
Simplicity Bias of Two-Layer Networks beyond Linearly Separable Data

Nikita Tsoy¹ Nikola Konstantinov¹

Abstract

Simplicity bias, the propensity of deep models to over-rely on simple features, has been identified as a potential reason for limited out-of-distribution generalization of neural networks (Shah et al., 2020). Despite the important implications, this phenomenon has been theoretically confirmed and characterized only under strong dataset assumptions, such as linear separability (Lyu et al., 2021). In this work, we characterize simplicity bias for general datasets in the context of two-layer neural networks initialized with small weights and trained with gradient flow. Specifically, we prove that in the early training phases, network features cluster around a few directions that do not depend on the size of the hidden layer. Furthermore, for datasets with an XOR-like pattern, we precisely identify the learned features and demonstrate that simplicity bias intensifies during later training stages. These results indicate that features learned in the middle stages of training may be more useful for OOD transfer. We support this hypothesis with experiments on image data.

1. Introduction

Out-of-distribution (OOD) generalization is a key challenge towards the widespread adoption of machine learning. Specifically, since training data may not always cover all possible test scenarios, networks often rely on shortcuts: spurious rules that hold on the training distribution but not in more complicated real-world situations (Geirhos et al., 2020). For example, convolutional networks often prioritize texture over shape (Geirhos et al., 2019), or transformers might rely on simplistic heuristics in natural language inference (McCoy et al., 2019).

One possible mechanism behind shortcuts is *simplicity bias*,

¹INSAIT, Sofia University, Bulgaria. Correspondence to: Nikita Tsoy <nikita.tsoy@insait.ai>.

the propensity of neural networks to rely only on “simple” features. As Shah et al. (2020) demonstrated, this bias might be persistent and hurt OOD generalization in image classification tasks. Simplicity bias is also a peculiar phenomenon from a theoretical perspective. Since many neural architectures are universal function approximators (Cybenko, 1989; Hornik et al., 1989), one might hope that models will learn other, more sophisticated patterns within the data.

Despite the importance of simplicity bias, a thorough theoretical understanding of this phenomenon is still lacking. To the best of our knowledge, existing works on simplicity bias (Lyu et al., 2021; Safran et al., 2022; Morwani et al., 2023) only demonstrate its emergence by employing stringent assumptions, which restrict training data to be linearly separable or one-dimensional.

Contributions We give the first proof of the existence of simplicity bias and a precise mathematical characterization of the features learned during training beyond linearly separable data. We do so for two-layer neural networks trained with gradient flow from a small initialization, a model popular in the theoretical literature (Luo et al., 2021a; Lyu et al., 2021; Boursier et al., 2022). We characterize simplicity bias as a property that only a small set of *prominent* neurons governs the behavior of the network. These prominent neurons cluster in several directions, which do not depend on the size of the hidden layer.

Specifically, our theoretical analysis divides training into three stages. During the first two stages (Section 4), in which the weights grow from small to constant scale, we prove for general datasets that the most prominent features recovered by the training dynamics cluster around the extrema of a data-dependent function, which does not depend on the number of neurons and disentangles the interactions between them. For the last stage, where the network converges to zero loss, we prove that simplicity bias can become extreme even in non-linearly-separable datasets. In this stage, we cover the case of XOR-like data under an assumption about the convergence of a 4-neuron network, which we validate experimentally. On a methodological level, our work generalizes the analysis of Lyu et al. (2021) beyond linearly separable data and provides an implicit description of the most prominent features for general datasets.

Our theoretical findings additionally lead to a hypothesis with potential practical implications: networks trained to a very small loss may be prone to stronger simplicity bias and, hence, may be harder to finetune to new tasks. We test this intuition with experiments on a domino dataset of MNIST-CIFAR10 pairs (Shah et al., 2020) and observe experimental support for our hypothesis.¹

2. Related Work

Simplicity bias of features The simplicity bias of neural networks was observed and linked to generalization and OOD performance in prior work (Valle-Perez et al., 2019; Shah et al., 2020). Several works show that two-layer networks provably learn a linear decision boundary on linearly separable datasets (Brutzkus et al., 2018; Pellegrini & Biroli, 2020; Sarussi et al., 2021; Phuong & Lampert, 2021; Lyu et al., 2021; Englert & Lazic, 2022; Frei et al., 2023b; Kou et al., 2023; Morwani et al., 2023; Wang & Ma, 2023; Chistikov et al., 2023; Min et al., 2024). While we use some of the techniques developed by these works, we focus on the non-linearly-separable case, which requires further theoretical analysis. Safran et al. (2022) prove that, on one-dimensional data, two-layer networks converge to a model with few linear regions. In contrast, we study datasets in \mathbb{R}^d . Brutzkus & Globerson (2019) also show simplicity bias for XOR-like data, but only for a 4-point dataset in \mathbb{R}^2 , while we focus on datasets of arbitrary size in \mathbb{R}^d .

Simplicity bias of training dynamics A direction related to our work is simplicity bias in terms of training dynamics, the propensity of neural networks to learn simple patterns first in training. This property was observed in several empirical works (Arpit et al., 2017; Xu et al., 2019; Rahaman et al., 2019; Kalimeris et al., 2019) and demonstrated theoretically in certain settings (Arora et al., 2019; Basri et al., 2020; Luo et al., 2021a; Bowman & Montufar, 2022). Another closely related topic is the distributional simplicity bias proposed by Refinetti et al. (2023). In contrast to these works, we seek to characterize the learned features instead of analyzing some complexity invariant.

Small initialization and initial condensation Our results in Section 4 provably demonstrate initial condensation, the propensity of neural networks to condense neurons in few directions during early stages of training, for two-layer networks with small initialization. Several theoretical works have analyzed this phenomenon previously. Maennel et al. (2018) study condensation in regression and classification problems with two-layer ReLU networks. While they also recognize the function G we use to describe the network

¹Replication files are available at <https://github.com/nikita-tsoy98/simplicity-bias-beyond-linear-replication>

dynamics as an important factor in learning, their derivations and discussions of the near-zero initialization regime, which we study in our paper, are informal. In addition, their bounds for the speed of neuron alignment are insufficient to differentiate between prominent and non-prominent neurons, as we do in Section 4. Zhou et al. (2022) give another theoretical description of condensation in regression tasks. However, their results do not apply to ReLU activation or our differentiable approximation of ReLU (due to irregularity at 0). Boursier et al. (2022) also identify the condensation phenomenon in regression tasks, but only for orthogonal data, while we put considerably milder assumptions on the data. Additionally, Xu & Du (2023) recognize the condensation phenomenon in a regression setting with a one-neuron teacher network, which is similar to a linearly separable case because the direction of the gradients will be correlated with the direction of the teacher neuron.

Finally, the concurrent work of Boursier & Flammarion (2024) proves an alignment result similar to that in our Theorem 4.1. However, their result does not quantify the differences in the neurons’ growth rates and, hence, can not differentiate between prominent and non-prominent neurons. Moreover, there are several technical differences. Boursier & Flammarion (2024) work directly with (leaky) ReLU and make milder assumptions on the initialization. However, they require the function G to have no saddle points and only show alignment results for neurons that satisfy a certain technical condition (Condition 1 in their manuscript).

Small initialization and mean-field regime Since we work in the small initialization regime of two-layer networks, our setting is similar to the mean-field regime (Chizat & Bach, 2018; Mei et al., 2018; Wei et al., 2019; Li et al., 2020; Ge et al., 2021). However, in our setting, we do not increase the number of neurons when we decrease the initialization scale (which corresponds to the condensed regime in the classification of Luo et al., 2021b). Thus, in our limit, at the beginning of training, the neurons evolve independently, while, in the mean-field limit, they interact via the velocity field (the derivative of a loss function). At the same time, some techniques in our work are similar to the mean-field techniques since both we and the mean-field works analyze the behavior of a loss function near zero.

3. Setting

Throughout the paper, we analyze feature learning on a binary classification problem with a two-layer network initialized with small random parameter values.

Notation We use the following notation: v^j is j th component of vector v , so that $v = (v^1, v^2, \dots, v^d)^T$, $\|v\| := \sqrt{\sum_{j=1}^d (v^j)^2}$ is the usual l_2 -norm of v , $\hat{v} := \frac{v}{\|v\|}$ is

the unit direction of \mathbf{v} , $\mathbf{P}_v := \mathbf{I} - \hat{\mathbf{v}}\hat{\mathbf{v}}^\top$ is the projector on the space orthogonal to \mathbf{v} , $[k] := \{1, \dots, k\}$, $\mathbb{S}^{d-1} := \{\mathbf{x} \in \mathbb{R}^d \mid \|\mathbf{x}\| = 1\}$ is the unit sphere, $\mathbb{D}^d := \{\mathbf{x} \in \mathbb{R}^d \mid \|\mathbf{x}\| < 1\}$ is the unit disk.

Objective Denote by $f(\boldsymbol{\theta}, \cdot)$ a network parameterized by $\boldsymbol{\theta}$. The sign of f stands for the classification result. Denote by $D := (\mathbf{x}_i \in \mathbb{R}^d, y_i \in \{-1, 1\})_{i=1}^n$ an arbitrary training dataset, such that $\forall i \|\mathbf{x}_i\| \leq 1$. We consider networks trained to minimize the cross-entropy loss

$$L(\boldsymbol{\theta}) := \frac{1}{n} \sum_{i=1}^n \ell(f(\boldsymbol{\theta}, \mathbf{x}_i) y_i), \text{ where } \ell(z) := \ln(1 + e^{-z}).$$

Architecture We consider two-layer networks

$$f(\boldsymbol{\theta}, \mathbf{x}) := \sum_{j=1}^m u_j \phi(\mathbf{v}_j, \mathbf{x}),$$

where $\boldsymbol{\theta} := (u_1, \dots, u_m, \mathbf{v}_1^\top, \dots, \mathbf{v}_m^\top)^\top$, $u_j \in \mathbb{R}$ and $\mathbf{v}_j \in \mathbb{R}^d$, are the network parameters and ϕ is an activation function. We denote $\forall A \subseteq [m] \|\boldsymbol{\theta}\|_A := \max_{j \in A} \max(|u_j|, \|\mathbf{v}_j\|)$.

One of the most commonly used activation functions is ReLU, for which $\phi(\mathbf{v}, \mathbf{x}) = (\mathbf{v}^\top \mathbf{x})_+ = \max(\mathbf{v}^\top \mathbf{x}, 0)$. In our paper, for the technical reasons, we need the activation to be smooth to avoid some degenerate cases in the dynamics of gradient flow. Since this property does not hold for ReLU, we consider a differentiable approximation,

$$\phi(\mathbf{v}, \mathbf{x}) := \phi_{Q, \xi}(\mathbf{v}, \mathbf{x}) := \int_{\mathbb{R}^d} (\mathbf{v}^\top (\mathbf{x} + \xi \mathbf{z}))_+ Q(d\mathbf{z}),$$

where $\xi > 0$, Q is the uniform measure on \mathbb{D}^d , and $(z)_+ := \max(0, z)$. For the purposes of our analysis, ξ could be set to be much smaller than machine precision. Thus, there is no practical difference between our activation and ReLU. We also note that in our experiments in Sections 6 and 7 we use the usual ReLU activation.

We call a function $f(x, y)$ k -positively homogeneous in x if $f(cx, y) = c^k f(x, y)$ for all vectors x, y and all $c > 0$. Notice that $\phi(\mathbf{v}, \mathbf{x})$ is 1-positively homogeneous in \mathbf{v} , i.e., $\forall c > 0 \phi(c\mathbf{v}, \mathbf{x}) = c\phi(\mathbf{v}, \mathbf{x})$. This property implies that $f(\boldsymbol{\theta}, \mathbf{x})$ is 2-positively homogeneous in $\boldsymbol{\theta}$.

Optimization We consider the training of f via gradient flow, $\frac{d\boldsymbol{\theta}}{dt} = -\nabla L(\boldsymbol{\theta})$, which implies dynamics

$$\begin{aligned} \frac{du_j}{dt} &= \frac{1}{n} \sum_{i=1}^n (-\ell'(f(\boldsymbol{\theta}, \mathbf{x}_i) y_i)) \phi(\mathbf{v}_j, \mathbf{x}_i) y_i, \\ \frac{d\mathbf{v}_j}{dt} &= \frac{u_j}{n} \sum_{i=1}^n (-\ell'(f(\boldsymbol{\theta}, \mathbf{x}_i) y_i)) \nabla_{\mathbf{v}} \phi(\mathbf{v}_j, \mathbf{x}_i) y_i. \end{aligned} \quad (1)$$

We initialize the system with a small weights $\boldsymbol{\theta}(0) = \sigma \boldsymbol{\theta}^0$, where $\sigma \approx 0$, similarly Lyu et al. (2021). For simplicity, we assume that $|u_j^0| = \|\mathbf{v}_j^0\|$. The results of Du et al. (2018) imply that $u_j^2 - \|\mathbf{v}_j\|^2 = \text{const}$. Then, $\text{sign}(u_j) =: s_j$ is constant (Lemma 1, Boursier et al., 2022).

4. Simplicity Bias for General Data

First, we analyze Equation (1) in the early (Section 4.1) and middle (Section 4.2) training phases, in which $\boldsymbol{\theta}$ remains small or grow to a constant scale, respectively. These phases correspond to Phases 1 and 2 of Lyu et al. (2021).

Key challenge The key challenge in characterizing the features learned throughout training is the lack of universal training invariants to trace. To compensate for this, most of the works that describe features of neural networks make structural assumptions about the dataset, such as linear separability (Lyu et al., 2021), orthogonality or near-orthogonality (e.g., Brutzkus & Globerson, 2019; Phuong & Lampert, 2021; Frei et al., 2023a) or high-dimensionality (e.g., Ba et al., 2022).

In contrast to these works, we do not make structural assumptions about the data. Instead, we exploit that, for small weights, a simpler data-dependent function, G , which does not depend on the number of neurons, can approximate the network dynamics. Specifically, we link the features learned by the original system to the global extrema of G . While our results do not explicitly characterize the learned features, they are sufficient to prove the presence of simplicity bias.

4.1. Feature Learning from a Small Initialization

Disentangling training dynamics First, we informally motivate our approximation of Equation (1) for a small initialization. By the mean value theorem,

$$\begin{aligned} \ell'(f(\boldsymbol{\theta}, \mathbf{x}) y) - \ell'(0) &= \ell''(\zeta) f(\boldsymbol{\theta}, \mathbf{x}) y \\ \implies |\ell'(f(\boldsymbol{\theta}, \mathbf{x}) y) - \ell'(0)| &\leq |f(\boldsymbol{\theta}, \mathbf{x})| \sup_{z \in \mathbb{R}} |\ell''(z)|, \end{aligned}$$

for some $\zeta \in [0, f(\boldsymbol{\theta}, \mathbf{x}) y]$. Since $f(\boldsymbol{\theta}, \mathbf{x})$ is 2-homogeneous, we get $|\ell'(f(\boldsymbol{\theta}, \mathbf{x}) y) - \ell'(0)| = \mathcal{O}(\|\boldsymbol{\theta}\|^2)$. Thus, when $\sigma \approx 0$, Equation (1) behaves similarly to the following system with linearized loss (Maennel et al., 2018),

$$\frac{du_j^l}{dt} := G(\mathbf{v}_j^l) := \frac{1}{n} \sum_{i=1}^n (-\ell'(0)) \phi(\mathbf{v}_j^l, \mathbf{x}_i) y_i, \quad (2)$$

$$\frac{d\mathbf{v}_j^l}{dt} := u_j^l \nabla G(\mathbf{v}_j^l),$$

where $u_j^l(0) = u_j(0)$ and $\mathbf{v}_j^l(0) = \mathbf{v}_j(0)$.

Note that the neurons in Equation (2) evolve independently of each other, while the neurons in Equation (1) interact

via $\ell'(f(\boldsymbol{\theta}, \mathbf{x}_i)y_i)$. This property significantly facilitates the analysis of Equation (2) compared to the original system.

The following theorem formalizes the link between the two systems and links the features learned by the original dynamics (1) to the global extrema of the function G .

Theorem 4.1 (Proof in Appendix B). *Assume that $\boldsymbol{\theta}$ follows Equation (1), $\forall i \|\mathbf{x}_i\| \leq 1$, $d \geq 2$, and d is odd. Then $\exists \kappa^* > 0$, $P \subseteq [m]$, $(\kappa_j > 0, u_j^* \in \mathbb{R}, \hat{\mathbf{v}}_j^* \in \mathbb{S}^{d-1})_{j=1}^m$ such that for $\sigma = r^{1+\kappa^*}$, $T_1 := \frac{1}{\lambda} \ln(\frac{r}{\sigma})$, and $r \rightarrow 0$, we get*

$$\begin{aligned} \forall j \in P \quad & |u_j(T_1) - ru_j^*| \leq O(r^{1+\kappa^*}), \\ & \|\hat{\mathbf{v}}_j(T_1) - \hat{\mathbf{v}}_j^*\| \leq O(r^{\kappa^*}), \quad s_j G(\hat{\mathbf{v}}_j^*) = \lambda, \\ \forall j \in R \quad & |u_j(T_1)| = \|v_j(T_1)\| \leq O(r^{1+\kappa_j}), \end{aligned}$$

where $R := [m] \setminus P$, $\lambda := \max_{\hat{\mathbf{v}} \in \mathbb{S}^{d-1}} |G(\hat{\mathbf{v}})|$ and G is defined by Equation (2).

Moreover, to ensure that a particular global extrema $\hat{\mathbf{v}}^*$ of $|G|$ ($|G(\hat{\mathbf{v}}^*)| = \lambda$) is captured ($\exists j \in P : \hat{\mathbf{v}}_j^* = \hat{\mathbf{v}}^*$) with probability at least $1 - \delta$ over isotropic initialization of $\boldsymbol{\theta}^0$, we need $m = O(-\ln(\delta))$ neurons, where the constants in the big- O notation depend only on the properties of data.

Remark 4.2. We expect that the result can be extended to the case of even d . However, this extension will require exploiting the concept of o-minimal structures and proving a corresponding Lojasiewicz inequality for the o-minimal structure containing arcsin function (see our proof, Ji & Telgarsky (2020), and Example 1.5 of Loi (2010)). Since such an analysis is not necessarily informative from a machine learning perspective, we stick to d being odd for simplicity.

Discussion Theorem 4.1 suggests that, when we start training from small initialization ($\sigma \rightarrow 0$), the neurons either align with the global extrema of G ($j \in P$) or grow very slowly ($j \in R$). In particular, for the neurons in R , $|u_j(T_1)| = \|v_j(T_1)\| = O(r^{1+\kappa_j})$. Therefore, their contribution to the decision boundary is negligible compared to the *prominent* neurons in P , for which $|u_j(T_1)| = \|v_j(T_1)\| = \Theta(r|u_j^*|) = \Theta(r)$. Thus, at the start of the training, the network exhibits simplicity bias: regardless of the number of neurons, m , only prominent ones contribute to the network’s decision boundary. Moreover, these prominent neurons are aligned with the global extrema directions of G , which do not depend on m . When the number of neurons is sufficiently large ($m = \Omega(-\ln(\delta))$), the network will learn all global extrema directions of G , which makes the characterization very precise for small σ .

4.2. Feature Growth to a Constant Scale

Next, we extend our analysis beyond the stage studied in Theorem 4.1 to the point when the network weights reach a constant scale. Specifically, we show that the prominent

features preserve their alignment and that the network essentially behaves like a smaller p -neuron network, where p is the number of extrema of G .

Embedding function To formalize our claim, we consider a specific smaller network that describes Equation (1) well. Since the prominent neurons ($j \in P$) cluster around the global extrema of G at the end of the first phase, we can divide them according to their direction

$$P = \sqcup_{k=1}^p P_k \text{ s.t. } \forall k \forall i, j \in P_k \hat{\mathbf{v}}_i^* = \hat{\mathbf{v}}_j^* \quad \wedge \forall k \neq k' \forall i \in P_k, j \in P_{k'} \hat{\mathbf{v}}_i^* \neq \hat{\mathbf{v}}_j^*.$$

We denote by $\hat{\mathbf{v}}_{P_k}^*$ the direction of neurons in P_k and by $s_{P_k} := \text{sign}(G(\hat{\mathbf{v}}_{P_k}^*))$. Now, consider the following auxiliary system

$$\begin{aligned} \frac{du_k^e}{dt} &= \frac{1}{n} \sum_{i=1}^n (-\ell'(f(\boldsymbol{\theta}^e, \mathbf{x}_i)y_i)) \phi(\mathbf{v}_k^e, \mathbf{x}_i)y_i, \\ \frac{d\mathbf{v}_k^e}{dt} &= \frac{u_k^e}{n} \sum_{i=1}^n (-\ell'(f(\boldsymbol{\theta}^e, \mathbf{x}_i)y_i)) \nabla_{\mathbf{v}} \phi(\mathbf{v}_k^e, \mathbf{x}_i)y_i, \end{aligned} \quad (3)$$

where

$$u_k^e(T_1) = s_{P_k} r \sqrt{\sum_{j \in P_k} (u_j^*)^2}, \quad \mathbf{v}_k^e(T_1) = r |u_k^e(T_1)| \hat{\mathbf{v}}_{P_k}^*.$$

These equations describe the dynamics of a p -neuron network $\boldsymbol{\theta}^e = (\mathbf{v}_1^e, \mathbf{v}_2^e, \dots, \mathbf{v}_p^e, u_1^e, u_2^e, \dots, u_p^e)$, initialized in a way that preserves the alignment and scale of $\boldsymbol{\theta}$.

Our goal will be to show that each neuron in the original network can be approximated by corresponding neuron in the p -neuron network. To do it, below, we define an embedding function (Lyu et al., 2021) that maps the p -neuron network $\boldsymbol{\theta}^e$ to a m -neuron network $\boldsymbol{\theta}^x := \chi(\boldsymbol{\theta}^e)$

$$\begin{aligned} \forall j \in P_k \quad & u_j^x = \frac{ru_j^*}{u_k^e(T_1)} u_k^e, \quad \mathbf{v}_j^x = \frac{ru_j^*}{u_k^e(T_1)} \mathbf{v}_k^e, \\ \forall j \in R \quad & u_j^x = 0, \quad \mathbf{v}_j^x = \mathbf{0}. \end{aligned} \quad (4)$$

Propagation of simplicity bias We prove that the original network behaves approximately as an image of the embedding above. We only require a mild assumption on the data that avoids degenerate cases in which a data point perfectly aligns with an extrema of G .

Definition 4.3. A direction $\hat{\mathbf{v}}$ is Δ -regular if $\forall i |\hat{\mathbf{v}}^\top \mathbf{x}_i| \geq \Delta$.

The following result holds under the assumption that all global extrema of G are regular. We note that this is a generalization of Assumption 4.5 of Lyu et al. (2021) for the case of non-linearly-separable data.

Theorem 4.4 (Proof in Appendix C). *In the setting of Theorem 4.1, consider Equation (3) and assume that $\forall j \in P \hat{\mathbf{v}}_j^*$*

is $(\xi + 2\Delta)$ -regular. Then, $\forall \varepsilon \leq \min\{\Delta, 1/2\}$, the following holds. First, $\exists \theta^{\varepsilon,*} := \lim_{r \rightarrow 0} \theta^\varepsilon(T_2^\varepsilon)$, where $T_2^\varepsilon := T_1 + t_4^\varepsilon$, $t_4^\varepsilon := \frac{1}{2\lambda} \ln\left(\frac{\lambda \varepsilon}{2ar^2 \|\theta^*\|_{[m]}^2}\right)$, and $a := \frac{m(1+\xi)^2}{4}$. Second, denote $\theta^{\chi,\varepsilon,*} := \chi(\theta^{\varepsilon,*})$, then, as $r \rightarrow 0$, we have

$$\begin{aligned} \forall j \in P \quad & |u_j(T_2^\varepsilon) - u_j^{\chi,\varepsilon,*}| = O(r^{\kappa_j} + r^2), \\ & \|\hat{v}_j(T_2^\varepsilon) - \hat{v}_j^{\chi,\varepsilon,*}\| = O(r^{\kappa_j} + r^2), \\ \forall j \in R \quad & |u_j(T_2^\varepsilon)| = \|\mathbf{v}_j(T_2^\varepsilon)\| = O(r^{\kappa_j}), \end{aligned}$$

where $\kappa := \min\{\kappa^*, \min_{j \in R} \kappa_j\}$.

Moreover, $|u_j^{\chi,\varepsilon,*}| = \Theta(\sqrt{\varepsilon})$ and $\forall j \in P \|\hat{v}_j^{\chi,\varepsilon,*} - \hat{v}_j^*\| \leq \varepsilon$.

Remark 4.5. The assumption that the critical directions \hat{v}^* are regular is essentially needed only to show that the Hessian matrix $\nabla^2|G|(z)$ is negative semi-definite in some $(\xi + 2\Delta)$ -neighborhood of \hat{v}^* . Since the Hessian matrix is semi-negative at $z = \hat{v}^*$ and G is twice continuously differentiable, this assumption appears rather mild.

Discussion Similarly to Theorem 4.1, the neurons in R have a negligible effect on the network since u_j and v_j are of smaller magnitude compared to the remaining weights. Theorem 4.4 suggests that the network experiences the simplicity bias not only at the start of the training but also until the weights grow to a constant scale $\Theta(\sqrt{\varepsilon})$. Even in the second phase, the prominent neurons in P stay near the extrema of G , which they learned initially. In addition, $\theta(T_2^\varepsilon) \approx \theta^{\chi,\varepsilon,*}$ implies that the original network is an approximate embedding of the p -neural network above.

Theorems 4.1 and 4.4 show an interesting separation of the training dynamics of two-layer networks trained from small initialization. First, the hidden layer features traverse the unit sphere until some become prominent, capturing a supremum direction of G and aligning with it. Then, the prominent features grow without much change in direction.

5. Extreme Simplicity Bias for Specific Data

Our results so far describe the features learned by two-layer networks in the early stages of training, as the parameters go from small to constant scale. However, as indicated by Lyu et al. (2021); Shah et al. (2020), the simplicity bias might persist not only in the initial stages of training but also when the network reaches perfect accuracy on a training dataset. To test to what extent our mechanism outlined in Section 4 can explain this empirical observation, we extend our analysis to the infinite time training limit.

A key challenge in this setup is the lack of convergence guarantees for non-convex models, which necessitates at least some assumptions on the train data. Given our focus on non-linearly-separable data, we focus on datasets in \mathbb{R}^d that feature an XOR-like pattern in a 2-dimensional subspace as a prime example that breaks the linear separability.

