

## A Technical background

### A.1 Hermite polynomials

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

$$\mathbb{E}\{\text{He}_j(G)\text{He}_k(G)\} = k! \delta_{jk}. \quad (18)$$

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

$$g(x) = \sum_{k=0}^{\infty} \frac{\mu_k(g)}{k!} \text{He}_k(x), \quad \mu_k(g) \equiv \mathbb{E}\{g(G)\text{He}_k(G)\}. \quad (19)$$

### A.2 Notations

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

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

and respectively:  $h_1(d) = o_{d,\mathbb{P}}(h_2(d))$ , if  $h_1(d)/h_2(d)$  converges to 0 in probability.

We will occasionally hide logarithmic factors using the  $\tilde{O}_d(\cdot)$  notation (resp.  $\tilde{o}_d(\cdot)$ ):  $h_1(d) = \tilde{O}_d(h_2(d))$  if there exists a constant  $C$  such that  $h_1(d) \leq C(\log d)^C h_2(d)$ . Similarly, we will denote  $\tilde{O}_{d,\mathbb{P}}(\cdot)$  (resp.  $\tilde{o}_{d,\mathbb{P}}(\cdot)$ ) when considering the big-O in probability notation up to a logarithmic factor.

## B Proofs for quadratic functions

Our results for quadratic functions (qf) assume  $\mathbf{x}_i \sim \mathcal{N}(0, \mathbf{I}_d)$  and  $y_i = f_*(\mathbf{x}_i)$  where

$$f_*(\mathbf{x}_i) \equiv b_0 + \langle \mathbf{x}, \mathbf{B}\mathbf{x} \rangle. \quad (20)$$

Throughout this section, we will denote  $\mathbb{E}_{\mathbf{x}}$  the expectation operator with respect to  $\mathbf{x} \sim \mathcal{N}(0, \mathbf{I}_d)$ , and  $\mathbb{E}_{\mathbf{w}}$  the expectation operator with respect to  $\mathbf{w} \sim \mathcal{N}(0, \mathbf{\Gamma})$ .

### B.1 Random Features model: proof of Theorem 1

Recall the definition

$$R_{\text{RF},N}(f_*) = \min_{\hat{f} \in \mathcal{F}_{\text{RF},N}(\mathbf{W})} \mathbb{E}\{(f_*(\mathbf{x}) - \hat{f}(\mathbf{x}))^2\},$$

where

$$\mathcal{F}_{\text{RF},N}(\mathbf{W}) = \left\{ f_N(\mathbf{x}) = \sum_{i=1}^N a_i \sigma(\langle \mathbf{w}_i, \mathbf{x} \rangle) : a_i \in \mathbb{R}, i \in [N] \right\}.$$

Note that it is easy to see from the proof that the result stays the same if we add an offset  $c$ .

#### B.1.1 Representation of the RF risk

**Lemma 1.** *Consider the RF model. We have*

$$R_{\text{RF},N}(f_*) = \mathbb{E}_{\mathbf{x}}[f_*(\mathbf{x})^2] - \mathbf{V}^\top \mathbf{U}^{-1} \mathbf{V}, \quad (21)$$

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

$$\begin{aligned} V_i &= \mathbb{E}_{\mathbf{x}}[f_*(\mathbf{x})\sigma(\langle \mathbf{w}_i, \mathbf{x} \rangle)], \\ U_{ij} &= \mathbb{E}_{\mathbf{x}}[\sigma(\langle \mathbf{w}_i, \mathbf{x} \rangle)\sigma(\langle \mathbf{w}_j, \mathbf{x} \rangle)]. \end{aligned}$$

*Proof of Lemma 1.* Simply write the KKT conditions. The optimum is achieved at  $\mathbf{a} = \mathbf{U}^{-1}\mathbf{V}$ .  $\square$

### B.1.2 Approximation of kernel matrix $\mathbf{U}$

**Lemma 2.** Let  $\sigma \in L^2(\mathbb{R}, \gamma)$  be an activation function. Denote  $\lambda_k = \mathbb{E}_{G \sim \mathcal{N}(0,1)}[\sigma(G)\text{He}_k(G)]$  the  $k$ -th Hermite coefficient of  $\sigma$  and assume  $\lambda_0 = 0$ . Let  $\mathbf{U} = (U_{ij})_{i,j \in [N]}$  be a random matrix with

$$U_{ij} = \mathbb{E}_{\mathbf{x}}[\sigma(\langle \mathbf{w}_i, \mathbf{x} \rangle) \sigma(\langle \mathbf{w}_j, \mathbf{x} \rangle)],$$

where  $(\mathbf{w}_i)_{i \in [N]} \sim \mathcal{N}(\mathbf{0}, \mathbf{\Gamma})$  independently. Assume conditions **A1** and **A2** hold.

Let  $\mathbf{W} = (\mathbf{w}_1, \dots, \mathbf{w}_N) \in \mathbb{R}^{d \times N}$ , and denote  $\mathbf{U}_0 = \{(U_0)_{ij}\}_{i,j \in [N]}$ , with

$$(U_0)_{ij} = \tilde{\lambda} \delta_{ij} + \lambda_1^2 \langle \mathbf{w}_i, \mathbf{w}_j \rangle + \kappa/d + \mu_i \mu_j,$$

where

$$\begin{aligned} \mu_i &= \lambda_2(\|\mathbf{w}_i\|_2^2 - 1)/2, \\ \tilde{\lambda} &= \mathbb{E}[\sigma(G)^2] - \lambda_1^2, \\ \kappa &= d\lambda_2^2 \text{Tr}(\mathbf{\Gamma}^2)/2. \end{aligned}$$

Then we have as  $N/d = \rho$  and  $d \rightarrow \infty$ ,

$$\|\mathbf{U} - \mathbf{U}_0\|_{\text{op}} = o_d(\mathbb{P}(1)).$$

*Proof of Lemma 2.*

**Step 1. Hermite expansion of  $\sigma$  for  $\|\mathbf{w}_i\|_2 \neq 1$ .** Denote  $\sigma_i(x) = \sigma(\|\mathbf{w}_i\|_2 \cdot x)$ . First notice that by a change of variables, we get

$$\mathbb{E}[\sigma(tG)] = \mathbb{E}[(\sigma(G)/t) \exp(G^2(1 - 1/t^2)/2)]. \quad (22)$$

By Assumption **A1**, there exists  $c_1 < 1$  such that

$$\sigma(u)^2 \exp(u^2(1 - 1/t^2)) \leq c_0 \exp(u^2(c_1/2 + 1 - 1/t^2)).$$

Hence for  $|t - 1|$  sufficiently small, we have  $\sigma_i \in L^2(\mathbb{R}, \gamma)$  and we can consider its Hermite expansion

$$\sigma_i(x) = \sum_{k=0}^{\infty} \frac{\zeta_k(\sigma_i)}{k!} \text{He}_k(x),$$

where

$$\zeta_k(\sigma_i) = \mathbb{E}_{G \sim \mathcal{N}(0,1)}[\sigma(\|\mathbf{w}_i\|_2 G) \text{He}_k(G)].$$

Denote the Hermite expansion of  $\sigma$  to be

$$\sigma(x) = \sum_{k=0}^{\infty} \lambda_k(\sigma) \text{He}_k(x)/k!,$$

where

$$\lambda_k(\sigma) = \mathbb{E}_{G \sim \mathcal{N}(0,1)}[\sigma(G) \text{He}_k(G)].$$

By dominated convergence theorem, we have

$$\lim_{t \rightarrow 1} \mathbb{E}_{G \sim \mathcal{N}(0,1)}[(\sigma(G) - \sigma(tG))^2] = 0.$$

In addition, by sub-Gaussianity of the norm of a multivariate Gaussian random variable (see [Ver10]), it is easy to show that

$$\sup_{i \in [N]} |\|\mathbf{w}_i\|_2 - 1| = o_d(\mathbb{P}(1)). \quad (23)$$

Hence we have

$$\begin{aligned} \sup_{i \in [N]} \|\sigma - \sigma_i\|_{L^2} &= o_{d, \mathbb{P}}(1), \\ \sup_{i \in [N]} |\zeta_k(\sigma_i) - \lambda_k(\sigma)| &\leq \sup_{i \in [N]} \|\sigma - \sigma_i\|_{L^2} \mathbb{E}[\text{He}_k(G)^2]^{1/2} = o_{d, \mathbb{P}}(1), \end{aligned} \quad (24)$$

for any fixed integer  $k$ .

**Step 2. Expansion of  $U$ .** Denote  $\mathbf{u}_i = \mathbf{w}_i / \|\mathbf{w}_i\|_2$ , then we have

$$U_{ij} = \underbrace{\zeta_0(\sigma_i)\zeta_0(\sigma_j)}_{T_{0,ij}} + \underbrace{\zeta_1(\sigma_i)\zeta_1(\sigma_j)\langle \mathbf{u}_i, \mathbf{u}_j \rangle}_{T_{1,ij}} + \underbrace{\zeta_2(\sigma_i)\zeta_2(\sigma_j)\frac{\langle \mathbf{u}_i, \mathbf{u}_j \rangle^2}{2}}_{T_{2,ij}} + \underbrace{\sum_{k \geq 3} \zeta_k(\sigma_i)\zeta_k(\sigma_j)\frac{\langle \mathbf{u}_i, \mathbf{u}_j \rangle^k}{k!}}_{T_{3,ij}}. \quad (25)$$

We define

$$\mathbf{T}_k = \left( \zeta_k(\sigma_i)\zeta_k(\sigma_j)\frac{\langle \mathbf{w}_i, \mathbf{w}_j \rangle^k}{k!} \right)_{i,j \in [N]}.$$

**Step 3. Term  $\mathbf{T}_0$ .** By definition of  $\mu_i$ , we have

$$\mathbf{T}_0 = (\zeta_0(\sigma_i)\zeta_0(\sigma_j))_{i,j \in [N]} = \mathbf{D}_0[(\lambda_2/2)(\|\mathbf{w}_i\|_2^2 - 1)(\|\mathbf{w}_j\|_2^2 - 1)]_{i,j \in [N]} \mathbf{D}_0,$$

where (by the assumption that  $\mathbb{E}_G[\sigma(G)] = 0$ )

$$(\mathbf{D}_0)_{ii} = \frac{\zeta_0(\sigma_i)}{\lambda_2(\|\mathbf{w}_i\|_2^2 - 1)/2} = \mathbb{E}\left[\frac{\sigma(\|\mathbf{w}_i\|_2 G) - \sigma(G)}{\|\mathbf{w}_i\|_2 - 1}\right] \cdot \frac{1}{\lambda_2(\|\mathbf{w}_i\|_2 + 1)/2}.$$

Let us show:

$$\lim_{t \rightarrow 1} \mathbb{E}\left[\frac{\sigma(tG) - \sigma(G)}{t - 1}\right] = \lambda_2(\sigma), \quad (26)$$

or equivalently:

$$\lim_{t \rightarrow 1} \mathbb{E}\left[\frac{\sigma(tG) - \sigma(G)}{t - 1} - (G^2 - 1)\sigma(G)\right] = 0$$

Recall the change of variable (22) and do a first order Taylor expansion of the exponential: there exists a function  $\xi(G) \in [0, G]$  such that

$$\begin{aligned} &\mathbb{E}\left[\frac{\sigma(tG) - \sigma(G)}{t - 1} - (G^2 - 1)\sigma(G)\right] \\ &= \mathbb{E}\left[\sigma(G)\left(\exp(G^2(1 - 1/t^2)/2) - t - t(t - 1)(G^2 - 1)\right)\right] \cdot \frac{1}{t(t - 1)} \\ &= \mathbb{E}\left[\sigma(G)(t - 1)\left(1 - G^2[2t + 1]/(2t^2) + G^4(t + 1)^2/(8t^4)\exp(\xi(G)^2(1 - 1/t^2)/2)\right)\right] \cdot \frac{1}{t}. \end{aligned}$$

We see that the integrand goes to zero as  $t \rightarrow 1$ . For  $|t - 1|$  sufficiently small, we have

$$\frac{\left|\exp(G^2(1 - 1/t^2)/2) - t - t(t - 1)(G^2 - 1)\right|}{|t - 1|} \leq 2 + 2G^2 + 2G^4 \exp(G^2/5),$$

which is squared integrable. Recalling that  $\sigma \in L^2(\mathbb{R}, \gamma)$ , we obtain (26) by dominated convergence.

Hence, combining (23) and (26) gives

$$\|\mathbf{D}_0 - \mathbf{I}_d\|_{\text{op}} = o_{d, \mathbb{P}}(1).$$

Furthermore, for  $\boldsymbol{\mu} = (\mu_i)_{i \in [N]}$  with  $\mu_i = \lambda_2(\|\mathbf{w}_i\|_2^2 - 1)/2$ , we have

$$\mathbb{E}[\|\boldsymbol{\mu}\boldsymbol{\mu}^\top\|_{\text{op}}] = \mathbb{E}[\|\boldsymbol{\mu}\|_2^2] = \frac{\lambda_2^2}{4} N \mathbb{E}[(\|\mathbf{w}_i\|_2^2 - 1)^2] = \frac{\lambda_2^2}{2} N \|\boldsymbol{\Gamma}\|_F^2 \leq \frac{\lambda_2^2}{2} N^2 \|\boldsymbol{\Gamma}\|_{\text{op}}^2 = O_{d, \mathbb{P}}(1),$$

where the last equality comes from assumption **A2**. We get

$$\|\mathbf{T}_0 - \boldsymbol{\mu}\boldsymbol{\mu}^\top\|_{\text{op}} \leq 2\|\mathbf{D}_0 - \mathbf{I}_d\|_{\text{op}}\|\boldsymbol{\mu}\boldsymbol{\mu}^\top\|_{\text{op}}(\|\mathbf{D}_0\|_{\text{op}} + 1) = o_{d,\mathbb{P}}(1). \quad (27)$$

**Step 4. Term  $\mathbf{T}_1$ .** For  $\mathbf{T}_1$ , we have

$$\mathbf{T}_1 = (\zeta_1(\sigma_i)\zeta_1(\sigma_j)\langle \mathbf{u}_i, \mathbf{u}_j \rangle)_{i,j \in [N]} = \mathbf{D}_1 \mathbf{W}^\top \mathbf{W} \mathbf{D}_1,$$

where

$$\mathbf{D}_1 = \text{diag}((\zeta_1(\sigma_i))/\|\mathbf{w}_i\|_2).$$

By the uniform convergence of  $\zeta_1(\sigma_i)$  to  $\lambda_1(\sigma)$ , cf Eq. (24), we have

$$\|\mathbf{D}_1 - \lambda_1(\sigma)\mathbf{I}_d\|_{\text{op}} = o_{d,\mathbb{P}}(1).$$

Moreover, we have

$$\|\mathbf{W}^\top \mathbf{W}\|_{\text{op}} = \|\mathbf{W} \mathbf{W}^\top\|_{\text{op}} \leq \|\sqrt{d}\boldsymbol{\Gamma}^{1/2}\|_{\text{op}}^2 \|\mathbf{G} \mathbf{G}^\top\|_{\text{op}} = O_{d,\mathbb{P}}(1),$$

where we denoted by  $\mathbf{G}$  the matrix with columns  $\mathbf{g}_i \sim \mathbf{N}(\mathbf{0}, \mathbf{I}_d/d)$ . Hence, we have

$$\|\mathbf{T}_1 - \lambda_1^2 \mathbf{W}^\top \mathbf{W}\|_{\text{op}} \leq \|\mathbf{D}_1 - \lambda_1 \mathbf{I}_d\|_{\text{op}} \|\mathbf{W}^\top \mathbf{W}\|_{\text{op}} (\|\mathbf{D}_1\|_{\text{op}} + 1) = o_{d,\mathbb{P}}(1). \quad (28)$$

**Step 5. Term  $\mathbf{T}_2$ .** We have

$$\mathbf{T}_2 = (\zeta_2(\sigma_i)\zeta_2(\sigma_j)\langle \mathbf{u}_i, \mathbf{u}_j \rangle^2/2)_{i,j \in [N]} = \mathbf{D}_2 (\langle \mathbf{w}_i, \mathbf{w}_j \rangle^2/2)_{i,j \in [N]} \mathbf{D}_2,$$

where

$$\mathbf{D}_2 = \text{diag}((\zeta_2(\sigma_i))/\|\mathbf{w}_i\|_2^2).$$

By the uniform convergence of  $\zeta_2(\sigma_i)$  to  $\lambda_2(\sigma)$ , we have

$$\|\mathbf{D}_2 - \lambda_2 \mathbf{I}_d\|_{\text{op}} = o_{d,\mathbb{P}}(1).$$

Moreover, we have (see below)

$$\|(\langle \mathbf{w}_i, \mathbf{w}_j \rangle^2)_{i,j \in [N]}\|_{\text{op}} = O_{d,\mathbb{P}}(1).$$

Hence, we have

$$\|\mathbf{T}_2 - \lambda_2^2 (\langle \mathbf{w}_i, \mathbf{w}_j \rangle^2/2)_{i,j \in [N]}\|_{\text{op}} \leq \|\mathbf{D}_2 - \lambda_2 \mathbf{I}_d\|_{\text{op}} \|(\langle \mathbf{w}_i, \mathbf{w}_j \rangle^2/2)_{i,j \in [N]}\|_{\text{op}} (\|\mathbf{D}_2\|_{\text{op}} + 1) = o_{d,\mathbb{P}}(1).$$

Moreover, by the estimates in proof of Theorem 2.1 in [EK<sup>+</sup>10], we have

$$\|(\langle \mathbf{w}_i, \mathbf{w}_j \rangle^2/2)_{i,j \in [N]} - [\text{Tr}(\boldsymbol{\Gamma}^2)/2] \mathbf{1}\mathbf{1}^\top - (1/2)\mathbf{I}_N\|_{\text{op}} = o_{d,\mathbb{P}}(1).$$

Hence, we get

$$\|\mathbf{T}_2 - \lambda_2^2 [\text{Tr}(\boldsymbol{\Gamma}^2)/2] \mathbf{1}\mathbf{1}^\top - [\lambda_2^2/2] \mathbf{I}_N\|_{\text{op}} = o_{d,\mathbb{P}}(1). \quad (29)$$

**Step 6. Term  $\sum_{k \geq 3} \text{ddiag}(\mathbf{T}_k)$ .** Denote  $\text{ddiag}(\mathbf{T}_k)$  the diagonal matrix composed of diagonal entries of  $\mathbf{T}_k$ . We have

$$\begin{aligned} \left| \sum_{k \geq 3} ((\mathbf{T}_k)_{ii} - \lambda_k(\sigma)^2/k!) \right| &= \left| \|\sigma_i\|_{L^2}^2 - \sum_{k=0}^2 \zeta_k(\sigma_i)^2/k! - \|\sigma\|_{L^2}^2 + \sum_{k=0}^2 \lambda_k(\sigma)^2/k! \right| \\ &\leq \|\sigma - \sigma_i\|_{L^2} [2\|\sigma\|_{L^2} + \|\sigma - \sigma_i\|_{L^2}] + \sum_{k=0}^2 |\zeta_k(\sigma_i)^2 - \lambda_k(\sigma)^2|/k!. \end{aligned}$$

Note that we have shown (cf Eq. (24))

$$\sup_{i \in [N]} \max \left\{ \|\sigma - \sigma_i\|_{L^2}, \max_{k=0,1,2} |\zeta_k(\sigma_i) - \lambda_k(\sigma)| \right\} = o_{d,\mathbb{P}}(1).$$

Therefore, we have

$$\left\| \sum_{k \geq 3} \text{ddiag}(\mathbf{T}_k) - (\tilde{\lambda} - \lambda_2^2/2) \mathbf{I}_N \right\|_{\text{op}} = o_{d, \mathbb{P}}(1). \quad (30)$$

