

# ON THE DYNAMICS OF LEARNING LINEAR FUNCTIONS WITH NEURAL NETWORKS

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

This paper studies the gradient descent training dynamics of fitting a one-hidden-layer network with multi-dimensional outputs to linear target functions. That is, we focus on a realizable model where the inputs are drawn i.i.d. from a Gaussian distribution and the labels are generated according to a planted linear model with multiple outputs. This framework serves as a good model for a variety of interesting problems including end-to-end training in inverse problems and various auto-encoder models in machine learning. Despite the seemingly simple formulation, understanding training dynamics is a challenging unresolved problem. This is in part due to the fact that the training landscape contains multiple non-strict saddle points and it is completely unclear why gradient descent from random initialization is able to escape such bad stationary points. In this work, we develop the first comprehensive analysis of the gradient descent dynamics for learning linear target functions with ReLU networks. We provide a comprehensive characterization of the optimization landscape. Furthermore, we show that gradient descent with moderately small random initialization converges to a global minimizer at a linear rate with an order-wise optimal sample complexity. To rigorously show that GD avoids non-strict saddle points, we develop intricate techniques to decompose the loss and control the GD trajectory, which may have broader implications for the analysis of non-convex optimization problems involving non-strict saddles. We corroborate our theoretical results with extensive experiments with various configurations.

## 1 INTRODUCTION

### 1.1 MOTIVATION

End-to-end training of neural networks (NNs) via Gradient Descent (GD) has recently achieved remarkable success on many tasks. Of particular interest, these models have been adopted to solve inverse problems by taking the measurements as input and mapping them directly to the desired signal with successful scientific applications in computer vision (Ledig et al., 2017; Wang et al., 2018), MRI reconstruction (Sriram et al., 2020; Fabian et al., 2022), sparse-view computed tomography (CT) (Jin et al., 2017b), and phase retrieval (Hand et al., 2018). These models not only fit the training data but also appear to capture useful features and nuanced priors that enable them to generalize to unseen test examples. Despite this empirical success, the reasons behind the success of NNs for end-to-end training and how they can extract useful features from data remain unclear.

Perhaps the most classical form of end-to-end training is those arising in autoencoder type problems where the goal is to teach a neural network to learn a linear mapping (e.g. identity for autoencoders). Surprisingly, the dynamics of training of such a model is not well understood for nonlinear models. For linear networks, a classical result by Baldi and Hornik (1989) provided a complete characterization, showing how gradient descent recovers the principal components of the data. In contrast, understanding the dynamics of non-linear encoders has remained an open and challenging problem even for simple target functions. In this paper we aim to take a step towards a systematic understanding of the training dynamics of such problems by addressing the following question:

How do the dynamics of training ReLU neural networks with gradient descent starting from random initialization facilitate learning simple priors and structures such as linear target functions?

Despite significant recent progress in understanding neural networks (especially shallow networks) (Chizat et al., 2019; Soltanolkotabi et al., 2018; Jacot et al., 2018; Du et al., 2018; Ongie et al., 2019), many aspects of the dynamics of GD and how it facilitates learning remain mysterious even in seemingly simple settings. A particularly simple one involves learning linear target functions via GD, that is, teaching a one-hidden-layer network to mimic the output of a simple linear model. Surprisingly, understanding the dynamics of GD in this simple setting has remained elusive. Although there are many results on learning specific target functions such as ReLUs (Xu and Du, 2023; Soltanolkotabi, 2017) and polynomials (Damian et al., 2022), these results typically exclude linear function classes. In fact, many of the existing papers use a pre-processing step or alter the early optimization trajectory to avoid complications arising from the dynamics of learning linear functions (Damian et al., 2022). This is in part due to the fact that the optimization landscape of learning linear target functions contains multiple non-strict saddle points (i.e. where the gradient vanishes and the Hessian is PSD but has a 0 eigenvalue) requiring a subtle trajectory analysis to ensure GD avoid these bad points (See Section 1.3 for further details). We note that despite the simple formulation, quite a few interesting scenarios, including autoencoder training dynamics, are captured in this framework.

Our main contributions are as follows:

- To gain insight into the inner working of nonlinear autoencoders, we focus our attention on learning linear target functions using one-hidden-layer Neural Networks (NNs) via Gradient Descent (GD). We empirically observe an interesting pattern in the weights of the neural network with exact parametrization (when the number of hidden units is exactly twice the number of target directions). We find that GD iterations converge to a solution where hidden units cluster into *pairs*: incoming and outgoing weights from these pairs are the negative of each other (Figure 3).
- Fixing the outer layer of the NN according to the said sign pattern, with exact parameterization, we develop theory for running GD on the NN with moderately small initialization, demonstrating exact convergence to the ground truth at a linear rate and with an optimal sample complexity that scales linearly in the number of parameters. That is, we show that the inner weights of the NN recover the target directions *exactly*.
- As detailed further in Section 1.3 the training landscape studied in this paper contains multiple non-strict saddles. To prove that the trajectory of GD from moderately small random initialization avoids these bad minima we develop new techniques to control the GD trajectory which we combine with intricate uniform concentration bounds. We believe our novel proof techniques may have broader implications for the analysis of non-convex optimization problems involving non-strict saddles.
- We further corroborate our results with various experimental investigations.

## 1.2 PROBLEM FORMULATION

We first state the general family of problems of interest in this paper.

**Data Model.** We assume there are  $n$  pairs of training data consisting of input features  $\mathbf{x}_i \in \mathbb{R}^d$  and corresponding targets  $\mathbf{y}_i \in \mathbb{R}^r$ . As mentioned before, we consider the class of linear models where the relationship between  $\mathbf{x}_i$  and  $\mathbf{y}_i$  is given by the equation:  $\mathbf{y}_i = \mathbf{A}\mathbf{x}_i$  where  $\mathbf{A} \in \mathbb{R}^{r \times d}$  is the labeling matrix. Conceptually,  $\mathbf{A}$  contains  $r$  target *directions* ( $\mathbf{a}_1, \dots, \mathbf{a}_r$ ) that our predictor should *learn*. In the special case when  $r = 1$  (there is a *single* direction to be learned), we use  $\mathbf{a}^T$  instead of  $\mathbf{A}$  to emphasize this single direction  $\mathbf{a}$ . For our theoretical analysis we assume the data points  $\mathbf{x}_i$  are drawn i.i.d. according to a standard normal distribution  $\mathcal{N}(\mathbf{0}, \mathbf{I}_d)$ .

**Network Model.** We consider one-hidden-layer neural networks of the form  $\hat{\mathbf{y}} = f(\boldsymbol{\theta}, \mathbf{x}) = f(\mathbf{V}, \mathbf{W}, \mathbf{x}) = \mathbf{V}^T \phi(\mathbf{W}\mathbf{x})$  as our predictor. Here  $k$  denotes the number of hidden-units,  $\mathbf{V} \in \mathbb{R}^{k \times r}$  is the outer layer of the neural network,  $\mathbf{W} \in \mathbb{R}^{k \times d}$  is the inner layer of the neural network, and  $\phi(z)$  is the activation function. We refer to individual rows of  $\mathbf{V}/\mathbf{W}$  as  $\mathbf{v}_i/\mathbf{w}_i$  respectively. In this paper,

we specifically consider neural networks with ReLU activation functions i.e.  $\phi(\mathbf{z}) = \text{ReLU}(\mathbf{z}) = \max(0, \mathbf{z})$ , where max is applied to the input vector  $\mathbf{z}$  element-wise.

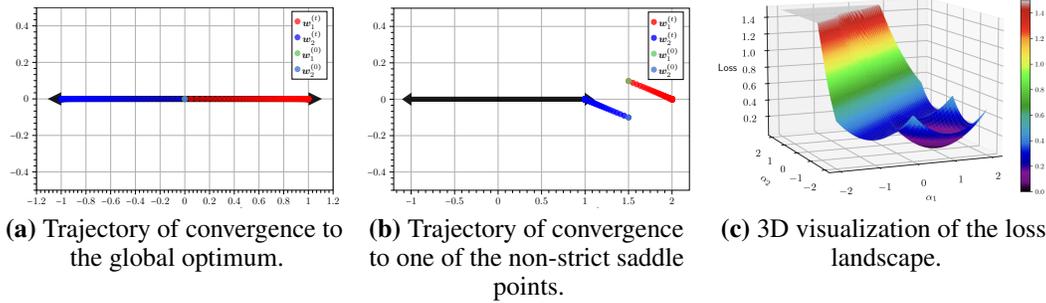
**Training Loss.** We minimize the squared loss between the target and the prediction

$$\mathcal{L}(\theta) = \frac{1}{2n} \sum_{i=1}^n \|\hat{\mathbf{y}}_i - \mathbf{y}_i\|^2 = \frac{1}{2n} \sum_{i=1}^n \left\| \mathbf{V}^T \phi(\mathbf{W}\mathbf{x}) - \mathbf{A}\mathbf{x} \right\|^2 \quad (1)$$

using gradient descent. For part of our theoretical analysis of GD, we also consider the population loss (i.e. infinite data asymptotics as  $n \rightarrow \infty$ ) with  $\mathbf{x}$  drawn randomly from an isotropic Gaussian distribution  $\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d)$ . Concretely, the population loss is given by

$$\mathcal{L}(\theta) = \frac{1}{2} \mathbb{E}_{\mathbf{x}} \left[ \left\| \sum_{i=1}^k \mathbf{v}_i \phi(\mathbf{w}_i^T \mathbf{x}) - \mathbf{A}\mathbf{x} \right\|^2 \right]. \quad (2)$$

### 1.3 WHY IS LEARNING LINEAR FUNCTIONS WITH RELU NETWORKS CHALLENGING?



**Figure 1: Characterization of the Loss landscape.** We run gradient descent updates on the population loss. Here, the fitted model is a one-hidden-layer ReLU network with two hidden units and the outer layer is fixed to  $\mathbf{v} = [1, -1]^T$ . We use two different choices of initializations for  $\mathbf{w}_1^{(0)}$  and  $\mathbf{w}_2^{(0)}$  in parts (a) and (b). Note that in both figures black arrows indicate the  $\pm \mathbf{a}$  direction and a randomly selected orthogonal direction to  $\mathbf{a}$  is shown for the y-axis in order to visualize the neurons in 2D. In part (a) on the left, we initialize the weights small, i.e. near the origin, and observe that the weights converge to the global optimum. However, when we initialize the weights near  $1.5\mathbf{a}$  as depicted in part (b) on the right, the weights converge to  $\mathbf{a}$  and  $2\mathbf{a}$ . This corresponds to one of the non-strict saddle points of the population loss (2). To corroborate this further, in part (c) we visualize the loss landscape when  $\mathbf{w}_1 = \alpha_1 \mathbf{a}$  and  $\mathbf{w}_2 = \alpha_2 \mathbf{a}$ .

Given the simple nature of the target function it is natural to wonder about what makes GD analysis in this setting challenging. The main challenge arises from the fact that the loss landscape of the problem with ReLU NNs has many non-strict saddles—in fact infinitely many! For instance, when the output is one dimensional, if we fix the outer layer to be  $\mathbf{v} = [v_1, -v_2]^T$  with  $v_1, v_2 > 0$ , and the target function is given by  $\mathbf{x} \mapsto \mathbf{a}^T \mathbf{x}$  it is easy to see that any  $\mathbf{w}_1$  and  $\mathbf{w}_2$  of the form  $\mathbf{w}_1 = \frac{(c+1)\mathbf{a}}{v_1}$  and  $\mathbf{w}_2 = \frac{c\mathbf{a}}{v_2}$  for any  $c > 0$  or  $c < -1$  is a non-strict saddle (See Appendix C for comprehensive characterization of the loss landscape). In Figure 1 (c) we visualize this by setting  $v_1 = v_2 = 1$  and drawing the loss as a 3D plot when we restrict student neurons to be  $\mathbf{w}_1 = \alpha_1 \mathbf{a}$  and  $\mathbf{w}_2 = \alpha_2 \mathbf{a}$ . The gradient vanishes around the  $\alpha_1 - \alpha_2 = 1$  valley, but the loss is greater than 0.

As a concrete example, in Figure 1 we show that the initialization of the network directly influences whether GD will converge to the global optimum or a non-strict saddle. This experiment gives us a hint that as long as  $\|\mathbf{w}_1 + \mathbf{w}_2\|$  ( $\|v_1 \mathbf{w}_1 + v_2 \mathbf{w}_2\|$  for general  $v_1, v_2$ ) is sufficiently small (which is satisfied for small initialization), we are far away from bad stationary points. This observation will play a crucial role in our analysis.

## 2 THEORETICAL RESULTS IN THE POPULATION SETTING

We begin by stating our main result when the output is one dimensional i.e.  $r = 1$  and defer the general case to the appendix.

**Theorem 1** *Suppose the feature vectors are distributed i.i.d. according to a Gaussian distribution  $\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d)$ . We assume the corresponding output are generated according to a linear target function of the form  $\mathbf{y} = \mathbf{a}^T \mathbf{x}$  where  $\mathbf{a} \in \mathbb{R}^d$  is an arbitrary weight vector. To learn this linear function we fit a one hidden layer ReLU network with two hidden nodes of the form*

$$\mathbf{x} \mapsto \mathbf{v}^T \text{ReLU}(\mathbf{W} \mathbf{x}) = v_1 \text{ReLU}(\mathbf{w}_1^T \mathbf{x}) - v_2 \text{ReLU}(\mathbf{w}_2^T \mathbf{x}).$$

Here, we fix  $\mathbf{v} = [v_1, -v_2]^T$  with  $v_1, v_2 > 0$  and define  $\mathbf{W} = [\mathbf{w}_1, \mathbf{w}_2]^T \in \mathbb{R}^{2 \times d}$ . Consider the population loss

$$\mathcal{L}(\mathbf{W}) = \frac{1}{2} \mathbb{E}_{\mathbf{x}} \left[ (\mathbf{v}^T \text{ReLU}(\mathbf{W} \mathbf{x}) - \mathbf{a}^T \mathbf{x})^2 \right].$$

To fit this model we run gradient updates of the form  $\mathbf{W}^{(\tau+1)} = \mathbf{W}^{(\tau)} - \mu \text{diag}(\mathbf{v})^{-2} \nabla \mathcal{L}(\mathbf{W}^{(\tau)})$ , starting from an initial estimate  $\mathbf{W}^{(0)} = [\mathbf{w}_1^{(0)} \quad \mathbf{w}_2^{(0)}]^T$  with step size obeying  $\mu \leq c_1$ .

Then, as long as the initialization obeys  $\|v_1 \mathbf{w}_1^{(0)} + v_2 \mathbf{w}_2^{(0)}\| \leq \frac{1}{2} \|\mathbf{a}\|$ , we have

$$\|\mathbf{W}^{(\tau)} - \mathbf{W}^*\|_F^2 \leq \frac{\max(v_1^2, v_2^2)}{\min(v_1^2, v_2^2)} (1 - c_2 \mu)^\tau \|\mathbf{W}^{(0)} - \mathbf{W}^*\|_F^2$$

for all iterations  $\tau$ . Here,  $\mathbf{W}^* = [\frac{\mathbf{a}}{v_1}, -\frac{\mathbf{a}}{v_2}]^T$  and all  $c_j$ 's are fixed numerical constants.

This result shows that one can indeed use GD to train a one-hidden layer network with two hidden nodes to learn a linear function. We note that any linear function of the form  $\mathbf{a}^T \mathbf{x}$  can also be written as a difference of two ReLUs:  $v_1 \text{ReLU}(\frac{1}{v_1} \mathbf{a}^T \mathbf{x}) - v_2 \text{ReLU}(\frac{-1}{v_2} \mathbf{a}^T \mathbf{x})$  for any  $v_1, v_2 > 0$ , so that two hidden nodes are necessary. We indeed show directly that the GD updates result in directional convergence:  $\mathbf{w}_1^{(\infty)} = \frac{\mathbf{a}}{v_1}$  and  $\mathbf{w}_2^{(\infty)} = -\frac{\mathbf{a}}{v_2}$ .

Stated differently, with exact parametrization, GD indeed finds the underlying structure in the data (we verify this experimentally in Section 4.1). It is also worth noting that the dependence on the initialization scale is moderate and is only through ensuring that  $\|v_1 \mathbf{w}_1 + v_2 \mathbf{w}_2\|$  is not *too large* at initialization. Finally, the training is rather fast enjoying a geometric (a.k.a. linear) rate of convergence. It is worth noting that this result is based on running GD with a *scaled* step size due to the  $\text{diag}(\mathbf{v})^{-2}$  term. While this scaling is absent in standard GD updates, it significantly accelerates convergence which we demonstrate empirically in Section 4.3.

As discussed previously in Section 1.3 the optimization landscape in this problem is rather complex involving multiple non-strict saddles. Nevertheless, the above theorem demonstrates geometric convergence to the global optimum. As we explain in the proofs this is possible in part an interesting control of the trajectory of the iterates where we show that throughout the training dynamics  $\|v_1 \mathbf{w}_1^{(\tau)} + v_2 \mathbf{w}_2^{(\tau)}\|$ , continuously decrease. This facilitates a more refined control of the trajectory of GD, ensuring that GD can in fact avoid the bad stationary points.

Finally, we note that when running GD on both layers from small random initialization the outer weight will have opposite sign with roughly the same absolute value ( $v_1 \approx v_2$ ) (see Section 4). This holds if the output layer weights have opposite sign at initialization. Perhaps to be expected if the outer weights are initialized with the same sign the model gets stuck at a local optimum.

## 3 THEORETICAL RESULTS IN THE EMPIRICAL SETTING

In the previous section, we stated our main results for the population setting. Now, we focus on the more practical empirical setting when the output is one dimensional i.e.  $r = 1$ .

**Theorem 2** Suppose we have  $n$  feature vectors  $\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$  that are sampled i.i.d. according to a Gaussian distribution  $\mathbf{x}_i \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d)$ . We assume the corresponding output are generated according to a linear target function of the form  $\mathbf{y}_i = \mathbf{a}^T \mathbf{x}_i$  where  $\mathbf{a} \in \mathbb{R}^d$  is an arbitrary weight vector. To learn this linear function we fit a one hidden layer ReLU network with two hidden nodes of the form

$$\mathbf{x} \mapsto \mathbf{v}^T \text{ReLU}(\mathbf{W} \mathbf{x}) = v_1 \text{ReLU}(\mathbf{w}_1^T \mathbf{x}) - v_2 \text{ReLU}(\mathbf{w}_2^T \mathbf{x}).$$

Here, we fix the outer weights  $\mathbf{v} = [v_1, -v_2]^T$  for  $v_1, v_2 > 0$  and optimize the loss over the inner weights  $\mathbf{W} = [\mathbf{w}_1, \mathbf{w}_2]^T \in \mathbb{R}^{2 \times d}$  on the empirical loss

$$\hat{\mathcal{L}}(\mathbf{W}) = \frac{1}{2n} \sum_{i=1}^n (\mathbf{v}^T \text{ReLU}(\mathbf{W} \mathbf{x}_i) - \mathbf{a}^T \mathbf{x}_i)^2.$$

To fit this model we run gradient updates of the form  $\mathbf{W}^{(\tau+1)} = \mathbf{W}^{(\tau)} - \mu_\tau \text{diag}(\mathbf{v})^{-2} \nabla \mathcal{L}(\mathbf{W}^{(\tau)})$ , with step size obeying  $\mu_1 = 2$ ,  $\mu_\tau = \mu \leq c_5$ ,  $\forall \tau \geq 2$ . Furthermore, we assume a sufficiently small random initial estimate  $\mathbf{W}^{(0)} = [\mathbf{w}_1^{(0)} \quad \mathbf{w}_2^{(0)}]^T$  of the form  $\mathbf{w}_1^{(0)} \sim \mathcal{N}(\mathbf{0}, \sigma_1^2 \mathbf{I}_d)$  and  $\mathbf{w}_2^{(0)} \sim \mathcal{N}(\mathbf{0}, \sigma_2^2 \mathbf{I}_d)$  with the standard deviations obeying  $\sqrt{v_1^2 \sigma_1^2 + v_2^2 \sigma_2^2} \sqrt{d} \leq c_6 \|\mathbf{a}\|$ . Then as long as the number of training samples satisfy  $n \geq c_7 d$ , we have

$$\left\| \mathbf{W}^{(\tau)} - \mathbf{W}^* \right\|_F^2 \leq \frac{\max(v_1^2, v_2^2)}{\min(v_1^2, v_2^2)} (1 - c_8 \mu)^\tau \left\| \mathbf{W}^{(0)} - \mathbf{W}^* \right\|_F^2$$

with probability at least  $1 - Ce^{-cd}$ . Here,  $\mathbf{W}^* = [\frac{\mathbf{a}}{v_1}, -\frac{\mathbf{a}}{v_2}]^T$  and all  $c_j$ 's are fixed numerical constants independent of any problem dimensions.

This result demonstrates that gradient descent can indeed be used to train a one-hidden-layer network with two hidden nodes on a finite number of training samples to learn a linear function. Similar to the population setting, we show that student neurons  $\mathbf{w}_1$  and  $\mathbf{w}_2$  directionally converge in direction to the ground truth with high probability. Notably, the convergence is fast, exhibiting a geometric rate. Moreover, the sample complexity is optimal, scaling linearly with the problem dimension  $d$ . Finally, we note that the theorem above allows sample reuse across all iterations requiring the development and use of intricate uniform concentration inequalities in the proof.

We also highlight that the initialization radius required in the empirical setting, given by  $\sqrt{v_1^2 \sigma_1^2 + v_2^2 \sigma_2^2} \sqrt{d} \leq c_6 \|\mathbf{a}\|$ , directly parallels the condition in the population setting, where  $\|v_1 \mathbf{w}_1^{(0)} + v_2 \mathbf{w}_2^{(0)}\| \leq \frac{1}{2} \|\mathbf{a}\|$ . However, our analysis for the empirical case relies on random initialization, whereas in the population setting, any initialization in the specified ball is sufficient.

## 4 EXPERIMENTS

In this section we show experimental results in various output dimension ( $r = 1$  vs.  $r > 1$ ), and initialization scale (small vs large). We use PyTorch for experiments and unless mentioned otherwise, network weights are initialized with Xavier Normal initialization (for a matrix  $\mathbf{W} \in \mathbb{R}^{r \times d}$ ,  $\mathbf{W}_{ij} \sim \mathcal{N}\left(0, \frac{2}{r+d}\right)$ ). For visualization purposes, we set  $v_1 = v_2 = 1$  (in fact the proof for arbitrary  $v_1, v_2$  setting can be reduced to  $v_1 = v_2 = 1$  which we discuss in detail in Appendix D). In order to change the initialization scale, we multiply the default initialization with a positive scalar  $\alpha$ . For small initialization experiments, we use  $\alpha = 10^{-8}$ , otherwise it is set to  $\alpha = 1$ . We set  $d = 100$  and  $\mu = 0.1$ . All experiments are run on a server with an Intel Xeon Gold 5220R CPU. We would like to stress that even though the visualizations in this paper are based on a single trial, we ran these experiments for different random seeds and the behavior of the visualizations did not change.

### 4.1 LEARNING LINEAR TARGETS (SINGLE OUTPUT: $r = 1$ )

In experiments w.l.o.g. we choose  $\mathbf{a} = \mathbf{e}_1$  where  $\mathbf{e}_1$  is the first standard basis in  $\mathbb{R}^d$ . This does not effect the results due to the rotational symmetry of isotropic Gaussian distribution of which  $\mathbf{x}$  are drawn from. Note that this implies  $\|\mathbf{a}\| = 1$  in our experiments. Finally, in this section we focus

our experiments on the population loss. Similar results continue to hold in the empirical case with moderate sample sizes i.e. when  $n \geq crd$  with  $c$  a sufficiently large constant.

When the model is exactly parameterized with two hidden nodes ( $k = 2$ ), and outer layer is fixed as  $\pm 1$  we show that the population loss decreases at a linear rate (Figure 5) and weights of the inner layer align themselves with  $\pm a$  (Figure 2 (a)). If  $v$  is also initialized randomly, we empirically see that the model cannot converge to the global optima consistently. When it does,  $v^{(\infty)}$  indeed becomes  $\pm 1$  and  $w_1^{(\infty)}$  and  $w_2^{(\infty)}$  recover  $\pm a$  exactly. For the remaining time, the GD iterates converge to one of the many local optima of this problem similar to the depiction in Figure 1 (part b). We further observe that iterates get stuck only when  $v_1^{(0)}$  and  $v_2^{(0)}$  both have the same signs which happens with probability  $\frac{1}{2}$ .

When  $k > 2$ , the probability of *all*  $v_i$ 's having the same sign decreases rapidly. Therefore, iterates typically converge to the global optima. However, in this case global minima is not unique anymore. To demonstrate this, consider the case where there are four hidden units ( $k = 4$ ) instead of two. We fix half of  $v$  as  $+1$  and remaining half as  $-1$ . The trajectory of the inner weights across GD iterations is depicted in Figure 2.

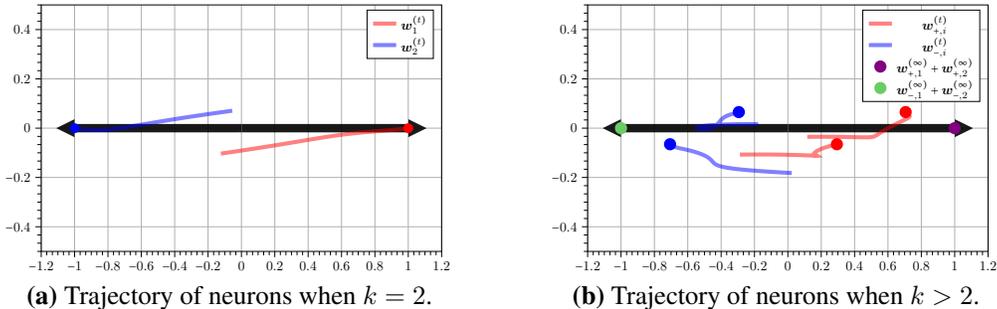


Figure 2: **Trajectory of neurons for different values of  $k$ .** We run gradient descent updates on the population loss after fixing  $v$  with half  $+1$ 's and half  $-1$ 's. A randomly selected orthogonal direction to  $a$  is shown for the y-axis in order to visualize the neurons in 2D. Black arrows indicate  $\pm a$  direction. We use colors red and blue to indicate whether  $v_i$  corresponding to  $w_i$  is  $1$  or  $-1$  respectively. Points at the end of each trajectory denotes the final weight GD converges to.

