
SHARPNESS OF MINIMA IN DEEP MATRIX FACTORIZATION: EXACT EXPRESSIONS

Anonymous authors

Paper under double-blind review

ABSTRACT

Understanding the geometry of the loss landscape near a minimum is key to explaining the implicit bias of gradient-based methods in non-convex optimization problems such as deep neural network training and deep matrix factorization. A central quantity to characterize this geometry is the maximum eigenvalue of the Hessian of the loss, which measures the sharpness of the landscape. Currently, its precise role has been obfuscated because no exact expressions for this sharpness measure were known in general settings. In this paper, we present the first exact expression for the maximum eigenvalue of the Hessian of the squared-error loss at *any* minimizer in general overparameterized deep matrix factorization (i.e., deep linear neural network training) problems, resolving an open question posed by Mulayoff & Michaeli (2020). This expression uncovers a fundamental property of the loss landscape of depth-2 matrix factorization problems: *a minimum is flat if and only if it is spectral-norm balanced*, which implies that *flat minima are not necessarily Frobenius-norm balanced*. Furthermore, to complement our theory, we empirically investigate an escape phenomenon observed during gradient-based training near a minimum that crucially relies on our exact expression of the sharpness.

1 INTRODUCTION

Decades of research in learning theory suggest limiting model complexity to prevent overfitting. However, modern deep learning is heavily overparameterized and has nonetheless achieved unprecedented success in practice over the past decade (Krizhevsky et al., 2012; Vaswani et al., 2017). Generally, in overparameterized settings, the loss function has infinitely many global minima that achieve zero training error (interpolation regime), yet these models still perform well. This phenomenon has been explored in various settings such as nonparametric regression, (Belkin et al., 2019), training two-layer neural networks with logistic loss (Frei et al., 2022), and linear regression (Bartlett et al., 2020).

The propensity of neural network training dynamics to converge to *good minima* is attributed to the ability of gradient-based optimization algorithms to avoid *bad minima* (Neyshabur et al., 2017; Zhang et al., 2017). This is related to the *implicit bias* of gradient descent (GD) (Neyshabur et al., 2014), and a large body of work has focused on its understanding (Gunasekar et al., 2017; 2018; Soudry et al., 2018; Arora et al., 2019; Ji & Telgarsky, 2020; Yun et al., 2021).

It has been observed that *dynamical stability* of GD near a minimum is a key factor in characterizing its implicit bias toward particular solutions (Wu et al., 2018; Nar & Sastry, 2018). Conceptually, dynamical stability refers to the ability of GD to *stably converge* to a minimum, and it is closely related to the sharpness of the loss landscape in its vicinity (Mulayoff et al., 2021). This topic has been investigated in numerous works (Nar & Sastry, 2018; Wu et al., 2018; Ma & Ying, 2021; Mulayoff et al., 2021; Nacson et al., 2023; Qiao et al., 2024; Liang et al., 2025) within the framework of the classical notion of *linear stability* in dynamical systems (Strogatz, 2024).

Ultimately, this understanding boils down to understanding the *geometry* of the loss landscape near a minimum. The maximum eigenvalue of the Hessian of the loss serves as a key measure to quantify the *sharpness* of the landscape near a minimum. Despite its significance, its precise role is not well-understood, particularly because closed-form expressions are generally unknown, outside a

Table 1: Closed-form expressions for the maximum Hessian eigenvalue in the literature. Ω denotes the set of *all* global minimizers, $\Omega_F \subseteq \Omega$ denotes the set of *flat* global minimizers, and $\Omega_B \subseteq \Omega$ denotes the set of *balanced* global minimizers.

Related Work	Depth	Input Dim.	Output Dim.	$\lambda_{\max}(\nabla^2 \mathcal{L}(\mathbf{w}))$	Layers
Mulayoff & Michaeli (2020, Theorem 1)	L	d_0	d_L	$\mathbf{w} \in \Omega_F$	$\mathbb{R}^{a \times b}$
Zhu et al. (2023, Appendix B.1)	2	1	1	$\mathbf{w} \in \mathbb{R}^N$	\mathbb{R}
Singh & Hofmann (2024, Theorem 1)	2	1	1	$\mathbf{w} \in \mathbb{R}^N$	\mathbb{R}^a
Ghosh et al. (2025, Lemma 1)	L	d_0	d_L	$\mathbf{w} \in \Omega_B$	$\mathbb{R}^{a \times b}$
Theorem 5 (This Paper)	L	d_0	d_L	$\mathbf{w} \in \Omega$	$\mathbb{R}^{a \times b}$

few particular cases. We summarize the current state of understanding as well as the contributions of our paper in Table 1.

Most notably, the seminal work of Mulayoff & Michaeli (2020) derives a closed-form expression for the maximum eigenvalue of the Hessian at *flat* global minima of deep linear networks (i.e., deep matrix factorization) with squared-error loss. However, obtaining a closed-form expression for *all* global minima in deep linear networks/deep matrix factorization was an open problem. In particular, Mulayoff & Michaeli (2020) claim that finding a closed-form expression for arbitrary global minima is intractable. In this paper, we refute this claim and positively answer the following fundamental question.

Does a closed-form expression for the maximum eigenvalue of the Hessian exist for overparameterized deep matrix factorization problems?

In particular, in Theorem 5, we provide a closed-form expression for the maximum Hessian eigenvalue at arbitrary minima of depth- L overparameterized deep matrix factorization. We also highlight that overparameterized deep matrix factorization and the deep linear neural network setting of Mulayoff & Michaeli (2020) are equivalent for investigating the sharpness measure $\lambda_{\max}(\nabla^2 \mathcal{L}(\mathbf{w}^*))$ when the data covariance matrix is identity, i.e., $\Sigma_x = \mathbf{I}$ (see Mulayoff & Michaeli, 2020, Equation 21). To the best of our knowledge, our analysis provides the first exact expression of the maximum eigenvalue for deep matrix factorization/deep linear neural network problems. In the case of deep overparameterized scalar factorization (Theorem 4) and depth-2 matrix factorization (Corollary 6), our closed-form expression simplifies considerably.

With our closed-form expression in hand, we then empirically explore in Section 7 the *escape phenomenon*, observed by Wu et al. (2018) (who only studied a one-dimensional setting). We find that this phenomenon also occurs for overparameterized deep matrix factorization problems. Therefore, we empirically observe the following.

GD escapes from a dynamically unstable minimum almost surely.

We explore this phenomenon through the lens of *dynamical stability* introduced by Wu et al. (2018).

Consider a twice continuously differentiable loss function $\mathcal{L} : \mathbb{R}^N \rightarrow \mathbb{R}$, and let \mathbf{x}^* be a minimizer of \mathcal{L} . For $\delta > 0$, denote the open ball $B_\delta(\mathbf{x}^*) := \{\mathbf{x} \in \mathbb{R}^N : \|\mathbf{x} - \mathbf{x}^*\|_2 < \delta\}$. There exists a $\delta > 0$ such that, for all $\mathbf{x} \in B_\delta(\mathbf{x}^*)$,

$$\mathcal{L}(\mathbf{x}) = \mathcal{L}(\mathbf{x}^*) + (\mathbf{x} - \mathbf{x}^*)^\top \nabla \mathcal{L}(\mathbf{x}^*) + \frac{1}{2}(\mathbf{x} - \mathbf{x}^*)^\top \nabla^2 \mathcal{L}(\mathbf{x}^*)(\mathbf{x} - \mathbf{x}^*) + o(\|\mathbf{x} - \mathbf{x}^*\|_2^2). \quad (1)$$

Definition 1. Let $\mathcal{L} : \mathbb{R}^N \rightarrow \mathbb{R}$ be a twice continuously differentiable loss function, and let \mathbf{x}^* be a minimizer of \mathcal{L} . The corresponding linearized GD update rule for \mathcal{L} is given by

$$\mathbf{x}_{t+1} = \mathbf{x}_t - \eta \nabla^2 \mathcal{L}(\mathbf{x}^*)(\mathbf{x}_t - \mathbf{x}^*), \quad (2)$$

where $\eta > 0$ is the step size. Given $\delta > 0$ such that (1) holds, if there exists an initial point $\mathbf{x}_0 \in B_\delta(\mathbf{x}^*)$ such that the residuals $\epsilon_t = \mathbf{x}_t - \mathbf{x}^*$ diverge, i.e.,

$$\lim_{t \rightarrow \infty} \|\epsilon_t\|_2 = \infty, \quad (3)$$

then \mathbf{x}^* is said to be dynamically unstable.

108 A necessary and sufficient condition for a minimum to be dynamically unstable is that
109 $\lambda_{\max}(\nabla^2 \mathcal{L}(\mathbf{x}^*)) > 2/\eta$ (cf. Wu et al. (2018); Mulayoff et al. (2021); Chemnitz & Engel (2025)).
110 Thus, we see that an exact expression for the maximum eigenvalue is necessary to explore the escape
111 phenomenon empirically.

113 1.1 CONTRIBUTIONS

114
115 In this paper, we present the first exact expression for the maximum eigenvalue of the Hessian
116 of the squared-error loss at any minimizer in general overparameterized deep matrix factorization
117 problems. Our results lead to the following remarkable observations:

- 118 • A minimizer of deep overparameterized scalar factorization loss is flat **if and only if** the
119 product of spectral norms of left and right intermediate factors is constant across layers
120 (Corollary 7).
- 121 • Flat minima are spectral-norm balanced in depth-2 matrix factorization (Corollary 8). This
122 implies that **flat minima are not necessarily balanced**, contrary to claims made in several
123 works (Ding et al., 2024; Ghosh et al., 2025). We further discuss this in Section 6.
- 124 • A minimizer of deep matrix factorization loss is flat **if** the product of spectral norms of left
125 and right intermediate factors is constant across layers (Corollary 10).
- 126 • A recent work (Anonymous, 2025) provides empirical evidence that, in two-layer linear
127 networks, the largest eigenvalue of the Hessian of the loss function is correlated with the
128 spectral norm of the weight matrices. Furthermore, for deep linear networks, they observe
129 that this largest Hessian eigenvalue correlates with the spectral norm of the products of the
130 weight matrices. Our results in Theorem 4, Theorem 5, and Corollary 6 provide the first
131 theoretical explanation for their empirical observations.

133 1.2 RELATED WORK

134
135 **Flat Minima.** Mulayoff & Michaeli (2020) derived a closed-form expression for the maximum
136 eigenvalue of the Hessian at flat minima for deep linear neural networks. They also showed that
137 the Hessian at a global minimum is rank-deficient by at least the order of $1 - 1/L$, where L is the
138 depth of the network. Moreover, they showed that the sharpness of the flattest minima increases
139 approximately linearly with L if $L \gg 1$. Singh & Hofmann (2024) provided a full characterization
140 of the Hessian spectrum at a point in parameter space for linear and ReLU networks in the *scalar*
141 *regression* case. They observed that the eigenvalues scale in proportion to the input variance within
142 one hidden-layer scalar linear networks. More recently, Josz (2025) has shown that locally flat
143 minima are globally flat in depth-2 matrix factorization problems.

144 **Balanced Minima.** Ghosh et al. (2025) provided a full characterization of the Hessian spectrum at
145 balanced minima in deep matrix factorization. Furthermore, they showed that the maximum eigen-
146 value of the Hessian at the flattest minima is equal to that of the balanced minima. Ding et al. (2024)
147 showed that *norm-minimal*, *balanced*, and *flat* solutions coincide in depth-2 matrix factorization,
148 where sharpness measured by the *scaled trace* of the Hessian matrix of the loss function. We further
149 discuss Ghosh et al. (2025) and Ding et al. (2024)’s results in Section 6. Finally, Wang et al. (2022)
150 showed that large step size GD training induces a *balancing effect* between factors in depth-2 matrix
151 factorization.

