

## A APPENDIX

## CONTENTS

<a href="#">A.1 Stronger Generalization bounds for polynomials</a>	15
<a href="#">A.2 The Doubling Trick</a>	16
<a href="#">A.3 Amount of Overparameterization</a>	17
<a href="#">A.4 PSD property of <math>K^\infty</math></a>	17
<a href="#">A.5 Complexity upper bound</a>	18
<a href="#">A.6 Structure of the proof, setting <math>m_3</math>, and further definitions</a>	18
<a href="#">A.7 Complexity measure and the <math>\zeta</math>-norm</a>	20
<a href="#">A.8 Core Generalization Result</a>	22
<a href="#">A.9 Proof of Theorem 2</a>	24
<a href="#">A.10 Optimization</a>	25
<a href="#">A.11 Rademacher Complexity</a>	28
<a href="#">A.12 Constructing <math>W^*, V^*</math></a>	34
<a href="#">A.12.1 First Layer, Construction of <math>W^*</math></a>	34
<a href="#">A.12.2 Second Layer, Construction of <math>V^*</math></a>	39
<a href="#">A.12.3 Construction of <math>V^*</math></a>	45
<a href="#">A.13 Existence of a good direction</a>	48
<a href="#">A.14 Existence of a good direction Helper Lemmas</a>	53
<a href="#">A.15 Bounding the worst-case Scenario</a>	65
<a href="#">A.16 Convergence</a>	73
<a href="#">A.17 Process from a higher view: definition of the <math>(X)</math> sequence</a>	77
<a href="#">A.18 Bounding the MGF of <math>X_i</math>'s</a>	79
<a href="#">A.19 Proof of Theorem 7</a>	80
<a href="#">A.20 Gaussian Smoothing</a>	81
<a href="#">A.20.1 Setting <math>\beta_1</math> and <math>\beta_2</math></a>	86
<a href="#">A.21 Basic Tools</a>	87
<a href="#">A.21.1 Defining the rare events <math>E_i</math></a>	91
<a href="#">A.21.2 Bounding the value of <math>f'</math></a>	93
<a href="#">A.21.3 Bounding the difference between Original and Smoothed Functions</a>	94

### A.1 STRONGER GENERALIZATION BOUNDS FOR POLYNOMIALS

In this section, we prove an explicit generalization bound for functions represented as a polynomial sum. Note that the bounds in [Arora et al \(2019a\)](#) for polynomials assume the monomials with degree larger than one to have even powers, while here we do not impose this restriction. In addition, note that despite [Arora et al \(2019a\)](#), our bounds remain meaningful in the noisy case (recall our [Theorem 2](#)).

More specifically, we bound the  $\zeta$  norm of such functions. Consider the target function  $s$  with the following power series formula:

$$y = s(x) = \sum_{p=1}^{\infty} a_p (w_p^T x)^p, \quad (26)$$

where  $a_p \in \mathbb{R}$  and  $w_p \in \mathbb{R}^d$ . We can write

$$s(x) = g_1(x) + \sum_{k=1}^d x_k g_2^k(x), \quad (27)$$

where  $x_k$  denotes the  $k$ th entry of vector  $x$  here and

$$g_1(x) = \sum_{p \in A_1 := \{p=1 \text{ or } p \text{ even}\}} a_p (w_p^T x)^p,$$

and for all  $k \in [d]$ :

$$g_2^k(x) = \sum_{p \in A_2 := \{p>2, p \text{ odd}\}} w_{pk} a_p (w_p^T x)^{p-1}.$$

Then, using the Taylor series of  $x(\frac{1}{4} + \frac{\arcsin(x)}{2\pi}) = \sum_{p=1}^{\infty} \gamma_p x^p$  for  $|x| \leq 1$ , the RKHS space  $\mathcal{H}(H^\infty)$  of NTK can be identified by square-summable sequences of reals  $(a_{p'})_{p'=1}^{\infty}$  with dot product

$$\langle (a_{p'})_{p'=1}^{\infty}, (b_{p'})_{p'=1}^{\infty} \rangle = \sum_{p'=1}^{\infty} \gamma_{\lambda(p')} a_{p'} b_{p'},$$

where  $\lambda(p') : \mathbb{Z}_{>0} \rightarrow \mathbb{Z}_{>0}$  such that it maps zero to zero, the first  $d$  positive integers are mapped to one, the next  $d^2$  ones are mapped to 2, etc. Moreover, the RKHS mapping  $\Psi : \mathbb{R}^d \rightarrow \mathcal{H}(H^\infty)$  from the Euclidean space is:

$$\Psi(x) = \|x\| \left( x'_1, \dots, x'_d, (x'_{k_1} x'_{k_2})_{k_1, k_2 \in [d]}, \dots, (x'_{k_1} x'_{k_2} \dots x'_{k_p})_{k_1, \dots, k_p \in [d]}, \dots \right),$$

where  $x' = x/\|x\|$  and in the notation above we are presenting a sequence of sequences, by which we mean the inner sequences simply unfold. Using this identification, one can see using the linear representations of  $g_1, g_2^k$  in  $\mathcal{H}$ :

$$\|g_1\|_{H^\infty}^2 = \sum_{p \in A_1} \gamma_p a_p^2 \|w_p\|_2^{2p}, \quad (28)$$

$$\|g_2^k\|_{H^\infty}^2 = \sum_{p \in A_2} \gamma_p w_{pk}^2 a_p^2 \|w_p\|_2^{2(p-1)}, \quad (29)$$

Summing above and noting the linear representation of  $g$ :

$$\|g_1\|_{H^\infty}^2 + \sum_{k=1}^d \|g_2^k\|_{H^\infty}^2 = \sum_{p \in A_1} \gamma_p a_p^2 \|w_p\|_2^{2p} + \sum_{p \in A_2} \gamma_p a_p^2 \|w_p\|_2^{2p} = \sum_{p=1}^{\infty} \gamma_p a_p^2 \|w_p\|_2^{2p} = \|g\|_{H^\infty}^2. \quad (30)$$

Now for  $\{g\} := \{g'_k\}_{k=1}^{d+1} := \{g_1\} \cup \{g_2^k\}_{k=1}^d$ , we consider the kernel  $K_{\{g\}}$ . Expanding the Taylor series of  $F_2(2F_3)(x) = \sum_{p=0}^{\infty} \mu_p x^p$ , we find the identification  $(h_{p'}^k)_{k \in [d+1], p'=0, \dots, \infty}$  with dot product

$$= \sum_{p'=0}^{\infty} \mu_{\lambda(p')} \sum_{k=1}^{d+1} h_{p'}^k q_{p'}^k,$$

with RKHS map

$$\Psi_2(x) = \left( g'_1(x), \dots, g'_{d+1}(x), x'_1 g'_1(x), \dots, x'_1 g'_{d+1}(x), \dots, x'_d g'_1(x), \dots, \right. \\ \left. x'_d g'_{d+1}(x), (x'_{k_1} x'_{k_2} g'_k(x))_{k_1, k_2 \in [d], k \in [d+1]}, \dots, (x'_{k_1} x'_{k_2} \dots x'_{k_p} g'_k(x))_{k_1, \dots, k_p \in [d], k \in [d+1]}, \dots \right).$$

Now, we compute the norm of function  $s$  with respect to  $K_{\{g\}}$ , combining the above representation and dot product with Equation (27) and the fact that we work with unit norm  $x$ , so  $x' = x$ :

$$\|s\|_{K_{\{g\}}}^2 = \mu_0 + (d+1)\mu_1. \quad (31)$$

Plugging the above and Equation (30) into the definition of  $\|\cdot\|_\zeta$  in (??), we conclude

$$\|s\|_\zeta^2 \leq \|s\|_{K_{\{g\}}}^2 (\|g_1\|_{H^\infty}^2 + \sum_{k=1}^d \|g_2^k\|_{H^\infty}^2) \leq (\mu_0 + (d+1)\mu_1) \sum_{p=1}^{\infty} \gamma_p a_p^2 \|w_p\|_2^{2p}. \quad (32)$$

Note that if the odd exponents (except possibly one) in the definition of  $s$  in (26) are zero, then we could consider only the function  $g_1$  and kernel  $K_{g_1}$ , which would have implied a bound of  $\mu_0 \sum_{p=1}^{\infty} \gamma_p a_p^2 \|w_p\|_2^{2p}$ .

## A.2 THE DOUBLING TRICK

For the SGD optimization, we set the regularization coefficients in the loss  $L_1$  as

$$\psi_1 = \nu/4, \quad \psi_2 = \nu/(4\zeta(f^*, G)), \quad (33)$$

with  $\nu := \max\{R_n(f^*)/2, B^2/n\}$ . This assumes we know the  $f^*$  and  $G$  that minimize the adaptation within the complexity measure (12). Here, we explain how to use a simple doubling trick to get over the fact that we might not know these optimal solutions  $f^*$  and  $G$ .

**Theorem 1** *Without explicitly knowing the exact value of the complexity measure, i.e., the optimal solution of Equation (12), one can still achieve the generalization bound in Theorem 1.*

### Proof of Theorem 1

The proof is simply adding a doubling trick on top of the argument of Theorem 3. In the rest of the proofs, for simplicity, we refer to  $(W_{\text{PSGD}}, V_{\text{PSGD}})$  by  $(W', V')$ . Let  $f^{**}, G^*$  be the optimal solution to (10). With a simple rescaling of  $G^*$ , we can assume  $\langle H^{\infty-1}, G^* \rangle = 1$ . (Note that the complexity does not change by such rescaling). Now one can exploit the condition  $\|y_i\|_\infty \leq B$ , and consider the setting  $f^* = 0$  to get the following trivial upper bound on the complexity measure:

$$\mathfrak{S}((x_i), (y_i)) \leq 2nB^2.$$

Therefore,

$$2nR_n(f^{**}) + f^{**} A^{-1} f^{**} (c' \varpi) \leq 2nB^2. \quad (34)$$

Using Equation (34) and the optimality of  $(f^{**}, G^*)$ :

$$R_n(f^{**}) \leq B^2, \\ \zeta = \zeta(f^{**}, G^*) \leq 2nB^2. \quad (35)$$

Combining the first equation above with the definition of  $\nu$  in Equation (31), we get

$$B^2/n \leq \nu \leq B^2. \quad (36)$$

To initialize  $\psi_1$  and  $\psi_2$ , we use Equations (33) for any  $f^*$  and  $G$ , and as a result we get a generalization bound as in Equation (35). However, to achieve the best possible rate characterized by our complexity measure in Theorem 1 without explicitly computing the answer of (12), we use a simple doubling trick; for every pair  $(\zeta', \nu')$  such that  $\zeta'$  is a power of 2 between  $B^2$  and  $2nB^2$ , and  $\nu'$  is a power of two between  $B^2/n$  and  $B^2$ , we initialize  $\psi_1, \psi_2$  as in Equation (33) and run the algorithm,

then return the network that minimizes the upper bound of the risk given by Equation (54). Note that we do not need to know  $R_n(f^{**})$  to compute this upper bound. Now let  $\nu'$  be the power of 2 within  $\nu(G^*, f^{**})/2 \leq \nu' < \nu(G^*, f^{**})$ . If we are in the case

$$f^{**T} A^{-1} f^{**} < B^2, \quad (37)$$

then for  $\zeta'$  equal to the smallest power of two larger than  $B^2$ , when we run PSGD with pair  $(\nu', \zeta')$ , by Theorem 3:

$$R(f_{W', V'}) \leq 2R_n(f^*) + c''\varpi \frac{2B^2 + B^2}{n} \leq \frac{\mathfrak{S}((x_i)_{i=1}^n, (y_i)_{i=1}^n)}{n} + c' \frac{B^2\varpi}{n}. \quad (38)$$

Because we return the minimum upper bound on the risk (the tighter lower bound of Equation (54)) among all such powers of two, we certainly achieve the above rate in (38). Otherwise, if  $f^{**T} A^{-1} f^{**} \geq B^2$ , let  $\zeta'$  be the power of two within  $f^{**T} A^{-1} f^{**} \leq \zeta' \leq 2f^{**T} A^{-1} f^{**}$ , then again it is easy to check that conditions of Theorem 3 are satisfied, hence we get the following generalization bound:

$$R(f_{W', V'}) \leq 2R_n(f^{**}) + c''\varpi \frac{(\zeta' + B^2)}{n} \leq 2R_n(f^{**}) + c''\varpi \frac{2\zeta(f^{**}, G^*) + B^2}{n} \quad (39)$$

$$\leq \frac{\mathfrak{S}((x_i)_{i=1}^n, (y_i)_{i=1}^n)}{n} + c' \frac{B^2\varpi}{n}, \quad (40)$$

which proves the bound of Theorem 3.

### A.3 AMOUNT OF OVERPARAMETERIZATION

In this section, to provide high-level insight, we indicate the right order of magnitude that our overparameterization should be in, with respect to one another. Note that the exact coefficients in these inequalities would depend on the basic parameters  $B, 1/\lambda_0, n$ , which we have avoided here for sake of simplicity. We refer the reader to our main proof (mostly Appendix parts A.12, A.13) for more details.

$$\begin{aligned} \kappa_1 \kappa_2 m_3 &\ll 1, \\ \kappa_2 \sqrt{m_2} &\gg 1, \\ \kappa_1 \sqrt{m_3} &\gg \kappa_2 \sqrt{m_2}, \\ m_1 &\gg m_3^4, \\ \kappa_1 \sqrt{m_1} &\gg m_3^{3/2}, \\ \kappa_2 &\ll 1/\sqrt{m_3} \\ \sqrt{m_3} \kappa_2 &\ll 1/\sqrt{m_3} \\ \sqrt{m_2} &\gg m_3^{3/2} \kappa_1 \kappa_2 \\ m_3^3 (\kappa_2 m_2) &\ll \kappa_1 m_1 \\ m_1, m_2, m_3, 1/\kappa_1, 1/\kappa_2 &= \text{poly}(n, B \vee 1/B, 1/\lambda_0). \end{aligned}$$

In addition, we set the smoothing parameters as

$$\begin{aligned} \beta_2 &:= \Theta_p \left( (\kappa_1 \sqrt{m_3})^{-1} (\sqrt{m_2} \kappa_2)^{-\frac{2}{3}} \right), \\ \beta_1 &:= \Theta_p \left( m_3 \sqrt{m_3} / (\kappa_1 \sqrt{m_1}) \right), \end{aligned}$$

where  $\Theta_p$  only shows polynomial dependencies on the overparameterization.

### A.4 PSD PROPERTY OF $K^\infty$

The Schur product theorem states that for PSD matrices  $A$  and  $B$ ,  $A \odot B$  is also PSD. Now given an analytic function  $F$  whose Taylor series coefficients are all nonnegative, Suppose we apply  $F$  on some PSD matrix  $A$  entrywise, denoted by  $F(A)$ , under the condition that the entries of  $A$  are in

the radius of convergence of  $F$ , then using Schur product theorem, it is straightforward that  $F(A)$  is also PSD.

Using the above property, one can then check that the Taylor series of the defined functions  $F_2$  and  $F_3$  are nonnegative, hence, the application of the function  $F_2(2F_3(x))$  on the gram matrix of  $(x_i)_{i=1}^n$  is a PSD matrix, (note that  $|\langle x_i, x_j \rangle| \leq 1$  is in the convergence radius of  $F_2(2F_3(x))$ .) thus  $K^\infty$  is indeed a kernel.

#### A.5 COMPLEXITY UPPER BOUND

First we mention a simple fact that hadamard product respects matrix orderings. Given PSD matrices  $A, B, C$  such that  $A \preceq B$ , the fact that  $A \odot C \preceq B \odot C$  is an easy consequence of the Schur Product Theorem; indeed,  $B - A$  is PSD by definition, so  $(B - A) \odot C = B \odot C - A \odot C$  is also PSD.

Next, it is easy to check that the Taylor series of  $\arcsin(x)$  has all nonnegative coefficients. Therefore, for a PSD matrix  $X$ , as we discussed in Appendix A.4, applying  $\arcsin$  entrywise on  $X$ , namely  $\arcsin X$ , is also PSD. Setting  $X$  equal to the entrywise application of  $2F_3$  to the gram matrix of datapoints  $(x_i)_{i=1}^n$ , we realize the matrix  $\arcsin\left(2F_3\left(\left(\langle x_i, x_j \rangle\right)_{1 \leq i, j \leq n}\right)\right)$  is also PSD. Noting the definition of  $K^\infty$  in Equation (8), we conclude that for the data kernel matrix  $K^\infty$  we have

$$K^\infty \geq \frac{1}{4} \mathbb{1} \mathbb{1}^T,$$

where  $\mathbb{1}$  is the all ones  $n$ -dimensional vector.

Combining the two mentioned facts, we can lower bound the matrix  $K = K^\infty \odot G$  for any matrix  $G$  as

$$K = K^\infty \odot G \geq \frac{1}{4} \mathbb{1} \mathbb{1}^T \odot G = \frac{1}{4} G.$$

Substituting the rank one matrix  $f^* f^{*T}$  for the  $n$ -dimensional vector  $f^*$  in Equation (13):

$$K_{\{f^*/\|f^*\|}}^{-1} = K^\infty \odot f^* f^{*T} / \|f^*\|^2 \geq \frac{1}{4} f^* f^{*T} / \|f^*\|^2. \quad (41)$$

The inequality used in (13) then follows from Equation (41).

#### A.6 STRUCTURE OF THE PROOF, SETTING $m_3$ , AND FURTHER DEFINITIONS

Almost all of our proofs in the rest are in the aim of proving Theorem B. Throughout the proof,  $(W', V')$  represents the pair of matrices of the current iteration of PSGD,  $(W^*, V^*)$  are the ‘‘ideal’’ matrices that we construct in Appendix A.12,  $(W^\rho, V^\rho)$  and refers to the gaussian smoothing matrices. Importantly, we assume the loss function  $\ell(f, y)$  is **zero** at  $f = y$  for any label.

There are four main parts in our proof:

1. we construct a ‘‘good’’ underlying network, Appendix A.12
2. we find a ‘‘good’’ random direction and study the landscape of the objective, Appendix A.13
3. we bound the Rademacher Complexity of the class  $\mathcal{G}_{\gamma_1, \gamma_2}$ , Appendix A.14
4. we prove the convergence, Appendix A.16

We have tried to make the lower level proofs into sub-lemmas and create a manageable hierarchy as much as we could, to make the document more clear and readable.

Similar to the conditions in Theorem B, through out most of the proofs we assume that we are given a pair  $(f^*, G)$  with a slightly more general setting of Theorem B:

$$\begin{aligned} f^{*T} A^{-1} f^* &\leq \zeta_2, \text{ for } A = G \odot K^\infty, \\ \langle G, H^\infty \rangle &\leq \zeta_1. \end{aligned}$$

Particularly,  $\zeta_1, \zeta_2$  appear in Appendix [A.13](#). Because we are allowed to rescale  $G$ , we do not really gain much by assuming this more general setting, though we pick to work with the general setting as the abstraction makes the proof more straightforward to understand.

We refer to the parameters  $B, 1/\lambda_0, n$  as the “basic parameters”,  $m_1, m_2, m_3, 1/\kappa_1, \kappa_2$  as the “overparameterization”, and  $\beta_1, \beta_2$  as the “smoothing parameters.” By the phrase “having enough overparameterization” we mean it suffices to pick the overparameterization  $m_1, m_2, m_3, 1/\kappa_1, 1/\kappa_2$  only **polynomially** large in the basic parameters.

Throughout the proof, we denote the change in the output of the first layer at  $W^{(0)} + W' + W^\rho$  compared to the initialization value by  $\phi^{(2)}(x_i)$ , i.e.

$$\phi^{(2)}(x_i) = \frac{1}{\sqrt{m_1}} W^s \sigma((W^{(0)} + W' + W^\rho)x_i) - \phi^{(0)}(x_i),$$

while recall that  $\phi'(x_i)$  has a similar definition except without the smoothing matrix  $W^\rho$ . Although our model is a three layer network, throughout the proof, we refer to the parts  $W^s \frac{1}{\sqrt{m_1}} \sigma((W^{(0)} + W')x)$  and  $\frac{1}{\sqrt{m_2}} a^T \sigma((V^{(0)} + V')(\cdot))$  as the “first layer” and “second layer,” respectively.

Also, we sometimes refer to the binary sign pattern of vector  $x$  multiplied to matrix  $W$  by  $D_{W,x}$  ( $D_{W,x} := \text{Sgn}(Wx)$ ), i.e. the  $j$ th diagonal entry of  $D_{W,x}$  is one if  $W_j^T x \geq 0$ , and is zero otherwise. To refer to the  $j$ th row of  $W$  as a vector, we sometimes drop the comma in  $W_j$ , and write it as  $W_j$ .

For brevity, we denote the Frobenius norm  $\|W\|_F$  of matrix by  $\|W\|$ , and the Euclidean norm of a vector  $x$  by  $\|x\|$ . For matrices  $W_1, W_2$  we denote their natural dot product by  $\langle W_1, W_2 \rangle := \text{tr}(W_1^T W_2)$ . We refer to the smallest eigenvalue of a matrix by  $\lambda_{\min}(\cdot)$ . We write  $\mathcal{R}(\cdot)$  for the Rademacher complexity of a function class. We refer to the smoothed version of the network by  $f'_{W', V'}(x)$ , defined by

$$f'_{W', V'}(x) = \mathbb{E}_{W^\rho, V^\rho} f_{W'+W^\rho, V'+V^\rho}(x).$$

In the proof, we mainly work with the loss over the smoothed network  $f'$ , defined as

$$L(W', V') = R_n(f'_{W', V'}) + \psi_1 \|W'\|^2 + \psi_2 \|V'\|^2. \quad (42)$$

Our algorithm, PSGD can be regarded roughly as an SGD over  $L$ .

Similar to what we discussed in section [3](#), let the functions  $\{g_k\}_{k=1}^{m_3}$  be some feature representation whose gram matrix is equal to  $G$  and  $\langle H^\infty, G \rangle = \sum_{k=1}^{m_3} \|g_k\|_{H^\infty}^2$ . In such setting, it is not hard to check that we can assume each  $g_k$  is the minimum norm NTK function which maps  $(x_i)_{i=1}^n$  to  $(g_k(x_i))_{i=1}^n$ 's. Indeed, if this is not the case for some  $g_k$ , we can project the RKHS representation of  $g_k$  onto the span of the representations of  $(x_i)_{i=1}^n$ , which can only decrease the complexity measure. Hence we can represent  $g_k \in \mathcal{H}_{H^\infty}$  as a linear combination of basic functions  $H^\infty(x_i, \cdot)$  on data points:

$$\forall k \in [m_3], g_k(x) := \sum_{i=1}^n \mathcal{V}_{k,i} H^\infty(x_i, x). \quad (43)$$

Here, the sum of squared- $H^\infty$  norms of  $\mathcal{V}_k$  is bounded as

$$\sum_k \|\mathcal{V}_k\|_{H^\infty}^2 = \sum_k \|g_k\|_{H^\infty}^2 = \langle H^\infty, G \rangle \leq \zeta_1. \quad (44)$$

For each  $i \in [n]$ , we refer to the feature representation vector  $(g_k(x_i))_{k=1}^{m_3}$  on  $x_i$  as  $\bar{x}_i$ . Note that we have the relationship

$$\bar{x}_i = (H_{i,\cdot}^\infty \mathcal{V}_k)_{k=1}^{m_3}, \quad (45)$$

where  $H_{i,\cdot}^\infty$  is the  $i$ th row of  $H^\infty$ . In the analysis, we work with a bound  $\xi$  on the quantity  $\max_k \|\mathcal{V}_k\|$  which should be bounded polynomially by other basic parameters; in particular, it is defined in Lemma [10](#) and is used to bound a cross term in Lemma [14](#). However,  $\max_k \|\mathcal{V}_k\|$  might not be effectively bounded for an arbitrary feature representation. Fortunately, we can remedy this by a simple trick; for every natural number  $s$ , one can substitute every  $g_k$  by  $s$  copies of  $g_k/\sqrt{s}$ , without changing the gram matrix  $G$ . Therefore, for any  $\delta$ , one can increase the multiset of functions  $(g_k)$

to a bigger set  $(\tilde{g}_k)$ , by adding at most  $O(\zeta_1/\delta)$  functions, making sure of the following for the new functions:

$$\forall k : \mathcal{V}_k^T H^\infty \mathcal{V}_k = \|\tilde{g}_k\|_{H^\infty}^2 \leq \delta. \quad (46)$$

(This is because  $\sum_k \|\tilde{g}_k\|_{H^\infty}^2 \leq 1$ ). Furthermore, observe that for each gram matrix  $G$ , we have an  $n$ -dimensional feature representation  $(g_k)_{k=1}^n$  for  $G$  according to the Cholesky factorization. Combining these facts, we conclude, to guarantee Equation (46), in the worst case, we need  $m_3$  to be as large as  $n + O(\zeta_1/\delta)$ .

Finally, observing the following inequality

$$\|\mathcal{V}_k\|^2 \leq \|\mathcal{V}_k\|_{H^\infty}^2 / \lambda_0 \leq \|\tilde{g}_k\|_{H^\infty}^2 / \lambda_0. \quad (47)$$

in order to guarantee  $\max_k \|\mathcal{V}_k\| \leq \xi$  we need to take  $m_3$  as large as  $n + O(\zeta_1/(\xi^2 \lambda_0))$ , which is indeed bounded polynomially by the basic parameters because of the same condition for  $1/\xi$ . This computation also brings into sight an important point:

“Although each gram matrix  $G$  is representable by  $n$  features, in order for the algorithm to be able to find a suitable network,  $m_3$  might need to be larger than  $n$ .”

Moreover, for every  $1 \leq k, i \leq n$ , we define the matrices  $Z_k^i \in \mathbb{R}^{m_1}$  as

$$Z_k^i = 1/\sqrt{m_1} \left( W_{k,j}^s \mathbb{1}\{W_j^{(0)T} x_i\} x_i \right)_{j=1}^{m_1}, \quad (48)$$

where in the above notation,  $j$  is enumerating the columns of the matrix. We also define the following matrices which we use in our construction later:

$$W_k^{+T} = W^{k+T} = \sum_{i=1}^n \mathcal{V}_{k,i} Z_k^i, \quad (49)$$

and  $W^+$  as

$$W^+ = \sum_{k=1}^{m_3} W^{k+}.$$

Finally, to avoid unnecessary complication, we often argue **high probability** bounds without an explicit representation of their dependency on the chance of failure (which is a negligible logarithmic factor). We also ignore all constants and log factors, and mainly work with the notation  $\lesssim$  which ignores constants; we write  $a \lesssim b \pm c$  as a short form for  $b - O(c) \leq a \leq b + O(c)$ . As there are several hierarchies of new parameters that are defined based on lower-level ones, we rename the new parameters and continue viewing them as black-box. This makes the proofs more readable, since we also do not care about the exact dependency of the underlying parameters most of the time, rather we are interested in their orders of magnitude, for example that a given parameter goes to zero polynomially fast with respect to the overparameterization, etc. Due to the large number of symbols that we have to work with, we might use a symbol more than once, of course when it is clear from the context which one we are referring to.

## A.7 COMPLEXITY MEASURE AND THE $\zeta$ -NORM

This is a brief section regarding some basic properties of  $\mathfrak{S}$  and  $\|\cdot\|_\zeta$ .

First, note that the two versions of the complexity measure in Equations (40) and (41) are equivalent, as for any finite set of functions  $\{g\}$ , we can define the gram matrix with respect to the feature vectors of these functions on data, and for an arbitrary nonzero PSD  $G$  we can consider a Cholesky factorization for  $G$  as  $G = \bar{X}^T \bar{X}$ , then define the functions  $\{g_k\}$  as the minimum-NTK norm functions which map the input to the features corresponding to  $\bar{X}$ . This observation further implies we can suppose the factor matrix  $\bar{X}$  is in  $\mathbb{R}^{n \times n}$ , and there is a set of at most  $n$  functions  $\{g_k\}_{k=1}^n$  which corresponds to this  $G$ .

Next, we show that for an arbitrary function  $f$ , its sup norm over the unit ball is bounded by its  $\zeta$  norm:

$$\sup_{\|x\|=1} |f(x)| \leq \|f\|_\zeta. \quad (50)$$

Note that for a kernel  $K$  which satisfies  $K(x, x) \leq 1$ , using Cauchy Schwarz we simply obtain

$$f(x) \leq \|f\|_K \sqrt{K(x, x)} \leq \|f\|_K,$$

where recall that  $\|\cdot\|_K$  is the norm corresponding to the RKHS space of  $K$ . Hence, to show (50), it suffice to show that for all kernels  $K \in \mathcal{K}$  and unit norm  $x$  we have  $K(x, x) \leq 1$ . To see this fact, note that the norm of each  $x \in \mathbb{R}^d$  in the NTK-space is  $H^\infty(x, x) = \frac{1}{2}$ . Therefore, for each function  $g$  with bounded-NTK norm, again using Cauchy Schwarz:

$$|g(x)| \leq \frac{1}{2} \|g\|_{H^\infty}.$$

As a result, for a family of functions  $\{g\}$  with  $\sum_{g \in \{g\}} \|g\|_{H^\infty}^2 \leq 1$ , we have on every unit norm  $x$ :

$$\sum_{g \in \{g\}} g(x)^2 \leq 1.$$

On the other hand, it is easy to check that for every unit norm  $x$ , we have  $K^\infty(x, x) \leq \frac{1}{2}$ , so for every such  $\{g\}$ , we have by definition

$$K_{\{g\}}(x) \leq 1,$$

which completes the proof of Equation (50).

## A.8 CORE GENERALIZATION RESULT

In this section, we prove our core generalization result for the trained network, Theorem B, which underlies our generalization bounds in Theorems D and E. Recall that in the rest of the proofs, we refer to the solution  $(W_{\text{PSGD}}, V_{\text{PSGD}})$  returned by PSGD simply by  $(W', V')$ .

**Theorem 3** Suppose we have a candidate pair  $(f^*, G)$  regarding our complexity measure in (17) that satisfies  $\langle H^{\infty-1}, G \rangle \leq 1$ ,  $f^{*T} A^{-1} f^* \leq \zeta$  (recall  $A = G \odot H^\infty$ ). Then, for

$$\nu = \max\{R_n(f^*)/4, B^2/n\}, \quad (51)$$

if we are given  $\nu/2 \leq \nu' \leq \nu$ , and we set

$$\psi_1 = \frac{\nu'}{4}, \quad (52)$$

$$\psi_2 = \frac{\nu'}{4\zeta}, \quad (53)$$

then for the solution  $(W', V')$  returned by PSGD we have the following generalization bound:

$$R(f_{W', V'}) \leq \frac{5}{4} R_n(f_{W', V'}) + c''' \varpi \frac{(\zeta + B^2)}{n} \quad (54)$$

$$\leq \frac{5}{3} R_n(f^*) + c'' \varpi \frac{(\zeta + B^2)}{n}, \quad (55)$$

for constants  $c''$ ,  $c'''$  and log factor  $\varpi = \log(n)^3 + \log(1/\lambda_0)$ .

**Proof of Theorem B**

To prove Theorem B, we use a generalization bound from Srebro et al (2010). Crucially, in order to apply this bound, we need to establish two big results:

1. We need to show that the final network has small training loss, and is within the class  $\mathcal{G}_{\gamma_1, \gamma_2}$  for some suitable  $\gamma_1, \gamma_2$ . This is handled by Theorem A in Appendix A.10. We define the class  $\mathcal{G}_{\gamma_1, \gamma_2}$  roughly as the class of networks with norm bounds  $\|W - W^{(0)}\| \leq \gamma_1$ ,  $\|V - V^{(0)}\| \leq \gamma_2$  where the rows of  $V - V^{(0)}$  are orthogonal to the subspace  $\Phi$ , plus an additional structure defined in Appendix A.11.
2. The Rademacher Complexity of the class  $\mathcal{G}_{\gamma_1, \gamma_2}$  needs to be suitably bounded. This is handled by Theorem B in Appendix A.11.

With access to these results, we show how Theorem B follows. For fixed constants  $z_1, z_3$  and every integer  $i \geq 0$ , we use Theorem 1 in Srebro et al (2010) for the class  $\mathcal{G}_{z_1, \gamma_i}$ ,  $\gamma_i = 2^i \times B/z_3$  with confidence probability  $1 - 2^{-i}\delta_3$ , which, with a union bound, implies that with probability at least  $1 - \delta_3$  for every  $i$  and  $f_{W', V'} \in \mathcal{G}_{z_1, \gamma_i}$ :

$$R(f_{W', V'}) \leq R_n(f_{W', V'}) + K(\sqrt{R_n(f_{W', V'})}(\sqrt{\log(n)^{1.5} \mathcal{R}(\mathcal{G}_{z_1, \gamma_i})} + \sqrt{\frac{b \log(1/(2^{-i}\delta_3))}{n}})) \quad (56)$$

$$+ \log(n)^3 \mathcal{R}(\mathcal{G}_{z_1, \gamma_i})^2 + \frac{b \log(1/(2^{-i}\delta_3))}{n}, \quad (57)$$

where  $\ell(f_{W', V'}(x), y)$  is a.s. bounded by  $b$  for function within the class  $\mathcal{G}_{z_1, \gamma_i}$ , and  $K$  is a universal constant. In the following, we aim to further bound the Rademacher complexity  $\mathcal{R}$  and parameter  $b$ .

Applying the AM-GM inequality with respect to ratio  $z_4 > 0$  for the second term:

$$\begin{aligned}
R(f_{W',V'}) &\leq (1+z_4)R_n(f_{W',V'}) + K^2/z_4 \left( \log(n)^{1.5} \mathcal{R}(\mathcal{G}_{z_1,\gamma_i}) + \sqrt{\frac{b \log(1/(2^{-i}\delta_3))}{n}} \right)^2 \\
&\quad + \log(n)^3 \mathcal{R}(\mathcal{G}_{z_1,\gamma_i})^2 + \frac{b \log(1/(2^{-i}\delta_3))}{n} \\
&\leq (1+z_4)R_n(f_{W',V'}) + K^2/z_4 \left( \log(n)^{1.5} \mathcal{R}(\mathcal{G}_{z_1,\gamma_i}) + \sqrt{\frac{b \log(1/(2^{-i}\delta_3))}{n}} \right)^2 \\
&\quad + \log(n)^3 \mathcal{R}(\mathcal{G}_{z_1,\gamma_i})^2 + \frac{b \log(1/(2^{-i}\delta_3))}{n} \\
&\leq (1+z_4)R_n(f_{W',V'}) + (2K^2/z_4 + 1) \log(n)^3 \mathcal{R}(\mathcal{G}_{z_1,\gamma_i})^2 + (2K^2/z_4 + 1) \frac{b \log(1/(2^{-i}\delta_3))}{n}.
\end{aligned} \tag{58}$$

Now let  $\gamma^*$  be the smallest number of the form  $2^i B/z_3$  (for some  $i$ ) which is not smaller than  $z_2 \sqrt{\zeta}$ . This definition implies

$$\gamma^* \leq \max\{2z_2 \sqrt{\zeta}, B/z_3\}. \tag{59}$$

Now Theorem 5 in Appendix A.1 bounds the Rademacher complexity:

$$\mathbb{R}(\mathcal{G}_{z_1,\gamma^*}) \leq \frac{2z_1 \gamma^*}{\sqrt{n}}. \tag{60}$$

On the other hand, from Theorem 4 by setting  $z_1, z_2 = \sqrt{40}$ , we get  $f_{W',V'} \in \mathcal{G}_{z_1,\gamma^*}$ .

Moreover, from Lemma 54, for  $f_{W',V'} \in \mathcal{G}_{z_1,\gamma^*}$ , we have for every  $\|x\| \leq 1$ :

$$|f_{W',V'}(x)| \leq 2z_1 \gamma^*, \tag{61}$$

so the loss  $\ell(f_{W',V'}(x), y)$  can be bounded by  $(B + 2z_1 \gamma^*)^2$  using the 1 smoothness property. Therefore, for the class  $\mathcal{G}_{z_1,\gamma^*}$  we can set  $b = (B + 2z_1 \gamma^*)^2$ . Combining this with Equation (60) and plugging into Equation (58):

$$R(f_{W',V'}) \leq (1+z_4)R_n(f_{W',V'}) + (2K^2/z_4 + 1) \log(n)^3 \frac{4z_1^2 \gamma^{*2}}{n} + (2K^2/z_4 + 1) \frac{(B + 2z_1 \gamma^*)^2 \log(1/(2^{-i^*} \delta_3))}{n}.$$

Furthermore, by definition of  $\gamma^*$ , we have  $2^i \leq 2z_2 z_3 \sqrt{\zeta}/B$ :

$$R(f_{W',V'}) \leq (1+z_4)R_n(f_{W',V'}) + (2K^2/z_4 + 1) 4z_1^2 (\log(n)^3 + 2 \log(2z_2 z_3 \sqrt{\zeta}/B)) \frac{4z_1^2 \gamma^{*2}}{n} \tag{62}$$

$$+ (2K^2/z_4 + 1) \frac{2B^2 \log(2z_2 z_3 \sqrt{\zeta}/B)}{n}. \tag{63}$$

Now applying the upper bound on  $\gamma^*$ :

$$R(f_{W',V'}) \leq (1+z_4)R_n(f_{W',V'}) + (2K^2/z_4 + 1) 4z_1^2 (\log(n)^3 + 2 \log(2z_2 z_3 \sqrt{\zeta}/B)) \frac{4z_1^2 (2z_2 \sqrt{\zeta} + 2B/z_3)^2}{n} \tag{64}$$

$$+ (2K^2/z_4 + 1) \frac{2B^2 \log(2z_2 z_3 \sqrt{\zeta}/B)}{n}. \tag{65}$$

If  $\zeta > B^2$ , in the third term above we substitute  $B$  by  $\sqrt{\zeta}$ . Finally, similar to the bound we stated in Equation (15), note that we have the following trivial bound for  $\zeta$ :

$$\zeta \leq y^T H^{\infty-1} y \leq 4nB^2/\lambda_0, \tag{66}$$

i.e. there is no point in considering larger  $\zeta$ 's, which implies  $\log(2z_2 z_3 \sqrt{\zeta}/B) = O(\log(n) + \log(1/\lambda_0))$ . Plugging this above and picking  $z_4 = 1/3$  show the proof of Equation (54). Furthermore, applying Equation (70) in Theorem 4 to the  $R_n(f_{W',V'})$  term in Equation (54) further gives the second Equation (55).

**Remark 1** In the same setting of Theorem 3, if we have  $\nu/2 \leq \nu'$  but not generally upper bounded by  $\nu$ , then PSGD leads to the following generalization bound:

$$R(f_{W',V'}) \leq R_n(f^*) + \nu' + c'''\varpi \frac{(\zeta + B^2)}{n},$$

using a similar argument as we did for Theorem 3.

#### A.9 PROOF OF THEOREM 2

In this section, we prove Theorem 2, stated below.

**Theorem 2** For any function  $f : \mathbb{R}^d \rightarrow \mathbb{R}$ , in the same setting as Theorem 1, the population risk of the trained network  $(W', V')$  can be bounded as

$$R(f_{W',V'}) \leq 2R(f) + O\left(\alpha\varpi \frac{\|f\|_\zeta^2 + B^2}{n}\right). \quad (67)$$

#### Proof of Theorem 2

Theorem 2 is a simple consequence of Theorem 3; for the given function  $f$ , we apply Theorem 3 with the smaller coefficient  $\gamma = \frac{4}{3}$  for  $R_n(f^*)$ , regarding the complexity upper bound, by setting  $f^* := (f(x_i))_{i=1}^n$ :

$$\begin{aligned} R(f_{W',V'}) &\leq \frac{4}{3}R_n(f^*) + (\alpha\varpi) \min_{K \in \mathcal{K}} \frac{f^{*T} K^{-1} f^*}{n} + \frac{B^2 \alpha \varpi}{n} \\ &= \frac{4}{3}R_n(f) + (\alpha\varpi) \min_{K \in \mathcal{K}} \frac{f^{*T} K^{-1} f^*}{n} + \frac{B^2 \alpha \varpi}{n}. \end{aligned}$$

On the other hand, because  $f^{*T} K^{-1} f^*$  is the minimum-RKHS norm of a function with respect to kernel  $K$  which maps  $x_i$ 's to  $f^*_i$  and  $f$  is one such function, we have  $f^{*T} K^{-1} f^* \leq \|f\|_K$ . This inequality implies

$$\min_{K \in \mathcal{K}} f^{*T} K^{-1} f^* \leq \|f\|_\zeta,$$

so we obtain

$$R(f_{W',V'}) \leq \frac{4}{3}R_n(f) + (\alpha\varpi) \frac{\|f\|_\zeta^2 + B^2}{n}. \quad (68)$$

Therefore, it remains to bound  $R_n(f)$  by  $R(f)$ .

As we showed in Appendix A.7, for every input  $x$  we have  $f(x) \leq \|f\|_\zeta$ , so for every data  $(x, y)$ , by the fact that  $|y| \leq B$  a.s. and  $\alpha$  smoothness of the loss, we have  $\ell(f(x), y) \leq \alpha(\|f\|_\zeta + B)^2$ . Moreover, note that the random variable  $\ell(f(x), y)$  has mean  $R(f)$ . It is easy to check that in this setting, the variance of  $\ell(f(x), y)$  is at most  $R(f)\alpha(\|f\|_\zeta + B)^2$ . Therefore, an application of the Bernstein inequality, we have with high probability over the dataset

$$R_n(f) \leq R(f) + O\left(\sqrt{\frac{R(f)\alpha(\|f\|_\zeta + B)^2}{n}} + \frac{\alpha(\|f\|_\zeta + B)^2}{n}\right) \leq \frac{3}{2}R(f) + O\left(\frac{\alpha(\|f\|_\zeta + B)^2}{n}\right).$$

Plugging this back to Equation (68) completes the proof. As a result, the learned network can compete with any function that has reasonably small  $\|f\|_\zeta$ :

$$R(f_{W',V'}) \leq \min_f \left\{ 2R(f) + O\left(\alpha\varpi \frac{\|f\|_\zeta^2 + B^2}{n}\right) \right\}.$$

## A.10 OPTIMIZATION

In this section, we glue together

- the existence of a good random direction that we prove in Appendix [A.13](#)
- the convergence analysis of PSGD that we do based on the work [Ge et al. \(2015b\)](#) in Appendix [A.16](#).

**Theorem 4** *In the same setting of Theorem [3](#), assume the network  $(W', V')$  returned by PSGD, has sufficient polynomially large “overparameterization”. Then, for the solution  $(W', V')$  returned by PSGD we have*

$$L(W', V') \leq R_n(f^*) + \nu, \quad (69)$$

which further implies

$$R_n(f_{W', V'}) \leq R_n(f^*) + 2\nu, \quad (70)$$

$$\|W'\|^2 \leq 40, \|V'\|^2 \leq 40\zeta. \quad (71)$$

Moreover, for every  $i \in [n], j \in [m_1], j \notin P$  for  $P$  defined in [4](#), we have that  $\text{sign}((W_j^{(0)} + W'_j)^T x_i)$  and  $\text{sign}(W_j^{(0)T} x_i)$  are the same.

**Proof of Theorem [4](#)**

Let  $\Upsilon \in \mathbb{R}^{m_2(m_3-n) \times m_2 m_3}$  be a matrix whose rows are an orthonormal basis for the space of matrices whose rows are orthogonal to  $\text{span}(\{\phi^{(0)}(x_i)\}_{i=1}^n)$ , i.e.  $\Phi^\perp$ , as defined in [\(24\)](#). Then, we consider a linear change of coordinates for the subspace  $\Phi^\perp$ , regarding the second layer weights, as  $v' = \Upsilon \text{vec}(V')$  where  $\text{vec}(\cdot)$  splits out the vectorized version of a matrix. For consistent notation, we also denote  $w' = W'$ , so we now have a new coordinate system  $(w', v') \in \mathbb{R}^{m_2(m_3-n) \times m_1 d}$  for pairs of weights  $(W', V')$  such that  $V' \in \Phi^\perp$ . We also define the loss function

$$L^\Pi(w := (w', v')) = L(W', V'),$$

with respect to the change of coordinate.

Now it is easy to see that running PSGD on  $L$  in the normal coordinates is equivalent to running stochastic gradient descent on  $L^\Pi$  with respect to  $(w', v')$ . Moreover, because multiplying to matrix  $\Upsilon$  is an orthonormal change of coordinates for  $\Phi^\perp$  and because  $V'$  is already in  $\Phi^\perp$  at each step of PSGD, then  $\|v'\| = \|V'\|$ , so the conditions  $\|W'\| \leq C_1, \|V'\| \leq C_2$  are equivalent to  $\|w'\| \leq C_1, \|v'\| \leq C_2$ . Furthermore, by our construction, the random matrix  $V_\Sigma^*$  is in the subspace  $\Phi^\perp$ , so the norm bounds  $\|W^*\| \leq \zeta_1, \|V^*\| \leq \zeta_2$  are equivalent to  $\|w^*\| \leq \zeta_1, \|v^*\| \leq \zeta_2$  for  $w^* := W_\Sigma^*$  and  $v^* := \Upsilon \text{vec}(V_\Sigma^*)$ .

Now we apply the result of Theorem [6](#) on  $L^\Pi$  with parameter  $\nu$  set as  $\nu', \zeta_2 := \zeta$  and  $\zeta_1 := 1$ , and  $\Delta := R_n(f^*)$ , as defined in Theorem [3](#). More specifically, based on our arguments above regarding the natural isometry in the change of coordinate, any pair  $(w', v')$  in the domain  $\|w'\| \leq C_1, \|v'\| \leq C_2$ ,  $L^\Pi(w', v') \geq R_n(f^*) + \nu'$  translates into a pair  $(W', V')$  in the domain  $\|W'\| \leq C_1, \|V'\| \leq C_2$ ,  $L(W', V') \geq R_n(f^*) + \nu'$ , for which by Theorem [6](#) there exists  $(W_\Sigma^*, V_\Sigma^*)$  such that

$$\mathbb{E}_\Sigma L(W' - \eta/2W' + \sqrt{\eta}W_\Sigma^*, V' - \eta/2V' + \sqrt{\eta}V_\Sigma^*) \leq L(W', V') - \eta\nu'/4. \quad (72)$$

Translating back to the change of coordinates:

$$\mathbb{E}_\Sigma L^\Pi(w' - \eta/2w' + \sqrt{\eta}w_\Sigma^*, v' - \eta/2v' + \sqrt{\eta}v_\Sigma^*) \leq L(w', v') - \eta\nu'/4. \quad (73)$$

Now we apply Lemma [44](#) to translate this into an argument about the landscape of  $L^\Pi$ . As a result, applying the bounds in Equations [\(107\)](#) and [\(127\)](#), we obtain that for  $(w', v')$  such that

$$L^\Pi(w', v') \geq R_n(f^*) + \nu',$$

we should either have

$$\begin{aligned}\|\nabla L^\Pi(w', v')\| &\geq \frac{\nu/4}{4\sqrt{\|w'\|^2 + \|v'\|^2}} \\ &= \frac{\nu/4}{4\sqrt{\|W'\|^2 + \|V'\|^2}} \\ &= \frac{\nu}{16\sqrt{C_1^2 + C_2^2}},\end{aligned}$$

or

$$\begin{aligned}\lambda_{\min}(\nabla^2 L^\Pi(w', v')) &\leq -\frac{\nu/4}{2\min_\Sigma(\|w^*\|^2 + \|v^*\|^2)} \\ &= -\frac{\nu}{2\min_\Sigma(\|W_\Sigma^*\|^2 + \|V_\Sigma^*\|^2)} \\ &\leq -\frac{\nu}{16(\zeta_1 + \zeta_2)} \\ &= -\frac{\nu}{16(1 + \zeta)}.\end{aligned}$$

Next, we want to apply Theorem [7](#) by setting

$$\begin{aligned}\gamma &= \frac{\nu}{16(1 + \zeta)}, \\ \aleph_\ell &= R_n(f^*) + \nu',\end{aligned}$$

and lipschitz parameters  $\rho_1, \rho_2, \rho_3 = \text{poly}(B, C_1, C_2, m_1, m_2, m_3)$  set as described in Appendix [B.1](#), Theorem [9](#). Also, note that as prescribed by Theorem [7](#), we set

$$\begin{aligned}C_1 &:= \frac{\aleph + 4l}{\psi_1}, \\ C_2 &:= \frac{\aleph + 4l}{\psi_2},\end{aligned}\tag{74}$$

where  $l = O(1)$  depends on our desired chance of success for the algorithm, specified in Theorem [7](#). Finally, note that Theorem [7](#) needs to work with a bounded noise on the gradient whose covariance matrix is bounded between two multipliers of identity. The point of injecting extra noise to SGD in  $\text{PSGD}$  is in fact because of this covariance condition that we need in Theorem [7](#). On the other hand, note that in general, because of the gaussian smoothing that we use, the noise vector is not supported on a bounded domain, which makes it a bit harder to apply Hoeffding type concentration. To remedy this, we introduce a coupling between our unbounded noise vector for  $L(W', V')$  and another noise random variable whose support is bounded, which with high probability is equal to the real noise, along all iterations. In Corollary, we further translate this coupling for the objective  $L^\Pi$  after change of coordinates, and write down the exact dependencies of the parameters  $Q, \sigma_1$  and  $\sigma_2$ , which are all polynomial in the basic parameters and the overparameterization.

Hence, the conditions of Theorem [7](#) are satisfied, so we conclude that after at most  $\text{poly}(\rho_1, \rho_2, \rho_3, Q, \aleph, C_1, C_2, 1/\gamma, \log(\sigma_1/\sigma_2)) = \text{poly}(B, m_1, m_2, m_3, C_1, C_2, \zeta_1, \zeta_2) = \text{poly}(n, B \vee 1/B, 1/\gamma_0)$  number of iterations,  $\text{PSGD}$  reach a point  $w_t$  in some iteration  $t$  with  $L^\Pi(w_t) \leq \aleph_\ell$ .

Translating back this  $w_t = (w'_t, v'_t)$  by multiplying the  $v'_t$  part to  $\Upsilon^T$ , we get a pair  $(W'_t, V'_t)$  with objective value bounded as

$$L(W'_t, V'_t) \leq R_n(f^*) + \nu'.$$

But note that we obviously have the condition  $\|W'\| \leq C_1, \|V'\| \leq C_2$  through the whole iterations, for the choice of  $C_1, C_2$  in Equation [\(74\)](#). Therefore, using Lemma [54](#), for every  $i \in [n]$ :

$$|f'_{W', V'}(x_i)| = O(C_1, C_2),\tag{75}$$

$$|f_{W', V'}(x_i)| = O(C_1, C_2)\tag{76}$$

From Equations (76), as also stated in Theorem 9, we know that for all  $i \in [n]$ ,  $\ell(\cdot, y_i)$  is  $O(C_1 C_2) + B^2$ -Lipschitz at points  $f_{W', V'}(x_i)$  and  $f'_{W', V'}(x_i)$ , so we can bound the difference  $|\ell(f'_{W', V'}(x_i), y_i) - \ell(f_{W', V'}(x_i), y_i)|$  by  $(O(C_1 C_2) + B^2)|f'_{W', V'}(x_i) - f_{W', V'}(x_i)|$ , which in turn can become arbitrarily small having enough overparameterization using Lemma 35, in particular, we force it to be smaller than  $O(\nu'/(B^2 + C_1 C_2))$  (recall  $\nu' \geq \nu/2 \geq B^2/(2n)$ ). As a result, we get  $|\ell(f'_{W', V'}(x_i), y_i) - \ell(f_{W', V'}(x_i), y_i)| = O(\nu)$  for every  $i \in [n]$ , which in turn implies  $|L(W', V') - L_1(W', V')| \leq \nu$  by picking small constants, where recall that the objective  $L_1$  is the same as  $L$  but without the smoothing. Now applying this bound to Equation (75), we get

$$L_1(W'_t, V'_t) \leq R_n(f^*) + 2\nu'.$$

Therefore, as PSGD check the values of  $L_1$  in the loop, it terminates at such pair  $(W_t, V_t)$ . From this point onward, we refer to the returned  $(W'_t, V'_t)$  as just  $(W', V')$ .

Opening the definition of  $L_1(W', V')$ , we clearly get

$$R_n(f_{W', V'}) \leq L_1(W', V') \leq R_n(f^*) + 2\nu' \leq R_n(f^*) + 2\nu. \quad (77)$$

Furthermore, noting the setting of  $\psi_1, \psi_2$  in Theorem 3 and the fact that  $\nu' \geq \nu/2 \geq R_n(f^*)/8$ , we get

$$\|W'\|^2 \leq \frac{4(R_n(f^*) + 2\nu')}{\nu'} \leq 40, \quad (78)$$

$$\|V'\|^2 \leq \frac{4\zeta(R_n(f^*) + 2\nu')}{\nu'} \leq 40\zeta, \quad (79)$$

which completes the proof. The fact that for every  $i \in [n], j \in [m_1], j \notin P$  we have that  $\text{sign}((W_j^{(0)} + W'_j)^T x_i)$  and  $\text{sign}(W_j^{(0)T} x_i)$  are the same follows from Lemma 11.

## A.11 RADEMACHER COMPLEXITY

In this section we show the proof for our Rademacher Complexity bound, which is used in Theorem 5.

**Theorem 5** *Let  $\mathcal{G}_{\gamma_1, \gamma_2}$  be the class of neural nets with weights  $(W, V)$  in our three layer setting, such that  $\|W - W^{(0)}\| \leq \gamma_1$ ,  $\|V - V^{(0)}\| \leq \gamma_2$ , where for every  $j \in [m_2], i \in [n]$ :  $V_j^\perp \perp \phi^{(0)}(x_i)$ , and for every  $i \in [n], j \in [m_1], j \notin P$  for  $P \subseteq [m_1]$  defined in Lemma 4, it satisfies  $\text{sign}((W_j^{(0)} + W_j')x_i) = \text{sign}(W^{(0)}x_i)$ . Then, for large enough overparameterization, we have the following bound on the Rademacher complexity:*

$$\mathcal{R}(\mathcal{G}_{\gamma_1, \gamma_2}) \leq \frac{2\gamma_1\gamma_2}{\sqrt{n}}.$$

**Proof of Theorem 5**

Here, we do not have the smoothing matrices  $W^\rho, V^\rho$  anymore. In this section, unlike the optimization section that we used  $\{x'_i\}_{i=1}^n$  to denote the output of the first layer by incorporating also the smoothing matrices, here we define it without them:

$$x'_i = \frac{1}{\sqrt{m_1}} W^s \sigma((W^{(0)} + W')x_i).$$

Now define the matrices

$$\begin{aligned} Z'_i &= 1/\sqrt{m_2} \left( a_j \mathbb{1}\{V_j^{(0)}x'_i \geq 0\} x'_i \right)_{j=1}^{m_2}, \\ Z_i^+ &= 1/\sqrt{m_2} \left( a_j (\mathbb{1}\{V_j x'_i \geq 0\} - \mathbb{1}\{V_j^{(0)}x'_i \geq 0\}) x'_i \right)_{j=1}^{m_2}. \end{aligned}$$

To bound the  $Z_i^+$  part, note that substituting  $C_1$  by  $\gamma_1$  in lemma 6 and assuming conditions

$$\begin{aligned} m_1 &= \tilde{\Omega}(m_3^4), \\ \frac{2C_1^{3/2}}{\sqrt{\kappa_1}} \left( \frac{n^3 m_3^3}{m_1 \lambda_0} \right)^{1/4} &\leq \gamma_1, \end{aligned}$$

(we can use this result because we do not have the smoothing matrix  $W^\rho$  here), we get with high probability over the initialization for every  $i \in [n]$ :

$$\|\phi'(x_i)\| = \|x'_i - \phi^{(0)}(x_i)\| \lesssim \gamma_1. \quad (80)$$

Therefore, we can write

$$\begin{aligned} |\text{trace}(V Z_i^+)| &= \frac{1}{\sqrt{m_2}} \left| \sum_j a_j \mathbb{1}\{\text{sign}(V_j x'_i) \neq \text{sign}(V_j^{(0)} x'_i)\} V_j x'_i \right| \\ &\leq \frac{1}{\sqrt{m_2}} \sum_j \mathbb{1}\{\text{sign}(V_j x'_i) \neq \text{sign}(V_j^{(0)} x'_i)\} |V_j x'_i| \\ &\leq \frac{1}{\sqrt{m_2}} \sum_j \mathbb{1}\{|V_j^{(0)} x'_i| \leq |(V_j - V_j^{(0)}) x'_i|\} \left( |(V_j - V_j^{(0)}) x'_i| + |V_j^{(0)} x'_i| \right) \\ &\leq \frac{1}{\sqrt{m_2}} \sum_j \mathbb{1}\{|V_j^{(0)} x'_i| \leq |(V_j - V_j^{(0)}) x'_i|\} \left( 2|(V_j - V_j^{(0)}) x'_i| \right). \\ &\leq \frac{1}{\sqrt{m_2}} \sum_j \mathbb{1}\{|V_j^{(0)} x'_i|, |(V_j - V_j^{(0)}) x'_i| \leq \gamma_2^{2/3} \left( \frac{\kappa_2}{m_2} \right)^{1/3} \min(\gamma_1, \|x'_i\|)\} \left( 2|(V_j - V_j^{(0)}) x'_i| \right) \\ &+ \frac{1}{\sqrt{m_2}} \sum_j \mathbb{1}\{|(V_j - V_j^{(0)}) x'_i| \geq \gamma_2^{2/3} \left( \frac{\kappa_2}{m_2} \right)^{1/3} \min(\gamma_1, \|x'_i\|)\} \left( 2|(V_j - V_j^{(0)}) x'_i| \right). \end{aligned}$$

Now using the fact that  $V_j - V_j^{(0)}$  is orthogonal to  $\phi^{(0)}(x_i)$ 's:

$$\begin{aligned} LHS &\leq \frac{2\gamma_2^{2/3}(\frac{\kappa_2}{m_2})^{1/3}\gamma_1}{\sqrt{m_2}} \sum_j \mathbb{1}\{|V_j^{(0)}x'_i| \leq \gamma_2^{2/3}(\frac{\kappa_2}{m_2})^{1/3}\|x'_i\|\} \\ &\quad + \frac{2}{\sqrt{m_2}} \sum_j \mathbb{1}\{\|V_j - V_j^{(0)}\| \|x'_i - \phi^{(0)}(x_i)\| \geq \gamma_2^{2/3}(\frac{\kappa_2}{m_2})^{1/3}\gamma_1\} \|V_j - V_j^{(0)}\| \|x'_i - \phi^{(0)}(x_i)\|. \end{aligned}$$

Next, using the upper bound on  $\|x'_i - \phi^{(0)}(x_i)\|$ :

$$\begin{aligned} LHS &\lesssim \frac{2\gamma_2^{2/3}(\frac{\kappa_2}{m_2})^{1/3}\gamma_1}{\sqrt{m_2}} \sum_j (\mathbb{1}\{|V_j^{(0)}x'_i| \leq \gamma_2^{2/3}(\frac{\kappa_2}{m_2})^{1/3}\|x'_i\|\}) \\ &\quad + \frac{\gamma_1}{\sqrt{m_2}} \sum_j \mathbb{1}\{\|V_j - V_j^{(0)}\| \gtrsim \gamma_2^{2/3}(\frac{\kappa_2}{m_2})^{1/3}\} \|V_j - V_j^{(0)}\| \\ &\lesssim \frac{\gamma_2^{2/3}(\frac{\kappa_2}{m_2})^{1/3}\gamma_1}{\sqrt{m_2}} \sum_j (\mathbb{1}\{|V_j^{(0)}x'_i| \leq \gamma_2^{2/3}(\frac{\kappa_2}{m_2})^{1/3}\|x'_i\|\}) \\ &\quad + \frac{\gamma_1}{\sqrt{m_2}} \sqrt{\sum_j \mathbb{1}\{\|V_j - V_j^{(0)}\|^2 \geq \gamma_2^{4/3}(\frac{\kappa_2}{m_2})^{2/3}\}} \sqrt{\sum_j \|V_j - V_j^{(0)}\|^2} \\ &\leq \frac{\gamma_2^{2/3}(\frac{\kappa_2}{m_2})^{1/3}\gamma_1}{\sqrt{m_2}} \left( \sum_j \mathbb{1}\{|V_j^{(0)}x'_i| \leq \gamma_2^{2/3}(\frac{\kappa_2}{m_2})^{1/3}\|x'_i\|\} \right) + \frac{\gamma_2^2\gamma_1}{\sqrt{m_2}} \times \left(\frac{m_2}{\kappa_2}\right)^{1/3} \frac{1}{\gamma_2^{2/3}}. \end{aligned}$$

Then, applying the first argument of Lemma [29](#), we have with high probability over the randomness of  $V^{(0)}$ :

$$\begin{aligned} &\leq \frac{\gamma_2^{2/3}(\frac{\kappa_2}{m_2})^{1/3}\gamma_1}{\sqrt{m_2}} \left(\frac{m_2}{\kappa_2} \left(\frac{\kappa_2}{m_2}\right)^{1/3} \gamma_2^{2/3}\right) + \frac{\gamma_2^2\gamma_1}{\sqrt{m_2}} \times \left(\frac{m_2}{\kappa_2}\right)^{1/3} \frac{1}{\gamma_2^{2/3}} \\ &\leq \frac{\gamma_2^{4/3}\gamma_1}{(\kappa_2\sqrt{m_2})^{1/3}} + \frac{\gamma_2^{4/3}\gamma_1}{(\kappa_2\sqrt{m_2})^{1/3}} \\ &\lesssim \frac{\gamma_2^{4/3}\gamma_1}{(\kappa_2\sqrt{m_2})^{1/3}}. \end{aligned}$$

Therefore, we can write:

$$\begin{aligned}
\mathcal{R}(\mathcal{G}_{\gamma_1, \gamma_2})|_{(x_i), (y_i)} &= \frac{1}{n} \mathbb{E}_\epsilon \sup_{V, W \in \mathcal{S}} \sum_{i=1}^n \epsilon_i f_{V, W}(x_i) \\
&= \frac{1}{n} \mathbb{E}_\epsilon \sup_{V, W \in \mathcal{S}} \sum_{i=1}^n \epsilon_i a^T \sigma(1/\sqrt{m_2} V W^s \sigma(1/\sqrt{m_1} W x_i)) \\
&= \frac{1}{n} \mathbb{E}_\epsilon \sup_{V, W \in \mathcal{S}} \sum_{i=1}^n \epsilon_i a^T \sigma(1/\sqrt{m_2} V x'_i) \\
&= \frac{1}{n} \mathbb{E}_\epsilon \sup_{W \in \mathcal{S}} \sup_{V \in \mathcal{S}} \sum_{i=1}^n \epsilon_i \text{trace}(V(Z'_i + Z_i'^+)) \\
&\lesssim \frac{1}{n} \mathbb{E}_\epsilon \sup_{W, V \in \mathcal{S}} \sum_{i=1}^n \epsilon_i \text{trace}(V Z'_i) + \frac{\gamma_2^{4/3} \gamma_1}{(\kappa_2 \sqrt{m_2})^{1/3}} \\
&\leq \frac{1}{n} \mathbb{E}_\epsilon \sup_{W \in \mathcal{S}} \sup_{V \in \mathcal{S}} \text{trace}(V (\sum_{i=1}^n \epsilon_i Z'_i)) + \frac{\gamma_2^{4/3} \gamma_1}{(\kappa_2 \sqrt{m_2})^{1/3}} \\
&= \frac{1}{n} \mathbb{E}_\epsilon \sup_{W \in \mathcal{S}} \sup_{V \in \mathcal{S}} \text{trace}((V - V^{(0)}) (\sum_{i=1}^n \epsilon_i Z'_i)) \\
&\quad + \frac{1}{n} \mathbb{E}_\epsilon \sup_{W \in \mathcal{S}} \text{trace}(V^{(0)} (\sum_{i=1}^n \epsilon_i Z'_i)) + \frac{\gamma_2^{4/3} \gamma_1}{(\kappa_2 \sqrt{m_2})^{1/3}}. \tag{81}
\end{aligned}$$

For the first term, for every  $j \in [m_2]$ , define  $H_j$  to be the set of  $i$ 's in  $[n]$  where the  $j$ th column of  $Z'_i$  is non-zero, i.e.

$$H_j = \{i \in [n] : V_j^{(0)} x'_i \geq 0\}.$$

Here, we use the crucial assumption that  $(V - V^{(0)})_j^T \phi^{(0)}(x_i) = 0$ , so we can drop the  $\phi^{(0)}(x_i)$  term when  $x'_i$  is multiplied to  $V - V^{(0)}$ . Using this trick and applying Cauchy Schwarz, we bound the first term as:

$$\begin{aligned}
&\frac{1}{n} \mathbb{E}_\epsilon \sup_{W, V \in \mathcal{S}} \text{trace}((V - V^{(0)}) (\sum_{i=1}^n \epsilon_i Z'_i)) \\
&\leq \frac{1}{n} \mathbb{E}_\epsilon \|V - V^{(0)}\| \sup_{W \in \mathcal{S}} \sqrt{\frac{1}{m_2} \sum_{j=1}^{m_2} \left\| \sum_{i \in H_j} \epsilon_i \phi^{(2)}(x_i) \right\|^2}.
\end{aligned}$$

Further using Jensen's inequality:

$$\begin{aligned}
&\frac{1}{n} \mathbb{E}_\epsilon \sup_{W, V \in \mathcal{S}} \text{trace}((V - V^{(0)}) (\sum_{i=1}^n \epsilon_i Z'_i)) \\
&\leq \frac{\|V - V^{(0)}\|}{n} \sqrt{\mathbb{E}_\epsilon \frac{1}{m_2} \sum_{j=1}^{m_2} \sup_{W \in \mathcal{S}} \left\| \sum_{i \in H_j} \epsilon_i \phi^{(2)}(x_i) \right\|^2}. \tag{82}
\end{aligned}$$

Using Equation (A11) of Lemma 6 (note that we do not have the smoothing matrix  $W^\rho$  here, so we are allowed to use this result), we obtain

$$\| \langle W - W^{(0)}, Z_i^k \rangle - \phi'(x_i) \| \lesssim \frac{2C_1^{3/2}}{\sqrt{\kappa_1}} (n^3 m_3^3 / m_1 \lambda_0)^{1/4},$$

where  $Z_i^k$ 's are defined in Equation (318).

