N_1}\|_2, & \text{if } i \in [1, N_1 - 1], \\ \frac{q_4 \mu_{1,1} \mu_{2,1} \eta_{i+1}}{d} \|\bar{\boldsymbol{\theta}}_{N_1}\|_2, & \text{if } i \in [N_1, N - 1], \\ \frac{1}{\sqrt{d}} \Phi_1\left(\frac{1}{\sqrt{d}} u_{i-N+1} \|\bar{\boldsymbol{\theta}}_{N_1}\|_2\right), & \text{if } i \geq N, \end{cases}$$

and

$$\mathbf{U} = \frac{1}{\sqrt{d}} \begin{bmatrix} \boldsymbol{\eta}_1 & & \\ & \boldsymbol{\eta}_2 & \\ & & \mathbf{u} \end{bmatrix} \in \mathbb{R}^{(P-1) \times 3}, \quad \mathbf{M} = \begin{bmatrix} q_4 \mu_{1,1}^2 & q_4 \mu_{1,1} \mu_{2,1} & \phi_{1,1} \\ q_4 \mu_{1,1} \mu_{2,1} & q_4 \mu_{2,1}^2 & \phi_{2,1} \\ \phi_{1,1} & \phi_{2,1} & q_5 \end{bmatrix}.$$

By direct verification, we have

$$\boldsymbol{\Delta} = \mathbf{U} \mathbf{M} \mathbf{U}^\top.$$

We now decompose  $w_2$  into the terms related with  $\tilde{\mathbf{B}}$ ,  $\mathbf{v}$  and  $\mathbf{U}$ . By Schur complement formula, we have

$$\begin{aligned} (\tilde{\mathbf{B}} + \mathbf{U} \mathbf{M} \mathbf{U}^\top - \xi \mathbf{I}_{P-1})^{-1} &= (\tilde{\mathbf{B}} - \xi \mathbf{I}_{P-1})^{-1} \\ &\quad - (\tilde{\mathbf{B}} - \xi \mathbf{I}_{P-1})^{-1} \mathbf{U} [\mathbf{M}^{-1} + \mathbf{U}^\top (\tilde{\mathbf{B}} - \xi \mathbf{I}_{P-1})^{-1} \mathbf{U}]^{-1} \mathbf{U}^\top (\tilde{\mathbf{B}} - \xi \mathbf{I}_{P-1})^{-1}. \end{aligned} \quad (\text{IV.15})$$

Then  $w_2$  can be rewritten as

$$\begin{aligned} w_2 &= \left( -\xi + q_2 \mu_{1,*}^2 + q_4 \mu_{1,1}^2 - \mathbf{v}^\top (\tilde{\mathbf{B}} + \mathbf{U} \mathbf{M} \mathbf{U}^\top - \xi \mathbf{I}_{P-1})^{-1} \mathbf{v} \right)^{-1} \\ &= \left[ -\xi + q_2 \mu_{1,*}^2 + q_4 \mu_{1,1}^2 - \mathbf{v}^\top (\tilde{\mathbf{B}} - \xi \mathbf{I}_{P-1})^{-1} \mathbf{v} \right. \\ &\quad \left. + \mathbf{v}^\top (\tilde{\mathbf{B}} - \xi \mathbf{I}_{P-1})^{-1} \mathbf{U} (\mathbf{M}^{-1} + \mathbf{U}^\top (\tilde{\mathbf{B}} - \xi \mathbf{I}_{P-1})^{-1} \mathbf{U})^{-1} \mathbf{U}^\top (\tilde{\mathbf{B}} - \xi \mathbf{I}_{P-1})^{-1} \mathbf{v} \right]^{-1}, \end{aligned} \quad (\text{IV.16})$$

where the first equation is the definition of  $w_2$ , and the second equation follows by (IV.15). To continue the calculation, we study the terms  $\mathbf{v}^\top (\tilde{\mathbf{B}} - \xi \mathbf{I}_{P-1})^{-1} \mathbf{v}$ ,  $\mathbf{v}^\top (\tilde{\mathbf{B}} - \xi \mathbf{I}_{P-1})^{-1} \mathbf{U}$  and  $\mathbf{U}^\top (\tilde{\mathbf{B}} - \xi \mathbf{I}_{P-1})^{-1} \mathbf{U}$  in the denominator of (IV.16). To do so, we note that  $\tilde{\mathbf{B}}$  is independent of  $\mathbf{v}$  and  $\mathbf{U}$ . Moreover, by the leave-one-out argument, the Stieltjes transform corresponding to  $\tilde{\mathbf{B}}$  shares the same asymptotics as the Stieltjes transform corresponding to  $\bar{\mathbf{A}}$ . Notice that  $\eta_i$  is independent on  $\tilde{\mathbf{B}}$  conditioned on  $\bar{\boldsymbol{\theta}}_{N_1}$ , and  $\tilde{\mathbf{B}}$  is independent on  $\bar{\boldsymbol{\theta}}_{N_1}$ . We have

$$\begin{aligned} \mathbb{E} \mathbf{v}^\top (\tilde{\mathbf{B}} - \xi \mathbf{I}_{P-1})^{-1} \mathbf{v} &= \mathbb{E} \text{tr}(\tilde{\mathbf{B}} - \xi \mathbf{I}_{P-1})^{-1} \mathbf{v} \mathbf{v}^\top = \text{tr}(\mathbb{E}(\tilde{\mathbf{B}} - \xi \mathbf{I}_{P-1})^{-1} \mathbb{E} \mathbf{v} \mathbf{v}^\top) \\ &= \text{tr} \left( \begin{bmatrix} \frac{d\bar{m}_{1,d}}{N_1} \mathbf{I}_{N_1-1} & * & * \\ * & \frac{d\bar{m}_{2,d}}{N_2} \mathbf{I}_{N_2} & * \\ * & * & \frac{d\bar{m}_{3,d}}{n} \mathbf{I}_n \end{bmatrix} \right. \\ &\quad \cdot \frac{1}{d} \begin{bmatrix} (q_4^2 \mu_{1,1}^4 + o_d(1)) \mathbf{I}_{N_1-1} & & \\ & (q_4^2 \mu_{1,1}^2 \mu_{2,1}^2 + o_d(1)) \mathbf{I}_{N_2} & \\ & & (\phi_{1,1}^2 + \phi_{1,*}^2 + o_d(1)) \mathbf{I}_n \end{bmatrix} \left. \right) \\ &= q_4^2 \mu_{1,1}^2 (\mu_{1,1}^2 \bar{m}_{1,d} + \mu_{2,1}^2 \bar{m}_{2,d}) + (\phi_{1,1}^2 + \phi_{1,*}^2) \bar{m}_{3,d} + o_d(1), \end{aligned}$$

where the second equality follows from the fact that  $\mathbb{E} \Phi_1^2 \left( \frac{1}{\sqrt{d}} u_{i-N+1} \|\bar{\boldsymbol{\theta}}_{N_1}\|_2 \right) = \phi_{1,1}^2 + \phi_{1,*}^2 + o_d(1)$ , and we have denoted by “\*” the blocks that are irrelevant to the calculation. By a concentration measure argument (see in Tao (2012) Section 2.4.3), we have

$$\mathbf{v}^\top (\tilde{\mathbf{B}} - \xi \mathbf{I}_{P-1})^{-1} \mathbf{v} = q_4^2 \mu_{1,1}^2 (\mu_{1,1}^2 \bar{m}_{1,d} + \mu_{2,1}^2 \bar{m}_{2,d}) + (\phi_{1,1}^2 + \phi_{1,*}^2) \bar{m}_{3,d} + o_{\mathbb{P}}(1). \quad (\text{IV.17})$$

After direct calculation with the same argument, we obtain that

$$\mathbf{v}^\top (\tilde{\mathbf{B}} - \xi \mathbf{I}_{P-1})^{-1} \mathbf{U} = \begin{bmatrix} q_4 \mu_{1,1}^2 \bar{m}_{1,d} \\ q_4 \mu_{1,1} \mu_{2,1} \bar{m}_{2,d} \\ \phi_{1,1} \bar{m}_{3,d} \end{bmatrix}^\top + o_{\mathbb{P}}(1), \quad (\text{IV.18})$$

$$\mathbf{U}^\top (\tilde{\mathbf{B}} - \xi \mathbf{I}_{P-1})^{-1} \mathbf{U} = \begin{bmatrix} \bar{m}_{1,d} & & \\ & \bar{m}_{2,d} & \\ & & \bar{m}_{3,d} \end{bmatrix} + o_{\mathbb{P}}(1). \quad (\text{IV.19})$$

Now since  $|w_2| \leq 1/\xi_0$  is deterministically bounded, by dominated convergence theorem, we have the  $L_1$  convergence of  $w_2$  by plugging the main terms of (IV.17), (IV.18) and (IV.19) into (IV.16). Further note that equation (IV.16) has a part in the form of  $(\mathbf{A}^{-1} + \mathbf{M}^{-1})^{-1}$ , where

$$\mathbf{A} = \begin{bmatrix} 1/\bar{m}_{1,d} & & \\ & 1/\bar{m}_{2,d} & \\ & & 1/\bar{m}_{3,d} \end{bmatrix}, \quad \mathbf{M} = \begin{bmatrix} q_4 \mu_{1,1}^2 & q_4 \mu_{1,1} \mu_{2,1} & \phi_{1,1} \\ q_4 \mu_{1,1} \mu_{2,1} & q_4 \mu_{2,1}^2 & \phi_{2,1} \\ \phi_{1,1} & \phi_{2,1} & q_5 \end{bmatrix}.$$

By the formula  $(\mathbf{A}^{-1} + \mathbf{M}^{-1})^{-1} = \mathbf{A} - \mathbf{A}(\mathbf{A} + \mathbf{M})^{-1} \mathbf{A}$ , we have

$$(\mathbf{M}^{-1} + \mathbf{A}^{-1})^{-1} = \begin{bmatrix} 1/\bar{m}_{1,d} & & \\ & 1/\bar{m}_{2,d} & \\ & & 1/\bar{m}_{3,d} \end{bmatrix} - \mathbf{A} \begin{bmatrix} q_4 \mu_{1,1}^2 + 1/\bar{m}_{1,d} & q_4 \mu_{1,1} \mu_{2,1} & \phi_{1,1} \\ q_4 \mu_{1,1} \mu_{2,1} & q_4 \mu_{2,1}^2 + 1/\bar{m}_{2,d} & \phi_{2,1} \\ \phi_{1,1} & \phi_{2,1} & q_5 + 1/\bar{m}_{3,d} \end{bmatrix}^{-1} \mathbf{A}.$$

Denote  $\mathbf{l} = [q_4\mu_{1,1}^2\bar{m}_{1,d} \quad q_4\mu_{1,1}\mu_{2,1}\bar{m}_{2,d} \quad \phi_{1,1}\bar{m}_{3,d}]^\top$ . Then by plugging the equation above into (IV.16), and combining it with (IV.17), (IV.18) and (IV.19), we finally get

$$\begin{aligned} \bar{m}_{1,d} &= \psi_1 \mathbb{E} w_2 \\ &= \psi_1 \left\{ -\xi + q_2\mu_{1,*}^2 + q_4\mu_{1,1}^2 - q_4^2\mu_{1,1}^2(\mu_{1,1}^2\bar{m}_{1,d} + \mu_{2,1}^2\bar{m}_{2,d}) \right. \\ &\quad \left. - (\phi_{1,1}^2 + \phi_{1,*}^2)\bar{m}_{3,d} + \mathbf{l}^\top \mathbf{A} \mathbf{l} - \mathbf{l}^\top \mathbf{A}(\mathbf{A} + \mathbf{M})^{-1} \mathbf{A} \mathbf{l} \right\}^{-1} + o_d(1) \\ &= \psi_1 \left\{ -\xi + q_2\mu_{1,*}^2 + q_4\mu_{1,1}^2 - \phi_{1,*}^2\bar{m}_{3,d} - \begin{bmatrix} q_4\mu_{1,1}^2 \\ q_4\mu_{1,1}\mu_{2,1} \\ \phi_{1,1} \end{bmatrix}^\top (\mathbf{A} + \mathbf{M})^{-1} \begin{bmatrix} q_4\mu_{1,1}^2 \\ q_4\mu_{1,1}\mu_{2,1} \\ \phi_{1,1} \end{bmatrix} \right\}^{-1} \\ &\quad + o_d(1). \end{aligned}$$

Now note that  $\phi_{j,1}^2 = \mu_{j,1}^2(1 + q_1)^2$ ,  $\phi_{j,*}^2 = \mu_{j,*}^2$ ,  $j = 1, 2$ . Therefore with direct calculation, we have

$$\bar{m}_{1,d} = \psi_1 \left\{ -\xi + q_2\mu_{1,*}^2 - \mu_{1,*}^2\bar{m}_{3,d} + \frac{H_{1,d}}{H_{D,d}} \right\}^{-1} + o_d(1), \quad (\text{IV.20})$$

where

$$\begin{aligned} H_{1,d} &= \mu_{1,1}^2 q_4(1 + \bar{m}_{3,d} q_5) - \mu_{1,1}^2(1 + q_1)^2 \bar{m}_{3,d}, \\ H_{D,d} &= (1 + \mu_{1,1}^2 \bar{m}_{1,d} q_4 + \mu_{2,1}^2 \bar{m}_{2,d} q_4)(1 + \bar{m}_{3,d} q_5) - \mu_{2,1}^2(1 + q_1)^2 \bar{m}_{2,d} \bar{m}_{3,d} \\ &\quad - \mu_{1,1}^2(1 + q_1)^2 \bar{m}_{1,d} \bar{m}_{3,d}. \end{aligned}$$

The equation above shows that the magnitude of the first element of  $\bar{\mathbf{m}}_d(\xi) - \mathbf{F}(\bar{\mathbf{m}}_d(\xi))$  is  $o_d(1)$ . With exactly the same proof, we also have

$$\begin{aligned} \bar{m}_{2,d} &= \psi_2 \left\{ -\xi + q_2\mu_{2,*}^2 - \mu_{2,*}^2\bar{m}_{3,d} + \frac{H_{2,d}}{H_{D,d}} \right\}^{-1} + o_d(1), \\ \bar{m}_{3,d} &= \psi_3 \left\{ -\xi + q_3 - \mu_{1,*}^2\bar{m}_{1,d} - \mu_{2,*}^2\bar{m}_{2,d} + \frac{H_{3,d}}{H_{D,d}} \right\}^{-1} + o_d(1), \end{aligned} \quad (\text{IV.21})$$

where

$$\begin{aligned} H_{2,d} &= \mu_{2,1}^2 q_4(1 + \bar{m}_{3,d} q_5) - \mu_{2,1}^2(1 + q_1)^2 \bar{m}_{3,d}, \\ H_{3,d} &= q_5(1 + \mu_{1,1}^2 \bar{m}_{1,d} q_4 + \mu_{2,1}^2 \bar{m}_{2,d} q_4) - \mu_{2,1}^2(1 + q_1)^2 \bar{m}_{2,d} - \mu_{1,1}^2(1 + q_1)^2 \bar{m}_{1,d}. \end{aligned}$$

This completes the proof of Lemma IV.5. ■

### IV.3 Proof for conclusions 1 and 2 in Proposition A.6

We first introduce an important lemma about the property of Stieltjes transforms, which is given in [Hastie et al. \(2022\)](#).

**Lemma IV.6** (Lemma 7 in [Hastie et al. \(2022\)](#)). *The functions  $\xi \rightarrow \bar{m}_{i,d}(\xi)$ ,  $i = 1, 2, 3$ , have the following properties:*

1.  $\bar{m}_{i,d}$ ,  $i = 1, 2, 3$  are analytical on  $\mathbb{C}_+$ , and map  $\mathbb{C}_+$  into  $\mathbb{C}_+$ .

2. Let  $\Omega \subset \mathbb{C}_+$  be a set with an accumulation point. If  $\bar{m}_{i,d} \rightarrow m_i(\xi)$  for all  $\xi \in \Omega$ , then  $m_i(\xi)$  has a unique analytic continuation to  $\mathbb{C}_+$  and  $\bar{m}_{i,d} \rightarrow m_i(\xi)$  for all  $\xi \in \mathbb{C}_+$ . Moreover, the convergence is uniform over compact sets  $\Omega \subset \mathbb{C}_+$ .

