

LEARNING DYNAMICS OF DEEP MATRIX FACTORIZATION BEYOND THE EDGE OF STABILITY

Anonymous authors

Paper under double-blind review

ABSTRACT

Deep neural networks trained using gradient descent with a fixed learning rate η often operate in the regime of “edge of stability” (EOS), where the largest eigenvalue of the Hessian equilibrates about the stability threshold $2/\eta$. Existing theoretical analyses of EOS focus on simple prototypes, such as scalar functions or second-order regression models, which limits our understanding of the phenomenon in deep networks. In this work, we present a fine-grained analysis of the learning dynamics of (deep) linear networks (DLN) within the deep matrix factorization loss beyond EOS. For DLNs, loss oscillations within EOS follow a period-doubling route to chaos. We theoretically analyze the regime of the 2-period orbit and show that the loss oscillations occur within a small subspace, with the dimension of the subspace precisely characterized by the learning rate. Our analysis contributes to explaining two key phenomena in deep networks: (i) shallow models and simple tasks do not always exhibit EOS (Cohen et al., 2021); and (ii) oscillations occur within top features (Zhu et al., 2023a). We present experiments to support our theory, along with examples demonstrating how these phenomena occur in nonlinear networks and how they differ from those in DLNs.

1 INTRODUCTION

Understanding generalization in deep neural networks requires an understanding of the optimization process in gradient descent (GD). In the literature, it has been empirically observed that the learning rate η plays a key role in driving generalization (Hayou et al., 2024; Lewkowycz et al., 2020). The “descent lemma” from classical optimization theory says that for a β -smooth loss $\mathcal{L}(\Theta)$ parameterized by Θ , gradient descent (GD) iterates satisfy

$$\mathcal{L}(\Theta(t+1)) \leq \mathcal{L}(\Theta(t)) - \frac{\eta(2-\eta\beta)}{2} \|\nabla\mathcal{L}(\Theta(t))\|_2^2,$$

and so the learning rate should be chosen as $\eta < 2/\beta$ to monotonically decrease the loss. However, many recent works have shown that the training loss decreases even for $\eta > 2/\beta$, albeit non-monotonically. Surprisingly, it has been observed that choosing such a learning rate often provides better generalization over smaller ones that lie within the stability threshold. This observation has led to a series of works analyzing the behavior of GD within a regime dubbed “the edge of stability” (EOS). By letting Θ be a deep network, we formally define EOS as follows:

Definition 1 (Edge of Stability (Cohen et al., 2021)). *During training, the sharpness of the loss, defined as $S(\Theta) := \|\nabla^2\mathcal{L}(\Theta)\|_2$, continues to grow until it reaches $2/\eta$ (progressive sharpening), after which it stabilizes around $2/\eta$. During this process, the training loss behaves non-monotonically over short timescales but consistently decreases over long timescales.*

Using a large learning rate to operate within the EOS regime is hypothesized to give better generalization performance by inducing “catapults” in the training loss (Zhu et al., 2023a). Intuitively, whenever the sharpness $S(\Theta)$ exceeds the local stability limit $2/\eta$, the GD iterates momentarily diverge (or catapults) out of a sharp region and self-stabilizes (Damian et al., 2023) to settle for a flatter region where the sharpness is below $2/\eta$, which has shown to correlate with better generalization (Keskar et al., 2017; Izmailov et al., 2019; Petzka et al., 2021; Foret et al., 2021; Gatmiry et al., 2023). Of course, the dynamics within EOS differ based on the loss landscape. When the loss

054
055
056
057
058
059
060
061
062
063
064
065
066
067
068
069
070
071
072
073
074
075
076
077
078
079
080
081
082
083
084
085
086
087
088
089
090
091
092
093
094
095
096
097
098
099
100
101
102
103
104
105
106
107

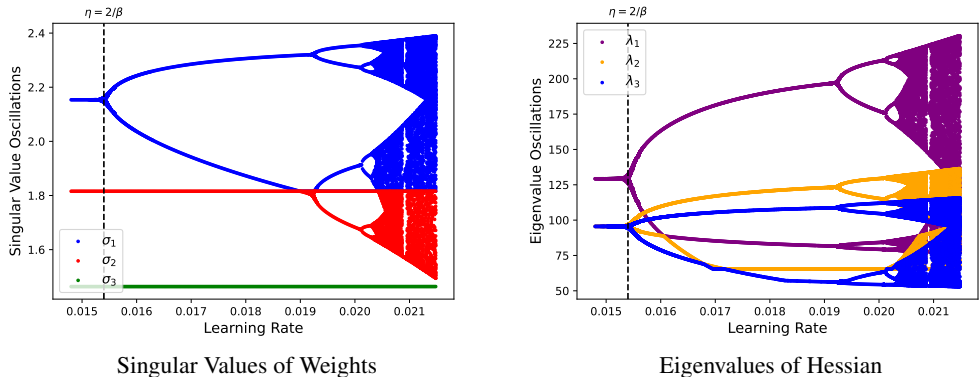


Figure 1: Bifurcation plot of the oscillations in the singular values (left) and the eigenvalues of the Hessian (right) of a 3-layer end-to-end DLN. The bifurcation plots indicate the existence of a period-doubling route to chaos in DLNs, which we analyze by examining the two-period orbit. Here, $\eta > 2/\beta$ corresponds to the EOS regime, where $\beta = L\sigma_{*,1}^{2-2/L}$ is the sharpness at the minima, L is the depth of the network and $\sigma_{*,1}$ is the first singular value of the target matrix \mathbf{M}_* .

landscape is highly non-convex with many local valleys, catapults may occur, whereas sustained oscillations may exist for other landscapes. It is of great interest to understand these behaviors within different network architectures to further our understanding of EOS.

From a theoretical perspective, there have been many recent efforts to understand EOS. These works generally focus on analyzing “simple” functions, examples including scalar losses (Zhu et al., 2023b; Wang et al., 2023; Kreisler et al., 2023), quadratic regression models (Agarwala et al., 2022), diagonal linear networks (Even et al., 2024) and two-layer matrix factorization (Chen & Bruna, 2023). However, the simplicity of these functions cannot fully capture the behaviors of deep neural networks within the EOS regime. Specifically, the following observations remain unexplained by existing analyses: (i) mild (or no) sharpening occurs when either networks are shallow or “simple” datasets are used for training (Caveat 2 from (Cohen et al., 2021)); and (ii) the oscillations and catapults in the weights occur within the top singular values of each weight matrix (Zhu et al., 2023a).

In this work, we present a fine-grained analysis of the learning dynamics of deep linear networks (DLNs) within the EOS regime, demonstrating that these phenomena can be replicated and effectively explained using DLNs. Generally, there are two lines of work for DLNs: (i) those that analyze the effects of depth and initialization scale, and how they implicitly bias the trajectory of gradient flow towards low-rank solutions when the learning rate is chosen to be stable (Saxe et al., 2014; Arora et al., 2018; 2019; You et al., 2020; Liu et al., 2022; Zhang et al., 2024a; Pesme & Flammarion, 2023; Jacot et al., 2022), and (ii) those that analyze the similarities in behavior between linear and nonlinear networks (Wang et al., 2024; Zhang et al., 2024b; Yaras et al., 2023). Our analysis builds upon these works to show that DLNs exhibit intricate and interesting behaviors outside the stability regime and to demonstrate how factors such as depth and initialization scale contribute to the EOS regime. Our main contributions can be summarized as follows:

- **Oscillations in Top Subspaces.** We show that there exist periodic oscillations within r -dimensional subspaces in DLNs, where r is precisely characterized by the learning rate. For DLNs, a period-doubling route to chaos (Ott, 2002) exists in both the singular values of the DLN and the eigenvalues of the Hessian, as shown in Figure 1. We rigorously characterize the case of the two-period orbit, aiming to contribute to explaining the empirical observations by Zhu et al. (2023a) and Cohen et al. (2021). We also prove that the learning rate needed to enter EOS is a function of the network depth, further revealing its role in deep networks.
- **Difference in DLNs and Diagonal Linear Networks.** While DLNs and diagonal linear networks exhibit similar behaviors under a stable learning rate, we demonstrate that their dynamics differ within the EOS regime (Gidel et al., 2019b). Near the global minima, we show that DLNs have additional curvature directions that influence behavior within the EOS regime, distinguish-

ing them from diagonal linear networks. This offers a unique perspective on how changes in the landscape affect behaviors within the EOS regime, depending on the learning rate and depth.

2 NOTATION AND PROBLEM SETUP

Notation. We denote vectors with bold lower-case letters (e.g., \mathbf{x}) and matrices with bold upper-case letters (e.g., \mathbf{X}). We use \mathbf{I}_n to denote an identity matrix of size $n \in \mathbb{N}$. We use $[L]$ to denote the set $\{1, 2, \dots, L\}$. We use the notation $\sigma_i(\mathbf{A})$ to denote the i -th singular value of the matrix \mathbf{A} . We also use the notation $\sigma_{i,\ell}$ to denote the i -th singular value of the matrix \mathbf{W}_ℓ .

Deep Matrix Factorization Loss. The objective in deep matrix factorization is to model a low-rank matrix $\mathbf{M}_* \in \mathbb{R}^{d \times d}$ with $\text{rank}(\mathbf{M}_*) = r$ via a DLN parameterized by a set of parameters $\Theta = (\mathbf{W}_1, \mathbf{W}_2, \dots, \mathbf{W}_L)$, which can be estimated by solving

$$\underset{\Theta}{\text{argmin}} f(\Theta) := \frac{1}{2} \left\| \underbrace{\mathbf{W}_L \cdots \mathbf{W}_1}_{=:\mathbf{W}_{L:1}} - \mathbf{M}_* \right\|_F^2, \quad (1)$$

where we adopt the abbreviation $\mathbf{W}_{j:i} = \mathbf{W}_j \cdots \mathbf{W}_i$ to denote the end-to-end DLN and is identity when $j < i$. We assume that each weight matrix has dimensions $\mathbf{W}_\ell \in \mathbb{R}^{d \times d}$ to observe the effects of overparameterization. We also assume that the singular values of \mathbf{M}_* are distinct such that $\sigma_{*,1} > \dots > \sigma_{*,r}$.

Optimization. We update each weight matrix $\mathbf{W}_\ell \in \mathbb{R}^{d \times d}$ using GD with iterations given by

$$\mathbf{W}_\ell(t) = \mathbf{W}_\ell(t-1) - \eta \cdot \nabla_{\mathbf{W}_\ell} f(\Theta(t-1)), \quad \forall \ell \in [L], \quad (2)$$

where $\eta > 0$ is the learning rate and $\nabla_{\mathbf{W}_\ell} f(\Theta(t))$ is the gradient of $f(\Theta)$ with respect to the ℓ -th weight matrix at the t -th GD iterate.

Initialization. To encompass a wide range of initialization schemes, we consider both a balanced and unbalanced initialization, respectively:

$$\mathbf{W}_\ell(0) = \alpha \mathbf{I}_d, \quad \forall \ell \in [L], \quad \text{and} \quad \mathbf{W}_L(0) = \mathbf{0}, \quad \mathbf{W}_\ell(0) = \alpha \mathbf{I}_d, \quad \forall \ell \in [L-1], \quad (3)$$

where $\alpha > 0$ is a small constant. We assume α is chosen small enough such that $\alpha \in (0, \sigma_{*,r})$, where $\sigma_{*,r}$ is the r -th singular value of \mathbf{M}_* . Generally, many existing works on both shallow and deep linear networks assume a zero-balanced initialization (i.e., $\mathbf{W}_i^\top(0)\mathbf{W}_i(0) = \mathbf{W}_j(0)\mathbf{W}_j^\top(0)$ for $i \neq j$). This introduces the invariant $\mathbf{W}_i^\top(t)\mathbf{W}_i(t) = \mathbf{W}_j(t)\mathbf{W}_j^\top(t)$ for all $t > 0$, ensuring two (degenerate) conditions throughout the training trajectory: (i) the singular vectors of each of the layers remain aligned and (ii) the singular values stay balanced. For the unbalanced initialization, the zero weight layer can be viewed as the limiting case of initializing the weights with a (very) small constant $\alpha' \ll \alpha$, and has been similarly explored by Varre et al. (2023); Xu et al. (2024), albeit for two-layer networks. The zero weight layer relieves the balancing condition of the singular values. Rather than staying balanced, we show that the singular values become increasingly balanced (see Lemma 2). This allows us to jointly analyze the singular values of the weights for either case.

Nevertheless, we also show that our analysis is not limited to either initialization but applies to *any initialization* that converges to a set we call the singular vector stationary set (see Proposition 1). To the best of our knowledge, it is common to assume that the singular vectors remain aligned, as many existing works make the same assumption (Varre et al., 2023; Arora et al., 2019; Saxe et al., 2014; Gidel et al., 2019a; Chou et al., 2024b; Min Kwon et al., 2024).

3 DEEP MATRIX FACTORIZATION BEYOND THE EDGE OF STABILITY

When using a large learning rate, the learning dynamics can typically be separated into two distinct stages: (i) progressive sharpening and (ii) the edge of stability. Within the progressive sharpening stage, the sharpness lies below $2/\eta$ and tends to continually rise. Our goal is to analyze the EOS stage under the deep matrix factorization formulation. Here, we observe that the training loss fluctuates due to layerwise singular value oscillations, as illustrated in Figure 2.

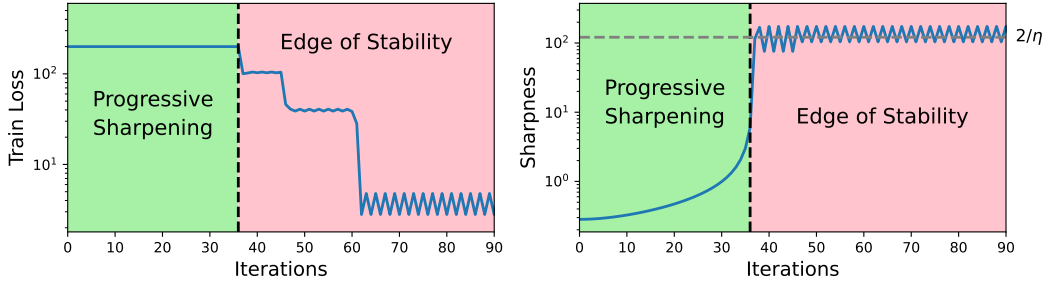


Figure 2: Depiction of the two phases of learning in the deep matrix factorization problem for a network of depth 3. It appears that upon escaping the first saddle point, the GD iterates enter the EOS regime, where the sharpness hovers just above $2/\eta$.

3.1 MAIN RESULTS

Before we present our main results, we provide a definition of what we refer to as a *strict balanced state* of the singular values for the weight matrices. If the parameters are said to be in a strict balanced state, then the singular values of each weight are balanced in the sense that they take the same values across all weight layers.

Definition 2 (Strict Balanced State). *The parameters Θ of the DLN from Equation (1) are said to be in a strict balanced state if for some $t \geq 0$*

$$\sigma_i(\mathbf{W}_\ell(t)) = \sigma_i(\mathbf{W}_k(t)), \quad \forall i \in [r], \quad \forall \ell, k \in [L],$$

where $\sigma_i(\mathbf{W}_\ell)$ denotes the i -th singular value of the ℓ -th layer and r is the rank of the matrix \mathbf{M}_* .

It is straightforward to show that the parameters are in a strictly balanced state for all $t \geq 0$ if we initialize the singular values to be the same across all weight matrices \mathbf{W}_ℓ . Hence, it immediately holds that the balanced initialization is in a strict balanced state. However, the one-zero initialization in Equation (3) sets the singular values of \mathbf{W}_L to zero, meaning the parameters are not initially in a strictly balanced state. Since we prove in Lemma 2 that the singular values across layers become increasingly balanced, we assume a strict balanced state throughout the rest of this paper. Next, we derive the eigenvalues of the Hessian at convergence, such that we can identify the learning rate needed to enter the EOS regime for DLNs.

Lemma 1 (Eigenvalues of Hessian at Convergence). *Consider running GD on the deep matrix factorization loss $f(\Theta)$ defined in Equation (1). Under strict balancing with any stationary point Θ with $\nabla_{\Theta} f(\Theta) = 0$, the set of all non-zero eigenvalues of the training loss Hessian are given by*

$$\lambda_{\Theta} = \underbrace{\left\{ L\sigma_{*,i}^{2-\frac{2}{L}}, \sigma_{*,i}^{2-\frac{2}{L}} \right\}_{i=1}^r}_{\text{self-interaction}} \cup \underbrace{\left\{ \sum_{\ell=0}^{L-1} \left(\sigma_{*,i}^{1-\frac{1}{L}-\frac{1}{L}\ell} \cdot \sigma_{*,j}^{\frac{1}{L}\ell} \right)^2 \right\}_{i \neq j}^r}_{\text{interaction with other singular values}} \cup \underbrace{\left\{ \sum_{\ell=0}^{L-1} \left(\sigma_{*,k}^{1-\frac{1}{L}-\frac{1}{L}\ell} \cdot \alpha^{\ell} \right)^2 \right\}_{k=1}^r}_{\text{interaction with initialization}}$$

where $\sigma_{*,i}$ is the i -th singular value of the target matrix $\mathbf{M}_* \in \mathbb{R}^{d \times d}$, $\alpha \in \mathbb{R}$ is the initialization scale, L is the depth of the network, and the second element of the set under “self-interaction” has a multiplicity of $d - r$.

We defer all of the proofs to Appendix C. By Lemma 1, we can see that the sharpness is exactly $\|\nabla^2 f(\Theta)\|_2 = L\sigma_{*,1}^{2-\frac{2}{L}}$ under strict balancing. Hence, as long as η is set to $\eta > 2/L\sigma_{*,1}^{2-\frac{2}{L}}$, we will observe oscillations in the loss. Interestingly, notice that all non-zero eigenvalues are a function of network depth. For a deeper network, the sharpness will be larger, implying that a smaller learning rate can be used to drive the DLN into EOS. This provides a unique perspective on how the learning rate should be chosen as networks become deeper and explains the observation made by Cohen et al. (2021), who observed that sharpness scales with the depth of the network. In Section 3.3, we show that these eigenvalues account for the primary difference between DLNs and diagonal linear networks – the eigenvalues that correspond to the “interaction with other singular values” are absent in diagonal linear networks, leading to oscillations in two or more dimensions occurring at different learning rates. Next, we present our result on the two-period orbit in the first singular value.

Theorem 1 (Rank-1 Oscillation). *Let $\mathbf{M}_* = \mathbf{U}_* \Sigma_* \mathbf{V}_*^\top$ denote the SVD of the target matrix and let $S := L\sigma_{*,1}^{2-\frac{2}{L}}$, $\alpha' := \left(\ln \left(\frac{2\sqrt{2}}{\eta L \sigma_{*,1}^{\frac{2}{L}}} \cdot \frac{\sigma_{*,1}^{\frac{4}{L}}}{L^2 \cdot 2^{\frac{2L-3}{L}}} \right) \right)^{\frac{1}{4}}$, and $K' := \max \left\{ \sum_{\ell=0}^{L-1} \left(\sigma_{*,1}^{1-\frac{1}{L}-\frac{1}{L}\ell} \cdot \sigma_{*,2}^{\frac{1}{L}\ell} \right)^2, \frac{S}{2\sqrt{2}} \right\}$. If we run GD on the deep matrix factorization loss with initialization scale $\alpha < \alpha'$ and learning rate $\eta = \frac{2}{K}$, where $K' < K < S$, then under strict balancing, each weight matrix $\mathbf{W}_\ell \in \mathbb{R}^{d \times d}$ oscillates around the minima in a 2-period fixed orbit ($i \in \{1, 2\}$) as follows:*

$$\mathbf{W}_L = \underbrace{\rho_i \cdot \mathbf{u}_{*,1} \mathbf{v}_{*,1}^\top}_{\text{oscillation subspace}} + \underbrace{\sum_{j=2}^r \sigma_{*,j} \mathbf{u}_{*,j} \mathbf{v}_{*,j}^\top}_{\text{stationary subspace}}, \quad i \in \{1, 2\},$$

$$\mathbf{W}_\ell = \underbrace{\rho_i \cdot \mathbf{v}_{*,1} \mathbf{v}_{*,1}^\top}_{\text{oscillation subspace}} + \underbrace{\sum_{j=2}^r \sigma_{*,j} \mathbf{v}_{*,j} \mathbf{v}_{*,j}^\top}_{\text{stationary subspace}}, \quad i \in \{1, 2\}, \quad \forall \ell \in [L-1],$$

where $\rho_1 \in (0, \sigma_{*,1}^{1/L})$ and $\rho_2 \in (\sigma_{*,1}^{1/L}, (2\sigma_{*,1})^{1/L})$ are the two real roots of the polynomial $g(\rho) = 0$, where

$$g(\rho) = \rho^L \cdot \frac{1 + (1 + \eta L (\sigma_{*,1} - \rho^L) \cdot \rho^{L-2})^{2L-1}}{1 + (1 + \eta L (\sigma_{*,1} - \rho^L) \cdot \rho^{L-2})^{L-1}} - \sigma_{*,1}.$$

Remarks. From Lemma 1, the second largest eigenvalue is given by $\sum_{\ell=0}^{L-1} \left(\sigma_{*,1}^{1-\frac{1}{L}-\frac{1}{L}\ell} \cdot \sigma_{*,2}^{\frac{1}{L}\ell} \right)^2$. If the learning rate is chosen such that K less than the second largest eigenvalue, then oscillation occurs within the top eigenvector direction of the Hessian, which is amenable to the oscillation in the first singular value as shown in Theorem 1. The oscillation amplitude is governed by ρ_1 and ρ_2 , which are below and above the minima, respectively. Hence, the oscillations will occur about the minima. This aims to theoretically show why (i) oscillations only occur within top subspaces of the network as observed by Zhu et al. (2023a) and (ii) oscillations are more prevalent in the direction of the stronger features (measured by the magnitude of the singular values). Finally, we remark that the additional bound on the learning rate as well as the initialization scale is an artifact of Lemma 2, which are needed to ensure that balancing occurs.

Our result also generalizes the recent theoretical findings of Chen & Bruna (2023), where they proved the existence of a certain class of scalar functions $f(x)$ for which GD does not diverge even when operating beyond the stability threshold $\eta > \frac{2}{f''(\hat{x})}$, where \hat{x} is a local minimum of $f(x)$. Specifically, they showed that for a function dependent $\epsilon > 0$, there exists a range $\eta \in \left(\frac{2}{f''(\hat{x})}, \frac{2}{f''(\hat{x})}(1 + \epsilon) \right)$, where the loss oscillates around the local minima with a certain periodicity. As η increases beyond $\frac{2}{f''(\hat{x})}$, the oscillations gradually enter higher periodic orbits (e.g., 2, 4, 8 periods), then transition into chaotic behavior, and ultimately lead to divergence. In our work, we prove that this oscillatory behavior beyond the stability threshold occurs even in DLNs.

To generalize Theorem 1 to show oscillations in two or more subspaces, we require a careful treatment of the derivation of the roots for ρ_1 and ρ_2 and the corresponding singular vectors that contribute to the oscillation subspace. However, in the following, we prove that as long as K is set to $\lambda_r \geq K > \lambda_{r+1}$, where λ_r is the r -th eigenvalue of the Hessian, we will observe oscillations in the top- r subspaces. This is proved using Lemma 11 (restated from Chen & Bruna (2023)), where the necessary condition for stable two-period orbit is that $f_{\Delta_i}^{(3)}(x) \neq 0$ and $3[f_{\Delta_i}^{(3)}(x)]^2 - f_{\Delta_i}^{(2)}(x)f_{\Delta_i}^{(4)}(x) > 0$ around a local minima x , where Δ_i is the i -th eigenvector of the Hessian and f_{Δ_i} is the loss function restricted to the line $\{y : y = x + t\Delta_i, t \in \mathbb{R}\}$. In Theorem 2, we prove that this condition holds for each eigenvector direction.

Theorem 2 (Stable Subspace Oscillations (Informal)). *Consider running GD on the loss in Equation (1) with initialization scale $\alpha < \alpha'$ from Theorem 1. If $\eta = \frac{2}{K}$ with $\lambda_i \leq K < \lambda_{i+1}$, then*

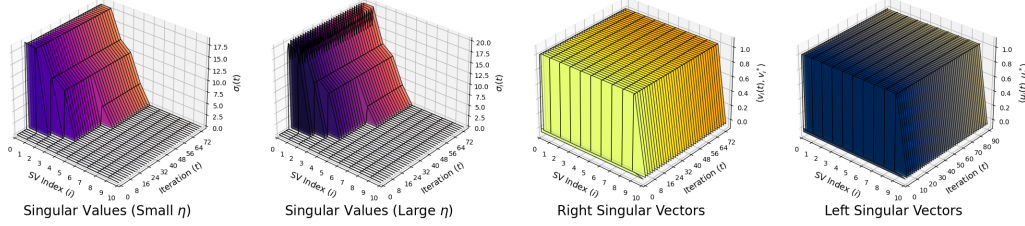


Figure 3: Illustrations of the singular vector and value evolution of the end-to-end DLN. The singular vectors of the network remain static across all iterations, as suggested by the singular vector stationary set, regardless of the learning rate. The angle between the true singular vectors and those of the network remains aligned throughout. The first singular values undergo oscillations in the large η regime, whereas they remain constant in the small η regime.

2-period orbit oscillation occurs in the direction of Δ_i , where λ_i and Δ_i denote the i -th largest eigenvalue and its corresponding eigenvector of the Hessian at convergence, respectively.

3.2 TOOLS USED IN THE ANALYSES

This section presents the two main tools used in our analyses: the singular vector stationary set and singular value balancedness. First, we present the singular vector stationary set, which allows us to encompass a wider range of initialization schemes. This set defines a broad class of initialization for which singular vector alignment occurs, simplifying the dynamics to only singular values.

Proposition 1 (Singular Vector Stationary Set). *Consider the deep matrix factorization loss in Equation (1). Let $\mathbf{M}_\star = \mathbf{U}_\star \Sigma_\star \mathbf{V}_\star^\top$ and $\mathbf{W}_\ell(t) = \mathbf{U}_\ell(t) \Sigma_\ell(t) \mathbf{V}_\ell^\top(t)$ denote the compact SVD for the target matrix and the ℓ -th layer weight matrix at time t , respectively. For any time $t \geq 0$, if $\dot{\mathbf{U}}_\ell(t) = \dot{\mathbf{V}}_\ell(t) = 0$ for all $\ell \in [L]$, then the singular vector stationary points for each weight matrix are given by*

$$\text{SVS}(f(\Theta)) = \begin{cases} (\mathbf{U}_L, \mathbf{V}_L) & = (\mathbf{U}_\star, \mathbf{Q}_L), \\ (\mathbf{U}_\ell, \mathbf{V}_\ell) & = (\mathbf{Q}_{\ell+1}, \mathbf{Q}_\ell), \quad \forall \ell \in [2, L-1], \\ (\mathbf{U}_1, \mathbf{V}_1) & = (\mathbf{Q}_2, \mathbf{V}_\star), \end{cases}$$

where $\{\mathbf{Q}_\ell\}_{\ell=2}^L$ are any set of orthogonal matrices.

The singular vector stationary set states that for any set of weights where the gradients with respect to the singular vectors become zero, the singular vectors become fixed points for subsequent iterations. Once the singular vectors become stationary, running GD further isolates the dynamics on the singular values. Hence, throughout our analysis, we re-write and consider the loss

$$\frac{1}{2} \|\mathbf{W}_{L:1}(t) - \mathbf{M}^\star\|_F^2 = \frac{1}{2} \|\Sigma_{L:1} - \Sigma^\star\|_F^2 = \frac{1}{2} \sum_{i=1}^r (\sigma_i(\Sigma_{L:1}(t)) - \sigma_{\star,i})^2, \quad (4)$$

where $\Sigma_{L:1}$ are the singular values of $\mathbf{W}_{L:1}$. This allows us to decouple the dynamics of the singular vectors and singular values, focusing on the periodicity that occurs in the singular values within the EOS regime. In Propositions 2 and 3, we prove that both the unbalanced and balanced initializations converge to this set respectively, with an illustration in Figure 3. Specifically, we show that the singular vectors belongs to the singular vector stationary set after GD iteration $t = 1$ (far before entering the EOS regime), allowing us to consider the loss in Equation (4). In Appendix B, we provide another example that belongs to this set. Next, we present a result to validate our use of the strictly balanced assumption on the unbalanced initialization case by showing that the singular values become increasingly balanced throughout the GD iterations.

Lemma 2 (Balancing). *Let $\sigma_{*,i}$ and $\sigma_{\ell,i}(t)$ denote the i -th singular value of $\mathbf{M}_* \in \mathbb{R}^{d \times d}$ and $\mathbf{W}_\ell(t)$, respectively and define $S := L\sigma_{*,1}^{2-\frac{2}{L}}$. Consider GD on the deep matrix factorization loss in Equation (1) with unbalanced initialization in (3) and learning rate $\eta < \frac{2\sqrt{2}}{S}$. If the initialization scale α satisfies $0 < \alpha < \left(\ln \left(\frac{2\sqrt{2}}{\eta L \sigma_{*,1}^{2-\frac{2}{L}}} \right) \cdot \frac{\sigma_{*,1}^{\frac{4}{L}}}{L^2 \cdot 2^{\frac{2L-3}{L}}} \right)^{\frac{1}{4}}$, then there exists a $c \in (0, 1]$ such that for all $i \in [r]$, we have $\left| \sigma_{L,i}^2(t+1) - \sigma_{\ell,i}^2(t+1) \right| < c \left| \sigma_{L,i}^2(t) - \sigma_{\ell,i}^2(t) \right|$.*

To summarize, Lemma 2 states that, provided α is chosen below a certain threshold, the top- r singular values of the weights across all layers become increasingly balanced during GD, even if they are unbalanced as in the initialization of Equation (3). This can be viewed as an implicit property of GD, which has been shown to hold for two-layer matrix factorization (Wang et al., 2021; Ye & Du, 2021; Chen & Bruna, 2023). Our result is an extension of these analyses, but to the deep matrix factorization case. Our analysis shows that the constant c changes for two different cases¹: (i) $0 < c < 1$ when the product of singular values across all layers $\sigma_i(\Sigma_{L:1}) < \sigma_{*,i}$ and (ii) $c = 1$ when $\sigma_i(\Sigma_{L:1}) > \sigma_{*,i}$. In the literature, it has been widely shown that the dynamics of DLNs (along with diagonal linear networks) exhibit an incremental learning phenomenon, where the singular values $\sigma_i(\mathbf{W}_\ell)$ start from α and increase to the target singular value one-by-one (Gissin et al., 2020; Berthier, 2023; Min Kwon et al., 2024; Jacot et al., 2021). Empirically, this implies that we often operate in the regime of $\sigma_i(\Sigma_{L:1}) < \sigma_{*,i}$, as the oscillations begin to occur once we reach and about the minima. Hence, throughout most of the learning trajectory, $0 < c < 1$ holds, and the balancing gap becomes infinitesimally small. In Figure 4, we plot the balancing gap between the top-3 singular values of a weight matrix initialized to zero and those initialized to α for a rank-3 matrix. This plot shows that the gap decreases and goes to zero empirically, and this is consistently the case across all of our experiments, with additional results provided in Appendix B. Interestingly, the iterations that correspond to the decrease in the balancing gap in Figure 4 corresponds to the case in which $\sigma_i(\Sigma_{L:1}) < \sigma_{*,i}$ and that the gap almost goes to zero before entering the EOS regime. To this end, we use this insights to assume that strict balancing holds for both initializations in Equation (3). This allows us to write the loss of the singular values into the form $\sigma_i(\Sigma_{L:1}(t)) = \sigma_i^L(t)$, which allows us to focus on the dynamics in the singular values.

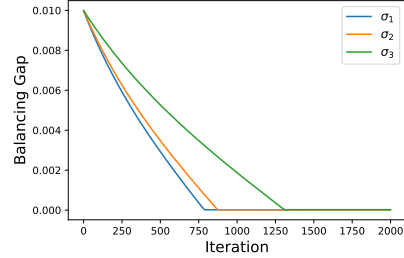


Figure 4: Plot of $|\sigma_{i,L}^2(t) - \sigma_{i,\ell}^2(t)|$ for initialization scale $\alpha = 0.01$ showing strict balancing.