We observe that while no individual  $w_i$  align itself with  $\pm a$  direction, grouping hidden units based on their corresponding signs in  $v$  and summing them recovers  $\pm a$  exactly (purple and green points in Figure 2). Although not depicted here, we have tried various values for  $k > 2$  and the observation that grouping weights recover  $\pm a$  was consistent. This suggests that combining node aggregation technique from (Li et al., 2024) with our proof strategy may extend our results for the  $k > 2$  setting. We leave this to future work.

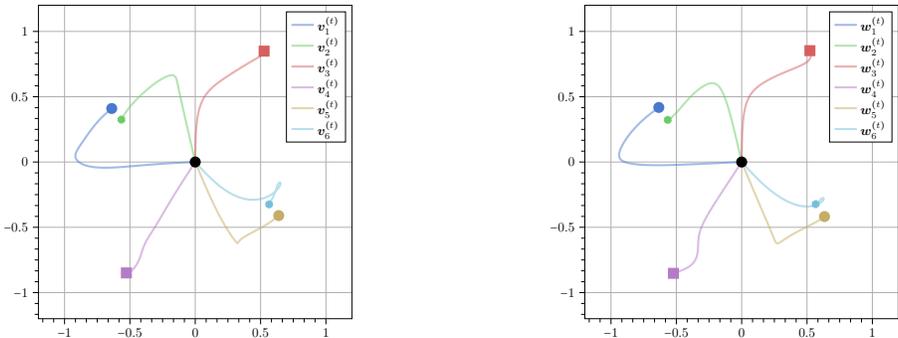
#### 4.2 LEARNING LINEAR TARGETS (MULTI-DIMENSIONAL OUTPUT: $r > 1$ )

We only consider the case where the model is exactly parameterized i.e.  $k = 2r$ . We first show that an interesting pattern arises if both the inner and outer layers of the neural network are initialized sufficiently small. For visualization purposes in Figure 3, we pick  $r = 3$  and  $k = 6$  (see Appendix G.1 for additional figures with  $r > 3$ ). As for the target function, we pick  $a_1, a_2, a_3$  to be  $e_1, e_2, e_3$  respectively which correspond to the standard basis vectors in  $\mathbb{R}^d$ . We plot the trajectory of both the inner and outer layer weights of the network across iterations and observe a peculiar pattern in both  $v_i$ 's and  $w_i$ 's. At convergence, weights can be grouped into pairs such that one of the weights is approximately negative of the other. As a concrete example, in Figure 3, we observe that  $v_3^{(\infty)} \approx -v_4^{(\infty)}$ ,  $v_1^{(\infty)} \approx -v_5^{(\infty)}$ , and  $v_2^{(\infty)} \approx -v_6^{(\infty)}$  which also holds similarly for  $w_i$ 's as well. This suggests that after a permutation of the hidden units, we get

$$V^{(\infty)} \approx [I_r, -I_r]^T \tilde{V}, \quad W^{(\infty)} \approx [I_r, -I_r]^T \tilde{W}.$$

This inspires us to fix  $V$  according to the pattern above for our theoretical results (see Theorem 6 in the Appendix) as a natural extension to the  $v_i = \pm 1$  pattern in the single output setting.

324  
325  
326  
327  
328  
329  
330  
331  
332  
333  
334



(a) Trajectory of  $v_i$ 's and their pairing behavior. (b) Trajectory of  $w_i$ 's and their pairing behavior.

335  
336  
337  
338  
339  
340  
341  
342  
343  
344

**Figure 3: Pairing pattern in multi-dimensional setting.** We train the network from small initialization when exactly parameterized ( $k = 6$  and  $r = 3$ ). On left (a), we depict the trajectories of individual weights in the outer layer ( $v_i$ 's) across iterations. We observe that the weights at convergence can be grouped into three pairs such that one of the weights is approximately negative of the other. For instance, we observe that  $v_3^{(\infty)} \approx -v_4^{(\infty)}$ . Which neurons end up pairing with each other is indicated by the usage of same symbol (square, circle, etc.). A similar pairing is observed for the inner layer weights as well (b). While these vectors all lie in a higher dimensional space, we pick an arbitrary two dimensional axis to plot them in 2D.

345

### 4.3 BENEFITS OF SCALING THE LEARNING RATE

346  
347  
348  
349  
350

In this section we illustrate why a scaled learning rate (i.e.  $\mu_i = \frac{\mu}{v_i}$ ), which appears in Theorems 1 and 2, is a natural choice. While this scaling is absent in standard GD updates, we show that it significantly accelerates convergence in the fixed- $v$  setting.

351  
352  
353  
354  
355

Concretely, in this experiment we focus on a single output setting ( $r = 1$ ) where  $v_i$ 's are fixed but with a big difference in their magnitudes. In particular, we set  $v_1 = 100$  and  $v_2 = 1$ . We initialize  $\mathbf{W}$  moderately small with  $\alpha = 0.01$  and compare the convergence speed of the loss in the population setting with and without the scaled step size. We separately find the best  $\mu$  (via binary search) that achieves the fastest convergence. The loss across iterations is depicted in Figure 4.

356  
357  
358  
359  
360  
361  
362

As evident in this figure, the difference in convergence speed is rather significant. Without our scaling, about  $10^4$  iterations are needed to bring the loss below  $10^{-5}$ , while with normalization it takes only **3** iterations. Our proposed rescaling, thus enables the use of a large  $\mu$  throughout training.

363

## 5 RELATED WORK

364  
365  
366  
367  
368  
369  
370

There is a large body of work on developing global convergence guarantees for nonconvex problems. We briefly review this literature and compare the differences with the setting discussed in this paper.

371  
372  
373  
374  
375  
376  
377

**Nonconvex low-rank matrix recovery:** In low-rank matrix recovery, numerous studies have shown that nonconvex gradient descent, when initiated with spectral initialization, can effectively solve low-rank reconstruction problems across various domains. This includes phase retrieval (Candès et al., 2015; Chen and Candès, 2017; Ma et al., 2020), matrix sensing Tu et al. (2016), blind deconvolution Li et al. (2019); Ling and Strohmer (2019), and matrix completion Chen et al. (2020). In practice, random initialization is frequently employed instead of specialized spectral initialization methods. As a result, more recent literature Sun et al. (2018); Ge et al. (2016); Zhang et al. (2019), have turned to analyzing the loss landscape. These studies demonstrate that, despite their non-convex nature,

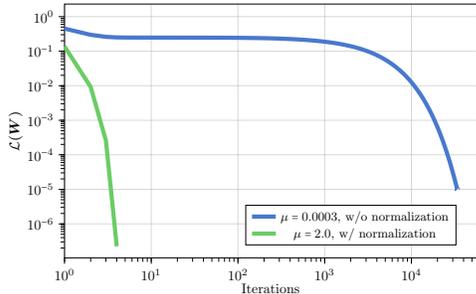


Figure 4: **Effect of scaled step-size.** The population loss decreases significantly faster with a scaled step size compared to its fixed step size counterpart. Step size  $\mu$  is tuned separately in both cases.

these loss landscapes remain well-behaved under certain assumptions. Specifically, they contain no spurious local minima (i.e., all minimizers are global minima), and saddle points exhibit a strict direction of negative curvature (also known as strict saddle points) Sun et al. (2015). Then specialized truncation or saddle escaping algorithms such as trust region, cubic regularization Nesterov and Polyak (2006); Nocedal and Wright (2006) or noisy (stochastic) gradient-based methods Jin et al. (2017a); Ge et al. (2015); Raginsky et al. (2017); Zhang et al. (2017) are deployed to provably find a global optimum. In contrast to the above literature, the landscape of our loss contain non-strict saddle points. Furthermore, we do not seek any modification to the initialization or the GD updates. Indeed, our result holds with moderately small initialization. As mentioned earlier, we are able to establish this result by developing intricate control of the GD updates throughout the trajectory.

**Gradient-based analysis for neural networks:** A recent line of work is concerned with connecting the analysis of neural network training with the so-called neural tangent kernel (NTK) Jacot et al. (2018); Oymak and Soltanolkotabi (2019; 2020); Du et al. (2019); Arora et al. (2019). The core idea is that with sufficiently large initialization, a neural network can be approximated by its linearization around the origin. This approximation facilitates linking neural network analysis to the well-established theory of kernel methods. This approach is sometimes referred to as lazy training since, under such initialization, the network parameters remain close to their initial values throughout training. However, some research suggests that NTK-based analysis alone may not fully account for the practical success of neural networks. For instance, Chizat et al. (2019) presents empirical evidence indicating that reducing the initialization size can lead to lower test error. Similarly, Ghorbani et al. (2020) observes a performance gap between neural networks and their NTK counterparts, with the gap widening when the covariance matrix is isotropic. We note that in an NTK analysis the parameters stay close to the initialization which is not the case in our setting. Furthermore, an NTK analysis that relies on linearization can not deal with trajectory analysis that avoids local optima. Indeed, an NTK analysis will not yield the directional convergence established in this paper. So in this sense our result can be viewed as going beyond the lazy training in NTK theory. We demonstrate the lack of directional convergence and lack of generalization in the NTK regime empirically in Appendix G.2.

#### Beyond NTK and learning of specific target functions.

Recent work carries out analysis of neural networks beyond NTK regime including Damian et al. (2022); Ba et al. (2022); Lee et al. (2024); Xu and Du (2023). Many of these results also focus on learning specific target functions such as ReLUs (Xu and Du, 2023), (Soltanolkotabi, 2017) and polynomials (Damian et al., 2022). These results however typically exclude linear function classes and do not directly involve analysis that requires avoiding bad stationary points explicitly. In fact, many of the existing papers use a pre-processing step or alter the early optimization trajectory to avoid complications arising from the dynamics of learning linear functions (Damian et al., 2022). In contrast, our focus is directly dealing with such intricacies.

Among these papers, perhaps the closest to ours in spirit is (Xu and Du, 2023) which studies the problem of fitting an overparameterized ReLU network to a single ReLU target function with a one dimensional output. Our one-dimensional result can be viewed as a generalization of this work (in particular their exact parametrization result) where the target function has two ReLUs with a particular pattern. This is due to the fact that any linear function of the form  $\mathbf{a}^T \mathbf{x}$  can also be written as a difference of two ReLUs:  $v_1 \text{ReLU}\left(\frac{1}{v_1} \mathbf{a}^T \mathbf{x}\right) - v_2 \text{ReLU}\left(\frac{-1}{v_2} \mathbf{a}^T \mathbf{x}\right)$  for any  $v_1, v_2 > 0$ . The addition of this new ReLU with a negative sign introduces non-strict saddle points and various intricacies in the landscape necessitating a completely different analysis. However, compared to (Xu and Du, 2023) we do not study the effect of overparameterization theoretically. Our empirical results in Section 4.1 suggest that such an extension may be possible.

We highlight that besides Xu and Du (2023), there are several other works on learning a single neuron Yehudai and Shamir (2022); Vardi et al. (2022); Chistikov et al. (2023) and variants Brutzkus and

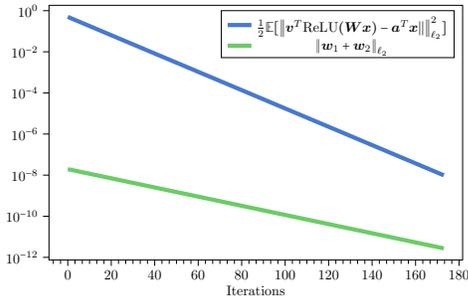


Figure 5:  $\|w_1 + w_2\|$  is monotonically decreasing. Both the population loss  $\mathcal{L}(\theta)$  and  $\|w_1 + w_2\|$  go to zero at a linear rate.

432 Globerson (2017). As explained before, such results cannot be used to analyze linear targets due to  
 433 the interaction terms between positive and negative ReLU neurons. Furthermore, we note that the  
 434 landscape for fitting a single ReLU is fundamentally different as it contains only a single basin of  
 435 attraction (albeit a non-convex one). In contrast, as discussed earlier the landscape in our problem  
 436 include non-strict saddle points significantly complicating gradient descent analysis.

437 We would also like to discuss the difference between our work and a few other papers Zhong et al.  
 438 (2017); Zhang et al. (2018); Zhu et al. (2025) that have planted one-hidden layer models. These  
 439 papers differ in at least one of three ways focusing on (1) local analysis, (2) have sub-optimal sample  
 440 complexity, and/or (3) assume non-negative outer layer weights. For instance, Zhong et al. (2017)  
 441 utilize tensor initialization, performing a local analysis rather than a global GD analysis. This local  
 442 analysis however can not be used to analyze the linear target setting. Indeed, as noted in Remark 4.3  
 443 of their work, their analysis requires  $W^*$  to be full-rank which does not hold in the linear setting  
 444 (where the rows of the weight matrix are negatives of each other leading to a minimum singular value  
 445 is zero). Furthermore, this result also requires resampling the data points at each iteration to ensure  
 446 convergence of gradient descent where as we use the same samples across all iterations. On a related  
 447 note, their sample complexity has polynomial dependency on many problem parameters (Theorem  
 448 4.2) whereas our proof only requires sample size linear in input dimension  $d$ .

449 Similarly, Zhang et al. (2018) provide a local analysis of GD when the outer layer weights are fixed  
 450 to be all ones. They also utilize results of Zhong et al. (2017) and share similar limitations in terms of  
 451 the rank requirement on  $W^*$ . Thus this result can not be used in the linear target setting even for a  
 452 local analysis. While they improve the sample complexity of (Zhang et al., 2018) by getting rid of  
 453 the resampling trick, they still end up with a sample complexity polynomial in width of the network.

454 Wu et al. (2018) consider the setting when student and teacher networks both have 2 neurons. In  
 455 particular, when the teachers are *orthogonal*, and the outer weights are all ones; they demonstrated  
 456 an interesting result that the landscape is benign and all saddles are strict. In contrast, the landscape  
 457 in our problem include non-strict saddle points significantly complicating gradient descent analysis.  
 458 More recently, Zhu et al. (2025) also consider learning multiple *orthogonal* ReLU neurons in a  
 459 teacher-student framework with outer layer weights fixed to all ones. As just discussed, having  
 460 orthogonal teacher weights leads to a much more benign landscape. Moreover, assumptions in the  
 461 aforementioned works strictly exclude the linear target setting, where the outer layer must contain  
 462 negative coefficients. Furthermore, they impose strong restrictions on the initialization. Specifically,  
 463 they look at the convergence after “weak alignment” where for each student neuron there exists  
 464 only one teacher neuron that is not near perpendicular. Our population results on the other hand can  
 465 handle initializations where both student neurons are perpendicular to the target direction as long  
 466 as we have  $\|w_1 + w_2\| \leq \frac{1}{2} \|\alpha\|$  (e.g.  $w_1 = w_2 = \frac{\|\alpha\|}{4} \bar{\alpha}_\perp$ ). That said, their analysis can handle  
 467 over-parametrization ( $k \gg k^*$ ) and teacher networks with more than 2 neurons.

468 In recent and independent work, Boursier and Flammarion (2025) also consider the problem of learn-  
 469 ing linear target functions. The authors demonstrate an interesting result: despite over-parametrization,  
 470 the sum of positive (resp. negative) neurons aligns with the OLS estimator obtained from the “positive”  
 471 (resp. negative) subset of the data. To prove this, the authors impose heavy restrictions on the data  
 472 distribution (in particular, Conditions 3 and 4 in their paper) to essentially align the data with the  
 473 target direction and avoid changes in the activation cone. We quote the authors:

474 “However, item 3 is quite restrictive: it is needed to ensure that the volume of the activation cone  
 475 containing  $\beta^*$  does not vanish when  $n \rightarrow \infty$ . A similar assumption is considered by Chistikov et al.  
 476 (2023); Tsoy and Konstantinov (2024), for similar reasons. Additionally, Condition 4 ensures that  
 477  $\mathbb{E}_x [xx^T] \beta^*$  and  $\beta^*$  are in the same activation cone. This assumption allows the training dynamics to  
 478 remain within a single cone after the early alignment phase, significantly simplifying our analysis.”

479 In contrast, we demonstrate feature learning in the linear target setting by performing a full char-  
 480 acterization of GD dynamics with a generic data distribution and initialization without any of the  
 481 restrictive assumptions mentioned above.

## 482 6 OVERVIEW AND KEY IDEAS OF THE PROOF

483 In this section we provide an overview of our proof strategy. In the entire proof we focus on the  
 484  $v_1 = v_2 = 1$  setting. As we demonstrate in Appendix D this is without any loss in generality.

**Key Proof Ideas in the Population Setting.**

Our proof technique relies on showing the following two inequalities:

1. **Correlation inequality:**  $\langle \mathbf{W} - \mathbf{W}^*, \nabla_{\mathbf{W}} \mathcal{L}(\mathbf{W}) \rangle \geq \alpha \|\mathbf{W} - \mathbf{W}^*\|_F^2$

2. **Gradient smoothness (towards the global optima):**  $\|\nabla_{\mathbf{W}} \mathcal{L}(\mathbf{W})\|_F \leq \beta \|\mathbf{W} - \mathbf{W}^*\|_F$

These two inequalities combined imply that the GD iterates converge to a global optima at a linear rate. However, as stated earlier in Section 1.3 the loss landscape has many non-strict saddle points. This immediately implies that the first correlation inequality can not hold over the entire domain. To circumvent this, we will utilize the following key observation.

**Key Observation:** A critical idea towards proving the correlation inequality is showing that along the GD trajectory  $\|\mathbf{w}_1 + \mathbf{w}_2\|$  decreases as also depicted in Figure 5. This is formalized below.

**Lemma 3** For all iterations  $\tau$ , we have  $\|\mathbf{w}_1^{(\tau+1)} + \mathbf{w}_2^{(\tau+1)}\| \leq \|\mathbf{w}_1^{(\tau)} + \mathbf{w}_2^{(\tau)}\|$  when  $\mu \leq 1$ .

With this key observation in place, we turn our attention to the proof which follows these steps:

- Step 1 (Correlation inequality): We define

$$h(\mathbf{w}_1, \mathbf{w}_2, \mathbf{a}) = \langle \mathbf{W} - \mathbf{W}^*, \nabla_{\mathbf{W}} \mathcal{L}(\mathbf{W}) \rangle - \alpha \|\mathbf{W} - \mathbf{W}^*\|_F^2, \quad (3)$$

and compute  $\tilde{h}(\mathbf{w}_1, \mathbf{w}_2) = \min_{\mathbf{a}} h(\mathbf{w}_1, \mathbf{w}_2, \mathbf{a})$  assuming that  $\|\mathbf{w}_1 + \mathbf{w}_2\| \leq \frac{1}{2} \|\mathbf{a}\|$ . We show that

$\frac{1}{\|\mathbf{w}_1\|^2} \tilde{h}(\mathbf{w}_1, \mathbf{w}_2)$  is only a function of  $\theta$  (angle between  $\mathbf{w}_1$  and  $\mathbf{w}_2$ ) and  $\frac{\|\mathbf{w}_2\|}{\|\mathbf{w}_1\|}$ . This allows us to draw  $\frac{1}{\|\mathbf{w}_1\|^2} \tilde{h}$  in 2D and prove that  $\tilde{h}$  is indeed non-negative. We provide the visualization in the Appendix. This leads us to the following lemma:

**Lemma 4**  $\langle \mathbf{W} - \mathbf{W}^*, \nabla_{\mathbf{W}} \mathcal{L}(\mathbf{W}) \rangle \geq \alpha \|\mathbf{W} - \mathbf{W}^*\|_F^2$  holds with  $\alpha = 0.3$  as long as  $\|\mathbf{w}_1 + \mathbf{w}_2\| \leq \frac{1}{2} \|\mathbf{a}\|$ .

- Step 2 (Gradient smoothness towards the global optima): This step is more straightforward and only requires algebraic manipulations. Concretely we have the following lemma.

**Lemma 5**  $\|\nabla_{\mathbf{W}} \mathcal{L}(\mathbf{W})\|_F \leq \beta \|\mathbf{W} - \mathbf{W}^*\|_F$  holds for all  $\mathbf{W} \in \mathbb{R}^{2 \times d}$  with  $\beta$  a fixed constant.

**Key Proof Ideas in the Empirical Setting.** To prove our results for the empirical case, we need only establish the empirical counterparts of Lemma 4 and Lemma 5. Here, we will focus on outlining the key ideas behind the proof for the counterpart to Lemma 4 which is more involved. A first idea may be to try to show that the empirical correlation concentrates around the population correlation. However, since we reuse samples across iterations and the correlation inequality involves complex and heavy tail nonlinear functions of the data points, this can be quite challenging. To circumvent this, our key idea is to observe that the empirical correlation can be decomposed into two terms

$$\langle \mathbf{W} - \mathbf{W}^*, \nabla_{\mathbf{W}} \hat{\mathcal{L}}(\mathbf{W}) \rangle = \frac{1}{n} \sum_{i=1}^n r^2(\mathbf{x}_i) + \frac{1}{n} \sum_{i=1}^n r(\mathbf{x}_i) (\mathbf{a}^T \mathbf{x}_i) (1 - \phi'(\mathbf{w}_1^T \mathbf{x}_i) - \phi'(\mathbf{w}_2^T \mathbf{x}_i))$$

where  $r(\mathbf{x}) = \mathbf{v}^T \text{ReLU}(\mathbf{W}\mathbf{x}) - \mathbf{a}^T \mathbf{x}$  is the residual function. Critically we can show that the second term is dominated by the first term allowing us to conclude that

$$\langle \mathbf{W} - \mathbf{W}^*, \nabla_{\mathbf{W}} \hat{\mathcal{L}}(\mathbf{W}) \rangle \geq c \left( \frac{1}{n} \sum_{i=1}^n r^2(\mathbf{x}_i) \right)$$

holds for a fixed numerical constant (see appendix for detail). Thus to show the empirical version of the correlation inequality with parameter  $\alpha$ , it suffices to show that

$$\frac{1}{n} \sum_{i=1}^n r^2(\mathbf{x}_i) \geq \frac{\alpha}{c} \|\mathbf{W} - \mathbf{W}^*\|_F^2, \quad (4)$$

holds with high probability. To prove this we will first show that  $\mathbb{E}_{\mathbf{x}} [r^2(\mathbf{x})] \geq \frac{\alpha}{c(1-\delta)} \|\mathbf{W} - \mathbf{W}^*\|_F^2$ .

Next, we show that as long as  $n \geq c \frac{d}{\delta^2}$ , then

$$\frac{1}{n} \sum_{i=1}^n r^2(\mathbf{x}_i) \geq (1 - \delta) \mathbb{E}_{\mathbf{x}} [r^2(\mathbf{x})],$$

holds with high probability. The combination of the latter two results immediately implies (4) finishing the proof. This reduction may seem naive as the sum of the residuals above also involve uniform concentration of heavy-tail stochastic processes. However, critically, the summands are now positive allowing us to utilize powerful one-sided uniform concentration results that hold despite their heavy-tail nature.

## REFERENCES

- 540  
541  
542 Sanjeev Arora, Simon S Du, Wei Hu, Zhiyuan Li, and Ruosong Wang. Fine-grained analysis  
543 of optimization and generalization for overparameterized two-layer neural networks. In *36th*  
544 *International Conference on Machine Learning, ICML 2019*, pages 477–502. International Machine  
545 Learning Society (IMLS), 2019.
- 546 Jimmy Ba, Murat A Erdogdu, Taiji Suzuki, Zhichao Wang, Denny Wu, and Greg Yang. High-  
547 dimensional asymptotics of feature learning: How one gradient step improves the representation.  
548 *Advances in Neural Information Processing Systems*, 35:37932–37946, 2022.
- 549 Pierre Baldi and Kurt Hornik. Neural networks and principal component analysis: Learning  
550 from examples without local minima. *Neural Networks*, 2(1):53–58, 1989. ISSN 0893-6080.  
551 doi: [https://doi.org/10.1016/0893-6080\(89\)90014-2](https://doi.org/10.1016/0893-6080(89)90014-2). URL <https://www.sciencedirect.com/science/article/pii/0893608089900142>.
- 552 Etienne Boursier and Nicolas Flammarion. Simplicity bias and optimization threshold in two-layer  
553 relu networks, 2025. URL <https://arxiv.org/abs/2410.02348>.
- 554 Alon Brutzkus and Amir Globerson. Globally optimal gradient descent for a convnet with gaussian  
555 inputs, 2017. URL <https://arxiv.org/abs/1702.07966>.
- 556 Emmanuel J. Candès, Xiaodong Li, and Mahdi Soltanolkotabi. Phase retrieval via Wirtinger flow:  
557 theory and algorithms. *IEEE Trans. Inf. Theory*, 61(4):1985–2007, 2015. ISSN 0018-9448.
- 558 Ji Chen, Dekai Liu, and Xiaodong Li. Nonconvex rectangular matrix completion via gradient descent  
559 without  $\ell_2, \infty$  regularization. *IEEE Trans. Inf. Theory*, 66(9):5806–5841, 2020. ISSN 0018-9448.
- 560 Yuxin Chen and Emmanuel J. Candès. Solving random quadratic systems of equations is nearly  
561 as easy as solving linear systems. *Commun. Pure Appl. Math.*, 70(5):822–883, 2017. ISSN  
562 0010-3640; 1097-0312/e.
- 563 Dmitry Chistikov, Matthias Englert, and Ranko Lazic. Learning a neuron by a shallow relu network:  
564 Dynamics and implicit bias for correlated inputs, 2023. URL <https://arxiv.org/abs/2306.06479>.
- 565 Lenaïc Chizat, Edouard Oyallon, and Francis Bach. On lazy training in differentiable programming.  
566 *Advances in neural information processing systems*, 32, 2019.
- 567 Alexandru Damian, Jason Lee, and Mahdi Soltanolkotabi. Neural networks can learn representations  
568 with gradient descent. In *Conference on Learning Theory*, pages 5413–5452. PMLR, 2022.
- 569 Amit Daniely, Roy Frostig, and Yoram Singer. Toward deeper understanding of neural networks: The  
570 power of initialization and a dual view on expressivity. *Advances in neural information processing*  
571 *systems*, 29, 2016.
- 572 Simon Du, Jason Lee, Haochuan Li, Liwei Wang, and Xiyu Zhai. Gradient descent finds global  
573 minima of deep neural networks. In *International Conference on Machine Learning*, pages  
574 1675–1685. PMLR, 2019.
- 575 Simon S Du, Xiyu Zhai, Barnabas Poczos, and Aarti Singh. Gradient descent provably optimizes  
576 over-parameterized neural networks. *arXiv preprint arXiv:1810.02054*, 2018.
- 577 Zalan Fabian, Berk Tinaz, and Mahdi Soltanolkotabi. Humus-net: Hybrid unrolled multi-scale  
578 network architecture for accelerated mri reconstruction. *Advances in Neural Information Processing*  
579 *Systems*, 35:25306–25319, 2022.
- 580 Rong Ge, Furong Huang, Chi Jin, and Yang. Yuan. Escaping from saddle points: online stochastic  
581 gradient for tensor decomposition. In *Proceedings of The 28th Conference on Learning Theory*,  
582 pages 797–842, 2015.
- 583 Rong Ge, Jason D Lee, and Tengyu Ma. Matrix completion has no spurious local minimum. *Advances*  
584 *in Neural Information Processing Systems*, 29:2973–2981, 2016.

