

## 5 Appendix

Let  $\angle(\mathbf{h}, \mathbf{h}_0) = \bar{\theta}_0^{(1)}$  and  $\angle(\mathbf{m}, \mathbf{m}_0) = \bar{\theta}_0^{(2)}$  for non-zero  $\mathbf{h}, \mathbf{h}_0 \in \mathbb{R}^n$  and  $\mathbf{m}, \mathbf{m}_0 \in \mathbb{R}^p$ . In order to understand how the operators  $\mathbf{h} \rightarrow \mathbf{W}_{+, \mathbf{h}}^{(1)} \mathbf{h}$  and  $\mathbf{m} \rightarrow \mathbf{W}_{+, \mathbf{m}}^{(2)} \mathbf{m}$  distort angles, we define

$$g(\theta) = \cos^{-1} \left( \frac{(\pi - \theta) \cos \theta + \sin \theta}{\pi} \right). \quad (8)$$

Also, for a fixed  $\mathbf{p}, \mathbf{q} \in \mathbb{R}^n$ , define

$$\tilde{\mathbf{t}}_{\mathbf{p}, \mathbf{q}}^{(k)} := \frac{1}{2^{a^{(k)}}} \left[ \left( \prod_{i=0}^{a^{(k)}-1} \frac{\pi - \bar{\theta}_i^{(k)}}{\pi} \right) \mathbf{q} + \sum_{i=0}^{a^{(k)}-1} \frac{\sin \bar{\theta}_i^{(k)}}{\pi} \left( \prod_{j=i+1}^{a^{(k)}-1} \frac{\pi - \bar{\theta}_j^{(k)}}{\pi} \right) \frac{\|\mathbf{q}\|_2}{\|\mathbf{p}\|_2} \mathbf{p} \right], \quad (9)$$

where  $\bar{\theta}_i^{(k)} = g(\bar{\theta}_{i-1}^{(k)})$  for  $g$  given by (8),  $\bar{\theta}_0^{(k)} = \angle(\mathbf{p}, \mathbf{q})$ ,  $a^{(1)} = d$ , and  $a^{(2)} = s$ .

### 5.1 Proof of Deterministic Theorem

**Theorem 4** (Also Theorem 3). *Fix  $\epsilon > 0$ ,  $0 < \alpha_1 \leq 1$  and  $0 < \alpha_2 \leq 1$  such that  $K_1(d^7 s^2 + d^2 s^7) \epsilon^{1/4} / (\alpha_1 \alpha_2) < 1$ ,  $d \geq 2$ , and  $s \geq 2$ . Let  $\mathcal{K} = \{(\mathbf{h}, \mathbf{0}) \in \mathbb{R}^{n \times p} | \mathbf{h} \in \mathbb{R}^n\} \cup \{(\mathbf{0}, \mathbf{m}) \in \mathbb{R}^{n \times p} | \mathbf{m} \in \mathbb{R}^p\}$ . Suppose that  $\mathbf{W}_i^{(1)} \in \mathbb{R}^{n_i \times n_{i-1}}$  for  $i = 1, \dots, d-1$  and  $\mathbf{W}_i^{(2)} \in \mathbb{R}^{p_i \times p_{i-1}}$  for  $i = 1, \dots, s-1$  satisfy the WDC with constant  $\epsilon$  and 1. Suppose  $\mathbf{W}_d^{(1)} \in \mathbb{R}^{\ell \times n_{d-1}}$  satisfy WDC with constants  $\epsilon$  and  $\alpha_1$ , and  $\mathbf{W}_s^{(2)} \in \mathbb{R}^{\ell \times p_{s-1}}$  satisfy WDC with constants  $\epsilon$  and  $\alpha_2$ . Also, suppose  $(\mathbf{W}_d^{(1)}, \mathbf{W}_s^{(2)})$  satisfy joint-WDC with constants  $\epsilon, \alpha = \alpha_1 \cdot \alpha_2$ . Let  $\mathcal{K} = \{(\mathbf{h}, \mathbf{0}) \in \mathbb{R}^{n \times p} | \mathbf{h} \in \mathbb{R}^n\} \cup \{(\mathbf{0}, \mathbf{m}) \in \mathbb{R}^{n \times p} | \mathbf{m} \in \mathbb{R}^p\}$  and  $\mathcal{A} = \mathcal{A}_{K_2 d^3 s^3 \epsilon^{\frac{1}{4}} \alpha^{-1}, (\mathbf{h}_0, \mathbf{m}_0)} \cup \mathcal{A}_{K_2 d^8 s^3 \epsilon^{\frac{1}{4}} \alpha^{-1}, (-\rho_d^{(1)} \mathbf{h}_0, \mathbf{m}_0)} \cup \mathcal{A}_{K_2 d^3 s^8 \epsilon^{\frac{1}{4}} \alpha^{-1}, (\rho_s^{(2)} \mathbf{h}_0, -\mathbf{m}_0)} \cup \mathcal{A}_{K_2 d^8 s^8 \epsilon^{\frac{1}{4}} \alpha^{-1}, (-\rho_d^{(1)} \rho_s^{(2)} \mathbf{h}_0, -\mathbf{m}_0)}$ . Then, for  $(\mathbf{h}_0, \mathbf{m}_0) \neq (\mathbf{0}, \mathbf{0})$ , and*

$$(\mathbf{h}, \mathbf{m}) \notin \mathcal{A} \cup \mathcal{K}$$

the one-sided directional derivative of  $f$  in the direction of  $\mathbf{g} = \mathbf{g}_{1, (\mathbf{h}, \mathbf{m})}$  or  $\mathbf{g} = \mathbf{g}_{2, (\mathbf{h}, \mathbf{m})}$  satisfy  $D_{-\mathbf{g}} f(\mathbf{h}, \mathbf{m}) < 0$ . Additionally, for all  $(\mathbf{h}, \mathbf{m}) \in \mathcal{K}$  and for all  $(\mathbf{x}, \mathbf{y})$

$$D_{(\mathbf{x}, \mathbf{y})} f(\mathbf{h}, \mathbf{m}) \leq 0.$$

Here,  $\rho_d^{(k)}$  are positive numbers that converge to 1 as  $d \rightarrow \infty$ , and  $K_1$ , and  $K_2$  are absolute constants.

*Proof.* Recall that

$$\begin{aligned} \mathbf{v}_{(\mathbf{h}, \mathbf{m}), (\mathbf{h}_0, \mathbf{m}_0)}^{(1)} &= \begin{cases} \nabla_{\mathbf{h}} f(\mathbf{h}, \mathbf{m}) & G \text{ is differentiable at } (\mathbf{h}, \mathbf{m}), \\ \lim_{\delta \rightarrow 0^+} \nabla_{\mathbf{h}} f((\mathbf{h}, \mathbf{m}) + \delta \mathbf{w}) & \text{otherwise} \end{cases} \\ \mathbf{v}_{(\mathbf{h}, \mathbf{m}), (\mathbf{h}_0, \mathbf{m}_0)}^{(2)} &= \begin{cases} \nabla_{\mathbf{m}} f(\mathbf{h}, \mathbf{m}) & G \text{ is differentiable at } (\mathbf{h}, \mathbf{m}), \\ \lim_{\delta \rightarrow 0^+} \nabla_{\mathbf{m}} f((\mathbf{h}, \mathbf{m}) + \delta \mathbf{w}) & \text{otherwise} \end{cases} \end{aligned}$$

where  $G(\mathbf{h}, \mathbf{m})$  is differentiable at  $(\mathbf{h}, \mathbf{m}) + \delta \mathbf{w}$  for sufficiently small  $\delta$ . Such a  $\delta$  exists because of piecewise linearity of  $G(\mathbf{h}, \mathbf{m})$  and any such  $\mathbf{w}$  can be arbitrarily selected. Also, recall that for any differentiable point  $(\mathbf{h}, \mathbf{m})$ , we have

$$\begin{aligned} \nabla_{\mathbf{h}} f(\mathbf{h}, \mathbf{m}) &= \left( \Lambda_{d,+, \mathbf{h}}^{(1)} \right)^{\top} \operatorname{diag} \left( \Lambda_{s,+, \mathbf{m}}^{(2)} \mathbf{m} \right)^2 \Lambda_{d,+, \mathbf{h}}^{(1)} \mathbf{h} \\ &\quad - \left( \Lambda_{d,+, \mathbf{h}}^{(1)} \right)^{\top} \operatorname{diag} \left( \Lambda_{s,+, \mathbf{m}}^{(2)} \mathbf{m} \odot \Lambda_{s,+, \mathbf{m}_0}^{(2)} \mathbf{m}_0 \right) \Lambda_{d,+, \mathbf{h}_0}^{(1)} \mathbf{h}_0, \\ \nabla_{\mathbf{m}} f(\mathbf{h}, \mathbf{m}) &= \left( \Lambda_{s,+, \mathbf{m}}^{(2)} \right)^{\top} \operatorname{diag} \left( \Lambda_{d,+, \mathbf{h}}^{(1)} \mathbf{h} \right)^2 \Lambda_{s,+, \mathbf{m}}^{(2)} \mathbf{m} \\ &\quad - \left( \Lambda_{s,+, \mathbf{m}}^{(2)} \right)^{\top} \operatorname{diag} \left( \Lambda_{d,+, \mathbf{h}}^{(1)} \mathbf{h} \odot \Lambda_{d,+, \mathbf{h}_0}^{(1)} \mathbf{h}_0 \right) \Lambda_{s,+, \mathbf{m}_0}^{(2)} \mathbf{m}_0. \end{aligned}$$

Let

$$\begin{aligned}\mathbf{g}_{1,(\mathbf{h},\mathbf{m})} &= \begin{bmatrix} \mathbf{v}_{(\mathbf{h},\mathbf{m}),(\mathbf{h}_0,\mathbf{m}_0)}^{(1)} \\ \mathbf{0} \end{bmatrix} \in \mathbb{R}^{n \times p}, \\ \mathbf{g}_{2,(\mathbf{h},\mathbf{m})} &= \begin{bmatrix} \mathbf{0} \\ \mathbf{v}_{(\mathbf{h},\mathbf{m}),(\mathbf{h}_0,\mathbf{m}_0)}^{(2)} \end{bmatrix} \in \mathbb{R}^{n \times p},\end{aligned}$$

$$\mathbf{t}_{(\mathbf{h},\mathbf{m}),(\mathbf{h}_0,\mathbf{m}_0)}^{(1)} = \frac{\alpha}{2^{d+s}\ell} \|\mathbf{m}\|_2^2 \mathbf{h} - \frac{\alpha}{\ell} \mathbf{m}^\top \tilde{\mathbf{t}}_{\mathbf{m},\mathbf{m}_0}^{(2)} \tilde{\mathbf{t}}_{\mathbf{h},\mathbf{h}_0}^{(1)}, \quad (10)$$

$$\mathbf{t}_{(\mathbf{h},\mathbf{m}),(\mathbf{h}_0,\mathbf{m}_0)}^{(2)} = \frac{\alpha}{2^{d+s}\ell} \|\mathbf{h}\|_2^2 \mathbf{m} - \frac{\alpha}{\ell} \mathbf{h}^\top \tilde{\mathbf{t}}_{\mathbf{h},\mathbf{h}_0}^{(1)} \tilde{\mathbf{t}}_{\mathbf{m},\mathbf{m}_0}^{(2)}, \quad (11)$$

$$\begin{aligned}S_{\epsilon,(\mathbf{h}_0,\mathbf{m}_0)}^{(1)} &= \left\{ (\mathbf{h}, \mathbf{m}) \in \mathbb{R}^{n \times p} \setminus \mathcal{K} \mid \frac{\|\mathbf{t}_{(\mathbf{h},\mathbf{m}),(\mathbf{h}_0,\mathbf{m}_0)}^{(1)}\|_2}{\|\mathbf{m}\|_2} \leq \frac{\epsilon \max(\|\mathbf{h}\|_2 \|\mathbf{m}\|_2, \|\mathbf{h}_0\|_2 \|\mathbf{m}_0\|_2)}{2^{d+s}\ell} \right\}, \\ S_{\epsilon,(\mathbf{h}_0,\mathbf{m}_0)}^{(2)} &= \left\{ (\mathbf{h}, \mathbf{m}) \in \mathbb{R}^{n \times p} \setminus \mathcal{K} \mid \frac{\|\mathbf{t}_{(\mathbf{h},\mathbf{m}),(\mathbf{h}_0,\mathbf{m}_0)}^{(2)}\|_2}{\|\mathbf{h}\|_2} \leq \frac{\epsilon \max(\|\mathbf{h}\|_2 \|\mathbf{m}\|_2, \|\mathbf{h}_0\|_2 \|\mathbf{m}_0\|_2)}{2^{d+s}\ell} \right\},\end{aligned}$$

where  $\tilde{\mathbf{t}}_{\mathbf{m},\mathbf{m}_0}^{(2)}$  and  $\tilde{\mathbf{t}}_{\mathbf{h},\mathbf{h}_0}^{(1)}$  is as defined in (9). For brevity of notation, write  $\mathbf{v}_{(\mathbf{h},\mathbf{m})}^{(1)} = \mathbf{v}_{(\mathbf{h},\mathbf{m}),(\mathbf{h}_0,\mathbf{m}_0)}^{(1)}$ ,  $\mathbf{v}_{(\mathbf{h},\mathbf{m})}^{(2)} = \mathbf{v}_{(\mathbf{h},\mathbf{m}),(\mathbf{h}_0,\mathbf{m}_0)}^{(2)}$ ,  $\mathbf{t}_{(\mathbf{h},\mathbf{m})} = \mathbf{t}_{(\mathbf{h},\mathbf{m}),(\mathbf{h}_0,\mathbf{m}_0)}$ ,  $\mathbf{t}_{(\mathbf{h},\mathbf{m})}^{(1)} = \mathbf{t}_{(\mathbf{h},\mathbf{m}),(\mathbf{h}_0,\mathbf{m}_0)}^{(1)}$  and  $\mathbf{t}_{(\mathbf{h},\mathbf{m})}^{(2)} = \mathbf{t}_{(\mathbf{h},\mathbf{m}),(\mathbf{h}_0,\mathbf{m}_0)}^{(2)}$ .

Since  $\mathbf{W}_i^{(1)} \in \mathbb{R}^{n_i \times n_{i-1}}$  for  $i = 1, \dots, d-1$  and  $\mathbf{W}_i^{(2)} \in \mathbb{R}^{p_i \times p_{i-1}}$  for  $i = 1, \dots, s-1$  satisfy the WDC with constant  $\epsilon$  and 1,  $\mathbf{W}_d^{(1)} \in \mathbb{R}^{\ell \times n_{d-1}}$  satisfy WDC with constants  $\epsilon$  and  $\alpha_1$ ,  $\mathbf{W}_s^{(2)} \in \mathbb{R}^{\ell \times p_{s-1}}$  satisfy WDC with constants  $\epsilon$  and  $\alpha_2$ , and  $(\mathbf{W}_d^{(1)}, \mathbf{W}_s^{(2)})$  satisfy joint-WDC with constants  $\epsilon, \alpha = \alpha_1 \cdot \alpha_2$ , we have lemma 2 implying for all nonzero  $\mathbf{h}, \mathbf{h}_0 \in \mathbb{R}^n$  and nonzero  $\mathbf{m}, \mathbf{m}_0 \in \mathbb{R}^p$

$$\|\nabla_{\mathbf{h}} f(\mathbf{h}, \mathbf{m}) - \mathbf{t}_{(\mathbf{h},\mathbf{m})}^{(1)}\|_2 \leq K \frac{d^3 s^3 \sqrt{\epsilon}}{2^{d+s}\ell} \max(\|\mathbf{h}\|_2 \|\mathbf{m}\|_2, \|\mathbf{h}_0\|_2 \|\mathbf{m}_0\|_2) \|\mathbf{m}\|_2, \quad (12)$$

$$\|\nabla_{\mathbf{m}} f(\mathbf{h}, \mathbf{m}) - \mathbf{t}_{(\mathbf{h},\mathbf{m})}^{(2)}\|_2 \leq K \frac{d^3 s^2 \sqrt{\epsilon}}{2^{d+s}\ell} \max(\|\mathbf{h}\|_2 \|\mathbf{m}\|_2, \|\mathbf{h}_0\|_2 \|\mathbf{m}_0\|_2) \|\mathbf{h}\|_2. \quad (13)$$

Thus, we have, for all nonzero  $\mathbf{h}, \mathbf{h}_0 \in \mathbb{R}^n$  and nonzero  $\mathbf{m}, \mathbf{m}_0 \in \mathbb{R}^p$ ,

$$\begin{aligned}\|\mathbf{v}_{(\mathbf{h},\mathbf{m})}^{(1)} - \mathbf{t}_{(\mathbf{h},\mathbf{m})}^{(1)}\|_2 &= \lim_{\delta \rightarrow 0^+} \|\nabla_{\mathbf{h}} f((\mathbf{h}, \mathbf{x}) + \delta \mathbf{w}) - \mathbf{t}_{(\mathbf{h},\mathbf{m})+\delta \mathbf{w}}^{(1)}\|_2 \\ &\leq K \frac{d^3 s^3 \sqrt{\epsilon}}{2^{d+s}\ell} \max(\|\mathbf{h}\|_2 \|\mathbf{m}\|_2, \|\mathbf{h}_0\|_2 \|\mathbf{m}_0\|_2) \|\mathbf{m}\|_2, \text{ and} \\ \|\mathbf{v}_{(\mathbf{h},\mathbf{m})}^{(2)} - \mathbf{t}_{(\mathbf{h},\mathbf{m})}^{(2)}\|_2 &= \lim_{\delta \rightarrow 0^+} \|\nabla_{\mathbf{m}} f((\mathbf{h}, \mathbf{x}) + \delta \mathbf{w}) - \mathbf{t}_{(\mathbf{h},\mathbf{m})+\delta \mathbf{w}}^{(2)}\|_2 \\ &\leq K \frac{d^3 s^2 \sqrt{\epsilon}}{2^{d+s}\ell} \max(\|\mathbf{h}\|_2 \|\mathbf{m}\|_2, \|\mathbf{h}_0\|_2 \|\mathbf{m}_0\|_2) \|\mathbf{h}\|_2,\end{aligned}$$

where the inequalities follow from (12) and (13).

Note that the one-sided directional derivative of  $f$  in the direction of  $(\mathbf{x}, \mathbf{y}) \neq \mathbf{0}$  at  $(\mathbf{h}, \mathbf{y})$  is  $D_{(\mathbf{x},\mathbf{y})} f(\mathbf{h}, \mathbf{x}) = \lim_{t \rightarrow 0^+} \frac{1}{t} (f((\mathbf{h}, \mathbf{x}) + t(\mathbf{x}, \mathbf{y})) - f(\mathbf{h}, \mathbf{x}))$ . Due to the continuity and piecewise linearity of the function

$$\mathcal{G}(\mathbf{h}, \mathbf{m}) = \Lambda_{d,+,\mathbf{h}}^{(1)} \mathbf{h} \odot \Lambda_{s,+,\mathbf{m}}^{(2)} \mathbf{m},$$

we have that for any  $(\mathbf{h}, \mathbf{m}) \neq (\mathbf{0}, \mathbf{0})$  and  $(\mathbf{x}, \mathbf{y}) \neq \mathbf{0}$  that there exists a sequence  $\{(\mathbf{h}_n, \mathbf{m}_n)\} \rightarrow (\mathbf{h}, \mathbf{m})$  such that  $f$  is differentiable at each  $(\mathbf{h}_n, \mathbf{m}_n)$  and  $D_{(\mathbf{x},\mathbf{y})} f(\mathbf{h}, \mathbf{m}) = \lim_{n \rightarrow \infty} \nabla f(\mathbf{h}_n, \mathbf{m}_n) \cdot (\mathbf{x}, \mathbf{y})$ . Thus, as  $\nabla f(\mathbf{h}_n, \mathbf{m}_n) = \begin{bmatrix} \mathbf{v}_{(\mathbf{h}_n,\mathbf{m}_n)}^{(1)} \\ \mathbf{v}_{(\mathbf{h}_n,\mathbf{m}_n)}^{(2)} \end{bmatrix}$ ,

$$D_{-\mathbf{g}_{1,(\mathbf{h},\mathbf{m})}} f(\mathbf{h}, \mathbf{m}) = \lim_{n \rightarrow \infty} \nabla f(\mathbf{h}_n, \mathbf{m}_n) \cdot \frac{-\mathbf{g}_{1,(\mathbf{h},\mathbf{m})}}{\|\mathbf{g}_{1,(\mathbf{h},\mathbf{m})}\|_2} = \frac{-1}{\|\mathbf{g}_{1,(\mathbf{h},\mathbf{m})}\|_2} \lim_{n \rightarrow \infty} \mathbf{v}_{(\mathbf{h}_n,\mathbf{m}_n)}^{(1)} \cdot \mathbf{v}_{(\mathbf{h},\mathbf{m})}^{(1)},$$

$$D_{-\mathbf{g}_{2,(\mathbf{h},\mathbf{m})}} f(\mathbf{h}, \mathbf{m}) = \lim_{n \rightarrow \infty} \nabla f(\mathbf{h}_n, \mathbf{m}_n) \cdot \frac{-\mathbf{g}_{2,(\mathbf{h},\mathbf{m})}}{\|\mathbf{g}_{2,(\mathbf{h},\mathbf{m})}\|_2} = \frac{-1}{\|\mathbf{g}_{2,(\mathbf{h},\mathbf{m})}\|_2} \lim_{n \rightarrow \infty} \mathbf{v}_{(\mathbf{h}_n, \mathbf{m}_n)}^{(2)} \cdot \mathbf{v}_{(\mathbf{h}, \mathbf{m})}^{(2)}.$$

Now, we write

$$\begin{aligned} & \mathbf{v}_{(\mathbf{h}_n, \mathbf{m}_n)}^{(1)} \cdot \mathbf{v}_{(\mathbf{h}, \mathbf{m})}^{(1)} \\ &= \mathbf{t}_{(\mathbf{h}_n, \mathbf{m}_n)}^{(1)} \cdot \mathbf{t}_{(\mathbf{h}, \mathbf{m})}^{(1)} + (\mathbf{v}_{(\mathbf{h}_n, \mathbf{m}_n)}^{(1)} - \mathbf{t}_{(\mathbf{h}_n, \mathbf{m}_n)}^{(1)}) \cdot \mathbf{t}_{(\mathbf{h}, \mathbf{m})}^{(1)} + \mathbf{t}_{(\mathbf{h}_n, \mathbf{m}_n)}^{(1)} \cdot (\mathbf{v}_{(\mathbf{h}, \mathbf{m})}^{(1)} - \mathbf{t}_{(\mathbf{h}, \mathbf{m})}^{(1)}) \\ &\quad + (\mathbf{v}_{(\mathbf{h}_n, \mathbf{m}_n)}^{(1)} - \mathbf{t}_{(\mathbf{h}_n, \mathbf{m}_n)}^{(1)}) \cdot (\mathbf{v}_{(\mathbf{h}, \mathbf{m})}^{(1)} - \mathbf{t}_{(\mathbf{h}, \mathbf{m})}^{(1)}) \\ &\geq \mathbf{t}_{(\mathbf{h}_n, \mathbf{m}_n)}^{(1)} \cdot \mathbf{t}_{(\mathbf{h}, \mathbf{m})}^{(1)} - \|\mathbf{v}_{(\mathbf{h}_n, \mathbf{m}_n)}^{(1)} - \mathbf{t}_{(\mathbf{h}_n, \mathbf{m}_n)}^{(1)}\|_2 \|\mathbf{t}_{(\mathbf{h}, \mathbf{m})}^{(1)}\|_2 - \|\mathbf{t}_{(\mathbf{h}_n, \mathbf{m}_n)}^{(1)}\|_2 \|\mathbf{v}_{(\mathbf{h}, \mathbf{m})}^{(1)} - \mathbf{t}_{(\mathbf{h}, \mathbf{m})}^{(1)}\|_2 \\ &\quad \|\mathbf{v}_{(\mathbf{h}_n, \mathbf{m}_n)}^{(1)} - \mathbf{t}_{(\mathbf{h}_n, \mathbf{m}_n)}^{(1)}\|_2 \|\mathbf{v}_{(\mathbf{h}, \mathbf{m})}^{(1)} - \mathbf{t}_{(\mathbf{h}, \mathbf{m})}^{(1)}\|_2 \\ &\geq \mathbf{t}_{(\mathbf{h}_n, \mathbf{m}_n)}^{(1)} \cdot \mathbf{t}_{(\mathbf{h}, \mathbf{m})}^{(1)} - K \frac{d^3 s^3 \sqrt{\epsilon}}{2^{d+s} \ell} \max(\|\mathbf{h}_n\|_2 \|\mathbf{m}_n\|_2, \|\mathbf{h}_0\|_2 \|\mathbf{m}_0\|_2) \|\mathbf{m}_n\|_2 \|\mathbf{t}_{(\mathbf{h}_n, \mathbf{m}_n)}^{(1)}\|_2 \\ &\quad - K \frac{d^3 s^3 \sqrt{\epsilon}}{2^{d+s} \ell} \max(\|\mathbf{h}\|_2 \|\mathbf{m}\|_2, \|\mathbf{h}_0\|_2 \|\mathbf{m}_0\|_2) \|\mathbf{m}\|_2 \|\mathbf{t}_{(\mathbf{h}, \mathbf{m})}^{(1)}\|_2 \\ &\quad - \left( K \frac{d^3 s^3 \sqrt{\epsilon}}{2^{d+s} \ell} \right)^2 \max(\|\mathbf{h}_n\|_2 \|\mathbf{m}_n\|_2, \|\mathbf{h}_0\|_2 \|\mathbf{m}_0\|_2) \\ &\quad \max(\|\mathbf{h}\|_2 \|\mathbf{m}\|_2, \|\mathbf{h}_0\|_2 \|\mathbf{m}_0\|_2) \|\mathbf{m}_n\|_2 \|\mathbf{m}\|_2. \end{aligned}$$