**Step 7. Term  $\sum_{k \geq 3} [\mathbf{T}_k - \text{ddiag}(\mathbf{T}_k)]$ .** We have

$$\begin{aligned} & \left\| \sum_{k \geq 3} [\mathbf{T}_k - \text{ddiag}(\mathbf{T}_k)] \right\|_F \leq \sum_{k \geq 3} \|\mathbf{T}_k - \text{ddiag}(\mathbf{T}_k)\|_F \\ & \leq \sum_{k \geq 3} \left[ \left( \sum_{i,j=1}^N \zeta_k(\sigma_i)^2 \zeta_k(\sigma_j)^2 \right) \left( \sup_{i \neq j} \langle \mathbf{u}_i, \mathbf{u}_j \rangle^{2k} / (k!)^2 \right) \right]^{1/2} \\ & \leq \left[ \sum_{k \geq 3} \sum_{i=1}^N \zeta_k(\sigma_i)^2 / k! \right] \max_{i \neq j} \langle \mathbf{u}_i, \mathbf{u}_j \rangle^3 \\ & \leq \|\sigma_i\|_{L^2}^2 \times N \max_{i \neq j} \langle \mathbf{u}_i, \mathbf{u}_j \rangle^3. \end{aligned}$$

Note we have  $\max_{i \in [N]} \|\sigma_i\|_{L^2}^2 = O_{d, \mathbb{P}}(1)$ . Moreover, we have (see for example Lemma 10 in [GMMM19])

$$\max_{i \neq j} \langle \mathbf{u}_i, \mathbf{u}_j \rangle^3 = \tilde{O}_{d, \mathbb{P}}(d^{-3/2}).$$

Therefore, we have

$$\left\| \sum_{k \geq 3} [\mathbf{T}_k - \text{ddiag}(\mathbf{T}_k)] \right\|_F = o_{d, \mathbb{P}}(1). \quad (31)$$

Combining the bounds (27), (28), (29), (30) and (31) into the decomposition (25) proves the lemma.  $\square$

### B.1.3 Approximation of the $\mathbf{V}$ vector

**Lemma 3.** *Under the assumptions of Theorem 1, define  $\mathbf{V} = (V_1, \dots, V_N)^\top$  with*

$$V_i = \mathbb{E}_{\mathbf{x}}[f_*(\mathbf{x})\sigma(\langle \mathbf{w}_i, \mathbf{x} \rangle)]$$

*where  $(\mathbf{w}_i)_{i \in [N]} \sim \mathcal{N}(\mathbf{0}, \mathbf{\Gamma})$  independently. Then as  $N/d = \rho$  with  $d \rightarrow \infty$ , we have*

$$\|\mathbf{V} - \tau \mathbf{1} / \sqrt{d}\|_2^2 = \|\mathbf{B}\|_F^2 \cdot o_{d, \mathbb{P}}(1),$$

*where*

$$\tau = \sqrt{d} \cdot \lambda_2 \text{Tr}(\mathbf{B}\mathbf{\Gamma}).$$

*Proof of Lemma 3.* Without loss of generality, we assume  $\|\mathbf{B}\|_F = 1$  in the proof (it suffices to divide  $V_i$  by  $\|\mathbf{B}\|_F$ ). Consider  $\mathbf{w}_i \in \mathbb{R}^d$ . Take  $\mathbf{R}$  to be an orthogonal matrix such that  $\mathbf{R}\mathbf{w}_i = \|\mathbf{w}_i\|_2 \mathbf{e}_1$ , then we have

$$\begin{aligned} V_i &= \mathbb{E}_{\mathbf{x}}[f_*(\mathbf{R}^\top \mathbf{x})\sigma(\|\mathbf{w}_i\|_2 x_1)] \\ &= \mathbb{E}_{\mathbf{x}}[(\langle \mathbf{x}, \mathbf{RBR}^\top \mathbf{x} \rangle - \text{Tr}(\mathbf{B}))\sigma(\|\mathbf{w}_i\|_2 x_1)] \\ &= \mathbb{E}_{x_1} \left[ \left( x_1^2 \frac{\langle \mathbf{w}_i, \mathbf{B}\mathbf{w}_i \rangle}{\|\mathbf{w}_i\|_2^2} + \text{Tr}(\mathbf{P}_{\perp \mathbf{w}_i} \mathbf{B}) - \text{Tr}(\mathbf{B}) \right) \sigma(\|\mathbf{w}_i\|_2 x_1) \right] \\ &= \mathbb{E}_{x_1} \left[ (x_1^2 - 1) \frac{\langle \mathbf{w}_i, \mathbf{B}\mathbf{w}_i \rangle}{\|\mathbf{w}_i\|_2^2} \sigma(\|\mathbf{w}_i\|_2 x_1) \right] \\ &\equiv \frac{\langle \mathbf{w}_i, \mathbf{B}\mathbf{w}_i \rangle}{\|\mathbf{w}_i\|_2^2} \zeta_2(\sigma_i), \end{aligned}$$

where  $\mathbf{P}_{\perp \mathbf{w}_i}$  is the projection on the hyperplane orthogonal to  $\mathbf{w}_i$ , and we recall the definition of  $\zeta_2(\sigma_i)$  of Lemma 2:

$$\zeta_2(\sigma_i) = \mathbb{E}_G[(G^2 - 1)\sigma(\|\mathbf{w}_i\|_2 G)],$$

with  $G$  a standard normal random variable.

We define the following interpolating variables:

$$V_i^{(1)} = \frac{\langle \mathbf{w}_i, \mathbf{B}\mathbf{w}_i \rangle}{\|\mathbf{w}_i\|_2^2} \lambda_2, \quad V_i^{(2)} = \langle \mathbf{w}_i, \mathbf{B}\mathbf{w}_i \rangle \lambda_2, \quad V_i^{(3)} = \text{Tr}(\mathbf{\Gamma}\mathbf{B}) \lambda_2,$$

and the associated vectors  $\mathbf{V}^{(1)}$ ,  $\mathbf{V}^{(2)}$  and  $\mathbf{V}^{(3)}$ . We bound successively the distance between these vectors. We will denote by  $\mathbf{P}_{\mathbf{w}_i}$  the projection onto vector  $\mathbf{w}_i$ . First, we consider:

$$\|\mathbf{V} - \mathbf{V}^{(1)}\|_2^2 = \sum_{i=1}^N \text{Tr}(\mathbf{P}_{\mathbf{w}_i} \mathbf{B})^2 (\zeta_2(\sigma_i) - \lambda_2)^2.$$

One can check, using a similar argument as for Eq. (26) and dominated convergence, that

$$\lim_{t \rightarrow 1} \frac{\mathbb{E}[(G^2 - 1)(\sigma(tG) - \sigma(G))]}{t - 1} = \lambda_4(\sigma) + 2\lambda_2(\sigma). \quad (32)$$

Hence, recalling (23), we have

$$\|\mathbf{V} - \mathbf{V}^{(1)}\|_2^2 = O_{d,\mathbb{P}} \left( \left( \sup_{i \in [N]} \text{Tr}(\mathbf{P}_{\mathbf{w}_i} \mathbf{B}) \right)^2 \sum_{i=1}^N (\|\mathbf{w}_i\|_2 - 1)^2 \right). \quad (33)$$

Let us first show that the sum is bounded with high probability: denoting  $\mathbf{g} \sim \mathbf{N}(\mathbf{0}, \mathbf{I}_d)$ , classical sub-Gaussian concentration inequalities (see for example Theorem 6.3.2 in [Ver10]) shows that

$$\left\| \|\mathbf{\Gamma}^{1/2} \mathbf{g}\|_2 - \|\mathbf{\Gamma}^{1/2}\|_F \right\|_{\psi_2} \leq C \|\mathbf{\Gamma}^{1/2}\|_{\text{op}}, \quad (34)$$

where  $\|\cdot\|_{\psi_2}$  denotes the sub-Gaussian Orlicz norm. By assumption, we have  $\|\mathbf{\Gamma}^{1/2}\|_{\text{op}} = \|\mathbf{\Gamma}\|_{\text{op}}^{1/2} = O_d(d^{-1/2})$ , and  $\|\mathbf{\Gamma}^{1/2}\|_F = \sqrt{\text{Tr} \mathbf{\Gamma}} = 1$ . Hence, for  $\mathbf{w}_i \sim \mathbf{N}(\mathbf{0}, \mathbf{\Gamma})$ , we have

$$\left\| \sqrt{d} \|\mathbf{w}_i\|_2 - \sqrt{d} \right\|_{\psi_2} \leq C. \quad (35)$$

Therefore, we have

$$\sum_{i=1}^N (\|\mathbf{w}_i\|_2 - 1)^2 = O_{d,\mathbb{P}}(1). \quad (36)$$

Furthermore, we readily have (for example from (23))

$$\sup_{i \in [N]} \|\mathbf{w}_i\|^{-4} = O_{d,\mathbb{P}}(1). \quad (37)$$

Noticing that  $\text{Tr}(\mathbf{w}_i \mathbf{w}_i^\top \mathbf{B}) = \|\mathbf{B}^{1/2} \mathbf{w}_i\|_2^2$  and by the same argument as for (34), we have:

$$\left\| \|\mathbf{B}^{1/2} \mathbf{\Gamma}^{1/2} \mathbf{g}\|_2 - \mathbb{E}[\|\mathbf{B}^{1/2} \mathbf{\Gamma}^{1/2} \mathbf{g}\|_2] \right\|_{\psi_2} \leq C \|\mathbf{B}^{1/2} \mathbf{\Gamma}^{1/2}\|_{\text{op}}. \quad (38)$$

By assumption **A2**, we have  $\|\mathbf{B}^{1/2} \mathbf{\Gamma}^{1/2}\|_{\text{op}} \leq \|\mathbf{B}^{1/2}\|_{\text{op}} \|\mathbf{\Gamma}^{1/2}\|_{\text{op}} = O_d(d^{-1/2})$  and

$$\mathbb{E}[\|\mathbf{B}^{1/2} \mathbf{\Gamma}^{1/2} \mathbf{g}\|_2] \leq (\mathbb{E}[\|\mathbf{B}^{1/2} \mathbf{\Gamma}^{1/2} \mathbf{g}\|_2^2])^{1/2} = \text{Tr}(\mathbf{\Gamma} \mathbf{B})^{1/2} \leq \|\mathbf{\Gamma}\|_F^{1/2} \|\mathbf{B}\|_F^{1/2} \leq \|\mathbf{\Gamma}\|_{\text{op}}^{1/4} \text{Tr}(\mathbf{\Gamma})^{1/4} = O_d(d^{-1/2}),$$

which combined with (38) yields

$$\sup_{i \in [N]} \|\mathbf{B}^{1/2} \mathbf{w}_i\|_2^2 = o_{d,\mathbb{P}}(1). \quad (39)$$

Combining the bounds (36), (37) and (39) into (33), we get

$$\|\mathbf{V} - \mathbf{V}^{(1)}\|_2^2 = o_{d,\mathbb{P}}(1). \quad (40)$$

Consider now

$$\begin{aligned} \|\mathbf{V}^{(1)} - \mathbf{V}^{(2)}\|_2^2 &= \sum_{i=1}^N \lambda_2^2 \langle \mathbf{w}_i, \mathbf{B}\mathbf{w}_i \rangle^2 \left( \frac{1}{\|\mathbf{w}_i\|_2^2} - 1 \right)^2 \\ &\leq \lambda_2^2 \left( \sup_{i \in [N]} \|\mathbf{B}^{1/2} \mathbf{w}_i\|_2^2 / \|\mathbf{w}_i\|_2^2 \right) \sum_{i=1}^N (\|\mathbf{w}_i\|_2^2 - 1)^2. \end{aligned} \quad (41)$$

We have

$$\mathbb{E}_{\mathbf{w}_i \sim \mathcal{N}(\mathbf{0}, \mathbf{\Gamma})} [(\|\mathbf{w}_i\|_2^2 - 1)^2] = \mathbb{E}_{\mathbf{g} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} [(\langle \mathbf{g}\mathbf{g}^\top, \mathbf{\Gamma} \rangle - \text{Tr}(\mathbf{\Gamma}))^2] = 2\|\mathbf{\Gamma}\|_F^2 = O_{d,\mathbb{P}}(d^{-1}).$$

Hence we must have

$$\sum_{i=1}^N (\|\mathbf{w}_i\|_2^2 - 1)^2 = O_{d,\mathbb{P}}(1),$$

which, combined with (39) and (41), yields

$$\|\mathbf{V}^{(1)} - \mathbf{V}^{(2)}\|_2^2 = o_{d,\mathbb{P}}(1). \quad (42)$$

Consider the last comparison:

$$\|\mathbf{V}^{(2)} - \mathbf{V}^{(3)}\|_2^2 = \sum_{i=1}^N \lambda_2^2 \left( \langle \mathbf{w}_i, \mathbf{B}\mathbf{w}_i \rangle - \text{Tr}(\mathbf{\Gamma}\mathbf{B}) \right)^2.$$

Taking the expectation:

$$\begin{aligned} \mathbb{E}_{\mathbf{w}_i \sim \mathcal{N}(\mathbf{0}, \mathbf{\Gamma})} [(\langle \mathbf{w}_i, \mathbf{B}\mathbf{w}_i \rangle - \text{Tr}(\mathbf{\Gamma}\mathbf{B}))^2] &= \mathbb{E}_{\mathbf{g} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} [(\langle \mathbf{g}\mathbf{g}^\top, \mathbf{\Gamma}^{1/2} \mathbf{B} \mathbf{\Gamma}^{1/2} \rangle - \text{Tr}(\mathbf{\Gamma}\mathbf{B}))^2] \\ &= 2\|\mathbf{\Gamma}^{1/2} \mathbf{B} \mathbf{\Gamma}^{1/2}\|_F^2 \\ &\leq 2\|\mathbf{\Gamma}\|_{\text{op}}^2 \|\mathbf{B}\|_F^2 = O_d(d^{-2}). \end{aligned}$$

We conclude that

$$\sum_{i=1}^N \left( \langle \mathbf{w}_i, \mathbf{B}\mathbf{w}_i \rangle - \text{Tr}(\mathbf{\Gamma}\mathbf{B}) \right)^2 = o_{d,\mathbb{P}}(1),$$

and therefore

$$\|\mathbf{V}^{(2)} - \mathbf{V}^{(3)}\|_2^2 = o_{d,\mathbb{P}}(1), \quad (43)$$

where  $\mathbf{V}^{(3)} = \lambda_2 \text{Tr}(\mathbf{\Gamma}\mathbf{B}) \mathbf{1}$ . Combining the above three bounds (33), (42) and (43) yields the desired result.  $\square$

#### B.1.4 Calculating $\mathbf{1}^\top \mathbf{U}_0^{-1} \mathbf{1} / d$

The following proposition is stated in slightly more general terms, in order to be used in both the proofs of Theorem 1 and Theorem 4.

**Proposition 2.** *Let  $(\mathbf{w}_i)_{i \in [N]} \sim \mathcal{N}(\mathbf{0}, \mathbf{\Gamma})$  independently, where  $\mathbf{\Gamma}$  satisfies assumption **A2** (resp. **B2**). Denote by  $\lambda_k = \mathbb{E}_{G \sim \mathcal{N}(0,1)} [\sigma(G) \text{He}_k(G)]$  the  $k$ -th Hermite coefficient of  $\sigma$ . Define  $\tilde{\lambda} = \mathbb{E}_{G \sim \mathcal{N}(0,1)} [\sigma(G)^2] - \lambda_1^2$ . Consider  $\kappa \equiv \kappa(d)$  positive constants that are uniformly upper bounded. Define*

$$\mathbf{U}_0 = \mathbf{A}_0 + \kappa \mathbf{1} \mathbf{1}^\top / d + \boldsymbol{\mu} \boldsymbol{\mu}^\top,$$

where

$$\begin{aligned}\mathbf{A}_0 &= \tilde{\lambda} \mathbf{I}_N + \lambda_1^2 \mathbf{W}^\top \mathbf{W}, \\ \mu_i &= \lambda_2 (\|\mathbf{w}_i\|_2^2 - 1)/2.\end{aligned}$$

Then we have

$$\langle \mathbf{1}, \mathbf{U}_0^{-1} \mathbf{1} \rangle / d = \psi / (1 + \kappa \psi) + o_{d, \mathbb{P}}(1),$$

where  $\psi > 0$  is the unique solution of

$$-\tilde{\lambda} = -\frac{\rho}{\psi} + \int \frac{\lambda_1^2 t}{1 + \lambda_1^2 t \psi} \mathcal{D}(dt), \quad (44)$$

where  $\mathcal{D}$  is the empirical distribution of eigenvalues of  $d \cdot \mathbf{\Gamma}$ .

The proof of Proposition 2 is a direct combination of Lemma 4, 5, and 6 below.

**Lemma 4.** *Let  $(\mathbf{w}_i)_{i \in [N]} \sim \mathcal{N}(\mathbf{0}, \mathbf{\Gamma})$  independently. Assume condition **A2** holds (resp. **B2**). Let  $\boldsymbol{\mu} = (\|\mathbf{w}_i\|_2^2 - 1)_{i \in [N]}$ , and  $\mathbf{A}_0 = c_1 \mathbf{I}_N + c_2 \mathbf{W}^\top \mathbf{W}$ , where  $c_1 \equiv c_1(d)$  and  $c_2 \equiv c_2(d)$  are constants that are asymptotically upper and lower bounded by strictly positive constants. Then as  $d \rightarrow \infty$  and  $N/d \rightarrow \rho$ , we have*

$$\langle \mathbf{1}, \mathbf{A}_0^{-1} \boldsymbol{\mu} \rangle / \sqrt{d} = o_{d, \mathbb{P}}(1). \quad (45)$$

*Proof.* We first prove the lemma under the following extra assumption on the covariance matrix: there exists a (fixed) integer  $K$  such that

$$\mathbf{\Gamma} = \mathbf{Q} \text{diag}(\gamma_1 \mathbf{I}_{d_1}, \dots, \gamma_K \mathbf{I}_{d_K}) \mathbf{Q}^\top, \quad (46)$$

for some orthogonal matrix  $\mathbf{Q}$  and  $d \cdot \gamma_i \leq C$ . Furthermore, there exists an  $\varepsilon > 0$  such that  $d_k/d \geq \varepsilon$  for  $d$  sufficiently large.

Without loss of generality, we assume  $\mathbf{\Gamma} = \text{diag}(\gamma_1 \mathbf{I}_{d_1}, \dots, \gamma_K \mathbf{I}_{d_K})$ , and we divide  $\mathbf{w}_i$  into vectors corresponding to each block

$$\mathbf{w}_i = (\mathbf{w}_{i,1}; \dots; \mathbf{w}_{i,K}) \in \mathbb{R}^d,$$

where  $\mathbf{w}_{i,k} \in \mathbb{R}^{d_k}$ , and we denote  $\mathbf{W}_k = [\mathbf{w}_{1,k}, \mathbf{w}_{2,k}, \dots, \mathbf{w}_{N,k}] \in \mathbb{R}^{d_k \times N}$  for  $k \in [K]$ .

**Step 1. Decouple the randomness.**

Let  $(\tilde{\mathbf{w}}_i)_{i \in [N]} \sim \mathcal{N}(\mathbf{0}, \mathbf{\Gamma})$  independently and independent of  $(\mathbf{w}_i)_{i \in [N]}$ . We divide  $\tilde{\mathbf{w}}_i$  into segments corresponding to each blocks

$$\tilde{\mathbf{w}}_i = (\tilde{\mathbf{w}}_{i,1}; \dots; \tilde{\mathbf{w}}_{i,K}),$$

where  $\tilde{\mathbf{w}}_{i,k} \in \mathbb{R}^{d_k}$ , and we denote  $\tilde{\mathbf{W}}_k = [\tilde{\mathbf{w}}_{1,k}, \tilde{\mathbf{w}}_{2,k}, \dots, \tilde{\mathbf{w}}_{N,k}] \in \mathbb{R}^{d_k \times N}$  for  $k \in [K]$ .