- 594 Behrooz Ghorbani, Song Mei, Theodor Misiakiewicz, and Andrea Montanari. When do neural  
595 networks outperform kernel methods? *arXiv preprint arXiv:2006.13409*, 2020.
- 596 Paul Hand, Oscar Leong, and Vlad Voroninski. Phase retrieval under a generative prior. *Advances in*  
597 *Neural Information Processing Systems*, 31, 2018.
- 599 Arthur Jacot, Franck Gabriel, and Clément Hongler. Neural tangent kernel: Convergence and  
600 generalization in neural networks. *Advances in neural information processing systems*, 31, 2018.
- 601 Chi Jin, Rong Ge, Praneeth Netrapalli, Sham M. Kakade, and Michael I. Jordan. How to escape  
602 saddle points efficiently. page 1724–1732, 2017a.
- 604 Kyong Hwan Jin, Michael T McCann, Emmanuel Froustey, and Michael Unser. Deep convolutional  
605 neural network for inverse problems in imaging. *IEEE transactions on image processing*, 26(9):  
606 4509–4522, 2017b.
- 607 Christian Ledig, Lucas Theis, Ferenc Huszár, Jose Caballero, Andrew Cunningham, Alejandro Acosta,  
608 Andrew Aitken, Alykhan Tejani, Johannes Totz, Zehan Wang, et al. Photo-realistic single image  
609 super-resolution using a generative adversarial network. In *Proceedings of the IEEE conference on*  
610 *computer vision and pattern recognition*, pages 4681–4690, 2017.
- 612 Jason D Lee, Kazusato Oko, Taiji Suzuki, and Denny Wu. Neural network learns low-dimensional  
613 polynomials with sgd near the information-theoretic limit. *arXiv preprint arXiv:2406.01581*, 2024.
- 614 Binghui Li, Zhixuan Pan, Kaifeng Lyu, and Jian Li. Feature averaging: An implicit bias of gradient  
615 descent leading to non-robustness in neural networks. *arXiv preprint arXiv:2410.10322*, 2024.
- 616 Xiaodong Li, Shuyang Ling, Thomas Strohmer, and Ke Wei. Rapid, robust, and reliable blind  
617 deconvolution via nonconvex optimization. *Appl. Comput. Harmon. Anal.*, 47(3):893–934, 2019.  
618 ISSN 1063-5203.
- 620 Shuyang Ling and Thomas Strohmer. Regularized gradient descent: a non-convex recipe for fast joint  
621 blind deconvolution and demixing. *Inf. Inference*, 8(1):1–49, 2019. ISSN 2049-8764; 2049-8772/e.
- 622 Cong Ma, Kaizheng Wang, Yuejie Chi, and Yuxin Chen. Implicit regularization in nonconvex  
623 statistical estimation: gradient descent converges linearly for phase retrieval, matrix completion,  
624 and blind deconvolution. *Found. Comput. Math.*, 20(3):451–632, 2020. ISSN 1615-3375; 1615-  
625 3383/e.
- 626 Yurii Nesterov and Boris T. Polyak. Cubic regularization of Newton method and its global perfor-  
627 mance. *Math. Program.*, 108(1 (A)):177–205, 2006. ISSN 0025-5610; 1436-4646/e.
- 629 Jorge Nocedal and Stephen J. Wright. Trust-region methods. *Numerical Optimization*, pages 66–100,  
630 2006.
- 631 Greg Ongie, Rebecca Willett, Daniel Soudry, and Nathan Srebro. A function space view of bounded  
632 norm infinite width relu nets: The multivariate case. *arXiv preprint arXiv:1910.01635*, 2019.
- 634 Samet Oymak and Mahdi Soltanolkotabi. Overparameterized nonlinear learning: Gradient descent  
635 takes the shortest path? In *International Conference on Machine Learning*, pages 4951–4960.  
636 PMLR, 2019.
- 637 Samet Oymak and Mahdi Soltanolkotabi. Towards moderate overparameterization: global con-  
638 vergence guarantees for training shallow neural networks. *IEEE Journal on Selected Areas in*  
639 *Information Theory*, 2020.
- 641 Maxim Raginsky, Alexander Rakhlin, and Matus Telgarsky. Non-convex learning via stochastic  
642 gradient langevin dynamics: a nonasymptotic analysis. pages 1674–1703, 2017.
- 643 Mahdi Soltanolkotabi. Learning relus via gradient descent. *Advances in neural information processing*  
644 *systems*, 30, 2017.
- 646 Mahdi Soltanolkotabi. Structured signal recovery from quadratic measurements: Breaking sample  
647 complexity barriers via nonconvex optimization. *IEEE Transactions on Information Theory*, 65(4):  
2374–2400, 2019.

- 648 Mahdi Soltanolkotabi, Adel Javanmard, and Jason D Lee. Theoretical insights into the optimization  
649 landscape of over-parameterized shallow neural networks. *IEEE Transactions on Information*  
650 *Theory*, 65(2):742–769, 2018.
- 651 Anuroop Sriram, Jure Zbontar, Tullie Murrell, Aaron Defazio, C Lawrence Zitnick, Nafissa Yakubova,  
652 Florian Knoll, and Patricia Johnson. End-to-end variational networks for accelerated mri recon-  
653 struction. In *Medical Image Computing and Computer Assisted Intervention–MICCAI 2020: 23rd*  
654 *International Conference, Lima, Peru, October 4–8, 2020, Proceedings, Part II 23*, pages 64–73.  
655 Springer, 2020.
- 656 Ju Sun, Qing Qu, and John Wright. When are nonconvex problems not scary? *arXiv preprint*  
657 *arXiv:1510.06096*, 2015.
- 658 Ju Sun, Qing Qu, and John Wright. A geometric analysis of phase retrieval. *Found. Comput. Math.*,  
659 18(5):1131–1198, 2018. ISSN 1615-3375; 1615-3383/e.
- 660 Stephen Tu, Ross Boczar, Max Simchowitz, Mahdi Soltanolkotabi, and Ben Recht. Low-rank  
661 solutions of linear matrix equations via procrustes flow. In *International Conference on Machine*  
662 *Learning*, pages 964–973. PMLR, 2016.
- 663 Gal Vardi, Gilad Yehudai, and Ohad Shamir. Learning a single neuron with bias using gradient  
664 descent, 2022. URL <https://arxiv.org/abs/2106.01101>.
- 665 Roman Vershynin. High-dimensional probability.
- 666 Yi Wang, Xin Tao, Xiaojuan Qi, Xiaoyong Shen, and Jiaya Jia. Image inpainting via generative  
667 multi-column convolutional neural networks. *Advances in neural information processing systems*,  
668 31, 2018.
- 669 Chenwei Wu, Jiajun Luo, and Jason D. Lee. No spurious local minima in a two hidden unit ReLU  
670 network, 2018. URL <https://openreview.net/forum?id=B14uJzW0b>.
- 671 Weihang Xu and Simon S. Du. Over-parameterization exponentially slows down gradient descent for  
672 learning a single neuron, 2023.
- 673 Gilad Yehudai and Ohad Shamir. Learning a single neuron with gradient methods, 2022. URL  
674 <https://arxiv.org/abs/2001.05205>.
- 675 Richard Y. Zhang, Somayeh Sojoudi, and Javad Lavaei. Sharp restricted isometry bounds for the  
676 inexistence of spurious local minima in nonconvex matrix recovery. *J. Mach. Learn. Res.*, 20(114):  
677 1–34, 2019. URL <http://jmlr.org/papers/v20/19-020.html>.
- 678 Xiao Zhang, Yaodong Yu, Lingxiao Wang, and Quanquan Gu. Learning one-hidden-layer relu  
679 networks via gradient descent, 2018. URL <https://arxiv.org/abs/1806.07808>.
- 680 Yuchen Zhang, Percy Liang, and Moses Charikar. A hitting time analysis of stochastic gradi-  
681 ent langevin dynamics. In Satyen Kale and Ohad Shamir, editors, *Proceedings of the 2017*  
682 *Conference on Learning Theory*, volume 65 of *Proceedings of Machine Learning Research*,  
683 pages 1980–2022. PMLR, 07–10 Jul 2017. URL [http://proceedings.mlr.press/](http://proceedings.mlr.press/v65/zhang17b.html)  
684 [v65/zhang17b.html](http://proceedings.mlr.press/v65/zhang17b.html).
- 685 Kai Zhong, Zhao Song, Prateek Jain, Peter L. Bartlett, and Inderjit S. Dhillon. Recovery guarantees for  
686 one-hidden-layer neural networks, 2017. URL <https://arxiv.org/abs/1706.03175>.
- 687 Zhenyu Zhu, Fanghui Liu, and Volkan Cevher. How gradient descent balances features: A dynamical  
688 analysis for two-layer neural networks. In *The Thirteenth International Conference on Learning*  
689 *Representations*, 2025. URL <https://openreview.net/forum?id=25j2ZEgwTj>.
- 690  
691  
692  
693  
694  
695  
696  
697  
698  
699  
700  
701

## APPENDIX

## A USEFUL CALCULATIONS

In this section we provide the derivation of several useful identities.

## A.1 POPULATION LOSS

Let  $\mathbf{a}, \mathbf{b} \in \mathbb{R}^d$  be two arbitrary vectors. Define

$$\begin{aligned} f(\mathbf{a}, \mathbf{b}) &= \mathbb{E}_{\mathbf{x}} \left[ [\mathbf{a}^T \mathbf{x}]_+ [\mathbf{b}^T \mathbf{x}]_+ \right] \\ &\stackrel{(a)}{=} \frac{1}{2\pi} \|\mathbf{a}\| \|\mathbf{b}\| (\sin(\theta_{\mathbf{a}, \mathbf{b}}) + (\pi - \theta_{\mathbf{a}, \mathbf{b}}) \cos(\theta_{\mathbf{a}, \mathbf{b}})) \end{aligned} \quad (5)$$

where  $\theta_{\mathbf{a}, \mathbf{b}} = \cos^{-1} \left( \frac{\mathbf{a}^T \mathbf{b}}{\|\mathbf{a}\| \|\mathbf{b}\|} \right)$ , expectation is over  $\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d)$  and inequality (a) follows from the Table 1 in (Daniely et al., 2016).

Using these we calculate the closed form for the population loss (2) as:

$$\begin{aligned} \mathcal{L}(\boldsymbol{\theta}) &= \frac{1}{2} \mathbb{E}_{\mathbf{x}} \left[ \left\| \sum_{i=1}^k \mathbf{v}_i \phi(\mathbf{w}_i^T \mathbf{x}) - \mathbf{A} \mathbf{x} \right\|^2 \right] \\ &= \frac{1}{2} \sum_{i=1}^k \sum_{j=1}^k \langle \mathbf{v}_i, \mathbf{v}_j \rangle \mathbb{E}_{\mathbf{x}} [\phi(\mathbf{w}_i^T \mathbf{x}) \phi(\mathbf{w}_j^T \mathbf{x})] - \sum_{i=1}^k \mathbf{v}_i^T \mathbf{A} \mathbb{E}_{\mathbf{x}} [\phi(\mathbf{w}_i^T \mathbf{x}) \mathbf{x}] + \frac{1}{2} \mathbb{E}_{\mathbf{x}} [\mathbf{x}^T \mathbf{A}^T \mathbf{A} \mathbf{x}] \\ &\stackrel{(a)}{=} \frac{1}{2} \sum_{i=1}^k \sum_{j=1}^k \langle \mathbf{v}_i, \mathbf{v}_j \rangle f(\mathbf{w}_i, \mathbf{w}_j) - \sum_{i=1}^k \mathbf{v}_i^T \mathbf{A} \mathbb{E}_{\mathbf{x}} [\phi(\mathbf{w}_i^T \mathbf{x}) \mathbf{x}] + \frac{1}{2} \mathbb{E}_{\mathbf{x}} [\mathbf{x}^T \mathbf{A}^T \mathbf{A} \mathbf{x}] \\ &\stackrel{(b)}{=} \frac{1}{2} \sum_{i=1}^k \sum_{j=1}^k \langle \mathbf{v}_i, \mathbf{v}_j \rangle f(\mathbf{w}_i, \mathbf{w}_j) - \sum_{i=1}^k \mathbf{v}_i^T \mathbf{A} \mathbb{E}_{\mathbf{x}} [\nabla_{\mathbf{x}} \phi(\mathbf{w}_i^T \mathbf{x})] + \frac{1}{2} \mathbb{E}_{\mathbf{x}} [\mathbf{x}^T \mathbf{A}^T \mathbf{A} \mathbf{x}] \\ &= \frac{1}{2} \sum_{i=1}^k \sum_{j=1}^k \langle \mathbf{v}_i, \mathbf{v}_j \rangle f(\mathbf{w}_i, \mathbf{w}_j) - \sum_{i=1}^k \mathbf{v}_i^T \mathbf{A} \mathbf{w}_i \mathbb{E}_{\mathbf{x}} [\phi'(\mathbf{w}_i^T \mathbf{x})] + \frac{1}{2} \mathbb{E}_{\mathbf{x}} [\mathbf{x}^T \mathbf{A}^T \mathbf{A} \mathbf{x}] \\ &\stackrel{(c)}{=} \frac{1}{2} \sum_{i=1}^k \sum_{j=1}^k \langle \mathbf{v}_i, \mathbf{v}_j \rangle f(\mathbf{w}_i, \mathbf{w}_j) - \frac{1}{2} \sum_{i=1}^k \mathbf{v}_i^T \mathbf{A} \mathbf{w}_i + \frac{1}{2} \mathbb{E}_{\mathbf{x}} [\mathbf{x}^T \mathbf{A}^T \mathbf{A} \mathbf{x}] \\ &\stackrel{(d)}{=} \frac{1}{2} \sum_{i=1}^k \sum_{j=1}^k \langle \mathbf{v}_i, \mathbf{v}_j \rangle f(\mathbf{w}_i, \mathbf{w}_j) - \frac{1}{2} \sum_{i=1}^k \mathbf{v}_i^T \mathbf{A} \mathbf{w}_i + \frac{1}{2} \text{Tr}(\mathbf{A} \mathbf{A}^T) \\ &= \frac{1}{4\pi} \sum_{i=1}^k \sum_{j=1}^k \langle \mathbf{v}_i, \mathbf{v}_j \rangle \|\mathbf{w}_i\| \|\mathbf{w}_j\| (\sin \theta_{ij} + (\pi - \theta_{ij}) \cos \theta_{ij}) - \frac{1}{2} \sum_{i=1}^k \mathbf{v}_i^T \mathbf{A} \mathbf{w}_i + \frac{1}{2} \text{Tr}(\mathbf{A} \mathbf{A}^T) \end{aligned} \quad (6)$$

where equation (a) follows from the definition of  $f(\mathbf{a}, \mathbf{b})$ , (b) follows from the Stein's Lemma, (c) follows from the fact that derivative of ReLU activation is the step function and  $\mathbf{w}_i^T \mathbf{x} > 0$  with probability  $\frac{1}{2}$ , and finally (d) follows from the cyclical property of the trace.

We also write this in a more compact matrix form as follows:

$$\mathcal{L}(\boldsymbol{\theta}) = \frac{1}{4\pi} \text{Tr}(\mathbf{V}^T \text{diag}(\boldsymbol{\omega}) (\sin(\boldsymbol{\Theta}) + (\pi \mathbb{1} \mathbb{1}^T - \boldsymbol{\Theta}) \odot \cos(\boldsymbol{\Theta})) \text{diag}(\boldsymbol{\omega}) \mathbf{V}) - \frac{1}{2} \text{Tr}((\mathbf{V}^T \mathbf{W} - \mathbf{A}) \mathbf{A}^T)$$

where  $\boldsymbol{\omega}_i = \|\mathbf{w}_i\|$  and  $\theta_{ij}$  is the angle between  $\mathbf{w}_i$  and  $\mathbf{w}_j$ . When  $r = 1$ , the expression simplifies further to

$$\mathcal{L}(\boldsymbol{\theta}) = \frac{1}{4\pi} \mathbf{u}^T (\sin(\boldsymbol{\Theta}) + (\pi \mathbb{1} \mathbb{1}^T - \boldsymbol{\Theta}) \odot \cos(\boldsymbol{\Theta})) \mathbf{u} - \frac{1}{2} \mathbf{a}^T \mathbf{W}^T \mathbf{v} + \frac{1}{2} \|\mathbf{a}\|^2$$

where  $\mathbf{u} = \text{diag}(\boldsymbol{\omega}) \mathbf{v}$ .

## A.2 POPULATION GRADIENT

Let us define,

$$\begin{aligned}
g(\mathbf{a}, \mathbf{b}) &= \frac{\partial}{\partial \mathbf{a}} f(\mathbf{a}, \mathbf{b}) \\
&= \frac{1}{2\pi} (\|\mathbf{b}\| \sin(\theta_{\mathbf{a}, \mathbf{b}}) \bar{\mathbf{a}} + (\pi - \theta_{\mathbf{a}, \mathbf{b}}) \mathbf{b}) \quad \left( \bar{\mathbf{a}} = \frac{\mathbf{a}}{\|\mathbf{a}\|}, \quad \bar{\mathbf{b}} = \frac{\mathbf{b}}{\|\mathbf{b}\|} \right) \\
&= \frac{\|\mathbf{b}\|}{2\pi} (\sin(\theta_{\mathbf{a}, \mathbf{b}}) \bar{\mathbf{a}} + (\pi - \theta_{\mathbf{a}, \mathbf{b}}) \bar{\mathbf{b}}). \tag{7}
\end{aligned}$$

Taking the derivative of (6) with respect to  $\mathbf{w}_i$ , we get

$$\begin{aligned}
\nabla_{\mathbf{w}_i} \mathcal{L}(\boldsymbol{\theta}) &= \frac{1}{2} \|\mathbf{v}_i\|^2 \mathbf{w}_i + \sum_{\substack{j=1 \\ i \neq j}}^k \langle \mathbf{v}_i, \mathbf{v}_j \rangle g(\mathbf{w}_i, \mathbf{w}_j) - \frac{1}{2} \mathbf{A}^T \mathbf{v}_i \\
&= \sum_{j=1}^k \langle \mathbf{v}_i, \mathbf{v}_j \rangle g(\mathbf{w}_i, \mathbf{w}_j) - \frac{1}{2} \mathbf{A}^T \mathbf{v}_i \\
&= \frac{1}{2\pi} \sum_{j=1}^k \langle \mathbf{v}_i, \mathbf{v}_j \rangle \|\mathbf{w}_j\| (\sin(\theta_{ij}) \bar{\mathbf{w}}_i + (\pi - \theta_{ij}) \bar{\mathbf{w}}_j) - \frac{1}{2} \mathbf{A}^T \mathbf{v}_i
\end{aligned}$$

In matrix form:

$$\nabla_{\mathbf{W}} \mathcal{L}(\boldsymbol{\theta}) = \frac{1}{2\pi} ((\mathbf{V}\mathbf{V}^T \odot (\pi \mathbb{1} \mathbb{1}^T - \boldsymbol{\Theta})) \text{diag}(\boldsymbol{\omega}) + \text{diag}((\mathbf{V}\mathbf{V}^T \odot \sin \boldsymbol{\Theta}) \boldsymbol{\omega})) \bar{\mathbf{W}} - \frac{1}{2} \mathbf{V} \mathbf{A}$$

where  $\boldsymbol{\omega}_i = \|\mathbf{w}_i\|$ . When  $r = 1$ , the expression simplifies further to

$$\nabla_{\mathbf{W}} \mathcal{L}(\boldsymbol{\theta}) = \frac{1}{2\pi} \text{diag}(\mathbf{v}) ((\pi \mathbb{1} \mathbb{1}^T - \boldsymbol{\Theta}) \text{diag}(\mathbf{u}) + \text{diag}(\sin(\boldsymbol{\Theta}) \mathbf{u})) \bar{\mathbf{W}} - \frac{1}{2} \mathbf{v} \mathbf{A}^T \tag{8}$$

where  $\mathbf{u} = \text{diag}(\boldsymbol{\omega}) \mathbf{v}$ .

## A.3 POPULATION HESSIAN

We calculate the population Hessian  $\nabla_{\text{vect}(\mathbf{W}), \text{vect}(\mathbf{W})}^2 \mathcal{L}(\boldsymbol{\theta})$  below. Define  $\bar{\mathbf{w}}_l = \frac{\mathbf{w}_l}{\|\mathbf{w}_l\|}$ ,  $\mathbf{P}_{\mathbf{w}_l^\perp} = (\mathbf{I} - \bar{\mathbf{w}}_l \bar{\mathbf{w}}_l^T)$ , and  $\mathbf{w}_{\ell, m^\perp} = \mathbf{P}_{\mathbf{w}_m^\perp} \mathbf{w}_\ell$ . Then,

$$\begin{aligned}
\nabla_{\mathbf{w}_\ell, \mathbf{w}_m}^2 \mathcal{L}(\boldsymbol{\theta}) &= \mathbb{E}_{\mathbf{x}} \left[ (f(\boldsymbol{\theta}; \mathbf{x}) - \mathbf{a}^T \mathbf{x}) \nabla_{\mathbf{w}_\ell, \mathbf{w}_m}^2 f(\boldsymbol{\theta}; \mathbf{x}) + \nabla_{\mathbf{w}_\ell} f(\boldsymbol{\theta}; \mathbf{x}) \nabla_{\mathbf{w}_m} f(\boldsymbol{\theta}; \mathbf{x})^T \right] \tag{9} \\
&= \begin{cases} \frac{\mathbf{v}_\ell^2}{2} \mathbf{I} + \frac{\mathbf{v}_\ell}{2\pi \|\mathbf{w}_\ell\|} \sum_{i=1}^k \mathbf{v}_i \|\mathbf{w}_i\| \sin(\theta_{\ell, i}) \left( \mathbf{P}_{\mathbf{w}_l^\perp} + \frac{\mathbf{P}_{\mathbf{w}_l^\perp} \bar{\mathbf{w}}_i \bar{\mathbf{w}}_i^T \mathbf{P}_{\mathbf{w}_l^\perp}}{\|\mathbf{P}_{\mathbf{w}_l^\perp} \bar{\mathbf{w}}_i\|^2} \right) & \ell = m \\ \frac{\mathbf{v}_\ell \mathbf{v}_m}{2\pi} \left( \bar{\mathbf{w}}_\ell \bar{\mathbf{w}}_{m, \ell^\perp}^T + \bar{\mathbf{w}}_m \bar{\mathbf{w}}_{\ell, m^\perp}^T + (\pi - \theta_{\ell, m}) \mathbf{I} \right) & \ell \neq m \end{cases}
\end{aligned}$$

for calculation of individual terms refer down below.