As  $\mathbf{t}_{(\mathbf{h}, \mathbf{m})}^{(1)}$  is continuous in  $(\mathbf{h}, \mathbf{m})$  for all  $(\mathbf{h}, \mathbf{m}) \notin \mathcal{K}$ , we have for all  $(\mathbf{h}, \mathbf{m}) \notin \mathcal{S}_{4Kd^3s^3\sqrt{\epsilon}, (\mathbf{h}_0, \mathbf{m}_0)}^{(1)} \cup \mathcal{K}$ ,

$$\begin{aligned} & \lim_{n \rightarrow \infty} \mathbf{v}_{(\mathbf{h}_n, \mathbf{m}_n)}^{(1)} \cdot \mathbf{v}_{(\mathbf{h}, \mathbf{m})}^{(1)} \\ &\geq \|\mathbf{t}_{(\mathbf{h}, \mathbf{m})}^{(1)}\|_2^2 - 2K \frac{d^3 s^3 \sqrt{\epsilon}}{2^{d+s} \ell} \max(\|\mathbf{h}\|_2 \|\mathbf{m}\|_2, \|\mathbf{h}_0\|_2 \|\mathbf{m}_0\|_2) \|\mathbf{m}\|_2 \|\mathbf{t}_{(\mathbf{h}, \mathbf{m})}^{(1)}\|_2 \\ &\quad - \left( K \frac{d^3 s^3 \sqrt{\epsilon}}{2^{d+s} \ell} \max(\|\mathbf{h}\|_2 \|\mathbf{m}\|_2, \|\mathbf{h}_0\|_2 \|\mathbf{m}_0\|_2) \|\mathbf{m}\|_2 \right)^2 \\ &\geq \frac{\|\mathbf{t}_{(\mathbf{h}, \mathbf{m})}^{(1)}\|_2^2}{2} \left[ \|\mathbf{t}_{(\mathbf{h}, \mathbf{m})}^{(1)}\|_2 - 4K \frac{d^3 s^3 \sqrt{\epsilon}}{2^{d+s} \ell} \max(\|\mathbf{h}\|_2 \|\mathbf{m}\|_2, \|\mathbf{h}_0\|_2 \|\mathbf{m}_0\|_2) \|\mathbf{m}\|_2 \right] + \\ &\quad \frac{1}{2} \left[ \|\mathbf{t}_{(\mathbf{h}, \mathbf{m})}^{(1)}\|_2^2 - 2 \left( K \frac{d^3 s^3 \sqrt{\epsilon}}{2^{d+s} \ell} \max(\|\mathbf{h}\|_2 \|\mathbf{m}\|_2, \|\mathbf{h}_0\|_2 \|\mathbf{m}_0\|_2) \|\mathbf{m}\|_2 \right)^2 \right] \\ &> 0. \end{aligned} \tag{14}$$

Similarly, we have for all  $(\mathbf{h}, \mathbf{m}) \notin \mathcal{S}_{4Kd^3s^3\sqrt{\epsilon}, (\mathbf{h}_0, \mathbf{m}_0)}^{(2)} \cup \mathcal{K}$ ,

$$\lim_{n \rightarrow \infty} \mathbf{v}_{(\mathbf{h}_n, \mathbf{m}_n)}^{(2)} \cdot \mathbf{v}_{(\mathbf{h}, \mathbf{m})}^{(2)} > 0. \tag{15}$$

So, for all  $(\mathbf{h}, \mathbf{m}) \notin \left( \mathcal{S}_{4Kd^3s^3\sqrt{\epsilon}, (\mathbf{h}_0, \mathbf{m}_0)}^{(1)} \cap \mathcal{S}_{4Kd^3s^3\sqrt{\epsilon}, (\mathbf{h}_0, \mathbf{m}_0)}^{(2)} \right) \cup \mathcal{K}$ , at least (14) or (15) holds. If (14) holds, then we have  $D_{-\mathbf{g}_{1,(\mathbf{h},\mathbf{m})}} f(\mathbf{h}, \mathbf{m}) < 0$  and if (15) holds, then we have  $D_{-\mathbf{g}_{2,(\mathbf{h},\mathbf{m})}} f(\mathbf{h}, \mathbf{m}) < 0$ . Let  $\mathcal{S} = \mathcal{S}_{4Kd^3s^3\sqrt{\epsilon}, (\mathbf{h}_0, \mathbf{m}_0)}^{(1)} \cap \mathcal{S}_{4Kd^3s^3\sqrt{\epsilon}, (\mathbf{h}_0, \mathbf{m}_0)}^{(2)}$ . Apply Lemma 3 and  $38(d^5 + s^5)\sqrt{4Kd^3s^3\sqrt{\epsilon}}/\alpha < 1$  to get

$$\begin{aligned} \mathcal{S} \subseteq & \mathcal{A}_{\tilde{K} \frac{d^3 s^3 \epsilon^{1/4}}{\alpha}, (\mathbf{h}_0, \mathbf{m}_0)} \cup \mathcal{A}_{\tilde{K} \frac{d^3 s^3 \epsilon^{1/4}}{\alpha}, (-\rho_d^{(1)} \mathbf{h}_0, \mathbf{m}_0)} \cup \mathcal{A}_{\tilde{K} \frac{d^3 s^8 \epsilon^{1/4}}{\alpha}, (\rho_s^{(2)} \mathbf{h}_0, -\mathbf{m}_0)} \\ & \cup \mathcal{A}_{\tilde{K} \frac{d^8 s^8 \epsilon^{1/4}}{\alpha}, (-\rho_d^{(1)} \rho_s^{(2)} \mathbf{h}_0, -\mathbf{m}_0)}, \end{aligned}$$

for some absolute constant  $\tilde{K}$ .

It remains to prove that elements of the set  $\mathcal{K}$  are local maximizers. We now show that elements of the set  $\{(\mathbf{0}, \mathbf{m}) \in \mathbb{R}^{n \times p} | \mathbf{m} \in \mathbb{R}^p\}$  are local maximizers. Fix a direction  $(\mathbf{x}, \mathbf{y}) \in \mathbb{R}^{n \times p}$  and let

$$\tilde{\mathbf{h}}_0 = \Lambda_{d-1,+, \mathbf{h}_0}^{(1)} \mathbf{h}_0, \quad \tilde{\mathbf{x}} = \Lambda_{d-1,+, \mathbf{x}}^{(1)} \mathbf{x}, \quad \tilde{\mathbf{m}} = \Lambda_{s-1,+, \mathbf{m}}^{(2)} \mathbf{m}, \quad \tilde{\mathbf{m}}_0 = \Lambda_{s-1,+, \mathbf{m}_0}^{(2)} \mathbf{m}_0,$$

$\bar{\theta}_i^{(k)} = g(\bar{\theta}_{i-1}^{(k)})$  for  $g$  given in (8),  $\bar{\theta}_0^{(1)} = \angle(\mathbf{x}, \mathbf{h}_0)$  and  $\bar{\theta}_0^{(2)} = \angle(\mathbf{m}, \mathbf{m}_0)$ . We compute

$$\begin{aligned}
& -D_{(\mathbf{x}, \mathbf{y})} f(\mathbf{h}, \mathbf{m}) \cdot \|(\mathbf{x}, \mathbf{y})\|_2 \\
&= \lim_{t \rightarrow 0^+} \frac{f((\mathbf{h}, \mathbf{m}) + t(\mathbf{x}, \mathbf{y})) - f(\mathbf{h}, \mathbf{m})}{t} \\
&= \left\langle \text{diag} \left( \Lambda_{s,+,\mathbf{m}}^{(2)} \mathbf{m} \right) \Lambda_{d,+,\mathbf{x}}^{(1)} \mathbf{x}, \text{diag} \left( \Lambda_{s,+,\mathbf{m}_0}^{(2)} \mathbf{m}_0 \right) \Lambda_{d,+,\mathbf{h}_0}^{(1)} \mathbf{h}_0 \right\rangle \\
&= \left\langle \tilde{\mathbf{x}}, \left( \mathbf{W}_{d-1,+,\mathbf{m}}^{(2)} \right)^\top \text{diag} \left( \mathbf{W}_{d-1,+,\mathbf{x}}^{(1)} \tilde{\mathbf{m}} \odot \mathbf{W}_{d-1,+,\mathbf{x}_0}^{(1)} \tilde{\mathbf{m}}_0 \right) \mathbf{W}_{d-1,+,\mathbf{h}_0}^{(1)} \tilde{\mathbf{h}}_0 \right\rangle \\
&= \left\langle \tilde{\mathbf{x}}, \left( \left( \mathbf{W}_{d-1,+,\mathbf{m}}^{(2)} \right)^\top \text{diag} \left( \mathbf{W}_{d-1,+,\mathbf{x}}^{(1)} \tilde{\mathbf{m}} \odot \mathbf{W}_{d-1,+,\mathbf{x}_0}^{(1)} \tilde{\mathbf{m}}_0 \right) \mathbf{W}_{d-1,+,\mathbf{h}_0}^{(1)} \right. \right. \\
&\quad \left. \left. - \frac{\alpha}{n} \mathbf{Q}_{\tilde{\mathbf{x}}, \tilde{\mathbf{h}}_0} \tilde{\mathbf{m}}^\top \mathbf{Q}_{\tilde{\mathbf{m}}, \tilde{\mathbf{m}}_0} \tilde{\mathbf{m}}_0 \right) \tilde{\mathbf{h}}_0 \right\rangle + \frac{\alpha}{n} \tilde{\mathbf{m}}^\top \mathbf{Q}_{\tilde{\mathbf{m}}, \tilde{\mathbf{m}}_0} \tilde{\mathbf{m}}_0 \cdot \tilde{\mathbf{x}}^\top \mathbf{Q}_{\tilde{\mathbf{x}}, \tilde{\mathbf{h}}_0} \tilde{\mathbf{h}}_0 \\
&\geq - \left\| \left( \mathbf{W}_{d-1,+,\mathbf{m}}^{(2)} \right)^\top \text{diag} \left( \mathbf{W}_{d-1,+,\mathbf{x}}^{(1)} \tilde{\mathbf{m}} \odot \mathbf{W}_{d-1,+,\mathbf{x}_0}^{(1)} \tilde{\mathbf{m}}_0 \right) \mathbf{W}_{d-1,+,\mathbf{h}_0}^{(1)} \right. \\
&\quad \left. - \frac{\alpha}{n} \mathbf{Q}_{\tilde{\mathbf{x}}, \tilde{\mathbf{h}}_0} \tilde{\mathbf{m}}^\top \mathbf{Q}_{\tilde{\mathbf{m}}, \tilde{\mathbf{m}}_0} \tilde{\mathbf{m}}_0 \right\| \|\tilde{\mathbf{x}}\|_2 \|\tilde{\mathbf{h}}_0\|_2 + \frac{\alpha}{4n} \left( \frac{(\pi - \bar{\theta}_{d-1}^{(1)}) \cos \bar{\theta}_{d-1}^{(1)} + \sin \bar{\theta}_{d-1}^{(1)}}{\pi} \right) \\
&\quad \left( \frac{(\pi - \bar{\theta}_{d-1}^{(2)}) \cos \bar{\theta}_{d-1}^{(2)} + \sin \bar{\theta}_{d-1}^{(2)}}{\pi} \right) \|\tilde{\mathbf{x}}\| \|\tilde{\mathbf{h}}_0\| \|\tilde{\mathbf{m}}\| \|\tilde{\mathbf{m}}_0\| \\
&\geq -\frac{4\epsilon}{n} \|\tilde{\mathbf{m}}\|_2 \|\tilde{\mathbf{m}}_0\|_2 \|\tilde{\mathbf{x}}\|_2 \|\tilde{\mathbf{h}}_0\|_2 + \frac{\alpha}{4n} \|\tilde{\mathbf{x}}\| \|\tilde{\mathbf{h}}_0\| \|\tilde{\mathbf{m}}\| \|\tilde{\mathbf{m}}_0\| \cos \bar{\theta}_d^{(1)} \cos \bar{\theta}_d^{(2)} \\
&\geq \left( -\frac{4\epsilon}{n} + \frac{\alpha}{4\pi^2 n} \right) \|\tilde{\mathbf{x}}\| \|\tilde{\mathbf{h}}_0\| \|\tilde{\mathbf{m}}\| \|\tilde{\mathbf{m}}_0\|.
\end{aligned}$$

By Lemma 1, we have  $\|\tilde{\mathbf{x}}\| \geq \frac{1}{4\pi} \frac{1}{2^{d-1}} \|\mathbf{x}\|$  and  $\|\tilde{\mathbf{m}}\| \geq \frac{1}{4\pi} \frac{1}{2^{s-1}} \|\mathbf{m}\|$ . So, if  $4\pi^2 \epsilon / \alpha < 1$ , then  $D_{(\mathbf{x}, \mathbf{y})} f(\mathbf{h}, \mathbf{m}) \cdot \|(\mathbf{x}, \mathbf{y})\|_2 < 0$  for all  $(\mathbf{x}, \mathbf{y}) \in \mathbb{R}^{n \times p}$  with  $\mathbf{x} \neq \mathbf{0}$  and  $D_{(\mathbf{x}, \mathbf{y})} f(\mathbf{h}, \mathbf{m}) \cdot \|(\mathbf{x}, \mathbf{y})\|_2 = 0$  for all  $(\mathbf{x}, \mathbf{y}) \in \mathbb{R}^{n \times p}$  with  $\mathbf{x} = \mathbf{0}$ . Thus, elements of the set  $\{(\mathbf{0}, \mathbf{m}) \in \mathbb{R}^{n \times p} | \mathbf{m} \in \mathbb{R}^p\}$  are local maximizers. Similarly, elements of the set  $\{(\mathbf{h}, \mathbf{0}) \in \mathbb{R}^{n \times p} | \mathbf{h} \in \mathbb{R}^n\}$  are local maximizers. This concludes the proof of Theorem 4.  $\square$

## 5.2 Concentration of terms in $\tilde{g}_{1,(\mathbf{h}, \mathbf{m})}$ and $\tilde{g}_{2,(\mathbf{h}, \mathbf{m})}$

**Lemma 1.** Fix  $0 < \epsilon < d^{-4}/(16\pi)^2$  and  $d \geq 2$ . Suppose that  $\mathbf{W}_i \in \mathbb{R}^{n_i \times n_{i-1}}$  satisfies the WDC with constant  $\epsilon$  and 1 for  $i = 1, \dots, d$ . Define

$$\tilde{t}_{\mathbf{p}, \mathbf{q}} = \frac{1}{2^d} \left[ \left( \prod_{i=0}^{d-1} \frac{\pi - \bar{\theta}_i}{\pi} \right) \mathbf{q} + \sum_{i=0}^{d-1} \frac{\sin \bar{\theta}_i}{\pi} \left( \prod_{j=i+1}^{d-1} \frac{\pi - \bar{\theta}_j}{\pi} \right) \frac{\|\mathbf{q}\|_2}{\|\mathbf{p}\|_2} \mathbf{p} \right],$$

where  $\bar{\theta}_i = g(\bar{\theta}_{i-1})$  for  $g$  given by (8) and  $\bar{\theta}_0 = \angle(\mathbf{p}, \mathbf{q})$ . For all  $\mathbf{p} \neq \mathbf{0}$  and  $\mathbf{q} \neq \mathbf{0}$ ,

$$\left\| \left( \prod_{i=d}^1 \mathbf{W}_{i,+,\mathbf{p}} \right)^\top \left( \prod_{i=d}^1 \mathbf{W}_{i,+,\mathbf{q}} \right) \mathbf{q} - \tilde{t}_{\mathbf{p}, \mathbf{q}} \right\|_2 \leq 24 \frac{d^3 \sqrt{\epsilon}}{2^d} \|\mathbf{q}\|_2, \quad (16)$$

$$\left\langle \left( \prod_{i=d}^1 \mathbf{W}_{i,+,\mathbf{p}} \right) \mathbf{p}, \left( \prod_{i=d}^1 \mathbf{W}_{i,+,\mathbf{q}} \right) \mathbf{q} \right\rangle \geq \frac{1}{4\pi} \frac{1}{2^d} \|\mathbf{p}\|_2 \|\mathbf{q}\|_2 \quad (17)$$

We refer the readers to Hand and Voroninski [2017] for proof of Lemma 1. We now state a related Lemma.

**Lemma 2.** Fix  $0 < \epsilon < 1/((d^4 + s^4)16\pi)^2$ ,  $d \geq 2$  and  $s \geq 2$ . Suppose that  $\mathbf{W}_i^{(1)} \in \mathbb{R}^{n_i \times n_{i-1}}$  for  $i = 1, \dots, d-1$  and  $\mathbf{W}_i^{(2)} \in \mathbb{R}^{p_i \times p_{i-1}}$  for  $i = 1, \dots, s-1$  satisfy the WDC with constant  $\epsilon$  and 1. Suppose  $\mathbf{W}_d^{(1)} \in \mathbb{R}^{\ell \times n_{d-1}}$  satisfy WDC with constants  $\epsilon$  and  $\alpha_1$ , and  $\mathbf{W}_s^{(2)} \in \mathbb{R}^{\ell \times p_{s-1}}$  satisfy WDC with constants  $\epsilon$  and  $\alpha_2$ . Also, suppose  $(\mathbf{W}_d^{(1)}, \mathbf{W}_s^{(2)})$  satisfy pair-WDC with constants  $\epsilon$ ,

$\alpha = \alpha_1 \cdot \alpha_2$ . Define

$$\tilde{t}_{\mathbf{p}, \mathbf{q}}^{(k)} = \frac{1}{2^{a^{(k)}}} \left[ \left( \prod_{i=0}^{a^{(k)}-1} \frac{\pi - \bar{\theta}_i^{(k)}}{\pi} \right) \mathbf{q} + \sum_{i=0}^{a^{(k)}-1} \frac{\sin \bar{\theta}_i^{(k)}}{\pi} \left( \prod_{j=i+1}^{a^{(k)}-1} \frac{\pi - \bar{\theta}_j^{(k)}}{\pi} \right) \frac{\|\mathbf{q}\|_2}{\|\mathbf{p}\|_2} \mathbf{p} \right],$$

where  $\bar{\theta}_i^{(k)} = g(\bar{\theta}_{i-1}^{(k)})$  for  $g$  given by (8),  $\bar{\theta}_0^{(k)} = \angle(\mathbf{p}, \mathbf{q})$ ,  $a^{(1)} = d$ , and  $a^{(2)} = s$ . For all  $\mathbf{h} \neq 0$ ,  $\mathbf{x} \neq 0$ ,  $\mathbf{m} \neq 0$  and  $\mathbf{y} \neq 0$ ,

$$\begin{aligned} & \left\| \left( \Lambda_{d,+,\mathbf{h}}^{(1)} \right)^T \text{diag} \left( \Lambda_{s,+,\mathbf{m}}^{(2)} \mathbf{m} \odot \Lambda_{s,+,\mathbf{y}}^{(2)} \mathbf{y} \right) \Lambda_{d,+,\mathbf{x}}^{(1)} \mathbf{x} - \frac{\alpha}{n} \left( \mathbf{m}^T \tilde{\mathbf{t}}_{\mathbf{m}, \mathbf{y}}^{(2)} \right) \tilde{\mathbf{t}}_{\mathbf{h}, \mathbf{x}}^{(1)} \right\|_2 \\ & \leq \frac{208d^3s^3\sqrt{\epsilon}}{2^{d+s}\ell} \|\mathbf{x}\|_2 \|\mathbf{m}\|_2 \|\mathbf{y}\|_2, \end{aligned} \quad (18)$$

$$\begin{aligned} & \left\| \left( \Lambda_{s,+,\mathbf{m}}^{(2)} \right)^T \text{diag} \left( \Lambda_{d,+,\mathbf{h}}^{(1)} \mathbf{h} \odot \Lambda_{d,+,\mathbf{x}}^{(1)} \mathbf{x} \right) \Lambda_{s,+,\mathbf{y}}^{(2)} \mathbf{y} - \frac{\alpha}{n} \left( \mathbf{h}^T \tilde{\mathbf{t}}_{\mathbf{h}, \mathbf{x}}^{(1)} \right) \tilde{\mathbf{t}}_{\mathbf{m}, \mathbf{y}}^{(2)} \right\|_2 \\ & \leq \frac{208d^3s^3\sqrt{\epsilon}}{2^{d+s}\ell} \|\mathbf{y}\|_2 \|\mathbf{h}\|_2 \|\mathbf{x}\|_2. \end{aligned} \quad (19)$$

*Proof.* We will prove (18). Proof of (19) is identical to proof of (18). Define  $\mathbf{h}_0 = \mathbf{h}$ ,  $\mathbf{x}_0 = \mathbf{x}$ ,  $\mathbf{m}_0 = \mathbf{m}$ ,  $\mathbf{y}_0 = \mathbf{y}$ ,

$$\begin{aligned} \mathbf{h}_d := \left( \prod_{i=d}^1 \mathbf{W}_{i,+,\mathbf{h}}^{(1)} \right) \mathbf{h} &= \left( \mathbf{W}_{d,+,\mathbf{h}}^{(1)} \mathbf{W}_{d-1,+,\mathbf{h}}^{(1)} \cdots \mathbf{W}_{1,+,\mathbf{h}}^{(1)} \right) \mathbf{h} \\ &= \mathbf{W}_{d,+,\mathbf{h}}^{(1)} \mathbf{h}_{d-1} \\ &= (\mathbf{W}_d^{(1)})_{+, \mathbf{h}_{d-1}} \mathbf{h}_{d-1}, \end{aligned}$$

and analogously  $\mathbf{x}_d = \left( \prod_{i=d}^1 \mathbf{W}_{i,+,\mathbf{x}}^{(1)} \right) \mathbf{x}$ ,  $\mathbf{m}_s = \left( \prod_{i=s}^1 \mathbf{W}_{i,+,\mathbf{m}}^{(2)} \right) \mathbf{m}$ , and  $\mathbf{y}_s = \left( \prod_{i=s}^1 \mathbf{W}_{i,+,\mathbf{y}}^{(2)} \right) \mathbf{y}$ . By the WDC, we have for all  $\mathbf{h} \neq 0$ ,  $\mathbf{m} \neq 0$ ,

$$\left\| \left( \mathbf{W}_i^{(1)} \right)_{+, \mathbf{h}}^T \left( \mathbf{W}_i^{(1)} \right)_{+, \mathbf{h}} - \frac{1}{2} \mathbf{I}_{n_{i-1}} \right\| \leq \epsilon \text{ for all } i = 1, \dots, d-1, \text{ and} \quad (20)$$

$$\left\| \left( \mathbf{W}_i^{(2)} \right)_{+, \mathbf{m}}^T \left( \mathbf{W}_i^{(2)} \right)_{+, \mathbf{m}} - \frac{1}{2} \mathbf{I}_{p_{i-1}} \right\| \leq \epsilon \text{ for all } i = 1, \dots, s-1. \quad (21)$$