5.1. Data with an XOR-pattern

We consider train data in \mathbb{R}^d that follows an XOR pattern in the 2-dimensional subspace generated by coordinate vectors e_1 and e_2 . Specifically, the points cluster around four vectors, $e_1, -e_1, e_2, -e_2$. They are symmetric w.r.t. the permutation of the first and second coordinate, the reflection of the first and second coordinate axes, and the reflection through the hyperplane generated by the first and second axes. Points that cluster around the directions e_1 and $-e_1$ have positive labels; others have negative labels. Finally, we make an additional assumption on the non-alignment of data points with coordinate axes (which would allow us to apply Theorem 4.4 to this dataset). Notice that similar assumptions appear in Lyu et al. (2021), but in our case the resulting dataset is not linearly separable.

To formalize these assumptions, denote $\mathbf{P} := \mathbf{I} + (e_2 - e_1)e_1^\top + (e_1 - e_2)e_2^\top$ (permutation of the first and second coordinates), $\forall a \in \{1, 2\} : \mathbf{R}_a := \mathbf{I} - 2e_a e_a^\top$ (reflection of a th coordinate), $\mathbf{R}_r := 2(e_1 e_1^\top + e_2 e_2^\top) - \mathbf{I}$ (reflection of the rest of the coordinates), and $D_x := \{\mathbf{x}_i\}_{i=1}^n$ for $\mathbf{x}_i \in \mathbb{R}^d$. Then, the formal assumptions will be the following.

Assumption 5.1. Denote $a(k) := k \bmod 2$, $b(k) := 2 \lfloor \frac{k-1}{2} \rfloor - 1$. There exists $\{S_1, S_2, S_3, S_4\} : \sqcup_k S_k = [n]$ such that the following properties hold.

1. $\exists \delta < \sqrt{2}/2 : \forall i \in S_k \|\mathbf{x}_i - b(k)e_{a(k)+1}\| \leq \delta$
 $\wedge y_i = 1 - 2a(k)$.
2. $\forall a \in \{1, 2, r\} \mathbf{R}_a D_x = D_x$.
3. $\mathbf{P} D_x = D_x$.
4. $\exists \Delta > 0 : \forall i, k |e_k^\top \mathbf{x}_i| \geq \xi + 2\Delta$.

5.2. Initial stages of training

We first use the results of Section 4 to analyze the behavior of a two-layer network trained from small initialization on our XOR-like dataset. To apply Theorem 4.1, we first describe the global extrema of G .

Lemma 5.2 (Proof in Appendix D.1). *If Assumption 5.1 holds and $\xi + \delta < 1/6$, the function $|G|$ have four extrema directions: $e_1, -e_1, e_2, -e_2$.*

Lemma 5.2 and Theorem 4.1 imply that, at the start of the training, the big neurons will converge in the four directions: $\pm e_1$ and $\pm e_2$. To evaluate the probability of the network capturing all extrema of G , we show the following fact.

Lemma 5.3 (Proof in Appendix D.2). *Assume the setting of Lemma 5.2. The probability of successful initialization that will capture all extrema of G is greater than $(1 - h^m)^4 \geq 1 - 4(3/4)^m (1 + O(\delta + \xi)) - O((9/16)^m)$ for $m \rightarrow \infty$ and $\delta + \xi \rightarrow 0$, where $h := 1 - \frac{1}{2} \frac{\text{Vol}(A)}{\text{Vol}(\mathbb{S}^{d-1})}$ and $A = \{\mathbf{x} \in \mathbb{S}^{d-1} \mid e_1^\top \mathbf{x} \geq \delta + \xi\}$.*

The lemmas above and Theorems 4.1 and 4.4 suggest that, with high probability, at the end of Phase 2, $\theta^{\varepsilon,*}$, the four-neuron approximation of the original network, has features aligned with the cluster directions of the data, $\|\hat{\mathbf{v}}_k^{\varepsilon,*} - b(k)e_{a(k)+1}\| \leq \varepsilon$.

5.3. Training Dynamics in the Infinite Time Limit

Now we are interested in the training dynamics on the XOR dataset, beyond the stage described by Theorem 4.4 and as $T \rightarrow \infty$. To study that, we use the fact that θ^e provides a good approximation of the original network θ and formalize this beyond the time T_2^ε used in Theorem 4.4.

The next result proves that the limit dynamics (1) of the original network converge to the same features as those learned by the simpler 4-neuron network from Section 5.2. Our result assumes that the features of the 4-neuron network remain aligned with the extrema of G from Lemma 5.2.

Assumption 5.4. In the setting of Lemma 5.2, the 4-neuron network initialized at $\theta^{\varepsilon,*}$ converges in direction to θ^{mm} , in which $\forall i, j$ $\|\mathbf{v}_i^{mm}\| = |u_i^{mm}| = |u_j^{mm}|$ and $\forall k$ $\hat{\mathbf{v}}_k^{mm} = b(k)e_{a(k)+1}$.

We provide experimental and theoretical evidence that this assumption holds in the next subsection. We also note that this assumption concerns the two-layer 4-neuron network and can be checked directly in an experimental manner. In contrast, the result below holds for any sufficiently big two-layer neural networks trained from small initialization.

Lemma 5.5 (Proof in Appendix D.5). *In the setting of Lemma 5.2 under Assumption 5.4, the original network θ converges in direction to $\chi(\theta^{mm})$ if the initialization scale σ is small:*

$$\exists \sigma^* : \forall \sigma < \sigma^* \lim_{t \rightarrow \infty} \left\| \frac{\theta(t)}{\|\theta(t)\|} - \frac{\chi(\theta^{mm})}{\|\chi(\theta^{mm})\|} \right\| = 0.$$

(Notice $\theta(0) = \sigma\theta^0$ and χ depends on θ^0 but not on σ .)

This lemma shows that the conclusions of Theorems 4.1 and 4.4 not only propagate to the later stages of training for the XOR data, but also exacerbate. The network forgets all features except those learned at the beginning of training, causing an extreme simplicity bias (Shah et al., 2020).

Proof sketch Our proof builds upon concepts studied in Lyu & Li 2020. Following their notation, we define the normalized margin $\gamma(\theta)$ of f on dataset D by

$$\gamma := \min_i f(\theta, \mathbf{x}_i)y_i / \|\theta\|^2.$$

Notice that, due to 2-homogeneity of f , γ depends only on the direction of θ : $\forall \lambda > 0$ $\gamma(\lambda\theta) = \gamma(\theta)$. We call a direction that is a (local) solution to the max-margin problem $\max_{\theta} \gamma(\theta)$ a (local)-max-margin direction. The main result we build upon in our proof is the following.

Theorem 5.6 (Theorem 5.6, Lyu et al. 2021). *Consider 2-positively-homogeneous network f trained with gradient flow on logistic loss. For any local-max-margin direction, $\hat{\theta}^*$, and $\zeta > 0$, $\exists \omega > 0, \rho \geq 1$ such that for any θ^0 with $\|\theta^0\| \geq \rho$ and $\|\hat{\theta}^0 - \hat{\theta}^*\| \leq \omega$, gradient flow starting with θ^0 directionally converges to some direction $\hat{\theta}$ with the same normalized margin γ as $\hat{\theta}^*$, and $\|\hat{\theta} - \hat{\theta}^*\| \leq \zeta$.*

We use this result in the following manner. First, we show that the assumed limit direction of θ^e is a local-max-margin direction. Second, we employ Theorem 4.4 and classical theorems about continuous dependency on initial conditions for initial value problems to show that our original system will satisfy the conditions of Theorem 5.6. Finally, we apply Theorem 5.6 to prove the desired result.

5.4. Convergence Behavior of Four-Neuron Network

Finally, we provide evidence for Assumption 5.4. Unfortunately, a precise characterization of this smaller network on our dataset is challenging. First, as far as we are aware, general results about convergence to perfect accuracy or zero loss of two-layer networks exist only for 1-neuron networks (Awasthi et al., 2023; Chistikov et al., 2023). Second, even if one can prove that the network converges to zero loss, the additional challenge to analyze the resulting limit direction remains. The only characterization of the limit directions of networks that we are aware of is in terms of the KKT conditions for the dual margin-maximization problem (Lyu & Li, 2020). However, in the general case a direct analysis of the KKT conditions remains intractable.

Thus, we give a theoretical motivation for the proposed direction θ^{mm} and validate our assumption empirically.

5.4.1. THEORETICAL EVIDENCE

Convergence to perfect accuracy First, we prove that the network converges to perfect accuracy.

Proposition 5.7 (Proof in Appendix D.3). *Assume that $\max_{i,j} |u_i^\varepsilon(0)|/|u_j^\varepsilon(0)| \leq 1000$, $\max_k |u_k^\varepsilon(0)| \in (0.001, 0.5)$, and $\delta + \varepsilon \leq 0.01$. Then there exist a time T^0 such that $\forall i$ $f(\theta^e(T^0), \mathbf{x}_i)y_i > 4.67 > 0$.*

To comment on the plausibility of the assumptions in this result, first notice that in the setting of Theorem 4.1 neurons evolve almost independently. Thus, due to the symmetry of our dataset, at the end of Phase 1 the scales u_j^* will be sums of m identically distributed values. Then by the law of large numbers they will converge to the same values in the limit $m \rightarrow \infty$. Thus, when m is big, at the end of Phase 1, the ratios $|u_i^\varepsilon(0)|/|u_j^\varepsilon(0)|$ will be close to one. Since in Phase 2 the scales of neurons and the deviation of features is small, we expect these ratios to increase by no more than $1 + O(\varepsilon)$, implying that these ratios again will be close to one.

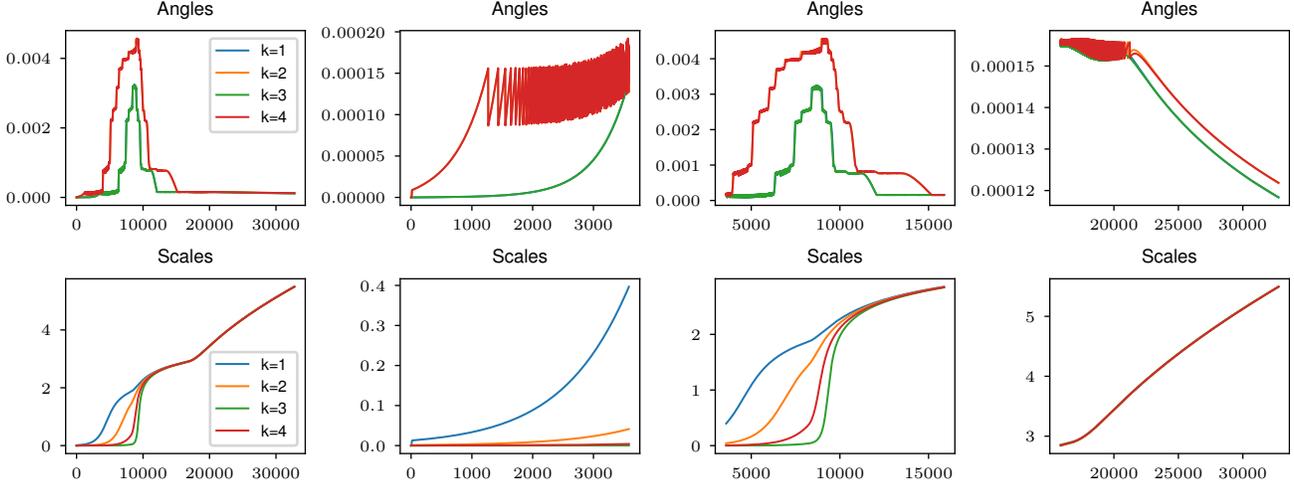


Figure 1. Evolution of 4-neuron network initialized at $(u_1^e(0), u_2^e(0), u_3^e(0), u_4^e(0)) = (10^{-4}, -10^{-5}, 10^{-7}, -10^{-6})$. The first column depicts the whole training process. We additionally depict different stages of training process for visual convenience: the second column depicts the first 3584 training epochs; the third column depicts the epochs from 3584 to 15872; the last column depicts training after the 15872th epoch. Notice that $\alpha_1 \approx \alpha_3$ and $\alpha_2 \approx \alpha_4$, where α_i are the angles between the network features and the cluster directions.

The assumption that $\max_k |u_k^e(0)| \in (0.001, 0.5)$ is implied by the property $u_j^{e,*} = \Theta(\sqrt{\varepsilon})$. If we assume that $\Delta \approx 0.001$ and choose $\varepsilon = \Delta$, we can expect that θ^e will have the desired scale. Finally, we assume $\delta + \xi \leq 0.01$ for technical reasons related to the proof technique. We expect that the property will hold for wider ranges of δ .

Implicit bias at the end of training The previous result about perfect accuracy also suggests that the network may achieve small loss at the end of training. Then, by Theorem 4.4 of Lyu & Li (2020) and Theorem 3.1 of Ji & Telgarsky (2020), the network converges to some KKT point of dual margin-maximization problem

$$\min_{\theta} \|\theta\|^2 \text{ s.t. } \forall i f(\theta, \mathbf{x}_i) y_i \geq 1.$$

Under the additional assumption that the features of this margin direction are $(\xi + 2\Delta)$ -regular for some $\Delta > 0$, a slight modification of the results of Vardi et al. (2022) implies that this direction is a local-max-margin direction. Finally, we show that θ^{mm} is one strict local-max-margin direction, which motivates our hypothesis.

Proposition 5.8 (Proof in Appendix D.4). *Assume the setting of Lemma 5.2. Then, the direction of θ^{mm} is a strict local-max-margin in the weight space of the 4-neuron network, while its embedding, $\chi(\theta^{mm})$ is a strict local-max-margin in the original weight space.*

5.4.2. EMPIRICAL EVIDENCE

Here we present experimental evidence for Assumption 5.4. To this end, we construct a random dataset in \mathbb{R}^2 that will satisfy Assumption 5.1 with $\delta \leq 0.01$ and train a

four neuron ReLU network using gradient descent for different initializations, in which the neurons are aligned with cluster directions, corresponding to the setting of Lemma 5.2. Specifically, after we pick initialization scales $(u_1^e(0), u_2^e(0), u_3^e(0), u_4^e(0))$, we initialize the first layer as $\mathbf{v}_k^e(0) = u_k^e(0) e_{a(k)+1}$ and train our network using plain gradient descent.

We used the following adaptive learning rate schedule η_t . During the first part of training, which corresponds to the evolution initialized at the limit point from Theorem 4.1 to the limit point in Theorem 4.4, the gradients are very stable. Therefore, we use constant-scale learning rates $\eta_t = 4$ for $t < 12$ and $\eta_t = 2^{-7}$ for $12 \leq t < 2^{13}$. The next part corresponds to the training initialized at the limit point from Theorem 4.4. At around $t = 2^{13}$ all features reach constant scales and the cross-entropy loss starts to dump gradients. This allows us to progressively increase the learning rate, in order to speed up the simulation, in the third part of schedule, setting $\eta_t = 2^{-7} (1 + 2^5 (\frac{t}{2^{13}} - 1))$ for $2^{13} \leq t < 2^{14}$. Finally, at the end of the third part, the cross-entropy loss causes gradients to decay exponentially and the training process almost stops. To combat this and simulate the later stages of training, we use an exponential learning rate $\eta_t = 2^{-7} (1 + 2^{5 + \frac{t-2^{14}}{2^9}})$ for $t \geq 2^{14}$.

Figure 1 depicts the evolution of the 4-neuron network. Here, plots titled “Angles” depict the signed angles between the network features and the cluster directions. And plots titled “Scales” depict $\|\mathbf{v}_k^e\|$. (See more experiments in Appendix E.) As we can see, eventually the scales of the neurons became almost identical and the network start to converge to the desired local-max-margin direction θ^{mm} (forth

column), empirically supporting our assumption.

6. Empirical Validation of Results

In this section, we empirically validate the predictions of Theorems 4.1, 4.4, and Lemma 5.5. We consider a skewed XOR-like dataset in \mathbb{R}^2 , similar to the dataset covered by Assumption 5.1, but in which the angle between cluster directions can be arbitrary. We consider two angles: $\alpha = \pi/2$ and $\alpha = \pi/3$. The experiments for $\alpha = \pi/2$ seek to verify our predictions from Section 5, while the experiments for $\alpha = \pi/3$ test if these predictions transfer to non-orthogonal cases. We train a two-layer neural network with $m = 2^{12}$ randomly initialized neurons with initialization scale $\sigma = 2^{-7}$. We train this network for 2^{13} epochs using gradient descent with $\text{lr} = 2^{-4}$. At that stage, we observed that the network essentially converged.

When $\alpha = \pi/2$, the scale of the first layer grows from 0.3 to 5.5 during the training (Figure 2, first row, left). At the same time, almost all neurons end up aligned in four directions (Figure 2, second row, left): $0, \pi/2, \pi, -\pi/2$, which are the global extrema of the function G . Similarly, when $\alpha = \pi/3$ the scale of the first layer grows from 0.3 to 5.9 during the training (Figure 2, first row, right). At the same time, almost all neurons end up aligned in four directions (Figure 2, second row, right): $-\pi/6, \pi/2, 5\pi/6, -\pi/2$, which again are the global extrema of the function G . In both cases, the

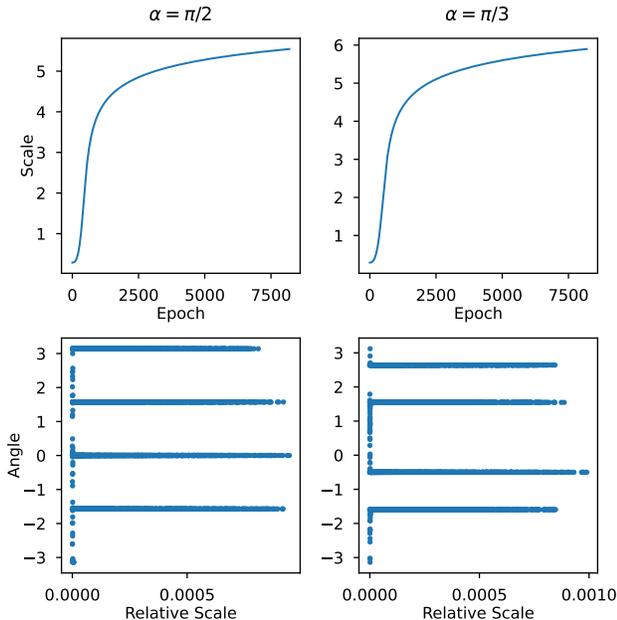


Figure 2. Empirical evidence of simplicity bias on XOR-like data. Relative scale in the second row defined as $\|v_i\|^2 / \sum_{j=1}^m \|v_j\|^2$.

few non-aligned neurons have smaller weights than aligned ones (Figure 2, second row). Therefore, our theory indeed predicts the alignment well in this setting.

7. Effects of Extreme Simplicity Bias

Finally, we test the possible implications of our theoretical characterization of simplicity bias on real-world datasets. Our previous results suggest that simplicity bias disproportionately amplifies “simple” features that are very informative about the target. Additionally, in the case of extreme simplicity bias, if the “simple” features are enough to classify the training set perfectly, the network effectively “forgets” all features except for the “simple” ones.

The latter fact suggests that if spurious features are enough to classify the target in the train distribution, then the network should progressively lose its ability to fine-tune out of distribution. We test this hypothesis on the MNIST-CIFAR-10 domino dataset proposed by Shah et al. (2020).



Figure 3. Examples of domino with a car (class 1 in CIFAR-10) in train (left) and test (right) dataset. Notice that the top MNIST image is an image of 1 only for the train data.

Setup The dataset contains vertical concatenations of MNIST and CIFAR-10 images and uses CIFAR-10 image labels. On the train distribution, the MNIST and CIFAR-10 labels of the images are perfectly correlated: digit i from MNIST will always be concatenated with the image from class i in CIFAR-10. On the test distribution, MNIST and CIFAR-10 labels are not correlated, i.e., the concatenations of images from different classes are random. Thus, MNIST images represent “simple” spurious features that can be used to perfectly classify the data. (See Figure 3.)

We fit a ResNet-18 model (He et al., 2016) on this domino dataset and track its parameter trajectory. We use the standard PyTorch initialization, multiplied by 2^{-5} to mimic the small initialization regime studied in the previous sections. We train this model for 2^8 epochs using the usual SGD optimizer with Nesterov momentum and linear schedule with warm-up (similar training recipe was used by Jain et al., 2023). Periodically during training, we apply the current model to our test set and extract its last layer features, X , on test data. We normalize these last layer features (to make feature scales comparable across different epochs), $X^{\text{norm}} = \frac{X}{\sqrt{\frac{1}{n} \sum_{i=1}^n \sum_{j=1}^k X_{ij}^2}}$, where n is the number of test

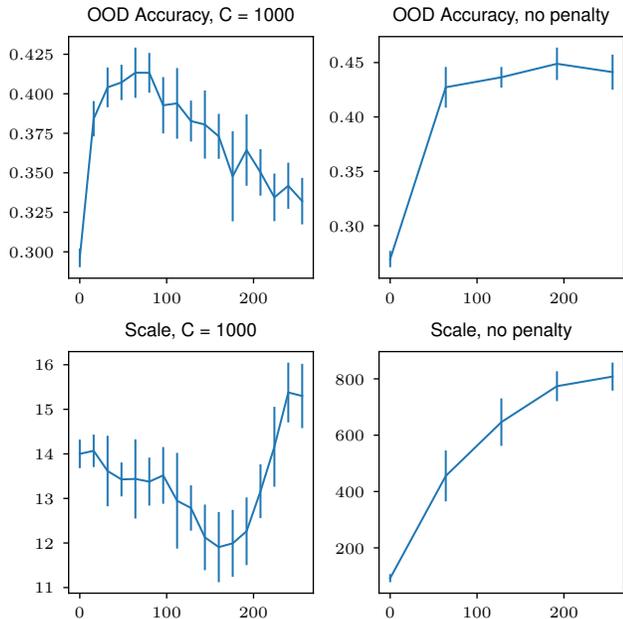


Figure 4. Accuracy and scale of the logistic regression on the validation part of the OOD test set (y -axis) vs. the training epoch at which the ResNet features are extracted (x -axis).

samples, k is the number of features. We then use these normalized last layer features to train simple logistic regression model. We train two types of linear models: with small regularization (with inverse regularization strength equal to 1000) and without regularization. We plot the accuracy on the validation part of the test distribution and the quadratic mean of the regression coefficients on Figure 4. (See training details in Appendix F.1.)

Analysis First, we can see that the model can not make reliable predictions on the test set, peaking in accuracy at about 45%. This is in contrast to training the model using the same recipe on plain CIFAR-10 data, where it achieves around 90% accuracy. At the same time, using the MNIST labels on the test set, we get around 99% accuracy (see Figure 8 in Appendix F.2; in line with the results of Shah et al., 2020; Hermann & Lampinen, 2020). These facts suggest that the network indeed experiences a simplicity bias and mainly relies on “simple” MNIST features for prediction. Second, as we can see, the OOD accuracy increases fast at the beginning of training. This may indicate that even simple MNIST features are better for classification of CIFAR-10 data compared to random features. Finally, we can see that the OOD accuracy does not increase in the latter stages of training, indicating that the simplicity bias persists even if we train longer. In Appendix F.2, we present additional setups, where we do not scale the initialization (Figure 9) or break the perfect correlation between MNIST and CIFAR-10 labels on the train set (Figure 10). In both

cases, we also observe the presence of simplicity bias.

Finally, we can see some evidence in favor of the extreme simplicity bias mechanism described in Section 5. For the small regularization setup, we observe a significant drop in OOD accuracy when using features from later training stages (the drop between epochs 64 and 256 is approximately $8.13\% \pm 0.80\%$, which gives a p-value around 10^{-6} according to the t-test). This result suggests a potential presence of extreme simplicity bias, which impedes the learning of complex features and makes the network forget the initial random features. For the no regularization setup, we can see that the network does not lose OOD accuracy, but the last-layer weights become much bigger. This result is again consistent with the presence of extreme simplicity bias, as our proposed mechanism does not force the network to forget “non-simple” features directly. Instead, it makes “simple” features grow faster, so the regression can only use “non-simple” features by applying huge weights (approximately $50\times$ bigger than the regularized regression).

Discussion We can draw three practical conclusions from these experiments. Suppose the training data can be classified using a simple heuristic. Then, it might be beneficial to train networks in a more “lazy regime” (Chizat et al., 2019), which forces the network to remember randomly initialized features. Similarly, one could apply early stopping of training. In this way, one could benefit from new features learned from the training data without losing access to potentially beneficial random features. Finally, since simplicity bias works by making “simple” features disproportionately large, the natural countermeasure is to use additional normalization layers, which have already been proven effective for few-shot transfer (e.g., Perez et al., 2018). At the same time, the “lazy regime” or early stopping might harm in-distribution generalization performance (Telgarsky, 2023; Lyu et al., 2024). Thus, a practical method would need to trade off the benefits of the “lazy regime” or early stopping for the OOD task to their potential harm to generalization.

8. Conclusion

We characterize simplicity bias beyond linearly separable datasets as the tendency of features to cluster in several directions, which do not depend on the network size. We also demonstrate that extreme simplicity bias may appear for non-linearly-separable datasets and observe it experimentally on image data. We see our results as an indication that the characterization of simplicity bias is a crucial step toward improving out-of-distribution generalization.

Acknowledgments

This research was partially funded from the Ministry of Education and Science of Bulgaria (support for INSAIT, part of the Bulgarian National Roadmap for Research Infrastructure). The authors thank Ivan Kirev and Kristian Minchev for their helpful feedback and discussions on this work.

Impact Statement

This paper presents work whose goal is to advance the field of Machine Learning. There are many potential societal consequences of our work, none which we feel must be specifically highlighted here.

