

When Both Layers Learn: Training Dynamics of Representing Linear Models via ReLU Networks

Berk Tinaz

Changzhi Xie

Mahdi Soltanolkotabi

University of Southern California, Los Angeles, CA.

TINAZ@USC.EDU

CHANGZHI@USC.EDU

SOLTANOL@USC.EDU

Editors: Steve Hanneke and Tor Lattimore

Abstract

In this paper, we study the gradient descent dynamics for jointly training *both* layers of a one-hidden-layer ReLU network to fit a linear target function. Concretely, we consider a realizable setting where inputs are drawn i.i.d. from a Gaussian distribution and labels follow a planted linear model. This stylized framework captures salient features of end-to-end training in inverse problems and certain auto-encoder models. Despite its apparent simplicity, the dynamics remain poorly understood, in part because the loss landscape contains multiple non-strict saddle points, making it unclear why gradient descent from random initialization reliably escapes bad stationary regions. We provide a detailed characterization of the optimization landscape and prove that gradient descent from a moderately small random initialization-*simultaneously training both layers*-converges to a global minimizer at a linear rate with order-wise optimal sample complexity. Our analysis tracks the trajectory through three phases: an *alignment phase* in which hidden weights progressively align with the planted direction while the output weights maintain the correct sign pattern; a *growth phase* in which the norms of both layers increase while preserving alignment; and a *local refinement phase* in which the aligned neurons rapidly converge to the planted direction, yielding fast local convergence. To rigorously show that GD avoids non-strict saddles, we develop trajectory-level control arguments for the end-to-end dynamics. In addition, we establish novel uniform concentration results that hold along the entire trajectory, and are essential for obtaining order-wise optimal sample complexity. We corroborate our theory with extensive experiments across a range of configurations.

Keywords: ReLU networks, gradient descent, learning theory, linear functions

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 that 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 such a model are 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?

Understanding this question requires reasoning not only about the final solution reached by GD, but about the entire trajectory of the optimization process. Recent empirical work suggests that several phenomena observed during neural network training, including grokking (or delayed generalization) ([Power et al., 2022](#)), are closely tied to the temporal evolution of gradient descent. In such settings, models may fit the training data well before exhibiting improved generalization, indicating that learning can unfold through distinct stages over the course of optimization. This perspective motivates a careful, trajectory-level analysis even in simple problem settings.

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](#))¹, 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 or genuinely training both layers ([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 2](#) 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:

- We present one of the first works that analyzes training dynamics of learning *both* layers in a one-hidden-layer ReLU network in a practical regime. That is, we do not use pre-processing or alter the early optimization trajectory to avoid complications that arise from non-linear training dynamics of optimizing both layers.

1. Due to space constraints, we refer the reader to [Appendix F](#) for in-depth discussion on related work.

- We develop a 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*, while the outer layer maintains the correct sign pattern.
- As detailed further in Section 2 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 stationary points, we develop new techniques to control the GD trajectory which we combine with intricate uniform concentration bounds. In particular, our refined analysis tracks the trajectory through three phases (*alignment*, *growth*, and *local refinement* phases). We believe our refined trajectory analysis may have broader implications for the analysis of non-convex optimization problems involving non-strict saddles.
- Since gradient descent repeatedly reuses the same finite dataset across all phases, the iterates become statistically dependent on the samples. We address this by proving new uniform concentration bounds for the gradient along the entire optimization trajectory, holding simultaneously for all iterates encountered by GD. A key component of our uniform concentration result is that the accuracy of the concentration increases as we get closer and closer to the global optima. These refined bounds are a key technical ingredient for achieving order-optimal sample complexity.
- 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 $y_i \in \mathbb{R}$. As mentioned before, we consider the class of linear models where the relationship between \mathbf{x}_i and y_i is given by the equation: $y_i = \mathbf{a}^T \mathbf{x}_i$ where $\mathbf{a} \in \mathbb{R}^d$ is the labeling vector. Conceptually, \mathbf{a} is the target *direction* that our predictor should *learn*. 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 $f(\mathbf{v}, \mathbf{W}, \mathbf{x}) := \mathbf{v}^T \phi(\mathbf{W}\mathbf{x})$ as our predictor. Here k denotes the number of hidden-units, $\mathbf{v} \in \mathbb{R}^k$ 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 $\phi(\mathbf{z})$ is the activation function. We refer to individual rows of \mathbf{v}/\mathbf{W} as v_i/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. Furthermore, we focus on the exact parametrized setting, i.e. $k = 2$, as a step towards understanding the behavior of the over-specified/parameterized setting with $k > 2$ neurons.

Training Loss – We minimize the squared loss between the target and the prediction

$$\widehat{\mathcal{L}}(\mathbf{v}, \mathbf{W}) = \frac{1}{2n} \sum_{i=1}^n (\mathbf{v}^T \phi(\mathbf{W} \mathbf{x}_i) - y_i)^2 \quad (1)$$

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

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

2. Landscape Analysis: Why is learning linear functions with ReLUs challenging?

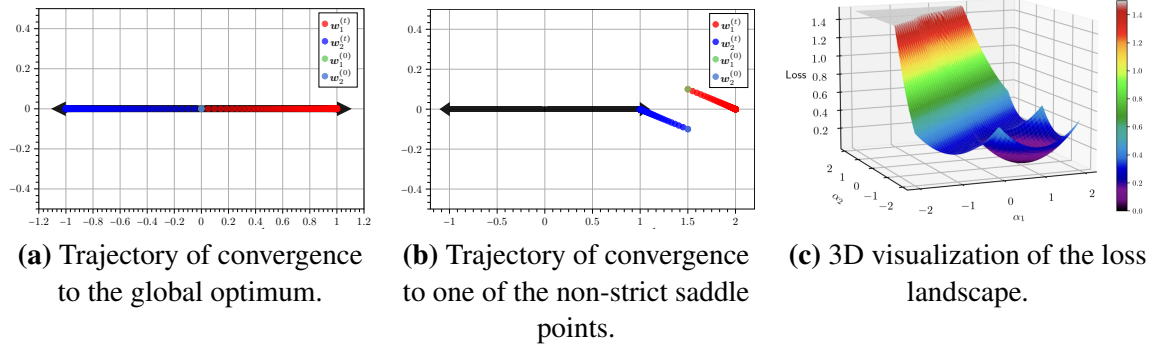


Figure 1: GD trajectories and population loss landscape. We run gradient descent on the population loss for a one-hidden-layer ReLU network with two hidden units and fixed output weights $\mathbf{v} = [1, -1]^T$. Panels (a) and (b) use two different initializations of $\mathbf{w}_1^{(0)}$ and $\mathbf{w}_2^{(0)}$. To visualize the dynamics in 2D, we plot each neuron in the plane spanned by \mathbf{a} (black arrows indicate $\pm \mathbf{a}$) and a randomly chosen direction orthogonal to \mathbf{a} (y-axis). In (a), a small initialization near the origin converges to the global optimum. In (b), initializing near $1.5\mathbf{a}$ leads GD to $(\mathbf{w}_1, \mathbf{w}_2) = (\mathbf{a}, 2\mathbf{a})$, a non-strict saddle of (2). Finally, (c) plots the population landscape in the reduced slice $\mathbf{w}_1 = \alpha_1 \mathbf{a}$, $\mathbf{w}_2 = \alpha_2 \mathbf{a}$.

Despite the simplicity of the target function, the gradient descent dynamics in this setting are surprisingly subtle. The difficulty is that the loss landscape is riddled with non-strict saddle points. Indeed, infinitely many of them—creating large flat directions where naive intuition about descent can fail. The next theorem makes this phenomenon precise for the population loss.

Theorem 1 (Landscape Characterization) *For $v_1, v_2 > 0$, the stationary points of the population loss (2) are either*

1. global optima: $v_1 \mathbf{w}_1 = \mathbf{a}, \quad v_2 \mathbf{w}_2 = -\mathbf{a}$,
2. or non-strict saddles: $v_1 \mathbf{w}_1 = (c + 1) \mathbf{a}, \quad v_2 \mathbf{w}_2 = c \mathbf{a}, \quad \text{where } c > 0 \text{ or } c < -1$.

We provide the proof of Theorem 1 in Appendix D.1.

Theorem 1 above shows that, beyond the global minima, the population loss contains a continuum of stationary points forming non-strict saddle manifolds parameterized by c . In particular, for every $c > 0$ and every $c < -1$, the equations $v_1 \mathbf{w}_1 = (c + 1)\mathbf{a}$ and $v_2 \mathbf{w}_2 = c\mathbf{a}$ define a stationary point with flat directions in the loss. Thus the landscape is highly degenerate: instead of isolated critical points, there are infinitely many saddle regions that gradient descent can enter and move along without encountering negative curvature. This proliferation of flat saddles is the primary geometric obstruction to analyzing the global behavior of gradient descent.

In Figure 1, we illustrate how the initialization determines whether GD converges to a global optimum or drifts toward a non-strict saddle. Figure 1(a) shows a trajectory that converging to the global optimum, while Figure 1(b) shows a trajectory that stalls near a saddle. To further visualize the landscape, Figure 1(c) fixes $v_1 = v_2 = 1$ and plots the loss in the reduced two-dimensional slice $\mathbf{w}_1 = \alpha_1 \mathbf{a}$, $\mathbf{w}_2 = \alpha_2 \mathbf{a}$. In this slice, the gradient vanishes along the $\alpha_1 - \alpha_2 = 1$ valley, even though the loss remains strictly positive.

3. Main Result: Convergence of the Gradient Descent Trajectory

We now present our main result, which characterizes the training dynamics when *both* layers of a ReLU network are trained in the practical empirical regime.

Theorem 2 (Convergence of GD Trajectory) *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 outputs are generated according to a linear target function of the form $y_i = \mathbf{a}^T \mathbf{x}_i$, 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*

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

by minimizing the empirical loss

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

over $\mathbf{v} = [v_1, -v_2]^T$ and $\mathbf{W} = [\mathbf{w}_1, \mathbf{w}_2]^T \in \mathbb{R}^{2 \times d}$ using gradient descent with step size $\mu := \frac{\bar{\mu}}{\|\mathbf{a}\|}$ with $\bar{\mu} \leq \frac{\mu_0}{\ln\left(\frac{\sqrt{\|\mathbf{a}\|}}{\sigma}\right)}$:

$$\mathbf{W}^{(\tau+1)} = \mathbf{W}^{(\tau)} - \mu \nabla_{\mathbf{W}} \widehat{\mathcal{L}}(\mathbf{v}^{(\tau)}, \mathbf{W}^{(\tau)}), \quad \mathbf{v}^{(\tau+1)} = \mathbf{v}^{(\tau)} - \mu \nabla_{\mathbf{v}} \widehat{\mathcal{L}}(\mathbf{v}^{(\tau)}, \mathbf{W}^{(\tau)})$$

Assume the initialization

$$\mathbf{w}_1^{(0)}, \mathbf{w}_2^{(0)} \sim \mathcal{N}\left(0, \frac{\sigma^2}{d} \mathbf{I}_d\right), \quad v_1^{(0)}, v_2^{(0)} \sim \frac{\sigma}{\sqrt{d}} \xi, \quad \xi^2 \sim \chi_d^2,$$

with $\sigma \leq \sigma_0 \sqrt{\|\mathbf{a}\|}$, χ_d^2 a chi-squared distribution with d degrees of freedom, and define $\mathbf{W}^* = [\mathbf{a}, -\mathbf{a}]^T$. As long as the number of training samples satisfies $n \geq Cd$, then with probability at

least $1 - Ce^{-cd}$ there exists $T \geq c' \frac{1}{\mu} \ln \left(\frac{\sqrt{\|\mathbf{a}\|}}{\sigma} \right)$ such that for all iterations $\tau > T$,

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

Here, $\mu_0, \sigma_0, c, \tilde{c}, c', C$ are fixed numerical constants independent of any problem dimensions.

This theorem shows that gradient descent can provably train a fully end-to-end one-hidden-layer ReLU network to learn a linear target from finitely many samples, despite the highly degenerate and saddle-rich optimization landscape. In particular, the result gives a global convergence guarantee for simultaneous optimization of the hidden and output weights; going beyond analyses that rely on effectively fixed features or only local perturbations around initialization. Starting from a small random initialization with the standard σ/\sqrt{d} scaling—consistent with common “default” initializations used in practice—the two student neurons w_1 and w_2 rapidly align with the ground-truth direction \mathbf{a} , after which the *effective parameters* $\text{diag} \left(\begin{bmatrix} v_1 \\ v_2 \end{bmatrix} \right) \mathbf{W}$ converge geometrically to the planted solution. This linear-rate convergence kicks in after only a short burn-in period of $T \geq c' \frac{1}{\mu} \ln \left(\frac{\sqrt{\|\mathbf{a}\|}}{\sigma} \right)$ iterations and requires only $n \geq Cd$ samples, which is information-theoretically optimal.

We note that the *same* finite dataset is reused across all iterations of gradient descent. Controlling the resulting dependence between the iterates and the samples requires uniform concentration arguments that hold *along the entire trajectory*, i.e., controlling population–empirical deviations simultaneously for all iterates visited by GD rather than at a fixed parameter value. Finally, the χ^2 -based initialization for the output weights is used for technical convenience. One can alternatively initialize v_1 and v_2 as Gaussians with variance $\sigma^2/2$; the same qualitative convergence behavior persists, but the success probability degrades to a fixed constant, rather than $1 - Ce^{-cd}$. This constant failure with Gaussian initialization is unavoidable for $k = 2$. For large k , scaling the initialization with $\frac{\sigma^2}{k}$ yields failure probability decaying as e^{-k} . We therefore use a slight non-Gaussian modification in the initialization to demonstrate that except for this $k = 2$ artifact our result holds with much higher probability.

4. Experiments

We run experiments on various output dimension (denoted with $r = 1$ vs. $r > 1$), and initialization scale (small vs large). In this section we show experimental results for single output $r = 1$ case and refer the reader to Appendix E.1 for multi-output $r > 1$ results. We use PyTorch for experiments and unless mentioned otherwise, network weights are initialized with Xavier Normal initialization (for a matrix $\mathbf{W} \in \mathbb{R}^{r \times d}$, $\mathbf{W}_{ij} \sim \mathcal{N} \left(0, \frac{2}{r+d} \right)$).

In order to change the initialization scale, we multiply the default initialization with a positive scalar σ . For *small* initialization experiments, we use $\sigma = 10^{-8}$, otherwise it is set to $\sigma = 1$. 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.

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 ($k = 2$), we empirically see that the model cannot converge to the global optima consistently. When it does, $\mathbf{v}^{(\infty)}$ indeed becomes ± 1 and $w_1^{(\infty)}$ and $w_2^{(\infty)}$ recover $\pm \mathbf{a}$ exactly. For the remaining time, the GD iterates converge to one of the many stationary points of this problem similar to the depiction in Figure 1 (part b). We further observe that iterates get stuck only when $\mathbf{v}_1^{(0)}$ and $\mathbf{v}_2^{(0)}$ both have the same signs which happens with probability $\frac{1}{2}$. This is also the reason for why we fix the correct sign pattern at initialization in Theorem 2.

When $k > 2$, the probability of *all* \mathbf{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. The trajectory of the inner weights across GD iterations is depicted in Figure 2.

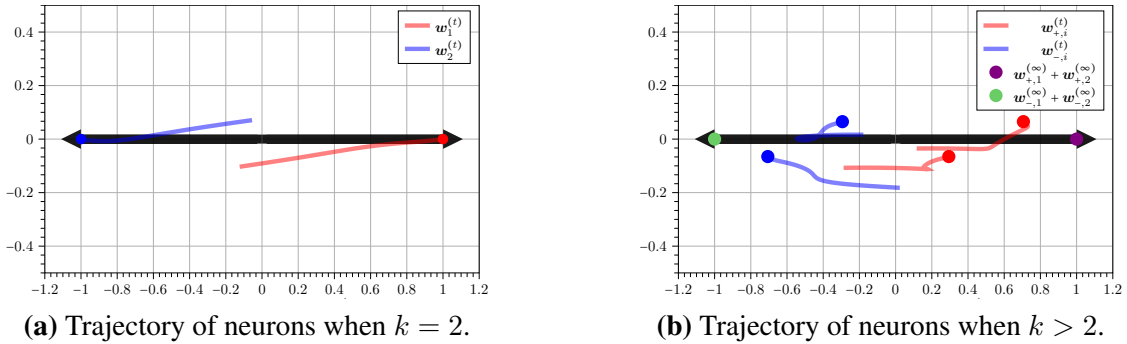


Figure 2: **Trajectory of neurons for different values of k .** We run gradient descent updates on the population loss. A randomly selected orthogonal direction to \mathbf{a} is shown for the y-axis in order to visualize the neurons in 2D. Black arrows indicate $\pm \mathbf{a}$ direction. We use colors **red** and **blue** to indicate whether \mathbf{v}_i corresponding to w_i is positive or negative 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 $\pm \mathbf{a}$ direction, grouping hidden units based on their corresponding signs in \mathbf{v} and summing them recovers $\pm \mathbf{a}$ *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 $\pm \mathbf{a}$ 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.

5. Overview and Key Ideas of the Proof

In this section, we outline the main ideas underlying our analysis. As mentioned previously, a major challenge is that the optimization landscape is riddled with non-strict saddle points that gradient descent can get stuck in. Thus, our analysis requires a very refined control of the trajectory to

guarantee that the iterates escape these saddle regions. We will show that the trajectory of full-batch gradient descent partitions into three distinct phases discussed below. Figure 3 illustrates the three phases and their interaction.

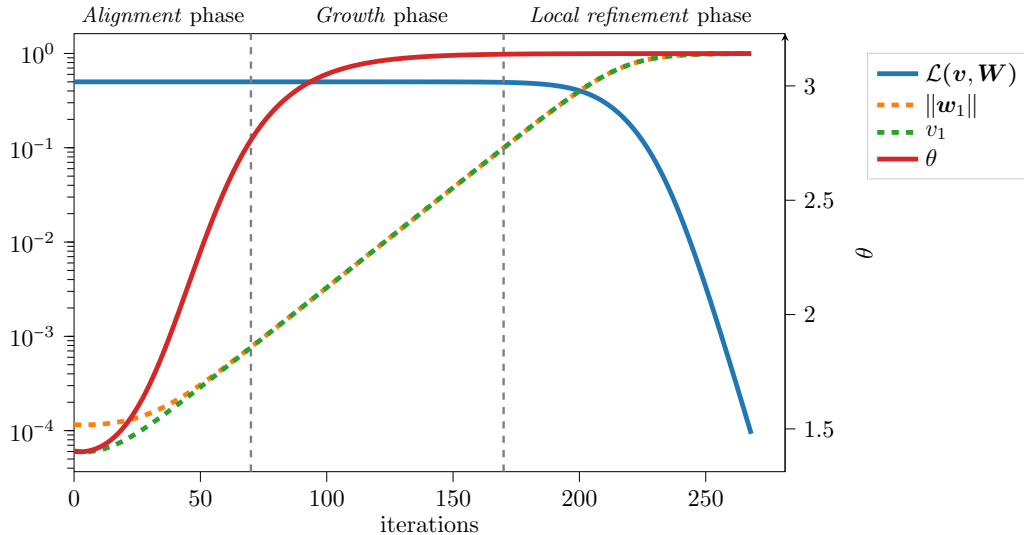


Figure 3: **Phases of GD Trajectory.** We run gradient descent updates on the population loss with small initialization $\sigma = 10^{-4}$. We track the population loss \mathcal{L} (blue), norms $v_1, \|w_1\|$ (green and yellow), and the θ – i.e. angle between w_1, w_2 – (red). For visualization purposes θ uses the right vertical axis. v_2 and $\|w_2\|$ behave similarly but omitted for clarity.

- (1) **Alignment phase (Section 5.4).** Starting from a small random initialization, we show that the hidden weights progressively align with the planted direction while the output weights maintain the correct sign pattern.
- (2) **Growth phase (Section 5.5).** Once sufficient alignment has been established, we prove that the norms of both the hidden and output layers grow in a coordinated fashion while preserving this alignment. This phase drives the effective parameters toward the correct scale and pushes the iterates away from flat saddle regions of the loss landscape. A key technical challenge here is to show that gradient descent does not drift into spurious stationary points despite the non-strict nature of these saddles.
- (3) **Local refinement phase (Section 5.6).** After the alignment and growth phases we enter a well-behaved region of the planted solution, where the dynamics become locally well-conditioned. In this phase, We show that the aligned neurons then converge rapidly to the ground-truth direction, and the effective parameters enjoy a linear rate of convergence to the global minimizer.

Throughout all three phases, the same finite dataset is reused across iterations. To control the resulting dependence between the iterates and the samples, we establish new trajectory-level uniform concentration bounds that hold simultaneously for all points visited by gradient descent. These results are crucial for obtaining order-wise optimal sample complexity. We give an overview of these

uniform concentration results in Section 5.7. Before we detail the specific phases of the trajectory, we also need to establish two sets of key identities. The first set demonstrates a specific property of balancedness between the inner and outer weights (Section 5.2). The second set concerns the stability of our training dynamics, ensuring that the evolution is monotonic in the sense that once the iterates enter a new phase, they do not revert to a previous one (Section 5.3). We begin with some quick notation used throughout our proofs.

5.1. Notation

In this section we gather some simple notation used in our proofs. As a reminder we use $\widehat{\mathcal{L}}$ and \mathcal{L} to denote the empirical and population losses, respectively. We use

$$\Delta\mathcal{G}_1 := \frac{2}{v_1} \left(\nabla_{\mathbf{w}_1} \widehat{\mathcal{L}} - \nabla_{\mathbf{w}_1} \mathcal{L} \right) \quad \text{and} \quad \Delta\mathcal{G}_2 := \frac{2}{v_2} \left(\nabla_{\mathbf{w}_2} \widehat{\mathcal{L}} - \nabla_{\mathbf{w}_2} \mathcal{L} \right)$$

to denote the scaled difference between the empirical and population gradients. Finally, we use θ to denote the angle between \mathbf{w}_1 and \mathbf{w}_2 . We also define θ_1 to be the angle between \mathbf{w}_1 and \mathbf{a} , θ_2 to be the angle between \mathbf{w}_2 and $-\mathbf{a}$. We note that all lemmas stated in this proof overview our under the assumptions of the main theorem, we avoid repeating these assumptions repeatedly for readability.

5.2. Controlling the imbalance term

A crucial identity used throughout our proofs is that from moderately small initialization the norms of the inner and outer weights remain close to each other. Concretely, we define the imbalance term as $b_1^{(\tau)} := \|\mathbf{w}_1^{(\tau)}\|^2 - (v_1^{(\tau)})^2$ and $b_2^{(\tau)} := \|\mathbf{w}_2^{(\tau)}\|^2 - (v_2^{(\tau)})^2$. A constant bound for the absolute value of these terms is required to prove that the norms remain bounded throughout the training process (see Lemma 5).

While the imbalance is invariant in gradient flow (Ji and Telgarsky, 2019), the discretization in gradient descent introduces a small drift given by:

$$b_1^{(\tau+1)} = b_1^{(\tau)} + \mu^2 \left(\|\nabla_{\mathbf{w}_1} \widehat{\mathcal{L}}\|^2 - (\nabla_{v_1} \widehat{\mathcal{L}})^2 \right).$$

A simple constant bound on this drift is insufficient for our analysis, as the errors could accumulate to infinity over an infinite number of iterations.

To address this, we prove a stronger result: the drift in each step is bounded by the distance between the effective weights $\text{diag} \left(\begin{bmatrix} v_1 \\ v_2 \end{bmatrix} \right) \mathbf{W}$ and the planted solution. Since $\text{diag} \left(\begin{bmatrix} v_1 \\ v_2 \end{bmatrix} \right) \mathbf{W}$ converges to the planted solution exponentially fast in Phase 3, the total accumulated drift remains finite even as $\tau \rightarrow \infty$. Concretely, we prove the following lemma.

Lemma 3 (Imbalance bound) *Assume that $v_1^{(\tau)}, v_2^{(\tau)} > 0$. For any $i \in \{1, 2\}$, we have*

$$\left| b_i^{(\tau+1)} - b_i^{(\tau)} \right| \leq c_6 \mu^2 \left((v_i^{(\tau)})^2 + \|\mathbf{w}_i^{(\tau)}\|^2 \right) \left(\left\| v_1^{(\tau)} \mathbf{w}_1^{(\tau)} - \mathbf{a} \right\|^2 + \left\| v_2^{(\tau)} \mathbf{w}_2^{(\tau)} + \mathbf{a} \right\|^2 \right).$$

Here, we set the constant as $c_6 = 6$.

This lemma is proven in Section C.1.

5.3. Stability of Training Dynamics

In this section, we establish two key stability properties that serve as the foundation for our proof. These results ensuring that the evolution is monotonic in the sense that once the iterates enter a new phase, they do not revert to a previous phase. The first lemma ensures that once the angle becomes small (at the end of the first phase) it continues to remain sufficiently small.

Lemma 4 (Angle stays small) *Assume that $\mu \leq \frac{c_0}{\|\mathbf{a}\|}$. For any iteration τ such that $0 < v_1^{(\tau)}, v_2^{(\tau)} \leq c_2 \sqrt{\|\mathbf{a}\|}$, $\theta_1^{(\tau)}, \theta_2^{(\tau)} \leq c_4$, $|b_1^{(\tau)}|, |b_2^{(\tau)}| \leq \gamma \|\mathbf{a}\|$ and $\|\Delta \mathcal{G}_1^{(\tau)}\|, \|\Delta \mathcal{G}_2^{(\tau)}\| \leq c_5 \|\mathbf{a}\|$, we have*

$$\theta_1^{(\tau+1)}, \theta_2^{(\tau+1)} \leq c_4.$$

Here, we set the constants as $c_0 \leq \frac{1}{2}$, $c_2 = 2$, $c_5 = \frac{1}{50}$, $c_4 = \frac{\pi}{20}$, $\gamma = \frac{1}{4}$.

This lemma is proven in Section C.2. The second result ensures the norms remain bounded.

Lemma 5 (Norms remain bounded) *Assume that $\mu \leq \frac{c_0}{\|\mathbf{a}\|}$. For any $0 \leq \beta \leq \frac{1}{4}$ and iteration τ such that $\beta \sqrt{\|\mathbf{a}\|} < v_1^{(\tau)}, v_2^{(\tau)} \leq c_2 \sqrt{\|\mathbf{a}\|}$, $\theta_1^{(\tau)}, \theta_2^{(\tau)} \leq c_4$, $|b_1^{(\tau)}|, |b_2^{(\tau)}| \leq \gamma \|\mathbf{a}\|$ and $\|\Delta \mathcal{G}_1^{(\tau)}\|, \|\Delta \mathcal{G}_2^{(\tau)}\| \leq c_5 \|\mathbf{a}\|$, we have*

$$\beta \sqrt{\|\mathbf{a}\|} < v_1^{(\tau+1)}, v_2^{(\tau+1)} \leq c_2 \sqrt{\|\mathbf{a}\|}.$$

Here, we set the constants as $c_0 \leq \frac{4}{25}$, $c_2 = 2$, $c_5 = \frac{1}{50}$, $c_4 = \frac{\pi}{20}$, $\gamma = \frac{1}{4}$.

This lemma is proven in Section C.3.

5.4. Overview of Alignment Phase

The primary objective of Phase 1 is to demonstrate that \mathbf{w}_1 becomes approximately aligned with \mathbf{a} (and \mathbf{w}_2 with $-\mathbf{a}$) within a constant number of steps, which is crucial for Phases 2 and 3. By symmetry, we focus on \mathbf{w}_1 .

Our key observation is that the gradient update is dominated by the signal direction. Specifically, the update can be decomposed as:

$$\mathbf{w}_1^{(\tau+1)} = \mathbf{w}_1^{(\tau)} + \mu \left(\frac{v_1}{2} \mathbf{a} + \zeta^{(\tau)} \right).$$

where the remainder term $\zeta^{(\tau)}$ consists of terms involving the weights and the empirical noise. Since we use small initialization, these weight-dependent terms are much smaller than the signal term. This implies that the projection of \mathbf{w}_1 onto the signal direction \mathbf{a} grows much faster than its projection onto the orthogonal subspace. Specifically, we have the following lemma proven in Section C.4:

Lemma 6 (Angle alignment) *Assume that $\mu \leq \frac{c_0}{\|\mathbf{a}\|}$, $\sigma \leq \sigma_0 \sqrt{\|\mathbf{a}\|}$. After $T_1 = \lceil \frac{c_9}{\mu \|\mathbf{a}\|} \rceil$ iterations, it holds that*

$$\theta_1^{(T_1)}, \theta_2^{(T_1)} \leq c_4$$

with probability at least $1 - Ce^{-cd}$. Moreover, this alignment is achieved while maintaining that:

$$c_3\sigma \leq v_1^{(T_1)}, v_2^{(T_1)} \leq c_2\sqrt{\|\mathbf{a}\|}, \quad \left|b_1^{(T_1)}\right|, \left|b_2^{(T_1)}\right| \leq c_{10}\|\mathbf{a}\|.$$

Here, we set the constants as $c_0 \leq 1, c_4 = \frac{\pi}{20}, c_9 = \frac{64}{\tan c_4}, c_2 = 2, c_3 = \frac{1}{4}$. With $\alpha = \frac{65}{\tan c_4}$, we have $c_{10} = \frac{1}{4e^{2\alpha}}, \sigma_0 = \frac{1}{8e^{2\alpha}}$.