We now give the proof of the conclusions in Proposition A.6 that  $\mathbf{m}(\xi; \mathbf{q}, \boldsymbol{\mu})$  is analytic on  $\{\xi : \Im(\xi) > \xi_0\}$  for some sufficiently large  $\xi_0$ , has unique analytic continuation to  $\mathbb{C}_+$  and maps  $\mathbb{C}_+$  to  $\mathbb{C}_+^3$ . Denote  $\bar{\mathbf{m}}_d = \bar{\mathbf{m}}_d(\xi) = [\bar{m}_{1,d}(\xi), \bar{m}_{2,d}(\xi), \bar{m}_{3,d}(\xi)]^\top$ . Then for any fixed  $\xi \in \mathbb{C}_+$ , Lemma IV.5 gives

$$\|\bar{\mathbf{m}}_d - \mathbf{F}(\bar{\mathbf{m}}_d)\|_2 = o_d(1). \quad (\text{IV.22})$$

By Lemma III.1, there exists a  $\xi_0 > 0$ , such that for all  $\xi$  with  $\Im(\xi) \geq \xi_0$ ,  $\mathbf{F}(\cdot)$  is 1/2-Lipschitz with respect to  $\ell_2$  norm. Moreover, for all  $\xi$  with  $\Im(\xi) \geq \xi_0$  we have

$$\begin{aligned} \|\bar{\mathbf{m}}_d - \mathbf{m}\|_2 &= \|\bar{\mathbf{m}}_d - \mathbf{F}(\mathbf{m})\|_2 \\ &\leq \|\bar{\mathbf{m}}_d - \mathbf{F}(\bar{\mathbf{m}}_d)\|_2 + \|\mathbf{F}(\bar{\mathbf{m}}_d) - \mathbf{F}(\mathbf{m})\|_2 \\ &\leq o_d(1) + \frac{1}{2} \cdot \|\bar{\mathbf{m}}_d - \mathbf{m}\|_2, \end{aligned}$$

where the equality is by the definition of  $\mathbf{m}$  as the unique fixed point of  $\mathbf{F}(\cdot)$ , the first inequality is by triangle inequality, the second inequality is by (IV.22) and the fact that  $\mathbf{F}(\cdot)$  is 1/2-Lipschitz with respect to  $\ell_2$  norm. Therefore we have  $\|\bar{\mathbf{m}}_d(\xi) - \mathbf{m}(\xi)\|_2 = o_d(1)$  for all  $\xi$  with  $\Im(\xi) \geq \xi_0$ . The properties of Stieljes transforms (see Lemma IV.6) then imply that  $\mathbf{m}(\xi; \mathbf{q}, \boldsymbol{\mu})$  is analytic in  $\{\xi : \Im(\xi) > \xi_0\}$ , and has a unique analytic continuation to  $\mathbb{C}_+$ . Moreover, the extended  $\mathbf{m}(\xi; \mathbf{q}, \boldsymbol{\mu})$  satisfies

$$\|\bar{\mathbf{m}}_d(\xi) - \mathbf{m}(\xi)\|_2 = o_d(1) \quad (\text{IV.23})$$

for any fixed  $\xi \in \mathbb{C}_+$ . This implies that  $\mathbf{m}(\xi; \mathbf{q}, \boldsymbol{\mu})$  is  $\mathbb{C}_+ \rightarrow \mathbb{C}_+^3$  by the definition of  $\bar{\mathbf{m}}_d$ . The proof of Conclusion 1 is complete.

To prove Conclusion 2, we first prove that  $\mathbf{m}(\xi)$  is a continuity point of  $\mathbf{F}(\cdot)$  for any fixed  $\xi \in \mathbb{C}_+$ . For any fixed  $\xi \in \mathbb{C}_+$ , assume that  $\mathbf{m}(\xi)$  is not a continuity point of  $\mathbf{F}(\cdot)$ , by the definition of  $\mathbf{F}(\cdot)$  we have  $\|\mathbf{F}(\mathbf{m}(\xi))\|_2 = +\infty$ . Therefore, for any  $M > 0$ , there exists  $\delta(\xi, M) > 0$  ( $\xi \in \mathbb{C}_+$  is fixed here), as long as  $\|\bar{\mathbf{m}}_d(\xi) - \mathbf{m}(\xi)\| < \delta(\xi, M)$ , the inequality  $\mathbf{F}(\bar{\mathbf{m}}_d(\xi)) > M$  holds. Moreover, for the  $\delta(\xi, M)$ , there always exists  $d_0$  such that  $\|\bar{\mathbf{m}}_d(\xi) - \mathbf{m}(\xi)\| < \delta(\xi, M)$  for all  $d > d_0$ . That is: for any fixed  $\xi \in \mathbb{C}_+$ , and any large constant  $M > 0$ , there always exists  $d_0$  such that  $\mathbf{F}(\bar{\mathbf{m}}_d(\xi)) > M$  for  $d > d_0$ . Combined with (IV.22), there exists  $d_1 > 0$ , such that  $\|\bar{\mathbf{m}}_d(\xi) - \mathbf{F}(\bar{\mathbf{m}}_d(\xi))\|_2 < 1$  for all  $d > d_1$ . Then for  $d > \max\{d_0, d_1\}$ , we have  $\mathbf{F}(\bar{\mathbf{m}}_d(\xi)) > M$  and  $\|\bar{\mathbf{m}}_d(\xi) - \mathbf{F}(\bar{\mathbf{m}}_d(\xi))\|_2 < 1$ . We have  $\|\bar{\mathbf{m}}_d(\xi)\|_2 > M - 1$  for  $d > \max\{d_0, d_1\}$ . On the other hand,  $\|\bar{\mathbf{m}}_d(\xi)\|_2 \leq 2(\psi_1 + \psi_2 + \psi_3)/\Im(\xi)$  from the definition of  $\bar{\mathbf{m}}_d(\xi)$ . Note that  $\xi \in \mathbb{C}_+$  is fixed here. Enlarging  $M$  leads to a contradiction. Therefore,  $\mathbf{m}(\xi)$  is the continuity point of  $\mathbf{F}(\cdot)$  for any fixed  $\xi \in \mathbb{C}_+$ .

For any fixed  $\xi \in \mathbb{C}_+$ , note that  $\mathbf{m}(\xi)$  is the continuity point of  $\mathbf{F}(\cdot)$ . Let  $d \rightarrow +\infty$ , (IV.22) and (IV.23) give us

$$\|\mathbf{m} - \mathbf{F}(\mathbf{m})\|_2 = 0.$$

This means that  $\mathbf{F}(\mathbf{m}(\xi; \mathbf{q}, \boldsymbol{\mu})) = \mathbf{m}(\xi; \mathbf{q}, \boldsymbol{\mu})$  for any fixed  $\xi \in \mathbb{C}_+$ . The proof of Conclusion 2 is complete.

#### IV.4 Point convergence for polynomial activation functions

We now give the proof of point convergence under the additional assumption that the activation functions are polynomials. We remind readers the “Gaussian version” of the problem defined in Sections IV.1 and IV.2, where the data inputs  $\bar{\mathbf{x}}_i$ ,  $i \in [n]$  and  $\bar{\boldsymbol{\theta}}_a$ ,  $a \in [N]$  are defined in Definition IV.1 and the activation functions  $\phi_1(x), \phi_2(x)$  are given in (IV.2). We also remind readers that the “Gaussian version” and “spherical version” Stieltjes transforms of the empirical eigenvalue distributions of linear pencil matrices are denoted as  $\bar{M}_d(\xi)$  and  $M_d(\xi)$ , respectively. Importantly, the expectation of  $\bar{M}_d(\xi)$  is denoted as  $\bar{m}_d$ , while  $m(\xi; \mathbf{q}, \boldsymbol{\mu}) = \sum_{i=1}^3 m_i(\xi)$ , where  $\mathbf{m} = \mathbf{m}(\xi) = (m_1(\xi), m_2(\xi), m_3(\xi))^\top$  is defined as the solution of (A.2) on  $\{\xi : \Im(\xi) \geq \xi_0\}$  and then extended to  $\mathbb{C}_+$  by analytic continuation.

By (IV.7), for all fixed  $\xi \in \mathbb{C}_+$ , we have

$$\mathbb{E} \left| \bar{M}_d(\xi) - \sum_{i=1}^3 \bar{m}_{i,d}(\xi) \right| = o_d(1). \quad (\text{IV.24})$$

In addition, by Lemma IV.3, when the activation functions are polynomials, we have

$$\mathbb{E} |M_d(\xi) - \bar{M}_d(\xi)| = o_d(1). \quad (\text{IV.25})$$

Combining (IV.23) (IV.24) and (IV.25) gives

$$\mathbb{E} |M_d(\xi) - m(\xi)| = o_d(1)$$

for any fixed  $\xi \in \mathbb{C}_+$ , which completes the proof of the point convergence for polynomial activation functions.

#### IV.5 Point convergence for general activation functions satisfying Assumption 3.2

We now extend the result for polynomial activation functions to general activation functions satisfying Assumption 3.2. Let  $\tau_d$  be the marginal distribution of  $\langle \mathbf{x}, \boldsymbol{\theta} \rangle / \sqrt{d}$  for  $\mathbf{x}, \boldsymbol{\theta} \sim_{\text{iid}} \text{Unif}(\sqrt{d} \cdot \mathbb{S}^{d-1})$ , and  $\bar{\tau}_d$  the marginal distribution of  $\langle \bar{\mathbf{x}}, \bar{\boldsymbol{\theta}} \rangle / \sqrt{d}$  for  $\bar{\mathbf{x}}, \bar{\boldsymbol{\theta}} \sim_{\text{iid}} \text{N}(0, \mathbf{I}_d)$ . For  $j = 1, 2$ , suppose that  $\sigma_j$  are activation functions satisfying Assumption 3.2. The idea here is to construct polynomial activation functions  $\tilde{\sigma}_j$  to approximate  $\sigma_j$ . To do so, we recall that  $\mathbf{m} = \mathbf{m}(\xi; \mathbf{q}, \boldsymbol{\mu})$  solves the implicit equations

$$\mathbf{m} = \mathbf{F}(\mathbf{m}; \xi, \mathbf{q}, \boldsymbol{\mu}),$$

where  $\mathbf{F}(\cdot; \xi, \mathbf{q}, \boldsymbol{\mu})$  is defined in Definition A.5. When  $\Im(\xi) > \xi_0$  for some large enough  $\xi_0$ , by the continuity of the solution of the fixed point equation with respect to  $\boldsymbol{\mu}$ , we have

$$\lim_{\tilde{\boldsymbol{\mu}} \rightarrow \boldsymbol{\mu}} \mathbf{m}(\xi; \mathbf{q}, \tilde{\boldsymbol{\mu}}) = \mathbf{m}(\xi; \mathbf{q}, \boldsymbol{\mu}).$$

According to our proof in Section IV.3, we can extend the definition of  $\mathbf{m}$  to  $\xi \in \mathbb{C}_+$  with analytic continuation. Then with the same proof as in Mei and Montanari (2022) (see equation (10.56) in Mei and Montanari (2022)), for any fixed  $\xi \in \mathbb{C}_+$  and any  $\varepsilon > 0$ , there exists  $\delta = \delta(\varepsilon, \xi, \mathbf{q}, \boldsymbol{\mu}) > 0$



such that

$$\|\mathbf{m}(\xi; \mathbf{q}, \tilde{\boldsymbol{\mu}}) - \mathbf{m}(\xi; \mathbf{q}, \boldsymbol{\mu})\|_2 \leq \varepsilon \quad (\text{IV.26})$$

for all  $\tilde{\boldsymbol{\mu}}$  with  $\|\boldsymbol{\mu} - \tilde{\boldsymbol{\mu}}\|_2 \leq \delta$ . Now by Assumption 3.2, for any fixed  $\varepsilon > 0$ , we can choose a sufficiently large integer  $\bar{k}$  and construct

$$\tilde{\sigma}_j(x) = \sum_{k=0}^{\bar{k}} \frac{\mu_{j,k}}{k!} H_k(x),$$

such that for  $G \sim N(0, 1)$ ,

$$\mathbb{E}[\sigma_j(G) - \tilde{\sigma}_j(G)]^2 \leq \varepsilon^2, \quad (\text{IV.27})$$

$$|\mathbb{E}\{\tilde{\sigma}_j(G)^2\} - \mathbb{E}\{\sigma_j(G)^2\}| \leq \delta^2/2. \quad (\text{IV.28})$$

Here,  $\{H_k(x)\}$  are the family of Hermite polynomials. Then by (IV.27) and Lemma 5 in Ghorbani et al. (2021), we have

$$\|\sigma_j - \tilde{\sigma}_j\|_{L^2(\tau_d)} \leq \varepsilon \quad (\text{IV.29})$$

for  $j = 1, 2$  and sufficiently large  $d$ , where we denote  $\|\sigma_j - \tilde{\sigma}_j\|_{L^2(\nu)} = \int (\sigma_j(x) - \tilde{\sigma}_j(x))^2 \nu(dx)$ .

Following Definition 3.1, we can also define the parameters  $\tilde{\mu}_{j,0}, \tilde{\mu}_{j,1}, \tilde{\mu}_{j,*}$  corresponding to the polynomial activation functions  $\tilde{\sigma}_j$  by

$$\tilde{\mu}_{j,0} \triangleq \mathbb{E}\{\tilde{\sigma}_j(G)\}, \quad \tilde{\mu}_{j,1} \triangleq \mathbb{E}\{G\tilde{\sigma}_j(G)\}, \quad \tilde{\mu}_{j,*}^2 \triangleq \mathbb{E}\{\tilde{\sigma}_j(G)^2\} - \tilde{\mu}_{j,0}^2 - \tilde{\mu}_{j,1}^2.$$

Then by the definition of  $\tilde{\sigma}_j$ , we have  $\mu_{j,0} = \tilde{\mu}_{j,0}$ ,  $\mu_{j,1} = \tilde{\mu}_{j,1}$  for  $j = 1, 2$ . Moreover, we also have

$$|\mu_{j,*} - \tilde{\mu}_{j,*}| \leq \sqrt{|\mu_{j,*}^2 - \tilde{\mu}_{j,*}^2|} = \sqrt{|\mathbb{E}\{\sigma_j(G)^2\} - \mathbb{E}\{\tilde{\sigma}_j(G)^2\}|} \leq \delta/\sqrt{2}$$

for  $j = 1, 2$ , where the first inequality follows from  $|a - b| \leq \sqrt{|a^2 - b^2|}$  for all  $a, b > 0$ , the equality follows by  $\mu_{j,0} = \tilde{\mu}_{j,0}$ ,  $\mu_{j,1} = \tilde{\mu}_{j,1}$  for  $j = 1, 2$ , and the last inequality follows by (IV.28). Therefore we have  $\|\tilde{\boldsymbol{\mu}} - \boldsymbol{\mu}\|_2 \leq \delta$ .