3.3 RELATION TO DIAGONAL LINEAR NETWORKS

Due to singular vector stationarity and balancing, DLNs may appear equivalent to diagonal linear networks at first glance. In this section, we characterize an explicit distinction between the two networks by deriving the eigenvalues of diagonal linear networks from the Hessian at convergence and explaining how they contribute to periodic oscillations in the EOS regime.

Theorem 3 (Subspace Oscillation for Diagonal Linear Networks). *Consider an L -layer diagonal linear network on the loss*

$$\mathcal{L}(\{\mathbf{s}_\ell\}_{\ell=1}^L) := \frac{1}{2} \|\mathbf{s}_1 \odot \dots \odot \mathbf{s}_L - \mathbf{s}_*\|_2^2, \quad (5)$$

where $\mathbf{s}_* \in \mathbb{R}^d$ be an r -sparse vector with ordered coordinates such that $s_{*,1} > \dots > s_{*,d}$ and define $S_p := Ls_{*,p}^{2-\frac{2}{L}}$ and $\alpha' := \left(\ln \left(\frac{2\sqrt{2}}{\eta L s_{*,1}^{2-\frac{2}{L}}} \right) \cdot \frac{s_{*,1}^{\frac{4}{L}}}{L^2 \cdot 2^{\frac{2L-3}{L}}} \right)^{\frac{1}{4}}$. For any $p < r - 1$ and $\alpha < \alpha'$, suppose we run GD on Equation (5) with learning rate $\eta = \frac{2}{K}$, where $S_p \geq K > S_{p+1}$ with

¹We exclude the case $\sigma_i(\Sigma_{L:1}) = \sigma_{*,i}$, as this occurs with probability 0 in the presence of oscillations in the EOS regime.

initialization $\mathbf{s}_\ell = \alpha \mathbf{1}_d$ for all $\ell \in [L - 1]$ and $\mathbf{s}_L = \mathbf{0}_d$. Then, under strict balancing, the top- p coordinates of \mathbf{s}_ℓ oscillate within a 2-period fixed orbit around the minima in the form

$$s_{\ell,i}(t) = \rho_{i,j}(t), \quad \forall i < p, \forall \ell \in [L],$$

where $\rho_{i,j}(t) \in \{\rho_{i,1}, \rho_{i,2}\}$, $\rho_{i,1} \in (0, s_{\star,i}^{1/L})$ and $\rho_{i,2} \in (s_{\star,i}^{1/L}, (2s_{\star,i})^{1/L})$ are two real roots of the polynomial $h(\rho) = 0$:

$$h(\rho) = \rho^L \cdot \frac{1 + (1 + \eta L(s_{\star,i} - \rho^L) \cdot \rho^{L-2})^{2L-1}}{1 + (1 + \eta L(s_{\star,i} - \rho^L) \cdot \rho^{L-2})^{L-1}} - s_{\star,i}.$$

Remarks. Similar to Theorem 1, each coordinate of the diagonal linear network undergoes periodic oscillations with an appropriately chosen learning rate. However, as shown in the proof of Theorem 3, the main difference lies in the eigenvalues themselves – the number non-zero eigenvalues for diagonal linear networks are much smaller than those of the DLN. The set of eigenvalues corresponding to the interaction with other singular values in Lemma 1 are missing for diagonal linear networks. Therefore, the top two eigenvalues of the diagonal linear network are $Ls_{\star,1}^{2-\frac{2}{L}}$ and $Ls_{\star,2}^{2-\frac{2}{L}}$, whereas the top two eigenvalues of the DLN are $Ls_{\star,1}^{2-\frac{2}{L}}$ and $\sum_{\ell=0}^{L-1} \left(s_{\star,1}^{1-\frac{1}{L}-\frac{1}{L}\ell} \cdot s_{\star,2}^{\frac{1}{L}\ell} \right)^2$. Hence, the dynamics in the EOS regime between the two networks differ significantly (see Figure 5). The primary difference in the landscape arises from the zero off-diagonal elements of the singular value diagonal matrix of the DLN which introduces additional curvature directions despite singular vector invariance.

4 EXPERIMENTAL RESULTS

Section 4.1 presents experiments that corroborate our theory. Section 4.2 discusses (i) phenomena in non-linear networks currently unexplained in the literature and how our theory can account for them in deep linear networks and (ii) how landscape in DLNs behave at EOS compared to more complicated non-convex landscapes.

4.1 SUBSPACE OSCILLATIONS IN DEEP NETWORKS

Firstly, we provide experimental results corroborating Theorem 1 and Theorem 3 to highlight their differences. We let the target matrix be $\mathbf{M}_\star \in \mathbb{R}^{50 \times 50}$ with rank 3, with dominant singular values $\sigma_\star = 10, 9, 6$. For the DLN, we consider a 3-layer network, with each layer as $\mathbf{W}_\ell \in \mathbb{R}^{50 \times 50}$ and use an initialization scale of $\alpha = 0.01$. For the diagonal linear network, we consider a similar setup, with initialization scale $\alpha = 0.01$ and the top-3 elements of $\mathbf{s}_\star \in \mathbb{R}^{50}$ to be 10, 9, 6. In Figure 5, we present the behaviors of both end-to-end networks under different learning rate regimes. By Theorem 1 and Theorem 3, the largest eigenvalue λ_1 is the same for both networks, and thus both undergo oscillations in the largest component when $\eta > 2/\lambda_1$. However, the difference in the loss landscape plays a role in oscillations for the second component. When $\eta > 2/\lambda_2$, where λ_2 is the second largest eigenvalue of the DLN, the diagonal linear network does not experience oscillations in the second component, as its second largest eigenvalue is much smaller. This highlights a distinction between the two networks and how the landscape changes between them. In Figure 6, we present an experiment demonstrating the relationship between the range of oscillations and the learning rate by plotting the amplitude of the singular value oscillations in the end-to-end network. We see that there exists no oscillations when $\eta < 2/\lambda_1$, but begins to occur at $\eta > 2/\lambda_1$ with increasing amplitude.

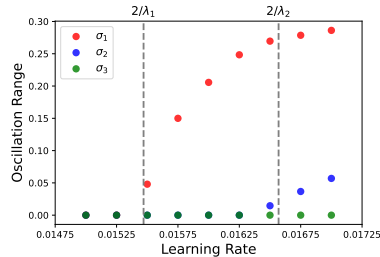


Figure 6: Oscillation range as a function of the learning rate.

4.2 SIMILARITIES AND DIFFERENCES BETWEEN LINEAR AND NONLINEAR NETS AT EOS

Mild Sharpening. “Mild” sharpening refers to the sharpness not rising to $2/\eta$ throughout learning, and generally occurs in tasks with low complexity as discussed in Caveat 2 of (Cohen et al., 2021).

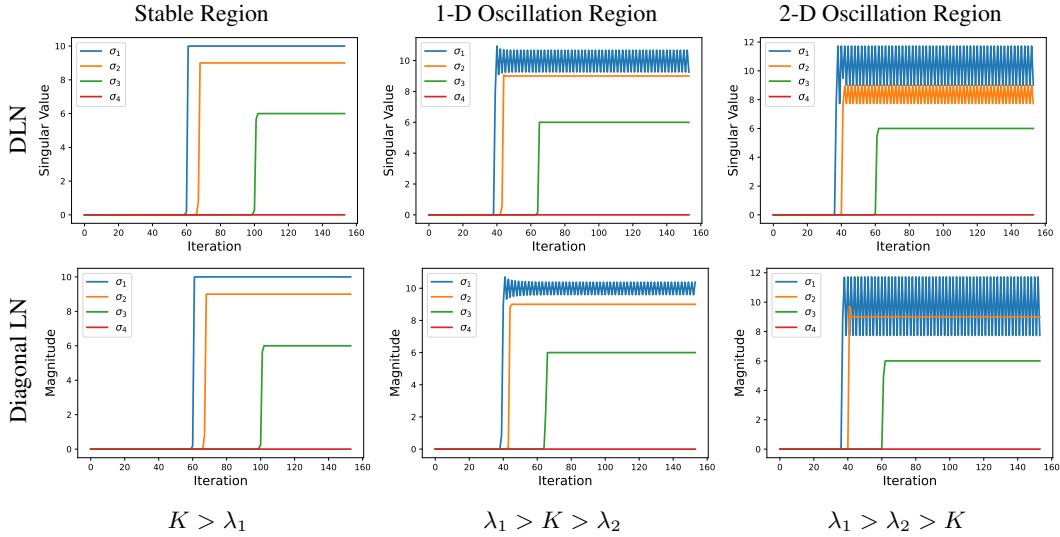


Figure 5: Dynamics of the singular values of the end-to-end DLN (top) and diagonal linear network (bottom) trained using a learning rate $\eta = 2/K$. $\lambda_1 = L\sigma_{\star,1}^{2-2/L}$ and $\lambda_2 = \sum_{\ell=0}^{L-1} \left(\sigma_{\star,1}^{1-\frac{1}{L}-\frac{1}{L}\ell} \cdot \sigma_{\star,2}^{\frac{1}{L}\ell} \right)^2$ are the top two eigenvalues of the Hessian of the training loss at convergence for the DLN. By Theorem 1, the DLN has corresponding oscillations for these two regimes, while the diagonal linear network does not due to the difference in the curvature of the landscape.

We illustrate mild sharpening in Figure 9, where we plot sharpness in two settings: (i) regression with simple images and (ii) classification with an MLP using a subset of the CIFAR-10 dataset.

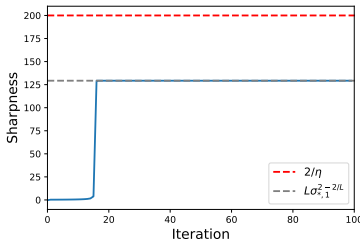


Figure 7: DLNs do not enter EOS regime if $L\sigma_{\star,1}^{2-\frac{2}{L}} < 2/\eta$.

For the regression task, we minimize the loss $\mathcal{L}(\Theta) = \|G(\Theta) - \mathbf{y}_{\text{image}}\|_2^2$, where $G(\Theta)$ is a UNet parameterized by Θ , and $\mathbf{y}_{\text{image}}$ denotes one of the images in Figure 9b. We observe that when $\mathbf{y}_{\text{image}}$ is a smooth, low-frequency image, the sharpness of the loss generally remains low. However, when $\mathbf{y}_{\text{image}}$ has higher frequency content, the sharpness increases and enters the EOS regime (Figure 9a). Similarly, for the classification task, we train a 2-layer fully connected neural network on N labeled training images from the CIFAR-10 dataset using MSE loss and plot the sharpness in Figure 9c. The sharpness links to N , the number of data points used for training. For small N values, such as 100 or 200, the network learns only a limited set of latent features, resulting in mild sharpening, and it does not reach the EOS threshold. However, when N exceeds 1000, the sharpness increases and reaches the EOS threshold. The intrinsic dimension update in neural networks for such low complexity tasks is usually smaller (Li et al., 2018) which could cause the sharpness to be small. Similar observations can also be seen in DLNs. In Figure 7, we show that the sharpness reaches $L\sigma_{\star,1}^{2-\frac{2}{L}}$, where $\sigma_{\star,1}$ is the singular value of the target matrix. Whenever $L\sigma_{\star,1}^{2-\frac{2}{L}} < 2/\eta$, the network will not enter the EOS regime. This can be viewed as low-complexity learning, as $\sigma_{\star,1}$ corresponds to the magnitude of the strongest feature of the target matrix. Hence, when $\sigma_{\star,1}$ is not large enough, the sharpness will not rise to $2/\eta$. While these observations do not fully explain mild sharpening, our experiments demonstrate that interpreting sharpness as a measure of complexity, combined with our findings from DLNs, marks an important first step toward fully understanding this phenomenon.

Difference in Oscillation Behaviors. Here, we discuss the differences in oscillations that arise in DLNs compared to catapults that occur in practical deep nonlinear networks. The main difference

486
487
488
489
490
491
492
493
494
495
496
497
498
499
500
501
502
503
504
505
506
507
508
509
510
511
512
513
514
515
516
517
518
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539

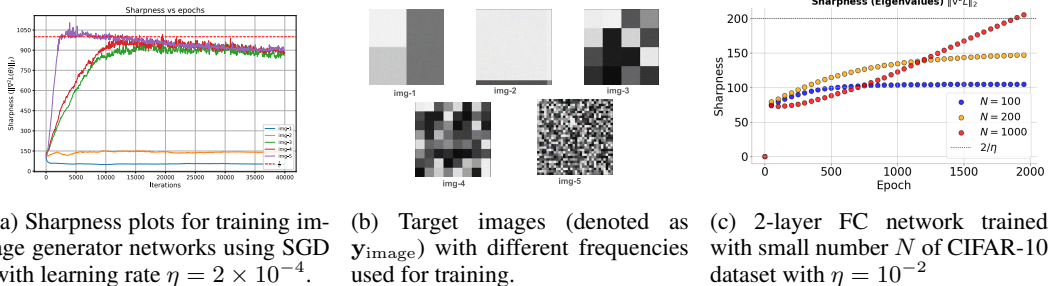
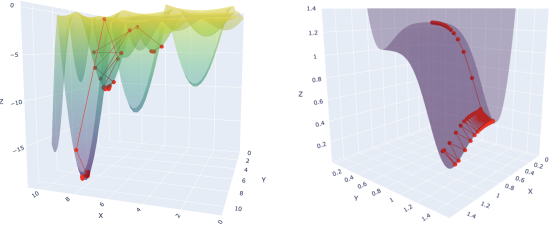


Figure 9: Illustration of Caveat 2 by Cohen et al. (2021) on how mild sharpening occurs on simple datasets and network. (a) Regression task showing the evolution of the sharpness when an UNet (with fixed initialization) is trained to fit a single image shown in (b). (c) Evolution of the minimal progressive sharpening on a classification task of a 2-layer MLP trained on a subset of CIFAR-10.

lies in the loss landscape—at convergence, the Hessian for DLNs is positive semi-definite, as shown in Lemma 1, meaning there are only directions of positive curvature and flat directions (in the null space of the Hessian). In this landscape, oscillations occur because the basin walls bounce off, without direction of escape. However, in deep nonlinear networks, it has been frequently observed that the Hessian at the minima has negative eigenvalues (Ghorbani et al., 2019; Sagun et al., 2016). This enables an escape direction along the negative curvature, preventing sustained oscillations. In Figure 8, we demonstrate these two differences by visualizing the loss landscapes and the iterates throughout GD marked in red. The Holder table function Figure 8 (left) exhibits numerous local minima, causing the loss to exhibit sharp “catapult” when a large learning rate is used. In contrast, for DLNs (shown in the right) the loss oscillates in a periodic orbit around the minima since there are no spurious local minima (Ge et al., 2016). In Appendix B.1, we provide experimental details.



Lastly, Damian et al. (2023) studies self-stabilization, where sharpness decreases below $2/\eta$ after initially exceeding $2/\eta$. Their analysis on self-stabilization requires certain assumptions such as $\nabla L(\theta) \cdot u(\theta) = 0$ and $\nabla S(\theta)$ lies in the null space of the Hessian, where $S(\theta)$ and $u(\theta)$ denotes the maximum eigenvalue and its corresponding eigenvector respectively. These assumptions do not hold exactly in DLNs. Rather, the sharpness oscillates about $2/\eta$ as shown in Figure 2 as the condition for stable oscillation holds along each eigenvector of the Hessian.

5 CONCLUSION, LIMITATIONS AND FUTURE WORK

In this paper, we presented a fine-grained analysis of the learning dynamics of deep matrix factorization with the aim of understanding unexplained phenomena in deep nonlinear networks within the EOS regime. Our analysis revealed that within EOS, DLNs exhibit periodic oscillations in small subspaces, where the subspace dimension is exactly characterized by the learning rate. There are two limitations to our work: we require (i) the dynamics converge to the singular vector stationary set, and (ii) strict balancing of the singular values. However, we provide thorough empirical evidence validating the use of these assumptions, along with more results in Appendix B. For the balancing assumption, we leave for future work on alleviating the assumption of strict balancing, and rigorously show that this holds before entering the EOS regime.

REFERENCES

- 540
541
542 Atish Agarwala, Fabian Pedregosa, and Jeffrey Pennington. Second-order regression models exhibit
543 progressive sharpening to the edge of stability. *arXiv preprint arXiv:2210.04860*, 2022.
- 544 Kwangjun Ahn, Sébastien Bubeck, Sinho Chewi, Yin Tat Lee, Felipe Suarez, and Yi Zhang. Learn-
545 ing threshold neurons via edge of stability. *Advances in Neural Information Processing Systems*,
546 36, 2024.
- 547 Sanjeev Arora, Nadav Cohen, and Elad Hazan. On the optimization of deep networks: Implicit
548 acceleration by overparameterization. In *International conference on machine learning*, pp. 244–
549 253. PMLR, 2018.
- 550
551 Sanjeev Arora, Nadav Cohen, Wei Hu, and Yiping Luo. Implicit regularization in deep
552 matrix factorization. In *Advances in Neural Information Processing Systems*, volume 32,
553 2019. URL [https://proceedings.neurips.cc/paper_files/paper/2019/
554 file/c0c783b5fc0d7d808f1d14a6e9c8280d-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2019/file/c0c783b5fc0d7d808f1d14a6e9c8280d-Paper.pdf).
- 555 Sanjeev Arora, Zhiyuan Li, and Abhishek Panigrahi. Understanding gradient descent on the edge
556 of stability in deep learning. In *International Conference on Machine Learning*, pp. 948–1024.
557 PMLR, 2022.
- 558
559 Raphaël Berthier. Incremental learning in diagonal linear networks. *Journal of Machine Learning*
560 *Research*, 24(171):1–26, 2023. URL <http://jmlr.org/papers/v24/22-1395.html>.
- 561
562 Lei Chen and Joan Bruna. Beyond the edge of stability via two-step gradient updates, 2023.
- 563 Xuxing Chen, Krishnakumar Balasubramanian, Promit Ghosal, and Bhavya Agrawalla. From sta-
564 bility to chaos: Analyzing gradient descent dynamics in quadratic regression. *arXiv preprint*
565 *arXiv:2310.01687*, 2023.
- 566
567 Hung-Hsu Chou, Carsten Gieshoff, Johannes Maly, and Holger Rauhut. Gradient descent for deep
568 matrix factorization: Dynamics and implicit bias towards low rank. *Applied and Computational*
569 *Harmonic Analysis*, 68:101595, 2024a. ISSN 1063-5203. doi: [https://doi.org/10.1016/j.acha.
570 2023.101595](https://doi.org/10.1016/j.acha.2023.101595). URL [https://www.sciencedirect.com/science/article/pii/
571 S1063520323000829](https://www.sciencedirect.com/science/article/pii/S1063520323000829).
- 572
573 Hung-Hsu Chou, Carsten Gieshoff, Johannes Maly, and Holger Rauhut. Gradient descent for deep
574 matrix factorization: Dynamics and implicit bias towards low rank. *Applied and Computational*
575 *Harmonic Analysis*, 68:101595, 2024b.
- 576
577 Jeremy Cohen, Simran Kaur, Yuanzhi Li, J Zico Kolter, and Ameet Talwalkar. Gradient descent on
578 neural networks typically occurs at the edge of stability. In *International Conference on Learning*
579 *Representations*, 2021. URL <https://openreview.net/forum?id=jh-rTtvkGeM>.
- 580
581 Alex Damian, Eshaan Nichani, and Jason D. Lee. Self-stabilization: The implicit bias of gradient
582 descent at the edge of stability. In *The Eleventh International Conference on Learning Represen-*
583 *tations*, 2023. URL <https://openreview.net/forum?id=nhKHA59gXz>.
- 584
585 Mathieu Even, Scott Pesme, Suriya Gunasekar, and Nicolas Flammarion. (S)GD over diagonal lin-
586 ear networks: Implicit bias, large stepsizes and edge of stability. *Advances in Neural Information*
587 *Processing Systems*, 36, 2024.
- 588
589 Pierre Foret, Ariel Kleiner, Hossein Mobahi, and Behnam Neyshabur. Sharpness-aware minimiza-
590 tion for efficiently improving generalization. In *International Conference on Learning Represen-*
591 *tations*, 2021. URL <https://openreview.net/forum?id=6TmlmposlrM>.
- 592
593 Khashayar Gatmiry, Zhiyuan Li, Tengyu Ma, Sashank J. Reddi, Stefanie Jegelka, and Ching-Yao
Chuang. What is the inductive bias of flatness regularization? A study of deep matrix factorization
models. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023. URL
<https://openreview.net/forum?id=2hQ7MBQApp>.
- Rong Ge, Jason D Lee, and Tengyu Ma. Matrix completion has no spurious local minimum. *Ad-*
vances in neural information processing systems, 29, 2016.

- 594 Behrooz Ghorbani, Shankar Krishnan, and Ying Xiao. An investigation into neural net optimization
595 via hessian eigenvalue density. In *International Conference on Machine Learning*, pp. 2232–
596 2241. PMLR, 2019.
- 597 Gauthier Gidel, Francis Bach, and Simon Lacoste-Julien. Implicit regularization of discrete gradient
598 dynamics in linear neural networks. In *Advances in Neural Information Processing Systems*, vol-
599 ume 32, 2019a. URL [https://proceedings.neurips.cc/paper_files/paper/
600 2019/file/f39ae9ff3a81f499230c4126e01f421b-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2019/file/f39ae9ff3a81f499230c4126e01f421b-Paper.pdf).
- 601 Gauthier Gidel, Francis Bach, and Simon Lacoste-Julien. Implicit regularization of discrete gradient
602 dynamics in linear neural networks. *Advances in Neural Information Processing Systems*, 32,
603 2019b.
- 604 Daniel Gissin, Shai Shalev-Shwartz, and Amit Daniely. The implicit bias of depth: How incremental
605 learning drives generalization. In *International Conference on Learning Representations*, 2020.
606 URL <https://openreview.net/forum?id=H11j0nNFwB>.
- 607 Soufiane Hayou, Nikhil Ghosh, and Bin Yu. LoRA+: Efficient low rank adaptation of large
608 models. In *Forty-first International Conference on Machine Learning*, 2024. URL <https://openreview.net/forum?id=NEv8YqBROO>.
- 609 Edward J Hu, yelong shen, Phillip Wallis, Zeyuan Allen-Zhu, Yanzhi Li, Shean Wang, Lu Wang,
610 and Weizhu Chen. LoRA: Low-rank adaptation of large language models. In *International Con-
611 ference on Learning Representations*, 2022. URL [https://openreview.net/forum?
612 id=nZeVKeeFYf9](https://openreview.net/forum?id=nZeVKeeFYf9).
- 613 Pavel Izmailov, Dmitrii Podoprikin, Timur Garipov, Dmitry Vetrov, and Andrew Gordon Wil-
614 son. Averaging weights leads to wider optima and better generalization. *arXiv preprint
615 arXiv:1803.05407*, 2019. URL <https://arxiv.org/abs/1803.05407>.
- 616 Arthur Jacot, François Ged, Berfin Şimşek, Clément Hongler, and Franck Gabriel. Saddle-to-saddle
617 dynamics in deep linear networks: Small initialization training, symmetry, and sparsity. *arXiv
618 preprint arXiv:2106.15933*, 2021.
- 619 Arthur Jacot, François Ged, Berfin Şimşek, Clément Hongler, and Franck Gabriel. Saddle-to-saddle
620 dynamics in deep linear networks: Small initialization training, symmetry, and sparsity. *arXiv
621 preprint arXiv:2106.15933*, 2022.
- 622 Stanisław Jastrzebski, Zachary Kenton, Nicolas Ballas, Asja Fischer, Yoshua Bengio, and Amos
623 Storkey. On the relation between the sharpest directions of dnn loss and the sgd step length. *arXiv
624 preprint arXiv:1807.05031*, 2018.
- 625 Stanislaw Jastrzebski, Maciej Szymczak, Stanislav Fort, Devansh Arpit, Jacek Tabor, Kyunghyun
626 Cho, and Krzysztof Geras. The break-even point on optimization trajectories of deep neural
627 networks. *arXiv preprint arXiv:2002.09572*, 2020.
- 628 Dayal Singh Kalra, Tianyu He, and Maissam Barkeshli. Universal sharpness dynamics in neural
629 network training: Fixed point analysis, edge of stability, and route to chaos. *arXiv preprint
630 arXiv:2311.02076*, 2023.
- 631 Nitish Shirish Keskar, Dheevatsa Mudigere, Jorge Nocedal, Mikhail Smelyanskiy, and Ping Tak Pe-
632 ter Tang. On large-batch training for deep learning: Generalization gap and sharp minima. In
633 *International Conference on Learning Representations*, 2017. URL [https://openreview.
634 net/forum?id=H1oyRlYgg](https://openreview.net/forum?id=H1oyRlYgg).
- 635 Itai Kreisler, Mor Shpigel Nacson, Daniel Soudry, and Yair Carmon. Gradient descent monoton-
636 ically decreases the sharpness of gradient flow solutions in scalar networks and beyond. In *Inter-
637 national Conference on Machine Learning*, pp. 17684–17744. PMLR, 2023.
- 638 Aitor Lewkowycz, Yasaman Bahri, Ethan Dyer, Jascha Sohl-Dickstein, and Guy Gur-Ari. The large
639 learning rate phase of deep learning: the catapult mechanism. *arXiv preprint arXiv:2003.02218*,
640 2020.

- 648 Chunyuan Li, Heerad Farkhoor, Rosanne Liu, and Jason Yosinski. Measuring the intrinsic dimension
649 of objective landscapes. *arXiv preprint arXiv:1804.08838*, 2018.
650
- 651 Sheng Liu, Zhihui Zhu, Qing Qu, and Chong You. Robust training under label noise by over-
652 parameterization. In *International Conference on Machine Learning*, pp. 14153–14172. PMLR,
653 2022.
- 654 Kaifeng Lyu, Zhiyuan Li, and Sanjeev Arora. Understanding the generalization benefit of normal-
655 ization layers: Sharpness reduction. *Advances in Neural Information Processing Systems*, 35:
656 34689–34708, 2022.
- 657 Pierre Marion and Lénaïc Chizat. Deep linear networks for regression are implicitly regularized
658 towards flat minima. *arXiv preprint arXiv:2405.13456*, 2024.
659
- 660 Soo Min Kwon, Zekai Zhang, Dogyoon Song, Laura Balzano, and Qing Qu. Efficient low-
661 dimensional compression of overparameterized models. In *Proceedings of The 27th International
662 Conference on Artificial Intelligence and Statistics*, volume 238 of *Proceedings of Machine Learn-
663 ing Research*, pp. 1009–1017. PMLR, 02–04 May 2024. URL [https://proceedings.
664 mlr.press/v238/min-kwon24a.html](https://proceedings.mlr.press/v238/min-kwon24a.html).
- 665 Leon Mirsky. A trace inequality of John von Neumann. *Monatshefte für mathematik*, 79(4):303–
666 306, 1975.
667
- 668 Edward Ott. *Chaos in Dynamical Systems*. Cambridge University Press, 2 edition, 2002.
- 669 Scott Pesme and Nicolas Flammarion. Saddle-to-saddle dynamics in diagonal linear networks. *Ad-
670 vances in Neural Information Processing Systems*, 36:7475–7505, 2023.
671
- 672 Henning Petzka, Michael Kamp, Linara Adilova, Cristian Sminchisescu, and Mario Boley. Relative
673 flatness and generalization. In *Advances in Neural Information Processing Systems*, 2021. URL
674 https://openreview.net/forum?id=sygvo7ctb_.
- 675 Levent Sagun, Leon Bottou, and Yann LeCun. Eigenvalues of the hessian in deep learning: Singu-
676 larity and beyond. *arXiv preprint arXiv:1611.07476*, 2016.
677
- 678 Andrew M. Saxe, James L. McClelland, and Surya Ganguli. Exact solutions to the nonlinear dy-
679 namics of learning in deep linear neural networks. In *2nd International Conference on Learning
680 Representations, ICLR, 2014*. URL <http://arxiv.org/abs/1312.6120>.
- 681 Minhak Song and Chulhee Yun. Trajectory alignment: understanding the edge of stability phe-
682 nomenon via bifurcation theory. *arXiv preprint arXiv:2307.04204*, 2023.
683
- 684 Aditya Vardhan Varre, Maria-Luiza Vladarean, Loucas PILLAUD-VIVIEN, and Nico-
685 las Flammarion. On the spectral bias of two-layer linear networks. In *Advances in
686 Neural Information Processing Systems*, volume 36, pp. 64380–64414, 2023. URL
687 [https://proceedings.neurips.cc/paper_files/paper/2023/file/
688 cad2fd66cf88226d868f90a7cbaa4a53-Paper-Conference.pdf](https://proceedings.neurips.cc/paper_files/paper/2023/file/cad2fd66cf88226d868f90a7cbaa4a53-Paper-Conference.pdf).
- 689 Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel R. Bowman.
690 GLUE: A multi-task benchmark and analysis platform for natural language understanding. In
691 *International Conference on Learning Representations*, 2019. URL [https://openreview.
692 net/forum?id=rJ4km2R5t7](https://openreview.net/forum?id=rJ4km2R5t7).
- 693 Peng Wang, Xiao Li, Can Yaras, Zhihui Zhu, Laura Balzano, Wei Hu, and Qing Qu. Understanding
694 deep representation learning via layerwise feature compression and discrimination. *arXiv preprint
695 arXiv:2311.02960*, 2024. URL <https://arxiv.org/abs/2311.02960>.
- 696 Yuqing Wang, Minshuo Chen, Tuo Zhao, and Molei Tao. Large learning rate tames homogeneity:
697 Convergence and balancing effect. *arXiv preprint arXiv:2110.03677*, 2021.
698
- 699 Yuqing Wang, Zhenghao Xu, Tuo Zhao, and Molei Tao. Good regularity creates large learning rate
700 implicit biases: edge of stability, balancing, and catapult. In *NeurIPS 2023 Workshop on Math-
701 ematics of Modern Machine Learning*, 2023. URL [https://openreview.net/forum?
id=6015A3h2yl](https://openreview.net/forum?id=6015A3h2yl).