### A.3.1 CALCULATING $\mathbb{E}_{\mathbf{x}} \left[ \nabla_{\mathbf{w}_\ell} f(\boldsymbol{\theta}; \mathbf{x}) \nabla_{\mathbf{w}_m} f(\boldsymbol{\theta}; \mathbf{x})^T \right]$

$$\begin{aligned}
\mathbb{E}_{\mathbf{x}} \left[ \nabla_{\mathbf{w}_\ell} f(\boldsymbol{\theta}; \mathbf{x}) \nabla_{\mathbf{w}_m} f(\boldsymbol{\theta}; \mathbf{x})^T \right] &= \mathbb{E}_{\mathbf{x}} \left[ \mathbf{v}_\ell \phi'(\mathbf{w}_\ell^T \mathbf{x}) \mathbf{x} \mathbf{x}^T \phi'(\mathbf{w}_m^T \mathbf{x}) \mathbf{v}_m \right] \\
&= \mathbf{v}_\ell \mathbf{v}_m \mathbb{E}_{\mathbf{x}} \left[ \phi'(\bar{\mathbf{w}}_\ell^T \mathbf{x}) \phi'(\bar{\mathbf{w}}_m^T \mathbf{x}) \mathbf{x} \mathbf{x}^T \right]
\end{aligned}$$

To tackle the expectation term, we use second order Stein's Lemma,  $\mathbb{E}_{\mathbf{x}} [g(\mathbf{x}) \mathbf{x} \mathbf{x}^T] = \mathbb{E}_{\mathbf{x}} [\nabla_{\mathbf{x}}^2 g(\mathbf{x})] + \mathbb{E}_{\mathbf{x}} [g(\mathbf{x})] \mathbf{I}$ .

$$\mathbb{E}_{\mathbf{x}} [\phi'(\bar{\mathbf{w}}_\ell^T \mathbf{x}) \phi'(\bar{\mathbf{w}}_m^T \mathbf{x}) \mathbf{x} \mathbf{x}^T] = \mathbb{E}_{\mathbf{x}} [\nabla_{\mathbf{x}}^2 (\phi'(\bar{\mathbf{w}}_\ell^T \mathbf{x}) \phi'(\bar{\mathbf{w}}_m^T \mathbf{x}))] + \mathbb{E}_{\mathbf{x}} [\phi'(\bar{\mathbf{w}}_\ell^T \mathbf{x}) \phi'(\bar{\mathbf{w}}_m^T \mathbf{x})] \mathbf{I}$$

First term is:

$$\begin{aligned} \mathbb{E}_{\mathbf{x}} [\nabla_{\mathbf{x}}^2 (\phi'(\bar{\mathbf{w}}_\ell^T \mathbf{x}) \phi'(\bar{\mathbf{w}}_m^T \mathbf{x}))] &= \mathbb{E}_{\mathbf{x}} [\delta'(\bar{\mathbf{w}}_\ell^T \mathbf{x}) \phi'(\bar{\mathbf{w}}_m^T \mathbf{x}) \bar{\mathbf{w}}_\ell \bar{\mathbf{w}}_\ell^T] \\ &\quad + \mathbb{E}_{\mathbf{x}} [\delta(\bar{\mathbf{w}}_\ell^T \mathbf{x}) \delta(\bar{\mathbf{w}}_m^T \mathbf{x}) (\bar{\mathbf{w}}_\ell \bar{\mathbf{w}}_m^T + \bar{\mathbf{w}}_m \bar{\mathbf{w}}_\ell^T)] \\ &\quad + \mathbb{E}_{\mathbf{x}} [\phi'(\bar{\mathbf{w}}_\ell^T \mathbf{x}) \delta'(\bar{\mathbf{w}}_m^T \mathbf{x}) \bar{\mathbf{w}}_m \bar{\mathbf{w}}_m^T] \end{aligned}$$

These terms can be grouped in two.

$$\begin{aligned} \mathbb{E}_{\mathbf{x}} [\delta'(\bar{\mathbf{w}}_\ell^T \mathbf{x}) \phi'(\bar{\mathbf{w}}_m^T \mathbf{x})] &= \mathbb{E}_{\mathbf{x},g} [\delta'(g) \phi'(\bar{\mathbf{w}}_m^T \mathbf{P}_{\mathbf{w}_\ell^\perp} \mathbf{x} + \bar{\mathbf{w}}_m^T \bar{\mathbf{w}}_\ell g)] \\ &= -\frac{\cos(\theta_{\ell,m})}{\sqrt{2\pi}} \mathbb{E}_{\mathbf{x}} [\delta(\bar{\mathbf{w}}_m^T \mathbf{P}_{\mathbf{w}_\ell^\perp} \mathbf{x})] \quad (\delta'(x) f(x) = -f'(0) \delta(x)) \\ &= -\frac{\cos(\theta_{\ell,m})}{2\pi \|\mathbf{P}_{\mathbf{w}_\ell^\perp} \bar{\mathbf{w}}_m\|} = -\frac{\cos(\theta_{\ell,m})}{2\pi \sin(\theta_{\ell,m})} \end{aligned}$$

and the other one is

$$\begin{aligned} \mathbb{E}_{\mathbf{x}} [\delta(\bar{\mathbf{w}}_\ell^T \mathbf{x}) \delta(\bar{\mathbf{w}}_m^T \mathbf{x})] &= \mathbb{E}_{\mathbf{x},g} [\delta(g) \delta(\bar{\mathbf{w}}_m^T \mathbf{P}_{\mathbf{w}_\ell^\perp} \mathbf{x} + \bar{\mathbf{w}}_m^T \bar{\mathbf{w}}_\ell g)] \\ &= \frac{1}{2\pi \|\mathbf{P}_{\mathbf{w}_\ell^\perp} \bar{\mathbf{w}}_m\|} = \frac{1}{2\pi \sin(\theta_{\ell,m})} \end{aligned}$$

Therefore we get:

$$\mathbb{E}_{\mathbf{x}} [\nabla_{\mathbf{x}}^2 (\phi'(\bar{\mathbf{w}}_\ell^T \mathbf{x}) \phi'(\bar{\mathbf{w}}_m^T \mathbf{x}))] = \frac{\bar{\mathbf{w}}_\ell \bar{\mathbf{w}}_m^T + \bar{\mathbf{w}}_m \bar{\mathbf{w}}_\ell^T}{2\pi \sin(\theta_{\ell,m})} - \frac{\cos(\theta_{\ell,m}) (\bar{\mathbf{w}}_\ell \bar{\mathbf{w}}_\ell^T + \bar{\mathbf{w}}_m \bar{\mathbf{w}}_m^T)}{2\pi \sin(\theta_{\ell,m})}$$

Second term is:

$$\mathbb{E}_{\mathbf{x}} [\phi'(\bar{\mathbf{w}}_\ell^T \mathbf{x}) \phi'(\bar{\mathbf{w}}_m^T \mathbf{x})] \mathbf{I} = \left( \frac{\pi - \theta_{\ell,m}}{2\pi} \right) \mathbf{I} \quad (\text{Dual activation of step function})$$

Combining everything:

$$\mathbb{E}_{\mathbf{x}} [\nabla_{\mathbf{w}_\ell} f(\boldsymbol{\theta}; \mathbf{x}) \nabla_{\mathbf{w}_m} f(\boldsymbol{\theta}; \mathbf{x})^T] = \mathbf{v}_\ell \mathbf{v}_m \left( \frac{\bar{\mathbf{w}}_\ell \bar{\mathbf{w}}_m^T + \bar{\mathbf{w}}_m \bar{\mathbf{w}}_\ell^T - \cos(\theta_{\ell,m}) (\bar{\mathbf{w}}_\ell \bar{\mathbf{w}}_\ell^T + \bar{\mathbf{w}}_m \bar{\mathbf{w}}_m^T)}{2\pi \sin(\theta_{\ell,m})} + \left( \frac{\pi - \theta_{\ell,m}}{2\pi} \right) \mathbf{I} \right)$$

or alternatively (by substituting  $\cos(\theta_{i,j}) = \bar{w}_i^T \bar{w}_j$ ):

$$\mathbb{E}_{\mathbf{x}} [\nabla_{\mathbf{w}_\ell} f(\boldsymbol{\theta}; \mathbf{x}) \nabla_{\mathbf{w}_m} f(\boldsymbol{\theta}; \mathbf{x})^T] = \frac{\mathbf{v}_\ell \mathbf{v}_m}{2\pi} \left( \bar{\mathbf{w}}_\ell \bar{\mathbf{w}}_{m,\ell^\perp}^T + \bar{\mathbf{w}}_m \bar{\mathbf{w}}_{\ell,m^\perp}^T + (\pi - \theta_{\ell,m}) \mathbf{I} \right)$$

### A.3.2 CALCULATING $\mathbb{E}_{\mathbf{x}} [(f(\boldsymbol{\theta}; \mathbf{x}) - \mathbf{a}^T \mathbf{x}) \nabla_{\mathbf{w}_\ell, \mathbf{w}_m}^2 f(\boldsymbol{\theta}; \mathbf{x})]$

Note that  $\nabla_{\mathbf{w}_\ell, \mathbf{w}_m}^2 f(\boldsymbol{\theta}; \mathbf{x}) = \text{diag}(\mathbf{v} \odot \phi''(\mathbf{W}\mathbf{x}))_{\ell,m} \mathbf{x} \mathbf{x}^T$ . This expectation is  $\mathbf{0}$  when  $\ell \neq m$ .

Define  $g = \bar{\mathbf{w}}_i^T \mathbf{x} \sim \mathcal{N}(0, 1)$ .

$$\begin{aligned} \mathbb{E}_{\mathbf{x}} [(f(\boldsymbol{\theta}; \mathbf{x}) - \mathbf{a}^T \mathbf{x}) \nabla_{\mathbf{w}_\ell, \mathbf{w}_\ell}^2 f(\boldsymbol{\theta}; \mathbf{x})] &= \mathbb{E}_{\mathbf{x}} [r(\mathbf{x}) \mathbf{v}_\ell \delta(\bar{\mathbf{w}}_\ell^T \mathbf{x}) \mathbf{x} \mathbf{x}^T] \\ &= \frac{\mathbf{v}_\ell}{\|\mathbf{w}_\ell\|} \mathbb{E}_{\mathbf{x}} [r(\mathbf{x}) \delta(\bar{\mathbf{w}}_\ell^T \mathbf{x}) \mathbf{x} \mathbf{x}^T] \\ &= \frac{\mathbf{v}_\ell}{\|\mathbf{w}_\ell\|} \mathbb{E}_{\mathbf{x},g} \left[ r(\mathbf{P}_{\mathbf{w}_\ell^\perp} \mathbf{x} + \bar{\mathbf{w}}_\ell g) \delta(g) (\mathbf{P}_{\mathbf{w}_\ell^\perp} \mathbf{x} + \bar{\mathbf{w}}_\ell g) (\mathbf{P}_{\mathbf{w}_\ell^\perp} \mathbf{x} + \bar{\mathbf{w}}_\ell g)^T \right] \\ &= \frac{\mathbf{v}_\ell}{\sqrt{2\pi} \|\mathbf{w}_\ell\|} \mathbf{P}_{\mathbf{w}_\ell^\perp} \mathbb{E}_{\mathbf{x}} [r(\mathbf{P}_{\mathbf{w}_\ell^\perp} \mathbf{x}) \mathbf{x} \mathbf{x}^T] \mathbf{P}_{\mathbf{w}_\ell^\perp} \end{aligned}$$

To tackle the expectation term, we use second order Stein's Lemma,  $\mathbb{E}_{\mathbf{x}} [g(\mathbf{x})\mathbf{x}\mathbf{x}^T] = \mathbb{E}_{\mathbf{x}} [\nabla_{\mathbf{x}}^2 g(\mathbf{x})] + \mathbb{E}_{\mathbf{x}} [g(\mathbf{x})] \mathbf{I}$ .

$$\mathbb{E}_{\mathbf{x}} \left[ r \left( \mathbf{P}_{\mathbf{w}_i^\perp} \mathbf{x} \right) \mathbf{x}\mathbf{x}^T \right] = \mathbb{E}_{\mathbf{x}} \left[ \nabla_{\mathbf{x}}^2 r \left( \mathbf{P}_{\mathbf{w}_i^\perp} \mathbf{x} \right) \right] + \mathbb{E}_{\mathbf{x}} \left[ r \left( \mathbf{P}_{\mathbf{w}_i^\perp} \mathbf{x} \right) \right] \mathbf{I}$$

First term is:

$$\begin{aligned} \mathbb{E}_{\mathbf{x}} \left[ \nabla_{\mathbf{x}}^2 r \left( \mathbf{P}_{\mathbf{w}_i^\perp} \mathbf{x} \right) \right] &= \mathbb{E}_{\mathbf{x}} \left[ \nabla_{\mathbf{x}}^2 \left( \sum_{i=1}^k \mathbf{v}_i \phi \left( \mathbf{w}_i^T \mathbf{P}_{\mathbf{w}_i^\perp} \mathbf{x} \right) \right) \right] \quad (\mathbf{a}^T \mathbf{x} \text{ vanishes.}) \\ &= \sum_{i=1}^k \mathbf{v}_i \mathbb{E}_{\mathbf{x}} \left[ \nabla_{\mathbf{x}}^2 \phi \left( \mathbf{w}_i^T \mathbf{P}_{\mathbf{w}_i^\perp} \mathbf{x} \right) \right] \\ &= \sum_{i=1}^k \mathbf{v}_i \mathbb{E}_{\mathbf{x}} \left[ \delta \left( \mathbf{w}_i^T \mathbf{P}_{\mathbf{w}_i^\perp} \mathbf{x} \right) \mathbf{P}_{\mathbf{w}_i^\perp} \mathbf{w}_i \mathbf{w}_i^T \mathbf{P}_{\mathbf{w}_i^\perp} \right] \\ &= \sum_{i=1}^k \frac{\mathbf{v}_i}{\left\| \mathbf{P}_{\mathbf{w}_i^\perp} \mathbf{w}_i \right\|} \mathbb{E}_u \left[ \delta(u) \right] \mathbf{P}_{\mathbf{w}_i^\perp} \mathbf{w}_i \mathbf{w}_i^T \mathbf{P}_{\mathbf{w}_i^\perp} \\ &= \frac{1}{\sqrt{2\pi}} \sum_{i=1}^k \frac{\mathbf{v}_i}{\left\| \mathbf{w}_i \right\| \sin(\theta_{\ell,i})} \mathbf{P}_{\mathbf{w}_i^\perp} \mathbf{w}_i \mathbf{w}_i^T \mathbf{P}_{\mathbf{w}_i^\perp} \quad (\text{Delta integration}) \\ &= \frac{1}{\sqrt{2\pi}} \sum_{i=1}^k \frac{\mathbf{v}_i \left\| \mathbf{w}_i \right\|}{\sin(\theta_{\ell,i})} \mathbf{P}_{\mathbf{w}_i^\perp} \bar{\mathbf{w}}_i \bar{\mathbf{w}}_i^T \mathbf{P}_{\mathbf{w}_i^\perp} \end{aligned}$$

Second term is:

$$\begin{aligned} \mathbb{E}_{\mathbf{x}} \left[ r \left( \mathbf{P}_{\mathbf{w}_i^\perp} \mathbf{x} \right) \right] &= \sum_{i=1}^k \mathbf{v}_i \mathbb{E}_{\mathbf{x}} \left[ \phi \left( \mathbf{w}_i^T \mathbf{P}_{\mathbf{w}_i^\perp} \mathbf{x} \right) \right] \\ &= \sum_{i=1}^k \mathbf{v}_i \left\| \mathbf{P}_{\mathbf{w}_i^\perp} \mathbf{w}_i \right\| \mathbb{E}_u \left[ \phi(u) \right] \quad (u \sim N(0,1)) \\ &= \frac{1}{\sqrt{2\pi}} \sum_{i=1}^k \mathbf{v}_i \left\| \mathbf{w}_i \right\| \sin(\theta_{\ell,i}) \quad \left( \text{Expectation of rectified Gaussian } f_{\mathbf{x}}(0) = \frac{1}{\sqrt{2\pi}} \right) \end{aligned}$$

Combining both terms we get

$$\mathbb{E}_{\mathbf{x}} \left[ r \left( \mathbf{P}_{\mathbf{w}_i^\perp} \mathbf{x} \right) \mathbf{x}\mathbf{x}^T \right] = \frac{1}{\sqrt{2\pi}} \sum_{i=1}^k \mathbf{v}_i \left\| \mathbf{w}_i \right\| \left( \sin(\theta_{\ell,i}) \mathbf{I} + \frac{\mathbf{P}_{\mathbf{w}_i^\perp} \bar{\mathbf{w}}_i \bar{\mathbf{w}}_i^T \mathbf{P}_{\mathbf{w}_i^\perp}}{\sin(\theta_{\ell,i})} \right).$$

Finally we plug this back to get:

$$\begin{aligned} \mathbb{E}_{\mathbf{x}} \left[ r(\mathbf{x}) \nabla_{\mathbf{w}_\ell}^2 f(\boldsymbol{\theta}; \mathbf{x}) \right] &= \frac{\mathbf{v}_\ell}{2\pi \left\| \mathbf{w}_\ell \right\|} \mathbf{P}_{\mathbf{w}_\ell^\perp} \left( \sum_{i=1}^k \mathbf{v}_i \left\| \mathbf{w}_i \right\| \left( \sin(\theta_{\ell,i}) \mathbf{I} + \frac{\mathbf{P}_{\mathbf{w}_i^\perp} \bar{\mathbf{w}}_i \bar{\mathbf{w}}_i^T \mathbf{P}_{\mathbf{w}_i^\perp}}{\sin(\theta_{\ell,i})} \right) \right) \mathbf{P}_{\mathbf{w}_\ell^\perp} \\ &= \frac{\mathbf{v}_\ell}{2\pi \left\| \mathbf{w}_\ell \right\|} \sum_{i=1}^k \mathbf{v}_i \left\| \mathbf{w}_i \right\| \left( \sin(\theta_{\ell,i}) \mathbf{P}_{\mathbf{w}_\ell^\perp} + \frac{\mathbf{P}_{\mathbf{w}_i^\perp} \bar{\mathbf{w}}_i \bar{\mathbf{w}}_i^T \mathbf{P}_{\mathbf{w}_i^\perp}}{\sin(\theta_{\ell,i})} \right) \\ &= \frac{\mathbf{v}_\ell}{2\pi \left\| \mathbf{w}_\ell \right\|} \sum_{i=1}^k \mathbf{v}_i \left\| \mathbf{w}_i \right\| \sin(\theta_{\ell,i}) \left( \mathbf{P}_{\mathbf{w}_\ell^\perp} + \frac{\mathbf{P}_{\mathbf{w}_i^\perp} \bar{\mathbf{w}}_i \bar{\mathbf{w}}_i^T \mathbf{P}_{\mathbf{w}_i^\perp}}{\left\| \mathbf{P}_{\mathbf{w}_\ell^\perp} \bar{\mathbf{w}}_i \right\|^2} \right) \end{aligned}$$

#### A.4 POPULATION CORRELATION

Using the population gradient in Eq. 8, for  $v_1 = v_2 = 1$  we calculate:

$$\begin{aligned}
\langle \mathbf{W} - \mathbf{W}^*, \nabla_{\mathbf{W}} \mathcal{L}(\mathbf{W}) \rangle &= \langle \mathbf{w}_1 - \mathbf{a}, \nabla_{\mathbf{w}_1} \mathcal{L}(\mathbf{W}) \rangle + \langle \mathbf{w}_2 + \mathbf{a}, \nabla_{\mathbf{w}_2} \mathcal{L}(\mathbf{W}) \rangle \\
&= \left\langle \mathbf{w}_1 - \mathbf{a}, -\frac{\mathbf{a}}{2} + \frac{1}{2\pi} (\pi \mathbf{w}_1 - \sin \theta \|\mathbf{w}_2\| \bar{\mathbf{w}}_1 - (\pi - \theta) \mathbf{w}_2) \right\rangle \\
&\quad + \left\langle \mathbf{w}_2 + \mathbf{a}, \frac{\mathbf{a}}{2} - \frac{1}{2\pi} ((\pi - \theta) \mathbf{w}_1 + \sin \theta \|\mathbf{w}_1\| \bar{\mathbf{w}}_2 - \pi \mathbf{w}_2) \right\rangle \\
&= \frac{\|\mathbf{w}_1 - \mathbf{a}\|^2}{2} - \frac{\|\mathbf{w}_2\|}{2\pi} \langle \mathbf{w}_1 - \mathbf{a}, \sin \theta \bar{\mathbf{w}}_1 + (\pi - \theta) \bar{\mathbf{w}}_2 \rangle \\
&\quad + \frac{\|\mathbf{w}_2 + \mathbf{a}\|^2}{2} - \frac{\|\mathbf{w}_1\|}{2\pi} \langle \mathbf{w}_2 + \mathbf{a}, \sin \theta \bar{\mathbf{w}}_2 + (\pi - \theta) \bar{\mathbf{w}}_1 \rangle \\
&= \frac{\|\mathbf{w}_1 - \mathbf{a}\|^2}{2} - \frac{\|\mathbf{w}_2\|}{2\pi} \|\mathbf{w}_1\| (\sin \theta + (\pi - \theta) \cos \theta) \\
&\quad + \frac{\|\mathbf{w}_2\|}{2\pi} \left( \frac{\sin \theta}{\|\mathbf{w}_1\|} \mathbf{w}_1^T \mathbf{a} + \frac{\pi - \theta}{\|\mathbf{w}_2\|} \mathbf{w}_2^T \mathbf{a} \right) + \frac{\|\mathbf{w}_2 + \mathbf{a}\|^2}{2} \\
&\quad - \frac{\|\mathbf{w}_1\|}{2\pi} \|\mathbf{w}_2\| (\sin \theta + (\pi - \theta) \cos \theta) - \frac{\|\mathbf{w}_1\|}{2\pi} \left( \frac{\sin \theta}{\|\mathbf{w}_2\|} \mathbf{w}_2^T \mathbf{a} + \frac{\pi - \theta}{\|\mathbf{w}_1\|} \mathbf{w}_1^T \mathbf{a} \right).
\end{aligned}$$

Rearranging the terms, we obtain:

$$\begin{aligned}
\langle \mathbf{W} - \mathbf{W}^*, \nabla_{\mathbf{W}} \mathcal{L}(\mathbf{W}) \rangle &= \frac{\|\mathbf{w}_1\|^2}{2} + \frac{\|\mathbf{w}_2\|^2}{2} - \left( \frac{(\pi - \theta) \cos \theta + \sin \theta}{\pi} \right) \|\mathbf{w}_1\| \|\mathbf{w}_2\| \\
&\quad - \left( \frac{3\pi - \sin \theta \frac{\|\mathbf{w}_2\|}{\|\mathbf{w}_1\|} - \theta}{2\pi} \right) \mathbf{w}_1^T \mathbf{a} + \left( \frac{3\pi - \sin \theta \frac{\|\mathbf{w}_1\|}{\|\mathbf{w}_2\|} - \theta}{2\pi} \right) \mathbf{w}_2^T \mathbf{a} + \|\mathbf{a}\|^2.
\end{aligned}$$

## B RELU NETWORKS WITH MULTI-DIMENSIONAL OUTPUTS

We now turn our attention to ReLU networks with multiple outputs. In this case running GD from small random initialization on both layers yields an interesting pattern. In particular, at convergence, weights can be grouped into pairs such that one of the weights in the pair is approximately negative of the other one (we discuss this further in Section 4.2).

Given the emergence of this interesting pattern in the outer layer, in our theory we fix the outer layer according to this pattern and focus on the trajectory of the input weights.

**Theorem 6** *Suppose the feature vectors are distributed i.i.d. according to a Gaussian distribution  $\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d)$ . We assume the corresponding output are generated according to a multi-dimensional linear target function of the form  $\mathbf{y} = \mathbf{A}\mathbf{x} \in \mathbb{R}^r$  where  $\mathbf{A} \in \mathbb{R}^{r \times d}$  is an arbitrary matrix. To learn this linear function we fit a one hidden layer ReLU network with  $2r$  hidden nodes of the form*

$$\mathbf{x} \mapsto \mathbf{V}^T \text{ReLU}(\mathbf{W}\mathbf{x})$$

Here,  $\mathbf{W} \in \mathbb{R}^{2r \times d}$  and we fix  $\mathbf{V}$  in the form of  $[\mathbf{I}_r, -\mathbf{I}_r]^T \tilde{\mathbf{V}}$ . Here, we assume  $\tilde{\mathbf{V}} \in \mathbb{R}^{r \times r}$  is of the form  $\tilde{\mathbf{V}} = \mathbf{\Sigma} \mathbf{R}$  where  $\mathbf{R} \in \mathbb{R}^{r \times r}$  is an orthonormal matrix and  $\mathbf{\Sigma}$  is a diagonal matrix with positive entries. Now consider the population loss

$$\mathcal{L}(\mathbf{W}) = \frac{1}{2} \mathbb{E}_{\mathbf{x}} \left[ \|\mathbf{V}^T \text{ReLU}(\mathbf{W}\mathbf{x}) - \mathbf{A}\mathbf{x}\|^2 \right].$$

To fit this model we run gradient updates of the form  $\mathbf{W}^{(\tau+1)} = \mathbf{W}^{(\tau)} - \mu \nabla \mathcal{L}(\mathbf{W}^{(\tau)})$ , starting from an initial estimate  $\mathbf{W}^{(0)}$  with step size obeying  $\mu \leq c_3$ . Furthermore, assume the initialization obeys

$$\|\mathbf{w}_\ell^{(0)} + \mathbf{w}_{\ell+r}^{(0)}\| \leq \frac{1}{2} \|\tilde{\mathbf{a}}_\ell\| \quad \text{for all } \ell = 1, 2, \dots, r,$$

where  $\tilde{\mathbf{a}}_\ell$  is the  $\ell$ th row of  $\Sigma^{-1}\mathbf{R}\mathbf{A}$ , we have

$$\left\| \mathbf{W}^{(\tau)} - \mathbf{W}^* \right\|_F^2 \leq (1 - c_4\mu)^\tau \left\| \mathbf{W}^{(0)} - \mathbf{W}^* \right\|_F^2$$

for all iterations  $\tau$ . Here,  $\mathbf{W}^* = \begin{bmatrix} \Sigma^{-1}\mathbf{R}\mathbf{A} \\ -\Sigma^{-1}\mathbf{R}\mathbf{A} \end{bmatrix}$  and all  $c_j$ 's are fixed numerical constants independent of any problem dimensions.

This result directly generalizes our one dimensional result. The initialization requirement is similar showing that sums of pairs of rows of the inner weights  $\mathbf{W}$  should be sufficiently small at initialization. Similarly, it shows that one can indeed use GD to train a one-hidden layer network with  $2r$  hidden nodes to learn any linear target function with  $r$  outputs. Our result also implies a directional convergence in the sense the rows of the first layer weights align themselves with the corresponding rows of  $\mathbf{A}$  or its negative direction.

## C LANDSCAPE CALCULATIONS

In this section, we verify that:

$$\mathbf{w}_1 = \frac{(c+1)\mathbf{a}}{v_1} \quad \text{and} \quad \mathbf{w}_2 = \frac{c\mathbf{a}}{v_2}$$

for arbitrary  $c > 0$  or  $c < -1$  are indeed stationary points of our optimization problem when  $k = 2$ ,  $r = 1$ . We first show that the gradient vanishes. Plugging the  $\mathbf{w}_1$  and  $\mathbf{w}_2$  into (2):

$$\begin{aligned} \nabla_{\mathbf{W}} \mathcal{L}(\boldsymbol{\theta}) &= \frac{1}{2\pi} \text{diag}(\mathbf{v}) \left( (\pi \mathbb{1}\mathbb{1}^T - \boldsymbol{\Theta}) \text{diag}(\mathbf{u}) + \text{diag}(\sin(\boldsymbol{\Theta})\mathbf{u}) \right) \bar{\mathbf{W}} - \frac{1}{2} \mathbf{v}\mathbf{a}^T \\ &\stackrel{(a)}{=} \frac{1}{2} \text{diag}(\mathbf{v}) \mathbb{1}\mathbb{1}^T \text{diag}(\mathbf{v}) \mathbf{W} - \frac{1}{2} \mathbf{v}\mathbf{a}^T \\ &= \frac{1}{2} \mathbf{v} (\mathbf{W}^T \mathbf{v} - \mathbf{a})^T \\ &= \frac{1}{2} \mathbf{v} (v_1 \mathbf{w}_1 - v_2 \mathbf{w}_2 - \mathbf{a})^T = \frac{1}{2} \mathbf{v} (\mathbf{a} - \mathbf{a})^T = \mathbf{0}. \end{aligned}$$

where (a) follows from the fact that  $\boldsymbol{\Theta} = \mathbf{0}$  at these points. Next we show that the Hessian at these points are PSD. Plugging the values into (9) we get:

$$\nabla_{\mathbf{w}_\ell, \mathbf{w}_m}^2 \mathcal{L}(\boldsymbol{\theta}) = \begin{cases} \frac{v_\ell^2}{2} \mathbf{I} & \ell = m \\ -\frac{v_\ell v_m}{2} \mathbf{I} & \ell \neq m \end{cases}$$

which follows from the fact that  $\theta_{\ell,i} = 0$  and  $\bar{\mathbf{w}}_{m,\ell^\perp} = \bar{\mathbf{w}}_{\ell,m^\perp} = \mathbf{0}$  for any choice of  $\ell, m, i \in [2]$ . Since  $k = 2$ , the resulting Hessian matrix have  $\frac{v_1^2}{2} \mathbf{I}$  and  $\frac{v_2^2}{2} \mathbf{I}$  on its diagonal blocks and  $-\frac{v_1 v_2}{2} \mathbf{I}$  on its off-diagonal blocks. Such a matrix have eigenvalues 0 and  $\frac{v_1^2 + v_2^2}{2}$  each with multiplicity  $d$ . Therefore, all the stationary points are in fact non-strict saddle points of the problem.