152 **Dynamical Stability.** In dynamical systems theory, it is well established that asymptotic conver-
153 gence to a critical point is determined solely by the local stability of that point (Strogatz, 2024).
154 In the seminal work of Wu et al. (2018) on the dynamical stability analysis of GD training, it was
155 shown that a global minimum is *dynamically stable* for GD if and only if the step size does not
156 exceed $2/\lambda_{\max}$, where λ_{\max} denotes the maximum eigenvalue of the Hessian of the loss at the
157 minimum. Mulayoff et al. (2021) investigated this mechanism in the space of learned functions for
158 two-layer overparameterized univariate ReLU networks in the interpolation regime. This was then
159 extended to multivariate ReLU networks by Nacson et al. (2023). The interpolation assumption was
160 then removed by Qiao et al. (2024); Liang et al. (2025).

161 **Edge-of-Stability.** Cohen et al. (2021) observed that neural networks trained with GD typically
operate in a regime called *edge of stability*, in which the maximum eigenvalue of the Hessian of

the loss function hovers just above the value $2/\eta$, where η is the step size, and argued that classical optimization theory fails to explain this phenomenon. Recently, Liang et al. (2025) empirically observed that explicit regularization seems to break the edge-of-stability phenomenon.

Sharpness and Generalization. It is widely recognized in the literature that flat minima are associated with better generalization (Hochreiter & Schmidhuber, 1997; Keskar et al., 2017). In a large-scale empirical investigation, Jiang et al. (2020) examined different measures for deep networks and found that a sharpness-based measure exhibited the strongest correlation with generalization. There is also theoretical evidence for this phenomenon in low-rank matrix recovery (Ding et al., 2024). On the other hand, Dinh et al. (2017) showed that *good minima* can be arbitrarily sharp in deep neural networks.

2 NOTATION, PRELIMINARIES, AND PROBLEM SETUP

We denote the *Kronecker product* by \otimes , the *Frobenius inner product* by $\langle \cdot, \cdot \rangle$, the *spectral norm* by $\sigma_{\max}(\cdot)$, and the *Frobenius norm* by $\|\cdot\|_F$. We denote by $[L]$ the set of natural numbers up to L , i.e., $[L] = \{1, 2, \dots, L\}$.

To simplify the notation for subsequent derivations, we define

$$\prod_{j=n}^m \mathbf{W}_j := \begin{cases} \mathbf{W}_m \mathbf{W}_{m-1} \cdots \mathbf{W}_n & \text{if } n \leq m, \\ \mathbf{I}_{d_m} & \text{otherwise, where } n, m \in [L], \end{cases} \quad (4)$$

where $\mathbf{W}_m \in \mathbb{R}^{d_m \times d_{m-1}}$.

Our analysis relies on matrix calculus and the formulation of directional second derivatives. Therefore, before proceeding to the technical details, we find it useful to first develop the intuition behind directional derivatives of real-valued functions of matrix variables.

Gâteaux Derivatives. Let $f : \mathbb{R}^{K \times L} \rightarrow \mathbb{R}$ be a differentiable function with continuous first- and second-order derivatives on $\mathbb{R}^{K \times L}$. Our objective is to derive closed-form expressions for the first- and second-order directional derivatives of f in the direction of $\mathbf{U} \in \mathbb{R}^{K \times L}$, where $\|\mathbf{U}\|_F < \infty$, denoted respectively by $D_{\mathbf{U}}f(\mathbf{X})$ and $D_{\mathbf{U}}^2f(\mathbf{X})$. By the limit definition of the derivative, the first derivative of $f(\mathbf{X})$ with respect to each entry of \mathbf{X} can be expressed as follows:

$$\frac{\partial f(\mathbf{X})}{\partial X_{ij}} = \lim_{\Delta t \rightarrow 0} \frac{f(\mathbf{X} + \Delta t \mathbf{e}_i \mathbf{e}_j^\top) - f(\mathbf{X})}{\Delta t}, \quad \forall (i, j) \in [K] \times [L], \quad (5)$$

where \mathbf{e}_i is the i^{th} standard basis vector of \mathbb{R}^K and \mathbf{e}_j is the j^{th} standard basis vector of \mathbb{R}^L . If the limit in (5) exists then by substitution of variables

$$\frac{\partial f(\mathbf{X})}{\partial X_{ij}} U_{ij} = \lim_{\Delta t \rightarrow 0} \frac{f(\mathbf{X} + \Delta t U_{ij} \mathbf{e}_i \mathbf{e}_j^\top) - f(\mathbf{X})}{\Delta t}, \quad \forall (i, j) \in [K] \times [L]. \quad (6)$$

By definition, the total change in $f(\mathbf{X})$ in the direction of \mathbf{U} is the sum of change due to each entry of \mathbf{X} . Then

$$D_{\mathbf{U}}f(\mathbf{X}) = \sum_{i,j \in [K] \times [L]} \frac{\partial f(\mathbf{X})}{\partial X_{ij}} U_{ij} \quad (7)$$

$$= \sum_{i,j \in [K] \times [L]} \lim_{\Delta t \rightarrow 0} \frac{f(\mathbf{X} + \Delta t U_{ij} \mathbf{e}_i \mathbf{e}_j^\top) - f(\mathbf{X})}{\Delta t} \quad (8)$$

$$= \lim_{\Delta t \rightarrow 0} \frac{f(\mathbf{X} + \Delta t \mathbf{U}) - f(\mathbf{X})}{\Delta t}. \quad (9)$$

We can rewrite (9) as follows:

$$\lim_{\Delta t \rightarrow 0} \frac{f(\mathbf{X} + (\Delta t + t)\mathbf{U}) - f(\mathbf{X} + t\mathbf{U})}{\Delta t} \Big|_{t=0} = \frac{\partial f(\mathbf{X} + t\mathbf{U})}{\partial t} \Big|_{t=0}. \quad (10)$$

This is known as the *Gâteaux derivative*, which represents the change in $f(\mathbf{X})$ under a perturbation in the direction of \mathbf{U} . By the same reasoning, we obtain the following result.

Lemma 2. The second directional derivative of f at \mathbf{X} in the direction $\mathbf{U} \in \mathbb{R}^{K \times L}$ is given by

$$D_{\mathbf{U}}^2 f(\mathbf{X}) = \left. \frac{\partial^2}{\partial t^2} f(\mathbf{X} + t\mathbf{U}) \right|_{t=0}. \quad (11)$$

The proof is deferred to Appendix A.1.

Directional Second Derivatives and Maximum Eigenvalue. Consider the following objective function for our real-valued matrix-variable function f .

$$f(\mathbf{W}_1, \mathbf{W}_2, \dots, \mathbf{W}_L) = \|\mathbf{M} - \mathbf{W}_L \mathbf{W}_{L-1} \cdots \mathbf{W}_1\|_F^2, \quad (12)$$

where $\mathbf{W}_i \in \mathbb{R}^{d_i \times d_{i-1}}$ and $\mathbf{M} \in \mathbb{R}^{d_L \times d_0}$ for all $i \in [L]$. In this setting, we can define the largest eigenvalue of the $\nabla^2 f(\mathbf{W}_1, \mathbf{W}_2, \dots, \mathbf{W}_L)$ at an arbitrary point in the parameter space as follows:

$$\lambda_{\max}(\nabla^2 f(\mathbf{W}_1, \mathbf{W}_2, \dots, \mathbf{W}_L)) = \max_{\substack{\mathbf{U}_1, \mathbf{U}_2, \dots, \mathbf{U}_L \\ \sum_{i=1}^L \|\mathbf{U}_i\|_F^2 = 1}} \left. \frac{d^2}{dt^2} f(\mathbf{W}_1 + t\mathbf{U}_1, \dots, \mathbf{W}_L + t\mathbf{U}_L) \right|_{t=0}. \quad (13)$$

This is the generalization of the Rayleigh quotient to the case where the Hessian is represented as a tensor and its eigenvectors take the form of matrices. This leads to the following lemma.

Lemma 3. For any $[\mathbf{W}_1^*, \mathbf{W}_2^*, \dots, \mathbf{W}_L^*]$ such that $\mathbf{M} = \prod_{j=1}^L \mathbf{W}_j^*$, the directional second derivative is given by

$$\nabla^2 f(\mathbf{W}_1^*, \dots, \mathbf{W}_L^*)[\mathbf{U}_1, \dots, \mathbf{U}_L] = 2 \left\| \sum_{i=1}^L \left[\left(\prod_{j=i+1}^L \mathbf{W}_j^* \right) \mathbf{U}_i \left(\prod_{j=1}^{i-1} \mathbf{W}_j^* \right) \right] \right\|_F^2. \quad (14)$$

The proof is deferred to Appendix A.2.

We study the sharpness of the loss landscape near any global minimum in deep matrix factorization problems. We consider the following optimization problem

$$\min_{\mathbf{w} \in \mathbb{R}^N} \mathcal{L}(\mathbf{w}) := \|\mathbf{M} - \mathbf{W}_L \mathbf{W}_{L-1} \cdots \mathbf{W}_1\|_F^2, \quad (15)$$

where $\mathbf{w} = \text{vec}([\mathbf{W}_1, \mathbf{W}_2, \dots, \mathbf{W}_L])$ denotes the collection of all parameters, and

$$N := \sum_{i=1}^L d_i \times d_{i-1} \quad (16)$$

is the total number of parameters in the model. $\mathbf{M} \in \mathbb{R}^{d_L \times d_0}$ denotes the matrix contains the parameters subject to factorization, $L \geq 2$ denotes the depth of factorization and $\mathbf{W}_i \in \mathbb{R}^{d_i \times d_{i-1}}$ is the i^{th} factor (layer). This objective is analogous to that of deep linear neural networks. To guarantee the feasibility of factorization at all points in $\mathbb{R}^{d_L \times d_0}$, we require

$$\min_i d_i \geq \min\{d_0, d_L\} \quad \forall i \in [L], \quad (17)$$

which follows directly from the fact that

$$\text{rank}(\mathbf{W}_L \mathbf{W}_{L-1} \cdots \mathbf{W}_1) \leq \min\{\text{rank}(\mathbf{W}_1), \text{rank}(\mathbf{W}_2), \dots, \text{rank}(\mathbf{W}_L)\}. \quad (18)$$

Define the set of global minima of $\mathcal{L}(\mathbf{w})$ as

$$\Omega := \arg \min_{\mathbf{w} \in \mathbb{R}^N} \mathcal{L}(\mathbf{w}) = \left\{ \mathbf{w} \in \mathbb{R}^N : \prod_{i=1}^L \mathbf{W}_i = \mathbf{M} \right\}. \quad (19)$$

Local Minima. Laurent & Brecht (2018) showed that all local minima of deep linear networks with convex and differentiable loss are global if the layers satisfy (17), i.e., hidden layers are at least as wide as input and output layers. In particular, they proved a more general theorem concerning real-valued functions that take as input a product of matrices. Thus, a corollary of Laurent & Brecht (2018, Theorem 1) is that overparameterized deep matrix factorization *does not have spurious minima*, i.e., all local minima are global.

3 DEEP OVERPARAMETERIZED SCALAR FACTORIZATION

Before we delve into our general results, we first investigate the deep overparameterized scalar factorization, i.e., a special case of deep matrix factorization in which the first and last layers are vectors. This simplified problem setup reveals the key proof techniques used to prove our general result in Section 4.