5.5. Overview of Growth Phase

In Phase 3, we show that the effective weights $\text{diag}\left(\begin{bmatrix} v_1 \\ v_2 \end{bmatrix}\right) \mathbf{W}$ converge to the planted solution at an exponential rate. A key ingredient is a Polyak–Lojasiewicz (PL) inequality for the population loss (Lemma 8), which lower-bounds the squared gradient norm in terms of the suboptimality gap. Importantly, the PL constant depends on the magnitudes of v_1 and v_2 .

At the end of Phase 1, $|v_1|$ and $|v_2|$ remain at their initialization scale, so the PL inequality only yields a weak contraction and therefore a slow convergence rate. The main goal of Phase 2 is to grow v_1 and v_2 to a sufficiently large scale, thereby strengthening the PL constant and enabling fast linear convergence in Phase 3.

Lemma 7 (Norm growth) *Assume that $\mu \leq \frac{c_0}{\|\mathbf{a}\| \ln\left(\frac{\sqrt{\|\mathbf{a}\|}}{\sigma}\right)}$. After $T_2 = \lceil \frac{c_8}{\mu\|\mathbf{a}\|} \ln\left(\frac{\sqrt{\|\mathbf{a}\|}}{\sigma}\right) \rceil$ iterations, we have*

$$v_1^{(T_1+T_2)}, v_2^{(T_1+T_2)} \geq c_7\sqrt{\|\mathbf{a}\|}, \quad \left|b_1^{(T_1+T_2)}\right|, \left|b_2^{(T_1+T_2)}\right| \leq c_{11}\|\mathbf{a}\|.$$

Here, we set the constants as $c_0 \leq \frac{1}{10^8}, c_7 = \frac{1}{4}, c_8 = 32, c_{11} = \frac{1}{50}$.

This lemma is proven in Section C.5.

5.6. Overview of the local Refinement Phase

In the final *local refinement* phase, we show that the effective weights $\text{diag}\left(\begin{bmatrix} v_1^{(\tau)} \\ v_2^{(\tau)} \end{bmatrix}\right) \mathbf{W}^{(\tau)}$ converge to the planted solution $\mathbf{W}^* = [\mathbf{a}, -\mathbf{a}]^T$. Rather than tracking parameters directly, we first prove that along empirical gradient descent the *population* loss decreases rapidly, and then convert this decay into the stated parameter convergence. The full argument (Proof of Theorem 2 in Section D.2) is technical: it couples a population-level gradient-descent analysis with trajectory-uniform concentration bounds (next section). For clarity, we sketch only the population argument here. This population reduction is essential because the empirical loss is not smooth and as discussed below, even the population loss is not uniformly smooth. The proof proceeds in two parts.

Part 1 (PL inequality) In this step we will show the following PL inequality proven in Section B.2.

Lemma 8 (PL Inequality for the population loss) *For $v_1, v_2 > 0$,*

$$\|\nabla_{\mathbf{w}_1} \mathcal{L}(\mathbf{v}, \mathbf{W})\|^2 + \|\nabla_{\mathbf{w}_2} \mathcal{L}(\mathbf{v}, \mathbf{W})\|^2 \geq \alpha \min(v_1^2, v_2^2) \mathcal{L}(\mathbf{v}, \mathbf{W})$$

holds with $\alpha = 0.05$ as long as $\theta > \frac{\pi}{2}$.

To prove this PL inequality we first show that it can be deduced by establishing the PL inequality when $v_1 = v_2 = 1$ via a clever reduction argument. To prove the latter we define

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

Note that since we set $v_1 = v_2 = 1$ the loss is now only a function of \mathbf{w}_1 and \mathbf{w}_2 . Also note that to prove the PL inequality it suffices to show that h is always positive. To do this, in our proof we show that $\tilde{h}(\mathbf{w}_1, \mathbf{w}_2) = \min_{\mathbf{a}} h(\mathbf{w}_1, \mathbf{w}_2, \mathbf{a})$ is always positive. The way we establish this is by showing that $\frac{1}{\|\mathbf{w}_2\|^2} \tilde{h}(\mathbf{w}_1, \mathbf{w}_2)$ is only a function of θ (the angle between \mathbf{w}_1 and \mathbf{w}_2) and $\frac{\|\mathbf{w}_1\|}{\|\mathbf{w}_2\|}$. Since this is now only a function of two variables θ and $\frac{\|\mathbf{w}_1\|}{\|\mathbf{w}_2\|}$ it is easy to establish non-negativity as long as $\theta > \frac{\pi}{2}$. The latter holds at the end of the growth phase and continues to remain large utilizing the stability of the dynamics established in Section 5.3 (concretely, Lemma 4 shows the angles with planted directions remain small which implies the angle between the weight vectors remain large).

Part 2 (Gradient smoothness) – While the population loss is not smooth in the entire domain, we show that in the region of the local refinement phase it is indeed smooth. Leaving the *growth* phase, we have lower/upper bounds on $v_1, v_2, \|\mathbf{w}_1\|, \|\mathbf{w}_2\|$. Additionally, due to the stability analysis in Section 5.3 we continue to have lower/upper bounds on these quantities. In the next lemma we show that assuming such lower/upper bounds the population loss is indeed smooth. This lemma is proven in Section B.3.

Lemma 9 (Smoothness of the population loss) $\|\nabla^2 \mathcal{L}(v, \mathbf{W})\|_F \leq L \|\mathbf{a}\|$ holds for all $v \in \mathbb{R}^2, \mathbf{W} \in \mathbb{R}^{2 \times d}$ such that $c_1 \sqrt{\|\mathbf{a}\|} \leq v_1, v_2, \|\mathbf{w}_1\|, \|\mathbf{w}_2\| \leq c_2 \sqrt{\|\mathbf{a}\|}$ holds. Here c_1, c_2, L are fixed constants.

Showing geometric decrease of the *population* loss under a PL inequality and smoothness is a classical optimization result. Our setting is more delicate because we run *empirical* gradient descent rather than its population counterpart. See the proof of Theorem 2 in Section D.2 for how we combine the trajectory-uniform concentration bounds (next section) with the PL and smoothness properties of the population loss stated above to obtain a geometric decrease in the population loss.

5.7. Uniform Concentration

In this section, we provide an overview of the novel uniform concentration result that we have established which is key to our near optimal sample complexity. In particular, the concentration holds along the entire trajectory of GD and is used across the three phases. We provide the setup next. Let $\mathbf{x}_1, \dots, \mathbf{x}_n \stackrel{\text{iid}}{\sim} \mathcal{N}(0, \mathbf{I}_d)$ in \mathbb{R}^d . For $\mathbf{w}, \mathbf{w}^* \in \mathbb{S}^{d-1}$ define

$$M(\mathbf{w}, \mathbf{w}^*) := \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{\{\langle \mathbf{x}_i, \mathbf{w} \rangle \geq 0\}} \mathbb{1}_{\{\langle \mathbf{x}_i, \mathbf{w}^* \rangle \geq 0\}} \mathbf{x}_i \mathbf{x}_i^T. \quad (4)$$

We prove a high-probability bound on

$$\sup_{\mathbf{w}, \mathbf{w}^* \in \mathbb{S}^{d-1}} \|M(\mathbf{w}, \mathbf{w}^*) - \mathbb{E}[M(\mathbf{w}, \mathbf{w}^*)]\|.$$

Lemma 10 (Uniform Concentration) Fix $\delta \in (0, 1/2)$. There exist universal constants $C, c > 0$ such that the following holds. If

$$n \geq C d \frac{\log^2(1/\delta)}{\delta^2}, \quad (5)$$

then with probability at least $1 - 3e^{-cd}$,

$$\sup_{\mathbf{w}, \mathbf{w}^* \in \mathbb{S}^{d-1}} \|M(\mathbf{w}, \mathbf{w}^*) - \mathbb{E}[M(\mathbf{w}, \mathbf{w}^*)]\| \leq \delta.$$

This lemma is proven in Section C.6. Notably, this lemma allows us to establish separate high-probability bounds for the deviations of the gradient components with respect to \mathbf{w}_1 and \mathbf{w}_2 . Concretely, it allows us to show the following lemma proven in Section C.7.

Lemma 11 (Component-wise Gradient Deviation Bounds) Fix $\delta \in (0, 1/2)$. Define the error vectors $\mathbf{h}_1 = v_1 \mathbf{w}_1 - \mathbf{a}$ and $\mathbf{h}_2 = v_2 \mathbf{w}_2 + \mathbf{a}$. Under the sample complexity $n \geq C d \frac{\log^2(1/\delta)}{\delta^2}$, with probability at least $1 - 3e^{-cd}$, the following bounds hold simultaneously:

$$\left\| \nabla_{\mathbf{w}_1} \widehat{\mathcal{L}} - \nabla_{\mathbf{w}_1} \mathcal{L} \right\| \leq v_1 \delta (\|\mathbf{h}_1\| + \|\mathbf{h}_2\|), \quad (6)$$

$$\left\| \nabla_{\mathbf{w}_2} \widehat{\mathcal{L}} - \nabla_{\mathbf{w}_2} \mathcal{L} \right\| \leq v_2 \delta (\|\mathbf{h}_1\| + \|\mathbf{h}_2\|). \quad (7)$$

This lemma highlights a particularly favorable self-regularizing property of the dynamics: the empirical–population gradient deviation scales *linearly* with the current error $\|\mathbf{h}_1\| + \|\mathbf{h}_2\|$. As the iterates approach the global optimum and the errors shrink, the concentration bounds automatically tighten, yielding increasingly accurate gradient estimates along the trajectory. In other words, concentration improves precisely when it is most needed in the local refinement regime, enabling stable geometric convergence in this region.

Acknowledgements

This work was partially supported by AWS credits through an Amazon Faculty Research Award, a NAIRR Pilot Award, and generous funding by Coefficient Giving. M. Soltanolkotabi is supported by the Packard Fellowship in Science and Engineering, a Sloan Research Fellowship in Mathematics, NSF CAREER Award #1846369, DARPA FastNICS program, NSF CIF Awards #1813877 and #2008443, and NIH Award DP2LM014564-01.

References

- Sanjeev Arora, Simon S Du, Wei Hu, Zhiyuan Li, and Ruosong Wang. Fine-grained analysis of optimization and generalization for overparameterized two-layer neural networks. In *36th International Conference on Machine Learning, ICML 2019*, pages 477–502. International Machine Learning Society (IMLS), 2019.
- Jimmy Ba, Murat A Erdogdu, Taiji Suzuki, Zhichao Wang, Denny Wu, and Greg Yang. High-dimensional asymptotics of feature learning: How one gradient step improves the representation. *Advances in Neural Information Processing Systems*, 35:37932–37946, 2022.
- Pierre Baldi and Kurt Hornik. Neural networks and principal component analysis: Learning from examples without local minima. *Neural Networks*, 2(1):53–58, 1989. ISSN 0893-6080. 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>.
- Etienne Boursier and Nicolas Flammarion. Simplicity bias and optimization threshold in two-layer relu networks, 2025. URL <https://arxiv.org/abs/2410.02348>.
- Olivier Bousquet. A Bennett concentration inequality and its application to suprema of empirical processes. *C. R. Acad. Sci. Paris, Ser. I*, 334(6):495–500, 2002.
- Alon Brutzkus and Amir Globerson. Globally optimal gradient descent for a convnet with gaussian inputs, 2017. URL <https://arxiv.org/abs/1702.07966>.
- Emmanuel J. Candès, Xiaodong Li, and Mahdi Soltanolkotabi. Phase retrieval via Wirtinger flow: theory and algorithms. *IEEE Trans. Inf. Theory*, 61(4):1985–2007, 2015. ISSN 0018-9448.
- Ji Chen, Dekai Liu, and Xiaodong Li. Nonconvex rectangular matrix completion via gradient descent without ℓ_2, ∞ regularization. *IEEE Trans. Inf. Theory*, 66(9):5806–5841, 2020. ISSN 0018-9448.
- Yuxin Chen and Emmanuel J. Candès. Solving random quadratic systems of equations is nearly as easy as solving linear systems. *Commun. Pure Appl. Math.*, 70(5):822–883, 2017. ISSN 0010-3640; 1097-0312/e.
- Dmitry Chistikov, Matthias Englert, and Ranko Lazic. Learning a neuron by a shallow relu network: Dynamics and implicit bias for correlated inputs, 2023. URL <https://arxiv.org/abs/2306.06479>.
- Lénaïc Chizat, Edouard Oyallon, and Francis Bach. On lazy training in differentiable programming. *Advances in Neural Information Processing Systems*, 32:2937–2947, 2019.
- Alexandru Damian, Jason Lee, and Mahdi Soltanolkotabi. Neural networks can learn representations with gradient descent. In *Conference on Learning Theory*, pages 5413–5452. PMLR, 2022.
- Amit Daniely, Roy Frostig, and Yoram Singer. Toward deeper understanding of neural networks: The power of initialization and a dual view on expressivity. *Advances in neural information processing systems*, 29, 2016.

- Simon Du, Jason Lee, Haochuan Li, Liwei Wang, and Xiyu Zhai. Gradient descent finds global minima of deep neural networks. In *International Conference on Machine Learning*, pages 1675–1685. PMLR, 2019.
- Simon S Du, Xiyu Zhai, Barnabas Poczos, and Aarti Singh. Gradient descent provably optimizes over-parameterized neural networks. *arXiv preprint arXiv:1810.02054*, 2018.
- Zalan Fabian, Berk Tinaz, and Mahdi Soltanolkotabi. Humus-net: Hybrid unrolled multi-scale network architecture for accelerated mri reconstruction. *Advances in Neural Information Processing Systems*, 35:25306–25319, 2022.
- Rong Ge, Furong Huang, Chi Jin, and Yang. Yuan. Escaping from saddle points: online stochastic gradient for tensor decomposition. In *Proceedings of The 28th Conference on Learning Theory*, pages 797–842, 2015.
- Rong Ge, Jason D Lee, and Tengyu Ma. Matrix completion has no spurious local minimum. *Advances in Neural Information Processing Systems*, 29:2973–2981, 2016.
- Behrooz Ghorbani, Song Mei, Theodor Misiakiewicz, and Andrea Montanari. When do neural networks outperform kernel methods? *arXiv preprint arXiv:2006.13409*, 2020.
- Paul Hand, Oscar Leong, and Vlad Voroninski. Phase retrieval under a generative prior. *Advances in Neural Information Processing Systems*, 31, 2018.
- Arthur Jacot, Franck Gabriel, and Clément Hongler. Neural tangent kernel: Convergence and generalization in neural networks. In *Advances in neural information processing systems*, pages 8571–8580, 2018.
- Ziwei Ji and Matus Telgarsky. Gradient descent aligns the layers of deep linear networks, 2019. URL <https://arxiv.org/abs/1810.02032>.
- Chi Jin, Rong Ge, Praneeth Netrapalli, Sham M. Kakade, and Michael I. Jordan. How to escape saddle points efficiently. page 1724–1732, 2017a.
- Kyong Hwan Jin, Michael T McCann, Emmanuel Froustey, and Michael Unser. Deep convolutional neural network for inverse problems in imaging. *IEEE transactions on image processing*, 26(9): 4509–4522, 2017b.
- Kevin Kögler, Alexander Shevchenko, Hamed Hassani, and Marco Mondelli. Compression of structured data with autoencoders: Provable benefit of nonlinearities and depth, 2024. URL <https://arxiv.org/abs/2402.05013>.
- Christian Ledig, Lucas Theis, Ferenc Huszár, Jose Caballero, Andrew Cunningham, Alejandro Acosta, Andrew Aitken, Alykhan Tejani, Johannes Totz, Zehan Wang, et al. Photo-realistic single image super-resolution using a generative adversarial network. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4681–4690, 2017.
- Jason D Lee, Kazusato Oko, Taiji Suzuki, and Denny Wu. Neural network learns low-dimensional polynomials with sgd near the information-theoretic limit. *arXiv preprint arXiv:2406.01581*, 2024.

- Binghui Li, Zhixuan Pan, Kaifeng Lyu, and Jian Li. Feature averaging: An implicit bias of gradient descent leading to non-robustness in neural networks. *arXiv preprint arXiv:2410.10322*, 2024.
- Xiaodong Li, Shuyang Ling, Thomas Strohmer, and Ke Wei. Rapid, robust, and reliable blind deconvolution via nonconvex optimization. *Appl. Comput. Harmon. Anal.*, 47(3):893–934, 2019. ISSN 1063-5203.
- Shuyang Ling and Thomas Strohmer. Regularized gradient descent: a non-convex recipe for fast joint blind deconvolution and demixing. *Inf. Inference*, 8(1):1–49, 2019. ISSN 2049-8764; 2049-8772/e.
- Cong Ma, Kaizheng Wang, Yuejie Chi, and Yuxin Chen. Implicit regularization in nonconvex statistical estimation: gradient descent converges linearly for phase retrieval, matrix completion, and blind deconvolution. *Found. Comput. Math.*, 20(3):451–632, 2020. ISSN 1615-3375; 1615-3383/e.
- Yurii Nesterov and Boris T. Polyak. Cubic regularization of Newton method and its global performance. *Math. Program.*, 108(1 (A)):177–205, 2006. ISSN 0025-5610; 1436-4646/e.
- Jorge Nocedal and Stephen J. Wright. Trust-region methods. *Numerical Optimization*, pages 66–100, 2006.
- Greg Ongie, Rebecca Willett, Daniel Soudry, and Nathan Srebro. A function space view of bounded norm infinite width relu nets: The multivariate case. *arXiv preprint arXiv:1910.01635*, 2019.
- Samet Oymak and Mahdi Soltanolkotabi. Overparameterized nonlinear learning: Gradient descent takes the shortest path? In *International Conference on Machine Learning*, pages 4951–4960. PMLR, 2019.
- Samet Oymak and Mahdi Soltanolkotabi. Towards moderate overparameterization: global convergence guarantees for training shallow neural networks. *IEEE Journal on Selected Areas in Information Theory*, 2020.
- Alethea Power, Yuri Burda, Harri Edwards, Igor Babuschkin, and Vedant Misra. Grokking: Generalization beyond overfitting on small algorithmic datasets, 2022. URL <https://arxiv.org/abs/2201.02177>.
- Maxim Raginsky, Alexander Rakhlin, and Matus Telgarsky. Non-convex learning via stochastic gradient langevin dynamics: a nonasymptotic analysis. pages 1674–1703, 2017.
- Yunwei Ren, Eshaan Nichani, Denny Wu, and Jason D. Lee. Emergence and scaling laws in sgd learning of shallow neural networks, 2025. URL <https://arxiv.org/abs/2504.19983>.
- Alexander Shevchenko, Kevin Kögler, Hamed Hassani, and Marco Mondelli. Fundamental limits of two-layer autoencoders, and achieving them with gradient methods, 2022.
- Mahdi Soltanolkotabi. Learning relus via gradient descent. *Advances in neural information processing systems*, 30, 2017.

- Mahdi Soltanolkotabi, Adel Javanmard, and Jason D Lee. Theoretical insights into the optimization landscape of over-parameterized shallow neural networks. *IEEE Transactions on Information Theory*, 65(2):742–769, 2018.
- Mahdi Soltanolkotabi, Dominik Stöger, and Changzhi Xie. Implicit balancing and regularization: Generalization and convergence guarantees for overparameterized asymmetric matrix sensing. In Gergely Neu and Lorenzo Rosasco, editors, *Proceedings of Thirty Sixth Conference on Learning Theory*, volume 195 of *Proceedings of Machine Learning Research*, pages 5140–5142. PMLR, 12–15 Jul 2023. URL <https://proceedings.mlr.press/v195/soltanolkotabi23a.html>.
- Anuroop Sriram, Jure Zbontar, Tullie Murrell, Aaron Defazio, C Lawrence Zitnick, Nafissa Yakubova, Florian Knoll, and Patricia Johnson. End-to-end variational networks for accelerated mri reconstruction. In *Medical Image Computing and Computer Assisted Intervention—MICCAI 2020: 23rd International Conference, Lima, Peru, October 4–8, 2020, Proceedings, Part II 23*, pages 64–73. Springer, 2020.
- Dominik Stöger and Mahdi Soltanolkotabi. Small random initialization is akin to spectral learning: Optimization and generalization guarantees for overparameterized low-rank matrix reconstruction. In M. Ranzato, A. Beygelzimer, Y. Dauphin, P.S. Liang, and J. Wortman Vaughan, editors, *Advances in Neural Information Processing Systems*, volume 34, pages 23831–23843. Curran Associates, Inc., 2021. URL https://proceedings.neurips.cc/paper_files/paper/2021/file/c82836ed448c41094025b4a872c5341e-Paper.pdf.
- Ju Sun, Qing Qu, and John Wright. When are nonconvex problems not scary? *arXiv preprint arXiv:1510.06096*, 2015.
- Ju Sun, Qing Qu, and John Wright. A geometric analysis of phase retrieval. *Found. Comput. Math.*, 18(5):1131–1198, 2018. ISSN 1615-3375; 1615-3383/e.
- Stephen Tu, Ross Boczar, Max Simchowitz, Mahdi Soltanolkotabi, and Ben Recht. Low-rank solutions of linear matrix equations via procrustes flow. In *International Conference on Machine Learning*, pages 964–973. PMLR, 2016.
- Aad W. Van der Vaart and Jon A. Wellner. *Weak Convergence and Empirical Processes: With Applications to Statistics*. Springer Series in Statistics. Springer, New York, NY, 1996. ISBN 978-0-387-94640-5.
- Gal Vardi, Gilad Yehudai, and Ohad Shamir. Learning a single neuron with bias using gradient descent, 2022. URL <https://arxiv.org/abs/2106.01101>.
- Martin J. Wainwright. *High-dimensional statistics*, volume 48 of *Cambridge Series in Statistical and Probabilistic Mathematics*. Cambridge University Press, Cambridge, 2019. ISBN 978-1-108-49802-9. doi: 10.1017/9781108627771. URL <https://doi.org/10.1017/9781108627771>. A non-asymptotic viewpoint.
- Yi Wang, Xin Tao, Xiaojuan Qi, Xiaoyong Shen, and Jiaya Jia. Image inpainting via generative multi-column convolutional neural networks. *Advances in neural information processing systems*, 31, 2018.

- Yuqing Wang, Minshuo Chen, Tuo Zhao, and Molei Tao. Large learning rate tames homogeneity: Convergence and balancing effect, 2022. URL <https://arxiv.org/abs/2110.03677>.
- Yuqing Wang, Zhenghao Xu, Tuo Zhao, and Molei Tao. Good regularity creates large learning rate implicit biases: edge of stability, balancing, and catapult, 2023. URL <https://arxiv.org/abs/2310.17087>.
- Chenwei Wu, Jiajun Luo, and Jason D. Lee. No spurious local minima in a two hidden unit reLU network, 2018. URL <https://openreview.net/forum?id=B14uJzW0b>.
- Weihang Xu and Simon S. Du. Over-parameterization exponentially slows down gradient descent for learning a single neuron, 2023.
- Gilad Yehudai and Ohad Shamir. Learning a single neuron with gradient methods, 2022. URL <https://arxiv.org/abs/2001.05205>.
- Richard Y. Zhang, Somayeh Sojoudi, and Javad Lavaei. Sharp restricted isometry bounds for the inexistence of spurious local minima in nonconvex matrix recovery. *J. Mach. Learn. Res.*, 20(114):1–34, 2019. URL <http://jmlr.org/papers/v20/19-020.html>.
- Xiao Zhang, Yaodong Yu, Lingxiao Wang, and Quanquan Gu. Learning one-hidden-layer relu networks via gradient descent, 2018. URL <https://arxiv.org/abs/1806.07808>.
- Yuchen Zhang, Percy Liang, and Moses Charikar. A hitting time analysis of stochastic gradient langevin dynamics. In Satyen Kale and Ohad Shamir, editors, *Proceedings of the 2017 Conference on Learning Theory*, volume 65 of *Proceedings of Machine Learning Research*, pages 1980–2022. PMLR, 07–10 Jul 2017. URL <http://proceedings.mlr.press/v65/zhang17b.html>.
- Kai Zhong, Zhao Song, Prateek Jain, Peter L. Bartlett, and Inderjit S. Dhillon. Recovery guarantees for one-hidden-layer neural networks, 2017. URL <https://arxiv.org/abs/1706.03175>.
- Zhenyu Zhu, Fanghui Liu, and Volkan Cevher. How gradient descent balances features: A dynamical analysis for two-layer neural networks. In *The Thirteenth International Conference on Learning Representations*, 2025. URL <https://openreview.net/forum?id=25j2ZEgwTj>.

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

$$\begin{aligned} f(\mathbf{a}, \mathbf{b}) &= \mathbb{E}_{\mathbf{x}} \left[[\mathbf{a}^T \mathbf{x}]_+ [\mathbf{b}^T \mathbf{x}]_+ \right] \\ &\stackrel{(a)}{=} \frac{1}{2\pi} \|\mathbf{a}\| \|\mathbf{b}\| (\sin(\theta_{\mathbf{a},\mathbf{b}}) + (\pi - \theta_{\mathbf{a},\mathbf{b}}) \cos(\theta_{\mathbf{a},\mathbf{b}})) \end{aligned} \quad (8)$$

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:

$$\begin{aligned} \mathcal{L}(\boldsymbol{\theta}) &= \frac{1}{2} \mathbb{E}_{\mathbf{x}} \left[\left\| \sum_{i=1}^k v_i \phi(\mathbf{w}_i^T \mathbf{x}) - \mathbf{a}^T \mathbf{x} \right\|^2 \right] \\ &= \frac{1}{2} \sum_{i=1}^k \sum_{j=1}^k v_i v_j \mathbb{E}_{\mathbf{x}} [\phi(\mathbf{w}_i^T \mathbf{x}) \phi(\mathbf{w}_j^T \mathbf{x})] - \sum_{i=1}^k v_i \mathbf{a}^T \mathbb{E}_{\mathbf{x}} [\phi(\mathbf{w}_i^T \mathbf{x}) \mathbf{x}] + \frac{1}{2} \mathbb{E}_{\mathbf{x}} [\mathbf{a}^T \mathbf{x} \mathbf{x}^T \mathbf{a}] \\ &= \frac{1}{2} \sum_{i=1}^k \sum_{j=1}^k v_i v_j \mathbb{E}_{\mathbf{x}} [\phi(\mathbf{w}_i^T \mathbf{x}) \phi(\mathbf{w}_j^T \mathbf{x})] - \sum_{i=1}^k v_i \mathbf{a}^T \mathbb{E}_{\mathbf{x}} [\phi(\mathbf{w}_i^T \mathbf{x}) \mathbf{x}] + \frac{\|\mathbf{a}\|^2}{2} \\ &\stackrel{(a)}{=} \frac{1}{2} \sum_{i=1}^k \sum_{j=1}^k v_i v_j f(\mathbf{w}_i, \mathbf{w}_j) - \sum_{i=1}^k v_i \mathbf{a}^T \mathbb{E}_{\mathbf{x}} [\phi(\mathbf{w}_i^T \mathbf{x}) \mathbf{x}] + \frac{\|\mathbf{a}\|^2}{2} \\ &\stackrel{(b)}{=} \frac{1}{2} \sum_{i=1}^k \sum_{j=1}^k v_i v_j f(\mathbf{w}_i, \mathbf{w}_j) - \sum_{i=1}^k v_i \mathbf{a}^T \mathbb{E}_{\mathbf{x}} [\nabla_{\mathbf{x}} \phi(\mathbf{w}_i^T \mathbf{x})] + \frac{\|\mathbf{a}\|^2}{2} \\ &= \frac{1}{2} \sum_{i=1}^k \sum_{j=1}^k v_i v_j f(\mathbf{w}_i, \mathbf{w}_j) - \sum_{i=1}^k v_i \mathbf{a}^T \mathbf{w}_i \mathbb{E}_{\mathbf{x}} [\phi'(\mathbf{w}_i^T \mathbf{x})] + \frac{\|\mathbf{a}\|^2}{2} \\ &\stackrel{(c)}{=} \frac{1}{2} \sum_{i=1}^k \sum_{j=1}^k v_i v_j f(\mathbf{w}_i, \mathbf{w}_j) - \frac{1}{2} \sum_{i=1}^k v_i \mathbf{a}^T \mathbf{w}_i + \frac{\|\mathbf{a}\|^2}{2} \\ &= \frac{1}{4\pi} \sum_{i=1}^k \sum_{j=1}^k v_i v_j \|\mathbf{w}_i\| \|\mathbf{w}_j\| (\sin \theta_{ij} + (\pi - \theta_{ij}) \cos \theta_{ij}) - \frac{1}{2} \sum_{i=1}^k v_i \mathbf{a}^T \mathbf{w}_i + \frac{\|\mathbf{a}\|^2}{2} \end{aligned} \quad (9)$$

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

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

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

where $\omega_i = \|\mathbf{w}_i\|$, $\mathbf{u} = \text{diag}(\boldsymbol{\omega}) \mathbf{v}$, and θ_{ij} is the angle between \mathbf{w}_i and \mathbf{w}_j .