Define

$$\begin{aligned}\mathbf{D}_{k,\mathbf{w}} &= \text{diag}(\|\mathbf{w}_{1,k}\|_2, \dots, \|\mathbf{w}_{N,k}\|_2) \in \mathbb{R}^{N \times N}, \\ \mathbf{D}_{k,\tilde{\mathbf{w}}} &= \text{diag}(\|\tilde{\mathbf{w}}_{1,k}\|_2, \dots, \|\tilde{\mathbf{w}}_{N,k}\|_2) \in \mathbb{R}^{N \times N}.\end{aligned}$$

Using the fact that  $\|\mathbf{g}\|_2$  is independent of  $\mathbf{g}/\|\mathbf{g}\|_2$  for  $\mathbf{g} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ , the following two sets of random variables have the same distribution:

$$\left\{ (\mathbf{W}_k^\top \mathbf{W}_k)_{k \in [K]}, (\|\mathbf{w}_{ik}\|_2)_{i \in [N], k \in [K]} \right\} \stackrel{d}{=} \left\{ (\mathbf{D}_{k,\mathbf{w}} \mathbf{D}_{k,\tilde{\mathbf{w}}}^{-1} \tilde{\mathbf{W}}_k^\top \tilde{\mathbf{W}}_k \mathbf{D}_{k,\tilde{\mathbf{w}}}^{-1} \mathbf{D}_{k,\mathbf{w}})_{k \in [K]}, (\|\mathbf{w}_{ik}\|_2)_{i \in [N], k \in [K]} \right\}.$$

Define

$$\bar{\mathbf{A}}_0 = c_1 \mathbf{I}_d + c_2 \sum_{k \in [K]} \mathbf{D}_{k,\mathbf{w}} \mathbf{D}_{k,\tilde{\mathbf{w}}}^{-1} \tilde{\mathbf{W}}_k^\top \tilde{\mathbf{W}}_k \mathbf{D}_{k,\tilde{\mathbf{w}}}^{-1} \mathbf{D}_{k,\mathbf{w}}.$$

Then we have

$$\langle \mathbf{1}, \mathbf{A}_0^{-1} \boldsymbol{\mu} \rangle / \sqrt{d} \stackrel{d}{=} \langle \mathbf{1}, \bar{\mathbf{A}}_0^{-1} \boldsymbol{\mu} \rangle / \sqrt{d}. \quad (47)$$

**Step 2. Bound the difference between  $\bar{\mathbf{A}}_0$  and  $\tilde{\mathbf{A}}_0$ .**

Define

$$\tilde{\mathbf{A}}_0 = c_1 \mathbf{I}_d + c_2 \sum_{k \in [K]} \tilde{\mathbf{W}}_k^\top \tilde{\mathbf{W}}_k.$$



Since  $d_k \rightarrow \infty$  as  $d \rightarrow \infty$ , we have

$$\|\mathbf{D}_{k,\tilde{\mathbf{w}}}^{-1}\mathbf{D}_{k,\mathbf{w}} - \mathbf{I}_N\|_{\text{op}} = o_{d,\mathbb{P}}(1),$$

and hence

$$\|\tilde{\mathbf{A}}_0 - \bar{\mathbf{A}}_0\|_{\text{op}} \leq 2c_2 \sum_{k \in [K]} \|\mathbf{D}_{k,\mathbf{w}}\mathbf{D}_{k,\tilde{\mathbf{w}}} - \mathbf{I}_d\|_{\text{op}} \|\tilde{\mathbf{W}}_k^T \tilde{\mathbf{W}}_k\|_{\text{op}} \|\mathbf{D}_{k,\mathbf{w}}\mathbf{D}_{k,\tilde{\mathbf{w}}}\|_{\text{op}} = o_{d,\mathbb{P}}(1).$$

By definition,  $\tilde{\mathbf{A}}_0, \bar{\mathbf{A}}_0 \succeq c_1 \mathbf{I}$  and therefore  $\|\tilde{\mathbf{A}}_0^{-1}\|_{\text{op}}, \|\bar{\mathbf{A}}_0^{-1}\|_{\text{op}} = O_{d,\mathbb{P}}(1)$ . We deduce

$$\|\tilde{\mathbf{A}}_0^{-1} - \bar{\mathbf{A}}_0^{-1}\|_{\text{op}} = \|\tilde{\mathbf{A}}_0^{-1}(\bar{\mathbf{A}}_0 - \tilde{\mathbf{A}}_0)\bar{\mathbf{A}}_0^{-1}\|_{\text{op}} = o_{d,\mathbb{P}}(1).$$

This gives (recalling that  $\|\boldsymbol{\mu}\|_2^2 = O_{d,\mathbb{P}}(1)$ )

$$\langle \mathbf{1}, \bar{\mathbf{A}}_0^{-1} \boldsymbol{\mu} \rangle / \sqrt{d} - \langle \mathbf{1}, \tilde{\mathbf{A}}_0^{-1} \boldsymbol{\mu} \rangle / \sqrt{d} = o_{d,\mathbb{P}}(1). \quad (48)$$

**Step 3. Calculating the second moment of  $\langle \mathbf{1}, \tilde{\mathbf{A}}_0^{-1} \boldsymbol{\mu} \rangle / \sqrt{d}$ .**

Since we have

$$\mathbb{E}_{\mathbf{W}}[(\langle \mathbf{1}, \tilde{\mathbf{A}}_0^{-1} \boldsymbol{\mu} \rangle / \sqrt{d})^2] = \langle \mathbf{1}, \tilde{\mathbf{A}}_0^{-2} \mathbf{1} \rangle / d \cdot \mathbb{E}_{\mathbf{w} \sim \mathbf{N}(\mathbf{0}, \mathbf{\Gamma})}[(\|\mathbf{w}\|_2^2 - 1)^2].$$

Note that

$$\mathbb{E}_{\mathbf{w} \sim \mathbf{N}(\mathbf{0}, \mathbf{\Gamma})}[(\|\mathbf{w}\|_2^2 - 1)^2] = O_{d,\mathbb{P}}(1/d),$$

and using that  $\|\tilde{\mathbf{A}}_0^{-1}\|_{\text{op}} = O_{d,\mathbb{P}}(1)$ ,

$$\langle \mathbf{1}, \tilde{\mathbf{A}}_0^{-2} \mathbf{1} \rangle / d = O_{d,\mathbb{P}}(1).$$

Therefore

$$\mathbb{E}_{\mathbf{W}}[(\langle \mathbf{1}, \tilde{\mathbf{A}}_0^{-1} \boldsymbol{\mu} \rangle / \sqrt{d})^2] = o_{d,\mathbb{P}}(1).$$

By Chebyshev inequality we have

$$\langle \mathbf{1}, \tilde{\mathbf{A}}_0^{-1} \boldsymbol{\mu} \rangle / \sqrt{d} = o_{d,\mathbb{P}}(1). \quad (49)$$

Combining (47), (48) and (49) proves the lemma in the case of a covariance of the form (46):

$$\langle \mathbf{1}, \mathbf{A}_0^{-1} \boldsymbol{\mu} \rangle / \sqrt{d} = o_{d,\mathbb{P}}(1). \quad (50)$$

**Step 4. From discrete to continuous spectrum.**

We consider  $\mathbf{\Gamma}$  a covariance matrix verifying assumption **A2**. For a given  $\varepsilon > 0$  and  $K$  sufficiently large, we consider  $\mathbf{\Gamma}_\varepsilon$  a matrix obtained from  $\mathbf{\Gamma}$  by binning its eigenvalues to at most  $K$  points of  $[0, C/d]$ , such that we have  $\text{Tr}(\mathbf{\Gamma}_\varepsilon) = 1$  and  $\lim_{d \rightarrow \infty} d \cdot \|\mathbf{\Gamma} - \mathbf{\Gamma}_\varepsilon\|_{\text{op}} \leq \varepsilon$  (recall that  $\|\mathbf{\Gamma}\|_{\text{op}} \leq C/d$  by assumption). Such a matrix always exists from the condition  $\text{Tr}(\mathbf{\Gamma}) = 1$  and the weak convergence of the spectrum of  $d \cdot \mathbf{\Gamma}$ .

By construction  $\mathbf{\Gamma}_\varepsilon$  is of the form (46). Consider  $\mathbf{G} = (\mathbf{g}_1, \dots, \mathbf{g}_N) \in \mathbb{R}^{d \times N}$  where  $\mathbf{g}_i \sim_{i.i.d.} \mathbf{N}(\mathbf{0}, \mathbf{I}_d)$ . We define:

$$\begin{aligned} \boldsymbol{\mu} &= (\|\mathbf{\Gamma}^{1/2} \mathbf{g}_i\|_2^2 - 1)_{i \in [N]}, & \boldsymbol{\mu}_\varepsilon &= (\|\mathbf{\Gamma}_\varepsilon^{1/2} \mathbf{g}_i\|_2^2 - 1)_{i \in [N]}, \\ \mathbf{A}_0 &= c_1 \mathbf{I}_d + c_2 \mathbf{G}^T \mathbf{\Gamma} \mathbf{G}, & \mathbf{A}_{0,\varepsilon} &= c_1 \mathbf{I}_d + c_2 \mathbf{G}^T \mathbf{\Gamma}_\varepsilon \mathbf{G}. \end{aligned}$$

We have for  $d$  sufficiently large,

$$\|\mathbf{A}_0 - \mathbf{A}_{0,\varepsilon}\|_{\text{op}} = \|\mathbf{G}^T (\mathbf{\Gamma} - \mathbf{\Gamma}_\varepsilon) \mathbf{G}\|_{\text{op}} \leq \|\mathbf{G}\|_{\text{op}}^2 \|\mathbf{\Gamma} - \mathbf{\Gamma}_\varepsilon\|_{\text{op}} \leq 2\varepsilon \|\mathbf{G}\|_{\text{op}}^2 / d.$$

Furthermore, using  $\text{Tr}(\mathbf{\Gamma} - \mathbf{\Gamma}_\varepsilon) = 0$ , we have

$$\mathbb{E}[\|\boldsymbol{\mu} - \boldsymbol{\mu}_\varepsilon\|_2^2] = N \mathbb{E}[(\langle \mathbf{g}_i \mathbf{g}_i^T, \mathbf{\Gamma} - \mathbf{\Gamma}_\varepsilon \rangle)^2] = 2N \|\mathbf{\Gamma} - \mathbf{\Gamma}_\varepsilon\|_F^2 \leq 2\rho \varepsilon^2.$$

Therefore

$$\begin{aligned} \left| \langle \mathbf{1}, \mathbf{A}_0^{-1} \boldsymbol{\mu} - \mathbf{A}_{0,\varepsilon}^{-1} \boldsymbol{\mu}_\varepsilon \rangle / \sqrt{d} \right| &\leq \left| \langle \mathbf{1}, \mathbf{A}_0^{-1} (\mathbf{A}_{0,\varepsilon} - \mathbf{A}_0) \mathbf{A}_{0,\varepsilon}^{-1} \boldsymbol{\mu} \rangle / \sqrt{d} \right| + \left| \langle \mathbf{1}, \mathbf{A}_{0,\varepsilon}^{-1} (\boldsymbol{\mu}_\varepsilon - \boldsymbol{\mu}) \rangle / \sqrt{d} \right| \\ &\leq \|\mathbf{A}_0^{-1}\|_{\text{op}} \|\mathbf{A}_0 - \mathbf{A}_{0,\varepsilon}\|_{\text{op}} \|\mathbf{A}_{0,\varepsilon}^{-1}\|_{\text{op}} \|\boldsymbol{\mu}\|_2 + \|\mathbf{A}_{0,\varepsilon}^{-1}\|_{\text{op}} \|\boldsymbol{\mu} - \boldsymbol{\mu}_\varepsilon\|_2. \end{aligned}$$

Noticing that  $\|\mathbf{A}_0^{-1}\|_{\text{op}}, \|\mathbf{A}_{0,\varepsilon}^{-1}\|_{\text{op}} \leq c_1^{-1}$ , and using (50) applied to  $\boldsymbol{\Gamma}_\varepsilon$ , we get for  $d$  sufficiently large:

$$\left| \langle \mathbf{1}, \mathbf{A}_0^{-1} \boldsymbol{\mu} \rangle / \sqrt{d} \right| \leq o_{d,\mathbb{P}}(1) + 2\varepsilon c_1^{-2} \|\boldsymbol{\mu}\|_2 \|\mathbf{G}\|_{\text{op}}^2 / d + c_1^{-1} \|\boldsymbol{\mu} - \boldsymbol{\mu}_\varepsilon\|_2. \quad (51)$$

We have  $\|\boldsymbol{\mu}\|_2 \|\mathbf{G}\|_{\text{op}}^2 / d = O_{d,\mathbb{P}}(1)$  hence for any  $\delta > 0$  there exists a constant  $C_\delta$  (which do not depend on  $\varepsilon$ ) such that:

$$\mathbb{P}(\varepsilon \|\boldsymbol{\mu}\|_2 \|\mathbf{G}\|_{\text{op}}^2 / d > \varepsilon C_\delta) \leq \delta.$$

Taking a sequence  $\delta \rightarrow 0$  and  $\varepsilon$  such that  $\varepsilon \propto C_\delta^{-1}$  shows that this is equivalent to

$$\varepsilon \|\boldsymbol{\mu}\|_2 \|\mathbf{G}\|_{\text{op}}^2 / d = o_{d,\mathbb{P}}(1). \quad (52)$$

By Markov inequality,

$$\lim_{d \rightarrow \infty} \mathbb{P}(\|\boldsymbol{\mu} - \boldsymbol{\mu}_\varepsilon\|_2 \geq \varepsilon \sqrt{2\rho/\delta}) \leq \delta.$$

Taking  $\varepsilon \propto \sqrt{\delta}$ , we deduce that this is equivalent to

$$\|\boldsymbol{\mu} - \boldsymbol{\mu}_\varepsilon\|_2 = o_{d,\mathbb{P}}(1). \quad (53)$$

Substituting (52) and (53) in (51) concludes the proof.  $\square$

**Lemma 5.** *Under the same setting as Proposition 2, we have*

$$\langle \mathbf{1}, \mathbf{U}_0^{-1} \mathbf{1} \rangle / d = \frac{\mathbf{1}^\top \mathbf{A}_0^{-1} \mathbf{1} / d}{1 + \kappa \mathbf{1}^\top \mathbf{A}_0^{-1} \mathbf{1} / d} + o_{d,\mathbb{P}}(1). \quad (54)$$

*Proof of Lemma 5.* Define  $\mathbf{z} = \sqrt{\kappa} \mathbf{1} / \sqrt{d}$ . Then we have

$$\mathbf{U}_0 = \mathbf{A}_0 + \mathbf{z} \mathbf{z}^\top + \boldsymbol{\mu} \boldsymbol{\mu}^\top.$$

By assumption, we have  $\kappa = O_{d,\mathbb{P}}(1)$  and therefore  $\|\mathbf{z}\|_2 = O_{d,\mathbb{P}}(1)$ . We have already seen that  $\|\mathbf{A}_0^{-1}\|_{\text{op}}, \|\mathbf{A}_0^{-1}\|_{\text{op}} = O_{d,\mathbb{P}}(1)$ . Furthermore

$$\|\mathbf{A}_0\|_{\text{op}} \leq \bar{\lambda} + \lambda_1^2 \lambda_{\max}(\mathbf{W}^\top \mathbf{W}) = O_{d,\mathbb{P}}(1).$$

By Sherman Morrison Woodbury formula, we have

$$\mathbf{1}^\top \mathbf{U}_0^{-1} \mathbf{1} / d = \mathbf{1}^\top \mathbf{A}_0^{-1} \mathbf{1} / d - \mathbf{1}^\top \mathbf{A}_0^{-1} [\mathbf{z}, \boldsymbol{\mu}] (\mathbf{I}_2 + [\mathbf{z}, \boldsymbol{\mu}]^\top \mathbf{A}_0^{-1} [\mathbf{z}, \boldsymbol{\mu}])^{-1} [\mathbf{z}, \boldsymbol{\mu}]^\top \mathbf{A}_0^{-1} \mathbf{1} / d.$$

Note that by

$$\|(\mathbf{I}_2 + [\mathbf{z}, \boldsymbol{\mu}]^\top \mathbf{A}_0^{-1} [\mathbf{z}, \boldsymbol{\mu}])^{-1}\|_F = O_{d,\mathbb{P}}(1),$$

and by Lemma 4, we have (since  $\mathbf{z}^\top \mathbf{A}_0^{-1} \boldsymbol{\mu}, \mathbf{1}^\top \mathbf{A}_0^{-1} \boldsymbol{\mu} / \sqrt{d} = o_{d,\mathbb{P}}(1)$ )

$$\begin{aligned} &\mathbf{1}^\top \mathbf{A}_0^{-1} [\mathbf{z}, \boldsymbol{\mu}] (\mathbf{I}_2 + [\mathbf{z}, \boldsymbol{\mu}]^\top \mathbf{A}_0^{-1} [\mathbf{z}, \boldsymbol{\mu}])^{-1} [\mathbf{z}, \boldsymbol{\mu}]^\top \mathbf{A}_0^{-1} \mathbf{1} / d \\ &= (\mathbf{1}^\top \mathbf{A}_0^{-1} \mathbf{z})^2 (1 + \mathbf{z}^\top \mathbf{A}_0^{-1} \mathbf{z})^{-1} / d + o_{d,\mathbb{P}}(1) = \kappa (\mathbf{1}^\top \mathbf{A}_0^{-1} \mathbf{1} / d)^2 (1 + \kappa \mathbf{1}^\top \mathbf{A}_0^{-1} \mathbf{1} / d)^{-1} + o_{d,\mathbb{P}}(1). \end{aligned}$$

This proves the lemma.  $\square$

In the following, we give an asymptotic expression for  $\langle \mathbf{1}, \mathbf{A}_0^{-1} \mathbf{1} \rangle / d$ .