Theorem 4. Consider the following objective function

$$\mathcal{L}(\mathbf{w}) := (m - \mathbf{w}_L \mathbf{W}_{L-1} \cdots \mathbf{W}_2 \mathbf{w}_1)^2, \quad (20)$$

where $m \in \mathbb{R}$, $d_0 = d_L = 1$, $\mathbf{w}_L \in \mathbb{R}^{1 \times d_{L-1}}$ and $\mathbf{w}_1 \in \mathbb{R}^{d_1 \times 1}$. For hidden factors (layers), i.e, for all $i \in \{2, 3, \dots, L-1\}$, we have $\mathbf{W}_i \in \mathbb{R}^{d_i \times d_{i-1}}$. Then, $\forall \mathbf{w}^* \in \Omega$,

$$\lambda_{\max}(\nabla^2 \mathcal{L}(\mathbf{w}^*)) = 2 \sum_{i=1}^L \sigma_{\max} \left(\prod_{j=i+1}^L \mathbf{W}_j^* \right)^2 \sigma_{\max} \left(\prod_{j=1}^{i-1} \mathbf{W}_j^* \right)^2. \quad (21)$$

The proof appears in Appendix B.

4 OVERPARAMETERIZED DEEP MATRIX FACTORIZATION

We now consider the general deep matrix factorization problem. In this section, we prove our main result, which is closed-form expression for the maximum eigenvalue of the Hessian for any global minimum to the objective (15).

Theorem 5. If $\mathbf{w}^* \in \Omega$ then

$$\lambda_{\max}(\nabla^2 \mathcal{L}(\mathbf{w}^*)) = 2\sigma_{\max} \left(\sum_{i=1}^L \mathbf{B}_i^\top \mathbf{B}_i \otimes \mathbf{A}_i \mathbf{A}_i^\top \right), \quad (22)$$

where $\mathbf{A}_k = \prod_{i=k+1}^L \mathbf{W}_i^*$ and $\mathbf{B}_k = \prod_{i=1}^{k-1} \mathbf{W}_i^*$.

The proof appears in Appendix C.1. Note that the deep overparameterized scalar factorization is a special case of deep matrix factorization where both $\mathbf{B}_i^\top \mathbf{B}_i$ and $\mathbf{A}_i \mathbf{A}_i^\top$ reduce to scalars. In that special case, we recover Theorem 4. Another corollary of Theorem 5 is the maximum Hessian eigenvalue for the classical (depth-2) matrix factorization problem. This result may be of independent interest as the expression simplifies considerably.

Corollary 6. Consider the following depth-2 matrix factorization objective

$$\mathcal{L}(\mathbf{L}, \mathbf{R}) = \|\mathbf{M} - \mathbf{L}\mathbf{R}^\top\|_F^2, \quad (23)$$

where $\mathbf{M} \in \mathbb{R}^{d_L \times d_0}$ is the target matrix and $\mathbf{L} \in \mathbb{R}^{d_L \times k}$, $\mathbf{R} \in \mathbb{R}^{d_0 \times k}$. To ensure the feasibility of the factorization every point in $\mathbb{R}^{d_L \times d_0}$, we choose $k \geq \min\{d_0, d_L\}$. We define the set of minimizers as follows:

$$\Omega := \arg \min_{\mathbf{L}, \mathbf{R}} \mathcal{L}(\mathbf{L}, \mathbf{R}) = \{(\mathbf{L}, \mathbf{R}) : \mathbf{M} = \mathbf{L}\mathbf{R}^\top\}. \quad (24)$$

If $(\mathbf{L}, \mathbf{R}) \in \Omega$ then

$$\lambda_{\max}(\nabla^2 \mathcal{L}(\mathbf{L}, \mathbf{R})) = 2(\sigma_{\max}(\mathbf{L})^2 + \sigma_{\max}(\mathbf{R})^2). \quad (25)$$

Proof. We have from Theorem 5 that

$$\lambda_{\max}(\nabla^2 \mathcal{L}(\mathbf{L}, \mathbf{R})) = 2\sigma_{\max}(\mathbf{I} \otimes \mathbf{L}\mathbf{L}^\top + \mathbf{R}\mathbf{R}^\top \otimes \mathbf{I}). \quad (26)$$

Using the fact from Horn & Johnson (1994, Theorem 4.4.5), we can write

$$2\sigma_{\max}(\mathbf{I} \otimes \mathbf{L}\mathbf{L}^\top + \mathbf{R}\mathbf{R}^\top \otimes \mathbf{I}) = 2\sigma_{\max}(\mathbf{I} \otimes \mathbf{L}\mathbf{L}^\top) + 2\sigma_{\max}(\mathbf{R}\mathbf{R}^\top \otimes \mathbf{I}). \quad (27)$$

Note that for any matrix \mathbf{A} and \mathbf{B} , $\sigma_{\max}(\mathbf{A} \otimes \mathbf{B}) = \sigma_{\max}(\mathbf{A})\sigma_{\max}(\mathbf{B})$, and $\sigma_{\max}(\mathbf{A}\mathbf{A}^\top) = \sigma_{\max}(\mathbf{A}^\top \mathbf{A}) = \sigma_{\max}(\mathbf{A})^2$. Hence,

$$2(\sigma_{\max}(\mathbf{I} \otimes \mathbf{L}\mathbf{L}^\top) + \sigma_{\max}(\mathbf{R}\mathbf{R}^\top \otimes \mathbf{I})) = 2(\sigma_{\max}(\mathbf{L})^2 + \sigma_{\max}(\mathbf{R})^2). \quad (28)$$

□

We also provide a self-contained proof of this corollary in Appendix C.2. This result was also recently observed by Josz (2025) independently.

5 FLATNESS

In this section, we reveal remarkable aspects of the loss landscape of general deep matrix factorization problems. First, we show that an optimal solution in deep overparameterized scalar factorization problem is flat if and only if the product of spectral norms of left and right intermediate networks is constant across layers. This is a direct consequence of Theorem 4.

Corollary 7. *Consider the following objective function*

$$\mathcal{L}(\mathbf{w}) := (m - \mathbf{w}_L \mathbf{W}_{L-1} \cdots \mathbf{W}_2 \mathbf{w}_1)^2, \quad (29)$$

where $m \in \mathbb{R}$, $d_0 = d_L = 1$, $\mathbf{w}_L \in \mathbb{R}^{1 \times d_{L-1}}$ and $\mathbf{w}_1 \in \mathbb{R}^{d_1 \times 1}$. For hidden factors (layers), i.e. for all $i \in \{2, 3, \dots, L-1\}$, we have $\mathbf{W}_i \in \mathbb{R}^{d_i \times d_{i-1}}$. Then, $\mathbf{w}^* \in \Omega_F$ if and only if

$$\sigma_{\max}(\mathbf{A}_k) \sigma_{\max}(\mathbf{B}_k) = |m|^{1-\frac{1}{L}} \quad \forall k \in [L], \quad (30)$$

where $\mathbf{A}_k = \prod_{i=k+1}^L \mathbf{W}_i^*$ and $\mathbf{B}_k = \prod_{i=1}^{k-1} \mathbf{W}_i^*$.

Proof. For a minimizer \mathbf{w}^* , let us assume that $\sigma_{\max}(\mathbf{A}_k) \sigma_{\max}(\mathbf{B}_k) = |m|^{1-\frac{1}{L}}$ for all $k \in [L]$. Then, by Theorem 4, $\lambda_{\max}(\nabla^2 \mathcal{L}(\mathbf{w}^*)) = 2L \times m^{2(1-1/L)}$. This implies that \mathbf{w}^* is flat (Mulayoff & Michaeli, 2020, Theorem 1). Now, assume that \mathbf{w}^* is flat. (Mulayoff & Michaeli, 2020, Theorem 2) showed that for any flat minimum \mathbf{w}^* ,

$$\sigma_{\max}(\mathbf{A}_k) = |m|^{1-\frac{k}{L}} \quad \text{and} \quad \sigma_{\max}(\mathbf{B}_k) = |m|^{\frac{k-1}{L}} \quad \forall k \in [L]. \quad (31)$$

Hence, $\sigma_{\max}(\mathbf{A}_k) \sigma_{\max}(\mathbf{B}_k) = |m|^{1-\frac{1}{L}}$ for all $k \in [L]$. \square

Second, we show that being flat in depth-2 matrix factorization is equivalent to being spectral-norm balanced. This result is a direct consequence of Corollary 6 and AM-GM inequality (Josz, 2025, see also Lemma 5).

Corollary 8. *Consider the following depth-2 matrix factorization objective*

$$\mathcal{L}(\mathbf{L}, \mathbf{R}) = \|\mathbf{M} - \mathbf{L}\mathbf{R}^\top\|_F^2, \quad (32)$$

where $\mathbf{M} \in \mathbb{R}^{d_L \times d_0}$ is the target matrix and $\mathbf{L} \in \mathbb{R}^{d_L \times k}$, $\mathbf{R} \in \mathbb{R}^{d_0 \times k}$ such that $k \geq \min\{d_0, d_L\}$. Then, $(\mathbf{L}^*, \mathbf{R}^*) \in \Omega_F$ if and only if

$$\sigma_{\max}(\mathbf{L}^*) = \sigma_{\max}(\mathbf{R}^*) = \sqrt{\sigma_{\max}(\mathbf{M})}. \quad (33)$$

Remark 9. *We discuss further implications of this equivalence in the subsequent section.*

Now, we show that a minimizer of deep matrix factorization loss is flat if the product of spectral norms of left and right intermediate factors is constant across layers.

Corollary 10. *Consider the optimization objective in (15). If $\mathbf{w}^* \in \Omega$ then*

$$\lambda_{\max}(\nabla^2 \mathcal{L}(\mathbf{w}^*)) \leq 2 \sum_{k=1}^L \sigma_{\max}(\mathbf{A}_k)^2 \sigma_{\max}(\mathbf{B}_k)^2, \quad (34)$$

where $\mathbf{A}_k = \prod_{i=k+1}^L \mathbf{W}_i^*$ and $\mathbf{B}_k = \prod_{i=1}^{k-1} \mathbf{W}_i^*$ (see Appendix B). This implies that if \mathbf{w}^* satisfies $\sigma_{\max}(\mathbf{A}_k) \sigma_{\max}(\mathbf{B}_k) = \sigma_{\max}(\mathbf{M})^{1-\frac{1}{L}}$ for all $k \in [L]$, then \mathbf{w}^* is flat.

Proof. Mulayoff & Michaeli (2020) showed that for any $\mathbf{w}^* \in \Omega$,

$$\lambda_{\max}(\nabla^2 \mathcal{L}(\mathbf{w}^*)) \geq 2L \sigma_{\max}(\mathbf{M})^{2(1-1/L)}. \quad (35)$$

Now, assume that $\sigma_{\max}(\mathbf{A}_k) \sigma_{\max}(\mathbf{B}_k) = \sigma_{\max}(\mathbf{M})^{1-\frac{1}{L}}$ for all $k \in [L]$. Then, by (34),

$$\lambda_{\max}(\nabla^2 \mathcal{L}(\mathbf{w}^*)) = 2L \sigma_{\max}(\mathbf{M})^{2(1-1/L)}. \quad (36)$$

This implies that \mathbf{w}^* is flat (Mulayoff & Michaeli, 2020, Theorem 1). \square

6 DISCUSSION

In this section, we thoroughly examine the claims made by Ding et al. (2024) and Ghosh et al. (2025) that flat minima coincide with Frobenius-norm balanced minima in depth-2 matrix factorization and then we discuss how misleading sharpness measures can be for the loss landscape analysis.