- 702 Zixuan Wang, Zhouzi Li, and Jian Li. Analyzing sharpness along gd trajectory: Progressive sharp-
703 ening and edge of stability. *Advances in Neural Information Processing Systems*, 35:9983–9994,
704 2022.
- 705
706 Jingfeng Wu, Vladimir Braverman, and Jason D Lee. Implicit bias of gradient descent for logistic
707 regression at the edge of stability. *Advances in Neural Information Processing Systems*, 36, 2024.
- 708
709 Zhenghao Xu, Yuqing Wang, Tuo Zhao, Rachel Ward, and Molei Tao. Provable acceleration of
710 nesterov’s accelerated gradient for rectangular matrix factorization and linear neural networks.
711 *arXiv preprint arXiv:2410.09640*, 2024.
- 712
713 Can Yaras, Peng Wang, Wei Hu, Zhihui Zhu, Laura Balzano, and Qing Qu. The law of parsimony
714 in gradient descent for learning deep linear networks. *arXiv preprint arXiv:2306.01154*, 2023.
- 715
716 Can Yaras, Peng Wang, Laura Balzano, and Qing Qu. Compressible dynamics in deep overpa-
717 rameterized low-rank learning & adaptation. In *Forty-first International Conference on Machine*
718 *Learning*, 2024. URL <https://openreview.net/forum?id=uDkXoZMzBv>.
- 719
720 Tian Ye and Simon S Du. Global convergence of gradient descent for asymmetric low-rank matrix
721 factorization. *Advances in Neural Information Processing Systems*, 34:1429–1439, 2021.
- 722
723 Chong You, Zhihui Zhu, Qing Qu, and Yi Ma. Robust recovery via implicit bias of discrepant learn-
724 ing rates for double over-parameterization. *Advances in Neural Information Processing Systems*,
725 33:17733–17744, 2020.
- 726
727 Xitong Zhang, Ismail R Alkhouri, and Rongrong Wang. Structure-preserving network com-
728 pression via low-rank induced training through linear layers composition. *arXiv preprint*
729 *arXiv:2405.03089*, 2024a.
- 730
731 Yedi Zhang, Andrew M Saxe, and Peter E. Latham. When are bias-free reLU networks like linear
732 networks? In *High-dimensional Learning Dynamics 2024: The Emergence of Structure and*
733 *Reasoning*, 2024b. URL <https://openreview.net/forum?id=LdYBMeWOG3>.
- 734
735 Libin Zhu, Chaoyue Liu, Adityanarayanan Radhakrishnan, and Mikhail Belkin. Quadratic models
736 for understanding neural network dynamics. *arXiv preprint arXiv:2205.11787*, 2022.
- 737
738 Libin Zhu, Chaoyue Liu, Adityanarayanan Radhakrishnan, and Mikhail Belkin. Catapults in SGD:
739 spikes in the training loss and their impact on generalization through feature learning. *arXiv*
740 *preprint arXiv:2306.04815*, 2023a.
- 741
742 Xingyu Zhu, Zixuan Wang, Xiang Wang, Mo Zhou, and Rong Ge. Understanding edge-of-
743 stability training dynamics with a minimalist example. In *The Eleventh International Confer-*
744 *ence on Learning Representations*, 2023b. URL <https://openreview.net/forum?id=p7EagBsMAEO>.
- 745
746
747
748
749
750
751
752
753
754
755

Appendix

CONTENTS

1	Introduction	1
2	Notation and Problem Setup	3
3	Deep Matrix Factorization Beyond the Edge of Stability	3
3.1	Main Results	4
3.2	Tools used in the Analyses	6
3.3	Relation to Diagonal Linear Networks	7
4	Experimental Results	8
4.1	Subspace Oscillations in Deep Networks	8
4.2	Similarities and Differences Between Linear and Nonlinear Nets at EOS	8
5	Conclusion, Limitations and Future Work	10
A	Discussion on Related Work	15
B	Additional Results	16
B.1	Experimental Details	16
B.2	Initialization Outside Singular Vector Invariant Set	18
B.3	Balancing of Singular Values	18
B.4	Additional Experiments for Balancing, Singular Vector Invariance, and Theory	19
B.5	Periodic and Free Oscillations	20
B.6	Investigation of Oscillations in Low-Rank Adaptors	22
C	Deferred Proofs	26
C.1	Deferred Proofs for Oscillations	26
C.2	Deferred Proofs for Singular Vector Invariance	49
C.3	Auxiliary Results	51

A DISCUSSION ON RELATED WORK

Implicit Bias of Edge of Stability. Edge of stability was first coined by Cohen et al. (2021), where they showed that the Hessian of the training loss plateaus around $2/\eta$ when deep models were trained using GD. However, Jastrzebski et al. (2020); Jastrzebski et al. (2018) previously demonstrated that the step size influences the sharpness along the optimization trajectory. Due to the important practical implications of the edge of stability, there has been an explosion of research dedicated to understanding this phenomenon and its implicit regularization properties. Here, we survey a few

of these works. Damian et al. (2023) explained edge of stability through a mechanism called “self-stabilization”, where they showed that during the momentary divergence of the iterates along the sharpest eigenvector direction of the Hessian, the iterates also move along the negative direction of the gradient of the curvature, which leads to stabilizing the sharpness to $2/\eta$. Agarwala et al. (2022) proved that second-order regression models (the simplest class of models after the linearized NTK model) demonstrate progressive sharpening of the NTK eigenvalue towards a slightly different value than $2/\eta$. Arora et al. (2022) mathematically analyzed the edge of stability, where they showed that the GD updates evolve along some deterministic flow on the manifold of the minima. Lyu et al. (2022) showed that the normalization layers had an important role in the edge of stability – they showed that these layers encouraged GD to reduce the sharpness of the loss surface and enter the EOS regime. Ahn et al. (2024) established the phenomenon in two-layer networks and find phase transitions for step-sizes in which networks fail to learn “threshold” neurons. Wang et al. (2022) also analyze a two-layer network, but provide a theoretical proof for the change in sharpness across four different phases. Even et al. (2024) analyzed the edge of stability in diagonal linear networks and found that oscillations occur on the sparse support of the vectors. Lastly, Wu et al. (2024) analyzed the convergence at the edge of stability for constant step size GD for logistic regression on linearly separable data.

Edge of Stability in Toy Functions. To analyze the edge of stability in slightly simpler settings, many works have constructed scalar functions to analyze the prevalence of this phenomenon. For example, Chen & Bruna (2023) studied a certain class of scalar functions and identified conditions in which the function enters the edge of stability through a two-step convergence analysis. Wang et al. (2023) showed that the edge of stability occurs in specific scalar functions, which satisfies certain regularity conditions and developed a global convergence theory for a family of non-convex functions without globally Lipschitz continuous gradients. Zhu et al. (2023b) analyzed local oscillatory behaviors for 4-layer scalar networks with balanced initialization. Song & Yun (2023); Kalra et al. (2023) provide analyses of learning dynamics at the EOS in simplified settings such as two-layer networks. Zhu et al. (2022); Chen et al. (2023) study GD dynamics for quadratic models in large learning rate regimes. Overall, all of these works showed that the necessary condition for the edge of stability to occur is that the second derivative of the loss function is non-zero, even though they assumed simple scalar functions. Our work takes one step further to analyze the prevalence of the edge of stability in DLNs. Although our loss simplifies to a loss in terms of the singular values, they precisely characterize the dynamics of the DLNs for the deep matrix factorization problem.

Deep Linear Networks. Over the past decade, many existing works have analyzed the learning dynamics of DLNs as a surrogate for deep nonlinear networks to study the effects of depth and implicit regularization (Saxe et al., 2014; Arora et al., 2018; 2019; Zhang et al., 2024a). Generally, these works focus on unveiling the dynamics of a phenomenon called “incremental learning”, where small initialization scales induce a greedy singular value learning approach (Min Kwon et al., 2024; Gissin et al., 2020; Saxe et al., 2014), analyzing the learning dynamics via gradient flow (Saxe et al., 2014; Chou et al., 2024a; Arora et al., 2019), or showing that the DLN is biased towards low-rank solution (Yaras et al., 2024; Arora et al., 2019; Min Kwon et al., 2024), amongst others. However, these works do not consider the occurrence of the edge of stability in such networks. On the other hand, while works such as those by Yaras et al. (2024) and Min Kwon et al. (2024) have similar observations in that the weight updates occur within an invariant subspace as shown by Proposition 2, they do not analyze the edge of stability regime.

B ADDITIONAL RESULTS

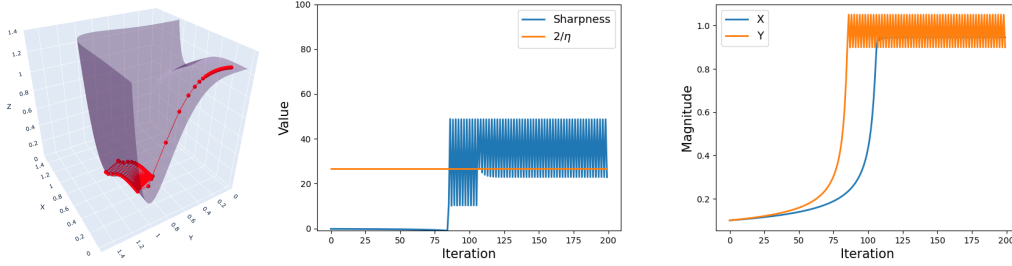
B.1 EXPERIMENTAL DETAILS

In this section, we provide additional details regarding the experiments used to generate the figures in the main text. For Figure 1, we consider a rank-3 target matrix $M_* \in \mathbb{R}^{5 \times 5}$ with ordered singular values 10, 6, 3. We use a 3-layer DLN to fit the target matrix. Since $\sigma_{*,1} = 10$, the network enters the EOS regime at

$$\eta = \frac{2}{L\sigma_{*,1}^{2-2/L}} = 0.0309.$$

864
865
866
867
868
869
870
871
872
873
874
875
876
877
878
879
880
881
882
883
884
885
886
887
888
889
890
891
892
893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
910
911
912
913
914
915
916
917

Oscillation along Y-axis: $2/\lambda_2 > \eta > 2/\lambda_1$



Oscillation along both X and Y-axis: $\eta > 2/\lambda_2$

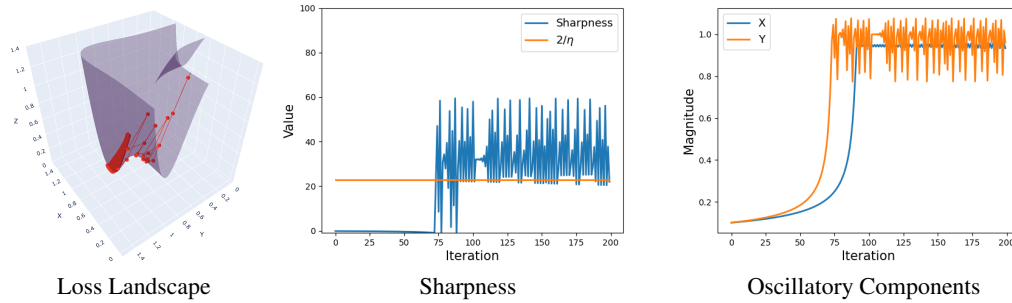


Figure 10: Demonstration of the EOS dynamics of a 2-dimensional depth-4 scalar network as shown in Equation (6). X, Y axes are the eigenvectors of the Hessian with eigenvalues λ_1 and λ_2 respectively. Top: when $\eta > 2/\lambda_1$, the X component remains fixed, while the Y component oscillates with a periodicity of 2. Bottom: for $\eta > 2/\lambda_2$, the iterates oscillation in both directions.

We show that there exists a two-period orbit after $0.0309/2 = 0.0154$, as we do not have a scaling of $1/2$ in the objective function for the code used to generate the figures.

In Figure 8 and 10, we compared the landscape of DLNs with that of a more complicated non-convex function such as the Holder table function. To mimic the DLN, we considered the loss function

$$z = L(x, y) = (x^4 - 0.8)^2 + (y^4 - 1)^2, \quad (6)$$

which corresponds for a 4-layer network. Here the eigenvector of the Hessian at the global minima coincides with the x, y -axis. We calculate the eigenvalues λ_1 and λ_2 at the minimum $(0.8^{0.25}, 1)$ and plot the dynamics of the iterate for step size range $\frac{2}{\lambda_2} > \eta > \frac{2}{\lambda_1}$ and $\eta > \frac{2}{\lambda_2}$. When $\frac{2}{\lambda_2} > \eta > \frac{2}{\lambda_1}$ the x -coordinate stays fixed at the minima $0.8^{0.25}$ and the y -coordinate oscillates around its minimum at $y = 1$. This is evident in the landscape figure. Similarly, when $\eta > \frac{2}{\lambda_2}$, oscillations occur in both the x and y direction. The loss landscape $z = L(x, y)$ does not have spurious local minima, so sustained oscillations take place in the loss basin.

For the non-convex landscape as shown in Figure 8 and 11, we consider the Holder table function:

$$f(x, y) = - \left| \sin(x) \cos(y) \exp \left(1 - \frac{\sqrt{x^2 + y^2}}{\pi} \right) \right|.$$

By observation, we initialize near a sharp minima and run GD with an increasing learning rate step size as shown in the lefthand side of Figure 11. When the learning rate is fixed, we observe that oscillations take place inside the local valley, but when learning rate is increased, it jumps out of the local valley to find a flatter basin. Similar to the observations by Cohen et al. (2021), the sharpness of the GD iterates are “regulated” by the threshold $2/\eta$, as it seems to closely follow this value as shown in Figure 11.

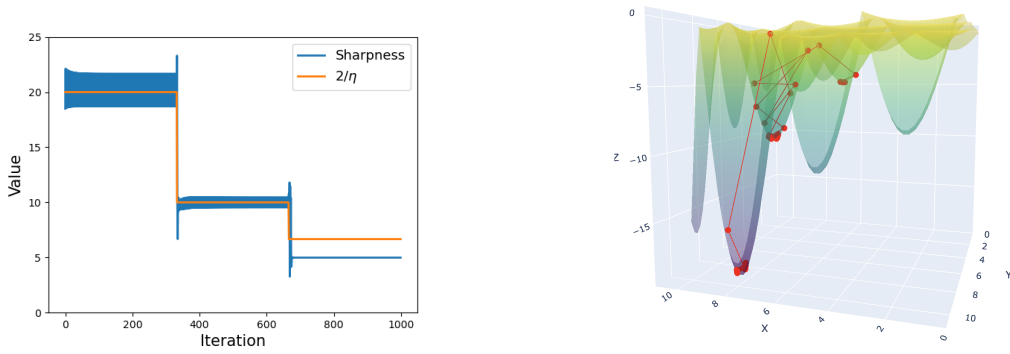


Figure 11: EOS dynamics at various step learning rates from the Holder table function. Left: plot of the learning rate steps and sharpness, showing that sharpness follows the EOS limit $2/\eta$. Right: Plot showing that the iterates catapult out of a local basin when the learning rate is increased and jumps out to a surface where the sharpness is about $2/\eta$.

Overall, these examples aim to highlight the difference in linear and complex loss landscapes. The former consists of *only* saddles and global minima, and hence (stably) oscillate about the global minimum. However, in more complicated non-convex landscapes, sharpness regularization due to large learning rates enable catapulting to flatter loss basins, where sharpness is smaller than $2/\eta$.

B.2 INITIALIZATION OUTSIDE SINGULAR VECTOR INVARIANT SET

In this section, we present an initialization example that is outside the Singular vector stationary set. We consider the following initialization:

$$\mathbf{W}_L(0) = \mathbf{0}, \quad \mathbf{W}_\ell(0) = \alpha \mathbf{P}_\ell, \quad \forall \ell \in [L-1], \quad (7)$$

where $\mathbf{P}_\ell \in \mathbb{R}^{d \times d}$ is an orthogonal matrix. Note that here for $\ell > 1$, the singular vectors do not align and lies outside the SVS set we defined in Proposition 2. We consider the deep matrix factorization problem with a target matrix $\mathbf{M}_* \in \mathbb{R}^{d \times d}$, where $d = 100$, $r = 5$, and $\alpha = 0.01$. We empirically obtain that the decomposition after convergence admits the form:

$$\mathbf{W}_L(t) = \mathbf{U}^* \begin{bmatrix} \boldsymbol{\Sigma}_L(t) & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \left[\left(\prod_{i=L-1}^1 \mathbf{P}_i \right) \mathbf{V}^* \right]^\top, \quad (8)$$

$$\mathbf{W}_\ell(t) = \left[\left(\prod_{i=\ell}^1 \mathbf{P}_i \right) \mathbf{V}^* \right] \begin{bmatrix} \boldsymbol{\Sigma}_\ell(t) & \mathbf{0} \\ \mathbf{0} & \alpha \mathbf{I}_{d-r} \end{bmatrix} \left[\left(\prod_{i=\ell-1}^1 \mathbf{P}_i \right) \mathbf{V}^* \right]^\top, \quad \forall \ell \in [2, L-1], \quad (9)$$

$$\mathbf{W}_1(t) = \mathbf{P}_1 \mathbf{V}^* \begin{bmatrix} \boldsymbol{\Sigma}_1(t) & \mathbf{0} \\ \mathbf{0} & \alpha \mathbf{I}_{d-r} \end{bmatrix} \mathbf{V}^{*\top}, \quad (10)$$

where $\mathbf{W}_L(0) = \mathbf{0}$ and $\mathbf{W}_\ell(0) = \alpha \mathbf{P}_\ell, \forall \ell \in [L-1]$. The decomposition after convergence lies in the SVS set as the singular vectors now align with each other. This demonstrates an example where even when the initialization is made outside the SVS set, GD aligns the singular vectors such that after certain iterations it lies in the SVS set.

B.3 BALANCING OF SINGULAR VALUES

In this section, we present additional experimental results on Lemma 2 and how close the iterates become for different initialization scales. To this end, we consider the same setup from the previous section, where we have a target matrix $\mathbf{M}_* \in \mathbb{R}^{d \times d}$, where $d = 100$, $r = 5$, and varying initialization α . In Figure 13, we observe that for larger values of α , the balancing quickly occurs, whereas for smaller values of α , the balancing is almost immediate. This is to also highlight that our bound on α in Lemma 2 may be an artifact of our analysis, and can choose larger values of α in practice.

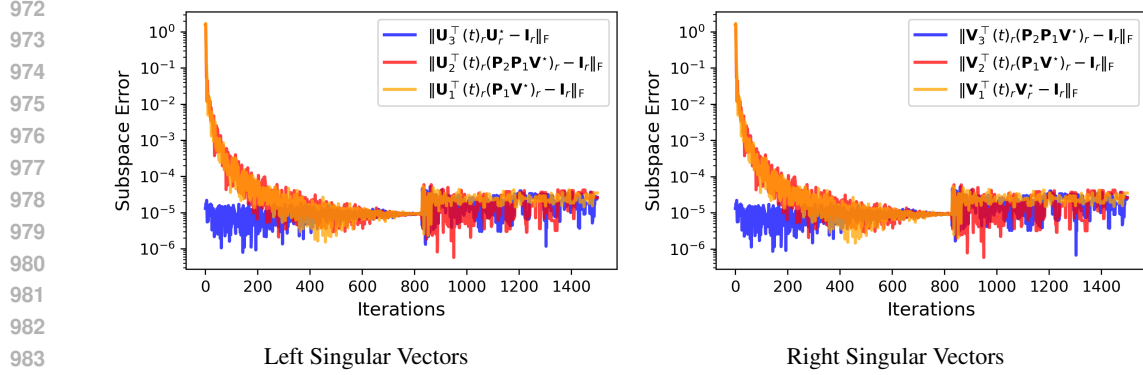


Figure 12: Empirical verification of the decomposition for initialization with orthogonal matrices (lying outside SVS set) in that after some GD iterations, the singular vectors of the intermediate matrices align to lie within SVS set, displaying singular vector invariance.

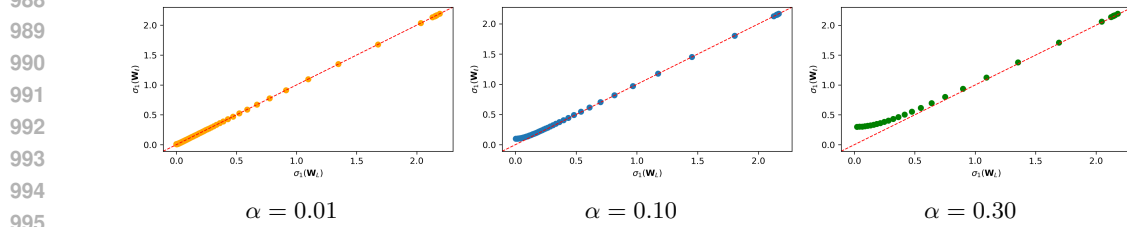


Figure 13: Observing the balancedness between the singular value initialized to 0 and a singular value initialized to α . The scattered points are successive GD iterations (going left to right). The initial gap between the two values is larger for a larger α , but quickly gets closer over more GD iterations.

B.4 ADDITIONAL EXPERIMENTS FOR BALANCING, SINGULAR VECTOR INVARIANCE, AND THEORY

Our theory relied on two tools and assumptions: balancing of singular values and stationarity of the singular vectors. In this section, we investigate how the dynamics at EOS are affected if these two assumptions do not hold.

Balancing. By Lemma 2, recall that balancing only holds as long as α chosen below a certain threshold. To this end, we consider the dynamics of a 3-layer DLN to fit a target matrix $\mathbf{M}_* \in \mathbb{R}^{10 \times 10}$ of rank-3 with ordered singular values 10, 8, 6. We use a learning rate of $\eta = 0.0166$, which corresponds to oscillations in the top-2 singular values. In Figure 14, we show the dynamics of when the initialization scale is $\alpha = 0.01$ and $\alpha = 0.5$, where balancing holds theoretically for the former but not for the latter. Clearly, we observe that balancing does not hold for $\alpha = 0.5$. However, examining the middle plots reveals that the oscillations in the singular values still have the same amplitude in both cases and for both singular values. This suggests that balancing is merely a tool for analysis, as the oscillations of interest remain prevalent in both scenarios.

Singular Vector Stationarity. Throughout this paper, we considered two initializations in Equation (3), where balancing holds immediately and one where balancing holds for a sufficiently small initialization scale. In this section, we investigate different initializations with aim to observe (i) if they do not converge to the SVS set and (ii) how they affect the oscillations if they do not belong to the SVS set. To this end, we consider the following:

$$\begin{aligned}
 \mathbf{W}_L(0) = \mathbf{0}, \quad \mathbf{W}_\ell(0) = \alpha \mathbf{I}_d, \quad \forall \ell \in [L-1], & \quad (\text{Original}) \\
 \mathbf{W}_L(0) = \mathbf{0}, \quad \mathbf{W}_\ell(0) = \alpha \mathbf{P}_\ell, \quad \forall \ell \in [L-1], & \quad (\text{Orthogonal}) \\
 \mathbf{W}_L(0) = \mathbf{0}, \quad \mathbf{W}_\ell(0) = \alpha \mathbf{H}_\ell, \quad \forall \ell \in [L-1], & \quad (\text{Random})
 \end{aligned}$$

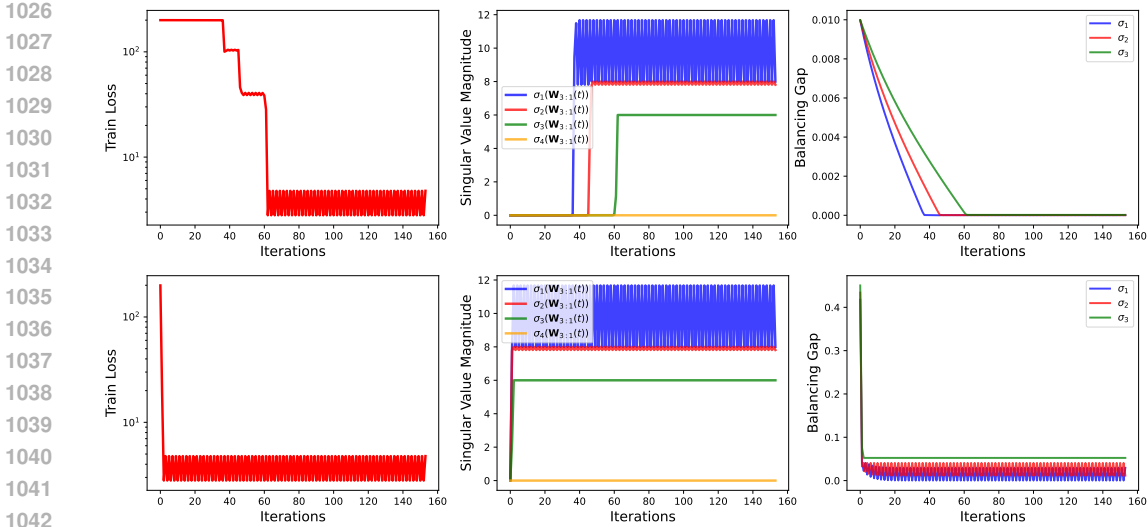


Figure 14: Top: EOS dynamics of a 3-layer DLN with initialization scale $\alpha = 0.01$, where balancing theoretically holds. Bottom: EOS dynamics of the DLN with initialization scale $\alpha = 0.5$. While the balancing does not hold for $\alpha = 0.5$, the oscillations in the singular values are still prevalent, with the same amplitude.

where \mathbf{P}_ℓ is an orthogonal matrix and \mathbf{H}_ℓ is a random matrix with Gaussian entries. For all of these initialization schemes, we consider the same setup as in the balancing case, with an initialization scale of $\alpha = 0.01$. To observe if singular vector stationarity holds, we consider the subspace distance as follows:

$$\text{Subspace Distance} = \|\mathbf{U}_{\ell-1,r}^\top \mathbf{V}_{\ell,r} - \mathbf{I}_r\|_F, \tag{11}$$

where $\mathbf{U}_{\ell,r}$ and $\mathbf{V}_{\ell,r}$ are the top- r left and right singular vectors of layer \mathbf{W}_ℓ , respectively. Since Proposition 1 implies that the intermediate singular vectors cancel, the initialization converges to the SVS set if the subspace distance goes to zero. In Figure 15, we plot the dynamics for all of the initializations. Generally, we observe that the subspace distance for all cases go to zero, validating the use of the SVS set for analysis purposes.

Additional Results. In this section, we provide more experimental results to corroborate our theory. Recall that in Lemma 1, we proved that the learning rate needed to enter the EOS is a function of the depth, and that deeper networks can enter EOS using a smaller learning rate. To verify this claim, we provide an additional experiment where the target matrix is $\mathbf{M}_* \in \mathbb{R}^{5 \times 5}$ with the top singular value set to $\sigma_{*,1} = 0.5$. We use an initialization scale of $\alpha = 0.01$. In Figure 16, we can clearly see that shallower networks need a larger learning rate, and vice versa to enter EOS. Here, black refers to stable learning and white refers to regions in which oscillations occur (EOS regime).

B.5 PERIODIC AND FREE OSCILLATIONS

In this section, we present additional experiments on oscillation and catapults in both deep linear and nonlinear networks to supplement the results in the main paper. First, we consider a 3-layer MLP without bias terms for the weights, with each hidden layer consisting of 1000 units. The network is trained using MSE loss with a learning rate of $\eta = 4$, along with random weights scaled by $\alpha = 0.01$ and full-batch gradient descent on a 5K subset of the MNIST dataset, following Cohen et al. (2021). The motivation for omitting bias terms comes from the findings of Zhang et al. (2024b), where they provably show that a ReLU network without bias terms behaves similarly to a linear network. With this in mind, we aimed to investigate how oscillations manifest in comparison to deep linear networks (DLNs). In Figure 17, we plot the training loss, top-5 singular values, and sharpness throughout training. Interestingly, despite the non-convexity of the loss landscape, the oscillations appear to be almost periodic across all three plots. It would be of great interest to theoretically study

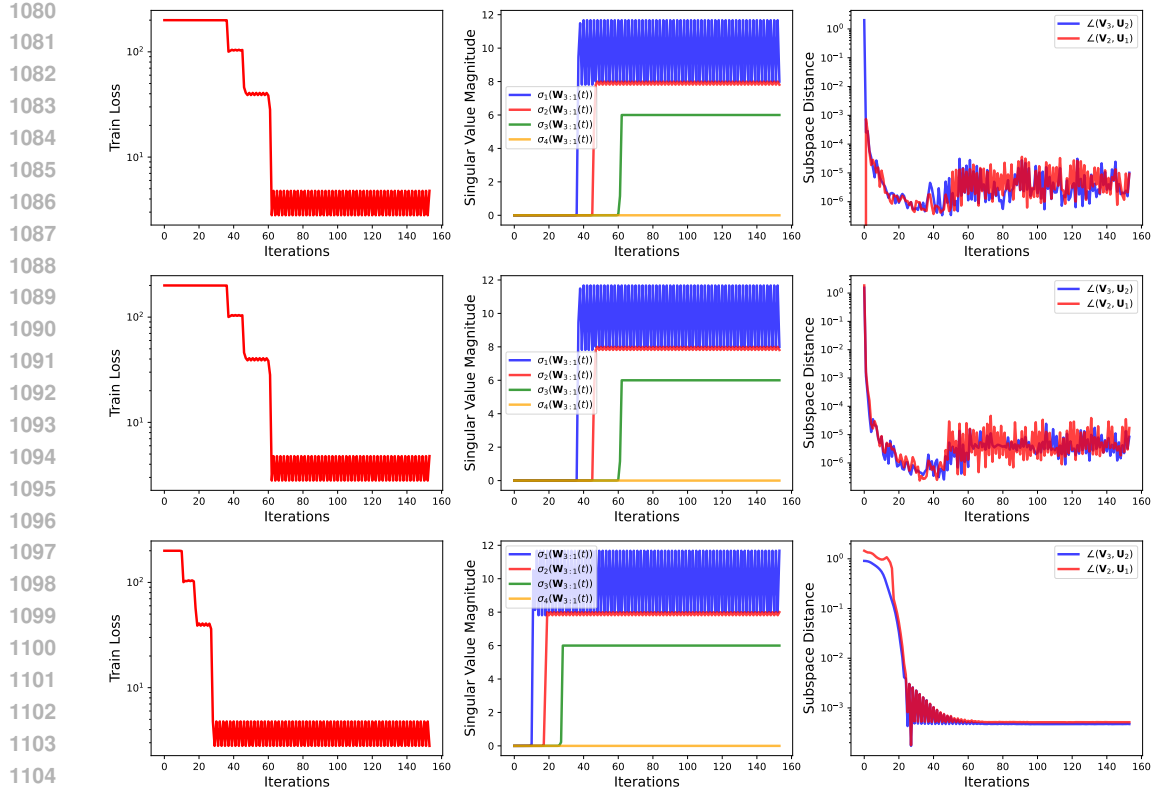


Figure 15: EOS dynamics of a 3-layer DLN for different initializations where it all converges to the SVD set. The subspace distance is defined in Equation (11). Top: Dynamics with the original identity initialization. Middle: Dynamics with orthogonal initialization. Bottom: Dynamics with random initialization.

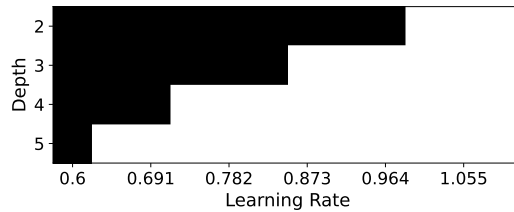


Figure 16: Demonstrating that deeper networks requires a smaller learning rate to enter the EOS regime for DLNs, as implied by Lemma 2, for a target matrix with top singular value $\sigma_{*,1} = 0.5$ and initialization $\alpha = 0.01$. Black refers to stable learning and white refers to regions in which oscillations in the loss and singular values occur. The EOS limit exactly matches $\eta = \frac{2}{L\sigma_{*,1}^{2-\frac{2}{L}}}$.

the behavior of EOS for this network architecture and determine whether our analyses extend to this case as well.

Next, we consider the DLN setting to corroborate our result from Theorem 1. We consider modeling rank-3 target matrix with singular values $\sigma_{*,i} = \{10, 9, 8\}$ with a 3-layer DLN with initialization scale $\alpha = 0.1$. By computing the sharpness under these settings, notice that $2/\lambda_1 = L\sigma_{*,1}^{2-\frac{2}{L}} \approx 0.01547$ and $2/\lambda_2 \approx 0.01657$. In Figure 18, we use learning rates near these values, and plot the oscillations in the singular values. Here, we can see that the oscillations follow exactly our theory.

Lastly, we provide additional experiments demonstrating stronger oscillation in feature directions as measured by the singular values. To this end, we consider a 4-layer MLP with ReLU activations

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

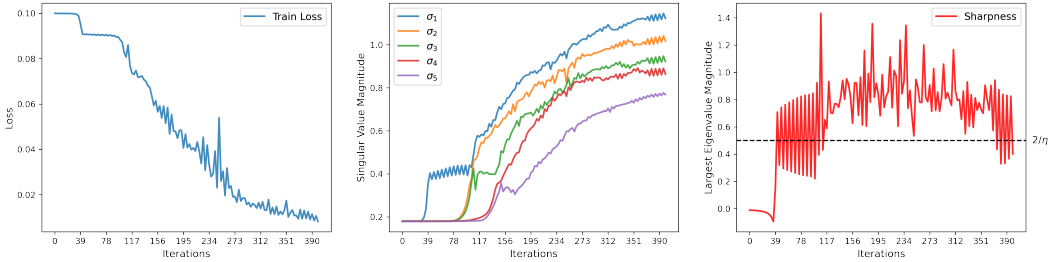


Figure 17: Plot of the training loss, singular values, and sharpness for an MLP network with no bias. Similar to the DLN case, there are oscillations in each of the plots throughout iterations.

with hidden layer size in each unit of 200 for classification on a subsampled 20K set on MNIST and CIFAR-10. In Figure 19, we show that the oscillations in the training loss are artifacts of jumps only in the top singular values, which is also what we observe in the DLN setting.