Let  $\tilde{\mathbf{m}}(\xi) = [\tilde{m}_1(\xi), \tilde{m}_2(\xi), \tilde{m}_3(\xi)]^\top$  be the solution of the implicit equations

$$\tilde{\mathbf{m}} = \mathbf{F}(\tilde{\mathbf{m}}; \xi, \mathbf{q}, \tilde{\boldsymbol{\mu}}),$$

and let  $\tilde{m}(\xi) = \tilde{m}_1(\xi) + \tilde{m}_2(\xi) + \tilde{m}_3(\xi)$ , where we drop the arguments  $\mathbf{q}, \tilde{\boldsymbol{\mu}}$  in  $\tilde{\mathbf{m}}(\xi; \mathbf{q}, \tilde{\boldsymbol{\mu}})$  for notation simplification. Then by (IV.26), we have

$$|\tilde{m}(\xi) - m(\xi)| \leq 3\varepsilon. \quad (\text{IV.30})$$

Let  $\tilde{\mathbf{A}}$  be the linear pencil matrix corresponding to  $\tilde{\sigma}$  in Definition A.3, and define  $\widetilde{M}_d(\xi) = (1/d) \cdot \text{tr}[(\tilde{\mathbf{A}} - \xi \mathbf{I})^{-1}]$ . Then we have

$$\mathbb{E}[|M_d(\xi) - \widetilde{M}_d(\xi)|] = \frac{1}{d} \mathbb{E}[|\text{tr}[(\mathbf{A} - \xi \mathbf{I})^{-1}(\tilde{\mathbf{A}} - \mathbf{A})(\tilde{\mathbf{A}} - \xi \mathbf{I})^{-1}]|]$$

$$\begin{aligned}
&\leq \frac{1}{d} \mathbb{E} [\|(\mathbf{A} - \xi \mathbf{I})^{-1}(\tilde{\mathbf{A}} - \xi \mathbf{I})^{-1}\|_{\text{op}} \|\tilde{\mathbf{A}} - \mathbf{A}\|_*] \\
&\leq [1/(\Im(\xi)^2) \cdot P^{-1/2} \cdot \mathbb{E}\{\|\tilde{\mathbf{A}} - \mathbf{A}\|_F^2\}]^{1/2} \\
&\leq C'(\xi, \boldsymbol{\psi}) \cdot [1/(\Im(\xi)^2) \cdot d^{-1/2} \cdot \mathbb{E}\{\|\tilde{\mathbf{A}} - \mathbf{A}\|_F^2\}]^{1/2} \\
&\leq C''(\xi, \mathbf{q}) \cdot (\|\sigma_1 - \tilde{\sigma}_1\|_{L^2(\tau_d)} + \|\sigma_2 - \tilde{\sigma}_2\|_{L^2(\tau_d)}), \tag{IV.31}
\end{aligned}$$

where  $C''(\xi, \boldsymbol{\psi}) > 0$  is a constant only depending on  $\xi$  and  $\boldsymbol{\psi}$ , and  $C''(\xi, \mathbf{q}) > 0$  only depends on  $\xi$ ,  $\mathbf{q}$  and  $\boldsymbol{\psi}$ . Here the second inequality above follows by Cauchy-Schwarz inequality, the third inequality follows by  $P = N_1 + N_2 + n$  and the assumption that  $N_1, N_2, n, d$  goes to infinity proportionally, and the last inequality follows by the definitions of  $\tilde{\mathbf{A}}$  and  $\mathbf{A}$ . Therefore, by (IV.29) and (IV.31), we have

$$\mathbb{E}|M_d(\xi) - \tilde{M}_d(\xi)| \leq 2C''(\xi, \mathbf{q}) \cdot \varepsilon \tag{IV.32}$$

for sufficiently large  $d$ . Moreover, since  $\tilde{\sigma}_j$ ,  $j = 1, 2$  are polynomial activation functions, by the results in Appendix IV.4, we have

$$\mathbb{E}|\tilde{M}_d(\xi) - \tilde{m}(\xi)| = o_d(1). \tag{IV.33}$$

Combining (IV.30), (IV.32) and (IV.33) and taking  $d \rightarrow \infty$ , we have

$$\limsup_{d \rightarrow +\infty} \mathbb{E}|M_d(\xi) - m(\xi)| \leq (2C''(\xi, \mathbf{q}) + 3) \cdot \varepsilon$$

for all fixed  $\xi \in \mathbb{C}_+$ . Taking  $\varepsilon \rightarrow 0^+$ , we conclude that  $\lim_{d \rightarrow \infty} \mathbb{E}|\tilde{M}_d(\xi) - \tilde{m}(\xi)| = 0$ , which proves the point convergence for general activation functions.

## IV.6 Uniform convergence on compact sets

In this section, we aim to prove that on compact sets the point convergence established above could be extended to uniform convergence. Consider a compact set  $\Omega \subset \mathbb{C}_+$ . From the proof above we have

$$\lim_{d \rightarrow +\infty} |\mathbb{E}M_d(\xi; \mathbf{q}, \boldsymbol{\mu}) - m(\xi; \mathbf{q}, \boldsymbol{\mu})| = 0.$$

Then from Lemma IV.6, we have

$$\lim_{d \rightarrow +\infty} \sup_{\xi \in \Omega} |\mathbb{E}M_d(\xi; \mathbf{q}, \boldsymbol{\mu}) - m(\xi; \mathbf{q}, \boldsymbol{\mu})| = 0. \tag{IV.34}$$

Moreover, by the definition of  $M_d(\xi; \mathbf{q}, \boldsymbol{\mu})$ , we have

$$\begin{aligned}
|M_d(\xi_1; \mathbf{q}, \boldsymbol{\mu}) - M_d(\xi_2; \mathbf{q}, \boldsymbol{\mu})| &= \frac{1}{d} |\text{tr}((\mathbf{A} - \xi_1 \mathbf{I})^{-1}(\xi_1 - \xi_2)(\mathbf{A} - \xi_2 \mathbf{I})^{-1})| \\
&\leq \frac{P}{d \cdot \Im(\xi_1) \Im(\xi_2)} \cdot |\xi_1 - \xi_2|.
\end{aligned}$$

Since  $P$  is proportional to  $d$ , there exists a constant  $L_0$  that only depends on  $\psi_1, \psi_2, \psi_3$  and  $\Omega$ , such that  $M_d(\xi; \mathbf{q}, \boldsymbol{\mu})$  is  $L_0$ -Lipschitz for all  $d \in \mathbb{N}$ . Then by the compactness of  $\Omega$ , for any  $\varepsilon > 0$ , there

exists a finite set  $\mathcal{N}_\varepsilon(\Omega) \subset \mathbb{C}_+$ , that is an  $\varepsilon/L_0$  covering of the compact set  $\Omega$ . Specifically, for any  $\xi \in \Omega$ , there exists a  $\xi_* \in \mathcal{N}_\varepsilon(\Omega)$  such that  $|\xi - \xi_*| < \varepsilon/L_0$ . Therefore

$$\begin{aligned} \sup_{\xi \in \Omega} \inf_{\xi_* \in \mathcal{N}_\varepsilon(\Omega)} |M_d(\xi; \mathbf{q}, \boldsymbol{\mu}) - M_d(\xi_*; \mathbf{q}, \boldsymbol{\mu})| &\leq \varepsilon, \\ \sup_{\xi \in \Omega} \inf_{\xi_* \in \mathcal{N}_\varepsilon(\Omega)} |\mathbb{E}M_d(\xi; \mathbf{q}, \boldsymbol{\mu}) - \mathbb{E}M_d(\xi_*; \mathbf{q}, \boldsymbol{\mu})| &\leq \varepsilon \end{aligned} \quad (\text{IV.35})$$

for all  $d \in \mathbb{N}$ . Moreover, since  $\mathcal{N}_\varepsilon(\Omega)$  is finite, the number of  $\xi_*$  is finite. Similar to the proof of (IV.7), we have

$$\sup_{\xi_* \in \mathcal{N}_\varepsilon(\Omega)} |M_d(\xi_*; \mathbf{q}, \boldsymbol{\mu}) - \mathbb{E}M_d(\xi_*; \mathbf{q}, \boldsymbol{\mu})| = o_{\mathbb{P}}(1).$$

Now since  $|M_d(\xi_*; \mathbf{q}, \boldsymbol{\mu})| \leq P/(d \cdot \Im(\xi_*)) \leq P/(d \cdot \inf_{\xi \in \Omega} \Im(\xi))$ ,  $|M_d(\xi_*; \mathbf{q}, \boldsymbol{\mu})|$  is bounded by some constant. By the dominated convergence theorem, we have

$$\mathbb{E} \sup_{\xi_* \in \mathcal{N}_\varepsilon(\Omega)} |M_d(\xi_*; \mathbf{q}, \boldsymbol{\mu}) - \mathbb{E}M_d(\xi_*; \mathbf{q}, \boldsymbol{\mu})| = o_d(1). \quad (\text{IV.36})$$

Combining (IV.34), (IV.35) and (IV.36), we obtain

$$\begin{aligned} &\mathbb{E} \left[ \sup_{\xi \in \Omega} |M_d(\xi; \mathbf{q}, \boldsymbol{\mu}) - m(\xi; \mathbf{q}, \boldsymbol{\mu})| \right] \\ &= \mathbb{E} \left\{ \sup_{\xi \in \Omega} \inf_{\xi_* \in \mathcal{N}_\varepsilon(\Omega)} \left| M_d(\xi; \mathbf{q}, \boldsymbol{\mu}) + M_d(\xi_*; \mathbf{q}, \boldsymbol{\mu}) - M_d(\xi_*; \mathbf{q}, \boldsymbol{\mu}) + \mathbb{E}M_d(\xi_*; \mathbf{q}, \boldsymbol{\mu}) \right. \right. \\ &\quad \left. \left. - \mathbb{E}M_d(\xi_*; \mathbf{q}, \boldsymbol{\mu}) + \mathbb{E}M_d(\xi; \mathbf{q}, \boldsymbol{\mu}) - \mathbb{E}M_d(\xi; \mathbf{q}, \boldsymbol{\mu}) - m(\xi; \mathbf{q}, \boldsymbol{\mu}) \right| \right\} \\ &\leq \mathbb{E} \left\{ \sup_{\xi \in \Omega} \inf_{\xi_* \in \mathcal{N}_\varepsilon(\Omega)} \left[ |M_d(\xi; \mathbf{q}, \boldsymbol{\mu}) - M_d(\xi_*; \mathbf{q}, \boldsymbol{\mu})| + |\mathbb{E}M_d(\xi; \mathbf{q}, \boldsymbol{\mu}) - \mathbb{E}M_d(\xi_*; \mathbf{q}, \boldsymbol{\mu})| \right. \right. \\ &\quad \left. \left. + |M_d(\xi_*; \mathbf{q}, \boldsymbol{\mu}) - \mathbb{E}M_d(\xi_*; \mathbf{q}, \boldsymbol{\mu})| + |\mathbb{E}M_d(\xi; \mathbf{q}, \boldsymbol{\mu}) - m(\xi; \mathbf{q}, \boldsymbol{\mu})| \right] \right\} \\ &\leq 2\varepsilon + o_d(1). \end{aligned}$$

Taking  $d \rightarrow +\infty$ , we have

$$\lim_{d \rightarrow +\infty} \mathbb{E} \left[ \sup_{\xi \in \Omega} |M_d(\xi; \mathbf{q}, \boldsymbol{\mu}) - m(\xi; \mathbf{q}, \boldsymbol{\mu})| \right] \leq 2\varepsilon.$$

Therefore, taking  $\varepsilon \rightarrow 0^+$  proves Conclusion 3 in Proposition A.6. The proof of Proposition A.6 is complete.

## V. Proof of Proposition A.7

We first present some lemmas in Section V.1, and then complete the proof in Section V.2. Recall that we assume  $\mathbf{q} \in \mathcal{Q}$  (see Definition A.3).

### V.1 Preliminary lemmas

The lemma below presents some additional properties of the function  $\mathbf{m}(\xi) = [m_1(\xi), m_2(\xi), m_3(\xi)]^\top$  defined in Proposition A.6.

**Lemma V.1.** *Let  $\mathbf{m}(\xi) = [m_1(\xi), m_2(\xi), m_3(\xi)]^\top$  defined on  $\xi \in \mathbb{C}_+$  be the analytic continuation of the solution of the implicit equations  $\mathbf{m} = \mathbf{F}(\mathbf{m}; \xi, \mathbf{q}, \boldsymbol{\mu})$  defined in Proposition A.6. Then for any fixed  $\xi_r \in \mathbb{R}$  and  $j = 1, 2, 3$ , we have*

$$\lim_{u \rightarrow +\infty} |m_j(\xi_r + iu) \cdot (\xi_r + iu) + \psi_j| = 0,$$

**Proof** [Proof of Lemma V.1] We denote  $\xi_u = \xi_r + iu$  for  $u > 0$ , and use the same definition of  $H_1(\mathbf{m}; \mathbf{q}, \boldsymbol{\mu})$  as in (III.1) that

$$H_1(\mathbf{m}; \mathbf{q}, \boldsymbol{\mu}) = -\mu_{1,*}^2 m_3 + \frac{1}{m_1 + \frac{-\mu_{2,1}^2(1+q_1)^2 m_2 m_3 + (1+\mu_{2,1}^2 m_2 q_4)(1+m_3 q_5)}{\mu_{1,1}^2 q_4(1+m_3 q_5) - \mu_{1,1}^2(1+q_1)^2 m_3}}.$$

Then we have

$$F_1(\mathbf{m}; \xi, \mathbf{q}, \boldsymbol{\mu}) = \frac{\psi_1}{-\xi + q_2 \mu_{1,*}^2 + H_1(\mathbf{m}; \mathbf{q}, \boldsymbol{\mu})}. \quad (\text{V.1})$$

Moreover, define  $\tilde{m}_1(\xi) = -\psi_1/\xi$ ,  $\tilde{m}_2(\xi) = -\psi_2/\xi$ , and  $\tilde{m}_3 = -\psi_3/\xi$ , and denote  $\check{\mathbf{m}}(\xi) = [\tilde{m}_1(\xi), \tilde{m}_2(\xi), \tilde{m}_3(\xi)]^\top$ . Then clearly we have  $\lim_{u \rightarrow +\infty} \check{\mathbf{m}}(\xi_u) = \mathbf{0}$ . By the definition of  $H_1(\mathbf{m}; \mathbf{q}, \boldsymbol{\mu})$ , with simple calculations, we can see that  $\lim_{u \rightarrow +\infty} H_1(\check{\mathbf{m}}(\xi_u); \mathbf{q}, \boldsymbol{\mu}) = q_4 \mu_{1,1}^2$ . Thus by (V.1), we have

$$|\xi_u \cdot [\tilde{m}_1 - F_1(\check{\mathbf{m}}(\xi_u); \xi_u, \mathbf{q}, \boldsymbol{\mu})]| = \psi_1 \cdot \left| \frac{q_2 \mu_{1,*}^2 + H_1(\check{\mathbf{m}}(\xi_u); \mathbf{q}, \boldsymbol{\mu})}{\xi_u - q_2 \mu_{1,*}^2 - H_1(\check{\mathbf{m}}(\xi_u); \mathbf{q}, \boldsymbol{\mu})} \right| = O_u\left(\frac{1}{u}\right).$$

Similarly, we can show that  $|\xi_u \cdot [\tilde{m}_j - F_j(\check{\mathbf{m}}(\xi_u); \xi_u, \mathbf{q}, \boldsymbol{\mu})]| = O_u(1/u)$ ,  $j = 2, 3$ . Therefore we have

$$\xi_u \cdot \|\check{\mathbf{m}}(\xi_u) - \mathbf{F}(\check{\mathbf{m}}(\xi_u); \xi_u, \mathbf{q}, \boldsymbol{\mu})\|_2 = O_u(u^{-1}). \quad (\text{V.2})$$

Moreover, by Lemma III.1, there exists a sufficiently large  $\xi_0$  such that for any  $K \geq \xi_0$ ,  $\mathbf{F}(\cdot; \xi_u, \mathbf{q}, \boldsymbol{\mu})$  is 1/2-Lipschitz on the domain  $\mathbb{D}(2\psi_1/\xi_0) \times \mathbb{D}(2\psi_2/\xi_0) \times \mathbb{D}(2\psi_3/\xi_0)$ . Therefore for sufficiently large  $K$ ,

$$\begin{aligned} & \|\check{\mathbf{m}}(\xi_u) - \mathbf{m}(\xi_u)\|_2 \\ &= \|\mathbf{F}(\check{\mathbf{m}}(\xi_u); \xi_u, \mathbf{q}, \boldsymbol{\mu}) - \mathbf{F}(\mathbf{m}(\xi_u); \xi_u, \mathbf{q}, \boldsymbol{\mu}) + \check{\mathbf{m}}(\xi_u) - \mathbf{F}(\check{\mathbf{m}}(\xi_u); \xi_u, \mathbf{q}, \boldsymbol{\mu})\|_2 \\ &\leq \|\mathbf{F}(\check{\mathbf{m}}(\xi_u); \xi_u, \mathbf{q}, \boldsymbol{\mu}) - \mathbf{F}(\mathbf{m}(\xi_u); \xi_u, \mathbf{q}, \boldsymbol{\mu})\|_2 + \|\check{\mathbf{m}}(\xi_u) - \mathbf{F}(\check{\mathbf{m}}(\xi_u); \xi_u, \mathbf{q}, \boldsymbol{\mu})\|_2 \\ &\leq \|\check{\mathbf{m}}(\xi_u) - \mathbf{m}(\xi_u)\|_2/2 + \|\check{\mathbf{m}}(\xi_u) - \mathbf{F}(\check{\mathbf{m}}(\xi_u); \xi_u, \mathbf{q}, \boldsymbol{\mu})\|_2, \end{aligned}$$

where the first equality follows by the definition of  $\mathbf{m}(\xi_u)$  as the fixed point of  $\mathbf{F}(\cdot; \xi_u, \mathbf{q}, \boldsymbol{\mu})$ , the first inequality follows by triangle inequality, and the second inequality follows by the 1/2-Lipschitz continuity of  $\mathbf{F}(\cdot; \xi_u, \mathbf{q}, \boldsymbol{\mu})$  on the domain  $\mathbb{D}(2\psi_1/\xi_0) \times \mathbb{D}(2\psi_2/\xi_0) \times \mathbb{D}(2\psi_3/\xi_0)$  (note that  $\mathbf{m}(\xi_u)$  is automatically in this domain according to Lemma III.1, and  $\check{\mathbf{m}}(\xi_u)$  is also in this domain by its

definition). Rearranging terms then gives

$$\|\check{\mathbf{m}}(\xi_u) - \mathbf{m}(\xi_u)\|_2 \leq 2\|\check{\mathbf{m}}(\xi_u) - \mathbf{F}(\check{\mathbf{m}}(\xi_u); \xi_u, \mathbf{q}, \boldsymbol{\mu})\|_2. \quad (\text{V.3})$$

Thus for  $j = 1, 2, 3$ , we have

$$\begin{aligned} |m_j(\xi_r + iu) \cdot (\xi_r + iu) + \psi_j| &= \xi_u \cdot |m_j(\xi_u) - \check{m}_j(\xi_u)| \\ &\leq 2\xi_u \cdot \|\check{\mathbf{m}}(\xi_u) - \mathbf{F}(\check{\mathbf{m}}(\xi_u); \xi_u, \mathbf{q}, \boldsymbol{\mu})\|_2 \\ &= O_u(u^{-1}), \end{aligned}$$

where the first inequality follows by (V.3), and the second equality follows by (V.2). This completes the proof.  $\blacksquare$

The following lemma shows the asymptotics of the functions  $G_d(iu; \mathbf{q}, \boldsymbol{\mu})$  and  $g(iu; \mathbf{q}, \boldsymbol{\mu})$  (defined in Definition A.3 and (A.3) respectively) as  $u$  goes to infinity.