## D REDUCTION OF $v_1, v_2 > 0$ TO $v_1 = v_2 = 1$ .

Let us define  $\tilde{\mathbf{w}}_i = v_i \mathbf{w}_i$  as the *simulated* student neuron. We will show that the iterates of the  $v_1, v_2 > 0$  with  $\mathbf{w}_i$  as student neurons is identical to having  $v_1 = v_2 = 1$  but with  $\tilde{\mathbf{w}}_i$  as student neurons. First we note that due to homogeneity of ReLU function, we can combine  $v_i$  and  $\mathbf{w}_i$  without changing the loss function. Furthermore, gradient with respect to  $\mathbf{w}_i$  and  $\tilde{\mathbf{w}}_i$  are related via the following identity:

$$\nabla_{\mathbf{w}_i} \mathcal{L} = \mathcal{D}_{\mathbf{w}_i} (v_i \mathbf{w}_i) \nabla_{\tilde{\mathbf{w}}_i} \mathcal{L} = v_i \mathbf{I} \nabla_{\tilde{\mathbf{w}}_i} \mathcal{L} = v_i \nabla_{\tilde{\mathbf{w}}_i} \mathcal{L}.$$

where  $\mathcal{D}$  denotes the derivative operation. Then, we can rewrite the GD iterates as:

$$\begin{aligned} \mathbf{w}_i^{(t+1)} &= \mathbf{w}_i^{(t)} - \frac{\mu}{v_i^2} \nabla_{\mathbf{w}_i} \mathcal{L} \\ v_i \mathbf{w}_i^{(t+1)} &= v_i \mathbf{w}_i^{(t)} - \mu \frac{\nabla_{\mathbf{w}_i} \mathcal{L}}{v_i} \\ \tilde{\mathbf{w}}_i^{(t+1)} &= \tilde{\mathbf{w}}_i^{(t)} - \mu \nabla_{\tilde{\mathbf{w}}_i} \mathcal{L}. \end{aligned}$$

Therefore, we first prove our theorems with  $v_1 = v_2 = 1$  which directly implies the proof for the general case by plugging in  $\tilde{\mathbf{w}}_i = v_i \mathbf{w}_i$  instead. Note that the direct substitution of  $\mathbf{w}_i \rightarrow v_i \mathbf{w}_i$  only implies:

$$\left\| \text{diag} \left( \begin{bmatrix} v_1 \\ v_2 \end{bmatrix} \right) \mathbf{W}^{(\tau)} - \mathbf{W}^* \right\|_F^2 \leq (1 - c\mu)^\tau \left\| \text{diag} \left( \begin{bmatrix} v_1 \\ v_2 \end{bmatrix} \right) \mathbf{W}^{(0)} - \mathbf{W}^* \right\|_F^2,$$

where  $\mathbf{W}^* = [\mathbf{a} \quad -\mathbf{a}]^T$  and  $c$  is picked appropriately for population and empirical settings separately. To complete the proof, we take the diag terms out, re-define the target to be  $\mathbf{W}^* = \left[ \frac{\mathbf{a}}{v_1} \quad -\frac{\mathbf{a}}{v_2} \right]^T$ , and use a simple inequality to obtain:

$$\min(v_1^2, v_2^2) \left\| \mathbf{W}^{(\tau)} - \mathbf{W}^* \right\|_F^2 \leq (1 - c\mu)^\tau \max(v_1^2, v_2^2) \left\| \mathbf{W}^{(0)} - \mathbf{W}^* \right\|_F^2.$$

Rearranging the terms completes the proof for the  $v_1, v_2 > 0$  setting.

## E PROOF OF KEY LEMMAS

### E.1 PROOF OF KEY LEMMAS IN THE POPULATION SETTING

#### E.1.1 PROOF OF THE MONOTONIC DECREASING PROPERTY (LEMMA 3)

For simplicity of notation, let  $\mathbf{w}_1 = \mathbf{w}_1^{(\tau)}$ ,  $\mathbf{w}_2 = \mathbf{w}_2^{(\tau)}$  and  $\theta$  to be the angle between  $\mathbf{w}_1$  and  $\mathbf{w}_2$ . To derive the formula for  $\mathbf{w}_1^{(\tau+1)} + \mathbf{w}_2^{(\tau+1)}$ , we compute that

$$\nabla_{\mathbf{w}_1} \mathcal{L}(\theta) + \nabla_{\mathbf{w}_2} \mathcal{L}(\theta) = \frac{1}{2} \left( \frac{\theta}{\pi} (\mathbf{w}_1 + \mathbf{w}_2) - \frac{\sin(\theta)}{\pi} (\|\mathbf{w}_2\| \bar{\mathbf{w}}_1 + \|\mathbf{w}_1\| \bar{\mathbf{w}}_2) \right),$$

which means

$$\mathbf{w}_1^{(\tau+1)} + \mathbf{w}_2^{(\tau+1)} = \mathbf{w}_1 + \mathbf{w}_2 - \frac{1}{2} \mu \left( \frac{\theta}{\pi} (\mathbf{w}_1 + \mathbf{w}_2) - \frac{\sin(\theta)}{\pi} (\|\mathbf{w}_2\| \bar{\mathbf{w}}_1 + \|\mathbf{w}_1\| \bar{\mathbf{w}}_2) \right),$$

The square norm of which is

$$\begin{aligned} & \left\| \mathbf{w}_1^{(\tau+1)} + \mathbf{w}_2^{(\tau+1)} \right\|^2 \\ &= \left\| \mathbf{w}_1 + \mathbf{w}_2 - \frac{1}{2} \mu \left( \frac{\theta}{\pi} (\mathbf{w}_1 + \mathbf{w}_2) - \frac{\sin(\theta)}{\pi} (\|\mathbf{w}_2\| \bar{\mathbf{w}}_1 + \|\mathbf{w}_1\| \bar{\mathbf{w}}_2) \right) \right\|^2 \\ &= \left\| \left( \left( 1 - \frac{\mu\theta}{2\pi} \right) \|\mathbf{w}_1\| + \frac{\mu \sin(\theta)}{2\pi} \|\mathbf{w}_2\| \right) \bar{\mathbf{w}}_1 + \left( \left( 1 - \frac{\mu\theta}{2\pi} \right) \|\mathbf{w}_2\| + \frac{\mu \sin(\theta)}{2\pi} \|\mathbf{w}_1\| \right) \bar{\mathbf{w}}_2 \right\|^2 \\ &= \left( \left( 1 - \frac{\mu\theta}{2\pi} \right)^2 + \left( \frac{\mu \sin(\theta)}{2\pi} \right)^2 + 2 \cos(\theta) \left( 1 - \frac{\mu\theta}{2\pi} \right) \frac{\mu \sin(\theta)}{2\pi} \right) (\|\mathbf{w}_1\|^2 + \|\mathbf{w}_2\|^2) \\ & \quad + 2 \left( 2 \left( 1 - \frac{\mu\theta}{2\pi} \right) \frac{\mu \sin(\theta)}{2\pi} + \cos \theta \left( \left( 1 - \frac{\mu\theta}{2\pi} \right)^2 + \left( \frac{\mu \sin(\theta)}{2\pi} \right)^2 \right) \right) \|\mathbf{w}_1\| \|\mathbf{w}_2\|. \end{aligned}$$

For simplicity of notation, call  $\alpha = \left( 1 - \frac{\mu\theta}{2\pi} \right)^2 + \left( \frac{\mu \sin(\theta)}{2\pi} \right)^2$  and  $\beta = 2 \left( 1 - \frac{\mu\theta}{2\pi} \right) \frac{\mu \sin(\theta)}{2\pi}$ . We have that

$$\left\| \mathbf{w}_1^{(\tau+1)} + \mathbf{w}_2^{(\tau+1)} \right\|^2 = (\alpha + \beta \cos(\theta)) (\|\mathbf{w}_1\|^2 + \|\mathbf{w}_2\|^2) + 2(\beta + \alpha \cos(\theta)) \|\mathbf{w}_1\| \|\mathbf{w}_2\|.$$

To show  $\left\| \mathbf{w}_1^{(\tau+1)} + \mathbf{w}_2^{(\tau+1)} \right\| \leq \left\| \mathbf{w}_1^{(\tau)} + \mathbf{w}_2^{(\tau)} \right\|$ , we note that

$$\begin{aligned}
& \left\| \mathbf{w}_1^{(\tau)} + \mathbf{w}_2^{(\tau)} \right\|^2 - \left\| \mathbf{w}_1^{(\tau+1)} + \mathbf{w}_2^{(\tau+1)} \right\|^2 \\
&= \left\| \mathbf{w}_1 + \mathbf{w}_2 \right\|^2 - \left\| \mathbf{w}_1^{(\tau+1)} + \mathbf{w}_2^{(\tau+1)} \right\|^2 \\
&= \left\| \mathbf{w}_1 \right\|^2 + \left\| \mathbf{w}_2 \right\|^2 + 2 \cos \theta \left\| \mathbf{w}_1 \right\| \left\| \mathbf{w}_2 \right\| \\
&\quad - (\alpha + \beta \cos(\theta)) \left( \left\| \mathbf{w}_1 \right\|^2 + \left\| \mathbf{w}_2 \right\|^2 \right) - 2(\beta + \alpha \cos(\theta)) \left\| \mathbf{w}_1 \right\| \left\| \mathbf{w}_2 \right\| \\
&= (1 - \alpha - \beta \cos(\theta)) \left( \left\| \mathbf{w}_1 \right\|^2 + \left\| \mathbf{w}_2 \right\|^2 \right) + 2(\cos(\theta) - \beta - \alpha \cos(\theta)) \left\| \mathbf{w}_1 \right\| \left\| \mathbf{w}_2 \right\|
\end{aligned}$$

To show this quantity is nonnegative, we only need to prove that  $1 - \alpha - \beta \cos(\theta) \geq |\cos(\theta) - \beta - \alpha \cos(\theta)|$ . Since  $\mu \leq 1$ , we have  $1 - \frac{\mu\theta}{2\pi} \geq 0$ , which means  $\beta \geq 0$ . Note that  $\alpha + \beta = \left(1 - \frac{\mu\theta}{2\pi} + \frac{\mu \sin(\theta)}{2\pi}\right)^2 \leq 1$ . We observe that

$$\begin{aligned}
& 1 - \alpha - \beta \cos(\theta) - (\cos(\theta) - \beta - \alpha \cos(\theta)) \\
&= (1 - \cos(\theta))(1 + \beta - \alpha) \geq 0,
\end{aligned}$$

and

$$\begin{aligned}
& 1 - \alpha - \beta \cos(\theta) + (\cos(\theta) - \beta - \alpha \cos(\theta)) \\
&= (1 + \cos(\theta))(1 - \beta - \alpha) \geq 0.
\end{aligned}$$

Combining both inequalities, we obtain that  $1 - \alpha - \beta \cos(\theta) \geq |\cos(\theta) - \beta - \alpha \cos(\theta)|$ . This finishes the proof.

### E.1.2 PROOF OF THE CORRELATION INEQUALITY IN THE POPULATION CASE (LEMMA 4)

As mentioned in the overview section, we define

$$h(\mathbf{w}_1, \mathbf{w}_2, \mathbf{a}) = \langle \mathbf{W} - \mathbf{W}^*, \nabla_{\mathbf{W}} \mathcal{L}(\mathbf{W}) \rangle - \alpha \left\| \mathbf{W} - \mathbf{W}^* \right\|_F^2.$$

Using the calculations in Appendix A.4, we can write it equivalently as

$$\begin{aligned}
h(\mathbf{w}_1, \mathbf{w}_2, \mathbf{a}) &= \underbrace{\left( \frac{1}{2} - \alpha \right) \left\| \mathbf{w}_1 \right\|^2 + \left( \frac{1}{2} - \alpha \right) \left\| \mathbf{w}_2 \right\|^2 - \left( \frac{(\pi - \theta) \cos \theta + \sin \theta}{\pi} \right) \left\| \mathbf{w}_1 \right\| \left\| \mathbf{w}_2 \right\|}_{:=g(\mathbf{w}_1, \mathbf{w}_2)} \\
&\quad - \underbrace{\left( \frac{3\pi - \sin \theta \frac{\left\| \mathbf{w}_2 \right\|}{\left\| \mathbf{w}_1 \right\|} - \theta}{2\pi} + 2\alpha \right)}_{:=\gamma_1} \mathbf{w}_1^T \mathbf{a} + \underbrace{\left( \frac{3\pi - \sin \theta \frac{\left\| \mathbf{w}_1 \right\|}{\left\| \mathbf{w}_2 \right\|} - \theta}{2\pi} - 2\alpha \right)}_{:= -\gamma_2} \mathbf{w}_2^T \mathbf{a} \\
&\quad + (1 - 2\alpha) \left\| \mathbf{a} \right\|^2 \\
&= g(\mathbf{w}_1, \mathbf{w}_2) - (\gamma_1 \mathbf{w}_1 + \gamma_2 \mathbf{w}_2)^T \mathbf{a} + (1 - 2\alpha) \left\| \mathbf{a} \right\|^2 \\
&\geq g(\mathbf{w}_1, \mathbf{w}_2) - \left\| \gamma_1 \mathbf{w}_1 + \gamma_2 \mathbf{w}_2 \right\| \left\| \mathbf{a} \right\| + (1 - 2\alpha) \left\| \mathbf{a} \right\|^2.
\end{aligned}$$

Noting that the expression above is quadratic in  $\left\| \mathbf{a} \right\|$ , we compute  $\tilde{h}(\mathbf{w}_1, \mathbf{w}_2) = \min_{\mathbf{a}} h(\mathbf{w}_1, \mathbf{w}_2, \mathbf{a})$  with the assumption that  $\left\| \mathbf{w}_1 + \mathbf{w}_2 \right\| \leq \frac{1}{2} \left\| \mathbf{a} \right\|$ . The choice of  $\left\| \mathbf{a} \right\|$  that minimizes the expression depends on the location of parabola vertex. Then, we have

$$\left\| \mathbf{a}_{\min} \right\| = \begin{cases} \frac{\left\| \gamma_1 \mathbf{w}_1 + \gamma_2 \mathbf{w}_2 \right\|}{2(1-2\alpha)}, & \text{if } 2 \left\| \mathbf{w}_1 + \mathbf{w}_2 \right\| \leq \frac{\left\| \gamma_1 \mathbf{w}_1 + \gamma_2 \mathbf{w}_2 \right\|}{2(1-2\alpha)}. \\ 2 \left\| \mathbf{w}_1 + \mathbf{w}_2 \right\|, & \text{otherwise} \end{cases}.$$

Plugging it in, we get

$$\tilde{h}(\mathbf{w}_1, \mathbf{w}_2) = \begin{cases} g(\mathbf{w}_1, \mathbf{w}_2) - \frac{\left\| \gamma_1 \mathbf{w}_1 + \gamma_2 \mathbf{w}_2 \right\|^2}{4(1-2\alpha)}, & \text{if } 2 \left\| \mathbf{w}_1 + \mathbf{w}_2 \right\| \leq \frac{\left\| \gamma_1 \mathbf{w}_1 + \gamma_2 \mathbf{w}_2 \right\|}{2(1-2\alpha)} \\ g(\mathbf{w}_1, \mathbf{w}_2) - 2 \left\| \gamma_1 \mathbf{w}_1 + \gamma_2 \mathbf{w}_2 \right\| \left\| \mathbf{w}_1 + \mathbf{w}_2 \right\| + 4(1 - 2\alpha) \left\| \mathbf{w}_1 + \mathbf{w}_2 \right\|^2, & \text{otherwise} \end{cases}.$$

1134  
1135  
1136  
1137  
1138  
1139  
1140  
1141  
1142  
1143  
1144  
1145  
1146  
1147  
1148  
1149  
1150  
1151  
1152  
1153  
1154  
1155  
1156  
1157  
1158  
1159  
1160  
1161  
1162  
1163  
1164  
1165  
1166  
1167  
1168  
1169  
1170  
1171  
1172  
1173  
1174  
1175  
1176  
1177  
1178  
1179  
1180  
1181  
1182  
1183  
1184  
1185  
1186  
1187

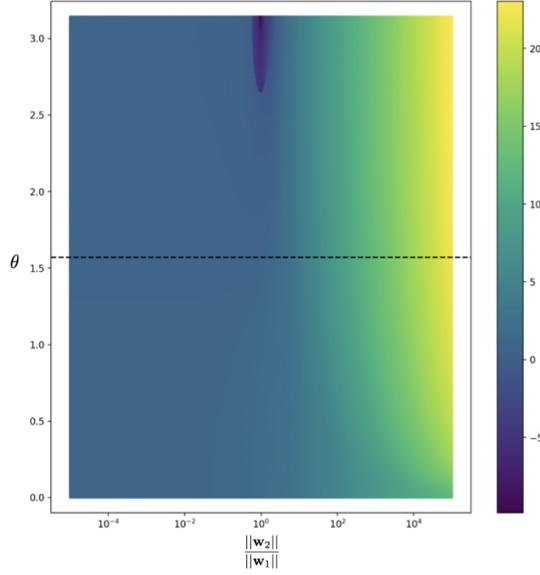


Figure 6:  $\langle \mathbf{W} - \mathbf{W}^*, \nabla_{\mathbf{W}} \mathcal{L}(\mathbf{W}) \rangle - \alpha \|\mathbf{W} - \mathbf{W}^*\|_F^2$  is non-negative. We set  $\alpha = 0.3$  and draw  $\frac{1}{\|\mathbf{w}_1\|} \tilde{h}(\mathbf{w}_1, \mathbf{w}_2)$  for  $\theta \in [0, \pi]$  and  $\frac{\|\mathbf{w}_2\|}{\|\mathbf{w}_1\|} \in [10^{-5}, 10^5]$ . Noting that the color bar is in log-scale, we show that  $\frac{1}{\|\mathbf{w}_1\|} \tilde{h}(\mathbf{w}_1, \mathbf{w}_2)$  is non-negative everywhere in the domain.

Since we are only interested in the sign of  $\tilde{h}$ , we can take  $\|\mathbf{w}_1\|$  outside. Defining the norm ratio  $r$  as  $r = \frac{\|\mathbf{w}_2\|}{\|\mathbf{w}_1\|}$ , it is clear that the expression  $\frac{1}{\|\mathbf{w}_1\|} \tilde{h}(\mathbf{w}_1, \mathbf{w}_2)$  is a function of only  $\theta$  and  $r$ . To complete the proof, in Figure 6, we set  $\alpha = 0.3$  and draw  $\frac{1}{\|\mathbf{w}_1\|} \tilde{h}(\mathbf{w}_1, \mathbf{w}_2)$  as a 2D plot for a wide range of  $r$  and  $\theta$  empirically. The plot demonstrates that  $\tilde{h}$  is non-negative everywhere. This finishes the proof of Lemma 4.

### E.1.3 PROOF OF GRADIENT SMOOTHNESS TOWARDS THE GLOBAL OPTIMA IN THE POPULATION CASE (LEMMA 5)

We first show that

$$\|\sin \theta \|\mathbf{w}_2\| \bar{\mathbf{w}}_1 + (\pi - \theta) \mathbf{w}_2\| \leq \pi \|\mathbf{w}_1 + \mathbf{w}_2\|. \quad (10)$$

Note that  $0 \leq \theta \leq \pi$  since it is the angle between  $\mathbf{w}_1$  and  $\mathbf{w}_2$ . We proceed by case analysis on the value of  $\theta$ . When  $0 \leq \theta < \frac{\pi}{2}$ , we have

$$\begin{aligned} \|\sin \theta \|\mathbf{w}_2\| \bar{\mathbf{w}}_1 + (\pi - \theta) \mathbf{w}_2\| &\stackrel{(a)}{\leq} \|\sin \theta \|\mathbf{w}_2\| \bar{\mathbf{w}}_1\| + \|(\pi - \theta) \mathbf{w}_2\| \\ &= \sin \theta \|\mathbf{w}_2\| + (\pi - \theta) \|\mathbf{w}_2\| \\ &\stackrel{(b)}{\leq} \pi \|\mathbf{w}_2\| \\ &\stackrel{(c)}{\leq} \pi \|\mathbf{w}_1 + \mathbf{w}_2\|. \end{aligned}$$

In Inequality (a) we use the triangle inequality. Inequality (b) follows from the fact that  $\sin \theta \leq \theta$  when  $\theta \geq 0$ . Inequality (c) follows from the fact that  $\theta \leq \frac{\pi}{2}$ .

When  $\theta \geq \frac{\pi}{2}$ , we observe that

$$\begin{aligned}
\|\sin \theta \|\mathbf{w}_2\| \bar{\mathbf{w}}_1 + (\pi - \theta) \mathbf{w}_2\| &\stackrel{(a)}{\leq} \|\sin \theta \|\mathbf{w}_2\| \bar{\mathbf{w}}_1\| + \|(\pi - \theta) \mathbf{w}_2\| \\
&= \sin \theta \|\mathbf{w}_2\| + (\pi - \theta) \|\mathbf{w}_2\| \\
&= \left(1 + \frac{\pi - \theta}{\sin \theta}\right) \sin \theta \|\mathbf{w}_2\| \\
&\stackrel{(b)}{\leq} \left(1 + \frac{\pi - \theta}{\sin \theta}\right) \|\mathbf{w}_1 + \mathbf{w}_2\| \\
&\stackrel{(c)}{\leq} \left(1 + \frac{\pi}{2}\right) \|\mathbf{w}_1 + \mathbf{w}_2\| \\
&\stackrel{(d)}{\leq} \pi \|\mathbf{w}_1 + \mathbf{w}_2\|.
\end{aligned}$$

In Inequality (a) we use the triangle inequality. Inequality (b) follows from the fact that  $\mathbf{w}_1 + \mathbf{w}_2$  has a component with magnitude  $\sin \theta \|\mathbf{w}_2\|$  perpendicular to  $\mathbf{w}_1$ . Inequality (c) follows since  $\frac{\pi - \theta}{\sin \theta}$  attains its maximum at  $\theta = \frac{\pi}{2}$  when restricted to the range  $\theta \geq \frac{\pi}{2}$ . Finally, (d) follows because  $1 \leq \frac{\pi}{2}$ . This finishes the proof of Ineq. 10. Note that due to symmetry we get the following as a corollary:

$$\|\sin \theta \|\mathbf{w}_1\| \bar{\mathbf{w}}_2 + (\pi - \theta) \mathbf{w}_1\| \leq \pi \|\mathbf{w}_1 + \mathbf{w}_2\|. \quad (11)$$

Now recall the population gradient in Eq. 8. The individual rows, i.e. gradient with respect to  $\mathbf{w}_1$  or  $\mathbf{w}_2$ , can be written separately as:

$$\begin{aligned}
\nabla_{\mathbf{w}_1} \mathcal{L} &= -\frac{\mathbf{a}}{2} + \frac{1}{2\pi} (\pi \mathbf{w}_1 - \sin \theta \|\mathbf{w}_2\| \bar{\mathbf{w}}_1 - (\pi - \theta) \mathbf{w}_2) \\
\nabla_{\mathbf{w}_2} \mathcal{L} &= \frac{\mathbf{a}}{2} - \frac{1}{2\pi} ((\pi - \theta) \mathbf{w}_1 + \sin \theta \|\mathbf{w}_1\| \bar{\mathbf{w}}_2 - \pi \mathbf{w}_2).
\end{aligned}$$

Using Ineq. 10, we can write

$$\begin{aligned}
\|\nabla_{\mathbf{w}_1} \mathcal{L}\|^2 &= \left\| -\frac{\mathbf{a}}{2} + \frac{1}{2\pi} (\pi \mathbf{w}_1 - \sin \theta \|\mathbf{w}_2\| \bar{\mathbf{w}}_1 - (\pi - \theta) \mathbf{w}_2) \right\|^2 \\
&= \left\| \frac{\mathbf{w}_1 - \mathbf{a}}{2} - \frac{1}{2\pi} (\sin \theta \|\mathbf{w}_2\| \bar{\mathbf{w}}_1 + (\pi - \theta) \mathbf{w}_2) \right\|^2 \\
&\leq 2 \left\| \frac{\mathbf{w}_1 - \mathbf{a}}{2} \right\|^2 + 2 \left\| \frac{1}{2\pi} (\sin \theta \|\mathbf{w}_2\| \bar{\mathbf{w}}_1 + (\pi - \theta) \mathbf{w}_2) \right\|^2 \\
&\leq \frac{1}{2} \|\mathbf{w}_1 - \mathbf{a}\|^2 + \frac{1}{2} \|\mathbf{w}_1 + \mathbf{w}_2\|^2 \\
&= \frac{1}{2} \|\mathbf{w}_1 - \mathbf{a}\|^2 + \frac{1}{2} \|\mathbf{w}_1 - \mathbf{a} + \mathbf{a} + \mathbf{w}_2\|^2 \\
&\leq \frac{3}{2} \|\mathbf{w}_1 - \mathbf{a}\|^2 + \|\mathbf{w}_2 + \mathbf{a}\|^2.
\end{aligned}$$

Similarly, using Eq. 11 on the gradient for  $\mathbf{w}_2$ , we get

$$\|\nabla_{\mathbf{w}_2} \mathcal{L}\|^2 \leq \frac{3}{2} \|\mathbf{w}_2 + \mathbf{a}\|^2 + \|\mathbf{w}_1 - \mathbf{a}\|^2.$$

Combining these, we obtain

$$\|\nabla_{\mathbf{w}_1} \mathcal{L}\|^2 + \|\nabla_{\mathbf{w}_2} \mathcal{L}\|^2 \leq \frac{5}{2} \left( \|\mathbf{w}_1 - \mathbf{a}\|^2 + \|\mathbf{w}_2 + \mathbf{a}\|^2 \right).$$

This completes the proof of Lemma 5 for  $\beta = \sqrt{\frac{5}{2}}$ .