A.2. Population Gradient

Gradient w.r.t. \mathbf{W} : Let us define,

$$\begin{aligned} g(\mathbf{a}, \mathbf{b}) &= \frac{\partial}{\partial \mathbf{a}} f(\mathbf{a}, \mathbf{b}) \\ &= \frac{1}{2\pi} (\|\mathbf{b}\| \sin(\theta_{\mathbf{a}, \mathbf{b}}) \bar{\mathbf{a}} + (\pi - \theta_{\mathbf{a}, \mathbf{b}}) \mathbf{b}) \quad \left(\bar{\mathbf{a}} = \frac{\mathbf{a}}{\|\mathbf{a}\|}, \quad \bar{\mathbf{b}} = \frac{\mathbf{b}}{\|\mathbf{b}\|} \right) \\ &= \frac{\|\mathbf{b}\|}{2\pi} (\sin(\theta_{\mathbf{a}, \mathbf{b}}) \bar{\mathbf{a}} + (\pi - \theta_{\mathbf{a}, \mathbf{b}}) \bar{\mathbf{b}}). \end{aligned} \quad (10)$$

Taking the derivative of (9) with respect to \mathbf{w}_i , we get

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

In matrix form:

$$\nabla_{\mathbf{W}} \mathcal{L}(\boldsymbol{\theta}) = \frac{1}{2\pi} \text{diag}(\mathbf{v}) ((\pi \mathbb{1} \mathbb{1}^T - \boldsymbol{\Theta}) \text{diag}(\mathbf{u}) + \text{diag}(\sin(\boldsymbol{\Theta}) \mathbf{u})) \bar{\mathbf{W}} - \frac{1}{2} \mathbf{v} \mathbf{a}^T \quad (11)$$

where $\omega_i = \|\mathbf{w}_i\|$, and $\mathbf{u} = \text{diag}(\boldsymbol{\omega}) \mathbf{v}$.

Gradient w.r.t. \mathbf{v} : Taking the derivative of (9) with respect to v_i , we get

$$\begin{aligned} \nabla_{v_i} \mathcal{L}(\boldsymbol{\theta}) &= v_i f(\mathbf{w}_i, \mathbf{w}_i) + \sum_{\substack{j=1 \\ i \neq j}}^k v_j f(\mathbf{w}_i, \mathbf{w}_j) - \frac{1}{2} \mathbf{a}^T \mathbf{w}_i \\ &= \sum_{j=1}^k v_j f(\mathbf{w}_i, \mathbf{w}_j) - \frac{1}{2} \mathbf{a}^T \mathbf{w}_i \\ &= \frac{1}{2\pi} \sum_{j=1}^k v_j \|\mathbf{w}_i\| \|\mathbf{w}_j\| (\sin \theta_{ij} + (\pi - \theta_{ij}) \cos \theta_{ij}) - \frac{1}{2} \mathbf{a}^T \mathbf{w}_i \end{aligned}$$

In matrix form:

$$\nabla_{\mathbf{v}} \text{Loss} = \frac{1}{2\pi} \text{diag}(\boldsymbol{\omega}) (\sin \boldsymbol{\Theta} + \cos \boldsymbol{\Theta} \odot (\pi \mathbb{1} \mathbb{1}^T - \boldsymbol{\Theta})) \text{diag}(\boldsymbol{\omega}) \mathbf{v} - \frac{1}{2} \mathbf{W} \mathbf{a} \quad (12)$$

where $\omega_i = \|\mathbf{w}_i\|$. Finally, we note that the gradient w.r.t \mathbf{w}_i and v_i are related with the following simple identity:

$$\mathbf{w}_i^T \nabla_{\mathbf{w}_i} \mathcal{L}(\boldsymbol{\theta}) = v_i \nabla_{v_i} \mathcal{L}(\boldsymbol{\theta}). \quad (13)$$

A.3. Population Hessian

The Hessian consists of four blocks (3 unique) due to interaction of \mathbf{v} and \mathbf{W} terms. We provide these individual blocks below and calculations in the following subsections. 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 we have,

$$\begin{aligned} \nabla_{v_\ell, v_m}^2 \mathcal{L}(\boldsymbol{\theta}) &= \frac{\|\mathbf{w}_\ell\| \|\mathbf{w}_m\|}{2\pi} ((\pi - \theta_{\ell, m}) \cos \theta_{\ell, m} + \sin \theta_{\ell, m}), \\ \nabla_{v_\ell, \mathbf{w}_m}^2 \mathcal{L}(\boldsymbol{\theta}) &= \begin{cases} \frac{v_\ell \mathbf{w}_\ell^T - \mathbf{a}^T}{2} + \sum_{i=1}^k \frac{v_i \|\mathbf{w}_i\|}{2\pi} ((\pi - \theta_{\ell, i}) \bar{\mathbf{w}}_i^T + \sin \theta_{\ell, i} \bar{\mathbf{w}}_\ell^T) & \ell = m \\ \frac{v_m \|\mathbf{w}_\ell\|}{2\pi} ((\pi - \theta_{\ell, m}) \bar{\mathbf{w}}_\ell^T + \sin \theta_{\ell, m} \bar{\mathbf{w}}_m^T) & \ell \neq m \end{cases}, \\ \nabla_{\mathbf{w}_\ell, \mathbf{w}_m}^2 \mathcal{L}(\boldsymbol{\theta}) &= \begin{cases} \frac{v_\ell^2}{2} \mathbf{I} + \frac{v_\ell}{2\pi \|\mathbf{w}_\ell\|} \sum_{i=1}^k v_i \|\mathbf{w}_i\| \sin(\theta_{\ell, i}) \left(\mathbf{P}_{\mathbf{w}_l^\perp} + \frac{\mathbf{P}_{\mathbf{w}_l^\perp} \bar{\mathbf{w}}_i \bar{\mathbf{w}}_i^T \mathbf{P}_{\mathbf{w}_l^\perp}}{\|\mathbf{P}_{\mathbf{w}_l^\perp} \bar{\mathbf{w}}_i\|^2} \right) & \ell = m \\ \frac{v_\ell v_m}{2\pi} (\bar{\mathbf{w}}_\ell \bar{\mathbf{w}}_{m, \ell^\perp}^T + \bar{\mathbf{w}}_m \bar{\mathbf{w}}_{\ell, m^\perp}^T + (\pi - \theta_{\ell, m}) \mathbf{I}) & \ell \neq m \end{cases} \end{aligned} \quad (14)$$

A.3.1. CALCULATING THE \mathbf{v}, \mathbf{v} BLOCK

We have,

$$\begin{aligned} \nabla_{v_\ell, v_m}^2 \mathcal{L}(\boldsymbol{\theta}) &= \mathbb{E}_{\mathbf{x}} \left[(f(\boldsymbol{\theta}; \mathbf{x}) - \mathbf{a}^T \mathbf{x}) \nabla_{v_\ell, v_m}^2 f(\boldsymbol{\theta}; \mathbf{x}) + \nabla_{v_\ell} f(\boldsymbol{\theta}; \mathbf{x}) \nabla_{v_m} f(\boldsymbol{\theta}; \mathbf{x})^T \right] \\ &= \mathbb{E}_{\mathbf{x}} \left[\nabla_{v_\ell} f(\boldsymbol{\theta}; \mathbf{x}) \nabla_{v_m} f(\boldsymbol{\theta}; \mathbf{x})^T \right] \\ &= \mathbb{E}_{\mathbf{x}} \left[\phi(\mathbf{w}_\ell^T \mathbf{x}) \phi(\mathbf{w}_m^T \mathbf{x}) \right] \\ &= \frac{\|\mathbf{w}_\ell\| \|\mathbf{w}_m\|}{2\pi} ((\pi - \theta_{\ell, m}) \cos \theta_{\ell, m} + \sin \theta_{\ell, m}) \end{aligned}$$

where the last step follows from the Table 1 in (Daniely et al., 2016).

A.3.2. CALCULATING THE $\mathbf{v}, \text{VECT}(\mathbf{W})$ BLOCK

We have,

$$\nabla_{v_\ell, \mathbf{w}_m}^2 \mathcal{L}(\boldsymbol{\theta}) = \mathbb{E}_{\mathbf{x}} \left[(f(\boldsymbol{\theta}; \mathbf{x}) - \mathbf{a}^T \mathbf{x}) \nabla_{v_\ell, \mathbf{w}_m}^2 f(\boldsymbol{\theta}; \mathbf{x}) + \nabla_{v_\ell} f(\boldsymbol{\theta}; \mathbf{x}) \nabla_{\mathbf{w}_m} f(\boldsymbol{\theta}; \mathbf{x})^T \right].$$

for calculation of individual terms refer down below.

Calculating the $\mathbb{E}_{\mathbf{x}} \left[\nabla_{v_\ell} f(\boldsymbol{\theta}; \mathbf{x}) \nabla_{\mathbf{w}_m} f(\boldsymbol{\theta}; \mathbf{x})^T \right]$ term: We have,

$$\begin{aligned}
 \mathbb{E}_{\mathbf{x}} \left[\nabla_{v_\ell} f(\boldsymbol{\theta}; \mathbf{x}) \nabla_{\mathbf{w}_m} f(\boldsymbol{\theta}; \mathbf{x})^T \right] &= \mathbb{E}_{\mathbf{x}} \left[\phi(\mathbf{w}_\ell^T \mathbf{x}) v_m \phi'(\mathbf{w}_m^T \mathbf{x}) \mathbf{x}^T \right] \\
 &= v_m \|\mathbf{w}_\ell\| \mathbb{E}_{\mathbf{x}} \left[\phi(\bar{\mathbf{w}}_\ell^T \mathbf{x}) \phi'(\bar{\mathbf{w}}_m^T \mathbf{x}) \mathbf{x}^T \right] \\
 &\stackrel{(a)}{=} v_m \|\mathbf{w}_\ell\| \mathbb{E}_{\mathbf{x}} \left[\nabla_{\mathbf{x}} \left(\phi(\bar{\mathbf{w}}_\ell^T \mathbf{x}) \phi'(\bar{\mathbf{w}}_m^T \mathbf{x}) \right) \right] \\
 &= v_m \|\mathbf{w}_\ell\| \mathbb{E}_{\mathbf{x}} \left[\phi'(\bar{\mathbf{w}}_\ell^T \mathbf{x}) \phi'(\bar{\mathbf{w}}_m^T \mathbf{x}) \bar{\mathbf{w}}_\ell^T + \phi(\bar{\mathbf{w}}_\ell^T \mathbf{x}) \delta(\bar{\mathbf{w}}_m^T \mathbf{x}) \bar{\mathbf{w}}_m^T \right] \\
 &\stackrel{(b)}{=} v_m \|\mathbf{w}_\ell\| \left(\left(\frac{\pi - \theta_{\ell,m}}{2\pi} \right) \bar{\mathbf{w}}_\ell^T + \mathbb{E}_{\mathbf{x}} \left[\phi(\bar{\mathbf{w}}_\ell^T \mathbf{x}) \delta(\bar{\mathbf{w}}_m^T \mathbf{x}) \bar{\mathbf{w}}_m^T \right] \right)
 \end{aligned}$$

where equation (a) follows from Stein's Lemma, and (b) follows from the dual activation of a step function. To handle the remaining expectation, we first define $g = \bar{\mathbf{w}}_m^T \mathbf{x} \sim \mathcal{N}(0, 1)$. Then,

$$\begin{aligned}
 \mathbb{E}_{\mathbf{x}} \left[\phi(\bar{\mathbf{w}}_\ell^T \mathbf{x}) \delta(\bar{\mathbf{w}}_m^T \mathbf{x}) \right] &= \mathbb{E}_{\mathbf{x}} \left[\phi \left(\bar{\mathbf{w}}_\ell^T \mathbf{P}_{\mathbf{w}_m^\perp} \mathbf{x} + \bar{\mathbf{w}}_\ell^T \bar{\mathbf{w}}_m g \right) \delta(g) \right] \\
 &= \frac{1}{\sqrt{2\pi}} \mathbb{E}_{\mathbf{x}} \left[\phi \left(\bar{\mathbf{w}}_\ell^T \mathbf{P}_{\mathbf{w}_m^\perp} \mathbf{x} \right) \right] \quad (\text{Delta integration}) \\
 &= \frac{\left\| \mathbf{P}_{\mathbf{w}_m^\perp} \bar{\mathbf{w}}_\ell \right\|}{\sqrt{2\pi}} \mathbb{E}_u \left[\phi(u) \right] \quad (u \sim N(0, 1)) \\
 &= \frac{\sin \theta_{\ell,m}}{2\pi} \left(\text{Expectation of rectified Gaussian } f_x(0) = \frac{1}{\sqrt{2\pi}} \right).
 \end{aligned}$$

Combining everything:

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

Calculating the $\mathbb{E}_{\mathbf{x}} \left[(f(\boldsymbol{\theta}; \mathbf{x}) - \mathbf{a}^T \mathbf{x}) \nabla_{v_\ell, \mathbf{w}_m}^2 f(\boldsymbol{\theta}; \mathbf{x}) \right]$ term: Note that

$$\nabla_{v_\ell, \mathbf{w}_m}^2 f(\boldsymbol{\theta}; \mathbf{x}) = \mathbb{1} \{ \ell = m \} \phi'(\bar{\mathbf{w}}_\ell^T \mathbf{x}) \mathbf{x}^T.$$

Hence we focus only on $\ell = m$ case.

$$\begin{aligned}
 \mathbb{E}_{\mathbf{x}} \left[(f(\boldsymbol{\theta}; \mathbf{x}) - \mathbf{a}^T \mathbf{x}) \nabla_{v_\ell, \mathbf{w}_m}^2 f(\boldsymbol{\theta}; \mathbf{x}) \right] &= \mathbb{E}_{\mathbf{x}} \left[r(\mathbf{x}) \phi'(\bar{\mathbf{w}}_\ell^T \mathbf{x}) \mathbf{x}^T \right] \\
 &\stackrel{(a)}{=} \mathbb{E}_{\mathbf{x}} \left[\nabla_{\mathbf{x}} r(\mathbf{x}) \phi'(\bar{\mathbf{w}}_\ell^T \mathbf{x}) + r(\mathbf{x}) \delta(\bar{\mathbf{w}}_\ell^T \mathbf{x}) \bar{\mathbf{w}}_\ell^T \right] \\
 &= \mathbb{E}_{\mathbf{x}} \left[\left(\sum_{i=1}^k v_i \phi'(\mathbf{w}_i^T \mathbf{x}) \mathbf{w}_i - \mathbf{a} \right)^T \phi'(\bar{\mathbf{w}}_\ell^T \mathbf{x}) + r(\mathbf{x}) \delta(\bar{\mathbf{w}}_\ell^T \mathbf{x}) \bar{\mathbf{w}}_\ell^T \right] \\
 &\stackrel{(b)}{=} \sum_{i=1}^k v_i \mathbf{w}_i^T \left(\frac{\pi - \theta_{\ell,i}}{2\pi} \right) - \frac{\mathbf{a}}{2} + \mathbb{E}_{\mathbf{x}} \left[r(\mathbf{x}) \delta(\bar{\mathbf{w}}_\ell^T \mathbf{x}) \right] \bar{\mathbf{w}}_\ell^T,
 \end{aligned}$$

where (a) follows from the Stein's identity, and (b) follows from the dual activation of step function. To handle the remaining expectation term, define $g = \bar{\mathbf{w}}_\ell^T \mathbf{x} \sim \mathcal{N}(0, 1)$. Then,

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

Combining everything,

$$\mathbb{E}_{\mathbf{x}} \left[(f(\boldsymbol{\theta}; \mathbf{x}) - \mathbf{a}^T \mathbf{x}) \nabla_{\mathbf{v}_\ell, \mathbf{w}_m}^2 f(\boldsymbol{\theta}; \mathbf{x}) \right] = \sum_{i=1}^k \frac{\mathbf{v}_i \|\mathbf{w}_i\|}{2\pi} \left((\pi - \theta_{\ell,i}) \bar{\mathbf{w}}_i^T + \sin \theta_{\ell,i} \bar{\mathbf{w}}_\ell^T \right) - \frac{\mathbf{a}^T}{2},$$

when $\ell = m$. Otherwise, this term is $\mathbf{0}$.

A.3.3. CALCULATING THE $\text{VECT}(\mathbf{W})$, $\text{VECT}(\mathbf{W})$ BLOCK

We have,

$$\nabla_{\mathbf{w}_\ell, \mathbf{w}_m}^2 \mathcal{L}(\boldsymbol{\theta}) = \mathbb{E}_{\mathbf{x}} \left[(f(\boldsymbol{\theta}; \mathbf{x}) - \mathbf{a}^T \mathbf{x}) \nabla_{\mathbf{w}_\ell, \mathbf{w}_m}^2 f(\boldsymbol{\theta}; \mathbf{x}) + \nabla_{\mathbf{w}_\ell} f(\boldsymbol{\theta}; \mathbf{x}) \nabla_{\mathbf{w}_m} f(\boldsymbol{\theta}; \mathbf{x})^T \right].$$

for calculation of individual terms refer down below.

Calculating the $\mathbb{E}_{\mathbf{x}} \left[\nabla_{\mathbf{w}_\ell} f(\boldsymbol{\theta}; \mathbf{x}) \nabla_{\mathbf{w}_m} f(\boldsymbol{\theta}; \mathbf{x})^T \right]$ term: We have,

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

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

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

First term is:

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

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 $\cos(\theta_{i,j}) = \bar{\mathbf{w}}_i^T \bar{\mathbf{w}}_j$):

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

Calculating the $\mathbb{E}_{\mathbf{x}} [(f(\boldsymbol{\theta}; \mathbf{x}) - \mathbf{a}^T \mathbf{x}) \nabla_{\mathbf{w}_{\ell}, \mathbf{w}_m}^2 f(\boldsymbol{\theta}; \mathbf{x})]$ term: 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 $g = \bar{\mathbf{w}}_{\ell}^T \mathbf{x} \sim \mathcal{N}(0, 1)$.

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

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

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

First term is:

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

Second term is:

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

Combining both terms we get

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

Finally we plug this back to get:

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

Appendix B. Proof of Key Lemmas in the Population Setting

For the simplicity of notation, let $\mathbf{w}_1 = \mathbf{w}_1^{(\tau)}$, $\mathbf{w}_2 = \mathbf{w}_2^{(\tau)}$.

B.1. Proof of Gradient Smoothness Towards the Global Optima in the Population Case (Lemma 12)

To establish the imbalance bound (Lemma 3), we first introduce a key lemma that characterizes the gradient smoothness toward the global optima in the population case. This result relates the norm of the population gradient to the relative distance between the current parameters and the global optima:

Lemma 12 *Under the constraint $v_1 = v_2 = 1$, the following inequality holds for all $\mathbf{w}_1, \mathbf{w}_2 \in \mathbb{R}^d$:*

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

Proof

We begin by demonstrating that:

$$\|\sin \theta \|\mathbf{w}_2\| \bar{\mathbf{w}}_1 + (\pi - \theta) \mathbf{w}_2\| \leq \pi \|\mathbf{w}_1 + \mathbf{w}_2\|. \quad (15)$$

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 θ . 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. 15. Note that due to symmetry we get the following as a corollary:

$$\|\sin \theta \|\mathbf{w}_1\| \bar{\mathbf{w}}_2 + (\pi - \theta) \mathbf{w}_1\| \leq \pi \|\mathbf{w}_1 + \mathbf{w}_2\|. \quad (16)$$

Under the constraint $v_1 = v_2 = 1$, the partial gradients with respect to \mathbf{w}_1 and \mathbf{w}_2 are given separately by:

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

Using Ineq. 15, 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. 16 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 12. ■

B.2. Proof of the PL Inequality in the Population Case (Lemma 8)

First, we show it is sufficient to analyze $v_1 = v_2 = 1$. For $v_1, v_2 > 0$, we define $\tilde{\mathbf{w}}_i = v_i \mathbf{w}_i$. Then,

$$\mathcal{L}(\boldsymbol{\theta}) = \mathbb{E}_{\mathbf{x}} \left[(v_1 \text{ReLU}(\mathbf{w}_1^T \mathbf{x}) - v_2 \text{ReLU}(\mathbf{w}_2^T \mathbf{x}) - \mathbf{a}^T \mathbf{x})^2 \right] = \mathbb{E}_{\mathbf{x}} \left[(\text{ReLU}(\tilde{\mathbf{w}}_1^T \mathbf{x}) - \text{ReLU}(\tilde{\mathbf{w}}_2^T \mathbf{x}) - \mathbf{a}^T \mathbf{x})^2 \right].$$

Let us focus on squared gradient norms:

$$\begin{aligned} \|\nabla_{\mathbf{w}_1} \mathcal{L}\|^2 + \|\nabla_{\mathbf{w}_2} \mathcal{L}\|^2 &= \|v_1 \nabla_{\tilde{\mathbf{w}}_1} \mathcal{L}\|^2 + \|v_2 \nabla_{\tilde{\mathbf{w}}_2} \mathcal{L}\|^2 \\ &\geq \min(v_1^2, v_2^2) \left(\|\nabla_{\tilde{\mathbf{w}}_1} \mathcal{L}\|^2 + \|\nabla_{\tilde{\mathbf{w}}_2} \mathcal{L}\|^2 \right) \end{aligned}$$

This suggests that proving $\|\nabla_{\mathbf{w}_1} \mathcal{L}\|^2 + \|\nabla_{\mathbf{w}_2} \mathcal{L}\|^2 \geq \alpha \mathcal{L}$ when $v_1 = v_2 = 1$ implies that

$$\|\nabla_{\mathbf{w}_1} \mathcal{L}\|^2 + \|\nabla_{\mathbf{w}_2} \mathcal{L}\|^2 \geq \min(v_1^2, v_2^2) \alpha \mathcal{L}$$

for arbitrary $v_1, v_2 > 0$. Now, we assume $v_1 = v_2 = 1$. We define

$$h(\mathbf{w}_1, \mathbf{w}_2, \mathbf{a}) = \|\nabla_{\mathbf{w}_1} \mathcal{L}\|^2 + \|\nabla_{\mathbf{w}_2} \mathcal{L}\|^2 - \alpha \mathcal{L}.$$

Using the gradient calculations in (11), we can write it equivalently as

$$\begin{aligned} h(\mathbf{w}_1, \mathbf{w}_2, \mathbf{a}) &= \frac{1}{4} \left(1 - \alpha + \frac{(\pi - \theta)^2 + 2(\pi - \theta) \sin \theta \cos \theta + (\sin \theta)^2}{\pi^2} \right) \left(\|\mathbf{w}_1\|^2 + \|\mathbf{w}_2\|^2 \right) \\ &\quad - \left(1 - \frac{\alpha}{2} \right) \frac{(\pi - \theta) \cos \theta + \sin \theta}{\pi} \|\mathbf{w}_1\| \|\mathbf{w}_2\| \\ &\quad - \frac{1}{2} \left(\left(\frac{\pi - \theta}{\pi} + 1 - \alpha \right) \|\mathbf{w}_1\| - \frac{\sin \theta}{\pi} \|\mathbf{w}_2\| \right) \tilde{\mathbf{w}}_1^T \mathbf{a} \\ &\quad + \frac{1}{2} \left(\left(\frac{\pi - \theta}{\pi} + 1 - \alpha \right) \|\mathbf{w}_2\| - \frac{\sin \theta}{\pi} \|\mathbf{w}_1\| \right) \tilde{\mathbf{w}}_2^T \mathbf{a} \\ &\quad + \left(\frac{1 - \alpha}{2} \right) \|\mathbf{a}\|^2 \\ &= \frac{1 - \alpha}{2} \|\mathbf{a}\|^2 + \mathbf{b}^T \mathbf{a} + c \end{aligned}$$

where \mathbf{b} and c are defined by the following terms for brevity,

$$\begin{aligned}\alpha_1 &= -\frac{1}{2} \left(\left(\frac{\pi - \theta}{\pi} + 1 - \alpha \right) \|\mathbf{w}_1\| - \frac{\sin \theta}{\pi} \|\mathbf{w}_2\| \right) \\ \alpha_2 &= +\frac{1}{2} \left(\left(\frac{\pi - \theta}{\pi} + 1 - \alpha \right) \|\mathbf{w}_2\| - \frac{\sin \theta}{\pi} \|\mathbf{w}_1\| \right) \\ c &= \frac{1}{4} \left(1 - \alpha + \frac{(\pi - \theta)^2 + 2(\pi - \theta) \sin \theta \cos \theta + (\sin \theta)^2}{\pi^2} \right) (\|\mathbf{w}_1\|^2 + \|\mathbf{w}_2\|^2) \\ &\quad - \left(1 - \frac{\alpha}{2} \right) \frac{(\pi - \theta) \cos \theta + \sin \theta}{\pi} \|\mathbf{w}_1\| \|\mathbf{w}_2\| \\ \mathbf{b} &= \alpha_1 \bar{\mathbf{w}}_1 + \alpha_2 \bar{\mathbf{w}}_2\end{aligned}$$

Noting that the expression above is quadratic in \mathbf{a} , we compute $\tilde{h}(\mathbf{w}_1, \mathbf{w}_2) = \min_{\mathbf{a}} h(\mathbf{w}_1, \mathbf{w}_2, \mathbf{a})$.