**Lemma V.2.** *Let  $G_d(\xi; \mathbf{q}, \boldsymbol{\mu})$  be defined in Definition A.3 and  $g(\xi; \mathbf{q}, \boldsymbol{\mu})$  defined in (A.3). The following limits hold:*

$$\begin{aligned} \lim_{u \rightarrow +\infty} \sup_{d \geq 1} \mathbb{E}|G_d(iu; \mathbf{q}, \boldsymbol{\mu}) - (\psi_1 + \psi_2 + \psi_3) \log(-iu)| &= 0, \\ \lim_{u \rightarrow +\infty} |g(iu; \mathbf{q}, \boldsymbol{\mu}) - (\psi_1 + \psi_2 + \psi_3) \log(-iu)| &= 0. \end{aligned}$$

**Proof** [Proof of Lemma V.2] The real and imaginary parts of  $G_d(iu; \mathbf{q}, \boldsymbol{\mu}) - (\psi_1 + \psi_2 + \psi_3) \log(-iu)$  are

$$\begin{aligned} \left| \Re \left[ \frac{1}{P} \sum (\log(\lambda_i(\mathbf{A}) - iu) - \log(-iu)) \right] \right| &= \frac{1}{2P} \sum_{i=1}^P \log(1 + \lambda_i(A)^2/u^2) \leq \frac{\|\mathbf{A}\|_F^2}{2Pu^2}, \\ \left| \Im \left[ \frac{1}{P} \sum (\log(\lambda_i(\mathbf{A}) - iu) - \log(-iu)) \right] \right| &= \frac{1}{P} \sum_{i=1}^P \arctan(\lambda_i(\mathbf{A})/u) \leq \frac{\|A\|_F}{P^{1/2}u}. \end{aligned}$$

By the definition of the linear pencil matrix  $\mathbf{A}$ , it is easy to see that  $\frac{1}{P} \mathbb{E}[\|\mathbf{A}\|_F^2] = O_d(1)$ , thus

$$\lim_{u \rightarrow +\infty} \sup_{d \geq 1} \mathbb{E}|G_d(iu; \mathbf{q}, \boldsymbol{\mu}) - (\psi_1 + \psi_2 + \psi_3) \log(-iu)| = 0.$$

For the asymptotics of  $g(iu; \mathbf{q}, \boldsymbol{\mu})$ , note that

$$L(\xi, z_1, z_2, z_3; \mathbf{q}, \boldsymbol{\mu}) = L_1(z_1, z_2, z_3; \mathbf{q}, \boldsymbol{\mu}) + L_2(\xi, z_1, z_2, z_3; \mathbf{q}, \boldsymbol{\mu}),$$

where

$$\begin{aligned} L_1(z_1, z_2, z_3; \mathbf{q}, \boldsymbol{\mu}) &= \log \left[ (1 + \mu_{1,1}^2 z_1 q_4 + \mu_{2,1}^2 z_2 q_4)(1 + z_3 q_5) - \mu_{1,1}^2 (1 + q_1)^2 z_1 z_3 - \mu_{2,1}^2 (1 + q_1)^2 z_2 z_3 \right] \\ &\quad - \mu_{1,*}^2 z_1 z_3 - \mu_{2,*}^2 z_2 z_3 + q_2 \mu_{1,*}^2 z_1 + q_2 \mu_{1,*}^2 z_2 + q_3 z_3, \\ L_2(\xi, z_1, z_2, z_3; \mathbf{q}, \boldsymbol{\mu}) &= -\psi_1 \log(z_1/\psi_1) - \psi_2 \log(z_2/\psi_2) - \psi_3 \log(z_3/\psi_3) - \xi(z_1 + z_2 + z_3) - \psi_1 - \psi_2 - \psi_3. \end{aligned}$$

We now calculate the limits of  $L_1(m_1(iu), m_2(iu), m_3(iu); \mathbf{q}, \boldsymbol{\mu})$  and  $L_2(iu, m_1(iu), m_2(iu), m_3(iu); \mathbf{q}, \boldsymbol{\mu})$  separately. For  $L_1$ , by Lemma V.1, we have

$$\lim_{u \rightarrow +\infty} m_1(iu) = 0, \quad \lim_{u \rightarrow +\infty} m_2(iu) = 0, \quad \lim_{u \rightarrow +\infty} m_3(iu) = 0,$$

which immediately implies that

$$\lim_{u \rightarrow +\infty} L_1(m_1(iu), m_2(iu), m_3(iu); \mathbf{q}, \boldsymbol{\mu}) = 0.$$

For  $L_2$ , note that by Lemma V.1 we also have

$$\lim_{u \rightarrow +\infty} |m_1(iu)iu + \psi_1| = 0, \quad \lim_{u \rightarrow +\infty} |m_2(iu)iu + \psi_2| = 0, \quad \lim_{u \rightarrow +\infty} |m_3(iu)iu + \psi_3| = 0.$$

Therefore,

$$\begin{aligned} & |L_2(iu, m_1(iu), m_2(iu), m_3(iu); \mathbf{q}, \boldsymbol{\mu}) - (\psi_1 + \psi_2 + \psi_3) \log(-iu)| \\ & \leq \psi_1 |\log(-ium_1(iu)/\psi_1)| + \psi_2 |\log(-ium_2(iu)/\psi_2)| + \psi_3 |\log(-ium_3(iu)/\psi_3)| \\ & \quad + |\psi_1 + ium_1(iu)| + |\psi_2 + ium_2(iu)| + |\psi_3 + ium_3(iu)| \rightarrow 0, \end{aligned}$$

which completes the proof. ■

The following lemma gives an important identity between  $g(\xi; \mathbf{q}, \boldsymbol{\mu})$  and  $m(\xi; \mathbf{q}, \boldsymbol{\mu})$ .

**Lemma V.3.** *For all  $\xi \in \mathbb{C}_+$ , it holds that*

$$\frac{\partial g}{\partial \xi}(\xi; \mathbf{q}, \boldsymbol{\mu}) = -(m_1 + m_2 + m_3)(\xi; \mathbf{q}, \boldsymbol{\mu}) = -m(\xi; \mathbf{q}, \boldsymbol{\mu}).$$

**Proof** [Proof of Lemma V.3] By the definition of  $L(\xi, z_1, z_2, z_3; \mathbf{q}, \boldsymbol{\mu})$ , it is easy to see that

$$\begin{aligned} & \partial_{z_1} L(\xi, z_1, z_2, z_3; \mathbf{q}, \boldsymbol{\mu}) \\ & = -\mu_{1,*}^2 z_3 + q_2 \mu_{1,*}^2 - \psi_1/z_1 - \xi \\ & \quad + \frac{\mu_{1,1}^2 q_4 (1 + z_3 q_5) - \mu_{1,1}^2 (1 + q_1)^2 z_3}{(1 + \mu_{1,1}^2 z_1 q_4 + \mu_{2,1}^2 z_2 q_4)(1 + z_3 q_5) - \mu_{2,1}^2 (1 + q_1)^2 z_2 z_3 - \mu_{1,1}^2 (1 + q_1)^2 z_1 z_3} \\ & = \psi_1 \left( \frac{1}{F_1(\mathbf{z})} - \frac{1}{z_1} \right), \end{aligned}$$

$$\begin{aligned} & \partial_{z_2} L(\xi, z_1, z_2, z_3; \mathbf{q}, \boldsymbol{\mu}) \\ & = -\mu_{2,*}^2 z_3 + q_2 \mu_{2,*}^2 - \psi_2/z_2 - \xi \\ & \quad + \frac{\mu_{2,1}^2 q_4 (1 + z_3 q_5) - \mu_{2,1}^2 (1 + q_1)^2 z_3}{(1 + \mu_{2,1}^2 z_2 q_4 + \mu_{1,1}^2 z_1 q_4)(1 + z_3 q_5) - \mu_{1,1}^2 (1 + q_1)^2 z_1 z_3 - \mu_{2,1}^2 (1 + q_1)^2 z_2 z_3} \\ & = \psi_2 \left( \frac{1}{F_2(\mathbf{z})} - \frac{1}{z_2} \right), \end{aligned}$$

$$\begin{aligned}
 \partial_{z_3} L(\xi, z_1, z_2, z_3; \mathbf{q}, \boldsymbol{\mu}) &= -\mu_{1,*}^2 z_1 - \mu_{2,*}^2 z_2 + q_3 - \psi_3/z_3 - \xi \\
 &\quad + \frac{q_5(1 + \mu_{1,1}^2 z_1 q_4 + \mu_{2,1}^2 z_2 q_4) - \mu_{2,1}^2(1 + q_1)^2 z_2 - \mu_{1,1}^2(1 + q_1)^2 z_1}{(1 + \mu_{2,1}^2 z_2 q_4 + \mu_{1,1}^2 z_1 q_4)(1 + z_3 q_5) - \mu_{1,1}^2(1 + q_1)^2 z_1 z_3 - \mu_{2,1}^2(1 + q_1)^2 z_2 z_3} \\
 &= \psi_3 \left( \frac{1}{\mathbf{F}_3(\mathbf{z})} - \frac{1}{z_3} \right),
 \end{aligned}$$

where we utilize the definition of  $\mathbf{F}$  in Definition A.5 and write  $\mathbf{z} = [z_1, z_2, z_3]$ . Then by Proposition A.6, we have

$$\nabla_{\mathbf{z}} L(\xi, \mathbf{z}; \mathbf{q}, \boldsymbol{\mu})|_{\mathbf{z}=\mathbf{m}} \equiv \mathbf{0}$$

for all  $\xi \in \mathbb{C}_+$ . By the formula of implicit differentiation, we have

$$\begin{aligned}
 \frac{\partial g(\xi; \mathbf{q}, \boldsymbol{\mu})}{\partial \xi} &= [\langle \nabla_{\mathbf{z}} L(\xi, \mathbf{z}; \mathbf{q}, \boldsymbol{\mu})|_{\mathbf{z}=\mathbf{m}}, \partial_{\xi} \mathbf{m} \rangle + \partial_{\xi} L(\xi, \mathbf{z}; \mathbf{q}, \boldsymbol{\mu})|_{\mathbf{z}=\mathbf{m}}] \\
 &= 0 + \frac{d L(\xi, \mathbf{z}; \mathbf{q}, \boldsymbol{\mu})}{d \xi} \Big|_{\mathbf{z}=\mathbf{m}} = -m(\xi; \mathbf{q}, \boldsymbol{\mu}).
 \end{aligned}$$

This completes the proof of Lemma V.3. ■

The following lemma further shows that the derivatives of  $G_d$  and  $g_d$  are asymptotically bounded.

**Lemma V.4.** *For fixed  $\xi \in \mathbb{C}_+$ , the following limits hold:*

$$\begin{aligned}
 \limsup_{d \rightarrow +\infty} \{ \mathbb{E} \sup_{\mathbf{q} \in \mathcal{Q}} \|\nabla_{\mathbf{q}} G_d(\xi; \mathbf{q}, \boldsymbol{\mu})\|_2 \} + \sup_{\mathbf{q} \in \mathcal{Q}} \|\nabla_{\mathbf{q}} g(\xi; \mathbf{q}, \boldsymbol{\mu})\|_2 &< +\infty, \\
 \limsup_{d \rightarrow +\infty} \{ \mathbb{E} \sup_{\mathbf{q} \in \mathcal{Q}} \|\nabla_{\mathbf{q}}^2 G_d(\xi; \mathbf{q}, \boldsymbol{\mu})\|_{\text{op}} \} + \sup_{\mathbf{q} \in \mathcal{Q}} \|\nabla_{\mathbf{q}}^2 g(\xi; \mathbf{q}, \boldsymbol{\mu})\|_{\text{op}} &< +\infty, \\
 \limsup_{d \rightarrow +\infty} \{ \mathbb{E} \sup_{\mathbf{q} \in \mathcal{Q}} \|\nabla_{\mathbf{q}}^3 G_d(\xi; \mathbf{q}, \boldsymbol{\mu})\|_{\text{op}} \} + \sup_{\mathbf{q} \in \mathcal{Q}} \|\nabla_{\mathbf{q}}^3 g(\xi; \mathbf{q}, \boldsymbol{\mu})\|_{\text{op}} &< +\infty.
 \end{aligned}$$

**Proof** [Proof of Lemma V.4] Let  $\xi = \xi_r + iu$ , where  $\xi_r \in \mathbb{R}$  and  $u \in \mathbb{R}_+$  are both fixed. We also denote

$$\begin{aligned}
 \mathbf{S}_1 &= \begin{bmatrix} \mathbf{0} & \mathbf{0} & \tilde{\mathbf{Z}}_1^\top \\ \mathbf{0} & \mathbf{0} & \tilde{\mathbf{Z}}_2^\top \\ \tilde{\mathbf{Z}}_1 & \tilde{\mathbf{Z}}_2 & \mathbf{0} \end{bmatrix} = \begin{bmatrix} \mathbf{0} & \tilde{\mathbf{Z}} \\ \tilde{\mathbf{Z}} & \mathbf{0} \end{bmatrix}, \quad \mathbf{S}_2 = \begin{bmatrix} \mathbf{M}_* \mathbf{M}_* & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix}, \\
 \mathbf{S}_3 &= \begin{bmatrix} \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{I}_n \end{bmatrix}, \quad \mathbf{S}_4 = \begin{bmatrix} \mathbf{M}_1 \frac{\boldsymbol{\Theta} \boldsymbol{\Theta}^\top}{d} \mathbf{M}_1 & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix}, \quad \mathbf{S}_5 = \begin{bmatrix} \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \frac{\mathbf{X} \mathbf{X}^\top}{d} \end{bmatrix}.
 \end{aligned}$$