In particular,  $\left\| \left( \mathbf{W}_{i,+,\mathbf{h}}^{(1)} \right)^T \mathbf{W}_{i,+,\mathbf{h}}^{(1)} - \frac{1}{2} \mathbf{I}_{n_{i-1}} \right\| \leq \epsilon$  and  $\left\| \left( \mathbf{W}_{i,+,\mathbf{m}}^{(2)} \right)^T \mathbf{W}_{i,+,\mathbf{m}}^{(2)} - \frac{1}{2} \mathbf{I}_{p_{i-1}} \right\| \leq \epsilon$ . and consequently,

$$\begin{aligned} \frac{1}{2} - \epsilon &\leq \left\| \mathbf{W}_{i,+,\mathbf{h}}^{(1)} \right\|^2 \leq \frac{1}{2} + \epsilon \\ \frac{1}{2} - \epsilon &\leq \left\| \mathbf{W}_{i,+,\mathbf{m}}^{(2)} \right\|^2 \leq \frac{1}{2} + \epsilon. \end{aligned}$$

Hence,

$$\left\| \prod_{i=d-1}^1 \mathbf{W}_{i,+,\mathbf{h}}^{(1)} \right\| \left\| \prod_{i=d-1}^1 \mathbf{W}_{i,+,\mathbf{x}}^{(1)} \right\| \leq \frac{1}{2^{d-1}} (1+2\epsilon)^{d-1} = \frac{1}{2^{d-1}} e^{(d-1)\log(1+2\epsilon)} \leq \frac{1+4\epsilon(d-1)}{2^{d-1}}, \quad (22)$$

where we used that  $\log(1+z) \leq z$ ,  $e^z \leq 1+2z$  for  $z < 1$ , and  $2(d-1)\epsilon \leq 1$ . Similarly,

$$\left\| \prod_{i=s-1}^1 \mathbf{W}_{i,+,\mathbf{m}}^{(2)} \right\| \left\| \prod_{i=s-1}^1 \mathbf{W}_{i,+,\mathbf{y}}^{(2)} \right\| \leq \frac{1+4\epsilon(s-1)}{2^{s-1}}. \quad (23)$$

Let

$$\tilde{\mathbf{h}} = \Lambda_{d-1,+,\mathbf{h}}^{(1)} \mathbf{h}, \quad \tilde{\mathbf{x}} = \Lambda_{d-1,+,\mathbf{x}}^{(1)} \mathbf{x}, \quad \tilde{\mathbf{m}} = \Lambda_{s-1,+,\mathbf{m}}^{(2)} \mathbf{m}, \quad \tilde{\mathbf{y}} = \Lambda_{s-1,+,\mathbf{y}}^{(2)} \mathbf{y},$$

and consider

$$\begin{aligned} & \left\| \left( \Lambda_{d,+,\mathbf{h}}^{(1)} \right)^\top \text{diag} \left( \Lambda_{s,+,\mathbf{m}}^{(2)} \mathbf{m} \odot \Lambda_{s,+,\mathbf{y}}^{(2)} \mathbf{y} \right) \Lambda_{d,+,\mathbf{x}}^{(1)} \mathbf{x} - \frac{\alpha}{\ell} \left( \mathbf{m}^\top \tilde{t}_{\mathbf{m},\mathbf{y}}^{(2)} \right) \tilde{t}_{\mathbf{h},\mathbf{x}}^{(1)} \right\|_2 \\ & \leq \left\| \left( \Lambda_{d-1,+,\mathbf{h}}^{(1)} \right)^\top \left( \left( \mathbf{W}_{d,+,\mathbf{h}}^{(1)} \right)^\top \text{diag} \left( \mathbf{W}_{d,+,\mathbf{m}}^{(2)} \tilde{\mathbf{m}} \odot \mathbf{W}_{d,+,\mathbf{y}}^{(2)} \tilde{\mathbf{y}} \right) \mathbf{W}_{d,+,\mathbf{x}}^{(1)} - \frac{\alpha}{\ell} \mathbf{Q}_{\tilde{\mathbf{h}},\tilde{\mathbf{x}}} \right. \right. \\ & \quad \left. \tilde{\mathbf{m}}^\top \mathbf{Q}_{\tilde{\mathbf{m}},\tilde{\mathbf{y}}} \tilde{\mathbf{y}} \right) \Lambda_{d-1,+,\mathbf{x}}^{(1)} \mathbf{x} \right\|_2 + \left\| \frac{\alpha}{\ell} \left( \Lambda_{d-1,+,\mathbf{h}}^{(1)} \right)^\top \mathbf{Q}_{\tilde{\mathbf{h}},\tilde{\mathbf{x}}} \tilde{\mathbf{m}}^\top \mathbf{Q}_{\tilde{\mathbf{m}},\tilde{\mathbf{y}}} \tilde{\mathbf{y}} \Lambda_{d-1,+,\mathbf{x}}^{(1)} \mathbf{x} \right. \\ & \quad \left. - \frac{\alpha}{\ell} \left( \mathbf{m}^\top \tilde{t}_{\mathbf{m},\mathbf{y}}^{(2)} \right) \tilde{t}_{\mathbf{h},\mathbf{x}}^{(1)} \right\|_2 \\ & \leq \left\| \left( \Lambda_{d-1,+,\mathbf{h}}^{(1)} \right)^\top \left( \left( \mathbf{W}_{d,+,\mathbf{h}}^{(1)} \right)^\top \text{diag} \left( \mathbf{W}_{d,+,\mathbf{m}}^{(2)} \tilde{\mathbf{m}} \odot \mathbf{W}_{d,+,\mathbf{y}}^{(2)} \tilde{\mathbf{y}} \right) \mathbf{W}_{d,+,\mathbf{x}}^{(1)} - \frac{\alpha}{\ell} \mathbf{Q}_{\tilde{\mathbf{h}},\tilde{\mathbf{x}}} \right. \right. \\ & \quad \left. \tilde{\mathbf{m}}^\top \mathbf{Q}_{\tilde{\mathbf{m}},\tilde{\mathbf{y}}} \tilde{\mathbf{y}} \right) \Lambda_{d-1,+,\mathbf{x}}^{(1)} \mathbf{x} \right\|_2 \\ & \quad + \frac{\alpha}{\ell} \left\| \tilde{\mathbf{m}}^\top \mathbf{Q}_{\tilde{\mathbf{m}},\tilde{\mathbf{y}}} \tilde{\mathbf{y}} \left( \Lambda_{d-1,+,\mathbf{h}}^{(1)} \right)^\top \mathbf{Q}_{\tilde{\mathbf{h}},\tilde{\mathbf{x}}} \Lambda_{d-1,+,\mathbf{x}}^{(1)} \mathbf{x} - \mathbf{m}^\top \tilde{t}_{\mathbf{m},\mathbf{y}}^{(2)} \left( \Lambda_{d-1,+,\mathbf{h}}^{(1)} \right)^\top \mathbf{Q}_{\tilde{\mathbf{h}},\tilde{\mathbf{x}}} \Lambda_{d-1,+,\mathbf{x}}^{(1)} \mathbf{x} \right\|_2 \\ & \quad + \frac{\alpha}{\ell} \left\| \mathbf{m}^\top \tilde{t}_{\mathbf{m},\mathbf{y}}^{(2)} \left( \Lambda_{d-1,+,\mathbf{h}}^{(1)} \right)^\top \mathbf{Q}_{\tilde{\mathbf{h}},\tilde{\mathbf{x}}} \Lambda_{d-1,+,\mathbf{x}}^{(1)} \mathbf{x} - \mathbf{m}^\top \tilde{t}_{\mathbf{m},\mathbf{y}}^{(2)} \tilde{t}_{\mathbf{h},\mathbf{x}}^{(1)} \right\|_2 \end{aligned} \tag{24}$$

where both the first and second inequality holds because of triangle inequality. We bound the terms in the inequality above separately. First consider

$$\begin{aligned} & \left\| \left( \Lambda_{d-1,+,\mathbf{h}}^{(1)} \right)^\top \left( \left( \mathbf{W}_{d,+,\mathbf{h}}^{(1)} \right)^\top \text{diag} \left( \mathbf{W}_{d,+,\mathbf{m}}^{(2)} \tilde{\mathbf{m}} \odot \mathbf{W}_{d,+,\mathbf{y}}^{(2)} \tilde{\mathbf{y}} \right) \mathbf{W}_{d,+,\mathbf{x}}^{(1)} - \frac{\alpha}{\ell} \mathbf{Q}_{\tilde{\mathbf{h}},\tilde{\mathbf{x}}} \right. \right. \\ & \quad \left. \tilde{\mathbf{m}}^\top \mathbf{Q}_{\tilde{\mathbf{m}},\tilde{\mathbf{y}}} \tilde{\mathbf{y}} \right) \Lambda_{d-1,+,\mathbf{x}}^{(1)} \mathbf{x} \right\|_2 \\ & \leq \left\| \left( \mathbf{W}_{d,+,\mathbf{h}}^{(1)} \right)^\top \text{diag} \left( \mathbf{W}_{d,+,\mathbf{m}}^{(2)} \tilde{\mathbf{m}} \odot \mathbf{W}_{d,+,\mathbf{y}}^{(2)} \tilde{\mathbf{y}} \right) \mathbf{W}_{d,+,\mathbf{x}}^{(1)} - \frac{\alpha}{\ell} \mathbf{Q}_{\tilde{\mathbf{h}},\tilde{\mathbf{x}}} \tilde{\mathbf{m}}^\top \mathbf{Q}_{\tilde{\mathbf{m}},\tilde{\mathbf{y}}} \tilde{\mathbf{y}} \right\| \\ & \quad \left\| \Lambda_{d-1,+,\mathbf{h}}^{(1)} \right\| \left\| \Lambda_{d-1,+,\mathbf{x}}^{(1)} \right\| \|\mathbf{x}\|_2 \\ & \leq \left( \frac{1+4\epsilon(d-1)}{2^{d-1}} \right) \frac{4\epsilon}{\ell} \|\mathbf{x}\|_2 \|\tilde{\mathbf{m}}\|_2 \|\tilde{\mathbf{y}}\|_2 \\ & = \frac{(1+4\epsilon(d-1))}{2^d} \frac{8\epsilon}{\ell} \|\mathbf{x}\|_2 \|\Lambda_{s-1,+,\mathbf{m}}^{(2)} \mathbf{m}\|_2 \|\Lambda_{s-1,+,\mathbf{y}}^{(2)} \mathbf{y}\|_2 \\ & \leq \frac{(1+4\epsilon(d-1))}{2^d} \frac{(1+4\epsilon(s-1))}{2^s} \frac{16\epsilon}{\ell} \|\mathbf{x}\|_2 \|\mathbf{m}\|_2 \|\mathbf{y}\|_2 \\ & \leq \frac{64\epsilon}{2^{d+s}\ell} \|\mathbf{x}\|_2 \|\mathbf{m}\|_2 \|\mathbf{y}\|_2. \end{aligned} \tag{25}$$

where the first inequality holds because spectral norm is a sub-multiplicative norm. The second inequality holds because of (22) and joint-WDC. The last inequality holds if  $4\epsilon(d-1) < 1$  and  $4\epsilon(s-1) < 1$ .

Second, consider

$$\begin{aligned} & \left\| \tilde{\mathbf{m}}^\top \mathbf{Q}_{\tilde{\mathbf{m}},\tilde{\mathbf{y}}} \tilde{\mathbf{y}} \left( \Lambda_{d-1,+,\mathbf{h}}^{(1)} \right)^\top \mathbf{Q}_{\tilde{\mathbf{h}},\tilde{\mathbf{x}}} \Lambda_{d-1,+,\mathbf{x}}^{(1)} \mathbf{x} - \mathbf{m}^\top \tilde{t}_{\mathbf{m},\mathbf{y}}^{(2)} \left( \Lambda_{d-1,+,\mathbf{h}}^{(1)} \right)^\top \mathbf{Q}_{\tilde{\mathbf{h}},\tilde{\mathbf{x}}} \Lambda_{d-1,+,\mathbf{x}}^{(1)} \mathbf{x} \right\|_2 \\ & = \left\| \left( \tilde{\mathbf{m}}^\top \mathbf{Q}_{\tilde{\mathbf{m}},\tilde{\mathbf{y}}} \tilde{\mathbf{y}} - \mathbf{m}^\top \tilde{t}_{\mathbf{m},\mathbf{y}}^{(2)} \right) \left( \Lambda_{d-1,+,\mathbf{h}}^{(1)} \right)^\top \mathbf{Q}_{\tilde{\mathbf{h}},\tilde{\mathbf{x}}} \Lambda_{d-1,+,\mathbf{x}}^{(1)} \mathbf{x} \right\|_2 \\ & \leq \frac{1+4\epsilon(d-1)}{2^{d-1}} \|\mathbf{Q}_{\tilde{\mathbf{h}},\tilde{\mathbf{x}}}\| \left\| \left( \Lambda_{s-1,+,\mathbf{m}}^{(2)} \right)^\top \mathbf{Q}_{\tilde{\mathbf{m}},\tilde{\mathbf{y}}} \Lambda_{s-1,+,\mathbf{y}}^{(2)} \mathbf{y} - \tilde{t}_{\mathbf{m},\mathbf{y}}^{(2)} \right\|_2 \|\mathbf{x}\|_2 \|\mathbf{m}\|_2 \\ & = \frac{1+4\epsilon(d-1)}{2^d} \left\| \left( \Lambda_{s-1,+,\mathbf{m}}^{(2)} \right)^\top \left( \mathbf{Q}_{\tilde{\mathbf{m}},\tilde{\mathbf{y}}} - \frac{1}{\alpha_2} \left( \mathbf{W}_{d,+,\mathbf{m}}^{(2)} \right)^\top \mathbf{W}_{d,+,\mathbf{y}}^{(2)} \right) \Lambda_{s-1,+,\mathbf{y}}^{(2)} \mathbf{y} \right. \\ & \quad \left. + \frac{1}{\alpha_2} \left( \Lambda_{s,+,\mathbf{m}}^{(2)} \right)^\top \Lambda_{s,+,\mathbf{y}}^{(2)} \mathbf{y} - \tilde{t}_{\mathbf{m},\mathbf{y}}^{(2)} \right\|_2 \|\mathbf{x}\|_2 \|\mathbf{m}\|_2 \\ & \leq \frac{1+4\epsilon(d-1)}{2^d} \left( \left\| \Lambda_{s-1,+,\mathbf{m}}^{(2)} \right\| \left\| \Lambda_{s-1,+,\mathbf{y}}^{(2)} \right\| \left\| \left( \mathbf{Q}_{\tilde{\mathbf{m}},\tilde{\mathbf{y}}} - \frac{1}{\alpha_2} \left( \mathbf{W}_{d,+,\mathbf{m}}^{(2)} \right)^\top \mathbf{W}_{d,+,\mathbf{y}}^{(2)} \right) \right\| \|\mathbf{y}\|_2 \right) \end{aligned}$$

$$\begin{aligned}
& + \left\| \frac{1}{\alpha_2} \left( \Lambda_{s,+,\mathbf{m}}^{(2)} \right)^\top \Lambda_{s,+,\mathbf{y}}^{(2)} \mathbf{y} - \tilde{t}_{\mathbf{m},\mathbf{y}}^{(2)} \right\|_2 \|\mathbf{x}\|_2 \|\mathbf{m}\|_2 \\
& \leq \frac{1+4\epsilon(d-1)}{2^{d+s}} \left( 6(1+4\epsilon(s-1))\epsilon/\alpha_2 + 24s^3\sqrt{\epsilon/\alpha_2} \right) \|\mathbf{x}\|_2 \|\mathbf{m}\|_2 \|\mathbf{y}\|_2 \\
& \leq \frac{2}{2^{d+s}} \left( 12\epsilon/\alpha_2 + 24s^3\sqrt{\epsilon/\alpha_2} \right) \|\mathbf{x}\|_2 \|\mathbf{m}\|_2 \|\mathbf{y}\|_2 \\
& \leq \frac{72s^3\sqrt{\epsilon}}{2^{2d}\alpha_2} \|\mathbf{x}\|_2 \|\mathbf{m}\|_2 \|\mathbf{y}\|_2.
\end{aligned} \tag{26}$$

where the first inequality holds because of (22). The second inequality holds because of triangle inequality. The third inequality holds because of (23),  $\frac{1}{\sqrt{\alpha_2}}\mathbf{W}_d^{(2)}$  satisfy WDC with constant  $\epsilon/\alpha_2$  and 1, and Lemma 1.

Third, consider

$$\begin{aligned}
& \left\| \mathbf{m}^\top \tilde{t}_{\mathbf{m},\mathbf{y}}^{(2)} \left( \Lambda_{d-1,+,\mathbf{h}}^{(1)} \right)^\top \mathbf{Q}_{\tilde{\mathbf{h}},\tilde{\mathbf{x}}} \Lambda_{d-1,+,\mathbf{x}}^{(1)} \mathbf{x} - \mathbf{m}^\top \tilde{t}_{\mathbf{m},\mathbf{y}}^{(2)} \tilde{t}_{\mathbf{h},\mathbf{x}}^{(1)} \right\|_2 \\
& = \left\| \mathbf{m}^\top \tilde{t}_{\mathbf{m},\mathbf{y}}^{(2)} \left\| \left( \Lambda_{d-1,+,\mathbf{h}}^{(1)} \right)^\top \mathbf{Q}_{\tilde{\mathbf{h}},\tilde{\mathbf{x}}} \Lambda_{d-1,+,\mathbf{x}}^{(1)} \mathbf{x} - \tilde{t}_{\mathbf{h},\mathbf{x}}^{(1)} \right\|_2 \right. \\
& \leq \left\| \tilde{t}_{\mathbf{m},\mathbf{y}}^{(2)} \right\|_2 \left\| \left( \Lambda_{d-1,+,\mathbf{h}}^{(1)} \right)^\top \left( \mathbf{Q}_{\tilde{\mathbf{h}},\tilde{\mathbf{x}}} - \frac{1}{\alpha_1} \left( \mathbf{W}_{d,+,\mathbf{h}}^{(1)} \right)^\top \mathbf{W}_{d,+,\mathbf{x}}^{(1)} \right) \Lambda_{d-1,+,\mathbf{x}}^{(1)} \mathbf{x} \right. \\
& \quad \left. + \frac{1}{\alpha_1} \left( \Lambda_{d,+,\mathbf{h}}^{(1)} \right)^\top \Lambda_{d,+,\mathbf{x}}^{(1)} - \tilde{t}_{\mathbf{h},\mathbf{x}}^{(1)} \right\|_2 \|\mathbf{m}\|_2 \\
& \leq \frac{1+s}{2^s} \left( \left\| \Lambda_{d-1,+,\mathbf{h}}^{(1)} \right\| \left\| \Lambda_{d-1,+,\mathbf{x}}^{(1)} \right\| \left\| \mathbf{Q}_{\tilde{\mathbf{h}},\tilde{\mathbf{x}}} - \frac{1}{\alpha_1} \left( \mathbf{W}_{d,+,\mathbf{h}}^{(1)} \right)^\top \mathbf{W}_{d,+,\mathbf{x}}^{(1)} \right\| \|\mathbf{x}\|_2 \right. \\
& \quad \left. + \left\| \frac{1}{\alpha_1} \left( \Lambda_{d,+,\mathbf{h}}^{(1)} \right)^\top \Lambda_{d,+,\mathbf{x}}^{(1)} - \tilde{t}_{\mathbf{h},\mathbf{x}}^{(1)} \right\|_2 \right) \|\mathbf{m}\|_2 \|\mathbf{y}\|_2 \\
& \leq \frac{2s}{2^{d+s}} \left( 6(1+4\epsilon(d-1))\epsilon/\alpha_1 + 24d^3\sqrt{\epsilon/\alpha_1} \right) \|\mathbf{x}\|_2 \|\mathbf{m}\|_2 \|\mathbf{y}\|_2 \\
& \leq \frac{72sd^3\sqrt{\epsilon}}{2^{d+s}\alpha_1} \|\mathbf{x}\|_2 \|\mathbf{m}\|_2 \|\mathbf{y}\|_2.
\end{aligned} \tag{27}$$

where the first inequality holds because of Cauchy-Schwartz inequality. The second inequality holds because of triangle inequality along with  $\|\tilde{t}_{\mathbf{m},\mathbf{y}}^{(2)}\|_2 \leq \frac{1+s}{2^s}$ . The third inequality holds because of (22),  $\frac{1}{\sqrt{\alpha_1}}\mathbf{W}_d^{(1)}$  satisfy WDC with constant  $\epsilon/\alpha_1$  and 1, and Lemma 1.