The choice of \mathbf{a} that minimizes the expression is $\mathbf{a} = -\frac{\mathbf{b}}{1-\alpha}$. Plugging this in back we get,

$$\tilde{h}(\mathbf{w}_1, \mathbf{w}_2) = c - \frac{\|\mathbf{b}\|^2}{2(1-\alpha)} = c - \frac{\alpha_1^2 + 2\alpha_1\alpha_2 \cos \theta + \alpha_2^2}{2(1-\alpha)}$$

Taking the norm out: Note that we are only interested in the positivity of \tilde{h} , therefore dividing it by $\|\mathbf{w}_2\|^2$ does not change the sign. Denote $\frac{\|\mathbf{w}_1\|}{\|\mathbf{w}_2\|} = r$. Then we still have

$$\frac{\tilde{h}(\mathbf{w}_1, \mathbf{w}_2)}{\|\mathbf{w}_2\|^2} = c - \frac{\alpha_1^2 + 2\alpha_1\alpha_2 \cos \theta + \alpha_2^2}{2(1-\alpha)}$$

but the variables are modified as

$$\begin{aligned}\alpha_1 &= -\frac{1}{2} \left(\left(\frac{\pi - \theta}{\pi} + 1 - \alpha \right) r - \frac{\sin \theta}{\pi} \right) \\ \alpha_2 &= +\frac{1}{2} \left(\left(\frac{\pi - \theta}{\pi} + 1 - \alpha \right) - \frac{\sin \theta}{\pi} r \right) \\ c &= \frac{1}{4} \left(1 - \alpha + \frac{(\pi - \theta)^2 + 2(\pi - \theta) \sin \theta \cos \theta + (\sin \theta)^2}{\pi^2} \right) (r^2 + 1) - \left(1 - \frac{\alpha}{2} \right) \frac{(\pi - \theta) \cos \theta + \sin \theta}{\pi} r\end{aligned}$$

We note that the expression is of the form $C_1(r^2 + 1) + C_2r$. Without changing the sign, we can take out r outside. Then we notice that the minima is achieved at $r = 1$. Therefore, it is sufficient for us to check the positivity of the expression at $r = 1$. That is, we draw:

$$\left. \frac{\tilde{h}(\mathbf{w}_1, \mathbf{w}_2)}{\|\mathbf{w}_2\|^2} \right|_{r=1} = c - \frac{(1 - \cos \theta)}{(1 - \alpha)} \tilde{\alpha}^2$$

where

$$\begin{aligned}\tilde{\alpha} &= \frac{1}{2} \left(\frac{\pi - \theta - \sin \theta}{\pi} + 1 - \alpha \right), \\ c &= \frac{1}{2} \left(1 - \alpha + \frac{(\pi - \theta)^2 + 2(\pi - \theta) \sin \theta \cos \theta + (\sin \theta)^2}{\pi^2} \right) - \left(1 - \frac{\alpha}{2} \right) \frac{(\pi - \theta) \cos \theta + \sin \theta}{\pi}.\end{aligned}$$

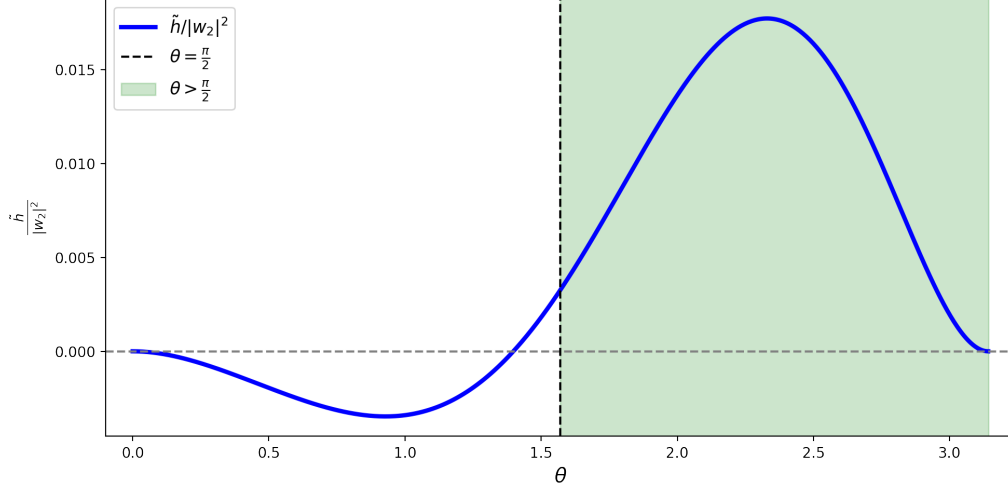


Figure 4: $\|\nabla_{w_1} \mathcal{L}\|^2 + \|\nabla_{w_2} \mathcal{L}\|^2 - \alpha \mathcal{L}$ is non-negative. We set $\alpha = 0.05$ and draw $\frac{1}{\|w_2\|} \tilde{h}(w_1, w_2)$ for $\theta \in [0, \pi]$. We show that $\frac{1}{\|w_2\|} \tilde{h}(w_1, w_2)$ is non-negative inside the shaded region ($\theta \geq \frac{\pi}{2}$).

To complete the proof, in Figure 4, we set $\alpha = 0.05$ and draw $\left. \frac{\tilde{h}(w_1, w_2)}{\|w_2\|^2} \right|_{r=1}$ as a 1D plot for $\theta \in [0, \pi]$. The plot demonstrates that \tilde{h} is non-negative for $\theta > \frac{\pi}{2}$. This finishes the proof.

B.3. Bound on the Smoothness of the Population Loss (Lemma 9)

We bound the population Hessian $\nabla^2 \mathcal{L}(v, \mathbf{W})$ in the local refinement phase. That is, we assume $c_1 \sqrt{\|\mathbf{a}\|} \leq v_1, v_2, \|w_1\|, \|w_2\| \leq c_2 \sqrt{\|\mathbf{a}\|}$. By the sub-additivity properties of the spectral norm, we have

$$\|\nabla^2 \mathcal{L}(v, \mathbf{W})\|_2 \leq \|\nabla_{v,v}^2 \mathcal{L}(v, \mathbf{W})\|_2 + 2 \left\| \nabla_{v, \text{vect}(\mathbf{W})}^2 \mathcal{L}(v, \mathbf{W}) \right\|_2 + \left\| \nabla_{\text{vect}(\mathbf{W}), \text{vect}(\mathbf{W})}^2 \mathcal{L}(v, \mathbf{W}) \right\|_2.$$

We bound each term separately below.

$\nabla_{v,v}^2 \mathcal{L}(\theta)$ term: We have

$$\|\nabla_{v,v}^2 \mathcal{L}(\theta)\|_2 \leq \sum_{\ell=1}^2 \sum_{m=1}^2 |\nabla_{v_\ell, v_m}^2 \mathcal{L}(\theta)|.$$

where,

$$\begin{aligned} |\nabla_{v_\ell, v_m}^2 \mathcal{L}(\theta)| &= \frac{\|w_\ell\| \|w_m\|}{2\pi} ((\pi - \theta_{\ell,m}) \cos \theta_{\ell,m} + \sin \theta_{\ell,m}) \\ &\leq \frac{c_2^2 \|\mathbf{a}\|}{2\pi} ((\pi - \theta_{\ell,m}) \cos \theta_{\ell,m} + \sin \theta_{\ell,m}) \\ &\leq \frac{c_2^2 \|\mathbf{a}\|}{2\pi} (\pi + 1) \\ &\leq c_2^2 \|\mathbf{a}\| \frac{1 + \pi}{2\pi}. \end{aligned}$$

Then,

$$\begin{aligned} \|\nabla_{\mathbf{v}, \mathbf{v}}^2 \mathcal{L}(\boldsymbol{\theta})\|_2 &\leq \sum_{\ell=1}^2 \sum_{m=1}^2 |\nabla_{v_\ell, v_m}^2 \mathcal{L}(\boldsymbol{\theta})| \\ &\leq 4c_2^2 \|\mathbf{a}\| \frac{1+\pi}{2\pi} = \left(2 + \frac{2}{\pi}\right) c_2^2 \|\mathbf{a}\|. \end{aligned}$$

$\nabla_{\mathbf{v}, \text{vect}(\mathbf{W})}^2 \mathcal{L}(\boldsymbol{\theta})$ term: We have

$$\left\| \nabla_{\mathbf{v}, \text{vect}(\mathbf{W})}^2 \mathcal{L}(\boldsymbol{\theta}) \right\|_2 \leq \sum_{\ell=1}^2 \sum_{m=1}^2 \|\nabla_{v_\ell, \mathbf{w}_m}^2 \mathcal{L}(\boldsymbol{\theta})\|.$$

For $\ell \neq m$, we have

$$\begin{aligned} \|\nabla_{v_\ell, \mathbf{w}_m}^2 \mathcal{L}(\boldsymbol{\theta})\| &= \left\| \frac{v_m \|\mathbf{w}_\ell\|}{2\pi} \left((\pi - \theta_{\ell, m}) \bar{\mathbf{w}}_\ell^T + \sin \theta_{\ell, m} \bar{\mathbf{w}}_m^T \right) \right\| \\ &= \frac{v_m \|\mathbf{w}_\ell\|}{2\pi} \|((\pi - \theta_{\ell, m}) \bar{\mathbf{w}}_\ell + \sin \theta_{\ell, m} \bar{\mathbf{w}}_m)\| \\ &\leq \frac{v_m \|\mathbf{w}_\ell\|}{2\pi} (\|(\pi - \theta_{\ell, m}) \bar{\mathbf{w}}_\ell\| + \|\sin \theta_{\ell, m} \bar{\mathbf{w}}_m\|) \\ &\leq \frac{v_m \|\mathbf{w}_\ell\|}{2\pi} (\pi + 1) \\ &\leq c_2^2 \|\mathbf{a}\| \frac{1+\pi}{2\pi}. \end{aligned}$$

For $\ell = m$, we have

$$\begin{aligned} \|\nabla_{v_\ell, \mathbf{w}_m}^2 \mathcal{L}(\boldsymbol{\theta})\| &= \left\| \frac{v_\ell \mathbf{w}_\ell^T - \mathbf{a}^T}{2} + \sum_{i=1}^2 \frac{v_i \|\mathbf{w}_i\|}{2\pi} \left((\pi - \theta_{\ell, i}) \bar{\mathbf{w}}_i^T + \sin \theta_{\ell, i} \bar{\mathbf{w}}_\ell^T \right) \right\| \\ &\leq \frac{\|v_\ell \mathbf{w}_\ell\|}{2} + \frac{\|\mathbf{a}\|}{2} + \sum_{i=1}^2 \left\| \frac{v_i \|\mathbf{w}_i\|}{2\pi} \left((\pi - \theta_{\ell, i}) \bar{\mathbf{w}}_i^T + \sin \theta_{\ell, i} \bar{\mathbf{w}}_\ell^T \right) \right\| \\ &\leq \frac{\|v_\ell \mathbf{w}_\ell\|}{2} + \frac{\|\mathbf{a}\|}{2} + \sum_{i=1}^2 \frac{v_i \|\mathbf{w}_i\|}{2\pi} (\|(\pi - \theta_{\ell, i}) \bar{\mathbf{w}}_i\| + \|\sin \theta_{\ell, i} \bar{\mathbf{w}}_\ell\|) \\ &\leq \frac{c_2^2 \|\mathbf{a}\|}{2} + \frac{\|\mathbf{a}\|}{2} + \sum_{i=1}^2 \frac{c_2^2 \|\mathbf{a}\|}{2\pi} (\pi + 1) = \left(\frac{1}{2} + \left(\frac{3}{2} + \frac{1}{\pi} \right) c_2^2 \right) \|\mathbf{a}\|. \end{aligned}$$

Combining both inequalities, we have

$$\begin{aligned} \left\| \nabla_{\mathbf{v}, \text{vect}(\mathbf{W})}^2 \mathcal{L}(\boldsymbol{\theta}) \right\|_2 &\leq \sum_{\ell=1}^2 \sum_{m=1}^2 \|\nabla_{v_\ell, \mathbf{w}_m}^2 \mathcal{L}(\boldsymbol{\theta})\| \\ &\leq 2 \left(c_2^2 \|\mathbf{a}\| \frac{1+\pi}{2\pi} + \left(\frac{1}{2} + \left(\frac{3}{2} + \frac{1}{\pi} \right) c_2^2 \right) \|\mathbf{a}\| \right) \\ &= \left(1 + \left(4 + \frac{3}{\pi} \right) c_2^2 \right) \|\mathbf{a}\|. \end{aligned}$$

$\nabla_{\text{vect}(\mathbf{W}), \text{vect}(\mathbf{W})}^2 \mathcal{L}(\boldsymbol{\theta})$ term: We have

$$\left\| \nabla_{\text{vect}(\mathbf{W}), \text{vect}(\mathbf{W})}^2 \mathcal{L}(\boldsymbol{\theta}) \right\|_2 \leq \sum_{\ell=1}^2 \sum_{m=1}^2 \left\| \nabla_{\mathbf{w}_\ell, \mathbf{w}_m}^2 \mathcal{L}(\boldsymbol{\theta}) \right\|.$$

For $\ell \neq m$, we have

$$\begin{aligned} \left\| \nabla_{\mathbf{w}_\ell, \mathbf{w}_m}^2 \mathcal{L}(\boldsymbol{\theta}) \right\|_2 &= \left\| \frac{v_\ell v_m}{2\pi} \left(\bar{\mathbf{w}}_\ell \bar{\mathbf{w}}_{m, \ell^\perp}^T + \bar{\mathbf{w}}_m \bar{\mathbf{w}}_{\ell, m^\perp}^T + (\pi - \theta_{\ell, m}) \mathbf{I} \right) \right\|_2 \\ &= \frac{v_\ell v_m}{2\pi} \left\| \bar{\mathbf{w}}_\ell \bar{\mathbf{w}}_{m, \ell^\perp}^T + \bar{\mathbf{w}}_m \bar{\mathbf{w}}_{\ell, m^\perp}^T + (\pi - \theta_{\ell, m}) \mathbf{I} \right\|_2 \\ &\leq \frac{v_\ell v_m}{2\pi} \left(\left\| \bar{\mathbf{w}}_\ell \bar{\mathbf{w}}_{m, \ell^\perp}^T \right\|_2 + \left\| \bar{\mathbf{w}}_m \bar{\mathbf{w}}_{\ell, m^\perp}^T \right\|_2 + \left\| (\pi - \theta_{\ell, m}) \mathbf{I} \right\|_2 \right) \\ &\leq \frac{v_\ell v_m}{2\pi} (1 + 1 + \pi) \\ &= v_\ell v_m \frac{2 + \pi}{2\pi} \\ &\leq c_2^2 \|\mathbf{a}\| \frac{2 + \pi}{2\pi}. \end{aligned}$$

For $\ell = m$, we have

$$\begin{aligned} \left\| \nabla_{\mathbf{w}_\ell, \mathbf{w}_m}^2 \mathcal{L}(\boldsymbol{\theta}) \right\|_2 &= \left\| \frac{v_\ell^2}{2} \mathbf{I} + \frac{v_\ell}{2\pi \|\mathbf{w}_\ell\|} \sum_{i=1}^2 v_i \|\mathbf{w}_i\| \sin(\theta_{\ell, i}) \left(\mathbf{P}_{\mathbf{w}_i^\perp} + \frac{\mathbf{P}_{\mathbf{w}_i^\perp} \bar{\mathbf{w}}_i \bar{\mathbf{w}}_i^T \mathbf{P}_{\mathbf{w}_i^\perp}}{\left\| \mathbf{P}_{\mathbf{w}_i^\perp} \bar{\mathbf{w}}_i \right\|^2} \right) \right\|_2 \\ &\leq \left\| \frac{v_\ell^2}{2} \mathbf{I} \right\|_2 + \frac{v_\ell}{2\pi \|\mathbf{w}_\ell\|} \sum_{i=1}^2 \left\| v_i \|\mathbf{w}_i\| \sin(\theta_{\ell, i}) \left(\mathbf{P}_{\mathbf{w}_i^\perp} + \frac{\mathbf{P}_{\mathbf{w}_i^\perp} \bar{\mathbf{w}}_i \bar{\mathbf{w}}_i^T \mathbf{P}_{\mathbf{w}_i^\perp}}{\left\| \mathbf{P}_{\mathbf{w}_i^\perp} \bar{\mathbf{w}}_i \right\|^2} \right) \right\|_2 \\ &\leq \left\| \frac{v_\ell^2}{2} \mathbf{I} \right\|_2 + \frac{v_\ell}{2\pi \|\mathbf{w}_\ell\|} \sum_{i=1}^2 v_i \|\mathbf{w}_i\| \sin(\theta_{\ell, i}) \left(\left\| \mathbf{P}_{\mathbf{w}_i^\perp} \right\| + \left\| \frac{\mathbf{P}_{\mathbf{w}_i^\perp} \bar{\mathbf{w}}_i \bar{\mathbf{w}}_i^T \mathbf{P}_{\mathbf{w}_i^\perp}}{\left\| \mathbf{P}_{\mathbf{w}_i^\perp} \bar{\mathbf{w}}_i \right\|^2} \right\|_2 \right) \\ &\leq \frac{v_\ell^2}{2} + \frac{v_\ell}{2\pi \|\mathbf{w}_\ell\|} \sum_{i=1}^2 v_i \|\mathbf{w}_i\| (1 + 1) \\ &\leq \frac{c_2^2 \|\mathbf{a}\|}{2} + \frac{c_2}{\pi c_1} (2c_2^2 \|\mathbf{a}\|) = c_2^2 \|\mathbf{a}\| \left(\frac{1}{2} + \frac{2c_2}{\pi c_1} \right). \end{aligned}$$

Combining both inequalities, we have

$$\begin{aligned} \left\| \nabla_{\text{vect}(\mathbf{W}), \text{vect}(\mathbf{W})}^2 \mathcal{L}(\boldsymbol{\theta}) \right\|_2 &\leq \sum_{\ell=1}^2 \sum_{m=1}^2 \left\| \nabla_{\mathbf{w}_\ell, \mathbf{w}_m}^2 \mathcal{L}(\boldsymbol{\theta}) \right\| \\ &\leq 2 \left(c_2^2 \|\mathbf{a}\| \frac{2 + \pi}{2\pi} + c_2^2 \|\mathbf{a}\| \left(\frac{1}{2} + \frac{2c_2}{\pi c_1} \right) \right) \\ &= \left(2 + \frac{2}{\pi} + \frac{4c_2}{\pi c_1} \right) c_2^2 \|\mathbf{a}\|. \end{aligned}$$

Combining the terms: Putting everything together,

$$\begin{aligned}
 \|\nabla^2 \mathcal{L}(\mathbf{v}, \mathbf{W})\|_2 &\leq \|\nabla_{\mathbf{v}, \mathbf{v}}^2 \mathcal{L}(\mathbf{v}, \mathbf{W})\|_2 + 2 \left\| \nabla_{\mathbf{v}, \text{vect}(\mathbf{W})}^2 \mathcal{L}(\mathbf{v}, \mathbf{W}) \right\|_2 + \left\| \nabla_{\text{vect}(\mathbf{W}), \text{vect}(\mathbf{W})}^2 \mathcal{L}(\mathbf{v}, \mathbf{W}) \right\|_2 \\
 &\leq \left(2 + \frac{2}{\pi}\right) c_2^2 \|\mathbf{a}\| + 2 \left(1 + \left(4 + \frac{3}{\pi}\right) c_2^2\right) \|\mathbf{a}\| + \left(2 + \frac{2}{\pi} + \frac{4c_2}{\pi c_1}\right) c_2^2 \|\mathbf{a}\| \\
 &= \left(2 + \left(12 + \frac{10}{\pi} + \frac{4c_2}{\pi c_1}\right) c_2^2\right) \|\mathbf{a}\| \\
 &:= L \|\mathbf{a}\|.
 \end{aligned}$$

This completes the proof of Lemma 9.

B.4. Population Loss Lower Bound (Lemma 13)

Lemma 13 (Population Loss Lower Bound) For $v_1, v_2 > 0$ and $\theta > \frac{\pi}{2}$. We have

$$\|v_1 \mathbf{w}_1 - \mathbf{a}\|^2 + \|v_2 \mathbf{w}_2 + \mathbf{a}\|^2 \leq 20 \mathcal{L}(\mathbf{v}, \mathbf{W}).$$

Proof Define $\tilde{\mathbf{w}}_i = v_i \mathbf{w}_i$. For $v_1, v_2 > 0$, both $\|v_1 \mathbf{w}_1 - \mathbf{a}\|^2 + \|v_2 \mathbf{w}_2 + \mathbf{a}\|^2$ and $\mathcal{L}(\mathbf{v}, \mathbf{W})$ are only functions of $\tilde{\mathbf{w}}_1, \tilde{\mathbf{w}}_2$. Next, we define

$$h(\tilde{\mathbf{w}}_1, \tilde{\mathbf{w}}_2, \mathbf{a}) = \mathcal{L}(\tilde{\mathbf{w}}_1, \tilde{\mathbf{w}}_2) - \tilde{\alpha} \left(\|\tilde{\mathbf{w}}_1 - \mathbf{a}\|^2 + \|\tilde{\mathbf{w}}_2 + \mathbf{a}\|^2 \right).$$

We can write it equivalently as

$$\begin{aligned}
 h(\tilde{\mathbf{w}}_1, \tilde{\mathbf{w}}_2, \mathbf{a}) &= \left(\frac{1}{4} - \tilde{\alpha}\right) \left(\|\tilde{\mathbf{w}}_1\|^2 + \|\tilde{\mathbf{w}}_2\|^2 \right) - \frac{(\pi - \theta) \cos \theta + \sin \theta}{2\pi} \|\tilde{\mathbf{w}}_1\| \|\tilde{\mathbf{w}}_2\| \\
 &\quad - \left(\frac{1}{2} - 2\tilde{\alpha}\right) \mathbf{a}^T (\tilde{\mathbf{w}}_1 - \tilde{\mathbf{w}}_2) + \left(\frac{1}{2} - 2\tilde{\alpha}\right) \|\mathbf{a}\|^2
 \end{aligned}$$

Noting that the expression above is quadratic in \mathbf{a} , we compute $\tilde{h}(\tilde{\mathbf{w}}_1, \tilde{\mathbf{w}}_2) = \min_{\mathbf{a}} h(\tilde{\mathbf{w}}_1, \tilde{\mathbf{w}}_2, \mathbf{a})$.

The choice of \mathbf{a} that minimizes the expression is $\mathbf{a} = \frac{\tilde{\mathbf{w}}_1 - \tilde{\mathbf{w}}_2}{2}$. Plugging this in back we get,

$$\begin{aligned}
 \tilde{h}(\tilde{\mathbf{w}}_1, \tilde{\mathbf{w}}_2) &= \left(\frac{1}{4} - \tilde{\alpha}\right) \left(\|\tilde{\mathbf{w}}_1\|^2 + \|\tilde{\mathbf{w}}_2\|^2 \right) - \frac{(\pi - \theta) \cos \theta + \sin \theta}{2\pi} \|\tilde{\mathbf{w}}_1\| \|\tilde{\mathbf{w}}_2\| \\
 &\quad - \left(\frac{1}{8} - \frac{\tilde{\alpha}}{2}\right) \|\tilde{\mathbf{w}}_1 - \tilde{\mathbf{w}}_2\|^2 \\
 &= \left(\frac{1}{8} - \frac{\tilde{\alpha}}{2}\right) \left(\|\tilde{\mathbf{w}}_1\|^2 + \|\tilde{\mathbf{w}}_2\|^2 \right) - \frac{\left(\pi \left(\frac{1}{2} + 2\tilde{\alpha}\right) - \theta\right) \cos \theta + \sin \theta}{2\pi} \|\tilde{\mathbf{w}}_1\| \|\tilde{\mathbf{w}}_2\|
 \end{aligned}$$

Taking the norm out: Note that we are only interested in the positivity of \tilde{h} , therefore dividing it by $\|\tilde{\mathbf{w}}_2\|^2$ does not change the sign. Denote $\frac{\|\tilde{\mathbf{w}}_1\|}{\|\tilde{\mathbf{w}}_2\|} = r$. Then,

$$\frac{\tilde{h}(\tilde{\mathbf{w}}_1, \tilde{\mathbf{w}}_2)}{\|\tilde{\mathbf{w}}_2\|^2} = \left(\frac{1}{8} - \frac{\tilde{\alpha}}{2}\right) (r^2 + 1) - \frac{\left(\pi \left(\frac{1}{2} + 2\tilde{\alpha}\right) - \theta\right) \cos \theta + \sin \theta}{2\pi} r$$

We note that the expression is of the form $C_1 (r^2 + 1) + C_2 r$. Note that the minima of this expression is achieved at $r = 1$. Therefore, it is sufficient for us to check the positivity of the expression at $r = 1$. To this aim we draw

$$\left. \frac{\tilde{h}(\tilde{\mathbf{w}}_1, \tilde{\mathbf{w}}_2)}{\|\tilde{\mathbf{w}}_2\|^2} \right|_{r=1} = \left(\frac{1}{4} - \tilde{\alpha} \right) - \frac{(\pi (\frac{1}{2} + 2\tilde{\alpha}) - \theta) \cos \theta + \sin \theta}{2\pi}$$

To complete the proof, in Figure 5, we set $\tilde{\alpha} = 0.05$ and draw $\left. \frac{\tilde{h}(\tilde{\mathbf{w}}_1, \tilde{\mathbf{w}}_2)}{\|\tilde{\mathbf{w}}_2\|^2} \right|_{r=1}$ as a 1D plot for

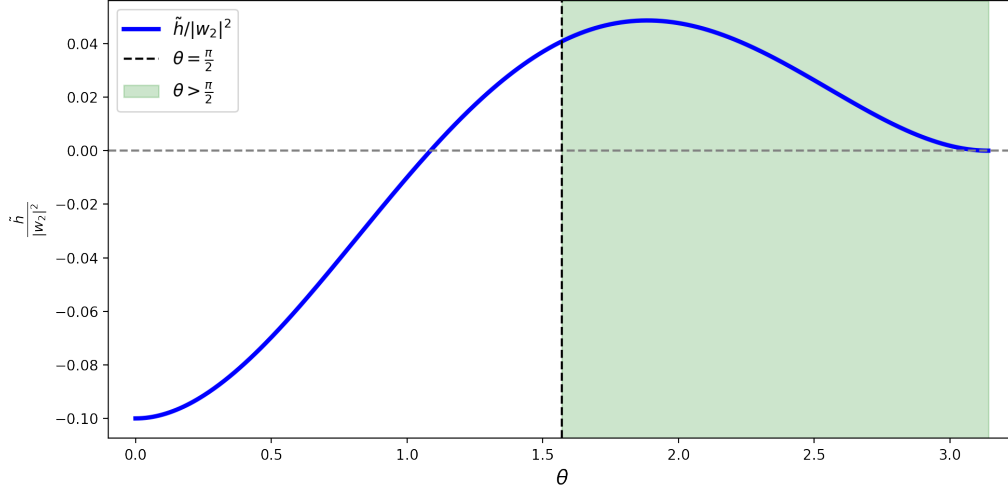


Figure 5: $\mathcal{L}(\mathbf{v}, \mathbf{W}) - \tilde{\alpha} \left(\|v_1 \mathbf{w}_1 - \mathbf{a}\|^2 + \|v_2 \mathbf{w}_2 + \mathbf{a}\|^2 \right)$ is non-negative. We set $\tilde{\alpha} = 0.05$ and draw $\frac{1}{\|\tilde{\mathbf{w}}_2\|^2} \tilde{h}(\tilde{\mathbf{w}}_1, \tilde{\mathbf{w}}_2)$ for $\theta \in [0, \pi]$. We show that $\frac{1}{\|\tilde{\mathbf{w}}_2\|^2} \tilde{h}(\tilde{\mathbf{w}}_1, \tilde{\mathbf{w}}_2)$ is non-negative inside the shaded region ($\theta \geq \frac{\pi}{2}$).