Then  $\mathbf{S}_1, \dots, \mathbf{S}_5$  are not related to  $\mathbf{q}$ , and it is easy to see that

$$\limsup_{d \rightarrow +\infty} \mathbb{E} \sup_{\mathbf{q} \in \mathcal{Q}} \|\mathbf{S}_i\|_{\text{op}}^{2k} < +\infty$$

for any fixed  $k \in \mathbb{N}$ . Moreover, define  $\mathbf{R} = \mathbf{R}(\mathbf{q}) = (\mathbf{A}(\mathbf{q}) - \xi_r \mathbf{I}_P - iu \mathbf{I}_P)^{-1}$ . Since  $\mathbf{A}(\mathbf{q})$  is a real symmetric matrix, the imaginary parts in the eigenvalues of  $\mathbf{A}(\mathbf{q}) - \xi_r \mathbf{I}_P - iu \mathbf{I}_P$  are all  $-iu$ , and

hence we deterministically have

$$\sup_{\mathbf{q}} \|\mathbf{R}\|_{\text{op}} \leq 1/u. \quad (\text{V.4})$$

Therefore, by (II.2), (II.3) and the definition of the linear pencil matrix  $\mathbf{A}$ , we have

$$\begin{aligned} \mathbb{E} \sup_{\mathbf{q} \in \mathcal{Q}} |\partial_{q_i} G_d(\xi; \mathbf{q})| &= \mathbb{E} \sup_{\mathbf{q} \in \mathcal{Q}} \frac{1}{d} |\text{tr}(\mathbf{R}\mathbf{S}_i)| \leq \mathbb{E} \sup_{\mathbf{q} \in \mathcal{Q}} \frac{1}{u} [\|\mathbf{S}_i\|_{\text{op}}] = O_d(1), \\ \mathbb{E} \sup_{\mathbf{q} \in \mathcal{Q}} |\partial_{q_i, q_j}^2 G_d(\xi; \mathbf{q})| &= \mathbb{E} \sup_{\mathbf{q} \in \mathcal{Q}} \frac{1}{d} |\text{tr}(\mathbf{R}\mathbf{S}_i \mathbf{R}\mathbf{S}_j)| \leq \left( \mathbb{E} \sup_{\mathbf{q} \in \mathcal{Q}} \frac{1}{u^2} [\|\mathbf{S}_i\|_{\text{op}}^2 \|\mathbf{S}_j\|_{\text{op}}^2] \right)^{\frac{1}{2}} = O_d(1). \end{aligned}$$

Similarly, for the third order derivatives, we also have

$$\begin{aligned} \mathbb{E} \sup_{\mathbf{q} \in \mathcal{Q}} |\partial_{q_i, q_j, q_l}^3 G_d(\xi; \mathbf{q})| &= \mathbb{E} \left\{ \sup_{\mathbf{q} \in \mathcal{Q}} \frac{1}{d} |\text{tr}(\mathbf{R}\mathbf{S}_i \mathbf{R}\mathbf{S}_j \mathbf{R}\mathbf{S}_l) + \text{tr}(\mathbf{R}\mathbf{S}_i \mathbf{R}\mathbf{S}_l \mathbf{R}\mathbf{S}_j)| \right\} \\ &\leq \frac{2}{u^3} \left( \mathbb{E} \sup_{\mathbf{q} \in \mathcal{Q}} [\|\mathbf{S}_i\|_{\text{op}}^4 \|\mathbf{S}_j\|_{\text{op}}^4 \|\mathbf{S}_l\|_{\text{op}}^4] \right)^{\frac{1}{4}} = O_d(1). \end{aligned}$$

This completes the Proof for  $G_d(\xi; \mathbf{q}, \boldsymbol{\mu})$ . As for  $g(\xi; \mathbf{q}, \boldsymbol{\mu})$ , we first show that if  $\mathbf{q}_1 \neq \mathbf{q}_2$ , the following property holds:

$$\begin{aligned} \frac{|m(\xi; \mathbf{q}_1, \boldsymbol{\mu}) - m(\xi; \mathbf{q}_2, \boldsymbol{\mu})|}{\|\mathbf{q}_1 - \mathbf{q}_2\|_2} &= \frac{\left| \lim_{d \rightarrow \infty} \mathbb{E}(M_d(\xi; \mathbf{q}_1, \boldsymbol{\mu}) - M_d(\xi; \mathbf{q}_2, \boldsymbol{\mu})) \right|}{\|\mathbf{q}_1 - \mathbf{q}_2\|_2} \\ &= \lim_{d \rightarrow \infty} \frac{\left| \mathbb{E}[\text{tr}(\mathbf{R}(\mathbf{q}_1) - \mathbf{R}(\mathbf{q}_2))] \right|}{d \|\mathbf{q}_1 - \mathbf{q}_2\|_2} \\ &= \lim_{d \rightarrow \infty} \frac{\left| \mathbb{E}[\text{tr}(\mathbf{R}(\mathbf{q}_1)(\mathbf{A}(\mathbf{q}_1) - \mathbf{A}(\mathbf{q}_2))\mathbf{R}(\mathbf{q}_2))] \right|}{d \|\mathbf{q}_1 - \mathbf{q}_2\|_2} \\ &\leq \lim_{d \rightarrow \infty} \frac{P}{d} \cdot \mathbb{E} \frac{\|\mathbf{A}(\mathbf{q}_1) - \mathbf{A}(\mathbf{q}_2)\|_{\text{op}}}{u^2 \cdot \|\mathbf{q}_1 - \mathbf{q}_2\|_2} < +\infty, \end{aligned}$$

where the first equality follows by Proposition A.6, the third equality follows by the identity  $\mathbf{A}^{-1} - \mathbf{B}^{-1} = \mathbf{A}^{-1}(\mathbf{B} - \mathbf{A})\mathbf{B}^{-1}$  for any invertible matrices  $\mathbf{A}, \mathbf{B}$ , the first inequality follows by  $|\text{tr}(\mathbf{A}\mathbf{B})| \leq P \cdot \|\mathbf{A}\|_{\text{op}} \|\mathbf{B}\|_{\text{op}}$  for all  $\mathbf{A}, \mathbf{B} \in \mathbb{C}^{P \times P}$  and (V.4), and the last inequality follows by the linearity of  $\mathbf{A}(\mathbf{q})$  in  $\mathbf{q}$ . Therefore we have  $\sup_{\mathbf{q} \in \mathcal{Q}} \|\nabla_{\mathbf{q}} m(\xi; \mathbf{q}, \boldsymbol{\mu})\|_2 < +\infty$ . Similarly, we can also show that  $\sup_{\mathbf{q} \in \mathcal{Q}} \|\nabla_{\mathbf{q}}^j m(\xi; \mathbf{q}, \boldsymbol{\mu})\|_{\text{op}} < +\infty$  for any fixed  $\xi \in \mathbb{C}_+$  and  $j = 2, 3$ . Moreover, by Lemma V.3, we have

$$\frac{d}{d\xi} g(\xi; \mathbf{q}, \boldsymbol{\mu}) = -m(\xi; \mathbf{q}, \boldsymbol{\mu}).$$

Then  $\sup_{\mathbf{q} \in \mathcal{Q}} \|\nabla_{\mathbf{q}}^j m(\xi; \mathbf{q}, \boldsymbol{\mu})\|_{\text{op}} < +\infty$  indicates that  $\sup_{\mathbf{q} \in \mathcal{Q}} \|\nabla_{\mathbf{q}}^j g(\xi; \mathbf{q}, \boldsymbol{\mu})\|_{\text{op}} < +\infty$ . This completes the proof of Lemma V.4.  $\blacksquare$



Finally, we present a classic result which shows that the derivatives of a function in a compact region can be upper bounded by the function value and the second derivatives of the function in the region.

**Lemma V.5** (lemma 11.4 in [Mei and Montanari \(2022\)](#)). *Let  $f \in C^2([a, b])$ . Then we have*

$$\sup_{x \in [a, b]} |f'(x)| \leq \left| \frac{f(a) - f(b)}{a - b} \right| + \frac{1}{2} \sup_{x \in [a, b]} |f''(x)| \cdot |a - b|.$$

Moreover, letting,  $f \in C^2(\mathbf{B}(\mathbf{x}_0, 2r))$  where  $\mathbf{B}(\mathbf{x}_0, 2r) = \{\mathbf{x} \in \mathbb{R}^d : \|\mathbf{x} - \mathbf{x}_0\|_2 \leq r\}$  with a point  $\mathbf{x}_0$ , we have

$$\sup_{\mathbf{x} \in \mathbf{B}(\mathbf{x}_0, 2r)} \|\nabla f(\mathbf{x})\|_2 \leq r^{-1} \sup_{\mathbf{x} \in \mathbf{B}(\mathbf{x}_0, 2r)} |f(\mathbf{x})| + 2r \sup_{\mathbf{x} \in \mathbf{B}(\mathbf{x}_0, 2r)} \|\nabla^2 f(\mathbf{x})\|_{\text{op}}.$$

## V.2 Completion of the proof

By Lemma V.3, we have  $\frac{\partial g}{\partial \xi}(\xi; \mathbf{q}, \boldsymbol{\mu}) = -m(\xi; \mathbf{q}, \boldsymbol{\mu})$ . Hence, for  $\xi \in \mathbb{C}_+$ ,  $u \in \mathbb{R}_+$ , and any compact continuous path  $c(\xi, iu)$  connecting  $\xi$  and  $iu$ , we have

$$g(\xi; \mathbf{q}, \boldsymbol{\mu}) - g(iu; \mathbf{q}, \boldsymbol{\mu}) = \int_{c(\xi, iu)} m(x; \mathbf{q}, \boldsymbol{\mu}) dx.$$

Moreover, from Definition A.3, we also have  $\frac{dG_d(\xi; \mathbf{q}, \boldsymbol{\mu})}{d\xi} = -M_d(\xi; \mathbf{q}, \boldsymbol{\mu})$ , and

$$G_d(\xi; \mathbf{q}, \boldsymbol{\mu}) - G_d(iu; \mathbf{q}, \boldsymbol{\mu}) = \int_{c(\xi, iu)} M_d(x; \mathbf{q}, \boldsymbol{\mu}) dx.$$

The two equations imply that

$$\begin{aligned} & \mathbb{E}[|G_d(\xi; \mathbf{q}, \boldsymbol{\mu}) - g(\xi; \mathbf{q}, \boldsymbol{\mu})|] \\ & \leq \mathbb{E}[|G_d(iu; \mathbf{q}, \boldsymbol{\mu}) - g(iu; \mathbf{q}, \boldsymbol{\mu})|] + \int_{c(\xi, iu)} \mathbb{E}[|M_d(x; \mathbf{q}, \boldsymbol{\mu}) - m(x; \mathbf{q}, \boldsymbol{\mu})|] dx. \end{aligned} \quad (\text{V.5})$$

Therefore, by taking supremum limit on both sides above and using Proposition A.6, we have

$$\begin{aligned} & \limsup_{d \rightarrow +\infty} \mathbb{E}[|G_d(\xi; \mathbf{q}, \boldsymbol{\mu}) - g(\xi; \mathbf{q}, \boldsymbol{\mu})|] \\ & \leq \limsup_{d \rightarrow +\infty} \mathbb{E}[|G_d(iu; \mathbf{q}, \boldsymbol{\mu}) - g(iu; \mathbf{q}, \boldsymbol{\mu})|] + \limsup_{d \rightarrow +\infty} \int_{c(\xi, iu)} \mathbb{E}[|M_d(x; \mathbf{q}, \boldsymbol{\mu}) - m(x; \mathbf{q}, \boldsymbol{\mu})|] dx \\ & = \limsup_{d \rightarrow +\infty} \mathbb{E}[|G_d(iu; \mathbf{q}, \boldsymbol{\mu}) - g(iu; \mathbf{q}, \boldsymbol{\mu})|]. \end{aligned}$$

Now the left hand side above does not depend on  $u$ . Moreover, by Lemma V.2, we have

$$\lim_{u \rightarrow +\infty} \limsup_{d \rightarrow +\infty} \mathbb{E}[|G_d(iu; \mathbf{q}, \boldsymbol{\mu}) - g(iu; \mathbf{q}, \boldsymbol{\mu})|] = 0.$$

Therefore

$$\lim_{d \rightarrow +\infty} \mathbb{E}[|G_d(\xi; \mathbf{q}, \boldsymbol{\mu}) - g(\xi; \mathbf{q}, \boldsymbol{\mu})|] = 0, \quad (\text{V.6})$$

which proves the first equality in Proposition A.7.

Next, we will prove the second and third equalities in Proposition A.7. Define  $V_d(\mathbf{q}) = G_d(\xi; \mathbf{q}, \boldsymbol{\mu}) - g(\xi; \mathbf{q}, \boldsymbol{\mu})$ . Then by Lemma V.5, we have

$$\sup_{\tilde{\mathbf{q}} \in \mathbf{B}(\mathbf{0}, \varepsilon)} \|\nabla V_d(\tilde{\mathbf{q}})\|_2 \leq \frac{\sup_{\tilde{\mathbf{q}} \in \mathbf{B}(\mathbf{0}, \varepsilon)} |V_d(\tilde{\mathbf{q}})|}{\varepsilon} + 2\varepsilon \sup_{\tilde{\mathbf{q}} \in \mathbf{B}(\mathbf{0}, \varepsilon)} \|\nabla^2 V_d(\tilde{\mathbf{q}})\|_{\text{op}}. \quad (\text{V.7})$$

By equation (V.6), Lemma V.4 and the covering number argument (similar to Section IV.6 and the proof in Section 11.2 in Mei and Montanari (2022)), we get that  $\lim_{d \rightarrow +\infty} \mathbb{E} \sup_{\tilde{\mathbf{q}} \in \mathcal{Q}_*} |V_d(\tilde{\mathbf{q}})| = 0$ . Again from Lemma V.4 and its proof, we already have

$$\lim_{d \rightarrow +\infty} \mathbb{E} \sup_{\tilde{\mathbf{q}} \in \mathbf{B}(\mathbf{0}, \varepsilon)} |\nabla^2 V_d(\tilde{\mathbf{q}})| < C,$$

for some absolute value  $C$ . Therefore, by (V.7), we have

$$\lim_{d \rightarrow +\infty} \mathbb{E} [\|\partial_{\mathbf{q}} G_d(iu; \mathbf{q}, \boldsymbol{\mu})|_{\mathbf{q}=\mathbf{0}} - \partial_{\mathbf{q}} g(iu; \mathbf{q}, \boldsymbol{\mu})|_{\mathbf{q}=\mathbf{0}}\|_2] \leq C\varepsilon.$$

Taking  $\varepsilon \rightarrow 0^+$ , we have

$$\lim_{d \rightarrow +\infty} \mathbb{E} [\|\partial_{\mathbf{q}} G_d(iu; \mathbf{q}, \boldsymbol{\mu})|_{\mathbf{q}=\mathbf{0}} - \partial_{\mathbf{q}} g(iu; \mathbf{q}, \boldsymbol{\mu})|_{\mathbf{q}=\mathbf{0}}\|_2] = 0.$$

This completes the proof of the second equality in Proposition A.7. The proof of the third equation in Proposition A.7 follows by a similar argument.

## VI. Proof of Proposition A.8

The existence result is obtained by directly checking that  $\mathbf{m}(\xi; \mathbf{0}, \boldsymbol{\mu})$  satisfies the two properties stated in Proposition A.8 as follows. The second property follows by the original definition of  $\mathbf{m}(\xi; \mathbf{q}, \boldsymbol{\mu})$  on  $\{\xi : \Im(\xi) \geq \xi_0\}$ . For the first property, by Proposition A.6, the analytic continuation of  $\mathbf{m}$  satisfies  $\mathbf{m}(\xi, \mathbf{0}, \boldsymbol{\mu}) \equiv \mathbf{F}[\mathbf{m}(\xi, \mathbf{0}, \boldsymbol{\mu}); \xi, \mathbf{0}, \boldsymbol{\mu}]$  for all  $\xi \in \mathbb{C}_+$ . This directly implies that  $\mathbf{m}(\xi, \mathbf{0}, \boldsymbol{\mu})$  solves the system (3.1) for all  $\xi \in \mathbb{C}_+$ , which verifies the first property in Proposition A.8. Moreover, by Proposition A.6,  $\mathbf{m}(\xi; \mathbf{0}, \boldsymbol{\mu})$  is analytic, and  $\mathbf{m}(\xi; \mathbf{0}, \boldsymbol{\mu}) \in \mathbb{C}_+^3$  for all  $\xi \in \mathbb{C}_+$ . Therefore  $\mathbf{m}(\xi; \mathbf{0}, \boldsymbol{\mu})$  is indeed an analytic function from  $\mathbb{C}_+$  to  $\mathbb{C}_+^3$ . This completes the proof of existence.