Plugging this back in (82):

$$\begin{aligned}
& \mathbb{E}_\epsilon \sup_{W \in S} \left\| \sum_{i \in H_j} \epsilon_i \phi'(x_i) \right\|^2 \\
& \leq \mathbb{E}_\epsilon \sup_{W \in S} \left( \left\| \sum_{i \in H_j} \epsilon_i \langle W - W^{(0)}, Z_i^k \rangle \right\| + \frac{2C_1^{3/2}}{\sqrt{\kappa_1}} \left( \frac{n^3 m_3^3}{m_1 \lambda_0} \right)^{1/4} \right)^2 \\
& \lesssim \mathbb{E}_\epsilon \sup_{W \in S} \left\| \sum_{i \in H_j} \epsilon_i \langle W - W^{(0)}, Z_i^k \rangle \right\|^2 + \frac{C_1^3}{\kappa_1} \left( \frac{n^3 m_3^3}{m_1 \lambda_0} \right)^{1/2} \\
& = \mathbb{E}_\epsilon \sup_{W \in S} \sum_{k=1}^{m_3} \left( \text{trace}((W - W^{(0)}) \left( \sum_{i \in H_j} \epsilon_i Z_i^k \right)) \right)^2 + \frac{C_1^3}{\kappa_1} \left( \frac{n^3 m_3^3}{m_1 \lambda_0} \right)^{1/2}. \tag{83}
\end{aligned}$$

Now for every fixed dataset, with high probability over the randomness of  $W^s$ , for every  $k_1 \neq k_2$ :

$$\begin{aligned}
& \left| \left\langle \sum_{i \in H_j} \epsilon_i Z_i^{k_1}, \sum_{i \in H_j} \epsilon_i Z_i^{k_2} \right\rangle \right| \leq \sum_{i_1, i_2 \in H_j} \left| \langle Z_{i_1}^{k_1}, Z_{i_2}^{k_2} \rangle \right| \\
& = \frac{1}{m_1} \sum_{i_1, i_2 \in H_j} \left| \sum_{j=1}^{m_1} W_{k_1, j}^s W_{k_2, j}^s \langle x_{i_1}, x_{i_2} \rangle \mathbb{1}\{W_j^{(0)} x_{i_1} \geq 0\} \mathbb{1}\{W_j^{(0)} x_{i_2} \geq 0\} \right|
\end{aligned}$$

But note that because  $\langle x_{i_1}, x_{i_2} \rangle \leq 1$ , the variables  $W_{k_1, j}^s W_{k_2, j}^s \langle x_{i_1}, x_{i_2} \rangle \mathbb{1}\{W_j^{(0)} x_{i_1} \geq 0\} \mathbb{1}\{W_j^{(0)} x_{i_2} \geq 0\}$  are subgaussian with parameter one with respect to the randomness of  $W^s$ . Hence, with high probability over the randomness of  $W^s$ , we get

$$\left| \left\langle \sum_{i \in H_j} \epsilon_i Z_i^{k_1}, \sum_{i \in H_j} \epsilon_i Z_i^{k_2} \right\rangle \right| \lesssim \frac{1}{m_1} \sum_{i_1, i_2 \in H_j} \sqrt{m_1} \leq \frac{n^2}{\sqrt{m_1}}. \tag{84}$$

Therefore, with high probability over the randomness of  $W^{(0)}$  and  $W'$  and the dataset, we get Equation (84). In order to get rid of the high probability argument on the dataset, we use the stronger Equation (319) in Lemma 46 which uniformly bounds  $\langle Z_{k_1}(x), Z_{k_2}(x') \rangle$  by  $\log(m_1)d/m_1$  for any  $x, x'$ , which in turn gives

$$\left| \left\langle \sum_{i \in H_j} \epsilon_i Z_i^{k_1}, \sum_{i \in H_j} \epsilon_i Z_i^{k_2} \right\rangle \right| \leq \sum_{i_1, i_2 \in H_j} \left| \langle Z_{i_1}^{k_1}, Z_{i_2}^{k_2} \rangle \right| \lesssim \frac{n^2 d \log(m_1)}{\sqrt{m_1}},$$

with high probability, independent of the choice of dataset. This bounds is slightly worse compared to (84), but still efficient for our purpose.

Furthermore, a similar bound to Equation (84) can be obtained in a more adversarial situation when we also take maximum against the choice of the dataset.

Note that the entries of  $\sum_{i \in H_j} \epsilon_i Z_i^k$  for  $1 \leq k \leq m_3$  can differ only in a sign. Therefore, their norms are all equal. Now suppose that  $\mathcal{C}_j$  is the random variable of the norm of these variables:

$$\mathcal{C}_j := \left\| \sum_{i \in H_j} \epsilon_i Z_i^k \right\|.$$

Then, by substituting  $r_k = \frac{1}{\mathcal{C}_j} \sum_{i \in H_j} \epsilon_i Z_i^k$  in Lemma 40, we get

$$\sum_{k=1}^{m_3} \left( \text{trace}((W - W^{(0)}) \left( \sum_{i \in H_j} \epsilon_i Z_i^k \right)) \right)^2 \leq \mathcal{C}_j^2 \left( 1 + m_3^2 O\left( \frac{n^2 d \log(m_1)}{\sqrt{m_1} \mathcal{C}_j^2} \right) \right) \|W - W^{(0)}\|_F^2 \tag{85}$$

$$= \left( \mathcal{C}_j^2 + \frac{n^2 m_3^2 d \log(m_1)}{\sqrt{m_1}} \right) \|W - W^{(0)}\|_F^2. \tag{86}$$

Now recall from Equation (80), we have

$$\|\phi'(x_i)\| \leq \gamma_1. \tag{87}$$

Hence, we can apply Corollary [84](#) with  $\phi^{(2)}(x_i)$  and  $C_1$  substituted by  $\phi'(x_i)$  and  $\gamma_1$  respectively, to argue with high probability over the initialization, there exists a set  $\tilde{P}_i$  such that for every  $i \in [n]$  and  $j \notin \tilde{P}_i$ , sign of  $V_j^T x'_i$  is the same as  $V_j^{(0)T} \phi^{(0)}(x_i)$ , and moreover,

$$|\tilde{P}_i| \lesssim \left(\frac{C_1^2}{m_3 \kappa_1^2}\right)^{1/3} m_2.$$

Now let

$$\bar{H}_j = \{i \in [n] : V_j^{(0)T} \phi^{(0)}(x_i) \geq 0\}.$$

Note that for  $j \notin \tilde{P} = \bigcup_i \tilde{P}_i$ , we have  $H_j = \bar{H}_j$ . Now note that the norm of each  $\sum_{i \in H_j} \epsilon_i Z_i^k$  is at most one. for each index  $1 \leq \ell \leq m_1 d$ , as the random variables  $\sum_{i \in \bar{H}_j} \epsilon_i (Z_i^k)_\ell$  are  $\sum_{i \in \bar{H}_j} (Z_i^k)_\ell^2 \leq \sum_{i \in [n]} (Z_i^k)_\ell^2$  subgaussian, we have with probability at least  $1 - \frac{1}{n}$  over the randomness of  $\epsilon_i$ 's, for every  $1 \leq \ell \leq m_1 d$  and every  $1 \leq j \leq m_2$ :

$$\left| \sum_{i \in \bar{H}_j} \epsilon_i (Z_i^k)_\ell \right| \leq \sqrt{\sum_{i \in [n]} (Z_i^k)_\ell^2 \log(m_1 d m_2 n)},$$

which implies for every  $j \in [m_2]$ :

$$\left\| \sum_{i \in \bar{H}_j} \epsilon_i Z_i^k \right\|^2 \leq \sum_{\ell} \sum_{i \in [n]} (Z_i^k)_\ell^2 \log(m_1 d) \leq n \log(m_1 d m_2 n).$$

Name this event  $\mathcal{B}$ , so

$$\mathbb{P}(\mathcal{B}) \leq \frac{1}{n}.$$

Note that although  $H_j$  might depend on the randomness of  $\epsilon_i$ 's,  $\bar{H}_j$  does not, and if  $j \notin \tilde{P}$ , we obtain

$$C_j = \left\| \sum_{i \in H_j} \epsilon_i Z_i^k \right\| \leq \sqrt{n \log(m_1 d m_2 n)}.$$

Moreover, note that we also have the following worse case bound:

$$C_j = \left\| \sum_{i \in H_j} \epsilon_i Z_i^k \right\| \leq \sum_{i \in H_j} \|Z_i^k\| \leq n.$$

Applying the last two inequalities into Equations [\(298\)](#) and [\(86\)](#):

$$\begin{aligned} & \mathbb{E}_\epsilon \frac{1}{m_2} \sum_{j=1}^{m_2} \sup_{W \in S} \left\| \sum_{i \in H_j} \epsilon_i \phi'(x_i) \right\|^2 \\ & \leq \frac{C_1^3}{\kappa_1} \left(\frac{n^3 m_3^3}{m_1 \lambda_0}\right)^{1/2} + \frac{1}{m_2} \sum_{j=1}^{m_2} \mathbb{E}_\epsilon \sup_{W \in S} \sum_{k=1}^{m_3} \left( \text{trace}((W - W^{(0)})(\sum_{i \in H_j} \epsilon_i Z_i^k)) \right)^2 \\ & \leq \frac{C_1^3}{\kappa_1} \left(\frac{n^3 m_3^3}{m_1 \lambda_0}\right)^{1/2} + \frac{1}{m_2} \mathbb{E}_\epsilon \mathbb{1}\{\mathcal{B}\} \sum_{j \in \tilde{P}} \left( C_j^2 + \frac{n^2 m_3^2 d \log(m_1)}{\sqrt{m_1}} \right) \|W - W^{(0)}\|_F^2 \\ & \quad + \frac{1}{m_2} \mathbb{E}_\epsilon \mathbb{1}\{\mathcal{B}\} \sum_{j \notin \tilde{P}} \left( C_j^2 + \frac{n^2 m_3^2 d \log(m_1)}{\sqrt{m_1}} \right) \|W - W^{(0)}\|_F^2 \\ & \quad + \frac{1}{m_2} \mathbb{E}_\epsilon \mathbb{1}\{\mathcal{B}^c\} \sum_{j=1}^{m_2} \left( C_j^2 + \frac{n^2 m_3^2 d \log(m_1)}{\sqrt{m_1}} \right) \|W - W^{(0)}\|_F^2 \\ & \leq \frac{C_1^3}{\kappa_1} \left(\frac{n^3 m_3^3}{m_1 \lambda_0}\right)^{1/2} + \|W - W^{(0)}\|_F^2 \left[ \frac{|\tilde{P}|}{m_2} \left( n^2 + \frac{n^2 m_3^2 d \log(m_1)}{\sqrt{m_1}} \right) + 2 \left( n + \frac{n^2 m_3^2 d \log(m_1)}{\sqrt{m_1}} \right) \right] \\ & \leq \frac{C_1^3}{\kappa_1} \left(\frac{n^3 m_3^3}{m_1 \lambda_0}\right)^{1/2} + \gamma_1^2 \left[ n^3 \left(\frac{C_1^2}{(m_3 \kappa_1^2)}\right)^{1/3} + \left(\frac{C_1^2}{(m_3 \kappa_1^2)}\right)^{1/3} \frac{n^3 m_3^2 d \log(m_1)}{\sqrt{m_1}} + 2 \frac{n^2 m_3^2 d \log(m_1)}{\sqrt{m_1}} + 2n \right]. \end{aligned} \tag{88}$$

Next, we analyze the term  $\frac{1}{n}\mathbb{E}_\epsilon \sup_{W \in \mathcal{S}} \text{trace}(V^{(0)}(\sum_{i=1}^n \epsilon_i Z'_i))$ . Noting that  $\|\phi^{(0)}(x_i)\| \lesssim \kappa_1 \sqrt{m_3}$  with high probability and using Equation (87):

$$\left\| \sum_{i=1}^n \epsilon_i Z'_i \right\|_F \leq \sum_{i=1}^n \|Z'_i\|_F \leq \sum_{i=1}^n \|x'_i\| \leq \sum_i (\|\phi'(x_i)\| + \|\phi^{(0)}(x_i)\|) \lesssim n(\sqrt{m_3}\kappa_1 + \gamma_1). \quad (89)$$

Hence

$$\begin{aligned} \frac{1}{n}\mathbb{E}_\epsilon \sup_{W \in \mathcal{S}} \text{trace}(V^{(0)}(\sum_{i=1}^n \epsilon_i Z'_i)) &= \frac{1}{n}\mathbb{E}_\epsilon \sup_{W \in \mathcal{S}} \sum_{i=1}^n \epsilon_i \text{trace}(V^{(0)} Z'_i) \\ &= \frac{1}{n}\mathbb{E}_\epsilon \sup_{W \in \mathcal{S}} \sum_{i=1}^n \epsilon_i \left( \frac{1}{\sqrt{m_2}} \sum_{j=1}^{m_2} a_j V_j^{(0)} x'_i \mathbb{1}\{V_j^{(0)} x'_i \geq 0\} \right). \\ &= \frac{1}{n}\mathbb{E}_\epsilon \sup_{W \in \mathcal{S}} \sum_{i=1}^n \epsilon_i \frac{1}{\sqrt{m_2}} a^T \sigma(V^{(0)} x'_i) \leq \sup_{W \in \mathcal{S}} \frac{1}{\sqrt{m_2}} a^T \sigma(V^{(0)} x'_i). \end{aligned}$$

But using Lemma 80:

$$LHS \lesssim \kappa_2 \sqrt{m_3} \|x'_i\|.$$

Applying a similar bound as we did in Equation (89) on  $\|x'_i\|$ :

$$\|x'_i\| \leq \|\phi^{(0)}(x_i)\| + \|\phi'(x_i)\| \lesssim \kappa_1 \sqrt{m_3} + \gamma_1.$$

Substituting above, we get

$$\frac{1}{n}\mathbb{E}_\epsilon \sup_{W \in \mathcal{S}} \text{trace}(V^{(0)}(\sum_{i=1}^n \epsilon_i Z'_i)) \leq \kappa_2 \sqrt{m_3} (\kappa_1 \sqrt{m_3} + \gamma_1). \quad (90)$$

Finally, Substituting Equations (88) into (82), then combining it with (90) into (81), we obtain a bound on Rademacher complexity which holds w.h.p over both the randomness of the initialization and the dataset:

$$\begin{aligned} \mathcal{R}(\mathcal{G}_{\gamma_1, \gamma_2})|_{x, y} &\lesssim \sqrt{\frac{C_1^3}{\kappa_1} \left( \frac{n^3 m_3^3}{m_1 \lambda_0} \right)^{1/2}} \quad (91) \\ &+ \frac{\gamma_1 \gamma_2}{n} \sqrt{n^3 \left( \frac{C_1^2}{(m_3 \kappa_1^2)} \right)^{1/3} + \left( \frac{C_1^2}{(m_3 \kappa_1^2)} \right)^{1/3} \frac{n^3 m_3^2 d \log(m_1)}{\sqrt{m_1}} + 2 \frac{n^2 m_3^2 d \log(m_1)}{\sqrt{m_1}} + 2n} \quad (92) \\ &+ \kappa_2 \sqrt{m_3} (\kappa_1 \sqrt{m_3} + \gamma_1) + \frac{\gamma_2^{4/3} \gamma_1}{(\kappa_2 \sqrt{m_2})^{1/3}}. \end{aligned}$$

Having enough overparameterization, we have for every dataset  $(x, y)$  (i.e. worst-case Rademacher complexity):

$$\mathcal{R}(\mathcal{G}_{\gamma_1, \gamma_2})|_{x, y} \leq 2\gamma_1 \gamma_2 / \sqrt{n}. \quad (93)$$

Note that for the bound (93) to hold, the overparameterization should be picked poly large in  $\gamma_1, \gamma_2$ , as well as in other basic parameters. However, noting Equations (59) and (66) in the proof of Theorem 3, we set  $\gamma_1 = 1, \gamma_2 \geq \Omega(B, n, 1/\gamma_0)$  in Theorem 3, so  $\gamma_1 \gamma_2$  is at most poly in the basic parameters. Therefore, again the overparameterization can be picked polynomially large in the basic only parameters (i.e. **independent** of  $\gamma_1, \gamma_2$  or  $\zeta$ ).

## A.12 CONSTRUCTING $W^*, V^*$

This section consists of two subsections; First, we prove a structural result for the first layer weights ( $W', V'$ ) that the algorithm visits, then construct a weight matrix  $W^*$  for the first layer with some good properties. Second, we do the same thing for the second layer (however, the structure of the first and second layers are completely different). Through out this section, we assume we have the norm bounds  $\|W'\| \leq C_1, \|V'\| \leq C_2$ .

Notably, we rely on a number of basic Lemmas more related to the representation power of the network, which we defer their proof into a later Appendix B.2 and refer to them here as needed.

### A.12.1 FIRST LAYER, CONSTRUCTION OF $W^*$

**Lemma 1** *Suppose  $m_1 \geq 16n^2m_3^2/\lambda_0^2$ . Let  $P_i = \{j \in [m_1] \mid |W_j^{(0)} x_i| \leq c_2/\sqrt{m_1}\}$  and  $P = \cup P_i$ . During SGD iterations, suppose we have  $\|W'\|_F \leq C_1$ . Then, for a value  $c_2$  satisfying*

$$2C_1\sqrt{nm_3}/\sqrt{\lambda_0} \leq c_2 \leq \kappa_1\lambda_0\sqrt{m_1}/(2n^2),$$

with high probability  $\forall i$ :

$$|P_i| \lesssim c_2\sqrt{m_1}/\kappa_1,$$

and for  $j \notin P$ , during the whole algorithm we have

$$\begin{aligned} \|W'_j\| &\leq \frac{\sqrt{nm_3}C_1}{\sqrt{m_1}\lambda_0} + c_2/(4\sqrt{m_1}) \leq c_2/(2\sqrt{m_1}), \\ c_2/\sqrt{m_2} &\leq |W_j^{(0)} x_i|. \end{aligned}$$

So the signs of neurons outside  $P$  never changes. In particular, we can set  $c_2$  as small as  $c_2 = C_1\sqrt{nm_3}/\sqrt{\lambda_0}$ . In the rest of the proof (i.e. other sections), we set  $c_2$  to this value.

#### Proof of Lemma 1

Define the matrix

$$\tilde{Z}_k^i = \frac{1}{\sqrt{m_1}} (W_{k,j}^s x_i \mathbb{1}\{\forall i : W_j^{(0)T} x_i \geq c_2/\sqrt{m_1}\})_{j=1}^n.$$

Let  $P_i$  be the set of indices  $j$  such that  $\mathbb{1}\{W_j^{(0)T} x_i \geq c_2/\sqrt{m_1}\}$  is zero. First of all, note that by Bernstein inequality:

$$|P_i| \leq c_2\sqrt{m_1}/\kappa_1 + O(\sqrt{c_2\sqrt{m_1}/\kappa_1} + 1) \lesssim c_2\sqrt{m_1}/\kappa_1.$$

Now suppose that until the current iteration of the algorithm the assumption has been true, i.e. the signs of the neurons outside of  $P$  have never changed. As a result, due to the specific update of the SGD for both of the terms  $\mathbb{E}_Z \ell(f_{V', W'}(x), y)$  and  $\|W'\|_F^2$ , if we define  $W'|_P$  to be the restriction of  $W'$  to indices that are not in  $P$  (i.e. the columns in  $P$  are equal to zero), then we can write

$$W'|_P^T = \sum_{k=1}^{m_3} \sum_{i=1}^n \alpha_{k,i} \tilde{Z}_k^i. \quad (94)$$

An issue here is that we also have some injected noise by PSGD into  $W'$  which violates Equation (94). To handle the injected noise as well, we define the subspace  $\Phi'$  of  $\mathbb{R}^{m_1 \times d}$  matrices to be the set of vectors with arbitrary rows for  $j \in [m_1]$  with  $j \in P$ , while restricted to the other rows  $j \notin P$  in should be in the span of  $(\tilde{Z}_k^i)_{i,k}$ . Then, we decompose  $W'$  into subspaces  $\Phi'$  and  $\Phi'^\perp$  respectively as  $W' = W'^{(1)} + W'^{(2)}$ , where  $W'^{(1)} \in \Phi', W'^{(2)} \in \Phi'^\perp$ . Here, we want to prove  $\|W_j^{(1)}\| \leq c_2/(4\sqrt{m_1})$ . We handle the  $\|W_j^{(2)}\|$  part in Appendix A.5. So instead of  $W'|_P$  in Equation (94) we consider  $W'^{(1)}|_P$ :

$$W'^{(1)}|_P^T = \sum_{k=1}^{m_3} \sum_{i=1}^n \alpha_{k,i} \tilde{Z}_k^i. \quad (95)$$

We handle the other part  $W'^{(2)}$  in Appendix [45](#). Now exactly similar to the derivation in Lemma [38](#), we can state with high probability

$$\begin{aligned} C_1^2 &\geq \|W'\|^2 \geq \|W'^{(1)}\|^2 \\ &\geq \|W'^{(1)}|_P\|^2 \geq \sum_{k=1}^{m_3} \left\| \sum_{i=1}^n \alpha_{k,i} \tilde{Z}_k^i \right\|^2 - O(nm_3/\sqrt{m_1}) \sum_k \|\alpha_k\|^2. \end{aligned} \quad (96)$$

Note that we are exploiting the fact that the norm  $\|W'\|$  remains bounded by  $C_1$ . Now using a Hoeffding bound for matrix  $H^{\infty'}$  defined below, we write:

$$\begin{aligned} H^{\infty'}_{i_1, i_2} &:= \mathbb{E}_{w: \mathcal{N}(0, \mathbb{R}^d)} \mathbb{1}\{\forall i: |w^T x_i| \geq c_2/\sqrt{m_1}\} x_{i_1}^T x_{i_2} (\mathbb{1}\{w^T x_{i_1} \geq 0\} \mathbb{1}\{w^T x_{i_2} \geq 0\}) \\ &= \mathbb{E}_{w: \mathcal{N}(0, \mathbb{R}^d)} (\mathbb{1}\{w^T x_{i_1} \geq 0\} \mathbb{1}\{w^T x_{i_2} \geq 0\}) x_{i_1}^T x_{i_2} \\ &\pm O(\mathbb{E} \mathbb{1}\{\exists i: |w^T x_i| \leq c_2/\sqrt{m_1}\} (\mathbb{1}\{w^T x_{i_1} \geq 0\} \mathbb{1}\{w^T x_{i_2} \geq 0\})) x_{i_1}^T x_{i_2} \\ &= H^{\infty}_{i_1, i_2} \pm O(nc_2/(\sqrt{m_1} \kappa_1) \|x_{i_1}\| \|x_{i_2}\|) \\ &= H^{\infty}_{i_1, i_2} \pm O(nc_2/(\sqrt{m_1} \kappa_1)). \end{aligned} \quad (97)$$

Now opening Equation [\(96\)](#) and using the property  $c_2 \leq k_1 \lambda_0 \sqrt{m_1}/(2n^2)$ , we get

$$\begin{aligned} LHS &= \sum_k \sum_{i_1, i_2} \alpha_{k, i_1} \alpha_{k, i_2} \langle \tilde{Z}_k^{i_1}, \tilde{Z}_k^{i_2} \rangle - O(nm_3/\sqrt{m_1}) \sum_k \|\alpha_k\|^2 \\ &= \sum_k \sum_{i_1, i_2} \alpha_{k, i_1} \alpha_{k, i_2} (H^{\infty'}_{i_1, i_2} \pm O(1/\sqrt{m_1})) - O(nm_3/\sqrt{m_1}) \sum_k \|\alpha_k\|^2 \\ &\geq \sum_k \sum_{i_1, i_2} \alpha_{k, i_1} \alpha_{k, i_2} H^{\infty}_{i_1, i_2} \pm \|\alpha_k\|_1^2 O(nc_2/\sqrt{m_1} \kappa_1) - O(nm_3/\sqrt{m_1}) \sum_k \|\alpha_k\|^2 \\ &\geq \sum_k \alpha_k^T H^{\infty} \alpha_k - O(nc_2/\sqrt{m_1} \kappa_1) \sum_k \|\alpha_k\|_1^2 - O(nm_3/\sqrt{m_1}) \sum_k \|\alpha_k\|^2 \\ &\geq \sum_k \alpha_k^T H^{\infty} \alpha_k - O(c_2 n^2/\sqrt{m_1} \kappa_1) \sum_k \|\alpha_k\|_2^2 - O(nm_3/\sqrt{m_1}) \sum_k \|\alpha_k\|^2 \\ &= (\lambda_0 - O(nm_3/\sqrt{m_1}) - O(c_2 n^2/\sqrt{m_1} \kappa_1)) \sum_k \|\alpha_k\|^2 \\ &\geq \lambda_0/2 \sum_k \|\alpha_k\|^2. \end{aligned}$$

For the last line to hold, we need enough overparameterization. This implies

$$\sum_k \|\alpha_k\|^2 \lesssim C_1^2/\lambda_0.$$

Now again, exactly similar to the derivation in Lemma [38](#), for  $j \notin P$  we have

$$\|W'_j{}^{(1)}\| \leq \sqrt{nm_3}/\sqrt{m_1} \sqrt{\sum_k \|\alpha_k\|^2} \lesssim \sqrt{m_3 n} C_1 / \sqrt{m_1} \lambda_0,$$

which completes most of the proof. For the rest, we are left to show that for the other part  $W'^{(2)}$ , we have  $\|W'_j{}^{(2)}\| \leq c_2/(4\sqrt{m_1})$ , which we do in Appendix [45](#).

**Lemma 2** Under condition  $m_3 n/\sqrt{m_1} \leq \lambda_0/4$ , there exist matrices  $\{W_k^*\}_{k=1}^{m_3} \in \mathbb{R}^{m_1 \times d}$  s.t. for every  $k \neq k' \in [m_3]$  and  $i \in [n]$ :

$$\begin{aligned} \langle W_k^*, Z_{k'}^i \rangle &= 0, \\ \|W_k^* - W_k^+\| &\lesssim \frac{n\sqrt{m_3}}{\lambda_0 \sqrt{m_1}} \|\mathcal{V}_k\|_{H^\infty}, \\ |\langle W_k^*, Z_k^i \rangle - \langle W_k^+, Z_k^i \rangle| &\lesssim \frac{n\sqrt{m_3}}{\lambda_0 \sqrt{m_1}} \|\mathcal{V}_k\|_{H^\infty}. \end{aligned}$$

Furthermore, for  $k_1 \neq k_2$ :

$$|\langle W_{k_1}^*, W_{k_2}^* \rangle| \leq \frac{n}{\lambda_0^2} \frac{\sqrt{m_3}}{\sqrt{m_1}} \left(1 + \frac{\sqrt{m_3}}{\sqrt{m_1}}\right) \|\mathcal{V}_{k_1}\|_{H^\infty} \|\mathcal{V}_{k_2}\|_{H^\infty}. \quad (98)$$

### Proof of Lemma 2

Let

$$W_k^+ = \sum_i \mathcal{V}_{k,i} Z_k^i.$$

we want to compute the norm of the projection  $P(W_k^+)$  of  $W_k^+$  onto the subspace spanned by all  $Z_{k'}^i$  for  $k' \neq k$  and  $i \in [n]$ :

$$\|P(W_k^+)\|^2 = (\langle W_k^+, Z_{k'}^i \rangle)_{k' \neq k, i \in [n]}^T \left( \langle Z_{k_1}^{i_1}, Z_{k_2}^{i_2} \rangle \right)_{(k_1, i_1), (k_2, i_2) \in [m_3] - \{k\} \times [n]}^{-1} (\langle W_k^+, Z_{k'}^i \rangle)_{k' \neq k, i \in [n]}, \quad (99)$$

where the first and third terms are vectors and the middle term is a matrix. Now note that for each  $k', k_1, k_2 \neq k$ , by Hoeffding inequality:

$$\left( \langle Z_{k'}^{i_1}, Z_{k'}^{i_2} \rangle \right)_{i_1, i_2 \in [n]} = H^\infty + (\pm 1 / \sqrt{m_1})_{i_1, i_2 \in [n]}, \quad (100)$$

$$\begin{aligned} \langle W_k^+, Z_{k'}^i \rangle &= \left\langle \sum_i \mathcal{V}_{k,i} Z_k^i, Z_{k'}^i \right\rangle \lesssim \frac{1}{\sqrt{m_1}} \sum_i |\mathcal{V}_{k,i}| \\ &\leq \frac{\sqrt{n}}{\sqrt{m_1}} \|\mathcal{V}_k\| \\ &\leq \frac{\sqrt{n}}{\sqrt{m_1} \lambda_0} \|\mathcal{V}_k\|_{H^\infty}. \end{aligned} \quad (101)$$

Therefore,

$$\|(\langle W_k^+, Z_{k'}^i \rangle)_{k' \neq k, i \in [n]}^T\| \leq n \sqrt{\frac{m_3}{m_1 \lambda_0}} \|\mathcal{V}_k\|_{H^\infty}. \quad (102)$$

Now Equation (100) implies for small enough  $m_1$

$$\lambda_{\min} \left( \left( \langle Z_{k'}^{i_1}, Z_{k'}^{i_2} \rangle \right)_{i_1, i_2 \in [n]} \right) \geq \lambda_0 / 2, \quad (103)$$

as long as  $\lambda_0 \geq 2n/m_1$ . Moreover, define  $\tilde{A}$  to be the block version of

$$A' = \left( \langle Z_{k_1}^{i_1}, Z_{k_2}^{i_2} \rangle \right)_{(k_1, i_1), (k_2, i_2) \in [m_3] - \{k\} \times [n]}^{-1},$$

i.e. for  $k_1 = k_2$  they are the same but for  $k_1 \neq k_2$   $\tilde{A}$  is zero. Then

$$\lambda_{\min}(\tilde{A}) \geq \lambda_0 / 2,$$

because the eigenvalues of each block is at least  $\lambda_0 / 2$  using Equation (103). But note that

$$\|A' - \tilde{A}\|_2 \leq \|A' - \tilde{A}\|_F \leq m_3 n / \sqrt{m_1}.$$

So as long as  $m_3 n / \sqrt{m_1} \leq \lambda_0 / 4$ , we have  $\lambda_{\min}(A) \geq \lambda_0 / 4$ . Combining this fact with Equation (102) and plugging it into Equation (99), we obtain

$$\|P(W_k^+)\|^2 \lesssim \frac{n^2 m_3}{m_1 \lambda_0^2} \|\mathcal{V}_k\|_{H^\infty}^2.$$

Now define  $W_k^* = W_k^+ - P(W_k^+)$ . Then

$$\begin{aligned}\|W_k^* - W_k^+\| &= \|P(W_k^+)\| \lesssim \frac{n\sqrt{m_3}}{\lambda_0\sqrt{m_1}} \|\mathcal{V}_k\|_{H^\infty}, \\ |\langle W_k^* - W_k^+, Z_k^i \rangle| &\leq \|P(W_k^+)\| \|Z_k^i\| \lesssim \frac{n\sqrt{m_3}}{\lambda_0\sqrt{m_1}} \|\mathcal{V}_k\|_{H^\infty}.\end{aligned}$$

Furthermore, note that  $W_{k_2}^*$  is orthogonal to  $W_{k_1}^+$  for  $k_1 \neq k_2$ , so

$$\begin{aligned}|\langle W_{k_1}^*, W_{k_2}^* \rangle| &= |\langle W_{k_1}^+ - P(W_{k_1}^+), W_{k_2}^* \rangle| \\ &= |\langle P(W_{k_1}^+), W_{k_2}^+ - P(W_{k_2}^+) \rangle| \\ &\leq \|P(W_{k_1}^+)\| \|W_{k_2}^+ - P(W_{k_2}^+)\| \\ &\leq \|P(W_{k_1}^+)\| (\|W_{k_2}^+\| + \|P(W_{k_2}^+)\|).\end{aligned}\tag{104}$$

But note that

$$\|W_k^+\| = \left\| \sum_i \mathcal{V}_{k,i} Z_k^i \right\| \leq \sum_i |\mathcal{V}_{k,i}| \|Z_k^i\| \leq \sum_i |\mathcal{V}_{k,i}| \leq \sqrt{n} \|\mathcal{V}_k\|_2 \leq \frac{\sqrt{n}}{\lambda_0} \|\mathcal{V}_k\|_{H^\infty}.$$

Therefore, we can bound Equation (104) as:

$$|\langle W_{k_1}^*, W_{k_2}^* \rangle| \leq \frac{n}{\lambda_0^2} \frac{\sqrt{m_3}}{\sqrt{m_1}} \left(1 + \frac{\sqrt{m_3}}{\sqrt{m_1}}\right) \|\mathcal{V}_{k_1}\|_{H^\infty} \|\mathcal{V}_{k_2}\|_{H^\infty}.$$

**Lemma 3** *There exists a matrix  $W_k^{+2}$  such that for every  $j \in P$ ,  $W_k^{+2} j = 0$ , and*

$$|\text{trace}(W_k^{+2} Z_i^k) - \bar{x}_{i,k}| \leq C_1 \sqrt{m_3} n^2 / (\lambda_0 \kappa_1 \sqrt{m_1}) \|\mathcal{V}_k\|_{H^\infty}.$$

**Proof of Lemma 3**

Define  $W_k^{+2}$  to be equal to  $W_k^+$  for  $j \notin P$  and equal to zero vector otherwise. Then, by Lemma 38: (note that  $|P_i| \leq C_1 \sqrt{n m_3} \sqrt{m_1} / (\sqrt{\lambda_0} \kappa_1)$ )

$$\begin{aligned}|\text{trace}(W_k^+ Z_i^k) - \text{trace}(W_k^{+2} Z_i^k)| &\leq 1/\sqrt{m_1} \sum_{j \in P} |W_k^+ j x_i| \\ &\leq \frac{|P|}{\sqrt{m_1}} \|W_k^+\| \\ &\leq \sqrt{n m_3} / (m_1 \sqrt{\lambda_0}) |P| \|\mathcal{V}_k\|_{H^\infty} \\ &\leq C_1 m_3 n^2 / ((\lambda_0 \kappa_1 \sqrt{m_1}) \|\mathcal{V}_k\|_{H^\infty}).\end{aligned}$$

Combining this with Lemma 37, the desired result follows.

**Lemma 4** *Under condition  $m_3 n / \sqrt{m_1} \leq \lambda_0 / 4$ , there exist matrix  $W_k^*$ 's exactly satisfying the same conditions in Lemma 2 but with respect to  $W_k^{+2}$  instead of  $W_k^+$ , and moreover, for  $j \in P$  we have  $W_k^* j = 0$ .*

**Proof of Lemma 4**

We can repeat the exact same procedure of Lemma 2 for  $W_k^{+2}$ . Using the bound in Equation (37), we have

$$\begin{aligned}\left(\langle \tilde{Z}_{k'}^{i_1}, \tilde{Z}_{k'}^{i_2} \rangle\right)_{i_1, i_2 \in [n]} &= H'^\infty + O(\pm 1 / \sqrt{m_1})_{i_1, i_2 \in [n]} \\ &= H^\infty + (\pm n c_2 / \sqrt{m_1} \kappa_1)_{i_1, i_2 \in [n]} + O(\pm 1 / \sqrt{m_1})_{i_1, i_2 \in [n]} \\ &= H^\infty + (\pm n c_2 / \sqrt{m_1} \kappa_1)_{i_1, i_2 \in [n]},\end{aligned}$$

so as long as

$$n^2 c_2 / \sqrt{m_1} \kappa_1 = n^2 C_1 \sqrt{nm_3} / (\kappa_1 \sqrt{m_1 \lambda_0}) \leq \lambda_0 / 2,$$

with similar argument as in Lemma [4](#), we get

$$\lambda_{\min}(\langle \tilde{Z}_{k'}^{i_1}, \tilde{Z}_{k'}^{i_2} \rangle_{i_1, i_2 \in [n]}) \geq \lambda_0 / 2.$$

Moreover,

$$\begin{aligned} \langle W_k^{+2}, \tilde{Z}_{k'}^i \rangle &= \langle \sum_i \mathcal{V}_{k,i} \tilde{Z}_k^i, \tilde{Z}_{k'}^i \rangle \lesssim \frac{1}{\sqrt{m_1}} \sum_i |\mathcal{V}_{k,i}| \\ &\leq \frac{\sqrt{n}}{\sqrt{m_1}} \|\mathcal{V}_k\| \\ &\leq \frac{\sqrt{n}}{\sqrt{m_1 \lambda_0}} \|\mathcal{V}_k\|_{H^\infty}. \end{aligned}$$

Thus, using the same argument as before the proof is complete.

**Lemma 5** *Suppose*

$$m_1 \geq n^7 m_3 / \lambda_0.$$

*During SGD, suppose we are currently at  $(V', W')$  with  $W' \leq C_1$ . For any matrix  $W_1$ , we denote the signs of the first layer imposed by  $W_1$  by  $D_{W_1, x_i}$ . Then with high probability, there exists  $W^* = \sum_{k \in [m_3]} W_k^*$  such that  $W_k^*$ 's is orthogonal to all other  $Z_{k'}^i$ 's for  $k' \neq k$ , and for every  $i \in [n]$ , we have:*

$$\left\| \frac{1}{\sqrt{m_1}} W^s D_{W^{(0)}+W', x_i} W^* x_i - \bar{x}_i \right\|_\infty \lesssim \frac{nm_3}{\sqrt{m_1} \lambda_0} \left[ 1 + \frac{nC_1}{\kappa_1} \right] \|\mathcal{V}_k\|_{H^\infty} := \mathfrak{R} \|\mathcal{V}_k\|_{H^\infty}.$$

*Moreover, we have*

$$\|W_j^*\| \leq \sqrt{nm_3} / (\sqrt{m_1} \lambda_0) \sqrt{\sum_k \|\mathcal{V}_k\|_{H^\infty}^2} + \frac{n\sqrt{m_3}}{\lambda_0 \sqrt{m_1}} \left( \sum_k \|\mathcal{V}_k\|_{H^\infty} \right) := \varrho \sqrt{\frac{m_3}{m_1}}, \quad (105)$$

*Particularly, for any diagonal sign matrix  $\Sigma \in \mathbb{R}^{m_3 \times m_3}$ , we have*

$$\|W_\Sigma^*\|_F^2 \leq \left( \frac{n\sqrt{n}}{\lambda_0^2} \frac{\sqrt{m_3}}{\sqrt{m_1}} \left( 1 + \frac{\sqrt{m_3}}{\sqrt{m_1}} \right) + \left( 1 + O(n/(\lambda_0 \sqrt{m_1})) + \frac{n^2 m_3}{\lambda_0^2 m_1} \right) \right) \sum_k \|\mathcal{V}_k\|_{H^\infty}^2. \quad (106)$$

*which, by having enough overparameterization, implies*

$$\|W_\Sigma^*\|_F \leq \sqrt{2 \sum_k \|\mathcal{V}_k\|_{H^\infty}^2} = \sqrt{2\zeta_1}, \quad (107)$$

*where*

$$W_\Sigma^* := \sum_{k=1}^{m_3} \Sigma_k W_k^*. \quad (108)$$

*Moreover, we have*

$$\frac{1}{\sqrt{m_1}} W^s D_{W^{(0)}+W', x_i} W_\Sigma^* x_i = \Sigma \frac{1}{\sqrt{m_1}} W^s D_{W^{(0)}+W', x_i} W^* x_i. \quad (109)$$

**Proof of Lemma [5](#)**

From Lemma [4](#), we have

$$|\bar{x}_{i,k} - \text{trace}(W_k^{+2} Z_k^i)| \leq C_1 m_3 n^2 / (\lambda_0 \kappa_1 \sqrt{m_1}) \|\mathcal{V}_k\|_{H^\infty}.$$

Combining this with the result of Lemma [4](#), we get:

$$\begin{aligned} |\bar{x}_{i,k} - \text{trace}(W_k^* Z_k^i)| &\lesssim \left[ \frac{n\sqrt{m_3}}{\lambda_0\sqrt{m_1}} + \frac{C_1 m_3 n^2}{\lambda_0 \kappa_1 \sqrt{m_1}} \right] \|\mathcal{V}_k\|_{H^\infty} \\ &= \frac{nm_3}{\sqrt{m_1}\lambda_0} \left[ 1 + \frac{nC_1}{\kappa_1} \right] \|\mathcal{V}_k\|_{H^\infty}. \end{aligned} \quad (110)$$

On the other hand, based on the property that  $W_{k_j}^* = 0$  for  $j \in P$  and its orthogonal property from Lemma [4](#), for  $j \in P$  we get

$$\begin{aligned} \frac{1}{\sqrt{m_1}} W_k^s D_{W^{(0)}+W',x_i} W^* x_i &= \frac{1}{\sqrt{m_1}} W_k^s D_{W^{(0)},x_i} W^* x_i \\ &= \text{trace}(W^* Z_k^i) = \text{trace}(W_k^* Z_k^i) \\ &= \frac{1}{\sqrt{m_1}} W_k^s D_{W^{(0)},x_i} W_k^* x_i, \end{aligned}$$

which combined with Equation ([110](#)) completes the proof. From the above, Equation ([109](#)) is also clear. Finally, note that by Lemma [58](#) we have

$$\|W^{+2}_j\| \leq \sqrt{nm_3}/(\sqrt{m_1}\lambda_0) \sqrt{\sum_k \|\mathcal{V}_k\|_{H^\infty}^2}.$$

which Combined with Lemma [4](#) implies

$$\|W_j^*\| \leq \sqrt{nm_3}/(\sqrt{m_1}\lambda_0) \sqrt{\sum_k \|\mathcal{V}_k\|_{H^\infty}^2} + \frac{n\sqrt{m_3}}{\lambda_0\sqrt{m_1}} (\sum_k \|\mathcal{V}_k\|_{H^\infty}) := \varrho \sqrt{\frac{m_3}{m_1}},$$

while the other claims follows from Lemma [59](#) and Lemma [4](#), combined with Equation ([98](#)):

$$\begin{aligned} \|W_\Sigma^*\|_F^2 &\leq \sum_k \|W_k^*\|^2 + \sum_{k_1 \neq k_2} |\langle W_{k_1}^*, W_{k_2}^* \rangle| \\ &\leq \sum_k \|W_k^*\|^2 + \frac{n}{\lambda_0^2} \frac{\sqrt{m_3}}{\sqrt{m_1}} (1 + \frac{\sqrt{m_3}}{\sqrt{m_1}}) (\sum_k \|\mathcal{V}_k\|_{H^\infty})^2 \\ &\leq \sum_k \|W_k^*\|^2 + \frac{n}{\lambda_0^2} \frac{\sqrt{m_3}}{\sqrt{m_1}} (1 + \frac{\sqrt{m_3}}{\sqrt{m_1}}) \sqrt{n} (\sum_k \|\mathcal{V}_k\|_{H^\infty}^2) \\ &\leq \sum_k \|W_k^*\|^2 + \frac{n\sqrt{n}}{\lambda_0^2} \frac{\sqrt{m_3}}{\sqrt{m_1}} (1 + \frac{\sqrt{m_3}}{\sqrt{m_1}}) \zeta_1 \\ &\leq \frac{n\sqrt{n}}{\lambda_0^2} \frac{\sqrt{m_3}}{\sqrt{m_1}} (1 + \frac{\sqrt{m_3}}{\sqrt{m_1}}) \zeta_1 + \sum_k \|W_k^{+2}\|^2 + \frac{n^2 m_3}{\lambda_0^2 m_1} \sum_k \|\mathcal{V}_k\|_{H^\infty}^2 \\ &\leq \frac{n\sqrt{n}}{\lambda_0^2} \frac{\sqrt{m_3}}{\sqrt{m_1}} (1 + \frac{\sqrt{m_3}}{\sqrt{m_1}}) \zeta_1 + \sum_k \|W_k^+\|^2 + \frac{n^2 m_3}{\lambda_0^2 m_1} \sum_k \|\mathcal{V}_k\|_{H^\infty}^2 \\ &\lesssim \frac{n\sqrt{n}}{\lambda_0^2} \frac{\sqrt{m_3}}{\sqrt{m_1}} (1 + \frac{\sqrt{m_3}}{\sqrt{m_1}}) \zeta_1 + (1 + O(n/(\lambda_0\sqrt{m_1}) + \frac{n^2 m_3}{\lambda_0^2 m_1})) \zeta_1. \end{aligned}$$

Next, we move on to construct  $V^*$  for the second layer.

#### A.12.2 SECOND LAYER, CONSTRUCTION OF $V^*$

In this section, we present a couple of lemmas that step by step lead to the construction of  $V^*$ . we remind the reader that  $\phi^{(0)}(x_i)$  is the output of the first layer at initialization weights,  $\phi'(x_i)$  and  $\phi^{(2)}(x_i)$  are the changes in the output of the first layer when  $W'$  and  $W' + W^\rho$  are added,

respectively, and finally  $\phi^*(x_i)$  is the optimal features that are generated by the matrix  $W^*$  but with the sign pattern of  $W^{(0)} + W'$ , i.e.

$$\phi^*(x_i) = \frac{1}{\sqrt{m_1}} W^s D_{W^{(0)}+W', x_i} W^* x_i.$$

We also define  $x'_i$  as

$$x'_i = \phi^{(0)}(x_i) + \phi^{(2)}(x_i) = \frac{1}{\sqrt{m_1}} W^s \sigma((W^{(0)} + W' + W^\rho)x_i).$$

To begin, we state a lemma to bound the magnitude of  $\|\phi'(x_i)\|$ , given that the norm of  $W'$  is bounded by  $C_1$  and the sign pattern  $\text{Sgn}((W^{(0)} + W')x_i)$  satisfies condition stated for the set of indices  $P$  in Lemma III. Later on, we exploit this Lemma in Lemma III to state bounds for  $\|\phi^{(2)}(x_i)\|$ .

**Lemma 6** *Let the matrix  $W'$  with norm bound  $\|W'\| \leq C_1$ , such that the signs of  $(W_j^{(0)} + W'_j)x_i$  and  $W_j^{(0)}x_i$  can be different only for  $j \in P$ , for  $P$  defined in Lemma III. (Note that for  $W'$  at every step of the algorithm, this is automatically satisfied by Lemma III) Then*

$$\|\phi'(x_i)\| \leq \frac{2C_1^{3/2}}{\sqrt{\kappa_1}} \left(\frac{n^3 m_3^3}{m_1 \lambda_0}\right)^{1/4} + (1 + O(m_3^2/\sqrt{m_1}))C_1.$$

Particularly for large enough  $m_1$  compared to  $n, m_3, \lambda_0, \kappa_1, C_1$ , we have

$$\|\phi'(x_i)\| \lesssim C_1.$$

### Proof of Lemma 6

We write

$$\begin{aligned} |\phi'_k(x_i) - \langle W', Z_k^i \rangle| &\leq 2/\sqrt{m_1} \sum_{j \in P} |W'_j x_i| \leq 2/\sqrt{m_1} \sum_{j \in P} \|W'_j\| \\ &\leq 2\sqrt{|P|}/\sqrt{m_1} \|W'\|_F \leq \frac{2C_1^{3/2}}{\sqrt{\kappa_1}} \left(\frac{n^3 m_3^3}{m_1 \lambda_0}\right)^{1/4}, \end{aligned} \quad (111)$$

where the last line follows from the bound on  $|P|$  from Lemma III.

On the other hand, because by Hoeffding we know that  $\langle Z_k^i, Z_{k'}^i \rangle \lesssim 1/\sqrt{m_1}$  by Lemma III, we get

$$\sum_{k=1}^{m_3} \langle W', Z_k^i \rangle^2 \leq (1 + O(m_3^2/\sqrt{m_1})) \|W'\|_F^2 \leq (1 + O(m_3^2/\sqrt{m_1})) C_1^2.$$

Combining this with Equation (III), we get

$$\begin{aligned} \|\phi'(x_i)\| &\leq \sqrt{\sum_k |\phi_k^{(2)}(x_i) - \langle W', Z_k^i \rangle|^2} + \sqrt{\sum_k \langle W', Z_k^i \rangle^2} \\ &\leq \frac{2C_1^{3/2}}{\sqrt{\kappa_1}} \left(\frac{n^3 m_3^3}{m_1 \lambda_0}\right)^{1/4} + (1 + O(m_3^2/\sqrt{m_1})) C_1. \end{aligned} \quad (112)$$

Next, we prove a structural lemma regarding the sign pattern in the second layer when we feed in  $x'_i$  to it, with the important message that the dominance of sign patterns are specified by  $\phi^{(0)}(x_i)$ .

**Lemma 7** *Suppose we have  $m_3 \kappa_1^2 \gtrsim C_1^2$ ,  $\kappa_2 \sqrt{m_2} \geq C_2$ , and  $m_1$  satisfies the condition on Lemma III. If we have the condition  $\|\phi^{(2)}(x_i)\| \lesssim C_1$ , which happens under the high probability event  $E^c$  defined in Lemma III, then for every  $i \in [n]$ , there exist a subset  $\tilde{P}_i$  which might depend on  $W^{(0)}, V^{(0)}, W', V'$ , such that*

$$|\tilde{P}_i| \lesssim \left( \left(\frac{C_1^2}{m_3 \kappa_1^2}\right)^{1/3} + \left(c_3 + \frac{C_1^2}{c_3^2 m_3 \kappa_1^2}\right) \left(\frac{C_2^2}{\kappa_2^2 m_2}\right)^{1/3} \right) m_2.$$

Moreover, for every  $i \in [n]$ , for  $j \notin \tilde{P}_i$ :

$$\begin{aligned} \frac{2}{3}|V_j^{(0)}\phi^{(0)}(x_i)| &\geq |V_j^{(0)}\phi^{(2)}(x_i)| + |V_j'(\phi^{(0)}(x_i) + \phi^{(2)}(x_i))|, \\ |V_j^{(0)}\phi^{(0)}(x_i)| &\gtrsim \left(\frac{\kappa_2}{m_2}\right)^{1/3} C_2^{2/3} c_3 \|\phi^{(0)}(x_i)\|, \\ |V_j^{(0)}\phi^{(0)}(x_i)| &\gtrsim \left(\frac{\kappa_2}{m_2}\right)^{1/3} C_2^{2/3} c_3 \|x_i'\|. \end{aligned}$$

### Proof of Lemma 7

By assumption, we know that during the algorithm, we have  $\|V'\| \leq C_2$ . Also, we know by Lemma 3 that under  $E^c$ :

$$\|\phi^{(2)}(x_i)\| \leq 2C_1.$$

Define the set

$$P'_i = \{j \mid |V_j^{(0)}\phi^{(0)}(x_i)| \leq c_3 \left(\frac{\kappa_2}{m_2}\right)^{1/3} C_2^{2/3} \|\phi^{(0)}(x_i)\|\} \quad (113)$$

and  $P' = \cup P'_i$ . We have

$$\mathbb{P}(|V_j^{(0)}\phi^{(0)}(x_i)| \leq c_3 \left(\frac{\kappa_2}{m_2}\right)^{1/3} C_2^{2/3} \|\phi^{(0)}(x_i)\|) \leq c_3 \left(\frac{\kappa_2}{m_2}\right)^{1/3} C_2^{2/3} / \kappa_2,$$

so by Bernstein, with high probability:

$$|P'_i| \leq m_2 C_2^{2/3} c_3 \left(\frac{\kappa_2}{m_2}\right)^{1/3} / \kappa_2 + \sqrt{m_2 C_2^{2/3} c_3 \left(\frac{\kappa_2}{m_2}\right)^{1/3} / \kappa_2 + 1} \lesssim c_3 m_2 C_2^{2/3} \left(\frac{\kappa_2}{m_2}\right)^{1/3} / \kappa_2,$$

so with high prob.

$$|P'_i| \lesssim c_3 C_2^{2/3} \left(\frac{m_2}{\kappa_2}\right)^{2/3}. \quad (114)$$

On the other hand, Note that

$$\phi_k^{(0)}(x_i) = \sum_{j=1}^{m_1} 1/\sqrt{m_1} W_{k,j}^s \sigma(W_j^{(0)} x_i) \quad (115)$$

is subGaussian with parameter  $\sigma^2 = O(1/m_1 \sum_j \sigma(W_j^{(0)} x_i)^2)$ . Furthermore, note that if we compute the variance of  $\phi_k^{(0)}(x_i)$  with respect to the randomness of  $W^s$ :

$$\mathbb{E} \phi_k^{(0)}(x_i)^2 = 1/m_1 \sum_{j=1}^{m_1} \sigma(W_j^{(0)} x_i)^2 := \aleph$$

which itself concentrates around  $1/2\kappa_1^2 \|x_i\|^2 = 1/2\kappa_1^2$  by another Bernstein, i.e.  $\aleph = 1/2\kappa_1^2(1 \pm O(1/\sqrt{m_1}))$ . Therefore, by concentration of subexponential variables (Bernstein), it is not hard to see that the squared norm of the vector  $\phi^{(0)}(x_i)$  is  $(m_3 \kappa_1^4, \kappa_2)$ -subexponential and concentrates around  $m_3 \aleph$ , i.e.

$$\|\phi^{(0)}(x_i)\|^2 = m_3 \aleph \pm O(\kappa_1^2 \sqrt{m_3}) = m_3 \kappa_1^2 / 2 \pm O(m_3 \kappa_1^2 / \sqrt{m_1}) \pm O(\kappa_1^2 \sqrt{m_3}), \quad (116)$$

with high probability. Combining this with the fact that  $\|\phi^{(2)}(x_i)\| \lesssim C_1$  implies with high probability:

$$\frac{\|\phi^{(0)}(x_i)\|}{\|\phi^{(2)}(x_i)\|} \gtrsim \frac{\sqrt{m_3} \kappa_1}{C_1}. \quad (117)$$

Now define  $P''_i = \{j \mid |V_j' x_i'| \geq |V_j^{(0)}\phi^{(0)}(x_i)|/3\}$ . If  $j \in P''_i - P'_i$ , then by Equation (117), with high probability

$$\|V_j'\| \|\phi^{(2)}(x_i)\| \geq |V_j'\phi^{(2)}(x_i)| = |V_j'(\phi^{(0)}(x_i) + \phi^{(2)}(x_i))| = |V_j' x_i'|$$

$$\geq |V_j^{(0)} \phi^{(0)}(x_i)|/3 \gtrsim c_3 \left(\frac{\kappa_2}{m_2}\right)^{1/3} C_2^{2/3} \|\phi^{(0)}(x_i)\|,$$

or

$$\|V_j'\|^2 \gtrsim c_3^2 \left(\frac{\kappa_2}{m_2}\right)^{2/3} C_2^{4/3} \frac{m_3 \kappa_1^2}{C_1^2}.$$

But note that  $\|V'\|_F^2 \leq C_2^2$  by our assumption, which implies

$$|P_i'' - P_i'| \lesssim C_2^2 / [c_3^2 \left(\frac{\kappa_2}{m_2}\right)^{2/3} C_2^{4/3} \frac{m_3 \kappa_1^2}{C_1^2}] = \frac{C_2^{2/3} C_1^2}{c_3^2 m_3 \kappa_1^2} \left(\frac{m_2}{\kappa_2}\right)^{2/3}. \quad (118)$$

Now combining Equations (114) and (118), we finally obtain

$$|P_i''| = |P_i'' - P_i'| + |P_i'| \lesssim \left(c_3 + \frac{C_1^2}{c_3^2 m_3 \kappa_1^2}\right) C_2^{2/3} \left(\frac{m_2}{\kappa_2}\right)^{2/3}.$$

Now define the set

$$P_i''' = \{j \mid |V_j^{(0)} \phi^{(2)}(x_i)| \geq |V_j^{(0)} \phi^{(0)}(x_i)|/3\}. \quad (119)$$

Note that for every  $j \in [m_2]$ ,  $V_j^{(0)} \phi^{(0)}(x_i)$  is gaussian with variance  $\|\phi^{(0)}(x_i)\|^2$  over the randomness of  $V_j^{(0)}$ , so

$$\mathbb{P}(V_j^{(0)} \phi^{(0)}(x_i) \leq \alpha \kappa_2 \|\phi^{(0)}(x_i)\|) \lesssim \alpha.$$

Therefore, if we define the set

$$Q_i = \{j \in [m_2] \mid |V_j^{(0)} \phi^{(0)}(x_i)| \leq \alpha \kappa_2 \|\phi^{(0)}(x_i)\|\},$$

then for large enough  $m_2$ , by Bernstein with high prob.:

$$|Q_i| \lesssim \alpha m_2. \quad (120)$$

Now note that  $\phi^{(0)}(x_i)$  is fixed during the algorithm. On the other hand, by random matrix theory, we know that with high probability, the eigenvalues of the matrix  $V^{(0)}$  are in  $(\kappa_2(\sqrt{m_2} - \sqrt{m_3}), \kappa_2(\sqrt{m_2} + \sqrt{m_3}))$ . Therefore, even if the vector  $\phi^{(2)}(x_i)$  is picked adversarially (because it keeps changing during the algorithm), we get that with high probability over the randomness of  $V^{(0)}$ :

$$\|V^{(0)} \phi^{(2)}(x_i)\|^2 \leq \kappa_2^2 (\sqrt{m_2} + \sqrt{m_3})^2 \|\phi^{(2)}(x_i)\|^2 \lesssim \kappa_2^2 m_2 \|\phi^{(2)}(x_i)\|^2. \quad (121)$$

Moreover, because  $\|\phi^{(2)}(x_i)\| \lesssim C_1$  and from Equation (116), with high probability over the randomness of  $W^{(0)}$ :

$$\frac{\|\phi^{(0)}(x_i)\|}{\|\phi^{(2)}(x_i)\|} \gtrsim \frac{\sqrt{m_3} \kappa_1}{C_1}.$$

This means that for  $j \in P_i''' - Q_i$ , combining these inequalities we conclude with high probability

$$|V_j^{(0)} \phi^{(2)}(x_i)| \geq |V_j^{(0)} \phi^{(0)}(x_i)|/3 \gtrsim \alpha \kappa_2 \|\phi^{(0)}(x_i)\| \geq \alpha \kappa_2 \frac{\sqrt{m_3} \kappa_1}{C_1} \|\phi^{(2)}(x_i)\|,$$

which combined with (121) implies

$$\|P_i'''\| \lesssim \frac{m_2 C_1^2}{m_3 \kappa_1^2 \alpha^2}.$$

Balancing this term with the one in Equation (120), we set

$$\alpha := \frac{C_1^{2/3}}{m_3^{1/3} \kappa_1^{2/3}},$$

which implies

$$|P_i'''\| \lesssim |P_i''' - Q_i| + |Q_i| \leq \left(\frac{C_1^2}{m_3 \kappa_1^2}\right)^{1/3} m_2.$$

Defining  $\tilde{P}_i = P_i'' \cup P_i'''$ , we finally get

$$|\tilde{P}_i| \lesssim \left(\left(\frac{C_1^2}{m_3 \kappa_1^2}\right)^{1/3} + \left(c_3 + \frac{C_1^2}{c_3^2 m_3 \kappa_1^2}\right) C_2^{2/3} \left(\frac{1}{\kappa_2^2 m_2}\right)^{1/3}\right) m_2.$$

Clearly by the definition of  $P_i''$  and  $P_i'''$  the proof is complete.