$\theta \in [0, \pi]$. The plot demonstrates that \tilde{h} is non-negative for $\theta > \frac{\pi}{2}$. This finishes the proof of Lemma 13. ■

Appendix C. Proof of Key Lemmas in the Empirical Setting

C.1. Bound for Imbalance Term (Lemma 3)

By symmetry, it suffices to prove the bound for $b_1^{(\tau+1)} - b_1^{(\tau)}$. We first evaluate the per-step change in the imbalance term $b_1^{(\tau)}$. By the update rule of gradient descent, we have

$$\begin{aligned}
 b_1^{(\tau+1)} &= \left\| \mathbf{w}_1^{(\tau+1)} \right\|^2 - \left(v_1^{(\tau+1)} \right)^2 \\
 &= \left\| \mathbf{w}_1 - \mu \nabla_{\mathbf{w}_1} \hat{\mathcal{L}} \right\|^2 - \left(v_1 - \mu \nabla_{v_1} \hat{\mathcal{L}} \right)^2 \\
 &= \left\| \mathbf{w}_1 \right\|^2 - 2\mu \mathbf{w}_1^T \nabla_{\mathbf{w}_1} \hat{\mathcal{L}} + \mu^2 \left\| \nabla_{\mathbf{w}_1} \hat{\mathcal{L}} \right\|^2 - \left(v_1^2 - 2\mu v_1 \nabla_{v_1} \hat{\mathcal{L}} + \mu^2 \left(\nabla_{v_1} \hat{\mathcal{L}} \right)^2 \right) \\
 &\stackrel{(a)}{=} \left\| \mathbf{w}_1 \right\|^2 - v_1^2 + \mu^2 \left\| \nabla_{\mathbf{w}_1} \hat{\mathcal{L}} \right\|^2 - \mu^2 \left(\nabla_{v_1} \hat{\mathcal{L}} \right)^2 \\
 &= b_1^{(\tau)} + \mu^2 \left(\left\| \nabla_{\mathbf{w}_1} \hat{\mathcal{L}} \right\|^2 - \left(\nabla_{v_1} \hat{\mathcal{L}} \right)^2 \right),
 \end{aligned}$$

where (a) follows from Eq. 13. It follows that

$$\begin{aligned}
 \left| b_1^{(\tau+1)} - b_1^{(\tau)} \right| &= \mu^2 \left| \left\| \nabla_{\mathbf{w}_1} \hat{\mathcal{L}} \right\|^2 - \left(\nabla_{v_1} \hat{\mathcal{L}} \right)^2 \right| \\
 &\leq \mu^2 \left(\left\| \nabla_{\mathbf{w}_1} \hat{\mathcal{L}} \right\|^2 + \left(\nabla_{v_1} \hat{\mathcal{L}} \right)^2 \right).
 \end{aligned} \tag{17}$$

To bound the drift, we decompose the empirical gradients into their population counterparts and the associated estimation errors:

$$\begin{aligned}
 &\left(\nabla_{v_1} \hat{\mathcal{L}} \right)^2 + \left\| \nabla_{\mathbf{w}_1} \hat{\mathcal{L}} \right\|^2 \\
 &= \left(\nabla_{v_1} \mathcal{L} + \nabla_{v_1} \hat{\mathcal{L}} - \nabla_{v_1} \mathcal{L} \right)^2 + \left\| \nabla_{\mathbf{w}_1} \mathcal{L} + \nabla_{\mathbf{w}_1} \hat{\mathcal{L}} - \nabla_{\mathbf{w}_1} \mathcal{L} \right\|^2 \\
 &\leq 2 \left(\nabla_{v_1} \mathcal{L} \right)^2 + 2 \left(\nabla_{v_1} \hat{\mathcal{L}} - \nabla_{v_1} \mathcal{L} \right)^2 + 2 \left\| \nabla_{\mathbf{w}_1} \mathcal{L} \right\|^2 + 2 \left\| \nabla_{\mathbf{w}_1} \hat{\mathcal{L}} - \nabla_{\mathbf{w}_1} \mathcal{L} \right\|^2 \\
 &= 2 \left(\left(\nabla_{v_1} \mathcal{L} \right)^2 + \left\| \nabla_{\mathbf{w}_1} \mathcal{L} \right\|^2 \right) + 2 \left(\nabla_{v_1} \hat{\mathcal{L}} - \nabla_{v_1} \mathcal{L} \right)^2 + 2 \left\| \nabla_{\mathbf{w}_1} \hat{\mathcal{L}} - \nabla_{\mathbf{w}_1} \mathcal{L} \right\|^2 \\
 &= 2 \left(\left(\nabla_{v_1} \mathcal{L} \right)^2 + \left\| \nabla_{\mathbf{w}_1} \mathcal{L} \right\|^2 \right) + 2 \left(\frac{\mathbf{w}_1^T}{v_1} \nabla_{\mathbf{w}_1} \hat{\mathcal{L}} - \frac{\mathbf{w}_1^T}{v_1} \nabla_{\mathbf{w}_1} \mathcal{L} \right)^2 + 2 \left\| \nabla_{\mathbf{w}_1} \hat{\mathcal{L}} - \nabla_{\mathbf{w}_1} \mathcal{L} \right\|^2 \\
 &\leq 2 \left(\left(\nabla_{v_1} \mathcal{L} \right)^2 + \left\| \nabla_{\mathbf{w}_1} \mathcal{L} \right\|^2 \right) + 2 \frac{\left\| \mathbf{w}_1 \right\|^2 \left\| \nabla_{\mathbf{w}_1} \hat{\mathcal{L}} - \nabla_{\mathbf{w}_1} \mathcal{L} \right\|^2}{v_1^2} + 2 \left\| \nabla_{\mathbf{w}_1} \hat{\mathcal{L}} - \nabla_{\mathbf{w}_1} \mathcal{L} \right\|^2 \\
 &= 2 \left(\left(\nabla_{v_1} \mathcal{L} \right)^2 + \left\| \nabla_{\mathbf{w}_1} \mathcal{L} \right\|^2 \right) + 2 \left(1 + \frac{\left\| \mathbf{w}_1 \right\|^2}{v_1^2} \right) \left\| \nabla_{\mathbf{w}_1} \hat{\mathcal{L}} - \nabla_{\mathbf{w}_1} \mathcal{L} \right\|^2.
 \end{aligned} \tag{18}$$

Regarding the population component, note that by considering a reparameterized set of weights $\tilde{v}_1 = \tilde{v}_2 = 1$, $\tilde{\mathbf{w}}_1 = v_1 \mathbf{w}_1$, $\tilde{\mathbf{w}}_2 = v_2 \mathbf{w}_2$, we can leverage the smoothness properties of the population loss (Lemma 12):

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

Given $v_1, v_2 > 0$, it holds that $\nabla_{\tilde{\mathbf{w}}_1} \mathcal{L} = \frac{1}{v_1} \nabla_{\mathbf{w}_1} \mathcal{L}$. It then follows that:

$$\begin{aligned}
 (\nabla_{v_1} \mathcal{L})^2 + \|\nabla_{\mathbf{w}_1} \mathcal{L}\|^2 &= \left(\frac{\mathbf{w}_1^T}{v_1} \nabla_{\mathbf{w}_1} \mathcal{L} \right)^2 + \|\nabla_{\mathbf{w}_1} \mathcal{L}\|^2 \\
 &\leq \frac{\|\mathbf{w}_1\|^2 \|\nabla_{\mathbf{w}_1} \mathcal{L}\|^2}{v_1^2} + \|\nabla_{\mathbf{w}_1} \mathcal{L}\|^2 \\
 &= (\|\mathbf{w}_1\|^2 + v_1^2) \|\nabla_{\tilde{\mathbf{w}}_1} \mathcal{L}\|^2 \\
 &\leq (\|\mathbf{w}_1\|^2 + v_1^2) \cdot \frac{5}{2} (\|\tilde{\mathbf{w}}_1 - \mathbf{a}\|^2 + \|\tilde{\mathbf{w}}_2 + \mathbf{a}\|^2) \\
 &= \frac{5}{2} (\|\mathbf{w}_1\|^2 + v_1^2) (\|v_1 \mathbf{w}_1 - \mathbf{a}\|^2 + \|v_2 \mathbf{w}_2 + \mathbf{a}\|^2). \tag{19}
 \end{aligned}$$

As for the second term, Lemma 11 provides the following concentration bound:

$$\left\| \nabla_{\mathbf{w}_1} \hat{\mathcal{L}} - \nabla_{\mathbf{w}_1} \mathcal{L} \right\| \leq v_1 \delta (\|v_1 \mathbf{w}_1 - \mathbf{a}\| + \|v_2 \mathbf{w}_2 + \mathbf{a}\|),$$

where $\delta \leq \frac{1}{2}$ is a constant. Squaring both sides and applying Jensen's inequality leads to:

$$\begin{aligned}
 \left\| \nabla_{\mathbf{w}_1} \hat{\mathcal{L}} - \nabla_{\mathbf{w}_1} \mathcal{L} \right\|^2 &\leq v_1^2 \delta^2 (\|v_1 \mathbf{w}_1 - \mathbf{a}\| + \|v_2 \mathbf{w}_2 + \mathbf{a}\|)^2 \\
 &\leq 2v_1^2 \delta^2 (\|v_1 \mathbf{w}_1 - \mathbf{a}\|^2 + \|v_2 \mathbf{w}_2 + \mathbf{a}\|^2). \tag{20}
 \end{aligned}$$

Substituting the bounds from (19) and (20) into (18), we obtain that

$$\begin{aligned}
 &(\nabla_{v_1} \hat{\mathcal{L}})^2 + \left\| \nabla_{\mathbf{w}_1} \hat{\mathcal{L}} \right\|^2 \\
 &\leq 2 \left((\nabla_{v_1} \mathcal{L})^2 + \|\nabla_{\mathbf{w}_1} \mathcal{L}\|^2 \right) + 2 \left(1 + \frac{\|\mathbf{w}_1\|^2}{v_1^2} \right) \left\| \nabla_{\mathbf{w}_1} \hat{\mathcal{L}} - \nabla_{\mathbf{w}_1} \mathcal{L} \right\|^2 \\
 &\leq 2 \cdot \frac{5}{2} (\|\mathbf{w}_1\|^2 + v_1^2) (\|v_1 \mathbf{w}_1 - \mathbf{a}\|^2 + \|v_2 \mathbf{w}_2 + \mathbf{a}\|^2) \\
 &\quad + 2 \left(1 + \frac{\|\mathbf{w}_1\|^2}{v_1^2} \right) \cdot 2v_1^2 \delta^2 (\|v_1 \mathbf{w}_1 - \mathbf{a}\|^2 + \|v_2 \mathbf{w}_2 + \mathbf{a}\|^2) \\
 &= (5 + 4\delta^2) (\|\mathbf{w}_1\|^2 + v_1^2) (\|v_1 \mathbf{w}_1 - \mathbf{a}\|^2 + \|v_2 \mathbf{w}_2 + \mathbf{a}\|^2).
 \end{aligned}$$

Putting the above inequality into (17), we complete the proof of the lemma with the constant $c_6 = 6$.

C.2. Proof of the Stability of Angles (Lemma 4)

We prove the lemma with the following constants: $c_0 \leq \frac{1}{2}$, $c_2 = 2$, $c_5 = \frac{1}{50}$, $c_4 = \frac{\pi}{20}$, $\gamma = \frac{1}{4}$.

By symmetry, it suffices to prove the bound for $\theta_1^{(\tau+1)}$. By the update rule of gradient descent, we have

$$\begin{aligned}
 & \mathbf{w}_1^{(\tau+1)} \\
 &= \mathbf{w}_1 - \mu \nabla_{\mathbf{w}_1} \widehat{\mathcal{L}} \\
 &= \mathbf{w}_1 - \mu \cdot \frac{v_1}{2} \left(v_1 \mathbf{w}_1 - \left(1 - \frac{\theta}{\pi}\right) v_2 \mathbf{w}_2 - \frac{\sin \theta}{\pi} v_2 \|\mathbf{w}_2\| \bar{\mathbf{w}}_1 - \mathbf{a} + \Delta \mathcal{G}_1 \right) \\
 &= \left(\left(1 - \frac{\mu v_1^2}{2}\right) \|\mathbf{w}_1\| + \frac{\mu v_1}{2} \cdot \frac{\sin \theta}{\pi} v_2 \|\mathbf{w}_2\| \right) \bar{\mathbf{w}}_1 + \frac{\mu v_1}{2} \left(\left(1 - \frac{\theta}{\pi}\right) v_2 \mathbf{w}_2 + \mathbf{a} - \Delta \mathcal{G}_1 \right).
 \end{aligned}$$

Given that $v_1, v_2 > 0$ and $1 - \frac{\mu v_1^2}{2} \geq 1 - \frac{c_0 c_2^2}{2} \geq 0$, the update vector $\mathbf{w}_1 - \mu \nabla_{\mathbf{w}_1} \widehat{\mathcal{L}}$ is a nonnegative linear combination of $\bar{\mathbf{w}}_1$ and the vector

$$\mathbf{q} := \left(1 - \frac{\theta}{\pi}\right) v_2 \mathbf{w}_2 + \mathbf{a} - \Delta \mathcal{G}_1.$$

Recall that the angle of a positive linear combination with a reference vector \mathbf{a} is bounded by the maximum angle of its components. Since we assume $\angle(\bar{\mathbf{w}}_1, \mathbf{a}) \leq c_4$, it suffices to show that $\angle(\mathbf{q}, \mathbf{a}) \leq c_4$ to conclude the proof.

Geometrically, to ensure this angle constraint on \mathbf{q} , we need to bound the perturbation magnitude $\|\Delta \mathcal{G}_1\|$ by the Euclidean distance to the cone boundary. Specifically, under the assumption that $0 < v_2 \leq c_2 \sqrt{\|\mathbf{a}\|}$ and $|b_2^{(\tau)}| \leq \gamma \|\mathbf{a}\|$, it follows that $\|v_2 \mathbf{w}_2\| \leq c_2 \sqrt{c_2^2 + \gamma} \|\mathbf{a}\|$. This implies that the condition:

$$\|\Delta \mathcal{G}_1\| \leq c_5 \|\mathbf{a}\| \leq \left(\frac{\|\mathbf{a}\|}{2 \cos c_4} - \frac{2c_4 c_2 \sqrt{c_2^2 + \gamma}}{\pi} \|\mathbf{a}\| \right) \sin(2c_4)$$

is sufficient to guarantee $\angle(\mathbf{q}, \mathbf{a}) \leq c_4$. Numerical verification confirms that the chosen constants $c_2 = 2$, $c_4 = \frac{\pi}{20}$, $c_5 = \frac{1}{50}$, and $\gamma = \frac{1}{4}$ satisfy the required inequality. Consequently, we have $\angle(\mathbf{w}_1^{(\tau+1)}, \mathbf{a}) \leq c_4$, which completes the proof.

C.3. Proof of the Stability of Norms (Lemma 5)

We prove the lemma with the following constants: $c_2 = 2$, $c_0 \leq \frac{4}{25}$, $c_4 \leq \frac{\pi}{10}$, $c_5 \leq \frac{1}{3}$, $\gamma \leq \frac{1}{2}$.

By symmetry, it suffices to prove the bound for $v_1^{(\tau+1)}$. We first show that $v_1^{(\tau+1)} \leq 2\sqrt{\|\mathbf{a}\|}$. Applying the gradient descent update rule, the partial derivative with respect to v_1 is bounded as

follows:

$$\begin{aligned}
 \nabla_{v_1} \widehat{\mathcal{L}} &= \frac{1}{2} \left(v_1 \|\mathbf{w}_1\|^2 - \frac{(\pi - \theta) \cos \theta + \sin \theta}{\pi} v_2 \|\mathbf{w}_1\| \|\mathbf{w}_2\| - \mathbf{w}_1^T \mathbf{a} + \mathbf{w}_1^T \Delta \mathcal{G}_1 \right) \\
 &\geq \frac{1}{2} \left(v_1 \|\mathbf{w}_1\|^2 - \frac{(\pi - \theta) \cos \theta + \sin \theta}{\pi} v_2 \|\mathbf{w}_1\| \|\mathbf{w}_2\| - \|\mathbf{w}_1\| \|\mathbf{a}\| - \|\mathbf{w}_1\| \|\Delta \mathcal{G}_1\| \right) \\
 &= \frac{\|\mathbf{w}_1\|}{2} \left(v_1 \|\mathbf{w}_1\| - \frac{(\pi - \theta) \cos \theta + \sin \theta}{\pi} v_2 \|\mathbf{w}_2\| - \|\mathbf{a}\| - \|\Delta \mathcal{G}_1\| \right) \\
 &\geq \frac{\|\mathbf{w}_1\|}{2} \left(v_1 \|\mathbf{w}_1\| - \frac{1}{3} \|\mathbf{a}\| - \|\mathbf{a}\| - \frac{1}{3} \|\mathbf{a}\| \right) \\
 &= \frac{\|\mathbf{w}_1\|}{2} \left(v_1 \|\mathbf{w}_1\| - \frac{5}{3} \|\mathbf{a}\| \right).
 \end{aligned}$$

In the penultimate line, we use the assumption $\|\Delta \mathcal{G}_1\| \leq c_5 \|\mathbf{a}\| \leq \frac{1}{3} \|\mathbf{a}\|$ and the fact that

$$\frac{(\pi - \theta) \cos \theta + \sin \theta}{\pi} v_2 \|\mathbf{w}_2\| \leq \frac{(2c_4) \cos(\pi - 2c_4) + \sin(\pi - 2c_4)}{\pi} c_2 \sqrt{c_2^2 + \gamma} \|\mathbf{a}\| \leq \frac{1}{3} \|\mathbf{a}\|.$$

To establish the bound, we consider the following two cases based on the magnitude of v_1 :

- **Case 1:** $v_1 \geq \frac{5}{3} \sqrt{\|\mathbf{a}\|}$.

Since $\left| \|\mathbf{w}_1\|^2 - v_1^2 \right| = |b_1| \leq \gamma \|\mathbf{a}\| \leq \|\mathbf{a}\|$, we have $\|\mathbf{w}_1\| \geq \sqrt{\|\mathbf{a}\|}$. Thus, we have

$$\begin{aligned}
 \nabla_{v_1} \widehat{\mathcal{L}} &\geq \frac{\|\mathbf{w}_1\|}{2} \left(v_1 \|\mathbf{w}_1\| - \frac{5}{3} \|\mathbf{a}\| \right) \\
 &\geq \frac{\|\mathbf{w}_1\|}{2} \left(\frac{5}{3} \sqrt{\|\mathbf{a}\|} \cdot \sqrt{\|\mathbf{a}\|} - \frac{5}{3} \|\mathbf{a}\| \right) \\
 &\geq 0,
 \end{aligned}$$

which means $v_1^{(\tau+1)} = v_1 - \mu \nabla_{v_1} \widehat{\mathcal{L}} \leq v_1 \leq 2\sqrt{\|\mathbf{a}\|}$.

- **Case 2:** $v_1 < \frac{5}{3} \sqrt{\|\mathbf{a}\|}$.

We have

$$\begin{aligned}
 \nabla_{v_1} \widehat{\mathcal{L}} &\geq \frac{\|\mathbf{w}_1\|}{2} \left(v_1 \|\mathbf{w}_1\| - \frac{5}{3} \|\mathbf{a}\| \right) \\
 &\geq -\frac{5}{6} \|\mathbf{w}_1\| \|\mathbf{a}\|,
 \end{aligned}$$

which means

$$\begin{aligned}
 v_1^{(\tau+1)} &= v_1 - \mu \nabla_{v_1} \widehat{\mathcal{L}} \\
 &\leq \frac{5}{3} \sqrt{\|\mathbf{a}\|} + \frac{5}{6} \mu \|\mathbf{w}_1\| \|\mathbf{a}\| \\
 &\leq 2\sqrt{\|\mathbf{a}\|}.
 \end{aligned}$$

Here we use the fact that $\mu \|\mathbf{w}_1\| \sqrt{\|\mathbf{a}\|} \leq c_0 \sqrt{c_2^2 + \gamma} \leq \frac{4}{25} \cdot \sqrt{5} \leq \frac{2}{5}$.

Next, we establish the lower bound $v_1^{(\tau+1)} > \beta\sqrt{\|\mathbf{a}\|}$. By the update rule of gradient descent, we have

$$\begin{aligned}\nabla_{v_1} \widehat{\mathcal{L}} &= \frac{1}{2} \left(v_1 \|\mathbf{w}_1\|^2 - \frac{(\pi - \theta) \cos \theta + \sin \theta}{\pi} v_2 \|\mathbf{w}_1\| \|\mathbf{w}_2\| - \mathbf{w}_1^T \mathbf{a} + \mathbf{w}_1^T \Delta \mathcal{G}_1 \right) \\ &\leq \frac{1}{2} \left(v_1 \|\mathbf{w}_1\|^2 - \mathbf{w}_1^T \mathbf{a} + \|\mathbf{w}_1\| \|\Delta \mathcal{G}_1\| \right) \\ &= \frac{\|\mathbf{w}_1\|}{2} (v_1 \|\mathbf{w}_1\| - \cos \theta_1 \|\mathbf{a}\| + \|\Delta \mathcal{G}_1\|) \\ &\leq \frac{\|\mathbf{w}_1\|}{2} \left(v_1 \|\mathbf{w}_1\| - \frac{5}{6} \|\mathbf{a}\| + \frac{1}{3} \|\mathbf{a}\| \right) \\ &= \frac{\|\mathbf{w}_1\|}{2} \left(v_1 \|\mathbf{w}_1\| - \frac{1}{2} \|\mathbf{a}\| \right).\end{aligned}$$

In the penultimate line, we use the assumption $\|\Delta \mathcal{G}_1\| \leq c_5 \|\mathbf{a}\| \leq \frac{1}{3} \|\mathbf{a}\|$ and $\theta_1 \leq c_4 \leq \frac{\pi}{10} \leq \arccos(\frac{5}{6})$.

To establish the bound, we consider the following two cases based on the magnitude of $v_1 \|\mathbf{w}_1\|$:

- **Case 1:** $v_1 \|\mathbf{w}_1\| \leq \frac{1}{2} \|\mathbf{a}\|$.
we have

$$\begin{aligned}\nabla_{v_1} \widehat{\mathcal{L}} &\leq \frac{\|\mathbf{w}_1\|}{2} \left(v_1 \|\mathbf{w}_1\| - \frac{1}{2} \|\mathbf{a}\| \right) \\ &\leq 0,\end{aligned}$$

which means $v_1^{(\tau+1)} = v_1 - \mu \nabla_{v_1} \widehat{\mathcal{L}} \geq \beta\sqrt{\|\mathbf{a}\|}$.

- **Case 2:** $v_1 \|\mathbf{w}_1\| > \frac{1}{2} \|\mathbf{a}\|$.

Since $\left| \|\mathbf{w}_1\|^2 - v_1^2 \right| = |b_1| \leq \gamma \|\mathbf{a}\| \leq \frac{1}{2} \|\mathbf{a}\|$, we have $v_1 > \frac{1}{2} \sqrt{\|\mathbf{a}\|}$ and $\|\mathbf{w}_1\| < \sqrt{\|\mathbf{a}\|}$. Thus, we have

$$\begin{aligned}\nabla_{v_1} \widehat{\mathcal{L}} &\leq \frac{\|\mathbf{w}_1\|}{2} \left(v_1 \|\mathbf{w}_1\| - \frac{1}{2} \|\mathbf{a}\| \right) \\ &\leq \frac{v_1 \|\mathbf{w}_1\|^2}{2} \\ &< \frac{v_1 \|\mathbf{a}\|}{2}.\end{aligned}$$

It follows that $v_1^{(\tau+1)} = v_1 - \mu \nabla_{v_1} \widehat{\mathcal{L}} \geq v_1 - \frac{\mu v_1 \|\mathbf{a}\|}{2} \geq \frac{1}{2} v_1 > \frac{1}{4} \sqrt{\|\mathbf{a}\|} \geq \beta\sqrt{\|\mathbf{a}\|}$. Here we use the assumption $\mu \leq \frac{c_0}{\|\mathbf{a}\|} \leq \frac{1}{\|\mathbf{a}\|}$.

This concludes the proof of Lemma 5.

C.4. Proof of Phase 1 (Lemma 6)

We prove the lemma with the following constants: $c_0 \leq 1, c_4 \leq \frac{\pi}{4}, c_9 = \frac{64}{\tan c_4}, c_2 = 2, c_3 = \frac{1}{4}, c_{10} = \frac{1}{4e^{2\alpha}}, \sigma_0 = \frac{1}{8e^{2\alpha}}$, where $\alpha = \frac{65}{\tan c_4}$.

For notational simplicity, we assume without loss of generality that $\|\mathbf{a}\| = 1$. Recall that our initialization scheme is given by

$$\mathbf{w}_1^{(0)}, \mathbf{w}_2^{(0)} \sim \mathcal{N}\left(0, \frac{\sigma^2}{d} \mathbf{I}_d\right), \quad v_1^{(0)}, v_2^{(0)} \sim \frac{\sigma}{\sqrt{d}} \xi, \quad \xi^2 \sim \chi_d^2,$$

where $\sigma \leq \sigma_0 \sqrt{\|\mathbf{a}\|}$ and χ_d^2 denotes the chi-squared distribution with d degrees of freedom.

By definition, both the squared scalar $\left(v_i^{(0)}\right)^2$ and the squared vector norm $\left\|\mathbf{w}_i^{(0)}\right\|^2$ follow the same Chi-squared distribution. By standard concentration inequalities, both variables concentrate sharply around the value σ^2 with probability at least $1 - O(e^{-cd})$. Furthermore, the projection $\mathbf{a}^T \mathbf{w}_i^{(0)}$ follows a Gaussian distribution $\mathcal{N}\left(0, \frac{\sigma^2}{d} \|\mathbf{a}\|^2\right)$, which concentrates around 0 with magnitude $O(1/\sqrt{d})$.

Specifically, applying the Laurent-Massart concentration bounds for the Chi-squared distribution and standard Gaussian tail bounds, we have that with probability at least $1 - 8 \exp\left(-\min\left(\frac{1}{16}, \frac{\hat{c}_4^2}{4}\right) d\right)$, the following inequalities hold simultaneously:

$$\begin{aligned} v_1^{(0)} &\geq \frac{1}{2} \left\|\mathbf{w}_1^{(0)}\right\|, \quad v_2^{(0)} \geq \frac{1}{2} \left\|\mathbf{w}_2^{(0)}\right\| \\ \frac{1}{2} \sigma &\leq v_1^{(0)}, v_2^{(0)} \leq 2\sigma \\ \mathbf{a}^T \mathbf{w}_1^{(0)} &\geq -\hat{c}_4 v_1^{(0)}, \quad \mathbf{a}^T \mathbf{w}_2^{(0)} \geq -\hat{c}_4 v_2^{(0)}. \end{aligned}$$

where \hat{c}_4 is a fixed positive constant.

Assume that $c_0 \leq 1$, $c_4 \leq \frac{\pi}{4}$ and $T_1 = \left\lceil \frac{64}{\mu \tan c_4} \right\rceil$. We aim to establish the following properties for all iterations $\tau \leq T$ via induction:

$$v_1^{(\tau)} \geq \frac{1}{2} v_1^{(0)}, \quad v_2^{(\tau)} \geq \frac{1}{2} v_2^{(0)} \tag{21}$$

$$\mathbf{a}^T \mathbf{w}_1^{(\tau)} \geq -\hat{c}_4 v_1^{(0)}, \quad \mathbf{a}^T \mathbf{w}_2^{(\tau)} \geq -\hat{c}_4 v_2^{(0)} \tag{22}$$

$$v_1^{(\tau)}, \left\|\mathbf{w}_1^{(\tau)}\right\| \leq 2(1 + \mu)^\tau v_1^{(0)} \leq \hat{c}_3 v_1^{(0)} \leq \frac{1}{\hat{c}_5} \tag{23}$$

$$v_2^{(\tau)}, \left\|\mathbf{w}_2^{(\tau)}\right\| \leq 2(1 + \mu)^\tau v_2^{(0)} \leq \hat{c}_3 v_2^{(0)} \leq \frac{1}{\hat{c}_5}$$

with constant $\hat{c}_4 = \frac{1}{2\alpha}$, $\hat{c}_3 = \hat{c}_5 = 2e^\alpha$, $\sigma_0 = \frac{1}{8e^{2\alpha}}$ where $\alpha = \frac{65}{\tan c_4}$.

We prove (21), (22), (23) by induction. At initialization it is true with probability at least $1 - \tilde{C}e^{-\tilde{c}d}$ as explained above. Assuming these hypotheses hold for some $\tau < T$, we proceed to show they remain valid for iteration $\tau + 1$. By symmetry, we focus on v_1 and \mathbf{w}_1 . For the sake of notation simplicity, we suppress the superscript (τ) where the context is clear.

We start with bound for $\mathbf{a}^T \mathbf{w}_1$. Specifically, the update rule for the alignment term yields

$$\begin{aligned} &\mathbf{a}^T \mathbf{w}_1^{(\tau+1)} \\ &= \mathbf{a}^T \mathbf{w}_1 - \mu \mathbf{a}^T \nabla_{\mathbf{w}_1} \hat{\mathcal{L}} \\ &= \mathbf{a}^T \mathbf{w}_1 + \frac{\mu}{2} v_1 - \frac{\mu v_1}{2} \left(v_1 \mathbf{a}^T \mathbf{w}_1 - \left(1 - \frac{\theta}{\pi}\right) v_2 \mathbf{a}^T \mathbf{w}_2 - \frac{\sin \theta}{\pi} v_2 \|\mathbf{w}_2\| \mathbf{a}^T \bar{\mathbf{w}}_1 + \mathbf{a}^T \Delta \mathcal{G}_1 \right). \end{aligned}$$

By Lemma 11, with probability at least $1 - 3e^{-cd}$, we have $\|\mathcal{G}_1\| \leq \frac{1}{\hat{c}_5^2} \|\mathbf{a}\|$. It follows that

$$\begin{aligned} & \left| v_1 \mathbf{a}^T \mathbf{w}_1 - \left(1 - \frac{\theta}{\pi}\right) v_2 \mathbf{a}^T \mathbf{w}_2 - \frac{\sin \theta}{\pi} v_2 \|\mathbf{w}_2\| \mathbf{a}^T \bar{\mathbf{w}}_1 + \mathbf{a}^T \Delta \mathcal{G}_1 \right| \\ & \leq v_1 \|\mathbf{w}_1\| + v_2 \|\mathbf{w}_2\| + v_2 \|\mathbf{w}_2\| + \|\mathbf{a}^T \Delta \mathcal{G}_1\| \\ & \leq \frac{4}{\hat{c}_5^2} \leq \frac{1}{2}. \end{aligned}$$

Combining the above bounds, the term inside the parenthesis is bounded by $\frac{1}{2}$ in absolute value. Consequently, the update satisfies

$$\mathbf{a}^T \mathbf{w}_1 + \frac{\mu}{4} v_1 \leq \mathbf{a}^T \mathbf{w}_1^{(\tau+1)} \leq \mathbf{a}^T \mathbf{w}_1 + \frac{3\mu}{4} v_1. \quad (24)$$

It follows that $\mathbf{a}^T \mathbf{w}_1^{(\tau+1)} \geq \mathbf{a}^T \mathbf{w}_1$. By the inductive hypothesis $\mathbf{a}^T \mathbf{w}_1 \geq -\hat{c}_4 v_1^{(0)}$, we conclude that (22) holds for iteration $\tau + 1$.