Hence, combining (24), (25), (26), and (27), we get

$$\begin{aligned}
& \left\| \left( \Lambda_{d,+,\mathbf{h}}^{(1)} \right)^\top \text{diag} \left( \Lambda_{s,+,\mathbf{m}}^{(2)} \mathbf{m} \odot \Lambda_{s,+,\mathbf{y}}^{(2)} \mathbf{y} \right) \Lambda_{d,+,\mathbf{x}}^{(1)} \mathbf{x} - \frac{\alpha_1\alpha_2}{\ell} \left( \mathbf{m}^\top \tilde{t}_{\mathbf{m},\mathbf{y}}^{(2)} \right) \tilde{t}_{\mathbf{h},\mathbf{x}}^{(1)} \right\|_2 \\
& \leq \left( \frac{64\epsilon}{2^{d+s}\ell} + \frac{72s^3\sqrt{\epsilon}}{2^{d+s}\ell} + \frac{72sd^3\sqrt{\epsilon}}{2^{2d}\ell} \right) \|\mathbf{x}\|_2 \|\mathbf{m}\|_2 \|\mathbf{y}\|_2 \\
& \leq \frac{208s^3d^3\sqrt{\epsilon}}{2^{d+s}\ell} \|\mathbf{x}\|_2 \|\mathbf{m}\|_2 \|\mathbf{y}\|_2.
\end{aligned}$$

□

### 5.3 Zeros of $t_{(\mathbf{h},\mathbf{m}),(\mathbf{h}_0,\mathbf{m}_0)}$

**Lemma 3.** Fix  $0 < \epsilon < 1$  and  $0 < \alpha \leq 1$  such that  $38(d^5 + s^5)\sqrt{\epsilon}/\alpha < 1$ . Let  $\mathcal{K} = \{(\mathbf{h}, \mathbf{0}) \in \mathbb{R}^{n \times p} \mid \mathbf{h} \in \mathbb{R}^n\} \cup \{(\mathbf{0}, \mathbf{m}) \in \mathbb{R}^{n \times p} \mid \mathbf{m} \in \mathbb{R}^p\}$ . Let

$$\begin{aligned}
S_{\epsilon,(\mathbf{h}_0,\mathbf{m}_0)}^{(1)} &= \left\{ (\mathbf{h}, \mathbf{x}) \in \mathbb{R}^{n \times p} \setminus \mathcal{K} \mid \frac{\|t_{(\mathbf{h},\mathbf{m}),(\mathbf{h}_0,\mathbf{m}_0)}^{(1)}\|_2}{\|\mathbf{m}\|_2} \leq \frac{\epsilon \max(\|\mathbf{h}\|_2 \|\mathbf{m}\|_2, \|\mathbf{h}_0\|_2 \|\mathbf{m}_0\|_2)}{2^{d+s}\ell} \right\}, \\
S_{\epsilon,(\mathbf{h}_0,\mathbf{m}_0)}^{(2)} &= \left\{ (\mathbf{h}, \mathbf{x}) \in \mathbb{R}^{n \times p} \setminus \mathcal{K} \mid \frac{\|t_{(\mathbf{h},\mathbf{m}),(\mathbf{h}_0,\mathbf{m}_0)}^{(2)}\|_2}{\|\mathbf{h}\|_2} \leq \frac{\epsilon \max(\|\mathbf{h}\|_2 \|\mathbf{m}\|_2, \|\mathbf{h}_0\|_2 \|\mathbf{m}_0\|_2)}{2^{d+s}\ell} \right\},
\end{aligned}$$

where  $d$  and  $s$  are integers greater than 1. Let

$$\tilde{\mathbf{t}}_{\mathbf{m}, \mathbf{y}}^{(k)} = \frac{1}{2^{a^{(k)}}} \left( \prod_{i=0}^{a^{(k)}-1} \frac{\pi - \bar{\theta}_i^{(k)}}{\pi} \mathbf{y} + \sum_{i=0}^{a^{(k)}-1} \frac{\sin \bar{\theta}_i^{(k)}}{\pi} \left( \prod_{j=i+1}^{a^{(k)}-1} \frac{\pi - \bar{\theta}_j^{(k)}}{\pi} \right) \frac{\|\mathbf{y}\|_2}{\|\mathbf{m}\|_2} \mathbf{m} \right), \quad (28)$$

where  $\bar{\theta}_i^{(k)} = g(\bar{\theta}_{i-1}^{(k)})$  for  $g$  given in (8),  $\bar{\theta}_0^{(k)} = \angle(\mathbf{m}, \mathbf{y})$ ,  $a^{(1)} = d$ , and  $a^{(2)} = s$ . Let

$$\mathbf{t}_{(\mathbf{h}, \mathbf{m}), (\mathbf{h}_0, \mathbf{m}_0)} = \begin{bmatrix} \mathbf{t}_{(\mathbf{h}, \mathbf{m}), (\mathbf{h}_0, \mathbf{m}_0)}^{(1)} \\ \mathbf{t}_{(\mathbf{h}, \mathbf{m}), (\mathbf{h}_0, \mathbf{m}_0)}^{(2)} \end{bmatrix} \text{ where}$$

$$\mathbf{t}_{(\mathbf{h}, \mathbf{m}), (\mathbf{h}_0, \mathbf{m}_0)}^{(1)} = \frac{\alpha}{2^{d+s} \ell} \|\mathbf{m}\|_2^2 \mathbf{h} - \frac{\alpha}{\ell} \mathbf{m}^\top \tilde{\mathbf{t}}_{\mathbf{m}, \mathbf{m}_0}^{(2)} \tilde{\mathbf{t}}_{\mathbf{h}, \mathbf{h}_0}^{(1)}, \quad (29)$$

$$\mathbf{t}_{(\mathbf{h}, \mathbf{m}), (\mathbf{h}_0, \mathbf{m}_0)}^{(2)} = \frac{\alpha}{2^{d+s} \ell} \|\mathbf{h}\|_2^2 \mathbf{m} - \frac{\alpha}{\ell} \mathbf{h}^\top \tilde{\mathbf{t}}_{\mathbf{h}, \mathbf{h}_0}^{(1)} \tilde{\mathbf{t}}_{\mathbf{m}, \mathbf{m}_0}^{(2)}. \quad (30)$$

Define

$$\rho_{a^{(k)}}^{(k)} := \sum_{i=1}^{a^{(k)}-1} \frac{\sin \bar{\theta}_i^{(k)}}{\pi} \left( \prod_{j=i+1}^{a^{(k)}-1} \frac{\pi - \bar{\theta}_j^{(k)}}{\pi} \right),$$

where  $\bar{\theta}_0^{(k)} = \pi$  and  $\bar{\theta}_i^{(k)} = g(\bar{\theta}_{i-1}^{(k)})$ . If  $(\mathbf{h}, \mathbf{m}) \in S_{\epsilon, (\mathbf{h}_0, \mathbf{m}_0)}^{(1)} \cap S_{\epsilon, (\mathbf{h}_0, \mathbf{m}_0)}^{(2)}$  then one of the following holds:

- $|\bar{\theta}_0^{(1)}| \leq 2\sqrt{\epsilon}$ ,  $|\bar{\theta}_0^{(2)}| \leq 2\sqrt{\epsilon}$  and

$$|\|\mathbf{h}\|_2 \|\mathbf{m}\|_2 - \|\mathbf{h}_0\|_2 \|\mathbf{m}_0\|_2| \leq 145 \frac{ds\sqrt{\epsilon}}{\alpha} \|\mathbf{h}_0\|_2 \|\mathbf{m}_0\|_2,$$

- $|\bar{\theta}_0^{(1)} - \pi| \leq 12\pi^2 d^3 \sqrt{\epsilon}/\alpha$ ,  $|\bar{\theta}_0^{(2)}| \leq 1.5\sqrt{\epsilon}$  and

$$|\|\mathbf{h}\|_2 \|\mathbf{m}\|_2 - \rho_d^{(1)} \|\mathbf{h}_0\|_2 \|\mathbf{m}_0\|_2| \leq 532 \frac{d^6 s \sqrt{\epsilon}}{\alpha} \|\mathbf{h}_0\|_2 \|\mathbf{m}_0\|_2,$$

- $|\bar{\theta}_0^{(1)}| \leq 2\sqrt{\epsilon}$ ,  $|\bar{\theta}_0^{(2)} - \pi| \leq 12\pi^2 s^3 \sqrt{\epsilon}/\alpha$  and

$$|\|\mathbf{h}\|_2 \|\mathbf{m}\|_2 - \rho_d^{(1)} \|\mathbf{h}_0\|_2 \|\mathbf{m}_0\|_2| \leq 532 \frac{ds^6 \sqrt{\epsilon}}{\alpha} \|\mathbf{h}_0\|_2 \|\mathbf{m}_0\|_2,$$

- $|\bar{\theta}_0^{(1)} - \pi| \leq 12\pi^2 d^3 \sqrt{\epsilon}/\alpha$ ,  $|\bar{\theta}_0^{(2)} - \pi| \leq 12\pi^2 d^3 \sqrt{\epsilon}/\alpha$  and

$$|\|\mathbf{h}\|_2 \|\mathbf{m}\|_2 - \rho_d^{(1)} \rho_s^{(2)} \|\mathbf{h}_0\|_2 \|\mathbf{m}_0\|_2| \leq 3915 \frac{d^6 s^6 \sqrt{\epsilon}}{\alpha} \|\mathbf{h}_0\|_2 \|\mathbf{m}_0\|_2.$$

In particular,

$$\begin{aligned} S_{\epsilon, (\mathbf{h}_0, \mathbf{m}_0)}^{(1)} \cap S_{\epsilon, (\mathbf{h}_0, \mathbf{m}_0)}^{(2)} &\subseteq \mathcal{A}_{437 \frac{ds\sqrt{\epsilon}}{\alpha}, (\mathbf{h}_0, \mathbf{m}_0)} \cup \mathcal{A}_{2436 \pi \frac{d^6 s \sqrt{\epsilon}}{\alpha}, (-\rho_d^{(1)} \mathbf{h}_0, \mathbf{m}_0)} \\ &\quad \cup \mathcal{A}_{2436 \pi \frac{ds^6 \sqrt{\epsilon}}{\alpha}, (\rho_s^{(2)} \mathbf{h}_0, -\mathbf{m}_0)} \cup \mathcal{A}_{16767 \pi^2 \frac{d^6 s^6 \sqrt{\epsilon}}{\alpha}, (-\rho_d^{(1)} \rho_s^{(2)} \mathbf{h}_0, -\mathbf{m}_0)}, \end{aligned}$$

where  $\mathcal{A}_{\epsilon, (\mathbf{h}_0, \mathbf{m}_0)}$  is defined in (3).

*Proof.* Without loss of generality, let  $\mathbf{h}_0 = \mathbf{e}_1$ ,  $\mathbf{m}_0 = \mathbf{e}_1$ ,  $\hat{\mathbf{h}} = \cos \bar{\theta}_0^{(1)} + \sin \bar{\theta}_0^{(1)}$  and  $\hat{\mathbf{m}} = \cos \bar{\theta}_0^{(2)} + \sin \bar{\theta}_0^{(2)}$  for some  $\bar{\theta}_i^{(1)}, \bar{\theta}_0^{(2)} \in [0, \pi]$ . First we introduce some notation for convenience. Let

$$\xi^{(k)} = \prod_{i=0}^{a^{(k)}-1} \frac{\pi - \bar{\theta}_i^{(k)}}{\pi}, \quad \zeta^{(k)} = \sum_{i=1}^{a^{(k)}-1} \frac{\sin \bar{\theta}_i^{(k)}}{\pi} \prod_{j=i+1}^{a^{(k)}-1} \frac{\pi - \bar{\theta}_j^{(k)}}{\pi},$$

$r^{(1)} = \|\mathbf{h}\|_2$ ,  $r^{(2)} = \|\mathbf{m}\|_2$ , and  $M = \max(r^{(1)} r^{(2)}, 1)$ . Using these notation, we can rewrite  $\mathbf{t}_{(\mathbf{h}, \mathbf{m}), (\mathbf{h}_0, \mathbf{m}_0)}^{(1)}$  as

$$\mathbf{t}_{(\mathbf{h}, \mathbf{m}), (\mathbf{h}_0, \mathbf{m}_0)}^{(1)}$$

$$\begin{aligned}
&= \frac{\alpha}{2^{d+s}\ell} \left( \|\mathbf{m}\|_2 \mathbf{h} - \left( \xi^{(2)} \cos \bar{\theta}_0^{(2)} + \zeta^{(2)} \right) \left( \xi^{(1)} \frac{\mathbf{h}_0}{\|\mathbf{h}_0\|_2} + \zeta^{(1)} \frac{\mathbf{h}}{\|\mathbf{h}\|_2} \right) \|\mathbf{h}_0\|_2 \|\mathbf{m}_0\|_2 \right) \|\mathbf{m}\|_2 \\
&= \frac{\alpha}{2^{d+s}\ell} \left( \|\mathbf{m}\|_2 \mathbf{h} - \cos \bar{\theta}_s^{(2)} \left( \xi^{(1)} \frac{\mathbf{h}_0}{\|\mathbf{h}_0\|_2} + \zeta^{(1)} \frac{\mathbf{h}}{\|\mathbf{h}\|_2} \right) \|\mathbf{h}_0\|_2 \|\mathbf{m}_0\|_2 \right) \|\mathbf{m}\|_2 \\
&= \frac{\alpha}{2^{d+s}\ell} \left( r^{(1)} r^{(2)} \left( \cos \bar{\theta}_0^{(1)} \mathbf{e}_1 + \sin \bar{\theta}_0^{(1)} \mathbf{e}_2 \right) \right. \\
&\quad \left. - \cos \bar{\theta}_s^{(2)} \left( \xi^{(1)} \mathbf{e}_1 + \zeta^{(1)} \left( \cos \bar{\theta}_0^{(1)} \mathbf{e}_1 + \sin \bar{\theta}_0^{(1)} \mathbf{e}_2 \right) \right) \right) r^{(2)}.
\end{aligned}$$

By inspecting the components of  $\mathbf{t}_{(\mathbf{h}, \mathbf{m}), (\mathbf{h}_0, \mathbf{m}_0)}^{(1)}$ , we have that  $(\mathbf{h}, \mathbf{m}) \in S_{\epsilon, (\mathbf{h}_0, \mathbf{m}_0)}^{(1)}$  implies

$$\left| r^{(1)} r^{(2)} \cos \bar{\theta}_0^{(1)} - \cos \bar{\theta}_s^{(2)} \left( \xi^{(1)} + \zeta^{(1)} \cos \bar{\theta}_0^{(1)} \right) \right| \leq \frac{\epsilon M}{\alpha} \quad (31)$$

$$\left| r^{(1)} r^{(2)} \sin \bar{\theta}_0^{(1)} - \cos \bar{\theta}_s^{(2)} \zeta^{(1)} \sin \bar{\theta}_0^{(1)} \right| \leq \frac{\epsilon M}{\alpha} \quad (32)$$

Similarly, by inspecting the components of  $\mathbf{t}_{(\mathbf{h}, \mathbf{m}), (\mathbf{h}_0, \mathbf{m}_0)}^{(2)}$ , we have that  $(\mathbf{h}, \mathbf{m}) \in S_{\epsilon, (\mathbf{h}_0, \mathbf{m}_0)}^{(2)}$  implies

$$\left| r^{(1)} r^{(2)} \cos \bar{\theta}_0^{(2)} - \cos \bar{\theta}_d^{(1)} \left( \xi^{(2)} + \zeta^{(2)} \cos \bar{\theta}_0^{(2)} \right) \right| \leq \frac{\epsilon M}{\alpha} \quad (33)$$

$$\left| r^{(1)} r^{(2)} \sin \bar{\theta}_0^{(2)} - \cos \bar{\theta}_d^{(1)} \zeta^{(2)} \sin \bar{\theta}_0^{(2)} \right| \leq \frac{\epsilon M}{\alpha} \quad (34)$$

Now, we record several properties. We have:

$$\bar{\theta}_i^{(k)} \in [0, \pi/2] \text{ for } i \geq 1 \quad (35)$$

$$\bar{\theta}_i^{(k)} \leq \bar{\theta}_{i-1}^{(k)} \text{ for } i \geq 1 \quad (36)$$

$$|\xi^{(k)}| \leq 1 \quad (37)$$

$$\check{\theta}_i^{(k)} \leq \frac{3\pi}{i+3} \text{ for } i \geq 0 \quad (38)$$

$$\check{\theta}_i^{(k)} \geq \frac{\pi}{i+1} \text{ for } i \geq 0 \quad (39)$$

$$\xi^{(k)} = \prod_{i=1}^{a^{(k)}-1} \frac{\pi - \bar{\theta}_i^{(k)}}{\pi} \geq \frac{\pi - \bar{\theta}_0^{(k)}}{\pi} a^{(k)-3} \quad (40)$$

$$\bar{\theta}_0^{(k)} = \pi + O_1(\delta) \Rightarrow \bar{\theta}_i^{(k)} = \check{\theta}_i^{(k)} + O_1(i\delta) \quad (41)$$

$$\bar{\theta}_0^{(k)} = \pi + O_1(\delta) \Rightarrow |\xi^{(k)}| \leq \frac{\delta}{\pi} \quad (42)$$

$$\bar{\theta}_0^{(k)} = \pi + O_1(\delta) \Rightarrow \zeta^{(k)} = \rho_d^{(k)} + O_1(3a^{(k)}\delta) \text{ if } \frac{a^{(k)}\delta}{\pi} \leq 1 \quad (43)$$

$$|\zeta^{(k)}| = |\xi^{(k)} \cos \bar{\theta}_0^{(k)} - \cos \bar{\theta}_{a^{(k)}}^{(k)}| \leq 2 \quad (44)$$

$$\cos \bar{\theta}_i^{(k)} \geq \frac{1}{\pi} \text{ for } i \geq 2 \quad (45)$$

For a proof of (37)-(43), we refer the readers to Lemma 8 of Hand and Voroninski [2017]. Also, we note that (45) follows directly from (44).

We first show that if  $(\mathbf{h}, \mathbf{m}) \in S_{\epsilon, (\mathbf{h}_0, \mathbf{m}_0)}^{(1)}$  then  $r^{(1)} r^{(2)} \leq 6$ , and thus  $M \leq 6$ . Suppose  $r^{(1)} r^{(2)} > 1$ . At least one of the following holds:  $|\sin \bar{\theta}_0^{(1)}| \geq 1/\sqrt{2}$  or  $|\cos \bar{\theta}_0^{(1)}| \geq 1/\sqrt{2}$ . If  $|\sin \bar{\theta}_0^{(1)}| \geq 1/\sqrt{2}$  then (32) implies that  $|r^{(1)} r^{(2)} - \cos \bar{\theta}_s^{(2)} \zeta^{(1)}| \leq \sqrt{2}\epsilon r^{(1)} r^{(2)}/\alpha$ . Using (44), we get  $r^{(1)} r^{(2)} \leq \frac{2}{1-\sqrt{2}\epsilon/\alpha} \leq 4$  if  $\epsilon/\alpha < 1/4$ . If  $|\cos \bar{\theta}_0^{(1)}| \geq 1/\sqrt{2}$ , then (31) implies  $|r^{(1)} r^{(2)} - \cos \bar{\theta}_s^{(2)} \zeta^{(1)}| \leq \sqrt{2}(\epsilon r^{(1)} r^{(2)}/\alpha + |\xi^{(1)}|)$ . Using (37), (44), and  $\epsilon/\alpha < 1/4$ , we get

$r^{(1)}r^{(2)} \leq \frac{\sqrt{2}|\xi^{(k)}| + \cos \bar{\theta}_s^{(2)}\zeta^{(1)}}{1 - \sqrt{2}\epsilon/\alpha} \leq \frac{2+\sqrt{2}}{1-\sqrt{2}\epsilon/\alpha} \leq 6$ . Thus, we have  $(\mathbf{h}, \mathbf{m}) \in S_{\epsilon, (\mathbf{h}_0, \mathbf{m}_0)}^{(1)} \Rightarrow r^{(1)}r^{(2)} \leq 6 \Rightarrow M \leq 6$ . Similarly, we have  $(\mathbf{h}, \mathbf{m}) \in S_{\epsilon, (\mathbf{h}_0, \mathbf{m}_0)}^{(2)} \Rightarrow r^{(1)}r^{(2)} \leq 6 \Rightarrow M \leq 6$ .

Next we establish that we only need to consider the small angle case and the large angle case (i.e.  $\bar{\theta}_0^{(k)} \approx 0$  or  $\pi$ ) if  $(\mathbf{h}, \mathbf{m}) \in S_{\epsilon, (\mathbf{h}_0, \mathbf{m}_0)}^{(1)} \cap S_{\epsilon, (\mathbf{h}_0, \mathbf{m}_0)}^{(2)}$ . Exactly one of the following holds:  $|r^{(1)}r^{(2)} - \cos \bar{\theta}_s^{(2)}\zeta^{(1)}| \geq \sqrt{\epsilon}M/\alpha$  or  $|r^{(1)}r^{(2)} - \cos \bar{\theta}_s^{(2)}\zeta^{(1)}| < \sqrt{\epsilon}M/\alpha$ . If  $|r^{(1)}r^{(2)} - \cos \bar{\theta}_s^{(2)}\zeta^{(1)}| \geq \sqrt{\epsilon}M/\alpha$ , then by (32), we have  $|\sin \bar{\theta}_0^{(1)}| \leq \sqrt{\epsilon}$ . Hence  $\bar{\theta}_0^{(1)} = O_1(2\sqrt{\epsilon})$  or  $\bar{\theta}_0^{(1)} = \pi + O_1(2\sqrt{\epsilon})$ , as  $\epsilon < 1$ . If  $|r^{(1)}r^{(2)} - \cos \bar{\theta}_s^{(2)}\zeta^{(1)}| < \sqrt{\epsilon}M/\alpha$ , then by (31) and (45) we have  $|\xi^{(1)}| \leq 2\pi\sqrt{\epsilon}M/\alpha$ . Using (40), we get  $\bar{\theta}_0^{(1)} = \pi + O_1(2\pi^2d^3\sqrt{\epsilon}M/\alpha)$ . Thus, we only need to consider the small angle case,  $\bar{\theta}_0^{(1)} = O_1(2\sqrt{\epsilon})$  and the large angle case  $\bar{\theta}_0^{(1)} = \pi + O_1(12\pi^2d^3\sqrt{\epsilon}/\alpha)$ , where we have used  $M \leq 6$ . Similarly, we only need to consider the small angle case,  $\bar{\theta}_0^{(2)} = O_1(2\sqrt{\epsilon})$  and the large angle case  $\bar{\theta}_0^{(2)} = \pi + O_1(12\pi^2s^3\sqrt{\epsilon}/\alpha)$ .

**Case 1:**  $\bar{\theta}_0^{(1)} \approx 0$  and  $\bar{\theta}_0^{(2)} \approx 0$ . Assume  $\bar{\theta}_0^{(k)} = O_1(2\sqrt{\epsilon})$ . As  $\bar{\theta}_i^{(k)} \leq \bar{\theta}_0^{(k)} \leq 2\sqrt{\epsilon}$  for all  $i$ , we have  $\xi^{(k)} \geq \left(1 - \frac{2\sqrt{\epsilon}}{\pi}\right)^{a^{(k)}} = 1 + O_1(\frac{4a^{(k)}\sqrt{\epsilon}}{\pi})$  provided  $2a^{(k)}\sqrt{\epsilon} \leq 1/2$ . By (44), we also have  $\zeta^{(k)} = O_1(\frac{a^{(k)}}{\pi}2\sqrt{\epsilon}) = O_1(a^{(k)}\sqrt{\epsilon})$ . By (31), we have

$$\left|r^{(1)}r^{(2)} \cos \bar{\theta}_0^{(1)} - (\xi^{(2)} \cos \bar{\theta}_0^{(2)} + \zeta^{(2)}) (\xi^{(1)} + \zeta^{(1)} \cos \bar{\theta}_0^{(1)})\right| \leq \frac{\epsilon M}{\alpha}$$

where we used  $\cos \bar{\theta}_{a^{(k)}}^{(k)} = \xi^{(k)} \cos \bar{\theta}_0^{(k)} + \zeta^{(k)}$ . As  $\cos \bar{\theta}_0^{(k)} = 1 + O_1((\bar{\theta}_0^{(k)})^2/2) = 1 + O_1(2\epsilon)$ ,

$$\xi^{(2)} \cos \bar{\theta}_0^{(2)} + \zeta^{(2)} = 1 + O_1(8s\epsilon\sqrt{\epsilon} + 4s\sqrt{\epsilon} + 2\epsilon + s\sqrt{\epsilon}) = 1 + O_1(15s\sqrt{\epsilon}),$$

$$\xi^{(1)} + \zeta^{(1)} \cos \bar{\theta}_0^{(1)} = 1 + O_1(4d\sqrt{\epsilon} + 2d\epsilon\sqrt{\epsilon} + d\sqrt{\epsilon}) = 1 + O_1(7d\sqrt{\epsilon}).$$

Thus,

$$r^{(1)}r^{(2)} = 1 + O_1(12\epsilon + 6\epsilon/\alpha + 105ds\epsilon + 7d\sqrt{\epsilon} + 15s\sqrt{\epsilon}) = 1 + O_1(145ds\sqrt{\epsilon}/\alpha). \quad (46)$$

We now show  $(\mathbf{h}, \mathbf{m})$  is close to  $(c\mathbf{h}_0, \frac{1}{c}\mathbf{m}_0)$ , where  $c = \frac{\|\mathbf{m}_0\|_2}{\|\mathbf{m}\|_2}$ . Consider

$$\begin{aligned} & \left\| \mathbf{h} - \frac{\|\mathbf{m}_0\|_2}{\|\mathbf{m}\|_2} \mathbf{h}_0 \right\|_2 \\ & \leq \frac{1}{\|\mathbf{m}\|_2} \left( |\|\mathbf{h}\|_2 \|\mathbf{m}\|_2 - \|\mathbf{m}_0\|_2 \|\mathbf{h}_0\|_2| + (\|\mathbf{h}_0\|_2 \|\mathbf{m}_0\|_2 + |\|\mathbf{h}\|_2 \|\mathbf{m}\|_2 - \|\mathbf{m}_0\|_2 \|\mathbf{h}_0\|_2|) \bar{\theta}_0^{(1)} \right) \\ & \leq \frac{1}{\|\mathbf{m}\|_2} (145ds\sqrt{\epsilon} \|\mathbf{h}_0\|_2 \|\mathbf{m}_0\|_2 / \alpha + (\|\mathbf{h}_0\|_2 \|\mathbf{m}_0\|_2 + 145ds\sqrt{\epsilon} \|\mathbf{h}_0\|_2 \|\mathbf{m}_0\|_2 / \alpha) 2\sqrt{\epsilon}) \\ & \leq 437 \frac{ds\sqrt{\epsilon}}{\alpha} \frac{\|\mathbf{h}_0\|_2 \|\mathbf{m}_0\|_2}{\|\mathbf{m}\|_2}. \end{aligned}$$

Similarly,

$$\left\| \mathbf{m} - \frac{\|\mathbf{m}\|_2}{\|\mathbf{m}_0\|_2} \mathbf{m}_0 \right\|_2 \leq \left( |\|\mathbf{m}\|_2 - \|\mathbf{m}_0\|_2| + (\|\mathbf{m}\|_2 + |\|\mathbf{m}\|_2 - \|\mathbf{m}_0\|_2|) \bar{\theta}_0^{(2)} \right) \leq 2\sqrt{\epsilon} \|\mathbf{m}\|_2.$$