**Corollary 5.1** *Under the condition  $\|\phi^{(2)}(x_i)\| \lesssim C_1$  (which happens under the event  $E^c$  defined in Lemma 3.3), setting  $c_3 := \frac{C_1^{2/3}}{m_3^{1/3} \kappa_1^{2/3}} (\frac{\kappa_2^2 m_2}{C_2^2})^{1/3}$  in the previous Lemma, we obtain  $\forall i \in [n]$ :*

$$|\tilde{P}_i| \lesssim \left(\frac{C_1^2}{m_3 \kappa_1^2}\right)^{1/3} m_2.$$

Also for  $j \notin \tilde{P}_i$ , the conditions in (113) and (119) becomes the same as

$$|W_j^{(0)} \phi^{(0)}(x_i)| \leq \kappa_2 \frac{C_1^{2/3}}{m_3^{1/3} \kappa_1^{2/3}} \|\phi^{(0)}(x_i)\|.$$

Hence, for every  $i \in [n]$  and for  $j \notin \tilde{P}_i$ , with high probability:

$$\frac{2}{3} |W_j^{(0)} \phi^{(0)}(x_i)| \geq |W_j^{(0)} \phi^{(2)}(x_i)| + |W_j'(\phi^{(0)}(x_i) + \phi^{(2)}(x_i))|, \quad (122)$$

$$|W_j^{(0)} \phi^{(0)}(x_i)| \gtrsim \kappa_2 \frac{C_1^{2/3}}{m_3^{1/3} \kappa_1^{2/3}} \|\phi^{(0)}(x_i)\| \gtrsim \kappa_2 (\sqrt{m_3} \kappa_1 C_1^2)^{1/3}, \quad (123)$$

$$|W_j^{(0)} \phi^{(0)}(x_i)| \gtrsim \frac{C_1^{2/3}}{m_3^{1/3} \kappa_2^{2/3}} \|x_i'\|. \quad (124)$$

Next, we state concentration result for the gram matrix of  $\phi^{(0)}(x_i)$ 's.

**Lemma 8** *For every  $i_1, i_2 \in [n]$ , with high probability over the randomness of  $W^{(0)}$  and  $V^{(0)}$  we have*

$$\langle \phi^{(0)}(x_{i_1}), \phi^{(0)}(x_{i_2}) \rangle = m_3 \mathbb{E} \sigma(W_j^{(0)} x_{i_1}) \sigma(W_j^{(0)} x_{i_2}) \pm O(m_3 \kappa_1^2 / \sqrt{m_1} + \sqrt{m_3} \kappa_1^2).$$

### Proof of Lemma 8

First, we compute the expectation:

$$\begin{aligned} \mathbb{E} \langle \phi^{(0)}(x_{i_1}), \phi^{(0)}(x_{i_2}) \rangle &= 1/m_1 \sum_{j_1, j_2 \in [m_1]} \sum_{k \in [m_3]} \mathbb{E} W_{k, j_1}^s W_{k, j_2}^s \sigma(W_{j_1}^{(0)} x_{i_1}) \sigma(W_{j_2}^{(0)} x_{i_2}) \\ &= 1/m_1 \sum_{j_1 \neq j_2} \mathbb{E} \sum_{k \in [m_3]} W_{k, j_1}^s W_{k, j_2}^s \sigma(W_{j_1}^{(0)} x_{i_1}) \sigma(W_{j_2}^{(0)} x_{i_2}) + m_3/m_1 \sum_{j \in [m_1]} \sigma(W_j^{(0)} x_{i_1}) \sigma(W_j^{(0)} x_{i_2}). \end{aligned}$$

But  $\sigma(W_{j_1}^{(0)} x_{i_1}) \sigma(W_{j_2}^{(0)} x_{i_2})$  is  $(m_1 \kappa_1^4, \kappa_1^2)$ -sub-exponential, so

$$\sum_{j \in [m_1]} \sigma(W_j^{(0)} x_{i_1}) \sigma(W_j^{(0)} x_{i_2}) = m_1 \mathbb{E} \sigma(W_j^{(0)} x_{i_1}) \sigma(W_j^{(0)} x_{i_2}) \pm O(\sqrt{m_1} \kappa_1^2),$$

which means with high probability:

$$\mathbb{E} \langle \phi^{(0)}(x_{i_1}), \phi^{(0)}(x_{i_2}) \rangle = m_3 \mathbb{E} \sigma(W_j^{(0)} x_{i_1}) \sigma(W_j^{(0)} x_{i_2}) \pm O(m_3 \kappa_1^2 / \sqrt{m_1}).$$

On the other side, we know that  $\phi_k^{(0)}(x_{i_1})$  is subgaussian with parameters  $\sigma^2 = 1/m_1 \sum_j (W_j^{(0)} x_{i_1})^2 := \aleph_1$  and  $\sigma^2 = 1/m_1 \sum_j (W_j^{(0)} x_{i_2})^2 := \aleph_2$  respectively. On the other hand, we know that by Bernstein w.h.p

$$\begin{aligned} \aleph_1 &= 1/2 \kappa_1^2 (1 \pm O(1/\sqrt{m_1})), \\ \aleph_2 &= 1/2 \kappa_1^2 (1 \pm O(1/\sqrt{m_1})). \end{aligned}$$

Hence,  $\phi_k^{(0)}(x_{i_1}) \phi_k^{(0)}(x_{i_2})$  is  $(\aleph_1 \aleph_2, \sqrt{\aleph_1 \aleph_2})$ -subexponential, and so  $\langle \phi^{(0)}(x_{i_1}), \phi^{(0)}(x_{i_2}) \rangle$  is  $(m_3 \aleph_1 \aleph_2, \sqrt{\aleph_1 \aleph_2})$ -subexponential. Therefore, applying another Bernstein on the top, we get

$$\begin{aligned} \langle \phi^{(0)}(x_{i_1}), \phi^{(0)}(x_{i_2}) \rangle &= \mathbb{E} \langle \phi^{(0)}(x_{i_1}), \phi^{(0)}(x_{i_2}) \rangle \pm O(\sqrt{m_3} \sqrt{\aleph_1 \aleph_2}) \\ &= m_3 \mathbb{E} \sigma(W_j^{(0)} x_{i_1}) \sigma(W_j^{(0)} x_{i_2}) \pm O(m_3 \kappa_1^2 / \sqrt{m_1}) \pm \frac{\sqrt{m_3} \kappa_1^2}{2} (1 \pm O(1/\sqrt{m_1})) \\ &= m_3 \mathbb{E} \sigma(W_j^{(0)} x_{i_1}) \sigma(W_j^{(0)} x_{i_2}) \pm O(m_3 \kappa_1^2 / \sqrt{m_1} + \sqrt{m_3} \kappa_1^2). \end{aligned}$$

Now we define the matrix  $L_i \in \mathbb{R}^{m_3 \times m_2}$ , with its  $j$ th column  $L_{i,j}$  equal to

$$\frac{a_j}{\sqrt{m_2}} \mathbb{1}\{V_j^{(0)} \phi^{(0)}(x_i) \geq 0\} \phi^*(x_i).$$

First, we state the following lemma which characterize a concentration result for the gram matrix of  $(L_i)_{i=1}^n$ .

**Lemma 9** *With high probability, we have the following approximation:*

$$\langle L_{i_1}, L_{i_2} \rangle = \langle \phi^*(x_{i_1}), \phi^*(x_{i_2}) \rangle \left[ F_2 \left( 2F_3(\langle x_{i_1}, x_{i_2} \rangle) \right) \pm O(m_1^{-1/4} + m_2^{-1/4} + m_3^{-1/4}) \right].$$

**Proof of Lemma 9**

By Hoeffding:

$$\begin{aligned} \langle L_{i_1}, L_{i_2} \rangle &= 1/m_2 \sum_{j \in m_2} \phi^*(x_{i_1})^T \phi^*(x_{i_2}) \mathbb{1}\{V_j^{(0)} \phi^{(1)}(x_{i_1}) \geq 0\} \mathbb{1}\{V_j^{(0)} \phi^{(1)}(x_{i_2}) \geq 0\} \\ &= \phi^*(x_{i_1})^T \phi^*(x_{i_2}) \left( \mathbb{E} \mathbb{1}\{V_j^{(0)} \phi^{(1)}(x_{i_1}) \geq 0\} \mathbb{1}\{V_j^{(0)} \phi^{(1)}(x_{i_2}) \geq 0\} \pm O(1/\sqrt{m_2}) \right) \\ &= \phi^*(x_{i_1})^T \phi^*(x_{i_2}) \left( F_2 \left( \langle \phi^{(1)}(x_{i_1}), \phi^{(1)}(x_{i_2}) \rangle / (\|\phi^{(1)}(x_{i_1})\| \|\phi^{(1)}(x_{i_2})\|) \right) \pm O(1/\sqrt{m_2}) \right), \end{aligned}$$

where recall

$$F_2(x) = 1/4 + \arcsin(x)/2\pi,$$

measures the angle between two unit vectors based on their dot product. Now notice that according to Lemma 8, with high probability:

$$\begin{aligned} &\langle L_{i_1}, L_{i_2} \rangle / \langle \phi^*(x_{i_1}), \phi^*(x_{i_2}) \rangle \\ &= F_2 \left( \frac{m_3 \mathbb{E} \sigma(W_j^{(0)} x_{i_1}) \sigma(W_j^{(0)} x_{i_2}) \pm O((m_3/\sqrt{m_1} + \sqrt{m_3}) \kappa_1^2)}{\sqrt{(m_3 \mathbb{E} \sigma(W_j^{(0)} x_{i_1})^2 \pm O((m_3/\sqrt{m_1} + \sqrt{m_3}) \kappa_1^2)) (m_3 \mathbb{E} \sigma(W_j^{(0)} x_{i_2})^2 \pm O(\dots))}} \pm O(1/\sqrt{m_2}) \right) \\ &= F_2 \left( \frac{F_3(\langle x_{i_1}, x_{i_2} \rangle) \pm O(1/\sqrt{m_1} + 1/\sqrt{m_3})}{1/2 \pm O(1/\sqrt{m_1} + 1/\sqrt{m_3})} \pm O(1/\sqrt{m_2}) \right), \end{aligned}$$

where recall  $F_3 : [-1, +1] \rightarrow [-1/2, 1/2]$  is defined as:

$$F_3(x) := \frac{\sqrt{1-x^2}}{2\pi} + \frac{x}{4} + \frac{x \arcsin x}{2\pi}.$$

It is easy to see  $F_3$  has the property that for unit vectors  $x_1, x_2$  and  $w$  sampled as standard normal:

$$F_3(\langle x_1, x_2 \rangle) = \mathbb{E} \sigma(w^T x_1) \sigma(w^T x_2).$$

But because  $|F_3(\cdot)| = O(1)$ , we have

$$\langle L_{i_1}, L_{i_2} \rangle / \langle \phi^*(x_{i_1}), \phi^*(x_{i_2}) \rangle = F_2 \left( 2F_3(\langle x_{i_1}, x_{i_2} \rangle) \pm O(1/\sqrt{m_1} + 1/\sqrt{m_2} + 1/\sqrt{m_3}) \right).$$

Now notice that the derivative of  $F_2$ , i.e.  $1/2\pi\sqrt{1-x^2}$  is increasing in the interval  $(0, 1)$ , so for a fixed  $\delta$ , the maximum of  $|F_2(x) - F_2(x - \delta)|$  happens at  $x = 1$ . On the other hand, by writing the first order approximation of  $\arcsin(1 - t^2)$  around  $t = 0$  and upper bounding its derivative in the interval  $[0, 1]$ , we get that for  $0 \leq \delta \leq 1$ :

$$\arcsin(1 - \delta) \geq \arcsin(1) - 2\sqrt{\delta}.$$

Therefore,  $F_2(x \pm \delta) = F_2(x) \pm O(\sqrt{\delta})$ . Hence:

$$\begin{aligned} \langle L_{i_1}, L_{i_2} \rangle / \langle \phi^*(x_{i_1}), \phi^*(x_{i_2}) \rangle &= F_2 \left( 2F_3(\langle x_{i_1}, x_{i_2} \rangle) \right) \pm O(\sqrt{1/\sqrt{m_1} + 1/\sqrt{m_2} + 1/\sqrt{m_3}}) \\ &= F_2 \left( 2F_3(\langle x_{i_1}, x_{i_2} \rangle) \right) \pm O(m_1^{-1/4} + m_2^{-1/4} + m_3^{-1/4}), \end{aligned}$$

which completes the proof.

Finally, we are ready to construct the weights  $V^*$  for the second layer.

A.12.3 CONSTRUCTION OF  $V^*$ 

**Lemma 10** *Let*

$$\mathfrak{R} = \frac{n\sqrt{m_3}}{\sqrt{m_1}\lambda_0} \left[ 1 + \frac{nC_1}{\kappa_1} \right].$$

*Suppose we have the condition that for every  $k \in [m_3]$ :*

$$\max_k \|\mathcal{V}_k\| \leq \xi, \quad (125)$$

*where recall the definition of  $\mathcal{V}_k$  in Equation (43). We assume enough overparameterization to make sure  $\mathfrak{R} < 1$ . Recall for the matrix  $A$  defined by*

$$A = \left( \langle \bar{x}_{i_1}, \bar{x}_{i_2} \rangle F_2(2F_3(\langle x_{i_1}, x_{i_2} \rangle)) \right)_{1 \leq i_1, i_2 \leq n}, \quad (126)$$

*we have*

$$(f^*(x_i))_{i=1}^n{}^T A^{-1} (f^*(x_i))_{i=1}^n \leq \zeta_2.$$

*Then, there exists weight matrix  $V^*$  which only depends on the random initializations  $W^{(0)}, V^{(0)}$  (e.g. not on  $V'$  and  $W'$ ) for the second layer, such that having enough overparameterization*

$$\|V^*\|_F^2 \leq 2\zeta_2, \quad (127)$$

*and for every  $j \in [m_2]$ :*

$$\|V_j^*\|_\infty \leq \frac{(1 + \mathfrak{R})n\sqrt{n\zeta_2}}{\sqrt{m_2}} \xi := \varrho_3 \xi / \sqrt{m_2}, \quad (128)$$

$$\|V_j^*\|_2 \leq \frac{1}{\sqrt{m_2}} n(1 + \mathfrak{R}) \sqrt{\frac{\zeta_2 \sum_k \|\mathcal{V}_k\|_{H^\infty}^2}{\lambda_0}} := \varrho_2 / \sqrt{m_2}, \quad (129)$$

*and further under the high probability event  $E^c$  defined in Lemma 33:*

$$\left| \frac{1}{\sqrt{m_2}} a^T D_{V^{(0)}+V', x_i} V^* \phi^*(x_i) - f^*(x_i) \right| \lesssim \left( \frac{C_1}{\sqrt{m_3}\kappa_1} \right)^{1/3} (1 + \mathfrak{R}) \sqrt{\zeta_2 \sum_k \|\mathcal{V}_k\|_{H^\infty}^2} := \mathfrak{R}_3. \quad (130)$$

### Proof of Lemma 11

Let

$$V^* = \sum_{i=1}^n \mathcal{V}_i^* L_i,$$

be the minimum norm vector which maps  $L_i$ 's to  $f^*(x_i)$ 's. As a result, for the matrix

$$L = \left( \langle L_{i_1}, L_{i_2} \rangle \right)_{i_1, i_2}$$

it is easy to see

$$\|V^*\|_F^2 = (f^*(x_i))_{i=1}^n{}^T L^{-1} (f^*(x_i))_{i=1}^n.$$

Now combining Lemmas 5 and 41, we get

$$\|\phi^*(x_i)\|_\infty \leq (1 + \mathfrak{R})\xi, \quad (131)$$

$$\|\phi^*(x_i)\| \leq (1 + \mathfrak{R}) \sqrt{\sum_k \|\mathcal{V}_k\|_{H^\infty}^2}, \quad (132)$$

and

$$\left| \langle \phi^*(x_{i_1}), \phi^*(x_{i_2}) \rangle - \langle \bar{x}_{i_1}, \bar{x}_{i_2} \rangle \right| \leq (2\mathfrak{R} + \mathfrak{R}^2) \sum_k \|\mathcal{V}_k\|_{H^\infty}^2.$$

Now by Lemma 9:

$$|\langle L_{i_1}, L_{i_2} \rangle - A_{i_1, i_2}| \lesssim (2\mathfrak{R} + \mathfrak{R}^2) \left( \sum_k \|\mathcal{V}_k\|_{H^\infty}^2 \right) \left| F_2 \left( 2F_3(\langle x_{i_1}, x_{i_2} \rangle) \right) \right| + \langle \bar{x}_{i_1}, \bar{x}_{i_2} \rangle (m_1^{-1/4} + m_2^{-1/4} + m_3^{-1/4}) \\ + \text{other cross term.}$$

By applying  $\langle \bar{x}_{i_1}, \bar{x}_{i_2} \rangle \leq \|\bar{x}_{i_1}\| \|\bar{x}_{i_2}\|$  we get

$$\begin{aligned} LHS &\leq \left( \sum_k \|\mathcal{V}_k\|_{H^\infty}^2 \right) \left( (2\mathfrak{R} + \mathfrak{R}^2) \left| F_2 \left( 2F_3(\langle x_{i_1}, x_{i_2} \rangle) \right) \right| + (m_1^{-1/4} + m_2^{-1/4} + m_3^{-1/4}) \right) \\ &\lesssim \left( \sum_k \|\mathcal{V}_k\|_{H^\infty}^2 \right) \left( \mathfrak{R} + m_1^{-1/4} + m_2^{-1/4} + m_3^{-1/4} \right). \end{aligned}$$

Therefore,

$$\begin{aligned} \|A - (\langle L_{i_1}, L_{i_2} \rangle)_{i_1, i_2}\|_2 &\leq \|A - (\langle L_{i_1}, L_{i_2} \rangle)_{i_1, i_2}\|_F \\ &\leq n \left( \sum_k \|\mathcal{V}_k\|_{H^\infty}^2 \right) \left( \mathfrak{R} + m_1^{-1/4} + m_2^{-1/4} + m_3^{-1/4} \right) := \mathfrak{R}_2. \end{aligned}$$

Note that  $\mathfrak{R}_2$  naturally goes to zero (with poly dependence) as  $\mathfrak{R} \rightarrow 0$  and  $m_1, m_2, m_3$  are large enough. Now if all of the eigenvalues of the matrix  $A$  are  $\Omega(1/n^2)$ , then if we overparameterize enough such that  $\mathfrak{R}_2 = O(1/n^2)$  with small enough constant so that  $\mathfrak{R}_2$  is less than half of the smallest eigenvalue of  $A$ , then for the  $i$ th eigenvalue  $\lambda_i$  of  $A$  and  $L$  we can write

$$\lambda_i(L) \geq \lambda_i(A) - \mathfrak{R}_2 \geq \lambda_i(A)/2,$$

so

$$\lambda_i(L^{-1}) \leq 2\lambda_i(A^{-1}),$$

which implies the property

$$\|V^*\|_F^2 \leq 2\zeta_2. \quad (133)$$

However,  $A$  might have very small eigenvalues. To remedie this, we use Lemma 42; we can substitute  $f^*$  with some  $\bar{f}^*$  such that

$$R_n(\bar{f}^*) \leq 2R_n(f^*) + \frac{B^2}{n}, \quad (134)$$

$$\bar{f}^{*T} A^{-1} \bar{f}^* \leq f^{*T} A^{-1} f^*, \quad (135)$$

where  $\bar{f}^*$  is on the subspace of eigenvectors of  $A$  whose eigenvalues are larger than  $\Omega(1/n^2)$ . But it is easy to check that in the context of Theorem 3, such substitution results in a  $\bar{f}^{*T} A^{-1} \bar{f}^* \leq f^{*T} A^{-1} f^* \leq \zeta$  and  $\bar{\nu}(f^*)$  parameter (as defined in (50)) with respect to  $\bar{f}^*$  which satisfies  $\bar{\nu}/2 \leq \nu$ . This enables us to use remark ?? with respect to  $\bar{f}^*$ . Note that the algorithm is with respect to the setting  $\nu$ , however we want to exploit generalization bound with respect to  $\bar{f}^*$  whose parameter is  $\bar{\nu}$  as it enables us to use our analysis in this Lemma. Furthermore, note that using Equation (134) we can further upper bound the empirical risk of  $\bar{f}^*$  with that of  $f^*$ , which makes it straightforward to derive a similar generalization bound as in (55) with respect to  $\bar{f}^*$ , of course with a change of constants. Therefore, without loss of generality we can use substitute  $f^*$  by  $\bar{f}^*$  and still obtain Equation (133).

On the other hand, the definition of  $V^*$  implies

$$\frac{1}{\sqrt{m_2}} a^T D_{V^{(0)}, x_i} V^* \phi^*(x_i) = f^*(x_i).$$

But note that by Corollary 51, under the high probability event  $E^c$  defined in Lemma 33,  $D_{V^{(0)}, x_i}$  and  $D_{V^{(0)+V', x_i}}$  can only be different in the index set  $\tilde{P}_i$  and

$$|\tilde{P}_i| \lesssim \left( \frac{C_1^2}{m_3 \kappa_1^2} \right)^{1/3} m_2,$$

Therefore, for all  $i \in [n]$ :

$$\begin{aligned}
& \left| \frac{1}{\sqrt{m_2}} a^T D_{V^{(0)}+V',x_i} V^* \phi^*(x_i) - \frac{1}{\sqrt{m_2}} a^T D_{V^{(0)},x_i} V^* \phi^*(x_i) \right| \\
& \leq 1/\sqrt{m_2} \sum_{j \in \tilde{P}} |V_j^* \phi^*(x_i)| \\
& \leq 1/\sqrt{m_2} \sum_{j \in \tilde{P}} \|V_j^*\| \|\phi^*(x_i)\| \\
& \leq \frac{\sqrt{|\tilde{P}|}}{\sqrt{m_2}} \|V^*\| (1 + \Re) \sqrt{\sum_k \|\mathcal{V}_k\|_{H^\infty}^2} \\
& \lesssim \left(\frac{C_1}{\sqrt{m_3 \kappa_1}}\right)^{1/3} \sqrt{\zeta_2} (1 + \Re) \sqrt{\sum_k \|\mathcal{V}_k\|_{H^\infty}^2},
\end{aligned}$$

which proves the first claim. On the other hand, we get:

$$2\zeta_2 \geq \|V^*\|_F^2 \geq \mathcal{V}^T L \mathcal{V}.$$

But because  $\lambda_{\min}(L) \gtrsim 1/n^2$ , we get

$$\mathcal{V}^T L \mathcal{V} \gtrsim \|\mathcal{V}\|_2^2/n^2,$$

which implies

$$\|\mathcal{V}\|_2 \lesssim n\sqrt{\zeta_2}.$$

But now using Equation (131), we can write

$$\begin{aligned}
|V_{j,k}^*| & \leq \sum_{i=1}^n |\mathcal{V}_i| |L_{i,j,k}| \leq \frac{1}{\sqrt{m_2}} \|\phi^*(x_i)\|_\infty \sum_i |\mathcal{V}_i| \\
& \lesssim \frac{(1 + \Re)\xi}{\sqrt{m_2}} \sum_i |\mathcal{V}_i| \leq (1 + \Re)\xi \sqrt{n} \|\mathcal{V}\|/\sqrt{m_2} \\
& \lesssim \frac{(1 + \Re)n\sqrt{n\zeta_2}}{\sqrt{m_2}} \xi,
\end{aligned}$$

which proves the other part. Moreover,

$$\|V_j^*\|^2 \leq \|\mathcal{V}\|^2 \frac{1}{\sqrt{m_2}} \left(\sum_i \|\phi^*(x_i)\|_2^2\right) \leq \frac{1}{\sqrt{m_2}} n(1 + \Re) \sqrt{\zeta_2 \sum_k \|\mathcal{V}_k\|_{H^\infty}^2/\lambda_0}. \quad (136)$$

## A.13 EXISTENCE OF A GOOD DIRECTION

Our aim in this section is to show that if the objective value is above certain threshold, there exists a good random direction which reduces the objective in expectation. Particularly our aim is to prove the following theorem (informal):

**Theorem 6** *For a given pair  $(f^*, G)$  with*

$$\begin{aligned} \langle H^\infty, G \rangle &\leq \zeta_1, \\ f^{*T} (K^\infty \odot G)^{-1} f^* &\leq \zeta_2, \\ R_n(f^*) &\leq \Delta, \end{aligned}$$

*recall the ideal random matrices  $(W_\Sigma^*, V_\Sigma^*)$  constructed in Appendix A.12, where  $\Sigma$  is a random diagonal sign matrix. Specifically,  $W_\Sigma^*$  is defined in Equation (108), and  $V_\Sigma^*$  is the projection of the rows of matrix  $V^*$  onto the orthogonal subspace spanned by  $(\phi^{(0)}(x_i))_{i=1}^n$ .*

*Using the parameter setting for  $i = 1, 2$*

$$\psi_i = \frac{\nu}{4\zeta_i}, \quad (137)$$

*with respect to an arbitrary parameter  $\nu > 0$ , then for every pair  $(W', V')$  such that  $\|W'\| \leq C_1$ ,  $\|V'\| \leq C_2$  and*

$$L(W', V') \geq \Delta + \nu, \quad (138)$$

*for parameters  $m_1, m_2, m_3, 1/\kappa_1, 1/\kappa_2$  polynomially large enough in  $B, 1/\lambda_0, n, C_1, C_2$  and small enough step size  $\eta$ , we have*

$$\mathbb{E}_\Sigma L(W' - \eta/2W' + \sqrt{\eta}W_\Sigma^*, V' - \eta/2V' + \sqrt{\eta}V_\Sigma^*) \leq L(W', V') - \eta\nu/4. \quad (139)$$

In order to prove the above theorem, we first state and prove the following lemma which is the core of Theorem 6.

**Lemma 11** *For matrices  $(W^*, V^*)$  constructed in Appendix A.12, specifically for their random coupling  $(W_\Sigma^*, V_\Sigma^*)$  as denoted above, we have:*

$$\mathbb{E}_\Sigma \ell(f'_{(1-\eta/2)W' + \sqrt{\eta}W_\Sigma^*}, (1-\eta/2)V' + \sqrt{\eta}V_\Sigma^*(x_i), y_i) \leq (1-\eta)\ell(f'_{W', V'}(x_i), y_i) + \eta\ell(f^*(x_i), y_i) \pm \eta\wp,$$

*where  $\wp$  goes to zero with polynomially large overparameterization (the exact dependence is revealed via the proof).*

**Proof of Lemma 11**

For brevity, we use the notation  $D', \rho$  here to refer to the diagonal binary sign matrix when the input is multiplied by the sum of weight and smoothing matrices. It will be clear in the context of the equation that what the ‘‘input’’ and the ‘‘weight’’ matrices are. This notation is also defined and used in Lemma 28. Here, we bound multiple cross terms that are created as a result of moving in the random direction. To simplify the presentation and avoid confusing recursions in the proof, we have made a sublemma for each of these cross terms and has deferred its proof to Appendix A.14. We use difference sub-indices of the symbol  $\mathfrak{R}$  to illustrate terms that go to zero by growing the overparameterization in our architecture.

We start by using Lemma [D8](#),

$$\begin{aligned}
& \mathbb{E}_{\Sigma} \ell(f'_{(1-\eta/2)W' + \sqrt{\eta}W_{\Sigma}^*, (1-\eta/2)V' + \sqrt{\eta}V_{\Sigma}^*}(x_i), y_i) \\
&= \mathbb{E}_{\Sigma} \ell(\mathbb{E}_{W^{\rho}, V^{\rho}} f_{(1-\eta/2)W' + \sqrt{\eta}W_{\Sigma}^* + W^{\rho}, (1-\eta/2)V' + \sqrt{\eta}V_{\Sigma}^* + V^{\rho}}(x_i), y_i) \\
&= \mathbb{E}_{\Sigma} \ell\left(\mathbb{E}_{W^{\rho}, V^{\rho}} a^T D_{\cdot, \rho}(V^{(0)} + (1-\eta/2)V' + V^{\rho} + \sqrt{\eta}V_{\Sigma}^*) W^s D_{\cdot, \rho}(W^{(0)} + (1-\eta/2)W' + W^{\rho} + \sqrt{\eta}W_{\Sigma}^*) x_i \right. \\
&\quad \left. + \Re_8 \eta, y_i\right). \\
&= \mathbb{E}_{\Sigma} \ell\left(\mathbb{E}_{W^{\rho}, V^{\rho}} \left[ a^T D_{\cdot, \rho}(V^{(0)} + (1-\eta/2)V' + V^{\rho}) W^s D_{\cdot, \rho}(W^{(0)} + (1-\eta/2)W' + W^{\rho}) x_i \right. \right. \\
&\quad \left. \left. + \eta a^T D_{\cdot, \rho} V_{\Sigma}^* W^s D_{\cdot, \rho} W_{\Sigma}^* x_i \right] + \sqrt{\eta} \mathbb{E}_{W^{\rho}, V^{\rho}} \left[ a^T D_{\cdot, \rho}(V^{(0)} + (1-\eta/2)V' + V^{\rho}) W^s D_{\cdot, \rho} W_{\Sigma}^* x_i \right. \right. \\
&\quad \left. \left. + a^T D_{\cdot, \rho} V_{\Sigma}^* W^s D_{\cdot, \rho}(W^{(0)} + (1-\eta/2)W' + W^{\rho}) x_i \right] \right. \\
&\quad \left. + \Re_8 \eta, y_i\right).
\end{aligned}$$

Now using the notation introduced in Lemma [D8](#), we have

$$W^s D_{\cdot, \rho}(W^{(0)} + (1-\eta/2)W' + W^{\rho}) x_i = \phi^{(0)}(x_i) + (1-\eta/2)\phi^{(2)}(x_i) + \frac{\eta}{2}\phi^{(2)'}(x_i).$$

By Lemma [D8](#), we have the following bound for  $\phi^{(2)'}(x_i)$ :

$$\begin{aligned}
& \mathbb{E}_{W^{\rho}, V^{\rho}} \left| \frac{1}{\sqrt{m_2}} a^T D_{V^{(0)} + V^{\rho} + V', x_i}(V^{(0)} + V^{\rho} + (1-\eta/2)V') \phi^{(2)'}(x_i) \right| \\
& \lesssim (\kappa_2 \sqrt{m_2 m_3} + \sqrt{m_3} \beta_2 + C_2) \Re_5.
\end{aligned}$$

Therefore, Combining this with Lemma [D1](#), we get

$$\begin{aligned}
&= \mathbb{E}_{\Sigma} \ell\left(\mathbb{E}_{W^{\rho}, V^{\rho}} \left[ a^T D_{\cdot, \rho}(V^{(0)} + (1-\eta/2)V' + V^{\rho}) W^s (\phi^{(0)}(x_i) + (1-\eta/2)\phi^{(2)}(x_i)) \right. \right. \\
&\quad \left. \left. + \eta a^T D_{\cdot, \rho} V_{\Sigma}^* W^s D_{\cdot, \rho} W_{\Sigma}^* x_i \right] + \sqrt{\eta} \mathbb{E}_{W^{\rho}, V^{\rho}} \left[ a^T D_{\cdot, \rho}(V^{(0)} + (1-\eta/2)V' + V^{\rho}) W^s D_{\cdot, \rho} W_{\Sigma}^* x_i \right. \right. \\
&\quad \left. \left. + a^T D_{\cdot, \rho} V_{\Sigma}^* W^s D_{\cdot, \rho}(W^{(0)} + (1-\eta/2)W' + W^{\rho}) x_i \right] \right. \\
&\quad \left. \pm O((\kappa_2 \sqrt{m_2 m_3} + \sqrt{m_3} \beta_2 + C_2) \Re_5 \eta) \pm O(\Re_8 \eta), y_i\right) \\
&= \mathbb{E}_{\Sigma} \ell\left(\mathbb{E}_{W^{\rho}, V^{\rho}} \left[ (1-\eta) a^T D_{\cdot, \rho}(V^{(0)} + V' + V^{\rho}) W^s (\phi^{(0)}(x_i) + \phi^{(2)}(x_i)) \right. \right. \\
&\quad \left. \left. + \eta a^T D_{\cdot, \rho} V_{\Sigma}^* W^s D_{\cdot, \rho} W_{\Sigma}^* x_i \right] + \sqrt{\eta} \mathbb{E}_{W^{\rho}, V^{\rho}} \left[ a^T D_{\cdot, \rho}(V^{(0)} + (1-\eta/2)V' + V^{\rho}) W^s D_{\cdot, \rho} W_{\Sigma}^* x_i \right. \right. \\
&\quad \left. \left. + a^T D_{\cdot, \rho} V_{\Sigma}^* W^s D_{\cdot, \rho}(W^{(0)} + (1-\eta/2)W' + W^{\rho}) x_i \right] \right. \\
&\quad \left. \pm O(\eta(\Re_6' + \Re_4 + (\sqrt{m_3} \kappa_2 + \beta_2)(C_1 + \sqrt{m_3} \beta_1))) \pm O((\kappa_2 \sqrt{m_2 m_3} + \sqrt{m_3} \beta_2 + C_2) \Re_5 \eta) \pm O(\Re_8 \eta), y_i\right). \\
&= \mathbb{E}_{\Sigma} \ell\left((1-\eta) f'_{W', V'}(x_i) + \eta \mathbb{E}_{W^{\rho}, V^{\rho}} a^T D_{\cdot, \rho} V_{\Sigma}^* W^s D_{\cdot, \rho} W_{\Sigma}^* x_i \right. \\
&\quad \left. + \sqrt{\eta} \mathbb{E}_{W^{\rho}, V^{\rho}} \left[ a^T D_{\cdot, \rho}(V^{(0)} + (1-\eta/2)V' + V^{\rho}) W^s D_{\cdot, \rho} W_{\Sigma}^* x_i \right. \right. \\
&\quad \left. \left. + a^T D_{\cdot, \rho} V_{\Sigma}^* W^s D_{\cdot, \rho}(W^{(0)} + (1-\eta/2)W' + W^{\rho}) x_i \right] \right. \\
&\quad \left. \pm O(\eta(\Re_6' + \Re_4 + (\sqrt{m_3} \kappa_2 + \beta_2)(C_1 + \sqrt{m_3} \beta_1))) \pm O((\kappa_2 \sqrt{m_2 m_3} + \sqrt{m_3} \beta_2 + C_2) \Re_5 \eta) \pm O(\Re_8 \eta), y_i\right).
\end{aligned}$$

Moreover, using the notation  $\bar{\phi}^{*'}(x_i)$  introduced in Lemma [15](#) and the bound in Lemma [17](#), we can rewrite the second term as:

$$\begin{aligned}
LHS &= \mathbb{E}_\Sigma \ell \left( (1-\eta) f'_{W',V'}(x_i) + \eta \mathbb{E}_{W^\rho, V^\rho} a^T D_{',\rho} V_\Sigma^* (\phi^*(x_i) + \bar{\phi}^{*'}(x_i)) \right. \\
&\quad + \sqrt{\eta} \mathbb{E}_{W^\rho, V^\rho} \left[ a^T D_{',\rho} (V^{(0)} + (1-\eta/2)V' + V^\rho) W^s D_{',\rho} W_\Sigma^* x_i \right. \\
&\quad \left. \left. + a^T D_{',\rho} V_\Sigma^* W^s D_{',\rho} (W^{(0)} + (1-\eta/2)W' + W^\rho) x_i \right] \right. \\
&\quad \left. \pm O(\eta(\mathfrak{R}'_6 + \mathfrak{R}_4 + (\sqrt{m_3}\kappa_2 + \beta_2)(C_1 + \sqrt{m_3}\beta_1))) \pm O((\kappa_2\sqrt{m_2m_3} + \sqrt{m_3}\beta_2 + C_2)\mathfrak{R}_5\eta) \pm O(\mathfrak{R}_8\eta), y_i \right) \\
&= \mathbb{E}_\Sigma \ell \left( (1-\eta) f'_{W',V'}(x_i) + \eta \mathbb{E}_{W^\rho, V^\rho} a^T D_{',\rho} V_\Sigma^* \phi^*(x_i) \right. \\
&\quad + \sqrt{\eta} \mathbb{E}_{W^\rho, V^\rho} \left[ a^T D_{',\rho} (V^{(0)} + (1-\eta/2)V' + V^\rho) W^s D_{',\rho} W_\Sigma^* x_i \right. \\
&\quad \left. \left. + a^T D_{',\rho} V_\Sigma^* W^s D_{',\rho} (W^{(0)} + (1-\eta/2)W' + W^\rho) x_i \right] \right. \\
&\quad \left. \pm O(\eta\mathfrak{R}_{10}) \pm O(\eta(\mathfrak{R}'_6 + \mathfrak{R}_4 + (\sqrt{m_3}\kappa_2 + \beta_2)(C_1 + \sqrt{m_3}\beta_1))) \pm O((\kappa_2\sqrt{m_2m_3} + \sqrt{m_3}\beta_2 + C_2)\mathfrak{R}_5\eta) \pm \right. \\
&\quad \left. O(\mathfrak{R}_8\eta), y_i \right). \tag{140}
\end{aligned}$$

Now we write the gradient-lipshitz inequality for  $\ell$  at point

$$p_\Sigma^{(1)} := (1-\eta) f'_{W',V'}(x_i) + \eta \mathbb{E}_{W^\rho, V^\rho} a^T D_{',\rho} V_\Sigma^* \phi^*(x_i) \pm \eta \wp_1,$$

and regarding the following vector, where  $\wp_1$  is the sum of all the noise terms above and goes to zero by over parameterization:

$$\begin{aligned}
p_\Sigma^{(2)} &:= \sqrt{\eta} \mathbb{E}_{W^\rho, V^\rho} \left[ a^T D_{',\rho} (V^{(0)} + (1-\eta/2)V' + V^\rho) W^s D_{',\rho} W_\Sigma^* x_i \right. \\
&\quad \left. + a^T D_{',\rho} V_\Sigma^* W^s D_{',\rho} (W^{(0)} + (1-\eta/2)W' + W^\rho) x_i \right].
\end{aligned}$$

Hence, using the 1 smoothness of  $\ell(\cdot, y_i)$ :

$$LHS \leq \mathbb{E}_\Sigma \ell \left( p_\Sigma^{(1)} \right) + \mathbb{E}_\Sigma \dot{\ell}(p) \sqrt{\eta} p_\Sigma^{(2)} + \frac{1}{2} \eta \mathbb{E}_\Sigma (p_\Sigma^{(2)})^2. \tag{141}$$

But note that

$$\mathbb{E}_\Sigma \dot{\ell}(p) \sqrt{\eta} p_\Sigma^{(2)} = \dot{\ell}(p) \sqrt{\eta} \mathbb{E}_\Sigma p_\Sigma^{(2)} = 0. \tag{142}$$

On the other hand, using the notation of Lemma [15](#) and the result of Lemma [16](#):

$$\mathbb{E}_\Sigma \left( \mathbb{E}_{W^\rho, V^\rho} a^T D_{',\rho} (V^{(0)} + (1-\eta/2)V' + V^\rho) W^s D_{',\rho} W_\Sigma^* x_i \right)^2 \tag{143}$$

$$= \mathbb{E}_\Sigma \left( \mathbb{E}_{W^\rho, V^\rho} a^T D_{',\rho} (V^{(0)} + (1-\eta/2)V' + V^\rho) (\Sigma \phi^*(x_i) + \bar{\phi}^{*'}(x_i)) \right)^2 \tag{144}$$

$$\leq 4 \mathbb{E}_\Sigma \left( \mathbb{E}_{W^\rho, V^\rho} a^T D_{',\rho} (V^{(0)} + (1-\eta/2)V' + V^\rho) \Sigma \phi^*(x_i) \right)^2 \tag{145}$$

$$+ 4 \mathbb{E}_\Sigma \left( \mathbb{E}_{W^\rho, V^\rho} a^T D_{',\rho} (V^{(0)} + (1-\eta/2)V' + V^\rho) \bar{\phi}^{*'}(x_i) \right)^2 \tag{146}$$

$$\lesssim \mathfrak{R}_{12}^2 + \mathfrak{R}_{11}^2. \tag{147}$$

Moreover, using again the result on  $\phi^{(2)'}(x_i)$  from Lemma [18](#) and the fact that  $\phi^{(0)}(x_i)$  is orthogonal to the rows of  $V_\Sigma^*$ :

$$\begin{aligned}
&\mathbb{E}_{W^\rho, V^\rho} a^T D_{',\rho} V_\Sigma^* W^s D_{',\rho} (W^{(0)} + (1-\eta/2)W' + W^\rho) x_i \\
&= a^T D_{',\rho} V_\Sigma^* (\phi^{(0)}(x_i) + (1-\eta/2)\phi^{(2)}(x_i)) \\
&\quad + \frac{\eta}{2} a^T D_{',\rho} V_\Sigma^* \phi^{(2)'}(x_i) \\
&= (1 - \frac{\eta}{2}) a^T D_{',\rho} V_\Sigma^* \phi^{(2)}(x_i) \\
&\quad + \frac{\eta}{2} a^T D_{',\rho} V_\Sigma^* W^s D_{',\rho} \phi^{(2)'}(x_i) \\
&\lesssim (1 - \frac{\eta}{2}) a^T D_{',\rho} V_\Sigma^* \phi^{(2)}(x_i) \pm \mathfrak{R}_5.
\end{aligned}$$

Combining the last Equation with Lemma [4](#):

$$\begin{aligned} & \mathbb{E}_\Sigma \left( \mathbb{E}_{W^\rho, V^\rho} a^T D_{\cdot, \rho} V_\Sigma^* W^s D_{\cdot, \rho} (W^{(0)} + (1 - \eta/2)W' + W^\rho) x_i \right)^2 \\ & \lesssim (1 - \frac{\eta}{2})^2 a^T D_{\cdot, \rho} V_\Sigma^* W^s D_{\cdot, \rho} \phi^{(2)}(x_i) + \mathfrak{R}_5^2 \\ & \lesssim \mathfrak{R}_0^2 + \mathfrak{R}_5^2. \end{aligned} \quad (148)$$

Combining Equations [47](#) and [48](#):

$$\mathbb{E}_\Sigma (p_\Sigma^{(2)})^2 \lesssim \mathfrak{R}_{11}^2 + \mathfrak{R}_{12}^2 + \mathfrak{R}_0^2 + \mathfrak{R}_5^2 := \wp_2. \quad (149)$$

Combining Equations [42](#) and [49](#), plugging into [41](#), and reopening the definition of  $p_\Sigma^{(1)}$ :

$$LHS \lesssim \mathbb{E}_\Sigma \ell \left( (1 - \eta) f'_{W', V'}(x_i) + \eta \mathbb{E}_{W^\rho, V^\rho} a^T D_{\cdot, \rho} V_\Sigma^* \phi^*(x_i) \pm \eta \wp_1, y_i \right) + \eta \mu_2 \wp_2. \quad (150)$$

Now note that we can easily bound the magnitude of the term  $\eta \mathbb{E}_{W^\rho, V^\rho} a^T D_{\cdot, \rho} V_\Sigma^* \phi^*(x_i)$  as:

$$\begin{aligned} & |\mathbb{E}_{W^\rho, V^\rho} a^T D_{\cdot, \rho} V_\Sigma^* \phi^*(x_i)| \leq \mathbb{E}_{W^\rho, V^\rho} |a^T D_{\cdot, \rho} V_\Sigma^* \phi^*(x_i)| \\ & \leq \|V_\Sigma^*\|_F \|\phi^*(x_i)\| \leq \|V^*\| \|\phi^*(x_i)\| \leq \sqrt{2\zeta_2} (1 + \mathfrak{R}) \sqrt{\sum_k \|\mathcal{V}_k\|_{H^\infty}^2}, \end{aligned}$$

while using Lemma [34](#):

$$|f'_{W', V'}(x_i)| \leq (\kappa_2 \sqrt{m_3} + \beta_2) \left( \sqrt{m_3} \kappa_1 + C_1 + \sqrt{m_3} \beta_1 \right) + C_2 (C_1 + \sqrt{m_3} \beta_1),$$

which is  $O(C_1 C_2)$  for enough overparameterization and smoothing parameters  $\beta_1, \beta_2$  as defined in [A.20](#). Furthermore, from Equations [32](#) and [27](#), we easily see that

$$\mathbb{E}_{W^\rho, V^\rho} a^T D_{\cdot, \rho} V_\Sigma^* \phi^*(x_i) \leq \sqrt{2\zeta_2} (1 + \mathfrak{R}) \sqrt{\sum_k \|\mathcal{V}_k\|_{H^\infty}^2}.$$

Now taking  $\eta$  small enough so that the bound  $\eta \sqrt{2\zeta_2} (1 + \mathfrak{R}) \sqrt{\sum_k \|\mathcal{V}_k\|_{H^\infty}^2}$  and  $\eta \wp_1$  both also be bounded of order  $O(C_1 C_2)$ , we observe that the term inside the argument of  $\ell(\cdot, y_i)$  Equation in [150](#) is  $O(C_1 C_2)$ . Hence, we can use the lipschitz parameter of  $\ell$  in the interval  $[-O(C_1 C_2), O(C_1 C_2)]$ , given by Lemma [9](#) to take out the noise term:

$$LHS \lesssim \mathbb{E}_\Sigma \ell \left( (1 - \eta) f'_{W', V'}(x_i) + \eta \mathbb{E}_{W^\rho, V^\rho} a^T D_{\cdot, \rho} V_\Sigma^* \phi^*(x_i), y_i \right) \pm \eta \wp_1 + \eta O(C_1 C_2 + B) \wp_2. \quad (151)$$

Now by applying Lemma [9](#) and writing the Lipschitz property of  $\ell$  at point  $(1 - \eta) f'_{W', V'}(x_i) = O(C_1 C_2)$ :

$$\begin{aligned} LHS & \lesssim \mathbb{E}_\Sigma \ell \left( (1 - \eta) f'_{W', V'}(x_i) + \eta f^*(x_i) \pm \eta \mathfrak{R}_9, y_i \right) \pm \eta \wp_1 + \eta O(C_1 C_2 + B) \wp_2 \\ & := \mathbb{E}_\Sigma \ell \left( (1 - \eta) f'_{W', V'}(x_i) + \eta f^*(x_i), y_i \right) \pm \eta \mathfrak{R}_9 \pm \eta \wp_1 + \eta O(C_1 C_2 + B) \wp_2 \\ & = \ell \left( (1 - \eta) f'_{W', V'}(x_i) + \eta f^*(x_i), y_i \right) \pm \eta \wp, \end{aligned}$$

where the last line is just definition. Now Convexity of  $\ell$  finishes the proof.

Next, using Lemma [11](#) we prove Theorem [6](#).

**Restating Theorem [6](#)** In the same setting as Theorem [6](#) and having enough overparameterization such that  $\wp \leq \frac{\nu}{8}$  ( $\wp$  defined in Lemma [11](#)) and polynomially small enough step size  $\eta$ , we have

$$\mathbb{E}_\Sigma L(W' - \eta W' + \sqrt{\eta} W_\Sigma^*, V' - \eta V' + \sqrt{\eta} V_\Sigma^*) \leq L(W', V') - \eta \nu / 4.$$

**Proof of Theorem 6**

First, note that taking expectation w.r.t  $\Sigma$ :

$$\begin{aligned}\mathbb{E}_\Sigma \|(1-\eta/2)W' + \sqrt{\eta}W_\Sigma^*\|^2 &= \mathbb{E}_\Sigma (1-\eta/2)^2 \|W'\|^2 + 2(1-\eta/2)\sqrt{\eta}\langle W', \sum_{k=1}^{m_3} \Sigma_k W_k^* \rangle + \eta \left\| \sum_{k=1}^{m_3} \Sigma_k W_k^* \right\|^2 \\ &= (1-\eta/2)^2 \|W'\|^2 + \eta \sum_k \|W_k^*\|^2,\end{aligned}$$

which by orthogonality of  $W_k^*$ 's:

$$LHS = (1-\eta/2)^2 \|W'\|^2 + \eta \|W^*\|^2 = (1-\eta) \|W'\|^2 + \eta \|W^*\|^2 + \eta^2 \|W'\|^2.$$

Similarly for  $V'$ :

$$\mathbb{E}_\Sigma \|(1-\eta/2)V' + \sqrt{\eta}V_\Sigma^*\|^2 = (1-\eta/2)^2 \|V'\|^2 + \eta \mathbb{E}_\Sigma \|V^*\|^2 = (1-\eta) \|V'\|^2 + \eta \|V^*\|^2 + \eta^2 \|V'\|^2.$$

Now using Lemma [□](#):

$$\begin{aligned}\mathbb{E}_\Sigma L(W' - \eta/2W' + \sqrt{\eta}W_\Sigma^*, V' - \eta/2V' + \sqrt{\eta}V_\Sigma^*) \\ \leq (1-\eta) \mathbb{E}_Z \ell(f'_{W', V'}(x), y) + \eta \mathbb{E}_Z \ell(f^*(x), y) \\ + (1-\eta) (\psi_1 \|W'\|^2 + \psi_2 \|V'\|^2) + \eta (\psi_1 \|W^*\|^2 + \psi_2 \|V^*\|^2) + \eta (\wp + \eta (\|W'\|^2 + \|V'\|^2)) \\ \leq L(W', V') - \eta (L(W', V') - \Delta - \psi_1 \zeta_1 - \psi_2 \zeta_2) + \eta (\wp + \eta (\|W'\|^2 + \|V'\|^2)),\end{aligned}$$

which by the choice of  $\psi_i$ 's is equal to

$$\begin{aligned}LHS &\leq L(W', V') - \eta (L(W', V') - \Delta - \nu/2) + \eta (\wp + \eta (\|W'\|^2 + \|V'\|^2)) \\ LHS &\leq L(W', V') - \eta \nu/2 + \eta (\wp + \eta (\|W'\|^2 + \|V'\|^2)).\end{aligned}$$

Moreover, using the condition

$$\wp \leq \nu/8,$$

and picking  $\eta$  as small as

$$\eta (\|W'\|^2 + \|V'\|^2) \leq \eta (C_1^2 + C_2^2) \leq \nu/8,$$

we finally get

$$LHS \leq L(W', V') - \eta \nu/4.$$

## A.14 EXISTENCE OF A GOOD DIRECTION HELPER LEMMAS

In this section, we state and prove the core lemmas that are used in the proof of Lemma [11](#). Notably, through all of this section, we assume the norm bounds  $\|W'\| \leq C_1$ ,  $\|V'\| \leq C_2$  and that as our usual assumption, the rows of  $V'$  are orthogonal to  $\phi^{(0)}(x_i)$ 's for all  $i \in [n]$ . A notation that we use throughout the proofs is  $V_\Sigma^*$  which refers to the projection of  $V^*\Sigma$  onto the orthogonal subspace to  $(\phi^{(0)}(x_i))_{i=1}^n$ .

**Lemma 12** *Let  $P(\cdot)$  be the projection operator onto the subspace spanned by  $(\phi^{(0)}(x_i))_{i=1}^n$ . Also, we denote the projection of rows of  $V^*\Sigma$  onto the orthogonal subspace to  $(\phi^{(0)}(x_i))_{i=1}^n$  by  $V_{\Sigma_j}^*$ . Then*

$$\mathbb{E}_\Sigma \|V_{\Sigma_j}^* - V_j^*\Sigma\|^2 \leq \varrho_3^2 \xi^2 n / m_2,$$

with high probability

$$\|V_{\Sigma_j}^* - V_j^*\Sigma\| \lesssim \frac{\varrho_3 \xi \sqrt{n}}{\sqrt{m_2}},$$

**Proof of Lemma [12](#)**

By Equation ([129](#)), we have  $\|V_j^*\|_\infty \leq \varrho_3 \xi / \sqrt{m_2}$ . Now suppose that  $u_1, \dots, u_n$  are an orthonormal basis for the subspace  $\text{span}(\phi^{(0)}(x_i))_{i=1}^n$ . Then

$$\mathbb{E}_\Sigma \|V_{\Sigma_j}^* - V_j^*\Sigma\|^2 = \mathbb{E}_\Sigma \|P(V_j^*\Sigma)\|^2 = \sum_i \sum_k V_j^{*2} u_{ik}^2 \leq \|V_j^*\|_\infty^2 n \leq \varrho_3^2 \xi^2 n / m_2.$$

Also, by Hoeffding, with high probability:

$$\|P(V_j^*\Sigma)\|^2 = \sum_i \left( \sum_{k=1}^{m_3} V_j^* u_{ik} \Sigma_k \right)^2 \lesssim n \|V_j^*\|_\infty^2,$$

which implies the second part.

**Lemma 13** *The first cross term goes away because of the definition of  $V_\Sigma^*$ . (inside the expectations is zero almost surely)*

$$\mathbb{E}_\Sigma \left( \mathbb{E}_{V^\rho, W^\rho} \left[ \frac{1}{\sqrt{m_2}} a^T D_{V^{(0)} + V^\rho + V', x_i} V_\Sigma^* \phi^{(0)}(x_i) \right] \right)^2 = 0.$$

**Lemma 14** *Second cross term:*

$$\mathbb{E}_\Sigma \left( \mathbb{E}_{V^\rho, W^\rho} \left[ \frac{1}{\sqrt{m_2}} a^T D_{V^{(0)} + V^\rho + V', x_i} V_\Sigma^* \phi^{(2)}(x_i) \right] \right)^2 \quad (152)$$

$$\lesssim \xi^2 ((1 + \mathfrak{R})^2 n \zeta_2 + \varrho_3^2 n) (C_1^2 + m_3 \beta_1^2) = \mathfrak{R}_0^2. \quad (153)$$

**Proof of Lemma [14](#)**

This time we use Equation (B29) in Lemma B1 and Lemma B2:

$$\begin{aligned}
& \mathbb{E}_\Sigma \left( \mathbb{E}_{V^\rho, W^\rho} \left[ \frac{1}{\sqrt{m_2}} a^T D_{V^{(0)}+V^\rho+V', x_i} V_\Sigma^* \phi^{(2)}(x_i) \right] \right)^2 \\
& \leq \mathbb{E}_\Sigma \left( \mathbb{E}_{V^\rho, W^\rho} \frac{1}{\sqrt{m_2}} \sum_j |V_{\Sigma_j}^* \phi^{(2)}(x_i)| \right)^2 \\
& \leq \frac{1}{m_2} \mathbb{E}_{\Sigma, V^\rho, W^\rho} \left( \sum_j |V_{\Sigma_j}^* \phi^{(2)}(x_i)| \right)^2 \\
& \leq \mathbb{E}_{V^\rho, W^\rho} E_\Sigma \sum_j |V_{\Sigma_j}^* \phi^{(2)}(x_i)|^2 \\
& \lesssim E_{V^\rho, W^\rho} E_\Sigma \sum_j |(V_{\Sigma_j}^* - V_j^* \Sigma) \phi^{(2)}(x_i)|^2 + \sum_j |V_j^* \Sigma \phi^{(2)}(x_i)|^2 \\
& \lesssim \mathbb{E}_{V^\rho, W^\rho} E_\Sigma \sum_j \|V_{\Sigma_j}^* - V_j^* \Sigma\|^2 \|\phi^{(2)}(x_i)\|^2 + \sum_j \|V_j^*\|_\infty^2 \|\phi^{(2)}(x_i)\|_2^2 \\
& \lesssim ((1 + \Re)^2 n \zeta_2 \xi^2 + \varrho_3^2 \xi^2 n) \mathbb{E}_{V^\rho, W^\rho} \|\phi^{(2)}(x_i)\|_2^2.
\end{aligned}$$

Now according to Lemma B3, we have

$$\mathbb{E}_{V^\rho, W^\rho} \|\phi^{(2)}(x_i)\|^2 \lesssim C_1^2 + m_3 \beta_1^2,$$

which completes the proof.