**Lemma 6.** Let  $(\mathbf{w}_i)_{i \in [N]} \sim \mathbf{N}(\mathbf{0}, \mathbf{\Gamma})$  independently, while  $\mathbf{\Gamma}$  satisfies assumption **A2** (resp. **B2**). Denote  $\mathbf{W} = (\mathbf{w}_1, \dots, \mathbf{w}_N) \in \mathbb{R}^{d \times N}$ . Let  $\tilde{\lambda}$  and  $\lambda_1$  be two positive constants. Define

$$\mathbf{A}_0 = \tilde{\lambda} \mathbf{I}_N + \lambda_1^2 \mathbf{W}^\top \mathbf{W}.$$

Let  $\rho \in (0, \infty)$ . We have almost surely

$$\lim_{N/d=\rho, d \rightarrow \infty} |\mathbf{1}^\top \mathbf{A}_0^{-1} \mathbf{1}/d - \text{Tr}(\mathbf{A}_0^{-1})/d| = 0. \quad (55)$$

In addition, assume  $\mathcal{D}$  is the limiting spectral distribution of  $d \cdot \mathbf{\Gamma}$ . Then, we have almost surely

$$\lim_{N/d=\rho, d \rightarrow \infty} \frac{1}{d} \text{Tr}(\mathbf{A}_0^{-1}) = m_{\mathcal{D}}(-\tilde{\lambda}), \quad (56)$$

where  $m_{\mathcal{D}}(\cdot) : \mathbb{C}^+ \rightarrow \mathbb{C}^+$  is the companion Stieltjes transform associated with  $\mathcal{D}$ . For any  $x \in \mathbb{C}^+$ ,  $m_{\mathcal{D}}(x)$  satisfies the so called Silverstein's equation:

$$x = -\frac{\rho}{m_{\mathcal{D}}(x)} + \int \frac{\lambda_1^2 t}{1 + \lambda_1^2 t m_{\mathcal{D}}(x)} \mathcal{D}(dt). \quad (57)$$

*Proof of Lemma 6.* Consider the event

$$A_N(t) := \{|\mathbf{1}^\top \mathbf{A}_0^{-1} \mathbf{1}/d - \text{Tr}(\mathbf{A}_0^{-1})/d| > t\}.$$

Let  $\mathbf{Q} \in \mathbb{R}^{N \times N}$  be an orthogonal matrix. By rotation invariance of Gaussian random variables,  $\mathbf{Q}\mathbf{W}^\top$  has the same distribution as  $\mathbf{W}$ . In fact, by Fubini's theorem, we can draw  $\mathbf{Q}$  uniformly (independent of  $\mathbf{A}_0$ ) from orthogonal matrices and the distribution would still be unchanged. Let

$$\tilde{A}_N(t) := \{|\mathbf{1}^\top (\mathbf{Q}\mathbf{A}_0^{-1}\mathbf{Q}^\top)^{-1} \mathbf{1}/d - \text{Tr}(\mathbf{Q}\mathbf{A}_0^{-1}\mathbf{Q}^\top)/d| > t\}.$$

By the argument above,

$$\mathbb{P}[A_N(t)] = \mathbb{P}[\tilde{A}_N(t)].$$

Since  $\mathbf{Q}$  is orthogonal,  $\tilde{A}_N(t)$  can be written as

$$\{|\mathbf{1}^\top \mathbf{Q}\mathbf{A}_0^{-1}\mathbf{Q}^\top \mathbf{1}/d - \text{Tr}(\mathbf{A}_0^{-1})/d| > t\}. \quad (58)$$

Since  $\mathbf{Q}$  is a uniformly chosen orthogonal matrix,  $\mathbf{Q}^\top \mathbf{1}/\sqrt{d}$  is uniformly distributed on  $\mathbb{S}^{N-1}(\sqrt{\rho})$ , independently of  $\mathbf{A}_0$ . Hence  $\mathbf{Q}^\top \mathbf{1}/\sqrt{d}$  has the same distribution as  $\sqrt{\rho} \mathbf{z}/\|\mathbf{z}\|_2$  where  $\mathbf{z} \sim \mathbf{N}(0, \mathbf{I}_N)$ . In particular,

$$\mathbb{P}[\tilde{A}_N(t)] = \mathbb{P}\left\{\left|\frac{1}{\|\mathbf{z}\|_2^2} \mathbf{z}^\top \mathbf{A}_0^{-1} \mathbf{z} - \text{Tr}(\mathbf{A}_0^{-1})/N\right| > \frac{t}{\rho}\right\} \quad (59)$$

$$\leq \mathbb{P}\left\{\left|\frac{N}{\|\mathbf{z}\|_2^2} - 1\right| \mathbf{z}^\top \mathbf{A}_0^{-1} \mathbf{z}/N + |\mathbf{z}^\top \mathbf{A}_0^{-1} \mathbf{z}/N - \text{Tr}(\mathbf{A}_0^{-1})/N| > \frac{t}{\rho}\right\} \quad (60)$$

$$\leq P_1 + P_2, \quad (61)$$

where

$$P_1 = \mathbb{P}\left\{\left|\frac{N}{\|\mathbf{z}\|_2^2} - 1\right| \mathbf{z}^\top \mathbf{A}_0^{-1} \mathbf{z}/N > \frac{t}{2\rho}\right\}, \quad P_2 = \mathbb{P}\left\{|\mathbf{z}^\top \mathbf{A}_0^{-1} \mathbf{z}/N - \text{Tr}(\mathbf{A}_0^{-1})/N| > \frac{t}{2\rho}\right\}.$$

Let's consider  $P_1$  first. Since  $\mathbf{A}_0^{-1} \preceq \mathbf{I}/\tilde{\lambda}$ , we have

$$\frac{\mathbf{z}^\top \mathbf{A}_0^{-1} \mathbf{z}}{N} \leq \frac{1}{\tilde{\lambda}} \frac{\|\mathbf{z}\|^2}{N},$$

which yields

$$P_1 \leq \mathbb{P}\left\{\left|\frac{N}{\|\mathbf{z}\|^2} - 1\right| \frac{\|\mathbf{z}\|^2}{N} > \frac{\tilde{\lambda}t}{2\rho}\right\} = \mathbb{P}\left\{\left|\frac{\|\mathbf{z}\|^2}{N} - 1\right| > \frac{\tilde{\lambda}t}{2\rho}\right\}. \quad (62)$$

We know due to fast concentration of  $\|\mathbf{z}\|^2/N$  around one (see e.g. [BLM13]),  $P_1$  vanish exponentially fast in  $N$  (equivalently in  $d$  since  $N/d$  is fixed to be  $\rho$ ).

Now, let's consider  $P_2$ .  $|\mathbf{z}^\top \mathbf{A}_0^{-1} \mathbf{z}/N - \text{Tr}(\mathbf{A}_0^{-1})/N|$ . By Hanson-Wright inequality (see e.g. [BLM13]), we have

$$\mathbb{P}\left(|\mathbf{z}^\top \mathbf{A}_0^{-1} \mathbf{z}/N - \text{Tr}(\mathbf{A}_0^{-1})/N| > \frac{t}{2\rho} \middle| \mathbf{A}_0\right) \leq 2 \exp\left\{-c \min\left(\frac{t^2}{\|\mathbf{A}_0^{-1}/N\|_F^2}, \frac{t}{\|\mathbf{A}_0^{-1}/N\|_{\text{op}}}\right)\right\} \quad (63)$$

$$\leq 2 \exp\left\{-c' \min\left(N\tilde{\lambda}^2 t^2, \tilde{\lambda} t N\right)\right\}. \quad (64)$$

Since the bound in (64) is independent of  $\mathbf{A}_0$ , it holds unconditionally. Therefore, we conclude  $P_2$  vanishes exponentially fast in  $N$  and  $d$ . We conclude that  $\Pr[\tilde{A}_N(t)]$  vanishes exponentially fast as  $d, N \rightarrow \infty$ . Therefore, by Borel-Cantelli lemma we recover (55).

Convergence of  $\text{Tr}(\mathbf{A}_0^{-1})/d$  to  $m_D(-\tilde{\lambda})$  is a standard result in random matrix theory. We refer the reader to [BS10] Chapters 3 and 6.  $\square$

### B.1.5 Proof of Theorem 1

By Lemma 1, the risk has a representation

$$R_{\text{RF},N}(f_*) = \mathbb{E}_{\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d)}[f_*(\mathbf{x})^2] - \mathbf{V}^\top \mathbf{U}^{-1} \mathbf{V}.$$

By Lemma 2, we have

$$\|\mathbf{U} - \mathbf{U}_0\|_{\text{op}} = o_{d,\mathbb{P}}(1).$$

By Lemma 3, we have

$$\|\mathbf{V} - \tau \mathbf{1}/\sqrt{d}\|_2 = \|\mathbf{B}\|_F \cdot o_{d,\mathbb{P}}(1),$$

where

$$\tau = \sqrt{d} \cdot \lambda_2 \text{Tr}(\mathbf{B}\mathbf{\Gamma}).$$

Hence, we have

$$|\mathbf{V}^\top \mathbf{U}^{-1} \mathbf{V} - \tau^2 \mathbf{1}^\top \mathbf{U}_0^{-1} \mathbf{1}/d| = \|\mathbf{B}\|_F^2 \cdot o_{d,\mathbb{P}}(1).$$

Proposition 2 gives the expression for

$$\mathbf{1}^\top \mathbf{U}_0^{-1} \mathbf{1}/d = \psi/(1 + \kappa\psi) + o_{d,\mathbb{P}}(1),$$

where

$$\kappa = d \cdot \lambda_2^2 \text{Tr}(\mathbf{\Gamma}^2)/2.$$

Hence we have

$$\mathbf{V}^\top \mathbf{U} \mathbf{V} = \tau^2 \psi/(1 + \kappa\psi) + \|\mathbf{B}\|_F^2 \cdot o_{d,\mathbb{P}}(1).$$

Recalling the assumption  $\mathbb{E}(f_*) = 0$ , we have  $\|f_*\|_{L^2}^2 = 2\|\mathbf{B}\|_F^2$ , which concludes the proof.

## B.2 Neural Tangent model: proof of Theorem 2

Recall the definition

$$R_{\text{NT},N}(f_*) = \min_{\hat{f} \in \mathcal{F}_{\text{NT},N}(\mathbf{W})} \mathbb{E}\{(f_*(\mathbf{x}) - \hat{f}(\mathbf{x}))^2\},$$

where

$$\mathcal{F}_{\text{NT},N}(\mathbf{W}) = \left\{f_N(\mathbf{x}) = c + \sum_{i=1}^N \sigma'(\langle \mathbf{w}_i, \mathbf{x} \rangle) \langle \mathbf{a}_i, \mathbf{x} \rangle : c \in \mathbb{R}, \mathbf{a}_i \in \mathbb{R}^d, i \in [N]\right\}.$$

*Proof of Theorem 2.* We can rewrite the neural tangent model with a squared non-linearity  $\sigma(x) = x^2$  as

$$\hat{f}(\mathbf{W}, \mathbf{A}, c) = 2 \sum_{i=1}^N \langle \mathbf{w}_i, \mathbf{x} \rangle \langle \mathbf{a}_i, \mathbf{x} \rangle + c = 2 \langle \mathbf{W} \mathbf{A}^\top, \mathbf{x} \mathbf{x}^\top \rangle + c.$$

with  $\mathbf{W} = [\mathbf{w}_1, \dots, \mathbf{w}_N] \in \mathbb{R}^{d \times N}$  and  $\mathbf{A} = [\mathbf{a}_1, \dots, \mathbf{a}_N] \in \mathbb{R}^{d \times N}$ . Note that we have

$$\begin{aligned} & \mathbb{E}_{\mathbf{x}}[\langle \mathbf{B} - 2\mathbf{W} \mathbf{A}^\top, \mathbf{x} \mathbf{x}^\top \rangle + b_0 - c]^2 \\ &= 2\|\mathbf{B} - \mathbf{W} \mathbf{A}^\top - \mathbf{A} \mathbf{W}^\top\|_F^2 + \text{Tr}(\mathbf{B} - 2\mathbf{W} \mathbf{A}^\top)^2 - 2\text{Tr}(\mathbf{B} - 2\mathbf{W} \mathbf{A}^\top)(c - b_0) + (c - b_0)^2, \end{aligned}$$

which, after minimizing over  $c \in \mathbb{R}$ , simplifies to:

$$\min_{c \in \mathbb{R}} \|f_* - \hat{f}(\mathbf{W}, \mathbf{A}, c)\|_{L^2}^2 = 2\|\mathbf{B} - \mathbf{W} \mathbf{A}^\top - \mathbf{A} \mathbf{W}^\top\|_F^2.$$

For  $\mathbf{w}_i \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d)$ , we have  $\text{rank}(\mathbf{W}) = \min(d, N) \equiv r$  with probability one. Let  $\mathbf{W} = \mathbf{P}_1 \mathbf{S} \mathbf{V}^\top$  be the singular value decomposition of  $\mathbf{W}$ , with  $\mathbf{P}_1 \in \mathbb{R}^{d \times r}$ ,  $\mathbf{S} \in \mathbb{R}^{r \times r}$  and  $\mathbf{V} \in \mathbb{R}^{N \times r}$ . Defining  $\mathbf{G} = \mathbf{S} \mathbf{V}^\top \mathbf{A} \in \mathbb{R}^{r \times d}$ , we get almost surely

$$\min_{\mathbf{A} \in \mathbb{R}^{d \times N}, c \in \mathbb{R}} \|f_* - \hat{f}(\mathbf{W}, \mathbf{A}, c)\|_{L^2}^2 = \min_{\mathbf{G} \in \mathbb{R}^{r \times d}} 2\|\mathbf{B} - \mathbf{P}_1 \mathbf{G} - \mathbf{G}^\top \mathbf{P}_1^\top\|_F^2.$$

In the case  $N \geq d$ , we can take  $\mathbf{G} = \mathbf{P}_1^\top \mathbf{B} / 2$  and we get almost surely over  $\mathbf{W} \in \mathbb{R}^{d \times N}$

$$R_{\text{NT}, N}(f_*) = 0.$$

Consider the case when  $N < d$ , we define  $\mathbf{P}_2 \in \mathbb{R}^{d \times (d-N)}$  the completion of  $\mathbf{P}_1$  to a full basis  $\mathbf{P} = [\mathbf{P}_1, \mathbf{P}_2] \in \mathbb{R}^{d \times d}$ . We define  $\mathbf{G}_1 = \mathbf{G} \mathbf{P}_1 \in \mathbb{R}^{N \times N}$  and  $\mathbf{G}_2 = \mathbf{G} \mathbf{P}_2 \in \mathbb{R}^{N \times (d-N)}$  and we perform our computation in the  $\mathbf{P}$  basis. We have

$$\mathbf{B} - \mathbf{P}_1 \mathbf{G} - \mathbf{G}^\top \mathbf{P}_1^\top = \begin{pmatrix} \mathbf{B}_{11} - \mathbf{G}_1 - \mathbf{G}_1^\top & \mathbf{B}_{12} - \mathbf{G}_2 \\ \mathbf{B}_{21} - \mathbf{G}_2^\top & \mathbf{B}_{22} \end{pmatrix},$$

where  $\mathbf{B}_{ij} = \mathbf{P}_i^\top \mathbf{B} \mathbf{P}_j$  for  $i, j = 1, 2$ . We readily deduce that

$$\min_{\mathbf{G} \in \mathbb{R}^{r \times d}} 2\|\mathbf{B} - \mathbf{P}_1 \mathbf{G} - \mathbf{G}^\top \mathbf{P}_1^\top\|_F^2 = 2\|\mathbf{P}_2^\top \mathbf{B} \mathbf{P}_2\|_F^2.$$

Let us compute its expectation over  $\mathbf{w}_i \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d)$ , i.e over  $\mathbf{P}_2 = [\mathbf{v}_1, \dots, \mathbf{v}_{d-N}]$  where the  $\mathbf{v}_i \in \mathbb{R}^d$  are  $(d-N)$  orthogonal vectors uniformly distributed on the unit sphere in  $\mathbb{R}^d$ . Let  $\mathbf{B} = \sum_{i=1}^s \lambda_i \mathbf{e}_i \mathbf{e}_i^\top$  with  $\mathbf{e}_i$  the orthonormal eigenvectors of  $\mathbf{B}$ . We get:

$$\begin{aligned} \mathbb{E}[\|\mathbf{P}_2^\top \mathbf{B} \mathbf{P}_2\|_F^2] &= \sum_{i,j=1}^s \sum_{k,l=1}^{d-N} \lambda_i \lambda_j \mathbb{E}[\langle \mathbf{v}_k, \mathbf{e}_i \rangle \langle \mathbf{v}_l, \mathbf{e}_j \rangle \langle \mathbf{v}_l, \mathbf{e}_i \rangle \langle \mathbf{v}_k, \mathbf{e}_j \rangle] \\ &= \|\mathbf{B}\|_F^2 (d-N) \mathbb{E}[\langle \mathbf{v}_1, \mathbf{e}_1 \rangle^4] + \|\mathbf{B}\|_F^2 (d-N)(d-N-1) \mathbb{E}[\langle \mathbf{v}_1, \mathbf{e}_1 \rangle^2 \langle \mathbf{v}_2, \mathbf{e}_1 \rangle^2] \\ &\quad + 2 \left( \sum_{i < j} \lambda_i \lambda_j \right) (d-N) \mathbb{E}[\langle \mathbf{v}_1, \mathbf{e}_1 \rangle^2 \langle \mathbf{v}_1, \mathbf{e}_2 \rangle^2] \\ &\quad + 2 \left( \sum_{i < j} \lambda_i \lambda_j \right) (d-N)(d-N-1) \mathbb{E}[\langle \mathbf{v}_1, \mathbf{e}_1 \rangle \langle \mathbf{v}_1, \mathbf{e}_2 \rangle \langle \mathbf{v}_2, \mathbf{e}_1 \rangle \langle \mathbf{v}_2, \mathbf{e}_2 \rangle]. \end{aligned} \tag{65}$$

We bound each term separately. For  $\mathbf{u} \sim \text{Unif}(\mathbb{S}^{d-1})$ , we have the convergence in distribution of the first two coordinates  $\sqrt{d}(u_1, u_2) \Rightarrow \mathcal{N}(\mathbf{0}, \mathbf{I}_2)$ , hence:

$$\lim_{d \rightarrow \infty} d^2 \mathbb{E}[\langle \mathbf{v}_1, \mathbf{e}_1 \rangle^4] = 3, \quad \lim_{d \rightarrow \infty} d^2 \mathbb{E}[\langle \mathbf{v}_1, \mathbf{e}_1 \rangle^2 \langle \mathbf{v}_1, \mathbf{e}_2 \rangle^2] = 1. \tag{66}$$

Furthermore, conditioned on  $\mathbf{v}_1$ ,  $\mathbf{v}_2$  is uniformly distributed over the sphere  $\mathbb{S}^{d-2}$  in the hyperplane orthogonal to  $\mathbf{v}_1$ . We get the uniform convergence

$$\lim_{d \rightarrow \infty} \sup_{\mathbf{v}_1 \in \mathbb{S}^{d-1}} |d\mathbb{E}[\langle \mathbf{v}_2, \mathbf{e}_1 \rangle^2 | \mathbf{v}_1] - (1 - \langle \mathbf{v}_1, \mathbf{e}_1 \rangle^2)| = 0.$$

By dominated convergence theorem, we get

$$\lim_{d \rightarrow \infty} d^2 \mathbb{E}[\langle \mathbf{v}_1, \mathbf{e}_1 \rangle^2 \langle \mathbf{v}_2, \mathbf{e}_1 \rangle^2] = 1. \quad (67)$$

The last term of the sum (65) is also derived by first conditioning on  $\mathbf{v}_1$ . Let us denote  $\mathbf{z}_1 = \mathbf{P}_{\perp \mathbf{v}_1} \mathbf{e}_1$  and  $\mathbf{z}_2 = \mathbf{P}_{\perp \mathbf{v}_1} \mathbf{e}_2$  the projections of  $(\mathbf{e}_1, \mathbf{e}_2)$  on the hyperplane perpendicular to  $\mathbf{v}_1$ , on which  $\mathbf{v}_2$  is uniformly distributed over the unit sphere. We decompose  $\mathbf{z}_2$  into two components: one along  $\mathbf{z}_1$  that we denote  $\mathbf{z}_2^{(1)} = \mathbf{P}_{\parallel \mathbf{z}_1} \mathbf{z}_2$  and one perpendicular to  $\mathbf{z}_1$ , denoted  $\mathbf{z}_2^{(2)} = \mathbf{P}_{\perp \mathbf{z}_1} \mathbf{z}_2$ . Then we have:

$$\begin{aligned} \mathbb{E}[\langle \mathbf{v}_2, \mathbf{e}_1 \rangle \langle \mathbf{v}_2, \mathbf{e}_2 \rangle | \mathbf{v}_1] &= \mathbb{E}[\langle \mathbf{v}_2, \mathbf{z}_1 \rangle \langle \mathbf{v}_2, \mathbf{z}_2 \rangle | \mathbf{v}_1] \\ &= \mathbb{E} \left[ \langle \mathbf{v}_2, \mathbf{z}_1 \rangle \left( \langle \mathbf{v}_2, \mathbf{z}_2^{(1)} \rangle + \langle \mathbf{v}_2, \mathbf{z}_2^{(2)} \rangle \right) \middle| \mathbf{v}_1 \right] \\ &= \langle \mathbf{z}_1, \mathbf{z}_2 \rangle \mathbb{E}[u_1^2] + \|\mathbf{z}_1\|_2 \|\mathbf{z}_2^{(2)}\|_2 \mathbb{E}[u_1 u_2] \\ &= \frac{\langle \mathbf{z}_1, \mathbf{z}_2 \rangle}{d-1}, \end{aligned}$$

where  $(u_1, u_2)$  are the first two coordinates of a uniform random variable on the sphere  $\mathbb{S}^{d-2}$ . Using that:

$$\langle \mathbf{z}_1, \mathbf{z}_2 \rangle = \langle \mathbf{e}_1 - \langle \mathbf{e}_1, \mathbf{v}_1 \rangle \mathbf{v}_1, \mathbf{e}_2 - \langle \mathbf{e}_2, \mathbf{v}_1 \rangle \mathbf{v}_1 \rangle = -\langle \mathbf{e}_1, \mathbf{v}_1 \rangle \langle \mathbf{e}_2, \mathbf{v}_1 \rangle,$$

we get

$$\mathbb{E}[\langle \mathbf{v}_1, \mathbf{e}_1 \rangle \langle \mathbf{v}_1, \mathbf{e}_2 \rangle \langle \mathbf{v}_2, \mathbf{e}_1 \rangle \langle \mathbf{v}_2, \mathbf{e}_2 \rangle] = -\frac{1}{d-1} \mathbb{E}[\langle \mathbf{v}_1, \mathbf{e}_1 \rangle^2 \langle \mathbf{v}_1, \mathbf{e}_2 \rangle^2] = -\frac{1}{d^3} + o_d(d^{-3}), \quad (68)$$

where we used the same argument as for (66). Plugging the above limits (66), (67) and (68) in the expansion (65), we get

$$\mathbb{E}[R_{\text{NT},N}(f_*)] = 2\|\mathbf{B}\|_F^2 \left[ (1-\rho)_+^2 + (1-\rho)_+ \frac{\text{Tr}(\mathbf{B})^2}{d\|\mathbf{B}\|_F^2} - (1-\rho)_+^2 \frac{\text{Tr}(\mathbf{B})^2}{d\|\mathbf{B}\|_F^2} + o_d(1) \right]. \quad (69)$$

Recalling the assumption  $\mathbb{E}(f_*) = 0$ , we have  $\|f_*\|_{L^2}^2 = 2\|\mathbf{B}\|_F^2$ , which concludes the proof.  $\square$

**Remark 1.** The above formula for the RF risk Eq. (69) has two terms that corresponds to the two limits  $\text{Tr}(\mathbf{B})/\|\mathbf{B}\|_F = o_d(\sqrt{d})$  (e.g. spiked matrix)

$$\mathbb{E}[R_{\text{NT},N}(f_*)] = 2(1-\rho)_+^2 \|\mathbf{B}\|_F^2 + o_d(\|\mathbf{B}\|_F^2),$$

and  $\text{Tr}(\mathbf{B})^2 = d\|\mathbf{B}\|_F^2$  (i.e.  $\mathbf{B} \propto \mathbf{I}$ )

$$\mathbb{E}[R_{\text{NT},N}(f_*)] = 2(1-\rho)_+ \|\mathbf{B}\|_F^2.$$

It is possible to show concentration of  $\|\mathbf{P}_2^\top \mathbf{B} \mathbf{P}_2\|_F^2$  on its mean  $\mathbb{E}[\|\mathbf{P}_2^\top \mathbf{B} \mathbf{P}_2\|_F^2]$  for  $\mathbf{B}$  that satisfies  $\|\mathbf{B}\|_{\text{op}} \|\mathbf{B}\|_F \leq C$  (see Theorem 5).