Strict assumptions. First, we examine the latter work, where Ghosh et al. (2025) examined the learning dynamics of deep linear networks within the deep matrix factorization loss beyond the edge of stability. Their analysis relies on two strict assumptions, which the latter follows the first one. First, they introduce the *singular vector stationary set*, i.e., for any initialization of GD from this set, *GD does not rotate the layers during training*. In Proposition 4, they prove that the balanced initialization they considered in the paper, i.e., $\mathbf{W}_L(0) = \mathbf{W}_{L-1}(0) = \dots = \mathbf{W}_1(0) = \alpha \mathbf{I}$ where $\alpha \in \mathbb{R}$, is a member of the singular vector stationary set. This is leveraged to decouple the dynamics of the singular vectors and singular values. Second, for this specific initialization, singular values remain balanced during training, which is not true for arbitrary initialization. Under these assumptions, their optimization objective becomes a deep scalar factorization problem in which the spectral norm equals the Frobenius norm. Therefore, they claim that *Frobenius-norm balanced minima correspond to flat minima* in deep matrix factorization problem. However, this is only valid from the perspective of GD that initialized at a specific set of points. Globally, we can find minimizers that are **flat but not Frobenius-norm balanced**. To see this, consider a depth-2 matrix factorization problem, i.e.,

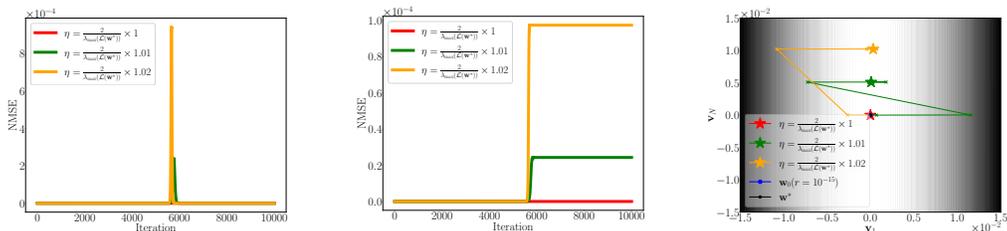
$$\min_{\mathbf{L}, \mathbf{R}} \mathcal{L}(\mathbf{L}, \mathbf{R}) \quad \text{such that} \quad \mathcal{L}(\mathbf{L}, \mathbf{R}) = \|\mathbf{M} - \mathbf{L}\mathbf{R}^\top\|_F^2, \quad (37)$$

where $\mathbf{M} \in \mathbb{R}^{m \times n}$, $\mathbf{L} \in \mathbb{R}^{m \times k}$ and $\mathbf{R} \in \mathbb{R}^{n \times k}$. Suppose $m = n = k = 3$ and $\mathbf{M} = \text{diag}(11, 2, 1)$. We know that for any flat minimum $(\mathbf{L}', \mathbf{R}')$, $\lambda_{\max}(\nabla^2(\mathcal{L}(\mathbf{L}', \mathbf{R}')))) = 4 \times \sigma_{\max}(\mathbf{M})$, which equals 44 (Mulayoff & Michaeli, 2020). By definition, any minimum whose sharpness is equal to 44 is a flat minimum. Let's investigate a specific minimizer $(\mathbf{L}^*, \mathbf{R}^*)$ such that $\mathbf{L}^* = \text{diag}(\sqrt{11}, \frac{2}{3}, 1)$ and $\mathbf{R}^* = \text{diag}(\sqrt{11}, 3, 1)$. Using Corollary 6, $\lambda_{\max}(\nabla^2(\mathcal{L}(\mathbf{L}^*, \mathbf{R}^*))) = 44$, which means that $(\mathbf{L}^*, \mathbf{R}^*)$ is a flat minimum. On the other hand, by definition of balancedness, i.e., $\mathbf{L}\mathbf{L}^\top = \mathbf{R}\mathbf{R}^\top$, $(\mathbf{L}^*, \mathbf{R}^*)$ is not a balanced minimizer. Therefore, any flat minimum is not necessarily a balanced minimizer.

Sharpness measures are delusive. Ding et al. (2024) showed that flat minima recover the groundtruth matrix in low rank matrix recovery. Furthermore, they showed that *norm-minimal*, *balanced*, and *flat* solutions coincide in depth-2 matrix factorization. They measure sharpness of a minimum by using *scaled trace*. Under this measure, flat minima coincide with Frobenius-norm balanced minima; however, as we showed above, this is not necessarily true when the worst-case sharpness, i.e., maximum Hessian eigenvalue, is used as the sharpness measure. In fact, we showed that flat minima are spectral-norm balanced when the maximum Hessian eigenvalue is used as the sharpness measure. This means that a loss landscape analysis in general matrix factorization problems using a specific sharpness measure might not lead to the same inferences as those made by using a different metric.

The necessity of an absolute sharpness measure. There is still no consensus in the literature regarding the definition of flatness. For instance, while several works define flat minima as global minimizers that minimize the maximum eigenvalue of the Hessian of the loss (Mulayoff & Michaeli, 2020; Liu et al., 2021; Marion & Chizat, 2024), others define them as global minimizers that minimize the trace of the Hessian (Dinh et al., 2017; Gatmiry et al., 2023). Even a scaled version of the Hessian trace designed to define flat minima in depth-2 matrix factorization problems (Ding et al., 2024), as we mentioned before. As we have shown in Section 5 and discussed above, different sharpness measures can lead to contradictory interpretations of balanced minima. Thus, the necessity of a robust sharpness measure that consolidates the interpretations from existing ones is a matter of urgency. Addressing this concern, Josz (2025) showed how minimizing the maximum Hessian eigenvalue over the solution set can disregard the higher-order variation near minimizers and result in misleading interpretations of flat minima. Therefore, Josz (2025) defined the flat minima of any smooth function \mathcal{L} as the local minima of $\lambda_{\max}(\nabla^2 \mathcal{L}(\mathbf{w}))$ under the constraint $\mathcal{L}(\mathbf{w}) = \mathcal{L}(\mathbf{w}^*)$ and demonstrated that this notion coincides with other notions of flatness in the depth-2 matrix factorization.

432
433
434
435
436
437
438
439
440
441
442
443

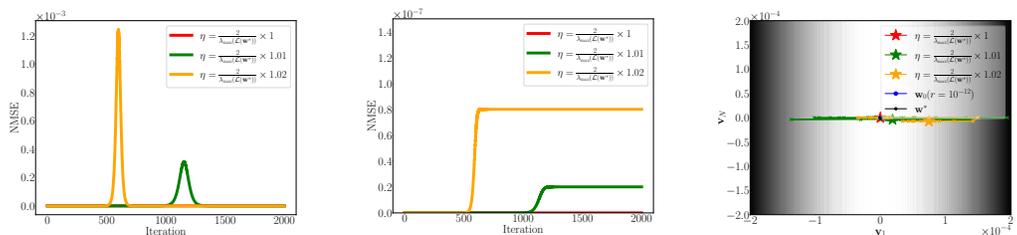


(a) Normalized training error across iterations, i.e., $\mathcal{L}(\mathbf{w}_k)/\|\mathbf{M}\|_F^2$. (b) Normalized ℓ^2 distance of \mathbf{w}_k from the minimum \mathbf{w}^* , i.e., $\|\mathbf{w}_k - \mathbf{w}^*\|_2 / \|\mathbf{w}^*\|_2$. (c) Trajectories of GD on the contour map of the loss landscape around the minimum.

444
445
446
447
448
449
450

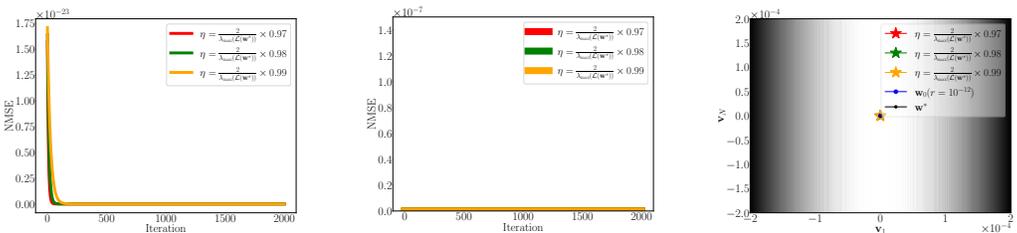
Figure 1: GD dynamics with different step sizes, $\eta \geq 2/\lambda_{\max}(\nabla^2 \mathcal{L}(\mathbf{w}^*))$, indicated by different colors, are initialized within a radius of 10^{-15} from the minimum in the direction of the Hessian eigenvector corresponding to the largest eigenvalue, for depth-2 matrix factorization, $\mathbf{M} = \mathbf{L}\mathbf{R}^\top$, of a random Gaussian matrix, where $\mathbf{L} \in \mathbb{R}^{10 \times 20}$ and $\mathbf{R} \in \mathbb{R}^{20 \times 20}$. The vector \mathbf{v}_1 denotes the eigenvector of the Hessian corresponding to the largest eigenvalue, while \mathbf{v}_N denotes the eigenvector corresponding to the smallest eigenvalue. The value of $\lambda_{\max}(\nabla^2 \mathcal{L}(\mathbf{w}^*))$ is computed using the closed-form expression derived in Corollary 6.

451
452
453
454
455
456
457
458
459
460
461



(a) Normalized training error across iterations, i.e., $\mathcal{L}(\mathbf{w}_k)/m^2$. (b) Normalized ℓ^2 distance of \mathbf{w}_k from the minimum \mathbf{w}^* . (c) Trajectories of GD on the contour map of the loss landscape.

462
463
464
465
466
467
468
469
470
471



(d) Normalized training error across iterations, i.e., $\mathcal{L}(\mathbf{w}_k)/m^2$. (e) Normalized ℓ^2 distance of \mathbf{w}_k from the minimum \mathbf{w}^* . (f) Trajectories of GD on the contour map of the loss landscape.

472
473
474
475
476

Figure 2: GD dynamics with different step sizes indicated by different colors are initialized within a radius of 10^{-12} from the minimum in the direction of the Hessian eigenvector corresponding to the largest eigenvalue, for a 15-layer overparameterized scalar factorization of a random scalar. The value of $\lambda_{\max}(\nabla^2 \mathcal{L}(\mathbf{w}^*))$ is computed using the closed-form expression derived in Theorem 4.

477 478 479 7 EXPERIMENTS

480
481
482
483
484
485

The dynamical stability analysis relies on how accurate the quadratic approximation of the loss function is in the δ -neighborhood of \mathbf{x}^* (Wu et al., 2018) (also see Definition 1). Therefore, we initialize GD in a δ -neighborhood with δ on the order of 10^{-15} to 10^{-9} to ensure validity of the quadratic approximation. To observe escape phenomenon clearly, we choose the perturbation direction for the initial point to be the eigenvector of $\nabla^2 \mathcal{L}(\mathbf{w}^*)$ corresponding to the largest eigenvalue, so as to avoid choosing a direction that is orthogonal to the eigenspace corresponding to the eigenvalues larger than $2/\eta$.

Note that, we could have also considered a random perturbation direction uniformly from \mathbb{S}^{N-1} . Mulayoff & Michaeli (2020) showed that for the deep matrix factorization problem (which has the same Hessian structure with deep linear networks), the Hessian is rank-deficient at all minima by at least the order of $1 - (1/L)$. This means that at least $1 - (1/L)$ of the eigenvalues are zero at a minimum. Thus, the probability of choosing a direction that is orthogonal to the eigenspace corresponding to eigenvalues larger than $2/\eta$ is 0. Therefore, random perturbations would lead to the escape phenomenon with probability 1.