Next, we bound the orthogonal component $\left\| (\mathbf{I} - \mathbf{a}\mathbf{a}^T) \mathbf{w}_1^{(\tau+1)} \right\|$. Observe that:

$$\begin{aligned} & (\mathbf{I} - \mathbf{a}\mathbf{a}^T) \mathbf{w}_1^{(\tau+1)} \\ & = (\mathbf{I} - \mathbf{a}\mathbf{a}^T) \mathbf{w}_1 - \mu (\mathbf{I} - \mathbf{a}\mathbf{a}^T) \nabla_{\mathbf{w}_1} \text{Loss} \\ & = (\mathbf{I} - \mathbf{a}\mathbf{a}^T) \mathbf{w}_1 - \mu \frac{v_1}{2} (\mathbf{I} - \mathbf{a}\mathbf{a}^T) \left(v_1 \mathbf{w}_1 - \left(1 - \frac{\theta}{\pi}\right) v_2 \mathbf{w}_2 - \frac{\sin \theta}{\pi} v_2 \|\mathbf{w}_2\| \bar{\mathbf{w}}_1 - \mathbf{a} + \Delta \mathcal{G}_1 \right). \end{aligned}$$

Since $(\mathbf{I} - \mathbf{a}\mathbf{a}^T) \mathbf{a} = 0$, the update simplifies to:

$$\begin{aligned} & (\mathbf{I} - \mathbf{a}\mathbf{a}^T) \mathbf{w}_1^{(\tau+1)} \\ & = (\mathbf{I} - \mathbf{a}\mathbf{a}^T) \mathbf{w}_1 - \mu \frac{v_1}{2} (\mathbf{I} - \mathbf{a}\mathbf{a}^T) \left(v_1 \mathbf{w}_1 - \left(1 - \frac{\theta}{\pi}\right) v_2 \mathbf{w}_2 - \frac{\sin \theta}{\pi} v_2 \|\mathbf{w}_2\| \bar{\mathbf{w}}_1 + \Delta \mathcal{G}_1 \right). \end{aligned}$$

Applying the triangle inequality and substituting the bounds for $\|\mathbf{w}_1\|$, $\|\mathbf{w}_2\|$, and $\|\Delta \mathcal{G}_1\|$, we have

$$\begin{aligned} & \left\| (\mathbf{I} - \mathbf{a}\mathbf{a}^T) \mathbf{w}_1^{(\tau+1)} \right\| \\ & \leq \left\| (\mathbf{I} - \mathbf{a}\mathbf{a}^T) \mathbf{w}_1 \right\| + \mu \left\| \frac{v_1}{2} (\mathbf{I} - \mathbf{a}\mathbf{a}^T) \left(v_1 \mathbf{w}_1 - \left(1 - \frac{\theta}{\pi}\right) v_2 \mathbf{w}_2 - \frac{\sin \theta}{\pi} v_2 \|\mathbf{w}_2\| \bar{\mathbf{w}}_1 + \Delta \mathcal{G}_1 \right) \right\| \\ & \leq \left\| (\mathbf{I} - \mathbf{a}\mathbf{a}^T) \mathbf{w}_1 \right\| + \frac{\mu v_1}{2} (v_1 \|\mathbf{w}_1\| + 2v_2 \|\mathbf{w}_2\| + \|\Delta \mathcal{G}_1\|) \\ & \leq \left\| (\mathbf{I} - \mathbf{a}\mathbf{a}^T) \mathbf{w}_1 \right\| + \frac{2\mu v_1}{\hat{c}_5^2}. \end{aligned} \quad (25)$$

Combining the upper bound for $\mathbf{a}^T \mathbf{w}_1^{(\tau+1)}$ (from inequality (24)) and the bound in (25), we derive the upper bound for $\|\mathbf{w}_1^{(\tau+1)}\|$ as follows:

$$\begin{aligned}
 \|\mathbf{w}_1^{(\tau+1)}\|^2 &= \left(\mathbf{a}^T \mathbf{w}_1^{(\tau+1)}\right)^2 + \left\|(\mathbf{I} - \mathbf{a}\mathbf{a}^T)\mathbf{w}_1^{(\tau+1)}\right\|^2 \\
 &\leq \left(\mathbf{a}^T \mathbf{w}_1 + \frac{3\mu}{4}v_1\right)^2 + \left(\|(\mathbf{I} - \mathbf{a}\mathbf{a}^T)\mathbf{w}_1\| + \frac{2}{c_5^2}\mu v_1^{(0)}\right)^2 \\
 &\leq \|\mathbf{w}_1\|^2 + 2\left(\frac{3\mu}{4}v_1 + \frac{2}{c_5^2}\mu v_1^{(0)}\right)\|\mathbf{w}_1\| + \left(\frac{3\mu}{4}v_1\right)^2 + \left(\frac{2}{c_5^2}\mu v_1\right)^2 \\
 &\leq \left(\|\mathbf{w}_1\| + \frac{3\mu}{4}v_1 + \frac{2}{c_5^2}\mu v_1^{(0)}\right)^2,
 \end{aligned}$$

Taking the square root and using the inductive hypothesis $v_1^{(\tau)}$, $\|\mathbf{w}_1^{(\tau)}\| \leq 2(1+\mu)^\tau v_1^{(0)}$, along with the fact that $\frac{2}{c_5^2} \leq \frac{1}{4}$, we conclude $\|\mathbf{w}_1^{(\tau+1)}\| \leq 2(1+\mu)^{\tau+1} v_1^{(0)}$, which establishes the upper bound for $\|\mathbf{w}_1\|$ in (23) for iteration $\tau + 1$.

We now turn to the evolution of v_1 . the update for v_1 is given by:

$$\begin{aligned}
 &v_1^{(\tau+1)} \\
 &= v_1 - \mu \nabla_{v_1} Loss \\
 &= v_1 - \mu \frac{\mathbf{w}_1^T}{v_1} \nabla_{\mathbf{w}_1} Loss \\
 &= v_1 - \frac{\mu}{2} \mathbf{w}_1^T \left(v_1 \mathbf{w}_1 - \left(1 - \frac{\theta}{\pi}\right) v_2 \mathbf{w}_2 - \frac{\sin \theta}{\pi} v_2 \|\mathbf{w}_2\| \bar{\mathbf{w}}_1 - \mathbf{a} + \Delta \mathcal{G}_1 \right) \\
 &= v_1 + \frac{\mu}{2} \mathbf{a}^T \mathbf{w}_1 - \frac{\mu}{2} \mathbf{w}_1^T \mathbf{w}_1 v_1 - \frac{\mu}{2} \mathbf{w}_1^T \Delta \mathcal{G}_1 + \frac{\mu}{2} v_2 \left(\left(1 - \frac{\theta}{\pi}\right) \mathbf{w}_1^T \mathbf{w}_2 + \frac{\sin \theta}{\pi} \|\mathbf{w}_1\| \|\mathbf{w}_2\| \right).
 \end{aligned}$$

By triangle inequalities, we have

$$\begin{aligned}
 &\left| -\frac{\mu}{2} \mathbf{w}_1^T \mathbf{w}_1 v_1 - \frac{\mu}{2} \mathbf{w}_1^T \Delta \mathcal{G}_1 + \frac{\mu}{2} v_2 \left(\left(1 - \frac{\theta}{\pi}\right) \mathbf{w}_1^T \mathbf{w}_2 + \frac{\sin \theta}{\pi} \|\mathbf{w}_1\| \|\mathbf{w}_2\| \right) \right| \\
 &\leq \frac{\mu}{2} \|\mathbf{w}_1\|^2 v_1 + \frac{\mu}{2} \|\mathbf{w}_1\| \|\mathcal{G}_1\| + \frac{\mu}{2} v_2 \|\mathbf{w}_1\| \|\mathbf{w}_2\| + \frac{\mu}{2} v_2 \|\mathbf{w}_1\| \|\mathbf{w}_2\| \\
 &\leq \frac{\mu}{2} \cdot 4 \frac{\hat{c}_3}{\hat{c}_5^2} v_1^{(0)}.
 \end{aligned}$$

For the upper bound, we have

$$\begin{aligned}
 v_1^{(\tau+1)} &\leq v_1 + \frac{\mu}{2} \mathbf{a}^T \mathbf{w}_1 + \frac{\mu}{2} \cdot 4 \frac{\hat{c}_3}{\hat{c}_5^2} v_1^{(0)} \\
 &\leq v_1 + \frac{\mu}{2} \mathbf{a}^T \mathbf{w}_1 + \frac{\mu}{2} v_1^{(0)} \\
 &\leq \left(1 + \frac{\mu}{2} + \frac{\mu}{2}\right) \cdot 2(1+\mu)^\tau v_1^{(0)} \\
 &= 2(1+\mu)^{\tau+1} v_1^{(0)},
 \end{aligned}$$

which shows the bound of v_1 in (23) for iteration $\tau + 1$.

For the lower bound, we have

$$\begin{aligned} v_1^{(\tau+1)} &\geq v_1 + \frac{\mu}{2} \left(-\hat{c}_4 v_1^{(0)} \right) - \frac{\mu}{2} \cdot 4 \frac{\hat{c}_3}{\hat{c}_5^2} v_1^{(0)} \\ &\geq v_1 - \frac{1}{2T_1} v_1^{(0)}. \end{aligned}$$

This implies that v_1 decrease at most one half in the first T_1 iterations, which shows (21) for iteration $\tau + 1$. Here we use the fact that $\hat{c}_4 + 4 \frac{\hat{c}_3}{\hat{c}_5^2} \leq \frac{1}{\mu T_1} \leq 1$.

Resuming the proof of Lemma 6. Inequalities (21) and (23) directly imply that at iteration T_1 :

$$v_1^{(T_1)} \geq \frac{1}{2} v_1^{(0)} \geq \frac{1}{4} \sigma,$$

and

$$v_1^{(T_1)} \leq \frac{1}{\hat{c}_5} \leq 2\sqrt{\|\mathbf{a}\|}.$$

Additionally, (23) yields an upper bound on the norm of the imbalance term $|b_1^{(T_1)}|$:

$$|b_1^{(T_1)}| = \left| \left(v_1^{(T_1)} \right)^2 - \|\mathbf{w}_1^{(T_1)}\|^2 \right| \leq \frac{1}{\hat{c}_5^2}$$

Next, we estimate the alignment angle $\theta_1^{(T_1)}$. Summing the updates in (24) and (25) over T_1 iterations, we have

$$\begin{aligned} \mathbf{a}^T \mathbf{w}_1^{(T_1)} &\geq \mathbf{a}^T \mathbf{w}_1^{(0)} + \frac{\mu}{8} v_1^{(0)} \cdot T_1 \\ &\geq -\hat{c}_4 v_1^{(0)} + \frac{\mu}{8} T_1 v_1^{(0)} \end{aligned}$$

and

$$\begin{aligned} \left\| (\mathbf{I} - \mathbf{a}\mathbf{a}^T) \mathbf{w}_1^{(T_1)} \right\| &\leq \left\| (\mathbf{I} - \mathbf{a}\mathbf{a}^T) \mathbf{w}_1^{(0)} \right\| + \frac{2c_3}{\hat{c}_5^2} \mu v_1^{(0)} \cdot T_1 \\ &\leq 2v_1^{(0)} + \frac{2\hat{c}_3}{\hat{c}_5^2} \mu T_1 v_1^{(0)}. \end{aligned}$$

It follows that $\mathbf{a}^T \mathbf{w}_1^{(T_1)} \tan c_4 \geq \left\| (\mathbf{I} - \mathbf{a}\mathbf{a}^T) \mathbf{w}_1^{(T_1)} \right\|$, which confirms

$$\theta_1^{(T_1)} = \angle(\mathbf{w}_1^{(T_1)}, \mathbf{a}) \leq c_4.$$

Here we use the fact that $\mu T \left(\frac{\tan c_4}{8} - \frac{2\hat{c}_3}{\hat{c}_5^2} \right) \geq \hat{c}_4 \tan c_4 + 2$. By symmetry, identical bounds hold for $v_2^{(T_1)}, b_2^{(T_1)}, \theta_2^{(T_1)}$. This completes the proof of Lemma 6 with constants $c_0 \leq 1, c_4 \leq \frac{\pi}{4}, c_9 = \frac{64}{\tan c_4}, c_2 = 2, c_3 = \frac{1}{4}, c_{10} = \frac{1}{4e^{2\alpha}}, \sigma_0 = \frac{1}{8e^{2\alpha}}$, where $\alpha = \frac{65}{\tan c_4}$.

C.5. Proof of Phase 2 (Lemma 7)

We prove the lemma with the following constants: $c_0 \leq \frac{1}{10^8}$, $c_7 = \frac{1}{4}$, $c_8 = 32$, $c_{11} = \frac{1}{50}$.

We first show that for any iteration $\tau \in [T_1, T_1 + T_2]$, the following bounds hold:

$$0 < v_1^{(\tau)}, v_2^{(\tau)} \leq 2\sqrt{\|\mathbf{a}\|} \quad (26)$$

$$\theta_1^{(\tau)}, \theta_2^{(\tau)} \leq c_4 \leq \frac{\pi}{10} \quad (27)$$

$$\left| b_1^{(\tau)} \right|, \left| b_2^{(\tau)} \right| \leq c_{10} \|\mathbf{a}\| + (\tau - T_1)\mu \|\mathbf{a}\|^2 \leq c_{11} \|\mathbf{a}\| \quad (28)$$

We proceed by induction. According to Lemma 6, equations (26), (27) and (28) hold for $\tau = T_1$ with probability at least $1 - Ce^{-cd}$. Now assume that we have (26), (27) and (28) for $\tau = T_1 + t$ with $t \in [0, T_2 - 1]$. For $\tau = T_1 + t + 1$, observe that

$$\left| b_1^{(\tau)} \right|, \left| b_2^{(\tau)} \right| \leq c_{10} \|\mathbf{a}\| + t\mu \|\mathbf{a}\|^2 \leq \gamma \|\mathbf{a}\|.$$

By invoking Lemma 4, Lemma 5 with $\beta = 0$, we establish that (26) and (27) hold for $\tau = T_1 + t + 1$. To prove (28), we apply the Lemma 3 with the constant $c_6 = 6$, which yields:

$$\begin{aligned} \left| b_1^{(\tau+1)} - b_1^{(\tau)} \right| &\leq c_6 \mu^2 \left(\left(v_1^{(\tau)} \right)^2 + \left\| \mathbf{w}_1^{(\tau)} \right\|^2 \right) \left(\left\| v_1^{(\tau)} \mathbf{w}_1^{(\tau)} - \mathbf{a} \right\|^2 + \left\| v_2^{(\tau)} \mathbf{w}_2^{(\tau)} + \mathbf{a} \right\|^2 \right) \\ &\leq c_6 \mu^2 (4 \|\mathbf{a}\| + 5 \|\mathbf{a}\|) \left((2\sqrt{5} + 1)^2 \|\mathbf{a}\|^2 + (2\sqrt{5} + 1)^2 \|\mathbf{a}\|^2 \right) \\ &\leq 3300 \mu^2 \|\mathbf{a}\|^3. \end{aligned}$$

Here, we use the assumptions $v_1^{(\tau)}, v_2^{(\tau)} \leq 2\sqrt{\|\mathbf{a}\|}$ and the fact that $\left\| \mathbf{w}_1^{(\tau)} \right\|^2 \leq \left(v_1^{(\tau)} \right)^2 + b_1^{(\tau)} \leq (4 + c_{11}) \|\mathbf{a}\| \leq 5 \|\mathbf{a}\|$, $\left\| \mathbf{w}_2^{(\tau)} \right\|^2 \leq \left(v_2^{(\tau)} \right)^2 + b_2^{(\tau)} \leq (4 + c_{11}) \|\mathbf{a}\| \leq 5 \|\mathbf{a}\|$. Substituting these into the recursive relation for b_1 and using the Inequality (28) for $\tau = T_1 + t$ we have:

$$\left| b_1^{(\tau+1)} \right| \leq \left| b_1^{(\tau)} \right| + \left| b_1^{(\tau+1)} - b_1^{(\tau)} \right| \leq c_{10} \|\mathbf{a}\| + (\tau + 1 - T_1) \cdot 3300 \mu^2 \|\mathbf{a}\|^3 \leq c_{11} \|\mathbf{a}\|.$$

Here we use the assumption that $\tau + 1 - T_1 \leq T_2 = \lceil \frac{c_8}{\mu \|\mathbf{a}\|} \ln \left(\frac{\sqrt{\|\mathbf{a}\|}}{\sigma} \right) \rceil$ and $\mu \leq \frac{c_0}{\|\mathbf{a}\| \ln \left(\frac{\sqrt{\|\mathbf{a}\|}}{\sigma} \right)}$.

By substituting the constants $c_8 = 32$, $c_{10} \leq \frac{1}{100}$, $c_{11} = \frac{1}{50}$, and $c_0 \leq \frac{1}{10^8}$, one can verify that the requirement for c_{11} is satisfied. By symmetry, we also have $\left| b_2^{(\tau+1)} \right| \leq c_{11} \|\mathbf{a}\|$. Thus, we have (28) for $\tau = T_1 + t + 1$.

By symmetry, we only need to focus on v_1, \mathbf{w}_1 . We first establish a lower bound on $\left\| \mathbf{w}_1^{(\tau+1)} \right\|$. Since $\mathbf{w}_1^T \nabla_{\mathbf{w}_1} \hat{\mathcal{L}} = v_1 \nabla_{v_1} \hat{\mathcal{L}}$ (due to Eq. 13), we have

$$\begin{aligned} \left\| \mathbf{w}_1^{(\tau+1)} \right\|^2 &= \left\| \mathbf{w}_1 - \mu \nabla_{\mathbf{w}_1} \hat{\mathcal{L}} \right\|^2 \\ &= \left\| \mathbf{w}_1 \right\|^2 - 2\mu \mathbf{w}_1^T \nabla_{\mathbf{w}_1} \hat{\mathcal{L}} + \mu^2 \left\| \nabla_{\mathbf{w}_1} \hat{\mathcal{L}} \right\|^2 \\ &= \left\| \mathbf{w}_1 \right\|^2 - 2\mu v_1 \nabla_{v_1} \hat{\mathcal{L}} + \mu^2 \left\| \nabla_{\mathbf{w}_1} \hat{\mathcal{L}} \right\|^2, \end{aligned}$$

and

$$\begin{aligned}\left\|\nabla_{\mathbf{w}_1}\widehat{\mathcal{L}}\right\| &\geq \frac{|\mathbf{w}_1^T\nabla_{\mathbf{w}_1}\widehat{\mathcal{L}}|}{\|\mathbf{w}_1\|} \\ &= \frac{|v_1\nabla_{v_1}\widehat{\mathcal{L}}|}{\|\mathbf{w}_1\|}.\end{aligned}$$

Thus, we have

$$\begin{aligned}\left\|\mathbf{w}_1^{(\tau+1)}\right\|^2 &= \|\mathbf{w}_1\|^2 - 2\mu v_1\nabla_{v_1}\widehat{\mathcal{L}} + \mu^2\left\|\nabla_{\mathbf{w}_1}\widehat{\mathcal{L}}\right\|^2 \\ &\geq \|\mathbf{w}_1\|^2 - 2\mu v_1\nabla_{v_1}\widehat{\mathcal{L}} + \mu^2\left(\frac{|v_1\nabla_{v_1}\widehat{\mathcal{L}}|}{\|\mathbf{w}_1\|}\right)^2 \\ &= \left(\|\mathbf{w}_1\| - \mu\frac{v_1\nabla_{v_1}\widehat{\mathcal{L}}}{\|\mathbf{w}_1\|}\right)^2.\end{aligned}$$

It follows that

$$\left\|\mathbf{w}_1^{(\tau+1)}\right\| \geq \|\mathbf{w}_1\| - \mu\frac{v_1\nabla_{v_1}\widehat{\mathcal{L}}}{\|\mathbf{w}_1\|}. \quad (29)$$

We continue by estimating $\nabla_{v_1}\widehat{\mathcal{L}}$. Note that

$$\begin{aligned}\nabla_{v_1}\widehat{\mathcal{L}} &= \frac{1}{2}\left(v_1\|\mathbf{w}_1\|^2 - v_2\frac{(\pi-\theta)\cos\theta + \sin\theta}{\pi}\|\mathbf{w}_1\|\|\mathbf{w}_2\| - \mathbf{w}_1^T\mathbf{a} + \mathbf{w}_1^T\Delta\mathcal{G}_1\right) \\ &\leq \frac{1}{2}\left(v_1\|\mathbf{w}_1\|^2 - \mathbf{w}_1^T\mathbf{a} + \mathbf{w}_1^T\Delta\mathcal{G}_1\right) \\ &\leq \frac{1}{2}\left(v_1\|\mathbf{w}_1\|^2 - \cos c_4\|\mathbf{w}_1\|\|\mathbf{a}\| + \|\mathbf{w}_1\|\|\Delta\mathcal{G}_1\|\right) \\ &\leq \frac{1}{2}\left(v_1\|\mathbf{w}_1\|^2 - \frac{5}{6}\|\mathbf{w}_1\|\|\mathbf{a}\| + \frac{1}{3}\|\mathbf{w}_1\|\|\mathbf{a}\|\right) \\ &\leq \frac{1}{2}\left(v_1\|\mathbf{w}_1\|^2 - \frac{1}{2}\|\mathbf{w}_1\|\|\mathbf{a}\|\right).\end{aligned}$$

In the penultimate line, we use the assumptions $c_4 = \frac{\pi}{10} \leq \arccos\left(\frac{5}{6}\right)$ and $\|\Delta\mathcal{G}_1\| \leq \frac{1}{3}\|\mathbf{a}\|$. To establish the lower bound for $v_1^{(T_1+T_2)}$, we will show that in phase 2, $v_1\|\mathbf{w}_1\|$ increases when it is small and will never decrease too much when it is large. Specifically, we consider the following three cases based on the magnitude of $v_1\|\mathbf{w}_1\|$:

- **Case 1:** $v_1\|\mathbf{w}_1\| \leq \frac{1}{4}\|\mathbf{a}\|$.

We have

$$\nabla_{v_1}\widehat{\mathcal{L}} \leq \frac{1}{2}\left(v_1\|\mathbf{w}_1\|^2 - \frac{1}{2}\|\mathbf{w}_1\|\|\mathbf{a}\|\right) \leq -\frac{1}{8}\|\mathbf{a}\|\|\mathbf{w}_1\|.$$

By Inequality (29), we have

$$\begin{aligned}
 \left\| \mathbf{w}_1^{(\tau+1)} \right\| &\geq \left\| \mathbf{w}_1 \right\| - \mu \frac{v_1 \nabla_{v_1} \widehat{\mathcal{L}}}{\left\| \mathbf{w}_1 \right\|} \\
 &\geq \left\| \mathbf{w}_1 \right\| - \mu \frac{v_1 \left(-\frac{1}{8} \left\| \mathbf{w}_1 \right\| \left\| \mathbf{a} \right\| \right)}{\left\| \mathbf{w}_1 \right\|} \\
 &= \left\| \mathbf{w}_1 \right\| + \frac{1}{8} \mu \left\| \mathbf{a} \right\| v_1.
 \end{aligned}$$

Summing the two inequalities we conclude that

$$v_1^{(\tau+1)} + \left\| \mathbf{w}_1^{(\tau+1)} \right\| \geq \left(1 + \frac{1}{8} \mu \left\| \mathbf{a} \right\| \right) \left(v_1^{(\tau)} + \left\| \mathbf{w}_1^{(\tau)} \right\| \right).$$

- **Case 2:** $\frac{1}{4} \left\| \mathbf{a} \right\| < v_1 \left\| \mathbf{w}_1 \right\| \leq \frac{1}{2} \left\| \mathbf{a} \right\|$.

We have

$$\begin{aligned}
 \nabla_{v_1} \widehat{\mathcal{L}} &\leq \frac{1}{2} \left(v_1 \left\| \mathbf{w}_1 \right\|^2 - \frac{1}{2} \left\| \mathbf{w}_1 \right\| \left\| \mathbf{a} \right\| \right) \\
 &\leq 0,
 \end{aligned}$$

which means $v_1^{(\tau+1)} = v_1 - \mu \nabla_{v_1} \widehat{\mathcal{L}} \geq v_1$. By Inequality (29), we have

$$\begin{aligned}
 \left\| \mathbf{w}_1^{(\tau+1)} \right\| &\geq \left\| \mathbf{w}_1 \right\| - \mu \frac{v_1 \nabla_{v_1} \widehat{\mathcal{L}}}{\left\| \mathbf{w}_1 \right\|} \\
 &\geq \left\| \mathbf{w}_1 \right\|.
 \end{aligned}$$

Thus, we have $v_1^{(\tau+1)} \left\| \mathbf{w}_1 \right\|^{(\tau+1)} \geq v_1 \left\| \mathbf{w}_1 \right\| \geq \frac{1}{4} \left\| \mathbf{a} \right\|$.

- **Case 3:** $v_1 \left\| \mathbf{w}_1 \right\| > \frac{1}{2} \left\| \mathbf{a} \right\|$.

We have

$$\begin{aligned}
 \nabla_{v_1} \widehat{\mathcal{L}} &\leq \frac{1}{2} \left(v_1 \left\| \mathbf{w}_1 \right\|^2 - \frac{1}{2} \left\| \mathbf{w}_1 \right\| \left\| \mathbf{a} \right\| \right) \\
 &\leq \frac{1}{2} v_1 \left\| \mathbf{w}_1 \right\|^2,
 \end{aligned}$$

which means $v_1^{(\tau+1)} = v_1 - \mu \nabla_{v_1} \widehat{\mathcal{L}} \geq v_1 - \frac{\mu v_1}{2} \left\| \mathbf{w}_1 \right\|^2 \geq v_1 - \frac{1}{4} v_1 = \frac{3}{4} v_1$. Here we use the fact that $\mu \left\| \mathbf{w}_1 \right\|^2 \leq c_0 (c_2^2 + \gamma) \leq \frac{1}{2}$. Using Inequality (29), we have

$$\begin{aligned}
 \left\| \mathbf{w}_1^{(\tau+1)} \right\| &\geq \left\| \mathbf{w}_1 \right\| - \mu \frac{v_1 \nabla_{v_1} \widehat{\mathcal{L}}}{\left\| \mathbf{w}_1 \right\|} \\
 &\geq \left\| \mathbf{w}_1 \right\| - \frac{\mu v_1^2}{2} \left\| \mathbf{w}_1 \right\| \\
 &\geq \left\| \mathbf{w}_1 \right\| - \frac{1}{4} \left\| \mathbf{w}_1 \right\| \\
 &= \frac{3}{4} \left\| \mathbf{w}_1 \right\|.
 \end{aligned}$$

Here we use the fact that $\mu v_1^2 \leq c_0 c_2^2 \leq \frac{1}{2}$. Thus, we have $v_1^{(\tau+1)} \|\mathbf{w}_1\|^{(\tau+1)} \geq \frac{3}{4} v_1 \cdot \frac{3}{4} \|\mathbf{w}_1\| \geq \frac{1}{2} v_1 \|\mathbf{w}_1\| \geq \frac{1}{4} \|\mathbf{a}\|$.

As a result, the sum $v_1 + \|\mathbf{w}_1\|$ will increase by a factor $(1 + \frac{1}{8}\mu \|\mathbf{a}\|)$ as long as $v_1 \|\mathbf{w}_1\| < \frac{1}{4} \|\mathbf{a}\|$. Once we have $v_1 \|\mathbf{w}_1\| \geq \frac{1}{4} \|\mathbf{a}\|$ at some iteration, it remains bounded below by $\frac{1}{4} \|\mathbf{a}\|$. By Lemma 6, $v_1^{(T_1)} + \|\mathbf{w}_1^{(T_1)}\| \geq c_3 \sigma \geq \frac{1}{4} \sigma$. Consequently, after $T_2 = \lceil \frac{c_8}{\mu \|\mathbf{a}\|} \ln \left(\frac{\sqrt{\|\mathbf{a}\|}}{\sigma} \right) \rceil$ iterations, where we have used the constant $c_8 = 32$, we have $v_1^{(T_1+T_2)} \|\mathbf{w}_1^{(T_1+T_2)}\| \geq \frac{1}{4} \|\mathbf{a}\|$. Combining this with the imbalance bound $\left| \left(v_1^{(T_1+T_2)} \right)^2 - \left\| \mathbf{w}_1^{(T_1+T_2)} \right\|^2 \right| \leq c_{11} \|\mathbf{a}\| \leq \frac{1}{2} \|\mathbf{a}\|$, we have $v_1^{(T_1+T_2)} \geq \frac{1}{4} \sqrt{\|\mathbf{a}\|}$. This completes the proof with the constant $c_{11} = \frac{1}{2}$.

C.6. Uniform Concentration (Lemma 10)

We prove the result in 5 steps.

Step 1: Standard net reduction for operator norm.