Hence,

$$\left\| (\mathbf{h}, \mathbf{m}) - \left( c\mathbf{h}_0, \frac{1}{c}\mathbf{m}_0 \right) \right\|_2 \leq 437 \frac{ds\sqrt{\epsilon}}{\alpha} \left\| \left( c\mathbf{h}_0, \frac{1}{c}\mathbf{m}_0 \right) \right\|_2.$$

**Case 2:**  $\bar{\theta}_0^{(1)} \approx \pi$  and  $\bar{\theta}_0^{(2)} \approx 0$ . Assume  $\bar{\theta}_0^{(1)} = \pi + O_1(\delta)$  where  $\delta = 12\pi^2d^3\sqrt{\epsilon}/\alpha$ . By (42) and (43), we have  $\xi^{(1)} = O_1(\delta/\pi)$ , and we have  $\zeta^{(1)} = \rho_d^{(1)} + O_1(3d^3\delta)$  if  $38d^5\sqrt{\epsilon}/\alpha \leq 1$ . Also, assume

$\bar{\theta}_0^{(2)} = O_1(2\sqrt{\epsilon})$ . As  $\bar{\theta}_i^{(2)} \leq \bar{\theta}_0^{(2)} \leq 2\sqrt{\epsilon}$  for all  $i$ , we have  $\xi^{(2)} \geq \left(1 - \frac{2\sqrt{\epsilon}}{\pi}\right)^s = 1 + O_1(\frac{4s\sqrt{\epsilon}}{\pi})$  provided  $2s\sqrt{\epsilon} \leq 1/2$ . By (44), we also have  $\zeta^{(2)} = O_1(\frac{s}{\pi}2\sqrt{\epsilon}) = O_1(s\sqrt{\epsilon})$ . By (33), we have

$$\left| r^{(1)}r^{(2)} \cos \bar{\theta}_0^{(2)} - \left( \xi^{(1)} \cos \bar{\theta}_0^{(1)} + \zeta^{(1)} \right) \left( \xi^{(2)} + \zeta^{(2)} \cos \bar{\theta}_0^{(2)} \right) \right| \leq \frac{\epsilon M}{\alpha}$$

where we used  $\cos \bar{\theta}_{a^{(k)}}^{(k)} = \xi^{(k)} \cos \bar{\theta}_0^{(k)} + \zeta^{(k)}$ . As  $\cos \bar{\theta}_0^{(1)} = -1 + O((\bar{\theta}_0^{(1)} - \pi)^2/2) = -1 + O_1(\delta^2/2)$  provided  $\delta < 1$  and  $\cos \bar{\theta}_0^{(2)} = 1 + O_1((\bar{\theta}_0^{(2)})^2/2) = 1 + O_1(2\epsilon)$ ,

$$\begin{aligned} \xi^{(1)} \cos \bar{\theta}_0^{(1)} + \zeta^{(1)} &= \rho_d^{(1)} + O_1\left(\frac{\delta^3}{2\pi} + \frac{\delta}{\pi} + 3d^3\delta\right) = \rho_d^{(1)} + O_1(4\delta d^3), \\ \xi^{(2)} + \zeta^{(2)} \cos \bar{\theta}_0^{(2)} &= 1 + O_1(4s\sqrt{\epsilon} + 2s\epsilon\sqrt{\epsilon} + s\sqrt{\epsilon}) = 1 + O_1(7s\sqrt{\epsilon}). \end{aligned}$$

Thus,

$$\begin{aligned} r^{(1)}r^{(2)} &= \rho_d^{(1)} + O_1(12\epsilon + 6\epsilon/\alpha + 4\delta d^3 + 7s\sqrt{\epsilon} + 28d^3s\delta\sqrt{\epsilon}) \\ &= \rho_d^{(1)} + O_1(30d\sqrt{\epsilon}/\alpha + 4\delta d^3 + 28d^3s\sqrt{\epsilon}) \\ &= \rho_d^{(1)} + O_1(532d^6s\sqrt{\epsilon}/\alpha). \end{aligned}$$

where, in the second equality, we use  $\delta < 1$ . We now show  $(\mathbf{h}, \mathbf{m})$  is close to  $(-\mathbf{c}\rho_d^{(1)}\mathbf{h}_0, \frac{1}{c}\mathbf{m}_0)$ , where  $c = \frac{\|\mathbf{m}_0\|_2}{\|\mathbf{m}\|_2}$ . Consider

$$\begin{aligned} &\left\| \mathbf{h} + \frac{\|\mathbf{m}_0\|_2}{\|\mathbf{m}\|_2} \rho_d^{(1)} \mathbf{h}_0 \right\|_2 \\ &\leq \frac{1}{\|\mathbf{m}\|_2} \left( \left| \|\mathbf{h}\|_2 \|\mathbf{m}\|_2 - \rho_d^{(1)} \|\mathbf{h}_0\|_2 \|\mathbf{m}_0\|_2 \right| + \left( \rho_d^{(1)} \|\mathbf{h}_0\|_2 \|\mathbf{m}_0\|_2 + \|\mathbf{h}\|_2 \|\mathbf{m}\|_2 \right. \right. \\ &\quad \left. \left. - \rho_d^{(1)} \|\mathbf{h}_0\|_2 \|\mathbf{m}_0\|_2 \right) \bar{\theta}_0^{(1)} \right) \\ &\leq \frac{1}{\|\mathbf{m}\|_2} \left( 532d^6s\sqrt{\epsilon} \|\mathbf{m}_0\|_2 \|\mathbf{h}_0\|_2 / \alpha + (2\|\mathbf{h}_0\|_2 \|\mathbf{m}_0\|_2 + 532d^6s\sqrt{\epsilon} \|\mathbf{m}_0\|_2 \|\mathbf{h}_0\|_2 / \alpha) 119d^3\sqrt{\epsilon}/\alpha \right) \\ &\leq \frac{1}{\|\mathbf{m}\|_2} \left( 532d^6s\sqrt{\epsilon} \|\mathbf{m}_0\|_2 \|\mathbf{h}_0\|_2 / \alpha + (2\|\mathbf{h}_0\|_2 \|\mathbf{m}_0\|_2 + 14s \|\mathbf{m}_0\|_2 \|\mathbf{h}_0\|_2) 119d^3\sqrt{\epsilon}/\alpha \right) \\ &\leq 2436\pi \frac{d^6s\sqrt{\epsilon}}{\alpha} \rho_d^{(1)} \frac{\|\mathbf{h}_0\|_2 \|\mathbf{m}_0\|_2}{\|\mathbf{m}\|_2}. \end{aligned}$$

Similarly,

$$\left\| \mathbf{m} - \frac{\|\mathbf{m}\|_2}{\|\mathbf{m}_0\|_2} \mathbf{m}_0 \right\|_2 \leq \left( \|\mathbf{m}\|_2 - \|\mathbf{m}_0\|_2 + (\|\mathbf{m}\|_2 + \|\mathbf{m}_0\|_2 - \|\mathbf{m}\|_2) \bar{\theta}_0^{(2)} \right) \leq 2\sqrt{\epsilon} \|\mathbf{m}\|_2.$$

Hence,

$$\left\| (\mathbf{h}, \mathbf{m}) - \left( -\mathbf{c}\rho_d^{(1)}\mathbf{h}_0, \frac{1}{c}\mathbf{m}_0 \right) \right\|_2 \leq 2436\pi \frac{d^6s\sqrt{\epsilon}}{\alpha} \left\| \left( -\mathbf{c}\rho_d^{(1)}\mathbf{h}_0, \frac{1}{c}\mathbf{m}_0 \right) \right\|_2.$$

**Case 3:**  $\bar{\theta}_0^{(1)} \approx 0$  and  $\bar{\theta}_0^{(2)} \approx \pi$ . The analysis is similar to case 2. Using (31), we get

$$r^{(1)}r^{(2)} = \rho_s^{(2)} + O_1(532ds^6\sqrt{\epsilon}/\alpha).$$

Again, similar to case 2, we can show  $(\mathbf{h}, \mathbf{m})$  is close to  $(\mathbf{c}\rho_s^{(2)}\mathbf{h}_0, -\frac{1}{c}\mathbf{m}_0)$ , where  $c = \frac{\|\mathbf{h}_0\|_2}{\|\mathbf{h}\|_2}$ . We get,

$$\left\| (\mathbf{h}, \mathbf{m}) - \left( \mathbf{c}\rho_s^{(2)}\mathbf{h}_0, -\frac{1}{c}\mathbf{m}_0 \right) \right\|_2 \leq 2436\pi \frac{ds^6\sqrt{\epsilon}}{\alpha} \left\| \left( \mathbf{c}\rho_s^{(2)}\mathbf{h}_0, -\frac{1}{c}\mathbf{m}_0 \right) \right\|_2.$$

**Case 4:**  $\bar{\theta}_0^{(1)} \approx \pi$  and  $\bar{\theta}_0^{(2)} \approx \pi$ . Assume  $\bar{\theta}_0^{(k)} = \pi + O_1(\delta^{(k)})$  where  $\delta^{(k)} = 12\pi^2 a^{(k)3} \sqrt{\epsilon}/\alpha$ . By (42) and (43), we have  $\xi^{(k)} = O_1(\delta^{(k)}/\pi)$ , and we have  $\zeta^{(k)} = \rho_d^{(k)} + O_1(3a^{(k)3}\delta^{(k)})$  if  $\frac{a^{(k)2}\delta}{\pi} \leq 1$ . By (31), we have

$$\left| r^{(1)}r^{(2)} \cos \bar{\theta}_0^{(1)} - \left( \xi^{(2)} \cos \bar{\theta}_0^{(2)} + \zeta^{(2)} \right) \left( \xi^{(1)} + \zeta^{(1)} \cos \bar{\theta}_0^{(1)} \right) \right| \leq \frac{\epsilon M}{\alpha}$$

where we used  $\cos \bar{\theta}_{a^{(k)}}^{(k)} = \xi^{(k)} \cos \bar{\theta}_0^{(k)} + \zeta^{(k)}$ . As  $\cos \bar{\theta}_0^{(k)} = -1 + O((\bar{\theta}_0^{(k)} - \pi)^2/2) = -1 + O_1((\delta^{(k)})^2/2)$ ,

$$\xi^{(2)} \cos \bar{\theta}_0^{(2)} + \zeta^{(2)} = \rho_s^{(2)} + O_1\left(\frac{(\delta^{(2)})^3}{2\pi} + \frac{\delta^{(2)}}{\pi} + 3s^3\delta^{(2)}\right) = \rho_s^{(2)} + O_1(4\delta^{(2)}s^3),$$

$$\xi^{(1)} + \zeta^{(1)} \cos \bar{\theta}_0^{(1)} = -\rho_d^{(1)} + O_1\left(\frac{\delta^{(1)}}{\pi} + \frac{3}{2}d^3(\delta^{(1)})^3 + 3\delta^{(1)}d^3\right) = -\rho_d^{(1)} + O_1(5\delta^{(1)}d^3).$$

Thus,

$$\begin{aligned} r^{(1)}r^{(2)} &= \rho_d^{(1)}\rho_s^{(2)} + O_1(6\epsilon/\alpha + 4(\delta^{(1)})^2 + 4\delta^{(2)}s^3 + 5\delta^{(1)}d^3 + 20\delta^{(1)}\delta^{(2)}d^3s^3) \\ &= \rho_d^{(1)}\rho_s^{(2)} + O_1(6\epsilon/\alpha + 4\delta^{(1)} + 4\delta^{(2)}s^3 + 5\delta^{(1)}d^3 + 20\delta^{(2)}d^3s^3) \\ &= \rho_d^{(1)}\rho_s^{(2)} + O_1(6\epsilon/\alpha + 3909d^6s^6\sqrt{\epsilon}/\alpha) \\ &= \rho_d^{(1)}\rho_s^{(2)} + O_1(3915d^6s^6\sqrt{\epsilon}/\alpha), \end{aligned}$$

where, in the second equality, we used  $\delta^{(1)} \leq \frac{\pi}{d^2} < 1$ . We now show  $(\mathbf{h}, \mathbf{m})$  is close to  $(-\rho_d^{(1)}\rho_s^{(2)}\mathbf{h}_0, -\frac{1}{c}\mathbf{m}_0)$ , where  $c = \frac{\|\mathbf{m}_0\|_2}{\|\mathbf{m}\|_2}$ . Consider

$$\begin{aligned} &\left\| \mathbf{h} + \frac{\|\mathbf{m}_0\|_2}{\|\mathbf{m}\|_2} \rho_d^{(1)}\rho_s^{(2)} \mathbf{h}_0 \right\|_2 \\ &\leq \frac{1}{\|\mathbf{m}\|_2} \left( \left| \|\mathbf{h}\|_2 \|\mathbf{m}\|_2 - \rho_d^{(1)}\rho_s^{(2)} \|\mathbf{m}_0\|_2 \|\mathbf{h}_0\|_2 \right| + \left( \rho_d^{(1)}\rho_s^{(2)} \|\mathbf{h}_0\|_2 \|\mathbf{m}_0\|_2 + \right. \right. \\ &\quad \left. \left. \left| \|\mathbf{h}\|_2 \|\mathbf{m}\|_2 - \rho_d^{(1)}\rho_s^{(2)} \|\mathbf{m}_0\|_2 \|\mathbf{h}_0\|_2 \right| \right) \bar{\theta}_0^{(1)} \right) \\ &\leq \frac{1}{\|\mathbf{m}\|_2} \left( 3915d^6s^6\sqrt{\epsilon} \|\mathbf{h}_0\|_2 \|\mathbf{m}_0\|_2 / \alpha + (4\|\mathbf{h}_0\|_2 \|\mathbf{m}_0\|_2 \right. \\ &\quad \left. + 3915d^6s^6\sqrt{\epsilon} \|\mathbf{h}_0\|_2 \|\mathbf{m}_0\|_2 / \alpha) 119d^3\sqrt{\epsilon}/\alpha \right) \\ &\leq \frac{1}{\|\mathbf{m}\|_2} \left( 3915d^6s^6\sqrt{\epsilon} \|\mathbf{h}_0\|_2 \|\mathbf{m}_0\|_2 / \alpha + (4\|\mathbf{h}_0\|_2 \|\mathbf{m}_0\|_2 + 104ds^6 \|\mathbf{h}_0\|_2 \|\mathbf{m}_0\|_2) 119d^3\sqrt{\epsilon}/\alpha \right) \\ &\leq 16767\pi^2 \frac{d^6s^6\sqrt{\epsilon}}{\alpha} \rho_d^{(1)}\rho_s^{(2)} \frac{\|\mathbf{h}_0\|_2 \|\mathbf{m}_0\|_2}{\|\mathbf{m}\|_2}. \end{aligned}$$

Similarly,

$$\left\| \mathbf{m} + \frac{\|\mathbf{m}\|_2}{\|\mathbf{m}_0\|_2} \mathbf{m}_0 \right\|_2 \leq \left( (\|\mathbf{m}\|_2 - \|\mathbf{m}_0\|_2) + (\|\mathbf{m}\|_2 + \|\mathbf{m}_0\|_2 - \|\mathbf{m}\|_2) \bar{\theta}_0^{(2)} \right) \leq 119s^3\sqrt{\epsilon} \|\mathbf{m}\|_2 / \alpha.$$

Hence,

$$\left\| (\mathbf{h}, \mathbf{m}) - \left( -\rho_d^{(1)}\rho_s^{(2)}\mathbf{h}_0, -\frac{1}{c}\mathbf{m}_0 \right) \right\|_2 \leq 16767\pi^2 \frac{d^6s^6\sqrt{\epsilon}}{\alpha} \left\| \left( c\rho_d^{(1)}\rho_s^{(2)}\mathbf{h}_0, \frac{1}{c}\mathbf{m}_0 \right) \right\|_2.$$

□

## 5.4 Proof of WDC condition

We first state a lemma that shows that the weight  $\mathbf{W} \in \mathbb{R}^{\ell \times n}$  of a layer of a neural network layer with i.i.d.  $\mathcal{N}(0, 1/\ell)$  entries satisfies the WDC with constant  $\epsilon$  and 1, and we refer the readers to Hand and Voroninski [2017] for a proof of the lemma.

**Lemma 4.** Fix  $0 < \epsilon < 1$ . Let  $\mathbf{W} \in \mathbb{R}^{\ell \times n}$  have i.i.d.  $\mathcal{N}(0, 1/\ell)$  entries. If  $\ell > cn \log n$ , then with probability at least  $1 - 8\ell e^{-\gamma n}$ ,  $\mathbf{W}$  satisfies the WDC with constant  $\epsilon$  and 1. Here  $c, \gamma^{-1}$  are constants that depend only polynomially on  $\epsilon^{-1}$ .

We now state a lemma similar to Lemma 4 which applies to truncated random variable. The proof follows the proof of lemma 4 in Hand and Voroninski [2017].

**Lemma 5.** Fix  $0 < \epsilon < 1$ . Let  $\mathbf{W} \in \mathbb{R}^{\ell \times n}$  where  $i$ th row of  $\mathbf{W}$  satisfy  $\mathbf{w}_i^\top = \mathbf{w}^\top \cdot \mathbf{1}_{\|\mathbf{w}\|_2 \leq 3\sqrt{n/\ell}}$  and  $\mathbf{w} \sim \mathcal{N}(\mathbf{0}, \frac{1}{\ell} \mathbf{I}_n)$ . If  $\ell > cn \log n$ , then with probability at least  $1 - 8ne^{-\gamma n}$ ,  $\mathbf{W}$  satisfies the WDC with constant  $\epsilon$  and  $\alpha$ . Here  $c, \gamma^{-1}$  are constants that depend only polynomially on  $\epsilon^{-1}$  and

$$\alpha = \frac{\Gamma\left(\frac{n+2}{2}\right) - \Gamma\left(\frac{n+2}{2}, \frac{9n}{2}\right)}{\Gamma\left(\frac{n+2}{2}\right)}, \quad (47)$$

where  $\Gamma$  is the Gamma function.

The WDC condition with constant  $\epsilon$  and  $\alpha$  can be written as

$$\|\mathbf{W}_{+,x}^\top \mathbf{W}_{+,y} - \alpha \mathbf{Q}_{x,y}\| \leq \epsilon$$

for all nonzero  $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$ . We note that

$$\mathbf{W}_{+,x}^\top \mathbf{W}_{+,y} = \sum_{i=1}^{\ell} \mathbf{1}_{\mathbf{w}_i^\top \mathbf{x}} \mathbf{1}_{\mathbf{w}_i^\top \mathbf{y}} \mathbf{w}_i \mathbf{w}_i^\top$$

and it is not continuous in  $\mathbf{x}$  and  $\mathbf{y}$ . So, we consider an arbitrarily good continuous approximation of  $\mathbf{W}_{+,x}^\top \mathbf{W}_{+,y}$ . Let

$$t_{-\epsilon}(z) = \begin{cases} 0 & z \leq -\epsilon, \\ 1 + \frac{z}{\epsilon} & -\epsilon \leq z \leq 0, \\ 1 & z \geq 0, \end{cases} \quad \text{and} \quad t_\epsilon(z) = \begin{cases} 0 & z \leq 0, \\ \frac{z}{\epsilon} & 0 \leq z \leq \epsilon, \\ 1 & z \geq \epsilon. \end{cases}$$

and define

$$\begin{aligned} H_{-\epsilon}(\mathbf{x}, \mathbf{y}) &:= \sum_{i=1}^{\ell} t_{-\epsilon}(\mathbf{w}_i^\top \mathbf{x}) t_{-\epsilon}(\mathbf{w}_i^\top \mathbf{y}) \mathbf{w}_i \mathbf{w}_i^\top, \\ H_\epsilon(\mathbf{x}, \mathbf{y}) &:= \sum_{i=1}^{\ell} t_\epsilon(\mathbf{w}_i^\top \mathbf{x}) t_\epsilon(\mathbf{w}_i^\top \mathbf{y}) \mathbf{w}_i \mathbf{w}_i^\top. \end{aligned}$$

The proof of Lemma 5 follows from the follow two lemmas. We first provide an upper bound on the singular values of  $H_{-\epsilon}(\mathbf{x}, \mathbf{y})$ .

**Lemma 6.** Fix  $0 < \epsilon < 1$ . Let  $\mathbf{W} \in \mathbb{R}^{\ell \times n}$  where  $i$ th row of  $\mathbf{W}$  satisfy  $\mathbf{w}_i^\top = \mathbf{w}^\top \cdot \mathbf{1}_{\|\mathbf{w}\|_2 \leq 3\sqrt{n/\ell}}$  and  $\mathbf{w} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_n)$ . If  $\ell > cn \log n$ , then with probability at least  $1 - 4ne^{-\gamma n}$ ,

$$\forall (\mathbf{x}, \mathbf{y}) \neq (\mathbf{0}, \mathbf{0}), \quad H_{-\epsilon}(\mathbf{x}, \mathbf{y}) \preceq \alpha \ell \mathbf{Q}_{x,y} + 3\ell \epsilon I_n.$$

Here,  $c$  and  $\gamma^{-1}$  are constants that depend only polynomially on  $\epsilon^{-1}$  and  $\alpha$  is

$$\alpha = \frac{\Gamma\left(\frac{n+2}{2}\right) - \Gamma\left(\frac{n+2}{2}, \frac{9n}{2}\right)}{\Gamma\left(\frac{n+2}{2}\right)}, \quad (48)$$

where  $\Gamma$  is the Gamma function.