B.6 INVESTIGATION OF OSCILLATIONS IN LOW-RANK ADAPTORS

Previously, we investigated the differences in oscillations (i.e., oscillations versus catapults) in deep linear and nonlinear networks, and how changes in the landscape present one behavior or the other. Low-rank adaptation (LoRA) (Hu et al., 2022) has arguably become one of the most popular methods for fine-tuning deep neural networks. By viewing the adaptations as individual low-rank matrix factorization problems, then this formulation closely aligns with our theoretical setup with a depth of 2. Here we pose the question (i) does oscillations may appear in such a setup and (ii) what these oscillations may imply in terms of generalization.

Briefly, the main idea behind LoRA is that rather than training from scratch, we can update two low-rank factor matrices to “append” onto an existing weight matrix. That is, give a pre-trained weight matrix $\mathbf{W}_0 \in \mathbb{R}^{d_1 \times d_2}$, LoRA involves updating two low-rank factors commonly referred to as “adaptors”:

$$\underbrace{\mathbf{W}_*}_{\text{new weight}} = \underbrace{\mathbf{W}_0}_{\text{pre-trained weight}} + \underbrace{\mathbf{A}\mathbf{B}^\top}_{\text{adaptors}}.$$

For a sufficiently small rank r , upon training only $\mathbf{A} \in \mathbb{R}^{d_1 \times r}$ and $\mathbf{B} \in \mathbb{R}^{d_2 \times r}$, $\mathbf{W}_* \in \mathbb{R}^{d_1 \times d_2}$ is used for inference.

For the experimental settings, we follow the setup used by Yaras et al. (2024) and consider a pre-trained BERT (Wang et al., 2019) base model and apply adaptation on all attention and feedforward weights in the transformer, resulting in 72 adapted layers in total. For initialization, we use random weights and scale them using an initialization scale of $\alpha = 10^{-3}$ for the adaptors and randomly sample 512 examples from the STS-B (Wang et al., 2019) dataset for fine-tuning. We choose a batch size of 64 with a maximum sequence length of 128 tokens. First, we experiment how large the rank of the adaptors must be to drive the entire network to EOS. Using a learning rate of $\eta = 10^{-4}$, Figure 20 shows oscillatory behavior across all ranks. However, this behavior may also be an artifact of the stochasticity induced by updating with only a batch of samples.

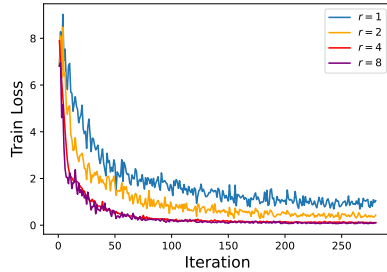


Figure 20: Catapults in the training loss for different ranks for LoRA.

In Figure 21, we present the training loss and Pearson correlation for different learning rates using rank $r = 8$. When $\eta = 10^{-4}$, the training loss catapults from a magnitude of 7 to 4 whereas other learning rates cannot decrease the loss with such a magnitude. Consequently, $\eta = 10^{-4}$ achieves the best Pearson correlation coefficient. This suggests that learning rate plays an important role when the optimization is restricted to a small subspace as used in LoRA. We leave this for future work a careful study of this observation, aiming to accurately select

1188 the learning rate to maximize the efficiency of LoRA.
1189
1190
1191
1192
1193
1194
1195
1196
1197
1198
1199
1200
1201
1202
1203
1204
1205
1206
1207
1208
1209
1210
1211
1212
1213
1214
1215
1216
1217
1218
1219
1220
1221
1222
1223
1224
1225
1226
1227
1228
1229
1230
1231
1232
1233
1234
1235
1236
1237
1238
1239
1240
1241



Figure 18: Depiction of the edge of stability progressively where oscillation is occurring on each singular value depending on the learning rate η . When $\eta = 2/L\sigma_{\star,1}^{2-\frac{2}{L}} \approx 0.0154$, oscillation occur on the first singular value. When $\eta = 2/\sum_{\ell=0}^{L-1} \left(\sigma_{\star,1}^{1-\frac{1}{L}-\frac{1}{L}\ell} \cdot \sigma_{\star,2}^{\frac{1}{L}\ell}\right)^2 \approx 0.0165$, oscillation occur on second singular value and so on.

1296
 1297
 1298
 1299
 1300
 1301
 1302
 1303
 1304
 1305
 1306
 1307
 1308
 1309
 1310
 1311
 1312
 1313
 1314
 1315
 1316
 1317
 1318
 1319
 1320
 1321
 1322
 1323
 1324
 1325
 1326
 1327
 1328
 1329
 1330
 1331
 1332
 1333
 1334
 1335
 1336
 1337
 1338
 1339
 1340
 1341
 1342
 1343
 1344
 1345
 1346
 1347
 1348
 1349

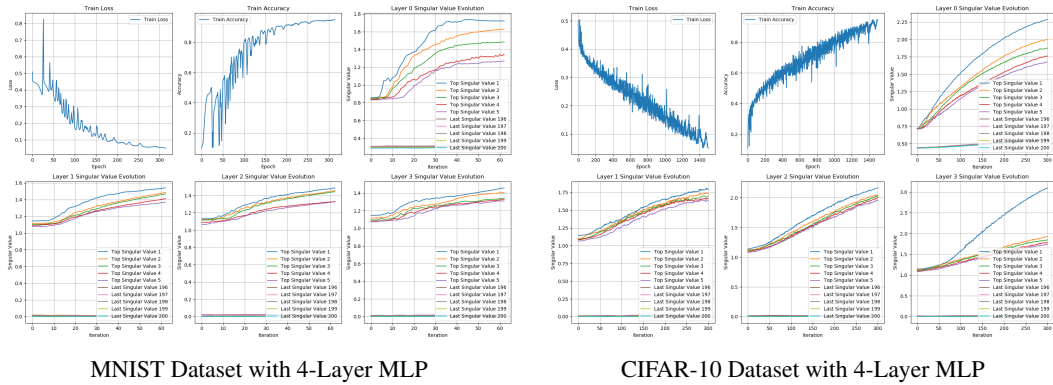


Figure 19: Prevalence of oscillatory behaviors in top subspaces in 4-layer networks with ReLU activations on two different datasets.

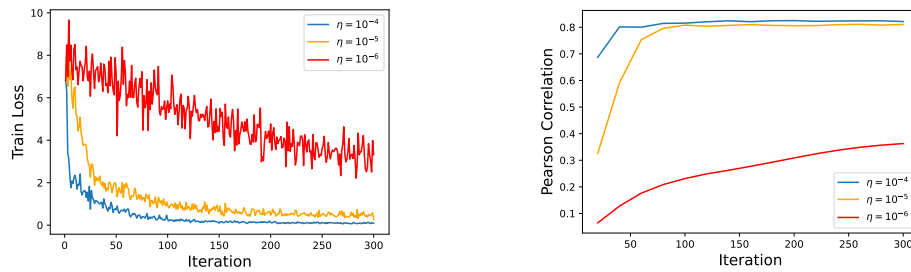


Figure 21: Illustration of different behaviors in the training loss for various learning rates with a fixed rank of $r = 8$ for fine-tuning BERT using LoRA. These plots indicate that larger learning rates lead to higher Pearson correlations. When $\eta = 10^{-4}$, the training loss catapults from a magnitude of 7 to 4 whereas other learning rates do not decrease the loss with such large magnitude.

C DEFERRED PROOFS

C.1 DEFERRED PROOFS FOR OSCILLATIONS

Proposition 1 (Singular Vector Stationary Set). *Consider the deep matrix factorization loss in Equation (1). Let $\mathbf{M}_\star = \mathbf{U}_\star \boldsymbol{\Sigma}_\star \mathbf{V}_\star^\top$ and $\mathbf{W}_\ell(t) = \mathbf{U}_\ell(t) \boldsymbol{\Sigma}_\ell(t) \mathbf{V}_\ell^\top(t)$ denote the compact SVD for the target matrix and the ℓ -th layer weight matrix at time t , respectively. For any time $t \geq 0$, if $\dot{\mathbf{U}}_\ell(t) = \dot{\mathbf{V}}_\ell(t) = 0$ for all $\ell \in [L]$, then the singular vector stationary points for each weight matrix are given by*

$$\text{SVS}(f(\boldsymbol{\Theta})) = \begin{cases} (\mathbf{U}_L, \mathbf{V}_L) &= (\mathbf{U}_\star, \mathbf{Q}_L), \\ (\mathbf{U}_\ell, \mathbf{V}_\ell) &= (\mathbf{Q}_{\ell+1}, \mathbf{Q}_\ell), \quad \forall \ell \in [2, L-1], \\ (\mathbf{U}_1, \mathbf{V}_1) &= (\mathbf{Q}_2, \mathbf{V}_\star), \end{cases}$$

where $\{\mathbf{Q}_\ell\}_{\ell=2}^L$ can be any orthogonal matrices.

Proof. Let us consider the dynamics of $\mathbf{W}_\ell(t)$ in terms of its SVD with respect to time:

$$\dot{\mathbf{W}}_\ell(t) = \dot{\mathbf{U}}_\ell(t) \boldsymbol{\Sigma}_\ell(t) \mathbf{V}_\ell^\top(t) + \mathbf{U}_\ell(t) \dot{\boldsymbol{\Sigma}}_\ell(t) \mathbf{V}_\ell^\top(t) + \mathbf{U}_\ell(t) \boldsymbol{\Sigma}_\ell(t) \dot{\mathbf{V}}_\ell^\top(t). \quad (12)$$

By left multiplying by $\mathbf{U}_\ell^\top(t)$ and right multiplying by $\mathbf{V}_\ell(t)$, we have

$$\mathbf{U}_\ell^\top(t) \dot{\mathbf{W}}_\ell(t) \mathbf{V}_\ell(t) = \mathbf{U}_\ell^\top(t) \dot{\mathbf{U}}_\ell(t) \boldsymbol{\Sigma}_\ell(t) + \dot{\boldsymbol{\Sigma}}_\ell(t) + \boldsymbol{\Sigma}_\ell(t) \dot{\mathbf{V}}_\ell^\top(t) \mathbf{V}_\ell(t), \quad (13)$$

where we used the fact that $\mathbf{U}_\ell(t)$ and $\mathbf{V}_\ell(t)$ have orthonormal columns. Now, note that we also have

$$\mathbf{U}_\ell^\top(t) \mathbf{U}_\ell(t) = \mathbf{I}_r \implies \dot{\mathbf{U}}_\ell^\top(t) \mathbf{U}_\ell(t) + \mathbf{U}_\ell^\top(t) \dot{\mathbf{U}}_\ell(t) = \mathbf{0},$$

which also holds for $\mathbf{V}_\ell(t)$. This implies that $\dot{\mathbf{U}}_\ell^\top(t) \mathbf{U}_\ell(t)$ is a skew-symmetric matrix, and hence have zero diagonals. Since $\boldsymbol{\Sigma}_\ell(t)$ is diagonal, $\mathbf{U}_\ell^\top(t) \dot{\mathbf{U}}_\ell(t) \boldsymbol{\Sigma}_\ell(t)$ and $\boldsymbol{\Sigma}_\ell(t) \dot{\mathbf{V}}_\ell^\top(t) \mathbf{V}_\ell(t)$ have zero diagonals as well. On the other hand, since $\dot{\boldsymbol{\Sigma}}_\ell(t)$ is a diagonal matrix, we can write

$$\hat{\mathbf{I}}_r \odot \left(\mathbf{U}_\ell^\top(t) \dot{\mathbf{W}}_\ell(t) \mathbf{V}_\ell(t) \right) = \mathbf{U}_\ell^\top(t) \dot{\mathbf{U}}_\ell(t) \boldsymbol{\Sigma}_\ell(t) + \boldsymbol{\Sigma}_\ell(t) \dot{\mathbf{V}}_\ell^\top(t) \mathbf{V}_\ell(t), \quad (14)$$

where \odot stands for the Hadamard product and $\hat{\mathbf{I}}_r$ is a square matrix holding zeros on its diagonal and ones elsewhere. Taking transpose of Equation (14), while recalling that $\mathbf{U}_\ell^\top(t) \dot{\mathbf{U}}_\ell(t)$ and $\mathbf{V}_\ell^\top(t) \dot{\mathbf{V}}_\ell(t)$ are skew-symmetric, we have

$$\hat{\mathbf{I}}_r \odot \left(\mathbf{V}_\ell^\top(t) \dot{\mathbf{W}}_\ell^\top(t) \mathbf{U}_\ell(t) \right) = -\boldsymbol{\Sigma}_\ell(t) \mathbf{U}_\ell^\top(t) \dot{\mathbf{U}}_\ell(t) - \dot{\mathbf{V}}_\ell^\top(t) \mathbf{V}_\ell(t) \boldsymbol{\Sigma}_\ell(t). \quad (15)$$

Then, by right multiplying Equation (14) by $\boldsymbol{\Sigma}_\ell(t)$, left-multiply Equation (15) by $\boldsymbol{\Sigma}_\ell(t)$, and by adding the two terms, we get

$$\begin{aligned} \hat{\mathbf{I}}_r \odot \left(\mathbf{U}_\ell^\top(t) \dot{\mathbf{W}}_\ell(t) \mathbf{V}_\ell(t) \boldsymbol{\Sigma}_\ell(t) + \boldsymbol{\Sigma}_\ell(t) \mathbf{V}_\ell^\top(t) \dot{\mathbf{W}}_\ell^\top(t) \mathbf{U}_\ell(t) \right) \\ = \mathbf{U}_\ell^\top(t) \dot{\mathbf{U}}_\ell(t) \boldsymbol{\Sigma}_\ell^2(t) - \boldsymbol{\Sigma}_\ell^2(t) \dot{\mathbf{V}}_\ell^\top(t) \mathbf{V}_\ell(t). \end{aligned}$$

Since we assume that the singular values of \mathbf{M}_\star are distinct, the top- r diagonal elements of $\boldsymbol{\Sigma}_\ell^2(t)$ are also distinct (i.e., $\Sigma_r^2(t) \neq \Sigma_{r'}^2(t)$ for $r \neq r'$). This implies that

$$\mathbf{U}_\ell^\top(t) \dot{\mathbf{U}}_\ell(t) = \mathbf{H}(t) \odot \left[\mathbf{U}_\ell^\top(t) \dot{\mathbf{W}}_\ell(t) \mathbf{V}_\ell(t) \boldsymbol{\Sigma}_\ell(t) + \boldsymbol{\Sigma}_\ell(t) \mathbf{V}_\ell^\top(t) \dot{\mathbf{W}}_\ell^\top(t) \mathbf{U}_\ell(t) \right],$$

where the matrix $\mathbf{H}(t) \in \mathbb{R}^{d \times d}$ is defined by:

$$H_{r,r'}(t) := \begin{cases} (\Sigma_{r'}^2(t) - \Sigma_r^2(t))^{-1}, & r \neq r', \\ 0, & r = r'. \end{cases} \quad (16)$$

1404 Then, multiplying from the left by $\mathbf{U}_\ell(t)$ yields

$$1405 \quad \mathbf{P}_{\mathbf{U}_\ell(t)} \dot{\mathbf{U}}_\ell(t) = \mathbf{U}_\ell(t) \left(\mathbf{H}(t) \odot \left[\mathbf{U}_\ell^\top(t) \dot{\mathbf{W}}_\ell(t) \mathbf{V}_\ell(t) \boldsymbol{\Sigma}_\ell(t) + \boldsymbol{\Sigma}_\ell(t) \mathbf{V}_\ell^\top(t) \dot{\mathbf{W}}_\ell^\top(t) \mathbf{U}_\ell(t) \right] \right), \quad (17)$$

1407 with $\mathbf{P}_{\mathbf{U}_\ell(t)} := \mathbf{U}_\ell(t) \mathbf{U}_\ell^\top(t)$ being the projection onto the subspace spanned by the (orthonormal) columns of $\mathbf{U}_\ell(t)$. Denote by $\mathbf{P}_{\mathbf{U}_{\ell\perp}(t)}$ the projection onto the orthogonal complement (i.e., $\mathbf{P}_{\mathbf{U}_{\ell\perp}(t)} := \mathbf{I}_r - \mathbf{U}_\ell(t) \mathbf{U}_\ell^\top(t)$). Apply $\mathbf{P}_{\mathbf{U}_{\ell\perp}(t)}$ to both sides of Equation (12):

$$1411 \quad \mathbf{P}_{\mathbf{U}_{\ell\perp}(t)} \dot{\mathbf{U}}_\ell(t) = \mathbf{P}_{\mathbf{U}_{\ell\perp}(t)} \dot{\mathbf{U}}_\ell(t) \boldsymbol{\Sigma}_\ell(t) \mathbf{V}_\ell^\top(t) + \mathbf{P}_{\mathbf{U}_{\ell\perp}(t)} \mathbf{U}_\ell(t) \dot{\boldsymbol{\Sigma}}_\ell(t) \mathbf{V}_\ell^\top(t) \quad (18)$$

$$1412 \quad + \mathbf{P}_{\mathbf{U}_{\ell\perp}(t)} \mathbf{U}_\ell(t) \boldsymbol{\Sigma}_\ell(t) \dot{\mathbf{V}}_\ell^\top(t). \quad (19)$$

1414 Note that $\mathbf{P}_{\mathbf{U}_{\ell\perp}(t)} \mathbf{U}_\ell(t) = 0$, and multiply from the right by $\mathbf{V}_\ell(t) \boldsymbol{\Sigma}_\ell^{-1}(t)$ (the latter is well-defined since we have the compact SVD and the top- r elements are non-zero):

$$1416 \quad \mathbf{P}_{\mathbf{U}_{\ell\perp}(t)} \dot{\mathbf{U}}_\ell(t) = \mathbf{P}_{\mathbf{U}_{\ell\perp}(t)} \dot{\mathbf{W}}_\ell(t) \mathbf{V}_\ell(t) \boldsymbol{\Sigma}_\ell^{-1}(t) = (\mathbf{I}_r - \mathbf{U}_\ell(t) \mathbf{U}_\ell^\top(t)) \dot{\mathbf{W}}_\ell(t) \mathbf{V}_\ell(t) \boldsymbol{\Sigma}_\ell^{-1}(t). \quad (20)$$

1418 Then by adding the two equations above, we obtain an expression for $\dot{\mathbf{U}}_\ell(t)$:

$$1419 \quad \begin{aligned} \dot{\mathbf{U}}_\ell(t) &= \mathbf{P}_{\mathbf{U}_\ell(t)} \dot{\mathbf{U}}_\ell(t) + \mathbf{P}_{\mathbf{U}_{\ell\perp}(t)} \dot{\mathbf{U}}_\ell(t) \\ &= \mathbf{U}_\ell(t) \left(\mathbf{H}(t) \odot \left[\mathbf{U}_\ell^\top(t) \dot{\mathbf{W}}_\ell(t) \mathbf{V}_\ell(t) \boldsymbol{\Sigma}_\ell(t) + \boldsymbol{\Sigma}_\ell(t) \mathbf{V}_\ell^\top(t) \dot{\mathbf{W}}_\ell^\top(t) \mathbf{U}_\ell(t) \right] \right) \\ &\quad + (\mathbf{I}_r - \mathbf{U}_\ell(t) \mathbf{U}_\ell^\top(t)) \dot{\mathbf{W}}_\ell(t) \mathbf{V}_\ell(t) \boldsymbol{\Sigma}_\ell^{-1}(t). \end{aligned} \quad (21)$$

1424 We can similarly derive the dynamics for $\dot{\mathbf{V}}_\ell(t)$ and $\dot{\boldsymbol{\Sigma}}_\ell(t)$:

$$1425 \quad \dot{\mathbf{V}}_\ell(t) = \mathbf{V}_\ell(t) \left(\mathbf{H}(t) \odot \left[\boldsymbol{\Sigma}_\ell(t) \mathbf{U}_\ell^\top(t) \dot{\mathbf{W}}_\ell(t) \mathbf{V}_\ell(t) + \mathbf{V}_\ell^\top(t) \dot{\mathbf{W}}_\ell^\top(t) \mathbf{U}_\ell(t) \boldsymbol{\Sigma}_\ell(t) \right] \right) \quad (22)$$

$$1426 \quad + (\mathbf{I}_r - \mathbf{V}_\ell(t) \mathbf{V}_\ell^\top(t)) \dot{\mathbf{W}}_\ell^\top(t) \mathbf{U}_\ell(t) \boldsymbol{\Sigma}_\ell^{-1}(t), \quad (23)$$

$$1428 \quad \dot{\boldsymbol{\Sigma}}_\ell(t) = \mathbf{I}_r \odot \left[\mathbf{U}_\ell^\top(t) \dot{\mathbf{W}}_\ell(t) \mathbf{V}_\ell(t) \right].$$

1431 Now, we will left multiply $\dot{\mathbf{U}}_\ell(t)$ and $\dot{\mathbf{V}}_\ell(t)$ with $\mathbf{U}_\ell^\top(t)$ and $\mathbf{V}_\ell^\top(t)$, respectively, to obtain

$$1432 \quad \mathbf{U}_\ell^\top(t) \dot{\mathbf{U}}_\ell(t) = -\mathbf{H}(t) \odot \left[\mathbf{U}_\ell^\top(t) \nabla_{\mathbf{W}_\ell} f(\boldsymbol{\Theta}) \mathbf{V}_\ell(t) \boldsymbol{\Sigma}_\ell(t) + \boldsymbol{\Sigma}_\ell(t) \mathbf{V}_\ell^\top(t) \nabla_{\mathbf{W}_\ell} f(\boldsymbol{\Theta}) \mathbf{U}_\ell(t) \right],$$

$$1433 \quad \mathbf{V}_\ell^\top(t) \dot{\mathbf{V}}_\ell(t) = -\mathbf{H}(t) \odot \left[\boldsymbol{\Sigma}_\ell(t) \mathbf{U}_\ell^\top(t) \nabla_{\mathbf{W}_\ell} f(\boldsymbol{\Theta}) \mathbf{V}_\ell(t) + \mathbf{V}_\ell^\top(t) \nabla_{\mathbf{W}_\ell} f(\boldsymbol{\Theta}) \mathbf{U}_\ell(t) \boldsymbol{\Sigma}_\ell(t) \right],$$

1435 where we replaced $\dot{\mathbf{W}}_\ell(t) := -\nabla_{\mathbf{W}_\ell} f(\boldsymbol{\Theta})$, as $\dot{\mathbf{W}}_\ell(t)$ is the gradient of $f(\boldsymbol{\Theta})$ with respect to \mathbf{W}_ℓ by definition. By rearranging and multiplying by $\boldsymbol{\Sigma}_\ell(t)$, we have

$$1437 \quad \mathbf{U}_\ell^\top(t) \dot{\mathbf{U}}_\ell(t) \boldsymbol{\Sigma}_\ell(t) - \boldsymbol{\Sigma}_\ell(t) \mathbf{V}^\top(t) \dot{\mathbf{V}}_\ell(t) = -\hat{\mathbf{I}}_r \odot \left[\mathbf{U}_\ell^\top(t) \nabla_{\mathbf{W}_\ell} f(\boldsymbol{\Theta}) \mathbf{V}_\ell(t) \right]. \quad (24)$$

1439 Hence, when $\dot{\mathbf{U}}_\ell(t) = 0$ and $\dot{\mathbf{V}}_\ell(t) = 0$, it must be that the left-hand side is zero and so $\mathbf{U}_\ell^\top(t) \nabla_{\mathbf{W}_\ell} f(\boldsymbol{\Theta}) \mathbf{V}_\ell(t)$ is a diagonal matrix.

1441 Now, notice that for the given loss function $f(\boldsymbol{\Theta})$, we have

$$1442 \quad -\dot{\mathbf{W}}_\ell(t) = \nabla_{\mathbf{W}_\ell} f(\boldsymbol{\Theta}(t)) = \mathbf{W}_{L:\ell+1}^\top(t) \cdot (\mathbf{W}_{L:1}(t) - \mathbf{M}_\star) \cdot \mathbf{W}_{\ell-1:1}^\top(t).$$

1444 Then, from Equation (24), when the singular vectors are stationary, we have

$$1445 \quad \mathbf{U}_\ell^\top(t) \mathbf{W}_{L:\ell+1}^\top(t) \cdot (\mathbf{W}_{L:1}(t) - \mathbf{M}_\star) \cdot \mathbf{W}_{\ell-1:1}^\top(t) \mathbf{V}_\ell(t)$$

1446 must be a diagonal matrix for all $\ell \in [L]$. The only solution to the above should be (since the intermediate singular vectors need to cancel to satisfy the diagonal condition), is the set

$$1448 \quad \text{SVS}(f(\boldsymbol{\Theta})) = \begin{cases} (\mathbf{U}_L, \mathbf{V}_L) &= (\mathbf{U}_\star, \mathbf{Q}_L), \\ (\mathbf{U}_\ell, \mathbf{V}_\ell) &= (\mathbf{Q}_{\ell+1}, \mathbf{Q}_\ell), \quad \forall \ell \in [2, L-1], \\ (\mathbf{U}_1, \mathbf{V}_1) &= (\mathbf{Q}_2, \mathbf{V}_\star), \end{cases}$$

1452 where $\{\mathbf{Q}_\ell\}_{\ell=2}^L$ are any set of orthogonal matrices. Then, notice that when the singular vectors are stationary, the dynamics become isolated on the singular values:

$$1454 \quad \dot{\boldsymbol{\Sigma}}_\ell(t) = \mathbf{I}_r \odot \left[\mathbf{U}_\ell^\top(t) \dot{\mathbf{W}}_\ell(t) \mathbf{V}_\ell(t) \right],$$

1456 since $\left[\mathbf{U}_\ell^\top(t) \dot{\mathbf{W}}_\ell(t) \mathbf{V}_\ell(t) \right]$ is diagonal. This completes the proof.

1457

□

Theorem 1 (Rank-1 Oscillation). Let $\mathbf{M}_* = \mathbf{U}_* \Sigma_* \mathbf{V}_*^\top$ denote the SVD of the target matrix and let $S := L\sigma_{*,1}^{2-\frac{2}{L}}$, $\alpha' := \left(\ln \left(\frac{2\sqrt{2}}{\eta L \sigma_{*,1}^{\frac{2}{L}}} \right) \cdot \frac{\sigma_{*,1}^{\frac{4}{L}}}{L^2 \cdot 2^{\frac{2L-3}{L}}} \right)^{\frac{1}{4}}$, and $K' := \max \left\{ \sum_{\ell=0}^{L-1} \left(\sigma_{*,1}^{1-\frac{1}{L}-\frac{1}{L}\ell} \cdot \sigma_{*,2}^{\frac{1}{L}\ell} \right)^2, \frac{S}{2\sqrt{2}} \right\}$. If we run GD on the deep matrix factorization loss with initialization scale $\alpha < \alpha'$ and learning rate $\eta = \frac{2}{K}$, where $K' < K < S$, then under strict balancing, each weight matrix $\mathbf{W}_\ell \in \mathbb{R}^{d \times d}$ oscillates around the minima in a 2-period fixed orbit ($i \in \{1, 2\}$) as follows:

$$\mathbf{W}_L(t) = \underbrace{\rho_i(t) \cdot \mathbf{u}_{*,1} \mathbf{v}_{*,1}^\top}_{\text{oscillation subspace}} + \underbrace{\sum_{j=2}^r \sigma_{*,i} \mathbf{u}_{*,j} \mathbf{v}_{*,j}^\top}_{\text{stationary subspace}}, \quad i = 1, 2,$$

$$\mathbf{W}_\ell(t) = \underbrace{\rho_i(t) \cdot \mathbf{v}_{*,1} \mathbf{v}_{*,1}^\top}_{\text{oscillation subspace}} + \underbrace{\sum_{j=2}^r \sigma_{*,i} \mathbf{v}_{*,j} \mathbf{v}_{*,j}^\top}_{\text{stationary subspace}}, \quad i = 1, 2, \quad \forall \ell \in [L-1],$$

where $\rho_i(t) \in \{\rho_1, \rho_2\}$ and $\rho_1 \in (0, \sigma_{*,1}^{1/L})$ and $\rho_2 \in (\sigma_{*,1}^{1/L}, (2\sigma_{*,1})^{1/L})$ are the two real roots of the polynomial $g(\rho) = 0$, where

$$g(\rho) = \rho^L \cdot \frac{1 + (1 + \eta L (\sigma_{*,1} - \rho^L) \cdot \rho^{L-2})^{2L-1}}{1 + (1 + \eta L (\sigma_{*,1} - \rho^L) \cdot \rho^{L-2})^{L-1}} - \sigma_{*,1}.$$

Proof. For ease of exposition, let us denote the first singular value as $\sigma_1 := \sigma_{\ell,1}$. Under balancing, consider the two-step GD update on the first singular value:

$$\sigma_1(t+1) = \sigma_1(t) + \eta L \cdot (\sigma_{*,1} - \sigma_1^L(t)) \cdot \sigma_1^{L-1}(t)$$

$$\sigma_1(t) = \sigma_1(t+2) = \sigma_1(t+1) + \eta L \cdot (\sigma_{*,1} - \sigma_1^L(t+1)) \cdot \sigma_1^{L-1}(t+1). \quad (\text{By 2-period orbit})$$

Define $z := (1 + \eta L \cdot (\sigma_{*,1} - \sigma_1^L(t)) \cdot \sigma_1^{L-2}(t))$ and by plugging in $\sigma_1(t+1)$ for $\sigma_1(t)$, we have

$$\sigma_1(t) = \sigma_1(t)z + \eta L \cdot (\sigma_{*,1} - \sigma_1^L(t)z^L) \cdot \sigma_1^{L-1}(t)z^{L-1}$$

$$\implies 1 = z + \eta L \cdot (\sigma_{*,1} - \sigma_1^L(t)z^L) \cdot \sigma_1^{L-2}(t)z^{L-1}$$

$$\implies 1 = (1 + \eta L \cdot (\sigma_{*,1} - \sigma_1^L(t)) \cdot \sigma_1^{L-2}(t)) + \eta L \cdot (\sigma_{*,1} - \sigma_1^L(t)z^L) \cdot \sigma_1^{L-2}(t)z^{L-1}$$

$$\implies 0 = (\sigma_{*,1} - \sigma_1^L(t)) + (\sigma_{*,1} - \sigma_1^L(t)z^L) \cdot z^{L-1}$$

Simplifying this expression further, we have

$$0 = \sigma_{*,1} - \sigma_1^L(t) + \sigma_{*,1}z^{L-1} - \sigma_1^L(t)z^{2L-1}$$

$$\implies \sigma_1^L(t) + \sigma_1^L(t)z^{2L-1} = \sigma_{*,1} + \sigma_{*,1}z^{L-1}$$

$$\implies \sigma_1^L(t) \cdot (1 + z^{2L-1}) = \sigma_{*,1} \cdot (1 + z^{L-1})$$

$$\implies \sigma_1^L(t) \frac{(1 + z^{2L-1})}{(1 + z^{L-1})} = \sigma_{*,1},$$

and by letting $\rho := \sigma_1(t)$, we obtain the polynomial

$$\sigma_{*,1} = \rho^L \frac{1 + z^{2L-1}}{1 + z^{L-1}}, \quad \text{where } z := (1 + \eta L (\sigma_{*,1} - \rho^L) \cdot \rho^{L-2}).$$

Next, we show the existence of the roots within the ranges for ρ_1 and ρ_2 . First, consider $\rho_1 \in (0, \sigma_{*,1}^{1/L})$. We will show that for two values within this range, there is a sign change for all $L \geq 2$.