## E.2 KEY IDENTITIES AND PROOFS IN THE EMPIRICAL SETTING

### E.2.1 GUARANTEES FOR THE FIRST ITERATION

In this section, we will show that with high probability, after the first step of gradient descent with step size 2,  $\mathbf{W}$  is sufficiently close to the global optimum. Specifically, we have the following Lemma:

**Lemma 7** Assume that  $\mathbf{w}_1^{(0)} \sim \mathcal{N}(\mathbf{0}, \sigma_1^2 \mathbf{I}_d)$  and  $\mathbf{w}_2^{(0)} \sim \mathcal{N}(\mathbf{0}, \sigma_2^2 \mathbf{I}_d)$  with the standard deviations obeying  $\sqrt{\sigma_1^2 + \sigma_2^2} \sqrt{d} \leq c_6 \|\mathbf{a}\|$ . Furthermore, assume  $n \geq \frac{C}{\epsilon^2} d$ . After one step of the gradient descent with step size  $\mu = 2$ , we have

$$\|\mathbf{w}_1^{(1)} - \mathbf{a}\|^2 + \|\mathbf{w}_2^{(1)} + \mathbf{a}\|^2 \leq \epsilon^2,$$

with probability at least  $1 - Ce^{-cd}$ . Here  $c, C$  are fixed positive numerical constants.

To prove the Lemma, it suffices to show  $\|\mathbf{w}_1^{(1)} - \mathbf{a}\|^2 \leq \frac{\epsilon^2}{2}$  as proving  $\|\mathbf{w}_2^{(1)} + \mathbf{a}\|^2 \leq \frac{\epsilon^2}{2}$  is the same by a symmetry argument. As a reminder the empirical gradient of  $\mathbf{w}_1$  is

$$\nabla_{\mathbf{w}_1} \widehat{\mathcal{L}} = \frac{1}{n} \sum_{i=1}^n r(\mathbf{x}_i) \phi'(\mathbf{x}_i^T \mathbf{w}_1) \mathbf{x}_i.$$

we have

$$\begin{aligned} \mathbf{w}_1^{(1)} &= \mathbf{w}_1^{(0)} - 2\nabla_{\mathbf{w}_1} \widehat{\mathcal{L}} \\ &= \mathbf{w}_1^{(0)} - \frac{2}{n} \sum_{i=1}^n r(\mathbf{x}_i) \phi'(\mathbf{x}_i^T \mathbf{w}_1^{(0)}) \mathbf{x}_i \\ &= \mathbf{w}_1^{(0)} - \frac{2}{n} \sum_{i=1}^n \left( \phi(\mathbf{x}_i^T \mathbf{w}_1^{(0)}) - \phi(\mathbf{x}_i^T \mathbf{w}_2^{(0)}) - \mathbf{a}^T \mathbf{x}_i \right) \phi'(\mathbf{x}_i^T \mathbf{w}_1^{(0)}) \mathbf{x}_i \\ &= \mathbf{w}_1^{(0)} - \frac{2}{n} \sum_{i=1}^n \left( \phi(\mathbf{x}_i^T \mathbf{w}_1^{(0)}) - \phi(\mathbf{x}_i^T \mathbf{w}_2^{(0)}) - \mathbf{a}^T \mathbf{x}_i \right) \mathbf{1}_{\{\mathbf{x}_i^T \mathbf{w}_1^{(0)} \geq 0\}} \mathbf{x}_i \\ &= \mathbf{w}_1^{(0)} - \frac{2}{n} \sum_{i=1}^n \left( \phi(\mathbf{x}_i^T \mathbf{w}_1^{(0)}) - \phi(\mathbf{x}_i^T \mathbf{w}_2^{(0)}) - \phi(\mathbf{x}_i^T \mathbf{a}) + \phi(\mathbf{x}_i^T (-\mathbf{a})) \right) \mathbf{1}_{\{\mathbf{x}_i^T \mathbf{w}_1^{(0)} \geq 0\}} \mathbf{x}_i. \end{aligned}$$

To prove this vector is close to  $\mathbf{a}$ , we first show the following Lemma:

**Lemma 8**  $\{\mathbf{x}_i\}_{i=1}^n$  are i.i.d. Gaussian random vectors distributed as  $\mathcal{N}(\mathbf{0}, \mathbf{I}_d)$ . Furthermore, assume

$$n \geq \frac{C}{\delta^2} d.$$

For any fixed vector  $\mathbf{p}, \mathbf{q}$ , we have

$$\left\| \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{\{\mathbf{x}_i^T \mathbf{p} \geq 0\}} \mathbf{x}_i \phi(\mathbf{x}_i^T \mathbf{q}) - \mathbb{E} [\mathbf{1}_{\{\mathbf{x}^T \mathbf{p} \geq 0\}} \mathbf{x} \phi(\mathbf{x}^T \mathbf{q})] \right\| \leq \delta \|\mathbf{q}\|$$

with probability at least  $1 - 6e^{-\gamma d}$ . Here,  $\phi: \mathbb{R} \mapsto \mathbb{R}$  is the ReLU function.

Without loss of generality, assume that  $\mathbf{p} = \mathbf{e}_1$ ,  $\text{span}(\mathbf{p}, \mathbf{q}) = (\mathbf{e}_1, \mathbf{e}_2)$  and  $\|\mathbf{q}\| = 1$ . Consider a random variable  $Y = \mathbf{1}_{\{\mathbf{x}^T \mathbf{p} \geq 0\}} \mathbf{x} \phi(\mathbf{x}^T \mathbf{q})$ . Assume that  $\mathbf{x} = [g_1, g_2, \dots, g_d]^T$ . The last  $d - 2$  entries of  $Y$ , defined as  $Y_2$ , is in the form of  $\mathbf{1}_{\{g_1 \geq 0\}} \phi(\mathbf{x}^T \mathbf{q}) [g_3, \dots, g_d]^T$ . We observe that  $[g_3, \dots, g_d]^T$  is independent with  $\mathbf{1}_{\{g_1 \geq 0\}} \phi(\mathbf{x}^T \mathbf{q})$ . For fixed  $\mathbf{1}_{\{\mathbf{x}_i^T \mathbf{p} \geq 0\}} \phi(\mathbf{x}_i^T \mathbf{q}), i = 1, 2, \dots, n$ ,  $\frac{1}{n} \sum_{i=1}^n \mathbf{1}_{\{\mathbf{x}_i^T \mathbf{p} \geq 0\}} \mathbf{x}_i \phi(\mathbf{x}_i^T \mathbf{q})$  is a weighted sum of independent Gaussian random vectors, which means it is also a Gaussian random vector, the variance of which is

$$\frac{1}{n^2} \sum_{i=1}^n (\mathbf{1}_{\{\mathbf{x}_i^T \mathbf{p} \geq 0\}} \phi(\mathbf{x}_i^T \mathbf{q}))^2 \cdot \mathbf{I}_{d-2}.$$

Denote  $\sigma = \frac{1}{n^2} \sum_{i=1}^n (\mathbf{1}_{\{\mathbf{x}_i^T \mathbf{p} \geq 0\}} \phi(\mathbf{x}_i^T \mathbf{q}))^2$ . we have that

$$\Pr \left[ \|Y_2\| \geq 2\sigma \sqrt{d} \right] \leq e^{-d/2}.$$

Let  $Z = \mathbf{1}_{\{\mathbf{x}^T \mathbf{p} \geq 0\}} \phi(\mathbf{x}^T \mathbf{q})$ . Note that  $Z$  is a sub-gaussian random variable, which implies that  $Z^2$  is sub-exponential. Indeed, the sub-exponential norm of  $Z^2$ , defined as

$$\|Z^2\|_{\psi_1} = \sup_{p \geq 1} p^{-1} (\mathbb{E}|Z^2|^p)^{1/p},$$

is upper bounded by some numerical constant  $\gamma_1$ . By the Bernstein-type inequality, we have that

$$\Pr \left[ \left| \sum_{i=1}^n Z_i^2 - \mathbb{E}[Z_i^2] \right| \geq t \right] \leq 2 \exp \left[ -c \min \left( \frac{t^2}{n\gamma_1^2}, \frac{t}{\gamma_1} \right) \right].$$

By union bound, we have that with probability at least  $1 - 4e^{-\gamma d}$ ,  $\|Y_2\| \leq \frac{\delta}{2}$ . Similarly, the first two entries of  $Y$  is also sub-exponential, the sub-exponential norm of which is upper bounded by some numerical constant  $\gamma_2$ . By the Bernstein-type inequality, we have that with probability at least  $1 - 2e^{-\gamma d}$ ,  $\|Y_1\| \leq \frac{\delta}{2}$ . By union bound, we have that with probability at least  $1 - 6e^{-\gamma d}$ ,  $\|Y\| \leq \|Y_1\| + \|Y_2\| \leq \delta$ , which finishes the proof of the Lemma.

We use this Lemma with  $\mathbf{p} = \mathbf{w}_1^{(0)}$  and  $\mathbf{q} = \mathbf{w}_1^{(0)}, \mathbf{w}_2^{(0)}, -\mathbf{a}, \mathbf{a}$  and sum the four inequalities together. Combining with the triangle inequality, we obtain that

$$\begin{aligned} & \left\| \frac{1}{n} \sum_{i=1}^n \left( \phi(\mathbf{x}_i^T \mathbf{w}_1^{(0)}) - \phi(\mathbf{x}_i^T \mathbf{w}_2^{(0)}) - \mathbf{a}^T \mathbf{x}_i \right) \mathbf{1}_{\{\mathbf{x}_i^T \mathbf{w}_1^{(0)} \geq 0\}} \mathbf{x}_i \right. \\ & \quad \left. - \mathbb{E} \left[ \left( \phi(\mathbf{x}^T \mathbf{w}_1^{(0)}) - \phi(\mathbf{x}^T \mathbf{w}_2^{(0)}) - \mathbf{a}^T \mathbf{x} \right) \mathbf{1}_{\{\mathbf{x}^T \mathbf{w}_1^{(0)} \geq 0\}} \mathbf{x} \right] \right\| \\ & \leq \delta (\|\mathbf{w}_1^{(0)}\| + \|\mathbf{w}_2^{(0)}\| + 2\|\mathbf{a}\|), \end{aligned} \quad (12)$$

with probability at least  $1 - 24e^{-\gamma d}$ . We compute that

$$\begin{aligned} & \mathbb{E} \left[ \left( \phi(\mathbf{x}^T \mathbf{w}_1^{(0)}) - \phi(\mathbf{x}^T \mathbf{w}_2^{(0)}) - \mathbf{a}^T \mathbf{x} \right) \mathbf{1}_{\{\mathbf{x}^T \mathbf{w}_1^{(0)} \geq 0\}} \mathbf{x} \right] \\ & = \mathbb{E} \left[ \left( \phi(\mathbf{x}^T \mathbf{w}_1^{(0)}) - \phi(\mathbf{x}^T \mathbf{w}_2^{(0)}) \right) \mathbf{1}_{\{\mathbf{x}^T \mathbf{w}_1^{(0)} \geq 0\}} \mathbf{x} \right] - \frac{\mathbf{a}}{2}. \end{aligned}$$

This implies that

$$\begin{aligned} \mathbf{w}_1^{(1)} - \mathbf{a} &= \mathbf{w}_1^{(0)} - \frac{2}{n} \sum_{i=1}^n \left( \phi(\mathbf{x}_i^T \mathbf{w}_1^{(0)}) - \phi(\mathbf{x}_i^T \mathbf{w}_2^{(0)}) - \mathbf{a}^T \mathbf{x}_i \right) \mathbf{1}_{\{\mathbf{x}_i^T \mathbf{w}_1^{(0)} \geq 0\}} \mathbf{x}_i - \mathbf{a} \\ &= \mathbf{w}_1^{(0)} - \mathbf{a} - 2 \left( \frac{1}{n} \sum_{i=1}^n \left( \phi(\mathbf{x}_i^T \mathbf{w}_1^{(0)}) - \phi(\mathbf{x}_i^T \mathbf{w}_2^{(0)}) - \mathbf{a}^T \mathbf{x}_i \right) \mathbf{1}_{\{\mathbf{x}_i^T \mathbf{w}_1^{(0)} \geq 0\}} \mathbf{x}_i \right. \\ & \quad \left. - \mathbb{E} \left[ \left( \phi(\mathbf{x}^T \mathbf{w}_1^{(0)}) - \phi(\mathbf{x}^T \mathbf{w}_2^{(0)}) - \mathbf{a}^T \mathbf{x} \right) \mathbf{1}_{\{\mathbf{x}^T \mathbf{w}_1^{(0)} \geq 0\}} \mathbf{x} \right] \right) \\ & \quad - 2 \mathbb{E} \left[ \left( \phi(\mathbf{x}^T \mathbf{w}_1^{(0)}) - \phi(\mathbf{x}^T \mathbf{w}_2^{(0)}) - \mathbf{a}^T \mathbf{x} \right) \mathbf{1}_{\{\mathbf{x}^T \mathbf{w}_1^{(0)} \geq 0\}} \mathbf{x} \right] \\ &= \mathbf{w}_1^{(0)} - 2 \left( \frac{1}{n} \sum_{i=1}^n \left( \phi(\mathbf{x}_i^T \mathbf{w}_1^{(0)}) - \phi(\mathbf{x}_i^T \mathbf{w}_2^{(0)}) - \mathbf{a}^T \mathbf{x}_i \right) \mathbf{1}_{\{\mathbf{x}_i^T \mathbf{w}_1^{(0)} \geq 0\}} \mathbf{x}_i \right. \\ & \quad \left. - \mathbb{E} \left[ \left( \phi(\mathbf{x}^T \mathbf{w}_1^{(0)}) - \phi(\mathbf{x}^T \mathbf{w}_2^{(0)}) - \mathbf{a}^T \mathbf{x} \right) \mathbf{1}_{\{\mathbf{x}^T \mathbf{w}_1^{(0)} \geq 0\}} \mathbf{x} \right] \right) \\ & \quad - 2 \mathbb{E} \left[ \left( \phi(\mathbf{x}^T \mathbf{w}_1^{(0)}) - \phi(\mathbf{x}^T \mathbf{w}_2^{(0)}) \right) \mathbf{1}_{\{\mathbf{x}^T \mathbf{w}_1^{(0)} \geq 0\}} \mathbf{x} \right]. \end{aligned}$$

Plugging in the Ineq.12 and using the triangle inequality, we obtain that

$$\begin{aligned} \|\mathbf{w}_1^{(1)} - \mathbf{a}\| &\leq \|\mathbf{w}_1^{(0)}\| + 2\delta (\|\mathbf{w}_1^{(0)}\| + \|\mathbf{w}_2^{(0)}\| + 2\|\mathbf{a}\|) \\ &\quad + 2 \left\| \mathbb{E} \left[ \left( \phi(\mathbf{x}^T \mathbf{w}_1^{(0)}) - \phi(\mathbf{x}^T \mathbf{w}_2^{(0)}) \right) \mathbf{1}_{\{\mathbf{x}^T \mathbf{w}_1^{(0)} \geq 0\}} \mathbf{x} \right] \right\|. \end{aligned}$$

By computation, we have that

$$\left\| \mathbb{E} \left[ \left( \phi \left( \mathbf{x}^T \mathbf{w}_1^{(0)} \right) - \phi \left( \mathbf{x}^T \mathbf{w}_2^{(0)} \right) \right) \mathbf{1}_{\{\mathbf{x}^T \mathbf{w}_1^{(0)} \geq 0\}} \mathbf{x} \right] \right\| \leq \|\mathbf{w}_1^{(0)}\| + \|\mathbf{w}_2^{(0)}\|,$$

which implies that

$$\|\mathbf{w}_1^{(1)} - \mathbf{a}\| \leq (3 + 2\delta)(\|\mathbf{w}_1^{(0)}\| + \|\mathbf{w}_2^{(0)}\|) + 4\delta\|\mathbf{a}\|.$$

Remind that  $\mathbf{w}_1^{(0)} \sim \mathcal{N}(\mathbf{0}, \sigma_1^2 \mathbf{I}_d)$  and  $\mathbf{w}_2^{(0)} \sim \mathcal{N}(\mathbf{0}, \sigma_2^2 \mathbf{I}_d)$  with the standard deviations obeying  $\sqrt{\sigma_1^2 + \sigma_2^2} \sqrt{d} \leq c_6 \|\mathbf{a}\|$ , it is well-known that with probability at least  $1 - 2e^{-\frac{d}{8}}$ ,  $\|\mathbf{w}_1^{(0)}\| + \|\mathbf{w}_2^{(0)}\| \leq 2c_6 \|\mathbf{a}\|$ . Let  $c_6 = \frac{1}{40}\epsilon$ ,  $\delta \leq \frac{1}{40}\epsilon$ . For any  $\epsilon \leq 1$ , we have

$$\begin{aligned} \|\mathbf{w}_1^{(1)} - \mathbf{a}\| &\leq (3 + 2\delta)(\|\mathbf{w}_1^{(0)}\| + \|\mathbf{w}_2^{(0)}\|) + 4\delta\|\mathbf{a}\| \\ &\leq (3 + 2\delta) \cdot 2c_6 \|\mathbf{a}\| + 4\delta\|\mathbf{a}\| \\ &\leq \frac{\epsilon}{\sqrt{2}}, \end{aligned}$$

with probability at least  $1 - 24e^{-\gamma d} - 2e^{-\frac{d}{8}}$ . Similarly, we have  $\|\mathbf{w}_2^{(1)} - \mathbf{a}\| \leq \frac{\epsilon}{\sqrt{2}}$  with probability at least  $1 - 24e^{-\gamma d} - 2e^{-\frac{d}{8}}$ . By union bound, we conclude that  $\|\mathbf{w}_1^{(1)} - \mathbf{a}\|^2 + \|\mathbf{w}_2^{(1)} - \mathbf{a}\|^2 \leq \epsilon^2$  with probability at least  $1 - Ce^{-cd}$  with some positive numerical constant  $c, C$ . This finishes the proof of the Lemma 7.

## E.2.2 PROOF OF THE CORRELATION INEQUALITY IN THE EMPIRICAL CASE

Without loss of generality, assume that  $\|\mathbf{a}\| = 1$ . For simplicity of notations, define  $\phi(x) = \text{ReLU}(x)$ . We observe that

$$\begin{aligned} &\langle \nabla_{\mathbf{w}_1} \widehat{\mathcal{L}}, \mathbf{w}_1 - \mathbf{a} \rangle + \langle \nabla_{\mathbf{w}_2} \widehat{\mathcal{L}}, \mathbf{w}_2 + \mathbf{a} \rangle \\ &= \frac{1}{n} \sum_{i=1}^n r(\mathbf{x}_i) \phi'(\mathbf{w}_1^T \mathbf{x}_i) \mathbf{x}_i^T (\mathbf{w}_1 - \mathbf{a}) - \frac{1}{n} \sum_{i=1}^n r(\mathbf{x}_i) \phi'(\mathbf{w}_2^T \mathbf{x}_i) \mathbf{x}_i^T (\mathbf{w}_2 + \mathbf{a}) \\ &= \frac{1}{n} \sum_{i=1}^n r(\mathbf{x}_i) (\phi(\mathbf{x}_i^T \mathbf{w}_1) - \phi(\mathbf{x}_i^T \mathbf{w}_2)) - \frac{1}{n} \sum_{i=1}^n r(\mathbf{x}_i) (\mathbf{a}^T \mathbf{x}_i) (\phi'(\mathbf{w}_1^T \mathbf{x}_i) + \phi'(\mathbf{w}_2^T \mathbf{x}_i)) \\ &= \frac{1}{n} \sum_{i=1}^n r^2(\mathbf{x}_i) + \frac{1}{n} \sum_{i=1}^n r(\mathbf{x}_i) (\mathbf{a}^T \mathbf{x}_i) (1 - (\phi'(\mathbf{w}_1^T \mathbf{x}_i) + \phi'(\mathbf{w}_2^T \mathbf{x}_i))). \end{aligned}$$

By the Cauchy-Schwarz inequality, we have

$$\begin{aligned} LHS &\geq \frac{1}{n} \sum_{i=1}^n r^2(\mathbf{x}_i) - \frac{1}{n} \sqrt{\sum_{i=1}^n r^2(\mathbf{x}_i) \sum_{i=1}^n (\mathbf{a}^T \mathbf{x}_i)^2 (1 - (\phi'(\mathbf{w}_1^T \mathbf{x}_i) + \phi'(\mathbf{w}_2^T \mathbf{x}_i)))^2} \\ &\geq \frac{1}{n} \sum_{i=1}^n r^2(\mathbf{x}_i) - \frac{1}{n} \sqrt{\sum_{i=1}^n r^2(\mathbf{x}_i) \sum_{i=1}^n (\mathbf{a}^T \mathbf{x}_i)^2 (1 - \text{sgn}(\mathbf{a}^T \mathbf{x}_i) \text{sgn}(\mathbf{w}_1^T \mathbf{x}_i) + 1 - \text{sgn}(-\mathbf{a}^T \mathbf{x}_i) \text{sgn}(\mathbf{w}_2^T \mathbf{x}_i))}. \end{aligned}$$

where the last inequality is derived by decomposing 1 as  $\phi'(\mathbf{a}^T \mathbf{x}_i) + \phi'(-\mathbf{a}^T \mathbf{x}_i)$ , using the identity  $\text{sgn}(\cdot) + 1 = 2\phi'(\cdot)$ , and the fact that  $(x + y)^2 \leq 2(x^2 + y^2)$ .

The following inequality is established within the proof of (Soltanolkotabi, 2019, Lemma 7.17):

**Lemma 9** Assume the measurement vectors  $\{\mathbf{x}_i\}_{i=1}^n$  are i.i.d. Gaussian random vectors distributed as  $\mathcal{N}(\mathbf{0}, \mathbf{I}_d)$ , where  $n \geq \frac{cd}{\delta}$ . Furthermore, assume  $\mathbf{a}$  is a fixed vector independent of the measurement vectors such that  $\|\mathbf{a}\| = 1$ . Define the set  $E(\epsilon) = \{\mathbf{w} \in \mathbb{R}^d \mid \|\mathbf{w} - \mathbf{a}\| \leq \epsilon\}$ . Assume that  $\delta \leq \frac{\epsilon}{100}$ .

We have

$$\begin{aligned}
\sqrt{\frac{1}{n} \sum_{i=1}^n (\mathbf{a}^T \mathbf{x}_i)^2 (1 - \text{sgn}(\mathbf{a}^T \mathbf{x}_i) \text{sgn}(\mathbf{w}^T \mathbf{x}_i))} &\leq \frac{\sqrt{2}}{1-\epsilon} \sqrt{\frac{1}{n} \sum_{i=1}^n \mathbf{1}_{\{(1-\epsilon)|\mathbf{x}_i^T \mathbf{a}| \leq |\mathbf{x}_i^T \mathbf{h}_\perp|\}} (\mathbf{h}_\perp^T \mathbf{x}_i)^2} \\
&\leq \frac{\sqrt{2}}{1-\epsilon} \left( \delta + \sqrt{\frac{21}{20}} \epsilon \right) \|\mathbf{w} - \mathbf{a}\| \\
&\leq \frac{3}{2} \cdot \frac{\epsilon}{1-\epsilon} \|\mathbf{w} - \mathbf{a}\|,
\end{aligned}$$

holds for all  $\mathbf{w} \in E(\epsilon)$  with probability at least  $1 - 3e^{-\gamma n}$ , where  $\mathbf{h}_\perp$  is the part of  $\mathbf{h} = \mathbf{w} + \mathbf{a}$  that is perpendicular to  $\mathbf{a}$  and  $\gamma$  is a fixed numerical constant.

This allows us to proceed as follows:

$$\begin{aligned}
LHS &\geq \frac{1}{n} \sum_{i=1}^n r^2(\mathbf{x}_i) - \sqrt{\frac{1}{n} \sum_{i=1}^n r^2(\mathbf{x}_i) \sqrt{\left(\frac{3}{2} \cdot \frac{\epsilon}{1-\epsilon} \|\mathbf{w}_1 - \mathbf{a}\|\right)^2 + \left(\frac{3}{2} \cdot \frac{\epsilon}{1-\epsilon} \|\mathbf{w}_2 + \mathbf{a}\|\right)^2}} \\
&\geq \frac{1}{n} \sum_{i=1}^n r^2(\mathbf{x}_i) - \frac{1}{2} \left( \frac{1}{n} \sum_{i=1}^n r^2(\mathbf{x}_i) + \left(\frac{3}{2} \cdot \frac{\epsilon}{1-\epsilon} \|\mathbf{w}_1 - \mathbf{a}\|\right)^2 + \left(\frac{3}{2} \cdot \frac{\epsilon}{1-\epsilon} \|\mathbf{w}_2 + \mathbf{a}\|\right)^2 \right) \\
&= \frac{1}{2n} \sum_{i=1}^n r^2(\mathbf{x}_i) - \frac{1}{2} \left(\frac{3}{2} \cdot \frac{\epsilon}{1-\epsilon}\right)^2 (\|\mathbf{w}_1 - \mathbf{a}\|^2 + \|\mathbf{w}_2 + \mathbf{a}\|^2), \tag{13}
\end{aligned}$$

holds for all  $\mathbf{W}$  such that  $\|\mathbf{W} - \mathbf{W}^*\|_F \leq \epsilon$  with probability at least  $1 - 6e^{-\gamma n}$ .