On the other hand, for stable minima, it is important to choose the direction as the eigenvector of $\nabla^2\mathcal{L}(\mathbf{w}^*)$ corresponding to the largest eigenvalue. The reasoning is the same. If you choose a direction that is not orthogonal to the eigenspace corresponding to the eigenvalues that are zero, then GD never converges to \mathbf{w}^* . This means that if we choose the perturbation direction randomly, then with probability 1, we choose a direction that is not orthogonal to the eigenspace corresponding to the eigenvalues that are zero. Therefore, for the experiment, it is convenient to choose perturbation direction as the eigenvector of $\nabla^2\mathcal{L}(\mathbf{w}^*)$ corresponding to the largest eigenvalue.

To measure the distance between the convergence point and the minimizer, we plot the normalized ℓ^2 -norm of $\mathbf{w}_k - \mathbf{w}^*$ at each iteration. Furthermore, as shown in Fig. 1 and Figs. 2a-2c if $\eta > 2/\lambda_{\max}$, where $\lambda_{\max} := \lambda_{\max}(\nabla^2\mathcal{L}(\mathbf{w}^*))$, GD always escapes from the minimum. On the other hand, if $\eta = 2/\lambda_{\max}$ then GD converges as shown in Figs. 2d-2f. A catapult in the training error indicates GD's escape from the basin of a minimum, after which it eventually converges to another minimum as observed by Wu et al. (2018). For the methodology used to generate contour maps of the loss landscape near a minimum, see Appendix D.1, and for additional experiments, see Appendix D.2.

8 CONCLUSION

In this paper, we derived an exact expression for the maximum eigenvalue of the Hessian of the squared-error loss for overparameterized deep matrix factorization problems at any minimizer. We also showed that this expression simplifies considerably for depth-2 matrix factorization and deep, overparameterized scalar factorization. To complement our theory, we conducted GD experiments that crucially rely on our exact expression of the sharpness to observe the escape phenomenon during training. Then, we showed how flat minima are spectral-norm balanced in depth-2 matrix factorization, and we showed that a minimizer of deep overparameterized scalar factorization loss is flat if and only if the product of spectral norms of left and right intermediate factors is constant across layers. Furthermore, we showed that a minimizer of deep matrix factorization loss is flat if the product of spectral norms of left and right intermediate factors is constant across layers. These results led to the following remarkable observation: flat minima are not necessarily Frobenius-norm balanced, which is contradictory to the claim made in (Ghosh et al., 2025). Finally, we discussed how delusive the sharpness measures can be for the loss landscape analysis. Therefore, there is an urgent need for future research to explore a robust sharpness measure that consolidates the interpretations from existing ones.

REFERENCES

- Anonymous. Cracking the hessian: Closed-form hessian spectra for fundamental neural networks. In *Submitted to The Fourteenth International Conference on Learning Representations*, 2025. under review.
- Sanjeev Arora, Nadav Cohen, Wei Hu, and Yuping Luo. Implicit regularization in deep matrix factorization. *Advances in neural information processing systems*, 32, 2019.
- Peter L Bartlett, Philip M Long, Gábor Lugosi, and Alexander Tsigler. Benign overfitting in linear regression. *Proceedings of the National Academy of Sciences*, 117(48):30063–30070, 2020.
- Mikhail Belkin, Alexander Rakhlin, and Alexandre B Tsybakov. Does data interpolation contradict statistical optimality? In *The 22nd international conference on artificial intelligence and statistics*, pp. 1611–1619. PMLR, 2019.

540 Dennis Chemnitz and Maximilian Engel. Characterizing dynamical stability of stochastic gradient
541 descent in overparameterized learning. *Journal of Machine Learning Research*, 26(134):1–46,
542 2025.

543 Jeremy Cohen, Simran Kaur, Yuanzhi Li, J Zico Kolter, and Ameet Talwalkar. Gradient descent on
544 neural networks typically occurs at the edge of stability. In *International Conference on Learning*
545 *Representations*, 2021.

546 Lijun Ding, Dmitriy Drusvyatskiy, Maryam Fazel, and Zaid Harchaoui. Flat minima generalize for
547 low-rank matrix recovery. *Information and Inference: A Journal of the IMA*, 13(2):iaae009, 2024.

548 Laurent Dinh, Razvan Pascanu, Samy Bengio, and Yoshua Bengio. Sharp minima can generalize
549 for deep nets. In *International Conference on Machine Learning*, pp. 1019–1028. PMLR, 2017.

550 Spencer Frei, Niladri S Chatterji, and Peter Bartlett. Benign overfitting without linearity: Neural
551 network classifiers trained by gradient descent for noisy linear data. In *Conference on Learning*
552 *Theory*, pp. 2668–2703. PMLR, 2022.

553 Khashayar Gatmiry, Zhiyuan Li, Tengyu Ma, Sashank Reddi, Stefanie Jegelka, and Ching-Yao
554 Chuang. What is the inductive bias of flatness regularization? a study of deep matrix factor-
555 ization models. *Advances in Neural Information Processing Systems*, 36:28040–28052, 2023.

556 Avrajit Ghosh, Soo Min Kwon, Rongrong Wang, Saiprasad Ravishankar, and Qing Qu. Learning dy-
557 namics of deep matrix factorization beyond the edge of stability. In *The Thirteenth International*
558 *Conference on Learning Representations*, 2025.

559 Ian J Goodfellow, Oriol Vinyals, and Andrew M Saxe. Qualitatively characterizing neural network
560 optimization problems. *arXiv preprint arXiv:1412.6544*, 2014.

561 Suriya Gunasekar, Blake E Woodworth, Srinadh Bhojanapalli, Behnam Neyshabur, and Nati Sre-
562 bro. Implicit regularization in matrix factorization. *Advances in neural information processing*
563 *systems*, 30, 2017.

564 Suriya Gunasekar, Jason D Lee, Daniel Soudry, and Nati Srebro. Implicit bias of gradient descent
565 on linear convolutional networks. *Advances in neural information processing systems*, 31, 2018.

566 Sepp Hochreiter and Jürgen Schmidhuber. Flat minima. *Neural computation*, 9(1):1–42, 1997.

567 Roger A Horn and Charles R Johnson. *Topics in matrix analysis*. Cambridge university press, 1994.

568 Ziwei Ji and Matus Telgarsky. Directional convergence and alignment in deep learning. *Advances*
569 *in Neural Information Processing Systems*, 33:17176–17186, 2020.

570 Yiding Jiang, Behnam Neyshabur, Hossein Mobahi, Dilip Krishnan, and Samy Bengio. Fant-
571 astic generalization measures and where to find them. In *International Conference on Learning*
572 *Representations*, 2020.

573 Cédric Jozz. On the geometry of flat minima, 2025.

574 Nitish Shirish Keskar, Dheevatsa Mudigere, Jorge Nocedal, Mikhail Smelyanskiy, and Ping Tak Pe-
575 ter Tang. On large-batch training for deep learning: Generalization gap and sharp minima. In
576 *International Conference on Learning Representations*, 2017.

577 Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convo-
578 lutional neural networks. *Advances in neural information processing systems*, 25, 2012.

579 Thomas Laurent and James Brecht. Deep linear networks with arbitrary loss: All local minima are
580 global. In *International conference on machine learning*, pp. 2902–2907. PMLR, 2018.

581 Hao Li, Zheng Xu, Gavin Taylor, Christoph Studer, and Tom Goldstein. Visualizing the loss land-
582 scape of neural nets. *Advances in neural information processing systems*, 31, 2018.

583 Tongtong Liang, Dan Qiao, Yu-Xiang Wang, and Rahul Parhi. Stable minima of relu neural net-
584 works suffer from the curse of dimensionality: The neural shattering phenomenon. *arXiv preprint*
585 *arXiv:2506.20779*, 2025.

-
- 594 Tianyi Liu, Yan Li, Song Wei, Enlu Zhou, and Tuo Zhao. Noisy gradient descent converges to flat
595 minima for nonconvex matrix factorization. In *International Conference on Artificial Intelligence*
596 *and Statistics*, 2021.
- 597 Chao Ma and Lexing Ying. On linear stability of sgd and input-smoothness of neural networks.
598 *Advances in Neural Information Processing Systems*, 34:16805–16817, 2021.
- 600 Pierre Marion and Lénaïc Chizat. Deep linear networks for regression are implicitly regularized
601 towards flat minima. *Advances in Neural Information Processing Systems*, 37:76848–76900,
602 2024.
- 603 Rotem Mulayoff and Tomer Michaeli. Unique properties of flat minima in deep networks. In
604 *International conference on machine learning*, pp. 7108–7118. PMLR, 2020.
- 605 Rotem Mulayoff, Tomer Michaeli, and Daniel Soudry. The implicit bias of minima stability: A
606 view from function space. *Advances in Neural Information Processing Systems*, 34:17749–17761,
607 2021.
- 608 Mor Shpigel Nacson, Rotem Mulayoff, Greg Ongie, Tomer Michaeli, and Daniel Soudry. The
609 implicit bias of minima stability in multivariate shallow relu networks. In *ICLR*, 2023.
- 611 Kamil Nar and Shankar Sastry. Step size matters in deep learning. *Advances in Neural Information*
612 *Processing Systems*, 31, 2018.
- 613 Behnam Neyshabur, Ryota Tomioka, and Nathan Srebro. In search of the real inductive bias: On the
614 role of implicit regularization in deep learning. *arXiv preprint arXiv:1412.6614*, 2014.
- 615 Behnam Neyshabur, Srinadh Bhojanapalli, David McAllester, and Nati Srebro. Exploring general-
616 ization in deep learning. *Advances in neural information processing systems*, 30, 2017.
- 617 Dan Qiao, Kaiqi Zhang, Esha Singh, Daniel Soudry, and Yu-Xiang Wang. Stable minima cannot
618 overfit in univariate relu networks: Generalization by large step sizes. *Advances in Neural Infor-*
619 *mation Processing Systems*, 37:94163–94208, 2024.
- 620 Sidak Pal Singh and Thomas Hofmann. Closed form of the hessian spectrum for some neural
621 networks. In *High-dimensional Learning Dynamics 2024: The Emergence of Structure and Rea-*
622 *soning*, 2024.
- 623 Daniel Soudry, Elad Hoffer, Mor Shpigel Nacson, Suriya Gunasekar, and Nathan Srebro. The im-
624 plicit bias of gradient descent on separable data. *Journal of Machine Learning Research*, 19(70):
625 1–57, 2018.
- 626 Steven H Strogatz. *Nonlinear dynamics and chaos: with applications to physics, biology, chemistry,*
627 *and engineering*. Chapman and Hall/CRC, 2024.
- 628 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez,
629 Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural informa-*
630 *tion processing systems*, 30, 2017.
- 631 Yuqing Wang, Minshuo Chen, Tuo Zhao, and Molei Tao. Large learning rate tames homogeneity:
632 Convergence and balancing effect. In *International Conference on Learning Representations*,
633 2022.
- 634 Lei Wu, Chao Ma, et al. How sgd selects the global minima in over-parameterized learning: A
635 dynamical stability perspective. *Advances in Neural Information Processing Systems*, 31, 2018.
- 636 Chulhee Yun, Shankar Krishnan, and Hossein Mobahi. A unifying view on implicit bias in training
637 linear neural networks. In *International Conference on Learning Representations*, 2021.
- 638 Chiyuan Zhang, Samy Bengio, Moritz Hardt, Benjamin Recht, and Oriol Vinyals. Understanding
639 deep learning requires rethinking generalization. In *International Conference on Learning Rep-*
640 *resentations*, 2017.
- 641 Xingyu Zhu, Zixuan Wang, Xiang Wang, Mo Zhou, and Rong Ge. Understanding edge-of-stability
642 training dynamics with a minimalist example. In *The Eleventh International Conference on*
643 *Learning Representations*, 2023.