More specifically, we show that there exists $\rho \in (0, \sigma_{*,1}^{1/L})$ such that

$$\rho^L \frac{1 + z^{2L-1}}{1 + z^{L-1}} - \sigma_{*,1} > 0 \quad \text{and} \quad \rho^L \frac{1 + z^{2L-1}}{1 + z^{L-1}} - \sigma_{*,1} < 0.$$

For the positive case, consider $\rho = (\frac{1}{2}\sigma_{*,1})^{1/L}$. We need to show that

$$\frac{1 + z^{2L-1}}{1 + z^{L-1}} = \frac{1 + \left(1 + \eta L \cdot \left(\frac{\sigma_{*,1}}{2}\right)^{\frac{1-\frac{2}{L}}{2^{1-\frac{2}{L}}}}\right)^{2L-1}}{1 + \left(1 + \eta L \cdot \left(\frac{\sigma_{*,1}}{2}\right)^{\frac{1-\frac{2}{L}}{2^{1-\frac{2}{L}}}}\right)^{L-1}} > 2.$$

To do this, we will plug in the smallest possible value of $\eta = \frac{2}{L\sigma_{*,1}^{\frac{2-\frac{2}{L}}{L}}}$ to show that the fraction is still greater than 2, which gives us

$$\frac{1 + \left(1 + \frac{1}{2^{1-\frac{2}{L}}}\right)^{2L-1}}{1 + \left(1 + \frac{1}{2^{1-\frac{2}{L}}}\right)^{L-1}}, \quad (25)$$

which is an increasing function of L for all $L \geq 2$ and so Equation (25) must be greater than 2. For the negative case, we can simply consider $\rho = 0$. Hence, since the polynomial is continuous, by the Intermediate Value Theorem (IVT), there must exist a root within the range $\rho \in (0, \sigma_{*,1}^{1/L})$.

Next, consider the range $\rho_2 \in (\sigma_{*,1}^{1/L}, (2\sigma_{*,1})^{1/L})$. Similarly, we will show sign changes for two values in ρ_2 . For the positive case, consider $\rho = (\frac{3}{2}\sigma_{*,1})^{1/L}$. For η , we can plug in the smallest possible value within the range to show that this value of ρ still provides a positive quantity. Specifically, we need to show that

$$\frac{1 + z^{2L-1}}{1 + z^{L-1}} - \sigma_{*,1} > \frac{2}{3} \implies \frac{1 + \left(1 + \frac{2}{\sigma_{*,1}^{\frac{2-\frac{2}{L}}{L}}} \cdot (\sigma_{*,1} - \frac{3}{2}\sigma_{*,1}) \cdot \left(\frac{3}{2}\sigma_{*,1}\right)^{1-\frac{2}{L}}\right)^{2L-1}}{1 + \left(1 + \frac{2}{\sigma_{*,1}^{\frac{2-\frac{2}{L}}{L}}} \cdot (\sigma_{*,1} - \frac{3}{2}\sigma_{*,1}) \cdot \left(\frac{3}{2}\sigma_{*,1}\right)^{1-\frac{2}{L}}\right)^{L-1}} > \frac{2}{3}.$$

We can simplify the fraction as follows:

$$\frac{1 + \left(1 + \frac{2}{\sigma_{*,1}^{\frac{2-\frac{2}{L}}{L}}} \cdot (\sigma_{*,1} - \frac{3}{2}\sigma_{*,1}) \cdot \left(\frac{3}{2}\sigma_{*,1}\right)^{1-\frac{2}{L}}\right)^{2L-1}}{1 + \left(1 + \frac{2}{\sigma_{*,1}^{\frac{2-\frac{2}{L}}{L}}} \cdot (\sigma_{*,1} - \frac{3}{2}\sigma_{*,1}) \cdot \left(\frac{3}{2}\sigma_{*,1}\right)^{1-\frac{2}{L}}\right)^{L-1}} = \frac{1 + \left(1 - \left(\frac{3}{2}\right)^{1-\frac{2}{L}}\right)^{2L-1}}{1 + \left(1 - \left(\frac{3}{2}\right)^{1-\frac{2}{L}}\right)^{L-1}}.$$

Then, since we are subtracting by $(\frac{3}{2})^{1-\frac{2}{L}}$, we can plug in its largest value for $L \geq 2$, which is $3/2$. This gives us

$$\frac{1 + (-0.5)^{2L-1}}{1 + (-0.5)^{L-1}} > \frac{2}{3},$$

as for odd values of L , the function increases to 1 starting from $L = 2$, and decreases to 1 for even L . To check negativity, let us define

$$h(\rho) := \frac{f(\rho)}{g(\rho)} := \frac{\rho^L (1 + z^{2L-1})}{1 + z^{L-1}}.$$

We will show that $h'(\sigma_{\star,1}^{1/L}) < 0$:

$$\begin{aligned} h'(\sigma_{\star,1}^{1/L}) &= \frac{f'(\sigma_{\star,1}^{1/L})g(\sigma_{\star,1}^{1/L}) - f(\sigma_{\star,1}^{1/L})g'(\sigma_{\star,1}^{1/L})}{g^2(\sigma_{\star,1}^{1/L})} \\ &= \frac{f'(\sigma_{\star,1}^{1/L}) - \sigma_{\star,1}g'(\sigma_{\star,1}^{1/L})}{2} \\ &= \frac{L\sigma_{\star,1}^{1-\frac{1}{L}} - \sigma_{\star,1}(2L-1)(\eta L^2\sigma_{\star,1}^{2-\frac{3}{L}}) - \sigma_{\star,1}(L-1)(\eta L^2\sigma_{\star,1}^{2-\frac{3}{L}})}{2} \\ &= \frac{L\sigma_{\star,1}^{1-\frac{1}{L}} - (3L-2)(\eta L^2\sigma_{\star,1}^{3-\frac{3}{L}})}{2} < 0, \end{aligned}$$

as otherwise we need $\eta \leq \frac{\sigma_{\star,1}^{2/L-2}}{3L^2-2L}$, which is out of the range of interest. Since $h'(\rho) < 0$, it follows that there exists a $\delta > 0$ such that $h(\rho) > h(x)$ for all x such that $\rho < x < \rho + \delta$. Lastly, since $h(\rho) - \sigma_{\star,1} = 0$ for $\rho = \sigma_{\star,1}^{1/L}$, it follows that $h(\rho) - \sigma_{\star,1}$ must be negative at $\rho + \delta$. Similarly, by IVT, there must exist a root within the range $\rho_2 \in (\sigma_{\star,1}^{1/L}, (2\sigma_{\star,1})^{1/L})$.

Then, by Proposition 2, notice that we can write the dynamics of the weight matrices as

$$\begin{aligned} \mathbf{W}_L(t) &= \sigma_{\star,1}\mathbf{u}_{\star,1}\mathbf{v}_{\star,1}^\top + \sum_{j=1}^r \sigma_{i,\star}\mathbf{u}_{\star,j}\mathbf{v}_{\star,j}^\top, \\ \mathbf{W}_\ell(t) &= \sigma_{\star,1}\mathbf{v}_{\star,1}\mathbf{v}_{\star,1}^\top + \sum_{j=1}^r \sigma_{i,\star}\mathbf{v}_{\star,j}\mathbf{v}_{\star,j}^\top, \quad \forall \ell \in [L-1]. \end{aligned}$$

By the oscillations, we can replace $\sigma_{\star,1}$ with ρ_i for $i = 1, 2$. This completes the proof. \square

Lemma 1 (Hessian Eigenvalues at Convergence). *Consider running GD on the deep matrix factorization loss $f(\Theta)$ defined in Equation (1). Under strict balancing, for any stationary point Θ such that $\nabla_{\Theta}f(\Theta) = 0$, the set of all non-zero eigenvalues of the Hessian of the training loss are given by*

$$\lambda_{\Theta} = \underbrace{\left\{ L\sigma_{\star,i}^{2-\frac{2}{L}}, \sigma_{\star,i}^{2-\frac{2}{L}} \right\}_{i=1}^r}_{\text{self-interaction}}, \cup \underbrace{\left\{ \sum_{\ell=0}^{L-1} \left(\sigma_{\star,i}^{1-\frac{1}{L}-\frac{1}{L}\ell} \cdot \sigma_{\star,j}^{\frac{1}{L}\ell} \right)^2 \right\}_{i \neq j}^r}_{\text{interaction with other singular values}}, \cup \underbrace{\left\{ \sum_{\ell=0}^{L-1} \left(\sigma_{\star,k}^{1-\frac{1}{L}-\frac{1}{L}\ell} \cdot \alpha^\ell \right)^2 \right\}_{k=1}^r}_{\text{interaction with initialization}}$$

where $\sigma_{\star,i}$ is the i -th singular value of the target matrix $\mathbf{M}_{\star} \in \mathbb{R}^{d \times d}$, $\alpha \in \mathbb{R}$ is the initialization scale, L is the depth of the network, and the second element of the set under ‘‘self-interaction’’ has a multiplicity of $d - r$.

Proof. By Proposition 2, notice that we can re-write the loss in Equation (1) as

$$\frac{1}{2} \|\mathbf{W}_{L:1} - \mathbf{M}_{\star}\|_{\text{F}}^2 = \frac{1}{2} \|\Sigma_{L:1} - \Sigma_{\star}\|_{\text{F}}^2,$$

where $\Sigma_{L:1}$ are the singular values of $\mathbf{W}_{L:1}$. We will first show that the eigenvalues of the Hessian with respect to the weight matrices \mathbf{W}_ℓ are equivalent to those of the Hessian taken with respect to its singular values Σ_ℓ . To this end, consider the vectorized form of the loss:

$$f(\Theta) := \frac{1}{2} \|\mathbf{W}_{L:1} - \mathbf{M}_{\star}\|_{\text{F}}^2 = \frac{1}{2} \|\text{vec}(\mathbf{W}_{L:1}) - \text{vec}(\mathbf{M}_{\star})\|_2^2.$$

Then, each block of the Hessian $\nabla_{\Theta}^2 f(\Theta) \in \mathbb{R}^{d^2 L \times d^2 L}$ with respect to the vectorized parameters is given as

$$[\nabla_{\Theta}^2 f(\Theta)]_{m,\ell} = \nabla_{\text{vec}(\mathbf{W}_m)} \nabla_{\text{vec}(\mathbf{W}_\ell)}^\top f(\Theta) \in \mathbb{R}^{d^2 \times d^2}.$$

By the vectorization trick, each vectorized layer matrix has an SVD of the form $\text{vec}(\mathbf{W}_\ell) = \text{vec}(\mathbf{U}_\ell \boldsymbol{\Sigma}_\ell \mathbf{V}_\ell^\top) = (\mathbf{V}_\ell \otimes \mathbf{U}_\ell) \cdot \text{vec}(\boldsymbol{\Sigma}_\ell)$. Then, notice that we have

$$\nabla_{\text{vec}(\mathbf{W}_\ell)} f(\boldsymbol{\Theta}(t)) = (\mathbf{V}_\ell \otimes \mathbf{U}_\ell) \cdot \nabla_{\text{vec}(\boldsymbol{\Sigma}_\ell)} f(\boldsymbol{\Theta}(t)),$$

which gives us that each block of the Hessian is given by

$$\begin{aligned} [\nabla_{\boldsymbol{\Theta}}^2 f(\boldsymbol{\Theta})]_{m,\ell} &= \nabla_{\text{vec}(\mathbf{W}_m)} \nabla_{\text{vec}(\mathbf{W}_\ell)}^\top f(\boldsymbol{\Theta}) \\ &= (\mathbf{V}_m \otimes \mathbf{U}_m) \cdot \underbrace{\nabla_{\text{vec}(\boldsymbol{\Sigma}_m)} \nabla_{\text{vec}(\boldsymbol{\Sigma}_\ell)}^\top f(\boldsymbol{\Theta})}_{=:\mathbf{H}_{m,\ell}} \cdot (\mathbf{V}_\ell \otimes \mathbf{U}_\ell)^\top. \end{aligned}$$

Then, since the Kronecker product of two orthogonal matrices is also an orthogonal matrix by Lemma 9, we can write the overall Hessian matrix as

$$\tilde{\mathbf{H}} = \begin{bmatrix} \mathbf{R}_1 \mathbf{H}_{1,1} \mathbf{R}_1 & \mathbf{R}_1 \mathbf{H}_{1,2} \mathbf{R}_2 & \dots & \mathbf{R}_1 \mathbf{H}_{1,L} \mathbf{R}_L \\ \mathbf{R}_2 \mathbf{H}_{2,1} \mathbf{R}_1 & \mathbf{R}_2 \mathbf{H}_{2,2} \mathbf{R}_2 & \dots & \mathbf{R}_2 \mathbf{H}_{2,L} \mathbf{R}_L \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{R}_L \mathbf{H}_{L,1} \mathbf{R}_1 & \mathbf{R}_L \mathbf{H}_{L,2} \mathbf{R}_2 & \dots & \mathbf{R}_L \mathbf{H}_{L,L} \mathbf{R}_L \end{bmatrix},$$

for orthogonal matrices $\{\mathbf{R}_\ell\}_{\ell=1}^L$. Then, by Lemma 8, the eigenvalues of $\tilde{\mathbf{H}}$ are the same as those of \mathbf{H} , where $\mathbf{H} \in \mathbb{R}^{d^2 L \times d^2 L}$ is the Hessian matrix with respect to the vectorized $\boldsymbol{\Sigma}_\ell$:

$$\mathbf{H} = \begin{bmatrix} \mathbf{H}_{1,1} & \mathbf{H}_{1,2} & \dots & \mathbf{H}_{L,1} \\ \mathbf{H}_{2,1} & \mathbf{H}_{2,2} & \dots & \mathbf{H}_{L,2} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{H}_{1,L} & \mathbf{H}_{2,L} & \dots & \mathbf{H}_{L,L} \end{bmatrix}.$$

Now, we can consider the following vectorized loss:

$$\begin{aligned} f(\boldsymbol{\Theta}) &= \frac{1}{2} \|\boldsymbol{\Sigma}_{L:1} - \boldsymbol{\Sigma}_\star\|_{\text{F}}^2 = \frac{1}{2} \|\text{vec}(\boldsymbol{\Sigma}_{L:1} - \boldsymbol{\Sigma}_\star)\|_2^2 \\ &= \frac{1}{2} \|\underbrace{(\boldsymbol{\Sigma}_{\ell-1:1}^\top \otimes \boldsymbol{\Sigma}_{L:\ell+1})}_{=:\mathbf{A}_\ell} \cdot \text{vec}(\boldsymbol{\Sigma}_\ell) - \text{vec}(\boldsymbol{\Sigma}_\star)\|_2^2. \end{aligned}$$

Then, the gradient with respect to $\text{vec}(\boldsymbol{\Sigma}_\ell)$ is given by

$$\nabla_{\text{vec}(\boldsymbol{\Sigma}_\ell)} f(\boldsymbol{\Theta}) = \mathbf{A}_\ell^\top (\mathbf{A}_\ell \cdot \text{vec}(\boldsymbol{\Sigma}_\ell) - \text{vec}(\boldsymbol{\Sigma}_\star)).$$

Then, for $m = \ell$, we have

$$\mathbf{H}_{\ell,\ell} = \nabla_{\text{vec}(\boldsymbol{\Sigma}_\ell)}^2 f(\boldsymbol{\Theta}) = \mathbf{A}_\ell^\top \mathbf{A}_\ell.$$

For $m \neq \ell$, we have

$$\begin{aligned} \mathbf{H}_{m,\ell} &= \nabla_{\text{vec}(\boldsymbol{\Sigma}_m)} \nabla_{\text{vec}(\boldsymbol{\Sigma}_\ell)} f(\boldsymbol{\Theta}) = \nabla_{\text{vec}(\boldsymbol{\Sigma}_m)} [\mathbf{A}_\ell^\top (\mathbf{A}_\ell \text{vec}(\boldsymbol{\Sigma}_\ell) - \text{vec}(\mathbf{M}^*))] \\ &= \nabla_{\text{vec}(\boldsymbol{\Sigma}_m)} \mathbf{A}_\ell^\top \cdot \underbrace{(\mathbf{A}_\ell \text{vec}(\boldsymbol{\Sigma}_\ell) - \text{vec}(\mathbf{M}^*))}_{=0 \text{ at convergence}} + \mathbf{A}_\ell^\top \cdot \nabla_{\text{vec}(\boldsymbol{\Sigma}_m)} (\mathbf{A}_\ell \text{vec}(\boldsymbol{\Sigma}_\ell) - \text{vec}(\mathbf{M}^*)) \\ &= \mathbf{A}_\ell^\top \mathbf{A}_m, \end{aligned}$$

where we have used the product rule along with the fact that $\mathbf{A}_\ell \text{vec}(\boldsymbol{\Sigma}_\ell) = \mathbf{A}_m \text{vec}(\boldsymbol{\Sigma}_m)$.

Overall, the Hessian at convergence for GD is given by

$$\mathbf{H} = \begin{bmatrix} \mathbf{A}_1^\top \mathbf{A}_1 & \mathbf{A}_1^\top \mathbf{A}_2 & \dots & \mathbf{A}_1^\top \mathbf{A}_L \\ \mathbf{A}_2^\top \mathbf{A}_1 & \mathbf{A}_2^\top \mathbf{A}_2 & \dots & \mathbf{A}_2^\top \mathbf{A}_L \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{A}_L^\top \mathbf{A}_1 & \mathbf{A}_L^\top \mathbf{A}_2 & \dots & \mathbf{A}_L^\top \mathbf{A}_L \end{bmatrix}$$

Now, we can derive an explicit expression for each $\mathbf{A}_{m,\ell}$ by considering the implicit balancing effect of GD in Lemma 2. Under balancing and Proposition 2, we have that at convergence,

$$\Sigma_{L:1} = \Sigma_\star \implies \Sigma_\ell = \begin{bmatrix} \Sigma_{\star,r}^{1/L} & \mathbf{0} \\ \mathbf{0} & \alpha \cdot \mathbf{I}_{d-r} \end{bmatrix}, \quad \forall \ell \in [L-1], \quad \text{and } \Sigma_L = \Sigma_\star^{1/L}.$$

Thus, we have

$$\mathbf{H}_{m,\ell} = \begin{cases} \Sigma_\ell^{2(\ell-1)} \otimes \Sigma_\star^{\frac{2(L-\ell)}{L}} & \text{for } m = \ell, \\ \Sigma_\ell^{m+\ell-2} \otimes \Sigma_\star^{2L-m-\ell} & \text{for } m \neq \ell. \end{cases}$$

Now, we are left with computing the eigenvalues of $\mathbf{H} \in \mathbb{R}^{d^2 L \times d^2 L}$. To do this, let us block diagonalize \mathbf{H} into $\mathbf{H} = \mathbf{P}\mathbf{C}\mathbf{P}^\top$, where \mathbf{P} is a permutation matrix and

$$\mathbf{C} = \begin{bmatrix} \mathbf{C}_1 & & \\ & \ddots & \\ & & \mathbf{C}_{d^2} \end{bmatrix} \in \mathbb{R}^{d^2 L \times d^2 L},$$

where each (i, j) -th entry of $\mathbf{C}_k \in \mathbb{R}^{L \times L}$ is the k -th diagonal element of $\mathbf{H}_{i,j}$. Since \mathbf{C} and \mathbf{H} are similar matrices, they have the same eigenvalues. Then, since \mathbf{C} is a block diagonal matrix, its eigenvalues (and hence the eigenvalues of \mathbf{H}) are the union of each of the eigenvalues of its blocks.

By observing the structure of $\mathbf{H}_{m,\ell}$, notice that each \mathbf{C}_k is a rank-1 matrix. Hence, when considering the top- r diagonal elements of $\mathbf{H}_{m,\ell}$ corresponding to each Kronecker product to construct \mathbf{C}_k , each \mathbf{C}_k can be written as an outer product $\mathbf{u}\mathbf{u}^\top$, where $\mathbf{u} \in \mathbb{R}^L$ is

$$\mathbf{u}^\top = \left[\sigma_{\star,i}^{1-\frac{1}{L}} \sigma_{\star,j}^0 \quad \sigma_{\star,i}^{1-\frac{2}{L}} \sigma_{\star,j}^{\frac{1}{L}} \quad \sigma_{\star,i}^{1-\frac{3}{L}} \sigma_{\star,j}^{\frac{2}{L}} \quad \dots \quad \sigma_{\star,i}^0 \sigma_{\star,j}^{1-\frac{1}{L}} \right]^\top. \quad (26)$$

Then, the non-zero eigenvalue of this rank-1 matrix is simply $\|\mathbf{u}\|_2^2$, which simplifies to

$$\|\mathbf{u}\|_2^2 = \sum_{\ell=0}^{L-1} \left(\sigma_{\star,i}^{1-\frac{1}{L}-\frac{1}{L}\ell} \cdot \sigma_{\star,j}^{\frac{1}{L}\ell} \right)^2.$$

Next, we can consider the remaining $d-r$ components of each Kronecker product of $\mathbf{H}_{m,\ell}$. Notice that for $m = \ell = L$, we have

$$\mathbf{H}_{L,L} = \begin{bmatrix} \sigma_{\star,1}^{\frac{2(L-1)}{L}} \cdot \mathbf{I}_d & & & \\ & \ddots & & \\ & & \sigma_{\star,r}^{\frac{2(L-1)}{L}} \cdot \mathbf{I}_d & \\ & & & \alpha^{2(L-1)} \mathbf{I}_{d-r} \otimes \mathbf{I}_d \end{bmatrix}.$$

This amounts to a matrix \mathbf{C}_k with a single element $\sigma_{\star,i}^{\frac{2(L-1)}{L}}$ and 0 elsewhere. This gives an eigenvalue $\sigma_{\star,i}^{\frac{2(L-1)}{L}}$ for all $i \in [r]$, with multiplicity $d-r$.

Lastly, we can consider the diagonal components of $\mathbf{H}_{m,\ell}$ that is a function of the initialization scale α . For this case, each \mathbf{C}_k can be written as an outer product $\mathbf{v}\mathbf{v}^\top$, where

$$\mathbf{v}^\top = \left[\sigma_{\star,i}^{1-\frac{1}{L}} \alpha^0 \quad \sigma_{\star,i}^{1-\frac{2}{L}} \alpha \quad \sigma_{\star,i}^{1-\frac{3}{L}} \alpha^2 \quad \dots \quad \sigma_{\star,i}^0 \alpha^{L-1} \right]^\top. \quad (27)$$

Similarly, the non-zero eigenvalue is simply $\|\mathbf{v}\|_2^2$, which corresponds to

$$\|\mathbf{v}\|_2^2 = \sum_{\ell=0}^{L-1} \left(\sigma_{\star,k}^{1-\frac{1}{L}-\frac{1}{L}\ell} \cdot \alpha^\ell \right)^2.$$

This completes the proof. \square

Lemma 2 (Balancing). *Suppose we run GD on the deep matrix factorization loss in Equation (1) with learning rate $\eta < \frac{2\sqrt{2}}{L\sigma_{*,1}^{\frac{2}{L-2}}}$, where $\sigma_{*,1}$ is the first singular value of $\mathbf{M}_* \in \mathbb{R}^{d \times d}$. Let $\sigma_{i,\ell}$ denote the i -th singular value of the ℓ -th layer matrix. If the initialization scale α satisfies $0 < \alpha < \left(\ln \left(\frac{2\sqrt{2}}{\eta L \sigma_{*,1}^{\frac{2}{L-2}}} \right) \cdot \frac{\sigma_{*,1}^{\frac{4}{L}}}{L^2 \cdot 2^{\frac{2L-3}{L}}} \right)^{\frac{1}{4}}$, we have $|\sigma_{i,L}^2(t+1) - \sigma_{i,\ell}^2(t+1)| < c |\sigma_{i,L}^2(t) - \sigma_{i,\ell}^2(t)|$ for some $0 < c \leq 1$.*

Proof. From Proposition 1, we can re-write the loss in terms of the singular values:

$$\frac{1}{2} \|\mathbf{W}_{L:1}(t) - \mathbf{M}^*\|_{\text{F}}^2 = \frac{1}{2} \sum_{k=1}^r (\sigma_k(\Sigma_{L:1}(t)) - \sigma_{*,k})^2 = \sum_{k=1}^r \frac{1}{2} \left(\prod_{\ell=1}^L \sigma_{\ell,k} - \sigma_{*,k} \right)^2. \quad (28)$$

We aim to prove that balancing occurs on each singular value scalar index k , and so we focus on the scalar loss

$$\frac{1}{2} \left(\prod_{\ell=1}^L \sigma_{\ell,k} - \sigma_{*,k} \right)^2 =: \frac{1}{2} \left(\prod_{\ell=1}^L \sigma_{\ell} - \sigma_* \right)^2,$$

and omit the dependency on k for ease of exposition. Then, let us define the balancing dynamics between σ_i and σ_j as $b_{i,j}^{(t+1)} := \left(\sigma_i^{(t+1)} \right)^2 - \left(\sigma_j^{(t+1)} \right)^2$ and $\pi^{(t)} := \prod_{\ell=1}^L \sigma_{\ell}(t)$ for the product of singular values at iteration t . We can simplify the balancing dynamics as follows:

$$\begin{aligned} b_{i,j}^{(t+1)} &= \left(\sigma_i^{(t+1)} \right)^2 - \left(\sigma_j^{(t+1)} \right)^2 \\ &= \left(\sigma_i^{(t)} - \eta \left(\pi^{(t)} - \sigma_* \right) \frac{\pi^{(t)}}{\sigma_i^{(t)}} \right)^2 - \left(\sigma_j^{(t)} - \eta \left(\pi^{(t)} - \sigma_* \right) \frac{\pi^{(t)}}{\sigma_j^{(t)}} \right)^2 \\ &= \left(\sigma_i^{(t)} \right)^2 - \left(\sigma_j^{(t)} \right)^2 + \eta^2 \left(\pi^{(t)} - \sigma_* \right)^2 \left(\frac{\left(\pi^{(t)} \right)^2}{\left(\sigma_i^{(t)} \right)^2} - \frac{\left(\pi^{(t)} \right)^2}{\left(\sigma_j^{(t)} \right)^2} \right) \\ &= \left(\left(\sigma_i^{(t)} \right)^2 - \left(\sigma_j^{(t)} \right)^2 \right) \left(1 - \eta^2 \left(\pi^{(t)} - \sigma_* \right)^2 \frac{\left(\pi^{(t)} \right)^2}{\left(\sigma_i^{(t)} \right)^2 \left(\sigma_j^{(t)} \right)^2} \right) \\ &= b_{i,j}^{(t)} \left(1 - \eta^2 \left(\pi^{(t)} - \sigma_* \right)^2 \frac{\left(\pi^{(t)} \right)^2}{\left(\sigma_i^{(t)} \right)^2 \left(\sigma_j^{(t)} \right)^2} \right). \end{aligned}$$

Then, in order to show that $|b_{i,j}^{(t+1)}| < c |b_{i,j}^{(t)}|$, we need to prove that

$$\left| 1 - \eta^2 \left(\pi^{(t)} - \sigma_* \right)^2 \frac{\left(\pi^{(t)} \right)^2}{\left(\sigma_i^{(t)} \right)^2 \left(\sigma_j^{(t)} \right)^2} \right| < c,$$

for all iterations t and for some $0 < c \leq 1$. Now we introduce a definition called gradient flow solution sharpness (GFS sharpness) before we proceed.

Definition 3 (GFS Sharpness). *The GFS sharpness denoted by $\psi(x)$ is the sharpness achieved by the global minima which lies in the same GF trajectory of x (i.e., $\|\nabla^2 L(z)\|$ such that $L(z) = 0$ and $z = GF(x)$, where $GF(\cdot)$ denotes the gradient flow solution).*

Then, we complete the following two steps:

- 1782 (i) We show that for all scalars σ in the trajectory, if $\psi(\sigma) < \frac{2\sqrt{1+c}}{\eta}$ and $\sigma > 0$, then it holds that
 1783
$$\sum_{i=1}^{\min\{2, L-1\}} \frac{\eta^2(\pi(\sigma) - \sigma_*)^2 \pi^2(\sigma)}{\sigma_{L-i}^2 \sigma_L^2} \leq 1 + c$$
, where $\pi(\sigma)$ denotes the product given the trajectory
 1784 of all σ_i . This case is analyzed when $\pi(\sigma) \in [0, \sigma_*)$ where $0 < c < 1$ and when $\pi(\sigma) > \sigma_*$
 1785 where $c = 1$.
 1786
 1787 (ii) If $\sum_{i=1}^{\min\{2, L-1\}} \frac{\eta^2(\pi(\sigma) - \sigma_*)^2 \pi^2(\sigma)}{\sigma_{L-i}^2 \sigma_L^2} \leq 1 + c$, then iterates become more balanced, i.e., $|b_{i,j}^{(t+1)}| <$
 1788 $c|b_{i,j}^{(t)}|$.
 1789

1791 We prove (i) in Lemma 3 and (ii) in Lemma 4. Both of the proofs are originally from Kreisler et al.
 1792 (2023), which we adapted using our notation for ease of the reader. Then, in Lemma 5, we show
 1793 that for each σ_* , as long as the initialization scale satisfies

$$\alpha < \left(\ln \left(\frac{2\sqrt{2}}{\eta L \sigma_*^{2-\frac{2}{L}}} \right) \cdot \frac{\sigma_*^{\frac{4}{L}}}{L^2 \cdot 2^{\frac{2L-3}{L}}} \right)^{\frac{1}{4}},$$

1798 then it holds that the GFS sharpness satisfies $\psi(\sigma) < \frac{2\sqrt{2}}{\eta}$, which is the necessary condition for
 1799 balancing. Then, to satisfy this condition for all singular values $\sigma_{*,i}$ for all $i \in [r]$, we need

$$\alpha < \left(\ln \left(\frac{2\sqrt{2}}{\eta L \sigma_{*,1}^{2-\frac{2}{L}}} \right) \cdot \frac{\sigma_{*,1}^{\frac{4}{L}}}{L^2 \cdot 2^{\frac{2L-3}{L}}} \right)^{\frac{1}{4}} \implies \eta < \frac{2\sqrt{2}}{L \sigma_{*,1}^{2-\frac{2}{L}}}, \quad (29)$$

1805 for the validity of the initialization scale. Thus, as long as the conditions in Equation (29) hold, we
 1806 will have balancing. This completes the proof. \square

1807 **Lemma 3.** *If the GFS sharpness $\psi(\sigma) \leq \frac{2\sqrt{1+c}}{\eta}$ and $\sigma > 0$, then $\sum_{i=1}^{\min\{2, L-1\}} \frac{\eta^2(\pi(\sigma) - \sigma_*)^2 \pi^2(\sigma)}{\sigma_{[L-i]\sigma_{[D]}}^2} \leq$
 1808 $(1 + c)$ for some $0 < c \leq 1$.*

1810 *Proof.* We will consider two cases: (i) $\pi(\sigma) \in [0, \sigma_*)$ and (ii) $\pi(\sigma) > \sigma_*^2$.

1813 **Case 1:** Let $\sigma \in \mathbb{R}^D$ and consider the case where $\pi(\sigma) \in [0, \sigma_*)$. Then, we have

$$\sum_{i=1}^{\min\{2, L-1\}} \frac{\eta^2(\pi(\sigma) - \sigma_*)^2 \pi^2(\sigma)}{\sigma_{L-i}^2 \sigma_L^2} \leq \frac{\eta^2 \pi^2(\sigma)}{\sigma_{L-i}^2 \sigma_L^2}.$$

1818 Our goal is to show that if $\psi(\sigma) \leq \frac{2\sqrt{1+c}}{\eta}$ for some $0 < c < 1$ then,

$$\sum_{i=1}^{\min\{2, L-1\}} \frac{\eta^2(\pi(\sigma) - 1)^2 \pi^2(\sigma)}{\sigma_{L-i}^2 \sigma_L^2} \leq \frac{\eta^2 \pi^2(\sigma)}{\sigma_{L-i}^2 \sigma_L^2} \leq 1 + c.$$

1824 Since the GFS sharpness is constant for all the weights on the gradient flow (GF) trajectory by
 1825 definition, we can focus on the singular values (or weights) at the global minima. Consider $z =$
 1826 $\text{GF}(\sigma)$, the GF solution of σ . In Lemma 6, we proved that GF preserves unbalancedness, such that
 1827 $\sigma_l^2 - \sigma_m^2 = z_l^2 - z_m^2$ for all layers. Hence, it is sufficient to show that $\sum_{i=1}^{\min\{2, L-1\}} \frac{\eta^2 \pi^2(z)}{z_{L-i}^2 z_L^2} \leq 1 + c$
 1828 in order to ensure $\sum_{i=1}^{\min\{2, L-1\}} \frac{\eta^2 \pi^2(\sigma)}{\sigma_{L-i}^2 \sigma_L^2} \leq 1 + c$. Note that $\pi(z) = \sigma_*$, since it lies on the global
 1829 minima. Then,
 1830

$$\sum_{i=1}^{\min\{2, L-1\}} \frac{\eta^2 \pi^2(z)}{z_{L-i}^2 z_L^2} = \sum_{i=1}^{\min\{2, L-1\}} \frac{\eta^2 \sigma_*^2}{z_{L-i}^2 z_L^2}. \quad (30)$$

1834
 1835 ²We ignore the case $\pi^{(t)} = \sigma_*$ when we get $b_{i,j}^{(t+1)} = b_{i,j}^{(t)}$. Since the occurrence $\pi^{(t)} = \sigma_*$ holds with a
 probability of zero where EOS ceases to exist.