*Proof.* First we bound  $\mathbb{E}[H_{-\epsilon}(\mathbf{x}, \mathbf{y})]$  for fixed  $\mathbf{x}, \mathbf{y} \in \mathcal{S}^{n-1}$ . Noting that  $t_{-\epsilon}(z) \leq \mathbf{1}_{z \geq -\epsilon}(z) = \mathbf{1}_{z>0}(z) + \mathbf{1}_{-\epsilon \leq z \leq 0}(z)$ , we have

$$\begin{aligned} &\mathbb{E}[H_{-\epsilon}(\mathbf{x}, \mathbf{y})] \\ &\leq \mathbb{E}\left[\sum_{i=1}^{\ell} \mathbf{1}_{\mathbf{w}_i^\top \mathbf{x} \geq -\epsilon} \mathbf{1}_{\mathbf{w}_i^\top \mathbf{y} \geq -\epsilon} \mathbf{w}_i \mathbf{w}_i^\top\right] \end{aligned}$$

$$\begin{aligned}
&= \ell \mathbb{E} \left[ \mathbf{1}_{\mathbf{w}_i^\top \mathbf{x} \geq -\epsilon} \mathbf{1}_{\mathbf{w}_i^\top \mathbf{y} \geq -\epsilon} \mathbf{w}_i \mathbf{w}_i^\top \right] \\
&= \ell \mathbb{E} \left[ \left( \mathbf{1}_{\mathbf{w}_i^\top \mathbf{x} \geq 0} \mathbf{1}_{\mathbf{w}_i^\top \mathbf{y} \geq 0} \mathbf{w}_i \mathbf{w}_i^\top \right) \right] + 2\ell \mathbb{E} \left[ \left( \mathbf{1}_{-\epsilon \leq \mathbf{w}_i^\top \mathbf{x} \leq 0} \mathbf{w}_i \mathbf{w}_i^\top \right) \right].
\end{aligned}$$

We first note that  $\mathbb{E} \left[ \mathbf{1}_{\mathbf{w}_i^\top \mathbf{x} \geq 0} \mathbf{1}_{\mathbf{w}_i^\top \mathbf{y} \geq 0} \mathbf{b}_i \mathbf{b}_i^\top \right] = \alpha \mathbf{Q}_{\mathbf{x}, \mathbf{y}}$  where  $\alpha$  satisfies  $0.97 < \alpha < 1$ . Also, we have  $\mathbb{E} \left[ \mathbf{1}_{-\epsilon \leq \mathbf{w}_i^\top \mathbf{x} \leq 0} \mathbf{w}_i \mathbf{w}_i^\top \right] \preceq \frac{\epsilon \alpha}{2} \mathbf{I}_n$ . Thus,

$$\begin{aligned}
\mathbb{E} [H_{-\epsilon}(\mathbf{x}, \mathbf{y})] &\preceq \alpha \ell \cdot \mathbf{m}^\top \mathbf{Q}_{\mathbf{x}, \mathbf{y}} \mathbf{y} + \epsilon \alpha \ell \mathbf{I}_n \\
&\preceq \alpha \ell \cdot \mathbf{m}^\top \mathbf{Q}_{\mathbf{x}, \mathbf{y}} + \epsilon \ell \mathbf{I}_n
\end{aligned} \tag{49}$$

Second, we show concentration of  $H_{-\epsilon}(\mathbf{x}, \mathbf{y})$  for fixed  $\mathbf{x}, \mathbf{y} \in \mathcal{S}^{n-1}$ . Let

$$\xi_i = \sqrt{t_{-\epsilon}(\mathbf{w}_i^\top \mathbf{x}) t_{-\epsilon}(\mathbf{w}_i^\top \mathbf{y}) \mathbf{w}_i}.$$

We have

$$\begin{aligned}
H_{-\epsilon}(\mathbf{x}, \mathbf{y}) - \mathbb{E} [H_{-\epsilon}(\mathbf{x}, \mathbf{y})] &= \sum_{i=1}^{\ell} \left( t_{-\epsilon}(\mathbf{w}_i^\top \mathbf{x}) t_{-\epsilon}(\mathbf{w}_i^\top \mathbf{y}) \mathbf{w}_i \mathbf{w}_i^\top - \mathbb{E} [t_{-\epsilon}(\mathbf{w}_i^\top \mathbf{x}) t_{-\epsilon}(\mathbf{w}_i^\top \mathbf{y}) \mathbf{w}_i \mathbf{w}_i^\top] \right) \\
&= \sum_{i=1}^{\ell} (\xi_i \xi_i^\top - \mathbb{E} [\xi_i \xi_i^\top]).
\end{aligned}$$

Note that  $\xi_i$  is sub-Gaussian for all  $i$  and that the sub-Gaussian norm of  $\xi_i$  is bounded from above by an absolute constant which we call  $K$ . By first part of Remark 5.40 in Vershynin [2012], there exists a  $c_K$  and  $\gamma_K$  such that for all  $t \geq 0$ , with probability at least  $1 - 2e^{-\gamma_K t^2}$ ,

$$\|H_{-\epsilon}(\mathbf{x}, \mathbf{y}) - \mathbb{E} [H_{-\epsilon}(\mathbf{x}, \mathbf{y})]\| \leq \max(\delta, \delta^2) \ell, \quad \text{where } \delta = c_K \sqrt{\frac{n}{\ell}} + \frac{t}{\sqrt{\ell}}.$$

If  $\ell > (2c_K/\epsilon)^2 n$ ,  $t = \epsilon \sqrt{\ell}/2$ , and  $\epsilon < 1$ , we have

$$\|H_{-\epsilon}(\mathbf{x}, \mathbf{y}) - \mathbb{E} [H_{-\epsilon}(\mathbf{x}, \mathbf{y})]\| \leq \epsilon \ell \tag{50}$$

with probability at least  $1 - 2e^{-\gamma_K \frac{\epsilon^2 \ell}{4}}$ .

Third, we bound the Lipschitz constant of  $H_{-\epsilon}$ . For  $\tilde{\mathbf{x}}, \tilde{\mathbf{y}} \in \mathbb{R}^n$  we have

$$\begin{aligned}
H_{-\epsilon}(\mathbf{x}, \mathbf{y}) - H_{-\epsilon}(\tilde{\mathbf{x}}, \tilde{\mathbf{y}}) &= \sum_{i=1}^{\ell} \left[ t_{-\epsilon}(\mathbf{w}_i^\top \mathbf{x}) t_{-\epsilon}(\mathbf{w}_i^\top \mathbf{y}) - t_{-\epsilon}(\mathbf{w}_i^\top \tilde{\mathbf{x}}) t_{-\epsilon}(\mathbf{w}_i^\top \tilde{\mathbf{y}}) \right] \mathbf{w}_i \mathbf{w}_i^\top \\
&= \sum_{i=1}^{\ell} \left[ t_{-\epsilon}(\mathbf{w}_i^\top \mathbf{x}) (t_{-\epsilon}(\mathbf{w}_i^\top \mathbf{y}) - t_{-\epsilon}(\mathbf{w}_i^\top \tilde{\mathbf{y}})) \right. \\
&\quad \left. + t_{-\epsilon}(\mathbf{w}_i^\top \tilde{\mathbf{y}}) (t_{-\epsilon}(\mathbf{w}_i^\top \mathbf{x}) - t_{-\epsilon}(\mathbf{w}_i^\top \tilde{\mathbf{x}})) \right] \mathbf{w}_i \mathbf{w}_i^\top \\
&= \mathbf{W}^\top \left[ \text{diag}(t_{-\epsilon}(\mathbf{W}\mathbf{x})) \text{diag}((\mathbf{W}\mathbf{y})_+ - (\mathbf{W}\tilde{\mathbf{y}})_+) \right. \\
&\quad \left. + \text{diag}(t_{-\epsilon}(\mathbf{W}\tilde{\mathbf{y}})) \text{diag}((\mathbf{W}\mathbf{x})_+ - (\mathbf{W}\tilde{\mathbf{x}})_+) \right] \mathbf{W}
\end{aligned}$$

Thus,

$$\begin{aligned}
&\|H_{-\epsilon}(\mathbf{x}, \mathbf{y}) - H_{-\epsilon}(\tilde{\mathbf{x}}, \tilde{\mathbf{y}})\| \\
&\leq \|\mathbf{W}\|^2 \left[ \|t_{-\epsilon}(\mathbf{W}\mathbf{y}) - t_{-\epsilon}(\mathbf{W}\tilde{\mathbf{y}})\|_\infty + \|t_{-\epsilon}(\mathbf{W}\mathbf{x}) - t_{-\epsilon}(\mathbf{W}\tilde{\mathbf{x}})\|_\infty \right] \\
&\leq \|\mathbf{W}\|^2 \left[ \max_{i \in [\ell]} |t_{-\epsilon}(\mathbf{w}_i^\top \mathbf{y}) - t_{-\epsilon}(\mathbf{w}_i^\top \tilde{\mathbf{y}})| + \max_{i \in [\ell]} |t_{-\epsilon}(\mathbf{w}_i^\top \mathbf{x}) - t_{-\epsilon}(\mathbf{w}_i^\top \tilde{\mathbf{x}})| \right]
\end{aligned}$$

$$\begin{aligned}
&\leq \|\mathbf{W}\|^2 \left[ \max_{i \in [\ell]} \frac{1}{\epsilon} |\mathbf{w}_i^\top (\mathbf{y} - \tilde{\mathbf{y}})| + \max_{i \in [\ell]} \frac{1}{\epsilon} |\mathbf{w}_i^\top (\mathbf{x} - \tilde{\mathbf{x}})| \right] \\
&\leq \|\mathbf{W}\|^2 \left[ \frac{1}{\epsilon} \max_{i \in [\ell]} \|\mathbf{w}_i\|_2 \|\mathbf{y} - \tilde{\mathbf{y}}\| + \frac{1}{\epsilon} \max_{i \in [\ell]} \|\mathbf{w}_i\|_2 \|\mathbf{x} - \tilde{\mathbf{x}}\| \right] \\
&\leq \|\mathbf{W}\|^2 \left[ \frac{9}{\epsilon} \sqrt{n} \|\mathbf{x} - \tilde{\mathbf{x}}\| + \frac{9}{\epsilon} \sqrt{n} \|\mathbf{h} - \tilde{\mathbf{h}}\| \right]
\end{aligned}$$

where the first inequality follows because  $|t_{-\epsilon}(z)| \leq 1$  for all  $z$ , and the third inequality follows because  $t_{-\epsilon}(z)$  is  $1/\epsilon$ -Lipschitz. Let  $E_1$  be the event that  $\|\mathbf{W}\| \leq 3\sqrt{\ell}$ . By Corollary 5.35 in Vershynin [2012], for  $\mathbf{A} \in \mathbb{R}^{\ell \times n}$  with rows of  $\mathbf{A}$  following  $\mathcal{N}(\mathbf{0}, \mathbf{I}_n)$ , we have  $\mathbb{P}(\|\mathbf{A}\| \leq 3\sqrt{\ell}) \geq 1 - 2e^{-\ell/2}$ , if  $\ell \geq n$ . As rows of  $\mathbf{W}$  are truncated, we have  $\mathbb{P}(E_1) \geq 1 - 2e^{-\ell/2}$ , if  $\ell \geq n$  as well. On  $E_1$ , we have

$$\begin{aligned}
&\|H_{-\epsilon}(\mathbf{x}, \mathbf{y}) - H_{-\epsilon}(\tilde{\mathbf{x}}, \tilde{\mathbf{y}})\| \\
&\leq \frac{27\ell\sqrt{n}}{\epsilon} [\|\mathbf{x} - \tilde{\mathbf{x}}\| + \|\mathbf{y} - \tilde{\mathbf{y}}\|]
\end{aligned} \tag{51}$$

for all  $\mathbf{x}, \mathbf{y}, \tilde{\mathbf{x}}, \tilde{\mathbf{y}} \in \mathcal{S}^{n-1}$ .

Finally, we complete the proof by a covering argument. Let  $\mathcal{N}_\delta$  be a  $\delta$ -net on  $\mathcal{S}^{n-1}$  such that  $|\mathcal{N}_\delta| \leq (3/\delta)^n$ . Take  $\delta = \frac{\epsilon^2}{54\sqrt{n}}$ . Combining (49) and (54), we have

$$\begin{aligned}
\forall (\mathbf{x}, \mathbf{y}), \in \mathcal{N}_\delta, \quad H_{-\epsilon}(\mathbf{x}, \mathbf{y}) &\preceq \mathbb{E} H_{-\epsilon}(\mathbf{x}, \mathbf{y}) + \ell\epsilon I_n \\
&\preceq \alpha\ell Q_{\mathbf{x}, \mathbf{y}} + 2\ell\epsilon I_n.
\end{aligned}$$

with probability at least

$$1 - 2|\mathcal{N}_\delta| e^{-\gamma_K \epsilon^2 \ell/4} \geq 1 - 2 \left( \frac{3}{\delta} \right)^n e^{-\gamma_K \epsilon^2 \ell/4} \geq 1 - 2e^{-\gamma_K \epsilon^2 \ell/4 + n \log(3 \cdot 54\sqrt{n}/\epsilon^2)}.$$

If  $\ell \geq \tilde{c}n \log(n)$  for some  $\tilde{c} = \Omega(\epsilon^2 \log \epsilon)$ , then this probability is at least  $1 - 2e^{-\tilde{\gamma}\ell}$  for some  $\tilde{\gamma} = O(\epsilon^2)$ . For  $\mathbf{x}, \mathbf{y} \in \mathcal{S}^{n-1}$ , let  $\tilde{\mathbf{x}}, \tilde{\mathbf{y}} \in \mathcal{N}_\delta$  be such that  $\|\mathbf{x} - \tilde{\mathbf{x}}\|_2 \leq \delta$ , and  $\|\mathbf{y} - \tilde{\mathbf{y}}\|_2 \leq \delta$ . By (51), we have that

$$\begin{aligned}
\forall \mathbf{x}, \mathbf{y} \neq \mathbf{0}, \quad H_{-\epsilon}(\mathbf{x}, \mathbf{y}) &\preceq H_{-\epsilon}(\tilde{\mathbf{x}}, \tilde{\mathbf{y}}) + \frac{27\ell\sqrt{n}}{\epsilon} 2\delta I_n \\
&\preceq \alpha\ell Q_{\mathbf{x}, \mathbf{y}} + 3\ell\epsilon I_n.
\end{aligned}$$

In conclusion, the result of this lemma holds if  $\ell > (2c_K/\epsilon)^2 n$  and  $\ell \geq \tilde{c}(n) \log n$ , with probability at least  $1 - 2e^{-\gamma_K \epsilon^2 \ell/4} - 2e^{-\ell/2} - 2e^{-\tilde{\gamma}\ell} > 1 - 6e^{-\gamma\ell}$  for some  $\gamma = O(\epsilon^2)$  and  $\tilde{c} = \Omega(\epsilon^2 \log \epsilon)$ .  $\square$

Next, we now provide an upper bound on the singular values of  $G_\epsilon(\mathbf{h}, \mathbf{x}, \mathbf{m}, \mathbf{y})$ .

**Lemma 7.** Fix  $0 < \epsilon < 1$ . Let  $\mathbf{W} \in \mathbb{R}^{\ell \times n}$  where  $i$ th row of  $\mathbf{W}$  satisfy  $\mathbf{w}_i^\top = \mathbf{w}^\top \cdot \mathbf{1}_{\|\mathbf{w}\|_2 \leq 3\sqrt{n}}$  and  $\mathbf{w} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_n)$ . If  $\ell > cn \log n$ , then with probability at least  $1 - 4ne^{-\gamma n}$ ,

$$\forall (\mathbf{x}, \mathbf{y}) \neq (\mathbf{0}, \mathbf{0}), \quad H_\epsilon(\mathbf{x}, \mathbf{y}) \succeq \alpha\ell Q_{\mathbf{x}, \mathbf{y}} - 3\ell\epsilon I_n.$$

Here,  $c$  and  $\gamma^{-1}$  are constants that depend only polynomially on  $\epsilon^{-1}$  and  $\alpha$  is

$$\alpha = \frac{\Gamma\left(\frac{n+2}{2}\right) - \Gamma\left(\frac{n+2}{2}, \frac{9n}{2}\right)}{\Gamma\left(\frac{n+2}{2}\right)}, \tag{52}$$

where  $\Gamma$  is the Gamma function.

*Proof.* First we bound  $\mathbb{E}[H_{-\epsilon}(\mathbf{x}, \mathbf{y})]$  for fixed  $\mathbf{x}, \mathbf{y} \in \mathcal{S}^{n-1}$ . Noting that  $t_\epsilon(z) \geq \mathbf{1}_{z>0}(z) - \mathbf{1}_{-\epsilon \leq z \leq 0}(z)$  for all  $z$ , we have

$$\mathbb{E}[H_\epsilon(\mathbf{x}, \mathbf{y})]$$

$$\succeq \ell \mathbb{E} \left[ \left( \mathbf{1}_{\mathbf{w}_i^\top \mathbf{x} \geq 0} \mathbf{1}_{\mathbf{w}_i^\top \mathbf{y} \geq 0} \mathbf{w}_i \mathbf{w}_i^\top \right) \right] - 2\ell \mathbb{E} \left[ \left( \mathbf{1}_{-\epsilon \leq \mathbf{w}_i^\top \mathbf{x} \leq 0} \mathbf{w}_i \mathbf{w}_i^\top \right) \right].$$

We first note that  $\mathbb{E} \left[ \mathbf{1}_{\mathbf{w}_i^\top \mathbf{x} \geq 0} \mathbf{1}_{\mathbf{w}_i^\top \mathbf{y} \geq 0} \mathbf{b}_i \mathbf{b}_i^\top \right] = \alpha \mathbf{Q}_{\mathbf{x}, \mathbf{y}}$  where  $\alpha$  satisfies  $0.97 < \alpha < 1$ . Also, we have  $\mathbb{E} \left[ \mathbf{1}_{-\epsilon \leq \mathbf{w}_i^\top \mathbf{x} \leq 0} \mathbf{w}_i \mathbf{w}_i^\top \right] \preceq \frac{\epsilon \alpha}{2} \mathbf{I}_n$ . Thus,

$$\begin{aligned} \mathbb{E} [H_{-\epsilon}(\mathbf{x}, \mathbf{y})] &\succeq \alpha \ell \cdot \mathbf{m}^\top \mathbf{Q}_{\mathbf{x}, \mathbf{y}} \mathbf{y} - \epsilon \alpha \ell \mathbf{I}_n \\ &\succeq \alpha \ell \cdot \mathbf{m}^\top \mathbf{Q}_{\mathbf{x}, \mathbf{y}} - \epsilon \ell \mathbf{I}_n \end{aligned} \quad (53)$$

Second, the same argument as in Lemma 6 provides that for fixed  $\mathbf{x}, \mathbf{y} \in \mathcal{S}^{n-1}$ , if  $\ell > (2c_K/\epsilon)^2 n$ , then we have with probability at least  $1 - 2e^{-\gamma_K \frac{\epsilon^2 \ell}{4}}$ ,

$$\|H_{-\epsilon}(\mathbf{x}, \mathbf{y}) - \mathbb{E} [H_{-\epsilon}(\mathbf{x}, \mathbf{y})]\| \leq \epsilon \ell \quad (54)$$

Third, same argument as in Lemma 6 provides on the event  $E_1$ , we have

$$\|H_{-\epsilon}(\mathbf{x}, \mathbf{y}) - H_{-\epsilon}(\tilde{\mathbf{x}}, \tilde{\mathbf{y}})\| \leq \frac{27\ell\sqrt{n}}{\epsilon} [\|\mathbf{x} - \tilde{\mathbf{y}}\| + \|\mathbf{y} - \tilde{\mathbf{y}}\|]$$

for all  $\mathbf{x}, \mathbf{y}, \tilde{\mathbf{x}}, \tilde{\mathbf{y}} \in \mathcal{S}^{n-1}$ .

Finally, we complete the proof by an identical covering argument as in Lemma 6. We have if  $\ell \geq c_0 n \log n$  then with probability at least  $1 - 6e^{-\gamma \ell}$ ,

$$\forall \mathbf{x}, \mathbf{y} \neq \mathbf{0}, \quad H_\epsilon(\mathbf{x}, \mathbf{y}) \succeq \alpha \ell \mathbf{Q}_{\mathbf{x}, \mathbf{y}} - 3\ell \epsilon \mathbf{I}_n.$$

□

## 5.5 Proof of joint-WDC condition

We now state a result that states random gaussian matrices with truncated rows satisfy joint-WDC.

**Lemma 8.** Fix  $0 < \epsilon < 1$ . Let  $\mathbf{B} \in \mathbb{R}^{\ell \times n}$  where  $i$ th row of  $\mathbf{B}$  satisfy  $\mathbf{b}_i^\top = \mathbf{b}^\top \cdot \mathbf{1}_{\|\mathbf{b}\|_2 \leq 3\sqrt{n}/\ell}$  and  $\mathbf{b} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_n/\ell)$ . Similarly, let  $\mathbf{C} \in \mathbb{R}^{\ell \times p}$  where  $i$ th row of  $\mathbf{C}$  satisfy  $\mathbf{c}_i^\top = \mathbf{c}^\top \cdot \mathbf{1}_{\|\mathbf{c}\|_2 \leq 3\sqrt{p}/\ell}$  and  $\mathbf{c} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_p/\ell)$ . If  $\ell > c((n \log n)^2 + (p \log p)^2)$ , then with probability at least  $1 - 8e^{-\gamma \ell/(n \log n)}$ ,  $\mathbf{B}$  and  $\mathbf{C}$  satisfy joint-WDC with constants  $\epsilon$  and  $\alpha = \alpha_1 \cdot \alpha_2$ . Here,  $c$  and  $\gamma^{-1}$  are constants that depend only polynomially on  $\epsilon^{-1}$  and

$$\alpha_1 = \frac{\Gamma\left(\frac{n+2}{2}\right) - \Gamma\left(\frac{n+2}{2}, \frac{9n}{2}\right)}{\Gamma\left(\frac{n+2}{2}\right)} \text{ and } \alpha_2 = \frac{\Gamma\left(\frac{p+2}{2}\right) - \Gamma\left(\frac{p+2}{2}, \frac{9p}{2}\right)}{\Gamma\left(\frac{p+2}{2}\right)}, \quad (55)$$

where  $\Gamma$  is the Gamma function.

The proof of Lemma 8 follows directly from Lemmas 9 and 10. Using Corollary 1, we provide a concentration result of  $\mathbf{B}_{+, h}^\top \text{diag}(\mathbf{C}_{+, m} \mathbf{m}) \text{diag}(\mathbf{C}_{+, y} \mathbf{y}) \mathbf{B}_{+, x}$ , which is part of the joint-WDC condition. We note that

$$\mathbf{B}_{+, h}^\top \text{diag}(\mathbf{C}_{+, m} \mathbf{m}) \text{diag}(\mathbf{C}_{+, y} \mathbf{y}) \mathbf{B}_{+, x} = \sum_{i=1}^{\ell} \mathbf{1}_{\mathbf{b}_i^\top \mathbf{h} > 0} \mathbf{1}_{\mathbf{b}_i^\top \mathbf{x} > 0} (\mathbf{c}_i^\top \mathbf{m})_+ (\mathbf{c}_i^\top \mathbf{y})_+ \mathbf{b}_i \mathbf{b}_i^\top$$

and it is not continuous in  $\mathbf{h}$  and  $\mathbf{x}$ . So, we consider an arbitrarily good continuous approximation of  $\mathbf{B}_{+, h}^\top \text{diag}(\mathbf{C}_{+, m} \mathbf{m}) \text{diag}(\mathbf{C}_{+, y} \mathbf{y}) \mathbf{B}_{+, x}$ . Let

$$t_{-\epsilon}(z) = \begin{cases} 0 & z \leq -\epsilon, \\ 1 + \frac{z}{\epsilon} & -\epsilon \leq z \leq 0, \\ 1 & z \geq 0, \end{cases} \quad \text{and} \quad t_\epsilon(z) = \begin{cases} 0 & z \leq 0, \\ \frac{z}{\epsilon} & 0 \leq z \leq \epsilon, \\ 1 & z \geq \epsilon. \end{cases}$$

and define

$$\begin{aligned} G_{-\epsilon}(\mathbf{h}, \mathbf{x}, \mathbf{m}, \mathbf{y}) &:= \sum_{i=1}^{\ell} t_{-\epsilon}(\mathbf{b}_i^\top \mathbf{h}) t_{-\epsilon}(\mathbf{b}_i^\top \mathbf{x}) (\mathbf{c}_i^\top \mathbf{m})_+ (\mathbf{c}_i^\top \mathbf{y})_+ \mathbf{b}_i \mathbf{b}_i^\top, \\ G_\epsilon(\mathbf{h}, \mathbf{x}, \mathbf{m}, \mathbf{y}) &:= \sum_{i=1}^{\ell} t_\epsilon(\mathbf{b}_i^\top \mathbf{h}) t_\epsilon(\mathbf{b}_i^\top \mathbf{x}) (\mathbf{c}_i^\top \mathbf{m})_+ (\mathbf{c}_i^\top \mathbf{y})_+ \mathbf{b}_i \mathbf{b}_i^\top. \end{aligned}$$

We now provide an upper bound on the singular values of  $G_{-\epsilon}(\mathbf{h}, \mathbf{x}, \mathbf{m}, \mathbf{y})$ .

**Lemma 9.** Fix  $0 < \epsilon < 1$ . Let  $\mathbf{B} \in \mathbb{R}^{\ell \times n}$  where  $i$ th row of  $\mathbf{B}$  satisfy  $\mathbf{b}_i^\top = \mathbf{b}^\top \cdot \mathbf{1}_{\|\mathbf{b}\|_2 \leq 3\sqrt{n}}$  and  $\mathbf{b} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_n)$ . Similarly, let  $\mathbf{C} \in \mathbb{R}^{\ell \times p}$  where  $i$ th row of  $\mathbf{C}$  satisfy  $\mathbf{c}_i^\top = \mathbf{c}^\top \cdot \mathbf{1}_{\|\mathbf{c}\|_2 \leq 3\sqrt{p}}$  and  $\mathbf{c} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_p)$ . If  $\ell > c((n \log n)^2 + (p \log p)^2)$ , then with probability at least  $1 - 4e^{-\gamma\ell/(n \log n)}$ ,

$$\begin{aligned} \forall (\mathbf{h}, \mathbf{x}) \neq (\mathbf{0}, \mathbf{0}) \text{ and } \mathbf{m}, \mathbf{y} \in \mathcal{S}^{p-1}, \\ G_{-\epsilon}(\mathbf{h}, \mathbf{x}, \mathbf{m}, \mathbf{y}) \preceq \alpha_1 \alpha_2 \ell \mathbf{Q}_{\mathbf{h}, \mathbf{x}} \mathbf{m}^\top \mathbf{Q}_{\mathbf{m}, \mathbf{y}} \mathbf{y} + 4\ell \epsilon \mathbf{I}_n. \end{aligned}$$

Here,  $c$  and  $\gamma^{-1}$  are constants that depend only polynomially on  $\epsilon^{-1}$  and  $\alpha_1$  and  $\alpha_2$  is as in (55).