**Lemma 15** We get an additional term  $\phi^{*'}(x_i)$  as a result of smoothing which we define as

$$\phi^{*'}(x_i) = \frac{1}{\sqrt{m_1}} W^s D_{W^{(0)}+W'+W^\rho, x_i} W_\Sigma^* x_i - \phi^*_\Sigma(x_i). \quad (154)$$

Then

$$\mathbb{P}(\phi^{*'}(x_i) \neq 0) \leq m_1 \exp\{-c_2^2/(8\beta_1^2)\}.$$

Moreover, we have the following inequality almost surely (over the randomness of  $W^\rho$ ):

$$\|\phi^{*'}(x_i)\|_\infty \lesssim \sqrt{\sum_k \|\mathcal{V}_k\|_{H^\infty}^2}.$$

### Proof of Lemma 15

According to Lemma B1, for  $j \notin P$ , for every  $i \in [n]$  we have

$$|(W_j^{(0)} + W_j')x_i| \geq c_2/2\sqrt{m_1}.$$

Now note that as long as the sign patterns for  $j \notin P$  does not change,  $\phi^{*'}(x_i)$  will be zero. Therefore by union bound

$$\begin{aligned}
\mathbb{P}(\phi^{*'}(x_i) \neq 0) & \leq \sum_{j=1}^{m_1} \mathbb{P}(\text{sign change in } j) \leq m_1 \mathbb{P}(|(W_j^{(0)} + W_j')x_i| \leq |W_j^\rho x_i|) \\
& \leq m_1 \mathbb{P}(|W_j^\rho x_i| \geq c_2/(2\sqrt{m_1})).
\end{aligned}$$

But  $(W_j^\rho)x_i$  is Gaussian with variance  $\beta_1^2/m_1$ . Hence

$$LHS \lesssim m_1 \exp\{-c_2^2/(8\beta_1^2)\},$$

which proves the first part. For the second part, according to Equation (B6) in Lemma B1, for every  $k \in [m_3]$ :

$$|\phi_k^{*'}(x_i)| \leq \left| \frac{1}{\sqrt{m_1}} W_k^s D_{W^{(0)}+W'+W^\rho, x_i} W_k^* x_i \right| + |\phi_k^*(x_i)| \quad (155)$$

$$\leq 2/\sqrt{m_1} \sum_j \|W_j^*\| \leq 2\|W^*\|_F \lesssim \sqrt{\sum_k \|\mathcal{V}_k\|_{H^\infty}^2}, \quad (156)$$

which implies the second part.

**Lemma 16** *Fourth Extra term:*

$$\begin{aligned} & \mathbb{E}_{W^\rho, V^\rho} \left[ \frac{1}{\sqrt{m_2}} a^T D_{V^{(0)}+V^\rho+V'}(V^{(0)} + V^\rho + (1-\eta/2)V') \phi^{*'}(x_i) \right] \\ & \lesssim (\kappa_2 \sqrt{m_2 m_3} + C_2 + \sqrt{m_3} \beta_2) m_1 \exp \{-c_2^2/(8\beta_1^2)\} \sqrt{\sum_k \|\mathcal{V}_k\|_{H^\infty}^2} := \mathfrak{R}_{11}. \end{aligned}$$

**Proof of Lemma 16**

Note that with high probability over the randomness of  $V^{(0)}$ , we have  $\|V^{(0)}\|_F \lesssim \sqrt{m_2 m_3} \kappa_2$ . Now according to Lemma 15 and using the fact that  $\|V'\|_F \leq C_2$ :

$$\begin{aligned} & \leq \mathbb{E}_{W^\rho, V^\rho} \frac{1}{\sqrt{m_2}} \|a\| \|V^{(0)} + V^\rho + (1-\eta)V'\|_2 \|\phi^{*'}(x_i)\| \\ & = \mathbb{E}_{W^\rho, V^\rho} \|V^{(0)} + V^\rho + (1-\eta/2)V'\|_F \|\phi^{*'}(x_i)\| \\ & \leq \sqrt{\mathbb{E}_{V^\rho} (\|V^{(0)} + (1-\eta/2)V'\|_F^2 + \|V^\rho\|_F^2)} m_1 \exp \{-c_2^2/(8\beta_1^2)\} \sqrt{\sum_k \|\mathcal{V}_k\|_{H^\infty}^2} \\ & \lesssim \sqrt{\|V^{(0)}\|_F^2 + \|V'\|_F^2 + m_3 \beta_2^2} m_1 \exp \{-c_2^2/(8\beta_1^2)\} \sqrt{\sum_k \|\mathcal{V}_k\|_{H^\infty}^2} \\ & \lesssim (\kappa_2 \sqrt{m_2 m_3} + C_2 + \sqrt{m_3} \beta_2) m_1 \exp \{-c_2^2/(8\beta_1^2)\} \sqrt{\sum_k \|\mathcal{V}_k\|_{H^\infty}^2}. \end{aligned}$$

**Lemma 17** *Fifth extra term:*

$$\left| \mathbb{E}_{W^\rho, V^\rho} \left[ \frac{1}{\sqrt{m_2}} a^T D_{V^{(0)}+V^\rho+V', x_i} V_\Sigma^* \phi^{*'}(x_i) \right] \right| \lesssim \sqrt{\zeta_2} m_1 \exp \{-c_2^2/(8\beta_1^2)\} \sqrt{\sum_k \|\mathcal{V}_k\|_{H^\infty}^2} := \mathfrak{R}_{10}.$$

**Proof of Lemma 17**

Similar to the previous Lemma, the inner expectation can be bounded as:

$$\leq \mathbb{E} \frac{1}{\sqrt{m_2}} \|a\| \|V_\Sigma^*\|_F \|\phi^{*'}(x_i)\| \leq \mathbb{E} \|V^*\|_F \|\phi^{*'}(x_i)\| \lesssim \sqrt{\zeta_2} m_1 \exp \{-c_2^2/(8\beta_1^2)\} \sqrt{\sum_k \|\mathcal{V}_k\|_{H^\infty}^2}.$$

**Lemma 18** *We have another extra term as a product of the movement  $-\eta/2W'$  in the first layer:*

$$\phi^{(2)'}(x_i) = \frac{2}{\eta \sqrt{m_1}} [W^s D_{W^{(0)}+W^\rho+W'}(W^{(0)}+W^\rho+(1-\eta/2)W') x_i - \phi^{(0)}(x_i) - (1-\eta/2)\phi^{(2)}(x_i)].$$

Then

$$\begin{aligned} & \left| \mathbb{E}_{W^\rho, V^\rho} \frac{1}{\sqrt{m_2}} a^T D_{V^{(0)}+V^\rho+V', x_i} V_\Sigma^* \phi^{(2)'}(x_i) \right| \\ & \lesssim \sqrt{\zeta_2 m_3 \left( \frac{\beta_1^2}{m_1} + \frac{c_2 C_1^2}{\sqrt{m_1} \kappa_1} + m_1 \exp \{-c_2^2/(8\beta_1^2)\} C_1^2 \right)} := \mathfrak{R}_5. \end{aligned} \quad (157)$$

$$\begin{aligned} & \mathbb{E}_{W^\rho, V^\rho} \left| \frac{1}{\sqrt{m_2}} a^T D_{V^{(0)}+V^\rho+V', x_i} (V^{(0)} + V^\rho + (1-\eta/2)V') \phi^{(2)'}(x_i) \right| \\ & \lesssim (\kappa_2 \sqrt{m_2 m_3} + \sqrt{m_3} \beta_2 + C_2) \mathfrak{R}_5. \end{aligned} \quad (158)$$

**Proof of Lemma 18**

First we prove the following approximation argument (for all  $k \in [m_3]$ ):

$$\mathbb{E}_{W^\rho} |\phi^{(2)'}(x_i)_k|^2 \leq \frac{\beta_1^2}{m_1} + \frac{c_2 C_1^2}{\sqrt{m_1} \kappa_1} + m_1 \exp \{-c_2^2/(8\beta_1^2)\} C_1^2. \quad (159)$$

We have

$$\begin{aligned} LHS &= \mathbb{E}_{W^\rho} \left| \frac{1}{\sqrt{m_1}} W_k^s D_{W^{(0)}+W^\rho+W',x_i} (W^{(0)} + W^\rho)x_i - \frac{1}{\sqrt{m_1}} W_k^s D_{W^{(0)},x_i} W^{(0)}x_i \right|^2 \\ &\lesssim \mathbb{E}_{W^\rho} \left| \frac{1}{\sqrt{m_1}} W_k^s D_{W^{(0)}+W^\rho,x_i} (W^{(0)} + W^\rho)x_i - \frac{1}{\sqrt{m_1}} W_k^s D_{W^{(0)},x_i} W^{(0)}x_i \right|^2 \\ &+ \mathbb{E}_{W^\rho} \left| \frac{1}{\sqrt{m_1}} W_k^s D_{W^{(0)}+W^\rho,x_i} (W^{(0)} + W^\rho)x_i - \frac{1}{\sqrt{m_1}} W_k^s D_{W^{(0)}+W^\rho+W',x_i} (W^{(0)} + W^\rho)x_i \right|^2 \end{aligned}$$

By the independence of  $W_j^\rho$ 's, the first term can be upper bounded as

$$\begin{aligned} &= \mathbb{E}_{W^\rho} \frac{1}{m_1} \sum_{j=1}^{m_1} \left( (W_j^{(0)} + W_j^\rho)x_i \mathbb{1}\{(W_j^{(0)} + W_j^\rho)x_i \geq 0\} - W_j^{(0)}x_i \mathbb{1}\{W_j^{(0)}x_i \geq 0\} \right)^2 = \\ &\leq \frac{1}{m_1} \sum_{j=1}^{m_1} \mathbb{E}_{W^\rho} |W_j^\rho x_i|^2 = \frac{1}{m_1} \sum \frac{\beta_1^2}{m_1} = \frac{\beta_1^2}{m_1}. \end{aligned}$$

For the second term, note that for every  $j \notin P$ , the  $j$ th entries of  $D_{W^{(0)}+W^\rho,x_i}$  and  $D_{W^{(0)}+W^\rho+W',x_i}$  are different only if  $W_j^\rho$  can make a sign change in the  $j$ th row, i.e.  $|(W_j^{(0)} + W_j^\rho)x_i| \leq |W_j^\rho x_i|$  should happen. We denote this event for every  $j \notin P$  by  $\tilde{E}_j$ . Furthermore, if this happens for some  $j$ , then the value of  $(W^{(0)} + W^\rho)_{j,x_i}$  is upper bounded by  $|W_j^\rho x_i|$ . Now similar to our discussion in Lemma [15](#) and using the result of Lemma [11](#):

$$\begin{aligned} \mathbb{P}(\cup_{j \notin P} \tilde{E}_j) &= \mathbb{P}(\text{sign change in some } j \notin P) \leq \sum_{j \notin P} \mathbb{P}(\text{sign change in } j) \\ &\leq m_1 \mathbb{P}(|(W_j^{(0)} + W_j^\rho)x_i| \leq |W_j^\rho x_i|) \leq m_1 \mathbb{P}(|W_j^\rho x_i| \geq c_2/(2\sqrt{m_1})). \end{aligned}$$

But note that  $(W_j^\rho)x_i$  is Gaussian with variance  $\beta_1^2/m_1$ . Hence

$$LHS \lesssim m_1 \exp\{-c_2^2/(8\beta_1^2)\},$$

So finally we can write

$$\begin{aligned} &\lesssim \frac{\beta_1^2}{m_1} + \mathbb{E}_{W^\rho} \frac{1}{m_1} \left( \sum_{j \in P} |W_j^\rho x_i|^2 + E_{W^\rho} \frac{1}{m_1} (\mathbb{1}\{\cup_{j \notin P} \tilde{E}_j\} \sum_{j \notin P} |W_j^\rho x_i|^2) \right) \\ &\lesssim \frac{\beta_1^2}{m_1} + \frac{|P|}{m_1} \|W^\rho\|^2 + \mathbb{P}(\cup_{j \notin P} \tilde{E}_j) \|W^\rho\|^2 \\ &\leq \frac{\beta_1^2}{m_1} + \frac{c_2 C_1^2}{\sqrt{m_1} \kappa_1} + m_1 \exp\{-c_2^2/(8\beta_1^2)\} C_1^2. \end{aligned}$$

which completes the proof for Equation [159](#). This immediately implies

$$\mathbb{E}_{W^\rho} \|\phi^{(2)'}(x_i)\| \leq \sqrt{\mathbb{E}_{W^\rho} \|\phi^{(2)'}(x_i)\|^2} \leq \sqrt{m_3 \left( \frac{\beta_1^2}{m_1} + \frac{c_2 C_1^2}{\sqrt{m_1} \kappa_1} + m_1 \exp\{-c_2^2/(8\beta_1^2)\} C_1^2 \right)}.$$

Now we first prove Equation [157](#):

$$\begin{aligned} \left| \mathbb{E}_{W^\rho, V^\rho} \left[ \frac{1}{\sqrt{m_2}} a^T D_{V^{(0)}+V^\rho+V',x_i} V_\Sigma^* \phi^{(2)'}(x_i) \right] \right| &\leq \mathbb{E}_{W^\rho} \frac{1}{\sqrt{m_2}} \|a\| \|D_{V^{(0)}+V^\rho+V',x_i} V_\Sigma^* \|_F \|\phi^{(2)'}(x_i)\| \\ &\leq \|V_\Sigma^*\|_F \mathbb{E}_{W^\rho} \|\phi^{(2)'}(x_i)\| \leq \|V^*\|_F \mathbb{E}_{W^\rho} \|\phi^{(2)'}(x_i)\| \\ &\lesssim \sqrt{\zeta_2 m_3 \left( \frac{\beta_1^2}{m_1} + \frac{c_2 C_1^2}{\sqrt{m_1} \kappa_1} + m_1 \exp\{-c_2^2/(8\beta_1^2)\} C_1^2 \right)}. \end{aligned}$$

To prove Equation (158):

$$\begin{aligned}
& \mathbb{E}_{W^\rho, V^\rho} \left| \frac{1}{\sqrt{m_2}} a^T D_{V^{(0)}+V^\rho+V', x_i} (V^{(0)} + V^\rho + (1 - \eta/2)V') \phi^{(2)'}(x_i) \right| \\
& \lesssim \mathbb{E}_{W^\rho, V^\rho} \frac{1}{\sqrt{m_2}} \|a\| \|D_{V^{(0)}+V^\rho+(1-\eta/2)V', x_i} (V^{(0)} + V^\rho + (1 - \eta/2)V')\|_F \|\phi^{(2)'}(x_i)\| \\
& \lesssim \mathbb{E}_{W^\rho, V^\rho} \frac{1}{\sqrt{m_2}} \|a\| \|V^{(0)} + V^\rho + (1 - \eta/2)V'\|_F \|\phi^{(2)'}(x_i)\| \\
& \lesssim \frac{1}{\sqrt{m_2}} \|a\| \sqrt{\mathbb{E}_{V^\rho} (\|V^{(0)}\|^2 + \|V^\rho\|_F^2 + \|(1 - \eta/2)V'\|_F^2)} \mathbb{E}_{W^\rho} \|\phi^{(2)'}(x_i)\| \\
& \lesssim \sqrt{(\kappa_2^2 m_2 m_3 + m_3 \beta_2^2 + C_2^2)} \mathfrak{R}_5 \lesssim (\kappa_2 \sqrt{m_2 m_3} + \sqrt{m_3} \beta_2 + C_2) \mathfrak{R}_5.
\end{aligned}$$

**Lemma 19** *Closeness condition:*

$$\mathbb{E}_\Sigma \left| \mathbb{E}_{W^\rho, V^\rho} \left[ \frac{1}{\sqrt{m_2}} a^T D_{V^{(0)}+V^\rho+V', x_i} V_\Sigma^* \phi_\Sigma^*(x_i) \right] - f^*(x_i) \right| \lesssim \mathfrak{R}_9,$$

where

$$\mathfrak{R}_9 := \varrho_3 \xi \sqrt{n} (1 + \mathfrak{R}) \sqrt{\sum_k \|\mathcal{V}_k\|_{H^\infty}^2} + \mathfrak{R}_3 \tag{160}$$

$$+ m_2 \left( \exp \{ -(m_2 \kappa_2^2 C_2^4)^{1/3} / (2\beta_2^2) \} + m_1 \exp \{ -C_1^2 / (8m_3 \beta_1^2) \} \right) \sqrt{\zeta_2} (1 + \mathfrak{R}) \sqrt{\sum_k \|\mathcal{V}_k\|_{H^\infty}^2}. \tag{161}$$

**Proof of Lemma 19**

Note that by Corollary 5.1 and according to the proof of Equation 130 in Lemma 10, if for every  $j \notin \tilde{P}$  we don't have a sign change in  $D_{V^{(0)}+V^\rho+V', x_i} V^* \phi^*(x_i)$ , then get

$$\left| \frac{1}{\sqrt{m_2}} a^T D_{V^{(0)}+V^\rho+V', x_i} V^* \phi^*(x_i) - f^*(x_i) \right| \leq \mathfrak{R}_3.$$

Also, note that we need the event  $E^c$  (defined in Lemma 3.3) to happen in order to be able to use Corollary 5.1. Hence, given a  $W^\rho$  for which  $E^c$  happens, we upper bound the probability of sign change with respect to the randomness of  $V^\rho$ . We define the following event with respect to the randomness of  $V^\rho$  when conditioned on a  $W^\rho$  for which  $E^c$  happens ( $P_i$ 's are defined in Lemma 10):

$$SC := \{ \exists j \notin P_i \text{ s.t. } |V_j^\rho x'_i| \gtrsim \left(\frac{\kappa_2}{m_2}\right)^{1/3} C_2^{2/3} \|x'_i\| \}.$$

Now from the result in Corollary 5.1 we have  $\leq \mathbb{1}\{\text{sign change in } j \notin P_i\} \leq \mathbb{1}\{SC\}$ . Therefore,

$$\begin{aligned}
E \mathbb{1}\{\text{sign change}\} & \leq \mathbb{1}\{SC\} \leq \sum_{j \notin P_i} \mathbb{P}(|V_j^\rho x'_i| \gtrsim \left(\frac{\kappa_2}{m_2}\right)^{1/3} C_2^{2/3} \|x'_i\|) \\
& \leq m_2 \mathbb{P}(|V_j^\rho x'_i| \gtrsim \left(\frac{\kappa_2}{m_2}\right)^{1/3} C_2^{2/3} \|x'_i\|).
\end{aligned}$$

But note that  $(V_j^\rho x'_i)$  is Gaussian with variance  $\beta_2^2 \|x'_i\|^2 / m_2$ . Hence

$$LHS \lesssim m_2 \exp \{ -(m_2 \kappa_2^2 C_2^4)^{1/3} / (2\beta_2^2) \}. \tag{162}$$

Now let  $D$  be a sign matrix random variable such that if  $E^c$  and  $SC^c$  both happens, then it is equal to the valid sign matrix  $D_{V^{(0)}+V^\rho+V', x_i}$ , and otherwise it is equal to an arbitrary valid sign matrix in the case when both  $E^c$  and  $SC^c$  happen. Now using Equation (117) we have with high probability

over the initialization:

$$\begin{aligned}
& \mathbb{E}_\Sigma \left| \mathbb{E}_{W^\rho, V^\rho} \left[ \frac{1}{\sqrt{m_2}} a^T D_{V^{(0)}+V^\rho+V', x_i} V_\Sigma^* \phi_\Sigma^*(x_i) \right] - f^*(x_i) \right| \\
& \leq \mathbb{E}_\Sigma \left| \mathbb{E}_{W^\rho, V^\rho} \left[ \frac{1}{\sqrt{m_2}} a^T D_{V^{(0)}+V^\rho+V', x_i} (V_\Sigma^* - V^* \Sigma) \phi_\Sigma^*(x_i) \right] \right| \\
& + \mathbb{E}_\Sigma \left| \mathbb{E}_{W^\rho, V^\rho} \left[ \frac{1}{\sqrt{m_2}} a^T D_{V^{(0)}+V^\rho+V', x_i} V^* \Sigma \phi_\Sigma^*(x_i) \right] - f^*(x_i) \right| \\
& \leq \mathbb{E}_{W^\rho, V^\rho} \mathbb{E}_\Sigma \left| \frac{1}{\sqrt{m_2}} a^T D_{V^{(0)}+V^\rho+V', x_i} (V_\Sigma^* - V^* \Sigma) \phi_\Sigma^*(x_i) \right| \\
& + \mathbb{E}_\Sigma \left| \mathbb{E}_{W^\rho, V^\rho} \left[ \frac{1}{\sqrt{m_2}} a^T D_{V^{(0)}+V^\rho+V', x_i} V^* \phi^*(x_i) \right] - f^*(x_i) \right| \\
& \leq \mathbb{E}_{W^\rho, V^\rho} \mathbb{E}_\Sigma \frac{1}{m_2} \sum_j \|V_{\Sigma^* j}^* - V_j^* \Sigma\| \|\phi_\Sigma^*(x_i)\| \\
& + \mathbb{E}_\Sigma \left| \mathbb{E}_{W^\rho, V^\rho} \left[ \left( \frac{1}{\sqrt{m_2}} a^T D V^* \phi^*(x_i) - f^*(x_i) \right) \right. \right. \\
& \left. \left. + \mathbb{1}\{SC \cup E\} \left( \frac{1}{\sqrt{m_2}} a^T D_{V^{(0)}+V', x_i} V^* \phi^*(x_i) - D \right) \right] \right| \\
& \leq \mathbb{E}_{W^\rho, V^\rho} \mathbb{E}_\Sigma \frac{1}{m_2} \sum_j \|V_{\Sigma^* j}^* - V_j^* \Sigma\| \|\phi_\Sigma^*(x_i)\| \\
& + \mathbb{E}_\Sigma \left| \mathbb{E}_{W^\rho, V^\rho} \left[ \left( \frac{1}{\sqrt{m_2}} a^T D V^* \phi^*(x_i) - f^*(x_i) \right) \right. \right. \\
& \left. \left. + \mathbb{E}_\Sigma \mathbb{E}_{W^\rho, V^\rho} \left[ \mathbb{1}\{SC \cup E\} \left( \frac{1}{\sqrt{m_2}} a^T D_{V^{(0)}+V', x_i} V^* \phi^*(x_i) - \frac{1}{\sqrt{m_2}} a^T D V^* \phi^*(x_i) \right) \right] \right] \right| \\
& \leq \mathbb{E}_{W^\rho, V^\rho} \frac{1}{\sqrt{m_2}} \sum_j \sqrt{\mathbb{E}_\Sigma \|V_{\Sigma^* j}^* - V_j^* \Sigma\|^2} \|\phi^*(x_i)\| \\
& + \mathfrak{R}_3 + 2\mathbb{P}(SC \cup E) \max_{D'} \left| \frac{1}{\sqrt{m_2}} a^T D' V^* \phi^*(x_i) \right|.
\end{aligned}$$

Now note that for any sign matrix  $D'$ , we have the following bound:

$$\left| \frac{1}{\sqrt{m_2}} a^T D' V^* \phi^*(x_i) \right| \leq \frac{1}{\sqrt{m_2}} \|a\| \|V^*\|_F \|\phi^*(x_i)\| \lesssim \sqrt{\zeta_2} (1 + \mathfrak{R}) \sqrt{\sum_k \|\mathcal{V}_k\|_{H^\infty}^2}.$$

Also, applying a union bound and using Lemmas [B3](#)

$$\begin{aligned}
\mathbb{P}(SC \cup E) & \leq \mathbb{P}(SC) + \mathbb{P}(E) \\
& \lesssim \exp\{-(m_2 \kappa_2^2 C_2^4)^{1/3} / (2\beta_2^2)\} + m_1 \exp\{-C_1^2 / (8m_3 \beta_1^2)\}.
\end{aligned}$$

Hence, also applying Lemma [A0](#), we further write

$$\begin{aligned}
LHS & \lesssim \varrho_3 \xi \sqrt{n} (1 + \mathfrak{R}) \sqrt{\sum_k \|\mathcal{V}_k\|_{H^\infty}^2} + \mathfrak{R}_3 \\
& + m_2 \left( \exp\{-(m_2 \kappa_2^2 C_2^4)^{1/3} / (2\beta_2^2)\} + m_1 \exp\{-C_1^2 / (8m_3 \beta_1^2)\} \right) \sqrt{\zeta_2} (1 + \mathfrak{R}) \sqrt{\sum_k \|\mathcal{V}_k\|_{H^\infty}^2}.
\end{aligned}$$

**Lemma 20** Suppose we have  $m_3 \kappa_1^2 \geq C_1^2$ . Then, for the following basic term we have:

$$\begin{aligned}
& \mathbb{E}_{W^\rho, V^\rho} \left[ \frac{1}{\sqrt{m_2}} a^T D_{V^{(0)}+V^\rho+V', x_i} (V^{(0)} + V^\rho + (1 - \eta/2)V') (\phi^{(0)}(x_i) + (1 - \eta/2)\phi^{(2)}(x_i)) \right] \\
& \lesssim (1 - \eta) \mathbb{E}_{W^\rho, V^\rho} \left[ \frac{1}{\sqrt{m_2}} a^T D_{V^{(0)}+V^\rho+V', x_i} (V^{(0)} + V^\rho + V') (\phi^{(0)}(x_i) + \phi^{(2)}(x_i)) \right]
\end{aligned}$$

$$\pm \eta \left( \mathfrak{R}'_6 + \mathfrak{R}_4 + (\sqrt{m_3} \kappa_2 + \beta_2)(C_1 + \sqrt{m_3} \beta_1) \right),$$

where

$$\mathfrak{R}_4 := C_2(C_1 + \sqrt{m_3} \beta_1) m_2 \exp\{-C_2^{4/3}(\sqrt{m_2} \kappa_2)^{2/3}/8\beta_2^2\} + \frac{C_2^{1/3}}{(\sqrt{m_2} \kappa_2)^{1/3}} C_2(C_1 + \sqrt{m_3} \beta_1),$$

$$\mathfrak{R}'_6 := m_1 \exp\{-C_1^2/(8m_3\beta_1^2)\} \sqrt{m_3} \kappa_1 (\sqrt{m_2} + \beta_2) + \mathfrak{R}_6,$$

and  $\mathfrak{R}_6$  is defined in Lemma 22.

### Proof of Lemma 20

First, note that by orthogonality of  $\phi^{(0)}(x_i)$  to the rows of  $V'$ :

$$\begin{aligned} & LHS - \eta/2 \mathbb{E}_{W^\rho, V^\rho} \left[ \frac{1}{\sqrt{m_2}} a^T D_{V^{(0)}+V^\rho+V', x_i} (V^{(0)} + V^\rho) \phi^{(0)}(x_i) \right] \\ &= LHS - \eta/2 \mathbb{E}_{W^\rho, V^\rho} \left[ \frac{1}{\sqrt{m_2}} a^T D_{V^{(0)}+V^\rho+V', x_i} (V^{(0)} + V^\rho + (1 - \eta/2)V') \phi^{(0)}(x_i) \right] \\ &= (1 - \eta/2) \mathbb{E}_{W^\rho, V^\rho} \left[ \frac{1}{\sqrt{m_2}} a^T D_{V^{(0)}+V^\rho+V', x_i} (V^{(0)} + V^\rho + (1 - \eta/2)V') (\phi^{(0)}(x_i) + \phi^{(2)}(x_i)) \right] \\ &= (1 - \eta/2)^2 \mathbb{E}_{V^\rho} \left[ \frac{1}{\sqrt{m_2}} a^T D_{V^{(0)}+V^\rho+V', x_i} (V^{(0)} + V^\rho + V') (\phi^{(0)}(x_i) + \phi^{(2)}(x_i)) \right] \\ &+ (1 - \eta/2)(\eta/2) \mathbb{E}_{V^\rho} \left[ \frac{1}{\sqrt{m_2}} a^T D_{V^{(0)}+V^\rho+V', x_i} (V^{(0)} + V^\rho) (\phi^{(0)}(x_i) + \phi^{(2)}(x_i)) \right] \end{aligned} \quad (163)$$

But note that for the second term:

$$\begin{aligned} & \mathbb{E}_{W^\rho, V^\rho} \left[ \frac{1}{\sqrt{m_2}} a^T D_{V^{(0)}+V^\rho+V', x_i} (V^{(0)} + V^\rho) (\phi^{(0)}(x_i) + \phi^{(2)}(x_i)) \right] \\ & \lesssim \mathbb{E}_{W^\rho, V^\rho} \left[ \frac{1}{\sqrt{m_2}} a^T D_{V^{(0)}+V^\rho, x_i} (V^{(0)} + V^\rho) (\phi^{(0)}(x_i) + \phi^{(2)}(x_i)) \right] \\ & \pm \frac{1}{\sqrt{m_2}} \sum_{j: \text{sign change}} |V'_j (\phi^{(0)}(x_i) + \phi^{(2)}(x_i))| \\ &= \mathbb{E}_{V^\rho} \left[ \frac{1}{\sqrt{m_2}} a^T D_{V^{(0)}+V^\rho, x_i} (V^{(0)} + V^\rho) (\phi^{(0)}(x_i) + \phi^{(2)}(x_i)) \right] \\ & \pm \frac{1}{\sqrt{m_2}} \sum_{j: \text{sign change}} |V'_j \phi^{(2)}(x_i)|. \end{aligned} \quad (164)$$

Now conditioned on  $x'_i$ , by the result of Lemma 21 we know there exists a set of indices  $O \subseteq [m_2]$ ,

s.t.  $|O| \leq \frac{C_2^{2/3}}{(\sqrt{m_2} \kappa_2)^{2/3}} m_2$  and for  $j \notin O$  we have

$$|V_j^{(0)} x'_i| \geq \frac{C_2^{2/3} (\sqrt{m_2} \kappa_2)^{1/3}}{\sqrt{m_2}} \|x'_i\|$$

and

$$|V'_j x'_i| \leq \frac{C_2^{2/3} (\sqrt{m_2} \kappa_2)^{1/3}}{2\sqrt{m_2}} \|x'_i\|.$$

Now for  $j \in [m_2]$ , define the event

$$R_j = \{|W_j^\rho x'_i| \geq \frac{C_2^{2/3} (\sqrt{m_2} \kappa_2)^{1/3}}{2\sqrt{m_2}} \|x'_i\|\},$$

and  $R = \cup_j R_j$ . First, note that using Gaussian tail bound,  $R$  is a rare event:

$$\mathbb{P}(R) \leq \sum_j \mathbb{P}(R_j) \leq m_2 \exp\{-C_2^{4/3} (\sqrt{m_2} \kappa_2)^{2/3}/8\beta_2^2\}.$$

Now for  $j \notin O$  and under  $R^c$ , clearly we have that the signs of  $(V_j^{(0)} + V_j^\rho)x'_i$  and  $(V_j^{(0)} + V_j^\rho + V_j')$  are the same. Therefore, applying Lemma B3, we can argue under  $R^c$ :

$$\begin{aligned} \mathbb{E}_{W^\rho, V^\rho} \frac{1}{\sqrt{m_2}} \sum_{j: \text{sign change}} |V_j' \phi^{(2)}(x_i)| &\leq \mathbb{E}_{W^\rho} \frac{1}{\sqrt{m_2}} \sum_{j \in O} |V_j' \phi^{(2)}(x_i)| \\ &\leq \sqrt{\frac{|O|}{m_2}} \|V'\| \mathbb{E}_{W^\rho} \|\phi^{(2)}(x_i)\| \leq \frac{C_2^{1/3}}{(\sqrt{m_2} \kappa_2)^{1/3}} C_2 (C_1 + \sqrt{m_3} \beta_1). \end{aligned}$$

Hence, overall, using Cauchy-Schwartz

$$\begin{aligned} \mathbb{E}_{W^\rho, V^\rho} \frac{1}{\sqrt{m_2}} \sum_{j: \text{sign change}} |V_j' \phi^{(2)}(x_i)| &\leq \|V'\| \mathbb{E}_{W^\rho} \|\phi^{(2)}(x_i)\| \mathbb{P}(R) + \frac{C_2^{1/3}}{(\sqrt{m_2} \kappa_2)^{1/3}} C_2 (C_1 + \sqrt{m_3} \beta_1) \\ &\leq C_2 (C_1 + \sqrt{m_3} \beta_1) m_2 \exp\{-C_2^{4/3} (\sqrt{m_2} \kappa_2)^{2/3} / 8 \beta_2^2\} + \frac{C_2^{1/3}}{(\sqrt{m_2} \kappa_2)^{1/3}} C_2 (C_1 + \sqrt{m_3} \beta_1) := \mathfrak{R}_4. \end{aligned} \quad (165)$$

On the other hand, using Lemma B0, we have with high probability over the randomness of initialization

$$\frac{1}{\sqrt{m_2}} a^T D_{V^{(0)}, x_i} V^{(0)} \phi^{(2)}(x_i) \leq \sqrt{m_3} \kappa_2 \|\phi^{(2)}(x_i)\|.$$

Hence:

$$\begin{aligned} &\mathbb{E}_{V^\rho} \left[ \frac{1}{\sqrt{m_2}} a^T D_{V^{(0)} + V^\rho, x_i} (V^{(0)} + V^\rho) \phi^{(2)}(x_i) \right] \\ &\leq \mathbb{E}_{W^\rho, V^\rho} \left[ \frac{1}{\sqrt{m_2}} a^T D_{V^{(0)}, x_i} V^{(0)} \phi^{(2)}(x_i) + \frac{1}{\sqrt{m_2}} \sum_j |V_j^\rho \phi^{(2)}(x_i)| \right] \\ &\leq \mathbb{E}_{W^\rho} \frac{1}{\sqrt{m_2}} a^T D_{V^{(0)}, x_i} V^{(0)} \phi^{(2)}(x_i) + \beta_2 \mathbb{E}_{W^\rho} \|\phi^{(2)}(x_i)\| \\ &\lesssim (\sqrt{m_3} \kappa_2 + \beta_2) \mathbb{E}_{W^\rho} \|\phi^{(2)}(x_i)\| \\ &\lesssim (\sqrt{m_3} \kappa_2 + \beta_2) (C_1 + \sqrt{m_3} \beta_1). \end{aligned} \quad (166)$$

Combining Equations (165) and (166) into Equation (164):

$$\left| \mathbb{E}_{V^\rho} \left[ \frac{1}{\sqrt{m_2}} a^T D_{V^{(0)} + V^\rho + V', x_i} (V^{(0)} + V^\rho) \phi^{(2)}(x_i) \right] \right| \lesssim \mathfrak{R}_4 + (\sqrt{m_3} \kappa_2 + \beta_2) (C_1 + \sqrt{m_3} \beta_1). \quad (167)$$

Moreover, for the first term in (163), using Equation (166) and Lemmas B3 and Lemma B0 we have

$$\begin{aligned} &|\mathbb{E}_{V^\rho} \left[ \frac{1}{\sqrt{m_2}} a^T D_{V^{(0)} + V^\rho + V', x_i} (V^{(0)} + V^\rho + V') \phi^{(2)}(x_i) \right]| \\ &\lesssim |\mathbb{E}_{W^\rho, V^\rho} \frac{1}{\sqrt{m_2}} a^T D_{V^{(0)}, x_i} V^{(0)} \phi^{(2)}(x_i)| + \mathbb{E}_{W^\rho, V^\rho} \frac{1}{\sqrt{m_2}} \sum_j |V_j^\rho \phi^{(2)}(x_i)| + \frac{1}{\sqrt{m_2}} \sum_j |V_j' \phi^{(2)}(x_i)| \\ &\lesssim \kappa_2 \sqrt{m_3} (C_1 + \beta_1 \sqrt{m_3}) + (C_2 + \beta_2) \mathbb{E}_{W^\rho} \|\phi^{(2)}(x_i)\| \\ &\lesssim \kappa_2 \sqrt{m_3} (C_1 + \beta_1 \sqrt{m_3}) + (C_2 + \beta_2) (C_1 + \sqrt{m_3} \beta_1). \end{aligned} \quad (168)$$

Substituting Equations (167) and (168) into Equation (163), we finally get

$$\begin{aligned}
LHS &= \eta/2 \mathbb{E}_{W^\rho, V^\rho} \left[ \frac{1}{\sqrt{m_2}} a^T D_{V^{(0)}+V^\rho+V', x_i} (V^{(0)} + V^\rho) \phi^{(0)}(x_i) \right] \\
&\lesssim (1 - \eta/2)^2 \mathbb{E}_{V^\rho} \left[ \frac{1}{\sqrt{m_2}} a^T D_{V^{(0)}+V^\rho+V', x_i} (V^{(0)} + V^\rho + V') \phi^{(2)}(x_i) \right] \\
&\quad \pm \frac{\eta}{2} \left( \mathfrak{R}_4 + (\sqrt{m_3} \kappa_2 + \beta_2) (C_1 + \sqrt{m_3} \beta_1) \right) \\
&\lesssim (1 - \eta) \mathbb{E}_{V^\rho} \left[ \frac{1}{\sqrt{m_2}} a^T D_{V^{(0)}+V^\rho+V', x_i} (V^{(0)} + V^\rho + V') \phi^{(2)}(x_i) \right] \\
&\quad \pm \frac{\eta^2}{4} \left| \mathbb{E}_{V^\rho} \left[ \frac{1}{\sqrt{m_2}} a^T D_{V^{(0)}+V^\rho+V', x_i} (V^{(0)} + V^\rho + V') \phi^{(2)}(x_i) \right] \right| \\
&\quad \pm \frac{\eta}{2} \left( \mathfrak{R}_4 + (\sqrt{m_3} \kappa_2 + \beta_2) (C_1 + \sqrt{m_3} \beta_1) \right) \\
&\lesssim (1 - \eta) \mathbb{E}_{V^\rho} \left[ \frac{1}{\sqrt{m_2}} a^T D_{V^{(0)}+V^\rho+V', x_i} (V^{(0)} + V^\rho + V') \phi^{(2)}(x_i) \right] \\
&\quad \pm \eta^2 (\kappa_2 \sqrt{m_3} (C_1 + \beta_1 \sqrt{m_3}) + (C_2 + \beta_2) (C_1 + \sqrt{m_3} \beta_1)) \\
&\quad \pm \eta \left( \mathfrak{R}_4 + (\sqrt{m_3} \kappa_2 + \beta_2) (C_1 + \sqrt{m_3} \beta_1) \right). \tag{169}
\end{aligned}$$

Now by picking  $\eta$  small enough so that the second term is dominated by the third term we get:

$$LHS - \eta/2 \mathbb{E}_{W^\rho, V^\rho} \left[ \frac{1}{\sqrt{m_2}} a^T D_{V^{(0)}+V^\rho+V', x_i} (V^{(0)} + V^\rho) \phi^{(0)}(x_i) \right] \tag{170}$$

$$\lesssim (1 - \eta) \mathbb{E}_{V^\rho} \left[ \frac{1}{\sqrt{m_2}} a^T D_{V^{(0)}+V^\rho+V', x_i} (V^{(0)} + V^\rho + V') \phi^{(2)}(x_i) \right] \tag{171}$$

$$\pm \eta \left( \mathfrak{R}_4 + (\sqrt{m_3} \kappa_2 + \beta_2) (C_1 + \sqrt{m_3} \beta_1) \right). \tag{172}$$

Now we aim to bound the term  $\mathbb{E}_{W^\rho, V^\rho} \left[ \frac{1}{\sqrt{m_2}} a^T D_{V^{(0)}+V^\rho+V', x_i} (V^{(0)} + V^\rho) \phi^{(0)}(x_i) \right]$ . First assume that we are in the event  $E^c$  defined in Lemma 53, i.e. we have  $\|\phi^{(2)}(x_i)\| \lesssim C_1$ . Conditioned on such  $W^\rho$ , we now work with the randomness of the initialization and  $V^\rho$ . Note that the random matrix  $V^{(0)} + V^\rho$  jointly over the randomness of  $V^\rho$  and the initialization is also Gaussian, and its variance is

$$\kappa_2^2 \leq \kappa_2^2 + \frac{\beta_2^2}{m_2} \leq 2\kappa_2^2, \tag{173}$$

where the inequality follows from the fact that  $\kappa_2 \geq \frac{1}{\sqrt{m_2}}$  and  $\beta_2 \leq 1$ . Therefore, applying Lemma 22 for the random matrix  $V^{(0)}$  in the Lemma as  $V^{(0)} + V^\rho$  here, the bound does not change up to constants because of the inequality (173). Hence, with high probability, lets say with prob.  $1 - \delta_1$  this time over both the randomness of initialization and  $V^\rho$ :

$$\mathcal{L} = \left| \frac{1}{\sqrt{m_2}} a^T D_{V^{(0)}+V^\rho+V', x_i} (V^{(0)} + V^\rho) \phi^{(0)}(x_i) \right| \leq \mathfrak{R}_6 \tag{174}$$

This means that with probability at least  $1 - \sqrt{\delta_1}$  over the random initialization, then we have (174) with prob. at least  $1 - \sqrt{\delta_1}$  over the randomness of  $V^\rho$ . We name the latter high probability statement as  $(\star)$ . Moreover, note that by Lemma 52 and assuming  $m_3 \log(m_2) < m_2$ , we have the following

almost surely bound (also note that  $V_j^\rho \phi^{(0)}(x_i)$  is Gaussian with std  $\frac{\beta_2}{\sqrt{m_2}} \|\phi^{(0)}(x_i)\|$ ):

$$\mathbb{E}_{V^\rho} \left| \frac{1}{\sqrt{m_2}} a^T D_{V^{(0)}+V^\rho+V', x_i} (V^{(0)} + V^\rho) \phi^{(0)}(x_i) \right| \quad (175)$$

$$= \mathbb{E}_{V^\rho} \left| \frac{1}{\sqrt{m_2}} a^T D_{V^{(0)}+V^\rho+V', x_i} (V^{(0)} + V^\rho) \phi^{(0)}(x_i) \right| \quad (176)$$

$$\leq \frac{1}{\sqrt{m_2}} \sum_{j=1}^{m_2} |V_j^{(0)} \phi^{(0)}(x_i)| + \mathbb{E}_{V^\rho} \frac{1}{\sqrt{m_2}} \sum_{j=1}^{m_2} |V_j^\rho \phi^{(0)}(x_i)| \quad (177)$$

$$\lesssim \|\phi^{(0)}(x_i)\| \sup_{\|x'\|=1} \frac{1}{\sqrt{m_2}} \sum_{j=1}^{m_2} |V_j^{(0)} x'| + \frac{1}{\sqrt{m_2}} \sum_{j=1}^{m_2} \frac{\beta_2}{\sqrt{m_2}} \|\phi^{(0)}(x_i)\| \quad (178)$$

$$\lesssim \|\phi^{(0)}(x_i)\| (\sqrt{m_2} + \beta_2). \quad (179)$$

Furthermore, note because each variable  $|V_j^\rho \phi^{(0)}(x_i)|$  is  $\frac{\beta_2}{\sqrt{m_2}} \|\phi^{(0)}(x_i)\|$  subGaussian. Therefore,  $\mathcal{L}$  is subGaussian with parameter  $\|\phi^{(0)}(x_i)\| \beta_2$  with respect to the randomness of  $V^\rho$ . Now the point is that the high probability argument in  $(\star)$  is much stronger than what one can get from the subGaussian inequality with parameter  $\|\phi^{(0)}(x_i)\| \beta_2$  (with the corresponding expectation term  $\|\phi^{(0)}(x_i)\| (\sqrt{m_2} + \beta_2)$ ). However, the disadvantage of  $(\star)$  is that it only works for a fixed  $\delta_1$ . In other words, at least it is not obvious from this argument that why for a fixed  $W^{(0)}$  in a high probability region of the random initialization, whether we can send  $\delta_1$  to zero by growing the constant behind  $\mathfrak{R}_6$  with logarithmic rate  $\log(1/\delta)$ . This makes our job hard for bounding the expectation with respect to  $V^\rho$  if we only wish to rely on  $(\star)$ . Therefore, we combine it with the inequality that we get from the subGaussian parameter that we introduced above. More rigorously, we define the thresholding parameter

$$\begin{aligned} \mathcal{U} &:= \|\phi^{(0)}(x_i)\| (\sqrt{m_2} + \beta_2) + \|\phi^{(0)}(x_i)\| \beta_2 \log(\|\phi^{(0)}(x_i)\| (\sqrt{m_2} + \beta_2) / \mathfrak{R}_6) \\ &= \Theta\left(\|\phi^{(0)}(x_i)\| (\sqrt{m_2} + \beta_2 \log(\|\phi^{(0)}(x_i)\| (\sqrt{m_2} + \beta_2) / \mathfrak{R}_6))\right), \end{aligned}$$

for which we have

$$\mathbb{E}[\mathcal{L} | \mathcal{U} \leq \mathcal{L}] \mathbb{P}(\mathcal{U} \leq \mathcal{L}) \lesssim \mathfrak{R}_6.$$

we divide the range of values for  $\mathcal{L}$  into three parts:

$$\begin{aligned} \mathbb{E}[\mathcal{L}] &= \mathbb{E}[\mathcal{L} | \mathcal{L} \leq \mathfrak{R}_6] \mathbb{P}(\mathcal{L} \leq \mathfrak{R}_6) \\ &\quad + \mathbb{E}[\mathcal{L} | \mathfrak{R}_6 \leq \mathcal{L} \leq \mathcal{U}] \mathbb{P}(\mathfrak{R}_6 \leq \mathcal{L} \leq \mathcal{U}) \\ &\quad + \mathbb{E}[\mathcal{L} | \mathcal{U} \leq \mathcal{L}] \mathbb{P}(\mathcal{U} \leq \mathcal{L}) \\ &\leq E[\mathcal{L} | \mathcal{L} \leq \mathfrak{R}_6] + \mathbb{P}(\mathfrak{R}_6 \leq \mathcal{L} \leq \mathcal{U}) + \mathfrak{R}_6 \\ &\lesssim \mathfrak{R}_6 + \sqrt{\delta_1} \mathcal{U}. \end{aligned}$$

Now by choosing  $\delta_1 \lesssim 1/\mathcal{U}$ , we conclude with high probability over initialization and conditioned on  $W^\rho$ 's such that  $E^c$  happens we have

$$\mathbb{E}_{V^\rho} \left| \frac{1}{\sqrt{m_2}} a^T D_{V^{(0)}+V^\rho+V', x_i} (V^{(0)} + V^\rho) \phi^{(0)}(x_i) \right| = \mathbb{E}[\mathcal{L}] \lesssim \mathfrak{R}_6.$$

Finally, we integrate also with respect to  $W^\rho$ . To control the random variable when  $E$  happens, we use the bound in (179) and the fact that  $E$  is a rare event due to Lemma B3:

$$\begin{aligned} \mathbb{E}_{W^\rho, V^\rho} \left| \frac{1}{\sqrt{m_2}} a^T D_{V^{(0)}+V^\rho+V', x_i} (V^{(0)} + V^\rho) \phi^{(0)}(x_i) \right| &\lesssim \mathbb{P}(E) \|\phi^{(0)}(x_i)\| (\sqrt{m_2} + \beta_2) + \mathbb{P}(E^c) \mathfrak{R}_6 \\ &\leq m_1 \exp\{-C_1^2/(8m_3\beta_1^2)\} \sqrt{m_3} \kappa_1 (\sqrt{m_2} + \beta_2) + \mathfrak{R}_6 := \mathfrak{R}'_6. \end{aligned}$$

Substituting this into (172) the proof is finally complete.

**Lemma 21** *Third cross term: with high probability over initialization, we have*

$$\begin{aligned} & \mathbb{E}_\Sigma \left( \mathbb{E}_{W^\rho, V^\rho} \left[ \frac{1}{\sqrt{m_2}} a^T D_{V^{(0)}+V^\rho+V', x_i} (V^{(0)} + V^\rho + (1-\eta)V') \Sigma \phi^*(x_i) \right] \right)^2 \\ & \lesssim \xi^2 (m_1 \exp\{-C_1^2/(8m_3\beta_1^2)\}) (\kappa_2^2 m_2 m_3 + \beta_1^2 m_3) + \mathfrak{R}_7^2 C_2^2 := \mathfrak{R}_{12}^2. \end{aligned}$$

**Proof of Lemma 21**

Note that the way we defined the matrix  $W^*$  and as a result  $\phi^*(x_i)$  only depends on the randomness of  $W^{(0)}$ , not on  $W'$  or the randomness of  $V^{(0)}$ . Now using Equation (131) and Jensen inequality we can write (for vector  $v$ , the notation  $v^{2\odot}$  is another vector with each entry as the second power of the corresponding entry in  $v$ ):

$$\begin{aligned} & = \mathbb{E}_\Sigma \left( \mathbb{E}_{W^\rho, V^\rho} \left[ \frac{1}{\sqrt{m_2}} a^T D_{V^{(0)}+V^\rho+V', x_i} (V^{(0)} + V^\rho + (1-\eta)V') \Sigma \phi^*(x_i) \right] \right)^2 \\ & \leq \mathbb{E}_\Sigma \mathbb{E}_{W^\rho, V^\rho} \left( \frac{1}{\sqrt{m_2}} a^T D_{V^{(0)}+V^\rho+V', x_i} (V^{(0)} + V^\rho + (1-\eta)V') \Sigma \phi^*(x_i) \right)^2 \\ & = \mathbb{E}_{W^\rho, V^\rho} \mathbb{E}_\Sigma \left( \left\langle \frac{1}{\sqrt{m_2}} a^T D_{V^{(0)}+V^\rho+V', x_i} (V^{(0)} + V^\rho + (1-\eta)V'), \Sigma \phi^*(x_i) \right\rangle \right)^2 \\ & = \mathbb{E}_{W^\rho, V^\rho} \left\langle \left( \frac{1}{\sqrt{m_2}} a^T D_{V^{(0)}+V^\rho+V', x_i} (V^{(0)} + V^\rho + (1-\eta)V') \right)^{2\odot}, \phi^*(x_i)^{2\odot} \right\rangle \\ & \leq \mathbb{E}_{W^\rho, V^\rho} \left\| \frac{1}{\sqrt{m_2}} a^T D_{V^{(0)}+V^\rho+V', x_i} (V^{(0)} + V^\rho + (1-\eta)V') \right\|_2^2 \left\| \phi^*(x_i) \right\|_\infty^2 \\ & \leq \xi^2 \mathbb{E}_{W^\rho, V^\rho} \left\| \frac{1}{\sqrt{m_2}} a^T D_{V^{(0)}+V^\rho+V', x_i} (V^{(0)} + V^\rho + (1-\eta)V') \right\|_2^2 \\ & \leq 2\xi^2 \mathbb{E}_{W^\rho, V^\rho} \left\| \frac{1}{\sqrt{m_2}} a^T D_{V^{(0)}+V^\rho+V', x_i} (V^{(0)} + V^\rho) \right\|_2^2 + 2\xi^2 \mathbb{E}_{W^\rho, V^\rho} \left\| \frac{1}{\sqrt{m_2}} a^T D_{V^{(0)}+V^\rho+V', x_i} (1-\eta)V' \right\|_2^2 \\ & \leq \xi^2 \mathbb{E}_{W^\rho, V^\rho} \left\| \frac{1}{\sqrt{m_2}} a^T D_{V^{(0)}+V^\rho+V', x_i} (V^{(0)} + V^\rho) \right\|_2^2 + \xi^2 (1-\eta)^2 \|V'\|_F^2 \\ & \leq \xi^2 \mathbb{E}_{W^\rho, V^\rho} \left\| \frac{1}{\sqrt{m_2}} a^T D_{V^{(0)}+V^\rho+V', x_i} (V^{(0)} + V^\rho) \right\|_2^2 + \xi^2 (1-\eta)^2 C_2^2. \end{aligned}$$

Now under the event  $E^c$  defined in Lemma 33 we get that  $\|\phi^{(2)}(x_i)\| \lesssim C_1$ , so we can bound the above as

$$\leq \xi^2 \mathbb{E}_{V^\rho} \sup_{\|V'\| \leq C_2, V' \perp \phi^{(0)}(x_i), \|x'\| \leq C_1} \frac{1}{m_2} \left\| \sum_j a_j \mathbb{1}\{(V_j^{(0)} + V_j^\rho + V_j')(\phi^{(0)}(x_i) + x') \geq 0\} (V_j^{(0)} + V_j^\rho) \right\|_2^2 \quad (180)$$

$$+ \xi^2 (1-\eta)^2 C_2^2. \quad (181)$$

Now defining

$$\mathcal{L}_2 := \frac{1}{m_2} \left\| \sum_j a_j \mathbb{1}\{(V_j^{(0)} + V_j^\rho + V_j')(\phi^{(0)}(x_i) + x') \geq 0\} (V_j^{(0)} + V_j^\rho) \right\|_2^2,$$

to bound the first term, we want to apply Lemma 23 using the same trick that we did in the proof of Lemma 20. Note that  $\mathcal{L}_2$  is the same term as  $\Gamma_{x', V'}^2$  in Lemma 23 except that it is defined with respect to  $V^{(0)} + V^\rho$  instead of  $V^{(0)}$ . On the other hand, note that  $V^{(0)} + V^\rho$  has Gaussian entries with variance  $\kappa_2^2 + \frac{\beta_2^2}{m_2}$  and we know  $\kappa_2^2 \leq \kappa_2^2 + \frac{\beta_2^2}{m_2} \leq 2\kappa_2^2$ , which means the argument of Lemma 23 holds true here up to constants:

$$\sup_{\|V'\| \leq C_2, V' \perp \phi^{(0)}(x_i), \|x'\| \leq C_1} \mathcal{L}_2 \lesssim \mathfrak{R}_7^2.$$

This holds with probability say  $1 - \delta_2$  over the randomness of both  $V^{(0)}$  and  $V^\rho$ . Therefore, with probability  $1 - \sqrt{\delta_2}$  over the initialization, then with probability at least  $1 - \sqrt{\delta_2}$  over the randomness

of  $V^\rho$  we have the above. Moreover, with a simple Cauchy-Swartz we get the following almost surely bound:

$$\mathcal{L}_2 \lesssim \|V^{(0)}\|_F^2 + \|V^\rho\|_F^2. \quad (182)$$

Now the variable  $\|V^\rho\|^2$  is subexponential with parameter  $(\beta_1^4 m_3^2, \beta_1^2 m_3)$ . Furthermore, with high probability we have  $\|V^{(0)}\|_F^2 \lesssim m_2 m_3 \kappa_2^2$ . Therefore, taking

$$\mathcal{U}_2 := \Theta\left(\kappa_2^2 m_2 m_3 + \beta_1^2 m_3 \log\left((\kappa_2^2 m_2 m_3 + \beta_1^2 m_3)/\mathfrak{R}_7\right)\right),$$

then one can easily see by the subexponential tail:

$$\begin{aligned} \mathbb{E}[\mathcal{L}_2 | \mathcal{L}_2 \geq \mathcal{U}_2] &= \Theta(\mathcal{U}_2), \\ \mathbb{P}(\mathcal{L}_2 \geq \mathcal{U}_2) &\leq \mathfrak{R}_7^2 / \mathcal{U}_2. \end{aligned}$$

Hence, we can apply the same trick as Lemma [20](#) as

$$\begin{aligned} \mathbb{E}[\mathcal{L}_2] &= \mathbb{E}\left[\mathcal{L}_2 | \mathcal{L}_2 \leq \mathfrak{R}_7^2\right] \mathbb{P}(\mathcal{L}_2 \leq \mathfrak{R}_7^2) \\ &\quad + \mathbb{E}\left[\mathcal{L}_2 | \mathfrak{R}_7^2 \leq \mathcal{L}_2 \leq \mathcal{U}_2\right] \mathbb{P}(\mathfrak{R}_7^2 \leq \mathcal{L}_2 \leq \mathcal{U}_2) \\ &\quad + \mathbb{E}\left[\mathcal{L}_2 | \mathcal{U}_2 \leq \mathcal{L}_2\right] \mathbb{P}(\mathcal{U}_2 \leq \mathcal{L}_2) \\ &\lesssim E\left[\mathcal{L}_2 | \mathcal{L}_2 \leq \mathfrak{R}_7^2\right] + \mathbb{P}(\mathfrak{R}_7^2 \leq \mathcal{L}_2 \leq \mathcal{U}_2) + \mathfrak{R}_7^2 \\ &\lesssim \mathfrak{R}_7^2 + \sqrt{\delta_2} \mathcal{U}_2. \end{aligned}$$

Now taking  $\delta_2 \lesssim \mathfrak{R}_7^4 / \mathcal{U}_2^2$ , we finally get that conditioned on  $W^\rho$ 's where  $E$  happens, then

$$\mathbb{E}_{V^\rho} \mathcal{L}_2 \leq \mathfrak{R}_7^2.$$

On the other hand, to handle the case when  $E$  happens, we can use the bound in [\(182\)](#) as it does not depend on the occurrence of  $E$  as well:

$$\begin{aligned} \mathbb{E}_{W^\rho, V^\rho} \mathcal{L}_2 &\leq \mathbb{P}(E) \mathbb{E}_{V^\rho} (\|V^{(0)}\|^2 + \|V^\rho\|^2) + \mathbb{P}(E^c) \mathfrak{R}_7^2 \\ &\lesssim m_1 \exp\{-C_1^2 / (8m_3 \beta_1^2)\} (\kappa_2^2 m_2 m_3 + \beta_1^2 m_3) + \mathfrak{R}_7^2. \end{aligned}$$

Plugging this back into [\(181\)](#) we finally get

$$LHS \lesssim \xi^2 (m_1 \exp\{-C_1^2 / (8m_3 \beta_1^2)\} (\kappa_2^2 m_2 m_3 + \beta_1^2 m_3) + \mathfrak{R}_7^2) + \xi^2 C_2^2.$$

## A.15 BOUNDING THE WORST-CASE SENARIO

**Lemma 22** Suppose  $m_3 \geq \log(m_2)$  and  $\sqrt{m_3}\kappa_1 \gtrsim C_1$ . We define the sign matrices  $D_{V^{(0)}+V',x_i}^{x'}$  and  $D_{V^{(0)},x_i}^{x'}$  with respect to the multiplications

$$(V^{(0)} + V')(\phi^{(0)}(x_i) + x'),$$

and

$$V^{(0)}(\phi^{(0)}(x_i) + x').$$

Then, with high probability:

$$\begin{aligned} & \sup_{\|x'\| \leq C_1, \|V'\|_F \leq C_2, V' \perp \phi^{(0)}} \frac{1}{\sqrt{m_2}} a^T D_{V^{(0)}+V',x_i}^{x'} V^{(0)} \phi^{(0)}(x_i) \\ & \lesssim \left( \frac{(C_1 C_2)^{4/3}}{(\sqrt{m_2} \kappa_2)^{1/3} (\sqrt{m_3} \kappa_1)^{1/3}} + \frac{(C_1 C_2)^{2/3} m_3^{2/3} (\kappa_1 \kappa_2)^{1/3} \sqrt{\log(m_2)}}{m_2^{1/3}} \right) \left( 1 + \log(m_2) \frac{C_1^{1/3} (\kappa_2 \sqrt{m_2})^{2/3}}{C_2^{2/3} (\kappa_1 \sqrt{m_3})^{1/3}} \right) \\ & + \frac{m_3^{3/2} \kappa_1 \kappa_2}{\sqrt{m_2}} \sqrt{\log(m_2)} (\log(m_3) + \log(\log(m_2))) + \kappa_1 \kappa_2 \sqrt{m_3 \log(m_2)} := \mathfrak{R}_6. \end{aligned}$$