### B.3 Neural Network model: proof of Theorem 3

We consider two-layers neural networks with quadratic activation function  $\sigma(x) = x^2$  and we fix the second layer weights to 1,

$$\hat{f}(\mathbf{x}; \mathbf{W}, c) = \sum_{i=1}^N \langle \mathbf{w}_i, \mathbf{x} \rangle^2 + c.$$

We consider the ground truth function  $f_*$  to be a quadratic function as per Eq. (20), and the risk function defined by

$$L(\mathbf{W}, c) = \mathbb{E}_{\mathbf{x}}[(f_*(\mathbf{x}) - \hat{f}(\mathbf{x}; \mathbf{W}, c))^2] = \mathbb{E}_{\mathbf{x}}\left[\left(\langle \mathbf{x}\mathbf{x}^\top, \mathbf{B} - \mathbf{W}\mathbf{W}^\top \rangle + b_0 - c\right)^2\right].$$

We consider running SGD dynamics upon the risk function for a fresh sample  $(\mathbf{x}_k, f_*(\mathbf{x}_k))$  for each iteration

$$(\mathbf{W}_{k+1}, c_{k+1}) = (\mathbf{W}_k, c_k) - \varepsilon \nabla_{\mathbf{W}, c} \left( f_*(\mathbf{x}_k) - \hat{f}(\mathbf{x}_k; \mathbf{W}, c) \right)^2,$$

and denote

$$R_{\text{NN}, N}(f_*; \ell, \varepsilon) = \mathbb{E}_{\mathbf{x}}[(f_*(\mathbf{x}) - \hat{f}(\mathbf{x}; \mathbf{W}_\ell, c_\ell))^2].$$

### B.3.1 Global minimum

**Lemma 7.** *Let  $f_* = \langle \mathbf{x}, \mathbf{B}\mathbf{x} \rangle + b_0$  for some  $\mathbf{B} \succeq 0$  and  $b_0 \in \mathbb{R}$ . Denote by  $(\lambda_i(\mathbf{B}))_{i \in [r]}$  the positive eigenvalues of  $\mathbf{B}$  in descending order. Then we have*

$$\inf_{\mathbf{W}, c} L(\mathbf{W}, c) = 2 \sum_{i=N+1}^r \lambda_i(\mathbf{B})^2.$$

*Proof of Lemma 7.* Note we have

$$\begin{aligned} L(\mathbf{W}, c) &= \mathbb{E}_{\mathbf{x}}[(\langle \mathbf{B} - \mathbf{W}\mathbf{W}^\top, \mathbf{x}\mathbf{x}^\top \rangle + b_0 - c)^2] \\ &= 2\|\mathbf{B} - \mathbf{W}\mathbf{W}^\top\|_F^2 + \text{Tr}(\mathbf{B} - \mathbf{W}\mathbf{W}^\top)^2 - 2\text{Tr}(\mathbf{B} - \mathbf{W}\mathbf{W}^\top)(c - b_0) + (c - b_0)^2, \end{aligned}$$

minimizing over  $c$  gives

$$\inf_c L(\mathbf{W}, c) = 2\|\mathbf{B} - \mathbf{W}\mathbf{W}^\top\|_F^2.$$

The infimum of  $L$  over  $\mathbf{W}$  is equivalent to the low-rank approximation problem of matrix  $\mathbf{B}$  in Frobenius norm, with rank less or equal to  $\max(d, N)$ , and is given by the Eckart-Young-Mirsky theorem (see [EY36]).  $\square$

### B.3.2 Landscape: proof of Proposition 1

Without loss of generality, throughout the proof, we assume that  $\mathbf{B}$  is diagonal and  $b_0 = 0$ . Our first proposition characterizes the critical points of  $L(\mathbf{W}, c)$ .

**Proposition 3.** *Let  $\mathbf{W} \in \mathbb{R}^{d \times N}$ , and  $\mathbf{B} \in \mathbb{R}^{d \times d}$  to be a positive semi-definite diagonal matrix. Define the risk function to be*

$$L(\mathbf{W}, c) = \mathbb{E}_{\mathbf{x}}[(\langle \mathbf{B} - \mathbf{W}\mathbf{W}^\top, \mathbf{x}\mathbf{x}^\top \rangle - c)^2].$$

*Then for any critical point  $(\mathbf{W}_0, c_0)$  of  $L(\mathbf{W}, c)$ , there exists a projection matrix  $\mathbf{P} = \sum_{i=1}^k \mathbf{e}_{\tau(i)} \mathbf{e}_{\tau(i)}^\top$  for some injection  $\tau: [k] \rightarrow [d]$ , such that  $\mathbf{\Gamma}_0 = \mathbf{W}_0 \mathbf{W}_0^\top$  is diagonal and satisfy*

$$\begin{aligned} \mathbf{\Gamma}_0 &= \mathbf{P} \mathbf{B} \mathbf{P}, \\ c_0 &= \text{Tr}(\mathbf{B} - \mathbf{\Gamma}_0). \end{aligned}$$

*Proof.* Calculating the risk function, we get

$$L(\mathbf{W}, c) = c^2 + 2c \cdot \text{Tr}(\mathbf{W}\mathbf{W}^\top - \mathbf{B}) + \text{Tr}(\mathbf{W}\mathbf{W}^\top - \mathbf{B})^2 + 2\|\mathbf{W}\mathbf{W}^\top - \mathbf{B}\|_F^2.$$

We consider the gradient of this function. We get:

$$\begin{aligned} \frac{\partial}{\partial c} L(\mathbf{W}, c) &= 2c + 2\text{Tr}(\mathbf{W}\mathbf{W}^\top - \mathbf{B}), \\ \nabla_{\mathbf{W}} L(\mathbf{W}, c) &= 2c\mathbf{W} + 2\text{Tr}(\mathbf{W}\mathbf{W}^\top - \mathbf{B})\mathbf{W} + 8(\mathbf{W}\mathbf{W}^\top - \mathbf{B})\mathbf{W}. \end{aligned}$$

By the stationary condition, at a critical point  $(\mathbf{W}_0, c_0)$ , we must have:

$$c_0 = -\text{Tr}(\mathbf{W}_0 \mathbf{W}_0^\top - \mathbf{B}), \quad (70)$$

$$\mathbf{B} \mathbf{W}_0 = \mathbf{W}_0 \mathbf{W}_0^\top \mathbf{W}_0. \quad (71)$$

Let us denote  $\mathbf{W}_0 = \mathbf{U} \mathbf{S} \mathbf{V}^\top$  the (extended) singular value decomposition of  $\mathbf{W}_0 \in \mathbb{R}^{d \times N}$  with  $\mathbf{U} \in \mathbb{R}^{d \times d}$ ,  $\mathbf{S} \in \mathbb{R}^{d \times N}$  and  $\mathbf{V} \in \mathbb{R}^{N \times N}$ . Then the stationary condition (71) gives

$$\mathbf{B} \mathbf{U} \mathbf{S} \mathbf{V}^\top = \mathbf{U} \mathbf{S}^3 \mathbf{V}^\top. \quad (72)$$

Let  $r$  be the rank of  $\mathbf{W}_0$  and  $\mathbf{S} = \text{diag}(\mathbf{S}_1, \mathbf{0})$ ,  $\mathbf{U} = (\mathbf{U}_1, \mathbf{U}_2)$  with  $\mathbf{S}_1 \in \mathbb{R}^{r \times r}$ ,  $\mathbf{U}_1 \in \mathbb{R}^{d \times r}$  and  $\mathbf{U}_2 \in \mathbb{R}^{d \times (d-r)}$ . Then we get:

$$\mathbf{B} \mathbf{U}_1 = \mathbf{U}_1 \mathbf{S}_1^2.$$

This is of the form of the eigenvalue equation of matrix  $\mathbf{B}$ . Hence we must have the columns of  $\mathbf{U}_1$  to be a set of eigenvectors and  $\mathbf{S}_1^2$  to be positive eigenvalues of  $\mathbf{B}$ . This proves the proposition.  $\square$

Note the global minimizers are attained for  $\mathbf{\Gamma}_0 = \mathbf{W}_0 \mathbf{W}_0^\top$  corresponding to the  $\min(N, d)$  directions of  $\mathbf{B}$  with the largest eigenvalues. We prove in the following proposition that stationary points that are not global minimizers are strict saddle points.

Define the spectral separation of  $\mathbf{B}$  as

$$\delta^{\text{sep}} = \min\{|\lambda_i(\mathbf{B}) - \lambda_j(\mathbf{B})| : i, j \in [d], \lambda_i(\mathbf{B}) \neq \lambda_j(\mathbf{B})\},$$

and  $\delta^{\text{eig}}$  the minimum strictly positive eigenvalue of  $\mathbf{B}$ .

**Proposition 4.** *Consider  $(\mathbf{W}_0, c_0)$  a stationary point of  $L(\mathbf{W}, c)$  but not a global minimizer. Then, we have*

$$\lambda_{\min}(\nabla_{\mathbf{W}}^2 L(\mathbf{W}_0, c_0)) \leq -4 \min\{\delta^{\text{eig}}, \delta^{\text{sep}}\} < 0.$$

*Proof.* Let us first compute the Hessian of the risk with respect to the  $\mathbf{W}$  variable. We have

$$\begin{aligned} \langle \mathbf{Z}, \nabla_{\mathbf{W}}^2 L(\mathbf{W}, c) \mathbf{Z} \rangle &= 2c \cdot \text{Tr}(\mathbf{Z} \mathbf{Z}^\top) + 2\text{Tr}(\mathbf{W} \mathbf{W}^\top - \mathbf{B}) \text{Tr}(\mathbf{Z} \mathbf{Z}^\top) + 4\text{Tr}(\mathbf{W} \mathbf{Z}^\top)^2 \\ &\quad + 4\|\mathbf{W} \mathbf{Z}^\top\|_F^2 + 4\text{Tr}(\mathbf{W} \mathbf{Z}^\top \mathbf{W} \mathbf{Z}^\top) + 4\langle \mathbf{W} \mathbf{W}^\top - \mathbf{B}, \mathbf{Z} \mathbf{Z}^\top \rangle. \end{aligned}$$

Plugging the value of  $c_0$  at a critical point (cf Eq. (70)), we get

$$\langle \mathbf{Z}, \nabla_{\mathbf{W}}^2 L(\mathbf{W}_0, c_0) \mathbf{Z} \rangle = 4\text{Tr}(\mathbf{W}_0 \mathbf{Z}^\top)^2 + 4\|\mathbf{W}_0 \mathbf{Z}^\top\|_F^2 + 4\text{Tr}(\mathbf{W}_0 \mathbf{Z}^\top \mathbf{W}_0 \mathbf{Z}^\top) + 4\langle \mathbf{W}_0 \mathbf{W}_0^\top - \mathbf{B}, \mathbf{Z} \mathbf{Z}^\top \rangle. \quad (73)$$

**Case 1:** Consider the case  $\text{rank}(\mathbf{W}_0) < \min\{\text{rank}(\mathbf{B}), N\}$ . Then there exists an  $i \in [d]$  such that  $\mathbf{B}_{ii} > 0$  (recall that we assumed  $\mathbf{B}$  diagonal, with diagonal elements given by the positive eigenvalues of  $\mathbf{B}$ ) and  $(\mathbf{W}_0 \mathbf{W}_0^\top)_{ii} = 0$ . For simplicity, let us permute the coordinates so that  $i = 1$ . The singular value decomposition of  $\mathbf{W}_0$  verifies

$$\mathbf{W}_0 = \mathbf{U}_0 \mathbf{S}_0 \mathbf{V}_0^\top = \begin{pmatrix} 0 & 0 & \dots & 0 \\ 0 & & & \\ \vdots & \tilde{\mathbf{U}}_0 \tilde{\mathbf{S}}_0 & & \\ 0 & & & \end{pmatrix} \mathbf{V}_0^\top,$$

where  $\tilde{\mathbf{U}}_0$  and  $\tilde{\mathbf{S}}_0$  are the sub-matrices corresponding respectively to the  $(d-1) \times (d-1)$  last coordinates of  $\mathbf{U}_0$  and  $(d-1) \times (N-1)$  last coordinates of  $\mathbf{S}_0$ . Let us consider

$$\mathbf{Z} = \begin{pmatrix} 1 & 0 & \dots & 0 \\ 0 & & & \\ \vdots & & \mathbf{0} & \\ 0 & & & \end{pmatrix} \mathbf{V}_0^\top.$$



We have  $\|\mathbf{Z}\|_F = 1$  and  $\mathbf{W}_0 \mathbf{Z}^\top = 0$ . Plugging these matrices in the above expression of the Hessian, see Eq. (73), we get

$$\langle \mathbf{Z}, \nabla_{\mathbf{W}}^2 L(\mathbf{W}_0, c_0) \mathbf{Z} \rangle = -4\mathbf{B}_{11} \leq -4\delta_{\text{eig}}.$$

**Case 2:** Consider the case when  $\text{rank}(\mathbf{W}_0 \mathbf{W}_0^\top) = N < \text{rank}(\mathbf{B})$  and  $\mathbf{W}_0 \mathbf{W}_0^\top$  does not correspond to the  $N$  largest eigenvalues of  $\mathbf{B}$ . Then there exists  $i \neq j \in [n]$ , such that  $\mathbf{B}_{ii} > \mathbf{B}_{jj}$ ,  $(\mathbf{W}_0 \mathbf{W}_0^\top)_{ii} = 0$  and  $(\mathbf{W}_0 \mathbf{W}_0^\top)_{jj} = \mathbf{B}_{jj}$ . For simplicity, let us permute the coordinates such that  $i = 1$  and  $j = 2$ . The SVD decomposition of  $\mathbf{W}_0$  now verifies:

$$\mathbf{W}_0 = \mathbf{U}_0 \mathbf{S}_0 \mathbf{V}_0^\top = \begin{pmatrix} 0 & 0 & \dots & 0 \\ \sqrt{\mathbf{B}_{22}} & 0 & \dots & 0 \\ 0 & & & \\ \vdots & & \tilde{\mathbf{U}}_0 \tilde{\mathbf{S}}_0 & \\ 0 & & & \end{pmatrix} \mathbf{V}_0^\top,$$

where  $\tilde{\mathbf{U}}_0 \tilde{\mathbf{S}}_0$  is the sub-matrix of the last  $(d-2) \times (N-1)$  coordinate of  $\mathbf{U}_0 \mathbf{S}_0$ . Let us consider again

$$\mathbf{Z} = \begin{pmatrix} 1 & 0 & \dots & 0 \\ 0 & & & \\ \vdots & & \mathbf{0} & \\ 0 & & & \end{pmatrix} \mathbf{V}_0^\top.$$

We have  $\|\mathbf{Z}\|_F = 1$ . Plugging these matrices in the above expression of the Hessian (73), note

$$\text{Tr}(\mathbf{W}_0 \mathbf{Z}^\top) = \text{Tr}(\mathbf{W}_0 \mathbf{Z}^\top \mathbf{W}_0 \mathbf{Z}^\top) = 0, \quad \|\mathbf{W}_0 \mathbf{Z}^\top\|_F^2 = \mathbf{B}_{22}, \quad \langle \mathbf{W}_0 \mathbf{W}_0^\top - \mathbf{B}, \mathbf{Z} \mathbf{Z}^\top \rangle = \mathbf{B}_{11},$$

we get

$$\langle \mathbf{Z}, \nabla_{\mathbf{W}}^2 L(\mathbf{W}_0, c_0) \mathbf{Z} \rangle = -4(\mathbf{B}_{11} - \mathbf{B}_{22}) \leq -4\delta^{\text{sep}}.$$

This proves the proposition.  $\square$

We can now prove Proposition 1.

*Proof of Proposition 1.* First, remark that  $L(\mathbf{W}, c)$  has compact sub-level sets. The proposition then follows from Proposition 4 and the continuity of the gradient  $\nabla L(\mathbf{x})$  and of the minimum eigenvalue of the Hessian  $\lambda_{\min}(\nabla^2 L(\mathbf{x}))$ .  $\square$

### B.3.3 Dynamics

The following lemma is a standard combination of Lojasiewicz inequality and center and stable manifold theorem. We prove it for completeness.

**Lemma 8.** *Let  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  be an analytic function that has compact level sets. Consider the gradient flow*

$$\dot{\mathbf{x}}_t = -\nabla f(\mathbf{x}_t).$$

*Then for (Lebesgue) almost all initialization  $\mathbf{x}_0$ , there exists a second order local minimizer  $\mathbf{x}_*$ , such that*

$$\lim_{t \rightarrow +\infty} \mathbf{x}_t = \mathbf{x}_*.$$

*Proof of Lemma 8.*

**Step 1. Show convergence to a critical point.** Since  $f$  is an analytic function, by Lojasiewicz inequality [Loj82], and the fact that the level set of  $f$  is compact, we have

$$\lim_{t \rightarrow +\infty} \mathbf{x}_t = \mathbf{x}_*$$

for  $\mathbf{x}_*$  some critical point of  $f$ .