*Proof.* First we bound  $\mathbb{E}[G_{-\epsilon}(\mathbf{h}, \mathbf{x}, \mathbf{m}, \mathbf{y})]$  for fixed  $\mathbf{h}, \mathbf{x} \in \mathcal{S}^{n-1}$  and  $\mathbf{m}, \mathbf{y} \in \mathcal{S}^{p-1}$ . Noting that  $t_{-\epsilon}(z) \leq \mathbf{1}_{z \geq -\epsilon}(z) = \mathbf{1}_{z > 0}(z) + \mathbf{1}_{-\epsilon \leq z \leq 0}(z)$ , we have

$$\begin{aligned} & \mathbb{E}[G_{-\epsilon}(\mathbf{h}, \mathbf{x}, \mathbf{m}, \mathbf{y})] \\ & \leq \mathbb{E} \left[ \sum_{i=1}^{\ell} \mathbf{1}_{\mathbf{b}_i^\top \mathbf{h} \geq -\epsilon} \mathbf{1}_{\mathbf{b}_i^\top \mathbf{x} \geq -\epsilon} (\mathbf{c}_i^\top \mathbf{m})_+ (\mathbf{c}_i^\top \mathbf{y})_+ \mathbf{b}_i \mathbf{b}_i^\top \right] \\ & = \ell \mathbb{E} \left[ \mathbf{1}_{\mathbf{b}_i^\top \mathbf{h} \geq -\epsilon} \mathbf{1}_{\mathbf{b}_i^\top \mathbf{x} \geq -\epsilon} (\mathbf{c}_i^\top \mathbf{m})_+ (\mathbf{c}_i^\top \mathbf{y})_+ \mathbf{b}_i \mathbf{b}_i^\top \right] \\ & \leq \ell \mathbb{E} \left[ \left( \mathbf{1}_{\mathbf{b}_i^\top \mathbf{h} \geq 0} \mathbf{1}_{\mathbf{b}_i^\top \mathbf{x} \geq 0} (\mathbf{c}_i^\top \mathbf{m})_+ (\mathbf{c}_i^\top \mathbf{y})_+ + \left( \mathbf{1}_{-\epsilon \leq \mathbf{b}_i^\top \mathbf{h} \leq 0} + \mathbf{1}_{-\epsilon \leq \mathbf{b}_i^\top \mathbf{x} \leq 0} \right) (\mathbf{c}_i^\top \mathbf{m})_+ (\mathbf{c}_i^\top \mathbf{y})_+ \right) \mathbf{b}_i \mathbf{b}_i^\top \right] \\ & = \ell \mathbb{E} \left[ \left( \mathbf{1}_{\mathbf{b}_i^\top \mathbf{h} \geq 0} \mathbf{1}_{\mathbf{b}_i^\top \mathbf{x} \geq 0} \mathbf{b}_i \mathbf{b}_i^\top \right) (\mathbf{c}_i^\top \mathbf{m})_+ (\mathbf{c}_i^\top \mathbf{y})_+ \right] + 2\ell \mathbb{E} \left[ \left( \mathbf{1}_{-\epsilon \leq \mathbf{b}_i^\top \mathbf{h} \leq 0} \mathbf{b}_i \mathbf{b}_i^\top \right) (\mathbf{c}_i^\top \mathbf{m})_+ (\mathbf{c}_i^\top \mathbf{y})_+ \right]. \end{aligned}$$

We first note that  $\mathbb{E} \left[ \mathbf{1}_{\mathbf{b}_i^\top \mathbf{h} \geq 0} \mathbf{1}_{\mathbf{b}_i^\top \mathbf{x} \geq 0} \mathbf{b}_i \mathbf{b}_i^\top \right] = \alpha_1 \mathbf{Q}_{\mathbf{h}, \mathbf{x}}$  and  $\mathbb{E}[(\mathbf{c}_i^\top \mathbf{m})_+ (\mathbf{c}_i^\top \mathbf{y})_+] = \alpha_2 \mathbf{m}^\top \mathbf{Q}_{\mathbf{m}, \mathbf{y}} \mathbf{y}$  where  $\alpha_i$  satisfies  $0.97 < \alpha_i < 1$ . Also, we have  $|\mathbf{m}^\top \mathbf{Q}_{\mathbf{m}, \mathbf{y}} \mathbf{y}| \leq \frac{1}{2}$  and  $\mathbb{E} \left[ \mathbf{1}_{-\epsilon \leq \mathbf{b}_i^\top \mathbf{h} \leq 0} \mathbf{b}_i \mathbf{b}_i^\top \right] \preceq \frac{\epsilon \alpha_1}{2} \mathbf{I}_n$ . Thus,

$$\begin{aligned} \mathbb{E}[G_{-\epsilon}(\mathbf{h}, \mathbf{x}, \mathbf{m}, \mathbf{y})] & \preceq \alpha_1 \alpha_2 \ell \cdot \mathbf{m}^\top \mathbf{Q}_{\mathbf{m}, \mathbf{y}} \mathbf{y} \cdot \mathbf{Q}_{\mathbf{h}, \mathbf{x}} + 2\ell \cdot \mathbb{E} \left[ \mathbf{1}_{-\epsilon \leq \mathbf{b}_i^\top \mathbf{h} \leq 0} \mathbf{b}_i \mathbf{b}_i^\top \right] \cdot \alpha_2 \mathbf{m}^\top \mathbf{Q}_{\mathbf{m}, \mathbf{y}} \mathbf{y} \\ & \preceq \alpha_1 \alpha_2 \ell \cdot \mathbf{m}^\top \mathbf{Q}_{\mathbf{m}, \mathbf{y}} \mathbf{y} \cdot \mathbf{Q}_{\mathbf{h}, \mathbf{x}} + \frac{\epsilon \alpha_1 \alpha_2 \ell}{2} \mathbf{I}_n \\ & \preceq \alpha_1 \alpha_2 \ell \cdot \mathbf{m}^\top \mathbf{Q}_{\mathbf{m}, \mathbf{y}} \mathbf{y} \cdot \mathbf{Q}_{\mathbf{h}, \mathbf{x}} + \frac{\epsilon \ell}{2} \mathbf{I}_n \end{aligned} \tag{56}$$

Second, we show concentration of  $G_{-\epsilon}(\mathbf{h}, \mathbf{x}, \mathbf{m}, \mathbf{y})$  for fixed  $\mathbf{h}, \mathbf{x} \in \mathcal{S}^{n-1}$  and  $\mathbf{m}, \mathbf{y} \in \mathcal{S}^{p-1}$ . Let  $\xi_i = \sqrt{t_{-\epsilon}(\mathbf{b}_i^\top \mathbf{h}) t_{-\epsilon}(\mathbf{b}_i^\top \mathbf{x})} (\mathbf{c}_i^\top \mathbf{m})_+ (\mathbf{c}_i^\top \mathbf{y})_+ + \mathbf{b}_i$ . We have

$$\begin{aligned} & G_{-\epsilon}(\mathbf{h}, \mathbf{x}, \mathbf{m}, \mathbf{y}) - \mathbb{E}[G_{-\epsilon}(\mathbf{h}, \mathbf{x}, \mathbf{m}, \mathbf{y})] \\ & = \sum_{i=1}^{\ell} \left( t_{-\epsilon}(\mathbf{b}_i^\top \mathbf{h}) t_{-\epsilon}(\mathbf{b}_i^\top \mathbf{x}) (\mathbf{c}_i^\top \mathbf{m})_+ (\mathbf{c}_i^\top \mathbf{y})_+ \mathbf{b}_i \mathbf{b}_i^\top - \mathbb{E}[t_{-\epsilon}(\mathbf{b}_i^\top \mathbf{h}) t_{-\epsilon}(\mathbf{b}_i^\top \mathbf{x}) (\mathbf{c}_i^\top \mathbf{m})_+ (\mathbf{c}_i^\top \mathbf{y})_+ \mathbf{b}_i \mathbf{b}_i^\top] \right) \\ & = \sum_{i=1}^{\ell} (\xi_i \xi_i^\top - \mathbb{E}[\xi_i \xi_i^\top]). \end{aligned}$$

Note that  $\xi_i$  is sub-Gaussian for all  $i$  and that the sub-Gaussian norm of  $\xi_i$  is bounded from above by  $K = \tilde{K} \sqrt{n}$ , where  $\tilde{K}$  is an absolute constant. By Corollary 1, there exists a  $c = \bar{c} \sqrt{n \log n}$  and  $\gamma = \frac{\bar{\gamma}}{n \log n}$  such that for all  $t \geq 0$ , with probability at least  $1 - 2e^{-\gamma t^2}$ ,

$$\|G_{-\epsilon}(\mathbf{h}, \mathbf{x}, \mathbf{m}, \mathbf{y}) - \mathbb{E}[G_{-\epsilon}(\mathbf{h}, \mathbf{x}, \mathbf{m}, \mathbf{y})]\| \leq \max(\delta, \delta^2) \ell, \quad \text{where } \delta = c \sqrt{\frac{n}{\ell}} + \frac{t}{\sqrt{\ell}}.$$

Here,  $\bar{c}$  and  $\bar{\gamma}$  are absolute constants. If  $\ell > (2\bar{c}/\epsilon)^2 n^2 (\log n)^2$ ,  $t = \epsilon \sqrt{\ell}/2$ , and  $\epsilon < 1$ , we have

$$\|G_{-\epsilon}(\mathbf{h}, \mathbf{x}, \mathbf{m}, \mathbf{y}) - \mathbb{E}[G_{-\epsilon}(\mathbf{h}, \mathbf{x}, \mathbf{m}, \mathbf{y})]\| \leq \epsilon \ell \tag{57}$$

with probability at least  $1 - 2e^{-\bar{\gamma} \frac{\epsilon^2}{4} \frac{\ell}{n \log n}}$ .

Third, we bound the Lipschitz constant of  $G_{-\epsilon}$ . For  $\tilde{\mathbf{h}}, \tilde{\mathbf{x}} \in \mathbb{R}^n$  and  $\tilde{\mathbf{m}}, \tilde{\mathbf{y}} \in \mathcal{S}^{p-1}$  we have

$$\begin{aligned}
& G_{-\epsilon}(\mathbf{h}, \mathbf{x}, \mathbf{m}, \mathbf{y}) - G_{-\epsilon}(\tilde{\mathbf{h}}, \tilde{\mathbf{x}}, \tilde{\mathbf{m}}, \tilde{\mathbf{y}}) \\
&= \sum_{i=1}^{\ell} \left[ t_{-\epsilon}(\mathbf{b}_i^\top \mathbf{h}) t_{-\epsilon}(\mathbf{b}_i^\top \mathbf{x}) (\mathbf{c}_i^\top \mathbf{m})_+ + (\mathbf{c}_i^\top \mathbf{y})_+ - t_{-\epsilon}(\mathbf{b}_i^\top \tilde{\mathbf{h}}) t_{-\epsilon}(\mathbf{b}_i^\top \tilde{\mathbf{x}}) (\mathbf{c}_i^\top \tilde{\mathbf{m}})_+ + (\mathbf{c}_i^\top \tilde{\mathbf{y}})_+ \right] \mathbf{b}_i \mathbf{b}_i^\top \\
&= \sum_{i=1}^{\ell} \left[ t_{-\epsilon}(\mathbf{b}_i^\top \mathbf{h}) t_{-\epsilon}(\mathbf{b}_i^\top \mathbf{x}) (\mathbf{c}_i^\top \mathbf{m})_+ ((\mathbf{c}_i^\top \mathbf{y})_+ - (\mathbf{c}_i^\top \tilde{\mathbf{y}})_+) \right. \\
&\quad + t_{-\epsilon}(\mathbf{b}_i^\top \mathbf{h}) t_{-\epsilon}(\mathbf{b}_i^\top \mathbf{x}) (\mathbf{c}_i^\top \tilde{\mathbf{y}})_+ ((\mathbf{c}_i^\top \mathbf{m})_+ - (\mathbf{c}_i^\top \tilde{\mathbf{m}})_+) \\
&\quad + t_{-\epsilon}(\mathbf{b}_i^\top \mathbf{h}) (\mathbf{c}_i^\top \tilde{\mathbf{m}})_+ (\mathbf{c}_i^\top \tilde{\mathbf{y}})_+ (t_{-\epsilon}(\mathbf{b}_i^\top \mathbf{x}) - t_{-\epsilon}(\mathbf{b}_i^\top \tilde{\mathbf{x}})) \\
&\quad \left. + t_{-\epsilon}(\mathbf{b}_i^\top \tilde{\mathbf{x}}) (\mathbf{c}_i^\top \tilde{\mathbf{m}})_+ (\mathbf{c}_i^\top \tilde{\mathbf{y}})_+ \left( t_{-\epsilon}(\mathbf{b}_i^\top \mathbf{h}) - t_{-\epsilon}(\mathbf{b}_i^\top \tilde{\mathbf{h}}) \right) \right] \mathbf{b}_i \mathbf{b}_i^\top \\
&= \mathbf{B}^\top \left[ \text{diag}(t_{-\epsilon}(\mathbf{B}\mathbf{h}) \odot t_{-\epsilon}(\mathbf{B}\mathbf{x}) \odot (\mathbf{C}\mathbf{m})_+) \text{diag}((\mathbf{C}\mathbf{y})_+ - (\mathbf{C}\tilde{\mathbf{y}})_+) \right. \\
&\quad + \text{diag}(t_{-\epsilon}(\mathbf{B}\mathbf{h}) \odot t_{-\epsilon}(\mathbf{B}\mathbf{x}) \odot (\mathbf{C}\tilde{\mathbf{y}})_+) \text{diag}((\mathbf{C}\mathbf{m})_+ - (\mathbf{C}\tilde{\mathbf{m}})_+) \\
&\quad + \text{diag}(t_{-\epsilon}(\mathbf{B}\mathbf{h}) \odot (\mathbf{C}\tilde{\mathbf{m}})_+ \odot (\mathbf{C}\tilde{\mathbf{y}})_+) \text{diag}(t_{-\epsilon}(\mathbf{B}\mathbf{x}) - t_{-\epsilon}(\mathbf{B}\tilde{\mathbf{x}})) \\
&\quad \left. + \text{diag}(t_{-\epsilon}(\mathbf{B}\tilde{\mathbf{x}}) \odot (\mathbf{C}\tilde{\mathbf{m}})_+ \odot (\mathbf{C}\tilde{\mathbf{y}})_+) \text{diag}(t_{-\epsilon}(\mathbf{B}\mathbf{h}) - t_{-\epsilon}(\mathbf{B}\tilde{\mathbf{h}})) \right] \mathbf{B}
\end{aligned}$$

Thus,

$$\begin{aligned}
& \|G_{-\epsilon}(\mathbf{h}, \mathbf{x}, \mathbf{m}, \mathbf{y}) - G_{-\epsilon}(\tilde{\mathbf{h}}, \tilde{\mathbf{x}}, \tilde{\mathbf{m}}, \tilde{\mathbf{y}})\| \\
&\leq \|\mathbf{B}\|^2 \left[ \|\mathbf{C}\mathbf{m}\|_\infty \|(\mathbf{C}\mathbf{y})_+ - (\mathbf{C}\tilde{\mathbf{y}})_+\|_\infty + \|\mathbf{C}\tilde{\mathbf{y}}\|_\infty \|(\mathbf{C}\mathbf{m})_+ - (\mathbf{C}\tilde{\mathbf{m}})_+\|_\infty \right. \\
&\quad \left. + \|\mathbf{C}\tilde{\mathbf{m}}\|_\infty \|\mathbf{C}\tilde{\mathbf{y}}\|_\infty \|t_{-\epsilon}(\mathbf{B}\mathbf{x}) - t_{-\epsilon}(\mathbf{B}\tilde{\mathbf{x}})\|_\infty + \|\mathbf{C}\tilde{\mathbf{m}}\|_\infty \|\mathbf{C}\tilde{\mathbf{y}}\|_\infty \|t_{-\epsilon}(\mathbf{B}\mathbf{h}) - t_{-\epsilon}(\mathbf{B}\tilde{\mathbf{h}})\|_\infty \right] \\
&\leq \|\mathbf{B}\|^2 \left[ \max_{i \in [\ell]} \|\mathbf{c}_i\|_2 \max_{i \in [\ell]} |(\mathbf{c}_i^\top \mathbf{y})_+ - (\mathbf{c}_i^\top \tilde{\mathbf{y}})_+| + \max_{i \in [\ell]} \|\mathbf{c}_i\|_2 \max_{i \in [\ell]} |(\mathbf{c}_i^\top \mathbf{m})_+ - (\mathbf{c}_i^\top \tilde{\mathbf{m}})_+| \right. \\
&\quad \left. + \left( \max_{i \in [\ell]} \|\mathbf{c}_i\|_2 \right)^2 \max_{i \in [\ell]} |t_{-\epsilon}(\mathbf{b}_i^\top \mathbf{x}) - t_{-\epsilon}(\mathbf{b}_i^\top \tilde{\mathbf{x}})| + \left( \max_{i \in [\ell]} \|\mathbf{c}_i\|_2 \right)^2 \max_{i \in [\ell]} |t_{-\epsilon}(\mathbf{b}_i^\top \mathbf{h}) - t_{-\epsilon}(\mathbf{b}_i^\top \tilde{\mathbf{h}})| \right] \\
&\leq \|\mathbf{B}\|^2 \left[ \max_{i \in [\ell]} \|\mathbf{c}_i\|_2 \max_{i \in [\ell]} |\mathbf{c}_i^\top (\mathbf{y} - \tilde{\mathbf{y}})| + \max_{i \in [\ell]} \|\mathbf{c}_i\|_2 \max_{i \in [\ell]} |\mathbf{c}_i^\top (\mathbf{m} - \tilde{\mathbf{m}})| \right. \\
&\quad \left. + \left( \max_{i \in [\ell]} \|\mathbf{c}_i\|_2 \right)^2 \max_{i \in [\ell]} \frac{1}{\epsilon} |\mathbf{b}_i^\top (\mathbf{x} - \tilde{\mathbf{x}})| + \left( \max_{i \in [\ell]} \|\mathbf{c}_i\|_2 \right)^2 \max_{i \in [\ell]} \frac{1}{\epsilon} |\mathbf{b}_i^\top (\mathbf{h} - \tilde{\mathbf{h}})| \right] \\
&\leq \|\mathbf{B}\|^2 \left[ \left( \max_{i \in [\ell]} \|\mathbf{c}_i\|_2 \right)^2 \|\mathbf{y} - \tilde{\mathbf{y}}\| + \left( \max_{i \in [\ell]} \|\mathbf{c}_i\|_2 \right)^2 \|\mathbf{m} - \tilde{\mathbf{m}}\| \right. \\
&\quad \left. + \frac{1}{\epsilon} \left( \max_{i \in [\ell]} \|\mathbf{c}_i\|_2 \right)^2 \max_{i \in [\ell]} \|\mathbf{b}_i\|_2 \|\mathbf{x} - \tilde{\mathbf{x}}\| + \frac{1}{\epsilon} \left( \max_{i \in [\ell]} \|\mathbf{c}_i\|_2 \right)^2 \max_{i \in [\ell]} \|\mathbf{b}_i\|_2 \|\mathbf{h} - \tilde{\mathbf{h}}\| \right] \\
&\leq \|\mathbf{B}\|^2 \left[ 9p \|\mathbf{y} - \tilde{\mathbf{y}}\| + 9p \|\mathbf{m} - \tilde{\mathbf{m}}\| + \frac{27}{\epsilon} \sqrt{np} \|\mathbf{x} - \tilde{\mathbf{x}}\| + \frac{27}{\epsilon} \sqrt{np} \|\mathbf{h} - \tilde{\mathbf{h}}\| \right]
\end{aligned}$$

where the first inequality follows because  $|t_{-\epsilon}(z)| \leq 1$  for all  $z$ , and the third inequality follows because  $t_{-\epsilon}(z)$  is  $1/\epsilon$ -Lipschitz and  $(z)_+$  is 1-Lipschitz. Let  $E_1$  be the event that  $\|\mathbf{B}\| \leq 3\sqrt{\ell}$ . By Corollary 5.35 in Vershynin [2012], for  $\mathbf{A} \in \mathbb{R}^{\ell \times n}$  with rows of  $\mathbf{A}$  following  $\mathcal{N}(\mathbf{0}, \mathbf{I}_n)$ , we have  $\mathbb{P}(\|\mathbf{A}\| \leq 3\sqrt{\ell}) \geq 1 - 2e^{-\ell/2}$ , if  $\ell \geq n$ . As rows of  $\mathbf{B}$  are truncated, we have  $\mathbb{P}(E_1) \geq 1 - 2e^{-\ell/2}$ , if  $\ell \geq n$  as well. On  $E_1$ , we have

$$\begin{aligned}
& \|G_{-\epsilon}(\mathbf{h}, \mathbf{x}, \mathbf{m}, \mathbf{y}) - G_{-\epsilon}(\tilde{\mathbf{h}}, \tilde{\mathbf{x}}, \tilde{\mathbf{m}}, \tilde{\mathbf{y}})\| \\
&\leq \frac{729\ell\sqrt{np}}{\epsilon} \left[ \|\mathbf{y} - \tilde{\mathbf{y}}\| + \|\mathbf{m} - \tilde{\mathbf{m}}\| + \|\mathbf{x} - \tilde{\mathbf{x}}\| + \|\mathbf{h} - \tilde{\mathbf{h}}\| \right]
\end{aligned} \tag{58}$$

for all  $\tilde{\mathbf{h}}, \tilde{\mathbf{x}} \in \mathcal{S}^{n-1}$  and  $\tilde{\mathbf{m}}, \tilde{\mathbf{y}} \in \mathcal{S}^{p-1}$ .