From Lemma 7, we know that the sharpness at the global minima is given as

$$\psi(\sigma) = \|\nabla^2 L(z)\| = \sum_{i=1}^L \frac{\sigma_*^2}{z_i^2}. \quad (31)$$

This immediately implies that $\frac{\sigma_*^2}{z_L^2} \leq \psi(\sigma)$ and equivalently, $\exists \alpha \in [0, 1]$ such that $\frac{\sigma_*^2}{z_L^2} = \alpha\psi(\sigma)$. Therefore, we have

$$\sum_{i=1}^{\min\{2, L-1\}} \frac{\sigma_*^2}{z_{L-i}^2} \leq (1 - \alpha)\psi(\sigma). \quad (32)$$

Substituting Equations (31) and (32) into the expression we aim to bound, we obtain

$$\sum_{i=1}^{\min\{2, L-1\}} \frac{\eta^2(\pi(\sigma) - \sigma_*^2)^2 \pi^2(\sigma)}{\sigma_{L-i}^2 \sigma_L^2} = \sum_{i=1}^{\min\{2, L-1\}} \frac{\eta^2 \sigma_*^2}{z_{L-i}^2 z_L^2} \leq \eta^2 \alpha (1 - \alpha) \psi^2(\sigma) \leq \frac{\eta^2}{4} \phi^2(\sigma) \leq 1 + c,$$

where we used the fact that the maximum of $\alpha(1 - \alpha)$ is $\frac{1}{4}$ when $\alpha = \frac{1}{2}$ and $\psi(\sigma) \leq \frac{2\sqrt{1+c}}{\eta}$. Thus, if $\psi(\sigma) \leq \frac{2\sqrt{1+c}}{\eta}$, then for every weight σ lying on its GF trajectory, we have

$$\sum_{i=1}^{\min\{2, L-1\}} \frac{\eta^2(\pi(\sigma) - \sigma_*)^2 \pi^2(\sigma)}{\sigma_{L-i}^2 \sigma_L^2} \leq 1 + c.$$

Case 2: Consider the case in which $\pi(\sigma) > \sigma_*$. We already have that $\sigma > 0$ throughout the trajectory (refer to Lemma 3.11 in Kreisler et al. (2023)) and so $\pi(\sigma) > 0$. So, the GD update from σ_i will also stay positive

$$\sigma_i - \eta(\pi(\sigma) - \sigma_*)\pi(\sigma) \frac{1}{\sigma_i} > 0.$$

From this, we get

$$2 > \frac{\eta(\pi(\sigma) - \sigma_*)\pi(\sigma)}{\sigma_i^2} > 0,$$

This implies $\sum_{i=1}^{\min\{2, L-1\}} \frac{\eta^2(\pi(\sigma) - \sigma_*)^2 \pi^2(\sigma)}{\sigma_{L-i}^2 \sigma_L^2} \leq (1 + c)$ with $c = 1$. This completes the proof. \square

Lemma 4. If $\sum_{i=1}^{\min\{2, L-1\}} \frac{\eta^2(\pi(\sigma) - \sigma_*)^2 \pi^2(\sigma)}{\sigma_{L-i}^2 \sigma_L^2} \leq 1 + c$ for $i, j \in [L]$ for some $0 < c \leq 1$, then

$$|b_{i,j}^{(t+1)}| < c |b_{i,j}^{(t)}|.$$

Proof. Recall that the condition for balancing was given by

$$b_{i,j}^{(t+1)} = b_{i,j}^{(t)} \left(1 - \eta^2 (\pi^{(t)} - \sigma_*)^2 \frac{\pi^{(t)2}}{(\sigma_i^{(t)})^2 (\sigma_j^{(t)})^2} \right). \quad (33)$$

WLOG, suppose that the σ are sorted such that $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_L$. We know that

$$\sum_{i=1}^{\min\{2, L-1\}} \frac{\eta^2(\pi(\sigma) - \sigma_*)^2 \pi^2(\sigma)}{\sigma_{L-i}^2 \sigma_L^2} \leq 1 + c,$$

which implies

$$\frac{\eta^2(\pi(\sigma) - \sigma_*)^2 \pi^2(\sigma)}{\sigma_{L-1}^2 \sigma_L^2} < 1 + c \quad \text{and} \quad \frac{\eta^2(\pi(\sigma) - \sigma_*)^2 \pi^2(\sigma)}{\sigma_i^2 \sigma_j^2} < \frac{1 + c}{2}, \quad (34)$$

for all $i \in [L]$, $j \in [L - 2]$ and $i < j$. Notice that the latter inequality comes from the fact that

$$\frac{\eta^2(\pi(\sigma) - \sigma_\star)^2 \pi^2(\sigma)}{\sigma_{L-2}^2 \sigma_L^2} + \frac{\eta^2(\pi(\sigma) - \sigma_\star)^2 \pi^2(\sigma)}{\sigma_{L-2}^2 \sigma_L^2} < \frac{\eta^2(\pi(\sigma) - \sigma_\star)^2 \pi^2(\sigma)}{\sigma_{L-1}^2 \sigma_L^2} + \frac{\eta^2(\pi(\sigma) - \sigma_\star)^2 \pi^2(\sigma)}{\sigma_{L-2}^2 \sigma_L^2} < 1 + c,$$

which implies that

$$2 \frac{\eta^2(\pi(\sigma) - \sigma_\star)^2 \pi^2(\sigma)}{\sigma_{L-2}^2 \sigma_L^2} < 1 + c \implies \frac{\eta^2(\pi(\sigma) - \sigma_\star)^2 \pi^2(\sigma)}{\sigma_{L-2}^2 \sigma_L^2} < \frac{1 + c}{2},$$

and since σ are sorted, it holds for all other σ . Therefore from Equation (33), we have for all $i \in [L - 2]$,

$$b_{i,i+1}^{(t+1)} < c b_{i,i+1}^{(t)} \quad \text{and} \quad b_{L-2,L}^{(t+1)} < c b_{L-2,L}^{(t)} \quad \text{and} \quad -c b_{L-1,L}^{(t)} < b_{L-1,L}^{(t+1)} < c b_{L-1,L}^{(t)}. \quad (35)$$

Then, notice that since we initialized all of the singular values σ_ℓ for $\ell \in [L - 1]$ to be the same, they follow the same dynamics. Since we already showed that $|b_{L-1,L}^{(t+1)}| < c |b_{L-1,L}^{(t)}|$, it must follow that

$$|b_{i,j}^{(t+1)}| < c |b_{i,j}^{(t)}| \quad \text{for } i, j \in [L].$$

This completes the proof. □

Lemma 5. Consider running GD with learning rate η in Equation (2) on the scalar loss

$$\mathcal{L}(\{\sigma_i\}_{i=1}^d) = \frac{1}{2} \left(\prod_{i=1}^L \sigma_i - \sigma_\star \right)^2,$$

with initialization $\sigma_L(0) = 0$ and $\sigma_\ell(0) = \alpha$ for all $\ell \in [L - 1]$. If $\alpha < \left(\ln \left(\frac{2\sqrt{2}}{\eta L \sigma_\star^{2-\frac{2}{L}}} \right) \cdot \frac{\sigma_\star^{\frac{4}{L}}}{L^{2 \cdot 2 \frac{L-3}{L}}} \right)^{\frac{1}{4}}$, then the GFS sharpness $\psi(\sigma) \leq \frac{2\sqrt{1+c}}{\eta}$ for some $0 < c < 1$.

Proof. Since the singular values σ_ℓ for all $\ell \in [L - 1]$ are initialized to α , note that they all follow the same dynamics. Then, let us define

$$y := \sigma_1 = \dots = \sigma_{L-1} \quad \text{and} \quad x := \sigma_L.$$

The gradient flow (GF) solution is the intersection between

$$xy^{L-1} = \sigma_\star \quad \text{and} \quad x^2 - y^2 = -\alpha^2,$$

where the first condition comes from convergence and the second comes from the conservation flow law of GF which we prove in Lemma 6. Then, if we can find a solution at the intersection such that

$$(\hat{x}(\alpha), \hat{y}(\alpha)) = \begin{cases} xy^{L-1} = \sigma_\star \\ x^2 - y^2 = -\alpha^2, \end{cases} \quad (36)$$

solely in terms of α , then we can plug in $(\hat{x}(\alpha), \hat{y}(\alpha))$ into the GFS³ from Lemma 7

$$\psi(\hat{x}(\alpha), \hat{y}(\alpha)) = \psi(\sigma) = \sum_{i=1}^L \frac{\sigma_\star^2}{\sigma_i^2} = \sigma_\star^2 \left(\frac{1}{\hat{x}(\alpha)^2} + \frac{L-1}{\hat{y}(\alpha)^2} \right) < \frac{2\sqrt{2}}{\eta}$$

and solve to find an upper bound on α . **The strict inequality ensures that we can find a c in $c \in [0, 1)$ such that $\psi(\alpha) < \frac{2\sqrt{1+c}}{\eta}$.** However, the intersection $(\hat{x}(\alpha), \hat{y}(\alpha))$ is a $2L$ -th order polynomial in $\hat{y}(\alpha)$ which does not have a straightforward closed-form solution solely in terms of α . Hence, we

³Note that throughout the proof $(\hat{x}(\alpha), \hat{y}(\alpha))$ denotes the gradient flow solution as function of α . It does not refer to the GF trajectory.

aim to find the upper bound on α by using a calculus of variations. By plugging in x , the solution $\hat{y}(\alpha)$ satisfies

$$y^{2L} - \alpha^2 y^{2L-2} = \sigma_*^2.$$

Then, by differentiating the relation with respect to α , we obtain the following variational relation:

$$\begin{aligned} 2Ly^{2L-1}dy - \alpha^2 2(L-1)y^{2L-3}dy - 2\alpha y^{2L-2}d\alpha &= 0 \\ \implies y^{2L-3}(y^2L - \alpha^2(L-1))dy &= \alpha y^{2(L-1)}d\alpha \\ \implies dy &= \frac{y\alpha}{(y^2L - \alpha^2(L-1))}d\alpha, \end{aligned} \quad (37)$$

where we used the fact that $y^{2L-2} > 0$ from Lemma 3 in the last line. Then, in order to have $\frac{dy}{d\alpha} > 0$, we need $y > \sqrt{\frac{L-1}{L}}\alpha$, which is always true since $y > \alpha$ from initialization. Then, since $\alpha \rightarrow 0$, $\lim_{\alpha \rightarrow 0} \hat{y}(\alpha) = \sigma_*^{\frac{1}{L}}$ and $\lim_{\alpha \rightarrow 0} \hat{x}(\alpha) = \sigma_*^{\frac{1}{L}}$, as it corresponds to exact balancing. Hence, $\frac{dy}{d\alpha} > 0$ implies as α increases from 0, $\hat{y}(\alpha)$ would increase from $\sigma_*^{\frac{1}{L}}$ and $\hat{y}(\alpha)$ is an increasing function of α . Similarly, the intersection at the global minima would satisfy the following relation for $\hat{x}(\alpha)$:

$$\begin{aligned} x^{(2+\frac{2}{L-1})} + x^{\frac{2}{L-1}}\alpha^2 &= \sigma_*^{\frac{2}{L-1}} \\ \implies \left(2 + \frac{2}{L-1}\right)x^{1+\frac{2}{L-1}}dx + \left(\frac{2}{L-1}\right)\alpha^2 x^{\frac{2}{L-1}-1}dx + x^{\frac{2}{L-1}}(2\alpha d\alpha) &= 0 \\ \implies dx &= \frac{-\alpha}{\left(\frac{L}{L-1}x + \frac{\alpha^2}{L-1}\frac{1}{x}\right)}d\alpha. \end{aligned} \quad (38)$$

Note that since $x > 0$, we will always have $\frac{dx}{d\alpha} < 0$. Then, since $\lim_{\alpha \rightarrow 0} \hat{x}(\alpha) = \sigma_*^{\frac{1}{L}}$, $\frac{dx}{d\alpha} < 0$ implies that as α increases, $\hat{x}(\alpha)$ would decrease from $\sigma_*^{\frac{1}{L}}$. Now, with the variational relations $\frac{d\hat{x}}{d\alpha}$ and $\frac{d\hat{y}}{d\alpha}$ in place, we aim to find $\frac{d\Psi}{d\alpha}$:

$$\begin{aligned} \Psi(\alpha) &:= \psi(\hat{x}(\alpha), \hat{y}(\alpha)) = \sigma_*^2 \left(\frac{1}{\hat{x}(\alpha)^2} + \frac{L-1}{\hat{y}(\alpha)^2} \right) \\ \implies d\Psi &= \sigma_*^2 \left(-\frac{2}{\hat{x}^3}d\hat{x} - \frac{2(L-1)}{\hat{y}^3}d\hat{y} \right) \\ \implies d\Psi &= \frac{1}{\hat{x}^3} \left[\frac{2\alpha\sigma_*^2}{\left(\frac{L}{L-1}\hat{x} + \left(\frac{\alpha^2}{L-1}\right)\frac{1}{\hat{x}}\right)} \right] d\alpha - \left[\frac{(L-1)}{\hat{y}^3} \frac{2\alpha\hat{y}\sigma_*^2}{(\hat{y}^2L - \alpha^2(L-1))} \right] d\alpha \\ \implies d\Psi &= \left[\frac{1}{\hat{x}^4 + \frac{\alpha^2}{L}\hat{x}^2} - \frac{1}{(\hat{y}^4 - \alpha^2\hat{y}^2\frac{(L-1)}{L})} \right] 2\frac{(L-1)\sigma_*^2}{L}\alpha d\alpha \\ \implies d\Psi &= G(\alpha)d\alpha, \end{aligned}$$

where we defined $G(\alpha) := \left[\frac{1}{\hat{x}^4 + \frac{\alpha^2}{L}\hat{x}^2} - \frac{1}{(\hat{y}^4 - \alpha^2\hat{y}^2\frac{(L-1)}{L})} \right] 2\frac{(L-1)\sigma_*^2}{L}\alpha$ and used the notation $\hat{x} = \hat{x}(\alpha)$ and $\hat{y} = \hat{y}(\alpha)$ for simplicity.

Next, we will show the three following steps:

- (i) Prove that $G(\alpha) > 0$ for all $\alpha > 0$ to show that the sharpness $\Psi(\alpha)$ is an increasing function of α .
- (ii) Solve the differential $d\Psi$ to find the relationship between $d\Psi$ and $\Psi(\alpha)$.
- (iii) Find an upper bound on a part of $\frac{d\Psi}{\Psi(\alpha)}$ found in Step 2.

1998
1999
2000
2001
2002
2003
2004
2005
2006
2007
2008
2009
2010
2011
2012
2013
2014
2015
2016
2017
2018
2019
2020
2021
2022
2023
2024
2025
2026
2027
2028
2029
2030
2031
2032
2033
2034
2035
2036
2037
2038
2039
2040
2041
2042
2043
2044
2045
2046
2047
2048
2049
2050
2051

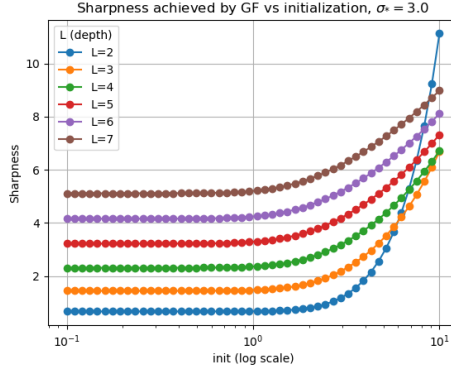


Figure 22: Sharpness $\Psi(\alpha)$ as a function of initialization α . The theoretical approximation bound $\Psi = \Psi_0 \exp\left(\frac{L^2 \cdot 2^{\frac{2(L-1)}{L}}}{2\sigma_*^{\frac{4}{L}}} \alpha^4\right)$ serves as proxy upper bound to this increasing function.

These series of steps comes from the fact that the intersection does not have a closed-form solution. The goal is to find a function in which we can upper bound $\frac{d\Psi}{\Psi(\alpha)}$ with a function with a closed-form solution to find a bound on α such that the sharpness $\psi(\alpha) < \frac{2\sqrt{2}}{\eta}$.

Step 1: Prove $G(\alpha) > 0$ to show sharpness $\Psi(\alpha)$ is an increasing function of α .

There have been several lines of work such as those by Kreisler et al. (2023) and Marion & Chizat (2024) which showed that GD would decrease the sharpness of the solution. The more balanced the solution (which corresponds to smaller α), the smaller the sharpness. We prove this again here:

$$\begin{aligned}
 G(\alpha) > 0 &\implies \hat{y}^4 - \alpha^2 \hat{y}^2 \frac{(L-1)}{L} > \hat{x}^4 + \frac{\alpha^2}{L} \hat{x}^2 \\
 &\implies (\hat{y}^4 - \hat{x}^4) > \alpha^2 \left(\frac{1}{L} \hat{x}^2 + \hat{y}^2 \frac{(L-1)}{L} \right) \\
 &\implies \underbrace{(\hat{y}^2 - \hat{x}^2)}_{=\alpha^2} (\hat{y}^2 + \hat{x}^2) > \alpha^2 \left(\frac{1}{L} \hat{x}^2 + \hat{y}^2 \frac{(L-1)}{L} \right) \\
 &\implies \hat{x}^2 \left(1 - \frac{1}{L}\right) + \hat{y}^2 \frac{1}{L} > 0,
 \end{aligned}$$

where the last inequality always holds since we have $L > 2$. This proves that Ψ is an increasing function of α since for $d\Psi = G(\alpha)d\alpha$, as it always holds that $G(\alpha) > 0$ for any $L > 2$ and $\alpha > 0$.

Step 2: Solve the differential to establish the relation between $\Psi(\alpha)$ and α .

Rewriting the expression for sharpness and establishing an equation we have

$$\Psi(\alpha) = \sigma_*^2 \left(\frac{1}{\hat{x}(\alpha)^2} + \frac{L-1}{\hat{y}(\alpha)^2} \right) \implies \Psi(\alpha) = \sigma_*^2 \left(\frac{\hat{y}^2 + (L-1)\hat{x}^2}{\hat{x}^2 \hat{y}^2} \right) \quad (39)$$

$$\implies \frac{\hat{y}^2}{L} + \left(1 - \frac{1}{L}\right) \hat{x}^2 = \frac{\Psi(\alpha) \hat{x}^2 \hat{y}^2}{L \sigma_*^2}. \quad (40)$$

2052
2053
2054
2055
2056
2057
2058
2059
2060
2061
2062
2063
2064
2065
2066
2067
2068
2069
2070
2071
2072
2073
2074
2075
2076
2077
2078
2079
2080
2081
2082
2083
2084
2085
2086
2087
2088
2089
2090
2091
2092
2093
2094
2095
2096
2097
2098
2099
2100
2101
2102
2103
2104
2105

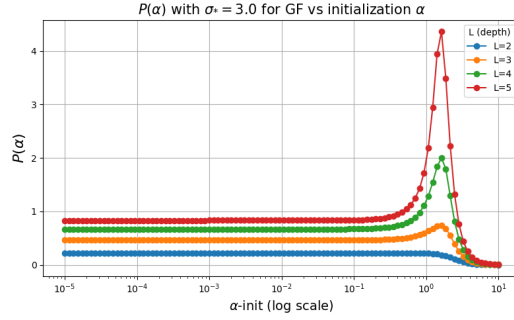


Figure 23: $P(\alpha)$ for GF has a unique maxima at $\alpha = \frac{\sigma_*^{\frac{1}{L}}}{\left(\frac{1}{\sqrt{L(L-2)}}(1 + \frac{1}{L(L-2)})^{\frac{L-1}{2}}\right)^{\frac{1}{L}}}$ for $L > 2$.

Now, we revisit the original differential between $\Psi(\alpha)$ and α :

$$\begin{aligned} d\Psi &= \left[\frac{1}{(\hat{x}^4 + \frac{\alpha^2}{L}\hat{x}^2)} - \frac{1}{(\hat{y}^4 - \alpha^2\hat{y}^2\frac{(L-1)}{L})} \right] 2\frac{(L-1)\sigma_*^2}{L}\alpha d\alpha \\ \implies d\Psi &= \frac{\hat{y}^4 - \hat{x}^4 - \alpha^2(\frac{\hat{x}^2}{L} + (1 - \frac{1}{L})\hat{y}^2)}{(\hat{x}^4 + \frac{\alpha^2}{L}\hat{x}^2)(\hat{y}^4 - \alpha^2\hat{y}^2\frac{(L-1)}{L})} 2\left(1 - \frac{1}{L}\right)\sigma_*^2\alpha d\alpha \\ \implies d\Psi &= \frac{\alpha^2\left(\frac{\hat{y}^2}{L} + (1 - \frac{1}{L})\hat{x}^2\right)}{(\hat{x}^4 + \frac{\alpha^2}{L}\hat{x}^2)(\hat{y}^4 - \alpha^2\hat{y}^2\frac{(L-1)}{L})} 2\left(1 - \frac{1}{L}\right)\sigma_*^2\alpha d\alpha \end{aligned} \quad (41)$$

Using the expression for $\frac{\hat{y}^2}{L} + (1 - \frac{1}{L})\hat{x}^2$ derived in Equation (40) and plugging it into Equation (41), we obtain

$$\begin{aligned} d\Psi &= \frac{\alpha^2\left(\frac{\Psi(\alpha)\hat{x}^2\hat{y}^2}{L\sigma_*^2}\right)}{(\hat{x}^4 + \frac{\alpha^2}{L}\hat{x}^2)(\hat{y}^4 - \alpha^2\hat{y}^2\frac{(L-1)}{L})} 2\left(1 - \frac{1}{L}\right)\sigma_*^2\alpha d\alpha \\ \implies \frac{d\Psi}{\Psi(\alpha)} &= \frac{2}{(\hat{x}^2 + \frac{\alpha^2}{L})(\hat{y}^2 - \alpha^2\frac{(L-1)}{L})} \left(\frac{1}{L} - \frac{1}{L^2}\right)\alpha^3 d\alpha \\ \implies \frac{d\Psi}{\Psi(\alpha)} &= \frac{2}{(\hat{x}^2 + \frac{\alpha^2}{L})^2} \left(\frac{1}{L} - \frac{1}{L^2}\right)\alpha^3 d\alpha \quad (42) \\ \implies \frac{d\Psi}{\Psi(\alpha)} &= P(\alpha)\alpha^3 d\alpha, \quad (43) \end{aligned} \quad (44)$$

where we have defined $P(\alpha) := \frac{2}{(\hat{x}^2 + \frac{\alpha^2}{L})^2} \left(\frac{1}{L} - \frac{1}{L^2}\right)$.

Solving the differential $\frac{d\Psi}{\Psi(\alpha)} = P(\alpha)\alpha^3 d\alpha$ in exact closed-form is difficult since \hat{x} is also an function of α . However, in Step 1, we proved that $\Psi(\alpha)$ is an increasing function of α , and so instead of solving exactly, we can find a differential equation $\frac{d\Psi}{\Psi(\alpha)} = F(\alpha)\alpha^3 d\alpha$ with $F(\alpha) > P(\alpha)$ such that $F(\alpha)$ is more increasing, and use it to solve the PDE instead. Though, note that the initialization limit on α that would be found after solving the surrogate PDE $\frac{d\Psi}{\Psi(\alpha)} = F(\alpha)\alpha^3 d\alpha$ would be smaller than the α if it was found using the original PDE $\frac{d\Psi}{\Psi(\alpha)} = P(\alpha)\alpha^3 d\alpha$.

Step 3: Finding an upper bound function and solving for initialization.

Note that the original coefficient in $\frac{d\Psi}{\Psi(\alpha)} = P(\alpha)\alpha^3 d\alpha$ is of the form

$$P(\alpha) = \frac{2}{(\hat{x}^2 + \frac{\alpha^2}{L})^2} \left(\frac{1}{L} - \frac{1}{L^2}\right) \quad (45)$$

Let us consider the two corner cases for α . We showed before that $\lim_{\alpha \rightarrow 0} \hat{x}(\alpha) = \sigma_*^{\frac{1}{L}}$, so

$$\lim_{\alpha \rightarrow 0} P(\alpha) = \frac{2 \left(\frac{1}{L} - \frac{1}{L^2} \right)}{\sigma_*^{\frac{4}{L}}}$$

As $\alpha \rightarrow \infty$, we have $\lim_{\alpha \rightarrow \infty} P(\alpha) \rightarrow 0$ since $\lim_{\alpha \rightarrow \infty} \hat{x} = 0$. Furthermore, we have

$$\begin{aligned} P'(\alpha) &= \frac{-4}{\left(\hat{x}^2 + \frac{\alpha^2}{L}\right)^3} \left(\frac{1}{L} - \frac{1}{L^2} \right) \left(2\hat{x} \frac{d\hat{x}}{d\alpha} + \frac{2\alpha}{L} \right) \\ &= \frac{-4 \left(\frac{1}{L} - \frac{1}{L^2} \right)}{\left(\hat{x}^2 + \frac{\alpha^2}{L}\right)^3} \left(2\hat{x} \left(\frac{-\alpha}{\left(\frac{L}{L-1}\hat{x} + \frac{\alpha^2}{L-1} \frac{1}{\hat{x}} \right)} \right) + \frac{2\alpha}{L} \right) \\ &= \frac{8\alpha \left(\frac{1}{L} - \frac{1}{L^2} \right)}{\left(\hat{x}^2 + \frac{\alpha^2}{L}\right)^3} \left(\frac{L-1}{L + \frac{\alpha^2}{\hat{x}^2}} - \frac{1}{L} \right) \end{aligned}$$

For $L = 2$, we always have $P'(\alpha) < 0$. Hence, choosing $F(\alpha) = \lim_{\alpha \rightarrow 0} P(\alpha) = \frac{2 \left(\frac{1}{L} - \frac{1}{L^2} \right)}{\sigma_*^{\frac{4}{L}}}$, will serve as the correct upper bound.

Then, let us consider $L > 2$. We can see that α at which $\hat{x}(\alpha) = \frac{\alpha}{\sqrt{L(L-2)}}$ is the critical point of $P(\alpha)$. Further, we note that when $\hat{x}(\alpha) < \frac{\alpha}{\sqrt{L(L-2)}}$, $P'(\alpha) < 0$ meaning $P(\alpha)$ is decreasing. Since $\hat{x}(\alpha)$ is itself decreasing in α , this states there for any $\alpha > \alpha_{crit}$, $P(\alpha)$ is decreasing. α_{crit} is the solution of $\hat{x}(\alpha) = \frac{\alpha}{\sqrt{L(L-2)}}$.

For any $\alpha < \alpha_{crit}$, $P'(\alpha) > 0$, so $P(\alpha)$ is increasing. So, $P(\alpha_{crit})$ corresponds to the maximum of P in α . Choosing $F(\alpha) = P(\alpha_{crit})$, a constant allows us to find an upper bound function for the function in Equation (45). Furthermore, note that since $\frac{d\hat{x}}{d\alpha} < 0$ and $\hat{x} > 0$, $\alpha > 0$, there must be only one critical point of $P(\alpha)$ which is at α_{crit} .

Hence, we have $\hat{x}(\alpha_{crit}) = \frac{\alpha_{crit}}{\sqrt{L(L-2)}}$. From Equation (36), we also get $\hat{y} = \alpha_{crit} \sqrt{1 + \frac{1}{L(L-2)}}$ and

$$\begin{aligned} \left(\frac{\alpha_{crit}}{\sqrt{L(L-2)}} \right) \left(\alpha_{crit} \sqrt{1 + \frac{1}{L(L-2)}} \right)^{L-1} &= \sigma_* \\ \implies \alpha_{crit} &= \frac{\sigma_*^{\frac{1}{L}}}{\left(\frac{1}{\sqrt{L(L-2)}} \left(1 + \frac{1}{L(L-2)} \right)^{\frac{L-1}{2}} \right)^{\frac{1}{L}}} \end{aligned}$$

Using α_{crit} , we obtain the maximum of $P(\alpha)$ to be

$$\begin{aligned} P(\alpha_{crit}) &= \frac{2}{\left(\hat{x}(\alpha_{crit})^2 + \frac{\alpha_{crit}^2}{L}\right)^2} \left(\frac{1}{L} - \frac{1}{L^2} \right) \\ &= \frac{2}{\left(\alpha_{crit}^2 \left(\frac{1}{L(L-2)} + \frac{1}{L} \right)\right)^2} \left(\frac{1}{L} - \frac{1}{L^2} \right) \\ &= \frac{2}{\sigma_*^{\frac{4}{L}}} g(L) \end{aligned}$$

where $g(L) = \frac{\left(\frac{1}{L} - \frac{1}{L^2} \right) \left[\frac{1}{\sqrt{L(L-2)}} \left(1 + \frac{1}{L(L-2)} \right)^{\frac{L-1}{2}} \right]^{\frac{4}{L}}}{\left(\frac{1}{L(L-2)} + \frac{1}{L} \right)^2}$.

Now, choosing $F(\alpha) = P(\alpha_{crit})$, we integrate the upper bound function as

$$\begin{aligned} \int \frac{d\Psi}{\Psi} &= \frac{2g(L)}{\sigma_*^{\frac{4}{L}}} \int \alpha^3 d\alpha \\ \implies \ln\left(\frac{\Psi}{\Psi_0}\right) &= \frac{g(L)}{2\sigma_*^{\frac{4}{L}}} (\alpha^4) \\ \implies \Psi &= \Psi_0 \exp\left(\frac{g(L)}{2\sigma_*^{\frac{4}{L}}} \alpha^4\right) \end{aligned}$$

where $\Psi_0 = \lim_{\alpha \rightarrow 0} \Psi = L\sigma_*^{2-\frac{2}{L}}$. We verify this upper bound empirically from Figure 22, where we see a near exponential growth in sharpness as function of α .

Now, note that the function $\Psi = \Psi_0 \exp\left(\frac{g(L)}{2\sigma_*^{\frac{4}{L}}} \alpha^4\right)$ acts an upper bound to the original sharpness function of α and both are increasing in α (Step 1). So, solving for an initialization α -upper limit with $\Psi = \frac{2\sqrt{2}}{\eta}$ would mean that the original sharpness with this initialization would be less than $\frac{2\sqrt{1+c}}{\eta}$ for some $0 < c < 1$. Hence, α is restricted to

$$\alpha < \left(\ln\left(\frac{\frac{2\sqrt{2}}{\eta}}{L\sigma_*^{2-\frac{2}{L}}}\right) \cdot \frac{2\sigma_*^{\frac{4}{L}}}{g(L)} \right)^{\frac{1}{4}}.$$

We can simplify the bound further by finding an upper bound on $g(L)$ ⁴:

$$g(L) \leq \frac{\left(\left(1 + \frac{1}{L(L-2)}\right)^{\frac{L-1}{2}}\right)^{\frac{4}{L}}}{\left(\frac{1}{L(L-2)} + \frac{1}{L}\right)^2} \leq L^2 \cdot \left(\left(1 + \frac{1}{L(L-2)}\right)^{\frac{L-1}{2}}\right)^{\frac{4}{L}} \leq L^2 \cdot 2^{\frac{2(L-1)}{L}}.$$

Then, we get obtain a lower bound on α :

$$\alpha < \left(\ln\left(\frac{2\sqrt{2}}{\eta L \sigma_*^{2-\frac{2}{L}}}\right) \cdot \frac{2\sigma_*^{\frac{4}{L}}}{L^2 \cdot 2^{\frac{2(L-1)}{L}}}\right)^{\frac{1}{4}} = \left(\ln\left(\frac{2\sqrt{2}}{\eta L \sigma_*^{2-\frac{2}{L}}}\right) \cdot \frac{\sigma_*^{\frac{4}{L}}}{L^2 \cdot 2^{\frac{2(L-1)}{L}}}\right)^{\frac{1}{4}}$$

Hence, as long as α satisfies this upper bound, we will have balancing. This completes the proof. \square

Lemma 6. Consider the minimizing the loss

$$\mathcal{L}(\{\sigma_\ell\}_{\ell=1}^L) = \frac{1}{2} \left(\prod_{\ell=1}^L \sigma_\ell - \sigma_* \right)^2,$$

using gradient flow. Then, the balancedness between two singular values defined by $\sigma_\ell^2(t) - \sigma_m^2(t)$ for all $m, \ell \in [L]$ is constant for all t .