To show the uniqueness, suppose that an analytic function  $\boldsymbol{\nu} : \mathbb{C}_+ \rightarrow \mathbb{C}_+^3$  satisfies the properties satisfied in Proposition A.6. It then suffices to show that  $\boldsymbol{\nu}(\xi; \boldsymbol{\mu}) \equiv \mathbf{m}(\xi; \mathbf{0}, \boldsymbol{\mu})$ . By Lemma III.1, we clearly have  $\boldsymbol{\nu}(\xi; \boldsymbol{\mu}) = \mathbf{m}(\xi; \mathbf{0}, \boldsymbol{\mu})$  for all  $\xi$  with  $\Im(\xi) > \xi_0$ . Now since both  $\boldsymbol{\nu}(\xi; \boldsymbol{\mu})$  and  $\mathbf{m}(\xi; \mathbf{0}, \boldsymbol{\mu})$  are analytic on  $\mathbb{C}_+$ , the result follows by the uniqueness of analytic continuation.

We denote by  $\boldsymbol{\nu}^* = \boldsymbol{\nu}(\sqrt{\lambda} \cdot i; \boldsymbol{\mu})$ . From the definition of  $\overline{\mathbf{m}}_d(\xi)$  in Lemma IV.5, we easily get that the elements in  $\overline{\mathbf{m}}_d(\xi)$  are purely imaginary in the upper half-plane of  $\mathbb{C}$  when  $\mathbf{q} = \mathbf{0}$  and  $\xi = \sqrt{\lambda} \cdot i$ . (IV.23) further indicates that elements in  $\boldsymbol{\nu}^*$  are purely imaginary. Based on the proof above, we have  $\nu_j^*/i \in \mathbb{R}_+$ .

## VII. Proofs of Lemmas and Propositions in Appendix C

### VII.1 Proof of Lemma C.6

When  $\Im(\xi) \geq \xi_0$  for some sufficiently large  $\xi_0$ , we prove the existence and uniqueness of the solution by the Banach fixed point theorem. To do so, we want to show that

1.  $\mathbf{F}(\cdot; \mathbf{q}, \boldsymbol{\mu})$  maps domain  $\mathbb{D}(2\psi_1/\xi_0) \times \cdots \times \mathbb{D}(2\psi_K/\xi_0) \times \mathbb{D}(2\psi_{K+1}/\xi_0)$  into itself.
2.  $\mathbf{F}(\cdot; \mathbf{q}, \boldsymbol{\mu})$  is Lipschitz continuous with a Lipschitz constant smaller than 1.

For  $\mathbf{F}_1(\cdot; \mathbf{q}, \boldsymbol{\mu})$ , by Definition C.5, we have

$$\mathbf{F}_1(\mathbf{m}; \xi, \mathbf{q}, \boldsymbol{\mu}) = \frac{\psi_1}{-\xi + q_2\mu_{1,*}^2 + H_1(\mathbf{m}; \mathbf{q}, \boldsymbol{\mu})},$$

where

$$H_1(\mathbf{m}; \mathbf{q}, \boldsymbol{\mu}) = -\mu_{1,*}^2 m_3 + \frac{1}{m_1 + \frac{-\sum_{j=2}^K \mu_{j,1}^2 (1+q_1)^2 m_j m_{K+1} + (1 + \sum_{j=2}^K \mu_{j,1}^2 m_j q_4)(1 + m_{K+1} q_5)}{\mu_{1,1}^2 q_4 (1 + m_3 q_5) - \mu_{1,1}^2 (1+q_1)^2 m_3}} \quad (\text{VII.1})$$

Note that  $q_4, q_5 \leq (1+q_1)/2$ , it is easy to see that for  $r_0$  small enough and  $\mathbf{m} \in \mathbb{D}(r_0) \times \mathbb{D}(r_0) \times \mathbb{D}(r_0)$ , we have

$$|H_1(\mathbf{m}; \mathbf{q}, \boldsymbol{\mu})| \leq 2 + 2|q_4|\mu_{1,1}^2. \quad (\text{VII.2})$$

Now as long as  $\xi_0 \geq 4 + 4|q_4|\mu_{1,1}^2$ , it is clear that for  $\xi$  with  $\Im(\xi) \geq \xi_0$  we have

$$\Im(\xi) \geq \xi_0/2 + \xi_0/2 \geq \xi_0/2 + 2 + 2|q_4|\mu_{1,1}^2 \geq \xi_0/2 + |H_1(\mathbf{m}; \mathbf{q}, \boldsymbol{\mu})|, \quad (\text{VII.3})$$

where the last inequality follows by (VII.2). Therefore we have

$$\begin{aligned} |\mathbf{F}_1(\mathbf{m}; \xi, \mathbf{q}, \boldsymbol{\mu})| &\leq \frac{\psi_1}{|\Im(\xi - q_2\mu_{1,*}^2 - H_1(\mathbf{m}; \mathbf{q}, \boldsymbol{\mu}))|} \\ &\leq \frac{\psi_1}{\Im(\xi) - |H_1(\mathbf{m}; \mathbf{q}, \boldsymbol{\mu})|} \leq \frac{2\psi_1}{\xi_0}, \end{aligned}$$

where the inequalities follow from (VII.3).

Similarly, for  $\mathbf{F}_j$ ,  $j = 2, \dots, K+1$ , we also have  $|\mathbf{F}_j(\mathbf{m}; \xi, \mathbf{q}, \boldsymbol{\mu})| \leq 2\psi_j/\xi_0$  provided  $\xi_0 \geq 4 + 4\max_j\{|q_4|\mu_{j,1}^2, |q_5|\}$ . Therefore if  $\xi_0$  satisfies  $2\max\{\psi_1, \dots, \psi_{K+1}\}/\xi_0 \leq r_0$  and  $\xi_0 \geq 4 + 4\max_j\{|q_4|\mu_{j,1}^2, |q_5|\}$ , it is clear that  $\mathbf{F}$  maps domain  $\mathbb{D}(2\psi_1/\xi_0) \times \cdots \times \mathbb{D}(2\psi_{K+1}/\xi_0)$  into itself.

As for the Lipschitz continuity of  $\mathbf{F}(\cdot; \mathbf{q}, \boldsymbol{\mu})$ , note that

$$\nabla_{\mathbf{m}} \mathbf{F}_1(\mathbf{m}; \xi, \mathbf{q}, \boldsymbol{\mu}) = -\frac{\psi_1}{(-\xi + q_2\mu_{1,*}^2 + H_1(\mathbf{m}; \mathbf{q}, \boldsymbol{\mu}))^2} \cdot \nabla_{\mathbf{m}} H_1(\mathbf{m}; \mathbf{q}, \boldsymbol{\mu}).$$

It is easy to see that when  $\xi_0$  is sufficiently large,  $\|\nabla_{\mathbf{m}} H_1(\mathbf{m}; \mathbf{q}, \boldsymbol{\mu})\|_2 \leq C(\mathbf{q}, \boldsymbol{\mu})$  for all  $\mathbf{m} \in \mathbb{D}(2\psi_1/\xi_0) \times \cdots \times \mathbb{D}(2\psi_{K+1}/\xi_0)$ , where  $C(\mathbf{q}, \boldsymbol{\mu})$  is a constant that only depends on  $\mathbf{q}$  and  $\boldsymbol{\mu}$ . Thus when  $\xi_0$  is sufficiently large, for  $\xi$  with  $\Im(\xi) \geq \xi_0$ ,

$$\|\nabla_{\mathbf{m}} \mathbf{F}_1(\mathbf{m}; \xi, \mathbf{q}, \boldsymbol{\mu})\|_2 \leq \frac{C(\mathbf{q}, \boldsymbol{\mu}) \cdot \psi_1}{\Im(\xi) - |H_1(\mathbf{m}; \mathbf{q}, \boldsymbol{\mu})|} \leq \frac{4C(\mathbf{q}, \boldsymbol{\mu}) \cdot \psi_1}{\xi_0} \leq \frac{1}{4K},$$

where we again utilize (VII.3). We can apply the same argument for  $F_2, \dots, F_{K+1}$ , and conclude that  $\mathbf{F}$  is  $\frac{1}{2}$ -Lipschitz on  $\mathbf{m} \in \mathbb{D}(2\psi_1/\xi_0) \times \dots \times \mathbb{D}(2\psi_{K+1}/\xi_0)$ . Therefore by Banach fixed point theorem, there exists a unique fixed point of  $\mathbf{F}$ . Thus the fixed point of the functions defined in Definition C.5 exists and is unique.

## VII.2 Proof of Proposition C.7

Following the same argument as in Lemma IV.3, we may assume that all the elements in  $\bar{\mathbf{X}}$  and  $\bar{\Theta}_j$  are independently generated from standard normal  $N(0, 1)$ , and the activation functions are polynomials and centralized as  $\phi_j(x) = \sigma_j(x) - \mu_{j,0}$ . The linear pencil matrix of this Gaussian version is defined as

$$\bar{\mathbf{A}}(\mathbf{q}, \boldsymbol{\mu}) = \begin{bmatrix} q_2\mu_{1,*}^2 \mathbf{I}_{N_1} + q_4\mu_{1,1}^2 \frac{\bar{\Theta}_1 \bar{\Theta}_1^\top}{d} & \cdots & q_4\mu_{1,1}\mu_{K,1} \frac{\bar{\Theta}_1 \bar{\Theta}_K^\top}{d} & \bar{\mathbf{Z}}_1^\top \\ \vdots & \ddots & \vdots & \vdots \\ q_4\mu_{K,1}\mu_{1,1} \frac{\bar{\Theta}_K \bar{\Theta}_1^\top}{d} & \cdots & q_2\mu_{K,*}^2 \mathbf{I}_{N_K} + q_4\mu_{K,1}^2 \frac{\bar{\Theta}_K \bar{\Theta}_K^\top}{d} & \bar{\mathbf{Z}}_K^\top \\ \bar{\mathbf{Z}}_1 & \cdots & \bar{\mathbf{Z}}_K & q_3 \mathbf{I}_n + q_5 \frac{\bar{\mathbf{X}} \bar{\mathbf{X}}^\top}{d} \end{bmatrix}.$$

Here  $\bar{\mathbf{Z}}_j = \Phi_j \left( \bar{\mathbf{X}} \bar{\Theta}_j^\top / \sqrt{d} \right) / \sqrt{d} \in \mathbb{R}^{n \times N_j}$ , and  $\Phi_j(x)$  is defined as  $\Phi_j(x) = \phi_j(x) + q_1\mu_{j,1}x$ . Moreover, for  $j = 1, \dots, K$  and with  $G \sim N(0, 1)$ , we denote  $\phi_{j,0} \triangleq \mathbb{E}\{\Phi_j(G)\}$ ,  $\phi_{j,1} \triangleq \mathbb{E}\{G\Phi_j(G)\}$ ,  $\phi_{j,*} \triangleq \mathbb{E}\{\Phi_j(G)^2\} - \phi_{j,0}^2 - \phi_{j,1}^2$ . It is easy to see  $\phi_{j,0} = 0$ ,  $\phi_{j,1}^2 = \mu_{j,1}^2(1 + q_1)^2$ ,  $\phi_{j,*}^2 = \mu_{j,*}^2$ .

We remind readers that  $\mathcal{N}_j$  is the index set of units that use the  $j$ -th activation function  $\sigma_j$ . Define the following terms:

$$\begin{aligned} \bar{m}_{j,d}(\xi; \mathbf{q}, \boldsymbol{\mu}) &= \mathbb{E}[\bar{M}_{j,d}(\xi; \mathbf{q}, \boldsymbol{\mu})], \quad \bar{M}_{j,d}(\xi; \mathbf{q}, \boldsymbol{\mu}) = \frac{1}{d} \text{tr}_{\mathcal{N}_j} [\bar{\mathbf{A}}(\mathbf{q}, \boldsymbol{\mu}) - \xi \mathbf{I}_P]^{-1}, \quad j = 1, \dots, K \\ \bar{m}_{K+1,d}(\xi; \mathbf{q}, \boldsymbol{\mu}) &= \mathbb{E}[\bar{M}_{K+1,d}(\xi; \mathbf{q}, \boldsymbol{\mu})], \quad \bar{M}_{K+1,d}(\xi; \mathbf{q}, \boldsymbol{\mu}) = \frac{1}{d} \text{tr}_{[N+1:P]} [\bar{\mathbf{A}}(\mathbf{q}, \boldsymbol{\mu}) - \xi \mathbf{I}_P]^{-1}. \end{aligned}$$

With the same argument as in Lemma IV.3, we obtain

$$\mathbb{E} \left| \sum_{j=1}^{K+1} \bar{M}_{j,d}(\xi; \mathbf{q}, \boldsymbol{\mu}) - M_d(\xi; \mathbf{q}, \boldsymbol{\mu}) \right| = o_d(1), \quad \text{for any fixed } \xi \in \mathbb{C}_+.$$

Next, by contraction properties we have

$$\mathbb{E} \left| \bar{M}_{j,d}(\xi; \mathbf{q}, \boldsymbol{\mu}) - \bar{m}_{j,d}(\xi; \mathbf{q}, \boldsymbol{\mu}) \right| = o_d(1), \quad \text{for any fixed } \xi \in \mathbb{C}_+.$$

To study  $\bar{M}_d(\xi; \mathbf{q}, \boldsymbol{\mu})$ , which is the Stieltjes transform of the empirical eigenvalue distribution of  $\bar{\mathbf{A}}(\mathbf{q}, \boldsymbol{\mu})$ , it suffices to derive the resolvent equations for  $\bar{m}_d(\xi; \mathbf{q}, \boldsymbol{\mu})$  here. This is done by the following lemma.

**Lemma VII.1.** *Let  $\bar{\mathbf{m}}_d(\xi) = [\bar{m}_{1,d}(\xi), \dots, \bar{m}_{K+1,d}(\xi)]^\top$ . Then for any fixed  $\xi \in \mathbb{C}_+$ , the following property holds:*

$$\|\bar{\mathbf{m}}_d(\xi) - \mathbf{F}(\bar{\mathbf{m}}_d(\xi))\|_2 = o_d(1).$$

**Proof** [Proof of Lemma VII.1] Since  $\bar{\mathbf{m}}_d(\xi), \mathbf{F}(\bar{\mathbf{m}}_d(\xi)) \in \mathbb{C}^{K+1}$ , Lemma VII.1 essentially contains results showing that each element of  $\bar{\mathbf{m}}_d(\xi) - \mathbf{F}(\bar{\mathbf{m}}_d(\xi))$  is asymptotically zero. Since the proofs of the results are almost the same, we mainly focus on the proof of the first element  $\bar{m}_{1,d}$ . The proof still consists of three main steps similar to the proof of Lemma IV.5.