References

- Arora, S., Du, S., Hu, W., Li, Z., and Wang, R. Fine-Grained Analysis of Optimization and Generalization for Over-parameterized Two-Layer Neural Networks. In Chaudhuri, K. and Salakhutdinov, R. (eds.), *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, pp. 322–332. PMLR, 09–15 Jun 2019.
- Arpit, D., Jastrzebski, S., Ballas, N., Krueger, D., Bengio, E., Kanwal, M. S., Maharaj, T., Fischer, A., Courville, A., Bengio, Y., and Lacoste-Julien, S. A Closer Look at Memorization in Deep Networks. In Precup, D. and Teh, Y. W. (eds.), *Proceedings of the 34th International Conference on Machine Learning*, volume 70 of *Proceedings of Machine Learning Research*, pp. 233–242. PMLR, 06–11 Aug 2017.
- Awasthi, P., Tang, A., and Vijayaraghavan, A. Agnostic Learning of General ReLU Activation Using Gradient Descent. In *The Eleventh International Conference on Learning Representations*, 2023.
- Ba, J., Erdogdu, M. A., Suzuki, T., Wang, Z., Wu, D., and Yang, G. High-dimensional Asymptotics of Feature Learning: How One Gradient Step Improves the Representation. In Koyejo, S., Mohamed, S., Agarwal, A., Belgrave, D., Cho, K., and Oh, A. (eds.), *Advances in Neural Information Processing Systems*, volume 35, pp. 37932–37946. Curran Associates, Inc., 2022.
- Basri, R., Galun, M., Geifman, A., Jacobs, D., Kasten, Y., and Kritchman, S. Frequency Bias in Neural Networks for Input of Non-Uniform Density. In III, H. D. and Singh, A. (eds.), *Proceedings of the 37th International Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning Research*, pp. 685–694. PMLR, 13–18 Jul 2020.
- Boursier, E. and Flammarion, N. Early alignment in two-layer networks training is a two-edged sword, 2024.
- Boursier, E., Pillaud-Vivien, L., and Flammarion, N. Gradient flow dynamics of shallow ReLU networks for square loss and orthogonal inputs. In Koyejo, S., Mohamed, S., Agarwal, A., Belgrave, D., Cho, K., and Oh, A. (eds.), *Advances in Neural Information Processing Systems*, volume 35, pp. 20105–20118. Curran Associates, Inc., 2022.
- Bowman, B. and Montufar, G. Implicit Bias of MSE Gradient Optimization in Underparameterized Neural Networks. In *International Conference on Learning Representations*, 2022.
- Brutzkus, A. and Globerson, A. Why do Larger Models Generalize Better? A Theoretical Perspective via the XOR Problem. In Chaudhuri, K. and Salakhutdinov, R. (eds.), *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, pp. 822–830. PMLR, 09–15 Jun 2019.
- Brutzkus, A., Globerson, A., Malach, E., and Shalev-Shwartz, S. SGD Learns Over-parameterized Networks that Provably Generalize on Linearly Separable Data. In *International Conference on Learning Representations*, 2018.
- Chistikov, D., Englert, M., and Lazic, R. Learning a Neuron by a Shallow ReLU Network: Dynamics and Implicit Bias for Correlated Inputs. In Oh, A., Neumann, T., Globerson, A., Saenko, K., Hardt, M., and Levine, S. (eds.), *Advances in Neural Information Processing Systems*, volume 36, pp. 23748–23760. Curran Associates, Inc., 2023.
- Chizat, L. and Bach, F. On the Global Convergence of Gradient Descent for Over-parameterized Models using Optimal Transport. In Bengio, S., Wallach, H., Larochelle, H., Grauman, K., Cesa-Bianchi, N., and Garnett, R. (eds.), *Advances in Neural Information Processing Systems*, volume 31. Curran Associates, Inc., 2018.
- Chizat, L., Oyallon, E., and Bach, F. On Lazy Training in Differentiable Programming. In Wallach, H., Larochelle, H., Beygelzimer, A., d'Alché-Buc, F., Fox, E., and Garnett, R. (eds.), *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019.
- Cybenko, G. Approximation by superpositions of a sigmoidal function. *Mathematics of control, signals and systems*, 2(4):303–314, 1989.
- Du, S. S., Hu, W., and Lee, J. D. Algorithmic Regularization in Learning Deep Homogeneous Models: Layers are Automatically Balanced. In Bengio, S., Wallach, H., Larochelle, H., Grauman, K., Cesa-Bianchi, N., and Garnett, R. (eds.), *Advances in Neural Information Processing Systems*, volume 31. Curran Associates, Inc., 2018.

- Englert, M. and Lazic, R. Adversarial Reprogramming Revisited. In Koyejo, S., Mohamed, S., Agarwal, A., Belgrave, D., Cho, K., and Oh, A. (eds.), *Advances in Neural Information Processing Systems*, volume 35, pp. 28588–28600. Curran Associates, Inc., 2022.
- Frei, S., Chatterji, N. S., and Bartlett, P. L. Random Feature Amplification: Feature Learning and Generalization in Neural Networks. *Journal of Machine Learning Research*, 24(303):1–49, 2023a.
- Frei, S., Vardi, G., Bartlett, P., Srebro, N., and Hu, W. Implicit Bias in Leaky ReLU Networks Trained on High-Dimensional Data. In *The Eleventh International Conference on Learning Representations*, 2023b.
- Ge, R., Ren, Y., Wang, X., and Zhou, M. Understanding Deflation Process in Over-parametrized Tensor Decomposition. In Ranzato, M., Beygelzimer, A., Dauphin, Y., Liang, P., and Vaughan, J. W. (eds.), *Advances in Neural Information Processing Systems*, volume 34, pp. 1299–1311. Curran Associates, Inc., 2021.
- Geirhos, R., Rubisch, P., Michaelis, C., Bethge, M., Wichmann, F. A., and Brendel, W. ImageNet-trained CNNs are biased towards texture; increasing shape bias improves accuracy and robustness. In *International Conference on Learning Representations*, 2019.
- Geirhos, R., Jacobsen, J.-H., Michaelis, C., Zemel, R., Brendel, W., Bethge, M., and Wichmann, F. A. Shortcut learning in deep neural networks. *Nature Machine Intelligence*, 2(11):665–673, 2020.
- Hartman, P. *Ordinary Differential Equations*. Society for Industrial and Applied Mathematics, second edition, 2002. doi: 10.1137/1.9780898719222.
- He, K., Zhang, X., Ren, S., and Sun, J. Deep Residual Learning for Image Recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2016.
- Hermann, K. and Lampinen, A. What shapes feature representations? exploring datasets, architectures, and training. In Larochelle, H., Ranzato, M., Hadsell, R., Balcan, M., and Lin, H. (eds.), *Advances in Neural Information Processing Systems*, volume 33, pp. 9995–10006. Curran Associates, Inc., 2020.
- Hornik, K., Stinchcombe, M., and White, H. Multilayer feedforward networks are universal approximators. *Neural Networks*, 2(5):359–366, 1989. ISSN 0893-6080. doi: 10.1016/0893-6080(89)90020-8.
- Jain, S., Lawrence, H., Moitra, A., and Madry, A. Distilling Model Failures as Directions in Latent Space. In *The Eleventh International Conference on Learning Representations*, 2023.
- Ji, Z. and Telgarsky, M. Directional convergence and alignment in deep learning. In Larochelle, H., Ranzato, M., Hadsell, R., Balcan, M., and Lin, H. (eds.), *Advances in Neural Information Processing Systems*, volume 33, pp. 17176–17186. Curran Associates, Inc., 2020.
- Kalimeris, D., Kaplun, G., Nakkiran, P., Edelman, B., Yang, T., Barak, B., and Zhang, H. SGD on Neural Networks Learns Functions of Increasing Complexity. In Wallach, H., Larochelle, H., Beygelzimer, A., d Alche-Buc, F., Fox, E., and Garnett, R. (eds.), *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019.
- Kou, Y., Chen, Z., and Gu, Q. Implicit Bias of Gradient Descent for Two-layer ReLU and Leaky ReLU Networks on Nearly-orthogonal Data. In Oh, A., Neumann, T., Globerson, A., Saenko, K., Hardt, M., and Levine, S. (eds.), *Advances in Neural Information Processing Systems*, volume 36, pp. 30167–30221. Curran Associates, Inc., 2023.
- Li, Y., Ma, T., and Zhang, H. R. Learning Over-Parametrized Two-Layer Neural Networks beyond NTK. In Abernethy, J. and Agarwal, S. (eds.), *Proceedings of Thirty Third Conference on Learning Theory*, volume 125 of *Proceedings of Machine Learning Research*, pp. 2613–2682. PMLR, 09–12 Jul 2020.
- Loi, T. L. Lecture 1: O-minimal structures. In *The Japanese-Australian Workshop on Real and Complex Singularities: JARCS III*, volume 43, pp. 19–31. Australian National University, Mathematical Sciences Institute, 2010.
- Luo, T., Ma, Z., Xu, Z.-Q. J., and Zhang, Y. Theory of the Frequency Principle for General Deep Neural Networks. *CSIAM Transactions on Applied Mathematics*, 2(3):484–507, 2021a. ISSN 2708-0579. doi: 10.4208/csiam-am.SO-2020-0005.
- Luo, T., Xu, Z.-Q. J., Ma, Z., and Zhang, Y. Phase Diagram for Two-layer ReLU Neural Networks at Infinite-width Limit. *Journal of Machine Learning Research*, 22(71): 1–47, 2021b.
- Lyu, K. and Li, J. Gradient Descent Maximizes the Margin of Homogeneous Neural Networks. In *International Conference on Learning Representations*, 2020.
- Lyu, K., Li, Z., Wang, R., and Arora, S. Gradient Descent on Two-layer Nets: Margin Maximization and Simplicity Bias. In Ranzato, M., Beygelzimer, A., Dauphin, Y., Liang, P., and Vaughan, J. W. (eds.), *Advances in Neural Information Processing Systems*, volume 34, pp. 12978–12991. Curran Associates, Inc., 2021.

- Lyu, K., Jin, J., Li, Z., Du, S. S., Lee, J. D., and Hu, W. Dichotomy of early and late phase implicit biases can provably induce grokking. In *The Twelfth International Conference on Learning Representations*, 2024.
- Maennel, H., Bousquet, O., and Gelly, S. Gradient Descent Quantizes ReLU Network Features, 2018.
- McCoy, T., Pavlick, E., and Linzen, T. Right for the Wrong Reasons: Diagnosing Syntactic Heuristics in Natural Language Inference. In Korhonen, A., Traum, D., and Márquez, L. (eds.), *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pp. 3428–3448, Florence, Italy, July 2019. Association for Computational Linguistics. doi: 10.18653/v1/P19-1334.
- Mei, S., Montanari, A., and Nguyen, P.-M. A mean field view of the landscape of two-layer neural networks. *Proceedings of the National Academy of Sciences*, 115(33): E7665–E7671, 2018. doi: 10.1073/pnas.1806579115.
- Min, H., Mallada, E., and Vidal, R. Early Neuron Alignment in Two-layer ReLU Networks with Small Initialization. In *The Twelfth International Conference on Learning Representations*, 2024.
- Morwani, D., Batra, J., Jain, P., and Netrapalli, P. Simplicity Bias in 1-Hidden Layer Neural Networks. In Oh, A., Neumann, T., Globerson, A., Saenko, K., Hardt, M., and Levine, S. (eds.), *Advances in Neural Information Processing Systems*, volume 36, pp. 8048–8075. Curran Associates, Inc., 2023.
- Pellegrini, F. and Biroli, G. An analytic theory of shallow networks dynamics for hinge loss classification. In Larochelle, H., Ranzato, M., Hadsell, R., Balcan, M., and Lin, H. (eds.), *Advances in Neural Information Processing Systems*, volume 33, pp. 5356–5367. Curran Associates, Inc., 2020.
- Perez, E., Strub, F., de Vries, H., Dumoulin, V., and Courville, A. FiLM: Visual Reasoning with a General Conditioning Layer. *Proceedings of the AAAI Conference on Artificial Intelligence*, 32(1), Apr. 2018. doi: 10.1609/aaai.v32i1.11671.
- Phuong, M. and Lampert, C. H. The inductive bias of ReLU networks on orthogonally separable data. In *International Conference on Learning Representations*, 2021.
- Rahaman, N., Baratin, A., Arpit, D., Draxler, F., Lin, M., Hamprecht, F., Bengio, Y., and Courville, A. On the Spectral Bias of Neural Networks. In Chaudhuri, K. and Salakhutdinov, R. (eds.), *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, pp. 5301–5310. PMLR, 09–15 Jun 2019.
- Refinetti, M., Ingrosso, A., and Goldt, S. Neural networks trained with SGD learn distributions of increasing complexity. In Krause, A., Brunskill, E., Cho, K., Engelhardt, B., Sabato, S., and Scarlett, J. (eds.), *Proceedings of the 40th International Conference on Machine Learning*, volume 202 of *Proceedings of Machine Learning Research*, pp. 28843–28863. PMLR, 23–29 Jul 2023.
- Safra, I., Vardi, G., and Lee, J. D. On the Effective Number of Linear Regions in Shallow Univariate ReLU Networks: Convergence Guarantees and Implicit Bias. In Koyejo, S., Mohamed, S., Agarwal, A., Belgrave, D., Cho, K., and Oh, A. (eds.), *Advances in Neural Information Processing Systems*, volume 35, pp. 32667–32679. Curran Associates, Inc., 2022.
- Sarussi, R., Brutzkus, A., and Globerson, A. Towards Understanding Learning in Neural Networks with Linear Teachers. In Meila, M. and Zhang, T. (eds.), *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*, pp. 9313–9322. PMLR, 18–24 Jul 2021.
- Shah, H., Tamuly, K., Raghunathan, A., Jain, P., and Netrapalli, P. The Pitfalls of Simplicity Bias in Neural Networks. In Larochelle, H., Ranzato, M., Hadsell, R., Balcan, M., and Lin, H. (eds.), *Advances in Neural Information Processing Systems*, volume 33, pp. 9573–9585. Curran Associates, Inc., 2020.
- Telgarsky, M. Feature selection and low test error in shallow low-rotation ReLU networks. In *The Eleventh International Conference on Learning Representations*, 2023.
- Valle-Perez, G., Camargo, C. Q., and Louis, A. A. Deep learning generalizes because the parameter-function map is biased towards simple functions. In *International Conference on Learning Representations*, 2019.
- Vardi, G., Shamir, O., and Srebro, N. On Margin Maximization in Linear and ReLU Networks. In Koyejo, S., Mohamed, S., Agarwal, A., Belgrave, D., Cho, K., and Oh, A. (eds.), *Advances in Neural Information Processing Systems*, volume 35, pp. 37024–37036. Curran Associates, Inc., 2022.
- Wang, M. and Ma, C. Understanding Multi-phase Optimization Dynamics and Rich Nonlinear Behaviors of ReLU Networks. In Oh, A., Neumann, T., Globerson, A., Saenko, K., Hardt, M., and Levine, S. (eds.), *Advances in Neural Information Processing Systems*, volume 36, pp. 35654–35747. Curran Associates, Inc., 2023.
- Wei, C., Lee, J. D., Liu, Q., and Ma, T. Regularization Matters: Generalization and Optimization of Neural Nets v.s. their Induced Kernel. In Wallach, H., Larochelle, H., Beygelzimer, A., d Alche-Buc, F., Fox, E., and Garnett,

R. (eds.), *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019.

Xu, W. and Du, S. Over-Parameterization Exponentially Slows Down Gradient Descent for Learning a Single Neuron. In Neu, G. and Rosasco, L. (eds.), *Proceedings of Thirty Sixth Conference on Learning Theory*, volume 195 of *Proceedings of Machine Learning Research*, pp. 1155–1198. PMLR, 12–15 Jul 2023.

Xu, Z.-Q. J., Zhang, Y., and Xiao, Y. Training Behavior of Deep Neural Network in Frequency Domain. In Gedeon, T., Wong, K. W., and Lee, M. (eds.), *Neural Information Processing*, pp. 264–274, Cham, 2019. Springer International Publishing. ISBN 978-3-030-36708-4.

Zhou, H., Qixuan, Z., Luo, T., Zhang, Y., and Xu, Z.-Q. Towards Understanding the Condensation of Neural Networks at Initial Training. In Koyejo, S., Mohamed, S., Agarwal, A., Belgrave, D., Cho, K., and Oh, A. (eds.), *Advances in Neural Information Processing Systems*, volume 35, pp. 2184–2196. Curran Associates, Inc., 2022.

Supplementary Material

The supplementary material is structured as follows.

- Appendix A explains some additional notation used in proofs.
- Appendix B contains the proof of Theorem 4.1.
- Appendix C contains the proof of Theorem 4.4.
- Appendix D contains the proofs of the results from Section 5.
- Appendix E contains the additional experimental results for Section 5.
- Appendix F.1 contains additional experimental results for Section 7.

A. Additional Notation

Denote $\mathbf{g}(\mathbf{v}) := \nabla_{\mathbf{v}} G(\mathbf{v})$.

B. Proof of Theorem 4.1

We will prove the theorem in several stages. First, we will analyze Equation (2). Then, we couple Equations (1) and (2). This will allow us to prove the desired result.

B.1. Analysis of Equation (2)

B.1.1. DIRECTIONAL CONVERGENCE

Now, we want to apply the results of Ji & Telgarsky (2020) about the directional convergence to Equation (2).

Dynamics of neuron direction Notice that

$$\frac{d\hat{\mathbf{v}}}{d\mathbf{v}} = \frac{d\mathbf{v}/\|\mathbf{v}\|}{d\mathbf{v}} = \frac{\|\mathbf{v}\|^2 - \mathbf{v}\mathbf{v}^\top}{\|\mathbf{v}\|^3}. \quad (5)$$

Thus,

$$\frac{d\hat{\mathbf{v}}_j^l}{dt} = \frac{d\hat{\mathbf{v}}_j^l}{d\mathbf{v}_j^l} \frac{d\mathbf{v}_j^l}{dt} = s_j \mathbf{P}_{\hat{\mathbf{v}}_j^l} \mathbf{g}(\hat{\mathbf{v}}_j^l)$$

and $\hat{\mathbf{v}}^l(0) = \frac{\mathbf{v}^0}{\|\mathbf{v}^0\|}$ does not depend on σ .

Consider the auxiliary system

$$\frac{d\mathbf{r}_j}{dt} = s_j \nabla_{\mathbf{r}} G(\hat{\mathbf{r}}_j) = \frac{s_j}{\|\mathbf{r}_j\|} \mathbf{P}_{\mathbf{r}_j} \mathbf{g}(\hat{\mathbf{r}}_j), \quad (6)$$

where $\mathbf{r}_j(0) = \hat{\mathbf{v}}_j^l(0)$. Notice that

$$\mathbf{r}_j \frac{d\mathbf{r}_j}{dt} = 0 \wedge \|\mathbf{r}_j(0)\| = 1 \implies \frac{d\mathbf{r}_j}{dt} = s_j \mathbf{P}_{\hat{\mathbf{r}}_j} \mathbf{g}(\hat{\mathbf{r}}_j).$$

Thus, $\forall t \mathbf{r}_j(t) = \hat{\mathbf{v}}_j^l(t)$: we could use \mathbf{r}_j instead of $\hat{\mathbf{v}}_j^l$ in further derivations.

We want to show the convergence of \mathbf{r}_j similarly to Theorem 3.1 of Ji & Telgarsky (2020). To do it, we want to use Lemma B.11 of Ji & Telgarsky (2020).

Lemma B.1 (Lemma B.11 of Ji & Telgarsky 2020). *Given a locally Lipschitz definable function $f: A \rightarrow \mathbb{R}$ with an open bounded domain A , there exists $\nu > 0$ and a definable desingularizing function ψ on $[0, \nu)$ such that*

$$\forall \mathbf{v} \in f^{-1}((0, \nu)) \psi'(f(\mathbf{v})) \|\nabla_{\mathbf{v}} f(\mathbf{v})\| \geq 1.$$

To apply it, we need to show that $G(\hat{r})$ is locally Lipschitz and definable. And since

$$G(\hat{v}) = \frac{-\ell'(0)}{n} \sum_{i=1}^n \phi(\hat{v}, \mathbf{x}_i) y_i,$$

we only need to show the desired properties for ϕ .

Proposition B.2. Denote $c := \max\left(-1, \min\left(1, \frac{\mathbf{v}^\top \mathbf{x}}{\xi \|\mathbf{v}\|}\right)\right)$ and the class of all polynomials of x_1, \dots, x_p by $\text{Poly}(x_1, \dots, x_p)$. We have

$$\phi(\mathbf{v}, \mathbf{x}) = \int_{-c}^1 (\mathbf{v}^\top \mathbf{x} + \xi \|\mathbf{v}\| a) (1 - a^2)^{\frac{d-1}{2}} \frac{\text{Vol}(\mathbb{D}^{d-1})}{\text{Vol}(\mathbb{D}^d)} da \in \xi \|\mathbf{v}\| \text{Poly}(c, \sqrt{1 - c^2}).$$

Proof. By the definition

$$\begin{aligned} \phi(\mathbf{v}, \mathbf{x}) &= \int_{\|\mathbf{z}\| \leq 1} (\mathbf{v}^\top (\mathbf{x} + \xi \mathbf{z}))_+ \frac{d\mathbf{z}}{\text{Vol}(\mathbb{D}^d)} = \int_{-1}^1 \int_{\|\mathbf{b}\|^2 \leq 1 - a^2} (\mathbf{v}^\top \mathbf{x} + \xi \|\mathbf{v}\| a)_+ \frac{d\mathbf{b} da}{\text{Vol}(\mathbb{D}^d)} \\ &= \int_{-1}^1 (\mathbf{v}^\top \mathbf{x} + \xi \|\mathbf{v}\| a)_+ (1 - a^2)^{\frac{d-1}{2}} \frac{\text{Vol}(\mathbb{D}^{d-1})}{\text{Vol}(\mathbb{D}^d)} da = \int_{-c}^1 (\mathbf{v}^\top \mathbf{x} + \xi \|\mathbf{v}\| a) (1 - a^2)^{\frac{d-1}{2}} \frac{\text{Vol}(\mathbb{D}^{d-1})}{\text{Vol}(\mathbb{D}^d)} da \\ &= \int_{-\arcsin(c)}^{\pi/2} (\mathbf{v}^\top \mathbf{x} + \xi \|\mathbf{v}\| \sin(\varphi)) \cos(\varphi)^d \frac{\text{Vol}(\mathbb{D}^{d-1})}{\text{Vol}(\mathbb{D}^d)} d\varphi, \end{aligned}$$

where $a := \hat{\mathbf{v}}^\top \mathbf{z}$ and $\mathbf{b} := \mathbf{P}_v \mathbf{z}$. To prove that the right hand part is a polynomial in c and $\sqrt{1 - c^2}$, notice the following identities.

$$\begin{aligned} \int_{-\arcsin(c)}^{\pi/2} \sin(\varphi) \cos(\varphi)^d d\varphi &= \frac{\cos(\arcsin(c))^{d+1}}{d+1} = \frac{(1 - c^2)^{\frac{d+1}{2}}}{d+1}. \\ \int_{-\arcsin(c)}^{\pi/2} \cos(\varphi)^d d\varphi &= \int_{-\arcsin(c)}^{\pi/2} \left(\frac{e^{i\varphi} + e^{-i\varphi}}{2} \right)^d d\varphi = \sum_{k=0}^{\frac{d-1}{2}} \binom{d}{k} \frac{e^{i(d-2k)\varphi} - e^{-i(d-2k)\varphi}}{i(d-2k)2^d} \Big|_{-\arcsin(c)}^{\pi/2} \\ &= \sum_{k=0}^{\frac{d-1}{2}} \binom{d}{k} \frac{\sin((d-2k)\arcsin(c)) + \sin(\frac{(d-2k)\pi}{2})}{(d-2k)2^{d-1}}. \end{aligned}$$

Finally, notice that

$$\sin(k \arcsin(c)) = \text{Im}(e^{ik \arcsin(c)}) = \text{Im}((\sqrt{1 - c^2} + ic)^k) \in \text{Poly}(c, \sqrt{1 - c^2}).$$

Thus, $\phi(\mathbf{v}, \mathbf{x}) \in \xi \|\mathbf{v}\| \text{Poly}(c, \sqrt{1 - c^2})$. \square

Remark B.3. Notice that, for even d , the activation function will have a term proportional to $\arcsin(c)$ in addition to the polynomial in c . While $\arcsin(c)$ is not definable in the smallest structure on $(\mathbb{R}, +, \times)$ (see discussion below), this function is definable on the structure $(\mathbb{R}, +, \times, \mathcal{A})$, where \mathcal{A} is the class of all restricted functions on $[-1, 1]^n$ (Example 1.5 of [Loi, 2010](#)).

Definability of ϕ To prove that ϕ is definable in the o-minimal structure on $(\mathbb{R}, +, \times)$, we will employ the following properties of definable functions ([Appendix B.1, Ji & Telgarsky, 2020](#)).

1. Let $f, g: D \rightarrow \mathbb{R}$ be definable functions. Then $\forall \alpha, \beta \in \mathbb{R} \alpha f + \beta g$ and $f g$ are definable. If $g \neq 0$, f/g is definable. If $f \geq 0$, $\forall l \in \mathbb{N} \sqrt[l]{f}$ is definable.
2. Let $f: D \rightarrow \mathbb{R}^d$. Then f is definable iff all coordinate projections of f are definable.
3. Composition of definable functions is definable.
4. Any coordinate permutation of a definable set is definable.

5. The image and pre-image of a definable set by a definable function is definable.
6. Any combination of finitely many definable functions with disjoint domains is definable. For example, the point-wise maximum and minimum of definable functions are definable.
7. Polynomial functions are definable.

Now, notice that $\mathbf{v}^\top \mathbf{x}$ is a linear function, $\|\mathbf{v}\|$ is a square root of a quadratic polynomial. Thus, their ratio $\frac{\mathbf{v}^\top \mathbf{x}}{\xi \|\mathbf{v}\|}$ is definable. Therefore, c is definable since it is a composition of the maxima and minima of definable functions. Finally, $\phi(\hat{\mathbf{v}}, \mathbf{x})$ is definable since it is a polynomial of two definable functions c and $\sqrt{1 - c^2}$.

Local Lipschitzness To prove that $\phi(\hat{\mathbf{v}}, \mathbf{x})$ is locally Lipschitz on A , we prove a stronger property that $\phi(\mathbf{v}, \mathbf{x})$ is twice continuously differentiable on A , which would be useful later. Notice that

$$\phi(\mathbf{v}, \mathbf{x}) = \int_{-c}^1 (\mathbf{v}^\top \mathbf{x} + \xi \|\mathbf{v}\| a) (1 - a^2)^{\frac{d-1}{2}} \frac{\text{Vol}(\mathbb{D}^{d-1})}{\text{Vol}(\mathbb{D}^d)} da.$$

Denote $H(z, \mathbf{v}, \mathbf{x}) := \int_{-z}^1 (\mathbf{v}^\top \mathbf{x} + \xi \|\mathbf{v}\| a) (1 - a^2)^{\frac{d-1}{2}} \frac{\text{Vol}(\mathbb{D}^{d-1})}{\text{Vol}(\mathbb{D}^d)} da$. Notice that $\phi(\mathbf{v}, \mathbf{x}) = H(c, \mathbf{v}, \mathbf{x})$.