A PROOFS FROM SECTION 2

A.1 PROOF OF LEMMA 2

Proof. We can express the derivative of $\frac{\partial f(\mathbf{X})}{\partial X_{ij}}$ with respect to each entry of \mathbf{X} , using the limit definition of the derivative, as follows:

$$\frac{\partial^2 f(\mathbf{X})}{\partial X_{kl} \partial X_{ij}} = \frac{\partial}{\partial X_{kl}} \left(\frac{\partial f(\mathbf{X})}{\partial X_{ij}} \right) = \lim_{\Delta t \rightarrow 0} \frac{\partial f(\mathbf{X} + \Delta t \mathbf{e}_i \mathbf{e}_j^\top) - \partial f(\mathbf{X})}{\partial X_{kl} \Delta t}, \quad \forall (k, l) \in [K] \times [L] \quad (38)$$

which is equal to

$$\lim_{\Delta h, \Delta t \rightarrow 0} \frac{f(\mathbf{X} + \Delta t \mathbf{e}_i \mathbf{e}_j^\top + \Delta h \mathbf{e}_k \mathbf{e}_l^\top) - f(\mathbf{X} + \Delta t \mathbf{e}_i \mathbf{e}_j^\top) - f(\mathbf{X} + \Delta h \mathbf{e}_k \mathbf{e}_l^\top) + f(\mathbf{X})}{\Delta h \Delta t}. \quad (39)$$

By using the substitution of variables as in (6),

$$\frac{\partial^2 f(\mathbf{X})}{\partial X_{kl} \partial X_{ij}} U_{ij} U_{kl} = \frac{\partial}{\partial X_{kl}} \left(\frac{\partial f(\mathbf{X})}{\partial X_{ij}} U_{ij} \right) U_{kl} = \lim_{\Delta t \rightarrow 0} \frac{\partial f(\mathbf{X} + \Delta t U_{ij} \mathbf{e}_i \mathbf{e}_j^\top) - \partial f(\mathbf{X})}{\partial X_{kl} \Delta t} U_{kl}. \quad (40)$$

Equivalently,

$$\lim_{\Delta h, \Delta t \rightarrow 0} \frac{f(\mathbf{X} + \Delta t U_{ij} \mathbf{e}_i \mathbf{e}_j^\top + \Delta h U_{kl} \mathbf{e}_k \mathbf{e}_l^\top) - f(\mathbf{X} + \Delta t U_{ij} \mathbf{e}_i \mathbf{e}_j^\top) - f(\mathbf{X} + \Delta h U_{kl} \mathbf{e}_k \mathbf{e}_l^\top) + f(\mathbf{X})}{\Delta h \Delta t}. \quad (41)$$

which can be proved by substitution of variables in (40). In turn, second order differential due to any $\mathbf{U} \in \mathbb{R}^{K \times L}$ is

$$D_{\mathbf{U}}^2 f(\mathbf{X}) = \sum_{i,j} \sum_{k,l} \frac{\partial^2 f(\mathbf{X})}{\partial X_{kl} \partial X_{ij}} U_{ij} U_{kl} = \left\langle \nabla \langle \nabla f(\mathbf{X}), \mathbf{U} \rangle, \mathbf{U} \right\rangle, \quad (42)$$

where

$$\nabla f(\mathbf{X}) = \begin{bmatrix} \frac{\partial f(\mathbf{X})}{\partial X_{11}} & \frac{\partial f(\mathbf{X})}{\partial X_{12}} & \cdots & \frac{\partial f(\mathbf{X})}{\partial X_{1L}} \\ \frac{\partial f(\mathbf{X})}{\partial X_{21}} & \frac{\partial f(\mathbf{X})}{\partial X_{22}} & \cdots & \frac{\partial f(\mathbf{X})}{\partial X_{2L}} \\ \vdots & \vdots & \cdots & \vdots \\ \frac{\partial f(\mathbf{X})}{\partial X_{K1}} & \frac{\partial f(\mathbf{X})}{\partial X_{K2}} & \cdots & \frac{\partial f(\mathbf{X})}{\partial X_{KL}} \end{bmatrix} \in \mathbb{R}^{K \times L}. \quad (43)$$

Equivalently,

$$D_{\mathbf{U}}^2 f(\mathbf{X}) = \sum_{k,l} \lim_{\Delta t \rightarrow 0} \frac{\partial f(\mathbf{X} + \Delta t \mathbf{U} \mathbf{e}_i \mathbf{e}_j^\top) - \partial f(\mathbf{X})}{\partial X_{kl} \Delta t} U_{kl} \quad (44)$$

$$= \lim_{\Delta t \rightarrow 0} \frac{f(\mathbf{X} + 2\Delta t \mathbf{U}) - 2f(\mathbf{X} + \Delta t \mathbf{U}) + f(\mathbf{X})}{\Delta t^2} \quad (45)$$

$$= \frac{\partial^2}{\partial t^2} f(\mathbf{X} + t\mathbf{U}) \Big|_{t=0}. \quad (46)$$

□

A.2 PROOF OF LEMMA 3

Proof. We can rewrite (12) through using the definition of *Frobenius inner product*

$$f(\mathbf{W}_1, \dots, \mathbf{W}_L) = \left\langle M - \prod_{i=1}^L \mathbf{W}_i, M - \prod_{i=1}^L \mathbf{W}_i \right\rangle. \quad (47)$$

702 Then

$$703 f(\mathbf{W}_1 + t\mathbf{U}_1, \dots, \mathbf{W}_L + t\mathbf{U}_L) = \left\langle M - \prod_{i=1}^L (\mathbf{W}_i + t\mathbf{U}_i), M - \prod_{i=1}^L (\mathbf{W}_i + t\mathbf{U}_i) \right\rangle. \quad (48)$$

704 Let's define $g : \mathbb{R} \rightarrow \mathbb{R}^{d_L \times d_0}$ such that

$$705 g(t) = M - \prod_{i=1}^L (\mathbf{W}_i + t\mathbf{U}_i), \quad f(\mathbf{W}_1 + t\mathbf{U}_1, \dots, \mathbf{W}_L + t\mathbf{U}_L) = \langle g(t), g(t) \rangle. \quad (49)$$

706 First, we need to differentiate $f(\mathbf{W}_1 + t\mathbf{U}_1, \dots, \mathbf{W}_L + t\mathbf{U}_L)$ w.r.t t . Using the fact that $\langle \mathbf{A}, \mathbf{B} \rangle = \text{tr}(\mathbf{A}^\top \mathbf{B})$, which simplifies the differentiation,

$$707 \frac{d}{dt} f(\mathbf{W}_1 + t\mathbf{U}_1, \dots, \mathbf{W}_L + t\mathbf{U}_L) = 2 \langle g(t), g'(t) \rangle. \quad (50)$$

$$708 \frac{d^2}{dt^2} f(\mathbf{W}_1 + t\mathbf{U}_1, \dots, \mathbf{W}_L + t\mathbf{U}_L) = 2 \langle g'(t), g'(t) \rangle + 2 \langle g(t), g''(t) \rangle. \quad (51)$$

709 Then, the directional second derivative $\nabla^2 f(\mathbf{W}_1, \dots, \mathbf{W}_L)[\mathbf{U}_1, \dots, \mathbf{U}_L]$ equals

$$710 \frac{d^2}{dt^2} f(\mathbf{W}_1 + t\mathbf{U}_1, \dots, \mathbf{W}_L + t\mathbf{U}_L) \Big|_{t=0} = 2 \langle g'(0), g'(0) \rangle + 2 \langle g(0), g''(0) \rangle. \quad (52)$$

711 It is straightforward to differentiate $g(t)$ such that

$$712 g(0) = M - \prod_{i=1}^L \mathbf{W}_i, \quad (53)$$

$$713 g'(0) = - \sum_{i=1}^L \left[\left(\prod_{j=i+1}^L \mathbf{W}_j \right) \mathbf{U}_i \left(\prod_{j=1}^{i-1} \mathbf{W}_j \right) \right], \quad (54)$$

$$714 g''(0) = -2 \sum_{1 \leq k < i \leq L} \left[\left(\prod_{j=i+1}^L \mathbf{W}_j \right) \mathbf{U}_i \left(\prod_{j=k+1}^{i-1} \mathbf{W}_j \right) \mathbf{U}_k \left(\prod_{j=1}^{k-1} \mathbf{W}_j \right) \right]. \quad (55)$$

715 Therefore, for any $[\mathbf{W}_1, \mathbf{W}_2, \dots, \mathbf{W}_L]$ in parameter space

$$716 \nabla^2 f(\mathbf{W}_1, \dots, \mathbf{W}_L)[\mathbf{U}_1, \dots, \mathbf{U}_L] = \quad (56)$$

$$717 2 \left\langle \sum_{i=1}^L \left[\left(\prod_{j=i+1}^L \mathbf{W}_j \right) \mathbf{U}_i \left(\prod_{j=1}^{i-1} \mathbf{W}_j \right) \right], \sum_{i=1}^L \left[\left(\prod_{j=i+1}^L \mathbf{W}_j \right) \mathbf{U}_i \left(\prod_{j=1}^{i-1} \mathbf{W}_j \right) \right] \right\rangle \quad (57)$$

$$718 -4 \left\langle M - \prod_{i=1}^L \mathbf{W}_i, \sum_{1 \leq k < i \leq L} \left[\left(\prod_{j=i+1}^L \mathbf{W}_j \right) \mathbf{U}_i \left(\prod_{j=k+1}^{i-1} \mathbf{W}_j \right) \mathbf{U}_k \left(\prod_{j=1}^{k-1} \mathbf{W}_j \right) \right] \right\rangle. \quad (58)$$

719 Note that for any minimizer $[\mathbf{W}_1^*, \mathbf{W}_2^*, \dots, \mathbf{W}_L^*]$, $M - \prod_{j=1}^L \mathbf{W}_j^* = 0$. Hence, for any global minimum

$$720 \nabla^2 f(\mathbf{W}_1^*, \dots, \mathbf{W}_L^*)[\mathbf{U}_1, \dots, \mathbf{U}_L] = \quad (59)$$

$$721 2 \left\langle \sum_{i=1}^L \left[\left(\prod_{j=i+1}^L \mathbf{W}_j^* \right) \mathbf{U}_i \left(\prod_{j=1}^{i-1} \mathbf{W}_j^* \right) \right], \sum_{i=1}^L \left[\left(\prod_{j=i+1}^L \mathbf{W}_j^* \right) \mathbf{U}_i \left(\prod_{j=1}^{i-1} \mathbf{W}_j^* \right) \right] \right\rangle. \quad (60)$$