Proof. Notice that the result holds specifically for gradient flow and not descent. The dynamics of each scalar factor for gradient flow can be written as

$$\dot{\sigma}_\ell(t) = - \left(\prod_{\ell=1}^L \sigma_\ell(t) - \sigma_* \right) \cdot \prod_{i \neq \ell} \sigma_i(t)$$

⁴This includes the case for $L = 2$ since $\left(\frac{1}{L} - \frac{1}{L^2}\right) < L^2 \cdot 2^{\frac{2(L-1)}{L}}$ for $L = 2$.

Then, the time derivative of balancing is given as

$$\begin{aligned} \frac{\partial}{\partial t}(\sigma_\ell^2(t) - \sigma_m^2(t)) &= \sigma_\ell(t)\dot{\sigma}_\ell(t) - \sigma_m(t)\dot{\sigma}_m(t) \\ &= -\sigma_\ell(t) \left(\prod_{\ell=1}^L \sigma_\ell(t) - \sigma_\star \right) \cdot \prod_{i \neq \ell}^L \sigma_i(t) + \sigma_m(t) \left(\prod_{m=1}^L \sigma_m(t) - \sigma_\star \right) \cdot \prod_{j \neq m}^L \sigma_j(t). \\ &= 0. \end{aligned}$$

Hence, the quantity $\sigma_\ell^2(t) - \sigma_m^2(t)$ remains constant for all time t , hence preserving unbalancedness. \square

Lemma 7. Consider the scalar loss

$$\mathcal{L}(\{\sigma_i\}_{i=1}^L) = \frac{1}{2} \left(\prod_{i=1}^L \sigma_i - \sigma_\star \right)^2,$$

The sharpness at the global minima is given as $\|\nabla^2 \mathcal{L}\|_2 = \sum_{i=1}^L \frac{\sigma_\star^2}{\sigma_i^2}$.

Proof. The gradient is given by

$$\nabla_{\sigma_i} \mathcal{L} = \left(\prod_{\ell=1}^L \sigma_\ell(t) - \sigma_\star \right) \prod_{j \neq i}^L \sigma_j(t).$$

Then,

$$\nabla_{\sigma_j} \nabla_{\sigma_i} \mathcal{L} = \prod_{\ell \neq i}^L \sigma_\ell(t) \prod_{\ell \neq j}^L \sigma_\ell(t) + \left(\prod_{\ell=1}^L \sigma_\ell(t) - \sigma_\star \right) \prod_{\ell \neq j, \ell \neq i}^L \sigma_\ell(t)$$

Let $\pi(t) = \prod_{i=1}^L \sigma_i(t)$. Then, at the global minima, we have

$$\nabla_{\sigma_j} \nabla_{\sigma_i} \mathcal{L} = \frac{\pi^2}{\sigma_i \sigma_j} = \frac{\sigma_\star^2}{\sigma_i \sigma_j}$$

Thus, the sharpness of the largest eigenvalue is given as $\|\nabla^2 \mathcal{L}\|_2 = \sum_{i=1}^L \frac{\sigma_\star^2}{\sigma_i^2}$. \square

Theorem 2 (Stable Subspace Oscillations). Consider running GD on the deep matrix factorization loss in Equation (1) and denote the SVD of the target matrix as $\mathbf{M}_\star = \mathbf{U}_\star \Sigma_\star \mathbf{V}_\star^\top$, with distinct singular values $\sigma_{\star,1} > \dots > \sigma_{\star,r}$. Let Δ_i denote the i -th eigenvector of the Hessian with unit norm, λ_i the corresponding eigenvalue after strict balancing occurs and denote f_{Δ_i} as the 1-D function at the cross section of the loss landscape and the line following the direction of Δ_i passing the minima. Then, if the minima of f_{Δ_i} satisfy $f_{\Delta_i}^{(3)} > 0$ and $3[f_{\Delta_i}^{(3)}]^2 - f_{\Delta_i}^{(2)} f_{\Delta_i}^{(4)} > 0$, then 2-period orbit oscillation occurs in direction of Δ_i if $\eta > \frac{2}{\lambda_i}$.

Proof. First, we derive the eigenvectors of the Hessian of the training loss at convergence (i.e., $\mathbf{M}_\star = \mathbf{W}_{L:1}$). To obtain the eigenvectors of the Hessian of parameters $(\mathbf{W}_L, \dots, \mathbf{W}_2, \mathbf{W}_1)$, consider a small perturbation of the parameters:

$$\Theta := (\Delta \mathbf{W}_\ell + \mathbf{W}_\ell)_{\ell=1}^L = (\mathbf{W}_L + \Delta \mathbf{W}_L, \dots, \mathbf{W}_2 + \Delta \mathbf{W}_2, \mathbf{W}_1 + \Delta \mathbf{W}_1).$$

Given that $\mathbf{W}_{L:1} = \mathbf{M}_\star$, consider and evaluate the loss function at this minima:

$$\mathcal{L}(\Theta) = \frac{1}{2} \left\| \sum_{\ell} \mathbf{W}_{L:\ell+1} \Delta \mathbf{W}_\ell \mathbf{W}_{\ell-1:1} \right\|_F^2 \quad (46)$$

$$+ \sum_{\ell < m} \mathbf{W}_{L:\ell+1} \Delta \mathbf{W}_\ell \mathbf{W}_{\ell-1:m+1} \Delta \mathbf{W}_m \mathbf{W}_{m-1:1} + \dots + \Delta \mathbf{W}_{L:1} \left\|_F^2. \quad (47)$$

By expanding each of the terms and splitting by the orders of $\Delta \mathbf{W}_\ell$ (perturbation), we get that the second-order term is equivalent to

$$\begin{aligned} \Theta \left(\sum_{\ell=1}^L \|\Delta \mathbf{W}_\ell\|^2 \right) &: \frac{1}{2} \left\| \sum_{\ell} \mathbf{W}_{L:\ell+1} \Delta \mathbf{W}_\ell \mathbf{W}_{\ell-1:1} \right\|_{\mathbb{F}}^2 \\ \Theta \left(\sum_{\ell=1}^L \|\Delta \mathbf{W}_\ell\|^3 \right) &: \text{tr} \left[\left(\sum_{\ell} \mathbf{W}_{L:\ell+1} \Delta \mathbf{W}_\ell \mathbf{W}_{\ell-1:1} \right)^\top \left(\sum_{\ell < m} \mathbf{W}_{L:\ell+1} \Delta \mathbf{W}_\ell \mathbf{W}_{\ell-1:m+1} \Delta \mathbf{W}_m \mathbf{W}_{m-1:1} \right) \right] \\ \Theta \left(\sum_{\ell=1}^L \|\Delta \mathbf{W}_\ell\|^4 \right) &: \frac{1}{2} \left\| \sum_{\ell < m} \mathbf{W}_{L:\ell+1} \Delta \mathbf{W}_\ell \mathbf{W}_{\ell-1:m+1} \Delta \mathbf{W}_m \mathbf{W}_{m-1:1} \right\|_{\mathbb{F}}^2 \\ &+ \text{tr} \left[\sum_l \left(\mathbf{W}_{L:\ell+1} \Delta \mathbf{W}_\ell \mathbf{W}_{\ell-1:1} \right)^\top \left(\sum_{l < m < p} \mathbf{W}_{L:\ell+1} \Delta \mathbf{W}_\ell \mathbf{W}_{\ell-1:m+1} \Delta \mathbf{W}_m \mathbf{W}_{m-1:p+1} \Delta \mathbf{W}_p \mathbf{W}_{p-1:1} \right) \right] \end{aligned}$$

The direction of the steepest change in the loss at the minima correspond to the largest eigenvector direction of the Hessian. Since higher order terms such as $\Theta \left(\sum_{\ell=1}^L \|\Delta \mathbf{W}_\ell\|^3 \right)$ are insignificant compared to the second order terms $\Theta \left(\sum_{\ell=1}^L \|\Delta \mathbf{W}_\ell\|^2 \right)$, finding the direction that maximizes the second order term leads to finding the eigenvector of the Hessian. Then, the eigenvector corresponding to the maximum eigenvalue of $\nabla^2 \mathcal{L}$ is the solution of

$$\Delta_1 := \text{vec}(\Delta \mathbf{W}_L, \dots, \Delta \mathbf{W}_1) = \underset{\|\Delta \mathbf{W}_L\|_{\mathbb{F}}^2 + \dots + \|\Delta \mathbf{W}_1\|_{\mathbb{F}}^2 = 1}{\text{argmax}} f(\Delta \mathbf{W}_L, \dots, \Delta \mathbf{W}_1), \quad (48)$$

where

$$f(\Delta \mathbf{W}_L, \dots, \Delta \mathbf{W}_1) := \frac{1}{2} \|\Delta \mathbf{W}_L \mathbf{W}_{L-1:1} + \dots + \mathbf{W}_{L:3} \Delta \mathbf{W}_2 \mathbf{W}_1 + \mathbf{W}_{L:2} \Delta \mathbf{W}_1\|_{\mathbb{F}}^2. \quad (49)$$

While the solution of Equation (48) gives the maximum eigenvector direction of the Hessian, Δ_1 , the other eigenvectors can be found by solving

$$\Delta_r := \underset{\substack{\|\Delta \mathbf{W}_L\|_{\mathbb{F}}^2 + \dots + \|\Delta \mathbf{W}_1\|_{\mathbb{F}}^2 = 1, \\ \Delta_r \perp \Delta_{r-1}, \dots, \Delta_r \perp \Delta_1}}{\text{argmax}} f(\Delta \mathbf{W}_L, \dots, \Delta \mathbf{W}_1). \quad (50)$$

By expanding $f(\cdot)$, we have that

$$\begin{aligned} f(\Delta \mathbf{W}_L, \dots, \Delta \mathbf{W}_1) &= \|\Delta \mathbf{W}_L \mathbf{W}_{L-1:1}\|_{\mathbb{F}}^2 + \dots + \|\mathbf{W}_{L:3} \Delta \mathbf{W}_2 \mathbf{W}_1\|_{\mathbb{F}}^2 + \|\mathbf{W}_{L:2} \Delta \mathbf{W}_1\|_{\mathbb{F}}^2 \\ &+ \text{tr} \left[(\Delta \mathbf{W}_L \mathbf{W}_{L-1:1})^\top (\mathbf{W}_{L:3} \Delta \mathbf{W}_2 \mathbf{W}_1 + \dots + \mathbf{W}_{L:2} \Delta \mathbf{W}_1) \right] + \dots + \\ &\text{tr} \left[(\mathbf{W}_{L:2} \Delta \mathbf{W}_1)^\top (\mathbf{W}_{L:3} \Delta \mathbf{W}_2 \mathbf{W}_1 + \dots + \mathbf{W}_{L:3} \Delta \mathbf{W}_2 \mathbf{W}_1) \right]. \end{aligned} \quad (51)$$

We can solve Equation (48) by maximizing each of the terms, which can be done in two steps:

- (i) Each Frobenius term in the expansion is maximized when the left singular vector of $\Delta \mathbf{W}_\ell$ aligns with $\mathbf{W}_{L:\ell+1}$ and the right singular vector aligns with $\mathbf{W}_{\ell-1:1}$. This is a result of Von Neumann's trace inequality (Mirsky, 1975). Similarly, each term in the trace is maximized when the singular vector of the perturbations align with the products.
- (ii) Due to the alignment, Equation (48) can be written in just the singular values. Let $\Delta s_{\ell,i}$ denote the i -th singular value of the perturbation matrix $\Delta \mathbf{W}_\ell$. Recall that all of the singular values of \mathbf{M}_* are distinct (i.e., $\sigma_{*,1} > \dots > \sigma_{*,r}$). Hence, it is easy to see that Equation (48) is maximized when $\Delta s_{\ell,i} = 0$ (i.e., all the weight goes to $\Delta s_{\ell,1}$). Thus, each perturbation matrix must be rank-1.

Now since each perturbation is rank-1, we can write each perturbation as

$$\Delta \mathbf{W}_\ell = \Delta s_\ell \Delta \mathbf{u}_\ell \Delta \mathbf{v}_\ell^\top, \quad \forall \ell \in [L], \quad (52)$$

for $\Delta s_\ell > 0$ and orthonormal vectors $\Delta \mathbf{u}_\ell \in \mathbb{R}^d$ and $\Delta \mathbf{v}_\ell \in \mathbb{R}^d$ with $\sum_{\ell=1}^L \Delta s_\ell^2 = 1$. Plugging this in each term, we obtain:

$$\|\mathbf{W}_{L:\ell+1} \Delta_1 \mathbf{W}_\ell \mathbf{W}_{\ell-1:1}\|_2^2 = \Delta_1 s_\ell^2 \cdot \left\| \underbrace{\mathbf{V}_* \sigma_{*,1}^{\frac{L-\ell}{L}} \mathbf{V}_*^\top \Delta \mathbf{u}_\ell}_{=: \mathbf{a}} \underbrace{\Delta \mathbf{v}_\ell^\top \mathbf{V}_* \sigma_{*,1}^{\frac{\ell-1}{L}} \mathbf{V}_*^\top}_{=: \mathbf{b}^\top} \right\|_2^2.$$

Since, allignment maximizes this expression as discussed in first point, we have:

$\mathbf{u}_\ell = \mathbf{v}_\ell = \mathbf{v}_{*,1}$ for all $\ell \in [2, L-1]$, then

$$\mathbf{a} = \sigma_{*,1}^{\frac{L-\ell}{L}} \mathbf{v}_{*,1} \quad \text{and} \quad \mathbf{b}^\top = \sigma_{*,1}^{\frac{\ell-1}{L}} \mathbf{v}_{*,1}^\top \implies \mathbf{a} \mathbf{b}^\top = \sigma_{*,1}^{1-\frac{1}{L}} \cdot \mathbf{v}_{*,1} \mathbf{v}_{*,1}^\top.$$

The very same argument can be made for the trace terms. Hence, in order to maximize $f(\cdot)$, we must have

$$\begin{aligned} \mathbf{v}_L &= \mathbf{v}_{*,1}, \quad \text{and} \quad \mathbf{u}_1 = \mathbf{v}_{*,1}, \\ \mathbf{u}_\ell &= \mathbf{v}_\ell = \mathbf{v}_{*,1}, \quad \forall \ell \in [2, L-1]. \end{aligned}$$

To determine \mathbf{u}_L and \mathbf{v}_1 , we can look at one of the trace terms:

$$\text{tr} \left[(\Delta_1 \mathbf{W}_L \mathbf{W}_{L-1:1})^\top (\mathbf{W}_{L:3} \Delta_1 \mathbf{W}_2 \mathbf{W}_1 + \dots + \mathbf{W}_{L:2} \Delta_1 \mathbf{W}_1) \right] \leq \left(\frac{L-1}{L} \right) \cdot \sigma_{*,1}^{2-\frac{2}{L}}.$$

To reach the upper bound, we require $\mathbf{u}_L = \mathbf{u}_{*,1}$ and $\mathbf{v}_1 = \mathbf{v}_{*,1}$. Finally, as the for each index, the singular values are balanced, we will have $\Delta_1 s_\ell = \frac{1}{\sqrt{L}}$ for all $\ell \in [L]$ to satisfy the constraint.

Finally, we get that the leading eigenvector is

$$\Delta_1 := \text{vec} \left(\frac{1}{\sqrt{L}} \mathbf{u}_1 \mathbf{v}_1^\top, \frac{1}{\sqrt{L}} \mathbf{v}_1 \mathbf{v}_1^\top, \dots, \frac{1}{\sqrt{L}} \mathbf{v}_1 \mathbf{v}_1^\top \right).$$

Notice that we can also verify that $f(\Delta_1) = L \sigma_{*,1}^{2-\frac{2}{L}}$, which is the leading eigenvalue (or sharpness) derived in Lemma 1.

To derive the remaining eigenvectors, we need to find all of the vectors in which $\Delta_i^\top \Delta_j = 0$ for $i \neq j$, where

$$\Delta_i = \text{vec}(\Delta_i \mathbf{W}_L, \dots, \Delta_i \mathbf{W}_1),$$

and $f(\Delta_i) = \lambda_i$, where λ_i is the i -th largest eigenvalue. By repeating the same process as above, we find that the eigenvector-eigenvalue pair as follows:

$$\begin{aligned} \Delta_1 &= \text{vec} \left(\frac{1}{\sqrt{L}} \mathbf{u}_1 \mathbf{v}_1^\top, \frac{1}{\sqrt{L}} \mathbf{v}_1 \mathbf{v}_1^\top, \dots, \frac{1}{\sqrt{L}} \mathbf{v}_1 \mathbf{v}_1^\top \right), \quad \lambda_1 = L \sigma_{*,1}^{2-\frac{2}{L}} \\ \Delta_2 &= \text{vec} \left(\frac{1}{\sqrt{L}} \mathbf{u}_1 \mathbf{v}_2^\top, \frac{1}{\sqrt{L}} \mathbf{v}_1 \mathbf{v}_2^\top, \dots, \frac{1}{\sqrt{L}} \mathbf{v}_1 \mathbf{v}_2^\top \right), \quad \lambda_2 = \left(\sum_{i=0}^{L-1} \sigma_{*,1}^{1-\frac{1}{L}-\frac{1}{L}i} \cdot \sigma_{*,2}^{\frac{1}{L}i} \right) \\ \Delta_3 &= \text{vec} \left(\frac{1}{\sqrt{L}} \mathbf{u}_2 \mathbf{v}_1^\top, \frac{1}{\sqrt{L}} \mathbf{v}_2 \mathbf{v}_1^\top, \dots, \frac{1}{\sqrt{L}} \mathbf{v}_2 \mathbf{v}_1^\top \right), \quad \lambda_3 = \left(\sum_{i=0}^{L-1} \sigma_{*,1}^{1-\frac{1}{L}-\frac{1}{L}i} \cdot \sigma_{*,2}^{\frac{1}{L}i} \right) \\ &\vdots \\ \Delta_d &= \text{vec} \left(\frac{1}{\sqrt{L}} \mathbf{u}_d \mathbf{v}_2^\top, \frac{1}{\sqrt{L}} \mathbf{v}_d \mathbf{v}_2^\top, \dots, \frac{1}{\sqrt{L}} \mathbf{v}_d \mathbf{v}_2^\top \right), \quad \lambda_d = L \sigma_{*,2}^{2-\frac{2}{L}} \\ &\vdots \\ \Delta_{dr+r} &= \text{vec} \left(\frac{1}{\sqrt{L}} \mathbf{u}_d \mathbf{v}_r^\top, \frac{1}{\sqrt{L}} \mathbf{v}_d \mathbf{v}_r^\top, \dots, \frac{1}{\sqrt{L}} \mathbf{v}_d \mathbf{v}_r^\top \right), \end{aligned}$$

2376 which gives a total of $dr + r$ eigenvectors.
2377

2378 Second, equipped with the eigenvectors, let us consider the 1-D function f_{Δ_i} generated by the cross-
2379 section of the loss landscape and each eigenvector Δ_i passing the minima:

$$\begin{aligned}
2380 & f_{\Delta_i}(\mu) = \mathcal{L}(\mathbf{W}_L + \mu\Delta\mathbf{W}_L, \dots, \mathbf{W}_2 + \mu\Delta\mathbf{W}_2, \mathbf{W}_1 + \mu\Delta\mathbf{W}_1), \\
2381 & = \mu^2 \cdot \frac{1}{2} \|\Delta\mathbf{W}_L \mathbf{W}_{L-1:1} + \dots + \mathbf{W}_{L:3} \Delta\mathbf{W}_2 \mathbf{W}_1 + \mathbf{W}_{L:2} \Delta\mathbf{W}_1\|_{\mathbb{F}}^2 \\
2382 & + \mu^3 \cdot \sum_{\ell=1, \ell < m}^L \text{tr} \left[(\mathbf{W}_{L:\ell+1} \Delta\mathbf{W}_\ell \mathbf{W}_{\ell-1:1})^\top (\mathbf{W}_{L:\ell+1} \Delta\mathbf{W}_\ell \mathbf{W}_{\ell-1:m+1} \Delta\mathbf{W}_m \mathbf{W}_{m-1:1}) \right] \\
2383 & + \mu^4 \cdot \frac{1}{2} \left\| \left(\sum_{\ell < m} \mathbf{W}_{L:\ell+1} \Delta\mathbf{W}_\ell \mathbf{W}_{\ell-1:m+1} \Delta\mathbf{W}_m \mathbf{W}_{m-1:1} \right) \right\|_{\mathbb{F}}^2 \\
2384 & + \mu^4 \cdot \sum_{\ell < m < p}^L \text{tr} \left[(\mathbf{W}_{L:\ell+1} \Delta\mathbf{W}_\ell \mathbf{W}_{\ell-1:1})^\top (\mathbf{W}_{L:\ell+1} \Delta\mathbf{W}_\ell \mathbf{W}_{\ell-1:m+1} \Delta\mathbf{W}_m \mathbf{W}_{m-1:p+1} \Delta\mathbf{W}_p \mathbf{W}_{p-1:1}) \right].
\end{aligned}$$

2394 Then, the several order derivatives of $f_{\Delta_i}(\mu)$ at $\mu = 0$ can be obtained from Taylor expansion as
2395

$$\begin{aligned}
2396 & f_{\Delta_i}^{(2)}(0) = \|\Delta_i \mathbf{W}_L \mathbf{W}_{L-1:1} + \dots + \mathbf{W}_{L:3} \Delta_i \mathbf{W}_2 \mathbf{W}_1 + \mathbf{W}_{L:2} \Delta_i \mathbf{W}_1\|_{\mathbb{F}}^2 = \lambda_i^2 \\
2397 & f_{\Delta_i}^{(3)}(0) = 6 \sum_{\ell=1}^L \text{tr} \left[(\mathbf{W}_{L:\ell+1} \Delta_i \mathbf{W}_\ell \mathbf{W}_{\ell-1:1})^\top (\mathbf{W}_{L:\ell+2} \Delta_i \mathbf{W}_{\ell+1} \mathbf{W}_\ell \Delta_i \mathbf{W}_{\ell-1} \mathbf{W}_{\ell-2:1}) \right] \\
2398 & = 6 \left\| \sum_{\ell} \mathbf{W}_{L:\ell+1} \Delta_i \mathbf{W}_\ell \mathbf{W}_{\ell-1:1} \right\|_{\mathbb{F}} \cdot \left\| \left(\sum_{\ell < m} \mathbf{W}_{L:\ell+1} \Delta\mathbf{W}_\ell \mathbf{W}_{\ell-1:m+1} \Delta\mathbf{W}_m \mathbf{W}_{m-1:1} \right) \right\|_{\mathbb{F}} \\
2399 & := 6\lambda_i \cdot \beta_i \\
2400 & f_{\Delta_i}^{(4)}(0) = 12 \|\Delta_i \mathbf{W}_L \Delta_i \mathbf{W}_{L-1} \mathbf{W}_{L-2:1} + \dots + \mathbf{W}_{L:4} \Delta_i \mathbf{W}_3 \mathbf{W}_2 \Delta_i \mathbf{W}_1 + \mathbf{W}_{L:3} \Delta_i \mathbf{W}_2 \Delta_i \mathbf{W}_1\|_{\mathbb{F}}^2 \\
2401 & + 24 \sum_{\ell=1}^L \text{tr} \left[(\mathbf{W}_{L:\ell+1} \Delta_i \mathbf{W}_\ell \mathbf{W}_{\ell-1:1})^\top \left(\sum_{\ell < m < p} \mathbf{W}_{L:\ell+1} \Delta\mathbf{W}_\ell \mathbf{W}_{\ell-1:m+1} \Delta\mathbf{W}_m \mathbf{W}_{m-1:p+1} \Delta\mathbf{W}_p \mathbf{W}_{p-1:1} \right) \right] \\
2402 & := 12\beta_i^2 + 24\lambda_i \cdot \delta_i,
\end{aligned}$$

2412 where we defined

$$\begin{aligned}
2413 & \lambda_i = \left\| \sum_{\ell} \mathbf{W}_{L:\ell+1} \Delta_i \mathbf{W}_\ell \mathbf{W}_{\ell-1:1} \right\|_{\mathbb{F}} \quad (\text{Total } \binom{L}{1} \text{ terms}) \\
2414 & \beta_i = \left\| \left(\sum_{\ell < m} \mathbf{W}_{L:\ell+1} \Delta\mathbf{W}_\ell \mathbf{W}_{\ell-1:m+1} \Delta\mathbf{W}_m \mathbf{W}_{m-1:1} \right) \right\|_{\mathbb{F}} \quad (\text{Total } \binom{L}{2} \text{ terms}) \\
2415 & \delta_i = \left\| \left(\sum_{\ell < m < p} \mathbf{W}_{L:\ell+1} \Delta\mathbf{W}_\ell \mathbf{W}_{\ell-1:m+1} \Delta\mathbf{W}_m \mathbf{W}_{m-1:p+1} \Delta\mathbf{W}_p \mathbf{W}_{p-1:1} \right) \right\|_{\mathbb{F}}, \\
2416 & \quad (\text{Total } \binom{L}{3} \text{ terms})
\end{aligned}$$

2426 and used the fact that $\text{tr}(\mathbf{A}^\top \mathbf{B}) = \|\mathbf{A}\|_{\mathbb{F}} \cdot \|\mathbf{B}\|_{\mathbb{F}}$ under singular vector alignment.
2427

2428 Then, since β_i has $\binom{L}{2}$ terms inside the sum, when the Frobenium term is expanded, it will have
2429 $\frac{\binom{L}{2} \binom{\binom{L}{2} + 1}{2}}$ number of terms. Under alignment and balancedness, $\beta_i^2 = \Delta s_i^2 \sigma_i^{2 - \frac{4}{L}} \times \frac{\binom{L}{2} \binom{\binom{L}{2} + 1}{2}}{2}$

and $\lambda_i \delta_i = \Delta s_\ell^2 \sigma_i^{2-\frac{4}{L}} \times \binom{L}{3} L$. Thus, we have the expression

$$\begin{aligned} 2\beta_i^2 - \lambda_i \delta_i &= \Delta s_\ell^2 \sigma_i^{2-\frac{4}{L}} \left(2 \frac{\binom{L}{2} \left(\binom{L}{2} + 1 \right)}{2} - \binom{L}{3} L \right) \\ &= \Delta s_\ell^2 \sigma_i^{2-\frac{4}{L}} \binom{L}{3} L \times \left(\frac{3 \left(\frac{L(L-1)}{2} + 1 \right)}{L(L-2)} - 1 \right) \\ &= \Delta s_\ell^2 \sigma_i^{2-\frac{4}{L}} \frac{2 \binom{L}{3} L}{L(L-2)} \times ((L-1)^2 + 5) > 0, \end{aligned}$$

for any depth $L > 2$. Finally, the condition of stable oscillation of 1-D function is

$$3[f_{\Delta_i}^{(3)}]^2 - f_{\Delta_i}^{(2)} f_{\Delta_i}^{(4)} = 108\lambda_i^2 \beta_i^2 - (\lambda_i^2)(12\beta_i^2 + 24(2\lambda_i)(\delta_i)) = 48\lambda_i^2(2\beta_i^2 - \lambda_i \delta_i) > 0,$$

which we have proven to be positive for any depth $L > 2$, for all the eigenvector directions corresponding to the non-zero eigenvalues. This completes the proof. \square

Theorem 3 (Subspace Oscillation for Diagonal Linear Networks). *Consider an L -layer diagonal linear network on the loss*

$$\mathcal{L}(\{\mathbf{s}_\ell\}_{\ell=1}^L) := \frac{1}{2} \|\mathbf{s}_1 \odot \dots \odot \mathbf{s}_L - \mathbf{s}_\star\|_2^2, \quad (53)$$

where $\mathbf{s}_\star \in \mathbb{R}^d$ be an r -sparse vector with ordered coordinates such that $s_{\star,1} > \dots > s_{\star,d}$ and

define $S_p := L s_{\star,p}^{2-\frac{2}{L}}$ and $\alpha' := \left(\ln \left(\frac{2\sqrt{2}}{\eta L s_{\star,1}^{\frac{2}{L}}} \cdot \frac{s_{\star,1}^{\frac{4}{L}}}{L^2 \cdot 2^{\frac{2L-3}{L}}} \right) \right)^{\frac{1}{4}}$. For any $p < r-1$ and $\alpha < \alpha'$,

suppose we run GD on Equation (5) with learning rate $\eta = \frac{2}{K}$, where $S_p \geq K > S_{p+1}$ with initialization $\mathbf{s}_\ell = \alpha \mathbf{1}_d$ for all $\ell \in [L-1]$ and $\mathbf{s}_L = \mathbf{0}_d$. Then, under strict balancing, the top- p coordinates of \mathbf{s}_ℓ oscillate within a 2-period fixed orbit around the minima in the form

$$s_{\ell,i}(t) = \rho_{i,j}(t), \quad \forall i < p, \forall \ell \in [L],$$

where $\rho_{i,j}(t) \in \{\rho_{i,1}, \rho_{i,2}\}$, $\rho_{i,1} \in (0, s_{\star,i}^{1/L})$ and $\rho_{i,2} \in (s_{\star,i}^{1/L}, (2s_{\star,i})^{1/L})$ are two real roots of the polynomial $h(\rho) = 0$:

$$h(\rho) = \rho^L \cdot \frac{1 + (1 + \eta L(s_{\star,i} - \rho^L) \cdot \rho^{L-2})^{2L-1}}{1 + (1 + \eta L(s_{\star,i} - \rho^L) \cdot \rho^{L-2})^{L-1}} - s_{\star,i}.$$

Proof. This proof essentially mimics that of the DLN proof from Theorem 1, in that we will

- (i) Compute the eigenvalues and eigenvectors of the flattened training loss Hessian of the diagonal linear network at convergence under the balancing assumption.
- (ii) Show that in 1D cross-section of the eigenvector, the stable condition oscillation $3[f_{\Delta_i}^{(3)}]^2 - f_{\Delta_i}^{(2)} f_{\Delta_i}^{(4)} > 0$ is satisfied, where f_{Δ_i} denotes the 1D cross-section function at the i -th eigenvector direction.