**Step 2. Show convergence to a local minimizer.** In this step, we proceed similarly to the proof of Theorem 3 in [PP16]. First, consider a sublevel set

$$\Omega(K) = \{\mathbf{x} : f(\mathbf{x}) \leq K\}.$$

Then we have  $\Omega(K)$  compact. Since  $f$  is an analytic function,  $\nabla f$  is Lipschitz in the compact set  $\Omega(K)$ . We define the map  $\phi_t : \Omega(K) \rightarrow \phi_t(\Omega(K))$ ,  $\mathbf{x} \mapsto \mathbf{x}_t$  where  $\mathbf{x}_t$  is defined as the solution of

$$\begin{aligned}\dot{\mathbf{x}}_t &= -\nabla f(\mathbf{x}_t), \\ \mathbf{x}_0 &= \mathbf{x}.\end{aligned}$$

By Picard's existence and uniqueness theorem, we have  $\phi_t$  is a diffeomorphism from  $\Omega(K)$  to  $\phi(\Omega(K))$  for any  $t > 0$ . Fix an  $\varepsilon_0 > 0$ , and we define  $g = \phi_{\varepsilon_0} : \Omega(K) \rightarrow \Omega(K)$ .

Let  $\mathbf{r}$  be a strict saddle point of  $f$ , then  $\mathbf{r}$  must be an unstable fixed point of the diffeomorphism  $g = \phi_{\varepsilon_0}$ . By center and stable manifold theorem (such as Theorem 9 in [PP16]), there exists a manifold  $W_{\text{loc}}^{\text{sc}}(\mathbf{r})$  of dimension at most  $d-1$ , and a ball  $B(\mathbf{r}, \varepsilon(\mathbf{r}))$  centered at  $\mathbf{r}$  with radius  $\varepsilon(\mathbf{r})$ , such that we have the following facts:

- (1)  $g(W_{\text{loc}}^{\text{sc}}(\mathbf{r}) \cap B(\mathbf{r}, \varepsilon(\mathbf{r}))) \subseteq W_{\text{loc}}^{\text{sc}}(\mathbf{r})$ ;
- (2) If  $g^n(\mathbf{x}) \in B(\mathbf{r}, \varepsilon(\mathbf{r}))$  for all  $n \geq 0$ , we have  $\mathbf{x} \in W_{\text{loc}}^{\text{sc}}(\mathbf{r})$  (here  $g^n$  means composition of  $g$  for  $n$  times).

We consider the union of the balls associated to all the strict saddle points of  $f$  in  $\Omega(K)$

$$A = \bigcup_{\mathbf{r} \in \Omega(K) : \mathbf{r} \text{ strict saddle}} B(\mathbf{r}, \varepsilon(\mathbf{r})).$$

Due to Lindelof's lemma, we can find a countable subcover for  $A$ , i.e., there exists fixed-points  $\mathbf{r}_1, \mathbf{r}_2, \dots$  such that  $A = \bigcup_{m=1}^{\infty} B(\mathbf{r}_m, \varepsilon(\mathbf{r}_m))$ . If gradient descent converges to a strict saddle point, starting from a point  $\mathbf{v} \in \Omega(K)$ , there must exist a  $t_0$  and  $m$  such that  $\phi_t(\mathbf{v}) \in B(\mathbf{r}_m, \varepsilon(\mathbf{r}_m))$  for all  $t \geq t_0$ . By center and stable manifold theorem, we get that  $\phi_t(\mathbf{v}) \in W_{\text{loc}}^{\text{sc}}(\mathbf{r}_m) \cap \Omega(K)$ . By setting  $D_1(\mathbf{r}_m) = g^{-1}(W_{\text{loc}}^{\text{sc}}(\mathbf{r}_m) \cap \Omega(K))$  and  $D_{i+1}(\mathbf{r}_m) = g^{-1}(D_i(\mathbf{r}_m) \cap \Omega(K))$  we get that  $\mathbf{v} \in D_k(\mathbf{r}_m)$  for all  $k \geq t_0$ . Hence the set of initial points in  $\Omega(K)$  such that gradient descent converges to a strict saddle point is a subset of

$$P = \bigcup_{m=1}^{\infty} \bigcup_{k \in \mathbb{N}} D_k(\mathbf{r}_m).$$

Note that the set  $W_{\text{loc}}^{\text{sc}}(\mathbf{r}_m) \cap \Omega(K)$  has Lebesgue measure zero in  $\mathbb{R}^d$ . Since  $g$  is a diffeomorphism,  $g^{-1}$  is continuously differentiable and thus it is locally Lipschitz. Therefore,  $g^{-1}$  preserves the null-sets and hence (by induction)  $D_i(\mathbf{r}_m)$  has measure zero for all  $i$ . Thereby we get that  $P$  is a countable union of measure zero sets. Hence  $P$  has measure 0.

Finally, note we have

$$\{\mathbf{x} \in \Omega(K) : \exists \mathbf{r}, \mathbf{r} \text{ is strict saddle}, \mathbf{r} = \lim_{t \rightarrow +\infty} \phi_t(\mathbf{x})\} \subseteq P.$$

Since  $P$  has measure 0, we have

$$\begin{aligned}& \{\mathbf{x} \in \mathbb{R}^d : \exists \mathbf{r}, \mathbf{r} \text{ is strict saddle}, \mathbf{r} = \lim_{t \rightarrow +\infty} \phi_t(\mathbf{x})\} \\ &= \bigcup_{K \in \mathbb{N}} \{\mathbf{x} \in \Omega(K) : \exists \mathbf{r}, \mathbf{r} \text{ is strict saddle}, \mathbf{r} = \lim_{t \rightarrow +\infty} \phi_t(\mathbf{x})\}\end{aligned}$$

has measure 0. This proves the lemma.  $\square$

The following lemma is standard, and a corollary of Theorem 2.11 in [Kur70].

**Lemma 9.** *Let*

$$F(\mathbf{x}) = \mathbb{E}_{\mathbf{z}}[f(\mathbf{x}; \mathbf{z})]$$

*be a  $C^2$  function on  $\Omega \subseteq \mathbb{R}^d$ . Assume*

$$\begin{aligned} \sup_{\mathbf{x} \in \Omega} \mathbb{E}_{\mathbf{z}}[\|\nabla_{\mathbf{x}} f(\mathbf{x}; \mathbf{z})\|_2] &< \infty, \\ \sup_{\mathbf{x} \in \Omega} \|\nabla^2 F(\mathbf{x})\|_{\text{op}} &< \infty. \end{aligned}$$

*Let  $\mathbf{x}_t$  be the trajectory of*

$$\dot{\mathbf{x}}_t = -\nabla F(\mathbf{x}_t),$$

*with initialization  $\mathbf{x}_0 \in \Omega$ . Further assume that there exists  $\eta > 0$ , such that  $\cup_{t \geq 0} \mathcal{B}(\mathbf{x}_t, \eta) \subseteq \Omega$ .*

*Consider the following Markov jump process  $\mathbf{x}_{t,\varepsilon}$  starting from  $\mathbf{x}_0$ , with jump time to be an exponential random variable with fixed mean  $\varepsilon$ , and jump direction  $-\varepsilon \nabla f(\mathbf{x}; \mathbf{z})$  where  $\mathbf{x}$  is the current state, and  $\mathbf{z}$  an independent sample. Then we have for any fixed  $T > 0$  and  $\delta > 0$ ,*

$$\lim_{\varepsilon \rightarrow 0^+} \mathbb{P}\left(\sup_{0 \leq t \leq T} \|\mathbf{x}_t - \mathbf{x}_{t,\varepsilon}\|_2 \geq \delta\right) = 0.$$

### B.3.4 Proof of Theorem 3

By Proposition 4, we know that for  $L(\mathbf{W}, c)$ , any critical point that is not a global minimizer is a strict saddle point. Consider the gradient flow

$$\frac{d}{dt}(\mathbf{W}_t, c_t) = -\nabla L(\mathbf{W}_t, c_t)$$

with random initialization  $(\mathbf{W}_0, c_0) \sim \nu_0$  where  $\nu_0$  is a distribution that is absolutely continuous with respect to Lebesgue measure. Since  $L(\mathbf{W}, c)$  is an analytic function, by Lemma 8, we have  $(\mathbf{W}_t, c_t)$  converges to a global minimizer of  $L(\mathbf{W}, c)$ . That is, we have almost surely (over  $\nu_0$ )

$$\lim_{t \rightarrow \infty} L(\mathbf{W}_t, c_t) = \inf_{\mathbf{W}, c} L(\mathbf{W}, c),$$

where  $\inf_{\mathbf{W}, c} L(\mathbf{W}, c)$  is calculated in Lemma 7.

Consider the following Markov jump process  $(\mathbf{W}_{t,\varepsilon}, c_{t,\varepsilon})$  starting from  $(\mathbf{W}_0, c_0) \sim \nu_0$ , with jump time to be an exponential random variable with fixed mean  $\varepsilon$ , and jump direction to be  $-\varepsilon \nabla L(\mathbf{W}, c; \mathbf{z})$  where

$$\nabla L(\mathbf{W}, c; \mathbf{z}) = \begin{pmatrix} \nabla_{\mathbf{W}} L(\mathbf{W}, c; \mathbf{z}) \\ \partial_c L(\mathbf{W}, c; \mathbf{z}) \end{pmatrix} = \begin{pmatrix} 2(c - b_0 + \langle \mathbf{z}\mathbf{z}^\top, \mathbf{W}\mathbf{W}^\top - \mathbf{B} \rangle) \mathbf{z}\mathbf{z}^\top \mathbf{W} \\ 2(c - b_0 + \langle \mathbf{z}\mathbf{z}^\top, \mathbf{W}\mathbf{W}^\top - \mathbf{B} \rangle) \end{pmatrix}$$

with  $(\mathbf{W}, c)$  the current state, and  $\mathbf{z}$  an independent sample. By Lemma 9, we have for any fixed  $T > 0$  and  $\delta > 0$ ,

$$\lim_{\varepsilon \rightarrow 0^+} \mathbb{P}\left(\sup_{0 \leq t \leq T} \|(\mathbf{W}_{t,\varepsilon}, c_{t,\varepsilon}) - (\mathbf{W}_t, c_t)\|_2 \geq \delta\right) = 0.$$

Note the sequence of Markov jump process at jump time is exactly the SGD iterates. Hence the SGD iterates with properly scaled number of iterations is uniformly close to  $(\mathbf{W}_t, c_t)$  over finite horizon as  $\varepsilon \rightarrow 0$ . This proves the Theorem.

## C Proofs for Mixture of Gaussians

In this section, we consider the mixture of Gaussian setting (mg):  $y_i = \pm 1$  with equal probability 1/2, and  $\mathbf{x}_i | y_i = +1 \sim \mathcal{N}(0, \Sigma^{(1)})$ ,  $\mathbf{x}_i | y_i = -1 \sim \mathcal{N}(0, \Sigma^{(2)})$  where  $\Sigma^{(1)} = \Sigma - \Delta$  and  $\Sigma^{(2)} = \Sigma + \Delta$ . With these notations,

$$\begin{aligned} \Sigma &= \frac{1}{2}(\Sigma^{(1)} + \Sigma^{(2)}), \\ \Delta &= \frac{1}{2}(\Sigma^{(2)} - \Sigma^{(1)}). \end{aligned}$$

Throughout this section, we will make the following assumptions:

**M1.** There exists constants  $0 < c_1 < c_2$  such that  $c_1 \mathbf{I}_d \preceq \Sigma \preceq c_2 \mathbf{I}_d$ ;

**M2.**  $\|\Delta\|_{\text{op}} = \Theta_d(1/\sqrt{d})$ .

Throughout this section, we will denote  $\mathbb{P}_{\Sigma, \Delta}$  the joint distribution of  $(y, \mathbf{x})$  under the **mg** model,  $\mathbb{E}_{\mathbf{x}, y}$  the expectation operator with respect to  $(y, \mathbf{x}) \sim \mathbb{P}_{\Sigma, \Delta}$  and  $\mathbb{E}_{\mathbf{x}}$  the expectation operator with respect to the marginal distribution  $\mathbf{x} \sim (1/2) \cdot \mathcal{N}(0, \Sigma^{(1)}) + (1/2) \cdot \mathcal{N}(0, \Sigma^{(2)})$ .

## C.1 Random Features model: proof of Theorem 4

Recall the definition

$$R_{\text{RF}, N}(\mathbb{P}) = \min_{\hat{f} \in \mathcal{F}_{\text{RF}, N}(\mathbf{W})} \mathbb{E}\{(y - \hat{f}(\mathbf{x}))^2\},$$

where

$$\mathcal{F}_{\text{RF}, N}(\mathbf{W}) = \left\{ f_N(\mathbf{x}) = \sum_{i=1}^N a_i \sigma(\langle \mathbf{w}_i, \mathbf{x} \rangle) : a_i \in \mathbb{R}, i \in [N] \right\}.$$

Note that it is easy to see from the proof that the result stays the same if we add an offset  $c$ .

**Remark 2.** We will state the lemmas for the case  $\Sigma = \mathbf{I}_d$ , which amounts to re-scaling  $\tilde{\Gamma} = \Sigma^{1/2} \Gamma \Sigma^{1/2}$  and  $\tilde{\Delta} = \Sigma^{-1/2} \Delta \Sigma^{-1/2}$ .

### C.1.1 Representation of the RF risk

**Lemma 10.** Consider the RF model introduced above. We have

$$R_{\text{RF}, N}(\mathbb{P}_{\mathbf{I}, \Delta}) = \mathbb{E}_{\mathbf{x}, y}[y^2] - \mathbf{V}^\top \mathbf{U}^{-1} \mathbf{V}, \quad (74)$$

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

$$\begin{aligned} V_i &= \mathbb{E}_{\mathbf{x}, y}[y \sigma(\langle \mathbf{w}_i, \mathbf{x} \rangle)], \\ U_{ij} &= \mathbb{E}_{\mathbf{x}, y}[\sigma(\langle \mathbf{w}_i, \mathbf{x} \rangle) \sigma(\langle \mathbf{w}_j, \mathbf{x} \rangle)]. \end{aligned}$$

*Proof.* Simply write the KKT conditions. The optimum is achieved at  $\mathbf{a} = \mathbf{U}^{-1} \mathbf{V}$ .  $\square$

### C.1.2 Approximation of kernel matrix $\mathbf{U}$

**Lemma 11.** Let  $\sigma \in L^2(\mathcal{N}(0, 1))$  be an activation function. Denote  $\lambda_k = \mathbb{E}_{G \sim \mathcal{N}(0, 1)}[\sigma(G) \text{He}_k(G)]$  the  $k$ -th Hermite coefficient of  $\sigma$  and assume  $\lambda_0 = 0$ . Let  $\mathbf{U} = (U_{ij})_{i, j \in [N]}$  be a random matrix with

$$U_{ij} = \mathbb{E}_{\mathbf{x}}[\sigma(\langle \mathbf{w}_i, \mathbf{x} \rangle) \sigma(\langle \mathbf{w}_j, \mathbf{x} \rangle)],$$

where  $(\mathbf{w}_i)_{i \in [N]} \sim \mathcal{N}(\mathbf{0}, \Gamma)$  independently. Assume conditions **A1** and **B2** hold.

Define  $\mathbf{W} = (\mathbf{w}_1, \dots, \mathbf{w}_N) \in \mathbb{R}^{d \times N}$ , and  $\mathbf{U}_0 = \{(U_0)_{ij}\}_{i, j \in [N]}$ , with

$$(U_0)_{ij} = \tilde{\lambda} \delta_{ij} + \lambda_1^2 \langle \mathbf{w}_i, \mathbf{w}_j \rangle + \kappa/d + \mu_i \mu_j,$$

where

$$\begin{aligned} \mu_i &= \lambda_2(\|\mathbf{w}_i\|_2^2 - 1)/2, \\ \tilde{\lambda} &= \mathbb{E}[\sigma(G)^2] - \lambda_1^2, \\ \kappa &= d \cdot \lambda_2^2[\text{Tr}(\Gamma^2)/2 + \text{Tr}(\Delta \Gamma)^2/4]. \end{aligned}$$

Then we have as  $N/d = \rho$  and  $d \rightarrow \infty$ , we have

$$\|\mathbf{U} - \mathbf{U}_0\|_{\text{op}} = o_d(\mathbb{P}(1)).$$

*Proof of Lemma 11.* Recalling that in the (mg) model, we have  $\mathbf{x} \sim (1/2) \cdot \mathbf{N}(\mathbf{0}, \mathbf{I} - \mathbf{\Delta}) + (1/2) \cdot \mathbf{N}(\mathbf{0}, \mathbf{I} + \mathbf{\Delta})$ , we have

$$\begin{aligned} U_{ij} &= \mathbb{E}_{\mathbf{x}}[\sigma(\langle \mathbf{w}_i, \mathbf{x} \rangle) \sigma(\langle \mathbf{w}_j, \mathbf{x} \rangle)] \\ &= \left\{ \mathbb{E}_{\mathbf{x} \sim \mathbf{N}(\mathbf{0}, \mathbf{I})}[\sigma(\langle (\mathbf{I} - \mathbf{\Delta})^{1/2} \mathbf{w}_i, \mathbf{x} \rangle) \sigma(\langle (\mathbf{I} - \mathbf{\Delta})^{1/2} \mathbf{w}_j, \mathbf{x} \rangle)] \right. \\ &\quad \left. + \mathbb{E}_{\mathbf{x} \sim \mathbf{N}(\mathbf{0}, \mathbf{I})}[\sigma(\langle (\mathbf{I} + \mathbf{\Delta})^{1/2} \mathbf{w}_i, \mathbf{x} \rangle) \sigma(\langle (\mathbf{I} + \mathbf{\Delta})^{1/2} \mathbf{w}_j, \mathbf{x} \rangle)] \right\} / 2. \end{aligned}$$

We can therefore readily use the result of Lemma 2 for  $\tilde{\mathbf{w}}_i \sim \mathbf{N}(\mathbf{0}, (\mathbf{I} - \mathbf{\Delta})^{1/2} \mathbf{\Gamma} (\mathbf{I} - \mathbf{\Delta})^{1/2})$  and  $\tilde{\mathbf{w}}_i \sim \mathbf{N}(\mathbf{0}, (\mathbf{I} + \mathbf{\Delta})^{1/2} \mathbf{\Gamma} (\mathbf{I} + \mathbf{\Delta})^{1/2})$ , to get

$$\|\mathbf{U} - \tilde{\mathbf{U}}_0\|_{\text{op}} = o_{d, \mathbb{P}}(1), \quad (75)$$

where  $\tilde{\mathbf{U}}_0 = (\tilde{U}_0)_{i,j \in [N]}$  with

$$(\tilde{U}_0)_{ij} = \tilde{\lambda} \delta_{ij} + \lambda_1^2 \langle \mathbf{w}_i, \mathbf{w}_j \rangle + \kappa/d + (\mu_i^+ \mu_j^+ + \mu_i^- \mu_j^-)/2,$$

and

$$\begin{aligned} \tilde{\lambda} &= \mathbb{E}[\sigma(G)^2] - \lambda_1^2, \\ \tilde{\kappa} &= d\lambda_2^2 [\text{Tr}((\mathbf{I} - \mathbf{\Delta})\mathbf{\Gamma}(\mathbf{I} - \mathbf{\Delta})\mathbf{\Gamma}) + \text{Tr}((\mathbf{I} + \mathbf{\Delta})\mathbf{\Gamma}(\mathbf{I} + \mathbf{\Delta})\mathbf{\Gamma})] / 4 \\ &= d\lambda_2^2 [\text{Tr}(\mathbf{\Gamma}^2) + \text{Tr}(\mathbf{\Delta}\mathbf{\Gamma}\mathbf{\Delta}\mathbf{\Gamma})] / 2, \\ \mu_i^+ &= \lambda_2(\|(\mathbf{I} + \mathbf{\Delta})^{1/2} \mathbf{w}_i\|_2^2 - 1)/2, \\ \mu_i^- &= \lambda_2(\|(\mathbf{I} - \mathbf{\Delta})^{1/2} \mathbf{w}_i\|_2^2 - 1)/2. \end{aligned}$$