**Proof of Lemma 22**

Consider a cover for the euclidean ball of radius  $C_1$  in  $\mathbb{R}^{m_3}$  with precision  $\epsilon$ , i.e.  $B_{C_1}(\epsilon)$ . So for every  $x' \in \mathbb{R}^{m_3}$ , there exists an  $x \in B_{C_1}(\epsilon)$  such that  $\|x - x'\| \leq \epsilon$ , and  $|B_{C_1}(\epsilon)| \lesssim (\frac{1}{\epsilon})^{m_3}$ . Now fix  $x'$  and  $x$ . We have

$$\Gamma_{x',V'} := \frac{1}{\sqrt{m_2}} a^T D_{V^{(0)}+V',x_i}^{x'} V^{(0)} \phi^{(0)}(x_i) = \frac{1}{\sqrt{m_2}} \sum_{j=1}^{m_2} a_j \mathbb{1}\{(V_j^{(0)}+V'_j)(\phi^{(0)}(x_i)+x') \geq 0\} V_j^{(0)} \phi^{(0)}(x_i).$$

Now by a union bound, because each variable  $V_j^{(0)} \phi^{(0)}(x_i)$  is Gaussian with parameter  $\kappa_2 \|\phi^{(0)}(x_i)\|$  and using Equation (116), with high probability we have for every  $j \in [m_2]$ :

$$V_j^{(0)} \phi^{(0)}(x_i) \lesssim \kappa_2 \|\phi^{(0)}(x_i)\| \sqrt{\log(m_2)} \lesssim \kappa_1 \kappa_2 \sqrt{m_3 \log(m_2)}. \quad (183)$$

Therefore, by Hoeffding over the randomness of the Bernoulli variables  $a_j$ , for a fixed  $x'$  with high probability:

$$\Gamma_{x'} := \frac{1}{\sqrt{m_2}} \sum_{j=1}^{m_2} a_j \mathbb{1}\{V_j^{(0)}(\phi^{(0)}(x_i) + x') \geq 0\} V_j^{(0)} \phi^{(0)}(x_i) \lesssim \kappa_1 \kappa_2 \sqrt{m_3 \log(m_2)}.$$

On the other hand, We know that the VC-dimension of the class of binary functions with respect to halfspaces in  $\mathbb{R}^{m_3}$  is  $m_3 + 1$ . Therefore, the set of different sign patterns in matrices  $D_{V^{(0)},x_i}^{x'}$  is bounded by  $m_2^{m_3+1}$ , i.e. for

$$\mathcal{D} = \{D_{V^{(0)},x_i}^{x'} \mid x' \in \mathbb{R}^{m_3}\},$$

we have

$$|\mathcal{D}| \lesssim m_2^{m_3+1}.$$

Therefore, by taking a union bound over all sign matrices in  $\mathcal{D}$ , we get with high probability

$$\sup_{x'} \Gamma_{x'} \lesssim \kappa_1 \kappa_2 \sqrt{m_3 \log(m_2)} \sqrt{\log(m_2^{m_3+1})} = \kappa_1 \kappa_2 m_3 \log(m_2). \quad (184)$$

Now for a threshold  $r$  which satisfies

$$r \geq 2\sqrt{m_3} \kappa_2 \epsilon, \quad (185)$$

we define

$$\mathcal{J}_{x,r} = \{j \in [m_2] \mid |V_j^{(0)}(\phi^{(0)}(x_i) + x)| \leq r\}.$$

Now by Equation (116) and the assumption of the Lemma  $\sqrt{m_3}\kappa_1 \gtrsim C_1$ , we have

$$\|\phi^{(0)}(x_i) + x\| \leq \|\phi^{(0)}(x_i)\| + \|x\| \lesssim \sqrt{m_3} \kappa_1 + C_1. \quad (186)$$

$$\|\phi^{(0)}(x_i) + x\| \geq \|\phi^{(0)}(x_i)\| - \|x\| \gtrsim \sqrt{m_3}\kappa_1 - C_1 \gtrsim \sqrt{m_3}\kappa_1. \quad (187)$$

Hence,  $V_j^{(0)}(\phi^{(0)}(x_i) + x)$  is Gaussian with standard deviation at least  $\Omega(\kappa_2\sqrt{m_3}\kappa_1)$ . Therefore,

$$\mathbb{P}(|V_j^{(0)}(\phi^{(0)}(x_i) + x)| \leq r) \lesssim \frac{r}{\sqrt{m_3}\kappa_1\kappa_2}.$$

This implies

$$\mathbb{E}[|\mathcal{J}_{x,r}|] \lesssim \frac{r}{\sqrt{m_3}\kappa_1\kappa_2} m_2. \quad (188)$$

On the other hand, note that  $|\mathcal{J}_{x,r}|$  is the sum of  $m_2$  Bernoulli random variables, so it is subGaussian with parameter  $m_2$ . Therefore, with high probability

$$|\mathcal{J}_{x,r}| \lesssim \frac{r}{\sqrt{m_3}\kappa_1\kappa_2} m_2 + \sqrt{m_2}.$$

Now taking maximum over all  $x \in B_{C_1}(\epsilon)$  and exploiting the subGaussian tail of the random variables, we get with high probability

$$\max_{x \in B_{C_1}(\epsilon)} |\mathcal{J}_{x,r}| \lesssim \frac{r}{\sqrt{m_3}\kappa_1\kappa_2} m_2 + \sqrt{m_2 \log(|B_{C_1}(\epsilon)|)} \lesssim \frac{r}{\sqrt{m_3}\kappa_1\kappa_2} m_2 + \sqrt{m_2 m_3 \log(1/\epsilon)}. \quad (189)$$

Moreover, consider a threshold  $1 < \theta$ , such that  $e^{-\theta^2/8} \leq m_2/m_3$ , and define the following set of indices

$$\mathcal{J}_{x',\theta}^{(2)} := \{j \in [m_2] \mid |V_j^{(0)}x'| \geq \theta\kappa_2 C_1\}.$$

Then, using Lemma 29 and noting the fact that the standard deviation of Gaussians in  $V^{(0)}$  is  $\kappa_2$  and that  $\|\phi^{(2)}(x_i)\| \leq C_1$ , with high probability:

$$\sup_{x': \|x'\|=1} |\mathcal{J}_{x',\theta}^{(2)}| \leq m_3(\log(m_3) + \log(\log(m_2))). \quad (190)$$

Now note that for each  $j \in [m_2]$ ,  $\|V_j^{(0)}\|^2$  is subexponential with parameters  $(m_3\kappa_2^4, \kappa_2^2)$ , which means that with high probability:

$$\max_j \|V_j^{(0)}\|^2 \lesssim m_3\kappa_2^2 + \sqrt{m_3}\kappa_2^2 \sqrt{\log(m_2)} + \kappa_2^2 \log(m_2).$$

But with condition  $m_3 \geq \log(m_2)$ , we can further upper bound it as

$$\max_j \|V_j^{(0)}\|^2 \lesssim m_3\kappa_2^2.$$

Now for fixed  $x, x'$ , for  $j \in \mathcal{J}_{x,r}$  we have

$$\begin{aligned} |V_j^{(0)}(\phi^{(0)}(x_i) + x')| &\leq |V_j^{(0)}(\phi^{(0)}(x_i) + x)| + |V_j^{(0)}(x' - x)| \\ &\leq |V_j^{(0)}(\phi^{(0)}(x_i) + x)| + \|V_j^{(0)}\| \|x' - x\| \\ &\lesssim r + \sqrt{m_3}\kappa_2\epsilon. \end{aligned}$$

On the other hand, for  $j \notin \mathcal{J}_{x',\theta}^{(2)}$ :

$$|V_j^{(0)}x'| \leq \theta\kappa_2 C_1. \quad (191)$$

Therefore, for  $j \in \mathcal{J}_{x,r} - \mathcal{J}_{x',\theta}^{(2)}$ :

$$|V_j^{(0)}\phi^{(0)}(x_i)| \leq |V_j^{(0)}(\phi^{(0)}(x_i) + x')| + |V_j^{(0)}x'| \lesssim r + \sqrt{m_3}\kappa_2\epsilon + \theta\kappa_2 C_1. \quad (192)$$

In a similar fashion, if  $j \notin \mathcal{J}_{x,r}$ , then using assumption (185):

$$|V_j^{(0)}(\phi^{(0)}(x_i) + x')| \geq |V_j^{(0)}(\phi^{(0)}(x_i) + x)| - |V_j^{(0)}(x - x')| \gtrsim r - \sqrt{m_3}\kappa_2\epsilon \geq r/2. \quad (193)$$

Hence, using the fact that  $\phi^{(0)}(x_i)$  is orthogonal to  $V_j'$ :

$$\begin{aligned}
& \left| \mathbb{1}\{(V_j^{(0)} + V_j')(\phi^{(0)}(x_i) + x') \geq 0\} - \mathbb{1}\{V_j^{(0)}(\phi^{(0)}(x_i) + x') \geq 0\} \right| \\
& \leq \mathbb{1}\{|V_j'(\phi^{(0)}(x_i) + x')| \gtrsim |V_j^{(0)}(\phi^{(0)}(x_i) + x')|\} \\
& \leq \mathbb{1}\{|V_j'x'| \gtrsim |V_j^{(0)}(\phi^{(0)}(x_i) + x')|\} \\
& \leq \mathbb{1}\{\|V_j'\| \|x'\| \gtrsim |V_j^{(0)}(\phi^{(0)}(x_i) + x')|\} \\
& \leq \mathbb{1}\{\|V_j'\| C_1 \gtrsim |V_j^{(0)}(\phi^{(0)}(x_i) + x')|\} \\
& \leq \mathbb{1}\{\|V_j'\| \gtrsim \frac{|V_j^{(0)}(\phi^{(0)}(x_i) + x')|}{2C_1} + \frac{r}{4C_1}\}.
\end{aligned} \tag{194}$$

Now by triangle inequality and Equations (194), (192), (191) and the fact that  $\|V'\|_F \leq C_2$ , we can write:

$$\begin{aligned}
& |\Gamma_{x'} - \Gamma_{x', V'}| \\
& \leq \frac{1}{\sqrt{m_2}} \sum_{j \in \mathcal{J}_{x,r} - \mathcal{J}_{x',\theta}^{(2)}} \left| \mathbb{1}\{(V_j^{(0)} + V_j')(\phi^{(0)}(x_i) + x') \geq 0\} - \mathbb{1}\{V_j^{(0)}(\phi^{(0)}(x_i) + x') \geq 0\} \right| |V_j^{(0)} \phi^{(0)}(x_i)| \\
& + \frac{1}{\sqrt{m_2}} \sum_{j \notin (\mathcal{J}_{x,r} \cup \mathcal{J}_{x',\theta}^{(2)})} \left| \mathbb{1}\{(V_j^{(0)} + V_j')(\phi^{(0)}(x_i) + x') \geq 0\} - \mathbb{1}\{V_j^{(0)}(\phi^{(0)}(x_i) + x') \geq 0\} \right| |V_j^{(0)} \phi^{(0)}(x_i)| \\
& + \frac{1}{\sqrt{m_2}} \sum_{j \in \mathcal{J}_{x',\theta}^{(2)}} \left| \mathbb{1}\{(V_j^{(0)} + V_j')(\phi^{(0)}(x_i) + x') \geq 0\} - \mathbb{1}\{V_j^{(0)}(\phi^{(0)}(x_i) + x') \geq 0\} \right| |V_j^{(0)} \phi^{(0)}(x_i)| \\
& \leq \frac{1}{\sqrt{m_2}} |\mathcal{J}_{x,r} - \mathcal{J}_{x',\theta}^{(2)}| \max_{j \in \mathcal{J}_{x,r} - \mathcal{J}_{x',\theta}^{(2)}} |V_j^{(0)} \phi^{(0)}(x_i)| \\
& + \frac{1}{\sqrt{m_2}} \sum_{j \notin (\mathcal{J}_{x,r} \cup \mathcal{J}_{x',\theta}^{(2)})} \mathbb{1}\{\|V_j'\| \gtrsim \frac{|V_j^{(0)}(\phi^{(0)}(x_i) + x')|}{2C_1} + \frac{r}{4C_1}\} |V_j^{(0)} \phi^{(0)}(x_i)| \\
& + \frac{1}{\sqrt{m_2}} |J_{x',\theta}^{(2)}| \max_{j \in m_2} |V_j^{(0)} \phi^{(0)}(x_i)| \\
& \leq \frac{1}{\sqrt{m_2}} |\mathcal{J}_{x,r} - \mathcal{J}_{x',\theta}^{(2)}| \max_{j \in \mathcal{J}_{x,r} - \mathcal{J}_{x',\theta}^{(2)}} |V_j^{(0)} \phi^{(0)}(x_i)| \\
& + \frac{1}{\sqrt{m_2}} \sum_{j \notin (\mathcal{J}_{x,r} \cup \mathcal{J}_{x',\theta}^{(2)})} \mathbb{1}\{\|V_j'\| \gtrsim \frac{|V_j^{(0)}(\phi^{(0)}(x_i) + x')|}{2C_1} + \frac{r}{4C_1}\} (|V_j^{(0)}(\phi^{(0)}(x_i) + x')| + \theta \kappa_2 C_1) \\
& + \frac{1}{\sqrt{m_2}} |J_{x',\theta}^{(2)}| \max_{j \in m_2} |V_j^{(0)} \phi^{(0)}(x_i)| \\
& \lesssim \frac{1}{\sqrt{m_2}} |\mathcal{J}_{x,r} - \mathcal{J}_{x',\theta}^{(2)}| \max_{j \in \mathcal{J}_{x,r} - \mathcal{J}_{x',\theta}^{(2)}} |V_j^{(0)} \phi^{(0)}(x_i)| \\
& + \frac{1}{\sqrt{m_2}} \sum_{j \notin (\mathcal{J}_{x,r} \cup \mathcal{J}_{x',\theta}^{(2)})} \mathbb{1}\{\|V_j'\| \gtrsim \frac{r}{C_1}\} (C_1 \|V_j'\| + \theta \kappa_2 C_1) \\
& + \frac{1}{\sqrt{m_2}} |J_{x',\theta}^{(2)}| \max_{j \in m_2} |V_j^{(0)} \phi^{(0)}(x_i)| \\
& \lesssim \frac{1}{\sqrt{m_2}} |\mathcal{J}_{x,r} - \mathcal{J}_{x',\theta}^{(2)}| (r + \sqrt{m_3} \kappa_2 \epsilon + \theta \kappa_2 C_1) + \frac{C_1}{\sqrt{m_2}} \sqrt{\#\left(j : \|V_j'\| \gtrsim \frac{r}{C_1}\right)} \|V'\|_F^2
\end{aligned}$$

$$\begin{aligned}
& + \frac{1}{\sqrt{m_2}} \# \left( j : \|V_j'\| \gtrsim \frac{r}{C_1} \right) \theta \kappa_2 C_1 + \frac{1}{\sqrt{m_2}} |\mathcal{J}_{x',\theta}^{(2)}| \max_{j \in m_2} |V_j^{(0)} \phi^{(0)}(x_i)| \\
& \lesssim \frac{1}{\sqrt{m_2}} |\mathcal{J}_{x,r} - \mathcal{J}_{x',\theta}^{(2)}| (r + \sqrt{m_3} \kappa_2 \epsilon + \theta \kappa_2 C_1) + \frac{C_1^2 C_2^2}{\sqrt{m_2} r} \\
& + \frac{C_1^3 C_2^2}{\sqrt{m_2} r^2} \theta \kappa_2 + \frac{1}{\sqrt{m_2}} |\mathcal{J}_{x',\theta}^{(2)}| \max_{j \in m_2} |V_j^{(0)} \phi^{(0)}(x_i)|.
\end{aligned}$$

Now using Equations (208), (190), (183), and (185), and the bound on  $|\mathcal{J}_{x',\theta}^{(2)}|$  from Lemma 49, we write

$$\begin{aligned}
& \lesssim \frac{1}{\sqrt{m_2}} |\mathcal{J}_{x,r}| (r + \theta \kappa_2 C_1) + \frac{1}{\sqrt{m_2}} |\mathcal{J}_{x',\theta}^{(2)}| \kappa_1 \kappa_2 \sqrt{m_3 \log(m_2)} + \frac{C_1^2 C_2^2}{\sqrt{m_2} r} + \frac{C_1^3 C_2^2}{\sqrt{m_2} r^2} \theta \kappa_2 \\
& \leq \frac{1}{\sqrt{m_2}} \left( \frac{r}{\sqrt{m_3} \kappa_1 \kappa_2} m_2 + \sqrt{m_2 m_3 \log(1/\epsilon)} \right) (r + \theta \kappa_2 C_1) \\
& + \frac{1}{\sqrt{m_2}} \left( m_3 (\log(m_3) + \log(\log(m_2))) \right) \kappa_1 \kappa_2 \sqrt{m_3 \log(m_2)} + \frac{C_1^2 C_2^2}{r \sqrt{m_2}} + \frac{C_1^3 C_2^2}{\sqrt{m_2} r^2} \theta \kappa_2. \\
& \leq \frac{1}{\sqrt{m_2}} \left( \frac{r}{\sqrt{m_3} \kappa_1 \kappa_2} m_2 + \sqrt{m_2 m_3 \log(1/\epsilon)} \right) (r + \theta \kappa_2 C_1) \\
& + \frac{m_3^{3/2} \kappa_1 \kappa_2}{\sqrt{m_2}} \sqrt{\log(m_2)} (\log(m_3) + \log(\log(m_2))) + \frac{C_1^2 C_2^2}{r \sqrt{m_2}} + \frac{C_1^3 C_2^2}{\sqrt{m_2} r^2} \theta \kappa_2. \tag{195}
\end{aligned}$$

Now setting

$$r^* := (C_1 C_2)^{2/3} \frac{m_3^{1/6} (\kappa_1 \kappa_2)^{1/3}}{m_2^{1/3}}.$$

By this choice, from (195) we obtain

$$\begin{aligned}
|\Gamma_{x'} - \Gamma_{x',V'}| & \leq \left( \frac{(C_1 C_2)^{4/3}}{(\sqrt{m_2} \kappa_2)^{1/3} (\sqrt{m_3} \kappa_1)^{1/3}} + \frac{(C_1 C_2)^{2/3} m_3^{2/3} (\kappa_1 \kappa_2)^{1/3} \sqrt{\log(1/\epsilon)}}{m_2^{1/3}} \right) \left( 1 + \theta \frac{C_1^{1/3} (\kappa_2 \sqrt{m_2})^{2/3}}{C_2^{2/3} (\kappa_1 \sqrt{m_3})^{1/3}} \right) \\
& + \frac{m_3^{3/2} \kappa_1 \kappa_2}{\sqrt{m_2}} \sqrt{\log(m_2)} (\log(m_3) + \log(\log(m_2))).
\end{aligned}$$

Now we set

$$\theta^* := 3 \log(m_2),$$

which also satisfies the condition of Lemma 49 and combining with Equation (184), we get that with high probability

$$\begin{aligned}
|\Gamma_{x',V'}| & \leq |\Gamma_{x',V'} - \Gamma_{x'}| + |\Gamma_{x'}| \\
& \lesssim \left( \frac{(C_1 C_2)^{4/3}}{(\sqrt{m_2} \kappa_2)^{1/3} (\sqrt{m_3} \kappa_1)^{1/3}} + \frac{(C_1 C_2)^{2/3} m_3^{2/3} (\kappa_1 \kappa_2)^{1/3} \sqrt{\log(1/\epsilon)}}{m_2^{1/3}} \right) \left( 1 + \log(m_2) \frac{C_1^{1/3} (\kappa_2 \sqrt{m_2})^{2/3}}{C_2^{2/3} (\kappa_1 \sqrt{m_3})^{1/3}} \right) \\
& + \frac{m_3^{3/2} \kappa_1 \kappa_2}{\sqrt{m_2}} \sqrt{\log(m_2)} (\log(m_3) + \log(\log(m_2))) + \kappa_1 \kappa_2 \sqrt{m_3 \log(m_2)},
\end{aligned}$$

where  $\|x'\| \lesssim C_2$  and  $\|V'\|_F \leq C_2$ ,  $\forall j : V_j' \phi^{(0)}(x_i) = 0$ . We also need to satisfy condition (185), which regarding this choice for  $\theta = \theta^*$  becomes

$$r^* = (C_1 C_2)^{2/3} \frac{m_3^{1/6} (\kappa_1 \kappa_2)^{1/3}}{m_2^{1/3}} \geq 2\sqrt{m_3} \kappa_2 \epsilon, \tag{196}$$

for which it suffices to set

$$\epsilon^* := (C_1 C_2)^{2/3} \frac{\kappa_1^{1/3}}{2(m_2 m_3)^{1/3} \kappa_2^{2/3}}, \tag{197}$$

Substituting this choice of  $\epsilon$  above and picking the overparameterization large enough to dominate the magnitude of  $C_1, C_2$  so that  $\log(1/\epsilon^*) \lesssim \log(m_2)$ , the proof is complete.

**Lemma 23** *Under the following condition*

$$(\sqrt{m_2}\kappa_2)^{1/3}(\sqrt{m_3}\kappa_1)^{2/3} \geq \log^{-7/6}(m_2)(C_1C_2)^{1/3},$$

with high probability we have

$$\begin{aligned} & \sup_{\|x'\| \leq C_1, \|V'\|_F \leq C_2, V' \perp \phi^{(0)}} \frac{1}{\sqrt{m_2}} \left\| \sum_j a_j \mathbb{1}\{(V_j^{(0)} + V_j')(\phi^{(0)}(x_i) + x') \geq 0\} V_j^{(0)} \right\| \\ & \lesssim \sqrt{m_3}\kappa_2 \log(m_2) + \frac{(\sqrt{m_2}\kappa_2)^{1/3}}{(\sqrt{m_3}\kappa_1)^{2/3}} (C_1C_2)^{2/3} \log^{1/6}(m_2) := \mathfrak{R}_7. \end{aligned}$$

**Proof of Lemma 23**

Similar to Lemma 22, define the helper functions  $\Gamma_{x'}$  and  $\Gamma_{x',V'}$  as

$$\Gamma_{x',V'} = \frac{1}{\sqrt{m_2}} \left\| \sum_j a_j \mathbb{1}\{(V_j^{(0)} + V_j')(\phi^{(0)}(x_i) + x') \geq 0\} V_j^{(0)} \right\|, \quad (198)$$

$$\Gamma_{x'} = \frac{1}{\sqrt{m_2}} \left\| \sum_j a_j \mathbb{1}\{V_j^{(0)}(\phi^{(0)}(x_i) + x') \geq 0\} V_j^{(0)} \right\|. \quad (199)$$

First we bound  $\sup_{x'} \Gamma_{x'}$ . To this end, note that because  $V_j^{(0)} \in \mathbb{R}^{m_3}$  and the VC-dimension of half-planes is  $m_3 + 1$ , then by Sauer's Lemma, the set

$$\mathcal{D} = \{D_{V^{(0)},x_i}^{x'} \mid x' \in \mathbb{R}^{m_3}, \|x'\| \lesssim C_1\}$$

of all sign pattern matrices has cardinality at most

$$|\mathcal{D}| \leq m_2^{m_3+1}.$$

Now note that with high probability, the entries of the matrix  $V^{(0)}$  are all less than  $O(\kappa_2 \sqrt{\log(m_2 m_3)})$ . On the other hand, for each fixed sign pattern  $D_{V^{(0)},x_i}^{x'}$ , we have for the sum with respect to this sign pattern:

$$\left\| \frac{1}{\sqrt{m_2}} \sum_j a_j \mathbb{1}\{V_j^{(0)}(\phi^{(0)}(x_i) + x') \geq 0\} V_j^{(0)} \right\|^2 \quad (200)$$

is  $(m_3\kappa_2^4 \log^2(m_2 m_3), \kappa_2^2 \log(m_2 m_3))$  sub-exponential with respect to the randomness of  $a$ , because each entry of the vector  $\frac{1}{\sqrt{m_2}} \sum_j a_j \mathbb{1}\{V_j^{(0)}(\phi^{(0)}(x_i) + x') \geq 0\} V_j^{(0)}$  is  $(\kappa_2 \sqrt{\log(m_2 m_3)})$ -subGaussian. Therefore, with high probability we have

$$\left\| \frac{1}{\sqrt{m_2}} \sum_j a_j \mathbb{1}\{V_j^{(0)}(\phi^{(0)}(x_i) + x') \geq 0\} V_j^{(0)} \right\|^2 \quad (201)$$

$$\leq \mathbb{E}_a \left[ \left\| \frac{1}{\sqrt{m_2}} \sum_j a_j \mathbb{1}\{V_j^{(0)}(\phi^{(0)}(x_i) + x') \geq 0\} V_j^{(0)} \right\|^2 \right] + \text{deviation} \quad (202)$$

$$\lesssim m_3\kappa_2^2 \log(m_2 m_3) + \sqrt{m_3}\kappa_2^2 \log(m_2 m_3) + \kappa_2^2 \log(m_2 m_3). \quad (203)$$

Similarly, if we take a union bound over all sign matrices in  $\mathcal{D}$  and using the fact that  $m_2 > m_3$ :

$$\sup_{x'} \Gamma_{x'}^2 = \sup_{x' \in \mathbb{R}^{m_3}} \left\| \frac{1}{\sqrt{m_2}} \sum_j a_j \mathbb{1}\{V_j^{(0)}(\phi^{(0)}(x_i) + x') \geq 0\} V_j^{(0)} \right\|^2 \quad (204)$$

$$\lesssim m_3\kappa_2^2 \log(m_2 m_3) + \sqrt{m_3}\kappa_2^2 \log(m_2 m_3) \sqrt{\log(m_2^{m_3+1})} + \kappa_2^2 \log(m_2 m_3) \log(m_2^{m_3+1}) \quad (205)$$

$$\lesssim m_3\kappa_2^2 \log^2(m_2), \quad (206)$$

which implies

$$\sup_{x'} \Gamma_{x'} \lesssim \sqrt{m_3}\kappa_2 \log(m_2). \quad (207)$$

Moreover, defining  $\mathcal{J}_{v,x}$  similar to Lemma [22](#) and using the similar approach we get with high probability

$$\max_{x \in B_{C_1}(\epsilon)} |\mathcal{J}_{x,r}| \lesssim \frac{r}{\sqrt{m_3} \kappa_1 \kappa_2} m_2 + \sqrt{m_2 \log(|B_{C_1}(\epsilon)|)} \lesssim \frac{r}{\sqrt{m_3} \kappa_1 \kappa_2} m_2 + \sqrt{m_2 m_3 \log(1/\epsilon)}. \quad (208)$$

Now for simplifying the analysis, we assume that for indices  $j \in \mathcal{J}_{x,r}$  we can change the sign pattern with no cost on  $V'$ , i.e. we can pick any subset of them. Therefore, we first compute a high probability upper bound on the following quantity:

$$\frac{1}{\sqrt{m_2}} \sup_{S \subset \mathcal{J}_{x,r}, \pm \text{signs}} \left\| \sum_{j \in S} \pm V_j^{(0)} \right\|. \quad (209)$$

If we form the matrix  $V^{(0)}(\mathcal{J}_{x,r})$  be the matrix which only keeps the rows with indices in  $\mathcal{J}_{x,r}$ , then the above quantity can be computed as

$$\frac{1}{\sqrt{m_2}} \sup_{S \subset \mathcal{J}_{x,r}, \pm \text{signs}} \left\| \sum_{j \in S} \pm V_j^{(0)} \right\| = \frac{1}{\sqrt{m_2}} \sup_{v \in \{0,1,-1\}^{|\mathcal{J}_{x,r}|}} \|v^T V^{(0)}(\mathcal{J}_{x,r})\| \quad (210)$$

$$\leq \frac{1}{\sqrt{m_2}} \lambda_{\max}(V^{(0)}(\mathcal{J}_{x,r})) \sup \|v\| \leq \frac{1}{\sqrt{m_2}} \lambda_{\max}(V^{(0)}(\mathcal{J}_{x,r})) |\mathcal{J}_{x,r}|, \quad (211)$$

where  $\lambda_{\max}$  is the maximum singular value of the matrix. Now by random matrix theory, we know for a fixed  $x$  and arbitrary  $t \geq 0$ , the following argument holds:

$$\mathbb{P}(\lambda_{\max}(V^{(0)}(\mathcal{J}_{x,r}))/\kappa_2 \gtrsim \sqrt{m_3} + \sqrt{|\mathcal{J}_{x,r}|} + t) \leq 2e^{-ct^2}. \quad (212)$$

Therefore, as  $|\mathcal{D}| \leq m_2^{m_3+1}$ , we get with high probability

$$\max_{x \in B_{C_1}(\epsilon)} \lambda_{\max}(V^{(0)}(\mathcal{J}_{x,r})) \lesssim \kappa_2 (\sqrt{m_3} + \sqrt{|\mathcal{J}_{x,r}|} + \sqrt{\log(m_2^{m_3+1})}) \quad (213)$$

$$\leq \kappa_2 \sqrt{\log(m_2) m_3} + \kappa_2 \sqrt{|\mathcal{J}_{x,r}|}. \quad (214)$$

Therefore, with high probability

$$\sup_{x \in B_{C_1}(\epsilon)} \frac{1}{\sqrt{m_2}} \sup_{S \subset \mathcal{J}_{x,r}} \left\| \sum_{j \in S} V_j^{(0)} \right\| \leq \frac{\kappa_2}{\sqrt{m_2}} (\sqrt{\log(m_2) m_3} |\mathcal{J}_{x,r}| + |\mathcal{J}_{x,r}|). \quad (215)$$

On the other hand, as in Equation [\(193\)](#) in the proof of Lemma [22](#), for  $j \notin \mathcal{J}_{x,r}$  we have:

$$|V_j^{(0)}(\phi^{(0)}(x_i) + x')| \geq |V_j^{(0)}(\phi^{(0)}(x_i) + x)| - |V_j^{(0)}(x - x')| \gtrsim r - \sqrt{m_3} \kappa_2 \epsilon. \quad (216)$$

Picking

$$\epsilon^* := \frac{r}{2\sqrt{m_3} \kappa_2},$$

we get for  $j \notin \mathcal{J}_{x,r}$

$$|V_j^{(0)}(\phi^{(0)}(x_i) + x')| \gtrsim r.$$

Now similar to the derivation in [\(194\)](#) we have

$$\left| \mathbb{1}\{(V_j^{(0)} + V_j')(\phi^{(0)}(x_i) + x') \geq 0\} - \mathbb{1}\{V_j^{(0)}(\phi^{(0)}(x_i) + x') \geq 0\} \right| \quad (217)$$

$$\leq \mathbb{1}\{\|V_j'\|_{C_1} \gtrsim |V_j^{(0)}(\phi^{(0)}(x_i) + x')|\} \quad (218)$$

$$\leq \mathbb{1}\{\|V_j'\| \gtrsim \frac{r}{C_1}\}. \quad (219)$$

Hence, because  $\|V'\|_F \leq C_2$ , the number of indices for which  $\mathbb{1}\{(V_j^{(0)} + V_j')(\phi^{(0)}(x_i) + x') \geq 0\} \neq \mathbb{1}\{V_j^{(0)}(\phi^{(0)}(x_i) + x') \geq 0\}$  is at most  $l = \frac{(C_1 C_2)^2}{r^2}$ . Therefore, we bound the following quantity to use in the analysis:

$$\sup_{S \subset [m_2] \ \& \ |S| \leq l, \pm \text{signs}} \left\| \sum_{j \in S} \pm V_j^{(0)} \right\|. \quad (220)$$

But if we define for  $m_2 < j \leq 2m_2$ ,

$$V_j^{(0)} = -V_{j-m_2}^{(0)},$$

then

$$\sup_{S \subset [m_2] \ \& \ |S| \leq l, \pm \text{signs}} \left\| \sum_{j \in S} \pm V_j^{(0)} \right\| \leq \sup_{S \subset [2m_2] \ \& \ |S| \leq l} \left\| \sum_{j \in S} V_j^{(0)} \right\|.$$

Now note that each entry of  $\sum_{j \in S} V_j^{(0)}$  is  $\sqrt{l}\kappa_2$  subGaussian. Hence, the quantity  $\left\| \sum_{j \in S} V_j^{(0)} \right\|^2$  is  $(m_3 l^2 \kappa_2^4, l \kappa_2^2)$  subexponential. Therefore, we have with high probability

$$\begin{aligned} \sup_{|S| \leq l} \left\| \sum_{j \in S} V_j^{(0)} \right\|^2 &\lesssim \mathbb{E} \left\| \sum_{j \in S} V_j^{(0)} \right\|^2 + \sqrt{m_3} l \kappa_2^2 \sqrt{\log \binom{2m_2}{l}} + l \kappa_2^2 \log \binom{2m_2}{l} \\ &\lesssim m_3 l \kappa_2^2 + \sqrt{m_3} l \kappa_2^2 \sqrt{\log \binom{2m_2}{l}} + l \kappa_2^2 \log \binom{2m_2}{l} \\ &\lesssim m_3 l \kappa_2^2 + \sqrt{m_3} l^{3/2} \kappa_2^2 \sqrt{\log(m_2)} + l^2 \kappa_2^2 \log(m_2). \end{aligned}$$

Hence

$$\sup_{|S| \leq l} \left\| \sum_{j \in S} V_j^{(0)} \right\| \lesssim \sqrt{m_3} \sqrt{l} \kappa_2 + l \kappa_2 \sqrt{\log(m_2)}.$$

Now using Equation , we can write

$$\begin{aligned} &|\Gamma_{x'} - \Gamma_{x', V'}| \\ &\leq \frac{1}{\sqrt{m_2}} \left\| \sum_{j \in \mathcal{J}_{x,r}} (\mathbb{1}\{(V_j^{(0)} + V'_j)(\phi^{(0)}(x_i) + x') \geq 0\} - \mathbb{1}\{V_j^{(0)}(\phi^{(0)}(x_i) + x') \geq 0\}) V_j^{(0)} \right\| \\ &+ \frac{1}{\sqrt{m_2}} \left\| \sum_{j \notin \mathcal{J}_{x,r}} (\mathbb{1}\{(V_j^{(0)} + V'_j)(\phi^{(0)}(x_i) + x') \geq 0\} - \mathbb{1}\{V_j^{(0)}(\phi^{(0)}(x_i) + x') \geq 0\}) V_j^{(0)} \right\| \\ &\leq \frac{1}{\sqrt{m_2}} \sup_{S \subset \mathcal{J}_{x,r}, \pm \text{signs}} \left\| \sum_{j \in S} \pm V_j^{(0)} \right\| \\ &+ \frac{1}{\sqrt{m_2}} \sup_{S \subset [m_2] \ \& \ |S| \leq (\frac{C_1 C_2}{r})^2, \pm \text{signs}} \left\| \sum_{j \in S} \pm V_j^{(0)} \right\| \\ &\leq \frac{\kappa_2}{\sqrt{m_2}} (\sqrt{\log(m_2)} m_3 |\mathcal{J}_{x,r}| + |\mathcal{J}_{x,r}|) + \frac{1}{\sqrt{m_2}} (\sqrt{m_3} \sqrt{l} \kappa_2 + l \kappa_2 \sqrt{\log(m_2)}) \\ &\leq \frac{\kappa_2}{\sqrt{m_2}} \left( \frac{r m_2}{\sqrt{m_3} \kappa_1 \kappa_2} + \sqrt{m_2 m_3 \log \left( \frac{\sqrt{m_3} \kappa_2}{r} \right)} \right)^{1/2} (\sqrt{\log(m_2)} m_3 + \sqrt{|\mathcal{J}_{x,r}|}) \\ &+ \frac{1}{\sqrt{m_2}} (\sqrt{m_3} (C_1 C_2 / r) \kappa_2 + (C_1 C_2 / r)^2 \kappa_2 \sqrt{\log(m_2)}) \\ &\lesssim \frac{\kappa_2}{\sqrt{m_2}} \left( \frac{r m_2}{\sqrt{m_3} \kappa_1 \kappa_2} + \sqrt{m_2 m_3 \log \left( \frac{\sqrt{m_3} \kappa_2}{r} \right)} \right) \\ &+ \frac{1}{\sqrt{m_2}} (\sqrt{m_3} (C_1 C_2 / r) \kappa_2 + (C_1 C_2 / r)^2 \kappa_2 \sqrt{\log(m_2)}). \end{aligned}$$

Combining this with (207):

$$|\Gamma_{x', V'}| \lesssim \sqrt{m_3} \kappa_2 \log(m_2) \tag{221}$$

$$+ \frac{r \sqrt{m_2}}{\sqrt{m_3} \kappa_1} + \kappa_2 \sqrt{m_3 \log \left( \frac{\sqrt{m_3} \kappa_2}{r} \right)} + \frac{1}{\sqrt{m_2}} (\sqrt{m_3} (C_1 C_2 / r) \kappa_2 + (C_1 C_2 / r)^2 \kappa_2 \sqrt{\log(m_2)}). \tag{222}$$

Now setting

$$r^* := \frac{m_3^{1/6} (\kappa_1 \kappa_2)^{1/3}}{m_2^{1/3}} (C_1 C_2)^{2/3} \log^{1/6}(m_2),$$

we get

$$LHS \lesssim \sqrt{m_3} \kappa_2 \log(m_2) + \frac{(\sqrt{m_2} \kappa_2)^{1/3}}{(\sqrt{m_3} \kappa_1)^{2/3}} (C_1 C_2)^{2/3} \log^{1/6}(m_2) \quad (223)$$

$$+ \kappa_2 m_3^{1/2} \log^{1/2}(m_2) + \frac{m_3^{1/3} \kappa_2^{2/3} (C_1 C_2)^{1/3}}{m_2^{1/6} \kappa_1^{1/3} \log^{1/6}(m_2)} \quad (224)$$

$$\lesssim \sqrt{m_3} \kappa_2 \log(m_2) + \frac{(\sqrt{m_2} \kappa_2)^{1/3}}{(\sqrt{m_3} \kappa_1)^{2/3}} (C_1 C_2)^{2/3} \log^{1/6}(m_2) \quad (225)$$

$$+ \frac{m_3^{1/2} \kappa_2 (C_1 C_2)^{1/3}}{(\sqrt{m_2} \kappa_2)^{1/3} (\sqrt{m_3} \kappa_1)^{1/3} \log^{1/6}(m_2)}. \quad (226)$$

Now under the condition

$$(\sqrt{m_2} \kappa_2)^{1/3} (\sqrt{m_3} \kappa_1)^{2/3} \geq \log^{-7/6}(m_2) (C_1 C_2)^{1/3},$$

The final term is dominated by the first term, which finally completes the proof.

## A.16 CONVERGENCE

The goal of this section is to prove Theorem [7](#).

**Theorem 7** *Letting  $\aleph = 4B^2$ , by Corollary [5](#), we have  $L^\Pi(0) \leq \aleph$ . We define the domain  $\mathcal{D}_l := \{\|w'\| \leq C_1 := \frac{\aleph+4l}{\psi_1}, \|v'\| \leq C_2 := \frac{\aleph+4l}{\psi_2}\}$ . For a large enough constant  $l = O(1)$  and function  $L^\Pi(w := (w', v')) : \mathbb{R}^N \rightarrow \mathbb{R}$ . Moreover, suppose  $L^\Pi$  is  $\rho_1$  lipschitz,  $\rho_2$  gradient lipschitz, and  $\rho_3$  hessian lipschitz in the domain  $\mathcal{D}_l$ , in the sense that their first, second, and third directional derivatives in an arbitrary unit direction is bounded by the corresponding parameters. Suppose we have access to the gradient of  $L^\Pi$  at each point in  $\mathcal{D}_l$  plus a zero mean noise vector  $\mathcal{L}$  such that  $\sigma_1^2 I \leq \mathbb{E}\mathcal{L}\mathcal{L}^T \leq \sigma_2^2 I$  and  $\|\mathcal{L}\| \leq Q$  almost surely. Also, suppose for a threshold  $\aleph_\ell \leq \aleph$ , if  $L^\Pi(w) \geq \aleph_\ell$  for any  $w \in \mathcal{D}_l$ , then we either have*

$$(1) \|\nabla L^\Pi(w)\| \geq \frac{\nu}{16\sqrt{C_1^2 + C_2^2}}, \quad (227)$$

$$(2) \lambda_{\min}(\nabla^2 L^\Pi(w)) \leq -\gamma. \quad (228)$$

Then starting from  $w_0 = 0$ , with probability at least 0.999 after at most  $\text{poly}(\rho_1, \rho_2, \rho_3, Q, \aleph, C_1, C_2, 1/\gamma, \log(\sigma_1/\sigma_2))$  number of iterations, we reach a point  $w_t$  such that  $L^\Pi(w_t) \leq \aleph_\ell$ .

**Proof** Our proof here is a refined version of that in [Ge et al. \(2015a\)](#). As we mentioned in section [5](#), the key fact that we are using in the other parts of our proof is a uniform upper bound  $\|w'\| \leq C_1, \|v'\| \leq C_2$  which is unjustified by only naively using [Ge et al. \(2015a\)](#). Here, first we restate a refined version of Lemmas 14 and 16 in [Ge et al. \(2015a\)](#) in Lemmas [24](#) and [26](#) respectively, and then use them to also bound the upward deviations of  $L^\Pi$ . Moreover, to avoid writing repeated proofs and overwhelm the reader, we mostly treat the arguments in Lemma 16 of [Ge et al. \(2015a\)](#) as blackbox and use them for our purpose here. A point to mention before we start, unlike Lemma 14 of [Ge et al. \(2015a\)](#) where the dependency on other parameters than the step size  $\eta$  is more explicit, Lemma 16 hides the dependencies on all the other parameters (which is polynomial). Here, we follow the same style.

We refer to the trajectory of the steps of algorithm by  $(w_t)_{t \geq 0}$ . In Lemmas of this section, To avoid introducing new notation and complicating things, we refer to the current point of the algorithm by  $w_0$ , while for the next point of the algorithm we use  $w_1$  (in Lemma [24](#)), and  $w_T$  (in Lemma [26](#)) respectively. Also, similar to [Ge et al. \(2015a\)](#),  $\tilde{O}$  and  $\tilde{\Omega}$  below means we are looking at the dependency on  $\eta$ .

**Lemma 24** *Suppose  $L^\Pi(w_0) \leq \aleph + 2l$ , and consider a parameter  $\chi > 1$  which can be set arbitrarily. For every point  $w_0$  such that  $L^\Pi(w_0) \leq \aleph + 2l, \|\nabla L^\Pi(w_0)\| \geq 2\sqrt{\eta(Q^2 + \sigma_2^2 N)}\rho_2\rho_1^2(2\chi + \frac{1}{2})$ , then for  $w_1 := w_0 - \eta(\nabla L^\Pi(w_0) + \mathcal{L})$  and random variable  $\mathfrak{R}_1$  (depending on  $w_0$ ) defined as*

$$\mathbb{E}L^\Pi(w_1) - L^\Pi(w_0) = -\eta^2\mathfrak{R}_1^2, \quad (229)$$

we have  $\mathfrak{R}_1 = \tilde{\Omega}(1)$  a.s., and almost surely:

$$\left|L^\Pi(w_1) - L^\Pi(w_0)\right| \leq \eta\mathfrak{R}_1/\sqrt{\rho_2\chi}.$$

(the expectation is over the randomness of  $\mathcal{L}$ ).

**Proof** This lemma is a tuned version of Lemma 14 in [Ge et al. \(2015a\)](#). First, note that the condition  $L^\Pi(w_0) \leq \aleph + 2l$  assures the smoothness coefficients  $\rho_1, \rho_2$  and  $\rho_3$  for  $L^\Pi$  by Corollary [??](#). We follow similar to [Ge et al. \(2015a\)](#) (picking  $\eta < 1/(2/\rho_2)$ ):

$$\begin{aligned} \mathbb{E}L^\Pi(w_1) - L^\Pi(w_0) &\leq -\frac{\eta}{2}\|\nabla L^\Pi(w_0)\|^2 + \frac{\eta^2\sigma_2^2\rho_2N}{2} \\ &\leq -\frac{\eta}{4}\|\nabla L^\Pi(w_0)\|^2 - \eta^2(\sigma_2^2N + Q^2)\rho_2\rho_1^2(2\chi + \frac{1}{2}) + \frac{\eta^2\sigma_2^2\rho_2N}{2} \\ &\leq -\frac{\eta}{4}\|\nabla L^\Pi(w_0)\|^2 - 2\eta^2Q^2\rho_2\rho_1^2\chi. \end{aligned} \quad (230)$$

On the other hand,  $L^\Pi$  is  $\rho_1$  Lipschitz, so we have almost surely

$$\begin{aligned} |L^\Pi(w_1) - L^\Pi(w_0)| &\leq \rho_1 \eta (\|\nabla L^\Pi(w_0) + \mathcal{L}\|) \leq \rho_1 \eta (\|\nabla L^\Pi(w_0)\| + \|\mathcal{L}\|) \\ &\leq \rho_1 \eta (\|\nabla L^\Pi(w_0)\| + Q). \end{aligned} \quad (231)$$

(To be completely precise, we should justify that we can write the Lipschitz inequality at point  $w_0$ , we also need to make sure that  $w_1$  remains in the domain that we have the Lipschitz parameter in, i.e.  $\mathcal{D}_l$ . To see why this is true, see the next Corollary).

Therefore

$$(\rho_2 \chi) \left| L^\Pi(w_1) - L^\Pi(w_0) \right|^2 \leq 2\eta^2 \rho_1^2 \rho_2 \chi \|\nabla L^\Pi(w_0)\|^2 + 2\rho_2 \chi \eta^2 \rho_1^2 Q^2. \quad (232)$$

Taking

$$\eta \leq (8\rho_1^2 \rho_2 \chi)^{-1}, \quad (233)$$

we get from Equation (230):

$$\mathbb{E}L^\Pi(w_1) - L^\Pi(w_0) \leq -2\eta^2 \rho_1^2 \rho_2 \chi \|\nabla L^\Pi(w_0)\|^2 - 2\rho_2 \chi \eta^2 \rho_1^2 Q^2. \quad (234)$$

Combining Equations (232) and (234), we see that

$$(\rho_2 \chi) \left| L^\Pi(w_1) - L^\Pi(w_0) \right|^2 \leq -(\mathbb{E}L^\Pi(w_1) - L^\Pi(w_0)).$$

Hence, if we define

$$\mathbb{E}L^\Pi(w_1) - L^\Pi(w_0) := -\eta^2 \mathfrak{R}_1^2,$$

we get

$$\left| L^\Pi(w_1) - L^\Pi(w_0) \right| \leq \eta \mathfrak{R}_1 / \sqrt{\rho_2 \chi}, \quad (235)$$

and from Equation (234), that

$$\mathfrak{R}_1^2 \geq 2\rho_1^2 \rho_2 \chi \|\nabla L^\Pi(w_0)\|^2 + 2\rho_2 \chi \rho_1^2 Q^2 \geq 2\rho_2 \chi \rho_1^2 Q^2 = \tilde{\Omega}(1).$$

Moreover, because the function is  $\rho_1$ -Lipshitz at the domain point  $w_0$ , we get from Equation (231):

$$-\eta_R^2 R R_1^2 = \mathbb{E}L^\Pi(w_1) - L^\Pi(w_0) \geq -\eta \rho_1 (\|\nabla L^\Pi(w_0)\| + Q) \geq -\eta \rho_1 (\rho_1 + Q) \geq -\frac{1}{\rho_2 \chi},$$

by taking

$$\eta \leq (\rho_1 (\rho_1 + Q) \rho_2 \chi)^{-1},$$

which implies

$$\eta \mathfrak{R}_1 \leq \frac{1}{\sqrt{\rho_2 \chi}}.$$

This combined with Equation (235) and a triangle inequality implies:

$$\left| L^\Pi(w_1) - \mathbb{E}L^\Pi(w_1) \right| \leq \eta \mathfrak{R}_1 / \sqrt{\rho_2 \chi} + \eta^2 \mathfrak{R}_1^2 \leq 2\eta \mathfrak{R}_1 / \sqrt{\rho_2 \chi}. \quad (236)$$

**Lemma 25** *As long as the value of the function at some  $w$  is bounded by  $\aleph + 2l$  ( $L^\Pi(w) \leq \aleph + 2l$ ), then  $\eta$  can be picked small enough (polynomially in other parameters) so that the change of the function by a step is a.s. bounded by  $l$ .*

**Proof** Let  $\psi = \min\{\psi_1, \psi_2\}$ . First, note that as the function is bounded by  $\aleph + 2l$ , we have the Lipschitz parameter  $\rho_1$ , hence  $\|\nabla L^\Pi(w)\| \leq \rho_1$ . Therefore, the change in  $w$  in a step is bounded as

$$\|\nabla L^\Pi(w) + \mathcal{L}\| \leq Q + \rho_1.$$

So by picking

$$\eta \leq \left( \sqrt{\frac{\aleph + 3l}{\psi}} - \sqrt{\frac{\aleph + 2l}{\psi}} \right) / (Q + \rho_1),$$

we guarantee that the value of  $w$  after a step remains in the ball of radius  $\sqrt{\frac{\aleph+3l}{\psi}}$ , hence we still have the smoothing parameters even after one step. Therefore, now we can use the Lipschitz parameter  $\rho_1$  to bound the value of the function after one step as it is written in Equation Equation (231). Using this Equation, it is enough to pick  $\eta$  as small as:

$$\eta \leq l/(\|\nabla L^\Pi(w_0)\| + Q), \quad (237)$$

so that the change in the function would be most  $l$  as desired.

**Lemma 26** For a fixed point  $w_0$  s.t.  $L^\Pi(w_0) \leq \aleph + 2l$ , suppose we pick  $\eta$  small enough such that

$$\xi(\eta) := 2\sqrt{\eta(Q^2 + \sigma_2^2 N)\rho_2\rho_1^2(2\chi + \frac{1}{2})} < \frac{\nu}{16\sqrt{C_1^2 + C_2^2}}.$$

Then, note that for  $\|\nabla L^\Pi(w_0)\| \leq \xi(\eta)$ , condition 228 implies:

$$\lambda_{min}(\nabla^2 L^\Pi(w_0)) \leq \gamma.$$

Then, using the notation  $\mathfrak{E}_T$  for the high probability event corresponding to Equations (36) and (44) in Ge et al (2015a), for small enough  $\eta$  (polynomially small w.r.t other parameters), for  $\mathfrak{R}_2$  defined as

$$\mathbb{E}[L^\Pi(w_T) - L^\Pi(w_0)]\mathbb{1}\{\mathfrak{E}_T\} = -\mathfrak{R}_2^2\eta^2, \quad (238)$$

we have almost surely

$$\left| [L^\Pi(w_T) - L^\Pi(w_0)]\mathbb{1}\{\mathfrak{E}_T\} \right| \leq \mathfrak{R}_2\eta/\sqrt{\rho_2\chi}. \quad (239)$$

Note that in the expectations above  $w_0$  is assumed fixed. Furthermore, we can assume  $\mathbb{P}(\mathfrak{E}_T) \geq 1 - \tilde{O}(\eta^5)$ .

**Proof** This Lemma is a tuned version of Lemma 16 in Ge et al (2015a). We change a couple of things here. First, we consider an implicit coupling that if  $w_t$  exits  $\mathcal{D}_l$  we do not move it anymore, i.e.  $w_{t'} = w_t, \forall t' \geq t$ , which means the noise vectors also becomes zero, i.e.  $\mathcal{L}_{t'} = 0, \forall t' \geq t$ . This way, the sequence of noise vectors remain bounded by  $Q$ , because if  $w_t$  is inside  $\mathcal{D}_l$ , then by assumption  $\|\mathcal{L}_t\| \leq Q$ , while otherwise  $\mathcal{L}_t = 0$ . We call this event  $\mathcal{E}_T$  for  $T$  defined in Lemma 16 of Ge et al (2015a). Note that we also have the smoothing parameters  $\rho_1, \rho_2, \rho_3$  for all  $(w_t)$  because of this coupling. In fact, we will use a more strict coupling; we consider the event  $\mathfrak{E}_T$  to be the high probability event corresponding to the bounds in Equations (44) and (36) of Ge et al (2015a) holding for all  $t \leq T$ ; We will see that  $\mathfrak{E}_T \subseteq \mathcal{E}_T$  at the end of this proof, but for now we assume it is true. An important point to note here is that in Ge et al (2015a),  $\mathbb{P}(\mathfrak{E}_T)$  is bounded by  $O(\eta^2)$ . However, the exponent dependency of  $\eta$  in this bound comes from Azuma-Hoeffding type inequalities, particularly used in Equations (60) and (42) in Ge et al (2015a), in which by considering larger constants one can easily get higher exponents. Unlike their analysis though, a bit stronger dependence of  $\eta^5$  is sufficient for our later use.

Also, because the distribution of our noise depends on the point  $w$ , our sequence of noise vectors  $(\mathcal{L}_t)$  is a martingale instead of being i.i.d, so we apply Azuma-Hoeffding inequality instead of the simple Hoeffdings in Lemma 16 of Ge et al (2015a). (because we are also sampling also a random  $(x_i, y_i)$  to compute the estimate of gradient, this could be simplified to the case where we compute the actual gradient and the injecting an i.i.d noise vector in each step, but it is an overhead to compute the actual gradient, so here we choose to analyze the more complicated case.)

Next, notice the definition of  $\Lambda$  and  $\tilde{\Lambda}$  right after Equation (66) in Ge et al (2015a), which in our notation translates to

$$\tilde{\Lambda} := \nabla L^\Pi(w_0)^T \tilde{\delta} + \frac{1}{2} \tilde{\delta}^T \mathcal{H} \tilde{\delta}, \quad \Lambda = \nabla L^\Pi(w_0)^T \delta + \frac{1}{2} \delta^T \mathcal{H} \delta. \quad (240)$$

where

$$\tilde{\delta} = \tilde{w}_T - w_0, \quad \delta = w_T - \tilde{w}_T,$$

for  $(\tilde{w}_t)$  which is a coupled sequence with  $(w_t)$  as defined in Ge et al (2015a). Note that we apply the coupling for the sequence  $\tilde{w}_t$  as well, i.e. if  $w_{t+1} = w_t$ , we also set  $\tilde{w}_{t+1} = \tilde{w}_t$ .

To show Equation (238), we want to use Equation (67) in Ge et al (2015a), though we only use the expansion for the first term which is under  $\mathbb{1}\{\mathcal{E}_T\}$ , i.e.

$$\mathbb{E}[L^\Pi(w_T) - L^\Pi(w_0)]\mathbb{1}_{\mathcal{E}_T} = \mathbb{E}\tilde{\Lambda}\mathbb{1}_{\mathcal{E}_T} + \mathbb{E}\Lambda\mathbb{1}_{\mathcal{E}_T}. \quad (241)$$

First of all, as it is mentioned in Lemma 16 of (Ge et al, 2015a), in the case where the noise vector  $\sigma_1^2 I \leq \mathbb{E}\mathcal{L}\mathcal{L}^T \leq \sigma_2^2 I$  instead of having  $\mathbb{E}\mathcal{L}\mathcal{L}^T = \sigma^2 I$  for a fixed  $\sigma$ , in order to still get a negative term of order  $\eta$  in Equation (68) of (Ge et al, 2015a), we just need the size of  $T_{\max}$  to be as large as  $O(\frac{1}{\gamma\eta}(\log d + \log \frac{\sigma_2}{\sigma_1}))$ , and it does not change the order of  $\eta$  in any other part of Lemma 16. Now similar to Equation (68) of (Ge et al, 2015a), if w.l.o.g we assume the smallest eigenvalue  $\gamma_0$  corresponds to  $i = 1$ :

$$\mathbb{E}\tilde{\Lambda}\mathbb{1}_{\mathcal{E}_T} \leq \frac{1}{2} \sum_{i=1}^N \lambda_i \sum_{\tau=0}^{T-1} \mathbb{1}_{\{\lambda_i < 0\}} (1 - \eta\lambda_i)^{2\tau} \eta^2 \sigma_1^2 \mathbb{P}(\mathcal{E}_T) + \quad (242)$$

$$\frac{1}{2} \sum_{i=1}^N \lambda_i \sum_{\tau=0}^{T-1} \mathbb{1}_{\{\lambda_i \geq 0\}} (1 - \eta\lambda_i)^{2\tau} \eta^2 \sigma_2^2 \quad (243)$$

$$\leq \frac{\eta^2}{2} \left[ \sigma_2^2 \frac{N-1}{\eta} - \gamma_0 \sigma_1^2 \mathbb{P}(\mathcal{E}_T) \sum_{\tau=0}^{T-1} (1 + \eta\gamma_0)^{2\tau} \right] \leq -\frac{\eta\sigma_1^2}{2}. \quad (244)$$

where in the last line we use the fact that  $\mathbb{P}(\mathcal{E}_T) \leq 1/2$  plus the additional  $\log(\sigma_2/\sigma_1)$  factor. Second, note that our threshold  $\mathfrak{s}(\eta)$  for the size of gradient in Lemmas 24 and 26 has the same order of  $\eta$  compared to that of Lemmas 14 and 16 in Ge et al (2015a). Therefore, the arguments in Lemma 16 that considers the order of  $\eta$  and treat the other parameters as constants is true here as well. Hence, we still have Equation (69) of Ge et al (2015a) which is under the event  $\mathcal{E}_T$ . Applying it to Equation (241),

Hence, finally by a similar derivation of Equation (67) in Ge et al (2015a):

$$\mathbb{E}[L^\Pi(w_T) - L^\Pi(w_0)]\mathbb{1}\{\mathcal{E}_T\} \leq -\tilde{\Omega}(\eta). \quad (245)$$

Furthermore, combining Equations (36) and (44) we get with high probability (we use the final high probability parameter of Lemma 16 which is the result of a union bound over all the high probability arguments which is equivalent to the occurrence of  $\mathcal{E}_T$ ), i.e. when  $\mathcal{E}_T$  happens,

$$\|w_T - w_0\| \leq \tilde{O}(\eta^{\frac{1}{2}} \log \frac{1}{\eta}). \quad (246)$$

Picking  $\eta$  small enough such that for the bound above we have

$$O(\eta^{\frac{1}{2}} \log \frac{1}{\eta}) \leq \sqrt{\frac{\aleph + 3l}{\psi}} - \sqrt{\frac{\aleph + 2l}{\psi}},$$

we get for every  $\bar{w}$  in the line connecting  $w_0$  to  $w_T$ :

$$\|\bar{w}\| \leq \sqrt{\frac{\aleph + 3l}{\psi}},$$

which implies that  $L^\Pi$  has the smoothing parameters  $\rho_1, \rho_2, \rho_3$  along  $w_0$  to  $w_T$ . Therefore, by the  $\rho_2$ -gradient smoothness property of  $L^\Pi$ :

$$\|\nabla L^\Pi(\bar{w}) - \nabla L^\Pi(w_0)\| \leq \rho_2 \|w_0 - \bar{w}\| \leq O(\rho_2 \eta^{\frac{1}{2}} \log \frac{1}{\eta}).$$

Combining the assumption of the Lemma  $\|\nabla L^\Pi(w_0)\| \leq \tilde{O}(\eta^{\frac{1}{2}})$ , we get

$$\|\nabla L^\Pi(\bar{w})\| = \tilde{O}(\eta^{1/2} \log(1/\eta)).$$

(The last  $\tilde{O}$  also hides the dependency on  $\rho_2$ ). Now integrating over the derivative along the direction from  $w_0$  to  $w_T$ :

$$L^\Pi(w_T) = L^\Pi(w_0) + \int_{t=0}^1 \nabla L^\Pi(tw_0 + (1-t)w_T)^T (w_T - w_0) dt.$$

Therefore, using (246) one more time, under the event  $\mathcal{E}_T$ :

$$\begin{aligned} |L^\Pi(w_T) - L^\Pi(w_0)| &\leq \int_0^1 \left| \nabla L^\Pi(tw_0 + (1-t)w_T)^T (w_T - w_0) \right| dt \\ &\leq \int_0^1 \|\nabla L^\Pi(tw_0 + (1-t)w_T)\| \|w_T - w_0\| dt \\ &\leq \tilde{O}(\eta^{1/2} \log 1/\eta) \|w_T - w_0\| = \tilde{O}(\eta \log^2 1/\eta). \end{aligned}$$

Hence

$$\left| [L^\Pi(w_T) - L^\Pi(w_0)] \mathbb{1}\{\mathcal{E}_T\} \right| \leq \tilde{O}(\eta \log^2 1/\eta), \quad (247)$$

which

Now comparing Equations (245) and (247), it is obvious that one can pick  $\eta$  small enough (again polynomially small in the other parameters) such that for some random variable  $\mathfrak{N}_2$ , which also depends on  $\eta$  Equations (238) and (239) hold. Furthermore, note that the bound in (247) is an a.s. upper bound on the change of the function value under the event  $\mathcal{E}_T$  for every  $1 \leq t \leq T$ . Therefore, by picking  $\eta$  small enough (polynomially) s.t. the quantity  $O(\eta \log(1/\eta)^2)$  in Equation (247) is bounded by  $l$ , we again make sure that the value of function during these steps changes by at most compared to  $w_0$ , i.e. for every  $1 \leq t \leq T$ :

$$\left| [L^\Pi(w_t) - L^\Pi(w_0)] \mathbb{1}\{\mathfrak{E}_t\} \right| \leq l \quad (248)$$

,therefore remains bounded by  $\aleph + 3l$ . This also implies that  $\mathfrak{E}_T \subseteq \mathcal{E}_T$  as promised.