**Step 1.** We first use a leave-one-out argument to calculate  $\bar{m}_{1,d}$ . Let  $\bar{\mathbf{A}}_{\cdot, N_1}$  be the  $N_1^{\text{th}}$  column of  $\bar{\mathbf{A}}$ , with the  $N_1^{\text{th}}$  entry removed. We further denote  $\bar{\mathbf{B}} \in \mathbb{R}^{(P-1) \times (P-1)}$  the matrix from  $\bar{\mathbf{A}}$  by removing the  $N_1^{\text{th}}$  row and  $N_1^{\text{th}}$  column. From the Schur complement formula, we get

$$\bar{m}_{1,d} = \psi_1 \mathbb{E} \left( -\xi + q_2 \mu_{1,*}^2 + q_4 \mu_{1,1}^2 \|\bar{\boldsymbol{\theta}}_{N_1}\|_2^2 / d - \bar{\mathbf{A}}_{\cdot, N_1}^\top (\bar{\mathbf{B}} - \xi \mathbf{I}_{P-1})^{-1} \bar{\mathbf{A}}_{\cdot, N_1} \right)^{-1}. \quad (\text{VII.4})$$

We decompose the vectors  $\bar{\boldsymbol{\theta}}_a, a \in [N]$  and  $\bar{\mathbf{x}}_i, i \in [n]$  into components along the direction of  $\bar{\boldsymbol{\theta}}_{N_1}$  and the other orthogonal directions:

$$\begin{aligned} \bar{\boldsymbol{\theta}}_a &= \eta_a \frac{\bar{\boldsymbol{\theta}}_{N_1}}{\|\bar{\boldsymbol{\theta}}_{N_1}\|} + \tilde{\boldsymbol{\theta}}_a, \quad \langle \bar{\boldsymbol{\theta}}_{N_1}, \tilde{\boldsymbol{\theta}}_a \rangle = 0, \quad a \in [N] \setminus \{N_1\}, \\ \bar{\mathbf{x}}_i &= u_i \frac{\bar{\boldsymbol{\theta}}_{N_1}}{\|\bar{\boldsymbol{\theta}}_{N_1}\|} + \tilde{\mathbf{x}}_i, \quad \langle \bar{\boldsymbol{\theta}}_{N_1}, \tilde{\mathbf{x}}_i \rangle = 0, \quad i \in [n]. \end{aligned} \quad (\text{VII.5})$$

Note that for any  $a \in [N] \setminus \{N_1\}$  and  $i \in [n]$ ,  $\eta_a, u_i$  are standard Gaussian and are independent of  $\tilde{\boldsymbol{\theta}}_a$  and  $\tilde{\mathbf{x}}_i$ . Moreover,  $\tilde{\boldsymbol{\theta}}_a$  and  $\tilde{\mathbf{x}}_i$  are conditionally independent on each other given  $\bar{\boldsymbol{\theta}}_{N_1}$ , with  $\tilde{\boldsymbol{\theta}}_a, \tilde{\mathbf{x}}_i \sim N(0, P_\perp)$ , where  $P_\perp$  is the projector orthogonal to  $\bar{\boldsymbol{\theta}}_{N_1}$ . We then have  $\bar{\mathbf{A}}_{\cdot, N_1} = (\bar{\mathbf{A}}_{1, N_1}, \dots, \bar{\mathbf{A}}_{P-1, N_1})^\top \in \mathbb{R}^{P-1}$  with

$$\bar{\mathbf{A}}_{i, N_1} = \begin{cases} \frac{q_4 \mu_{1,1}^2 \eta_i}{d} \|\bar{\boldsymbol{\theta}}_{N_1}\|_2, & \text{if } i \in [1, N_1 - 1], \\ \frac{q_4 \mu_{1,1} \mu_{j,1} \eta_{i+1}}{d} \|\bar{\boldsymbol{\theta}}_{N_1}\|_2, & \text{if } i+1 \in \mathcal{N}_j, \quad j \geq 2, \\ \frac{1}{\sqrt{d}} \Phi_1 \left( \frac{1}{\sqrt{d}} u_{i-N+1} \|\bar{\boldsymbol{\theta}}_{N_1}\|_2 \right), & \text{if } i \geq N. \end{cases}$$

To calculate the resolvent equations, we need to further represent the matrix  $\bar{\mathbf{B}}$  in (VII.4) with  $\eta_a, \tilde{\boldsymbol{\theta}}_a, u_i$ , and  $\tilde{\mathbf{x}}_i$  for  $a \in [N] \setminus \{N_1\}$  and  $i \in [n]$ . Below we first list some additional notations for easier reference. Write  $\boldsymbol{\eta}_1 = [\eta_1, \dots, \eta_{N_1-1}] \in \mathbb{R}^{N_1-1}$ ,  $\boldsymbol{\eta}_j = (\eta_{\mathcal{N}_j}) \in \mathbb{R}^{N_j}, j = 2, \dots, K$ ,  $\boldsymbol{\eta} = [\boldsymbol{\eta}_1^\top, \dots, \boldsymbol{\eta}_K^\top]^\top \in \mathbb{R}^{N-1}$ ,  $\mathbf{u} = (u_1, \dots, u_n)^\top \in \mathbb{R}^n$ ,  $\tilde{\boldsymbol{\Theta}}_1 = [\tilde{\boldsymbol{\theta}}_1, \dots, \tilde{\boldsymbol{\theta}}_{N_1-1}]^\top$ ,  $\tilde{\boldsymbol{\Theta}}_j = [\tilde{\boldsymbol{\theta}}_{\mathcal{N}_j}]^\top$ ,

$$\tilde{\boldsymbol{\Theta}} = \begin{bmatrix} \tilde{\boldsymbol{\Theta}}_1 \\ \vdots \\ \tilde{\boldsymbol{\Theta}}_K \end{bmatrix} \in \mathbb{R}^{(N-1) \times d}, \quad \tilde{\mathbf{M}}_1 = \begin{bmatrix} \mu_{1,1} \mathbf{I}_{N_1-1} & & \\ & \ddots & \\ & & \mu_{K,1} \mathbf{I}_{N_K} \end{bmatrix}, \quad \tilde{\mathbf{M}}_* = \begin{bmatrix} \mu_{1,*} \mathbf{I}_{N_1-1} & & \\ & \ddots & \\ & & \mu_{K,*} \mathbf{I}_{N_K} \end{bmatrix}.$$

With  $\bar{\mathbf{B}}$  defined previously and (VII.5),  $\bar{\mathbf{B}}_{[1:N-1], [1:N-1]}$  is decomposed into

$$\bar{\mathbf{B}}_{[1:N-1], [1:N-1]} = q_2 \tilde{\mathbf{M}}_* \tilde{\mathbf{M}}_* + \frac{q_4}{d} \tilde{\mathbf{M}}_1 \tilde{\boldsymbol{\Theta}} \tilde{\boldsymbol{\Theta}}^\top \tilde{\mathbf{M}}_1 + \frac{q_4}{d} \tilde{\mathbf{M}}_1 \boldsymbol{\eta} \boldsymbol{\eta}^\top \tilde{\mathbf{M}}_1. \quad (\text{VII.6})$$

Moreover, for  $i, j \in [n]$  and  $a \in [N] \setminus \{N_1\}$ , we define

$$(\tilde{\mathbf{H}})_{ij} = \frac{1}{d} \langle \tilde{\mathbf{x}}_i, \tilde{\mathbf{x}}_j \rangle.$$

Then we could decompose  $\bar{\mathbf{B}}_{[N:P-1],[N:P-1]}$  into

$$\bar{\mathbf{B}}_{[N:P-1],[N:P-1]} = q_3 \mathbf{I}_n + q_5 \tilde{\mathbf{H}} + \frac{q_5}{d} \mathbf{u} \mathbf{u}^\top. \quad (\text{VII.7})$$

$\bar{\mathbf{B}}_{[N:P-1],[1:N-1]} = \bar{\mathbf{B}}_{[1:N-1],[N:P-1]}^\top$  holds due to the symmetry of  $\bar{\mathbf{B}}$ . For  $i, j \in [n]$  and  $a \in \mathcal{N}_j \setminus \{N_1\}$ , elementally we have

$$\begin{aligned} (\bar{\mathbf{Z}})_{i,a} &= \frac{1}{\sqrt{d}} \Phi_j \left( \frac{1}{\sqrt{d}} \langle \tilde{\mathbf{x}}_i, \tilde{\boldsymbol{\theta}}_a \rangle \right) = \frac{1}{\sqrt{d}} \Phi_j \left( \frac{1}{\sqrt{d}} \langle \tilde{\mathbf{x}}_i, \tilde{\boldsymbol{\theta}}_a \rangle + \frac{1}{d} u_i \eta_a \right) \\ &= \frac{1}{\sqrt{d}} \Phi_j \left( \frac{1}{\sqrt{d}} \langle \tilde{\mathbf{x}}_i, \tilde{\boldsymbol{\theta}}_a \rangle \right) + \frac{\phi_{j,1}}{d} u_i \eta_a + \frac{1}{\sqrt{d}} \left[ \Phi_{j,\perp} \left( \frac{1}{\sqrt{d}} \langle \tilde{\mathbf{x}}_i, \tilde{\boldsymbol{\theta}}_a \rangle + \frac{1}{\sqrt{d}} u_i \eta_a \right) - \Phi_{j,\perp} \left( \frac{1}{\sqrt{d}} \langle \tilde{\mathbf{x}}_i, \tilde{\boldsymbol{\theta}}_a \rangle \right) \right], \end{aligned}$$

where  $\Phi_{j,\perp}(x) = \Phi_j(x) - \phi_{j,1}x$ . By the symmetry of  $\bar{\mathbf{B}}$ , we can then decompose  $\bar{\mathbf{B}}_{[N:P-1],[1:N-1]}$  into

$$\bar{\mathbf{B}}_{[N:P-1],[1:N-1]} = \tilde{\mathbf{Z}} + \frac{1}{d} \mathbf{u} \boldsymbol{\eta} \mathbf{M}_\phi + [\mathbf{E}_1, \mathbf{E}_2]. \quad (\text{VII.8})$$

Here, we define

$$\begin{aligned} \tilde{\mathbf{Z}} &= [\tilde{\mathbf{Z}}_1, \dots, \tilde{\mathbf{Z}}_K], \quad (\tilde{\mathbf{Z}}_j)_{i,a} = \frac{1}{\sqrt{d}} \Phi_j \left( \frac{1}{\sqrt{d}} \langle \tilde{\mathbf{x}}_i, \tilde{\boldsymbol{\theta}}_a \rangle \right), \quad \mathbf{M}_\phi = \begin{bmatrix} \phi_{1,1} \mathbf{I}_{N_1-1} & & \\ & \ddots & \\ & & \phi_{K,1} \mathbf{I}_{N_K} \end{bmatrix} \\ (\mathbf{E}_j)_{i,a} &= \frac{1}{\sqrt{d}} \left[ \Phi_{j,\perp} \left( \frac{1}{\sqrt{d}} \langle \tilde{\mathbf{x}}_i, \tilde{\boldsymbol{\theta}}_a \rangle + \frac{1}{\sqrt{d}} u_i \eta_a \right) - \Phi_{j,\perp} \left( \frac{1}{\sqrt{d}} \langle \tilde{\mathbf{x}}_i, \tilde{\boldsymbol{\theta}}_a \rangle \right) \right]. \end{aligned}$$

Combined (VII.6), (VII.7) and (VII.8), we decompose  $\bar{\mathbf{B}}$  into

$$\bar{\mathbf{B}} = \tilde{\mathbf{B}} + \boldsymbol{\Delta} + \mathbf{E} \in \mathbb{R}^{(P-1) \times (P-1)},$$

where

$$\begin{aligned} \tilde{\mathbf{B}} &= \begin{bmatrix} q_2 \tilde{\mathbf{M}}_* \tilde{\mathbf{M}}_* + \frac{q_4}{d} \tilde{\mathbf{M}}_1 \tilde{\boldsymbol{\Theta}} \tilde{\boldsymbol{\Theta}}^\top \tilde{\mathbf{M}}_1 & \tilde{\mathbf{Z}}^\top \\ \tilde{\mathbf{Z}} & q_3 \mathbf{I}_n + q_5 \tilde{\mathbf{H}} \end{bmatrix} \\ &= \begin{bmatrix} q_2 \mu_{1,*}^2 \mathbf{I}_{N_1} + q_4 \mu_{1,1}^2 \frac{\tilde{\boldsymbol{\Theta}}_1 \tilde{\boldsymbol{\Theta}}_1^\top}{d} & \cdots & q_4 \mu_{1,1} \mu_{K,1} \frac{\tilde{\boldsymbol{\Theta}}_1 \tilde{\boldsymbol{\Theta}}_K^\top}{d} & \tilde{\mathbf{Z}}_1^\top \\ \vdots & \ddots & \vdots & \vdots \\ q_4 \mu_{K,1} \mu_{1,1} \frac{\tilde{\boldsymbol{\Theta}}_K \tilde{\boldsymbol{\Theta}}_1^\top}{d} & \cdots & q_2 \mu_{K,*}^2 \mathbf{I}_{N_K} + q_4 \mu_{K,1}^2 \frac{\tilde{\boldsymbol{\Theta}}_K \tilde{\boldsymbol{\Theta}}_K^\top}{d} & \tilde{\mathbf{Z}}_K^\top \\ \tilde{\mathbf{Z}}_1 & \cdots & \tilde{\mathbf{Z}}_K & q_3 \mathbf{I}_n + q_5 \tilde{\mathbf{H}} \end{bmatrix}, \\ \boldsymbol{\Delta} &= \begin{bmatrix} \frac{q_4}{d} \tilde{\mathbf{M}}_1 \boldsymbol{\eta} \boldsymbol{\eta}^\top \tilde{\mathbf{M}}_1 & \frac{1}{d} \mathbf{M}_\phi \boldsymbol{\eta} \mathbf{u}^\top \\ \frac{1}{d} \mathbf{u} \boldsymbol{\eta} \mathbf{M}_\phi & \frac{q_5}{d} \mathbf{u} \mathbf{u}^\top \end{bmatrix} \\ &= \begin{bmatrix} \frac{q_4 \mu_{1,1}^2}{d} \boldsymbol{\eta}_1 \boldsymbol{\eta}_1^\top & \cdots & \frac{q_4 \mu_{1,1} \mu_{K,1}}{d} \boldsymbol{\eta}_1 \boldsymbol{\eta}_K^\top & \frac{\phi_{1,1}}{d} \boldsymbol{\eta}_1 \mathbf{u}^\top \\ \vdots & \ddots & \vdots & \vdots \\ \frac{q_4 \mu_{K,1} \mu_{1,1}}{d} \boldsymbol{\eta}_K \boldsymbol{\eta}_1^\top & \cdots & \frac{q_4 \mu_{K,1}^2}{d} \boldsymbol{\eta}_K \boldsymbol{\eta}_K^\top & \frac{\phi_{K,1}}{d} \boldsymbol{\eta}_K \mathbf{u}^\top \\ \frac{\phi_{1,1}}{d} \mathbf{u} \boldsymbol{\eta}_1^\top & \cdots & \frac{\phi_{K,1}}{d} \mathbf{u} \boldsymbol{\eta}_K^\top & \frac{q_5}{d} \mathbf{u} \mathbf{u}^\top \end{bmatrix}, \quad \mathbf{E} = \begin{bmatrix} \mathbf{0} & \cdots & \mathbf{0} & \mathbf{E}_1^\top \\ \vdots & \vdots & \vdots & \vdots \\ \mathbf{0} & \cdots & \mathbf{0} & \mathbf{E}_K^\top \\ \mathbf{E}_1 & \cdots & \mathbf{E}_K & \mathbf{0} \end{bmatrix}. \end{aligned}$$

Clearly, by the definition of  $\tilde{\mathbf{B}}$ , the Stieltjes transform corresponding to  $\tilde{\mathbf{B}}$  shares the same asymptotics as the Stieltjes transform corresponding to  $\bar{\mathbf{A}}$ .

**Step 2.** Define  $w_2 = \left(-\xi + q_2\mu_{1,*}^2 + q_4\mu_{1,1}^2 - \bar{\mathbf{A}}_{\cdot,N_1}^\top (\tilde{\mathbf{B}} + \mathbf{\Delta} - \xi\mathbf{I}_{P-1})^{-1}\bar{\mathbf{A}}_{\cdot,N_1}\right)^{-1}$ . Similar to the argument in Section IV.2, we have  $\bar{m}_{1,d} = \psi_1 \mathbb{E}w_2 + o_d(1)$ .