Note that we have

$$(\mu_i^+ \mu_j^+ + \mu_i^- \mu_j^-)/2 = \mu_i \mu_j + \lambda_2^2 \langle \mathbf{w}_i, \mathbf{\Delta} \mathbf{w}_i \rangle \langle \mathbf{w}_j, \mathbf{\Delta} \mathbf{w}_j \rangle / 4,$$

where

$$\mu_i = \lambda_2(\|\mathbf{w}_i\|_2^2 - 1)/2.$$

The matrix  $(\langle \mathbf{w}_i, \mathbf{\Delta} \mathbf{w}_i \rangle \langle \mathbf{w}_j, \mathbf{\Delta} \mathbf{w}_j \rangle)_{i,j \in [N]}$  is simply  $\mathbf{s} \mathbf{s}^\top$  with  $\mathbf{s} = (\langle \mathbf{w}_i, \mathbf{\Delta} \mathbf{w}_i \rangle)_{i \in [N]}$ . Defining  $\nu = \mathbb{E}[\langle \mathbf{w}_i, \mathbf{\Delta} \mathbf{w}_i \rangle] = \text{Tr}(\mathbf{\Gamma} \mathbf{\Delta})$ , we have

$$\mathbf{s} \mathbf{s}^\top = (\mathbf{s} - \nu \mathbf{1}) \nu \mathbf{1}^\top + \nu \mathbf{1} (\mathbf{s} - \nu \mathbf{1})^\top + \nu^2 \mathbf{1} \mathbf{1}^\top + (\mathbf{s} - \nu \mathbf{1}) (\mathbf{s} - \nu \mathbf{1})^\top.$$

Furthermore:

$$\|\mathbf{s} - \nu \mathbf{1}\|_2^2 = \sum_{i=1}^d \text{Tr}((\mathbf{w}_i \mathbf{w}_i^\top - \mathbf{\Gamma}) \mathbf{\Delta})^2.$$

Note that by assumptions **M2** and **B2**, we have  $\mathbb{E}[\text{Tr}((\mathbf{w}_i \mathbf{w}_i^\top - \mathbf{\Gamma}) \mathbf{\Delta})^2] = 2\|\mathbf{\Delta} \mathbf{\Gamma}\|_F^2 = o_{d, \mathbb{P}}(d^{-1})$ . We deduce that  $\|\mathbf{s} - \nu \mathbf{1}\|_2 = o_{d, \mathbb{P}}(1)$ , and therefore

$$\begin{aligned} \|(\mathbf{s} - \nu \mathbf{1}) \nu \mathbf{1}^\top\|_{\text{op}} &= o_{d, \mathbb{P}}(1), \\ \|(\mathbf{s} - \nu \mathbf{1}) (\mathbf{s} - \nu \mathbf{1})^\top\|_{\text{op}} &= o_{d, \mathbb{P}}(1). \end{aligned}$$

Hence, we get

$$\|(\boldsymbol{\mu}^+ \boldsymbol{\mu}^{+\top} + \boldsymbol{\mu}^- \boldsymbol{\mu}^{-\top})/2 - \boldsymbol{\mu} \boldsymbol{\mu}^\top - \text{Tr}(\mathbf{\Gamma} \mathbf{\Delta})^2 \mathbf{1} \mathbf{1}^\top\|_{\text{op}} = o_{d, \mathbb{P}}(1). \quad (76)$$

We also have  $\text{Tr}(\mathbf{\Delta} \mathbf{\Gamma} \mathbf{\Delta} \mathbf{\Gamma})^2 = o_d(d^{-1})$  by assumptions **M2** and **B2**, hence

$$\|\text{Tr}(\mathbf{\Delta} \mathbf{\Gamma} \mathbf{\Delta} \mathbf{\Gamma}) \mathbf{1} \mathbf{1}^\top\|_{\text{op}} = o_{d, \mathbb{P}}(1). \quad (77)$$

Therefore, combining (76) and (77), we get:

$$\|\tilde{\mathbf{U}}_0 - \mathbf{U}_0\|_{\text{op}} = o_{d, \mathbb{P}}(1). \quad (78)$$

Combining (75) and (78) concludes the proof.  $\square$

### C.1.3 Approximation of the $\mathbf{V}$ vector

**Lemma 12.** Under the assumption of Theorem 4, define  $\mathbf{V} = (V_1, \dots, V_N)^\top$  with

$$V_i = \mathbb{E}_{\mathbf{x}, y} [y \sigma(\langle \mathbf{w}_i, \mathbf{x} \rangle)]$$

where  $(\mathbf{w}_i)_{i \in [N]} \sim \mathcal{N}(\mathbf{0}, \mathbf{\Gamma})$  independently. Then as  $N/d = \rho$  with  $d \rightarrow \infty$ , we have

$$\|\mathbf{V} - \tau \mathbf{1} / \sqrt{d}\|_2 = o_{d, \mathbb{P}}(1),$$

where

$$\tau = -\sqrt{d} \cdot \lambda_2 \text{Tr}(\mathbf{\Delta} \mathbf{\Gamma}) / 2.$$

*Proof of Lemma 12.* We have

$$\begin{aligned} V_i &= \{\mathbb{E}_{\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{I} - \mathbf{\Delta})} [\sigma(\langle \mathbf{w}_i, \mathbf{x} \rangle)] - \mathbb{E}_{\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{I} + \mathbf{\Delta})} [\sigma(\langle \mathbf{w}_i, \mathbf{x} \rangle)]\} / 2 \\ &= \{\mathbb{E}_{\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} [\sigma(\langle (\mathbf{I} - \mathbf{\Delta})^{1/2} \mathbf{w}_i, \mathbf{x} \rangle)] - \mathbb{E}_{\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} [\sigma(\langle (\mathbf{I} + \mathbf{\Delta})^{1/2} \mathbf{w}_i, \mathbf{x} \rangle)]\} / 2 \\ &= \mathbb{E}_{G \sim \mathcal{N}(0, 1)} [\sigma(\|(\mathbf{I} - \mathbf{\Delta})^{1/2} \mathbf{w}_i\|_2 G) - \sigma(\|(\mathbf{I} + \mathbf{\Delta})^{1/2} \mathbf{w}_i\|_2 G)] / 2. \end{aligned}$$

We define three interpolating variables:

$$\begin{aligned} V_i^{(1)} &= \lambda_2 \{\|(\mathbf{I} - \mathbf{\Delta})^{1/2} \mathbf{w}_i\|_2 - \|(\mathbf{I} + \mathbf{\Delta})^{1/2} \mathbf{w}_i\|_2\} / 2, \\ V_i^{(2)} &= -\lambda_2 \{\text{Tr}(\mathbf{\Delta} \mathbf{w}_i \mathbf{w}_i^\top)\} / 2, \\ V_i^{(3)} &= -\lambda_2 \text{Tr}(\mathbf{\Delta} \mathbf{\Gamma}) / 2. \end{aligned}$$

We begin by bounding the difference between  $\mathbf{V}$  and  $\mathbf{V}^{(1)}$ . For convenience, we will define  $\tilde{\mathbf{w}}_i = (\mathbf{I} - \mathbf{\Delta})^{1/2} \mathbf{w}_i$ . We have:

$$\begin{aligned} &\mathbb{E}[\sigma(\|\tilde{\mathbf{w}}_i\|_2 G) - \sigma(G)] - \lambda_2(\|\tilde{\mathbf{w}}_i\|_2 - 1) \\ &= \mathbb{E} \left[ \frac{\sigma(\|\tilde{\mathbf{w}}_i\|_2 G) - \sigma(G) - (\|\tilde{\mathbf{w}}_i\|_2 - 1) G \sigma'(G)}{(\|\tilde{\mathbf{w}}_i\|_2 - 1)^2} \right] (\|\tilde{\mathbf{w}}_i\|_2 - 1)^2. \end{aligned} \quad (79)$$

Using dominated convergence theorem and arguments similar to those used to prove (26), one can check that

$$\lim_{t \rightarrow 1} \mathbb{E} \left[ \frac{\sigma(tG) - \sigma(G) - (t-1)G \sigma'(G)}{(t-1)^2} \right] = (\lambda_4(\sigma) + \lambda_2(\sigma)) / 2. \quad (80)$$

The same arguments as in the proofs of Lemma 2 and Lemma 3 show

$$\begin{aligned} \sup_{i \in [N]} |\|(\mathbf{I} - \mathbf{\Delta})^{1/2} \mathbf{w}_i\|_2 - 1| &= o_{d, \mathbb{P}}(1), \\ \sum_{i=1}^N (\|(\mathbf{I} - \mathbf{\Delta})^{1/2} \mathbf{w}_i\|_2 - 1)^2 &= O_{d, \mathbb{P}}(1). \end{aligned} \quad (81)$$

Combining (80) with (81) in (79), we get:

$$\begin{aligned} &\sum_{i=1}^N \left( \mathbb{E}[\sigma(\|\tilde{\mathbf{w}}_i\|_2 G) - \sigma(G)] - \lambda_2(\|\tilde{\mathbf{w}}_i\|_2 - 1) \right)^2 \\ &= \sum_{i=1}^N \left( \frac{\mathbb{E}[\sigma(\|\tilde{\mathbf{w}}_i\|_2 G) - \sigma(G)] - \lambda_2(\|\tilde{\mathbf{w}}_i\|_2 - 1)}{(\|\tilde{\mathbf{w}}_i\|_2 - 1)^2} \right) (\|\tilde{\mathbf{w}}_i\|_2 - 1)^4 \\ &= O_{d, \mathbb{P}}(1) \cdot \left( \sup_{i \in [N]} |\|(\mathbf{I} - \mathbf{\Delta})^{1/2} \mathbf{w}_i\|_2 - 1|^2 \right) \sum_{i=1}^N (\|(\mathbf{I} - \mathbf{\Delta})^{1/2} \mathbf{w}_i\|_2 - 1)^2 = o_{d, \mathbb{P}}(1). \end{aligned}$$

Bounding similarly the term depending on  $(\mathbf{I} + \mathbf{\Delta})^{1/2}\mathbf{w}_i$  in  $V_i^{(1)}$ , we get

$$\|\mathbf{V} - \mathbf{V}^{(1)}\|_2 = o_{d,\mathbb{P}}(1). \quad (82)$$

Now, consider the difference between  $\mathbf{V}^{(1)}$  and  $\mathbf{V}^{(2)}$ . We use the fact for  $x$  on a neighborhood of 0, there exists  $c$  such that

$$|\sqrt{1-x} - \sqrt{1+x} + x| \leq c|x|^3.$$

Hence, with high probability

$$| \|(\mathbf{I} - \mathbf{\Delta})^{1/2}\mathbf{w}_i\|_2 - \|(1 + \mathbf{\Delta})^{1/2}\mathbf{w}_i\|_2 + \langle \mathbf{w}_i, \mathbf{\Delta}\mathbf{w}_i \rangle | \leq c \frac{|\langle \mathbf{w}_i, \mathbf{\Delta}\mathbf{w}_i \rangle|^3}{\|\mathbf{w}_i\|_2^2}.$$

Furthermore, we have:

$$\begin{aligned} \mathbb{E}_{\mathbf{w}_i \sim \mathcal{N}(\mathbf{0}, \mathbf{\Gamma})} \left[ (\langle \mathbf{w}_i, \mathbf{\Delta}\mathbf{w}_i \rangle)^6 / \|\mathbf{w}_i\|_2^4 \right] &\leq \|\mathbf{\Delta}\|_{\text{op}}^2 \mathbb{E}[(\langle \mathbf{w}_i, \mathbf{\Delta}\mathbf{w}_i \rangle)^4] \\ &\leq C \|\mathbf{\Delta}\|_{\text{op}}^2 (\text{Tr}[\mathbf{\Gamma}^{1/2} \mathbf{\Delta} \mathbf{\Gamma}^{1/2}]^4 + \|\mathbf{\Gamma}^{1/2} \mathbf{\Delta} \mathbf{\Gamma}^{1/2}\|_F^4) = o_d(d^{-1}), \end{aligned}$$

where the last equality is due to assumptions **M2** and **B2**. We conclude that

$$\|\mathbf{V}^{(1)} - \mathbf{V}^{(2)}\|_2 = o_{d,\mathbb{P}}(1). \quad (83)$$

For the last comparison between  $\mathbf{V}^{(2)}$  and  $\mathbf{V}^{(3)}$ , we take the expectation:

$$\begin{aligned} \mathbb{E}_{\mathbf{w}_i \sim \mathcal{N}(\mathbf{0}, \mathbf{\Gamma})} [(\langle \mathbf{w}_i, \mathbf{\Delta}\mathbf{w}_i \rangle - \text{Tr}(\mathbf{\Gamma}\mathbf{\Delta}))^2] &= \mathbb{E}_{\mathbf{g} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} [(\langle \mathbf{g}\mathbf{g}^\top, \mathbf{\Gamma}^{1/2} \mathbf{\Delta} \mathbf{\Gamma}^{1/2} \rangle - \text{Tr}(\mathbf{\Gamma}\mathbf{\Delta}))^2] \\ &= 2\|\mathbf{\Gamma}^{1/2} \mathbf{\Delta} \mathbf{\Gamma}^{1/2}\|_F^2 \\ &\leq 2\|\mathbf{\Gamma}\|_{\text{op}}^2 \|\mathbf{\Delta}\|_F^2 = O_d(d^{-2}). \end{aligned}$$

We get

$$\|\mathbf{V}^{(3)} - \mathbf{V}^{(2)}\|_2 = o_{d,\mathbb{P}}(1). \quad (84)$$

Combining the above three bounds (82), (83) and (84) yields the desired result.  $\square$

#### C.1.4 Proof of Theorem 4

By Lemma 10, the risk has a representation

$$R_{\text{RF},N}(f_*) = 1 - \mathbf{V}^\top \mathbf{U}^{-1} \mathbf{V}.$$

By Lemma 11, we have

$$\|\mathbf{U} - \mathbf{U}_0\|_{\text{op}} = o_{d,\mathbb{P}}(1).$$

By Lemma 12, we have

$$\|\mathbf{V} - \tau \mathbf{1} / \sqrt{d}\|_2 = o_{d,\mathbb{P}}(1),$$

where

$$\tau = -\sqrt{d} \cdot \lambda_2 \text{Tr}(\mathbf{\Delta}\mathbf{\Gamma}) / 2.$$

Hence, we have

$$|\mathbf{V}^\top \mathbf{U}^{-1} \mathbf{V} - \tau^2 \mathbf{1}^\top \mathbf{U}_0^{-1} \mathbf{1} / d| = o_{d,\mathbb{P}}(1).$$

Proposition 2 gives the expression

$$\mathbf{1}^\top \mathbf{U}_0^{-1} \mathbf{1} / d = \psi / (1 + \kappa \psi) + o_{d,\mathbb{P}}(1),$$

where

$$\kappa = d \cdot \lambda_2^2 [\text{Tr}(\mathbf{\Gamma}^2) / 2 + \text{Tr}(\mathbf{\Delta}\mathbf{\Gamma})^2 / 4].$$

Hence we have

$$\mathbf{V}^\top \mathbf{U} \mathbf{V} = \tau^2 \psi / (1 + \kappa \psi) + o_{d,\mathbb{P}}(1).$$

This proves the theorem.

## C.2 Neural Tangent model: proof of Theorem 5

Recall the definition (note  $R_{\text{NT},N}(\mathbb{P})$  is a function of  $\mathbf{W}$ )

$$R_{\text{NT},N}(\mathbb{P}) = \min_{\hat{f} \in \mathcal{F}_{\text{NT},N}(\mathbf{W})} \mathbb{E}\{(y - \hat{f}(\mathbf{x}))^2\},$$

where

$$\mathcal{F}_{\text{NT},N}(\mathbf{W}) = \left\{ f_N(\mathbf{x}) = c + \sum_{i=1}^N \sigma'(\langle \mathbf{w}_i, \mathbf{x} \rangle) \langle \mathbf{a}_i, \mathbf{x} \rangle : c \in \mathbb{R}, \mathbf{a}_i \in \mathbb{R}^d, i \in [N] \right\}.$$

### C.2.1 A representation lemma

**Lemma 13.** *Assume conditions **M1** and **M2** hold. Consider the function*

$$\hat{f}(\mathbf{x}; \mathbf{\Gamma}, a, c) = a \langle \mathbf{\Gamma}, \mathbf{x} \mathbf{x}^\top \rangle + c. \quad (85)$$

*Define the risk function optimized over  $a, c$  while  $\mathbf{\Gamma}$  is fixed*

$$L(\mathbf{\Gamma}) = \inf_{a,c} \mathbb{E}_{\mathbf{x},y} [(y - \hat{f}(\mathbf{x}; \mathbf{\Gamma}, a, c))^2]. \quad (86)$$

*Then we have*

$$\sup_{\mathbf{\Gamma} \succeq 0} \left| L(\mathbf{\Gamma}) - \frac{2}{2 + \langle \mathbf{\Gamma}, \mathbf{\Delta} \rangle^2 / \|\mathbf{\Sigma}^{1/2} \mathbf{\Gamma} \mathbf{\Sigma}^{1/2}\|_F^2} \right| = o_d(1). \quad (87)$$

*Proof of Lemma 13.* Note we have

$$\begin{aligned} L(\mathbf{\Gamma}, a, c) &\equiv \mathbb{E}_{\mathbf{x},y} [(y - \hat{f}(\mathbf{x}; \mathbf{\Gamma}, a, c))^2] \\ &= 1 + c^2 + 2ac \langle \mathbf{\Gamma}, \mathbf{\Sigma} \rangle + 2a \langle \mathbf{\Gamma}, \mathbf{\Delta} \rangle \\ &\quad + a^2 [\langle \mathbf{\Gamma}, \mathbf{\Sigma} \rangle^2 + 2\text{Tr}(\mathbf{\Sigma} \mathbf{\Gamma} \mathbf{\Sigma} \mathbf{\Gamma}) + \langle \mathbf{\Gamma}, \mathbf{\Delta} \rangle^2 + 2\text{Tr}(\mathbf{\Delta} \mathbf{\Gamma} \mathbf{\Delta} \mathbf{\Gamma})]. \end{aligned}$$

Minimizing successively over  $c$  and  $a$ , we get the following formula:

$$L(\mathbf{\Gamma}) \equiv \min_{c,a \in \mathbb{R}} L(\mathbf{\Gamma}, a, c) = \frac{2}{2 + \langle \mathbf{\Gamma}, \mathbf{\Delta} \rangle^2 / [\text{Tr}(\mathbf{\Gamma} \mathbf{\Sigma} \mathbf{\Gamma} \mathbf{\Sigma}) + \text{Tr}(\mathbf{\Gamma} \mathbf{\Delta} \mathbf{\Gamma} \mathbf{\Delta})]}.$$

By Assumptions **M1** and **M2**, we have  $\mathbf{\Sigma} \succeq c \mathbf{I}_d$  and  $\|\mathbf{\Delta}\|_{\text{op}} \leq C/\sqrt{d}$  for some constants  $c$  and  $C$ . We get

$$\frac{\text{Tr}(\mathbf{\Gamma} \mathbf{\Delta} \mathbf{\Gamma} \mathbf{\Delta})}{\text{Tr}(\mathbf{\Gamma} \mathbf{\Sigma} \mathbf{\Gamma} \mathbf{\Sigma})} \leq \frac{C^2}{dc^2}.$$

We deduce that

$$\sup_{\mathbf{\Gamma} \succeq 0} \left| L(\mathbf{\Gamma}) - \frac{2}{2 + \langle \mathbf{\Gamma}, \mathbf{\Delta} \rangle^2 / \|\mathbf{\Sigma}^{1/2} \mathbf{\Gamma} \mathbf{\Sigma}^{1/2}\|_F^2} \right| \leq \left| \frac{1}{1 + C^2/(dc^2)} - 1 \right| = o_d(1).$$