By the chain rule for total derivative and the Leibniz integral rule (notice that $\mathbf{v}^\top \mathbf{x} + \xi \|\mathbf{v}\| a$ is smooth for $\mathbf{v} \neq \mathbf{0}$), when $\frac{\mathbf{v}^\top \mathbf{x}}{\xi \|\mathbf{v}\|} \notin \{-1, 1\}$, we get

$$\begin{aligned} \nabla_{\mathbf{v}} \phi(\mathbf{v}, \mathbf{x}) &= \frac{\partial H}{\partial c} \nabla_{\mathbf{v}} c + \nabla_{\mathbf{v}} H \\ &= (\mathbf{v}^\top \mathbf{x} - \xi \|\mathbf{v}\| c) (1 - c^2)^{\frac{d-1}{2}} \frac{\text{Vol}(\mathbb{D}^{d-1})}{\text{Vol}(\mathbb{D}^d)} \nabla_{\mathbf{v}} c + \int_{-c}^1 \left(\mathbf{x} + \xi \frac{\mathbf{v}}{\|\mathbf{v}\|} a \right) (1 - a^2)^{\frac{d-1}{2}} \frac{\text{Vol}(\mathbb{D}^{d-1})}{\text{Vol}(\mathbb{D}^d)} da \\ &= \int_{-c}^1 \left(\mathbf{x} + \xi \frac{\mathbf{v}}{\|\mathbf{v}\|} a \right) (1 - a^2)^{\frac{d-1}{2}} \frac{\text{Vol}(\mathbb{D}^{d-1})}{\text{Vol}(\mathbb{D}^d)} da. \end{aligned}$$

Thus, $\nabla_{\mathbf{v}} \phi(\mathbf{v}, \mathbf{x})$ exists and is continuous when $\frac{\mathbf{v}^\top \mathbf{x}}{\xi \|\mathbf{v}\|} \notin \{-1, 1\}$. When $\frac{\mathbf{v}^\top \mathbf{x}}{\xi \|\mathbf{v}\|} \in \{-1, 1\}$ (which implies $c \in \{-1, 1\}$), we could find the derivative by the definition. Notice that c is locally Lipschitz in \mathbf{v} when $\mathbf{v} \neq \mathbf{0}$ since it is a clipping of a continuously differentiable function. Now, using the Leibniz rule again, we get

$$\begin{aligned} \lim_{d\mathbf{v} \rightarrow \mathbf{0}} \frac{\phi(\mathbf{v} + d\mathbf{v}, \mathbf{x}) - \phi(\mathbf{v}, \mathbf{x})}{\|d\mathbf{v}\|} &= \lim_{d\mathbf{v} \rightarrow \mathbf{0}} \frac{H(c(\mathbf{v} + d\mathbf{v}), \mathbf{v} + d\mathbf{v}, \mathbf{x}) - H(c(\mathbf{v} + d\mathbf{v}), \mathbf{v}, \mathbf{x}) + H(c(\mathbf{v} + d\mathbf{v}), \mathbf{v}, \mathbf{x}) - H(c(\mathbf{v}), \mathbf{v}, \mathbf{x})}{\|d\mathbf{v}\|} \\ &= \lim_{d\mathbf{v} \rightarrow \mathbf{0}} \int_{-c(\mathbf{v} + d\mathbf{v})}^1 \left(\mathbf{x} + \xi \frac{\mathbf{v} + d\mathbf{v}}{\|\mathbf{v} + d\mathbf{v}\|} a \right) (1 - a^2)^{\frac{d-1}{2}} \frac{\text{Vol}(\mathbb{D}^{d-1})}{\text{Vol}(\mathbb{D}^d)} da + o(1) \\ &\quad + \frac{\int_{-c(\mathbf{v} + d\mathbf{v})}^{-c(\mathbf{v})} (\mathbf{v}^\top \mathbf{x} + \xi \|\mathbf{v}\| a) (1 - a^2)^{\frac{d-1}{2}} \frac{\text{Vol}(\mathbb{D}^{d-1})}{\text{Vol}(\mathbb{D}^d)} da}{\|d\mathbf{v}\|}. \end{aligned}$$

Since $|c(\mathbf{v})| = 1$, the integrand in the last expression satisfy

$$\begin{aligned} |(\mathbf{v}^\top \mathbf{x} + \xi \|\mathbf{v}\| a) (1 - a^2)^{\frac{d-1}{2}}| &\leq (|\mathbf{v}^\top \mathbf{x}| + \xi \|\mathbf{v}\|) (c(\mathbf{v})^2 - c(\mathbf{v} + d\mathbf{v})^2)^{\frac{d-1}{2}} \\ &\leq (|\mathbf{v}^\top \mathbf{x}| + \xi \|\mathbf{v}\|) 2^{\frac{d-1}{2}} |c(\mathbf{v}) - c(\mathbf{v} + d\mathbf{v})|^{\frac{d-1}{2}} \\ &= O(\|d\mathbf{v}\|^{\frac{d-1}{2}}). \end{aligned}$$

Thus, using Lebesgue's dominated convergence theorem, we get

$$\begin{aligned} \lim_{d\mathbf{v} \rightarrow \mathbf{0}} \frac{\phi(\mathbf{v} + d\mathbf{v}, \mathbf{x}) - \phi(\mathbf{v}, \mathbf{x})}{\|d\mathbf{v}\|} &= \lim_{d\mathbf{v} \rightarrow \mathbf{0}} \int_{-c(\mathbf{v}+d\mathbf{v})}^1 \left(\mathbf{x} + \xi \frac{\mathbf{v} + d\mathbf{v}}{\|\mathbf{v} + d\mathbf{v}\|} a \right) (1 - a^2)^{\frac{d-1}{2}} \frac{\text{Vol}(\mathbb{D}^{d-1})}{\text{Vol}(\mathbb{D}^d)} da + o(1) + O(\|d\mathbf{v}\|^{\frac{d-1}{2}}) \\ &= \int_{-c}^1 \left(\mathbf{x} + \xi \frac{\mathbf{v}}{\|\mathbf{v}\|} a \right) (1 - a^2)^{\frac{d-1}{2}} \frac{\text{Vol}(\mathbb{D}^{d-1})}{\text{Vol}(\mathbb{D}^d)} da. \end{aligned}$$

Therefore, $\phi(\mathbf{v}, \mathbf{x})$ is continuously differentiable in \mathbf{v} when $\mathbf{v} \neq \mathbf{0}$.

Similarly, we get the following Hessian

$$\begin{aligned} \nabla_{\mathbf{v}}^2 \phi(\mathbf{v}, \mathbf{x}) &= \int_{-c}^1 \frac{\xi a}{\|\mathbf{v}\|} \mathbf{P}_{\mathbf{v}} (1 - a^2)^{\frac{d-1}{2}} \frac{\text{Vol}(\mathbb{D}^{d-1})}{\text{Vol}(\mathbb{D}^d)} da + \left(\mathbf{x} - \xi \frac{\mathbf{v}}{\|\mathbf{v}\|} c \right) (1 - c^2)^{\frac{d-1}{2}} \frac{\text{Vol}(\mathbb{D}^{d-1})}{\text{Vol}(\mathbb{D}^d)} (\nabla_{\mathbf{v}} c)^{\top} \\ &= \mathbf{P}_{\mathbf{v}} \frac{\xi(1 - c^2)^{\frac{d+1}{2}} \text{Vol}(\mathbb{D}^{d-1})}{(d+1)\|\mathbf{v}\| \text{Vol}(\mathbb{D}^d)} + \left(\mathbf{x} - \xi \frac{\mathbf{v}}{\|\mathbf{v}\|} c \right) (1 - c^2)^{\frac{d-1}{2}} \frac{\text{Vol}(\mathbb{D}^{d-1})}{\text{Vol}(\mathbb{D}^d)} (\nabla_{\mathbf{v}} c)^{\top}. \end{aligned}$$

Notice that when $c \in (-1, 1)$

$$\nabla_{\mathbf{v}} c = \mathbf{P}_{\mathbf{v}} \frac{\mathbf{x}}{\xi}, \quad \mathbf{x} - \xi \frac{\mathbf{v}}{\|\mathbf{v}\|} c = \mathbf{P}_{\mathbf{v}} \mathbf{x}.$$

Thus, the Hessian exists, is continuous, and has the form $\nabla_{\mathbf{v}}^2 \phi(\mathbf{v}, \mathbf{x}) = \mathbf{P}_{\mathbf{v}} \mathbf{A}(\mathbf{v}, \mathbf{x}) \mathbf{P}_{\mathbf{v}}$.

Convergence of \mathbf{r}_j First, notice that $\exists \lim_{t \rightarrow \infty} s_j G(\mathbf{r}_j) =: G_j^*$ because $s_j G(\mathbf{r}_j)$ is monotonically increasing, $\mathbf{r}_j \in \mathbb{S}^{d-1}$, and G is continuous. Now, consider function $\zeta_j(t) := \int_0^t \left\| \frac{d\mathbf{r}_j}{dt} \right\| d\tau$, function $f_j(\mathbf{r}_j) := G_j^* - s_j G(\hat{\mathbf{r}}_j)$, and desingularizing function ψ_j , corresponding to f_j . Then for big enough t we have

$$\frac{df_j}{d\zeta_j} = -\|\nabla_{\mathbf{r}} G(\hat{\mathbf{r}}_j)\| = -\|\nabla_{\mathbf{r}} f_j(\mathbf{r}_j)\| \leq -\frac{1}{\psi_j'(f_j(\mathbf{r}_j))} \implies -\frac{d\psi_j}{dt} \geq \frac{d\zeta_j}{dt} \implies \zeta_j(t) \leq \zeta_j(t_0) + \psi_j(t_0) - \psi_j(t).$$

Thus, ζ_j is bounded, and hence $\lim_{t \rightarrow \infty} \zeta_j$ exists. Therefore, the trajectory of \mathbf{r}_j has a finite length, and hence $\lim_{t \rightarrow \infty} \mathbf{r}_j =: \hat{\mathbf{v}}_j^*$ exists. Also, since

$$\frac{df_j}{dt} = -\|\nabla_{\mathbf{r}} G(\hat{\mathbf{r}}_j)\|^2,$$

$\nabla_{\mathbf{r}} G(\hat{\mathbf{r}})$ is differentiable, and $\hat{\mathbf{r}}_j$ converges. Moreover, $\lim_{t \rightarrow \infty} \|\nabla_{\mathbf{r}} G(\hat{\mathbf{r}}_j)\| = 0$ (otherwise f would increase to infinity). Thus, $\hat{\mathbf{v}}_j^*$ is critical: $\mathbf{P}_{\hat{\mathbf{v}}_j^*} \mathbf{g}(\hat{\mathbf{v}}_j^*) = \mathbf{0}$.

B.1.2. DYNAMICS AROUND CRITICAL DIRECTIONS

We divide all neurons into two categories: prominent and non-prominent. Prominent neurons achieve the optimum of function G : $s_j G(\hat{\mathbf{v}}_j^*) = \lambda := \max_{\hat{\mathbf{v}} \in \mathbb{S}^{d-1}} |G(\hat{\mathbf{v}})|$. Denote $\lambda_j := s_j G(\hat{\mathbf{v}}_j^*)$, $P := \{j \mid \lambda_j = \lambda\}$ (prominent neurons), and $R := [m] \setminus P$ (non-prominent neurons).

Dynamics of small neurons First, we describe the dynamics of neurons from R . We have

$$u_j^l = u_j(0) \exp\left(\int_0^t s_j G(\hat{\mathbf{v}}_j^l) d\tau\right) \leq u_j(0) e^{\lambda_j t}.$$

Notice that the norm growth of these neurons is slower than $e^{\lambda t}$.

Dynamics of the big neurons' directions Now, we describe the dynamics of neurons from P . Near critical directions, the Hessian of G restricted on the unit sphere is non-positive if $s_j = 1$ and non-negative if $s_j = -1$ due to the local characterization of extremum. Consider the case $s_j = 1$ (the case $s_j = -1$ is similar). Near critical direction, we have

$$(\hat{v}_j^*)^\top \frac{d\hat{v}_j^l}{dt} = (\hat{v}_j^*)^\top \mathbf{P}_{\hat{v}_j^l} \mathbf{g}(\hat{v}_j^l) = (\hat{v}_j^*)^\top \mathbf{P}_{\hat{v}_j^l} (\mathbf{g}(\hat{v}_j^*) + \nabla^2 G(\hat{v}_j^*) (\hat{v}_j^l - \hat{v}_j^*) + O(\|\hat{v}_j^l - \hat{v}_j^*\|^2))$$

Denote $\hat{v}_j^l = \cos(\alpha)\hat{v}_j^* + \sin(\alpha)\epsilon$, where $\cos(\alpha) = (\hat{v}_j^*)^\top \hat{v}_j^l$ and $\epsilon := \frac{\mathbf{P}_{\hat{v}_j^*} \hat{v}_j^l}{\|\mathbf{P}_{\hat{v}_j^*} \hat{v}_j^l\|}$. We get

$$\mathbf{P}_{\hat{v}_j^l} \hat{v}_j^* = (1 - \cos(\alpha)^2)\hat{v}_j^* - \sin(\alpha)\cos(\alpha)\epsilon.$$

Since G is homogeneous, we get

$$\mathbf{g}(\hat{v}_j^*) = \lambda_j \hat{v}_j^*.$$

Also, we know that

$$\nabla^2 G(\hat{v}) = \mathbf{P}_{\hat{v}} \nabla^2 G(\hat{v}) \mathbf{P}_{\hat{v}} \implies \nabla^2 G(\hat{v}_j^*) (\hat{v}_j^l - \hat{v}_j^*) = \nabla^2 G(\hat{v}_j^*) \sin(\alpha)\epsilon.$$

These equations give

$$\begin{aligned} (\hat{v}_j^*)^\top \frac{d\hat{v}_j^l}{dt} &= ((1 - \cos(\alpha)^2)\hat{v}_j^* - \sin(\alpha)\cos(\alpha)\epsilon)^\top (\lambda \hat{v}_j^* + \nabla^2 G(\hat{v}_j^*) \sin(\alpha)\epsilon + O(\alpha^2)) \\ &= \lambda(1 - \cos(\alpha)^2) - \epsilon^\top \nabla^2 G(\hat{v}_j^*) \epsilon \sin(\alpha)^2 \cos(\alpha) + O(\alpha^3) \\ &\geq \lambda \|\hat{v}_j^* - \hat{v}_j^l\|^2 + O(\alpha^3). \end{aligned}$$

This equation implies

$$\begin{aligned} \frac{d\|\hat{v}_j^l - \hat{v}_j^*\|^2}{dt} &\leq -2\lambda \|\hat{v}_j^l - \hat{v}_j^*\|^2 + 2\delta \|\hat{v}_j^l - \hat{v}_j^*\|^3 \implies \frac{\|\hat{v}_j^l - \hat{v}_j^*\|}{1 - \frac{\delta}{\lambda} \|\hat{v}_j^l - \hat{v}_j^*\|} \leq \frac{\|\hat{v}_j^l(t_{0,j}) - \hat{v}_j^*(t_{0,j})\| e^{-\lambda(t-t_{0,j})}}{1 - \frac{\delta}{\lambda} \|\hat{v}_j^l(t_{0,j}) - \hat{v}_j^*(t_{0,j})\|} \\ \implies \|\hat{v}_j^l - \hat{v}_j^*\| &\leq \frac{\|\hat{v}_j^l(t_{0,j}) - \hat{v}_j^*(t_{0,j})\| e^{-\lambda(t-t_{0,j})}}{1 - \frac{\delta}{\lambda} \|\hat{v}_j^l(t_{0,j}) - \hat{v}_j^*(t_{0,j})\| (1 - e^{-\lambda(t-t_{0,j})})}, \end{aligned}$$

where $2\delta \|\hat{v}_j^l - \hat{v}_j^*\|^3$ comes from the term $O(\alpha^3)$ and $t_{0,j}$ is chosen such that $\|\hat{v}_j^l - \hat{v}_j^*\|$ is sufficiently small to bound our big-O term and, at the same time, $\frac{\delta}{\lambda} \|\hat{v}_j^l(t_0) - \hat{v}_j^*(t_0)\| \leq \frac{1}{2}$. Thus, we get an exponentially fast convergence near the critical direction.

Denote

$$t_0 := \max_j t_{0,j}, \quad c_0 := \max_j 2\|\hat{v}_j^l(t_{0,j}) - \hat{v}_j^*(t_{0,j})\| e^{\lambda t_{0,j}}.$$

We get

$$\forall j \in P, t \geq t_0 \quad \|\hat{v}_j^l - \hat{v}_j^*\| \leq c_0 e^{-\lambda t}.$$

Dynamics of the big neurons' scales Now, we will describe the dynamics of u_j^l . Notice

$$\frac{du_j^l}{dt} = u_j^l G(\hat{v}_j^l) \implies u_j^l = u_j^l(0) \exp\left(\int_0^t G(\hat{v}_j^l) d\tau\right).$$

This equality motivates us to consider a limit

$$u_j^* := \lim_{t \rightarrow \infty} u_j^l e^{-\lambda t}.$$

This limit exists since the right hand part is monotonically decreasing. We get

$$u_j^l = u_j^l(0) \exp\left(\int_0^t G(\hat{v}_j^l) d\tau\right) = u_j^* \exp\left(\lambda t + \int_t^\infty (\lambda - G(\hat{v}_j^l)) d\tau\right).$$

Thus,

$$u_j^* e^{\lambda t} \leq u_j^l \leq u_j^* \exp\left(\lambda t + \int_t^\infty c_0 b_0 e^{-\lambda \tau} d\tau\right) = u_j^* \exp\left(\lambda t + \frac{c_0 b_0}{\lambda} e^{-\lambda t}\right),$$

where $b_0 := \sup_{\mathbf{v} \neq \mathbf{v}'} \frac{|G(\mathbf{v}) - G(\mathbf{v}')|}{\|\mathbf{v} - \mathbf{v}'\|}$. Notice that this derivation also implies that neurons in P grow as $\Theta(e^{\lambda t})$.

Remark B.4. It is easy to see that $b_0 \leq (-\ell'(0))(1 + \xi)$.

Neuron capture Notice that all global extrema are attractive. Therefore, for isotropic initialization, with probability at least $1 - (f_j)^m$, some point will be attracted to the global extremum $\hat{\mathbf{v}}_j^*$, where $f_j := 1 - \frac{\text{volume of attraction region}}{\text{volume of sphere}}$.

B.2. Coupling Equations (1) and (2)

B.2.1. NORM GROWTH RATE

First, we prove the following proposition and lemma.

Proposition B.5. *Assuming $\|\mathbf{x}_i\| \leq 1$, the following identities hold*

$$\left| \frac{1}{n} \sum_{i=1}^n (-\ell'(f(\boldsymbol{\theta}, \mathbf{x}_i)y_i) + \ell'(0))\phi(\hat{\mathbf{v}}, \mathbf{x}_i)y_i \right| \leq a\|\boldsymbol{\theta}\|_{[m]}^2,$$

$$\left| \frac{1}{n} \sum_{i=1}^n (-\ell'(f(\boldsymbol{\theta}, \mathbf{x}_i)y_i) + \ell'(0))\nabla_{\mathbf{v}}\phi(\mathbf{z}, \mathbf{x}_i)|_{\mathbf{z}=\hat{\mathbf{v}}y_i} \right| \leq a\|\boldsymbol{\theta}\|_{[m]}^2,$$

where $a := \frac{m(1+\xi)^2}{4}$.

Proof. Notice that

$$|f(\boldsymbol{\theta}, \mathbf{x}_i)| \leq \sum_{j=1}^m u_j^2 |\phi(\hat{\mathbf{v}}_j, \mathbf{x}_i)| \leq (1 + \xi) \sum_{j=1}^m u_j^2 \leq m(1 + \xi)\|\boldsymbol{\theta}\|_{[m]}^2,$$

and

$$|-\ell'(f(\boldsymbol{\theta}, \mathbf{x}_i)y_i) + \ell'(0)| \leq |f(\boldsymbol{\theta}, \mathbf{x}_i)| \sup_z |\ell''(z)| \leq \frac{m(1 + \xi)}{4} \|\boldsymbol{\theta}\|_{[m]}^2.$$

Thus,

$$\left| \frac{1}{n} \sum_{i=1}^n (-\ell'(f(\boldsymbol{\theta}, \mathbf{x}_i)y_i) + \ell'(0))\phi(\hat{\mathbf{v}}, \mathbf{x}_i)y_i \right| \leq \frac{m(1 + \xi)}{4} \|\boldsymbol{\theta}\|_{[m]}^2 \max_{i, \hat{\mathbf{v}}} |\phi(\hat{\mathbf{v}}, \mathbf{x}_i)| \leq a\|\boldsymbol{\theta}\|_{[m]}^2.$$

Similarly for the gradients of activation. □

Lemma B.6. *Assume that $\boldsymbol{\theta}$ follows Equation (1) and $|u_j(0)| = \|\mathbf{v}_j(0)\|$. Then*

$$\forall t \leq t_1 \quad \|\boldsymbol{\theta}\|_{[m]}^2 \leq 2\|\boldsymbol{\theta}(0)\|_{[m]}^2 e^{2\lambda t},$$

where $t_1 := \frac{1}{2\lambda} \ln\left(\frac{\lambda}{2a\|\boldsymbol{\theta}(0)\|_{[m]}^2}\right)$.

Proof. Thus,

$$\frac{du_j}{dt} = G(\mathbf{v}_j) + \frac{1}{n} \sum_{i=1}^n (p_i(\boldsymbol{\theta}) + \ell'(0))\phi(\mathbf{v}_j, \mathbf{x}_i)y_i \implies \left| \frac{du_j}{dt} \right| \leq |u_j| |G(\hat{\mathbf{v}}_j)| + a\|\boldsymbol{\theta}\|_{[m]}^2 |u_j| \leq \lambda\|\boldsymbol{\theta}\|_{[m]} + a\|\boldsymbol{\theta}\|_{[m]}^3$$

$$\implies \int \frac{d\|\boldsymbol{\theta}\|_{[m]}^2}{\|\boldsymbol{\theta}\|_{[m]}^2 + \frac{a}{\lambda}\|\boldsymbol{\theta}\|_{[m]}^4} \leq 2\lambda t \implies \|\boldsymbol{\theta}\|_{[m]}^2 \leq \frac{\|\boldsymbol{\theta}(0)\|_{[m]}^2 e^{2\lambda t}}{1 - \frac{a}{\lambda}\|\boldsymbol{\theta}(0)\|_{[m]}^2 (e^{2\lambda t} - 1)}.$$

Therefore,

$$\|\boldsymbol{\theta}\|_{[m]}^2 \leq 2\|\boldsymbol{\theta}(0)\|_{[m]}^2 e^{2\lambda t} \quad \forall t \leq \frac{1}{2\lambda} \ln\left(\frac{\lambda}{2a\|\boldsymbol{\theta}(0)\|_{[m]}^2}\right).$$

□

B.2.2. COUPLING DIRECTIONS

We have

$$\begin{aligned}\frac{d\hat{\mathbf{v}}_j}{dt} &= s_j \mathbf{P}_{\mathbf{v}} \left(\mathbf{g}(\hat{\mathbf{v}}_j) + \frac{1}{n} \sum_{i=1}^n (-\ell'(f(\boldsymbol{\theta}, \mathbf{x}_i) y_i) + \ell'(0)) \nabla_{\mathbf{v}} \phi(\mathbf{v}_j, \mathbf{x}_i) y_i \right), \\ \frac{d\hat{\mathbf{v}}_j^l}{dt} &= s_j \mathbf{P}_{\mathbf{v}^l} \mathbf{g}(\hat{\mathbf{v}}_j^l).\end{aligned}$$

Denote $b_1 := \sup_{\hat{\mathbf{v}}, \hat{\mathbf{v}}' \in S^{d-1}, \hat{\mathbf{v}} \neq \hat{\mathbf{v}}'} \frac{\|\mathbf{P}_{\hat{\mathbf{v}}} \mathbf{g}(\hat{\mathbf{v}}) - \mathbf{P}_{\hat{\mathbf{v}}'} \mathbf{g}(\hat{\mathbf{v}}')\|}{\|\hat{\mathbf{v}} - \hat{\mathbf{v}}'\|}$ (since ϕ is twice continuously differentiable, this constant is defined).

Remark B.7. It is easy to see that $b_1 \leq 2 \sup_{\hat{\mathbf{v}}} \|\mathbf{g}(\hat{\mathbf{v}})\| + (-\ell'(0)) \sup_{i, \hat{\mathbf{v}}} \|\nabla_{\mathbf{v}}^2 \phi(\mathbf{v}, \mathbf{x}_i)\| \implies b_1 = O(1/\xi)$, when $\xi \rightarrow 0$.

We get

$$\begin{aligned}\frac{d\|\hat{\mathbf{v}}_j - \hat{\mathbf{v}}_j^l\|}{dt} &\leq \left\| \frac{d\hat{\mathbf{v}}_j}{dt} - \frac{d\hat{\mathbf{v}}_j^l}{dt} \right\| \leq \left\| \frac{1}{n} \sum_{i=1}^n (-\ell'(f(\boldsymbol{\theta}, \mathbf{x}_i) y_i) + \ell'(0)) \nabla_{\mathbf{v}} \phi(\mathbf{v}_j, \mathbf{x}_i) y_i \right\| + \|\mathbf{P}_{\hat{\mathbf{v}}} \mathbf{g}(\hat{\mathbf{v}}) - \mathbf{P}_{\hat{\mathbf{v}}^l} \mathbf{g}(\hat{\mathbf{v}}^l)\| \\ &\leq a \|\boldsymbol{\theta}\|_{[m]}^2 + b_1 \|\hat{\mathbf{v}}_j - \hat{\mathbf{v}}_j^l\|.\end{aligned}$$

Consider function $h := \|\hat{\mathbf{v}}_j - \hat{\mathbf{v}}_j^l\| e^{-b_1 t}$, we get

$$\begin{aligned}\forall t \leq t_1 \quad \frac{dh}{dt} &\leq a \|\boldsymbol{\theta}\|_{[m]}^2 e^{-b_1 t} \implies \forall t \leq t_1 \quad h \leq \int_0^t 2a \|\boldsymbol{\theta}(0)\|_{[m]}^2 e^{(2\lambda - b_1)\tau} d\tau \\ \implies \forall t \leq t_1 \quad h &\leq \frac{2a}{b_1 - 2\lambda} \|\boldsymbol{\theta}(0)\|_{[m]}^2 (1 - e^{(2\lambda - b_1)t}) \implies \forall t \leq t_1 \quad \|\hat{\mathbf{v}}_j - \hat{\mathbf{v}}_j^l\| \leq \frac{2a}{b_1 - 2\lambda} \|\boldsymbol{\theta}(0)\|_{[m]}^2 e^{b_1 t},\end{aligned}$$

where we have assumed that $b_1 > 2\lambda$ (which holds for small enough ξ).