**Step 3.** We calculate  $\mathbb{E}w_2$  by mathematical induction. Similar to Section IV.2, we give some notations which will be used in the following calculation on  $\mathbb{E}w_2$ . Let

$$\mathbf{v} = \bar{\mathbf{A}}_{\cdot,N_1}, \quad \mathbf{v}_i = \bar{\mathbf{A}}_{i,N_1} = \begin{cases} \frac{q_4\mu_{1,1}^2\eta_i}{d}\|\bar{\boldsymbol{\theta}}_{N_1}\|_2, & \text{if } i \in [1, N_1 - 1], \\ \frac{q_4\mu_{1,1}\mu_{j,1}\eta_{i+1}}{d}\|\bar{\boldsymbol{\theta}}_{N_1}\|_2, & \text{if } i \in \mathcal{N}_j - 1, j \geq 2, \\ \frac{1}{\sqrt{d}}\Phi_1\left(\frac{1}{\sqrt{d}}u_{i-N+1}\|\bar{\boldsymbol{\theta}}_{N_1}\|_2\right), & \text{if } i \geq N, \end{cases}$$

and

$$\mathbf{U} = \frac{1}{\sqrt{d}} \begin{bmatrix} \eta_1 & & & \\ & \eta_2 & & \\ & & \ddots & \\ & & & \eta_K \\ & & & & \mathbf{u} \end{bmatrix} \in \mathbb{R}^{(P-1) \times (K+1)}, \quad \mathbf{M} = \begin{bmatrix} q_4\mu_{1,1}^2 & \cdots & q_4\mu_{1,1}\mu_{K,1} & \phi_{1,1} \\ \vdots & \ddots & \vdots & \vdots \\ q_4\mu_{1,1}\mu_{K,1} & \cdots & q_4\mu_{K,1}^2 & \phi_{K,1} \\ \phi_{1,1} & \cdots & \phi_{K,1} & q_5 \end{bmatrix},$$

respectively. Then after direct calculation, we have the decomposition of  $\mathbf{\Delta}$  as

$$\mathbf{\Delta} = \mathbf{U}\mathbf{M}\mathbf{U}^\top.$$

Similar to (IV.16), we again get that

$$\begin{aligned} w_2 = & \left(-\xi + q_2\mu_{1,*}^2 + q_4\mu_{1,1}^2 - \mathbf{v}^\top (\tilde{\mathbf{B}} - \xi\mathbf{I}_{P-1})^{-1}\mathbf{v} \right. \\ & \left. + \mathbf{v}^\top (\tilde{\mathbf{B}} - \xi\mathbf{I}_{P-1})^{-1}\mathbf{U}(\mathbf{M}^{-1} + \mathbf{U}^\top (\tilde{\mathbf{B}} - \xi\mathbf{I}_{P-1})^{-1}\mathbf{U})^{-1}\mathbf{U}^\top (\tilde{\mathbf{B}} - \xi\mathbf{I}_{P-1})^{-1}\mathbf{v} \right)^{-1}. \end{aligned} \quad (\text{VII.9})$$

To continue the calculation, we still require to study the terms  $\mathbf{v}^\top (\tilde{\mathbf{B}} - \xi\mathbf{I}_{P-1})^{-1}\mathbf{v}$ ,  $\mathbf{v}^\top (\tilde{\mathbf{B}} - \xi\mathbf{I}_{P-1})^{-1}\mathbf{U}$  and  $\mathbf{U}^\top (\tilde{\mathbf{B}} - \xi\mathbf{I}_{P-1})^{-1}\mathbf{U}$  in the denominator of (VII.9). To do so, we note that  $\tilde{\mathbf{B}}$  is independent on  $\mathbf{v}$  and  $\mathbf{U}$ . Moreover, by the leave-one-out argument, the Stieltjes transform corresponding to  $\tilde{\mathbf{B}}$  shares the same asymptotics as the Stieltjes transform corresponding to  $\bar{\mathbf{A}}$ . Notice that  $\eta_i$  is independent on  $\tilde{\mathbf{B}}$  conditioned on  $\bar{\boldsymbol{\theta}}_{N_1}$ , and  $\tilde{\mathbf{B}}$  is independent on  $\bar{\boldsymbol{\theta}}_{N_1}$ . Similar to (IV.17)-(IV.19), we have

$$\mathbf{v}^\top (\tilde{\mathbf{B}} - \xi\mathbf{I}_{P-1})^{-1}\mathbf{v} = q_4^2\mu_{1,1}^2 \left( \sum_{j=1}^K \mu_{j,1}^2 \bar{m}_{j,d} \right) + (\phi_{1,1}^2 + \phi_{1,*}^2) \bar{m}_{K+1,d} + o_{\mathbb{P}}(1), \quad (\text{VII.10})$$

$$\mathbf{v}^\top (\tilde{\mathbf{B}} - \xi\mathbf{I}_{P-1})^{-1}\mathbf{U} = [q_4\mu_{1,1}^2 \bar{m}_{1,d} \quad \cdots \quad q_4\mu_{1,1}\mu_{K,1} \bar{m}_{K,d} \quad \phi_{1,1} \bar{m}_{K+1,d}] + o_{\mathbb{P}}(1), \quad (\text{VII.11})$$

$$\mathbf{U}^\top (\tilde{\mathbf{B}} - \xi\mathbf{I}_{P-1})^{-1}\mathbf{U} = \begin{bmatrix} \bar{m}_{1,d} & & \\ & \ddots & \\ & & \bar{m}_{K+1,d} \end{bmatrix} + o_{\mathbb{P}}(1). \quad (\text{VII.12})$$

Since  $|w_2| \leq \Im(\xi)$  is deterministically bounded, by dominated convergence theorem, we have the  $L_1$  convergence of  $w_2$  by plugging (VII.10)-(VII.12) into (VII.9). We have

$$\bar{m}_{1,d} = \psi_1 \{ -\xi + q_2 \mu_{1,*}^2 + q_4 \mu_{1,1}^2 - \phi_{1,*}^2 \bar{m}_{K+1,d} - \mathbf{l}_K^\top \mathbf{M}_K^{-1} \mathbf{l}_K \}^{-1} + o_d(1). \quad (\text{VII.13})$$

Here we define  $\mathbf{l}_K = [q_4 \mu_{1,1}^2 \quad \cdots \quad q_4 \mu_{1,1} \mu_{K,1} \quad \phi_{1,1}]^\top \in \mathbb{R}^{(K+1) \times 1}$ , and

$$\mathbf{M}_K = \begin{bmatrix} q_4 \mu_{1,1}^2 + \frac{1}{\bar{m}_{1,d}} & \cdots & q_4 \mu_{1,1} \mu_{K,1} & \phi_{1,1} \\ \vdots & \ddots & \vdots & \vdots \\ q_4 \mu_{1,1} \mu_{K,d} & \cdots & q_4 \mu_{K,1}^2 + \frac{1}{\bar{m}_{K,d}} & \phi_{K,1} \\ \phi_{1,1} & \cdots & \phi_{K,1} & q_5 + \frac{1}{\bar{m}_{K+1,d}} \end{bmatrix}.$$

Note that  $\phi_{j,1} = \mu_{j,1}(1 + q_1)$ ,  $\phi_{j,*} = \mu_{j,*}$ , we aim to prove the following equality:

$$q_4 \mu_{1,1}^2 - \mathbf{l}_K^\top \mathbf{M}_K^{-1} \mathbf{l}_K = \frac{\mu_{1,1}^2 q_4 (1 + q_5 \bar{m}_{K+1,d}) - \mu_{1,1}^2 (1 + q_1)^2 \bar{m}_{K+1,d}}{\left(1 + q_4 \sum_{j=1}^K \mu_{j,1}^2 \bar{m}_{j,d}\right) (1 + q_5 \bar{m}_{K+1,d}) - (1 + q_1)^2 \sum_{j=1}^K \mu_{j,1}^2 \bar{m}_{j,d} \bar{m}_{K+1,d}}. \quad (\text{VII.14})$$

We prove (VII.14) by mathematical induction. For  $K = 2$ , (VII.14) holds from Section IV.2. We assume that

$$q_4 \mu_{1,1}^2 - \mathbf{l}_{K-1}^\top \mathbf{M}_{K-1}^{-1} \mathbf{l}_{K-1} = \frac{\mu_{1,1}^2 q_4 (1 + q_5 \bar{m}_{K,d}) - \mu_{1,1}^2 (1 + q_1)^2 \bar{m}_{K,d}}{\left(1 + q_4 \sum_{j=1}^{K-1} \mu_{j,1}^2 \bar{m}_{j,d}\right) (1 + q_5 \bar{m}_{K,d}) - (1 + q_1)^2 \sum_{j=1}^{K-1} \mu_{j,1}^2 \bar{m}_{j,d} \bar{m}_{K,d}} \quad (\text{VII.15})$$

holds under the case  $K - 1$ . We aim to prove (VII.14) for general  $K$  under the assumption that (VII.15) holds. To prove so, define  $\boldsymbol{\mu}_K = [\mu_{1,1}, \dots, \mu_{K,1}]^\top$ . The vector  $\mathbf{l}_K$  could be separated into  $\mathbf{l}_K = [q_4 \mu_{1,1} \cdot \boldsymbol{\mu}_K^\top \quad (1 + q_1) \mu_{1,1}]^\top$ . If we further define

$$\mathbf{V}_0 = \begin{bmatrix} q_4 \mu_{1,1}^2 + \frac{1}{\bar{m}_{1,d}} & \cdots & q_4 \mu_{1,1} \mu_{K,1} \\ \vdots & \ddots & \vdots \\ q_4 \mu_{1,1} \mu_{K,1} & \cdots & q_4 \mu_{K,1}^2 + \frac{1}{\bar{m}_{K,d}} \end{bmatrix},$$

the target equation (VII.14) could be rewritten as

$$q_4 \mu_{1,1}^2 - \mathbf{l}_K^\top \mathbf{M}_K^{-1} \mathbf{l}_K = q_4 \mu_{1,1}^2 - \begin{bmatrix} q_4 \mu_{1,1} \cdot \boldsymbol{\mu}_K \\ (1 + q_1) \mu_{1,1} \end{bmatrix}^\top \begin{bmatrix} \mathbf{V}_0 & (1 + q_1) \boldsymbol{\mu}_K^\top \\ (1 + q_1) \boldsymbol{\mu}_K & q_5 + \frac{1}{\bar{m}_{K+1,d}} \end{bmatrix}^{-1} \begin{bmatrix} q_4 \mu_{1,1} \cdot \boldsymbol{\mu}_K \\ (1 + q_1) \mu_{1,1} \end{bmatrix}. \quad (\text{VII.16})$$

Clearly, the formula in (VII.16) requires us to investigate  $\boldsymbol{\mu}_K^\top \mathbf{V}_0^{-1} \boldsymbol{\mu}_K$  first. Under the case  $K - 1$ , (VII.15) holds from the induction hypothesis. Thus if we set  $(1 + q_1) = q_4 \mu_{K,1}$ ,  $q_5 = q_4 \mu_{K,1}^2$ , we have  $\mathbf{l}_{K-1} = q_4 \mu_{1,1} \cdot \boldsymbol{\mu}_K$ . Plugging  $\mathbf{l}_{K-1} = q_4 \mu_{1,1} \cdot \boldsymbol{\mu}_K$  into (VII.15) we obtain that



$$\begin{aligned}
 & \boldsymbol{\mu}_K^\top \mathbf{V}_0^{-1} \boldsymbol{\mu}_K \tag{VII.17} \\
 &= \frac{1}{q_4^2 \mu_{1,1}^2} \left( q_4 \mu_{1,1}^2 - \frac{\mu_{1,1}^2 q_4 (1 + \mu_{K,1}^2 q_4 \bar{m}_{K,d}) - \mu_{1,1}^2 q_4^2 \mu_{K,1}^2 \bar{m}_{K,d}}{\left(1 + q_4 \sum_{j=1}^{K-1} \mu_{j,1}^2 \bar{m}_{j,d}\right) (1 + \mu_{K,1}^2 q_4 \bar{m}_{K,d}) - q_4^2 \mu_{K,1}^2 \sum_{j=1}^{K-1} \mu_{j,1}^2 \bar{m}_{j,d} \bar{m}_{K,d}} \right) \\
 &= \frac{1}{q_4 \mu_{1,1}^2} \left( \mu_{1,1}^2 - \frac{\mu_{1,1}^2}{1 + q_4 \sum_{j=1}^K \mu_{j,1}^2 \bar{m}_{j,d}} \right) = \frac{\sum_{j=1}^K \mu_{j,1}^2 \bar{m}_{j,d}}{1 + q_4 \sum_{j=1}^K \mu_{j,1}^2 \bar{m}_{j,d}}.
 \end{aligned}$$

Therefore, for the case  $K$ , we have

$$\begin{aligned}
 q_4 \mu_{1,1}^2 - \mathbf{l}_K^\top \mathbf{M}_K^{-1} \mathbf{l}_K &= q_4 \mu_{1,1}^2 - \begin{bmatrix} q_4 \mu_{1,1} \cdot \boldsymbol{\mu}_K \\ (1 + q_1) \mu_{1,1} \end{bmatrix}^\top \begin{bmatrix} \mathbf{V}_0 & (1 + q_1) \boldsymbol{\mu}_K^\top \\ (1 + q_1) \boldsymbol{\mu}_K & q_5 + \frac{1}{\bar{m}_{K+1,d}} \end{bmatrix}^{-1} \begin{bmatrix} q_4 \mu_{1,1} \cdot \boldsymbol{\mu}_K \\ (1 + q_1) \mu_{1,1} \end{bmatrix} \\
 &= \frac{\mu_{1,1}^2 q_4 (1 + q_5 \bar{m}_{K,d}) - \mu_{1,1}^2 (1 + q_1)^2 \bar{m}_{K,d}}{\left(1 + q_4 \sum_{j=1}^{K-1} \mu_{j,1}^2 \bar{m}_{j,d}\right) (1 + q_5 \bar{m}_{K,d}) - (1 + q_1)^2 \sum_{j=1}^{K-1} \mu_{j,1}^2 \bar{m}_{j,d} \bar{m}_{K,d}}.
 \end{aligned}$$

Here, the first equality directly comes from (VII.16), and the second equality comes from Schur complement and (VII.17) after direct calculation. We completed the mathematical induction for the general case  $K$  and equation (VII.14) is proved. Then we have

$$\bar{m}_{1,d} = \psi_1 \left\{ -\xi + q_2 \mu_{1,*}^2 - \mu_{1,*}^2 \bar{m}_{K+1,d} + \frac{H_{1,d}}{H_{D,d}} \right\}^{-1} + o_d(1),$$

where

$$\begin{aligned}
 H_{1,d} &= \mu_{j,1}^2 q_4 (1 + q_5 \bar{m}_{K+1,d}) - \mu_{j,1}^2 (1 + q_1)^2 \bar{m}_{K+1,d}, \quad j = 1, \dots, K, \\
 H_{D,d} &= \left(1 + q_4 \sum_{j=1}^K \mu_{j,1}^2 \bar{m}_{j,d}\right) (1 + q_5 \bar{m}_{K+1,d}) - (1 + q_1)^2 \sum_{j=1}^K \mu_{j,1}^2 \bar{m}_{j,d} \bar{m}_{K+1,d}.
 \end{aligned}$$

After similar argument, we conclude that

$$\begin{aligned}
 \bar{m}_{j,d} &= \psi_j \left\{ -\xi + s_j \mu_{j,*}^2 - \mu_{j,*}^2 \bar{m}_{K+1,d} + \frac{H_{j,d}}{H_{D,d}} \right\}^{-1} + o_d(1), \quad j = 1, \dots, K+1, \\
 \bar{m}_{K+1,d} &= \psi_{K+1} \left\{ -\xi + q_3 - \sum_{j=1}^K \mu_{j,*}^2 \bar{m}_{j,d} + \frac{H_{K+1,d}}{H_{D,d}} \right\}^{-1} + o_d(1),
 \end{aligned}$$

where

$$\begin{aligned}
 H_{j,d} &= \mu_{j,1}^2 q_4 (1 + q_5 \bar{m}_{K+1,d}) - \mu_{j,1}^2 (1 + q_1)^2 \bar{m}_{K+1,d}, \quad j = 1, \dots, K, \\
 H_{K+1,d} &= q_5 \left(1 + q_4 \sum_{j=1}^K \mu_{j,1}^2 \bar{m}_{j,d}\right) - (1 + q_1)^2 \sum_{j=1}^K \mu_{j,1}^2 \bar{m}_{j,d}.
 \end{aligned}$$

We get that each element of  $\bar{\mathbf{m}}_d(\xi) - \mathbf{F}(\bar{\mathbf{m}}_d(\xi))$  is asymptotically zero. Therefore  $\|\bar{\mathbf{m}}_d(\xi) - \mathbf{F}(\bar{\mathbf{m}}_d(\xi))\|_2 = o_d(1)$ . The remaining arguments are similar to those in Section IV and the details are skipped. Wrapping all together, we complete the proof of Proposition C.7. ■

## References

- Xiuyuan Cheng and Amit Singer. The spectrum of random inner-product kernel matrices. *Random Matrices: Theory and Applications*, 2(04), 2013.
- Behrooz Ghorbani, Song Mei, Theodor Misiakiewicz, and Andrea Montanari. Linearized two-layers neural networks in high dimension. *The Annals of Statistics*, 49(2):1029–1054, 2021.
- Trevor Hastie, Andrea Montanari, Saharon Rosset, and Ryan J Tibshirani. Surprises in high-dimensional ridgeless least squares interpolation. *The Annals of Statistics*, 50(2):949–986, 2022.
- Luogeng Hua. *Harmonic Analysis of Functions of Several Complex Variables in the Classical Domains*. American Mathematical Soc., 1963.
- Song Mei and Andrea Montanari. The generalization error of random features regression: Precise asymptotics and the double descent curve. *Communications on Pure and Applied Mathematics*, 75(4):667–766, 2022.
- Terence Tao. *Topics in Random Matrix Theory*. American Mathematical Soc., 2012.
- Roman Vershynin. Introduction to the non-asymptotic analysis of random matrices. *arXiv preprint arXiv:1011.3027*, 2010.