722 □

B WARM-UP: DEEP OVERPARAMETERIZED SCALAR FACTORIZATION

Proof. According to (13) and (14),

$$\lambda_{\max}(\nabla^2 \mathcal{L}(\mathbf{w}^*)) = \max_{\substack{\mathbf{U}_1, \mathbf{U}_2, \dots, \mathbf{U}_L: \\ \sum_{i=1}^L \|\mathbf{U}_i\|_F^2 = 1}} 2 \left\| \sum_{i=1}^L \left[\left(\prod_{j=i+1}^L \mathbf{W}_j^* \right) \mathbf{U}_i \left(\prod_{j=1}^{i-1} \mathbf{W}_j^* \right) \right] \right\|_F^2 \quad (61)$$

$$\leq \max_{\substack{\mathbf{U}_1, \mathbf{U}_2, \dots, \mathbf{U}_L: \\ \sum_{i=1}^L \|\mathbf{U}_i\|_F^2 = 1}} 2 \left(\sum_{i=1}^L \left\| \left(\prod_{j=i+1}^L \mathbf{W}_j^* \right) \mathbf{U}_i \left(\prod_{j=1}^{i-1} \mathbf{W}_j^* \right) \right\|_F \right)^2 \quad (62)$$

$$= \max_{\substack{\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_L: \\ \sum_{i=1}^L \|\mathbf{u}_i\|_2^2 = 1}} 2 \left(\sum_{i=1}^L \left\| \left[\left(\prod_{j=1}^{i-1} \mathbf{W}_j^* \right)^\top \otimes \left(\prod_{j=i+1}^L \mathbf{W}_j^* \right) \right] \mathbf{u}_i \right\|_2 \right)^2 \quad (63)$$

$$\leq \max_{\substack{\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_L: \\ \sum_{i=1}^L \|\mathbf{u}_i\|_2^2 = 1}} 2 \left(\sum_{i=1}^L \sigma_{\max} \left(\left(\prod_{j=1}^{i-1} \mathbf{W}_j^* \right)^\top \otimes \left(\prod_{j=i+1}^L \mathbf{W}_j^* \right) \right) \|\mathbf{u}_i\|_2 \right)^2 \quad (64)$$

$$= \max_{\substack{\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_L: \\ \sum_{i=1}^L \|\mathbf{u}_i\|_2^2 = 1}} 2 \left(\sum_{i=1}^L \sigma_{\max} \left(\prod_{j=i+1}^L \mathbf{W}_j^* \right) \sigma_{\max} \left(\prod_{j=1}^{i-1} \mathbf{W}_j^* \right) \|\mathbf{u}_i\|_2 \right)^2, \quad (65)$$

where $\mathbf{U}_i \in \mathbb{R}^{d_i \times d_{i-1}}$ and $\mathbf{u}_i \in \mathbb{R}^{d_i d_{i-1}}$. We can upper bound the right-hand side of (61) by using the *triangle inequality*. By applying the *vectorization trick* of the Kronecker product, we can rewrite (62). Then, noting the fact that for any matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$ and vector $\mathbf{x} \in \mathbb{R}^n$, $\|\mathbf{A}\mathbf{x}\|_2 \leq \sigma_{\max}(\mathbf{A})\|\mathbf{x}\|_2$, we can upper bound the right-hand side of (63). Note that for any matrix \mathbf{A} and \mathbf{B} , $\sigma_{\max}(\mathbf{A} \otimes \mathbf{B}) = \sigma_{\max}(\mathbf{A})\sigma_{\max}(\mathbf{B})$. Hence, we can rewrite (64). Then, by using the Cauchy–Schwarz inequality,

$$\lambda_{\max}(\nabla^2 \mathcal{L}(\mathbf{w}^*)) \leq 2 \sum_{i=1}^L \sigma_{\max} \left(\prod_{j=i+1}^L \mathbf{W}_j^* \right)^2 \sigma_{\max} \left(\prod_{j=1}^{i-1} \mathbf{W}_j^* \right)^2. \quad (66)$$

Then, it suffices to show that there exists a direction $[\mathbf{U}_1^*, \mathbf{U}_2^*, \dots, \mathbf{U}_L^*]$ along which the bound in (66) is achieved.

Consider decomposition of $\prod_{j=i+1}^L \mathbf{W}_j^*$ by SVD, and denote by \mathbf{u}_{L_i} and \mathbf{v}_{L_i} the left and right singular vectors of $\prod_{j=i+1}^L \mathbf{W}_j^*$ corresponding to the largest singular value, respectively. Note that since \mathbf{W}_L is a vector, we have $\mathbf{u}_{L_i} = 1$ for all $i \in [L]$. Moreover, decompose $\prod_{j=1}^{i-1} \mathbf{W}_j^*$ by SVD, and denote by \mathbf{u}_{R_i} and \mathbf{v}_{R_i} the left and right singular vectors of $\prod_{j=1}^{i-1} \mathbf{W}_j^*$ corresponding to the largest singular value, respectively. Note that since \mathbf{W}_1 is a vector, we have $\mathbf{v}_{R_i} = 1$ for all $i \in [L]$. Now, we determine a particular direction $[\mathbf{U}_1^*, \mathbf{U}_2^*, \dots, \mathbf{U}_L^*]$ such that they achieve the upper bound while satisfying the constraint $\sum_{i=1}^L \|\mathbf{U}_i^*\|_F^2 = 1$. Choose

$$\mathbf{U}_i^* = \frac{\sigma_{\max} \left(\prod_{j=i+1}^L \mathbf{W}_j^* \right) \sigma_{\max} \left(\prod_{j=1}^{i-1} \mathbf{W}_j^* \right)}{\sqrt{\sum_{i=1}^L \sigma_{\max} \left(\prod_{j=i+1}^L \mathbf{W}_j^* \right)^2 \sigma_{\max} \left(\prod_{j=1}^{i-1} \mathbf{W}_j^* \right)^2}} \mathbf{v}_{L_i} \mathbf{u}_{R_i}^\top. \quad (67)$$

Then,

$$2 \left\| \sum_{i=1}^L \left[\left(\prod_{j=i+1}^L \mathbf{W}_j^* \right) \mathbf{U}_i^* \left(\prod_{j=1}^{i-1} \mathbf{W}_j^* \right) \right] \right\|_F^2 = 2 \sum_{i=1}^L \sigma_{\max} \left(\prod_{j=i+1}^L \mathbf{W}_j^* \right)^2 \sigma_{\max} \left(\prod_{j=1}^{i-1} \mathbf{W}_j^* \right)^2. \quad (68)$$

Since the upper bound is achieved, it implies

$$\lambda_{\max}(\nabla^2 \mathcal{L}(\mathbf{w}^*)) = 2 \sum_{i=1}^L \sigma_{\max} \left(\prod_{j=i+1}^L \mathbf{W}_j^* \right)^2 \sigma_{\max} \left(\prod_{j=1}^{i-1} \mathbf{W}_j^* \right)^2. \quad (69)$$

□

810 C PROOFS FROM SECTION 4

811 C.1 PROOF OF THEOREM 5

812 *Proof.* By definition,

$$813 \lambda_{\max}(\nabla^2 \mathcal{L}(\mathbf{w}^*)) = \max_{\substack{\mathbf{U}_1, \mathbf{U}_2, \dots, \mathbf{U}_L: \\ \sum_{i=1}^L \|\mathbf{U}_i\|_F^2 = 1}} 2 \left\| \sum_{i=1}^L \left[\left(\prod_{j=i+1}^L \mathbf{W}_j^* \right) \mathbf{U}_i \left(\prod_{j=1}^{i-1} \mathbf{W}_j^* \right) \right] \right\|_F^2 \quad (70)$$

$$814 = \max_{\substack{\mathbf{U}_1, \mathbf{U}_2, \dots, \mathbf{U}_L: \\ \sum_{i=1}^L \|\mathbf{U}_i\|_F^2 = 1}} 2 \left\| \text{vec} \left(\sum_{i=1}^L \left[\left(\prod_{j=i+1}^L \mathbf{W}_j^* \right) \mathbf{U}_i \left(\prod_{j=1}^{i-1} \mathbf{W}_j^* \right) \right] \right) \right\|_2^2 \quad (71)$$

$$815 = \max_{\substack{\mathbf{U}_1, \mathbf{U}_2, \dots, \mathbf{U}_L: \\ \sum_{i=1}^L \|\mathbf{U}_i\|_F^2 = 1}} 2 \left\| \sum_{i=1}^L \text{vec} \left(\left(\prod_{j=i+1}^L \mathbf{W}_j^* \right) \mathbf{U}_i \left(\prod_{j=1}^{i-1} \mathbf{W}_j^* \right) \right) \right\|_2^2 \quad (72)$$

$$816 = \max_{\substack{\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_L: \\ \sum_{i=1}^L \|\mathbf{u}_i\|_2^2 = 1}} 2 \left\| \sum_{i=1}^L \left[\left(\prod_{j=1}^{i-1} \mathbf{W}_j^* \right)^\top \otimes \left(\prod_{j=i+1}^L \mathbf{W}_j^* \right) \right] \mathbf{u}_i \right\|_2^2. \quad (73)$$

817 Note that vec is a linear operator. Therefore, (71) can be rewritten as (72). Then, by using the
818 vectorization trick of the Kronecker product, we can obtain (73). Let's define a block matrix and a
819 vector such that

$$820 \mathbf{K} = \left[\mathbf{I} \otimes \prod_{j=2}^L \mathbf{W}_j^* | \mathbf{W}_1^{*\top} \otimes \left(\prod_{j=3}^L \mathbf{W}_j^* \right) | \dots | \left(\prod_{j=1}^{L-1} \mathbf{W}_j^* \right)^\top \otimes \mathbf{I} \right], \quad (74)$$

$$821 \mathbf{u} = [\mathbf{u}_1^\top \quad \mathbf{u}_2^\top \quad \dots \quad \mathbf{u}_L^\top]^\top. \quad (75)$$

822 Then,

$$823 \lambda_{\max}(\nabla^2 \mathcal{L}(\mathbf{w}^*)) = \max_{\mathbf{u}: \|\mathbf{u}\|_2 = 1} 2 \|\mathbf{K} \mathbf{u}\|_2^2 \quad (76)$$

$$824 = \sigma_{\max}(\mathbf{K}^\top \mathbf{K}). \quad (77)$$

825 Note that $\sigma_{\max}(\mathbf{K}^\top \mathbf{K}) = \sigma_{\max}(\mathbf{K} \mathbf{K}^\top)$, and for any two block matrices \mathbf{A} and \mathbf{B} such that

$$826 \mathbf{A} = [\mathbf{A}_1 \quad \mathbf{A}_2 \quad \dots \quad \mathbf{A}_L] \in \mathbb{R}^{M_1 \times d}, \quad \mathbf{B} = \begin{bmatrix} \mathbf{B}_1 \\ \vdots \\ \mathbf{B}_L \end{bmatrix} \in \mathbb{R}^{d \times M_2} \quad (78)$$

$$827 \mathbf{A} \mathbf{B} = \sum_{i=1}^L \mathbf{A}_i \mathbf{B}_i, \quad \mathbf{A} \mathbf{B} \in \mathbb{R}^{M_1 \times M_2}. \quad (79)$$

828 Furthermore, for any matrices $\mathbf{A}, \mathbf{B}, \mathbf{C}, \mathbf{D}$ such that the matrix products $\mathbf{A} \mathbf{B}$ and $\mathbf{C} \mathbf{D}$ are well
829 defined, we have

$$830 (\mathbf{A} \otimes \mathbf{C})(\mathbf{B} \otimes \mathbf{D}) = \mathbf{A} \mathbf{B} \otimes \mathbf{C} \mathbf{D}. \quad (80)$$

831 Using the fact that $(\mathbf{A} \otimes \mathbf{C})^\top = \mathbf{A}^\top \otimes \mathbf{C}^\top$ together with the previous property, it follows that

$$832 \lambda_{\max}(\nabla^2 \mathcal{L}(\mathbf{w}^*)) = 2 \sigma_{\max} \left(\sum_{i=1}^L \mathbf{B}_i^\top \mathbf{B}_i \otimes \mathbf{A}_i \mathbf{A}_i^\top \right), \quad (81)$$

833 where $\mathbf{A}_k = \prod_{i=k+1}^L \mathbf{W}_i^*$ and $\mathbf{B}_k = \prod_{i=1}^{k-1} \mathbf{W}_i^*$.