B.2.3. COUPLING SCALES

We have

$$\begin{aligned}\frac{du_j}{dt} &= G(\mathbf{v}_j) + \frac{1}{n} \sum_{i=1}^n (-\ell'(f(\boldsymbol{\theta}, \mathbf{x}_i) y_i) + \ell'(0)) \phi(\mathbf{v}_j, \mathbf{x}_i) y_i, \\ \frac{du_j^l}{dt} &= G(\mathbf{v}_j^l).\end{aligned}$$

It gives

$$\begin{aligned}\forall t \leq t_1 \quad \frac{d|u_j - u_j^l|}{dt} &\leq \left| \frac{du_j}{dt} - \frac{du_j^l}{dt} \right| \leq a \|\boldsymbol{\theta}\|_{[m]}^2 |u_j| + |u_j^l| |G(\hat{\mathbf{v}}_j) - G(\hat{\mathbf{v}}_j^l)| + |u_j - u_j^l| |G(\hat{\mathbf{v}}_j)| \\ &\leq 2\sqrt{2}a \|\boldsymbol{\theta}(0)\|_{[m]}^3 e^{3\lambda t} + \frac{2ab_0 \|\boldsymbol{\theta}(0)\|_{[m]}^3}{b_1 - 2\lambda} e^{(\lambda_j + b_1)t} + \lambda |u_j - u_j^l|.\end{aligned}$$

As previously, consider function $h := |u_j - u_j^l| e^{-\lambda t}$, we get

$$\begin{aligned}\forall t \leq t_1 \quad h &\leq 2\sqrt{2}a \|\boldsymbol{\theta}(0)\|_{[m]}^3 e^{2\lambda t} + \frac{2ab_0 \|\boldsymbol{\theta}(0)\|_{[m]}^3}{b_1 - 2\lambda} e^{(\lambda_j + b_1 - \lambda)t} \\ \implies \forall t \leq t_1 \quad |u_j - u_j^l| &\leq \frac{a\sqrt{2}}{\lambda} \|\boldsymbol{\theta}(0)\|_{[m]}^3 (e^{3\lambda t} - e^{\lambda t}) + \frac{2ab_0}{(b_1 - 2\lambda)(\lambda_j + b_1 - \lambda)} \|\boldsymbol{\theta}(0)\|_{[m]}^3 (e^{(\lambda_j + b_1)t} - e^{\lambda t}) \\ \implies \forall t \leq t_1 \quad |u_j - u_j^l| &\leq c_1 \|\boldsymbol{\theta}(0)\|_{[m]}^3 e^{(\lambda_j + b_1)t},\end{aligned}$$

where we have assumed that $b_1 \geq 4\lambda$ and denoted $c_1 := \frac{a\sqrt{2}}{\lambda} + \frac{2ab_0}{(b_1 - 2\lambda)^2}$.

B.3. Final Bound

The result of the previous sections show

$$\begin{aligned} \forall j \in P, t \in [t_0, t_1] \quad & \|\hat{\mathbf{v}}_j - \hat{\mathbf{v}}_j^*\| \leq \|\hat{\mathbf{v}}_j - \hat{\mathbf{v}}_j^l\| + \|\hat{\mathbf{v}}_j^l - \hat{\mathbf{v}}_j^*\| \leq \frac{2a}{b_1 - 2\lambda} \|\boldsymbol{\theta}(0)\|_{[m]}^2 e^{c^g t} + c_0 e^{-\lambda t}, \\ \forall j \in P, t \in [t_0, t_1] \quad & |u_j - u_j^* e^{\lambda t}| \leq |u_j - u_j^l| + |u_j^l - u_j^* e^{\lambda t}| \\ & \leq c_1 \|\boldsymbol{\theta}(0)\|_{[m]}^3 e^{(\lambda + b_1)t} + |u_j^*| e^{\lambda t} \left(\exp\left(\frac{c_0 b_0}{\lambda} e^{-\lambda t}\right) - 1 \right), \\ \forall j \in R, t \in [0, t_1] \quad & |u_j| \leq c_1 \|\boldsymbol{\theta}(0)\|_{[m]}^3 e^{(\lambda_j + c^g)t} + |u_j(0)| e^{\lambda_j t}. \end{aligned}$$

We want to control these errors until the point when the fastest growing directions will have a predefined scale, r . To do this, we will choose $T_1 = \frac{1}{\lambda} \ln\left(\frac{r}{\sigma}\right)$ and $\sigma = r^{\kappa+1}$. (Notice that t_0 does not depend on r , and $t_1 = -\ln(\sigma)/\lambda + O(1)$. Thus, for sufficiently small r , we have $t_0 < T_1 < t_1$.)

This functional form will give us the following errors:

$$\begin{aligned} \forall j \in P \quad & \|\hat{\mathbf{v}}_j - \hat{\mathbf{v}}_j^*\| \leq O(r^{2-\kappa(b_1/\lambda-2)} + r^\kappa), \\ \forall j \in P \quad & |u_j - u_j^* e^{\lambda T_1}| \leq O(r^{3-\kappa(b_1/\lambda-2)} + r^{1+\kappa}), \\ \forall j \in R \quad & |u_j| \leq O\left(r^{3-\kappa\left(\frac{b_1+\lambda_j}{\lambda}-3\right)} + r^{1+\kappa\left(1-\frac{\lambda_j}{\lambda}\right)}\right). \end{aligned}$$

We choose $\kappa^* = \frac{2}{\frac{b_1}{\lambda}-1}$. It will give

$$\begin{aligned} \forall j \in P \quad & \|\hat{\mathbf{v}}_j - \hat{\mathbf{v}}_j^*\| \leq O(r^{\kappa^*}), \\ \forall j \in P \quad & |u_j - u_j^* e^{\lambda T_1}| \leq O(r^{1+\kappa^*}), \\ \forall j \in R \quad & |u_j| \leq O(r^{1+\kappa_j}), \end{aligned}$$

where $\kappa_j := \kappa^*(1 - \lambda_j/\lambda)$.

C. Proof of Theorem 4.4

By Lemma 5.3 of Lyu et al. (2021), we could write the following dynamics for embedding (4)

$$\begin{aligned} \frac{du_j^x}{dt} &= \frac{1}{n} \sum_{i=1}^n (-\ell'(f(\boldsymbol{\theta}^x, \mathbf{x}_i) y_i)) \phi(\mathbf{v}_j^x, \mathbf{x}_i) y_i, \\ \frac{d\mathbf{v}_j^x}{dt} &= \frac{u_j^x}{n} \sum_{i=1}^n (-\ell'(f(\boldsymbol{\theta}^x, \mathbf{x}_i) y_i)) \nabla_{\mathbf{v}} \phi(\mathbf{v}_j^x, \mathbf{x}_i) y_i. \end{aligned} \tag{7}$$

Denote $\boldsymbol{\theta}^* := \boldsymbol{\theta}^x(T_1)/r$.

We will prove the theorem in two stages. First, we will couple Equations (1) and (7). Then we will investigate the sensitivity of Equation (7) to initialization scale. This analysis will allow us to analyze the solution to Equation (1).

C.1. Coupling Equations (1) and (7)

C.1.1. NORM GROWTH

Since $\|\boldsymbol{\theta}^x(T_1)\|_{[m]} = \Theta(r)$ and $\|\boldsymbol{\theta}^x(T_1) - \boldsymbol{\theta}(T_1)\|_{[m]} = o(r)$, we have that $\|\boldsymbol{\theta}(T_1)\| \leq \sqrt{2} \|\boldsymbol{\theta}^x(T_1)\|$ for sufficiently small r . Using this fact and Lemma B.6, we get

$$\begin{aligned} \forall t \in [T_1, T_1 + t_3] \quad & \|\boldsymbol{\theta}(T_1)\|_{[m]}^2 \leq 2q^2 e^{2\lambda(t-T_1)}, \\ \forall t \in [T_1, T_1 + t_3] \quad & \|\boldsymbol{\theta}^x(T_1)\|_{[m]}^2 \leq 2q^2 e^{2\lambda(t-T_1)}, \end{aligned}$$

where $t_3 := \frac{1}{2\lambda} \ln\left(\frac{\lambda}{2aq^2}\right)$ and $q^2 := 2\|\boldsymbol{\theta}^x(T_1)\|_{[m]}^2 = 2r^2 \|\boldsymbol{\theta}^*\|_{[m]}^2$.

C.1.2. BOUNDING $\|\hat{\mathbf{v}}_j^X - \hat{\mathbf{v}}_j^*\|$

W.l.o.g. assume that $s_j = 1$. Additionally, assume $\|\hat{\mathbf{v}}_j^X - \hat{\mathbf{v}}_j^*\| \leq \varepsilon$ and denote $\hat{\mathbf{v}}_j^X = \cos(\alpha')\hat{\mathbf{v}}_j^* + \sin(\alpha')\boldsymbol{\epsilon}$. First, notice that $\|\hat{\mathbf{v}}_j^X - \hat{\mathbf{v}}_j^*\| \leq \varepsilon \leq \Delta$ implies

$$\phi(\hat{\mathbf{v}}_j^X, \mathbf{x}_i) = ((\hat{\mathbf{v}}_j^X)^\top \mathbf{x}_i)_+ \implies \mathbf{g}(\hat{\mathbf{v}}_j^X) = \mathbf{g}(\hat{\mathbf{v}}_j^*).$$

We have

$$\begin{aligned} \forall j \in P(\hat{\mathbf{v}}_j^*)^\top \frac{d\hat{\mathbf{v}}_j^X}{dt} &= (\mathbf{P}_{\hat{\mathbf{v}}_j^X \hat{\mathbf{v}}_j^*})^\top \left(\mathbf{g}(\hat{\mathbf{v}}_j^X) + \frac{1}{n} \sum_{i=1}^n (-\ell'(f(\boldsymbol{\theta}^X, \mathbf{x}_i))y_i + \ell'(0)) \nabla_{\mathbf{v}} \phi(\mathbf{v}_j^X, \mathbf{x}_i) y_i \right) \\ &\geq ((1 - \cos(\alpha'))^2 \hat{\mathbf{v}}_j^* - \sin(\alpha') \cos(\alpha') \boldsymbol{\epsilon})^\top \mathbf{g}(\hat{\mathbf{v}}_j^X) - \sin(\alpha') a \|\boldsymbol{\theta}^X\|_{[m]}^2 \\ &= \lambda \sin(\alpha')^2 - \sin(\alpha') a \|\boldsymbol{\theta}^X\|_{[m]}^2 \end{aligned}$$

where $\mathbf{z} \in [\hat{\mathbf{v}}_j^*, \hat{\mathbf{v}}_j^X]$. It implies

$$\forall t \in [T_1, T_1 + t_3] \frac{d\alpha'}{dt} \leq 2aq^2 e^{2\lambda(t-T_1)} - \lambda \sin(\alpha') \implies \alpha' \leq \frac{a}{\lambda} q^2 (e^{2\lambda(t-T_1)} - 1).$$

So, to ensure $\|\hat{\mathbf{v}}_j^X - \hat{\mathbf{v}}_j^*\| \leq \varepsilon$, it is sufficient to have

$$\frac{a}{\lambda} q^2 e^{2\lambda(t-T_1)} \leq \varepsilon \iff t - T_1 \leq \frac{1}{2\lambda} \ln \left(\frac{\lambda \varepsilon}{aq^2} \right) =: t_4^\varepsilon.$$

 C.1.3. BOUNDING $\|\boldsymbol{\theta} - \boldsymbol{\theta}^X\|$

Coupling Directions First, we want to bound $\|\hat{\mathbf{v}}_j - \hat{\mathbf{v}}_j^X\|$. Assuming that $\|\hat{\mathbf{v}}_j - \hat{\mathbf{v}}_j^X\| \leq \varepsilon$, we get $\|\hat{\mathbf{v}}_j - \hat{\mathbf{v}}_j^*\| \leq 2\varepsilon$ and

$$\begin{aligned} \frac{d\|\hat{\mathbf{v}}_j - \hat{\mathbf{v}}_j^X\|^2}{dt} &= -\hat{\mathbf{v}}_j^\top \frac{d\hat{\mathbf{v}}_j^X}{dt} - (\hat{\mathbf{v}}_j^X)^\top \frac{d\hat{\mathbf{v}}_j}{dt} \\ &= -(\mathbf{P}_{\hat{\mathbf{v}}_j^X \hat{\mathbf{v}}_j})^\top \left(\mathbf{g}(\hat{\mathbf{v}}_j^X) + \frac{1}{n} \sum_{i=1}^n (-\ell'(f(\boldsymbol{\theta}^X, \mathbf{x}_i))y_i + \ell'(0)) \nabla_{\mathbf{v}} \phi(\mathbf{v}_j^X, \mathbf{x}_i) y_i \right) \\ &\quad - (\mathbf{P}_{\hat{\mathbf{v}}_j \hat{\mathbf{v}}_j^X})^\top \left(\mathbf{g}(\hat{\mathbf{v}}_j) + \frac{1}{n} \sum_{i=1}^n (-\ell'(f(\boldsymbol{\theta}, \mathbf{x}_i))y_i + \ell'(0)) \nabla_{\mathbf{v}} \phi(\mathbf{v}_j, \mathbf{x}_i) y_i \right) \\ &= -(\mathbf{P}_{\hat{\mathbf{v}}_j^X \hat{\mathbf{v}}_j})^\top (\mathbf{g}(\hat{\mathbf{v}}_j^X) - \mathbf{g}(\hat{\mathbf{v}}_j)) \\ &\quad - (\mathbf{P}_{\hat{\mathbf{v}}_j \hat{\mathbf{v}}_j^X})^\top \left(\frac{1}{n} \sum_{i=1}^n (-\ell'(f(\boldsymbol{\theta}^X, \mathbf{x}_i))y_i + \ell'(0)) (\nabla_{\mathbf{v}} \phi(\mathbf{v}_j^X, \mathbf{x}_i) - \nabla_{\mathbf{v}} \phi(\mathbf{v}_j, \mathbf{x}_i)) y_i \right) \\ &\quad - (\mathbf{P}_{\hat{\mathbf{v}}_j^X \hat{\mathbf{v}}_j})^\top \left(\frac{1}{n} \sum_{i=1}^n (-\ell'(f(\boldsymbol{\theta}^X, \mathbf{x}_i))y_i + \ell'(f(\boldsymbol{\theta}, \mathbf{x}_i))y_i) \nabla_{\mathbf{v}} \phi(\mathbf{v}_j, \mathbf{x}_i) y_i \right) \\ &\quad - (\mathbf{P}_{\hat{\mathbf{v}}_j^X \hat{\mathbf{v}}_j} + \mathbf{P}_{\hat{\mathbf{v}}_j \hat{\mathbf{v}}_j^X})^\top \left(\mathbf{g}(\hat{\mathbf{v}}_j) + \frac{1}{n} \sum_{i=1}^n (-\ell'(f(\boldsymbol{\theta}, \mathbf{x}_i))y_i + \ell'(0)) \nabla_{\mathbf{v}} \phi(\mathbf{v}_j, \mathbf{x}_i) y_i \right) \\ &= -(\mathbf{P}_{\hat{\mathbf{v}}_j^X \hat{\mathbf{v}}_j})^\top \left(\frac{1}{n} \sum_{i=1}^n (-\ell'(f(\boldsymbol{\theta}^X, \mathbf{x}_i))y_i + \ell'(f(\boldsymbol{\theta}, \mathbf{x}_i))y_i) \nabla_{\mathbf{v}} \phi(\mathbf{v}_j, \mathbf{x}_i) y_i \right) \\ &\quad - (\mathbf{P}_{\hat{\mathbf{v}}_j^X \hat{\mathbf{v}}_j} + \mathbf{P}_{\hat{\mathbf{v}}_j \hat{\mathbf{v}}_j^X})^\top \left(\mathbf{g}(\hat{\mathbf{v}}_j) + \frac{1}{n} \sum_{i=1}^n (-\ell'(f(\boldsymbol{\theta}, \mathbf{x}_i))y_i + \ell'(0)) \nabla_{\mathbf{v}} \phi(\mathbf{v}_j, \mathbf{x}_i) y_i \right). \end{aligned}$$

Denote $\cos(\beta) := \hat{\mathbf{v}}_j^\top \hat{\mathbf{v}}_j^X$. We get

$$\|\mathbf{P}_{\hat{\mathbf{v}}_j \hat{\mathbf{v}}_j^X}\|^2 = \|\hat{\mathbf{v}}_j^X - \cos(\beta)\hat{\mathbf{v}}_j\|^2 = 1 - \cos(\beta)^2 \leq \|\hat{\mathbf{v}}_j - \hat{\mathbf{v}}_j^X\|^2.$$

Now, notice

$$\begin{aligned}
 |f(\boldsymbol{\theta}, \mathbf{x}_i) - f(\boldsymbol{\theta}^X, \mathbf{x}_i)| &= \left| \sum_{j=1}^m (u_j^2 \phi(\hat{\mathbf{v}}_j, \mathbf{x}_i) - (u_j^X)^2 \phi(\hat{\mathbf{v}}_j^X, \mathbf{x}_i)) \right| \\
 &= \left| \sum_{j=1}^m ((u_j^2 - (u_j^X)^2) \phi(\hat{\mathbf{v}}_j, \mathbf{x}_i) + (u_j^X)^2 (\phi(\hat{\mathbf{v}}_j, \mathbf{x}_i) - \phi(\hat{\mathbf{v}}_j^X, \mathbf{x}_i))) \right| \\
 &= \left| \sum_{j=1}^m ((u_j^2 - (u_j^X)^2) \phi(\hat{\mathbf{v}}_j, \mathbf{x}_i) + (u_j^X)^2 [(\hat{\mathbf{v}}_j^*)^\top \mathbf{x}_i \geq 0] (\hat{\mathbf{v}}_j - \hat{\mathbf{v}}_j^X)^\top \mathbf{x}_i) \right| \\
 &\leq |P|(1 + \xi)(\|\boldsymbol{\theta}\|_P + \|\boldsymbol{\theta}^X\|_P)U + |R|(1 + \xi)\|\boldsymbol{\theta}\|_R^2 + |P|\|\boldsymbol{\theta}\|_{[m]}^2 V,
 \end{aligned}$$

where $U := \max_{j \in P} |u_j - u_j^X|$, $V := \max_{j \in P} \|\hat{\mathbf{v}}_j - \hat{\mathbf{v}}_j^X\|$. This formula implies

$$\begin{aligned}
 &\left\| (\mathbf{P}_{\hat{\mathbf{v}}_j^X} \hat{\mathbf{v}}_j)^\top \left(\frac{1}{n} \sum_{i=1}^n (-\ell'(f(\boldsymbol{\theta}^X, \mathbf{x}_i) y_i) + \ell'(f(\boldsymbol{\theta}, \mathbf{x}_i) y_i)) \nabla_{\mathbf{v}} \phi(\mathbf{v}_j, \mathbf{x}_i) y_i \right) \right\| \\
 &\leq aU(\|\boldsymbol{\theta}\|_{[m]} + \|\boldsymbol{\theta}^X\|_{[m]}) \|\hat{\mathbf{v}}_j - \hat{\mathbf{v}}_j^X\| + a\|\boldsymbol{\theta}\|_R^2 \|\hat{\mathbf{v}}_j - \hat{\mathbf{v}}_j^X\| + a\|\boldsymbol{\theta}^X\|_{[m]}^2 V \|\hat{\mathbf{v}}_j - \hat{\mathbf{v}}_j^X\|.
 \end{aligned}$$

Finally,

$$\mathbf{P}_{\hat{\mathbf{v}}_j^X} \hat{\mathbf{v}}_j + \mathbf{P}_{\hat{\mathbf{v}}_j} \hat{\mathbf{v}}_j^X = (1 - \cos(\beta))(\hat{\mathbf{v}}_j + \hat{\mathbf{v}}_j^X) = \|\hat{\mathbf{v}}_j - \hat{\mathbf{v}}_j^X\|^2 \frac{\hat{\mathbf{v}}_j + \hat{\mathbf{v}}_j^X}{2},$$

which gives

$$\begin{aligned}
 (\mathbf{P}_{\hat{\mathbf{v}}_j^X} \hat{\mathbf{v}}_j + \mathbf{P}_{\hat{\mathbf{v}}_j} \hat{\mathbf{v}}_j^X)^\top \left(\mathbf{g}(\hat{\mathbf{v}}_j) + \frac{1}{n} \sum_{i=1}^n (-\ell'(f(\boldsymbol{\theta}, \mathbf{x}_i) y_i) + \ell'(0)) \nabla_{\mathbf{v}} \phi(\mathbf{v}_j, \mathbf{x}_i) y_i \right) \\
 \geq \lambda \|\hat{\mathbf{v}}_j - \hat{\mathbf{v}}_j^X\|^2 \frac{\cos(\alpha) + \cos(\alpha')}{2} - a\|\boldsymbol{\theta}\|_{[m]}^2 \|\hat{\mathbf{v}}_j - \hat{\mathbf{v}}_j^X\|^2,
 \end{aligned}$$

where $\cos(\alpha) = \hat{\mathbf{v}}_j^\top \hat{\mathbf{v}}_j^*$. Together, we get

$$\begin{aligned}
 \frac{dV^2}{dt} &\leq a(\|\boldsymbol{\theta}\|_{[m]}^2 + \|\boldsymbol{\theta}^X\|_{[m]}^2) V^2 + a\|\boldsymbol{\theta}\|_R^2 V + a(\|\boldsymbol{\theta}\|_{[m]} + \|\boldsymbol{\theta}^X\|_{[m]}) UV \\
 &\leq 4aq^2 e^{2\lambda(t-T_1)} V^2 + a\|\boldsymbol{\theta}\|_R^2 V + 2\sqrt{2}aq e^{\lambda(t-T_1)} UV \\
 \implies \frac{dV}{dt} &\leq 2aq^2 e^{2\lambda(t-T_1)} V + \frac{a}{2} \|\boldsymbol{\theta}\|_R^2 + \sqrt{2}aq e^{\lambda(t-T_1)} U.
 \end{aligned}$$

Bounding $\|\boldsymbol{\theta}\|_R$ We get

$$\frac{du_j}{dt} = G(\mathbf{v}_j) + \frac{1}{n} \sum_{i=1}^n (-\ell'(f(\boldsymbol{\theta}, \mathbf{x}_i) y_i) + \ell'(0)) \phi(\mathbf{v}_j, \mathbf{x}_i) y_i \implies \left| \frac{du_j}{dt} \right| \leq |u_j| (\lambda + a\|\boldsymbol{\theta}\|_{[m]}^2).$$

Thus,

$$|u_j| \leq |u_j(T_1)| \exp\left(\int_{T_1}^t \lambda + a\|\boldsymbol{\theta}\|_{[m]}^2 \tau\right) \leq |u_j(T_1)| \exp\left(\lambda(t - T_1) + \frac{a}{\lambda} q^2 e^{2\lambda(t-T_1)}\right).$$

If $t - T_1 \leq t_4^\varepsilon$, we get

$$|u_j| \leq e^{\lambda(t-T_1)+\varepsilon} |u_j(T_1)| \implies \|\boldsymbol{\theta}\|_R \leq e^{\lambda(t-T_1)+\varepsilon} \|\boldsymbol{\theta}(T_1)\|_R.$$

Bounding U Similarly to the previous cases, we have

$$\begin{aligned} \frac{d|u_j - u_j^x|}{dt} &\leq \left| \frac{du_j}{dt} - \frac{du_j^x}{dt} \right| \\ &\leq |u_j G(\hat{\mathbf{v}}_j) - u_j^x G(\hat{\mathbf{v}}_j^x)| + \left| \frac{1}{n} \sum_{i=1}^n ((\ell'(0) - \ell'(f(\boldsymbol{\theta}, \mathbf{x}_i) y_i)) \phi(\mathbf{v}_j, \mathbf{x}_i) - (\ell'(0) - \ell'(f(\boldsymbol{\theta}^x, \mathbf{x}_i) y_i)) \phi(\mathbf{v}_j^x, \mathbf{x}_i)) y_i \right| \leq \\ &\lambda U + \|\boldsymbol{\theta}^x\|_{[m]} |(\hat{\mathbf{v}}_j - \hat{\mathbf{v}}_j^x)^\top \mathbf{g}(\hat{\mathbf{v}}_j)| + a \|\boldsymbol{\theta}^x\|_{[m]} (U(\|\boldsymbol{\theta}\|_{[m]} + \|\boldsymbol{\theta}^x\|_{[m]}) + \|\boldsymbol{\theta}\|_R^2 + \|\boldsymbol{\theta}^x\|_{[m]}^2) V + a \|\boldsymbol{\theta}\|_{[m]}^2 \|\mathbf{v}_j - \mathbf{v}_j^x\|. \end{aligned}$$

Notice

$$\begin{aligned} |(\hat{\mathbf{v}}_j - \hat{\mathbf{v}}_j^x)^\top \mathbf{g}(\hat{\mathbf{v}}_j)| &= \lambda |\cos(\alpha) - \cos(\alpha')| = 2\lambda \sin\left(\frac{\alpha + \alpha'}{2}\right) \sin\left(\frac{|\alpha - \alpha'|}{2}\right) \leq 2\lambda \sin\left(\frac{\beta}{2}\right) \sin\left(\frac{\beta}{2} + \alpha'\right) \\ &\leq \lambda \|\hat{\mathbf{v}}_j - \hat{\mathbf{v}}_j^x\| \left(\frac{\|\hat{\mathbf{v}}_j - \hat{\mathbf{v}}_j^x\|}{2} + \alpha' \right). \end{aligned}$$

Also, notice

$$\|\mathbf{v}_j - \mathbf{v}_j^x\| = \|u_j \hat{\mathbf{v}}_j - u_j^x \hat{\mathbf{v}}_j^x\| \leq |u_j - u_j^x| + |u_j^x| \|\hat{\mathbf{v}}_j - \hat{\mathbf{v}}_j^x\|.$$

It will give

$$\frac{dU}{dt} \leq \lambda U + 6aq^2 e^{2\lambda(t-T_1)} U + \lambda \sqrt{2} q e^{\lambda(t-T_1)} \frac{V^2}{2} + \sqrt{2} a q \|\boldsymbol{\theta}(T_1)\|_R^2 e^{3\lambda(t-T_1)+\varepsilon} + 5\sqrt{2} a q^3 e^{3\lambda(t-T_1)} V$$

Bounding U and V Consider $h(t) := \exp\left(-\int_{T_1}^t 2aq^2 e^{2\lambda(\tau-T_1)} d\tau\right) \leq 1$, $\tilde{U} := \frac{U}{q} e^{-\lambda(t-T_1)} h(t)^3$ and $\tilde{V} := Vh(t)$.