□

### C.2.2 Proof of Theorem 5

We consider the re-scaled matrices  $\tilde{\mathbf{\Gamma}} = \mathbf{\Sigma}^{1/2} \mathbf{\Gamma} \mathbf{\Sigma}^{1/2}$  and  $\tilde{\mathbf{\Delta}} = \mathbf{\Sigma}^{-1/2} \mathbf{\Delta} \mathbf{\Sigma}^{-1/2}$ . We consider the NT model with a squared non-linearity:

$$\hat{f}(\mathbf{W}, \mathbf{A}) = 2 \sum_{i=1}^N \langle \mathbf{w}_i, \mathbf{x} \rangle \langle \mathbf{a}_i, \mathbf{x} \rangle + c = 2 \langle \mathbf{W} \mathbf{A}^\top, \mathbf{x} \mathbf{x}^\top \rangle + c.$$



with  $\mathbf{W} = [\mathbf{w}_1, \dots, \mathbf{w}_N] \in \mathbb{R}^{d \times N}$  and  $\mathbf{A} = [\mathbf{a}_1, \dots, \mathbf{a}_N] \in \mathbb{R}^{d \times N}$ . For  $\mathbf{w}_i \sim \mathcal{N}(\mathbf{0}, \Sigma)$ , we have with probability one  $\text{rank}(\mathbf{W}) = \min(d, N) \equiv r$ . We consider  $\mathbf{W} = \mathbf{P}_1 \mathbf{S} \mathbf{V}^\top$  the SVD decomposition of  $\mathbf{W}$ , with  $\mathbf{P}_1 \in \mathbb{R}^{d \times r}$ ,  $\mathbf{S} \in \mathbb{R}^{r \times r}$  and  $\mathbf{V} \in \mathbb{R}^{N \times r}$ . Define  $\mathbf{G} = \mathbf{S} \mathbf{V}^\top \mathbf{A} \in \mathbb{R}^{r \times d}$ , we obtain almost surely that the minimum over  $\mathbf{A}$  is the same as the minimum over  $\mathbf{G}$ . From Lemma 13, we deduce that almost surely

$$R_{\text{NT}, N}(\mathbb{P}_{\Sigma, \Delta}) = \min_{\mathbf{G} \in \mathbb{R}^{d \times d}} \left\{ \frac{2}{2 + \text{Tr}[(\mathbf{P}_1 \mathbf{G} + \mathbf{G}^\top \mathbf{P}_1^\top) \Delta]^2 / \|\mathbf{P}_1 \mathbf{G} + \mathbf{G}^\top \mathbf{P}_1^\top\|_F^2} \right\} + o_d(1) \quad (88)$$

**Case  $N/d \rightarrow \rho \geq 1$ .** In the case  $N \geq d$ , we can take  $\mathbf{G} = \mathbf{P}_1^\top \tilde{\mathbf{G}} / 2$  and we get almost surely over  $\mathbf{W} \in \mathbb{R}^{d \times N}$

$$R_{\text{NT}, N}(\mathbb{P}_{\Sigma, \Delta}) = \min_{\mathbf{G} \in \mathbb{R}^{d \times d}} \left\{ \frac{2}{2 + \langle \mathbf{G}, \Delta \rangle^2 / \|\mathbf{G}\|_F^2} \right\} + o_d(1) = \frac{2}{2 + \|\Delta\|_F^2} + o_d(1),$$

where the minimizer  $\mathbf{G} = \Delta$  is obtained by Cauchy-Schwarz inequality.

**Case  $N/d \rightarrow \rho < 1$ .** Consider now the case when  $N < d$ . From (88), the optimal  $\mathbf{G}$  is the one maximizing

$$\max_{\mathbf{G} \in \mathbb{R}^{N \times d}} \frac{\text{Tr}[(\mathbf{P}_1 \mathbf{G} + \mathbf{G}^\top \mathbf{P}_1^\top) \Delta]^2}{\|\mathbf{P}_1 \mathbf{G} + \mathbf{G}^\top \mathbf{P}_1^\top\|_F^2},$$

which we rewrite as the following convex problem

$$\max_{\mathbf{G} \in \mathbb{R}^{N \times d}} \text{Tr}[\mathbf{P}_1 \mathbf{G} \Delta], \quad \text{s.t.} \quad \|\mathbf{P}_1 \mathbf{G} + \mathbf{G}^\top \mathbf{P}_1^\top\|_F^2 \leq 1. \quad (89)$$

We define  $\mathbf{P}_2 \in \mathbb{R}^{d \times (d-N)}$  the completion of  $\mathbf{P}_1$  to a full basis  $\mathbf{P} = [\mathbf{P}_1, \mathbf{P}_2] \in \mathbb{R}^{d \times d}$ , and denote  $\mathbf{G}_1 = \mathbf{G} \mathbf{P}_1 \in \mathbb{R}^{N \times N}$  and  $\mathbf{G}_2 = \mathbf{G} \mathbf{P}_2 \in \mathbb{R}^{N \times (d-N)}$ . We can form the Lagrangian of problem (89):

$$\mathcal{L}(\mathbf{G}, \lambda) = \text{Tr}(\mathbf{P}_1 \mathbf{G} \Delta) + \lambda(1 - \|\mathbf{P}_1 \mathbf{G} + \mathbf{G}^\top \mathbf{P}_1^\top\|_F^2).$$

The stationary condition implies:

$$\nabla_{\mathbf{G}} \mathcal{L}(\mathbf{G}, \lambda) = \mathbf{P}_1^\top \Delta - 4\lambda(\mathbf{P}_1^\top \mathbf{G}^\top \mathbf{P}_1^\top + \mathbf{P}_1^\top \mathbf{P}_1 \mathbf{G}) = 0,$$

which yields, using  $\mathbf{P}_1^\top \mathbf{P}_1 = \mathbf{I}_N$ ,

$$\Delta_{12} = 4\lambda \mathbf{G}_2, \quad \Delta_{11} = 4\lambda(\mathbf{G}_1 + \mathbf{G}_1^\top), \quad (90)$$

where  $\Delta_{ij} = \mathbf{P}_i^\top \Delta \mathbf{P}_j$  for  $i, j = 1, 2$ . The constraint reads in the  $\mathbf{P}$  basis

$$\|\mathbf{P}_1 \mathbf{G} + \mathbf{G}^\top \mathbf{P}_1^\top\|_F^2 = \|\mathbf{G}_1 + \mathbf{G}_1^\top\|_F^2 + 2\|\mathbf{G}_2\|_F^2 = 1. \quad (91)$$

Substituting (90) in (91) yields:

$$4\lambda = \sqrt{\|\Delta_{11}\|_F^2 + 2\|\Delta_{12}\|_F^2}. \quad (92)$$

Considering the (unique) symmetric optimizer  $\mathbf{G}_1$  and substituting (92) in (90), we get the minimizer

$$\begin{aligned} \mathbf{G}_1^* &= \frac{1}{8\lambda} \Delta_{11} = \frac{1}{2\sqrt{\|\Delta_{11}\|_F^2 + 2\|\Delta_{12}\|_F^2}} \Delta_{11}, \\ \mathbf{G}_2^* &= \frac{1}{4\lambda} \Delta_{12} = \frac{1}{\sqrt{\|\Delta_{11}\|_F^2 + 2\|\Delta_{12}\|_F^2}} \Delta_{12}. \end{aligned} \quad (93)$$

Let's consider the objective function:

$$\begin{aligned} \text{Tr}(\mathbf{P}_1 \mathbf{G}^* \Delta) &= \text{Tr}(\mathbf{G}_1^* \Delta_{11} + \mathbf{G}_2^* \Delta_{21}) \\ &= \frac{1}{2\sqrt{\|\Delta_{11}\|_F^2 + 2\|\Delta_{12}\|_F^2}} \text{Tr}(\Delta_{11}^2 + 2\Delta_{12} \Delta_{21}) \\ &= \frac{1}{2} \sqrt{\|\Delta_{11}\|_F^2 + 2\|\Delta_{12}\|_F^2} \\ &= \frac{1}{2} \sqrt{\|\Delta\|_F^2 - \|\Delta_{22}\|_F^2}. \end{aligned} \quad (94)$$

Substituting (94) in (88), we then obtain

$$R_{\text{NT},N}(\mathbb{P}_{\Sigma,\Delta}) = \frac{2}{2 + \|\Delta\|_F^2 - \|\Delta_{22}\|_F^2} + o_d(1), \quad (95)$$

where  $\Delta_{22} = P_{W^\perp} \Delta P_{W^\perp}$  with  $P_{W^\perp} = \mathbf{I}_d - W(W^\top W)^{-1}W^\top$  is the random projection along the orthogonal subspace to the columns of  $W$ . From Theorem 2, we know that

$$\mathbb{E}[\|\Delta_{22}\|_F^2] = \|\Delta\|_F^2 \left[ (1-\rho)_+^2 \left( 1 - \frac{\text{Tr}(\Delta)^2}{d\|\Delta\|_F^2} \right) + (1-\rho)_+ \frac{\text{Tr}(\Delta)^2}{d\|\Delta\|_F^2} + o_d(1) \right]. \quad (96)$$

Let  $\mathbb{W}_d^N$  be the Stiefel manifold, i.e. the collection of all the sets of  $N$  orthonormal vectors in  $\mathbb{R}^d$  endowed with the Frobenius distance. In matrix representation, we have

$$\mathbb{W}_d^N = \{P \in \mathbb{R}^{d \times N} : P^\top P = \mathbf{I}_N\}.$$

By Theorem 2.4 in [Led01], the volume measure on  $\mathbb{W}_d^N$  has normal concentration. In particular, denote by  $F : \mathbb{W}_d^N \mapsto \mathbb{R}$ , the function  $F(P) = \|P^\top \Delta P\|_F^2$ . We upper bound the gradient of  $F$ :

$$\|\nabla F(P)\|_F = 4\|\Delta P P^\top \Delta P\|_F \leq 4\|\Delta P P^\top\|_{\text{op}} \|\Delta P\|_F \leq \|\Delta\|_{\text{op}} \|\Delta\|_F \leq C,$$

by assumption **M2** on  $\Delta$ . We deduce that there exists a constant  $c$  (that depends on  $\rho$  and  $C$ ) such that:

$$\mathbb{P}(|F(P) - \mathbb{E}[F(P)]| > t) \leq e^{-cdt^2}.$$

Therefore, we have

$$\mathbb{P}(|\|\Delta_{22}\|_F^2 - \mathbb{E}[\|\Delta_{22}\|_F^2]| > t) \leq e^{-cdt^2}. \quad (97)$$

Using (97) and (95), we deduce the final high probability formula for the risk of the NT model:

$$R_{\text{NT},N}(\mathbb{P}_{\Sigma,\Delta}) = \frac{2}{2 + \|\Delta\|_F^2 - \mathbb{E}[\|\Delta_{22}\|_F^2]} + o_{d,\mathbb{P}}(1).$$

Substituting  $\mathbb{E}[\|\Delta_{22}\|_F^2]$  by its expression (96) concludes the proof.

### C.3 Neural Network model: proof of Theorem 6

Recall the definition

$$R_{\text{NN},N}(\mathbb{P}) = \min_{\hat{f} \in \mathcal{F}_{\text{NN},N}(W)} \mathbb{E}\{(y - \hat{f}(x))^2\},$$

where we consider the function class of two-layers neural networks (with  $N$  neurons) with quadratic activation function and general offset and coefficients

$$\mathcal{F}_{\text{NN},N}(W) = \left\{ f_N(x) = c + \sum_{i=1}^N a_i (\langle w_i, x \rangle)^2 : c, a_i \in \mathbb{R}, i \in [N] \right\}.$$

We define the risk function for a given set of parameters as

$$L(W, a, c) = \mathbb{E}_{x,y}[(y - \hat{f}(x; W, a, c))^2].$$

The risk is optimized over  $(a_i, w_i)_{i \leq N}$  and  $c$ .

*Proof of Theorem 6.* Without loss of generality, we assume  $\Sigma = \mathbf{I}_d$  (it suffices to consider the re-scaled matrices  $\tilde{\Gamma} = \Sigma^{1/2} \Gamma \Sigma^{1/2}$  and  $\tilde{\Delta} = \Sigma^{-1/2} \Delta \Sigma^{-1/2}$ ). We rewrite the neural network function in a compact form:

$$\hat{f}(x; W, a, c) = \sum_{i=1}^N a_i \langle w_i, x \rangle^2 + c = \langle W A W^\top, x x^\top \rangle + c,$$

where  $\mathbf{A} = \text{diag}(\mathbf{a})$ . Define  $\mathbf{\Gamma} = \mathbf{W}\mathbf{A}\mathbf{W}^\top$  and using Eq. (87) in Lemma 13, the minimizer  $\mathbf{\Gamma}^*$  is the solution of

$$\max_{\mathbf{\Gamma} \in \mathcal{S}(\mathbb{R}^{d \times d})} \frac{\langle \mathbf{\Gamma}, \mathbf{\Delta} \rangle^2}{\|\mathbf{\Gamma}\|_F^2}, \quad \text{s.t.} \quad \text{rank}(\mathbf{\Gamma}) \leq \min(N, d) \equiv r.$$

where  $\mathcal{S}(\mathbb{R}^{d \times d})$  is the set of symmetric matrices in  $\mathbb{R}^{d \times d}$ .

Let us denote the eigendecomposition of  $\mathbf{\Gamma}$  by  $\mathbf{\Gamma} = \mathbf{U}\mathbf{S}\mathbf{U}^\top$  with  $\mathbf{U} \in \mathbb{R}^{d \times r}$  and  $\mathbf{S} = \text{diag}(\mathbf{s}) \in \mathbb{R}^{r \times r}$ . We have by Cauchy-Schwartz inequality

$$\frac{\langle \mathbf{\Gamma}, \mathbf{\Delta} \rangle^2}{\|\mathbf{\Gamma}\|_F^2} = \frac{\text{Tr}(\mathbf{S}\mathbf{U}^\top \mathbf{\Delta} \mathbf{U})^2}{\|\mathbf{S}\|_F^2} \leq \|\text{diag}(\mathbf{U}^\top \mathbf{\Delta} \mathbf{U})\|_2^2,$$

with equality if and only if  $\mathbf{S}_* = \text{ddiag}(\mathbf{U}^\top \mathbf{\Delta} \mathbf{U})$  where  $\text{ddiag}(\mathbf{U}^\top \mathbf{\Delta} \mathbf{U})$  is the vector of the diagonal elements of  $\mathbf{U}^\top \mathbf{\Delta} \mathbf{U}$ . Denoting  $\mathcal{D}(\mathbb{R}^{d \times d})$  the set of diagonal matrices in  $\mathbb{R}^{d \times d}$ , we get

$$\max_{\mathbf{S} \in \mathcal{D}(\mathbb{R}^{d \times d})} \frac{\langle \mathbf{U}\mathbf{S}\mathbf{U}^\top, \mathbf{\Delta} \rangle^2}{\|\mathbf{U}\mathbf{S}\mathbf{U}^\top\|_F^2} = \frac{\langle \mathbf{S}_*, \mathbf{U}^\top \mathbf{\Delta} \mathbf{U} \rangle^2}{\|\mathbf{S}_*\|_F^2} = \frac{\|\mathbf{S}_*\|_F^4}{\|\mathbf{S}_*\|_F^2} = \|\mathbf{S}_*\|_F^2.$$

Hence, the problem reduces to finding  $\mathbf{U} \in \mathbb{R}^{d \times r}$  with orthonormal columns which maximizes  $\|\text{ddiag}(\mathbf{U}^\top \mathbf{\Delta} \mathbf{U})\|_F^2$ . The maximizer is easily found as the eigendirections corresponding to the  $r$  largest singular values. We conclude that at the optimum

$$\frac{\langle \mathbf{\Gamma}_*, \mathbf{\Delta} \rangle^2}{\|\mathbf{\Gamma}_*\|_F^2} = \sum_{i=1}^r \lambda_i^2,$$

where the  $\lambda_i$ 's are the singular values of  $\mathbf{\Delta}$  in descending order. Plugging this expression in Eq. (87) concludes the proof.  $\square$

## D Additional Experiments

For the sake of theoretical analysis, we focused on the case of quadratic activations for NT and NN in the main text. However, the phenomena we presented persist (qualitatively) even when other activation functions are used. For example, figures 3 and 4 examine the performance of our models when ReLU non-linearity is used. These experiments suggest that when  $d$  is larger than  $N$  there is a significant performance gap between NN and NT. Moreover, similar to what was presented in the paper, we observe that the gap between RF( $I$ ) and NN does not vanish unless  $\frac{N}{d} \rightarrow \infty$ .

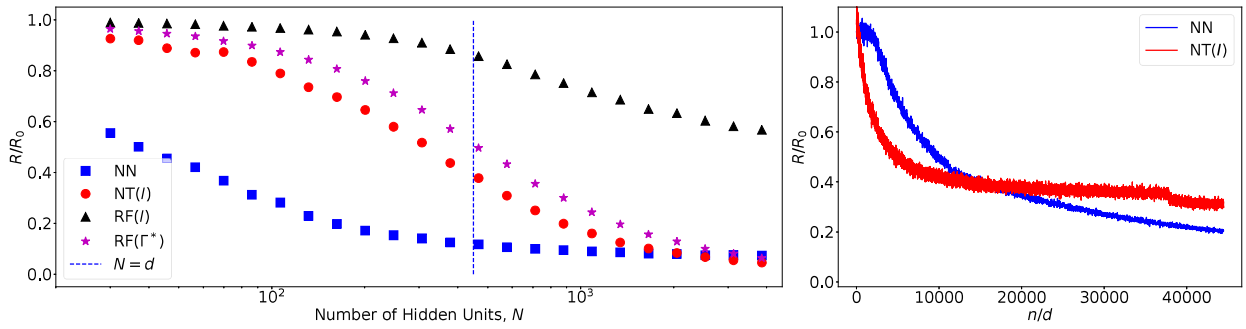


Figure 3: Left frame: Prediction (test) error of a two-layer neural networks in fitting a quadratic function in  $d = 450$  dimensions, as a function of the number of neurons  $N$ . We consider the large sample (population) limit  $n \rightarrow \infty$  and compare three training regimes: random features (RF), neural tangent (NT), and fully trained neural networks (NN). All models use **ReLU** activations. Right frame: Evolution of the risk for NT and NN with the number of samples.

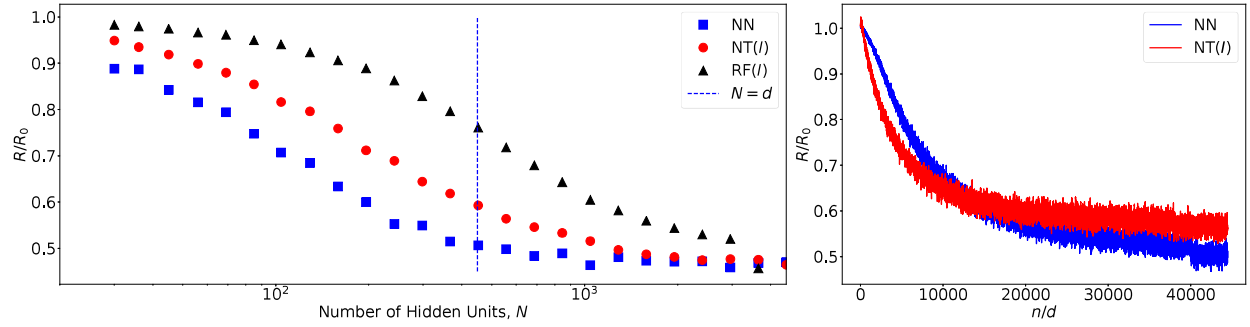


Figure 4: Left frame: Prediction (test) error of a two-layer neural networks in fitting a mixture of Gaussians in  $d = 450$  dimensions, as a function of the number of neurons  $N$ . We consider the large sample (population) limit  $n \rightarrow \infty$  and compare three training regimes: random features (RF), neural tangent (NT), and fully trained neural networks (NN). All models use **ReLU** activations. Right frame: Evolution of the risk for NT and NN with the number of samples.