Lemma 14 (Operator norm on a net) *Let $A \in \mathbb{R}^{d \times d}$ be symmetric and let $U \subset \mathbb{S}^{d-1}$ be an ε -net with $\varepsilon \in (0, 1/2)$. Then*

$$\|A\| \leq \frac{1}{1 - 2\varepsilon} \max_{u \in U} |u^\top A u|.$$

Moreover, there exists a $1/4$ -net U with $|U| \leq 9^d$.

Hence it suffices to control, uniformly over \mathbf{w}, \mathbf{w}^* ,

$$\max_{u \in U} \left| u^\top \left(M(\mathbf{w}, \mathbf{w}^*) - \mathbb{E}M(\mathbf{w}, \mathbf{w}^*) \right) u \right|.$$

Step 2: Truncation decomposition. Fix $\tau \geq 1$ and define truncation and remainder for $t \geq 0$:

$$T_\tau(t) := \min\{t, \tau\}, \quad R_\tau(t) := (t - \tau)_+.$$

Fix $u \in \mathbb{S}^{d-1}$ and define $W_u(x) := \langle x, u \rangle^2$ and

$$W_u^{(\tau)}(x) := T_\tau(W_u(x)), \quad G_u^{(\tau)}(x) := R_\tau(W_u(x)).$$

For any $(\mathbf{w}, \mathbf{w}^*)$ and any x ,

$$W_u(x) \mathbb{1}_{\{\langle x, \mathbf{w} \rangle \geq 0\}} \mathbb{1}_{\{\langle x, \mathbf{w}^* \rangle \geq 0\}} \leq W_u^{(\tau)}(x) \mathbb{1}_{\{\langle x, \mathbf{w} \rangle \geq 0\}} \mathbb{1}_{\{\langle x, \mathbf{w}^* \rangle \geq 0\}} + G_u^{(\tau)}(x), \quad (30)$$

since $W_u = W_u^{(\tau)} + G_u^{(\tau)}$ and indicators are ≤ 1 .

For fixed u , define the *bounded* function class

$$\mathcal{F}_{u, \tau} := \left\{ f_{\mathbf{w}, \mathbf{w}^*}(x) := W_u^{(\tau)}(x) \mathbb{1}_{\{\langle x, \mathbf{w} \rangle \geq 0\}} \mathbb{1}_{\{\langle x, \mathbf{w}^* \rangle \geq 0\}} : \mathbf{w}, \mathbf{w}^* \in \mathbb{S}^{d-1} \right\}.$$

Then every $f \in \mathcal{F}_{u, \tau}$ satisfies $0 \leq f \leq \tau$ pointwise.

Step 3: Tail remainder bound

Lemma 15 (Chi-square tail moments) *Let $Z \sim \mathcal{N}(0, 1)$ and $Y = (Z^2 - \tau)_+$. There exist universal constants $c, C > 0$ such that for all $\tau \geq 1$,*

$$\mathbb{E}Y \leq Ce^{-c\tau}, \quad \mathbb{E}Y^2 \leq Ce^{-c\tau}.$$

Proof Standard: $\mathbb{P}(Z^2 > t) \leq 2e^{-t/2}$ and $\mathbb{E}(Z^2 - \tau)_+ = \int_{\tau}^{\infty} \mathbb{P}(Z^2 > t) dt$, $\mathbb{E}(Z^2 - \tau)_+^2 = 2 \int_{\tau}^{\infty} (t - \tau) \mathbb{P}(Z^2 > t) dt$. \blacksquare

Lemma 16 (Uniform control of the truncation tail over a 1/4-net) *Let U be a 1/4-net with $|U| \leq 9^d$. There exist universal $c, C > 0$ such that for all $\tau \geq 1$, with probability at least $1 - 2e^{-cd}$,*

$$\max_{u \in \mathcal{U}} \left| \frac{1}{n} \sum_{i=1}^n G_u^{(\tau)}(x_i) - \mathbb{E}G_u^{(\tau)}(X) \right| \leq C \sqrt{\frac{e^{-c\tau}d}{n}} + C \frac{d}{n}.$$

Moreover, $\max_{u \in \mathcal{U}} \mathbb{E}G_u^{(\tau)}(X) \leq Ce^{-c\tau}$.

Proof Fix $u \in \mathcal{U}$. Then $G_u^{(\tau)}(X) \stackrel{d}{=} (Z^2 - \tau)_+$ with $Z \sim \mathcal{N}(0, 1)$. By Lemma 15, $\mathbb{E}G_u^{(\tau)} \lesssim e^{-c\tau}$ and $\mathbb{E}(G_u^{(\tau)})^2 \lesssim e^{-c\tau}$, and $G_u^{(\tau)}$ is sub-exponential (dominated by Z^2). More specifically, it is a (σ^2, b) subexponential with $\sigma^2 \leq Ce^{-c\tau}$ and b a constant.

Definition 17 (Sub-exponential random variable with parameters (σ^2, b)) *A real-valued random variable X is said to be sub-exponential with parameters (σ^2, b) if, for all $|\lambda| \leq \frac{1}{b}$,*

$$\mathbb{E} \exp(\lambda(X - \mathbb{E}X)) \leq \exp\left(\frac{\lambda^2 \sigma^2}{2}\right). \quad (31)$$

For such subexponential random variables we have the following standard refined Bernstein-type inequality.

Theorem 18 (Bernstein's-type inequality for sub-exponential sums) *Let X_1, \dots, X_n be independent mean-zero random variables, where each X_i are sub-exponential with parameters (σ^2, b) in the sense of (31). Then, for every $s \geq 0$,*

$$\mathbb{P}\left(\frac{1}{n} \sum_{i=1}^n X_i \geq s\right) \leq 2 \exp\left(-\frac{1}{2}n \cdot \min\left\{\frac{s^2}{\sigma^2}, \frac{s}{b}\right\}\right). \quad (32)$$

Thus using using this Bernstein's inequality for sub-exponential variables above with $s = C \sqrt{\frac{e^{-c\tau}t}{n}} + C \frac{t}{n}$ yields for all $t \geq 1$,

$$\mathbb{P}\left(\frac{1}{n} \sum_{i=1}^n G_u^{(\tau)}(x_i) - \mathbb{E}G_u^{(\tau)}(x) \geq C \sqrt{\frac{e^{-c\tau}t}{n}} + C \frac{t}{n}\right) \leq e^{-t}.$$

Set $t = c_1 d$ and union bound over $|\mathcal{U}| \leq 9^d$; choosing c_1 large enough makes the union-bound failure probability $\leq e^{-cd}$. \blacksquare

Step 4: Uniform control of the truncated term. We will use the following result on covering numbers for VC-subgraph classes.

Theorem 19 (Theorem 2.6.7 in Van der Vaart and Wellner (1996)) *Let \mathcal{H} be a VC-subgraph class of real-valued functions on \mathbb{R}^d with VC-subgraph dimension at most V and envelope bound $|h| \leq B$ pointwise. Then there exist absolute constants $A, C > 0$ such that for every probability measure Q and every $0 < \eta \leq B$,*

$$\log N(\eta, \mathcal{H}, L_2(Q)) \leq CV \log\left(\frac{AB}{\eta}\right).$$

We will also use the following Dudley-type bound for Rademacher averages.

Theorem 20 (Equation (5.48) in Wainwright (2019)) *There exists an absolute constant $C > 0$ such that, for any function class \mathcal{F} , conditionally on x_1, \dots, x_n ,*

$$\mathbb{E}_\varepsilon \left[\sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n \varepsilon_i f(x_i) \right| \right] \leq \frac{C}{\sqrt{n}} \int_0^{2r} \sqrt{\log N(\eta, \mathcal{F}, L_2(\mathbb{P}_n))} d\eta,$$

where \mathbb{P}_n is the empirical measure of x_1, \dots, x_n and $r := \sup_{f \in \mathcal{F}} \|f\|_{L_2(\mathbb{P}_n)}$.

Lemma 21 (Expected supremum for the truncated class) *There exists a universal constant $C > 0$ such that for each fixed u ,*

$$\mathbb{E} \left[\sup_{f \in \mathcal{F}_{u,\tau}} |(\mathbb{P}_n - \mathbb{P})f| \right] \leq C\tau \sqrt{\frac{d}{n}}.$$

Proof By symmetrization,

$$\mathbb{E} \sup_{f \in \mathcal{F}_{u,\tau}} |(\mathbb{P}_n - \mathbb{P})f| \leq 2\mathbb{E} \sup_{f \in \mathcal{F}_{u,\tau}} \left| \frac{1}{n} \sum_{i=1}^n \varepsilon_i f(x_i) \right|.$$

Let $r := \sup_{f \in \mathcal{F}_{u,\tau}} \|f\|_{L_2(\mathbb{P}_n)}$. Since $0 \leq f \leq \tau$ we have $r \leq \tau$. Applying Dudley's entropy integral bound for Rademacher averages per Theorem 20 and the covering bound from Theorem 19 with $Q = \mathbb{P}_n$ yields

$$\mathbb{E} \left[\sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n \varepsilon_i f(X_i) \right| \right] \leq \frac{C_3}{\sqrt{n}} \int_0^{2r} \sqrt{\log N(\eta, \mathcal{F}, L_2(\mathbb{P}_n))} d\eta \leq \frac{C_3}{\sqrt{n}} \int_0^{2r} \sqrt{V \log\left(\frac{A\tau}{\eta}\right)} d\eta.$$

Using the change of variables $\eta = re^{-s}$ and the elementary bound

$$\int_0^{2r} \sqrt{\log\left(\frac{A\tau}{\eta}\right)} d\eta \leq C_4 r \sqrt{\log\left(\frac{eA\tau}{r}\right)},$$

we obtain

$$\mathbb{E} \left[\sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n \varepsilon_i f(X_i) \right| \right] \leq C_5 \frac{r}{\sqrt{n}} \sqrt{V} \sqrt{\log\left(\frac{eA\tau}{r}\right)}.$$

Substituting $r \leq \tau$ and $V \leq c_2 d$, and absorbing constants (including A and C_0) into a single absolute constant C , gives

$$\mathbb{E} \left[\sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n \varepsilon_i f(X_i) \right| \right] \leq C\tau \sqrt{\frac{d}{n}}.$$

Finally, multiplying by 2 from symmetrization proves the claim. \blacksquare

Now we focus on establishing a high-probability bound. For this we will use Bousquet's concentration inequality for suprema of bounded empirical processes.

Theorem 22 (Theorem 2.3 in Bousquet (2002)) *Let \mathcal{F} be a class of measurable functions with $0 \leq f \leq b$. Let*

$$Z := \sup_{f \in \mathcal{F}} \sum_{i=1}^n (f(X_i) - \mathbb{E}f(X)), \quad \sigma^2 := \sup_{f \in \mathcal{F}} \text{Var}(f(X)).$$

Then for all $t \geq 0$, with probability at least $1 - e^{-t}$,

$$Z \leq \mathbb{E}Z + \sqrt{2t(n\sigma^2 + 2b\mathbb{E}Z)} + \frac{bt}{3}. \quad (33)$$

In particular, we will further upper bound the RHS of (33) as

$$\mathbb{E}Z + \sqrt{2t(n\sigma^2 + 2b\mathbb{E}Z)} + \frac{bt}{3} \leq 2\mathbb{E}Z + \sqrt{2tn\sigma^2} + \frac{4bt}{3}, \quad (34)$$

where we have used that $\sqrt{a+b} \leq \sqrt{a} + \sqrt{b}$.

Lemma 23 (Uniform high-probability deviation for $\mathcal{F}_{u,\tau}$) *There exist universal constants $C, c > 0$ such that for each fixed u and each $t \geq 1$, with probability at least $1 - e^{-t}$,*

$$\sup_{f \in \mathcal{F}_{u,\tau}} |(\mathbb{P}_n - \mathbb{P})f| \leq C\tau \sqrt{\frac{d}{n}} + C\sqrt{\frac{\tau t}{n}} + C\frac{\tau t}{n}.$$

Proof Apply (34) to $\mathcal{F} = \mathcal{F}_{u,\tau}$ with $b = \tau$. We have $\sigma^2 \leq \sup_f \mathbb{E}f^2 \leq \tau/2$ as above. Also, $\mathbb{E}Z = n \cdot \mathbb{E} \sup_{f \in \mathcal{F}_{u,\tau}} |(\mathbb{P}_n - \mathbb{P})f|$ which is bounded by Lemma 21. Divide by n to conclude. \blacksquare

Step 5: Complete the operator-norm bound Fix $\tau > 0$ and recall the decomposition

$$W_u(x) = W_u^{(\tau)}(x) + G_u^{(\tau)}(x), \quad W_u^{(\tau)}(x) := \min\{\langle x, u \rangle^2, \tau\}, \quad G_u^{(\tau)}(x) := (\langle x, u \rangle^2 - \tau)_+.$$

For $(\mathbf{w}, \mathbf{w}^*)$ define the ReLU sign indicator

$$\mathbb{1}_{\mathbf{w}, \mathbf{w}^*}(x) := \mathbb{1}\{\langle x, \mathbf{w} \rangle \geq 0\} \mathbb{1}\{\langle x, \mathbf{w}^* \rangle \geq 0\} \in \{0, 1\}.$$

Then for every x ,

$$W_u(x) \mathbb{1}_{\mathbf{w}, \mathbf{w}^*}(x) = W_u^{(\tau)}(x) \mathbb{1}_{\mathbf{w}, \mathbf{w}^*}(x) + G_u^{(\tau)}(x) \mathbb{1}_{\mathbf{w}, \mathbf{w}^*}(x),$$

and hence, by the triangle inequality,

$$\sup_{\mathbf{w}, \mathbf{w}^*} |(\mathbb{P}_n - \mathbb{P})(W_u \mathbb{1}_{\mathbf{w}, \mathbf{w}^*})| \leq \sup_{\mathbf{w}, \mathbf{w}^*} |(\mathbb{P}_n - \mathbb{P})(W_u^{(\tau)} \mathbb{1}_{\mathbf{w}, \mathbf{w}^*})| + \sup_{\mathbf{w}, \mathbf{w}^*} |(\mathbb{P}_n - \mathbb{P})(G_u^{(\tau)} \mathbb{1}_{\mathbf{w}, \mathbf{w}^*})|. \quad (35)$$

We control the second term in (35) uniformly over $(\mathbf{w}, \mathbf{w}^*)$ without any VC/covering argument. Using the identity

$$G_u^{(\tau)} \mathbb{1}_{\mathbf{w}, \mathbf{w}^*} = G_u^{(\tau)} - G_u^{(\tau)}(1 - \mathbb{1}_{\mathbf{w}, \mathbf{w}^*}),$$

we have for every $(\mathbf{w}, \mathbf{w}^*)$,

$$(\mathbb{P}_n - \mathbb{P})(G_u^{(\tau)} \mathbb{1}_{\mathbf{w}, \mathbf{w}^*}) = (\mathbb{P}_n - \mathbb{P})(G_u^{(\tau)}) - (\mathbb{P}_n - \mathbb{P})(G_u^{(\tau)}(1 - \mathbb{1}_{\mathbf{w}, \mathbf{w}^*})),$$

so

$$|(\mathbb{P}_n - \mathbb{P})(G_u^{(\tau)} \mathbb{1}_{\mathbf{w}, \mathbf{w}^*})| \leq |(\mathbb{P}_n - \mathbb{P})(G_u^{(\tau)})| + |(\mathbb{P}_n - \mathbb{P})(G_u^{(\tau)}(1 - \mathbb{1}_{\mathbf{w}, \mathbf{w}^*}))|. \quad (36)$$

Since $G_u^{(\tau)} \geq 0$ and $0 \leq 1 - \mathbb{1}_{\mathbf{w}, \mathbf{w}^*} \leq 1$, we have pointwise

$$0 \leq G_u^{(\tau)}(1 - \mathbb{1}_{\mathbf{w}, \mathbf{w}^*}) \leq G_u^{(\tau)},$$

and therefore

$$\mathbb{P}_n(G_u^{(\tau)}(1 - \mathbb{1}_{\mathbf{w}, \mathbf{w}^*})) \leq \mathbb{P}_n(G_u^{(\tau)}), \quad \mathbb{P}(G_u^{(\tau)}(1 - \mathbb{1}_{\mathbf{w}, \mathbf{w}^*})) \leq \mathbb{P}(G_u^{(\tau)}).$$

Consequently,

$$|(\mathbb{P}_n - \mathbb{P})(G_u^{(\tau)}(1 - \mathbb{1}_{\mathbf{w}, \mathbf{w}^*}))| \leq \mathbb{P}_n(G_u^{(\tau)}) + \mathbb{P}(G_u^{(\tau)}) = (\mathbb{P}_n - \mathbb{P})(G_u^{(\tau)}) + 2\mathbb{P}(G_u^{(\tau)}),$$

and hence

$$\sup_{\mathbf{w}, \mathbf{w}^*} |(\mathbb{P}_n - \mathbb{P})(G_u^{(\tau)} \mathbb{1}_{\mathbf{w}, \mathbf{w}^*})| \leq 2|(\mathbb{P}_n - \mathbb{P})(G_u^{(\tau)})| + 2\mathbb{P}(G_u^{(\tau)}). \quad (37)$$

Plugging (37) into (35) yields the corrected completion bound:

$$\sup_{\mathbf{w}, \mathbf{w}^*} |(\mathbb{P}_n - \mathbb{P})(W_u \mathbb{1}_{\mathbf{w}, \mathbf{w}^*})| \leq \sup_{\mathbf{w}, \mathbf{w}^*} |(\mathbb{P}_n - \mathbb{P})(W_u^{(\tau)} \mathbb{1}_{\mathbf{w}, \mathbf{w}^*})| + 2|(\mathbb{P}_n - \mathbb{P})(G_u^{(\tau)})| + 2\mathbb{P}(G_u^{(\tau)}). \quad (38)$$

Now apply Lemma 23 and Lemma 16 and union bound over $u \in \mathcal{U}$. Taking $t = c_0 d$ with c_0 large enough absorbs the 9^d factor, giving with probability $\geq 1 - 3e^{-cd}$:

$$\max_{u \in \mathcal{U}} \sup_{\mathbf{w}, \mathbf{w}^*} |u^\top (M - \mathbb{E}M)u| \leq C\tau \sqrt{\frac{d}{n}} + C\sqrt{\frac{\tau d}{n}} + C\frac{\tau d}{n} + C\sqrt{\frac{e^{-c\tau} d}{n}} + Ce^{-c\tau} + C\frac{d}{n}.$$

Multiplying by 2 from the net lemma gives the same bound for $\sup_{\mathbf{w}, \mathbf{w}^*} \|M - \mathbb{E}M\|$.

Finally choose $\tau = C \log(e/\delta)$ so that $e^{-c\tau} \ll \delta$. If $n \geq Cd \frac{\log^2(1/\delta)}{\delta^2}$, then each term on the right is $\leq \delta$ (after increasing the universal constants), which proves Theorem 10.

C.7. Concentration of Gradient Component Deviations (Lemma 11)

For the proof, we utilize the variational characterization of the Euclidean norm: $\|z\| = \sup_{\|\mathbf{u}\|=1} \langle \mathbf{u}, z \rangle$.

Bound for w_1 : Let $\mathbf{u} \in \mathbb{S}^{d-1}$ be an arbitrary unit vector. Using the Fundamental Theorem of Calculus, the definition of the ReLU gradient, and the residual $r(\mathbf{x}_i) = v_1\phi(\mathbf{w}_1^\top \mathbf{x}_i) - v_2\phi(\mathbf{w}_2^\top \mathbf{x}_i) - (\phi(\mathbf{a}^\top \mathbf{x}_i) - \phi(-\mathbf{a}^\top \mathbf{x}_i))$, we observe:

$$\begin{aligned} \left\langle \mathbf{u}, \nabla_{w_1} \widehat{\mathcal{L}} - \nabla_{w_1} \mathcal{L} \right\rangle &= v_1 \int_0^1 \mathbf{u}^\top (M(t(v_1 \mathbf{w}_1) + (1-t)\mathbf{a}, \mathbf{w}_1) - \mathbb{E}M(\dots)) \mathbf{h}_1 dt \\ &\quad + v_1 \int_0^1 \mathbf{u}^\top (M(-t\mathbf{a} + (1-t)(v_2 \mathbf{w}_2), \mathbf{w}_1) - \mathbb{E}M(\dots)) (-\mathbf{h}_2) dt, \end{aligned}$$

where $M(\tilde{\mathbf{u}}, \tilde{\mathbf{v}})$ is defined as in Eq. 4. The term v_1 appears due to the chain rule derivative with respect to w_1 . By Lemma 10, given the sample complexity $n \geq Cd \frac{\log(1/\delta)^2}{\delta^2}$, the following spectral deviation bound holds with probability at least $1 - 3e^{-cd}$:

$$\sup_{\tilde{\mathbf{u}}, \tilde{\mathbf{v}} \in \mathbb{S}^{d-1}} \|M(\tilde{\mathbf{u}}, \tilde{\mathbf{v}}) - \mathbb{E}[M(\tilde{\mathbf{u}}, \tilde{\mathbf{v}})]\| \leq \delta. \quad (39)$$

Since the indicator functions $\mathbb{1}_{\{\langle \mathbf{x}, \mathbf{u} \rangle \geq 0\}}$ are scale-invariant, this uniform bound applies to all directions $t\mathbf{u} + (1-t)\mathbf{v}$ appearing in the integrals. Then,

$$\begin{aligned} \left| \left\langle \mathbf{u}, \nabla_{w_1} \widehat{\mathcal{L}} - \nabla_{w_1} \mathcal{L} \right\rangle \right| &\leq v_1 \delta \|\mathbf{u}\| \|\mathbf{h}_1\| + v_1 \delta \|\mathbf{u}\| \|\mathbf{h}_2\| \\ &= v_1 \delta (\|\mathbf{h}_1\| + \|\mathbf{h}_2\|). \end{aligned}$$

Taking the supremum over \mathbf{u} proves the first bound.

Bound for w_2 : Similarly, for an arbitrary unit vector $\mathbf{v} \in \mathbb{S}^{d-1}$:

$$\begin{aligned} \left\langle \mathbf{v}, \nabla_{w_2} \widehat{\mathcal{L}} - \nabla_{w_2} \mathcal{L} \right\rangle &= v_2 \int_0^1 \mathbf{v}^\top (M(t(v_2 \mathbf{w}_2) - (1-t)\mathbf{a}, \mathbf{w}_2) - \mathbb{E}M(\dots)) \mathbf{h}_2 dt \\ &\quad + v_2 \int_0^1 \mathbf{v}^\top (M(t\mathbf{a} + (1-t)(v_1 \mathbf{w}_1), \mathbf{w}_2) - \mathbb{E}M(\dots)) (-\mathbf{h}_1) dt. \end{aligned}$$

Applying the same uniform spectral bound yields:

$$\left| \left\langle \mathbf{v}, \nabla_{w_2} \widehat{\mathcal{L}} - \nabla_{w_2} \mathcal{L} \right\rangle \right| \leq v_2 \delta (\|\mathbf{h}_1\| + \|\mathbf{h}_2\|).$$

This completes the proof of Lemma 11.

Corollary 24 (Bounds in terms of $\|\mathbf{a}\|$) *Further assume that $\|v_1 \mathbf{w}_1\| \leq C \|\mathbf{a}\|$ and $\|v_2 \mathbf{w}_2\| \leq C \|\mathbf{a}\|$. Then:*

$$\left\| \nabla_{w_1} \widehat{\mathcal{L}} - \nabla_{w_1} \mathcal{L} \right\| \leq \tilde{c} v_1 \|\mathbf{a}\|, \quad (40)$$

$$\left\| \nabla_{w_2} \widehat{\mathcal{L}} - \nabla_{w_2} \mathcal{L} \right\| \leq \tilde{c} v_2 \|\mathbf{a}\|, \quad (41)$$

where $\tilde{c} = 2\delta(C + 1)$.

Proof We bound the norms of the error vectors \mathbf{h}_1 and \mathbf{h}_2 using the triangle inequality:

$$\begin{aligned}\|\mathbf{h}_1\| &= \|v_1\mathbf{w}_1 - \mathbf{a}\| \leq \|v_1\mathbf{w}_1\| + \|\mathbf{a}\| \leq (C+1)\|\mathbf{a}\|, \\ \|\mathbf{h}_2\| &= \|v_2\mathbf{w}_2 + \mathbf{a}\| \leq \|v_2\mathbf{w}_2\| + \|\mathbf{a}\| \leq (C+1)\|\mathbf{a}\|.\end{aligned}$$

Substituting these into the theorem's bounds:

$$\begin{aligned}\left\|\nabla_{\mathbf{w}_1}\widehat{\mathcal{L}} - \nabla_{\mathbf{w}_1}\mathcal{L}\right\| &\leq v_1\delta(2(C+1)\|\mathbf{a}\|) = 2\delta(C+1)v_1\|\mathbf{a}\|, \\ \left\|\nabla_{\mathbf{w}_2}\widehat{\mathcal{L}} - \nabla_{\mathbf{w}_2}\mathcal{L}\right\| &\leq v_2\delta(2(C+1)\|\mathbf{a}\|) = 2\delta(C+1)v_2\|\mathbf{a}\|.\end{aligned}$$

■

Appendix D. Proof of Main Theorems

D.1. Proof of Theorem 1 for Landscape Characterization

To prove this theorem, we first show that $v_1\mathbf{w}_1 = \mathbf{a}$, $v_2\mathbf{w}_2 = -\mathbf{a}$ is the global optima. Since $v_1, v_2 > 0$,

$$v_1\phi(\mathbf{w}_1^T\mathbf{x}) - v_2\phi(\mathbf{w}_2^T\mathbf{x}) = \phi(v_1\mathbf{w}_1^T\mathbf{x}) - \phi(v_2\mathbf{w}_2^T\mathbf{x}) = \phi(\mathbf{a}^T\mathbf{x}) - \phi(-\mathbf{a}^T\mathbf{x}) = \mathbf{a}^T\mathbf{x}.$$

Hence, the given weights implement the planted model exactly. Next, we verify that all $v_1, v_2 > 0$, and $\mathbf{w}_1, \mathbf{w}_2$ that satisfy

$$v_1\mathbf{w}_1 - v_2\mathbf{w}_2 = \mathbf{a}, \quad \text{and} \quad \theta = 0$$

are indeed non-strict saddle points of our optimization problem when $k = 2$. We first show that the gradient vanishes. Plugging such $v_1, v_2, \mathbf{w}_1, \mathbf{w}_2$ 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. Furthermore, due to (13), $\nabla_{\mathbf{v}}\mathcal{L}(\boldsymbol{\theta})$ is also $\mathbf{0}$. Next we show that the Hessian at these points are PSD. Plugging the values into (14) we get:

$$\nabla^2\mathcal{L}(\boldsymbol{\theta}) = \frac{1}{2}\begin{bmatrix} \|\mathbf{w}_1\|^2 & -\|\mathbf{w}_1\|\|\mathbf{w}_2\| & v_1\mathbf{w}_1^T & -v_2\mathbf{w}_1^T \\ -\|\mathbf{w}_1\|\|\mathbf{w}_2\| & \|\mathbf{w}_2\|^2 & -v_1\mathbf{w}_2^T & v_2\mathbf{w}_2^T \\ v_1\mathbf{w}_1 & -v_1\mathbf{w}_2 & v_1^2\mathbf{I} & -v_1v_2\mathbf{I} \\ -v_2\mathbf{w}_1 & v_2\mathbf{w}_2 & -v_1v_2\mathbf{I} & v_2^2\mathbf{I} \end{bmatrix}$$

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 $\ell, m, i \in [2]$. This $(2d+2) \times (2d+2)$ matrix has eigenvalues 0 , $\frac{v_1^2+v_2^2}{2}$, and $\frac{\|\mathbf{w}_1\|^2+\|\mathbf{w}_2\|^2+v_1^2+v_2^2}{2}$ (all non-negative) with multiplicities $d+2$, $d-1$, and 1 respectively. Therefore, all the stationary points are in fact non-strict saddle points of the problem.