Notice that

$$h(t) = \exp\left(-\frac{aq^2}{\lambda} e^{2\lambda(t-T_1)}\right) \geq e^{-\varepsilon} \forall t \in [T_1, T_2^\varepsilon],$$

where $T_2^\varepsilon := T_1 + t_4^\varepsilon$. For $t \in [T_1, T_2^\varepsilon]$, we get

$$\begin{aligned} \frac{d\tilde{U}}{dt} &\leq \frac{\sqrt{2}\lambda}{2} \tilde{V}^2 + \sqrt{2} a e^{2\lambda(t-T_1)+\varepsilon} \|\boldsymbol{\theta}(T_1)\|_R^2 + 5\sqrt{2} a q^2 e^{2\lambda(t-T_1)} \tilde{V}, \\ \frac{d\tilde{V}}{dt} &\leq \frac{a}{2} e^{2\lambda(t-T_1)+\varepsilon} \|\boldsymbol{\theta}(T_1)\|_R^2 + \sqrt{2} a q^2 e^{2\lambda(t-T_1)+2\varepsilon} \tilde{U}. \end{aligned}$$

Consider $\tilde{W} := \max\{\tilde{V}, \tilde{U}\}$. We get

$$\frac{d\tilde{W}}{dt} \leq \lambda \tilde{W}^2 + 2a e^{2\lambda(t-T_1)+\varepsilon} \|\boldsymbol{\theta}(T_1)\|_R^2 + 8aq^2 e^{2\lambda(t-T_1)+2\varepsilon} \tilde{W}.$$

Again, consider $W := \tilde{W}h(t)^4$, we get

$$\frac{dW}{dt} \leq \lambda e^{4\varepsilon} W^2 + 2a e^{2\lambda(t-T_1)+5\varepsilon} \|\boldsymbol{\theta}(T_1)\|_R^2.$$

Now, consider the following system

$$\frac{d\bar{W}}{dt} = \lambda e^{4\varepsilon} \bar{W}^2 + 2a e^{2\lambda(t-T_1)+5\varepsilon} \|\boldsymbol{\theta}(T_1)\|_R^2.$$

where $\bar{W}(T_1) = W(T_1)$. It is easy to see that $\bar{W} \geq W$ since the right-hand part is the same function in both cases, this function is increasing in W , and the initial conditions are the same. Also, notice that \bar{W} is increasing. It implies

$$\bar{W}(t) \leq \bar{W}(T_1) + \lambda e^{4\varepsilon} \bar{W}(t)^2 (t - T_1) + \frac{a}{\lambda} e^{2\lambda(t-T_1)+5\varepsilon} \|\boldsymbol{\theta}(T_1)\|_R^2.$$

Assuming that $\lambda e^{4\varepsilon} \bar{W}(t)(t - T_1) \leq \frac{1}{2}$, we get

$$\bar{W} \leq 2\bar{W}(T_1) + \frac{2a}{\lambda} e^{2\lambda(t-T_1)+5\varepsilon} \|\boldsymbol{\theta}(T_1)\|_R^2 \leq 2\bar{W}(T_1) + \frac{2\varepsilon e^{5\varepsilon}}{q^2} \|\boldsymbol{\theta}(T_1)\|_R^2 = O(r^{\kappa^*} + r^{\kappa_R}),$$

where $\kappa_R := \min_{j \in R} \kappa_j$. Thus, $\lambda e^{4\varepsilon} \bar{W}(t)(t - T_1) = O(-\ln(r)(r^{\kappa^*} + r^{\kappa_R}))$. So, for sufficiently small r , our assumption holds.

Therefore, we get the following bounds

$$U = O(r^{\kappa^*} + r^{\kappa_R}), \quad V = O(r^{\kappa^*} + r^{\kappa_R}), \quad \forall j \in R \quad |u_j| = O(r^{\kappa_j}).$$

C.2. Finding a limit of Equation (7)

Now, we want to find $\lim_{r \rightarrow 0} \boldsymbol{\theta}^X(T_2^\varepsilon)$ to couple $\boldsymbol{\theta}(T_2^\varepsilon)$ with the vector that does not depend on r . To do this, we want to shift time $t \mapsto t - T_2^\varepsilon$ and consider Equation (7) for different r .

$$\begin{aligned} \frac{du_j^{X,r}}{dt} &= \frac{1}{n} \sum_{i=1}^n (-\ell'(f(\boldsymbol{\theta}^{X,r}, \mathbf{x}_i)y_i)) \phi(\mathbf{v}_j^{X,r}, \mathbf{x}_i)y_i, \\ \frac{d\mathbf{v}_j^{X,r}}{dt} &= \frac{1}{n} \sum_{i=1}^n (-\ell'(f(\boldsymbol{\theta}^{X,r}, \mathbf{x}_i)y_i)) u_j^{X,r} \nabla_{\mathbf{v}} \phi(\mathbf{v}_j^{X,r}, \mathbf{x}_i)y_i, \end{aligned} \quad (8)$$

where

$$u_j^{X,r}(-t_4^{\varepsilon,r}) = \begin{cases} ru_j^*, & j \in P, \\ 0, & j \in R, \end{cases} \quad \mathbf{v}_j^{X,r}(-t_4^{\varepsilon,r}) = \begin{cases} r|u_j^*| \hat{\mathbf{v}}_j^*, & j \in P, \\ 0, & j \in R, \end{cases}$$

$$\text{and } t_4^{\varepsilon,r} := \frac{1}{2\lambda} \ln \left(\frac{\lambda \varepsilon}{2ar^2 \|\boldsymbol{\theta}^*\|_{[m]}^2} \right).$$

Similarly to the previous subsections, we get

$$\forall t \in [-t_4^{\varepsilon,r}, 0] \quad \|\boldsymbol{\theta}^{X,r}\|_{[m]}^2 \leq 2r^2 \|\boldsymbol{\theta}^*\|_{[m]}^2 e^{2\lambda(t+t_4^{\varepsilon,r})},$$

and

$$\alpha^r(t) := \arccos((\hat{\mathbf{v}}_j^*)^\top \hat{\mathbf{v}}_j^{X,r}) \leq \frac{a}{\lambda} r^2 \|\boldsymbol{\theta}^*\|_{[m]}^2 (e^{2\lambda(t+t_4^{\varepsilon,r})} - 1).$$

Therefore, for $r' \geq r$

$$\alpha^{r'}(-t_4^{\varepsilon,r'}) \leq \frac{a}{\lambda} r'^2 \|\boldsymbol{\theta}^*\|_{[m]}^2 (e^{2\lambda(-t_4^{\varepsilon,r'}+t_4^{\varepsilon,r'})} - 1) = \frac{a}{\lambda} \|\boldsymbol{\theta}^*\|_{[m]}^2 ((r')^2 - r^2) = O((r')^2).$$

For $u_j^{X,r}$, we get

$$\begin{aligned} \frac{d(s_j u_j^{X,r} e^{-\lambda(t+t_4^{\varepsilon,r})})}{dt} &= |u_j^{X,r}| e^{-\lambda(t+t_4^{\varepsilon,r})} \left(\lambda(\cos(\alpha^r) - 1) + \frac{1}{n} \sum_{i=1}^n (-\ell'(f(\boldsymbol{\theta}^{X,r}, \mathbf{x}_i)y_i) + \ell'(0)) |u_j^X| \phi(\hat{\mathbf{v}}_j^X, \mathbf{x}_i)y_i \right) \\ &\leq a \|\boldsymbol{\theta}^{X,r}\|_{[m]}^2 |u_j^{X,r}| e^{-\lambda(t+t_4^{\varepsilon,r})} \\ \implies \left| \frac{u_j^{X,r} e^{-\lambda(t+t_4^{\varepsilon,r})}}{ru^*} \right| &\leq \exp \left(a \int_{-t_4^{\varepsilon,r}}^t \|\boldsymbol{\theta}^{X,r}\|_{[m]}^2 \right) \leq \exp \left(\frac{a}{\lambda} r^2 \|\boldsymbol{\theta}^*\|_{[m]}^2 (e^{2\lambda(t+t_4^{\varepsilon,r})} - 1) \right). \end{aligned}$$

Notice that since $\varepsilon \leq 1/2$, we get $\alpha^r/4 \leq \varepsilon/4 \leq 1$. Thus,

$$1 - \cos(\alpha^r) \leq \frac{(\alpha^r)^2}{2} \leq 2\alpha^r.$$

It gives

$$\begin{aligned} \frac{d(s_j u_j^{X,r} e^{-\lambda(t+t_4^{\varepsilon,r})})}{dt} &\geq |u_j^{X,r}| e^{-\lambda(t+t_4^{\varepsilon,r})} (\lambda(\cos(\alpha^r) - 1) - a \|\boldsymbol{\theta}^{X,r}\|_{[m]}^2) \geq |u_j^{X,r}| e^{-\lambda(t+t_4^{\varepsilon,r})} (-2\lambda\alpha^r - a \|\boldsymbol{\theta}^{X,r}\|_{[m]}^2) \\ \implies \left| \frac{u_j^{X,r} e^{-\lambda(t+t_4^{\varepsilon,r})}}{r u^*} \right| &\geq \exp\left(\frac{-2a}{\lambda} r^2 \|\boldsymbol{\theta}^*\|_{[m]}^2 (e^{2\lambda(t+t_4^{\varepsilon,r})} - 1)\right). \end{aligned}$$

Therefore,

$$\left| \ln\left(\frac{u_j^{X,r}(-t_4^{\varepsilon,r'})}{u_j^{X,r'}(-t_4^{\varepsilon,r'})}\right) \right| \leq \frac{2a}{\lambda} \|\boldsymbol{\theta}^*\|_{[m]}^2 ((r')^2 - r^2) \implies |u_j^{X,r}(-t_4^{\varepsilon,r'}) - u_j^{X,r'}(-t_4^{\varepsilon,r'})| \leq O((r')^3).$$

Notice that the above derivations also imply that

$$u_j^{X,r}(0) = \Theta(r u_j^* e^{\lambda t_4^{\varepsilon,r}}) = \Theta(\sqrt{\varepsilon}).$$

Thus, we can use the results of previous subsection with $\kappa = 2$ and get that

$$\|\hat{\mathbf{v}}_j^{X,r}(0) - \hat{\mathbf{v}}_j^{X,r'}(0)\| = O((r')^2), |u_j^{X,r}(0) - u_j^{X,r'}(0)| = O((r')^2).$$

This property ensures that the sequences $\hat{\mathbf{v}}_j^{X,r}(0)$ and $u_j^{X,r}(0)$ are fundamental. Since $u_j^{X,r}(0)$ is also bounded, $\exists \lim_{r \rightarrow 0} \boldsymbol{\theta}^{X,r}(0) = \boldsymbol{\theta}^{X,\varepsilon,*}$. Since $\boldsymbol{\theta}^{X,r}(0)$ is an image of embedding (4), we could find $\boldsymbol{\theta}^{\varepsilon,*}$ such that $\boldsymbol{\theta}^{X,\varepsilon,*} = \chi(\boldsymbol{\theta}^{\varepsilon,*})$. Also, since we have $\forall j \in P \|\hat{\mathbf{v}}_j^{X,\varepsilon,r}(0) - \hat{\mathbf{v}}^*\| \leq \varepsilon$, this property will also hold for the limit parameters.

C.3. Final Bound

Finally, using the results of both subsections, we have

$$\begin{aligned} \forall j \in P \quad |u_j(T_2^\varepsilon) - u_j^{X,\varepsilon,*}| &= O(r^2 + r^{\kappa^*} + r^{\kappa_R}), \\ \forall j \in P \quad \|\hat{\mathbf{v}}_j(T_2^\varepsilon) - \hat{\mathbf{v}}_j^{X,\varepsilon,*}\| &= O(r^2 + r^{\kappa^*} + r^{\kappa_R}), \\ \forall j \in R \quad |u_j(T_2^\varepsilon)| &= O(r^{\kappa_j}). \end{aligned}$$

D. Proofs for Section 5

D.1. Proof of Lemma 5.2

We only proof that \mathbf{e}_1 and $-\mathbf{e}_1$ are the maxima of function G . The proof that \mathbf{e}_2 and $-\mathbf{e}_2$ are the minima is similar.

First, we have

$$G(\mathbf{e}_1) = \frac{-\ell'(0)}{n} \left(\sum_{i \in S_1} \phi(\mathbf{e}_1, \mathbf{x}_i) + \sum_{i \in S_3} \phi(\mathbf{e}_1, \mathbf{x}_i) - \sum_{i \in S_2} \phi(\mathbf{e}_1, \mathbf{x}_i) - \sum_{i \in S_4} \phi(\mathbf{e}_1, \mathbf{x}_i) \right).$$

Notice

$$\begin{aligned} \sum_{i \in S_1} \phi(\mathbf{e}_1, \mathbf{x}_i) + \sum_{i \in S_3} \phi(\mathbf{e}_1, \mathbf{x}_i) &= \sum_{i \in S_1} \mathbf{x}_i^1 \geq \frac{n}{4}(1 - \delta), \\ \sum_{i \in S_2} \phi(\mathbf{e}_1, \mathbf{x}_i) + \sum_{i \in S_4} \phi(\mathbf{e}_1, \mathbf{x}_i) &\leq \sum_{i \in S_2} (\mathbf{x}_i^1 + \xi)_+ + \sum_{i \in S_4} (\mathbf{x}_i^1 + \xi)_+ \leq \frac{n}{2}(\delta + \xi). \end{aligned}$$

Thus,

$$G(\mathbf{e}_1) \geq \frac{-\ell'(0)}{4} (1 - 3\delta - 2\xi).$$

Let $\hat{\mathbf{v}}^*$ be a maximum direction of G . W.l.o.g., assume that $\mathbf{e}_1^\top \hat{\mathbf{v}}^* \geq 0$. Then we have

$$\begin{aligned} \sum_{i \in S_1} \phi(\hat{\mathbf{v}}^*, \mathbf{x}_i) + \sum_{i \in S_3} \phi(\hat{\mathbf{v}}^*, \mathbf{x}_i) &\leq \frac{n}{4} (\mathbf{e}_1^\top \hat{\mathbf{v}}^* + 2\xi + 2\delta), \\ \sum_{i \in S_2} \phi(\hat{\mathbf{v}}^*, \mathbf{x}_i) + \sum_{i \in S_4} \phi(\hat{\mathbf{v}}^*, \mathbf{x}_i) &\geq 0. \end{aligned}$$

Therefore,

$$G(\hat{\mathbf{v}}^*) \leq \frac{-\ell'(0)}{4} (\mathbf{e}_1^\top \hat{\mathbf{v}}^* + 2\xi + 2\delta).$$

These properties imply

$$\mathbf{e}_1^\top \hat{\mathbf{v}}^* \geq 1 - 5\delta - 4\xi > \delta + \xi.$$

Hence, $\hat{\mathbf{v}}^*$ should satisfy

$$\begin{aligned} \forall i \in S_3 \quad (\hat{\mathbf{v}}^*)^\top \mathbf{x}_i \leq -\xi &\implies \phi(\hat{\mathbf{v}}^*, \mathbf{x}_i) = 0, \\ \forall i \in S_1 \quad (\hat{\mathbf{v}}^*)^\top \mathbf{x}_i \geq \xi &\implies \phi(\hat{\mathbf{v}}^*, \mathbf{x}_i) = (\hat{\mathbf{v}}^*)^\top \mathbf{x}_i. \end{aligned}$$

Thus,

$$\sum_{i \in S_3} \phi(\hat{\mathbf{v}}^*, \mathbf{x}_i) + \sum_{i \in S_1} \phi(\hat{\mathbf{v}}^*, \mathbf{x}_i) = (\hat{\mathbf{v}}^*)^\top \left(\sum_{i \in S_1} \mathbf{x}_i \right).$$

Similarly,

$$\sum_{i \in S_3} \phi(\mathbf{e}_1, \mathbf{x}_i) + \sum_{i \in S_1} \phi(\mathbf{e}_1, \mathbf{x}_i) = \mathbf{e}_1^\top \left(\sum_{i \in S_1} \mathbf{x}_i \right) \geq \frac{n}{4} (1 - \delta).$$

Notice that, due to symmetry $\mathbf{R}_2 \mathbf{R}_r D_{\mathbf{x}} = D_{\mathbf{x}}$, $\sum_{i \in S_1} \mathbf{x}_i$ is the multiple of \mathbf{e}_1 . Thus,

$$(\hat{\mathbf{v}}^*)^\top \left(\sum_{i \in S_1} \mathbf{x}_i \right) = \hat{\mathbf{v}}^{*,1} \mathbf{e}_1^\top \left(\sum_{i \in S_1} \mathbf{x}_i \right).$$

Now, consider the points from S_2 and S_4 . Define \mathbf{x}_i^r by the following equality $(x_i^1, (\mathbf{x}_i^r)^\top)^\top = \mathbf{x}_i$. Since $\mathbf{R}_2 \mathbf{R}_r D_{\mathbf{x}} = D_{\mathbf{x}}$, all points $i \in S_2$ have corresponding point $i^- \in S_4$ such that $\mathbf{x}_{i^-} = \mathbf{R}_2 \mathbf{R}_r \mathbf{x}_i$. Since $\mathbf{R}_1 D_{\mathbf{x}} = D_{\mathbf{x}}$, all $i \in S_2$ have a corresponding point $i^+ \in S_2$ such that $\mathbf{x}_{i^+} = \mathbf{R}_1 \mathbf{x}_i$. Thus,

$$\begin{aligned} \sum_{i \in S_2} \phi(\hat{\mathbf{v}}, \mathbf{x}_i) + \sum_{i \in S_4} \phi(\hat{\mathbf{v}}, \mathbf{x}_i) &= \sum_{i \in S_2} \phi(\hat{\mathbf{v}}, \mathbf{x}_i) + \phi(\hat{\mathbf{v}}, \mathbf{x}_{i^-}) \\ &= \sum_{\substack{i \in S_2 \\ x_i^1 > 0}} \phi(\hat{\mathbf{v}}, \mathbf{x}_i) + \phi(\hat{\mathbf{v}}, \mathbf{x}_{i^-}) + \phi(\hat{\mathbf{v}}, \mathbf{x}_{i^+}) + \phi(\hat{\mathbf{v}}, \mathbf{x}_{(i^+)^-}) + \sum_{\substack{i \in S_2 \\ x_i^1 = 0}} \phi(\hat{\mathbf{v}}, \mathbf{x}_i) + \phi(\hat{\mathbf{v}}, \mathbf{x}_{i^-}). \end{aligned}$$

Now,

$$\phi(\hat{\mathbf{v}}, \mathbf{x}_i) + \phi(\hat{\mathbf{v}}, \mathbf{x}_{i^-}) = \int_{\mathbf{z} \in \mathbb{D}^d} ((\hat{\mathbf{v}})^\top (\mathbf{x}_i + \xi \mathbf{z}))_+ + ((\hat{\mathbf{v}})^\top \mathbf{R}_2 \mathbf{R}_r (\mathbf{x}_i + \xi \mathbf{z}))_+ Q(d\mathbf{z}).$$

Denote $\mathbf{x} := \mathbf{x}_i + \xi \mathbf{z}$ and consider integrand

$$((\hat{\mathbf{v}})^\top \mathbf{x})_+ + ((\hat{\mathbf{v}})^\top \mathbf{R}_2 \mathbf{R}_r \mathbf{x})_+ = (\hat{\mathbf{v}}^1 x^1 + (\hat{\mathbf{v}}^r)^\top \mathbf{x}^r)_+ + (\hat{\mathbf{v}}^1 x^1 - (\hat{\mathbf{v}}^r)^\top \mathbf{x}^r)_+ \geq (\hat{\mathbf{v}}^1 x^1)_+.$$

Therefore,

$$\phi(\hat{\mathbf{v}}, \mathbf{x}_i) + \phi(\hat{\mathbf{v}}, \mathbf{x}_{i^-}) \geq \int_{\mathbf{z} \in \mathbb{D}^d} (\hat{\mathbf{v}}^1 (x_i^1 + \xi z^1))_+ Q(d\mathbf{z}).$$

So,

$$\sum_{i \in S_2} \phi(\hat{\mathbf{v}}, \mathbf{x}_i) + \sum_{i \in S_4} \phi(\hat{\mathbf{v}}, \mathbf{x}_i) \geq \sum_{\substack{i \in S_2 \\ x_i^1 > 0}} \int_{\mathbf{z} \in \mathbb{D}^d} |\hat{\mathbf{v}}^1| |x_i^1 + \xi z^1| Q(d\mathbf{z}) + \sum_{\substack{i \in S_2 \\ x_i^1 = 0}} \int_{\mathbf{z} \in \mathbb{D}^d} |\hat{\mathbf{v}}^1| (\xi z^1)_+ Q(d\mathbf{z}).$$

Similarly,

$$\sum_{i \in S_2} \phi(\mathbf{e}_1, \mathbf{x}_i) + \sum_{i \in S_4} \phi(\mathbf{e}_1, \mathbf{x}_i) = \sum_{\substack{i \in S_2 \\ x_i^1 > 0}} \int_{\mathbf{z} \in \mathbb{D}^d} |x_i^1 + \xi z^1| Q(d\mathbf{z}) + \sum_{\substack{i \in S_2 \\ x_i^1 = 0}} \int_{\mathbf{z} \in \mathbb{D}^d} (\xi z^1)_+ Q(d\mathbf{z}) \leq \frac{n}{2}(\delta + \xi).$$

Thus,

$$\begin{aligned} \sum_{i \in S_2} \phi(\hat{\mathbf{v}}, \mathbf{x}_i) + \sum_{i \in S_4} \phi(\hat{\mathbf{v}}, \mathbf{x}_i) - \sum_{i \in S_2} \phi(\mathbf{e}_1, \mathbf{x}_i) - \sum_{i \in S_4} \phi(\mathbf{e}_1, \mathbf{x}_i) \\ \geq (|\hat{\mathbf{v}}^1| - 1) \left(\sum_{\substack{i \in S_2 \\ x_i^1 > 0}} \int_{\mathbf{z} \in \mathbb{D}^d} |x_i^1 + \xi z^1| Q(d\mathbf{z}) + \sum_{\substack{i \in S_2 \\ x_i^1 = 0}} \int_{\mathbf{z} \in \mathbb{D}^d} (\xi z^1)_+ Q(d\mathbf{z}) \right) \\ \geq \frac{n}{2} (|\hat{\mathbf{v}}^1| - 1) (\delta + \xi). \end{aligned}$$

So, we have

$$0 \geq G(\mathbf{e}_1) - G(\hat{\mathbf{v}}^*) \geq \frac{(-\ell'(0))}{4} (1 - \hat{\mathbf{v}}^{*,1}) (1 - \delta) - \frac{(-\ell'(0))}{2} (1 - \hat{\mathbf{v}}^{*,1}) (\delta + \xi) \geq \frac{(-\ell'(0))}{4} (1 - \hat{\mathbf{v}}^{*,1}) (1 - 3\delta - 2\xi).$$

Therefore, $\hat{\mathbf{v}}^{*,1} = 1 \implies \hat{\mathbf{v}}^* = \mathbf{e}_1$.

D.2. Proof of Lemma 5.3

W.l.o.g., consider direction \mathbf{e}_1 . We will show that this direction is attractive for positive neurons ($s_j = 1$) at the beginning of training in the region $\mathbf{e}_1^\top \hat{\mathbf{v}}^* \geq \delta + \xi$. Thus, the direction \mathbf{e}_1 captures the random neuron with probability greater than $h := \frac{1}{2} \frac{\text{Vol}(A)}{\text{Vol}(S^{d-1})} \geq \frac{1}{2} \left(\frac{1}{2} - \frac{\text{Vol}(\mathbb{D}^{d-2})}{\text{Vol}(S^{d-1})} \arcsin(\delta + \xi) \right) = \frac{1}{4} (1 - O(\delta + \xi))$, where $A = \{\mathbf{x} \in S^{d-1} \mid \mathbf{e}_1^\top \mathbf{x} \geq \delta + \xi\}$. This bound implies the following probability of success

$$\Pr(\text{success}) \geq (1 - (1 - h)^m)^4 = 1 - 4(1 - h)^m - O((1 - h)^{2m}) = 1 - 4 \left(\frac{3}{4} \right)^m (1 + O(\delta + \xi)) - O\left(\left(\frac{9}{16} \right)^m \right).$$

Now, we will proof that A is attractive region for \mathbf{e}_1 . Consider positive neuron with $\hat{\mathbf{v}} \in S^{d-1}$ such that $\mathbf{e}_1^\top \hat{\mathbf{v}} \geq \delta + \xi$. Equation (5) gives

$$\frac{d\hat{\mathbf{v}}}{dt} = \mathbf{P}_{\hat{\mathbf{v}}} \mathbf{g}(\hat{\mathbf{v}}).$$

We have

$$\frac{d\|\hat{\mathbf{v}} - \mathbf{e}_1\|^2}{dt} = -2\mathbf{e}_1^\top \frac{d\hat{\mathbf{v}}}{dt}.$$

So, we only need to show that $\mathbf{e}_1^\top \frac{d\hat{\mathbf{v}}}{dt}$ is positive.