#### A.17 PROCESS FROM A HIGHER VIEW: DEFINITION OF THE $(X)$ SEQUENCE

The goal is to find a  $w^*$  with  $L^\Pi(w^*) \leq \aleph_\ell$  using Lemmas 24 and 26. For this purpose we define a coupled sequence this way: First, as done in Ge et al (2015a), define a sequence of times  $\tau_i$  inductively in the following way: To define  $\tau_{i+1}$  based on  $\tau_i$ , if the condition

$$\aleph_\ell \leq L^\Pi(w_{\tau_{i+1}}) \leq \aleph + 2l \quad (249)$$

does not hold, then just set  $\tau_{i+1} = \tau_i \star (1)$ . Otherwise, using the conditions (228), we are either in the situation of Lemma 24 or Lemma 26, setting  $w_0 = w_{\tau_i}$ . In the first case, define  $\tau_{i+1} = \tau_i + 1 \star (2)$ . In the latter case, Let  $\mathfrak{E}_T$  be the same high probability event that we consider in Lemma (26), which happens when the aggregate behavior of the noise vectors does not behave oddly, so that we remain close to the starting point  $w_0$ . Note that from Lemma 26, we know  $\mathbb{P}(\mathfrak{E}_T) \geq 1 - O(\eta^2)$ . Now if the event  $\mathfrak{E}_T$  happens, define  $\tau_{i+1} := \tau_i + T \star (3)$ , for  $T$  also from Lemma 26 and defined originally in Lemma 16 of Ge et al (2015a), while otherwise define  $\tau_{i+1} = \tau_i \star (4)$  and, moreover, define all the rest of  $\tau_{i'}$ 's equal to  $\tau_i$ :  $\tau_{i'} = \tau_i$  for every  $i' \geq i$ . At the same time, we define the event  $\mathcal{G}_i$ , where  $\mathcal{G}_i$  happens in the case  $\star(4)$ , and  $\mathcal{G}_{i+1}$  happens in case  $\star(4)$ . Also,  $\mathcal{G}_i$  happens if any of the previous  $\mathcal{G}_{i'}$ 's happen for  $i' < i$ ; in other words,  $\mathcal{G}_i$  is included in  $\mathcal{G}_{i+1}$ . We use these events to bound the probability that the process remains above  $\aleph_\ell$ . Moreover, define the sequence of random variables  $(X_i)$  as  $X_i := L^\Pi(w_{\tau_i})$ . Note that by Lemma 25 and Equation (248) in Lemma 26, we have  $\aleph_\ell - l \leq X_i \leq \aleph + 3l$ . The key idea behind defining  $X_i$ 's is that we want to bound the MGF of  $L^\Pi(w_t)$ , without worrying about falling out of the assumptions of Lemmas 24 and 26. With the definition of  $(X_i)$  and  $\mathcal{G}_i$ , we are ready to state the theorem which roughly says the sequence  $\tau_i$  will most likely stop after a number of steps.

**Lemma 27** Let  $\mathcal{Q}_R := \bigcup_{i=1}^{\infty} (\{\tau_i \leq R\} \cap \bar{\mathcal{G}}_i)$ . Then, for some

$$R = \frac{O(\log(1/\delta_1))(\aleph + 3l)}{\theta\eta^2}, \quad (250)$$

we have  $\mathbb{P}(\mathcal{Q}_R) \leq \delta_1$ . In other words, after  $R$  iterations of PSGD, the defined sequence  $(X_i)$  above has either been in situation  $\star(1)$  or  $\star(4)$ .

**Proof** By Equations (229) and (245) in Lemmas 24 and 26, there exist a constant  $\theta$  depending polynomially on all parameters except  $\eta$  such that

$$\mathbb{E}[X_{i+1} - X_i | \bar{\mathcal{G}}_i] \leq -(\tau_{i+1} - \tau_i)\eta^2. \quad (251)$$

Now for some constant  $C$  that we specify later, define the random time  $\iota$  as the largest  $i$  where  $\tau_i \leq C/\eta^2$ . Using the fact that  $\mathcal{G}_{i-1} \subseteq \mathcal{G}_i$ , for every  $i$  we have a.s.:

$$X_{i+1}\mathbb{1}\{\bar{\mathcal{G}}_i\} - X_i\mathbb{1}\{\bar{\mathcal{G}}_{i-1}\} = \mathbb{1}\{\mathcal{G}_i - \mathcal{G}_{i-1}\}(-X_i) + (X_{i+1} - X_i)\mathbb{1}\{\bar{\mathcal{G}}_i\}.$$

Now summing this for  $i = 1$  to  $\iota$ , taking expectation from both sides and using (251):

$$\begin{aligned} \mathbb{E}X_{\iota+1}\mathbb{1}\{\bar{\mathcal{G}}_\iota\} - X_0 &= \sum_{i=1}^{\infty} \mathbb{E}\mathbb{1}\{\mathcal{G}_i - \mathcal{G}_{i-1}\}(-X_i)\mathbb{1}\{\iota \geq i\} \\ &\quad + \sum_{i=1}^{\infty} (X_{i+1} - X_i)\mathbb{1}\{\bar{\mathcal{G}}_i \cap \{\iota \geq i\}\} \\ &\leq \sum_{i=1}^{\infty} \mathbb{E}(\mathbb{1}\{\mathcal{G}_i\} - \mathbb{1}\{\mathcal{G}_{i-1}\})(-X_i) \\ &\quad + \sum_{i=1}^{\infty} \mathbb{E}(X_{i+1} - X_i \mid \bar{\mathcal{G}}_i \cap \{\iota \geq i\})\mathbb{P}(\bar{\mathcal{G}}_i \cap \{\iota \geq i\}) \\ &\leq \sup_i \sup |X_i| \\ &\quad + \sum_{i=1}^{\infty} \mathbb{E}(-(\tau_{i+1} - \tau_i)\eta^2 \mid \bar{\mathcal{G}}_i \cap \{\iota \geq i\})\mathbb{P}(\bar{\mathcal{G}}_i \cap \{\iota \geq i\}) \\ &= \sup_i \sup |X_i| - \eta^2 \sum_{i=1}^{\infty} \mathbb{E}(\tau_{i+1} - \tau_i)\mathbb{1}\{\bar{\mathcal{G}}_i \cap \{\iota \geq i\}\}. \end{aligned}$$

Now using Lemma 25, we know that in except when  $\mathcal{E}_T$  happens (in which we stop the time sequence  $\tau_i$ ), the increments of  $X_i$  are at most  $l$ . Therefore, the value of  $X_i$ 's always remain bounded by  $\aleph + 3l$ , hence:

$$LHS \leq \aleph + 3l - \theta\eta^2 \sum_{i=1}^{\infty} \mathbb{E}(\tau_{i+1} - \tau_i)\mathbb{1}\{\bar{\mathcal{G}}_i \cap \{\iota \geq i\}\}.$$

Also, by restricting the integration of the second term to the part  $\bigcup_{i=1}^{\infty} (\bar{\mathcal{G}}_i \cap \{\tau_i \geq 2C/\eta^2\})$  of the sample space, we know that under the event  $\{\iota \geq i\}$ ,  $\bar{\mathcal{G}}_i$  automatically happens (it is easy to check). Therefore:

$$\begin{aligned} LHS &\leq \aleph + 3l - \theta\eta^2 \mathbb{E}\mathbb{1}\left\{\bigcup_{i=1}^{\infty} (\bar{\mathcal{G}}_i \cap \{\tau_i \geq C/\eta^2\})\right\} \sum_{i=1}^{\infty} (\tau_{i+1} - \tau_i)\mathbb{1}\{\bar{\mathcal{G}}_i \cap \{\iota \geq i\}\} \\ &= \aleph + 3l - \theta\eta^2 \mathbb{E}\mathbb{1}\left\{\bigcup_{i=1}^{\infty} (\bar{\mathcal{G}}_i \cap \{\tau_i \geq C/\eta^2\})\right\} \sum_{i=1}^{\infty} (\tau_{i+1} - \tau_i)\mathbb{1}\{\iota \geq i\} \\ &= \aleph + 3l - \theta\eta^2 \mathbb{E}\mathbb{1}\left\{\bigcup_{i=1}^{\infty} (\bar{\mathcal{G}}_i \cap \{\tau_i \geq C/\eta^2\})\right\} \sum_{i=1}^{\iota} (\tau_{i+1} - \tau_i) \\ &= \aleph + 3l - \theta\eta^2 \mathbb{E}\mathbb{1}\left\{\bigcup_{i=1}^{\infty} (\bar{\mathcal{G}}_i \cap \{\tau_i \geq C/\eta^2\})\right\} \tau_{\iota+1}. \end{aligned}$$

Now by the definition of  $\iota$ ,  $\tau_{\iota+1} \geq C/\eta^2$ :

$$\begin{aligned} LHS &\leq \aleph + 3l - \theta\eta^2 \mathbb{E}\mathbb{1}\left\{\bigcup_{i=1}^{\infty} (\bar{\mathcal{G}}_i \cap \{\tau_i \geq C/\eta^2\})\right\} (C/\eta^2) \\ &\quad \aleph + 3l - C\theta\mathbb{P}\left(\bigcup_{i=1}^{\infty} (\bar{\mathcal{G}}_i \cap \{\tau_i \geq C/\eta^2\})\right). \end{aligned}$$

But note that  $X_i$ 's are a.s. bounded between 0 and  $\aleph + 3l$ , which implies the LHS above is at least  $-(\aleph + 3l)$ . Therefore, we finally get:

$$\mathbb{P}\left(\bigcup_{i=1}^{\infty} \left(\bar{\mathcal{G}}_i \cap \{\tau_i \geq C/\eta^2\}\right)\right) \leq \frac{2\aleph + 6l}{C\theta}.$$

Therefore, by picking  $C^* = 2(2\aleph + 6l)/\theta$ :

$$\mathbb{P}\left(\bigcup_{i=1}^{\infty} \left(\bar{\mathcal{G}}_i \cap \{\tau_i \geq C^*/\eta^2\}\right)\right) \leq \frac{1}{2}. \quad (252)$$

Note that the differences between  $\tau_i$ 's is at most  $T_{\max} = \tilde{O}(1/\eta)$  [Ge et al \(2015a\)](#). Hence, again for  $\eta$  polynomially small in other parameters, (252) implies that for  $R = \frac{2C^*}{\eta^2}$ , there exist  $\tilde{R} = \text{poly}(\cdot)$  such that after  $\tilde{R}$  iterations on the main sequence  $(w_t)$ , the corresponding sequence  $(\tau_i)$  has either been in  $\star(1)$  or  $\star(4)$  with chance at least  $1/2$ . Repeating this argument  $\log(1/\delta_1)$  times (using the markov property of the process) we conclude the proof.

#### A.18 BOUNDING THE MGF OF $X_i$ 'S

Next, we want to exploit  $X_i$ 's to bound the upward deviation of  $L^\Pi(w_t)$ . For a fix  $\theta$  the goal here is to bound  $\mathbb{E}[\exp\{\theta X_i\}]$  (this is a different  $\theta$ !). More precisely, let  $F_t$  be the sub-sigma field generated by variables  $w_i$  from time zero to  $t$ , and  $\mathcal{F}_i := F_{\tau_i}$  be the sigma field of the stop time  $\tau_i$ . Then, obviously,  $X_i$  is measurable w.r.t  $\mathcal{F}_i$ . We prove the following theorem:

**Theorem 8** *For any  $\theta > 0$ , the sequence  $(\mathbb{E}e^{\theta(X_i - X_0)})_{i=1}^{\infty}$  is a supermartingale with respect to the filtration  $(\mathcal{F}_i)$ ,*

**Proof** We proceed inductively by jointly conditioning on the previous  $X_i$  and whether  $\mathcal{G}_i$  has happened or not, and whether we are in situation  $\star(2)$  or  $\star(3)$ . We have

$$\begin{aligned} & \mathbb{E}[\exp\{\theta(X_{i+1} - X_0)\} | \mathcal{F}_i] \\ &= \mathbb{E}[\exp\{\theta(X_{i+1} - X_i + X_i - X_0)\} \mathbb{1}\{\mathcal{G}_i\} | \mathcal{F}_i] \\ &+ \mathbb{E}[\exp\{\theta(X_{i+1} - X_i + X_i - X_0)\} \mathbb{1}\{\bar{\mathcal{G}}_i \cap \star(2)\} | \mathcal{F}_i] \\ &+ \mathbb{E}[\exp\{\theta(X_{i+1} - X_i + X_i - X_0)\} \mathbb{1}\{\bar{\mathcal{G}}_i \cap \star(3)\} | \mathcal{F}_i]. \end{aligned}$$

Now by the a.s. bounds of Lemmas 24 and 26:

$$\begin{aligned} \mathbb{E}[X_{i+1} - X_i | \bar{\mathcal{G}}_i, w_{\tau_i} \text{ s.t. } \star(2)] &= -\mathfrak{R}_1^2 \eta^2, \\ \mathbb{E}[X_{i+1} - X_i | \bar{\mathcal{G}}_i, w_{\tau_i} \text{ s.t. } \star(3)] &= -\mathfrak{R}_2^2 \eta^2, \\ (X_{i+1} - X_i) \mathbb{1}\{\mathcal{G}_i\} &= 0. \text{ (a.s.)} \end{aligned}$$

Where  $\mathfrak{R}_1$  and  $\mathfrak{R}_2$  are r.v. defined in Lemmas 24 and 26 and are clearly  $\mathcal{F}_i$  measurable. This implies

$$\begin{aligned} \mathbb{E}[(X_{i+1} - X_i) \mathbb{1}\{\bar{\mathcal{G}}_i, w_{\tau_i} \text{ s.t. } \star(2)\} | \mathcal{F}_i] &= -\mathfrak{R}_1^2 \eta^2 \mathbb{1}\{\bar{\mathcal{G}}_i, w_{\tau_i} \text{ s.t. } \star(2)\}, \\ \mathbb{E}[(X_{i+1} - X_i) \mathbb{1}\{\bar{\mathcal{G}}_i, w_{\tau_i} \text{ s.t. } \star(3)\} | \mathcal{F}_i] &= -\mathfrak{R}_2^2 \eta^2 \mathbb{1}\{\bar{\mathcal{G}}_i, w_{\tau_i} \text{ s.t. } \star(3)\}. \end{aligned}$$

Now we mention the following fact:

**Fact** For a  $\sigma$  subGaussian random variable  $X$  we have  $\mathbb{E}[\exp\{\theta X\}] \leq \exp\{\theta^2 \sigma^2\}$ .

Using the a.s. bounds of Lemmas 24 and 26, we get that conditioned on  $\{\bar{\mathcal{G}}_i, w_{\tau_i} \text{ s.t. } \star(2)\}$ ,  $X_{i+1} - \mathbb{E}X_{i+1}$  is a.s. bounded by  $2\eta\theta\mathfrak{R}_1/(\rho_2\chi)$ , and conditioned on  $\{\bar{\mathcal{G}}_i, w_{\tau_i} \text{ s.t. } \star(3)\}$ ,  $X_{i+1} - \mathbb{E}X_{i+1}$  is bounded by  $2\eta\theta\mathfrak{R}_2/(\rho_2\chi)$ . Therefore, using the above fact

$$\mathbb{E}\left[\exp\{\theta(X_{i+1} - \mathbb{E}(X_{i+1} | \bar{\mathcal{G}}_i, w_{\tau_i} \text{ s.t. } \star(2)))\} \middle| \bar{\mathcal{G}}_i, w_{\tau_i} \text{ s.t. } \star(2)\right] \leq \exp\{4\eta^2\theta^2\mathfrak{R}_1^2/(\rho_2\chi)\}, \quad (253)$$

$$\mathbb{E}\left[\exp\{\theta(X_{i+1} - \mathbb{E}(X_{i+1} | \bar{\mathcal{G}}_i, w_{\tau_i} \text{ s.t. } \star(3)))\} \middle| \bar{\mathcal{G}}_i, w_{\tau_i} \text{ s.t. } \star(3)\right] \leq \exp\{4\eta^2\theta^2\mathfrak{R}_2^2/(\rho_2\chi)\}, \quad (254)$$

which implies in the notation of conditional expectation on sigma field:

$$\mathbb{E}\left[\exp\{\theta(X_{i+1} - \mathbb{E}(X_{i+1}|\mathcal{F}_i))\}\mathbb{1}\{\bar{\mathcal{G}}_i, \star(2)\} \middle| \mathcal{F}_i\right] \leq \exp\{4\eta^2\theta^2\mathfrak{R}_1^2/(\rho_2\chi)\}\mathbb{1}\{\bar{\mathcal{G}}_i, \star(2)\}, \quad (255)$$

$$\mathbb{E}\left[\exp\{\theta(X_{i+1} - \mathbb{E}(X_{i+1}|\mathcal{F}_i))\}\mathbb{1}\{\bar{\mathcal{G}}_i, \star(3)\} \middle| \mathcal{F}_i\right] \leq \exp\{4\eta^2\theta^2\mathfrak{R}_2^2/(\rho_2\chi)\}\mathbb{1}\{\bar{\mathcal{G}}_i, \star(3)\}. \quad (256)$$

Now we write:

$$\begin{aligned} LHS &\leq \mathbb{E}[\exp\{\theta(X_i - X_0)\}\mathbb{1}\{\bar{\mathcal{G}}_{i-1}\} | \mathcal{F}_i] \\ &+ \mathbb{E}[\exp\{\theta(X_{i+1} - \mathbb{E}[X_{i+1} | \mathcal{F}_i])\}\exp\{\theta(\mathbb{E}[X_{i+1} | \mathcal{F}_i] - X_i)\}\exp\{\theta(X_i - X_0)\}\mathbb{1}\{\bar{\mathcal{G}}_i \cap \star(2)\} | \mathcal{F}_i] \\ &+ \mathbb{E}[\exp\{\theta(X_{i+1} - \mathbb{E}[X_{i+1} | \mathcal{F}_i])\}\exp\{\theta(\mathbb{E}[X_{i+1} | \mathcal{F}_i] - X_i)\}\exp\{\theta(X_i - X_0)\}\mathbb{1}\{\bar{\mathcal{G}}_i \cap \star(3)\} | \mathcal{F}_i] \\ &\leq \exp\{\theta(X_i - X_0)\}\mathbb{1}\{\bar{\mathcal{G}}_i\} \\ &\quad + \exp\{\theta(X_i - X_0)\}\mathbb{E}[\exp\{\theta^2\mathfrak{R}_1^2\eta^2/(\rho_2\chi)\}\exp\{-\theta(\mathfrak{R}_1^2\eta^2)\}\mathbb{1}\{\bar{\mathcal{G}}_i \cap \star(2)\} | \mathcal{F}_i] \\ &\quad + \exp\{\theta(X_i - X_0)\}\mathbb{E}[\exp\{\theta^2\mathfrak{R}_2^2\eta^2/(\rho_2\chi)\}\exp\{-\theta(\mathfrak{R}_2^2\eta^2)\}\mathbb{1}\{\bar{\mathcal{G}}_i \cap \star(3)\} | \mathcal{F}_i] \\ &\leq \exp\{\theta(X_i - X_0)\}\mathbb{E}\left[\left(\mathbb{1}\{\bar{\mathcal{G}}_i\} \right. \right. \\ &\quad \left. \left. + \exp\{\theta^2\mathfrak{R}_1^2\eta^2/(\rho_2\chi) - \theta(\mathfrak{R}_1^2\eta^2)\}\mathbb{1}\{\bar{\mathcal{G}}_i \cap \star(2)\} \right. \right. \\ &\quad \left. \left. + \exp\{\theta^2\mathfrak{R}_2^2\eta^2/(\rho_2\chi) - \theta(\mathfrak{R}_2^2\eta^2)\}\mathbb{1}\{\bar{\mathcal{G}}_i \cap \star(3)\}\right) \middle| \mathcal{F}_i\right]. \end{aligned}$$

Now setting  $\theta := 1$  and picking  $\chi \geq 1/\rho_2$ :

$$\begin{aligned} LHS &\leq \exp\{(X_i - X_0)\}\mathbb{E}\left[\mathbb{1}\{\bar{\mathcal{G}}_{i-1}\} + \mathbb{1}\{\bar{\mathcal{G}}_{i-1} \cap \star(2)\} + \mathbb{1}\{\bar{\mathcal{G}}_{i-1} \cap \star(3)\} \middle| \mathcal{F}_i\right]. \\ &= \exp\{X_i - X_0\}. \end{aligned}$$

Now by hypothesis of Induction we have

$$\mathbb{E}[\exp\{X_{i+1} - X_0\}] = \mathbb{E}[\mathbb{E}[\exp\{X_{i+1} - X_0\} | \mathcal{F}_i]] \leq \mathbb{E}[\exp\{X_i - X_0\}] \leq 1,$$

which finishes the proof of step of induction.

Now using Doob's Maximal inequality for positive supermartingales and  $R$  defined in (250):

$$\begin{aligned} &\mathbb{P}\left(\sup_{1 \leq i \leq R} (X_i - X_0) \geq z\right) \\ &= \mathbb{P}\left(\sup_{1 \leq i \leq R} \exp\{X_i - X_0\} \geq \exp\{z\}\right) \leq \mathbb{E}[\exp\{X_R - X_0\}]/\exp\{z\} \leq e^{-z}. \end{aligned} \quad (257)$$

#### A.19 PROOF OF THEOREM 7

Finally with the developed tools, we are ready to prove Theorem 7.

**Proof of Theorem 7** Starting from  $w_0 = 0$  with  $L^\Pi(w_0) \leq \aleph$ , we use Equation (257) to get  $\mathbb{P}(\sup_{1 \leq i \leq R} \exp\{X_i - X_0\} \geq \Omega(\log(1/\delta_1))) \leq \delta_1$ . Therefore, setting  $l = \Theta(\log(1/\delta_1))$  and a union bound implies with probability at least  $1 - 2\delta_1$  we should have gotten into situation  $\star(1)$  or  $\star(4)$  without the value of  $X_i$  exceeding  $\aleph + 2l$ . On the other hand, using Lemma 26 we know that  $\mathfrak{E}_T$  happens with probability at least  $1 - \tilde{O}(\eta^5)$  for every  $1 \leq t \leq R$  which is equal to  $\tau_i$  for some  $i$  and when we are in the situation of Lemma 26. As a result, the chance that even one of  $\mathfrak{E}_T$ 's happen along  $R$  iterations is at most

$$R\tilde{O}(\eta^5) = \eta^3 \frac{O(\log(1/\delta_1))(\aleph + 3l)}{\theta}.$$

But picking  $\eta$  small enough with respect to  $\log(\delta_1)$  and other parameters, we conclude that with probability at least  $1 - 3\delta$ , after  $R$  rounds, we should have got into situation  $\star(1)$  and not  $\star(4)$  and not exceeding  $\aleph + 2l$ , which means that  $X_i = L^\Pi(w_{\tau_i})$  has gotten under the threshold  $\aleph_\ell$ . Note that as soon as that happens, we terminate the algorithm. We elaborate on this more in Appendix A.10.

## A.20 GAUSSIAN SMOOTHING

In this section, we describe our smoothing scheme and the approximation that it provides which enables us to keep the signs from the case  $\eta = 0$ . Recall that we use Gaussian smoothing matrices  $V_{j,k}^\rho \sim \mathcal{N}(0, \beta_1^2/m_1)$  and  $W_{j,k}^\rho \sim \mathcal{N}(0, \beta_2^2/m_2)$ . Here, we will particularly specify lower bounds for  $\beta_1$  and  $\beta_2$  in order for our sign approximation to be precise. On the other hand, we normally prefer the smoothing noise to be as low as possible so the primary and smoothed functions are close, so we set  $\beta_1, \beta_2$  equal to their lower bounds, and use this setting in the other parts.

To begin fix one of the inputs  $x_i$ . In order to reduce and simplify the amount of notations, we refer to the sign pattern matrix (diagonal sign matrix) of both the first and second layers by  $D$  with the appropriate indices. More specifically, for the first layer, we refer to  $\text{Sgn}(W^{(0)} + W' + W^\rho)x_i$  by  $D_{\cdot,\rho}$  and  $\text{Sgn}(W^{(0)} + (1 - \eta/2)W' + W^\rho + \sqrt{\eta}W^*)x_i$  by  $D'_{\cdot,\rho}$ . Similarly, for the second layer, of course depending on the input vector, we refer to the sign matrix with respect to the matrices  $V^{(0)} + V' + V^\rho$  and  $V^{(0)} + (1 - \eta/2)V' + V^\rho + \sqrt{\eta}V^*$  by  $D_{\cdot,\rho}$  and  $D'_{\cdot,\rho,\eta}$ , respectively. We introduce two new notations as well for the output of the first layer with respect to different matrix and sign patterns:

$$x'^{(1)} := W^s D_{\cdot,\rho}(W^{(0)} + (1 - \eta)W' + W^\rho + \sqrt{\eta}V^*)x_i, \quad (258)$$

$$x'^{(2)} := W^s D'_{\cdot,\rho,\eta}(W^{(0)} + (1 - \eta)W' + W^\rho + \sqrt{\eta}V^*)x_i. \quad (259)$$

For further brevity, we sometimes refer to  $x'^{(2)}$  by  $x'$ .

Now we are ready to mention our approximation theorem regarding the smoothing and the sign changes.

**Lemma 28** *Under the conditions  $\kappa_1\sqrt{m_3} \gtrsim C_1 + \beta_1\sqrt{m_3}$  and  $m_2 \geq m_3 \log(m_2)$ , then for every  $i \in [n]$ :*

$$\begin{aligned} & \left| \mathbb{E}_{W^\rho, V^\rho} a^T D_{\cdot,\rho,\eta}(V^{(0)} + (1 - \eta)V' + V^\rho + \sqrt{\eta}V^*)W^s D_{\cdot,\rho,\eta}(W^{(0)} + (1 - \eta)W' + W^\rho + \sqrt{\eta}V^*)x_i \right. \\ & \quad \left. - a^T D_{\cdot,\rho}(V^{(0)} + (1 - \eta)V' + V^\rho + \sqrt{\eta}V^*)W^s D_{\cdot,\rho}(W^{(0)} + (1 - \eta)W' + W^\rho + \sqrt{\eta}V^*)x_i \right| \\ & \leq \eta \varrho_2^2 \beta_2^{-1} \left[ (C_1 + \sqrt{m_3}\beta_1)^2 / (\kappa_1\sqrt{m_3}) + \left[ \sqrt{m_3}m_1\beta_1 + C_1 \right] \exp\{-C_1^2/(8m_3\beta_1^2)\} \right] \\ & \quad \times \left[ \exp\{-C_2^{4/3}(\sqrt{m_2}\kappa_2)^{2/3}/(8\beta_2^2)\} + \frac{C_2^{2/3}}{(\sqrt{m_2}\kappa_2)^{2/3}} \right] \\ & \quad + \eta \left( \kappa_2\sqrt{m_2} + C_1 \right) \left( \exp\{-c_2^2/(32\beta_1^2)\} + \frac{c_2}{\kappa_1\sqrt{m_1}} \right) \frac{\varrho^2 m_3 \sqrt{m_3}}{\beta_1} := \eta \mathfrak{R}_8. \end{aligned} \quad (260)$$

**Proof of Lemma 28**

We can bound the Left hand side above as

$$\begin{aligned} & LHS \leq \\ & \left| \mathbb{E}_{W^\rho, V^\rho} a^T D_{\cdot,\rho,\eta}(V^{(0)} + (1 - \eta)V' + V^\rho + \sqrt{\eta}V^*)W^s D_{\cdot,\rho,\eta}(W^{(0)} + (1 - \eta)W' + W^\rho + \sqrt{\eta}W^*)x_i \right. \\ & \quad \left. - a^T D_{\cdot,\rho}(V^{(0)} + (1 - \eta)V' + V^\rho + \sqrt{\eta}V^*)W^s D_{\cdot,\rho,\eta}(W^{(0)} + (1 - \eta)W' + W^\rho + \sqrt{\eta}W^*)x_i \right| \\ & \quad + \left| \mathbb{E}_{W^\rho, V^\rho} a^T D_{\cdot,\rho}(V^{(0)} + (1 - \eta)V' + V^\rho + \sqrt{\eta}V^*)W^s D_{\cdot,\rho,\eta}(W^{(0)} + (1 - \eta)W' + W^\rho + \sqrt{\eta}W^*)x_i \right. \\ & \quad \left. - a^T D_{\cdot,\rho}(V^{(0)} + (1 - \eta)V' + V^\rho + \sqrt{\eta}V^*)W^s D_{\cdot,\rho}(W^{(0)} + (1 - \eta)W' + W^\rho + \sqrt{\eta}W^*)x_i \right|. \\ & := A_1 + A_2. \end{aligned} \quad (261)$$

We bound  $A_1$  and  $A_2$  separately. First, we start with  $A_1$ .

Let  $\hat{P}_i$  be the set of indices  $j$  for which  $\mathbb{1}\{|(V_j^{(0)} + V_j')x'| \leq R^* \kappa_2 \|x'\|\}$  happens. Then, from Lemma [11](#), we have  $|\hat{P}_i| \lesssim R^* m_2$ . Now for  $j \in [m_2]$ , we write

$$\begin{aligned} & \mathbb{1}\{\text{sign change in the } j\text{th neuron}\} \times |\text{amount of change}| \tag{262} \\ & \leq \mathbb{1}\{V_j^\rho x' \in (-V_j^{(0)} x' - V_j' x' + \eta V_j' x' - \sqrt{\eta} V_j^* x', -V_j^{(0)} x' - V_j' x')\} \times \frac{1}{\sqrt{m_2}} (|\eta V_j' x'| + |\sqrt{\eta} V_j^* x'|) \tag{263} \end{aligned}$$

Moreover, note that

$$\begin{aligned} |V_j' x'| &= |V_j'(x' - \phi^{(0)}(x_i))|, \\ |V_j^* x'| &= |V_j^*(x' - \phi^{(0)}(x_i))|. \end{aligned}$$

Also, because  $\|V_j'\| \leq \|V'\| \leq 2C_2$  plus using Equation [\(129\)](#), we can further upper bound the above indicator as:

$$\begin{aligned} & \leq \mathbb{1}\{V_j^\rho x' \in (-V_j^{(0)} x' - V_j' x' - (\eta \|V_j'\| + \sqrt{\eta} \|V_j^*\|) \min\{\|x'\|, \|x' - \phi^{(0)}(x_i)\|\}, -V_j^{(0)} x' - V_j' x')\} \\ & \quad \times \frac{1}{\sqrt{m_2}} (\eta \|V_j'\| + \sqrt{\eta} \|V_j^*\|) \|x' - \phi^{(0)}(x_i)\| \\ & \leq \mathbb{1}\{V_j^\rho x' \in (-V_j^{(0)} x' - V_j' x' - (2\eta C_2 + \sqrt{\eta} \rho_2 / \sqrt{m_2}) \min\{\|x'\|, \|x' - \phi^{(0)}(x_i)\|\}, -V_j^{(0)} x' - V_j' x')\} \\ & \quad \times \frac{1}{\sqrt{m_2}} (2\eta C_2 + \sqrt{\eta} \rho_2 / \sqrt{m_2}) \|x' - \phi^{(0)}(x_i)\|. \end{aligned}$$

Taking  $\sqrt{\eta} \leq \rho_2 / (2C_2 \sqrt{m_2})$ , we can further upper bound as

$$\begin{aligned} & \lesssim \mathbb{1}\{V_j^\rho x' \in (-V_j^{(0)} x' - V_j' x' - 2\sqrt{\eta} \rho_2 \min\{\|x' - \phi^{(0)}(x_i)\|, \|x'\|\} / \sqrt{m_2}, -V_j^{(0)} x' - V_j' x')\} \\ & \quad \times (\sqrt{\eta} \rho_2 / m_2) \|x' - \phi^{(0)}(x_i)\|. \end{aligned}$$

Therefore, conditioned on  $x'$ :

$$\begin{aligned} & \mathbb{E}_{V^\rho} [\mathbb{1}\{\text{sign change in the } j\text{th neuron}\} \times |\text{amount of change}| \mid x'] \leq \\ & \mathbb{P}(V_j^\rho x' \in (-V_j^{(0)} x' - V_j' x' - 2\sqrt{\eta} \rho_2 \min\{\|x' - \phi^{(0)}(x_i)\|, \|x'\|\} / \sqrt{m_2}, -V_j^{(0)} x' - V_j' x')) \\ & \quad \times (\sqrt{\eta} \rho_2 / m_2) \|x' - \phi^{(0)}(x_i)\|. \end{aligned}$$

Now notice that for  $j \notin \hat{P}_i$ , we have

$$|-V_j^{(0)} x' - V_j' x'| \geq R^* \kappa_2 \|x'\|.$$

Also, note that the variable  $V_j^\rho x'$  is gaussian with variance  $\|x'\| \beta_2 / \sqrt{m_2}$ . Therefore, conditioned on  $x'$ , for  $j \notin \hat{P}_i$ , we have (note that  $x'$  does not depend on the randomness of  $V^\rho$ ):

$$\begin{aligned} & \mathbb{P}(V_j^\rho x' \in (-V_j^{(0)} x' - V_j' x' - 2\sqrt{\eta} \rho_2 \min\{\|x' - \phi^{(0)}(x_i)\|, \|x'\|\} / \sqrt{m_2}, -V_j^{(0)} x' - V_j' x')) \\ & \lesssim \exp\{\min\{|-V_j^{(0)} x' - V_j' x' - 2\sqrt{\eta} \rho_2 \|x'\| / \sqrt{m_2}|, |-V_j^{(0)} x' - V_j' x'|\} / (\sqrt{2} \|x'\| \beta_2 / \sqrt{m_2})\}^2 \\ & \quad \times (\|x'\| \beta_2 / \sqrt{m_2})^{-1} \times (\sqrt{\eta} \rho_2 \min\{\|x' - \phi^{(0)}(x_i)\|, \|x'\|\} / \sqrt{m_2}). \end{aligned}$$

This equation follows from the fact that

$$\mathbb{P}(a \leq \mathcal{N} \leq b) \lesssim \frac{|a - b|}{\sigma} e^{\min\{a, b\} / \sigma^2}. \tag{264}$$

On the other hand, note that for  $\sqrt{\eta} \lesssim \frac{C_2^{2/3} (\sqrt{m_2} \kappa_2)^{1/3}}{\rho_2}$  we have:

$$R^* \kappa_2 \|x'\| / 2 = \frac{C_2^{2/3} (\sqrt{m_2} \kappa_2)^{1/3}}{2\sqrt{m_2}} \|x'\| \geq 2\sqrt{\eta} \rho_2 \|x'\| / \sqrt{m_2}$$

which implies

$$\begin{aligned} &\lesssim \exp\{R^* \kappa_2 \|x'\| / (2\sqrt{2} \|x'\| \beta_2 / \sqrt{m_2})\}^2 (\sqrt{\eta} \varrho_2 / \beta_2) \\ &\leq \exp\{-C_2^{4/3} (\sqrt{m_2} \kappa_2)^{2/3} / (8\beta_2^2)\} (\sqrt{\eta} \varrho_2 / \beta_2). \end{aligned}$$

On the other side, for  $j \in \hat{P}_i$ , we can write

$$\begin{aligned} &\mathbb{P}(V_j^\rho x' \in (-V_j^{(0)} x' - V_j' x' - 2\sqrt{\eta} \varrho_2 \min\{\|x' - \phi^{(0)}(x_i)\|, \|x'\|\} / \sqrt{m_2}, -V_j^{(0)} x' - V_j' x')) \\ &\lesssim (\|x'\| \beta_2 / \sqrt{m_2})^{-1} (\sqrt{\eta} \varrho_2 \min\{\|x' - \phi^{(0)}(x_i)\|, \|x'\|\} / \sqrt{m_2}) = \min\{\|x' - \phi^{(0)}(x_i)\| / \|x'\|, 1\} \sqrt{\eta} \varrho_2 / \beta_2. \end{aligned}$$

Therefore, overall using the fact that  $\|V_j^*\| \leq \varrho_2 / \sqrt{m_2}$ , we can write

$$\begin{aligned} A_1 &\lesssim \sum_{j \notin \hat{P}} \exp\{-C_2^{4/3} (\sqrt{m_2} \kappa_2)^{2/3} / (8\beta_2^2)\} (\sqrt{\eta} \varrho_2 / \beta_2) \min\{\|x' - \phi^{(0)}(x_i)\|^2 / \|x'\|, \|x' - \phi^{(0)}(x_i)\|\} \\ &\quad + \sum_{j \in \hat{P}} (\sqrt{\eta} \varrho_2 / \beta_2) \min\{\|x' - \phi^{(0)}(x_i)\| / \|x'\|, 1\} \times (\sqrt{\eta} \varrho_2 / m_2) \|x' - \phi^{(0)}(x_i)\| \\ &\lesssim \left[ m_2 \times \exp\{-C_2^{4/3} (\sqrt{m_2} \kappa_2)^{2/3} / (8\beta_2^2)\} (\sqrt{\eta} \varrho_2 / \beta_2) \times (\sqrt{\eta} \varrho_2 / m_2) \right. \\ &\quad \left. + \eta \varrho_2^2 \beta_2^{-1} \frac{|\hat{P}_i|}{m_2} \right] \min\{\|x' - \phi^{(0)}(x_i)\|^2 / \|x'\|, \|x' - \phi^{(0)}(x_i)\|\} \\ &\leq \eta \varrho_2^2 \beta_2^{-1} \min\{\|x' - \phi^{(0)}(x_i)\|^2 / \|x'\|, \|x' - \phi^{(0)}(x_i)\|\} \left[ \exp\{-C_2^{4/3} (\sqrt{m_2} \kappa_2)^{2/3} / (8\beta_2^2)\} + \frac{C_2^{2/3}}{(\sqrt{m_2} \kappa_2)^{2/3}} \right]. \end{aligned} \tag{265}$$

Next, we bound  $A_2$ . First we bound  $\mathbb{E}_{W^\rho} \|x^{(1)} - x^{(2)}\|$ . Recalling the setting  $c_2 = 2\sqrt{nm_3} C_1 / \sqrt{\lambda_0}$  and the definition of in  $P$  from Lemma [4](#), we obtain that for  $j \notin P$ , we have for all  $i \in [n]$ :

$$\begin{aligned} |W_j^{(0)} x_i| &\geq c_2 / \sqrt{m_1}, \\ |W_j' x_i| &\leq c_2 / (2\sqrt{m_1}), \end{aligned}$$

which means for  $j \notin P$ :

$$|(W_j^{(0)} + W_j') x_i| \geq c_2 / (2\sqrt{m_1}). \tag{266}$$

Also, we have

$$|P| \leq c_2 \sqrt{m_1} / \kappa_1. \tag{267}$$

Now using Equation [\(105\)](#) in Lemma [5](#), we can write for every  $i \in [n]$ :

$$\begin{aligned} val_j &:= \mathbb{1}\{\text{sign change in the } j\text{th neuron}\} \times \left| \text{amount of change} \right| \\ &\leq \mathbb{1}\{W_j^\rho x_i \in (-W_j^{(0)} x_i - W_j' x_i + \eta W_j' x_i - \sqrt{\eta} W_j^* x_i, -W_j^{(0)} x_i - W_j' x_i)\} \\ &\quad \times \frac{1}{\sqrt{m_1}} (|\sqrt{\eta} W_j^* x_i + \eta W_j' x_i|). \end{aligned}$$

Using the fact that  $\|W_j'\| \leq \|W'\|_F \leq C_1$ , and Equation [\(105\)](#) ( $\|W_j^*\| \leq \varrho \sqrt{\frac{m_3}{m_1}}$ ) and picking  $\sqrt{\eta} \leq \frac{\varrho \sqrt{m_3}}{C_1 \sqrt{m_1}}$ , we obtain

$$\begin{aligned} &\leq \mathbb{1}\{W_j^\rho x_i \in (-W_j^{(0)} x_i - W_j' x_i - \eta \|W_j'\| - \sqrt{\eta} \|W_j^*\|, -W_j^{(0)} x_i - W_j' x_i)\} \\ &\quad \times \frac{1}{\sqrt{m_1}} (\sqrt{\eta} \|W_j^*\| + \eta \|W_j'\|) \\ &\leq \mathbb{1}\{W_j^\rho x_i \in (-W_j^{(0)} x_i - W_j' x_i - 2\sqrt{\eta} \varrho \frac{\sqrt{m_3}}{\sqrt{m_1}}, -W_j^{(0)} x_i - W_j' x_i)\} \\ &\quad \times \frac{2}{\sqrt{m_1}} (\sqrt{\eta} \varrho \sqrt{\frac{m_3}{m_1}}). \end{aligned}$$

Now for  $j \notin P$ , because  $W_j^\rho x_i$  is Gaussian with std  $\frac{\beta_1}{\sqrt{m_1}}$ :

$$\begin{aligned} \mathbb{E}_{W^\rho}[\text{val}_j] &\leq \mathbb{P}(W_j^\rho x_i \in (-W_j^{(0)} x_i - W_j' x_i - 2\sqrt{\eta}\rho \frac{\sqrt{m_3}}{\sqrt{m_1}}, -W_j^{(0)} x_i - W_j' x_i)) \times \frac{1}{\sqrt{m_1}} \sqrt{\eta}\rho \frac{\sqrt{m_3}}{\sqrt{m_1}} \\ &\lesssim \exp\{-\{\min\{|-W_j^{(0)} x_i - W_j' x_i - 2\sqrt{\eta}\rho \frac{\sqrt{m_3}}{\sqrt{m_1}}|, |-W_j^{(0)} x_i - W_j' x_i|\}/(\sqrt{2}\beta_1/\sqrt{m_1})\}^2\} \\ &\quad \times (\beta_1/\sqrt{m_1})^{-1} \times (\sqrt{\eta}\rho \frac{\sqrt{m_3}}{\sqrt{m_1}}) \times \frac{1}{\sqrt{m_1}} \sqrt{\eta}\rho \frac{\sqrt{m_3}}{\sqrt{m_1}}. \end{aligned}$$

Now from Equation (266) and by picking  $\sqrt{\eta} \lesssim \frac{c_2}{\rho\sqrt{m_3}}$  so that

$$2\sqrt{\eta}\rho \frac{\sqrt{m_3}}{\sqrt{m_1}} \leq c_2/(4\sqrt{m_1}),$$

then

$$LHS \lesssim \exp\{-c_2^2/(32\beta_1^2)\} \eta \frac{\rho^2 m_3}{\beta_1 m_1}. \quad (268)$$

On the other hand, for  $j \in P$  we have

$$\begin{aligned} \text{Eval}_j &\leq \mathbb{P}(W_j^\rho x_i \in (-W_j^{(0)} x_i - W_j' x_i - 2\sqrt{\eta}\rho \frac{\sqrt{m_3}}{\sqrt{m_1}}, -W_j^{(0)} x_i - W_j' x_i)) \times \frac{1}{\sqrt{m_1}} \sqrt{\eta}\rho \frac{\sqrt{m_3}}{\sqrt{m_1}} \\ &\lesssim (\beta_1/\sqrt{m_1})^{-1} \sqrt{\eta}\rho \frac{\sqrt{m_3}}{\sqrt{m_1}} \times \sqrt{\eta}\rho \frac{\sqrt{m_3}}{m_1} = \eta \frac{\rho^2 m_3}{\beta_1 m_1}. \end{aligned} \quad (269)$$

Now define the following random variable with respect to the randomness of  $W^\rho$ :

$$\text{Val} := \sum_{j=1}^{m_1} \mathbb{1}\{\text{sign change in the } j\text{th neuron}\} \times |\text{amount of change}|$$

then for every  $k \in [m_3]$ , we have

$$|x_k^{(1)} - x_k^{(2)}| \leq \text{Val},$$

which implies

$$\|x^{(1)} - x^{(2)}\| \leq \sqrt{m_3} \text{Val}.$$

But Combining Equations (268) and (269):

$$\begin{aligned} \mathbb{E} \text{Val} &\leq \left( \exp\{-c_2^2/(32\beta_1^2)\} + \frac{|P|}{m_1} \right) \eta \frac{\rho^2 m_3}{\beta_1} \\ &\leq \left( \exp\{-c_2^2/(32\beta_1^2)\} + \frac{c_2}{\kappa_1 \sqrt{m_1}} \right) \eta \frac{\rho^2 m_3}{\beta_1}, \end{aligned}$$

which implies

$$\mathbb{E}_{W^\rho} \|x^{(1)} - x^{(2)}\| \leq \left( \exp\{-c_2^2/(32\beta_1^2)\} + \frac{c_2}{\kappa_1 \sqrt{m_1}} \right) \eta \frac{\rho^2 m_3 \sqrt{m_3}}{\beta_1}. \quad (270)$$

Now we can write

$$\begin{aligned} &\left| a^T D_{\cdot, \rho}(V^{(0)} + (1-\eta)V' + V^\rho + \sqrt{\eta}V^*)x^{(2)} - a^T D_{\cdot, \rho}(V^{(0)} + (1-\eta)V' + V^\rho + \sqrt{\eta}V^*)x^{(1)} \right| \\ &\leq \frac{1}{\sqrt{m_2}} \sum_{j=1}^{m_2} |(V_j^{(0)} + (1-\eta)V_j' + V_j^\rho + \sqrt{\eta}V_j^*)(x^{(2)} - x^{(1)})| \\ &\leq \frac{1}{\sqrt{m_2}} \sum_{j=1}^{m_2} |V_j^{(0)}(x^{(2)} - x^{(1)})| + (1-\eta)|V_j'(x^{(2)} - x^{(1)})| + |V_j^\rho(x^{(2)} - x^{(1)})| + \sqrt{\eta}|V_j^*(x^{(2)} - x^{(1)})| \\ &= \frac{1}{\sqrt{m_2}} \sum_{j=1}^{m_2} |V_j^{(0)}(x^{(2)} - x^{(1)})| + |V_j^\rho(x^{(2)} - x^{(1)})| + \frac{1}{\sqrt{m_2}} \sum_{j=1}^{m_2} ((1-\eta)\|V_j'\| + \sqrt{\eta}\|V_j^*\|) \|x^{(2)} - x^{(1)}\| \\ &\leq \frac{1}{\sqrt{m_2}} \sum_{j=1}^{m_2} |V_j^{(0)}(x^{(2)} - x^{(1)})| + |V_j^\rho(x^{(2)} - x^{(1)})| + \left( (1-\eta)\|V'\|_F + \sqrt{\eta}\|V^*\|_F \right) \|x^{(2)} - x^{(1)}\|. \end{aligned}$$

Now by Equation (B27) in Lemma B10 (i.e.  $\|V^*\|_F \lesssim \sqrt{\zeta_2}$ ) and the fact that  $\|V'\|_F \leq C_2$ , and by taking

$$\sqrt{\eta} \leq \frac{C_2}{\sqrt{\zeta_2}},$$

we have

$$\left( (1 - \eta)\|V'\|_F + \sqrt{\eta}\|V^*\|_F \right) \lesssim C_1, \quad (271)$$

so we can bound the above as

$$LHS \lesssim \frac{1}{\sqrt{m_2}} \sum_{j=1}^{m_2} \left( |V_j^{(0)}(x'^{(2)} - x'^{(1)})| + |V_j^\rho(x'^{(2)} - x'^{(1)})| \right) + C_1 \|x'^{(2)} - x'^{(1)}\|.$$

Furthermore, using Lemma B2 and noting the fact that the entries of  $V^{(0)}$  are normal with standard deviation  $\kappa_2$ , we get with high probability over the randomness of  $V^{(0)}$ :

$$\lesssim \frac{1}{\sqrt{m_2}} \sum_{j=1}^{m_2} |V_j^\rho(x'^{(2)} - x'^{(1)})| + \left( \kappa_2 \sqrt{m_2} + \kappa_2 \sqrt{m_3 (\log(m_3) + \log(\log(m_2)))} + C_1 \right) \|x'^{(2)} - x'^{(1)}\|.$$

Now note that  $V_j^\rho(x'^{(2)} - x'^{(1)})$  is normal with standard deviation  $\frac{\beta_2}{\sqrt{m_2}} \|x'^{(2)} - x'^{(1)}\|$ . Hence, taking expectation with respect to  $V^\rho$ :

$$\begin{aligned} \mathbb{E}_{V^\rho} \left| a^T D_{\cdot, \rho, \eta} (V^{(0)} + (1 - \eta)V' + V^\rho + \sqrt{\eta}V^*) x'^{(2)} - a^T D_{\cdot, \rho, \eta} (V^{(0)} + (1 - \eta)V' + V^\rho + \sqrt{\eta}V^*) x'^{(1)} \right| \\ \lesssim \left( \kappa_2 \sqrt{m_2} + \kappa_2 \sqrt{m_3 (\log(m_3) + \log(\log(m_2)))} + C_1 \right) \|x'^{(2)} - x'^{(1)}\|. \end{aligned} \quad (272)$$

Finally, Combining Equation (B63) and (B72) and applying it to Equation (B61) implies with high probability over the random initialization:

$$\begin{aligned} & \left| \mathbb{E}_{W^\rho, V^\rho} a^T D_{\cdot, \rho, \eta} (V^{(0)} + (1 - \eta)V' + V^\rho + \sqrt{\eta}V^*) W^s x'^{(2)} \right. \\ & \left. - a^T D_{\cdot, \rho} (V^{(0)} + (1 - \eta)V' + V^\rho + \sqrt{\eta}V^*) W^s x'^{(1)} \right| \\ & \leq A_1 + A_2 \\ & \lesssim \mathbb{E}_{W^\rho} \eta \varrho_2^2 \beta_2^{-1} \min\{\|x' - \phi^{(0)}(x_i)\|^2 / \|x'\|, \|x' - \phi^{(0)}(x_i)\|\} \left[ \exp\{-C_2^{4/3} (\sqrt{m_2} \kappa_2)^{2/3} / (8\beta_2^2)\} + \frac{C_2^{2/3}}{(\sqrt{m_2} \kappa_2)^{2/3}} \right] \\ & + \mathbb{E}_{W^\rho} \left( \kappa_2 \sqrt{m_2} + \kappa_2 \sqrt{m_3 (\log(m_3) + \log(\log(m_2)))} + C_1 \right) \|x'^{(2)} - x'^{(1)}\| \end{aligned}$$

Now notice that under  $E^c$ , using the assumption  $\kappa_1 \sqrt{m_3} \gtrsim C_1 + \sqrt{m_3} \beta_1$  and Lemma B3 we have

$$\begin{aligned} \|x'\| & \geq \|\phi^{(0)}(x_i)\| - \|x' - \phi^{(0)}(x_i)\| \\ & \geq \kappa_1 \sqrt{m_3} - (C_1 + \sqrt{m_3} \beta_1) \\ & \gtrsim \kappa_1 \sqrt{m_3}, \end{aligned}$$

and

$$\|x' - \phi^{(0)}(x_i)\|^2 \leq (C_1 + \sqrt{m_3} \beta_1)^2, \quad (273)$$

which implies:

$$\begin{aligned} & \mathbb{E}_{W^\rho} \min\{\|x' - \phi^{(0)}(x_i)\|^2 / \|x'\|, \|x' - \phi^{(0)}(x_i)\|\} \\ & \leq \mathbb{E}_{W^\rho} \mathbb{1}\{E^c\} \|x' - \phi^{(0)}(x_i)\|^2 / \|x'\| + \mathbb{1}\{E\} \|x' - \phi^{(0)}(x_i)\| \\ & \lesssim (C_1 + \sqrt{m_3} \beta_1)^2 / (\kappa_1 \sqrt{m_3}) + \mathbb{E}_{W^\rho} \mathbb{1}\{E\} \|x' - \phi^{(0)}(x_i)\| \\ & \lesssim (C_1 + \sqrt{m_3} \beta_1)^2 / (\kappa_1 \sqrt{m_3}) + \left[ \sqrt{m_3} m_1 \beta_1 + C_1 \right] \exp\{-C_1^2 / (8m_3 \beta_1^2)\}. \end{aligned}$$

Substituting this above and further applying the result of Lemma B3 and Equation (270) and the assumption that  $m_2 \geq m_3 \log(m_2)$ :

$$\begin{aligned} A_1 + A_2 &\lesssim \eta \varrho_2^2 \beta_2^{-1} \left[ (C_1 + \sqrt{m_3} \beta_1)^2 / (\kappa_1 \sqrt{m_3}) + \left[ \sqrt{m_3} m_1 \beta_1 + C_1 \right] \exp\{-C_1^2 / (8m_3 \beta_1^2)\} \right] \\ &\times \left[ \exp\{-C_2^{4/3} (\sqrt{m_2} \kappa_2)^{2/3} / (8\beta_2^2)\} + \frac{C_2^{2/3}}{(\sqrt{m_2} \kappa_2)^{2/3}} \right] \\ &+ \eta \left( \kappa_2 \sqrt{m_2} + C_1 \right) \left( \exp\{-c_2^2 / (32\beta_1^2)\} + \frac{c_2}{\kappa_1 \sqrt{m_1}} \right) \frac{\varrho^2 m_3 \sqrt{m_3}}{\beta_1}, \end{aligned}$$

which completes the proof.

#### A.20.1 SETTING $\beta_1$ AND $\beta_2$

As we mentioned, to minimize the amount of deviation of the smoothed function compared to the original one, we prefer to choose  $\beta_1, \beta_2$  as small as possible. (The benefit of such choice, indeed, can be observed more explicitly in other parts of the proof, e.g. Appendix A.14.) Observing the bound in Equation (260) and noting that we can easily make the exponential terms orders of magnitude smaller than the poly terms, it is easy to find the following optimal setting for the smoothing parameters:

$$\begin{aligned} \beta_2 &:= \Theta_p \left( (\kappa_1 \sqrt{m_3})^{-1} (\sqrt{m_2} \kappa_2)^{-\frac{2}{3}} \right), \\ \beta_1 &:= \Theta_p \left( m_3 \sqrt{m_3} / (\kappa_1 \sqrt{m_1}) \right). \end{aligned}$$

Using this setting, we still can make  $\mathfrak{R}_8$  arbitrarily small. Here, we remind the reader that  $O_p$  only cares about the non-logarithmic dependencies on the overparameterization, i.e.  $m_1, m_2, m_3, \kappa_1, \kappa_2$ .

## A.21 BASIC TOOLS

In this section, we introduce and prove some lemmas that we use in our analysis as basic tools.

**Lemma 29** *Suppose  $V^{(0)} \in \mathbb{R}^{m_2 \times m_3}$  has standard normal entries and  $a$  is a random sign vector. Suppose  $\theta > 1, R < 1$  are given thresholds, such that*

$$m_2 R \gtrsim m_3 (\log(1/R) + \log(m_3) + \log(\log(m_2))),$$

$$e^{-\theta^2/8} \lesssim m_3/m_2.$$

Then, for the following quantities:

$$\begin{aligned} N_R^1(x) &= \#\left(j \in [m] : |V_j^{(0)} x| \leq R\right) \\ N_\theta^2(x) &= \#\left(j \in [m] : |V_j^{(0)} x| \geq \theta\right), \end{aligned}$$

with high probability we have

$$\begin{aligned} \sup_{\|x'\|=1} N_R^1(x') &\lesssim m_2 R, \\ \sup_{\|x'\|=1} N_\theta^2(x') &\lesssim m_3 (\log(m_3) + \log(\log(m_2))). \end{aligned}$$

**Proof of Lemma 29**

Suppose  $B_1(\epsilon)$  is a cover for the Euclidean ball in  $\mathbb{R}^{m_3}$  with precision  $\epsilon$ . We know

$$|B_1(\epsilon)| \lesssim (1/\epsilon)^{m_3}.$$

Now for a fixed  $\|x\| = 1$ , we have

$$\mathbb{P}(W_j^{(0)} x \leq 2R) \lesssim R.$$

Therefore, using Bernstein, with high probability we have

$$\#\left(j \in [m_2] : |V_j^{(0)} x| \leq 2R\right) \lesssim m_2 R + \sqrt{m_2 R} + 1.$$

Hence, using union bound, we have with high probability

$$\begin{aligned} \sup_{x \in B_1(\epsilon)} \#\left(j \in [m] : |V_j^{(0)} x| \leq 2R\right) &\lesssim m_2 R + \sqrt{\log |B_1(\epsilon)|} \sqrt{m_2 R} + \log |B_1(\epsilon)| \\ &= m_2 R + \sqrt{m_2 R m_3 \log(1/\epsilon)} + m_3 \log(1/\epsilon). \end{aligned}$$

By picking

$$\epsilon \lesssim R / (\sqrt{m_3 \log(m_2 m_3)}),$$

The assumption implies  $m_2 R \geq m_3 \log(1/\epsilon)$ , which implies

$$LHS \lesssim m_2 R.$$

On the other hand, note that with high probability we have

$$\sup_{j \in [m_2], k \in [m_3]} |V_{j,k}^{(0)}| \leq \sqrt{\log(m_2 m_3)}. \quad (274)$$

Now for  $\|x'\| = 1$  which is not in the cover, if  $x$  is the closest point to it in the cover, i.e.  $x \in B_1(\epsilon)$  and  $\|x - x'\| \leq \epsilon$ , then for every  $j \in [m_2]$  we have

$$\left| |V_j^{(0)} x| - |V_j^{(0)} x'| \right| \leq \|V_j^{(0)}\| \|x - x'\| \leq \sqrt{m_3 \log(m_2 m_3)} \epsilon \leq R,$$

by picking a small enough constant. Therefore, for a  $j$  that  $|V_j^{(0)} x| \geq 2R$ , then

$$|V_j^{(0)} x'| \geq 2R - R = R.$$

Therefore, we get that with high probability, for every  $\|x'\| = 1$ :

$$\sup_{\|x'\|=1} \#(j \in [m_2] : |V_j^{(0)} x'| \leq R) \lesssim m_2 R.$$

For the second part, note that for  $\|x\| = 1$ , by the tail bound for normal vars:

$$\mathbb{P}(W_j^{(0)} x \geq \theta/2) \lesssim e^{-\theta^2/8}.$$

Hence, again using Bernstein, we have with high probability

$$\begin{aligned} \sup_{x \in B_1(\epsilon)} \#(j \in [m_2] : |V_j^{(0)} x| \geq \theta/2) &\lesssim m_2 e^{-\theta^2/8} + \sqrt{\log |B_1(\epsilon)|} \sqrt{m_2 e^{-\theta^2/8}} + \log |B_1(\epsilon)| \\ &\lesssim m_2 e^{-\theta^2/8} + \sqrt{m_3 \log(1/\epsilon)} \sqrt{m_2 e^{-\theta^2/8}} + m_3 \log(1/\epsilon). \end{aligned}$$

By picking

$$\epsilon \lesssim 1/(\sqrt{m_3 \log(m_2 m_3)}),$$

and using the assumption  $m_2 e^{-\theta^2/8} \lesssim m_3$ , all terms are dominated by the third term so we can bound the above as

$$\sup_{x \in B_1(\epsilon)} \#(j \in [m_2] : |V_j^{(0)} x| \geq \theta/2) \lesssim m_3 (\log(m_3) + \log(\log(m_2))).$$

Now for  $\|x'\| = 1$  not in the cover, for the new  $\epsilon$  we can write

$$\|V_j^{(0)} x - |V_j^{(0)} x'\| \leq \|V_j^{(0)}\| \|x - x'\| \leq \sqrt{m_3 \log(m_2 m_3)} \epsilon \leq 1/2 \leq \theta/2.$$

Hence, with high prob.

$$\sup_{\|x'\|=1} \#(j \in [m_2] : |V_j^{(0)} x| \geq \theta) \lesssim m_3 (\log(m_3) + \log(\log(m_2))).$$

**Lemma 30** For  $x \in \mathbb{R}^d$  and  $W^{(0)} \in \mathbb{R}^{m \times d}$  which has standard normal entries (and  $a$  is a random sign vector), we have with high probability:

$$\sup_{\|x\|=1} f(x) := \frac{1}{\sqrt{m}} a^T \sigma(W^{(0)} x) \leq \sqrt{d}.$$

### Proof of Lemma 30

For the first part, we first compute an upper bound on

$$\mathbb{E} \sup_{\|x\|=1} \frac{1}{\sqrt{m}} a^T \sigma(W^{(0)} x).$$

To do so, we use Dudley's chaining. Note that for  $x_1, x_2 \in \mathbb{R}^d$ , the variable  $\sigma(W_j^{(0)} x_1) - \sigma(W_j^{(0)} x_2)$  is subGaussian with parameter  $\|x_1 - x_2\|$ , so the variable  $f(x_1) - f(x_2)$  is also subGaussian with parameter  $\|x_1 - x_2\|$ . Hence, by Dudley's integral:

$$\mathbb{E} \sup_{\|x\|=1} \frac{1}{\sqrt{m}} a^T \sigma(W^{(0)} x) \leq \int_0^1 \sqrt{\log(\mathcal{N}(\mathcal{B}_1^d, \epsilon))} \lesssim \sqrt{d}.$$

Now for a fixed  $x$ , note that

$$\frac{1}{\sqrt{m}} a^T \sigma(W_1 x) - \frac{1}{\sqrt{m}} a^T \sigma(W_2 x) \leq \frac{1}{\sqrt{m}} \sum_{j=1}^m \|W_{1j} - W_{2j}\| \leq \|W_1 - W_2\|_F.$$

Hence, the function  $f(x)$  is 1-lipchitz with respect to  $W$  and  $l_2$  norm, so is the function  $\sup f(x)$ . Hence, by Gaussian concentration,  $\sup f(x)$  is 1-subGaussian around its mean, so we finally get with high probability

$$\sup f(x) \lesssim \sqrt{d} + 1 \lesssim \sqrt{d}.$$

**Lemma 31** For

$$R^* := \frac{C_2^{2/3}}{(\sqrt{m_2}\kappa_2)^{2/3}},$$

we have with high probability over the randomness of  $V^{(0)}$ :

$$\sup_{x', V': \|V'\| \leq C_2} \# \left( j \in [m_2] : |(V_j^{(0)} + V_j')x'| \leq R^* \kappa_2 \|x'\| \right) \lesssim R^* m_2.$$

**Proof of Lemma 31**

Note that obviously the condition of Lemma 29 is satisfied with this choice of  $R = R^*$ . Therefore, with high probability we have for an arbitrary  $x'$ :

$$\# \left( |V_j^{(0)}x'| \leq 2R^* \kappa_2 \|x'\| \right) \leq m_2 R^*.$$

On the other hand, note that for  $j \in [m_2]$  such that  $|V_j'x'| \geq R\kappa_2 \|x'\|$ , we have

$$\|V_j'\| \|x'\| \geq |V_j'x'| \geq R\kappa_2 \|x'\|,$$

which implies

$$\|V_j'\| \geq R\kappa_2.$$

Therefore, there are at most  $\frac{C_2^2}{R^2\kappa_2^2}$ . Therefore, setting aside  $m_2 R + \frac{C_2^2}{R^2\kappa_2^2}$  of  $j$ 's, for the rest we have

$$|(V_j^{(0)} + V_j')x'| \geq |V_j^{(0)}x'| - |V_j'x'| \geq 2R\kappa_2 \|x'\| - R\kappa_2 \|x'\| = R\kappa_2 \|x'\|.$$

Setting  $R^*$  as defined above balances the terms  $m_2 R$  and  $\frac{C_2^2}{R^2\kappa_2^2}$ , which completes the proof.