With a slight abuse in notation, let $\mathbf{s} := \{\mathbf{s}_\ell\}_{\ell=1}^L$. Let us first derive the Hessian at convergence by considering each block of the flattened Hessian matrix denoted by $\mathbf{H}_{m,\ell}$:

$$\mathbf{H}_{\ell,\ell} = \begin{bmatrix} \prod_{k \neq \ell} s_{k,1} & & \\ & \ddots & \\ & & \prod_{k \neq \ell} s_{k,d} \end{bmatrix} \quad \text{if } m = \ell \quad (54)$$

$$\mathbf{H}_{m,\ell} = \begin{bmatrix} \gamma_1 & & \\ & \ddots & \\ & & \gamma_d \end{bmatrix} \quad \text{if } m \neq \ell, \quad (55)$$

where

$$\gamma_i := \left(\prod_{k \neq \ell} s_{k,i} \right) \left(\prod_{k \neq m} s_{k,i} \right) + \left(\prod_k s_{k,i} - s_{\star,i} \right) \left(\prod_{k \neq \ell, k \neq m} s_{k,i} \right).$$

Then, under Lemma 2, at convergence (i.e. the gradient is zero), we have

$$\left(\prod_k s_{k,i} - s_{\star,i} \right) = 0 \implies s_{k,i} = s_{\star,i}^{\frac{1}{L}},$$

which means that at convergence, the Hessian is given by

$$\mathbf{H} = \begin{bmatrix} \mathbf{A} & \dots & \mathbf{A} & \mathbf{A} \\ \vdots & \ddots & \vdots & \vdots \\ \mathbf{A} & \dots & \mathbf{A} & \mathbf{A} \\ \mathbf{A} & \dots & \mathbf{A} & \mathbf{B} \end{bmatrix} \in \mathbb{R}^{dL \times dL},$$

where

$$\mathbf{A} := \begin{bmatrix} s_{\star,1}^{2-\frac{2}{L}} & & & \\ & \ddots & & \\ & & s_{\star,r}^{2-\frac{2}{L}} & \\ & & & \mathbf{0}_{d-r} \end{bmatrix} \in \mathbb{R}^{d \times d}, \quad \mathbf{B} := \begin{bmatrix} s_{\star,1}^{2-\frac{2}{L}} & & & \\ & \ddots & & \\ & & s_{\star,r}^{2-\frac{2}{L}} & \\ & & & \alpha^{2(L-1)} \cdot \mathbf{I}_{d-r} \end{bmatrix} \in \mathbb{R}^{d \times d}.$$

To compute the eigenvalues of \mathbf{H} , we can block diagonalize \mathbf{H} into the form $\mathbf{C} = \mathbf{P}\mathbf{H}\mathbf{P}^\top$, where \mathbf{P} is a permutation matrix and

$$\mathbf{C} = \begin{bmatrix} \mathbf{C}_1 & & \\ & \ddots & \\ & & \mathbf{C}_d \end{bmatrix} \in \mathbb{R}^{dL \times dL},$$

where each (i, j) -th entry of $\mathbf{C}_k \in \mathbb{R}^{L \times L}$ is the k -th diagonal element of $\mathbf{H}_{i,j}$. Then, since \mathbf{C} is a block diagonal matrix, its eigenvalues are the union of each of the eigenvalues of its blocks. Then, notice that

$$\mathbf{C}_j = s_{\star,k}^{2-\frac{2}{L}} \cdot \mathbf{1}_L \mathbf{1}_L^\top, \quad \forall j \in [r] \quad \mathbf{C}_k = \begin{bmatrix} 0 & \dots & 0 & \alpha^{2(L-1)} \\ \vdots & \ddots & \vdots & \vdots \\ 0 & \dots & 0 & \alpha^{2(L-1)} \\ \alpha^{2(L-1)} & \dots & \alpha^{2(L-1)} & \alpha^{2(L-1)} \end{bmatrix}, \quad \forall k \in [r+1, d].$$

Hence, the eigenvalues of \mathbf{C} (and the eigenvalues of \mathbf{H}) are given by

$$\lambda_{\mathbf{H}} = \left\{ L s_{\star,i}^{2-\frac{2}{L}}, \underbrace{0}_{\text{multiplicity } L-1} \right\}_{i=1}^r \cup \left\{ \underbrace{\frac{-\alpha^{2(L-1)} \pm \sqrt{(4L-3) \cdot \alpha^{4(L-1)^2}}}{-2}}_{\text{multiplicity } d-r}, \underbrace{0}_{\text{multiplicity } (d-r)(L-2)} \right\},$$

which can be computed using co-factor expansion. For the eigenvectors, notice that we can write

$$\mathbf{C}\mathbf{v} = \mathbf{P}\mathbf{H}\mathbf{P}^\top \mathbf{v} = \lambda \mathbf{v} \implies \mathbf{H}\mathbf{P}^\top \mathbf{v} = \lambda \mathbf{P}^\top \mathbf{v}.$$

2538 Hence, we can find the eigenvectors of the block diagonal matrix \mathbf{C} , and left multiply them by \mathbf{P}^\top
 2539 to obtain the eigenvectors of the Hessian \mathbf{H} . This yields the eigenvector and eigenvalue pairs

$$2541 \Delta_1 = \mathbf{P}^\top \text{vec} \left(\frac{1}{\sqrt{L}} \mathbf{1}_L, \mathbf{0}, \dots, \mathbf{0} \right), \quad \lambda_1 = L s_{\star,1}^{2-\frac{2}{L}} \quad (56)$$

$$2543 \Delta_2 = \mathbf{P}^\top \text{vec} \left(\mathbf{0}, \frac{1}{\sqrt{L}} \mathbf{1}_L, \dots, \mathbf{0} \right), \quad \lambda_2 = L s_{\star,2}^{2-\frac{2}{L}} \quad (57)$$

$$2546 \vdots \quad \vdots \quad (58)$$

$$2547 \Delta_r = \mathbf{P}^\top \text{vec} \left(\mathbf{0}, \dots, \frac{1}{\sqrt{L}} \mathbf{1}_L, \dots, \mathbf{0} \right), \quad \lambda_r = L s_{\star,r}^{2-\frac{2}{L}} \quad (59)$$

$$2549 \vdots \quad \vdots \quad (60)$$

$$2551 \Delta_{r+j} = \mathbf{P}^\top \text{vec} (\mathbf{0}, \dots, \mathbf{e}_{r+j}, \dots, \mathbf{0}), \quad \lambda_{r+j} = \frac{-\alpha^{2(L-1)} \pm \sqrt{(4L-3) \cdot \alpha^{4(L-1)^2}}}{-2}, \quad (61)$$

2552 where \mathbf{e}_i is an i -th elementary basis vector.

2553 Then, in each 1-D eigenvector direction, we can analyze the loss and verify if it satisfies the stability
 2554 condition. Notice that we can consider the scalar loss

$$2555 \mathcal{L}_i(\mathbf{s}) = \frac{1}{2} (s_{1,i} \odot \dots \odot s_{L,i} - s_{\star,i})^2 = \frac{1}{2} (s_i^L - s_{\star,i})^2. \quad (\text{By Lemma 2})$$

2556 Using Corollary 5 by Chen & Bruna (2023) or restated Lemma 11 on the 1D scalar function, this 1D
 2557 loss is amenable to stable oscillation when learning rate $\eta > \frac{2}{\lambda_i}$. Finally, to prove the uniqueness
 2558 and existence of two period orbit fixed point for $\eta > \frac{2}{\lambda_i}$, we show that the polynomial obtained by
 2559 solving two step fixed point has a real root. This is the same loss we analyzed in Theorem 2, where
 2560 we showed that the oscillations are real roots of the polynomial

$$2561 \sigma_{\star,1} = \rho^L \frac{1+z^{2L-1}}{1+z^{L-1}}, \quad \text{where } z := (1 + \eta L (\sigma_{\star,1} - \rho^L) \cdot \rho^{L-2}).$$

2562 and $\rho_1 \in (0, \sigma_{\star,1}^{1/L})$ and $\rho_2 \in (\sigma_{\star,1}^{1/L}, (2\sigma_{\star,1})^{1/L})$ are the two real roots of the polynomial which
 2563 exists and are unique. Hence, whenever the learning rate η lies between $[2/\lambda_p, 2/\lambda_{p+1}]$, we will
 2564 have oscillations in all of the p eigenvector directions. This completes the proof.

2565 □

2566

2567

2568

2569

2570

2571

2572

2573

2574

2575

2576

2577

2578

2579

2580

2581

2582

2583

2584

2585

2586

2587

2588

2589

2590

2591

2592 C.2 DEFERRED PROOFS FOR SINGULAR VECTOR INVARIANCE
 2593

2594 **Proposition 2.** Let $\mathbf{M}_\star = \mathbf{U}_\star \boldsymbol{\Sigma}_\star \mathbf{V}_\star^\top$ denote the SVD of the target matrix. The initialization in
 2595 Equation (3) is a member of the singular vector stationary set in Proposition 1, where $\mathbf{Q}_L = \dots =$
 2596 $\mathbf{Q}_2 = \mathbf{V}_\star$.

2597 *Proof.* Recall that the initialization is given by

$$2598 \mathbf{W}_L(0) = \mathbf{0} \quad \text{and} \quad \mathbf{W}_\ell(0) = \alpha \mathbf{I}_d \quad \forall \ell \in [L-1].$$

2600 We will show that under this initialization, each weight matrix admits the following decomposition
 2601 for all $t \geq 1$:

$$2602 \mathbf{W}_L(t) = \mathbf{U}_\star \begin{bmatrix} \tilde{\boldsymbol{\Sigma}}_L(t) & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \mathbf{V}_\star^\top, \quad \mathbf{W}_\ell(t) = \mathbf{V}_\star \begin{bmatrix} \tilde{\boldsymbol{\Sigma}}(t) & \mathbf{0} \\ \mathbf{0} & \alpha \mathbf{I}_{d-r} \end{bmatrix} \mathbf{V}_\star^\top, \quad \forall \ell \in [L-1], \quad (62)$$

2605 where

$$2606 \tilde{\boldsymbol{\Sigma}}_L(t) = \tilde{\boldsymbol{\Sigma}}_L(t-1) - \eta \cdot \left(\tilde{\boldsymbol{\Sigma}}_L(t-1) \cdot \tilde{\boldsymbol{\Sigma}}^{L-1}(t-1) - \boldsymbol{\Sigma}_{\star,r} \right) \cdot \tilde{\boldsymbol{\Sigma}}^{L-1}(t-1)$$

$$2607 \tilde{\boldsymbol{\Sigma}}(t) = \tilde{\boldsymbol{\Sigma}}(t-1) \cdot \left(\mathbf{I}_r - \eta \cdot \tilde{\boldsymbol{\Sigma}}_L(t-1) \right) \cdot \left(\tilde{\boldsymbol{\Sigma}}_L(t-1) \cdot \tilde{\boldsymbol{\Sigma}}^{L-1}(t-1) - \boldsymbol{\Sigma}_{\star,r} \right) \cdot \tilde{\boldsymbol{\Sigma}}^{L-3}(t-1),$$

2610 where $\tilde{\boldsymbol{\Sigma}}_L(t), \tilde{\boldsymbol{\Sigma}}(t) \in \mathbb{R}^{r \times r}$ is a diagonal matrix with $\tilde{\boldsymbol{\Sigma}}_L(1) = \eta \alpha^{L-1} \cdot \boldsymbol{\Sigma}_{r,\star}$ and $\tilde{\boldsymbol{\Sigma}}(1) = \alpha \mathbf{I}_r$.

2611 This will prove that the singular vectors are stationary with $\boldsymbol{\Sigma}_L = \dots = \boldsymbol{\Sigma}_2 = \mathbf{V}_\star$. We proceed
 2612 with mathematical induction.

2614 **Base Case.** For the base case, we will show that the decomposition holds for each weight matrix
 2615 at $t = 1$. The gradient of $f(\boldsymbol{\Theta})$ with respect to \mathbf{W}_ℓ is

$$2616 \nabla_{\mathbf{W}_\ell} f(\boldsymbol{\Theta}) = \mathbf{W}_{L:\ell+1}^\top \cdot (\mathbf{W}_{L:1} - \mathbf{M}_\star) \cdot \mathbf{W}_{\ell-1:1}^\top.$$

2617 For $\mathbf{W}_L(1)$, we have

$$2618 \mathbf{W}_L(1) = \mathbf{W}_L(0) - \eta \cdot \nabla_{\mathbf{W}_L} f(\boldsymbol{\Theta}(0))$$

$$2619 = \mathbf{W}_L(0) - \eta \cdot (\mathbf{W}_{L:1}(0) - \mathbf{M}_\star) \cdot \mathbf{W}_{L-1:1}^\top(0)$$

$$2620 = \eta \alpha^{L-1} \boldsymbol{\Sigma}_\star$$

$$2621 = \mathbf{U}_\star \cdot (\eta \alpha^{L-1} \cdot \boldsymbol{\Sigma}_\star) \cdot \mathbf{V}_\star^\top$$

$$2622 = \mathbf{U}_\star \begin{bmatrix} \tilde{\boldsymbol{\Sigma}}_L(1) & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \mathbf{V}_\star^\top.$$

2623 Then, for each $\mathbf{W}_\ell(1)$ in $\ell \in [L-1]$, we have

$$2624 \mathbf{W}_\ell(1) = \mathbf{W}_\ell(0) - \eta \cdot \nabla_{\mathbf{W}_\ell} f(\boldsymbol{\Theta}(0))$$

$$2625 = \alpha \mathbf{I}_d,$$

2626 where the last equality follows from the fact that $\mathbf{W}_L(0) = \mathbf{0}$. Finally, we have

$$2627 \mathbf{W}_\ell(1) = \alpha \mathbf{V}_\star \mathbf{V}_\star^\top = \mathbf{V}_\star \begin{bmatrix} \tilde{\boldsymbol{\Sigma}}(1) & \mathbf{0} \\ \mathbf{0} & \alpha \mathbf{I}_{d-r} \end{bmatrix} \mathbf{V}_\star^\top, \quad \forall \ell \in [L-1].$$

2628 **Inductive Step.** By the inductive hypothesis, suppose that the decomposition holds. Then, notice
 2629 that we can simplify the end-to-end weight matrix to

$$2630 \mathbf{W}_{L:1}(t) = \mathbf{U}_\star \begin{bmatrix} \tilde{\boldsymbol{\Sigma}}_L(t) \cdot \tilde{\boldsymbol{\Sigma}}^{L-1}(t) & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \mathbf{V}_\star^\top,$$

2631 for which we can simplify the gradients to

$$2632 \nabla_{\mathbf{W}_L} f(\boldsymbol{\Theta}(t)) = \left(\mathbf{U}_\star \begin{bmatrix} \tilde{\boldsymbol{\Sigma}}_L(t) \cdot \tilde{\boldsymbol{\Sigma}}^{L-1}(t) - \boldsymbol{\Sigma}_{\star,r} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \mathbf{V}_\star^\top \right) \cdot \mathbf{V}_\star \begin{bmatrix} \tilde{\boldsymbol{\Sigma}}^{L-1}(t) & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \mathbf{V}_\star^\top$$

$$2633 = \mathbf{U}_\star \begin{bmatrix} \left(\tilde{\boldsymbol{\Sigma}}_L(t) \cdot \tilde{\boldsymbol{\Sigma}}^{L-1}(t) - \boldsymbol{\Sigma}_{\star,r} \right) \cdot \tilde{\boldsymbol{\Sigma}}^{L-1}(t) & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \mathbf{V}_\star^\top,$$

for the last layer matrix, and similarly,

$$\nabla_{\mathbf{W}_\ell} f(\Theta(t)) = \mathbf{V}_\star \begin{bmatrix} \tilde{\Sigma}_L(t) \cdot (\tilde{\Sigma}_L(t) \cdot \tilde{\Sigma}^{L-1}(t) - \Sigma_{\star,r}) \cdot \tilde{\Sigma}^{L-2}(t) & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \mathbf{V}_\star^\top, \quad \ell \in [L-1],$$

for all other layer matrices. Thus, for the next GD iteration, we have

$$\begin{aligned} \mathbf{W}_L(t+1) &= \mathbf{W}_L(t) - \eta \cdot \nabla_{\mathbf{W}_L} f(\Theta(t)) \\ &= \mathbf{U}_\star \begin{bmatrix} \tilde{\Sigma}_L(t) - \eta \cdot (\tilde{\Sigma}_L(t) \cdot \tilde{\Sigma}^{L-1}(t) - \Sigma_{\star,r}) \cdot \tilde{\Sigma}^{L-1}(t) & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \mathbf{V}_\star^\top \\ &= \mathbf{U}_\star \begin{bmatrix} \tilde{\Sigma}_L(t+1) & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \mathbf{V}_\star^\top. \end{aligned}$$

Similarly, we have

$$\begin{aligned} \mathbf{W}_\ell(t+1) &= \mathbf{W}_\ell(t) - \eta \cdot \nabla_{\mathbf{W}_\ell} f(\Theta(t)) \\ &= \mathbf{V}_\star \begin{bmatrix} \tilde{\Sigma}(t) - \eta \cdot \tilde{\Sigma}_L(t) \cdot (\tilde{\Sigma}_L(t) \cdot \tilde{\Sigma}^{L-1}(t) - \Sigma_{\star,r}) \cdot \tilde{\Sigma}^{L-2}(t) & \mathbf{0} \\ \mathbf{0} & \alpha \mathbf{I}_{d-r} \end{bmatrix} \mathbf{V}_\star^\top \\ &= \mathbf{V}_\star \begin{bmatrix} \tilde{\Sigma}(t) \cdot (\mathbf{I}_r - \eta \cdot \tilde{\Sigma}_L(t) \cdot (\tilde{\Sigma}_L(t) \cdot \tilde{\Sigma}^{L-1}(t) - \Sigma_{\star,r}) \cdot \tilde{\Sigma}^{L-3}(t)) & \mathbf{0} \\ \mathbf{0} & \alpha \mathbf{I}_{d-r} \end{bmatrix} \mathbf{V}_\star^\top \\ &= \mathbf{V}_\star \begin{bmatrix} \tilde{\Sigma}(t+1) & \mathbf{0} \\ \mathbf{0} & \alpha \mathbf{I}_{d-r} \end{bmatrix} \mathbf{V}_\star^\top, \end{aligned}$$

for all $\ell \in [L-1]$. This completes the proof. \square

Proposition 3. Let $\mathbf{M}_\star = \mathbf{V}_\star \Sigma_\star \mathbf{V}_\star^\top \in \mathbb{R}^{d \times d}$ denote the SVD of the target matrix. The balanced initialization in Equation (3) is a member of the singular vector stationary set in Proposition 1, where $\mathbf{U}_L = \mathbf{Q}_L = \dots = \mathbf{Q}_2 = \mathbf{V}_1 = \mathbf{V}_\star$.

Proof. Using mathematical induction, we will show that with the balanced initialization in Equation (3), each weight matrix admits a decomposition of the form

$$\mathbf{W}_\ell(t) = \mathbf{V}_\star \Sigma_\ell(t) \mathbf{V}_\star^\top, \quad (63)$$

which implies that the singular vectors are stationary for all t such that $\mathbf{U}_L = \mathbf{Q}_L = \dots = \mathbf{Q}_2 = \mathbf{V}_1 = \mathbf{V}_\star$.

Base Case. Consider the weights at iteration $t = 0$. By the initialization scheme, we can write each weight matrix as

$$\mathbf{W}_\ell(0) = \alpha \mathbf{I}_d \implies \mathbf{W}_\ell(0) = \alpha \mathbf{V}_\star \mathbf{V}_\star^\top,$$

which implies that $\mathbf{W}_\ell(0) = \mathbf{V}_\star \Sigma_\ell(0) \mathbf{V}_\star^\top$ with $\Sigma_\ell(0) = \alpha \mathbf{I}_d$.

Inductive Step. By the inductive hypothesis, assume that the decomposition holds for all $t \geq 0$. We will show that it holds for all iterations $t + 1$. Recall that the gradient of $f(\Theta)$ with respect to \mathbf{W}_ℓ is

$$\nabla_{\mathbf{W}_\ell} f(\Theta) = \mathbf{W}_{L:\ell+1}^\top \cdot (\mathbf{W}_{L:1} - \mathbf{M}_\star) \cdot \mathbf{W}_{\ell-1:1}^\top.$$

Then, for $\mathbf{W}_\ell(t+1)$, we have

$$\begin{aligned} \mathbf{W}_\ell(t+1) &= \mathbf{W}_\ell(t) - \eta \cdot \nabla_{\mathbf{W}_\ell} f(\Theta(t)) \\ &= \mathbf{V}_\star \Sigma_\ell(t) \mathbf{V}_\star^\top - \eta \mathbf{W}_{L:\ell+1}^\top(t) \cdot (\mathbf{W}_{L:1}(t) - \mathbf{M}_\star) \cdot \mathbf{W}_{\ell-1:1}^\top(t) \\ &= \mathbf{V}_\star \Sigma_\ell(t) \mathbf{V}_\star^\top - \eta \mathbf{V}_\star \cdot (\Sigma_\ell^{L-\ell}(t) \cdot (\Sigma_\ell^L(t) - \Sigma_\star) \cdot \Sigma_\ell^{\ell-1}(t)) \cdot \mathbf{V}_\star^\top \\ &= \mathbf{V}_\star \cdot (\Sigma_\ell(t) - \eta \cdot \Sigma_\ell^{L-\ell}(t) \cdot (\Sigma_\ell^L(t) - \Sigma_\star) \cdot \Sigma_\ell^{\ell-1}(t)) \cdot \mathbf{V}_\star^\top \\ &= \mathbf{V}_\star \Sigma(t) \mathbf{V}_\star^\top, \end{aligned}$$

where $\Sigma(t) = \Sigma_\ell(t) - \eta \cdot \Sigma_\ell^{L-\ell}(t) \cdot (\Sigma_\ell^L(t) - \Sigma_\star) \cdot \Sigma_\ell^{\ell-1}(t)$. This completes the proof. \square

C.3 AUXILIARY RESULTS

Lemma 8. Let $\{\mathbf{R}_\ell\}_{\ell=1}^L \in \mathbb{R}^{n \times n}$ be orthogonal matrices and $\mathbf{H}_{i,j} \in \mathbb{R}^{n^2 \times n^2}$ be diagonal matrices. Consider the two following block matrices:

$$\mathbf{H} = \begin{bmatrix} \mathbf{H}_{1,1} & \mathbf{H}_{1,2} & \cdots & \mathbf{H}_{L,1} \\ \mathbf{H}_{2,1} & \mathbf{H}_{2,2} & \cdots & \mathbf{H}_{L,2} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{H}_{1,L} & \mathbf{H}_{2,L} & \cdots & \mathbf{H}_{L,L} \end{bmatrix}$$

$$\tilde{\mathbf{H}} = \begin{bmatrix} \mathbf{R}_L \mathbf{H}_{1,1} \mathbf{R}_L^\top & \mathbf{R}_L \mathbf{H}_{1,2} \mathbf{R}_{L-1}^\top & \cdots & \mathbf{R}_L \mathbf{H}_{1,L} \mathbf{R}_1^\top \\ \mathbf{R}_{L-1} \mathbf{H}_{2,1} \mathbf{R}_L^\top & \mathbf{R}_{L-1} \mathbf{H}_{2,2} \mathbf{R}_{L-1}^\top & \cdots & \mathbf{R}_{L-1} \mathbf{H}_{2,L} \mathbf{R}_1^\top \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{R}_1 \mathbf{H}_{L,1} \mathbf{R}_L^\top & \mathbf{R}_1 \mathbf{H}_{L,2} \mathbf{R}_{L-1}^\top & \cdots & \mathbf{R}_1 \mathbf{H}_{L,L} \mathbf{R}_1^\top \end{bmatrix}.$$

Then, the two matrices \mathbf{H} and $\tilde{\mathbf{H}}$ are similar, in the sense that they have the same eigenvalues.

Proof. It suffices to show that \mathbf{H} and $\tilde{\mathbf{H}}$ have the same characteristic polynomials. Let us define

$$\tilde{\mathbf{H}} := \begin{bmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{C} & \mathbf{D} \end{bmatrix},$$

where

$$\mathbf{A} := \mathbf{R}_L \mathbf{H}_{1,1} \mathbf{R}_L^\top \quad \mathbf{B} := [\mathbf{R}_L \mathbf{H}_{1,2} \mathbf{R}_{L-1}^\top \quad \cdots \quad \mathbf{R}_L \mathbf{H}_{1,L} \mathbf{R}_1^\top] \quad (64)$$

$$\mathbf{C} := \begin{bmatrix} \mathbf{R}_{L-1} \mathbf{H}_{2,1} \mathbf{R}_L^\top \\ \vdots \\ \mathbf{R}_1 \mathbf{H}_{L,1} \mathbf{R}_L^\top \end{bmatrix} \quad \mathbf{D} := \begin{bmatrix} \mathbf{R}_{L-1} \mathbf{H}_{2,2} \mathbf{R}_{L-1}^\top & \cdots & \mathbf{R}_{L-1} \mathbf{H}_{2,L} \mathbf{R}_1^\top \\ \vdots & \ddots & \vdots \\ \mathbf{R}_1 \mathbf{H}_{L,2} \mathbf{R}_{L-1}^\top & \cdots & \mathbf{R}_1 \mathbf{H}_{L,L} \mathbf{R}_1^\top \end{bmatrix}. \quad (65)$$

Then, we have

$$\begin{aligned} \det(\tilde{\mathbf{H}} - \lambda \mathbf{I}) &= \det \left(\begin{bmatrix} \mathbf{A} - \lambda \mathbf{I} & \mathbf{B} \\ \mathbf{C} & \mathbf{D} - \lambda \mathbf{I} \end{bmatrix} \right) \\ &= \det(\mathbf{A} - \lambda \mathbf{I}) \cdot \det((\mathbf{D} - \lambda \mathbf{I}) - \mathbf{C}(\mathbf{A} - \lambda \mathbf{I})^{-1} \mathbf{B}), \end{aligned}$$

where the second equality is by the Schur complement. Notice that

$$\begin{aligned} (\mathbf{A} - \lambda \mathbf{I})^{-1} &= (\mathbf{R}_L \mathbf{H}_{1,1} \mathbf{R}_L^\top - \lambda \mathbf{I})^{-1} = (\mathbf{R}_L \mathbf{H}_{1,1} \mathbf{R}_L^\top - \lambda \mathbf{R}_L \mathbf{R}_L^\top)^{-1} \\ &= \mathbf{R}_L \cdot (\mathbf{H}_{1,1} - \lambda \mathbf{I})^{-1} \cdot \mathbf{R}_L^\top. \end{aligned}$$

Then, we also see that,

$$\mathbf{C}(\mathbf{A} - \lambda \mathbf{I})^{-1} \mathbf{B} = \underbrace{\begin{bmatrix} \mathbf{R}_{L-1} & & \\ & \ddots & \\ & & \mathbf{R}_1 \end{bmatrix}}_{=:\hat{\mathbf{V}}} \cdot \mathbf{E} \cdot \underbrace{\begin{bmatrix} \mathbf{R}_{L-1}^\top & & \\ & \ddots & \\ & & \mathbf{R}_1^\top \end{bmatrix}}_{=:\hat{\mathbf{V}}^\top}.$$

where

$$\mathbf{E} := \begin{bmatrix} \mathbf{H}_{2,1} \cdot (\mathbf{H}_{1,1} - \lambda \mathbf{I})^{-1} \cdot \mathbf{H}_{1,2} & \cdots & \mathbf{H}_{2,1} \cdot (\mathbf{H}_{1,1} - \lambda \mathbf{I})^{-1} \cdot \mathbf{H}_{1,L} \\ \vdots & \ddots & \vdots \\ \mathbf{H}_{L,1} \cdot (\mathbf{H}_{1,1} - \lambda \mathbf{I})^{-1} \cdot \mathbf{H}_{1,2} & \cdots & \mathbf{H}_{L,1} \cdot (\mathbf{H}_{1,1} - \lambda \mathbf{I})^{-1} \cdot \mathbf{H}_{1,L} \end{bmatrix}.$$

Similarly, we can write \mathbf{D} as

$$\mathbf{D} = \hat{\mathbf{V}} \underbrace{\begin{bmatrix} \mathbf{H}_{2,2} & \cdots & \mathbf{H}_{2,L} \\ \vdots & \ddots & \vdots \\ \mathbf{H}_{L,2} & \cdots & \mathbf{H}_{L,L} \end{bmatrix}}_{=\mathbf{F}} \hat{\mathbf{V}}^\top.$$

2754 Then, we have
2755

$$\begin{aligned} 2756 \det(\tilde{\mathbf{H}} - \lambda \mathbf{I}) &= \det(\mathbf{R}_L \cdot (\mathbf{H}_{1,1} - \lambda \mathbf{I}) \cdot \mathbf{R}_L^\top) \cdot \det(\hat{\mathbf{V}} \cdot (\mathbf{E} - \mathbf{F}) \cdot \hat{\mathbf{V}}^\top) \\ 2757 &= \det(\mathbf{H}_{1,1} - \lambda \mathbf{I}) \cdot \det(\mathbf{E} - \mathbf{F}), \end{aligned}$$

2759 which is not a function of $\mathbf{U}, \mathbf{V}, \{\mathbf{R}_\ell\}_{\ell=1}^L$. By doing the same for \mathbf{H} , we can show that both $\tilde{\mathbf{H}}$ and
2760 \mathbf{H} have the same characteristic polynomials, and hence the same eigenvalues. This completes the
2761 proof.
2762

□

2764 **Lemma 9.** Let $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{d \times d}$ be two orthogonal matrices. Then, the Kronecker product of \mathbf{A} and
2765 \mathbf{B} is also an orthogonal matrix:

$$2766 (\mathbf{A} \otimes \mathbf{B})^\top (\mathbf{A} \otimes \mathbf{B}) = (\mathbf{A} \otimes \mathbf{B})(\mathbf{A} \otimes \mathbf{B})^\top = \mathbf{I}_{d^2}.$$

2768 *Proof.* We prove this directly by using properties of Kronecker products:
2769

$$\begin{aligned} 2770 (\mathbf{A} \otimes \mathbf{B})^\top (\mathbf{A} \otimes \mathbf{B}) &= \mathbf{A}^\top \mathbf{A} \otimes \mathbf{B}^\top \mathbf{B} \\ 2771 &= \mathbf{I}_d \otimes \mathbf{I}_d = \mathbf{I}_{d^2}. \end{aligned}$$

2773 Similarly, we have

$$\begin{aligned} 2774 (\mathbf{A} \otimes \mathbf{B})(\mathbf{A} \otimes \mathbf{B})^\top &= \mathbf{A} \mathbf{A}^\top \otimes \mathbf{B} \mathbf{B}^\top \\ 2775 &= \mathbf{I}_d \otimes \mathbf{I}_d = \mathbf{I}_{d^2}. \end{aligned}$$

2777 This completes the proof. □

2778 **Lemma 10.** Let $\{a(t)\}_{t=1}^N$ be a sequence such that $a(t) \geq 0$ for all t . If there exists a constant
2779 $c \in (0, 1)$ such that $a(t+1) < c \cdot a(t)$ for all t , then $\lim_{t \rightarrow \infty} a(t) = 0$.
2780

2781 *Proof.* We prove this by direct reasoning. From the assumption $a(t+1) < c \cdot a(t)$ for some
2782 $c \in (0, 1)$, we can iteratively expand this inequality:
2783

$$2784 a(t+1) < c \cdot a(t), \quad a(t+2) < c \cdot a(t+1) < c^2 \cdot a(t),$$

2785 and, more generally, by induction:
2786

$$2787 a(t+k) < c^k \cdot a(t), \quad \text{for all } k \geq 0.$$

2788 Since $c \in (0, 1)$, the sequence $\{c^k\}_{k=0}^\infty$ converges to 0 as $k \rightarrow \infty$. Hence:
2789

$$2790 0 \leq \lim_{k \rightarrow \infty} a(t+k) \leq \lim_{k \rightarrow \infty} c^k \cdot a(t) = 0.$$

2792 Therefore, by the squeeze theorem, the sequence $\{a(t)\}$ converges to 0 as $t \rightarrow \infty$. □
2793

2794 **Lemma 11** (Chen & Bruna (2023)). Consider any 1-D differentiable function $f(x)$ around a local
2795 minima \bar{x} , satisfying (i) $f^{(3)}(\bar{x}) \neq 0$, and (ii) $3[f^{(3)}]^2 - f'' f^{(4)} > 0$ at \bar{x} . Then, there exists ϵ with
2796 sufficiently small $|\epsilon|$ and $\epsilon \cdot f^{(3)} > 0$ such that: for any point x_0 between \bar{x} and $\bar{x} - \epsilon$, there exists a
2797 learning rate η such that $F_\eta^2(x_0) = x_0$, and

$$2798 \frac{2}{f''(\bar{x})} < \eta < \frac{2}{f''(\bar{x}) - \epsilon \cdot f^{(3)}(\bar{x})}.$$

2800
2801
2802
2803
2804
2805
2806
2807