Denote $\hat{\mathbf{v}} =: \mathbf{e}_1 \cos(\alpha) + \boldsymbol{\epsilon} \sin(\alpha)$. Notice

$$\mathbf{P}_{\hat{\mathbf{v}}} \mathbf{e}_1 = \mathbf{e}_1 - \hat{\mathbf{v}} \cos(\alpha) = \mathbf{e}_1 \sin(\alpha)^2 - \boldsymbol{\epsilon} \sin(\alpha) \cos(\alpha).$$

Similarly to the proof of Lemma 5.2, we have

$$\mathbf{g}(\hat{\mathbf{v}}) = \frac{-\ell'(0)}{n} \left(\sum_{i \in S_1} \mathbf{x}_i - \sum_{i \in S_2} \nabla_{\mathbf{v}} (\phi(\mathbf{z}, \mathbf{x}_i) + \phi(\mathbf{z}, \mathbf{R}_2 \mathbf{R}_r \mathbf{x}_i)) \Big|_{\mathbf{z}=\hat{\mathbf{v}}} \right).$$

Notice

$$\begin{aligned} \nabla_{\mathbf{v}} (\phi(\mathbf{z}, \mathbf{x}_i) + \phi(\mathbf{z}, \mathbf{R}_2 \mathbf{R}_r \mathbf{x}_i)) \Big|_{\mathbf{z}=\hat{\mathbf{v}}} \\ = \int_{\mathbf{z} \in \mathbb{D}^d} ([\hat{\mathbf{v}}^\top (\mathbf{x}_i + \xi \mathbf{z}) \geq 0] (\mathbf{x}_i + \xi \mathbf{z}) + [\hat{\mathbf{v}}^\top \mathbf{R}_2 \mathbf{R}_r (\mathbf{x}_i + \xi \mathbf{z}) \geq 0] \mathbf{R}_2 \mathbf{R}_r (\mathbf{x}_i + \xi \mathbf{z})) Q(d\mathbf{z}). \end{aligned}$$

Consider the integrand. Denote $\mathbf{x} := \mathbf{x}_i + \xi \mathbf{z}$. We have four potential cases:

1. $\hat{\mathbf{v}}^\top \mathbf{x} \leq 0 \wedge \hat{\mathbf{v}}^\top \mathbf{R}_2 \mathbf{R}_r \mathbf{x} \leq 0$,
2. $\hat{\mathbf{v}}^\top \mathbf{x} > 0 \wedge \hat{\mathbf{v}}^\top \mathbf{R}_2 \mathbf{R}_r \mathbf{x} \leq 0$,
3. $\hat{\mathbf{v}}^\top \mathbf{x} \leq 0 \wedge \hat{\mathbf{v}}^\top \mathbf{R}_2 \mathbf{R}_r \mathbf{x} > 0$,
4. $\hat{\mathbf{v}}^\top \mathbf{x} > 0 \wedge \hat{\mathbf{v}}^\top \mathbf{R}_2 \mathbf{R}_r \mathbf{x} > 0$.

In the first case, the contribution of the integrand to expression $\mathbf{e}_1 \frac{d\hat{\mathbf{v}}}{dt}$ is zero. In the fourth case, the contribution of the integrand is

$$-(\mathbf{P}_{\hat{\mathbf{v}}} \mathbf{e}_1)^\top (\mathbf{x} + \mathbf{R}_2 \mathbf{R}_r \mathbf{x}) = -2x^1 \sin(\alpha)^2 \geq -2(\delta + \xi) \sin(\alpha)^2.$$

The second and third cases are symmetric. So, we consider only the second case. In the second case, we have

$$\hat{\mathbf{v}}^\top \mathbf{x} = x^1 \cos(\alpha) + \boldsymbol{\epsilon}^\top \mathbf{x}^r \sin(\alpha) \geq 0.$$

Thus, $\boldsymbol{\epsilon}^\top \mathbf{x}^r \geq 0$ (otherwise $\hat{\mathbf{v}}^\top \mathbf{R}_2 \mathbf{R}_r \mathbf{x} \geq \hat{\mathbf{v}}^\top \mathbf{x} > 0$). It implies

$$-(\mathbf{P}_{\hat{\mathbf{v}}} \mathbf{e}_1)^\top \mathbf{x} = -x^1 \sin(\alpha)^2 + \boldsymbol{\epsilon}^\top \mathbf{x} \cos(\alpha) \sin(\alpha) \geq -x^1 \sin(\alpha)^2.$$

Therefore, we get

$$\mathbf{e}_1 \frac{d\hat{\mathbf{v}}}{dt} \geq \frac{-\ell'(0)}{n} \left(\frac{n}{4}(1 - \delta) - \frac{n}{4}(2(\delta + \xi)) \right) \sin(\alpha)^2 \geq \frac{-\ell'(0)}{8} \sin(\alpha)^2.$$

It gives

$$\frac{d\alpha}{dt} \leq \frac{-\ell'(0)}{8} \sin(\alpha) \leq \frac{-\ell'(0)}{4\pi} \alpha,$$

which implies exponentially fast convergence to \mathbf{e}_1 .

D.3. Proof of Proposition 5.7

First, notice that, due to symmetry $\mathbf{R}_r D_{\mathbf{x}} = D_{\mathbf{x}}$, the features of 4-neuron network belong to a sub-space generated by vectors \mathbf{e}_1 and \mathbf{e}_2 and clustered around these directions.

Lemma D.1. $\mathbf{v}_j^{*,\varepsilon}$ form Theorem 4.4 belong to a sub-space generated by vectors \mathbf{e}_1 and \mathbf{e}_2 ($\mathbf{R}_r \mathbf{v}_j^{*,\varepsilon} = \mathbf{v}_j^{*,\varepsilon}$). Moreover, this symmetry is preserved under gradient flow dynamics (1).

Proof. Consider Equation (7). Notice that at the start of function $h(\mathbf{x}) := f(\boldsymbol{\theta}^\times(T_1), \mathbf{x})$ depends only on x^1 and x^2 . Due to symmetry $\mathbf{R}_r D_{\mathbf{x}} = D_{\mathbf{x}}$, this property will hold for the right-hand part of Equation (7). Therefore, features will not be able to escape the sub-space generated by \mathbf{e}_1 and \mathbf{e}_2 . Thus, limit features $\mathbf{v}_j^{*,\varepsilon}$ will also lie in the sub-space generated by \mathbf{e}_1 and \mathbf{e}_2 . Since the limit and dataset respect symmetry \mathbf{R}_r , it will propagate to further stages of training. \square

Consider dynamics

$$\begin{aligned} \frac{du_k^e}{dt} &= \sum_{i=1}^n (-\ell'(f(\boldsymbol{\theta}^e, \mathbf{x}_i) y_i)) \phi(\mathbf{v}_k^e, \mathbf{x}_i) y_i, \\ \frac{d\mathbf{v}_k^e}{dt} &= \sum_{i=1}^n (-\ell'(f(\boldsymbol{\theta}^e, \mathbf{x}_i) y_i)) u_k^e \nabla_{\mathbf{v}} \phi(\mathbf{v}_k^e, \mathbf{x}_i) y_i, \end{aligned}$$

where $u_k^e(0) = s_{P_k} \sqrt{\sum_{j \in P_k} (u_j^{\varepsilon,*})^2}$ and $\mathbf{v}_k^e(0) = |u_k^e(0)| \hat{\mathbf{v}}_{P_k}^{\varepsilon,*}$. Notice that, by Lemma D.1, $\mathbf{v}_k^e \in \text{span}\{\mathbf{e}_1, \mathbf{e}_2\}$.

Proposition D.2. Denote a signed angle between $\mathbf{R}_1^{\lfloor \frac{k-1}{2} \rfloor} \mathbf{P}^{a(k)} \mathbf{v}_k^e$ and $\mathbf{e}_{a(k)+1}$ by α_k . Then we have $|\sin(\alpha_k)| \leq \delta + \xi$.

Proof. W.l.o.g., consider \mathbf{v}_1^e and assume $\alpha_1 \geq 0$. We have

$$\frac{d\hat{\mathbf{v}}_1^e}{dt} = \sum_{i=1}^n (-\ell'(f(\boldsymbol{\theta}^e, \mathbf{x}_i) y_i)) \mathbf{P}_{\mathbf{v}_1^e} \nabla_{\mathbf{v}} \phi(\mathbf{v}_1^e, \mathbf{x}_i) y_i.$$

Notice

$$\mathbf{e}_1^\top \mathbf{P}_{\mathbf{v}_1^e} \nabla_{\mathbf{v}} \phi(\mathbf{v}_1^e, \mathbf{x}_i) y_i = \int_{\mathbf{z} \in \mathbb{D}^d} [(\mathbf{v}_1^e)^\top (\mathbf{x}_i + \xi \mathbf{z}) \geq 0] \mathbf{e}_1^\top \mathbf{P}_{\mathbf{v}_1^e} (\mathbf{x}_i + \xi \mathbf{z}) y_i Q(d\mathbf{z}).$$

Consider the integrand and denote $\mathbf{z}_i := \mathbf{x}_i + \xi \mathbf{z}$. Assuming that $\alpha_1 \leq \pi/4$, we get

$$\begin{aligned} \forall i \in S_1 [(\mathbf{v}_1^e)^\top \mathbf{z}_i \geq 0] \mathbf{e}_1^\top \mathbf{P}_{\mathbf{v}_1^e} \mathbf{z}_i y_i &= (\mathbf{e}_1 \sin(\alpha_1) - \mathbf{e}_2 \cos(\alpha_1))^\top \sin(\alpha_1) (\mathbf{e}_1 + \mathbf{z}_i - \mathbf{e}_1) \geq (\sin(\alpha_1) - (\delta + \xi)) \sin(\alpha_1), \\ \forall i \in S_2 [(\mathbf{v}_1^e)^\top \mathbf{z}_i \geq 0] \mathbf{e}_1^\top \mathbf{P}_{\mathbf{v}_1^e} \mathbf{z}_i y_i &= -[(\mathbf{v}_1^e)^\top \mathbf{z}_i \geq 0] (\mathbf{e}_1 \sin(\alpha_1) - \mathbf{e}_2 \cos(\alpha_1))^\top \sin(\alpha_1) (\mathbf{e}_2 + \mathbf{z}_i - \mathbf{e}_2) \\ &\geq (\cos(\alpha_1) - (\delta + \xi)) \sin(\alpha_1) \\ &\geq 0, \\ \forall i \in S_3 [(\mathbf{v}_1^e)^\top \mathbf{z}_i \geq 0] \mathbf{e}_1^\top \mathbf{P}_{\mathbf{v}_1^e} \mathbf{z}_i y_i &= 0, \\ \forall i \in S_4 [(\mathbf{v}_1^e)^\top \mathbf{z}_i \geq 0] \mathbf{e}_1^\top \mathbf{P}_{\mathbf{v}_1^e} \mathbf{z}_i y_i &\geq -[\sin(\alpha_1) \leq \delta + \xi] (1 + \delta + \xi). \end{aligned}$$

Thus, the function $h(\alpha_1) = \max(\alpha_1, \arcsin(\delta + \xi))$ is always decreasing (when α_1 is less than $\arcsin(\delta + \xi)$), the time derivative of $h(\alpha_1)$ is zero, when α_1 is greater than $\arcsin(\delta + \xi)$, the time derivative of $h(\alpha_1)$ is determined only by the points from S_1 and S_2 and hence negative). Therefore, α_1 never exceeds $\arcsin(\delta + \xi)$. \square

Proof of Proposition 5.7. Denote $a := 0.01$, $b := 0.001$, $c = 1000$, $q := 8/3$, $T := \inf\{t > 0 \mid \max_k |u_k^e(t)| \geq q\}$.

From the previous proposition, we know that $|\sin(\alpha_k)| \leq \delta + \xi \leq a$. Thus, for $t \in [0, T]$ and $i \in S_1$, we get

$$f(\boldsymbol{\theta}^e, x_i) = (u_1^e)^2 (\hat{\mathbf{v}}_1^e)^\top \mathbf{x}_i - (u_2^e)^2 \phi(\hat{\mathbf{v}}_2^e, \mathbf{x}_i) - (u_4^e)^2 \phi(\hat{\mathbf{v}}_4^e, \mathbf{x}_i).$$

Therefore,

$$\begin{aligned} f(\boldsymbol{\theta}^e, x_i) &\leq (u_1^e)^2 (1 + \delta) \leq (u_1^e)^2 (1 + a), \\ f(\boldsymbol{\theta}^e, x_i) &\geq (u_1^e)^2 (\cos(\alpha_k) - \delta) - q^2 (\phi(\hat{\mathbf{v}}_2^e, \mathbf{x}_i) + \phi(\hat{\mathbf{v}}_4^e, \mathbf{x}_i)) \geq (u_1^e)^2 (1 - (\delta + \xi)^2 - \delta) - q^2 (4\delta + 2\xi) \\ &\geq (u_1^e)^2 (1 - 2a) - 4q^2 a, \end{aligned}$$

where the last inequality follows from the following property:

$$\begin{aligned} \forall j \in \{2, 4\} \phi(\hat{\mathbf{v}}_j^e, \mathbf{x}_i) &= \int_{-1}^1 ((\hat{\mathbf{v}}_j^e)^\top \mathbf{x}_i + \xi z)_+ (1 - z^2)^{\frac{d-1}{2}} \frac{\text{Vol}(\mathbb{D}^{d-1})}{\text{Vol}(\mathbb{D}^d)} dz \\ &\leq \int_{-1}^1 (|\sin(\alpha_j)| + \delta + \xi z)_+ (1 - z^2)^{\frac{d-1}{2}} \frac{\text{Vol}(\mathbb{D}^{d-1})}{\text{Vol}(\mathbb{D}^d)} dz \\ &\leq \int_{-1}^1 (2\delta + \xi + \xi z) (1 - z^2)^{\frac{d-1}{2}} \frac{\text{Vol}(\mathbb{D}^{d-1})}{\text{Vol}(\mathbb{D}^d)} dz \\ &= 2\delta + \xi. \end{aligned}$$

Similarly, we could derive the same inequalities for S_2, \dots, S_4 .

These inequalities imply

$$\begin{aligned} \frac{du_1^e}{dt} &= \frac{1}{n} \sum_{i=1}^n (-\ell'(f(\boldsymbol{\theta}^e, \mathbf{x}_i) y_i)) \phi(\mathbf{v}_1^e, \mathbf{x}_i) y_i \\ &\geq \frac{u_1^e}{4} \left(\frac{1 - 2a}{1 + \exp((u_1^e)^2 (1 + a))} - \frac{2a}{1 + \exp((u_2^e)^2 (1 - 2a) - 4q^2 a)} - \frac{2a}{1 + \exp((u_4^e)^2 (1 - 2a) - 4q^2 a)} \right) \end{aligned}$$

and

$$\frac{du_1^e}{dt} \leq \frac{u_1^e (1 + a)}{4(1 + \exp((u_1^e)^2 (1 - 2a) - 4q^2 a))}.$$

Denote $x := \min_k |u_k^e|$ and $y := \max_k |u_k^e|$. The property above implies

$$\begin{aligned} \frac{dx}{dt} &\geq \frac{x}{4} \left(\frac{1-2a}{1 + \exp(x^2(1+a))} - \frac{4a}{1 + \exp(x^2(1-2a) - 4q^2a)} \right) \geq \frac{x((1-2a)e^{-7q^2a} - 4a)}{4(1 + \exp(x^2(1-2a) - 4q^2a))}, \\ \frac{dy}{dt} &\leq \frac{y(1+a)}{4(1 + \exp(y^2(1-2a) - 4q^2a))}. \end{aligned}$$

Now, denote $X := x^2(1-2a)$, $Y := y^2(1-2a)$, $A := \frac{(1-2a)e^{-7q^2a} - 4a}{2}$, $B := \frac{1+a}{2}$, and $C := e^{-4q^2a}$. We get

$$\begin{aligned} \frac{dX}{dt} &\geq \frac{AX}{1 + Ce^X} \implies \int_{X(0)}^{X(t)} \frac{dX(1 + Ce^X)}{X} \geq At, \\ \frac{dY}{dt} &\leq \frac{BY}{1 + Ce^Y} \implies \int_{Y(0)}^{Y(t)} \frac{dY(1 + Ce^Y)}{Y} \leq Bt. \end{aligned}$$

Denote $h(x) := \int_1^x \frac{dz(1+Ce^z)}{z}$. We get

$$h(X(t)) - h(X(0)) \geq At, \quad h(Y(t)) - h(Y(0)) \leq Bt.$$

We have two possible cases: $T < \infty$ and $T = \infty$. In the second case, we notice that our lower bound for the derivative of x holds on the whole timeline. In this case, $\lim_{t \rightarrow \infty} x(t) = \infty$. Thus, at some time point we would get $x > q$, which contradicts definition of T . Therefore, T is finite.

Now, notice that $Y(T) = (1-2a)q^2$ and we have

$$\begin{aligned} B(h(X(T)) - h(X(0))) &\geq ABT \geq A(h(Y(T)) - h(Y(0))) \\ \implies h(X(T)) &\geq \frac{A}{B}h(Y(T)) + \frac{B-A}{B}h(Y(0)) - h(Y(0)) + h(X(0)). \end{aligned}$$

By the definition of h ,

$$\begin{aligned} h(Y(T)) &= \ln(Y(T)) + C \int_1^{Y(T)} \frac{e^z}{z} dz, \\ h(Y(0)) &= - \int_{Y(0)}^1 \frac{1 + Ce^z}{z} dz \geq (1 + Ce) \ln(Y(0)) = (1 + Ce) \ln((1-2a)b^2), \\ h(Y(0)) - h(X(0)) &= \int_{X(0)}^{Y(0)} \frac{1 + Ce^z}{z} dz \leq (1 + Ce^{Y(0)}) \ln\left(\frac{Y(0)}{X(0)}\right) \leq (1 + Ce^{\frac{1-2a}{4}}) \ln(c^2). \end{aligned}$$

Using numerical integration, we get $h(X(T)) \geq 31.52$, which implies that $x(T) \geq 9/4$.

Therefore, at time T , we have

$$f(\boldsymbol{\theta}^e(T), \mathbf{x}_i) y_i \geq (1-2a)x(T)^2 - 4q^2a > 4.67 > 0,$$

the network classifies all points correctly. □

D.4. Proof of Proposition 5.8

First, consider direction $\boldsymbol{\theta}^{mm}$ and w.l.o.g. assume that $\|\boldsymbol{\theta}^{mm}\|^2 = 8$, which implies $\forall k |u_k^{mm}| = 1$. Consider some orbit, M , of data points under the group generated by $P, \mathbf{R}_1, \mathbf{R}_2, \mathbf{R}_r$.

Choose $i^* \in M \cap S_1 : e_2^\top \mathbf{x}_{i^*} \geq 0$. For this point, we have

$$f(\boldsymbol{\theta}^{mm}, \mathbf{x}_{i^*}) = \sum_{k=1}^4 (u_k^{mm})^2 \phi(\hat{\mathbf{v}}_k^{mm}, \mathbf{x}_{i^*}) = \phi(\hat{\mathbf{v}}_1^{mm}, \mathbf{x}_{i^*}) - \phi(\hat{\mathbf{v}}_2^{mm}, \mathbf{x}_{i^*}).$$

Since $\forall a |e_a^\top \mathbf{x}_{i^*}| > \xi$, we get

$$f(\boldsymbol{\theta}^{mm}, \mathbf{x}_{i^*}) = x_{i^*}^1 - x_{i^*}^2.$$

Due to symmetry of the network, all points in M will have the same margin.

Now, consider a small perturbation, $\boldsymbol{\theta}$, of $\chi(\boldsymbol{\theta}^{mm})$, which does not change activation patterns for points in M . Notice

$$\min_{i \in M} \frac{f(\boldsymbol{\theta}, \mathbf{x}_i) y_i}{\|\boldsymbol{\theta}\|^2} \leq \frac{1}{16\|\boldsymbol{\theta}\|^2} \sum_{i \in M} f(\boldsymbol{\theta}, \mathbf{x}_i) y_i.$$

However,

$$\begin{aligned} \sum_{i \in M} f(\boldsymbol{\theta}, \mathbf{x}_i) y_i &= \sum_{k=1}^4 \sum_{j \in P_k} \sum_{l=1}^4 \sum_{i \in S_l} u_j \phi(\mathbf{v}_j, \mathbf{x}_i) y_i = \sum_k \sum_{j \in P_k} 4|u_j| \|\mathbf{v}_j^{a(k)+1}\| (x_{i^*}^1 - x_{i^*}^2) \\ &\leq \sum_k \sum_{j \in P_k} 2(|u_j|^2 + \|\mathbf{v}_j\|^2) (x_{i^*}^1 - x_{i^*}^2). \end{aligned}$$

Thus,

$$\min_{i \in M} \frac{f(\boldsymbol{\theta}, \mathbf{x}_i) y_i}{\|\boldsymbol{\theta}\|^2} \leq \frac{x_{i^*}^1 - x_{i^*}^2}{8} = \min_{i \in M} \frac{f(\boldsymbol{\theta}^{mm}, \mathbf{x}_i) y_i}{\|\boldsymbol{\theta}^{mm}\|^2}.$$

It implies that, for sufficiently small perturbation,

$$\min_i \frac{f(\boldsymbol{\theta}, \mathbf{x}_i) y_i}{\|\boldsymbol{\theta}\|^2} = \min_M \min_{i \in M} \frac{f(\boldsymbol{\theta}, \mathbf{x}_i) y_i}{\|\boldsymbol{\theta}\|^2} \leq \min_i \frac{f(\boldsymbol{\theta}^{mm}, \mathbf{x}_i) y_i}{\|\boldsymbol{\theta}^{mm}\|^2},$$

where the equality is achieved only if $|u_j| = \|\mathbf{v}_j\|$ and \mathbf{v}_j is aligned with $b(k)e_{a(k)+1}$. However, in this case,

$$f(\boldsymbol{\theta}, \mathbf{x}_{i^*}) = U_1^2 x_{i^*}^1 - U_2^2 x_{i^*}^2,$$

where $U_k = s_{P_k} \sqrt{\sum_{j \in P_k} u_j^2}$. W.l.o.g., we could assume that $U_1^2 = \min_k U_k^2$, then, we get

$$\frac{f(\boldsymbol{\theta}, \mathbf{x}_{i^*})}{\|\boldsymbol{\theta}\|^2} \leq \frac{U_1^2 x_{i^*}^1 - U_1^2 x_{i^*}^2}{2 \sum_k U_k^2}.$$

But this margin should be equal to $\frac{x_{i^*}^1 - x_{i^*}^2}{8}$. Thus, $4U_1^2 \geq \sum_k U_k^2$. However, this implies that $U_1^2 = U_2^2 = U_3^2 = U_4^2$, i.e., $\boldsymbol{\theta}$ is proportional to $\chi(\boldsymbol{\theta}^{mm})$. Therefore, $\chi(\boldsymbol{\theta}^{mm})$ is indeed local extremum. Similarly, we can show that $\boldsymbol{\theta}^{mm}$ is an isolated local extremum.

D.5. Proof of Lemma 5.5

Denote 4-neuron network initialized at $\boldsymbol{\theta}^{\varepsilon,*}$ as $\boldsymbol{\theta}^f$. We get

$$\begin{aligned} \frac{du_j^f}{dt} &= \frac{1}{n} \sum_{i=1}^n (-\ell'(f(\boldsymbol{\theta}^f, \mathbf{x}_i) y_i)) \phi(\mathbf{v}_j^f, \mathbf{x}_i) y_i, \\ \frac{d\mathbf{v}_j^f}{dt} &= \frac{1}{n} \sum_{i=1}^n (-\ell'(f(\boldsymbol{\theta}^f, \mathbf{x}_i) y_i)) u_j^f \nabla_{\mathbf{v}} \phi(\mathbf{v}_j^f, \mathbf{x}_i) y_i, \end{aligned}$$

where $\boldsymbol{\theta}^f(0) = \boldsymbol{\theta}^{\varepsilon,*}$.

First, we want to show that $\|\boldsymbol{\theta}^f\| \xrightarrow{t \rightarrow \infty} \infty$. Similarly to the proof of Proposition 5.7, using Proposition D.2, we get that for all $i \in S_1$

$$\begin{aligned} f(\boldsymbol{\theta}^f, \mathbf{x}_i) &\leq (u_1^e)^2 (1 + a), \\ f(\boldsymbol{\theta}^f, \mathbf{x}_i) &\geq (u_1^e)^2 (1 - 2a) - 2a((u_2^e)^2 + (u_4^e)^2), \end{aligned}$$

where $a := \delta + \xi$. Since θ^f converges in direction to θ^{mm} , we could find a moment when ratio $\frac{8(u_k^e)^2}{\|\theta^f\|^2}$ will lie in interval $(1 - \epsilon, 1 + \epsilon)$, where $\frac{3\epsilon}{1-\epsilon} < \frac{1-6a}{4a}$. Then,

$$f(\theta^f, \mathbf{x}_i) \geq b\|\theta^f\|^2,$$

where $b := \frac{(1-\epsilon)(1-2a)-(1+\epsilon)2a}{8} > 0$. Similarly, for other clusters. Also define $c := \frac{(1+\epsilon)(1+a)}{8}$. Now, notice that

$$\begin{aligned} \frac{1}{4} \frac{d\|\theta\|^2}{dt} &= \sum_{k=1}^4 u_k^e \frac{du_k^e}{dt} = \frac{1}{n} \sum_{i=1}^n \sum_{k=1}^4 u_k^e (-\ell'(f(\theta^f, \mathbf{x}_i)y_i)) \phi(\mathbf{v}_k^e, \mathbf{x}_i)y_i = \frac{1}{n} \sum_{i=1}^n \frac{f(\theta^f, \mathbf{x}_i)y_i}{1 + \exp(f(\theta^f, \mathbf{x}_i)y_i)} \\ &\geq \frac{b\|\theta^f\|^2}{1 + \exp(c\|\theta^f\|)}. \end{aligned}$$

Similarly to the proof of Proposition 5.7, this differential inequality means that $\|\theta^f\| \xrightarrow{t \rightarrow \infty} \infty$.