**Lemma 32** If  $V^{(0)} \in \mathbb{R}^{m_2 \times m_3}$  is a matrix with standard normal entries, then with high probability

$$\sup_{\|x'\|=1} \frac{1}{\sqrt{m_2}} \sum_{j=1}^{m_2} |V_j^{(0)}x'| \lesssim \sqrt{m_2} + \sqrt{m_3(\log(m_3) + \log(\log(m_2)))}.$$

**Proof of Lemma 32**

Let  $B_1(\epsilon)$  be a cover for the unit Euclidean ball with precision  $\epsilon$ , for which we have  $|B_1(\epsilon)| \lesssim (\frac{1}{\epsilon})^{m_3}$ .

Now for a fixed  $x \in B_1(\epsilon)$ , note that because  $V_j^{(0)}x$  is a standard normal variable, the random variable  $|V_j^{(0)}x'| - \mathbb{E}|V_j^{(0)}x'|$  is  $O(1)$ -subGaussian, which means  $\frac{1}{\sqrt{m_2}} \sum_{j=1}^{m_2} (|V_j^{(0)}x'| - \mathbb{E}|V_j^{(0)}x'|)$  is also  $O(1)$ -subGaussian. Now from the tail of maximum of subGaussian variables:

$$\sup_{x \in B_1(\epsilon)} \frac{1}{\sqrt{m_2}} \sum_{j=1}^{m_2} (|V_j^{(0)}x'| - \mathbb{E}|V_j^{(0)}x'|) \lesssim \sqrt{\log(|B_1(\epsilon)|)} = \sqrt{m_3 \log(1/\epsilon)}.$$

On the other hand, note that  $\mathbb{E}|V_j^{(0)}x'| = O(1)$ , which implies w.h.p:

$$\sup_{x \in B_1(\epsilon)} \frac{1}{\sqrt{m_2}} \sum_{j=1}^{m_2} |V_j^{(0)}x'| \lesssim \sqrt{m_2} + \sqrt{m_3 \log(1/\epsilon)}.$$

Moreover, note that again by the tail of subGaussian variables, we have w.h.p:

$$\max_{j \in [m_2], k \in [m_3]} |V_{j,k}^{(0)}| \lesssim \sqrt{\log(m_2 m_3)},$$

which implies with high prob for every  $j \in [m_2]$ :

$$\|V_j^{(0)}\| \lesssim \sqrt{m_3 \log(m_2 m_3)}.$$

Now by picking

$$\epsilon := \left( \sqrt{m_3 \log(m_2 m_3)} \right)^{-1},$$

we get with high probability

$$\sup_{x \in B_1(\epsilon)} \frac{1}{\sqrt{m_2}} \sum_{j=1}^{m_2} |V_j^{(0)} x| \lesssim \sqrt{m_2} + \sqrt{m_3 (\log(m_3) + \log(\log(m_2)))}. \quad (275)$$

On the other hand, for an arbitrary  $x'$  with  $\|x'\| = 1$ , if  $x \in B_1(\epsilon)$  is the representative of  $x'$ , we have by definition  $\|x' - x\| \leq \epsilon$ , which combined with (A.21) implies

$$\begin{aligned} \left| |V_j^{(0)} x'| - |V_j^{(0)} x| \right| &\leq |V_j^{(0)}(x' - x)| \leq \|V_j^{(0)}\| \|x' - x\|, \\ &\lesssim \sqrt{m_3 \log(m_2 m_3)} \left( \sqrt{m_3 \log(m_2 m_3)} \right)^{-1} \leq 1. \end{aligned}$$

Therefore

$$\left| \frac{1}{\sqrt{m_2}} \sum_{j=1}^{m_2} |V_j^{(0)} x'| - \frac{1}{\sqrt{m_2}} \sum_{j=1}^{m_2} |V_j^{(0)} x| \right| \leq \sqrt{m_2}. \quad (276)$$

Combining Equations (275) and (276), we conclude the result.

A.21.1 DEFINING THE RARE EVENTS  $E_j$ 

**Lemma 33** For  $x^{(2)}$  defined in Equation (259) we have

$$\begin{aligned}\mathbb{E}_{W^\rho} \|x^{(2)} - \phi^{(0)}(x_i)\|, \mathbb{E}_{W^\rho} \|\phi^{(2)}(x_i)\| &\leq C_1 + \sqrt{m_3}\beta_1, \\ \mathbb{E}_{W^\rho} \|\phi^{(2)}(x_i)\|^2 &\lesssim C_1^2 + m_3\beta_1^2.\end{aligned}$$

Moreover, for the events

$$E_j = \{|W_j^\rho x_i| \geq C_1/\sqrt{m_3 m_1}\}, \quad E = \cup_j E_j,$$

we have under  $E^c$ :

$$\|x^{(2)} - \phi^{(0)}(x_i)\|, \|\phi^{(2)}(x_i)\| \lesssim C_1.$$

Furthermore,  $E$  happens rarely:

$$\begin{aligned}\mathbb{P}(E) &\lesssim m_1 \exp\{-C_1^2/(8m_3\beta_1^2)\}, \\ \mathbb{E}_{W^\rho} \mathbb{1}\{E\} \|\phi^{(2)}(x_i)\| &\leq \left[\sqrt{m_3}m_1\beta_1 + C_1\right] \exp\{-C_1^2/(8m_3\beta_1^2)\}. \\ \mathbb{E}_{W^\rho} \mathbb{1}\{E\} \|x^{(2)} - \phi^{(0)}(x_i)\| &\leq \left[\sqrt{m_3}m_1\beta_1 + C_1\right] \exp\{-C_1^2/(8m_3\beta_1^2)\}.\end{aligned}$$

Finally, we have the following almost surely bound:

$$\|\phi^{(2)}(x_i)\| \leq C_1 + \sqrt{m_3} \sum_{j=1}^{m_1} \frac{1}{\sqrt{m_1}} |W_j^\rho x_i|.$$

**Proof of Lemma 33**

We start by writing

$$|x_k^{(2)} - \frac{1}{\sqrt{m_1}} W_k^s \sigma(W^{(0)} + (1-\eta)W' + \sqrt{\eta}W^*)x_i| \leq \sum_{j=1}^{m_1} \frac{1}{\sqrt{m_1}} |W_j^\rho x_i|. \quad (277)$$

Now notice that by Lemma 11, we know for every  $j \notin \tilde{P}$ :

$$|W_j^{(0)} x_i| \geq c_2/\sqrt{m_2}, \quad (278)$$

$$|(1-\eta)W_j' x_i| \leq c_2/(2\sqrt{m_1}). \quad (279)$$

In addition, by Equations in (105) from Lemma 5, for every  $j \in [m_1]$ :

$$\|W_j^*\| \leq \varrho_1 \sqrt{\frac{m_3}{m_1}},$$

so by picking

$$\eta \leq c_2/(4\varrho\sqrt{m_3})$$

we obtain

$$\eta |W_j^* x_i| \leq \frac{c_2}{4\sqrt{m_1}}. \quad (280)$$

Combining this with Equations in (279), we see that the signs of  $(W_j^{(0)} + (1-\eta)W_j' + \sqrt{\eta}W_j^*)x_i$  and  $W_j^{(0)} x_i$  are the same for  $j \notin \tilde{P}$ .

Moreover, the matrix  $(1-\eta)W' + \sqrt{\eta}W^*$  satisfies

$$\|(1-\eta)W' + \sqrt{\eta}W^*\| \leq (1-\eta)C_1 + \sqrt{\eta}\sqrt{2\zeta_2} \leq C_1,$$

by picking  $\sqrt{\eta} \leq C_1/\sqrt{\zeta_2}$ . Hence, the conditions of Lemma 6 are satisfied and we get:

$$\left\| \frac{1}{\sqrt{m_1}} W^s \sigma(W^{(0)} + (1-\eta)W' + \sqrt{\eta}W^*)x_i - \phi^{(0)}(x_i) \right\| \leq C_1. \quad (281)$$

Combining Equations (277) and (281):

$$\|x'^{(2)} - \phi^{(0)}(x_i)\| \leq C_1 + \sqrt{m_3} \sum_{j=1}^{m_1} \frac{1}{\sqrt{m_1}} |W_j^\rho x_i|. \quad (282)$$

In exactly similar fashion, one can derive

$$\|\phi^{(2)}(x_i)\| \leq C_1 + \sqrt{m_3} \sum_{j=1}^{m_1} \frac{1}{\sqrt{m_1}} |W_j^\rho x_i|. \quad (283)$$

Now first of all, note

$$\mathbb{E}_{W^\rho} \sum_{j=1}^{m_1} \frac{1}{\sqrt{m_1}} |W_j^\rho x_i| \leq \beta_1,$$

which proves the first part of the claims. For the second part, note that by the Gaussian tail bound

$$\mathbb{P}(|W_j^\rho x_i| \geq C_1/\sqrt{m_3 m_1}) \lesssim \exp\{-C_1^2/(8m_3\beta_1^2)\}.$$

Therefore,

$$\mathbb{P}(E) \leq \sum_j \mathbb{P}(E_j) \leq m_1 \exp\{-C_1^2/(8m_3\beta_1^2)\}.$$

Moreover

$$\begin{aligned} \mathbb{E}_{W^\rho} \mathbb{1}\{E\} \sum_{j=1}^{m_1} \frac{1}{\sqrt{m_1}} |W_j^\rho x_i| &\leq \mathbb{E}_{W^\rho} \frac{1}{\sqrt{m_1}} \sum_{j_2} \sum_{j \neq j_2} \mathbb{1}\{E_{j_2}\} |W_j^\rho x_i| + \frac{1}{\sqrt{m_1}} \sum_j \mathbb{E}_{W^\rho} \mathbb{1}\{E_j\} |W_j^\rho x_i|, \\ &= \left[ \frac{1}{\sqrt{m_1}} \sum_{j_2} \sum_{j \neq j_2} \mathbb{E}_{W^\rho} |W_j^\rho x_i| + \frac{1}{\sqrt{m_1}} \sum_j \mathbb{E}_{W^\rho} [|W_j^\rho x_i| \mid E_j] \right] \mathbb{P}(E_j) \\ &\lesssim \left[ m_1 \beta_1 + C_1/\sqrt{m_3} \right] \exp\{-C_1^2/(8m_3\beta_1^2)\}. \end{aligned}$$

Plugging this into (282) finishes the proof. Also, under  $E^c$  by Equation (282) we have

$$\|x'^{(2)} - \phi^{(0)}(x_i)\|, \|\phi^{(2)}(x_i)\| \lesssim C_1.$$

Finally, exploiting Equation (283):

$$\begin{aligned} \mathbb{E}_{W^\rho} \|\phi^{(2)}(x_i)\|^2 &\lesssim C_1^2 + m_3 \frac{1}{m_1} \mathbb{E}_{W^\rho} \left( \sum_j |W_j^\rho x_i| \right)^2 \leq C_1^2 + m_3 \mathbb{E}_{W^\rho} \sum_j |W_j^\rho x_i|^2 \\ &\leq C_1^2 + m_3 \beta_1^2. \end{aligned}$$

A.21.2 BOUNDING THE VALUE OF  $f'$ 

The following Lemma provides a reasonable bound on the value of the smoothed function.

**Lemma 34** *We have the following general bound on the values of the smoothed function: With high probability over the initialization, for  $\|W'\| \leq C_1$ ,  $\|V'\| \leq C_2$  and  $\forall i \in [n]$  (having small enough choices of  $\beta_1, \beta_2$  described in Appendix A.20):*

$$|f'_{W',V'}(x_i)| \lesssim (\kappa_2\sqrt{m_3} + \beta_2)\left(\sqrt{m_3}\kappa_1 + C_1 + \sqrt{m_3}\beta_1\right) + C_2(C_1 + \sqrt{m_3}\beta_1),$$

which is  $O(C_1C_2)$  for large enough overparameterization as described in Appendix A.3. Moreover, we have the following almost surely bound (with respect to the randomness of  $W^\rho$  and  $V^\rho$ ):

$$\begin{aligned} & |f_{W'+W^\rho, V'+V^\rho}(x_i)| \\ & \lesssim (\kappa_2\sqrt{m_3} + \|V^\rho\|_F)\left(\sqrt{m_3}\kappa_1 + C_1 + \sqrt{m_3}\left(\frac{1}{\sqrt{m_1}} \sum_j |W_j^\rho x_i|\right)\right) + C_2(C_1 + \sqrt{m_3}\left(\frac{1}{\sqrt{m_1}} \sum_j |W_j^\rho x_i|\right)). \end{aligned}$$

Notably, with slightly higher overparameterization, the high probability bound in (34) holds even if we take supremum over  $x$ .

**Proof of Lemma 34**

Using Lemmas 30 and 33 and using the fact that  $\phi^{(0)}(x_i)$  is orthogonal to the rows of  $V'$  (recall  $x'_i = \phi^{(0)}(x_i) + \phi^{(2)}(x_i)$ ):

$$\begin{aligned} & |f_{W'+W^\rho, V'+V^\rho}(x_i)| \\ & \leq \frac{1}{\sqrt{m_2}} a^T \sigma(V^{(0)} x'_i) + \frac{1}{\sqrt{m_2}} \sum_j |(V_j^\rho + V_j') x'_i| \\ & \leq (\kappa_2\sqrt{m_3}) \|x'_i\| + \frac{1}{\sqrt{m_2}} \sum_j |V_j^\rho x'_i| + C_2 \|\phi^{(2)}(x_i)\| \\ & \leq (\kappa_2\sqrt{m_3} + \|V^\rho\|_F) \|x'_i\| + C_2 \|\phi^{(2)}(x_i)\| \\ & \lesssim (\kappa_2\sqrt{m_3} + \|V^\rho\|_F)\left(\sqrt{m_3}\kappa_1 + C_1 + \sqrt{m_3}\left(\frac{1}{\sqrt{m_1}} \sum_j |W_j^\rho x_i|\right)\right) + C_2(C_1 + \sqrt{m_3}\left(\frac{1}{\sqrt{m_1}} \sum_j |W_j^\rho x_i|\right)). \end{aligned} \tag{284}$$

Note that above, if we apply the stronger worst-case norm bound of the first layer's output presented in Lemma 46, we would get  $\sup_{x, \|x\|=1} |f_{W'+W^\rho, V'+V^\rho}(x)|$  is bounded by the RHS, which in turn proves a stronger uniform bound on  $f'$ .

Similarly, this time by taking expectation with respect to  $W^\rho$  and  $V^\rho$ :

$$\begin{aligned} |f'_{W',V'}(x_i)| &= |\mathbb{E}_{W^\rho, V^\rho} f_{W',V'}(x_i)| \\ &\leq \mathbb{E}_{W^\rho, V^\rho} |f_{W',V'}(x_i)| \\ &= \mathbb{E}_{W^\rho} (\kappa_2\sqrt{m_3} + \beta_2) \|x'_i\| + C_2 \|x'_i - \phi^{(0)}(x_i)\| \\ &\lesssim (\kappa_2\sqrt{m_3} + \beta_2)\left(\sqrt{m_3}\kappa_1 + C_1 + \sqrt{m_3}\beta_1\right) + C_2(C_1 + \sqrt{m_3}\beta_1). \end{aligned}$$

**Corollary 8.1** *If we set  $C_1 = C_2 = 0$  above, we get*

$$|f'_{0,0}(x_i)| \leq (\kappa_2\sqrt{m_3} + \beta_2)\left(\sqrt{m_3}\kappa_1 + \sqrt{m_3}\beta_1\right),$$

the point being these terms go to zero by an order of  $O((\sqrt{m_2}\kappa_2)^{-\frac{2}{3}})$ . Therefore, taking  $(\sqrt{m_2}\kappa_2)^{-\frac{2}{3}} \ll B$ , we make sure that  $|f'_{0,0}| < B$ , so by the 1 smoothness of  $\ell$  and  $B$  boundedness of the labels we get  $\ell(f'_{0,0}(x_i), y_i) < 4B^2$ .

## A.21.3 BOUNDING THE DIFFERENCE BETWEEN ORIGINAL AND SMOOTHED FUNCTIONS

The following Lemma bounds the difference between the smoothed function and original function of the network.

**Lemma 35** *Bound on the smoothing change under the assumption  $m_2 \geq m_3 \log(m_2)$ : with high probability over the initialization, for any  $(W', V')$  with  $\|W'\| \leq C_1, \|V'\| \leq C_2$ :*

$$\begin{aligned} & |f_{W',V'}(x_i) - f'_{W',V'}(x_i)| \\ & \leq \beta_2(\kappa_1\sqrt{m_3} + C_1 + \sqrt{m_3}\beta_1) + (C_2 + \kappa_2\sqrt{m_2})\sqrt{m_3}\beta_1. \end{aligned}$$

**Proof of Lemma 35**

We write

$$\begin{aligned} & |f_{W',V'}(x_i) - f'_{W',V'}(x_i)| = |f_{W',V'}(x_i) - \mathbb{E}_{W^\rho, V^\rho} f_{W'+W^\rho, V'+V^\rho}(x_i)| \\ & = \left| \mathbb{E}_{W^\rho, V^\rho} \left( f_{W',V'}(x_i) - f_{W'+W^\rho, V'+V^\rho}(x_i) \right) \right| \\ & \leq \mathbb{E}_{W^\rho, V^\rho} \left| f_{W',V'}(x_i) - f_{W'+W^\rho, V'+V^\rho}(x_i) \right|. \end{aligned}$$

In the following,  $\sigma$  means we apply Relu activation to the vector in front of it (entrywise):

$$\begin{aligned} LHS & \leq \mathbb{E}_{W^\rho, V^\rho} \left| \frac{1}{\sqrt{m_2}} a^T \sigma(V^{(0)} + V' + V^\rho) \frac{1}{\sqrt{m_1}} W^s \sigma(W^{(0)} + W' + W^\rho) x_i \right. \\ & \quad \left. - \frac{1}{\sqrt{m_2}} a^T \sigma(V^{(0)} + V') \frac{1}{\sqrt{m_1}} W^s \sigma(W^{(0)} + W' + W^\rho) x_i \right| \\ & + \mathbb{E}_{W^\rho, V^\rho} \left| \frac{1}{\sqrt{m_2}} a^T \sigma(V^{(0)} + V') \frac{1}{\sqrt{m_1}} W^s \sigma(W^{(0)} + W' + W^\rho) x_i \right. \\ & \quad \left. - \frac{1}{\sqrt{m_2}} a^T \sigma(V^{(0)} + V') \frac{1}{\sqrt{m_1}} W^s \sigma(W^{(0)} + W') x_i \right|. \end{aligned}$$

Now for the first term above, using the previous notation of  $x'_i$  representing the output of the first layer and using Lemma 33:

$$\begin{aligned} & \mathbb{E}_{W^\rho, V^\rho} \left| \frac{1}{\sqrt{m_2}} a^T \sigma(V^{(0)} + V' + V^\rho) \frac{1}{\sqrt{m_1}} W^s \sigma(W^{(0)} + W' + W^\rho) x_i \right. \\ & \quad \left. - \frac{1}{\sqrt{m_2}} a^T \sigma(V^{(0)} + V') \frac{1}{\sqrt{m_1}} W^s \sigma(W^{(0)} + W' + W^\rho) x_i \right| \\ & \leq \mathbb{E}_{W^\rho, V^\rho} \frac{1}{\sqrt{m_2}} \sum_j |V_j^\rho x'_i| \\ & \lesssim \beta_2 \mathbb{E}_{W^\rho} \|x'_i\| \\ & \leq \beta_2 \mathbb{E}_{W^\rho, V^\rho} (\|\phi^{(0)}(x_i)\| + \|\phi^{(2)}(x_i)\|) \\ & \lesssim \beta_2(\kappa_1\sqrt{m_3} + C_1 + \sqrt{m_3}\beta_1). \end{aligned} \tag{285}$$

For the second term, by starting off with a simple triangle inequality:

$$\begin{aligned} & \mathbb{E}_{W^\rho, V^\rho} \left| \frac{1}{\sqrt{m_2}} a^T \sigma(V^{(0)} + V') x'_i - \frac{1}{\sqrt{m_2}} a^T \sigma(V^{(0)} + V') (\phi^{(0)}(x_i) + \phi^{(2)}(x_i)) \right| \\ & \leq \mathbb{E}_{W^\rho, V^\rho} \frac{1}{\sqrt{m_2}} \sum_{j=1}^{m_2} \left| (V_j^{(0)} + V_j')(x'_i - \phi^{(0)}(x_i) - \phi^{(2)}(x_i)) \right| \\ & \leq \mathbb{E}_{W^\rho, V^\rho} \frac{1}{\sqrt{m_2}} \sum_{j=1}^{m_2} \left| V_j^{(0)}(x'_i - \phi^{(0)}(x_i) - \phi^{(2)}(x_i)) \right| + \left| V_j'(x'_i - \phi^{(0)}(x_i) - \phi^{(2)}(x_i)) \right| \\ & \leq C_2 \mathbb{E}_{W^\rho} \|x'_i - \phi^{(0)}(x_i) - \phi^{(2)}(x_i)\| + \mathbb{E}_{W^\rho} \frac{1}{\sqrt{m_2}} \sum_{j=1}^{m_2} \left| V_j^{(0)}(x'_i - \phi^{(0)}(x_i) - \phi^{(2)}(x_i)) \right|. \end{aligned}$$

Now using Lemma 32:

$$\begin{aligned}
&\lesssim \left(C_2 + \kappa_2 \sqrt{m_2}\right) \mathbb{E}_{W^\rho} \|x'_i - \phi^{(0)}(x_i) - \phi^{(2)\prime}(x_i)\| \\
&\lesssim \left(C_2 + \kappa_2 \sqrt{m_2}\right) \sqrt{m_3} E_{W^\rho} \sum_j \frac{1}{\sqrt{m_1}} |W_j^\rho x_i| \\
&\lesssim \left(C_2 + \kappa_2 \sqrt{m_2}\right) \sqrt{m_3} \beta_1.
\end{aligned} \tag{286}$$

Combining Equations (285) and (286) we conclude the proof.

## B APPENDIX

## CONTENTS

<b>B.1 Smoothness coefficients</b> . . . . .	96
<b>B.1.1 Computing the Lipschitz Coefficients of <math>f'_{W',V'}</math></b> . . . . .	98
<b>B.2 Representation Lemmas</b> . . . . .	104
<b>B.2.1 Representation Toolbox</b> . . . . .	104
<b>B.2.2 Some Linear Algebra</b> . . . . .	105
<b>B.2.3 Bound on the norm of <math>\tilde{x}_i</math>'s</b> . . . . .	106
<b>B.3 Coupling for <math>\tilde{\nabla}_W, \tilde{\nabla}_V</math></b> . . . . .	108
<b>B.4 Handling the Injected Noise by PSGD</b> . . . . .	113
<b>B.4.1 Bounding the Norm of the first Layer's Output in the Worst Case</b> . . . . .	114

## B.1 SMOOTHNESS COEFFICIENTS

Recall that for a function  $f \in C^3$  on  $\mathbb{R}^d$ , we say it is  $\mu_1$  lipschitz,  $\mu_2$  gradient lipschitz, and  $\mu_3$  hessian lipschitz at point  $x$  if for every unit direction  $v$ ,  $|\frac{d}{d\lambda} f(x + \lambda v)| \leq \mu_1$ ,  $|\frac{d^2}{d\lambda^2} f(x + \lambda v)| \leq \mu_2$ , and  $|\frac{d^3}{d\lambda^3} f(x + \lambda v)| \leq \mu_3$ .

The aim of this section is to bound the lipschitz coefficients of the loss  $\ell(\cdot, y)$  and objective  $L(W', V')$  in a bounded domain  $\|W'\| \leq C_1, \|V'\| \leq C_2$ . The following is our main Theorem in this regard:

**Theorem 9** *For given values  $C_1, C_2 > 0$ , in the domain  $\|W'\| \leq C_1, \|V'\| \leq C_2$ , for any label  $|y| \leq B$ , the loss function  $\ell(\cdot, y)$  is  $O((C_1 C_2 + B^2))$ -lipschitz (having enough overparameterization) and 1 gradient-lipschitz  $x = f'_{W', V'}$ . Moreover, the loss function  $L(W', V')$  is  $(O(C_1 C_2) + B)\Psi_1 + 2(C_1 + C_2)$  lipschitz,  $\Psi_1^2 + (O(C_1 C_2) + B)\Psi_2 + 4$  gradient lipschitz, and  $3\Psi_2\Psi_1 + (O(C_1 C_2) + B)\Psi_3$  hessian lipschitz, where  $\Psi_1, \Psi_2, \Psi_3$  are defined in Lemma [34](#).*

**Proof of Lemma 9**

As in the proof of Lemma [36](#), let  $(\tilde{W}, \tilde{V})$  be a unit direction, i.e.  $\|\tilde{W}\|^2 + \|\tilde{V}\|^2 = 1$ . Then, using Lemma [34](#), we know that for every  $i \in [n]$ :  $|f'_{W', V'}(x_i)| = O(C_1 C_2)$ , so by 1-smoothness of the loss and  $B$ -boundedness of the labels, we get that  $\ell(\cdot, y)$  is  $(O(C_1 C_2) + B)$  lipschitz at point  $f'_{W', V'}$ . The gradient smoothness parameter of the square loss  $\ell$  is bounded by 1 and its third derivative is zero. Now using these coefficients, we can easily compute the coefficients for  $L$  as well by simple differentiation:

$$\begin{aligned} \left| \frac{d}{d\lambda} \ell(f'_{W'+\lambda\tilde{W}, V'+\lambda\tilde{V}}(x_i), y_i) \right| &= |\dot{\ell}(f', y_i) \frac{d}{d\lambda} f'| \leq (O(C_1 C_2) + B)\Psi_1. \\ \left| \frac{d}{d\lambda^2} \ell(f'_{W'+\lambda\tilde{W}, V'+\lambda\tilde{V}}(x_i), y_i) \right| &= |\ddot{\ell}(f', y_i) (\frac{d}{d\lambda} f')^2 + \dot{\ell}(f', y_i) \frac{d^2}{d\lambda^2} f'| \leq \Psi_1^2 + (O(C_1 C_2) + B)\Psi_2. \\ \left| \frac{d}{d\lambda^3} \ell(f'_{W'+\lambda\tilde{W}, V'+\lambda\tilde{V}}(x_i), y_i) \right| &= |\ddot{\ell}(f', y_i) (\frac{d}{d\lambda} f')^3 + 3\dot{\ell}(f', y_i) \frac{d^2}{d\lambda^2} f' \frac{d}{d\lambda} f' + \dot{\ell}(f', y_i) \frac{d^3}{d\lambda^3} f'| \\ &\leq \Psi_1^3 + 3\Psi_2\Psi_1 + (O(C_1 C_2) + B)\Psi_3. \end{aligned}$$

Moreover, note that

$$\begin{aligned} \frac{d}{d\lambda} \|W' + \lambda\tilde{W}\|^2 &= 2\langle W' + \lambda\tilde{W}, \tilde{W} \rangle \Big|_{\lambda=0} = 2\langle W', \tilde{W} \rangle \leq 2\|W'\| = 2C_1, \\ \frac{d^2}{d\lambda^2} \|W' + \lambda\tilde{W}\|^2 &= \langle \tilde{W}, \tilde{W} \rangle = 2, \\ \frac{d^3}{d\lambda^3} \|W' + \lambda\tilde{W}\|^2 &= 0, \end{aligned}$$

and similarly for  $\|V' + \lambda\tilde{V}\|^2$ . Combining these results finishes the proof.

Above, we used parameters  $\Psi_1, \Psi_2, \Psi_3$ , the lipschitz coefficients of  $f'$  in domain  $\|W'\| \leq C_1, \|V'\| \leq C_2$ , which we bound in Lemma [B6](#) below.

### B.1.1 COMPUTING THE LIPSCHITZ COEFFICIENTS OF $f'_{W',V'}$

In this section, we bound the lipschitz coefficients of  $f'_{W',V'}$  in the domain  $\|W'\| \leq C_1, \|V'\| \leq C_2$  by  $poly(m_1, m_2, m_3, \beta_1, \beta_2)$  functions.

**Lemma 36** *For every point  $(W', V')$  in the domain  $\|W'\| \leq C_1, \|V'\| \leq C_2$ , we have the following bounds on the lipschitz coefficients of  $f'_{W',V'}$  ( $(\tilde{W}, \tilde{V})$  is a unit direction with  $\|\tilde{W}\|^2 + \|\tilde{V}\|^2 = 1$ ):*

$$\begin{aligned}
& \left| \frac{d}{d\lambda} f'_{W'+\lambda\tilde{W}, V'+\lambda\tilde{V}}(x_i) \Big|_{\lambda=0} \right| \\
& \lesssim \frac{m_2}{\beta_2^2} \left( \frac{\beta_2}{\sqrt{m_2}} \sqrt{m_3} (\kappa_1 + C_1 + \beta_1) (\kappa_2 \sqrt{m_3} + C_2) + \frac{\beta_2^2}{\sqrt{m_2}} \sqrt{m_3} (\kappa_1 + C_1 + \beta_1) \right) \\
& + \frac{m_1}{\beta_1^2} \left( \sqrt{m_3} \left( \frac{\beta_1}{\sqrt{m_1}} (\kappa_1 + C_1) + \frac{\beta_1^2}{\sqrt{m_1}} \right) (\kappa_2 \sqrt{m_3} + C_2) + \beta_2 \sqrt{m_3} \left( \frac{\beta_1}{\sqrt{m_1}} (\kappa_1 + C_1) + \frac{\beta_1^2}{\sqrt{m_1}} \right) \right) := \Psi_1, \\
& \frac{d^2}{d\lambda^2} f'_{W'+\lambda\tilde{W}, V'+\lambda\tilde{V}}(x_i) \Big|_{\lambda=0} \\
& \lesssim \left( \frac{m_1}{\beta_1^2} \|\tilde{W}\|^2 + \frac{m_2}{\beta_2^2} \|\tilde{V}\|^2 \right) \sqrt{m_3 \left( (\kappa_2 \sqrt{m_3} + C_2)^2 + \beta_2^2 \right) \left[ (\kappa_1 + C_1)^2 + \beta_1^2 \right]} := \Psi_2 \\
& \left| \frac{d^3}{d\lambda^3} f'_{W'+\lambda\tilde{W}, V'+\lambda\tilde{V}}(x_i) \Big|_{\lambda=0} \right| \\
& \lesssim \left( \frac{m_1}{\beta_1^2} \|\tilde{W}\|^2 + \frac{m_2}{\beta_2^2} \|\tilde{V}\|^2 \right)^{3/2} \sqrt{m_3 \left( (\kappa_2 \sqrt{m_3} + C_2)^2 + \beta_2^2 \right) \left[ (\kappa_1 + C_1)^2 + \beta_1^2 \right]} := \Psi_3
\end{aligned} \tag{287}$$

### Proof of Lemma 36

Let

$$\rho(W^\rho, V^\rho) := \frac{1}{(\sqrt{2\pi})^{m_2 m_3 + m_1 d} (\beta_1 / \sqrt{m_1})^{m_1 d} (\beta_2 / \sqrt{m_2})^{m_2 m_3}} \exp\left\{-\frac{\|W^\rho\|^2}{2\beta_1^2/m_1} - \frac{\|V^\rho\|^2}{2\beta_2^2/m_2}\right\},$$

be the density function of the law of  $W^\rho$  and  $V^\rho$  which is a joint Gaussian. Then to compute the derivative and second derivative of the function in the unit direction  $(\tilde{W}, \tilde{V})$ , s.t.  $\|\tilde{W}\|_F^2 + \|\tilde{V}\|_F^2 = 1$ , we can write the value of the smoothed function as an integration with density  $\rho$ , change variable, and then take derivatives:

$$\begin{aligned}
& \frac{d}{d\lambda} f'_{W'+\lambda\tilde{W}, V'+\lambda\tilde{V}}(x_i) \Big|_{\lambda=0} \\
& = \frac{d}{d\lambda} \mathbb{E}_{W^\rho, V^\rho} f_{W'+\lambda\tilde{W}+W^\rho, V'+\lambda\tilde{V}+V^\rho}(x_i) \\
& = \frac{d}{d\lambda} \int f_{W'+\lambda\tilde{W}+W^\rho, V'+\lambda\tilde{V}+V^\rho}(x_i) \rho(W^\rho, V^\rho) d(W^\rho, V^\rho) \\
& = \frac{d}{d\lambda} \int f_{W^+, V^+}(x_i) \rho(W^+ - (W' + \lambda\tilde{W}), V^+ - (V' + \lambda\tilde{V})) d(W^+, V^+).
\end{aligned}$$

But one can easily see that for fixed  $V'$  and  $\tilde{V}$ , the set of functions  $f_{W^+, V^+}(x_i) \rho(W^+ - (W' + \lambda\tilde{W}), V^+ - (V' + \lambda\tilde{V}))$  for a small neighborhood of  $\lambda$  can simultaneously be upper bounded by an integrable function. Hence, the Leibnitz rule holds here because of dominated convergence theorem,

and we can change the order of integration and derivation:

$$\begin{aligned}
&= \int f_{W^+, V^+}(x_i) \frac{d}{d\lambda} \rho(W^+ - (W' + \lambda\tilde{W}), V^+ - (V' + \lambda\tilde{V})) d(W^+, V^+) \\
&= - \int f_{W^+, V^+}(x_i) \left\langle (\tilde{W}, \tilde{V}), \left( \left( \frac{\beta_1^2}{m_1} \right)^{-1} (W^+ - (W' + \lambda\tilde{W})), \left( \frac{\beta_2^2}{m_2} \right)^{-1} (V^+ - (V' + \lambda\tilde{V})) \right) \right\rangle \\
&\quad \rho(W^+ - (W' + \lambda\tilde{W}), V^+ - (V' + \lambda\tilde{V})) d(W^+, V^+) \\
&= \int f_{W^+, V^+}(x_i) \left\langle (\tilde{W}, \tilde{V}), \left( \frac{m_1}{\beta_1^2} W^\rho, \frac{m_2}{\beta_2^2} V^\rho \right) \right\rangle \rho(W^\rho, V^\rho) d(W^\rho, V^\rho) \\
&= \mathbb{E}_{W^\rho, V^\rho} \left( \frac{m_1}{\beta_1^2} \langle \tilde{W}, W^\rho \rangle + \frac{m_2}{\beta_2^2} \langle \tilde{V}, V^\rho \rangle \right) f_{W'+\lambda\tilde{W}+W^\rho, V'+\lambda\tilde{V}+V^\rho}(x_i) \Big|_{\lambda=0} \\
&= \mathbb{E}_{W^\rho, V^\rho} \left( \frac{m_1}{\beta_1^2} \langle \tilde{W}, W^\rho \rangle + \frac{m_2}{\beta_2^2} \langle \tilde{V}, V^\rho \rangle \right) f_{W'+W^\rho, V'+V^\rho}(x_i). \tag{288}
\end{aligned}$$

Similarly we can compute the second derivative:

$$\begin{aligned}
&\frac{d^2}{d\lambda^2} f'_{W'+\lambda\tilde{W}, V'+\lambda\tilde{V}}(x_i) \Big|_{\lambda=0} \\
&= \frac{d}{d\lambda} \int f_{W^+, V^+}(x_i) \left\langle (\tilde{W}, \tilde{V}), \left( \left( \frac{\beta_1^2}{m_1} \right)^{-1} (W^+ - (W' + \lambda\tilde{W})), \left( \frac{\beta_2^2}{m_2} \right)^{-1} (V^+ - (V' + \lambda\tilde{V})) \right) \right\rangle \\
&\quad \rho(W^+ - (W' + \lambda\tilde{W}), V^+ - (V' + \lambda\tilde{V})) d(W^+, V^+) \\
&= \int f_{W^+, V^+}(x_i) \left[ \left\langle (\tilde{W}, \tilde{V}), \left( \left( \frac{\beta_1^2}{m_1} \right)^{-1} (W^+ - (W' + \lambda\tilde{W})), \left( \frac{\beta_2^2}{m_2} \right)^{-1} (V^+ - (V' + \lambda\tilde{V})) \right) \right\rangle^2 - \right. \\
&\quad \left. \left\langle (\tilde{W}, \tilde{V}), \left( \left( \frac{\beta_1^2}{m_1} \right)^{-1} \tilde{W}, \left( \frac{\beta_2^2}{m_2} \right)^{-1} \tilde{V} \right) \right\rangle \right] \rho(W^+ - (W' + \lambda\tilde{W}), V^+ - (V' + \lambda\tilde{V})) d(W^+, V^+) \\
&= \int f_{W^+, V^+}(x_i) \left[ \left\langle (\tilde{W}, \tilde{V}), \left( \frac{m_1}{\beta_1^2} W^\rho, \frac{m_2}{\beta_2^2} V^\rho \right) \right\rangle^2 + \left\langle (\tilde{W}, \tilde{V}), \left( \frac{m_1}{\beta_1^2} \tilde{W}, \frac{m_2}{\beta_2^2} \tilde{V} \right) \right\rangle \right] \rho(W^\rho, V^\rho) d(W^\rho, V^\rho) \\
&= \mathbb{E}_{W^\rho, V^\rho} \left[ \left( \frac{m_1}{\beta_1^2} \langle \tilde{W}, W^\rho \rangle + \frac{m_2}{\beta_2^2} \langle \tilde{V}, V^\rho \rangle \right)^2 - \left( \frac{m_1}{\beta_1^2} \|\tilde{W}\|^2 + \frac{m_2}{\beta_2^2} \|\tilde{V}\|^2 \right) \right] f_{W'+\lambda\tilde{W}+W^\rho, V'+\lambda\tilde{V}+V^\rho}(x_i) \Big|_{\lambda=0} \\
&= \mathbb{E}_{W^\rho, V^\rho} \left[ \left( \frac{m_1}{\beta_1^2} \langle \tilde{W}, W^\rho \rangle + \frac{m_2}{\beta_2^2} \langle \tilde{V}, V^\rho \rangle \right)^2 - \left( \frac{m_1}{\beta_1^2} \|\tilde{W}\|^2 + \frac{m_2}{\beta_2^2} \|\tilde{V}\|^2 \right) \right] f_{W'+W^\rho, V'+V^\rho}(x_i). \tag{289}
\end{aligned}$$

Similarly for the third derivative:

$$\begin{aligned}
&\frac{d^3}{d\lambda^3} f'_{W'+\lambda\tilde{W}, V'+\lambda\tilde{V}}(x_i) \Big|_{\lambda=0} \\
&= \frac{d}{d\lambda} \int f_{W^+, V^+}(x_i) \left[ \left\langle (\tilde{W}, \tilde{V}), \left( \left( \frac{\beta_1^2}{m_1} \right)^{-1} (W^+ - (W' + \lambda\tilde{W})), \left( \frac{\beta_2^2}{m_2} \right)^{-1} (V^+ - (V' + \lambda\tilde{V})) \right) \right\rangle^2 - \right. \\
&\quad \left. \left\langle (\tilde{W}, \tilde{V}), \left( \left( \frac{\beta_1^2}{m_1} \right)^{-1} \tilde{W}, \left( \frac{\beta_2^2}{m_2} \right)^{-1} \tilde{V} \right) \right\rangle \right] \rho(W^+ - (W' + \lambda\tilde{W}), V^+ - (V' + \lambda\tilde{V})) d(W^+, V^+) \\
&= \int f_{W^+, V^+}(x_i) \left[ \left\langle (\tilde{W}, \tilde{V}), \left( \left( \frac{\beta_1^2}{m_1} \right)^{-1} (W^+ - (W' + \lambda\tilde{W})), \left( \frac{\beta_2^2}{m_2} \right)^{-1} (V^+ - (V' + \lambda\tilde{V})) \right) \right\rangle^3 \right. \\
&\quad \left. - 3 \left\langle (\tilde{W}, \tilde{V}), \left( \left( \frac{\beta_1^2}{m_1} \right)^{-1} (W^+ - (W' + \lambda\tilde{W})), \left( \frac{\beta_2^2}{m_2} \right)^{-1} (V^+ - (V' + \lambda\tilde{V})) \right) \right\rangle \left\langle (\tilde{W}, \tilde{V}), \left( \left( \frac{\beta_1^2}{m_1} \right)^{-1} \tilde{W}, \left( \frac{\beta_2^2}{m_2} \right)^{-1} \tilde{V} \right) \right\rangle \right] \\
&\quad \rho(W^+ - (W' + \lambda\tilde{W}), V^+ - (V' + \lambda\tilde{V})) d(W^+, V^+) \\
&= \mathbb{E}_{W^\rho, V^\rho} \left[ \left( \frac{m_1}{\beta_1^2} \langle \tilde{W}, W^\rho \rangle + \frac{m_2}{\beta_2^2} \langle \tilde{V}, V^\rho \rangle \right)^3 - 3 \left( \frac{m_1}{\beta_1^2} \langle \tilde{W}, W^\rho \rangle + \frac{m_2}{\beta_2^2} \langle \tilde{V}, V^\rho \rangle \right) \left( \frac{m_1}{\beta_1^2} \|\tilde{W}\|^2 + \frac{m_2}{\beta_2^2} \|\tilde{V}\|^2 \right) \right] \\
&\quad f_{W'+W^\rho, V'+V^\rho}(x_i).
\end{aligned}$$

Now for first derivative, exactly similar to the derivation in (284), we can write

$$\begin{aligned}
& \left| \mathbb{E}_{W^\rho, V^\rho} \left( \frac{m_1}{\beta_1^2} \langle \tilde{W}, W^\rho \rangle + \frac{m_2}{\beta_2^2} \langle \tilde{V}, V^\rho \rangle \right) f_{W'+W^\rho, V'+V^\rho}(x_i) \right| \\
& \leq \mathbb{E}_{W^\rho, V^\rho} \left| \frac{m_1}{\beta_1^2} \langle \tilde{W}, W^\rho \rangle + \frac{m_2}{\beta_2^2} \langle \tilde{V}, V^\rho \rangle \right| \left| f_{W'+W^\rho, V'+V^\rho}(x_i) \right| \\
& \leq \mathbb{E}_{W^\rho, V^\rho} \left| \frac{m_1}{\beta_1^2} \langle \tilde{W}, W^\rho \rangle + \frac{m_2}{\beta_2^2} \langle \tilde{V}, V^\rho \rangle \right| \left( \kappa_2 \sqrt{m_3} \|x_i^{(2)}\| + \|V'\|_F \|x_i^{(2)}\| + \mathbb{E}_{V^\rho} \frac{1}{\sqrt{m_2}} \sum_j |V_j^\rho x_i^{(2)}| \right) \\
& \leq \mathbb{E}_{W^\rho, V^\rho} \left( \left| \frac{m_1}{\beta_1^2} \langle \tilde{W}, W^\rho \rangle \right| + \left| \frac{m_2}{\beta_2^2} \langle \tilde{V}, V^\rho \rangle \right| \right) \left( \kappa_2 \sqrt{m_3} \|x_i^{(2)}\| + \|V'\|_F \|x_i^{(2)}\| + \frac{1}{\sqrt{m_2}} \sum_j |V_j^\rho x_i^{(2)}| \right) \\
& = \frac{m_2}{\beta_2^2} \left( \mathbb{E}_{V^\rho} \left| \langle \tilde{V}, V^\rho \rangle \right| \mathbb{E}_{W^\rho} \|x_i^{(2)}\| \left( \kappa_2 \sqrt{m_3} + \|V'\|_F \right) + \mathbb{E}_{W^\rho} \mathbb{E}_{V^\rho} \frac{1}{\sqrt{m_2}} \left| \langle \tilde{V}, V^\rho \rangle \right| \sum_j |V_j^\rho x_i^{(2)}| \right) \\
& + \frac{m_1}{\beta_1^2} \left( \left( \mathbb{E}_{W^\rho} \|x_i^{(2)}\| \left| \langle \tilde{W}, W^\rho \rangle \right| \right) \left( \kappa_2 \sqrt{m_3} + \|V'\|_F \right) + \mathbb{E}_{W^\rho} \left| \langle \tilde{W}, W^\rho \rangle \right| \mathbb{E}_{V^\rho} \frac{1}{\sqrt{m_2}} \sum_j |V_j^\rho x_i^{(2)}| \right) \\
& \lesssim \frac{m_2}{\beta_2^2} \left( \mathbb{E}_{V^\rho} \left| \langle \tilde{V}, V^\rho \rangle \right| \mathbb{E}_{W^\rho} \|x_i^{(2)}\| \left( \kappa_2 \sqrt{m_3} + C_2 \right) + \mathbb{E}_{W^\rho} \mathbb{E}_{V^\rho} \frac{1}{\sqrt{m_2}} \left| \langle \tilde{V}, V^\rho \rangle \right| \sum_j |V_j^\rho x_i^{(2)}| \right) \\
& + \frac{m_1}{\beta_1^2} \left( \left( \mathbb{E}_{W^\rho} \|x_i^{(2)}\| \left| \langle \tilde{W}, W^\rho \rangle \right| \right) \left( \kappa_2 \sqrt{m_3} + C_2 \right) + \mathbb{E}_{W^\rho} \left| \langle \tilde{W}, W^\rho \rangle \right| \mathbb{E}_{V^\rho} \frac{1}{\sqrt{m_2}} \sum_j |V_j^\rho x_i^{(2)}| \right),
\end{aligned}$$

where the last line follows because  $\|V'\| \lesssim C_2$ . But notice that because  $\|\tilde{V}\|_F \leq 1$ ,  $\|\tilde{W}\|_F \leq 1$ , then  $\langle \tilde{V}, V^\rho \rangle$  and  $\langle \tilde{W}, W^\rho \rangle$  are Gaussian variables with variances at most  $\beta_2^2/m_2$  and  $\beta_1^2/m_1$ . Hence

$$\mathbb{E}_{V^\rho} \left| \langle \tilde{V}, V^\rho \rangle \right| \lesssim \beta_2 / \sqrt{m_2}, \quad (290)$$

$$\mathbb{E}_{W^\rho} \left| \langle \tilde{W}, W^\rho \rangle \right| \lesssim \beta_1 / \sqrt{m_1}. \quad (291)$$

Similarly, using the same derivation as in (B3), one can also get the following a.s. bound (over the randomness of  $W^\rho$ ):

$$\|x_i^{(2)}\| \lesssim \sqrt{m_3} \left( \kappa_1 + C_1 + \frac{1}{\sqrt{m_1}} \sum_j |W_j^\rho x_i| \right), \quad (292)$$

therefore

$$\mathbb{E}_{W^\rho} \|x_i^{(2)}\| \lesssim \sqrt{m_3} \left( \kappa_1 + C_1 + \mathbb{E}_{W^\rho} \frac{1}{\sqrt{m_1}} \sum_j |W_j^\rho x_i| \right) \quad (293)$$

$$\lesssim \sqrt{m_3} \left( \kappa_1 + C_1 + \beta_1 \right). \quad (294)$$

Moreover, for every  $j \in [m_2]$ :

$$\begin{aligned}
\mathbb{E}_{V^\rho} \left| \langle \tilde{V}, V^\rho \rangle \right| |V_j^\rho x_i^{(2)}| & \leq \sqrt{\mathbb{E}_{V^\rho} \left| \langle \tilde{V}, V^\rho \rangle \right|^2} \sqrt{\mathbb{E}_{V^\rho} |V_j^\rho x_i^{(2)}|^2} \\
& = \frac{\beta_2}{\sqrt{m_2}} \frac{\beta_2}{\sqrt{m_2}} \|x_i^{(2)}\| = \frac{\beta_2^2}{m_2} \|x_i^{(2)}\|, \\
\mathbb{E}_{W^\rho} \left| \langle \tilde{W}, W^\rho \rangle \right| |W_j^\rho x_i| & \leq \frac{\beta_1^2}{m_1}.
\end{aligned} \quad (295)$$

Similarly, using Equation (292) we bound

$$\begin{aligned}
\mathbb{E}_{W^\rho} \|x_i^{(2)}\| \left| \langle \tilde{W}, W^\rho \rangle \right| & \lesssim \sqrt{m_3} \left( \mathbb{E}_{W^\rho} \left| \langle \tilde{W}, W^\rho \rangle \right| \left( \kappa_1 + C_1 \right) + \mathbb{E}_{W^\rho} \frac{1}{\sqrt{m_1}} \sum_j \left| \langle \tilde{W}, W^\rho \rangle \right| |W_j^\rho x_i| \right) \\
& \leq \sqrt{m_3} \left( \frac{\beta_1}{\sqrt{m_1}} \left( \kappa_1 + C_1 \right) + \frac{\beta_1^2}{\sqrt{m_1}} \right).
\end{aligned} \quad (296)$$

Now applying these bounds (291), (294), (295), and (296) to (290) and using the fact that

$$\mathbb{E}_{V^\rho} \frac{1}{\sqrt{m_2}} \sum_j |V_j^\rho x_i'^{(2)}| \leq \beta_2 \|x_i'^{(2)}\|,$$

we get

$$\begin{aligned} LHS &\lesssim \frac{m_2}{\beta_2^2} \left( \frac{\beta_2}{\sqrt{m_2}} \sqrt{m_3} (\kappa_1 + C_1 + \beta_1) (\kappa_2 \sqrt{m_3} + C_2) + \mathbb{E}_{W^\rho} \frac{\beta_2^2}{\sqrt{m_2}} \|x_i'^{(2)}\| \right) \\ &\quad + \frac{m_1}{\beta_1^2} \left( \sqrt{m_3} \left( \frac{\beta_1}{\sqrt{m_1}} (\kappa_1 + C_1) + \frac{\beta_1^2}{\sqrt{m_1}} \right) (\kappa_2 \sqrt{m_3} + C_2) + \mathbb{E}_{W^\rho} \left| \langle \tilde{W}, W^\rho \rangle \right| \beta_2 \|x_i'^{(2)}\| \right) \\ &\lesssim \frac{m_2}{\beta_2^2} \left( \frac{\beta_2}{\sqrt{m_2}} \sqrt{m_3} (\kappa_1 + C_1 + \beta_1) (\kappa_2 \sqrt{m_3} + C_2) + \frac{\beta_2^2}{\sqrt{m_2}} \sqrt{m_3} (\kappa_1 + C_1 + \beta_1) \right) \\ &\quad + \frac{m_1}{\beta_1^2} \left( \sqrt{m_3} \left( \frac{\beta_1}{\sqrt{m_1}} (\kappa_1 + C_1) + \frac{\beta_1^2}{\sqrt{m_1}} \right) (\kappa_2 \sqrt{m_3} + C_2) + \beta_2 \sqrt{m_3} \left( \frac{\beta_1}{\sqrt{m_1}} (\kappa_1 + C_1) + \frac{\beta_1^2}{\sqrt{m_1}} \right) \right). \end{aligned}$$

To make it easier for handling the second and third derivatives, we first bound the expectations of  $f_{\tilde{W}'+W^\rho, V'+V^\rho(x_i)}^2(x_i)$  which enables us to use Cauchy-schwartz. Again using similar derivation as in (284) and Equation (292):

$$\begin{aligned} &\mathbb{E}_{W^\rho, V^\rho} f_{\tilde{W}'+W^\rho, V'+V^\rho(x_i)}^2(x_i) \\ &\leq \mathbb{E}_{W^\rho, V^\rho} \left( \kappa_2 \sqrt{m_3} \|x_i'^{(2)}\| + C_2 \|x_i'^{(2)}\| + \frac{1}{\sqrt{m_2}} \sum_j |V_j^\rho x_i'^{(2)}| \right)^2 \\ &\lesssim (\kappa_2 \sqrt{m_3} + C_2)^2 \mathbb{E}_{W^\rho} \|x_i'^{(2)}\|^2 + E_{W^\rho} E_{V^\rho} \frac{1}{m_2} \left( \sum_j |V_j^\rho x_i'^{(2)}| \right)^2 \\ &\lesssim (\kappa_2 \sqrt{m_3} + C_2)^2 \mathbb{E}_{W^\rho} \|x_i'^{(2)}\|^2 + E_{W^\rho} \frac{1}{m_2} \sum_j E_{V_j^\rho} |V_j^\rho x_i'^{(2)}|^2 + E_{W^\rho} \frac{1}{m_2} \sum_{j_1 \neq j_2} E_{V_{j_1}^\rho} |V_{j_1}^\rho x_i'^{(2)}| E_{V_{j_2}^\rho} |V_{j_2}^\rho x_i'^{(2)}| \\ &\lesssim (\kappa_2 \sqrt{m_3} + C_2)^2 \mathbb{E}_{W^\rho} \|x_i'^{(2)}\|^2 + E_{W^\rho} \frac{\beta_2^2}{m_2} \|x_i'^{(2)}\|^2 + E_{W^\rho} \beta_2^2 \|x_i'^{(2)}\|^2 \frac{m(m-1)}{m^2} \\ &\lesssim \left( (\kappa_2 \sqrt{m_3} + C_2)^2 + \beta_2^2 \right) \mathbb{E}_{W^\rho} \|x_i'^{(2)}\|^2 \\ &\lesssim \left( (\kappa_2 \sqrt{m_3} + C_2)^2 + \beta_2^2 \right) E_{W^\rho} m_3 \left( \kappa_1 + C_1 + \frac{1}{\sqrt{m_1}} \sum_j |W_j^\rho x_i| \right)^2 \\ &\lesssim \left( (\kappa_2 \sqrt{m_3} + C_2)^2 + \beta_2^2 \right) m_3 \left[ (\kappa_1 + C_1)^2 + \frac{1}{m_1} \sum_j \mathbb{E}_{W_j^\rho} |W_j^\rho x_i|^2 + \frac{1}{m_1} \sum_{j_1 \neq j_2} \mathbb{E}_{W_{j_1}^\rho} |W_{j_1}^\rho x_i| \mathbb{E}_{W_{j_2}^\rho} |W_{j_2}^\rho x_i| \right] \\ &\lesssim \left( (\kappa_2 \sqrt{m_3} + C_2)^2 + \beta_2^2 \right) m_3 \left[ (\kappa_1 + C_1)^2 + \frac{\beta_1^2}{m_1} + \beta_1^2 \frac{m(m-1)}{m^2} \right] \\ &= m_3 \left( (\kappa_2 \sqrt{m_3} + C_2)^2 + \beta_2^2 \right) \left[ (\kappa_1 + C_1)^2 + \beta_1^2 \right]. \tag{297} \end{aligned}$$

Now for the second derivative, we can proceed by applying Cauchy-Swartz:

$$\begin{aligned}
& \left. \frac{d^2}{d\lambda^2} f'_{W'+\lambda\tilde{W}, V'+\lambda\tilde{V}}(x_i) \right|_{\lambda=0} \\
& \leq \mathbb{E}_{W^\rho, V^\rho} \left| \left( \frac{m_1}{\beta_1^2} \langle \tilde{W}, W^\rho \rangle + \frac{m_2}{\beta_2^2} \langle \tilde{V}, V^\rho \rangle \right)^2 - \left( \frac{m_1}{\beta_1^2} \|\tilde{W}\|^2 + \frac{m_2}{\beta_2^2} \|\tilde{V}\|^2 \right) \right| \left| f_{W'+W^\rho, V'+V^\rho}(x_i) \right| \\
& \leq \mathbb{E}_{W^\rho, V^\rho} \left| \frac{m_1^2}{\beta_1^4} \langle \tilde{W}, W^\rho \rangle^2 - \frac{m_1}{\beta_1^2} \|\tilde{W}\|^2 + \frac{m_2^2}{\beta_2^4} \langle \tilde{V}, V^\rho \rangle^2 - \frac{m_2}{\beta_2^2} \|\tilde{V}\|^2 + \frac{2m_1m_2}{\beta_1^2\beta_2^2} \langle \tilde{W}, W^\rho \rangle \langle \tilde{V}, V^\rho \rangle \right| \\
& \left| f_{W'+W^\rho, V'+V^\rho}(x_i) \right| \\
& \leq \sqrt{\mathbb{E}_{W^\rho, V^\rho} \left| \frac{m_1^2}{\beta_1^4} \langle \tilde{W}, W^\rho \rangle^2 - \frac{m_1}{\beta_1^2} \|\tilde{W}\|^2 + \frac{m_2^2}{\beta_2^4} \langle \tilde{V}, V^\rho \rangle^2 - \frac{m_2}{\beta_2^2} \|\tilde{V}\|^2 + \frac{2m_1m_2}{\beta_1^2\beta_2^2} \langle \tilde{W}, W^\rho \rangle \langle \tilde{V}, V^\rho \rangle \right|^2} \\
& \sqrt{\mathbb{E}_{W^\rho, V^\rho} \left| f_{W'+W^\rho, V'+V^\rho}(x_i) \right|^2}.
\end{aligned}$$

Now note that the cross terms have expectation zero, so we get

$$\begin{aligned}
& = \sqrt{\mathbb{E}_{W^\rho, V^\rho} \left( \frac{m_1^2}{\beta_1^4} \langle \tilde{W}, W^\rho \rangle^2 - \frac{m_1}{\beta_1^2} \|\tilde{W}\|^2 \right)^2 + \left( \frac{m_2^2}{\beta_2^4} \langle \tilde{V}, V^\rho \rangle^2 - \frac{m_2}{\beta_2^2} \|\tilde{V}\|^2 \right)^2 + \frac{4m_1^2m_2^2}{\beta_1^4\beta_2^4} \mathbb{E}_{W^\rho, V^\rho} \langle \tilde{W}, W^\rho \rangle^2 \langle \tilde{V}, V^\rho \rangle^2} \\
& \sqrt{\mathbb{E}_{W^\rho, V^\rho} \left| f_{W'+W^\rho, V'+V^\rho}(x_i) \right|^2} \\
& \lesssim \sqrt{\frac{m_1^2}{\beta_1^4} \|\tilde{W}\|^4 + \frac{m_2^2}{\beta_2^4} \|\tilde{V}\|^4 + \frac{4m_1m_2}{\beta_1^2\beta_2^2} \|\tilde{W}\|^2 \|\tilde{V}\|^2} \sqrt{\mathbb{E}_{W^\rho, V^\rho} \left| f_{W'+W^\rho, V'+V^\rho}(x_i) \right|^2} \\
& \lesssim \left( \frac{m_1}{\beta_1^2} \|\tilde{W}\|^2 + \frac{m_2}{\beta_2^2} \|\tilde{V}\|^2 \right) \sqrt{\mathbb{E}_{W^\rho, V^\rho} \left| f_{W'+W^\rho, V'+V^\rho}(x_i) \right|^2}.
\end{aligned}$$

Now applying Cauchy-shwartz and Equation (297) to above and combining it with Equations (289):

$$\begin{aligned}
& \left. \frac{d^2}{d\lambda^2} f'_{W'+\lambda\tilde{W}, V'+\lambda\tilde{V}}(x_i) \right|_{\lambda=0} \\
& \lesssim \left( \frac{m_1}{\beta_1^2} \|\tilde{W}\|^2 + \frac{m_2}{\beta_2^2} \|\tilde{V}\|^2 \right) \sqrt{m_3 \left( (\kappa_2 \sqrt{m_3} + C_2)^2 + \beta_2^2 \right) \left[ (\kappa_1 + C_1)^2 + \beta_1^2 \right]}.
\end{aligned}$$

Similarly for the third derivative:

$$\begin{aligned}
& \left| \frac{d^3}{d\lambda^3} f'_{W'+\lambda\tilde{W}, V'+\lambda\tilde{V}}(x_i) \right|_{\lambda=0} \\
& = \mathbb{E}_{W^\rho, V^\rho} \left| \left( \frac{m_1}{\beta_1^2} \langle \tilde{W}, W^\rho \rangle + \frac{m_2}{\beta_2^2} \langle \tilde{V}, V^\rho \rangle \right)^3 - 3 \left( \frac{m_1}{\beta_1^2} \langle \tilde{W}, W^\rho \rangle + \frac{m_2}{\beta_2^2} \langle \tilde{V}, V^\rho \rangle \right) \left( \frac{m_1}{\beta_1^2} \|\tilde{W}\|^2 + \frac{m_2}{\beta_2^2} \|\tilde{V}\|^2 \right) \right| \\
& \left| f_{W'+W^\rho, V'+V^\rho}(x_i) \right| \\
& \leq \sqrt{\mathbb{E}_{W^\rho, V^\rho} \left[ \left( \frac{m_1}{\beta_1^2} \langle \tilde{W}, W^\rho \rangle + \frac{m_2}{\beta_2^2} \langle \tilde{V}, V^\rho \rangle \right)^3 - 3 \left( \frac{m_1}{\beta_1^2} \langle \tilde{W}, W^\rho \rangle + \frac{m_2}{\beta_2^2} \langle \tilde{V}, V^\rho \rangle \right) \left( \frac{m_1}{\beta_1^2} \|\tilde{W}\|^2 + \frac{m_2}{\beta_2^2} \|\tilde{V}\|^2 \right) \right]^2} \\
& \sqrt{\mathbb{E}_{W^\rho, V^\rho} \left| f_{W'+W^\rho, V'+V^\rho}(x_i) \right|^2}. \tag{298}
\end{aligned}$$

But note that

$$\begin{aligned}
& \mathbb{E}_{W^\rho, V^\rho} \left[ \left( \frac{m_1}{\beta_1^2} \langle \tilde{W}, W^\rho \rangle + \frac{m_2}{\beta_2^2} \langle \tilde{V}, V^\rho \rangle \right)^3 - 3 \left( \frac{m_1}{\beta_1^2} \langle \tilde{W}, W^\rho \rangle + \frac{m_2}{\beta_2^2} \langle \tilde{V}, V^\rho \rangle \right) \left( \frac{m_1}{\beta_1^2} \|\tilde{W}\|^2 + \frac{m_2}{\beta_2^2} \|\tilde{V}\|^2 \right) \right]^2 \\
& \leq 2 \mathbb{E}_{W^\rho, V^\rho} \left( \frac{m_1}{\beta_1^2} \langle \tilde{W}, W^\rho \rangle + \frac{m_2}{\beta_2^2} \langle \tilde{V}, V^\rho \rangle \right)^6 \\
& + 18 \left( \frac{m_1}{\beta_1^2} \|\tilde{W}\|^2 + \frac{m_2}{\beta_2^2} \|\tilde{V}\|^2 \right)^2 \mathbb{E}_{W^\rho, V^\rho} \left( \frac{m_1}{\beta_1^2} \langle \tilde{W}, W^\rho \rangle + \frac{m_2}{\beta_2^2} \langle \tilde{V}, V^\rho \rangle \right)^2
\end{aligned}$$

Now note that  $\frac{m_1}{\beta_1^2} \langle \tilde{W}, W^\rho \rangle + \frac{m_2}{\beta_2^2} \langle \tilde{V}, V^\rho \rangle$  is a normal variable with variance  $\frac{m_1}{\beta_1^2} \|\tilde{W}\|^2 + \frac{m_2}{\beta_2^2} \|\tilde{V}\|^2$ . Therefore, by the bound on the moments of normal random variables:

$$LHS \lesssim \left( \frac{m_1}{\beta_1^2} \|\tilde{W}\|^2 + \frac{m_2}{\beta_2^2} \|\tilde{V}\|^2 \right)^3. \quad (299)$$

Plugging this into Equation (298) and also using Equation (297):

$$\begin{aligned} & \left| \frac{d^3}{d\lambda^3} f'_{W'+\lambda\tilde{W}, V'+\lambda\tilde{V}}(x_i) \Big|_{\lambda=0} \right| \\ & \lesssim \left( \frac{m_1}{\beta_1^2} \|\tilde{W}\|^2 + \frac{m_2}{\beta_2^2} \|\tilde{V}\|^2 \right)^{3/2} \sqrt{m_3 \left( (\kappa_2 \sqrt{m_3} + C_2)^2 + \beta_2^2 \right) \left[ (\kappa_1 + C_1)^2 + \beta_1^2 \right]}. \end{aligned}$$

## B.2 REPRESENTATION LEMMAS

In this section, we prove lemmas mostly related to the representation power of the network, which we mainly use in Appendix [A.12](#).