Finally, we complete the proof by a covering argument. Let  $\mathcal{N}_\delta$  be a  $\delta$ -net on  $\mathcal{S}^{n-1} \times \mathcal{S}^{p-1}$  such that  $|\mathcal{N}_\delta| \leq (3/\delta)^{n+p}$ . Take  $\delta = \frac{\epsilon^2}{2916\sqrt{np}}$ . Combining (56) and (57), we have

$$\begin{aligned} \forall (\mathbf{h}, \mathbf{m}), (\mathbf{x}, \mathbf{y}) \in \mathcal{N}_\delta, \quad G_{-\epsilon}(\mathbf{h}, \mathbf{x}, \mathbf{m}, \mathbf{y}) &\preceq \mathbb{E}G_{-\epsilon}(\mathbf{h}, \mathbf{x}, \mathbf{m}, \mathbf{y}) + \ell\epsilon I_n \\ &\preceq \alpha^2 \ell Q_{\mathbf{h}, \mathbf{x}} \mathbf{m}^\top Q_{\mathbf{m}, \mathbf{y}} \mathbf{y} + 3\ell\epsilon I_n. \end{aligned}$$

with probability at least

$$1 - 2|\mathcal{N}_\delta| e^{-\gamma_K \epsilon^2 \frac{\ell}{4n \log n}} \geq 1 - 2 \left( \frac{3}{\delta} \right)^{n+p} e^{-\gamma_K \epsilon^2 \frac{\ell}{4n \log n}} \geq 1 - 2e^{-\gamma_K \epsilon^2 \frac{\ell}{4n \log n} + (n+p) \log(3 \cdot 2916\sqrt{np}/\epsilon^2)}.$$

If  $\ell \geq \tilde{c}n(n+p)\log(n)\log(np)$  for some  $\tilde{c} = \Omega(\epsilon^2 \log \epsilon^{-1})$ , then this probability is at least  $1 - 2e^{-\tilde{\gamma}\ell/(n \log(n))}$  for some  $\tilde{\gamma} = O(\epsilon^2)$ . For  $(\mathbf{h}, \mathbf{m}), (\mathbf{x}, \mathbf{y}) \in \mathcal{S}^{n-1} \times \mathcal{S}^{p-1}$ , let  $(\tilde{\mathbf{h}}, \tilde{\mathbf{m}}), (\tilde{\mathbf{x}}, \tilde{\mathbf{y}}) \in \mathcal{N}_\delta$  be such that  $\|\mathbf{h} - \tilde{\mathbf{h}}\|_2 \leq \delta$ ,  $\|\mathbf{x} - \tilde{\mathbf{x}}\|_2 \leq \delta$ ,  $\|\mathbf{m} - \tilde{\mathbf{m}}\|_2 \leq \delta$  and  $\|\mathbf{y} - \tilde{\mathbf{y}}\|_2 \leq \delta$ . By (58), we have that

$$\begin{aligned} \forall (\mathbf{h}, \mathbf{x}) \neq (\mathbf{0}, \mathbf{0}) \text{ and } \mathbf{m}, \mathbf{y} \in \mathcal{S}^{p-1}, \quad G_{-\epsilon}(\mathbf{h}, \mathbf{x}, \mathbf{m}, \mathbf{y}) &\\ &\preceq G_{-\epsilon}(\tilde{\mathbf{h}}, \tilde{\mathbf{x}}, \tilde{\mathbf{m}}, \tilde{\mathbf{y}}) + \frac{729\ell\sqrt{np}}{\epsilon} 4\delta I_n \\ &\preceq \alpha_1 \alpha_2 \ell Q_{\mathbf{h}, \mathbf{x}} \mathbf{m}^\top Q_{\mathbf{m}, \mathbf{y}} \mathbf{y} + 4\ell\epsilon I_n. \end{aligned}$$

In conclusion, the result of this lemma holds if  $\ell > (2\bar{c}/\epsilon)^2 n^2 (\log n)^2$  and  $\ell \geq \tilde{c}n(n+p)\log(n)\log(np)$ , with probability at least  $1 - 2e^{-\ell/2} - 2e^{-\tilde{\gamma}\ell/(n \log(n))} > 1 - 4e^{-\gamma\ell/(n \log(n))}$  for some  $\gamma = O(\epsilon^2)$  and  $\tilde{c} = \Omega(\epsilon^2 \log \epsilon)$ .  $\square$

Next, we now provide an upper bound on the singular values of  $G_\epsilon(\mathbf{h}, \mathbf{x}, \mathbf{m}, \mathbf{y})$ .

**Lemma 10.** Fix  $0 < \epsilon < 1$ . Let  $\mathbf{B} \in \mathbb{R}^{\ell \times n}$  where  $i$ th row of  $\mathbf{B}$  satisfy  $\mathbf{b}_i^\top = \mathbf{b}^\top \cdot \mathbf{1}_{\|\mathbf{b}\|_2 \leq 3\sqrt{n}}$  and  $\mathbf{b} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_n)$ . Similarly, let  $\mathbf{C} \in \mathbb{R}^{\ell \times p}$  where  $i$ th row of  $\mathbf{C}$  satisfy  $\mathbf{c}_i^\top = \mathbf{c}^\top \cdot \mathbf{1}_{\|\mathbf{c}\|_2 \leq 3\sqrt{p}}$  and  $\mathbf{c} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_p)$ . If  $\ell > c((n \log n)^2 + (p \log p)^2)$ , then with probability at least  $1 - 4e^{-\gamma\ell/(n \log n)}$ ,

$$\begin{aligned} \forall (\mathbf{h}, \mathbf{x}) \neq (\mathbf{0}, \mathbf{0}) \text{ and } \mathbf{m}, \mathbf{y} \in \mathcal{S}^{p-1}, \quad G_\epsilon(\mathbf{h}, \mathbf{x}, \mathbf{m}, \mathbf{y}) &\succeq \alpha_1 \alpha_2 \ell Q_{\mathbf{h}, \mathbf{x}} \mathbf{m}^\top Q_{\mathbf{m}, \mathbf{y}} \mathbf{y} - 4\ell\epsilon I_n. \end{aligned}$$

Here,  $c$  and  $\gamma^{-1}$  are constants that depend only polynomially on  $\epsilon^{-1}$  and  $\alpha_1$  and  $\alpha_2$  as in (55).

*Proof.* First we bound  $\mathbb{E}[G_\epsilon(\mathbf{h}, \mathbf{x}, \mathbf{m}, \mathbf{y})]$  for fixed  $\mathbf{h}, \mathbf{x} \in \mathcal{S}^{n-1}$  and  $\mathbf{m}, \mathbf{y} \in \mathcal{S}^{p-1}$ . Noting that  $t_\epsilon(z) \geq \mathbf{1}_{z>0}(z) - \mathbf{1}_{0 \leq z \leq \epsilon}(z)$ , we have

$$\begin{aligned} &\mathbb{E}[G_\epsilon(\mathbf{h}, \mathbf{x}, \mathbf{m}, \mathbf{y})] \\ &\succeq \ell \mathbb{E} \left[ \left( \mathbf{1}_{\mathbf{b}_i^\top \mathbf{h} \geq 0} \mathbf{1}_{\mathbf{b}_i^\top \mathbf{x} \geq 0} (\mathbf{c}_i^\top \mathbf{m})_+ (\mathbf{c}_i^\top \mathbf{y})_+ - \left( \mathbf{1}_{0 \leq \mathbf{b}_i^\top \mathbf{h} \leq \epsilon} + \mathbf{1}_{0 \leq \mathbf{b}_i^\top \mathbf{x} \leq \epsilon} \right) (\mathbf{c}_i^\top \mathbf{m})_+ (\mathbf{c}_i^\top \mathbf{y})_+ \right) \mathbf{b}_i \mathbf{b}_i^\top \right] \\ &= \ell \mathbb{E} \left[ \left( \mathbf{1}_{\mathbf{b}_i^\top \mathbf{h} \geq 0} \mathbf{1}_{\mathbf{b}_i^\top \mathbf{x} \geq 0} \mathbf{b}_i \mathbf{b}_i^\top \right) (\mathbf{c}_i^\top \mathbf{m})_+ (\mathbf{c}_i^\top \mathbf{y})_+ \right] - 2\ell \mathbb{E} \left[ \left( \mathbf{1}_{0 \leq \mathbf{b}_i^\top \mathbf{h} \leq \epsilon} \mathbf{b}_i \mathbf{b}_i^\top \right) (\mathbf{c}_i^\top \mathbf{m})_+ (\mathbf{c}_i^\top \mathbf{y})_+ \right]. \end{aligned}$$

We first note that  $\mathbb{E}[\mathbf{1}_{\mathbf{b}_i^\top \mathbf{h} \geq 0} \mathbf{1}_{\mathbf{b}_i^\top \mathbf{x} \geq 0} \mathbf{b}_i \mathbf{b}_i^\top] = \alpha_1 Q_{\mathbf{h}, \mathbf{x}}$  and  $\mathbb{E}[(\mathbf{c}_i^\top \mathbf{m})_+ (\mathbf{c}_i^\top \mathbf{y})_+] = \alpha_2 \mathbf{m}^\top Q_{\mathbf{m}, \mathbf{y}} \mathbf{y}$  where  $\alpha_i$  satisfies  $0.97 < \alpha_i < 1$ . Also, we have  $|\mathbf{m}^\top Q_{\mathbf{m}, \mathbf{y}} \mathbf{y}| \leq \frac{1}{2}$  and  $\mathbb{E}[\mathbf{1}_{0 \leq \mathbf{b}_i^\top \mathbf{h} \leq \epsilon} \mathbf{b}_i \mathbf{b}_i^\top] \preceq \frac{\alpha_1}{2} I_n$ . Thus,

$$\begin{aligned} \mathbb{E}[G_\epsilon(\mathbf{h}, \mathbf{x}, \mathbf{m}, \mathbf{y})] &\succeq \alpha_1 \alpha_2 \ell \cdot \mathbf{m}^\top Q_{\mathbf{m}, \mathbf{y}} \mathbf{y} \cdot Q_{\mathbf{h}, \mathbf{x}} - 2\ell \cdot \mathbb{E}[\mathbf{1}_{0 \leq \mathbf{b}_i^\top \mathbf{h} \leq \epsilon} \mathbf{b}_i \mathbf{b}_i^\top] \cdot \alpha_2 \mathbf{m}^\top Q_{\mathbf{m}, \mathbf{y}} \mathbf{y} \\ &\succeq \alpha_1 \alpha_2 \ell \cdot \mathbf{m}^\top Q_{\mathbf{m}, \mathbf{y}} \mathbf{y} \cdot Q_{\mathbf{h}, \mathbf{x}} - \frac{\epsilon \alpha_1 \alpha_2 \ell}{2} I_n \\ &\succeq \alpha_1 \alpha_2 \ell \cdot \mathbf{m}^\top Q_{\mathbf{m}, \mathbf{y}} \mathbf{y} \cdot Q_{\mathbf{h}, \mathbf{x}} - \ell I_n \end{aligned} \tag{59}$$

Second, we show concentration of  $G_\epsilon(\mathbf{h}, \mathbf{x}, \mathbf{m}, \mathbf{y})$  for fixed  $\mathbf{h}, \mathbf{x} \in \mathcal{S}^{n-1}$  and  $\mathbf{m}, \mathbf{y} \in \mathcal{S}^{p-1}$  and is similar to the steps shown in proof of Lemma 9. Let  $\xi_i = \sqrt{t_\epsilon(\mathbf{b}_i^\top \mathbf{h}) t_\epsilon(\mathbf{b}_i^\top \mathbf{x}) (\mathbf{c}_i^\top \mathbf{m})_+ (\mathbf{c}_i^\top \mathbf{y})_+} \mathbf{b}_i$ . If  $\ell > (2\bar{c}/\epsilon)^2 n^2 (\log n)^2$ , we have

$$\|G_\epsilon(\mathbf{h}, \mathbf{x}, \mathbf{m}, \mathbf{y}) - \mathbb{E}[G_\epsilon(\mathbf{h}, \mathbf{x}, \mathbf{m}, \mathbf{y})]\| \leq \epsilon n \tag{60}$$

with probability at least  $1 - 2e^{-\bar{\gamma} \frac{\epsilon^2}{4} \frac{\ell}{n \log n}}$ . Here,  $\bar{c}$  and  $\bar{\gamma}$  are absolute constants.

Third, we bound the Lipschitz constant of  $G_\epsilon$ , and is again similar to the steps shown in proof of Lemma 9. If  $\ell \geq n$  then we have

$$\begin{aligned} & \|G_\epsilon(\mathbf{h}, \mathbf{x}, \mathbf{m}, \mathbf{y}) - G_\epsilon(\tilde{\mathbf{h}}, \tilde{\mathbf{x}}, \tilde{\mathbf{m}}, \tilde{\mathbf{y}})\| \\ & \leq \frac{729\ell\sqrt{np}}{\epsilon} [\|\mathbf{y} - \tilde{\mathbf{y}}\| + \|\mathbf{m} - \tilde{\mathbf{m}}\| + \|\mathbf{x} - \tilde{\mathbf{x}}\| + \|\mathbf{h} - \tilde{\mathbf{h}}\|] \end{aligned} \quad (61)$$

for all  $\tilde{\mathbf{h}}, \tilde{\mathbf{x}} \in \mathcal{S}^{n-1}$  and  $\tilde{\mathbf{m}}, \tilde{\mathbf{y}} \in \mathcal{S}^{p-1}$  with probability at least  $1 - 2e^{-\ell/2}$ .

Finally, we complete the proof by a covering argument. Let  $\mathcal{N}_\delta$  be a  $\delta$ -net on  $\mathcal{S}^{n-1} \times \mathcal{S}^{p-1}$  such that  $|\mathcal{N}_\delta| \leq (3/\delta)^{n+p}$ . Take  $\delta = \frac{\epsilon^2}{2916\sqrt{np}}$ . Combining (59) and (60), we have

$$\begin{aligned} \forall (\mathbf{h}, \mathbf{m}), (\mathbf{x}, \mathbf{y}) \in \mathcal{N}_\delta, \quad & G_\epsilon(\mathbf{h}, \mathbf{x}, \mathbf{m}, \mathbf{y}) \succeq \mathbb{E}G_\epsilon(\mathbf{h}, \mathbf{x}, \mathbf{m}, \mathbf{y}) - \ell\epsilon I_n \\ & \succeq \alpha^2 \ell Q_{\mathbf{h}, \mathbf{x}} \mathbf{m}^\top Q_{\mathbf{m}, \mathbf{y}} \mathbf{y} - 3\ell\epsilon I_n \end{aligned}$$

with probability at least

$$1 - 2|\mathcal{N}_\delta| e^{-\gamma_K \epsilon^2 \frac{\ell}{4n \log n}} \geq 1 - 2 \left( \frac{3}{\delta} \right)^{n+p} e^{-\gamma_K \epsilon^2 \frac{\ell}{4n \log n}} \geq 1 - 2e^{-\gamma_K \epsilon^2 \frac{\ell}{4n \log n} + (n+p) \log(108\sqrt{np}/\epsilon^2)}.$$

If  $\ell \geq \tilde{c}n(n+p) \log(n) \log(np)$  for some  $\tilde{c} = \Omega(\epsilon^2 \log \epsilon)$ , then this probability is at least  $1 - 2e^{-\tilde{\gamma}\ell/(n \log n)}$  for some  $\tilde{\gamma} = O(\epsilon^2)$ . For  $(\mathbf{h}, \mathbf{m}), (\mathbf{x}, \mathbf{y}) \in \mathcal{S}^{n-1} \times \mathcal{S}^{p-1}$ , let  $(\tilde{\mathbf{h}}, \tilde{\mathbf{m}}), (\tilde{\mathbf{x}}, \tilde{\mathbf{y}}) \in \mathcal{N}_\delta$  be such that  $\|\mathbf{h} - \tilde{\mathbf{h}}\|_2 \leq \delta$ ,  $\|\mathbf{x} - \tilde{\mathbf{x}}\|_2 \leq \delta$ ,  $\|\mathbf{m} - \tilde{\mathbf{m}}\|_2 \leq \delta$  and  $\|\mathbf{y} - \tilde{\mathbf{y}}\|_2 \leq \delta$ . By (61), we have that

$$\begin{aligned} \forall (\mathbf{h}, \mathbf{x}) \neq (\mathbf{0}, \mathbf{0}) \text{ and } \mathbf{m}, \mathbf{y} \in \mathcal{S}^{p-1}, \quad & G_\epsilon(\mathbf{h}, \mathbf{x}, \mathbf{m}, \mathbf{y}) \\ & \succeq \alpha_1 \alpha_2 \ell Q_{\mathbf{h}, \mathbf{x}} \mathbf{m}^\top Q_{\mathbf{m}, \mathbf{y}} \mathbf{y} - 4\ell\epsilon I_n. \end{aligned}$$

In conclusion, the result of this lemma holds if  $\ell > (2\bar{c}/\epsilon)^2 n^2 (\log n)^2$  and  $\ell \geq \tilde{c}n(n+p) \log(n) \log(np)$ , with probability at least  $1 - 2e^{-\ell/2} - 2e^{-\tilde{\gamma}\ell/(n \log n)} > 1 - 4e^{-\gamma\ell/(n \log n)}$  for some  $\gamma = O(\epsilon^2)$  and  $\tilde{c} = \Omega(\epsilon^2 \log \epsilon)$ .  $\square$

## 5.6 Concentration of matrices with sub-gaussian rows

The proof of Lemmas 6 and 7 require results from concentration of sub-exponential random variables that has a better dependence on the sub-exponential parameters. To this end, we use the following Bernstein inequality and refer the readers to Jeong et al. [2019] for a proof of the theorem.

**Theorem 5.** *Let  $\mathbf{a} = (a_1, \dots, a_n)$  be a fixed non-zero vector and let  $y_1, \dots, y_m$  be independent, mean zero sub-exponential random variables satisfying  $\mathbb{E}|y_i| \leq 2$  and  $\|y_i\|_{\psi_1} \leq K_i^2$  ( $K_i \geq 2$ ). Then for every  $u \geq 0$ , we have*

$$\mathbb{P} \left( \left| \sum_{i=1}^m a_i y_i \right| \geq u \right) \leq 2 \exp \left[ -c \min \left( \frac{u^2}{\sum_{i=1}^m a_i^2 K_i^2 \log K_i}, \frac{u}{\|\mathbf{a}\|_\infty K^2 \log K} \right) \right],$$

where  $K = \max_i K_i$  and  $c$  is an absolute constant.

We now state a theorem that controls the singular values of a random matrix  $\mathbf{A}$ . The Theorem is exactly the same as Theorem 5.39 in Vershynin [2012] with the notable difference in the dependence of the constants to the sub-gaussian parameters. We use Theorem 5 to get this improved dependence.

**Theorem 6.** *Let  $\mathbf{A}$  be a  $N \times n$  matrix whose rows  $\mathbf{a}_i$  are independent sub-gaussian isotropic random vectors in  $\mathbb{R}^n$ . Then for every  $t \geq 0$ , with probability at least  $1 - 2 \exp(-ct^2)$  one has*

$$\sqrt{N} - C\sqrt{n} - t \leq s_{\min}(\mathbf{A}) \leq s_{\max}(\mathbf{A}) \leq \sqrt{N} + C\sqrt{n} + t. \quad (62)$$

Here  $C = C_K = K\sqrt{\log K} \sqrt{\frac{\log 9}{c_1}}$ ,  $c = c_K = \frac{c_1}{K^2 \log K} > 0$  with  $c_1$  is an absolute constant and  $K = \max_i \|\mathbf{a}_i\|_{\psi_2}$ .

The proof structure of Theorem 6 is exactly the same as the proof of Theorem 5.39 in Vershynin [2012], and so we provide the proof presented in Vershynin [2012] below.

*Proof.* The proof is a basic version of a covering argument, and it has three steps. We need to control  $\|\mathbf{A}\mathbf{x}\|_2$  for all vectors on the unit sphere. To this end, we discretize the sphere using a  $\mathcal{N}$  (the approximation step), establish a tight control of  $\|\mathbf{A}\mathbf{x}\|_2$  for every fixed vector  $\mathbf{x} \in \mathcal{N}$  with high probability (the concentration step), and finish off by taking a union bound over all  $\mathbf{x}$  in the net.

**Step 1: Approximation.** Using Lemma 5.36 in Vershynin [2012] for the matrix  $\mathbf{B} = \mathbf{A}/\sqrt{N}$  we see that the conclusion of the theorem is equivalent to

$$\left\| \frac{1}{N} \mathbf{A}^\top \mathbf{A} - \mathbf{I} \right\| \leq \max(\delta, \delta^2) =: \epsilon \text{ where } \delta = C \sqrt{\frac{n}{N}} + \frac{t}{\sqrt{N}}. \quad (63)$$

Using Lemma 5.34 in Vershynin [2012], we can evaluate the operator norm in (63) on a  $\frac{1}{4}$ -net  $\mathcal{N}$  of unit sphere  $\mathcal{S}^{n-1}$ :

$$\left\| \frac{1}{N} \mathbf{A}^\top \mathbf{A} - \mathbf{I} \right\| \leq 2 \max_{\mathbf{x} \in \mathcal{N}} \left| \left\langle \left( \frac{1}{N} \mathbf{A}^\top \mathbf{A} - \mathbf{I} \right) \mathbf{x}, \mathbf{x} \right\rangle \right| = 2 \max_{\mathbf{x} \in \mathcal{N}} \left| \frac{1}{N} \|\mathbf{A}\mathbf{x}\|_2^2 - 1 \right|.$$

So to complete the proof it suffices to show that, with high probability,

$$\max_{\mathbf{x} \in \mathcal{N}} \left| \frac{1}{N} \|\mathbf{A}\mathbf{x}\|_2^2 - 1 \right| \leq \frac{\epsilon}{2}. \quad (64)$$

By Lemma 5.2 in Vershynin [2012], we can choose the net  $\mathcal{N}$  so that it has cardinality  $|\mathcal{N}| \leq 9^n$ .

**Step 2: Concentration** Let us fix any vector  $\mathbf{x} \in \mathcal{S}^{n-1}$ . We can express  $\|\mathbf{A}\mathbf{x}\|_2^2$  as a sum of independent random variables

$$\|\mathbf{A}\mathbf{x}\|_2^2 = \sum_{i=1}^N \langle \mathbf{a}_i, \mathbf{x} \rangle =: \sum_{i=1}^N z_i^2 \quad (65)$$

where  $\mathbf{a}_i$  denote the rows of the matrix  $\mathbf{A}$ . By assumption,  $z_i = \langle \mathbf{a}_i, \mathbf{x} \rangle$  are independent sub-gaussian random variables with  $\mathbb{E} z_i^2 = 1$  and  $\|z_i\|_{\psi_2} \leq K$ . Therefore, by Remark 5.18 and Lemma 5.14 in Vershynin [2012],  $z_i^2 - 1$  are independent centered sub-exponential random variables with  $\|z_i^2 - 1\|_{\psi_1} \leq 2\|z_i^2\|_{\psi_1} \leq 4\|z_i\|_{\psi_2}^2 = 4K^2$ .

We can therefore use an exponential deviation inequality, Theorem 5, to control the sum (65).

$$\begin{aligned} \mathbb{P} \left\{ \left| \frac{1}{N} \|\mathbf{A}\mathbf{x}\|_2^2 - 1 \right| \geq \frac{\epsilon}{2} \right\} &= \mathbb{P} \left\{ \left| \frac{1}{N} \sum_{i=1}^N z_i^2 - 1 \right| \geq \frac{\epsilon}{2} \right\} \\ &\leq 2 \exp \left[ -\tilde{c}_1 \min \left( \frac{\epsilon^2 N^2 / 4}{\sum_{i=1}^N 4K_i^2 \log 2K_i}, \frac{\epsilon N / 2}{4K \log 2K} \right) \right] \\ &\leq 2 \exp \left[ -\frac{\tilde{c}_1}{4K^2 \log 2K} \min(\epsilon^2, \epsilon) N \right] \\ &= 2 \exp \left[ -\frac{\tilde{c}_1}{4K^2 \log 2K} \delta N \right] \\ &\leq 2 \exp \left[ -\frac{c_1}{K^2 \log K} (C^2 n + t^2) \right], \end{aligned}$$

where the last inequality follows by the definition of  $\delta$  and using the inequality  $(a+b)^2 \geq a^2 + b^2$  for  $a, b \geq 0$ .

**Step 3: Union bound.** Taking the union bound over all vectors  $\mathbf{x}$  in the net  $\mathcal{N}$  of cardinality  $|\mathcal{N}| \leq 9^n$ , we obtain

$$\mathbb{P} \left\{ \max_{\mathbf{x} \in \mathcal{N}} \left| \frac{1}{N} \|\mathbf{A}\mathbf{x}\|_2^2 - 1 \right| \geq \frac{\epsilon}{2} \right\} \leq 9^n \cdot 2 \exp \left[ -\frac{c_1}{K^2 \log K} (C^2 n + t^2) \right] \leq 2 \exp \left[ -\frac{c_1}{K^2 \log K} t^2 \right],$$

where the second inequality follows for  $C = C_K$  sufficiently large, e.g.  $C = K \sqrt{\log K} \sqrt{\frac{\log 9}{c_1}}$ .  $\square$

We now state a corollary of Theorem 6 that applies to general, non-isotropic sub-gaussian distribution.

**Corollary 1.** *Let  $\mathbf{A}$  be a  $N \times n$  matrix whose rows  $\mathbf{a}_i$  are independent sub-gaussian random vectors in  $\mathbb{R}^n$  with second moment matrix  $\Sigma$ . Then for every  $t \geq 0$ , with probability at least  $1 - 2 \exp(-ct^2)$  one has*

$$\left\| \frac{1}{N} \mathbf{A}^\top \mathbf{A} - \Sigma \right\| \leq \max(\delta, \delta^2) \text{ where } \delta = C \sqrt{\frac{n}{N}} + \frac{t}{\sqrt{N}}. \quad (66)$$

Here  $C = C_K = K \sqrt{\log K} \sqrt{\frac{\log 9}{c_1}}$ ,  $c = c_K = \frac{c_1}{K^2 \log K} > 0$  with  $c_1$  is an absolute constant and  $K = \max_i \|\mathbf{a}_i\|_{\psi_2}$ .