Finally, we show that there are no other stationary points besides the ones identified above. A necessary condition for $\nabla_{\mathbf{W}} \mathcal{L}(\boldsymbol{\theta}) = \mathbf{0}$ is that any linear combination of the gradient rows must vanish. Specifically, for $v_1, v_2 > 0$, we have:

$$\begin{bmatrix} \frac{1}{v_1} & \frac{1}{v_2} \end{bmatrix} \nabla_{\mathbf{W}} \mathcal{L}(\boldsymbol{\theta}) = \mathbf{0}.$$

By substituting the gradient expression, this implies:

$$\theta (v_1 \mathbf{w}_1 + v_2 \mathbf{w}_2) = \sin \theta (v_1 \|\mathbf{w}_1\| \bar{\mathbf{w}}_2 + v_2 \|\mathbf{w}_2\| \bar{\mathbf{w}}_1).$$

Note that the vectors $v_1 \mathbf{w}_1 + v_2 \mathbf{w}_2$ and $v_1 \|\mathbf{w}_1\| \bar{\mathbf{w}}_2 + v_2 \|\mathbf{w}_2\| \bar{\mathbf{w}}_1$ have identical norms. Taking the norm of both sides, the equality holds only if $|\theta| = |\sin \theta|$, which implies $\theta = 0$, or if the vectors themselves are zero ($v_1 \mathbf{w}_1 + v_2 \mathbf{w}_2 = \mathbf{0}$).

The case $\theta = 0$ corresponds to the non-strict saddle points previously identified. The case $v_1 \mathbf{w}_1 + v_2 \mathbf{w}_2 = \mathbf{0}$ corresponds to the global optima where the two neurons are anti-aligned ($\theta = \pi$) such that their combined contribution exactly implements the target \mathbf{a} . Consequently, there are no other stationary points in the optimization landscape. This completes the proof of the theorem.

D.2. Proof of Theorem 2 for Convergence of the GD Trajectory

To prove this theorem first we note that after $T_1 = \lceil \frac{c_9}{\mu \|\mathbf{a}\|} \rceil$ iterations of GD (i.e. *alignment* phase), using Lemma 6 from Section 5.4 we have with high probability

$$\theta_1^{(T_1)}, \theta_2^{(T_1)} \leq c_4, \quad \text{and} \quad c_1 \sigma \leq v_1^{(T_1)}, v_2^{(T_1)} \leq c_2 \sqrt{\|\mathbf{a}\|}.$$

Using Lemma 7, after $T = T_1 + T_2$ iterations, we have

$$\theta_1^{(T)}, \theta_2^{(T)} \leq c_4, \quad v_1^{(T)}, v_2^{(T)} \geq c_7 \sqrt{\|\mathbf{a}\|}, \quad \text{and} \quad |b_1^{(T)}|, |b_2^{(T)}| \leq c_{11} \|\mathbf{a}\|.$$

Using the definition of the imbalance term, $\|\mathbf{w}_i^{(T)}\|^2 = (v_i^{(T)})^2 + b_i^{(T)}$, we evaluate the weights at the end of the growth phase ($T = T_1 + T_2$). From Lemma 7, we have $v_i^{(T)} \geq c_7 \sqrt{\|\mathbf{a}\|}$ and $|b_i^{(T)}| \leq c_{11} \|\mathbf{a}\|$. Noting that $c_7^2 > c_{11}$ for $c_7 = \frac{1}{4}$, $c_{11} = \frac{1}{50}$; this implies:

$$(c_7^2 - c_{11}) \|\mathbf{a}\| \leq \|\mathbf{w}_1^{(T)}\|^2, \|\mathbf{w}_2^{(T)}\|^2 \leq (c_7^2 + c_{11}) \|\mathbf{a}\|. \quad (42)$$

We now establish that these bounds hold uniformly for all $\tau \geq T$. Lemma 4 ensures that since the angles $\theta_1^{(T)}, \theta_2^{(T)}$ are small, they remain bounded by c_4 for all subsequent iterations. Lemma 5 ensures that if the norms v_i start in the interval $[c_7 \sqrt{\|\mathbf{a}\|}, c_2 \sqrt{\|\mathbf{a}\|}]$, they remain within a fixed range $[c_7 \sqrt{\|\mathbf{a}\|}, c_2 \sqrt{\|\mathbf{a}\|}]$ for all $\tau > T$. Finally, as established in the convergence analysis below, the error $\|\mathbf{h}_1\|^2 + \|\mathbf{h}_2\|^2$ decays geometrically. By Lemma 3, the total drift in the imbalance terms is summable, keeping $|b_i^{(\tau)}|$ uniformly bounded by a constant $\gamma \|\mathbf{a}\|$ for all $\tau \geq T$.

Consequently, there exist universal constants c_{\min} and c_{\max} such that for all $\tau \geq T$:

$$c_{\min} \sqrt{\|\mathbf{a}\|} \leq v_1^{(\tau)}, v_2^{(\tau)}, \|\mathbf{w}_1^{(\tau)}\|, \|\mathbf{w}_2^{(\tau)}\| \leq c_{\max} \sqrt{\|\mathbf{a}\|}. \quad (43)$$

From (43), the conditions for the PL inequality (Lemma 8) and smoothness (Lemma 9) hold uniformly for all $\tau \geq T$, where the constants α and L now depend on c_{\min} and c_{\max} . Specifically, we have:

$$\left\| \nabla_{\mathbf{w}_1} \mathcal{L} \left(\mathbf{v}^{(\tau)}, \mathbf{W}^{(\tau)} \right) \right\|^2 + \left\| \nabla_{\mathbf{w}_2} \mathcal{L} \left(\mathbf{v}^{(\tau)}, \mathbf{W}^{(\tau)} \right) \right\|^2 \geq \alpha c_{\min}^2 \|\mathbf{a}\| \mathcal{L} \left(\mathbf{v}^{(\tau)}, \mathbf{W}^{(\tau)} \right),$$

and

$$\left\| \nabla^2 \mathcal{L} \left(\mathbf{v}^{(\tau)}, \mathbf{W}^{(\tau)} \right) \right\|_F \leq L \|\mathbf{a}\|.$$

Next, we keep track of the *population* loss while performing GD updates on the *empirical* loss. Let $\boldsymbol{\theta}$ denote $\begin{bmatrix} \mathbf{v} \\ \text{vect}(\mathbf{W}) \end{bmatrix}$, and also define errors vectors $\mathbf{h}_1 = v_1 \mathbf{w}_1 - \mathbf{a}$, $\mathbf{h}_2 = v_2 \mathbf{w}_2 + \mathbf{a}$. We note that $\left\| \text{diag} \left(\begin{bmatrix} v_1 \\ v_2 \end{bmatrix} \right) \mathbf{W} - \mathbf{W}^* \right\|_F^2 = \|\mathbf{h}_1\|^2 + \|\mathbf{h}_2\|^2$. For all $\tau \geq T$ we have

$$\begin{aligned} \mathcal{L} \left(\boldsymbol{\theta}^{(\tau+1)} \right) &= \mathcal{L} \left(\boldsymbol{\theta}^{(\tau)} - \mu \nabla \widehat{\mathcal{L}} \left(\boldsymbol{\theta}^{(\tau)} \right) \right) \\ &\stackrel{(a)}{\leq} \mathcal{L} \left(\boldsymbol{\theta}^{(\tau)} \right) - \mu \left\langle \nabla \mathcal{L} \left(\boldsymbol{\theta}^{(\tau)} \right), \nabla \widehat{\mathcal{L}} \left(\boldsymbol{\theta}^{(\tau)} \right) \right\rangle + \frac{L \|\mathbf{a}\|}{2} \mu^2 \left\| \nabla \widehat{\mathcal{L}} \left(\boldsymbol{\theta}^{(\tau)} \right) \right\|^2 \\ &= \mathcal{L} \left(\boldsymbol{\theta}^{(\tau)} \right) - \mu \left\| \nabla \mathcal{L} \left(\boldsymbol{\theta}^{(\tau)} \right) \right\|^2 - \mu \left\langle \nabla \mathcal{L} \left(\boldsymbol{\theta}^{(\tau)} \right), \nabla \widehat{\mathcal{L}} \left(\boldsymbol{\theta}^{(\tau)} \right) - \nabla \mathcal{L} \left(\boldsymbol{\theta}^{(\tau)} \right) \right\rangle \\ &\quad + \frac{L \|\mathbf{a}\|}{2} \mu^2 \left\| \nabla \widehat{\mathcal{L}} \left(\boldsymbol{\theta}^{(\tau)} \right) \right\|^2 \\ &\stackrel{(b)}{\leq} \mathcal{L} \left(\boldsymbol{\theta}^{(\tau)} \right) - \mu \left\| \nabla \mathcal{L} \left(\boldsymbol{\theta}^{(\tau)} \right) \right\|^2 + \eta \mu \left\| \nabla \mathcal{L} \left(\boldsymbol{\theta}^{(\tau)} \right) \right\|^2 + \frac{\mu}{\eta} \left\| \nabla \widehat{\mathcal{L}} \left(\boldsymbol{\theta}^{(\tau)} \right) - \nabla \mathcal{L} \left(\boldsymbol{\theta}^{(\tau)} \right) \right\|^2 \\ &\quad + \frac{L \|\mathbf{a}\|}{2} \mu^2 \left\| \nabla \widehat{\mathcal{L}} \left(\boldsymbol{\theta}^{(\tau)} \right) \right\|^2 \\ &\stackrel{(c)}{\leq} \mathcal{L} \left(\boldsymbol{\theta}^{(\tau)} \right) - \mu (1 - \eta) \left\| \nabla \mathcal{L} \left(\boldsymbol{\theta}^{(\tau)} \right) \right\|^2 + \frac{\mu}{\eta} \left\| \nabla \widehat{\mathcal{L}} \left(\boldsymbol{\theta}^{(\tau)} \right) - \nabla \mathcal{L} \left(\boldsymbol{\theta}^{(\tau)} \right) \right\|^2 \\ &\quad + L \|\mathbf{a}\| \mu^2 \left(\left\| \nabla \mathcal{L} \left(\boldsymbol{\theta}^{(\tau)} \right) \right\|^2 + \left\| \nabla \widehat{\mathcal{L}} \left(\boldsymbol{\theta}^{(\tau)} \right) - \nabla \mathcal{L} \left(\boldsymbol{\theta}^{(\tau)} \right) \right\|^2 \right) \\ &= \mathcal{L} \left(\boldsymbol{\theta}^{(\tau)} \right) - \mu (1 - \eta - \mu L \|\mathbf{a}\|) \left\| \nabla \mathcal{L} \left(\boldsymbol{\theta}^{(\tau)} \right) \right\|^2 \\ &\quad + \left(\frac{\mu}{\eta} + \mu^2 L \|\mathbf{a}\| \right) \left\| \nabla \widehat{\mathcal{L}} \left(\boldsymbol{\theta}^{(\tau)} \right) - \nabla \mathcal{L} \left(\boldsymbol{\theta}^{(\tau)} \right) \right\|^2 \\ &\stackrel{(d)}{\leq} \mathcal{L} \left(\boldsymbol{\theta}^{(\tau)} \right) - \alpha c_1^2 \|\mathbf{a}\| \mu (1 - \eta - \mu L \|\mathbf{a}\|) \mathcal{L} \left(\boldsymbol{\theta}^{(\tau)} \right) \\ &\quad + \left(\frac{\mu}{\eta} + \mu^2 L \|\mathbf{a}\| \right) \left\| \nabla \widehat{\mathcal{L}} \left(\boldsymbol{\theta}^{(\tau)} \right) - \nabla \mathcal{L} \left(\boldsymbol{\theta}^{(\tau)} \right) \right\|^2 \\ &\stackrel{(e)}{\leq} (1 - \alpha c_{\min}^2 \|\mathbf{a}\| \mu (1 - \eta - \mu L \|\mathbf{a}\|)) \mathcal{L} \left(\boldsymbol{\theta}^{(\tau)} \right) \end{aligned}$$

$$\begin{aligned}
 & + \left(\frac{\mu}{\eta} + \mu^2 L \|\mathbf{a}\| \right) \left(\left(v_1^{(\tau)} \right)^2 + \left(v_2^{(\tau)} \right)^2 + \left\| \mathbf{w}_1^{(\tau)} \right\|^2 + \left\| \mathbf{w}_2^{(\tau)} \right\|^2 \right) \delta^2 \left(\left\| \mathbf{h}_1^{(\tau)} \right\|^2 + \left\| \mathbf{h}_2^{(\tau)} \right\|^2 \right) \\
 & \stackrel{(f)}{\leq} \left(1 - \alpha c_{\min}^2 \|\mathbf{a}\| \mu (1 - \eta - \mu L \|\mathbf{a}\|) \right) \mathcal{L} \left(\boldsymbol{\theta}^{(\tau)} \right) \\
 & \quad + \left(\frac{\mu}{\eta} + \mu^2 L \|\mathbf{a}\| \right) 4c_{\max}^2 \|\mathbf{a}\| \delta^2 \left(\left\| \mathbf{h}_1^{(\tau)} \right\|^2 + \left\| \mathbf{h}_2^{(\tau)} \right\|^2 \right) \\
 & \stackrel{(g)}{\leq} \left(1 - \alpha c_{\min}^2 \|\mathbf{a}\| \mu (1 - \eta - \mu L \|\mathbf{a}\|) + 20 \left(\frac{\mu}{\eta} + \mu^2 L \|\mathbf{a}\| \right) 4c_{\max}^2 \|\mathbf{a}\| \delta^2 \right) \mathcal{L} \left(\boldsymbol{\theta}^{(\tau)} \right)
 \end{aligned}$$

where (a) follows from the quadratic upper bound and smoothness of \mathcal{L} , (b) follows from $\left\langle \sqrt{\eta} \mathbf{a}, \frac{1}{\sqrt{\eta}} \mathbf{b} \right\rangle \leq \eta \|\mathbf{a}\|^2 + \frac{1}{\eta} \|\mathbf{b}\|^2$ for any $\eta > 0$, (c) follows from triangle inequality, (d) follows from PL inequality, (e) follows from Lemma 10 and gradient identity (13), (f) follows from upper bounds on the norms, and finally (g) follows from applying the population loss lower bound (Lemma 13). Set $\mu = \frac{\bar{\mu}}{\|\mathbf{a}\|}$, $\eta = \frac{1}{4}$. We have

$$\begin{aligned}
 \mathcal{L} \left(\boldsymbol{\theta}^{(\tau+1)} \right) & \leq \left(1 - \alpha c_{\min}^2 \bar{\mu} \left(\frac{3}{4} - \bar{\mu} L \right) + \left(\frac{4\bar{\mu}}{3} + \bar{\mu}^2 L \right) 80c_{\max}^2 \delta^2 \right) \mathcal{L} \left(\boldsymbol{\theta}^{(\tau)} \right) \\
 & = \left(1 - \left(\frac{3\alpha c_{\min}^2}{4} - \frac{320c_{\max}^2 \delta^2}{3} \right) \bar{\mu} + (\alpha c_{\min}^2 + 80c_{\max}^2 \delta^2) L \bar{\mu}^2 \right) \mathcal{L} \left(\boldsymbol{\theta}^{(\tau)} \right).
 \end{aligned}$$

We now choose δ and $\bar{\mu}$ so that the quadratic factor above yields a strict contraction. First require that the linear coefficient is positive, i.e.

$$\delta^2 \leq \frac{9\alpha c_{\min}^2}{1280c_{\max}^2}.$$

Fix any such δ (this is ensured by Lemma 10 by taking n sufficiently large). Next, define,

$$a := \frac{3\alpha c_{\min}^2}{4} - \frac{320c_{\max}^2 \delta^2}{3} \quad \text{and} \quad b := (\alpha c_{\min}^2 + 80c_{\max}^2 \delta^2) L.$$

If we further choose

$$\bar{\mu} \leq \frac{a}{2b},$$

then the quadratic term is dominated by the linear term, and we have

$$1 - a\bar{\mu} + b\bar{\mu}^2 \leq 1 - \frac{a}{2}\bar{\mu}.$$

Consequently, for all $\tau \geq T$,

$$\mathcal{L} \left(\boldsymbol{\theta}^{(\tau+1)} \right) \leq (1 - c\bar{\mu}) \mathcal{L} \left(\boldsymbol{\theta}^{(\tau)} \right),$$

where $c := \frac{a}{2} > 0$ is a numerical constant depending only on $\alpha, c_{\min}, c_{\max}$. Iterating this inequality for all $\tau > T$ yields geometric decrease of the population loss

$$\mathcal{L} \left(\boldsymbol{\theta}^{(\tau)} \right) \leq (1 - c\bar{\mu})^{(\tau-T)} \mathcal{L} \left(\boldsymbol{\theta}^{(T)} \right).$$

Finally, we apply the population loss lower bound (Lemma 13) to lower bound the left-hand side:

$$\mathcal{L}(\boldsymbol{\theta}^{(\tau)}) \geq \frac{1}{\tilde{c}} \left\| \text{diag} \left(\begin{bmatrix} v_1^{(\tau)} \\ v_2^{(\tau)} \end{bmatrix} \right) \mathbf{W}^{(\tau)} - \mathbf{W} \right\|_F^2.$$

By further upper bounding the right-hand side using the fact that the ReLU activation is 1-Lipschitz, we have $\mathcal{L}(\boldsymbol{\theta}^{(T)}) \leq \left\| \text{diag} \left(\begin{bmatrix} v_1^{(T)} \\ v_2^{(T)} \end{bmatrix} \right) \mathbf{W}^{(T)} - \mathbf{W}^* \right\|_F^2$. Combining these yields:

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

for some numerical constant $\tilde{c} > 0$. This shows geometric convergence of the GD iterates, completing the proof of Theorem 2.

Appendix E. Additional Experimental Results

E.1. Pairing-up Behavior for $r \geq 3$

In this section, we present additional results on the pairing behavior of \mathbf{w}_i and \mathbf{v}_i for different values of r . Although our theoretical analysis is limited to the scalar output setting, for our experiments we also consider multi-dimensional outputs. 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 6, we pick $r = 3$ and $k = 6$. As for the target function, we pick $\mathbf{a}_1, \mathbf{a}_2, \mathbf{a}_3$ to be $\mathbf{e}_1, \mathbf{e}_2, \mathbf{e}_3$ 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 \mathbf{v}_i 's and \mathbf{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 6, we observe that $\mathbf{v}_3^{(\infty)} \approx -\mathbf{v}_4^{(\infty)}$, $\mathbf{v}_1^{(\infty)} \approx -\mathbf{v}_5^{(\infty)}$, and $\mathbf{v}_2^{(\infty)} \approx -\mathbf{v}_6^{(\infty)}$ which also holds similarly for \mathbf{w}_i 's as well. This suggests that after a permutation of the hidden units, we get

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

which can be considered as a natural extension to the $\mathbf{v}_i = \pm 1$ pattern in the single output setting.

Beyond the $r = 3$ case, 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.

Appendix F. Related Work

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

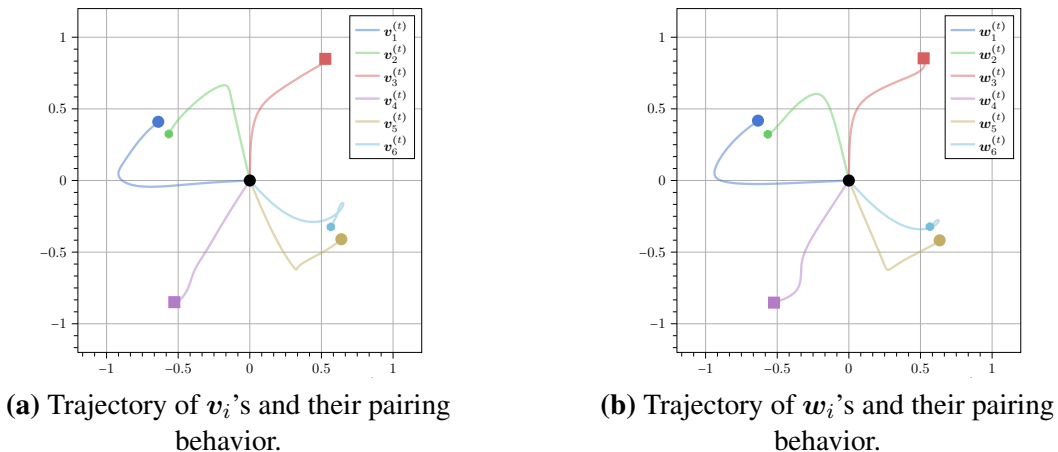


Figure 6: **Pairing pattern in multi-dimensional setting.** We train the network from small initialization when exactly parameterized ($k = 6$ and $r = 3$). On left (a), we depict the trajectories of individual weights in the outer layer (v_i 's) across iterations. We observe that the weights at convergence 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.

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, 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. This trajectory-level perspective (i.e. multi-phase analysis) is also explored in recent works on gradient descent dynamics and implicit bias under large learning rates (Wang et al., 2022, 2023), see also additional prior work (Stöger and Soltanolkotabi, 2021; Soltanolkotabi et al., 2023) on this topic. However, these works focus on matrix factorization and more general nonconvex objectives

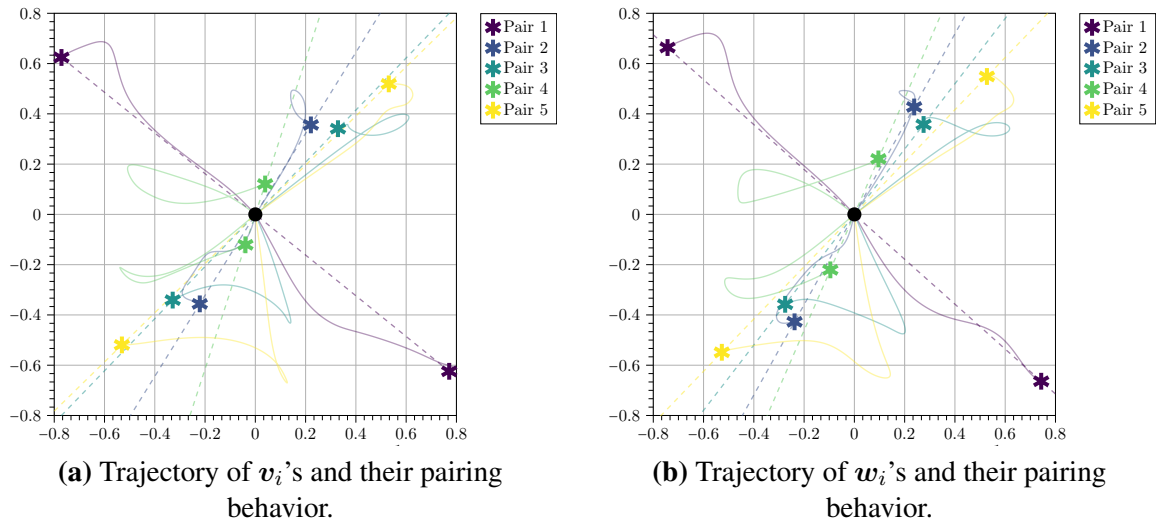


Figure 7: **Pairing pattern for $r = 5$.** We train the network from small initialization when exactly parameterized ($k = 10$ and $r = 5$). On left (a), we depict the trajectories of individual weights in the outer layer (v_i 's) across iterations. Each pair is indicated by the same color and the dashed 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.

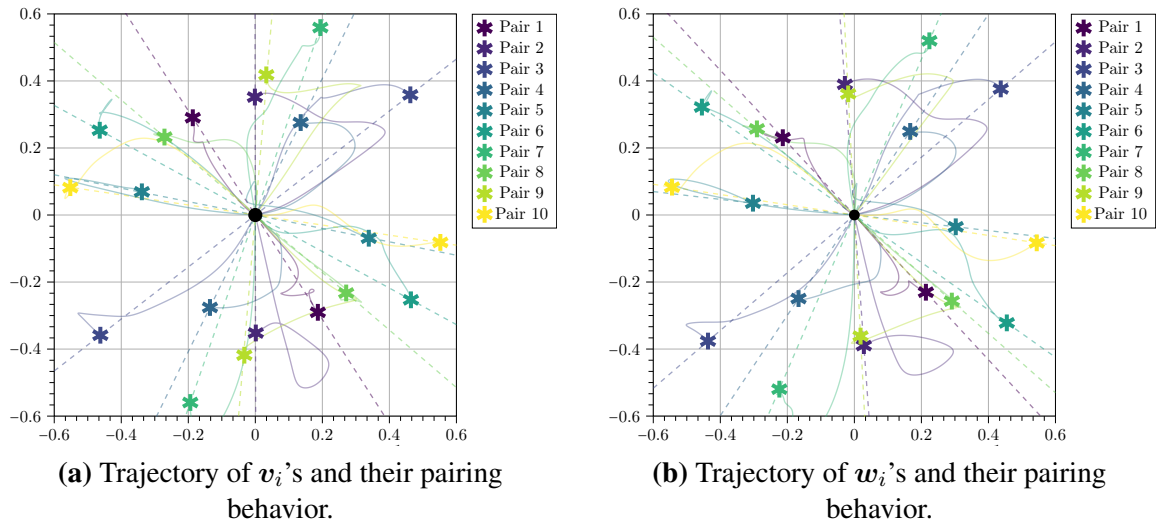


Figure 8: **Pairing pattern for $r = 10$.** We train the network from small initialization when exactly parameterized ($k = 20$ and $r = 10$). On left (a), we depict the trajectories of individual weights in the outer layer (v_i 's) across iterations. Each pair is indicated by the same color and the 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.

rather than neural network training.

Gradient-based analysis for neural networks: A recent line of work is concerned with con-

necting 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.

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

We would also like to discuss the difference between our work and a few other papers [Zhong et al. \(2017\)](#); [Zhang et al. \(2018\)](#); [Zhu et al. \(2025\)](#); [Ren et al. \(2025\)](#) that have planted one-hidden layer models. These papers differ in at least one of three ways focusing on (1) local analysis,

(2) have sub-optimal sample complexity, and/or (3) assume non-negative outer layer weights. For instance, [Zhong et al. \(2017\)](#) utilize tensor initialization, performing a local analysis rather than a global GD analysis. This local analysis however can not be used to analyze the linear target setting. Indeed, as noted in Remark 4.3 of their work, their analysis requires \mathbf{W}^* to be full-rank which does not hold in the linear setting (where the rows of the weight matrix are negatives of each other leading to a minimum singular value is zero). Furthermore, this result also requires resampling the data points at each iteration to ensure convergence of gradient descent where as we use the same samples across all iterations. On a related note, their sample complexity has polynomial dependency on many problem parameters (Theorem 4.2) whereas our proof only requires sample size linear in input dimension d .

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

[Wu et al. \(2018\)](#) consider the setting when student and teacher networks both have 2 neurons. In particular, when the teachers are *orthogonal*, and the outer weights are all ones; they demonstrated an interesting result that the landscape is benign and all saddles are strict. In contrast, the landscape in our problem include non-strict saddle points significantly complicating gradient descent analysis. In more recent work, [Ren et al. \(2025\)](#) study the complexity of learning a planted model with *orthogonal* planted directions, quadratic activations, and non-negative outer weights. They obtain interesting results on the scaling laws of the MSE loss via a multi-phase analysis. However, this problem setting is substantially different due to the difference between the activation and the orthogonal weights in the planted model that makes the landscape benign per above discussion. More recently, [Zhu et al. \(2025\)](#) also consider learning multiple *orthogonal* ReLU neurons in a teacher-student framework with outer layer weights fixed to all ones. As just discussed, having orthogonal teacher weights leads to a much more benign landscape. Moreover, assumptions in the aforementioned works strictly exclude the linear target setting, where the outer layer must contain negative coefficients. Furthermore, they impose strong restrictions on the initialization. Specifically, they look at the convergence after “weak alignment” where for each student neuron there exists only one teacher neuron that is not near perpendicular. Our results on the other hand can handle random initializations where student neurons *could* be perpendicular to the target direction. That said, their analysis can handle over-parametrization ($k \gg k^*$) and teacher networks with more than 2 neurons.

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

“However, item 3 is quite restrictive: it is needed to ensure that the volume of the activation cone containing β^* does not vanish when $n \rightarrow \infty$. A similar assumption is considered by [Chistikov et al. \(2023\)](#); [Tsoy and Konstantinov \(2024\)](#), for similar reasons. Additionally, Condition 4 ensures that

$\mathbb{E}_x[xx^T]\beta^*$ and β^* are in the same activation cone. This assumption allows the training dynamics to remain within a single cone after the early alignment phase, significantly simplifying our analysis.”

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

Finally, we note that recent interesting works by [Shevchenko et al. \(2022\)](#) and [Kögler et al. \(2024\)](#) provide a complementary perspective on nonlinear autoencoders, with the philosophical goal of understanding what happens when the autoencoder departs from PCA-like or linear-function learning. To isolate this genuinely nonlinear regime, these works study sign activations, and more generally odd activations, and characterize the optima of the corresponding population loss. They further show that these optima are achieved by suitable gradient-based procedures, including gradient flow with tied encoder/decoder weights and a projected/alternating gradient method in which the decoder is optimized to completion at each update of the encoder. In contrast, our setting studies standard empirical training of a ReLU network for a planted linear supervised target, where both layers are optimized simultaneously throughout gradient descent and the main challenge is controlling the full finite-sample trajectory through non-strict saddle regions.