### B.2.1 REPRESENTATION TOOLBOX

**Lemma 37** Recall the definitions of  $W_k^+$ , and  $\bar{x}_{i,k}$  from Equations [\(49\)](#), and [\(45\)](#). For all  $k, i \in [n]$ :

$$|\bar{x}_{i,k} - \text{trace}(W_k^+ Z_k^i)| \leq \sqrt{n/(m_1 \lambda_0)} \|\mathcal{V}_k\|_{H^\infty}.$$

**Proof of Lemma 37**

$$\text{trace}(W_k^+ Z_k^i) = \text{trace}(Z_k^i \sum_{j=1}^n \mathcal{V}_{k,j} Z_k^j) = \sum_{i'=1}^n \mathcal{V}_{k,i'} \langle Z_k^i, Z_k^{i'} \rangle. \quad (300)$$

But note

$$\begin{aligned} \langle Z_k^i, Z_k^{i'} \rangle &= 1/m_1 \sum_{j=1}^{m_1} W_{k,j}'^2 \mathbb{1}\{W_j^{(0)T} x_i\} \mathbb{1}\{W_j^{(0)T} x_{i'}\} \langle x_{i'}, x_i \rangle \\ &= \frac{\langle x_{i'}, x_i \rangle}{m_1} \sum_{j=1}^{m_1} \mathbb{1}\{W_j^{(0)T} x_i\} \mathbb{1}\{W_j^{(0)T} x_{i'}\}. \end{aligned}$$

Now note that  $\langle x_{i'}, x_i \rangle \leq 1$ , and  $\mathbb{1}\{W_j^{(0)T} x_i\} \mathbb{1}\{W_j^{(0)T} x_{i'}\}$  is a Bernoulli with

$$\mathbb{E} \mathbb{1}\{W_j^{(0)T} x_i\} \mathbb{1}\{W_j^{(0)T} x_{i'}\} = 1/4 + \arcsin(\langle x, y \rangle)/2\pi.$$

Therefore, by Hoeffding inequality we get

$$|\langle Z_k^i, Z_k^{i'} \rangle - H_{i,i'}^\infty| = O(1/\sqrt{m_1}).$$

Hence, because obviously  $\|H^\infty\|_2 \leq 1$ , we get

$$\text{trace}(W_k^+ Z_k^i) = \sum_{i'=1}^n \mathcal{V}_{k,i'} H_{i,i'}^\infty + O(1/\sqrt{m_1}) \sum_{i'=1}^n \mathcal{V}_{k,i'} = \bar{x}_{i,k} + O(1/\sqrt{m_1}) \sum_{i'=1}^n \mathcal{V}_{k,i'},$$

which implies

$$\begin{aligned} |\text{trace}(W_k^+ Z_k^i) - \bar{x}_{i,k}| &\leq O(1/\sqrt{m_1}) \sqrt{n} \|\mathcal{V}_k\|_2 \\ &\lesssim O(\sqrt{n}/(\sqrt{m_1} \sqrt{\lambda_0})) \|\mathcal{V}_k\|_{H^\infty} \\ &\leq \sqrt{n/(m_1 \lambda_0)} \|\mathcal{V}_k\|_{H^\infty}. \end{aligned}$$

**Lemma 38** (Bounding the rows norm) For every  $1 \leq j \leq m_1$ , we have

$$\|W_j^+\| \leq \sqrt{nm_3}/(\sqrt{m_1 \lambda_0}) \sqrt{\sum_k \|\mathcal{V}_k\|_{H^\infty}^2}.$$

Furthermore, for every  $k \in [m_3]$ , we have

$$\|W_j^{k+}\| \leq \sqrt{n}/\sqrt{\lambda_0 m_1} \|\mathcal{V}_k\|_{H^\infty}. \quad (301)$$

For the ease of notation, because here we want to work with row sub indices of the matrix  $W_k^+$ , we refer to it by  $W^{k+}$ . **Proof of Lemma 38**

For a fixed  $1 \leq j \leq m_1$  we have with high probability over the randomness of the sign matrix  $W^s$ :

$$\begin{aligned}
\|W_j^+\| &= \left\| \sum_{k=1}^{m_3} W_j^{k+} \right\| = \left\| \sum_{k=1}^{m_3} W_{k,j}^s \frac{1}{\sqrt{m_1}} \sum_{i=1}^n \mathcal{V}_{k,i} x_i \mathbb{1}\{W_j^{(0)} x_i \geq 0\} \right\| \\
&\leq \sqrt{m_3}/\sqrt{m_1} \sqrt{\sum_{k=1}^{m_3} \left\| \sum_{i=1}^n \mathcal{V}_{k,i} x_i \mathbb{1}\{W_j^{(0)} x_i \geq 0\} \right\|^2} \\
&\leq \sqrt{m_3}/\sqrt{m_1} \sqrt{\sum_k \left( \sum_i |\mathcal{V}_{k,i}| \right)^2} \\
&\leq \sqrt{m_3}/\sqrt{m_1} \sqrt{n \sum_k \|\mathcal{V}_k\|^2} \\
&\leq \sqrt{m_3 n}/(\sqrt{m_1} \lambda_0) \sqrt{\sum_k \|\mathcal{V}_k\|_{H^\infty}^2}.
\end{aligned}$$

Furthermore, for every  $k \in [m_3]$ , we have

$$\|W_j^{k+}\| \leq \frac{1}{\sqrt{m_1}} \sum_i |\mathcal{V}_{k,i}| \leq (\sqrt{n}/\sqrt{m_1}) \|\mathcal{V}_k\|_2 \leq (\sqrt{n}/\sqrt{\lambda_0 m_1}) \|\mathcal{V}_k\|_{H^\infty}.$$

**Lemma 39** *With high probability over the initialization, we have*

$$\|W^{k+}\|_F^2 \leq (1 \pm O(n/(\lambda_0 \sqrt{m_1}))) \|\mathcal{V}_k\|_{H^\infty}^2.$$

**Proof of Lemma 39**

Recall from the definition of  $W^{k+}$  in Equation (49):

$$\begin{aligned}
\|W^{k+}\|_F^2 &= \left\| \sum_{i=1}^n \mathcal{V}_{k,i} Z_k^i \right\|_F^2 = \sum_{i=1}^n \sum_{i'=1}^n \mathcal{V}_{k,i} \mathcal{V}_{k,i'} \langle Z_k^i, Z_k^{i'} \rangle \\
&= \sum_{i=1}^n \sum_{i'=1}^n \mathcal{V}_{k,i} \mathcal{V}_{k,i'} (H_{i,i'}^\infty \pm O(1/\sqrt{m_1})) \lesssim \mathcal{V}_k^T H^\infty \mathcal{V}_k \pm O((\|\mathcal{V}_k\|_1)^2/\sqrt{m_1}) \\
&= \|\mathcal{V}_k\|_{H^\infty}^2 \pm \|\mathcal{V}_k\|_{H^\infty}^2 O(n/(\lambda_0 \sqrt{m_1})) = (1 \pm O(n/(\lambda_0 \sqrt{m_1}))) \|\mathcal{V}_k\|_{H^\infty}^2
\end{aligned}$$

## B.2.2 SOME LINEAR ALGEBRA

**Lemma 40** *For  $n \leq s$ , let  $r_1, \dots, r_n$  be  $s$ -dimensional vectors that are approximately normalized and orthogonal to one another, i.e. given some  $\delta > 0$ , for every  $1 \leq i \neq j \leq n$ :*

$$-\delta \leq \langle r_i, r_j \rangle \leq \delta, \quad \|r_i\|^2 \leq 1 + \delta.$$

*Then, for any vector  $v$  we have*

$$\sum_{i=1}^n \langle v, r_i \rangle^2 \leq (1 + \delta + n(n-1)\delta(1+\delta)^2) \|v\|^2.$$

**Proof of Lemma 40**

Define

$$v_1 = \sum_{i=1}^n \langle v, r_i \rangle r_i, \quad v_2 = v - v_1.$$

First, note that

$$\begin{aligned}
\sum_{i=1}^n \langle v, r_i \rangle^2 &\leq (1 + \delta) \sum_{i=1}^n \langle v, r_i \rangle^2 \|r_i\|^2 \\
&= (1 + \delta) \left\| \sum_i \langle v, r_i \rangle r_i \right\|^2 - 2(1 + \delta) \sum_{1 \leq i \neq j \leq n} \langle v, r_i \rangle \langle v, r_j \rangle \langle r_i, r_j \rangle \\
&= (1 + \delta) \|v_1\|^2 - 2(1 + \delta) \sum_{1 \leq i \neq j \leq n} \langle v, r_i \rangle \langle v, r_j \rangle \langle r_i, r_j \rangle.
\end{aligned}$$

Next, we write

$$\begin{aligned}
\langle v_1, v - v_1 \rangle &= \langle v, \sum_{i=1}^n \langle v, r_i \rangle r_i \rangle - \left\langle \sum_{i=1}^n \langle v, r_i \rangle r_i, \sum_{i=1}^n \langle v, r_i \rangle r_i \right\rangle \\
&= \sum_i \langle v, r_i \rangle^2 - \sum_i \langle v, r_i \rangle^2 - 2 \sum_{i \neq j} \langle v, r_i \rangle \langle v, r_j \rangle \langle r_i, r_j \rangle \\
&= -2 \sum_{i \neq j} \langle v, r_i \rangle \langle v, r_j \rangle \langle r_i, r_j \rangle.
\end{aligned}$$

Therefore

$$\begin{aligned}
\sum_{i=1}^n \langle v, r_i \rangle^2 &\leq (1 + \delta) (\|v_1\|^2 + \|v - v_1\|^2) - 2(1 + \delta) \sum_{1 \leq i \neq j \leq n} \langle v, r_i \rangle \langle v, r_j \rangle \langle r_i, r_j \rangle. \\
&\leq (1 + \delta) (\|v_1\|^2 + \|v - v_1\|^2 + 2\langle v_1, v - v_1 \rangle) \\
&\quad + 2(1 + \delta) \sum_{1 \leq i \neq j \leq n} \langle v, r_i \rangle \langle v, r_j \rangle \langle r_i, r_j \rangle \\
&\leq (1 + \delta) \|v\|^2 + 2(1 + \delta) \sum_{1 \leq i \neq j \leq n} \|v\|^2 \|r_i\| \|r_j\| \delta \\
&\leq (1 + \delta) \|v\|^2 + 2(1 + \delta) \sum_{1 \leq i \neq j \leq n} \|v\|^2 (1 + \delta) \delta \\
&= (1 + \delta) \|v\|^2 + n(n - 1) \delta (1 + \delta)^2 \|v\|^2 \\
&= (1 + \delta + n(n - 1) \delta (1 + \delta)^2) \|v\|^2,
\end{aligned}$$

which completes the proof.

In the following lemma, we state a trivial bound on the norm of  $\bar{x}_i$  based on  $\zeta_1$ .

### B.2.3 BOUND ON THE NORM OF $\bar{x}_i$ 'S

**Lemma 41** *For every  $i \in [n]$ , we have*

$$\|\bar{x}_i\| \leq \sqrt{\sum_k \|\mathcal{V}_k\|_{H^\infty}^2} = \sqrt{\zeta_1}.$$

#### Proof of Lemma 41

By definition:

$$\bar{x}_i^T = \left( H_i^\infty \mathcal{V}_k \right)_{k=1}^{m_3}.$$

Now consider the Cholskey factorization  $H^\infty = K K^T$ . Because of the assumption  $\|x_i\| = 1$ , we know that the diagonal of  $H^\infty$  is all  $1/2$ . Hence, for the  $i$ th row of  $K$  we have  $\|K_i\| = 1/2$ . Now by Cauchy-Swartz, we have

$$x_{i_1}^2 = \left\langle \sum_i \mathcal{V}_{k,i} K_i, K_{i_1} \right\rangle^2 \leq \left\| \sum_i \mathcal{V}_{k,i} K_i \right\|^2 \|K_{i_1}\|^2 = 1/2 \|\mathcal{V}_k\|_{H^\infty}^2.$$

Summing over  $i$  and noting Equation (44) completes the proof.

**Lemma 42** In the context of Lemma III, for  $\zeta_2 \leq 2nB^2$ , one can substitute  $f^*$  by  $\bar{f}_1^*$  such that

$$\begin{aligned} R_n(\bar{f}^*) &\leq 2R_n(f^*) + \frac{B^2}{n}, \\ \bar{f}^{*T} A^{-1} \bar{f}^* &\leq f^{*T} A^{-1} f^*, \end{aligned}$$

and furthermore,  $\bar{f}^*$  is in the subspace of eigenvectors of  $A$  with eigenvalue larger than  $\Omega(\frac{1}{n^2})$ .

**Proof of Lemma 42**

For an arbitrary  $i \in [n]$ , define

$$\delta = |f_i^* - \bar{f}_i^*|,$$

and suppose the slope of  $\ell(\cdot, y_i)$  at point  $f_i^*$  is equal to  $c$ . Then, using the convexity, the fact that  $\ell(y_i, y_i) = 0$ , and the 1-smoothness of  $\ell(\cdot, y_i)$ , it is not hard to see the following poincare inequality between the value and derivative of  $\ell(\cdot, y_i)$  at point  $f_i^*$ :

$$c \leq \sqrt{2\ell(f_i^*, y_i)} := 2\ell. \quad (302)$$

where from now on, for brevity, we refer to  $\ell(f_i^*, y_i)$  by  $\ell$ . Also, from the definition of  $\delta$  and again using 1 smoothness property, it is easy to see that

$$\ell(\bar{f}_i^*, y_i) \leq (c + \delta)\delta + \ell(f_i^*, y_i) = (c + \delta)\delta + \ell, \quad (303)$$

Plugging Equation (302) into (303) and using AM-GM inequality:

$$\begin{aligned} \ell(\bar{f}_i^*, y_i) &\leq \delta^2 + c\delta + \ell \leq \delta^2 + \sqrt{2\ell}\delta + \ell \\ &\leq \delta^2 + \ell + \delta^2/2 + \ell \\ &\leq 2\ell + 3\delta^2/2. \end{aligned}$$

Summing above for  $i \in [n]$ , we obtain

$$R_n(\bar{f}^*) \leq 2R_n(f^*) + 3\|f^* - \bar{f}^*\|_2^2/2. \quad (304)$$

Now we write an eigendecomposition for  $A$  as  $A = \sum_{i=1}^n \lambda_i u_i u_i^T$  for orthonormal basis  $\{u_i\}$ , and let  $f^* = \sum_i \gamma_i u_i$  be the representation of  $f^*$  in this basis. Then, from our assumption, for arbitrary  $\omega > 0$

$$\sum_i \gamma_i^2 \lambda_i^{-1} = f^{*T} A^{-1} f^* \leq 4nB^2,$$

which implies

$$\omega^{-1} \sum_{i: \lambda_i \leq \omega} \gamma_i^2 \leq 4nB^2,$$

or equivalently

$$\sum_{i: \lambda_i \leq \omega} \gamma_i^2 \leq 4nB^2\omega, \quad (305)$$

where notice that  $\sum_{i: \lambda_i \leq \omega} \gamma_i^2$  is the squared norm of the projection of  $f^*$  onto the directions whose eigenvalue is at most  $\omega$ . Now taking  $\omega = \frac{1}{12n^2}$  and defining  $\bar{f}^*$  by keeping only the directions for which  $\lambda_i > \omega$  completes the proof.

### B.3 COUPLING FOR $\hat{\nabla}_W, \hat{\nabla}_V$

In general, because the gaussian smoothing matrices  $(W^\rho, V^\rho)$  can become unbounded, the gradient estimates  $(\hat{\nabla}_W, \hat{\nabla}_V) = \nabla_{W,V} \ell(f_{W'+W^\rho, V'+V^\rho}(x_i), y_i)$  also become unbounded. However, in analyzing the stochastic behavior of SGD and showing that it can escape saddle points, it is convenient to assume the gradient's noise vector is almost surely bounded. The goal of this section is to introduce a coupling between  $(W^\rho, V^\rho)$  and another random variable that is a.s. bounded polynomially in other parameters. As that the coupled random variables take different values is exponentially small while the number of iterations in our algorithm is only polynomially large, without any concern we instead work with this new random variable, and with an overload of notation we also denote it by  $(W^\rho, V^\rho)$ .

**Lemma 43** *For an arbitrary parameter  $\chi \gg 1$ , On any pair for  $(W', V')$  with  $\|W'\| \leq C_1$ ,  $\|V'\| \leq C_2$ , there exist a mean zero random vector  $\bar{\Lambda}$  with respect to the randomness of the uniformly picked data point  $(x_i, y_i)$  and the smoothing matrices  $W^{\rho,1}, W^{\rho,2}, V^{\rho,1}$ , and  $V^{\rho,2}$  which define  $\hat{\nabla}_{W',V'}$  (meaning it is a function of those variables), such that with probability at least*

$$1 - 2 \exp\{-(\chi^2 - 1)dm_1/4\} - 2 \exp\{-(\chi^2 - 1)m_3m_2/4\} := 1 - \delta_1,$$

we have

$$\hat{\nabla}_{W',V'} = \nabla_{W',V'} L(W', V') + \bar{\Lambda}, \quad (306)$$

and finally  $\bar{\Lambda}$  is a.s. polynomially bounded, i.e. almost surely we have

$$\|\bar{\Lambda}\| \leq \text{poly}(m_1, m_2, m_3, C_1, C_2, B, \chi).$$

#### Proof of Lemma 43

Remember that  $x'_i$  was the output of the first layer (by considering the smoothing matrix  $W^\rho$ ). Now with high probability over the initialization,

$$\begin{aligned} \|\nabla_{W'} f_{W'+W^\rho, V'+V^\rho}(x_i)\|_F &= \|\nabla_{x'_i} f_{W'+W^\rho, V'+V^\rho}(x_i)^T \frac{D(x'_i)}{dW'}\| \\ &= \|\nabla_{x'_i} f_{W'+W^\rho, V'+V^\rho}(x_i)\| \left\| \frac{D(x'_i)}{dW'} \right\| \\ &= \left\| \frac{1}{\sqrt{m_2}} a^T D_{V'+V^\rho}(V^{(0)} + V' + V^\rho) \right\| \left\| \frac{1}{\sqrt{m_1}} \sum_{k=1}^{m_3} \text{diag}(W_k^s) D_{W'+W^\rho, x_i} x_i^T \right\|_F \\ &\leq (\|V^{(0)}\|_F + \|V'\|_F + \|V^\rho\|_F) \left( \frac{1}{\sqrt{m_1}} \sum_k \|\text{diag}(W_k^s) D_{W'+W^\rho, x_i} x_i^T\|_F \right) \\ &\leq (\kappa_2 \sqrt{m_2 m_3} + C_2 + \|V^\rho\|_F) \left( \frac{1}{\sqrt{m_1}} \sum_k \|\text{diag}(W_k^s) D_{W'+W^\rho, x_i} x_i^T\|_F \right) \\ &\leq (\kappa_2 \sqrt{m_2 m_3} + C_2 + \|V^\rho\|_F). \end{aligned} \quad (307)$$

On the other hand, using the final bound in Lemma 33:

$$\begin{aligned} \|\nabla_{V'} f_{W'+W^\rho, V'+V^\rho}(x_i)\|_F &= \left\| \frac{1}{\sqrt{m_2}} \text{diag}(a) D_{V'+V^\rho} x_i'^T \right\|_F \leq \|x'_i\| \\ &\leq \kappa_1 \sqrt{m_3} + C_1 + \sqrt{m_3} \sum_j \frac{1}{\sqrt{m_1}} |W_j^\rho x_i| \\ &\leq \kappa_1 \sqrt{m_3} + C_1 + \sqrt{m_3} \|W^\rho\|_F. \end{aligned} \quad (308)$$

Denoting  $\hat{\ell}(f_{W'+W^\rho, V'+V^\rho}, y_i) \nabla_{W',V'} \ell(f_{W'+W^\rho, V'+V^\rho}, y_i)$  by  $\hat{\nabla}_{W',V'}$ , then combining Equations (307) and (308) and using the 1 gradient lipschitzness property of the square loss,

$$\begin{aligned} &: \\ &\|\hat{\nabla}_{W',V'}\|_F \\ &= \left| \frac{d(\ell(f, y_i))}{df} \right| \sqrt{\|\nabla_{W'} f_{W'+W^\rho, V'+V^\rho}(x_i)\|_F^2 + \|\nabla_{V'} f_{W'+W^\rho, V'+V^\rho}(x_i)\|_F^2} \\ &\leq \left( |f_{W'+W^\rho, V'+V^\rho}(x_i)| + |B| \right) \left( \kappa_1 \sqrt{m_3} + C_1 + \sqrt{m_3} \|W^{\rho,2}\|_F + \kappa_2 \sqrt{m_2 m_3} + C_2 + \|V^{\rho,2}\|_F \right). \end{aligned} \quad (309)$$

Finally, applying Cauchy-Swartz to the second a.s. bound in Lemma B4 we have:

$$\begin{aligned} & \left| f_{W'+W^\rho, V'+V^\rho}(x_i) \right| \\ & \leq (\kappa_2\sqrt{m_3} + \|V^\rho\|_F) \left( \sqrt{m_3}\kappa_1 + C_1 + \sqrt{m_3}\|W^\rho\| \right) + C_2(C_1 + \sqrt{m_3}\|W^\rho\|). \end{aligned}$$

Combining this with (B09):

$$\begin{aligned} & \|\tilde{\nabla}_{W', V'}\|_F \\ & \leq \left[ B + (\kappa_2\sqrt{m_3} + \|V^{\rho,1}\|_F) \left( \sqrt{m_3}\kappa_1 + C_1 + \sqrt{m_3}\|W^{\rho,1}\| \right) + C_2(C_1 + \sqrt{m_3}\|W^{\rho,1}\|) \right] \\ & \quad \times \left( \kappa_1\sqrt{m_3} + C_1 + \sqrt{m_3}\|W^{\rho,2}\|_F + \kappa_2\sqrt{m_2m_3} + C_2 + \|V^{\rho,2}\|_F \right). \end{aligned}$$

Therefore, using the lipschitz bound in Theorem 9:

$$\begin{aligned} & \|\tilde{\nabla}_{W', V'} - \nabla_{W', V'} \mathbb{E}_{(x_i, y_i) \sim \mathcal{Z}} \ell(f'_{W', V'}(x_i), y_i)\|_F \\ & \leq \|\hat{\nabla}_{W', V'}\|_F + \|\mathbb{E}_{(x_i, y_i) \sim \mathcal{Z}} \nabla_{W', V'} \ell(f_{W'+W^\rho, V'+V^\rho}(x_i), y_i)\|_F \\ & \leq \left[ B + (\kappa_2\sqrt{m_3} + \|V^{\rho,1}\|_F) \left( \sqrt{m_3}\kappa_1 + C_1 + \sqrt{m_3}\|W^{\rho,1}\| \right) + C_2(C_1 + \sqrt{m_3}\|W^{\rho,1}\|) \right] \\ & \quad \times \left( \kappa_1\sqrt{m_3} + C_1 + \sqrt{m_3}\|W^{\rho,2}\|_F + \kappa_2\sqrt{m_2m_3} + C_2 + \|V^{\rho,2}\|_F \right) + (O(C_1C_2) + B)\Psi_1. \end{aligned} \tag{310}$$

Now we define the following events

$$\begin{aligned} \Xi_1 & := \{ \|W^{\rho,1}\|_F \geq \chi\sqrt{d}\beta_1 \vee \|W^{\rho,2}\|_F \geq \chi\sqrt{d}\beta_1 \}, \\ \Xi_2 & := \{ \|V^{\rho,1}\|_F \geq \chi\sqrt{m_3}\beta_2 \vee \|V^{\rho,2}\|_F \geq \chi\sqrt{m_3}\beta_2 \}, \end{aligned}$$

where recall we assume  $\chi \gg 1$ . Then, as we know the variable  $\|W^\rho\|_F^2$  has mean  $d\beta_1^2$  and is subexponential with parameters  $(d\beta_1^4/m_1, \beta_1^2/m_1)$ . Hence, by a union bound and Bernstein (Note that  $W^{\rho,1}, W^{\rho,2}$  are independent):

$$\begin{aligned} & \mathbb{P}(\Xi_1) \\ & \leq 2\mathbb{P}(\|W^\rho\|_F \geq \chi\beta_1\sqrt{d}) \\ & = 2\mathbb{P}(\|W^\rho\|_F^2 \geq \chi^2\beta_1^2d) \\ & \leq 4 \max \left( \exp\{-(\chi^2 - 1)^2\beta_1^4d^2/(4d\beta_1^4/m_1)\}, \exp\{-(\chi^2 - 1)\beta_1^2d/(4\beta_1^2/m_1)\} \right) \\ & = 4 \max \left( \exp\{-(\chi^2 - 1)^2dm_1/4\}, 2 \exp\{-(\chi^2 - 1)dm_1/4\} \right) \leq 2 \exp\{-(\chi^2 - 1)dm_1/4\}. \end{aligned}$$

Similarly for  $\|V^\rho\|_F$ :

$$\mathbb{P}(\Xi_2) = \mathbb{P}(\|V^\rho\|_F \geq \chi\beta_2\sqrt{m_3}) \leq 4 \exp\{-(\chi^2 - 1)m_3m_2/4\}.$$

Moreover, because of the subexponential tails of  $\|W^\rho\|_F^2$  and  $\|V^\rho\|_F^2$ , for each of  $W^{\rho,1}$  or  $W^{\rho,2}$ ,  $V^{\rho,1}$ , or  $V^{\rho,2}$ :

$$\begin{aligned} \mathbb{E}(\|W^\rho\|_F | \Xi_1) & \lesssim \chi\sqrt{d}\beta_1, \quad \mathbb{E}(\|W^\rho\|_F^2 | \Xi_1) \lesssim \chi^2d\beta_1^2. \\ \mathbb{E}(\|V^\rho\|_F | \Xi_2) & \lesssim \chi\sqrt{m_3}\beta_2, \quad \mathbb{E}(\|V^\rho\|_F^2 | \Xi_2) \lesssim \chi^2m_3\beta_2^2. \end{aligned}$$

Now Defining  $\Xi = \Xi_1 \cup \Xi_2$  and combining the above equations:

$$\begin{aligned} \mathbb{E}(\|W^\rho\| \mathbb{1}\{\Xi\}) & \leq \mathbb{E}\|W^\rho\| (\mathbb{1}\{\Xi_1\} + \mathbb{1}\{\Xi_2\}) = \mathbb{E}(\|W^\rho\| | \Xi_1) \mathbb{P}(\Xi_1) + \mathbb{E}(\|W^\rho\| | \Xi_2) \mathbb{P}(\Xi_2) \\ & \lesssim \chi\sqrt{d}\beta_1 2 \exp\{-(\chi^2 - 1)dm_1/4\} + \sqrt{d}\beta_1 2 \exp\{-(\chi^2 - 1)m_3m_2/4\} \\ & = 2\sqrt{d}\beta_1 (\chi \exp\{-(\chi^2 - 1)dm_1/4\} + \exp\{-(\chi^2 - 1)m_3m_2/4\}), \end{aligned}$$

and

$$\begin{aligned}\mathbb{E}\|W^\rho\|^2\mathbb{1}\{\Xi_1\} &= \mathbb{E}(\|W^\rho\|^2 | \Xi_1)\mathbb{P}(\Xi_1) \\ &\leq 2d\beta_1^2\chi^2 \exp\{-(\chi^2 - 1)dm_1/4\}.\end{aligned}$$

Similarly

$$\mathbb{E}(\|V^\rho\|\mathbb{1}\{\Xi\}) \leq 2\sqrt{m_3}\beta_2(\exp\{-(\chi^2 - 1)dm_1/4\} + \chi \exp\{-(\chi^2 - 1)m_3m_2/4\}),$$

and

$$\mathbb{E}(\|V^\rho\|^2\mathbb{1}\{\Xi_2\}) \leq 2m_3\beta_2^2\chi^2 \exp\{-(\chi^2 - 1)m_3m_2/4\}.$$

Applying these equations to [\(B10\)](#) with Cauchy-Schwartz to get the upper bounds  $\mathbb{E}\mathbb{1}\{\Xi_1\}\|W^{\rho,1}\|\|W^{\rho,2}\| \leq \mathbb{E}\mathbb{1}\{\Xi_1\}\|W^\rho\|^2$  and  $(\mathbb{E}_{W^\rho}\|W^\rho\|)^2 \leq \mathbb{E}_{W^\rho}\|W^\rho\|^2$  (for terms with only one  $W^{\rho,i}$  or  $V^{\rho,i}$ , we simply write them as  $W^\rho$  and  $V^\rho$ ):

$$\begin{aligned}\mathbb{E}_{W^\rho, V^\rho} \|\tilde{\nabla}_{W', V'} - \nabla_{W', V'} \mathbb{E}_{(x_i, y_i) \sim \mathcal{Z}} \ell(f'_{W', V'}(x_i), y_i)\|_F \mathbb{1}\{\Xi\} \\ \leq \mathbb{E}_{W^\rho, V^\rho, (x_i, y_i)} \left[ B + (\kappa_2\sqrt{m_3} + \|V^{\rho,1}\|_F) \left( \sqrt{m_3}\kappa_1 + C_1 + \sqrt{m_3}\|W^\rho\| \right) + C_2(C_1 + \sqrt{m_3}\|W^\rho\|) \right] \\ \times \left( \kappa_1\sqrt{m_3} + C_1 + \sqrt{m_3}\|W^\rho\|_F + \kappa_2\sqrt{m_2m_3} + C_2 + \|V^\rho\|_F \right) + \alpha(O(C_1C_2) + B)\Psi_1 \\ = \left[ B + (\kappa_2\sqrt{m_3})(\sqrt{m_3}\kappa_1 + C_1) + C_1C_2 \right] \\ \times \left( \kappa_1\sqrt{m_3} + C_1 + \kappa_2\sqrt{m_2m_3} + C_2 + \sqrt{m_3}\mathbb{E}_{W^\rho}\mathbb{1}\{\Xi\}\|W^\rho\|_F + \mathbb{E}_{V^\rho}\mathbb{1}\{\Xi\}\|V^\rho\|_F \right) \\ + \left( B + (\sqrt{m_3}\kappa_1 + C_1)\sqrt{m_3} + (\kappa_2m_3 + C_2\sqrt{m_3}) + \sqrt{m_3}(\kappa_1\sqrt{m_3} + C_1 + \kappa_2\sqrt{m_2m_3} + C_2) \right) \\ \times (\mathbb{E}_{W^\rho}\mathbb{1}\{\Xi_1\}\|W^\rho\|_F \mathbb{E}_{V^\rho}\|V^\rho\|_F + \mathbb{E}_{W^\rho}\|W^\rho\|_F \mathbb{E}_{V^\rho}\mathbb{1}\{\Xi_2\}\|V^\rho\|_F) \\ + (B + \sqrt{m_3}\kappa_1 + C_1)(\mathbb{E}_{V^\rho}\mathbb{1}\{\Xi_2\}\|V^\rho\|^2 + \mathbb{P}(\Xi_1)\mathbb{E}\|V^\rho\|^2) \\ + C_2m_3(\mathbb{E}_{W^\rho}\mathbb{1}\{\Xi_1\}\|W^\rho\|^2 + \mathbb{P}(\Xi_2)\mathbb{E}\|W^\rho\|^2) \\ + m_3(\mathbb{E}_{W^\rho}\mathbb{1}\{\Xi_1\}\|W^\rho\|^2\mathbb{E}_{V^\rho}\|V^\rho\| + \mathbb{E}_{W^\rho}\|W^\rho\|^2\mathbb{E}_{V^\rho}\mathbb{1}\{\Xi_2\}\|V^\rho\|) \\ + \sqrt{m_3}(\mathbb{E}_{V^\rho}\mathbb{1}\{\Xi_2\}\|V^\rho\|^2\mathbb{E}_{W^\rho}\|W^\rho\| + \mathbb{E}_{V^\rho}\|V^\rho\|^2\mathbb{E}_{W^\rho}\mathbb{1}\{\Xi_1\}\|W^\rho\|) \\ + (O(C_1C_2) + B)\Psi_1\mathbb{P}(\Xi) \\ \leq (\exp\{-(\chi^2 - 1)dm_1/4\} + \chi \exp\{-(\chi^2 - 1)m_3m_2/4\})\text{poly}(m_1, m_2, m_3) = \text{negligible}.\end{aligned}\tag{311}$$

But note that

$$\begin{aligned}\hat{\nabla}_{W', V'} &= \tilde{\nabla}_{W', V'} + \nabla_{W', V'}(\psi_1\|W'\|^2 + \psi_2\|V'\|^2), \\ \nabla_{W', V'}L(W', V') &= \nabla_{W', V'}\mathbb{E}_{(x_i, y_i) \sim \mathcal{Z}} \ell(f'_{W', V'}(x_i), y_i) + \nabla_{W', V'}(\psi_1\|W'\|^2 + \psi_2\|V'\|^2).\end{aligned}$$

Applying this to Equation [\(B11\)](#), we get that if we define

$$\Lambda := \hat{\nabla}_{W', V'} - \nabla_{W', V'}L(W', V'),$$

then

$$\mathbb{E}\|\Lambda\mathbb{1}\{\Xi\}\| \leq (\exp\{-(\chi^2 - 1)dm_1/4\} + \chi \exp\{-(\chi^2 - 1)m_3m_2/4\})\text{poly}(m_1, m_2, m_3).$$

On the other hand, note that using again Equation [\(B10\)](#), we have the following a.s. bound:

$$\|\Lambda\mathbb{1}\{\Xi^c\}\| = \text{poly}(m_1, m_2, m_3, C_1, C_2, \chi).$$

Defining

$$\begin{aligned}\Lambda_1 &= \Lambda\mathbb{1}\{\Xi^c\}, \\ \Lambda_2 &= \mathbb{1}\{\Xi\}\mathbb{E}(\Lambda|\Xi), \\ \bar{\Lambda} &= \Lambda_1 + \Lambda_2,\end{aligned}$$

we get that with probability at least  $1 - \mathbb{P}(\Xi)$ :

$$\hat{\nabla}_{W',V'} = \nabla_{W',V'} L(W', V') + \bar{\Lambda}$$

and also note that

$$\mathbb{E}\bar{\Lambda} = \mathbb{E}\Lambda = 0.$$

Finally by the a.s. bound for  $\Lambda_1$ , we have a.s.:

$$\begin{aligned} \|\bar{\Lambda}\| &\leq \mathbb{E}_{W^\rho, V^\rho, (x_i, y_i)} \|\Lambda \mathbb{1}\{\Xi\}\| + \|\Lambda \mathbb{1}\{\Xi^c\}\| \\ &\leq (\exp\{-(\chi^2 - 1)dm_1/4\} + \chi \exp\{-(\chi^2 - 1)m_3m_2/4\}) \text{poly}(m_1, m_2, m_3) + \text{poly}(m_1, m_2, m_3) \\ &= \text{poly}(m_1, m_2, m_3, C_1, C_2, B, \chi), \end{aligned}$$

which completes the proof.

**Corollary 9.1** *It is easy to check that running PSGD with unbiased gradient estimate  $\hat{\nabla}_{W',V'}$  is equivalent to running SGD after our change of coordinates, with unbiased gradient estimate  $\hat{\nabla}_{w',v'} := \Upsilon \nabla_{W',V'}$ , where  $\Upsilon$  is the matrix for our change of coordinate, which is equal to  $\Upsilon$  defined in Appendix A.10 for the coordinates in  $V'$  and simply identity for the coordinates in  $W'$ . Therefore, projecting both sides in Equation (B.2) of Lemma B.3 onto  $\Phi^\perp$  by multiplying  $\Upsilon$  implies that with high probability for all iterations of the algorithm*

$$\begin{aligned} \hat{\nabla}_{w',v'} &= \Upsilon \nabla_{W',V'} L^\Pi(W', V') + \Upsilon \bar{\Lambda} \\ &= \nabla_{w',v'} L^\Pi(w', v') + \Upsilon \bar{\Lambda}, \end{aligned} \quad (312)$$

where  $\bar{\mathcal{L}} := \Upsilon \bar{\Lambda}$  (using the properties of  $\bar{\Lambda}$  in Lemma B.3) is a mean zero noise vector with almost surely bounded norm, i.e.  $\|\bar{\mathcal{L}}\| \leq Q'$  for some  $Q' = \text{poly}(m_1, m_2, m_3, C_1, C_2)$ . (we dropped the  $\chi$  parameters by considering constant high probability argument).

Finally, note that injecting noise  $(\Xi_1/\|\Xi_1\|, \Xi_2/(\sqrt{m_1}\|\Xi_2\|))$  by PSGD results in adding an extra zero mean noise  $(\tilde{\Xi}_1, \tilde{\Xi}_2) := (\Upsilon \Xi_1/\|\Xi_1\|, \Upsilon \Xi_2/(\sqrt{m_1}\|\Xi_2\|))$  to the gradient  $\nabla_{w',v'} L^\Pi(w', v')$ . Therefore, overall running SGD on  $L^\Pi$  (which is equivalent to PSGD on  $L$ ) observe an unbiased noise vector defined as  $\mathcal{L} := \bar{\mathcal{L}} + (\tilde{\Xi}_1, \tilde{\Xi}_2)$ . Now it is easy to check that the moment matrix of  $\tilde{\Xi}_1$  and  $\tilde{\Xi}_2$  are  $\sigma_1'^2 I$  and  $\sigma_2'^2 I$  for

$$\sigma_1'^2 := \frac{1}{m_1^2 d}, \quad (313)$$

$$\sigma_2'^2 := \frac{m_2(m_3 - n)}{m_2^m m_3}, \quad (314)$$

which implies the moment matrix of  $\mathcal{L}$  is upper bounded by

$$\sigma_2^2 I := (Q'/(m_2 m_3 + m_1 d) + \max\{\sigma_1'^2, \sigma_2'^2\}) I,$$

and lower bounded by

$$\sigma_2^2 I := \min\{\sigma_1'^2, \sigma_2'^2\} I,$$

i.e.

$$\sigma_1^2 I \mathbb{E} \mathcal{L} \mathcal{L}^T \leq \sigma_2^2 I.$$

(Note that we look at the new coordinates  $(w', v')$  as a vectors, so the term  $\mathbb{E} \mathcal{L} \mathcal{L}^T$  makes sense.)

Moreover,  $\|\tilde{\Xi}_1\| = 1/\sqrt{m_1}$ ,  $\|\tilde{\Xi}_2\| \leq 1$  almost surely, which implies the following almost surely bound for  $\mathcal{L}$ :

$$\|\mathcal{L}\| \leq Q := Q' + 1 + 1/\sqrt{m_1}.$$

**Lemma 44** Let  $g(x)$  be a second order differentiable function over  $\mathbb{R}^N$  such that at point  $x$ , there exist a random direction  $y$  and deterministic direction  $z$  and fixed positive real  $r$  with:

$$\begin{aligned}\mathbb{E}y &= 0, \\ \mathbb{E}_y g(x + \eta z + \sqrt{\eta}y) &\leq g(x) - \eta r.\end{aligned}$$

Then, for the gradient and Hessian at point  $x$ , we have either

$$\|\nabla g(x)\| \geq \frac{r}{4\|z\|},$$

or

$$\lambda_{\min}(\nabla^2 g(x)) \leq -\frac{r}{2\|y\|^2}.$$

**Proof of Lemma 44**

We write the second order Taylor approximation of  $g$  around  $x$ :

$$g(x + w) = g(x) + \nabla g(x)^T w + \frac{1}{2} w^T \nabla^2 g(x) w + o(\|w\|^2).$$

Now substituting  $w$  with  $\eta z + \sqrt{\eta}y$  and taking expectation with respect to  $y$ , as we send  $\eta \rightarrow 0$  and using the fact that  $\mathbb{E}y = 0$ :

$$\begin{aligned}\mathbb{E}_y g(x + \eta z + \sqrt{\eta}y) &= g(x) + \mathbb{E}_y \nabla g(x)^T (\eta z + \sqrt{\eta}y) + \frac{1}{2} (\eta z + \sqrt{\eta}y)^T \nabla^2 g(x) (\eta z + \sqrt{\eta}y) + o(\|\eta z + \sqrt{\eta}y\|^2) \\ &= \mathbb{E}_y g(x) + \eta \nabla g(x)^T z + \frac{1}{2} \eta^2 z^T \nabla^2 g(x) z + \eta \frac{1}{2} y^T \nabla^2 g(x) y + o(\eta \|y\|^2) \\ &= \mathbb{E}_y g(x) + \eta \nabla g(x)^T z + \eta \frac{1}{2} y^T \nabla^2 g(x) y + o(\eta).\end{aligned}$$

Combining the assumption with the above Equation, we get that for small enough  $\eta$ , we have

$$\eta \nabla g(x)^T z + \eta \frac{1}{2} \mathbb{E}_y y^T \nabla^2 g(x) y \leq -\eta r/2,$$

i.e.

$$\nabla g(x)^T z + \frac{1}{2} \mathbb{E}_y y^T \nabla^2 g(x) y \leq -r/2,$$

which means we should either have

$$\nabla g(x)^T z \leq -r/4,$$

which implies

$$\|\nabla g(x)\| \geq \frac{r}{4\|z\|},$$

or

$$\mathbb{E}_y y^T \nabla^2 g(x) y \leq -r/2,$$

which implies

$$\lambda_{\min}(\nabla^2 g(x)) \leq -\frac{r}{2 \max_{\tilde{y} \in \text{support}(y)} \|\tilde{y}\|^2}.$$

#### B.4 HANDLING THE INJECTED NOISE BY PSGD

In this section, we prove that having SGD injecting noise into our gradient estimates mostly does not change the sign pattern of the first layer, namely among the set of rows in  $P$  defined in Lemma [4](#).

**Lemma 45** *Having enough overparameterization, with high probability, at every iteration of the PSGD for  $W^{(2)}$  defined in Lemma [4](#), we have for every  $j \in [m_1]$ :*

$$\|W_j^{(2)}\| \leq c_2/(4\sqrt{m_1}).$$

#### Proof of Lemma [45](#)

Let  $\Phi'$  be the subspace of the first layer weight matrices which is zero in rows  $j \in P$  ( $P$  is defined in Lemma [4](#)), while in other rows it is the span of  $Z_i^k$ 's, i.e. using our notation  $\tilde{Z}_k^i$  introduced in the proof of Lemma [4](#), we can write  $\Phi'$  is  $\text{span}(Z_i^k)_{i,k}$ .

Recall from Lemma [4](#) that we decompose the first layer weight  $W'$  as  $W'^{(1)} + W'^{(2)}$ , namely the parts in the subspace  $\Phi'$  and subspace  $\Phi'^{\perp}$  respectively. Moreover, let  $\Xi_1/(\sqrt{m_1}\|\Xi_1\|) = \Xi^{(1)} + \Xi^{(2)}$  be the decomposition of the injected noise at some iteration of PSGD into subspaces  $\Phi'$  and  $\Phi'^{\perp}$  respectively.

Now recall that the current  $W'$  is the value of the previous iteration moved by the gradient plus the injected noise:

$$\begin{aligned} W' &= W' - \eta(\hat{\nabla}_{W'} + \Xi^{(1)} + \Xi^{(2)}) \\ &= W' - \eta\left(\tilde{\nabla}_{W'} + 2\psi_1 W'^{\cdot,(1)} + 2\psi_1 W'^{\cdot,(2)} + \Xi^{(1)} + \Xi^{(2)}\right), \end{aligned}$$

where  $W'$  is the weight of the previous iteration and  $W'^{\cdot,(1)}, W'^{\cdot,(2)}$  are again its decomposition to  $\Phi'$  and  $\Phi'^{\perp}$ , where  $\tilde{\nabla}_{W',V'}$  is defined in Lemma [43](#). Applying Lemma [24](#) for the previous iteration of the algorithm, we get  $\tilde{\nabla}_{W'} \in \Phi'$  since the bad events  $E$  defined in Lemma [43](#) occurs only with probability exponentially small (hence union bound across all the iterations rules it out). Hence, the decomposition for the current iteration becomes

$$W'^{(1)} = W'^{\cdot,(1)} - \eta(\hat{\nabla}_{W'} + 2\psi_1 W'^{\cdot,(1)} + \Xi^{(1)}), \quad (315)$$

$$W'^{(2)} = (1 - 2\eta\psi_1)W'^{\cdot,(2)} + \eta\Xi^{(2)}. \quad (316)$$

We handle the  $W'^{(1)}$  part in Lemma [4](#) and prove that as long as  $\|W'^{(1)}\|^2 \leq \|W'\|^2$  remains bounded by  $C_1^2$ , then the sign pattern of the first layer, when only considering the  $W'^{(1)}$  part, is specified by the initialization except within set  $P$ ; here we handle the  $W'^{(2)}$  part as well.

Note that for every row  $j \in [m_1]$ , the variable  $\left\|\left(\Xi_1/(\sqrt{m_1}\|\Xi_1\|)\right)_j\right\|^2$  is  $(O(1/(m_1^4 d)), O(1/(m_1^2 d))$ -subexponential with mean  $1/m_1$ . Therefore, with probability that is exponentially small in  $m_1$ ,  $\left\|\left(\Xi_1/(\sqrt{m_1}\|\Xi_1\|)\right)_j\right\|$  is bounded by  $O(1/m_1)$ . It is not hard to see the same argument holds for the projection of  $\Xi_1/(\sqrt{m_1}\|\Xi_1\|)$  onto  $\Phi^{\perp}$ , i.e.  $\Xi^{(2)}$ . Applying a union bound for all iterations, again using the fact that we run PSGD for poly iterations while the chance of error is exponentially small in  $m_1$ , we can then argue that with high probability over the noise of gradients, at every iteration and for every  $j \in [m_1]$ :

$$\|\Xi_j^{(2)}\| = \tilde{O}(1/m_1). \quad (317)$$

But applying triangle inequality to Equation [\(316\)](#) and writing it in a telescope form, particularly for the  $j$ th row, and further using the assumption in [317](#), we get that  $\|W_j'^{(2)}\|$  grows at most to  $O(1/(m_1\psi_1))$ ; as we set  $1/\psi_1 = O(\text{poly}(n))$ , assuming polynomially large enough  $m_1$  concludes the claim.

#### B.4.1 BOUNDING THE NORM OF THE FIRST LAYER'S OUTPUT IN THE WORST CASE

**Lemma 46** Suppose  $W'$  satisfies the assumption of Lemma 4, i.e.  $\|W'\| \leq C_1$ , and  $\|W'_j\| \leq c_2/(2\sqrt{m_1})$  except possible for indices in  $P$ , also defined in 4. Then, with high probability over initialization

$$\sup_{x, \|x\|=1} \|\phi'(x)\| \lesssim (1 + O(m_3^2 d^2 \log(m_1)^2 / m_1)) C_1 + \sqrt{m_3} \frac{C_1^{3/2}}{\sqrt{\kappa_1}} \left(\frac{n^3 m_3}{m_1 \lambda_0}\right)^{1/4},$$

which is  $O(C_1)$  for large enough overparameterization.

#### Proof of Lemma 46

Note that because the VC-dimension of the class of binary functions with respect to half-spaces in  $\mathbb{R}^d$  is  $d + 1$ , the number of different sign patterns  $D_{W^{(0)}, x}$  for different  $x$  can be at most  $m_1^{d+1}$ . Now similar to Equation (318), for  $k \in [m_3]$  define

$$Z_k(x) = 1/\sqrt{m_1} \left( W_{k,j}^s \mathbb{1}\{W_j^{(0)T} x\} \right)_{j=1}^{m_1}. \quad (318)$$

Then, for  $k_1 \neq k_2$ , as  $\|x\| = 1$ :

$$\begin{aligned} \langle Z_{k_1}(x), Z_{k_2}(x) \rangle &= \frac{1}{m_1} \sum_{j=1}^{m_1} W_{k_1,j}^s W_{k_2,j}^s \mathbb{1}\{W_j^{(0)T} x\} \\ &\leq \frac{1}{m_1} \sup_x \sum_{j=1}^{m_1} D_{W^{(0)}, x, j, j} W_{k_1,j}^s W_{k_2,j}^s. \end{aligned}$$

But for each fixed  $D_{W^{(0)}, x}$ , using Hoeffding bound, we have with probability  $1 - \delta$ :

$$\frac{1}{\sqrt{m_1}} \sum_{j=1}^{m_1} D_{W^{(0)}, x, j, j} W_{k_1,j}^s W_{k_2,j}^s \lesssim \sqrt{\frac{\log(1/\delta)}{m_1}}.$$

Applying the above for all possible sign patterns with  $\delta < O(1/m_1^{d+1})$  and a union bound, we have with high probability

$$\sup_{x, \|x\|=1} \langle Z_{k_1}(x), Z_{k_2}(x) \rangle \leq \frac{1}{m_1} \sup_x \sum_{j=1}^{m_1} D_{W^{(0)}, x, j, j} W_{k_1,j}^s W_{k_2,j}^s \lesssim d \log(m_1) / \sqrt{m_1}.$$

We can even state the following stronger bound with respect to two adversarially picked vectors  $x, x'$ :

$$\sup_{\|x\|=1, \|x'\|=1} \langle Z_{k_1}(x), Z_{k_2}(x') \rangle \leq \frac{1}{m_1} \sup_x \sum_{j=1}^{m_1} D_{W^{(0)}, x, j, j} D_{W^{(0)}, x', j, j} W_{k_1,j}^s W_{k_2,j}^s \lesssim d \log(m_1) / \sqrt{m_1}, \quad (319)$$

because each  $D_{W^{(0)}, x', j, j}$  has at most  $m_1^{d+1}$  cases as we discussed above, then  $D_{W^{(0)}, x, j, j} D_{W^{(0)}, x', j, j}$  has at most  $m_1^{d+1}$  possible cases, and applying a similar Hoeffding bound for each of them and a union bound as we did will imply (319). We will use this generalized version in another section.

Now combining Equation (319) with the fact that  $\|W'\| \leq C_1$  and applying Lemma 40:

$$\sup_{x, \|x\|=1} \sum_{k=1}^{m_3} \langle W', Z_k(x) \rangle^2 \leq (1 + O(m_3^2 d^2 \log(m_1)^2 / m_1)) C_1^2. \quad (320)$$

On the other hand, setting  $m_2 = m_1$ ,  $m_3 = d$ , and  $R = c_2/(2\sqrt{m_1} \kappa_1)$  in Lemma 29, we get with high probability

$$\#\left(j \in [m] : |V_j^{(0)} x| \leq c_2/(2\sqrt{m_1})\right) \leq m_1 c_2 / (2\sqrt{m_1}) = \sqrt{m_1} c_2 / (2\kappa_1).$$

Noting that  $\|W'_j\| \leq c_2/(2\sqrt{m_1})$  for  $j \notin P$ , we conclude that with high probability, for any  $x$ ,  $\mathbb{1}\{(W^{(0)} + W')^T_j x \geq 0\}$  and  $\mathbb{1}\{W_j^{(0)T} x \geq 0\}$  can be different in at most  $\sqrt{m_1}c_2/(2\kappa_1)$  of the  $j$ 's outside of  $[m_1] \setminus P$ . Therefore, as we have also  $|P| \lesssim nc_2\sqrt{m_1}/\kappa_1$  from Lemma III, we conclude that with high probability, for any  $x$ , there are at most  $O(nc_2\sqrt{m_1}/\kappa_1)$  sign changes by adding  $W'$  to  $W^{(0)}$ . This further implies:

$$\begin{aligned} |\phi'_k(x) - \langle W', Z_k(x) \rangle| &\leq 2/\sqrt{m_1} \sum_{j: \text{Sgn}(W_j^{(0)T} x) \neq \text{Sgn}((W^{(0)} + W')^T_j x)} |W'_j x| \\ &\leq \|W'\| 2\sqrt{|\{j | \text{Sgn}(W_j^{(0)T} x) \neq \text{Sgn}((W^{(0)} + W')^T_j x)\}|} / \sqrt{m_1} \\ &\lesssim C_1 \sqrt{nc_2\sqrt{m_1}/\kappa_1} / \sqrt{m_1} \\ &= \frac{C_1^{3/2}}{\sqrt{\kappa_1}} \left(\frac{n^3 m_3}{m_1 \lambda_0}\right)^{1/4}. \end{aligned}$$

Combining this with (B20), we conclude with high probability:

$$\sup_{x, \|x\|=1} \|\phi'(x)\| \lesssim (1 + O(m_3^2 d^2 \log(m_1)^2 / m_1)) C_1 + \sqrt{m_3} \frac{C_1^{3/2}}{\sqrt{\kappa_1}} \left(\frac{n^3 m_3}{m_1 \lambda_0}\right)^{1/4},$$

which completes the proof.

You may include other additional sections here.