Next we will lower bound

$$\frac{1}{n} \sum_{i=1}^n r^2(\mathbf{x}_i)$$

To do this note that

$$\frac{1}{n} \sum_{i=1}^n r^2(\mathbf{x}_i) = \frac{1}{n} \sum_{i=1}^n (\mathbf{1}_{\{\mathbf{a}^T \mathbf{x}_i \geq 0\}} r^2(\mathbf{x}_i) + \mathbf{1}_{\{\mathbf{a}^T \mathbf{x}_i < 0\}} r^2(\mathbf{x}_i)).$$

We proceed by bounding the first term. The bound of the second term is the same by a symmetry argument. To bound the first term we proceed by the following chain of inequalities

$$\begin{aligned}
\frac{1}{n} \sum_{i=1}^n \mathbf{1}_{\{\mathbf{a}^T \mathbf{x}_i \geq 0\}} r^2(\mathbf{x}_i) &\geq \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{\{\mathbf{a}^T \mathbf{x}_i \geq 0\}} r^2(\mathbf{x}_i) \\
&\stackrel{(a)}{\geq} \frac{1}{2n} \sum_{i=1}^n \mathbf{1}_{\{\mathbf{a}^T \mathbf{x}_i \geq 0\}} (\phi(\mathbf{w}_1^T \mathbf{x}_i) - \mathbf{a}^T \mathbf{x}_i)^2 - \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{\{\mathbf{a}^T \mathbf{x}_i \geq 0\}} \phi^2(\mathbf{w}_2^T \mathbf{x}_i) \\
&\stackrel{(b)}{\geq} \frac{1}{2n} \sum_{i=1}^n \mathbf{1}_{\{\mathbf{a}^T \mathbf{x}_i \geq 0\}} \mathbf{1}_{\{\mathbf{w}_1^T \mathbf{x}_i \geq 0\}} (\phi(\mathbf{w}_1^T \mathbf{x}_i) - \mathbf{a}^T \mathbf{x}_i)^2 - \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{\{\mathbf{a}^T \mathbf{x}_i \geq 0\}} \phi^2(\mathbf{w}_2^T \mathbf{x}_i) \\
&= \frac{1}{2n} \sum_{i=1}^n \mathbf{1}_{\{\mathbf{a}^T \mathbf{x}_i \geq 0\}} \mathbf{1}_{\{\mathbf{w}_1^T \mathbf{x}_i \geq 0\}} (\mathbf{h}_1^T \mathbf{x}_i)^2 - \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{\{\mathbf{a}^T \mathbf{x}_i \geq 0\}} \phi^2(\mathbf{w}_2^T \mathbf{x}_i) \\
&\stackrel{(c)}{\geq} c \|\mathbf{h}_1\|^2 - \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{\{\mathbf{a}^T \mathbf{x}_i \geq 0\}} \phi^2(\mathbf{w}_2^T \mathbf{x}_i). \tag{14}
\end{aligned}$$

Here, (a) follows from  $(x - y)^2 \geq \frac{1}{2}x^2 - y^2$ , (b) from  $1 \geq \mathbf{1}_{\{\mathbf{w}_1^T \mathbf{x}_i \geq 0\}}$ , and (c) from the Lemma 11 proven in Appendix E.2.4 below. To bound the second term note that we can rewrite  $\mathbf{w}_2$  as follows

$$\mathbf{w}_2 = -(-\mathbf{w}_2^T \mathbf{a}) \mathbf{a} + \mathbf{h}_{2,\perp},$$

where  $\mathbf{h}_{2,\perp}$  is the part of  $\mathbf{h}_2 = \mathbf{w}_2 + \mathbf{a}$  that is perpendicular to  $\mathbf{a}$ . Now note that when  $\mathbf{a}^T \mathbf{x}_i \geq 0$  and  $\|\mathbf{h}_2\| \leq \epsilon$  we have

$$\mathbf{w}_2^T \mathbf{x}_i = (-\mathbf{w}_2^T \mathbf{a})(-\mathbf{a}^T \mathbf{x}_i) + \mathbf{h}_{2,\perp}^T \mathbf{x}_i \leq \mathbf{h}_{2,\perp}^T \mathbf{x}_i,$$

and

$$\sum_{i=1}^n \mathbf{1}_{\{\mathbf{a}^T \mathbf{x}_i \geq 0\}} \phi^2(\mathbf{w}_2^T \mathbf{x}_i) = \sum_{i=1}^n \mathbf{1}_{\{\mathbf{a}^T \mathbf{x}_i \geq 0\}} \mathbf{1}_{\{\mathbf{w}_2^T \mathbf{x}_i \geq 0\}} (\mathbf{w}_2^T \mathbf{x}_i)^2 \leq \sum_{i=1}^n \mathbf{1}_{\{\mathbf{a}^T \mathbf{x}_i \geq 0\}} \mathbf{1}_{\{\mathbf{w}_2^T \mathbf{x}_i \geq 0\}} (\mathbf{h}_{2,\perp}^T \mathbf{x}_i)^2. \quad (15)$$

Next note that when  $\mathbf{a}^T \mathbf{x}_i \geq 0$  and  $\mathbf{w}_2^T \mathbf{x}_i \geq 0$  we have

$$\mathbf{x}_i^T \mathbf{h}_{2,\perp} = \mathbf{x}_i^T \mathbf{w}_2 + (-\mathbf{w}_2^T \mathbf{a})(\mathbf{a}^T \mathbf{x}_i) \geq (-\mathbf{w}_2^T \mathbf{a})(\mathbf{a}^T \mathbf{x}_i) \geq (1 - \epsilon)(\mathbf{a}^T \mathbf{x}_i) \geq 0,$$

which implies that

$$|\mathbf{x}_i^T \mathbf{h}_{2,\perp}| \geq (1 - \epsilon)|\mathbf{a}^T \mathbf{x}_i|.$$

Thus we have

$$\mathbf{1}_{\{\mathbf{a}^T \mathbf{x}_i \geq 0\}} \mathbf{1}_{\{\mathbf{w}_2^T \mathbf{x}_i \geq 0\}} \leq \mathbf{1}_{\{(1-\epsilon)|\mathbf{x}_i^T \mathbf{a}| \leq |\mathbf{x}_i^T \mathbf{h}_{2,\perp}|\}}.$$

Plugging this inequality into (15) we conclude that

$$\frac{1}{n} \sum_{i=1}^n \mathbf{1}_{\{\mathbf{a}^T \mathbf{x}_i \geq 0\}} \phi^2(\mathbf{w}_2^T \mathbf{x}_i) \leq \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{\{(1-\epsilon)|\mathbf{x}_i^T \mathbf{a}| \leq |\mathbf{x}_i^T \mathbf{h}_{2,\perp}|\}} (\mathbf{h}_{2,\perp}^T \mathbf{x}_i)^2.$$

Use Lemma 9 with  $\mathbf{w} = \mathbf{w}_2$ , we have that

$$\sqrt{\frac{1}{n} \sum_{i=1}^n \mathbf{1}_{\{(1-\epsilon)|\mathbf{x}_i^T \mathbf{a}| \leq |\mathbf{x}_i^T \mathbf{h}_{2,\perp}|\}} (\mathbf{h}_{2,\perp}^T \mathbf{x}_i)^2} \leq \left( \delta + \sqrt{\frac{21}{20}} \epsilon \right) \|\mathbf{w}_2 + \mathbf{a}\| \leq 2\epsilon \|\mathbf{w}_2 + \mathbf{a}\|$$

holds for all  $\|\mathbf{w}_2 - \mathbf{a}\| \leq \epsilon$  with probability at least  $1 - 3e^{-\gamma n}$ . Thus we have

$$\begin{aligned} \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{\{\mathbf{a}^T \mathbf{x}_i \geq 0\}} \phi^2(\mathbf{w}_2^T \mathbf{x}_i) &\leq \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{\{(1-\epsilon)|\mathbf{x}_i^T \mathbf{a}| \leq |\mathbf{x}_i^T \mathbf{h}_{2,\perp}|\}} (\mathbf{h}_{2,\perp}^T \mathbf{x}_i)^2 \\ &\leq 4\epsilon^2 \|\mathbf{w}_2 + \mathbf{a}\|^2. \end{aligned}$$

Plugging the latter into (14) we conclude that

$$\frac{1}{n} \sum_{i=1}^n \mathbf{1}_{\{\mathbf{a}^T \mathbf{x}_i \geq 0\}} r^2(\mathbf{x}_i) \geq c \|\mathbf{h}_1\|^2 - 4\epsilon^2 \|\mathbf{h}_2\|^2$$

Similarly, we have

$$\frac{1}{n} \sum_{i=1}^n \mathbf{1}_{\{\mathbf{a}^T \mathbf{x}_i \leq 0\}} r^2(\mathbf{x}_i) \geq c \|\mathbf{h}_2\|^2 - 4\epsilon^2 \|\mathbf{h}_1\|^2$$

Summing the latter two we conclude that

$$\frac{1}{n} \sum_{i=1}^n r^2(\mathbf{x}_i) \geq (c - 4\epsilon^2) \left( \|\mathbf{h}_1\|^2 + \|\mathbf{h}_2\|^2 \right)$$

Finally, plugging in the above into (13) we conclude that

$$\langle \nabla \widehat{\mathcal{L}}(\mathbf{W}), \mathbf{W} - \mathbf{W}^* \rangle \geq \left( \frac{c}{2} - 2\epsilon^2 - \frac{9\epsilon^2}{8(1-\epsilon)^2} \right) \|\mathbf{W} - \mathbf{W}^*\|_F^2$$

Thus, for  $\epsilon$  a sufficiently small constant we have

$$\langle \nabla \widehat{\mathcal{L}}(\mathbf{W}), \mathbf{W} - \mathbf{W}^* \rangle \geq \alpha \|\mathbf{W} - \mathbf{W}^*\|_F^2$$

holds with high probability with  $\alpha = c/3$ . This concludes the proof of the correlation inequality in the empirical case.

### E.2.3 PROOF OF GRADIENT SMOOTHNESS TOWARDS THE GLOBAL OPTIMA IN THE EMPIRICAL CASE

In this section we show that the empirical gradient also obeys a smoothness bound. Concretely we prove the following Lemma.

**Lemma 10 (smoothness of empirical gradient)** *Assume  $n \geq d$ . Then,*

$$\|\nabla \widehat{\mathcal{L}}(\mathbf{W})\|_F \leq \beta \|\mathbf{W} - \mathbf{W}^*\|_F$$

with  $\beta = 4\sqrt{2}$  holds for all  $\mathbf{W}$  simultaneously with probability at least  $1 - 2e^{-\gamma n}$ .

To prove this lemma note that

$$\|\nabla \widehat{\mathcal{L}}(\mathbf{W})\|_F^2 = \|\nabla_{\mathbf{w}_1} \widehat{\mathcal{L}}\|^2 + \|\nabla_{\mathbf{w}_2} \widehat{\mathcal{L}}\|^2$$

To continue note that

$$\|\nabla_{\mathbf{w}_1} \widehat{\mathcal{L}}\| = \sup_{\mathbf{u} \in \mathbb{S}^{d-1}} \langle \nabla_{\mathbf{w}_1} \widehat{\mathcal{L}}, \mathbf{u} \rangle = \sup_{\mathbf{u} \in \mathbb{S}^{d-1}} \frac{1}{n} \sum_{i=1}^n r(\mathbf{x}_i) \phi'(\mathbf{x}_i^T \mathbf{w}_1) (\mathbf{x}_i^T \mathbf{u})$$

Now applying Cauchy-Schwarz we conclude that for  $n \geq d$  with probability at least  $1 - e^{-\gamma d}$  we have

$$\begin{aligned} \|\nabla_{\mathbf{w}_1} \widehat{\mathcal{L}}\| &\leq \sqrt{\frac{1}{n} \sum_{i=1}^n r^2(\mathbf{x}_i)} \sqrt{\sup_{\mathbf{u} \in \mathbb{S}^{d-1}} \frac{1}{n} \sum_{i=1}^n (\mathbf{x}_i^T \mathbf{u})^2} \\ &\leq \sqrt{\frac{1}{n} \sum_{i=1}^n r^2(\mathbf{x}_i)} \sqrt{\frac{2(d+n)}{n}} \\ &\leq 2 \sqrt{\frac{1}{n} \sum_{i=1}^n r^2(\mathbf{x}_i)}, \end{aligned}$$

where in the penultimate step we used the fact that the spectral norm of a Gaussian matrix is bounded by  $\sqrt{2(d+n)}$  with probability at least  $1 - 2e^{-\gamma d}$  and in the last step we used the fact that  $n \geq d$ .

Using a similar identity for  $\|\nabla_{\mathbf{w}_2} \widehat{\mathcal{L}}\|$  we thus conclude that with probability at least  $1 - 2e^{-\gamma d}$

$$\|\nabla \widehat{\mathcal{L}}(\mathbf{W})\|_F^2 \leq \frac{4}{n} \sum_{i=1}^n r^2(\mathbf{x}_i) \quad (16)$$

To proceed note that we have

$$\begin{aligned} \frac{1}{n} \sum_{i=1}^n r^2(\mathbf{x}_i) &\stackrel{(a)}{\leq} \frac{2}{n} \sum_{i=1}^n (\phi(\mathbf{w}_1^T \mathbf{x}_i) - \phi(\mathbf{a}^T \mathbf{x}_i))^2 + \frac{2}{n} \sum_{i=1}^n (\phi(\mathbf{w}_2^T \mathbf{x}_i) - \phi(-\mathbf{a}^T \mathbf{x}_i))^2 \\ &\stackrel{(b)}{\leq} \frac{2}{n} \sum_{i=1}^n (\mathbf{x}_i^T \mathbf{h}_1)^2 + \frac{2}{n} \sum_{i=1}^n (\mathbf{x}_i^T \mathbf{h}_2)^2 \\ &\stackrel{(c)}{\leq} 8(\|\mathbf{h}_1\|^2 + \|\mathbf{h}_2\|^2) \end{aligned}$$

where (a) follows from the simple identity  $(a+b)^2 \leq 2(a^2+b^2)$ , (b) from the fact that ReLU is 1-Lipschitz, and (c) from the fact that the spectral norm of a Gaussian matrix is bounded by  $\sqrt{2(d+n)}$  with probability at least  $1 - 2e^{-\gamma d}$  and  $n \geq d$ . Plugging the latter into (16) we conclude that for  $n \geq d$ ,

$$\|\nabla \widehat{\mathcal{L}}(\mathbf{W})\|_F^2 \leq 32 \|\mathbf{W} - \mathbf{W}^*\|_F^2$$

holds with probability at least  $1 - 2e^{-\gamma d}$  concluding the proof with  $\beta = 4\sqrt{2}$ .

## E.2.4 PROOF OF LOWER-BOUND LEMMA VIA UNIFORM CONCENTRATION

We prove the following lemma

**Lemma 11** *Assume  $\mathbf{x}_i$  are distributed i.i.d.  $\mathcal{N}(\mathbf{0}, \mathbf{I}_d)$ . Furthermore, assume*

$$n \geq Cd.$$

*Then assuming  $\mathbf{a}$  a unit norm vector, for all  $\mathbf{w}$  obeying  $\|\mathbf{w} - \mathbf{a}\| \leq \epsilon$  we have*

$$\frac{1}{n} \sum_{i=1}^n \mathbf{1}_{\{\mathbf{a}^T \mathbf{x}_i \geq 0\}} \mathbf{1}_{\{\mathbf{w}^T \mathbf{x}_i \geq 0\}} (\mathbf{x}_i^T (\mathbf{w} - \mathbf{a}))^2 \geq \frac{1}{100\pi} \|\mathbf{w} - \mathbf{a}\|^2.$$

*holds with probability at least  $1 - 4e^{-\gamma n}$ .*

To prove this lemma, for all  $\mathbf{w}$  obeying  $\|\mathbf{w} - \mathbf{a}\| \leq \epsilon$  we divide the region

$$\mathcal{H} = \{\mathbf{h} : \|\mathbf{h}\| \leq \epsilon\}$$

into the following two regions

$$\mathcal{H}_1 = \left\{ \mathbf{h} : \|\mathbf{h}\| \leq \epsilon \quad \text{and} \quad \frac{\mathbf{a}^T \mathbf{h}}{\|\mathbf{h}\|} \geq -\rho \right\}$$

and

$$\mathcal{H}_2 = \left\{ \mathbf{h} : \|\mathbf{h}\| \leq \epsilon \quad \text{and} \quad \frac{\mathbf{a}^T \mathbf{h}}{\|\mathbf{h}\|} \leq -\rho \right\}$$

with  $\rho = \frac{1}{\sqrt{2}}$ . To prove the lemma it suffices to prove

$$\frac{1}{n} \sum_{i=1}^n \mathbf{1}_{\{\mathbf{a}^T \mathbf{x}_i \geq 0\}} \mathbf{1}_{\{\mathbf{h}^T \mathbf{x}_i + \mathbf{a}^T \mathbf{x}_i \geq 0\}} (\mathbf{x}_i^T \mathbf{h})^2 \geq c \|\mathbf{h}\|^2.$$

when  $\mathbf{h}$  belongs to each of the regions  $\mathcal{H}_1$  and  $\mathcal{H}_2$  separately.

**Case I:  $\mathbf{h} \in \mathcal{H}_1$ :**

In this case first note that when  $\mathbf{h}^T \mathbf{x}_i \geq 0$  and  $\mathbf{a}^T \mathbf{x}_i \geq 0$  we have  $\mathbf{h}^T \mathbf{x}_i + \mathbf{a}^T \mathbf{x}_i \geq 0$ . Thus,

$$\frac{1}{n} \sum_{i=1}^n \mathbf{1}_{\{\mathbf{a}^T \mathbf{x}_i \geq 0\}} \mathbf{1}_{\{\mathbf{h}^T \mathbf{x}_i + \mathbf{a}^T \mathbf{x}_i \geq 0\}} (\mathbf{x}_i^T \mathbf{h})^2 \geq \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{\{\mathbf{a}^T \mathbf{x}_i \geq 0\}} \mathbf{1}_{\{\mathbf{h}^T \mathbf{x}_i \geq 0\}} (\mathbf{x}_i^T \mathbf{h})^2.$$

In this case we will prove that for all  $\mathbf{u} \in \mathcal{U}_1 := \{\mathbf{u} \in \mathbb{S}^{d-1} : \mathbf{a}^T \mathbf{u} \geq -\rho\}$  we have

$$\frac{1}{n} \sum_{i=1}^n \mathbf{1}_{\{\mathbf{a}^T \mathbf{x}_i \geq 0\}} \phi^2(\mathbf{x}_i^T \mathbf{u}) \geq \frac{1}{10\sqrt{\pi}} \tag{17}$$

holds with probability at least  $1 - 2e^{-\gamma n}$ . By taking  $\mathbf{u} = \frac{\mathbf{h}}{\|\mathbf{h}\|}$  this immediately implies that

$$\frac{1}{n} \sum_{i=1}^n \mathbf{1}_{\{\mathbf{a}^T \mathbf{x}_i \geq 0\}} \mathbf{1}_{\{\mathbf{h}^T \mathbf{x}_i + \mathbf{a}^T \mathbf{x}_i \geq 0\}} (\mathbf{x}_i^T \mathbf{h})^2 \geq \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{\{\mathbf{a}^T \mathbf{x}_i \geq 0\}} \mathbf{1}_{\{\mathbf{h}^T \mathbf{x}_i \geq 0\}} (\mathbf{x}_i^T \mathbf{h})^2 \geq \frac{1}{100\pi} \|\mathbf{h}\|^2,$$

holds for all  $\mathbf{h} \in \mathcal{H}_1$  with the same probability completing the proof of Case I. We thus turn our attention to proving (17). To this aim, first note that using Jensen's inequality we have

$$\frac{1}{n} \sum_{i=1}^n \mathbf{1}_{\{\mathbf{a}^T \mathbf{x}_i \geq 0\}} \phi^2(\mathbf{x}_i^T \mathbf{u}) \geq \left( \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{\{\mathbf{a}^T \mathbf{x}_i \geq 0\}} \phi(\mathbf{x}_i^T \mathbf{u}) \right)^2$$

Next, define the stochastic process

$$\mathcal{X}_i(\mathbf{u}) := \mathbf{1}_{\{\mathbf{a}^T \mathbf{x}_i \geq 0\}} \phi(\mathbf{x}_i^T \mathbf{u}) - \mathbb{E}_{\mathbf{x}_i} \left[ \mathbf{1}_{\{\mathbf{a}^T \mathbf{x}_i \geq 0\}} \phi(\mathbf{x}_i^T \mathbf{u}) \right].$$

We next prove that this stochastic process has sub-Gaussian increments. In particular using a well-known centering argument for the  $\|\cdot\|_{\Psi_2}$  (denoting the Orlicz or sub-Gaussian norm) we have

$$\begin{aligned} \|\mathcal{X}_i(\mathbf{u}) - \mathcal{X}_i(\mathbf{v})\|_{\Psi_2} &\leq 2\|\mathbf{1}_{\{\mathbf{a}^T \mathbf{x}_i \geq 0\}} (\phi(\mathbf{x}_i^T \mathbf{u}) - \phi(\mathbf{x}_i^T \mathbf{v}))\|_{\Psi_2} \\ &\leq 2\|(\phi(\mathbf{x}_i^T \mathbf{u}) - \phi(\mathbf{x}_i^T \mathbf{v}))\|_{\Psi_2} \\ &\leq 2\|\mathbf{x}_i^T (\mathbf{u} - \mathbf{v})\|_{\Psi_2} \\ &\leq c\|\mathbf{u} - \mathbf{v}\| \end{aligned}$$

where in the penultimate step we used the fact that ReLU is 1-Lipschitz. As a result using the fact that weighted sums of sub-Gaussians are also sub-Gaussian, the stochastic process

$$\begin{aligned} \mathcal{X}(\mathbf{u}) &:= \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{\{\mathbf{a}^T \mathbf{x}_i \geq 0\}} \phi(\mathbf{x}_i^T \mathbf{u}) - \mathbb{E}_x \left[ \mathbf{1}_{\{\mathbf{a}^T \mathbf{x} \geq 0\}} \phi(\mathbf{x}^T \mathbf{u}) \right] \\ &= \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{\{\mathbf{a}^T \mathbf{x}_i \geq 0\}} \phi(\mathbf{x}_i^T \mathbf{u}) - \frac{1 + \mathbf{a}^T \mathbf{u}}{\sqrt{8\pi}} \end{aligned}$$

also has sub-Gaussian increments i.e.

$$\|\mathcal{X}(\mathbf{u}) - \mathcal{X}(\mathbf{v})\|_{\Psi_2} \leq \frac{c}{\sqrt{n}} \|\mathbf{u} - \mathbf{v}\|$$

Thus using Talagrand's majorizing theorem (e.g. see (Vershynin, Exercise 8.6.5)) we have

$$\sup_{\mathbf{u} \in \mathcal{U}_1} |\mathcal{X}(\mathbf{u})| \leq c \frac{\sqrt{d}}{\sqrt{n}} + t$$

holds with probability at least  $1 - 2e^{-\gamma t^2 n}$ . Thus for all  $\mathbf{u} \in \mathcal{U}_1$  using  $t = \delta/2$  and picking  $n$  such that  $c \frac{\sqrt{d}}{\sqrt{n}} \leq \frac{\delta}{2}$  we have

$$\frac{1}{n} \sum_{i=1}^n \mathbf{1}_{\{\mathbf{a}^T \mathbf{x}_i \geq 0\}} \phi(\mathbf{x}_i^T \mathbf{u}) \geq \frac{1 - \rho}{\sqrt{8\pi}} - \delta,$$

holds with probability at least  $1 - 2e^{-\gamma \delta^2 n}$  as long as  $n \geq C \frac{d}{\delta^2}$ . Using

$$\delta = \frac{5\sqrt{2} - 7}{20\sqrt{\pi}}$$

concludes the proof of Case I.

**Case II:  $\mathbf{h} \in \mathcal{H}_2$ :**

In this case first note that

$$\begin{aligned} &\frac{1}{n} \sum_{i=1}^n \mathbf{1}_{\{\mathbf{a}^T \mathbf{x}_i \geq 0\}} \mathbf{1}_{\{\mathbf{h}^T \mathbf{x}_i + \mathbf{a}^T \mathbf{x}_i \geq 0\}} (\mathbf{x}_i^T \mathbf{h})^2 \\ &= \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{\{\mathbf{a}^T \mathbf{x}_i \geq 0\}} \mathbf{1}_{\{\mathbf{h}^T \mathbf{x}_i \geq -\mathbf{a}^T \mathbf{x}_i\}} (\mathbf{x}_i^T \mathbf{h})^2 \\ &= \left( \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{\{\mathbf{a}^T \mathbf{x}_i \geq 0\}} \mathbf{1}_{\{\mathbf{h}^T \mathbf{x}_i \geq -\mathbf{a}^T \mathbf{x}_i\}} \left( \mathbf{x}_i^T \frac{\mathbf{h}}{\|\mathbf{h}\|} \right)^2 \right) \|\mathbf{h}\|^2 \\ &\geq \left( \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{\{\mathbf{a}^T \mathbf{x}_i \geq 0\}} \mathbf{1}_{\{\mathbf{x}_i^T \frac{\mathbf{h}}{\|\mathbf{h}\|} \geq -\frac{\mathbf{a}^T \mathbf{x}_i}{\epsilon}\}} \left( \mathbf{x}_i^T \frac{\mathbf{h}}{\|\mathbf{h}\|} \right)^2 \right) \|\mathbf{h}\|^2 \\ &\geq \left( \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{\{\mathbf{a}^T \mathbf{x}_i \geq 0\}} \mathbf{1}_{\{0 \geq \mathbf{x}_i^T \frac{\mathbf{h}}{\|\mathbf{h}\|} \geq -\frac{\mathbf{a}^T \mathbf{x}_i}{\epsilon}\}} \left( \mathbf{x}_i^T \frac{\mathbf{h}}{\|\mathbf{h}\|} \right)^2 \right) \|\mathbf{h}\|^2 \end{aligned}$$

where in the penultimate line we used the fact that  $\|\mathbf{h}\| \leq \epsilon$  and in the last line we added the indicator of  $\mathbf{x}_i^T \mathbf{h} \leq 0$ . Therefore, to complete the proof of this part it suffices to show that

$$\frac{1}{n} \sum_{i=1}^n \mathbf{1}_{\{\mathbf{a}^T \mathbf{x}_i \geq 0\}} \mathbf{1}_{\{0 \leq \mathbf{x}_i^T \mathbf{u} \leq -\frac{\mathbf{a}^T \mathbf{x}_i}{\epsilon}\}} (\mathbf{x}_i^T \mathbf{u})^2 \geq \frac{1}{100\pi}$$

holds with high probability for all  $\mathbf{u} \in \mathcal{U}_2 := \{\mathbf{u} \in \mathbb{S}^{d-1} : \mathbf{u}^T \mathbf{a} \leq -\rho\}$ . Or equivalently by flipping the sign of  $\mathbf{u}$

$$\frac{1}{n} \sum_{i=1}^n \mathbf{1}_{\{\mathbf{a}^T \mathbf{x}_i \geq 0\}} \mathbf{1}_{\{0 \leq \mathbf{x}_i^T \mathbf{u} \leq \frac{\mathbf{a}^T \mathbf{x}_i}{\epsilon}\}} (\mathbf{x}_i^T \mathbf{u})^2 \geq \frac{1}{100\pi}$$

holds with high probability for all  $\mathbf{u} \in \mathcal{U}_2 := \{\mathbf{u} \in \mathbb{S}^{d-1} : \mathbf{u}^T \mathbf{a} \geq \rho\}$ .