Now, we want to apply Theorem 5.6. To do it, choose ζ small enough so that $\chi(\theta^{mm})$ becomes the biggest local-max-margin direction in ζ -neighborhood around its image in the original weight space. Then, apply Theorem 5.6 for $\chi(\theta^{mm})$, which gives us parameters ω and ρ . Consider time T when $\chi(\theta^f)$ converged to the desired local-max-margin direction closer than $\omega/2$ and its scale became bigger than 2ρ . After that, apply Theorem 2.1, Chapter 5 from Hartman (2002) (notice that our activation function is twice continuously differentiable) and choose the initial scale σ to be sufficiently small so that original system at time $\theta(T + T_2^\epsilon)$ is bigger than ρ and $\omega/2$ -close in direction to θ^f . Then, θ will converge to some direction with normalized margin bigger than that of $\chi(\theta^{mm})$. However, since $\chi(\theta^{mm})$ is a strict local-max-margin direction, this would mean that θ will converge in direction to $\chi(\theta^{mm})$.

E. Additional Experiments for Section 5

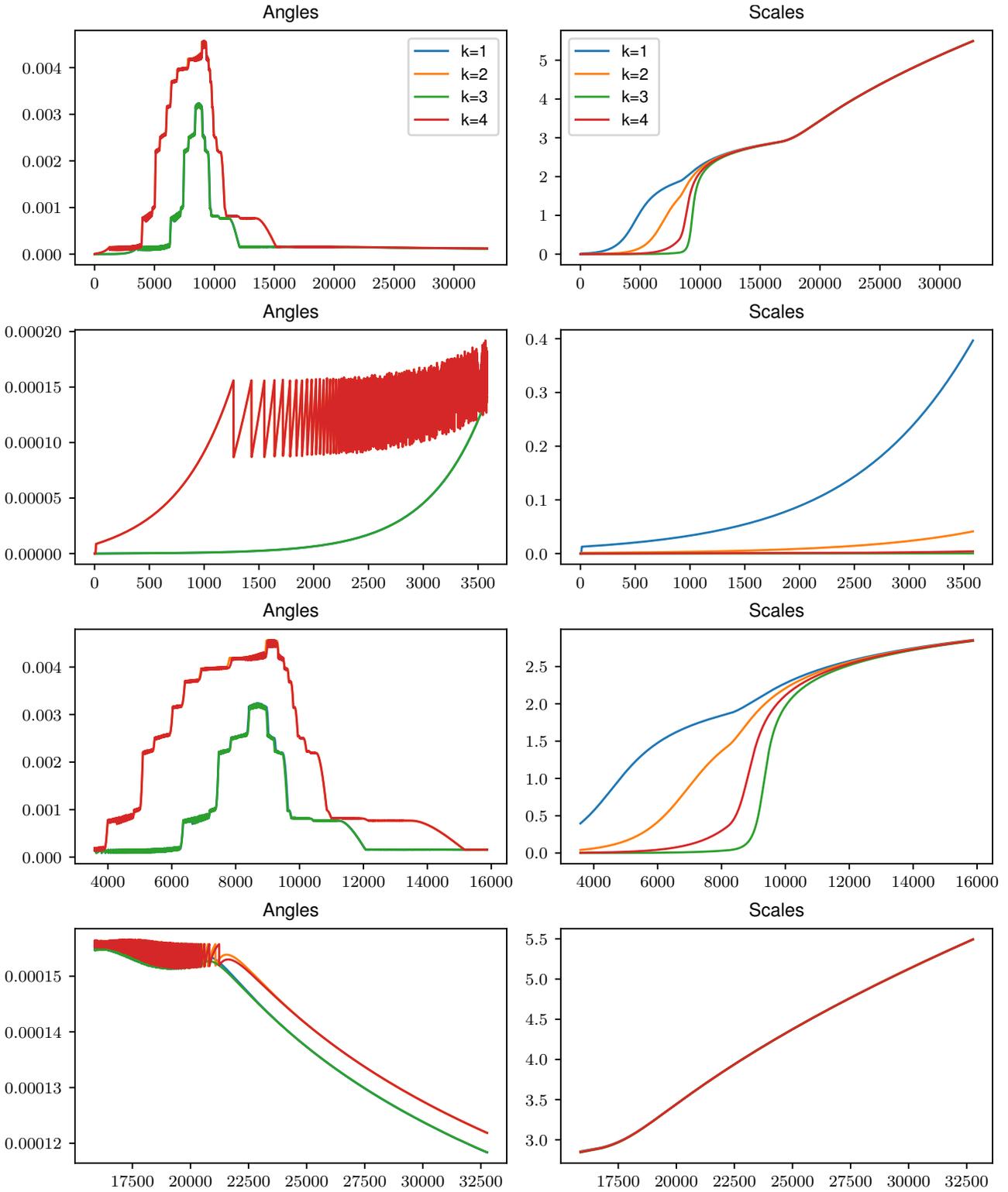


Figure 5. Evolution of 4-neuron network initialized at $(u_1^e(0), u_2^e(0), u_3^e(0), u_4^e(0)) = (10^{-4}, -10^{-5}, 10^{-7}, -10^{-6})$. The first row depicts the whole training process; the second row depicts the first 3584 training epochs; the third row depicts the epochs from 3584 to 15872; the last row depicts training after the 15872th epoch. Notice that $\alpha_1 \approx \alpha_3$ and $\alpha_2 \approx \alpha_4$.

Simplicity Bias of Two-Layer Networks beyond Linearly Separable Data

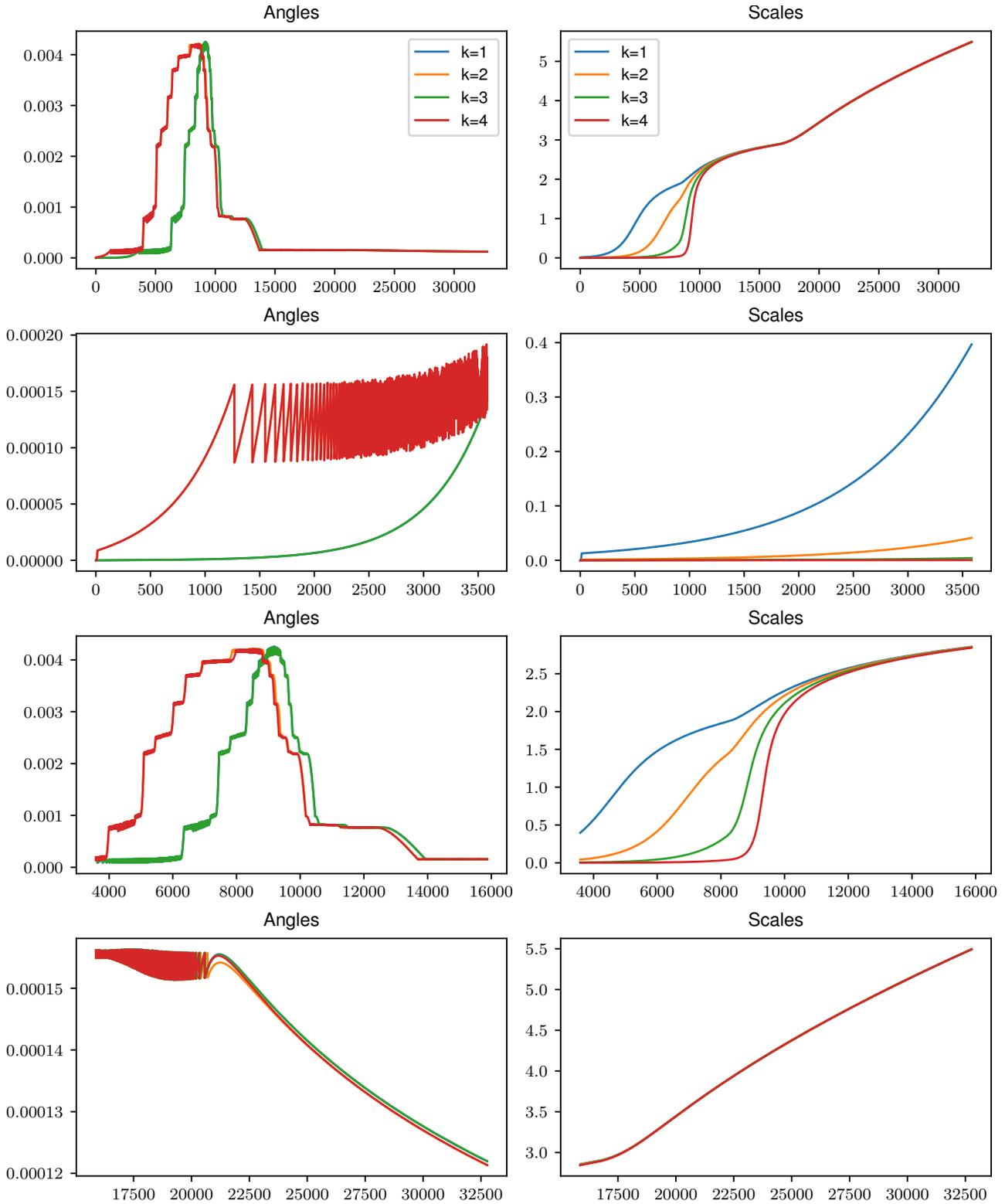


Figure 6. Evolution of 4-neuron network initialized at $(u_1^e(0), u_2^e(0), u_3^e(0), u_4^e(0)) = (10^{-4}, -10^{-5}, 10^{-6}, -10^{-7})$. The first row depicts the whole training process; the second row depicts the first 3584 training epochs; the third row depicts the epochs from 3584 to 15872; the last row depicts training after the 15872th epoch. Notice that $\alpha_1 \approx \alpha_3$ and $\alpha_2 \approx \alpha_4$.

Simplicity Bias of Two-Layer Networks beyond Linearly Separable Data

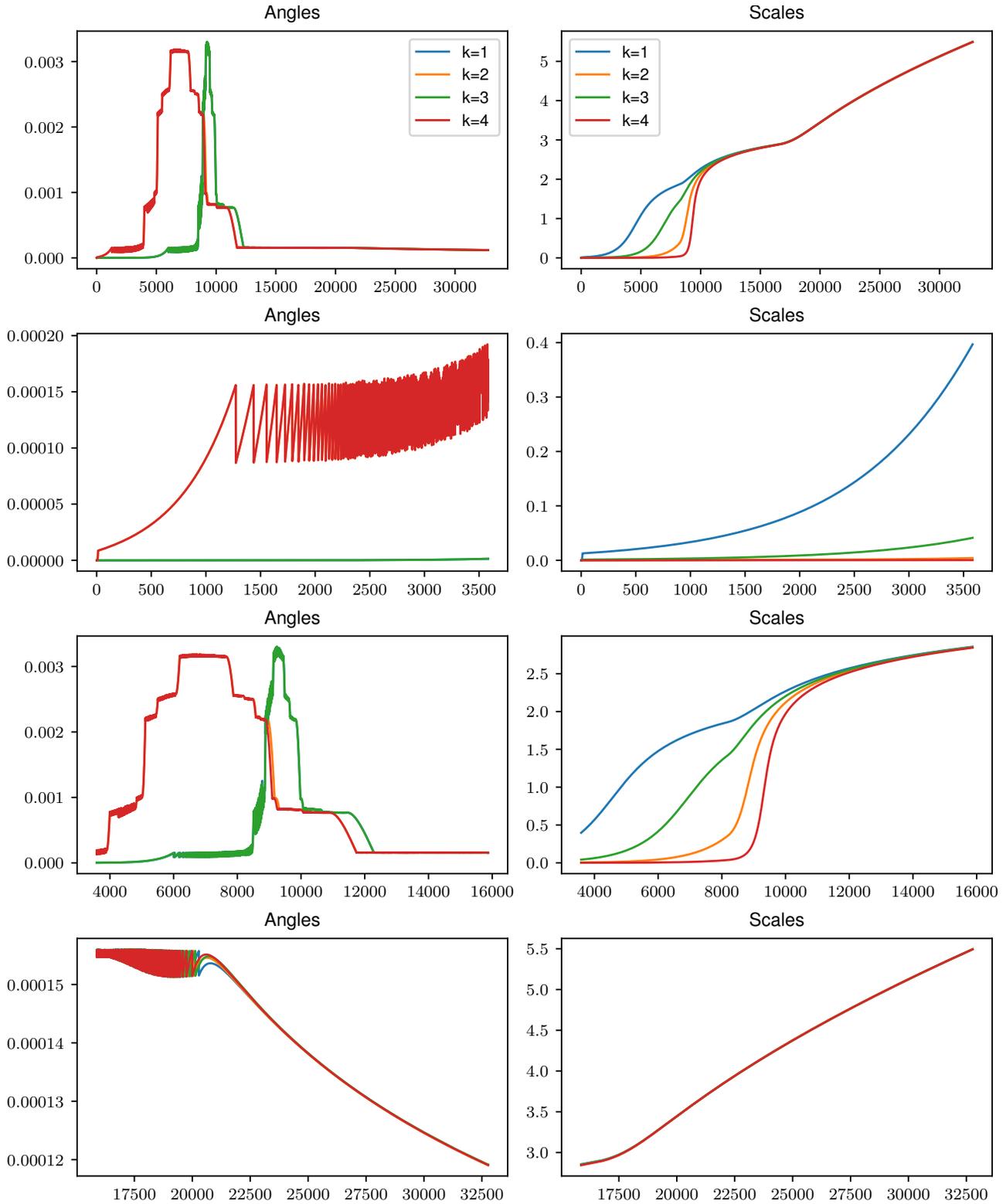


Figure 7. Evolution of 4-neuron network initialized at $(u_1^e(0), u_2^e(0), u_3^e(0), u_4^e(0)) = (10^{-4}, -10^{-6}, 10^{-5}, -10^{-7})$. The first row depicts the whole training process; the second row depicts the first 3584 training epochs; the third row depicts the epochs from 3584 to 15872; the last row depicts training after the 15872th epoch. Notice that $\alpha_1 \approx \alpha_3$ and $\alpha_2 \approx \alpha_4$.

F. Additional Experimental Results

F.1. Experimental Details

Data We used the usual MNIST and CIFAR-10 datasets for creation of dominos. For creating train and test data, we used the default train-test split of these datasets, resulting in 50000 images in train set and 10000 images in test set. We further devoted 25% train and test data for validation, giving us four datasets: train-train, train-validation, test-train, and test-validation. We also normalize images in these datasets using the default values for MNIST and CIFAR-10. During training, we also apply random horizontal flip augmentation.

Model We used the standard model from Torchvision library, but changed the first layer to 3×3 convolutions instead of the default 7×7 convolutions. Additionally, after initialization we multiplied all parameters of the model by a factor 2^{-5} to capture the desired simplicity bias mechanism.

Optimization procedure We use the standard SGD optimizer from PyTorch and linear learning scheduler with warm-up from Transformers library. The parameters of data and optimizer are listed below.

batch_size	128
lr	0.125
momentum	0.9
nesterov	True
weight_decay	0.0005
Share of warm-up steps	12.5%

Parameters of logistic regression We used the standard implementation of the logistic regression from scikit-learn library. By default, we use the following parameters.

penalty	l2
C	1000
max_iter	20000

Notice that effectively the current version of scikit-learn library does not allow to change the parameter `maxfun` parameter of the `lbfgs` optimizer. Thus, to ensure convergence we rerun fitting procedure 50 times using warm start.

F.2. Additional Experiments

Figures 8, 9, and 10 depict additional results for Section 7.

For Figure 8, we repeated the experiment but used MNIST labels on the test set. As we can see, the learned features are sufficient to achieve almost perfect accuracy on MNIST label, indicating that the network learned “simple” MNIST features.

For Figure 9, we repeated the experiment but did not scale the model at initialization. As we can see, the drop in OOD accuracy is less pronounced for the model initialized from the normal scale (approximately $6.43\% \pm 2.43\%$, implying p-value around 0.00062). This experiment indicates that closeness to more lazy training regime is indeed beneficial for OOD generalization in presence of simplicity bias.

For Figure 10, we repeated the experiment on different train data, on which the correlation between MNIST and CIFAR-10 classes is not perfect and with 5% probability MNIST class might not match CIFAR-10 class. As we can see, the model still experience simplicity bias. However, the drop in OOD accuracy at the end of training disappears.

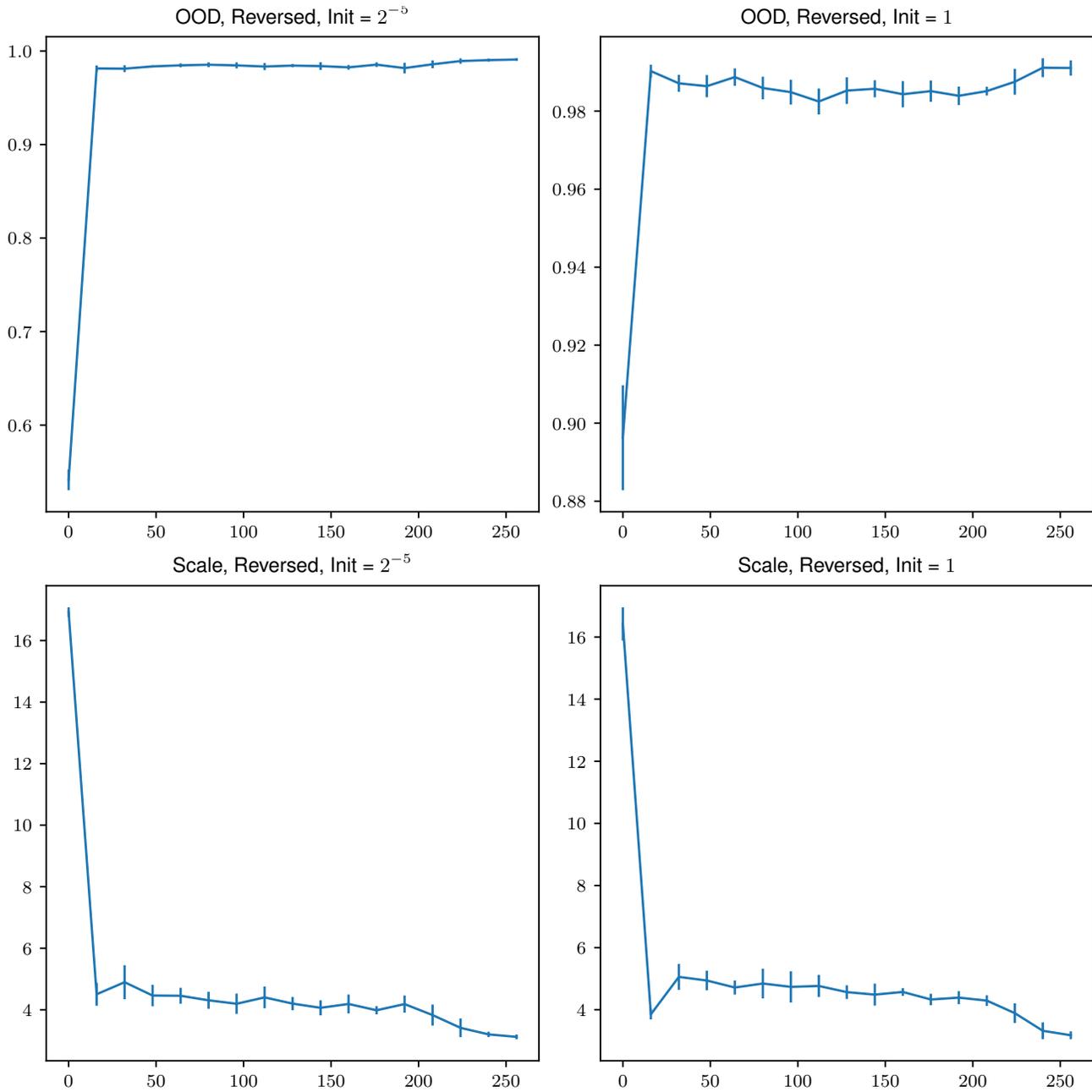


Figure 8. Accuracy and scale of the logistic regression on the validation part of the OOD test set (y -axis) vs. the training epoch at which the ResNet features are extracted (x -axis).

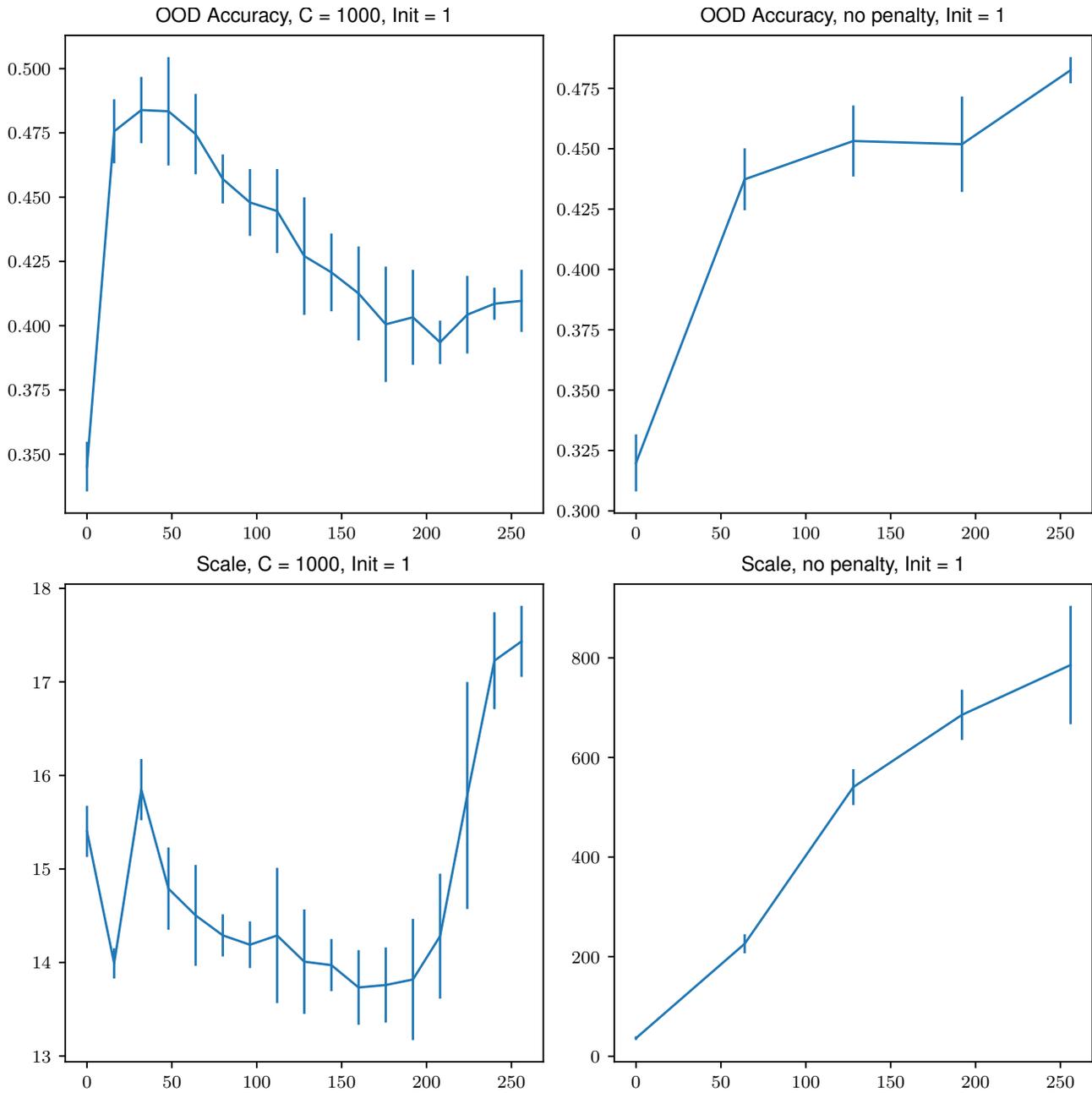


Figure 9. Accuracy and scale of the logistic regression on the validation part of the OOD test set (y -axis) vs. the training epoch at which the ResNet features are extracted (x -axis).

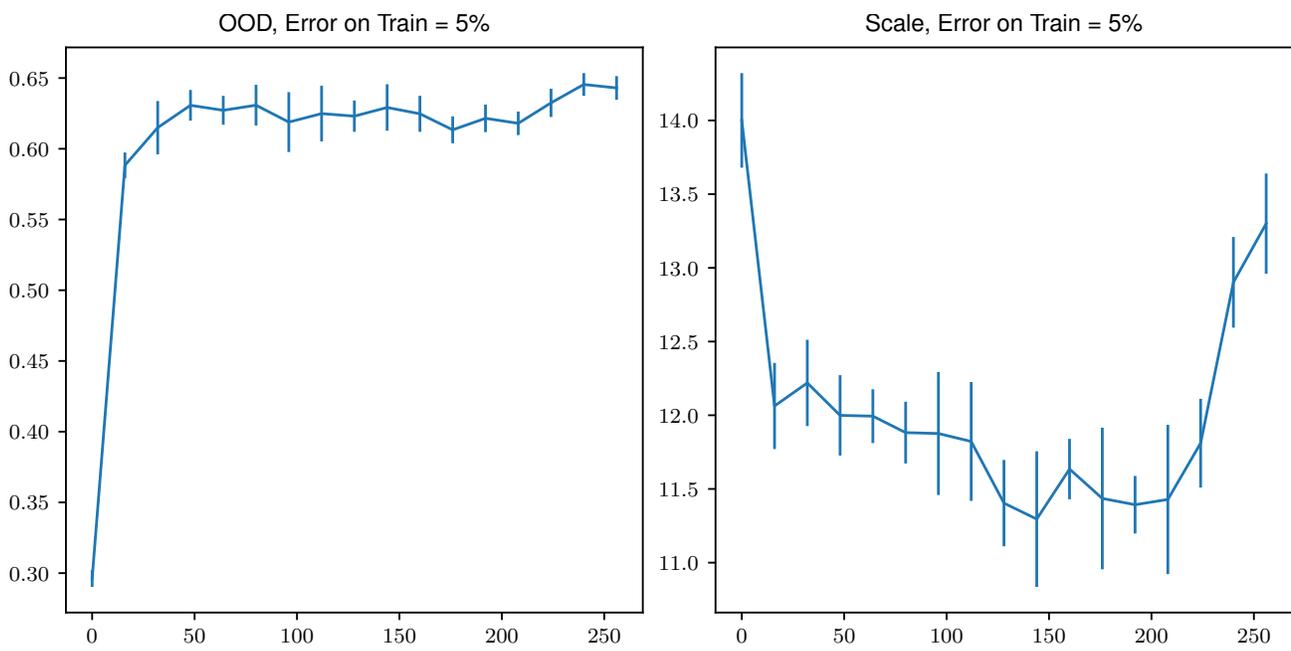


Figure 10. Accuracy and scale of the logistic regression on the validation part of the OOD test set (y -axis) vs. the training epoch at which the ResNet features are extracted (x -axis).