To do this we once again applying Jensen's inequality

$$\begin{aligned} \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{\{\mathbf{a}^T \mathbf{x}_i \geq 0\}} \mathbf{1}_{\{0 \leq \mathbf{x}_i^T \mathbf{u} \leq \frac{\mathbf{a}^T \mathbf{x}_i}{\epsilon}\}} (\mathbf{x}_i^T \mathbf{u})^2 &\geq \left( \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{\{\mathbf{a}^T \mathbf{x}_i \geq 0\}} \mathbf{1}_{\{0 \leq \mathbf{x}_i^T \mathbf{u} \leq \frac{\mathbf{a}^T \mathbf{x}_i}{\epsilon}\}} \mathbf{x}_i^T \mathbf{u} \right)^2 \\ &\geq \left( \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{\{\mathbf{a}^T \mathbf{x}_i \geq 0\}} \mathcal{S} \left( \mathbf{x}_i^T \mathbf{u}, \frac{\mathbf{x}_i^T \mathbf{a}}{\epsilon} \right) \right)^2 \end{aligned}$$

where

$$S(v; w) := \begin{cases} 0, & v < 0, \\ v, & 0 \leq v \leq \frac{w}{2}, \\ w - v, & \frac{w}{2} \leq v \leq w, \\ 0, & v \geq w. \end{cases}$$

Next, define the stochastic process

$$\mathcal{X}_i(\mathbf{u}) := \mathbf{1}_{\{\mathbf{a}^T \mathbf{x}_i \geq 0\}} \mathcal{S} \left( \mathbf{x}_i^T \mathbf{u}, \frac{\mathbf{x}_i^T \mathbf{a}}{\epsilon} \right) - \mathbb{E}_{\mathbf{x}_i} \left[ \mathbf{1}_{\{\mathbf{a}^T \mathbf{x}_i \geq 0\}} \mathcal{S} \left( \mathbf{x}_i^T \mathbf{u}, \frac{\mathbf{x}_i^T \mathbf{a}}{\epsilon} \right) \right].$$

We next prove that this stochastic process has sub-Gaussian increments. In particular using a well-known centering argument for the  $\|\cdot\|_{\Psi_2}$  (denoting the Orlicz or sub-Gaussian norm) we have

$$\begin{aligned} \|\mathcal{X}_i(\mathbf{u}) - \mathcal{X}_i(\mathbf{v})\|_{\Psi_2} &\leq 2 \mathbf{1}_{\{\mathbf{a}^T \mathbf{x}_i \geq 0\}} \left\| \mathcal{S} \left( \mathbf{x}_i^T \mathbf{u}, \frac{\mathbf{x}_i^T \mathbf{a}}{\epsilon} \right) - \mathcal{S} \left( \mathbf{x}_i^T \mathbf{v}, \frac{\mathbf{x}_i^T \mathbf{a}}{\epsilon} \right) \right\|_{\Psi_2} \\ &\leq 2 \left\| \mathcal{S} \left( \mathbf{x}_i^T \mathbf{u}, \frac{\mathbf{x}_i^T \mathbf{a}}{\epsilon} \right) - \mathcal{S} \left( \mathbf{x}_i^T \mathbf{v}, \frac{\mathbf{x}_i^T \mathbf{a}}{\epsilon} \right) \right\|_{\Psi_2} \\ &\leq 2 \|\mathbf{x}_i^T (\mathbf{u} - \mathbf{v})\|_{\Psi_2} \\ &\leq c \|\mathbf{u} - \mathbf{v}\| \end{aligned}$$

where in the penultimate step we used the fact that  $S(v; w)$  is 1-Lipschitz in its first argument. As a result using the fact that weighted sums of sub-Gaussians are also sub-Gaussian, the stochastic process

$$\begin{aligned} \mathcal{X}(\mathbf{u}) &:= \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{\{\mathbf{a}^T \mathbf{x}_i \geq 0\}} \mathcal{S} \left( \mathbf{x}_i^T \mathbf{u}, \frac{\mathbf{x}_i^T \mathbf{a}}{\epsilon} \right) - \mathbb{E}_{\mathbf{x}} \left[ \mathbf{1}_{\{\mathbf{a}^T \mathbf{x} \geq 0\}} \mathcal{S} \left( \mathbf{x}_i^T \mathbf{u}, \frac{\mathbf{x}_i^T \mathbf{a}}{\epsilon} \right) \right] \\ &= \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{\{\mathbf{a}^T \mathbf{x}_i \geq 0\}} \mathcal{S} \left( \mathbf{x}_i^T \mathbf{u}, \frac{\mathbf{x}_i^T \mathbf{a}}{\epsilon} \right) - f(\mathbf{a}^T \mathbf{u}, \epsilon) \end{aligned}$$

where  $f(\mathbf{a}^T \mathbf{u}, \epsilon) := \mathbb{E}_{\mathbf{x}} \left[ \mathcal{S} \left( \mathbf{x}_i^T \mathbf{u}, \frac{\mathbf{x}_i^T \mathbf{a}}{\epsilon} \right) \right]$ , also has sub-Gaussian increments i.e.

$$\|\mathcal{X}(\mathbf{u}) - \mathcal{X}(\mathbf{v})\|_{\Psi_2} \leq \frac{c}{\sqrt{n}} \|\mathbf{u} - \mathbf{v}\|$$

Thus using Talagrand’s majorizing theorem (e.g. see (Vershynin, Exercise 8.6.5)) we have

$$\sup_{\mathbf{u} \in \mathcal{U}_2} |\mathcal{X}(\mathbf{u})| \leq c \frac{\sqrt{d}}{\sqrt{n}} + t$$

holds with probability at least  $1 - 2e^{-\gamma t^2 n}$ . Thus for all  $\mathbf{u} \in \mathcal{U}_2$  using  $t = \delta/2$  and picking  $n$  such that  $c \frac{\sqrt{d}}{\sqrt{n}} \leq \frac{\delta}{2}$  we have

$$\frac{1}{n} \sum_{i=1}^n \mathbf{1}_{\{\mathbf{a}^T \mathbf{x}_i \geq 0\}} \mathcal{S} \left( \mathbf{x}_i^T \mathbf{u}, \frac{\mathbf{x}_i^T \mathbf{a}}{\epsilon} \right) \geq f(\mathbf{a}^T \mathbf{u}, \epsilon) - \delta,$$

holds with probability at least  $1 - 2e^{-\gamma \delta^2 n}$  as long as  $n \geq C \frac{d}{\delta^2}$ . It is easy to check that  $f(\mathbf{a}^T \mathbf{u}, \epsilon)$  is non-decreasing in its first argument so that for  $\mathbf{u} \in \mathcal{U}_2$  we have  $f(\mathbf{a}^T \mathbf{u}, \epsilon) \geq f(\rho, \epsilon)$ . Furthermore, for  $\rho = \frac{1}{\sqrt{2}}$  we have

$$f(\rho, \epsilon) \geq \frac{1}{8\sqrt{\pi}}$$

so that we can conclude that

$$\frac{1}{n} \sum_{i=1}^n \mathbf{1}_{\{\mathbf{a}^T \mathbf{x}_i \geq 0\}} \mathcal{S} \left( \mathbf{x}_i^T \mathbf{u}, \frac{\mathbf{x}_i^T \mathbf{a}}{\epsilon} \right) \geq \frac{1}{8\sqrt{\pi}} - \delta,$$

holds with probability at least  $1 - 2e^{-\gamma \delta^2 n}$  as long as  $n \geq C \frac{d}{\delta^2}$ . Picking  $\delta = \frac{1}{40\sqrt{\pi}}$  concludes the proof of Case II.

## F PROOF OF MAIN THEOREMS

### F.1 PROOF OF THEOREM 1 FOR ONE DIMENSIONAL OUTPUTS

Using the update rule  $\mathbf{W}^{(\tau+1)} = \mathbf{W}^{(\tau)} - \mu \nabla_{\mathbf{W}} \mathcal{L}(\mathbf{W}^{(\tau)})$ ,

$$\begin{aligned} \left\| \mathbf{W}^{(\tau+1)} - \mathbf{W}^* \right\|_F^2 &= \left\| \mathbf{W}^{(\tau)} - \mu \nabla_{\mathbf{W}} \mathcal{L}(\mathbf{W}^{(\tau)}) - \mathbf{W}^* \right\|_F^2 \\ &= \left\| \mathbf{W}^{(\tau)} - \mathbf{W}^* \right\|_F^2 - 2\mu \left\langle \mathbf{W}^{(\tau)} - \mathbf{W}^*, \nabla_{\mathbf{W}} \mathcal{L}(\mathbf{W}^{(\tau)}) \right\rangle + \mu^2 \left\| \nabla_{\mathbf{W}} \mathcal{L}(\mathbf{W}^{(\tau)}) \right\|_F^2 \\ &\stackrel{(a)}{\leq} \left\| \mathbf{W}^{(\tau)} - \mathbf{W}^* \right\|_F^2 - 2\mu\alpha \left\| \mathbf{W}^{(\tau)} - \mathbf{W}^* \right\|_F^2 + \mu^2 \left\| \nabla_{\mathbf{W}} \mathcal{L}(\mathbf{W}^{(\tau)}) \right\|_F^2 \\ &\stackrel{(b)}{\leq} \left\| \mathbf{W}^{(\tau)} - \mathbf{W}^* \right\|_F^2 - 2\mu\alpha \left\| \mathbf{W}^{(\tau)} - \mathbf{W}^* \right\|_F^2 + \mu^2 \beta^2 \left\| \mathbf{W}^{(\tau)} - \mathbf{W}^* \right\|_F^2 \\ &= (1 - 2\mu\alpha + \mu^2 \beta^2) \left\| \mathbf{W}^{(\tau)} - \mathbf{W}^* \right\|_F^2, \end{aligned}$$

where (a) follows from Lemma 4, and (b) follows from Lemma 5. Picking  $\mu \leq \frac{\alpha}{\beta^2}$  yields

$$\left\| \mathbf{W}^{(\tau+1)} - \mathbf{W}^* \right\|_F^2 \leq (1 - \mu\alpha) \left\| \mathbf{W}^{(\tau)} - \mathbf{W}^* \right\|_F^2.$$

Repeating the steps above  $\tau$  times completes the proof. The constants in the theorem statement are  $c_1 = \frac{\alpha}{\beta^2}$  and  $c_2 = \alpha$ .

### F.2 PROOF OF THEOREM 6 FOR MULTIDIMENSIONAL OUTPUTS

Our proof for multiple outputs can be derived by repeated application of the single output case. To see this we remind the reader that the loss takes the form

$$\mathcal{L}(\mathbf{W}) = \frac{1}{2} \mathbb{E}_{\mathbf{x}} \left[ \left\| \mathbf{V}^T \text{ReLU}(\mathbf{W}\mathbf{x}) - \mathbf{A}\mathbf{x} \right\|^2 \right].$$

1782 Plugging in the particular pattern for  $\mathbf{V}$  this is equal to

$$\begin{aligned}
1783 \mathcal{L}(\mathbf{W}) &= \frac{1}{2} \mathbb{E}_{\mathbf{x}} \left[ \left\| \tilde{\mathbf{V}}^T [\mathbf{I}_r \quad -\mathbf{I}_r] \text{ReLU}(\mathbf{W}\mathbf{x}) - \mathbf{A}\mathbf{x} \right\|^2 \right] \\
1784 &= \frac{1}{2} \mathbb{E}_{\mathbf{x}} \left[ \left\| \mathbf{R}^T \Sigma [\mathbf{I}_r \quad -\mathbf{I}_r] \text{ReLU}(\mathbf{W}\mathbf{x}) - \mathbf{A}\mathbf{x} \right\|^2 \right] \\
1785 &= \frac{1}{2} \mathbb{E}_{\mathbf{x}} \left[ \left\| \Sigma [\mathbf{I}_r \quad -\mathbf{I}_r] \text{ReLU}(\mathbf{W}\mathbf{x}) - \mathbf{R}\mathbf{A}\mathbf{x} \right\|^2 \right] \\
1786 &= \frac{1}{2} \sum_{\ell=1}^r \sigma_{\ell}^2 \mathbb{E}_{\mathbf{x}} \left[ \left( \text{ReLU}(\mathbf{w}_{\ell}^T \mathbf{x}) - \text{ReLU}(\mathbf{w}_{\ell+r}^T \mathbf{x}) - \tilde{\mathbf{a}}_{\ell}^T \mathbf{x} \right)^2 \right]
\end{aligned}$$

1793 where  $\tilde{\mathbf{a}}_{\ell}$  is the  $\ell$ th row of  $\Sigma^{-1} \mathbf{R}\mathbf{A}$  and  $\sigma_{\ell} = \Sigma_{\ell\ell}$ . Note that the loss decomposes into  $r$  optimization  
1794 problems of the form in the Theorem 1. Thus the result follows from applying Theorem 1 to the  
1795 summands of the above loss.

### 1797 F.3 PROOF OF THEOREM 2 FOR THE EMPIRICAL SETTING

1798 To prove this theorem first we note that after the first iteration using Lemma 7 from Section E.2.1 we  
1799 have that with high probability.

$$1802 \|\mathbf{W}^{(1)} - \mathbf{W}^*\|_F \leq \epsilon \|\mathbf{a}\|$$

1803 We show inductively below that assuming  $\|\mathbf{W}^{(\tau)} - \mathbf{W}^*\|_F \leq \epsilon \|\mathbf{a}\|$  the next iteration also obeys  
1804 this inequality staying in the local neighborhood of the global optima. In this local neighborhood as  
1805 shown in Sections E.2.2 and E.2.3 we have the correlation inequality

$$1806 \left\langle \mathbf{W} - \mathbf{W}^*, \nabla_{\mathbf{W}} \hat{\mathcal{L}}(\mathbf{W}) \right\rangle \geq \alpha \|\mathbf{W} - \mathbf{W}^*\|_F^2 \quad (18)$$

1807 and the smoothness bound

$$1808 \|\nabla \hat{\mathcal{L}}(\mathbf{W})\|_F \leq \beta \|\mathbf{W} - \mathbf{W}^*\|_F \quad (19)$$

1809 holds with high probability simultaneously for all  $\mathbf{W}$  obeying  $\|\mathbf{W}^{(1)} - \mathbf{W}^*\|_F \leq \epsilon \|\mathbf{a}\|$  with  $\alpha$   
1810 and  $\beta$  fixed numerical constants. We would like to emphasize that the proof (18) in the empirical  
1811 setting is quite intricate necessitating clever algebraic manipulations combined with sophisticated  
1812 empirical processing theory tools. With these identities in place using the update rule  $\mathbf{W}^{(\tau+1)} =$   
1813  $\mathbf{W}^{(\tau)} - \mu \nabla_{\mathbf{W}} \hat{\mathcal{L}}(\mathbf{W}^{(\tau)})$ , we have

$$\begin{aligned}
1814 \|\mathbf{W}^{(\tau+1)} - \mathbf{W}^*\|_F^2 &= \left\| \mathbf{W}^{(\tau)} - \mu \nabla_{\mathbf{W}} \hat{\mathcal{L}}(\mathbf{W}^{(\tau)}) - \mathbf{W}^* \right\|_F^2 \\
1815 &= \left\| \mathbf{W}^{(\tau)} - \mathbf{W}^* \right\|_F^2 - 2\mu \left\langle \mathbf{W}^{(\tau)} - \mathbf{W}^*, \nabla_{\mathbf{W}} \hat{\mathcal{L}}(\mathbf{W}^{(\tau)}) \right\rangle + \mu^2 \left\| \nabla_{\mathbf{W}} \hat{\mathcal{L}}(\mathbf{W}^{(\tau)}) \right\|_F^2 \\
1816 &\stackrel{(a)}{\leq} \left\| \mathbf{W}^{(\tau)} - \mathbf{W}^* \right\|_F^2 - 2\mu\alpha \left\| \mathbf{W}^{(\tau)} - \mathbf{W}^* \right\|_F^2 + \mu^2 \left\| \nabla_{\mathbf{W}} \hat{\mathcal{L}}(\mathbf{W}^{(\tau)}) \right\|_F^2 \\
1817 &\stackrel{(b)}{\leq} \left\| \mathbf{W}^{(\tau)} - \mathbf{W}^* \right\|_F^2 - 2\mu\alpha \left\| \mathbf{W}^{(\tau)} - \mathbf{W}^* \right\|_F^2 + \mu^2 \beta^2 \left\| \mathbf{W}^{(\tau)} - \mathbf{W}^* \right\|_F^2 \\
1818 &= (1 - 2\mu\alpha + \mu^2 \beta^2) \left\| \mathbf{W}^{(\tau)} - \mathbf{W}^* \right\|_F^2,
\end{aligned}$$

1819 where (a) follows from inequality 18, and (b) follows from inequality 19. Picking  $\mu \leq \frac{\alpha}{\beta^2}$  yields

$$1820 \left\| \mathbf{W}^{(\tau+1)} - \mathbf{W}^* \right\|_F^2 \leq (1 - \mu\alpha) \left\| \mathbf{W}^{(\tau)} - \mathbf{W}^* \right\|_F^2.$$

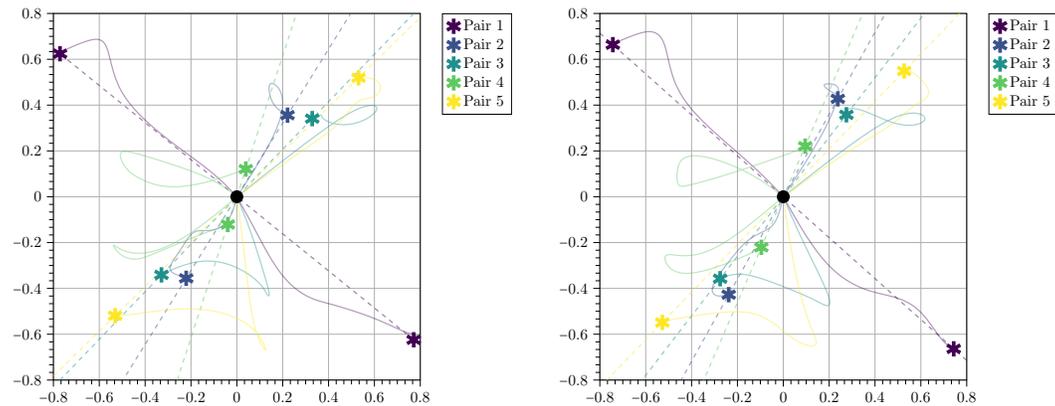
1821 Repeating the steps above  $\tau$  times completes the proof. The constants in the theorem statement are  
1822  $c_5 = \frac{\alpha}{\beta^2}$  and  $c_8 = \alpha$ .

## G ADDITIONAL EXPERIMENTAL RESULTS

### G.1 PAIRING-UP BEHAVIOR FOR $r > 3$

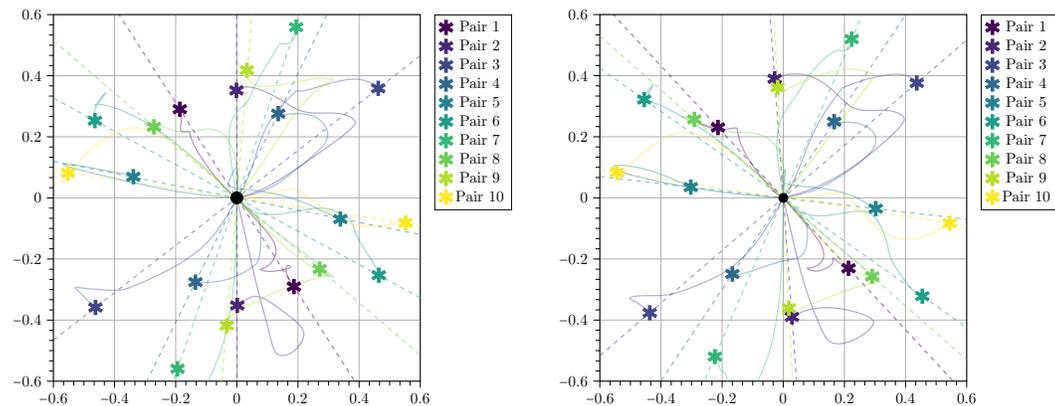
In this section, we present additional results on the pairing behavior of  $w_i$  and  $v_i$  for different values of  $r$ .

Beyond the  $r = 3$  case shown in Section 4.2, we illustrate the same behavior for  $r = 5$  in Figure 7 and for  $r = 10$  in Figure 8. While we also observe the pairing for  $r > 10$ , we omit those results here for visual clarity. In general, we note that the weights at convergence (indicated with *star* symbol in Figures 7 and 8) can be grouped into  $r$  pairs such that one of the weights is approximately negative of the other. To aid with detecting the pairs visually, we draw the line determined by each pair with dashed lines.



(a) Trajectory of  $v_i$ 's and their pairing behavior. (b) Trajectory of  $w_i$ 's and their pairing behavior.

Figure 7: **Pairing pattern for  $r = 5$ .** We train the network from small initialization when exactly parameterized ( $k = 10$  and  $r = 5$ ). On left (a), we depict the trajectories of individual weights in the outer layer ( $v_i$ 's) across iterations. Each pair is indicated by the same color and the dashed line. A similar pairing is observed for the inner layer weights as well (b). While these vectors all lie in a higher dimensional space, we pick an arbitrary two dimensional axis to plot them in 2D.



(a) Trajectory of  $v_i$ 's and their pairing behavior. (b) Trajectory of  $w_i$ 's and their pairing behavior.

Figure 8: **Pairing pattern for  $r = 10$ .** We train the network from small initialization when exactly parameterized ( $k = 20$  and  $r = 10$ ). On left (a), we depict the trajectories of individual weights in the outer layer ( $v_i$ 's) across iterations. Each pair is indicated by the same color and the dashed line. A similar pairing is observed for the inner layer weights as well (b). While these vectors all lie in a higher dimensional space, we pick an arbitrary two dimensional axis to plot them in 2D.

## G.2 LAZY VS. RICH REGIME COMPARISON

In this section, we compare training and generalization dynamics of, so called, "lazy" and "rich" regimes. In our experiments, we fix the learning rate and initialization scale. Therefore, the free parameter that controls the lazy vs. rich learning regime manifests itself through the model width. More specifically, we consider input dimension  $d = 100$ , fix the learning rate as  $\mu = 0.001$  and the initialization scale as  $\alpha = 1.0$  (default PyTorch initialization). We train exactly parametrized ( $k = 2$ ), mildly over-parametrized ( $k = 6$ ) and heavily over-parametrized ( $k = 100$ ) models. Using  $n = 1000$  data samples, we perform GD iterations until the training loss is below  $10^{-4}$ . The trajectory of neurons for all cases is given in Figure 9.

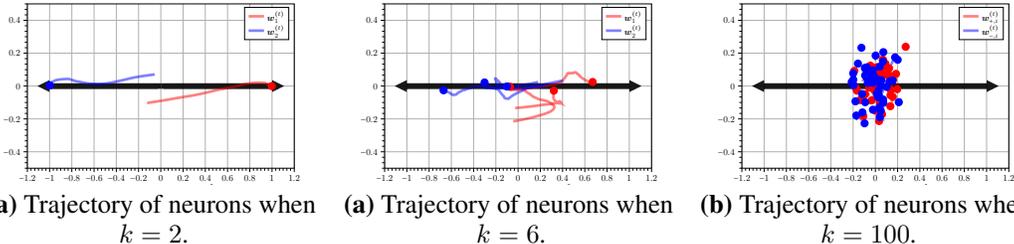


Figure 9: **Trajectory of neurons for different values of  $k$ .** We run gradient descent updates on the empirical loss after fixing  $v$  with half  $+1$ 's and half  $-1$ 's. Black arrows indicate  $\pm a$  direction. A randomly selected orthogonal direction to  $a$  is shown in y-axis in order to visualize the neurons in 2D. We use colors red and blue to indicate whether  $v_i$  corresponding to  $w_i$  is 1 or  $-1$  respectively. Points at the end of each trajectory denotes the final weight GD converges to.

We observe two completely different behaviors. The exactly parametrized model recovers the target directions using  $n \approx cd$  data samples as predicted by our theory. Similarly, the mildly over-parametrized model approximately aligns with the target direction. On the other hand, the heavily over-parametrized model has very little neuron movement from its initialization (as predicted by NTK theory). Consequently, the neurons are not necessarily aligned with the corresponding target directions. We measure the generalization performance by calculating the population loss given the final model weights. The exactly parametrized ( $k = 2$ ) and mildly over-parametrized ( $k = 6$ ) model achieves  $\approx 2 \times 10^{-4}$  and  $\approx 11 \times 10^{-4}$  test loss respectively that is in the same order as the training loss. The Heavily over-parametrized model ( $k = 100$ ) on the other hand, does not generalize, achieving  $\approx 4$  test loss. This experiment demonstrates that the existing theory of analyzing heavily over-parametrized models (such as Jacot et al. (2018); Oymak and Soltanolkotabi (2019; 2020); Du et al. (2019); Arora et al. (2019)) cannot be used to explain the generalization phenomenon present in exactly parametrized and mildly over-parametrized models.

## H LIMITATIONS

While our work provides the first comprehensive analysis of the gradient descent dynamics for learning linear target functions with ReLU networks, it does come with a few limitations. First, our analysis is restricted to shallow networks with a single hidden layer. Extending these results to deeper architectures remains an important direction for future research. Second, our theoretical guarantees rely on the assumption that input data are drawn i.i.d from a Gaussian distribution. Although our techniques can potentially be adapted to broader classes of distributions, the i.i.d assumption is central to the current analysis. That said, we believe this work serves as a stepping stone toward understanding more complex data models, architectures, and training dynamics.