□

864 C.2 PROOF OF COROLLARY 6

865 *Proof.* According to (14), for any $(\mathbf{L}^*, \mathbf{R}^*) \in \Omega$,

$$866 \quad \nabla^2 \mathcal{L}(\mathbf{L}^*, \mathbf{R}^*)[\mathbf{U}, \mathbf{V}] = 2 \|\mathbf{L}^* \mathbf{U}^\top + \mathbf{V} \mathbf{R}^{*\top}\|_F^2. \quad (82)$$

867 Then,

$$870 \quad \lambda_{\max}(\nabla^2 \mathcal{L}(\mathbf{L}^*, \mathbf{R}^*)) = \max_{\substack{\mathbf{U}, \mathbf{V} \\ \|\mathbf{U}\|_F^2 + \|\mathbf{V}\|_F^2 = 1}} 2 \|\mathbf{L}^* \mathbf{U}^\top + \mathbf{V} \mathbf{R}^{*\top}\|_F^2 \quad (83)$$

$$871 \quad \leq \max_{\substack{\mathbf{U}, \mathbf{V} \\ \|\mathbf{U}\|_F^2 + \|\mathbf{V}\|_F^2 = 1}} 2 (\|\mathbf{L}^* \mathbf{U}^\top\|_F + \|\mathbf{V} \mathbf{R}^{*\top}\|_F)^2 \quad (84)$$

$$872 \quad = \max_{\substack{\mathbf{u}, \mathbf{v} \\ \|\mathbf{u}\|_2^2 + \|\mathbf{v}\|_2^2 = 1}} 2 \left(\|(\mathbf{I} \otimes \mathbf{L}^*) \mathbf{u} \|_2 + \| (\mathbf{R}^* \otimes \mathbf{I}) \mathbf{v} \|_2 \right)^2 \quad (85)$$

$$873 \quad \leq \max_{\substack{\mathbf{u}, \mathbf{v} \\ \|\mathbf{u}\|_2^2 + \|\mathbf{v}\|_2^2 = 1}} 2 \left(\sigma_{\max}(\mathbf{I} \otimes \mathbf{L}^*) \|\mathbf{u}\|_2 + \sigma_{\max}(\mathbf{R}^* \otimes \mathbf{I}) \|\mathbf{v}\|_2 \right)^2. \quad (86)$$

874 We can upper bound the right-hand side of (83) using the *triangle inequality*. By applying the *vectorization trick* of the Kronecker product again, we can rewrite (84). Then, noting that for any matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$ and vector $\mathbf{x} \in \mathbb{R}^n$, $\|\mathbf{A}\mathbf{x}\|_2 \leq \sigma_{\max}(\mathbf{A})\|\mathbf{x}\|_2$, we can upper bound the right-hand side of (85). Note that for any matrix \mathbf{A} and \mathbf{B} , $\sigma_{\max}(\mathbf{A} \otimes \mathbf{B}) = \sigma_{\max}(\mathbf{A})\sigma_{\max}(\mathbf{B})$. Hence,

$$875 \quad \lambda_{\max}(\nabla^2 \mathcal{L}(\mathbf{L}^*, \mathbf{R}^*)) \leq \max_{\substack{\mathbf{u}, \mathbf{v} \\ \|\mathbf{u}\|_2^2 + \|\mathbf{v}\|_2^2 = 1}} 2 (\sigma_{\max}(\mathbf{L}^*) \|\mathbf{u}\|_2 + \sigma_{\max}(\mathbf{R}^*) \|\mathbf{v}\|_2)^2. \quad (87)$$

876 Then, by using Cauchy-Schwarz inequality,

$$877 \quad \lambda_{\max}(\nabla^2 \mathcal{L}(\mathbf{L}^*, \mathbf{R}^*)) \leq 2 (\sigma_{\max}(\mathbf{L}^*)^2 + \sigma_{\max}(\mathbf{R}^*)^2). \quad (88)$$

878 Now, we will show that this upper bound is achievable. Let us decompose \mathbf{L}^* as $\mathbf{U}_L \Sigma_L \mathbf{V}_L^\top$ by SVD, and denote by \mathbf{u}_L and \mathbf{v}_L the left and right singular vectors corresponding to the largest singular value, respectively. Moreover, decompose \mathbf{R}^* as $\mathbf{U}_R \Sigma_R \mathbf{V}_R^\top$ by SVD, and denote by \mathbf{u}_R and \mathbf{v}_R the left and right singular vectors corresponding to the largest singular value, respectively. We determine a particular $(\mathbf{U}^*, \mathbf{V}^*)$ such that it achieves the upper bound while satisfying the constraint $\|\mathbf{U}^*\|_F^2 + \|\mathbf{V}^*\|_F^2 = 1$. Choose

$$879 \quad \mathbf{U}^{*\top} = \frac{\sigma_{\max}(\mathbf{L}^*)}{\sqrt{\sigma_{\max}(\mathbf{L}^*)^2 + \sigma_{\max}(\mathbf{R}^*)^2}} \mathbf{v}_L \mathbf{u}_R^\top, \quad \mathbf{V}^* = \frac{\sigma_{\max}(\mathbf{R}^*)}{\sqrt{\sigma_{\max}(\mathbf{L}^*)^2 + \sigma_{\max}(\mathbf{R}^*)^2}} \mathbf{u}_L \mathbf{v}_R^\top. \quad (89)$$

880 Using the fact that, for any vectors \mathbf{x} and \mathbf{y} $\|\mathbf{x}\mathbf{y}^\top\|_F^2 = \|\mathbf{x}\|_2^2 \|\mathbf{y}\|_2^2$,

$$881 \quad 2 \left\| \frac{(\sigma_{\max}(\mathbf{L}^*)^2 + \sigma_{\max}(\mathbf{R}^*)^2)}{\sqrt{\sigma_{\max}(\mathbf{L}^*)^2 + \sigma_{\max}(\mathbf{R}^*)^2}} \mathbf{u}_L \mathbf{u}_R^\top \right\|_F^2 = 2 (\sigma_{\max}(\mathbf{L}^*)^2 + \sigma_{\max}(\mathbf{R}^*)^2). \quad (90)$$

882 □

883 D ADDITIONAL EXPERIMENTAL RESULTS

884 D.1 VISUALIZATION OF THE CONTOUR MAP OF THE LOSS LANDSCAPE

885 To study the dynamics of deep matrix factorization, we analyze the trajectories of GD. Previous works have visualized neural network loss landscapes to explore their highly non-convex and non-Euclidean structure (Goodfellow et al., 2014; Li et al., 2018). However, the high-dimensionality prevents full visualization. As a result, only 1-D (line) or 2-D (surface) visualizations are available. In this paper, we focus on contour maps of the loss landscape in the vicinity of a global minimum and a methodology employed in prior studies to generate them.

Contour Plots with Random Projections. We want to visualize the loss landscape around a global minimum $\mathbf{w}^* \in \mathbb{R}^N$. We select two random vectors, ζ and γ , from \mathbb{R}^N . Then, for any $K \subset \mathbb{R}^2$, we can define the function $p : K \rightarrow \mathbb{R}$:

$$p(x, y) = \mathcal{L}(\mathbf{w}^* + x\zeta + y\gamma), \quad \forall (x, y) \in K, \quad (91)$$

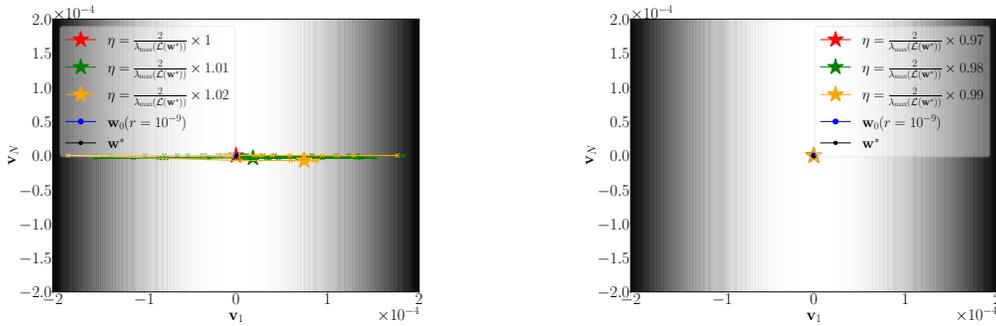
and plot p with the desired resolution.

Scale Invariance and Manifolds. Note that our loss function is *scale-invariant*, which means that for any nonzero scalar $c \in \mathbb{R}$, multiplying one layer by c and the next layer by $1/c$, or vice versa, yields the same end-to-end function. This phenomenon forms a manifold for global minimizers in the loss landscape (Dinh et al., 2017). Furthermore, we know that $\nabla^2 \mathcal{L}(\mathbf{w}^*)$ is rank-deficient by at least the order of $1 - 1/L$; that is, at least $1 - 1/L$ of the eigenvalues values of $\nabla^2 \mathcal{L}(\mathbf{w}^*)$ are zero (Mulayoff & Michaeli, 2020). This means that the ratio of the manifold dimension to the ambient space dimension increases as L grows.

Projection onto the Hessian Eigenvectors. If we use random projections in visualizations, plots might not be informative to track the optimization dynamics of GD due to the phenomenon caused by the scale invariance. To make contour maps as informative as possible, we choose ζ and γ to be \mathbf{v}_1 and \mathbf{v}_N , respectively — the eigenvectors of the largest and smallest eigenvalues of $\nabla^2 \mathcal{L}(\mathbf{w}^*)$.

D.2 EXPERIMENT DETAILS AND ADDITIONAL EXPERIMENTS

For the experiment, we first generate the layer dimensions randomly and then construct the optimal layers $[\mathbf{W}_1^*, \mathbf{W}_2^*, \dots, \mathbf{W}_L^*]$ as Gaussian random matrices, with each entry sampled from $N(0, 1)$ according to the generated dimensions. Then, we compute M or m by $\prod_{j=1}^L \mathbf{W}_j^*$. We then perform the same experiments as in Figs. 1–2, varying the depth, dimensions, and initialization distance r (as shown in Figs. 3-7). We note that oscillations occur along the eigenvector corresponding to the maximum eigenvalue of the Hessian. The dimensions of the factors, i.e., d_0, d_1, \dots, d_L , in Fig. 6 are given by 1, 9, 4, 8, 24, 16, 17, 11, 21, 3, 22, 3, 3, 15, 3, 18, 17, 16, 5, 12, 1, which implies $N = 2421$, while the dimensions of the factors in Fig. 7 are given by 1, 9, 4, 8, 1, which implies $N = 293$.

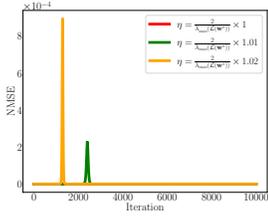


(a) Trajectories of GD initialized at \mathbf{w}_0 with step sizes $\geq 2/\lambda_{\max}$ are depicted by colored lines, and their corresponding convergence points are marked by colored \star symbols.

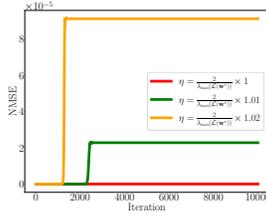
(b) Trajectories of GD initialized at \mathbf{w}_0 with step sizes $< 2/\lambda_{\max}$ are depicted by colored lines, and their corresponding convergence points are marked by colored \star symbols.

Figure 3: Contour map of the loss landscape around a minimum in a 15-layer overparameterized scalar factorization of a random scalar. GD with different step sizes η , indicated by different colors, is initialized within a radius of 10^{-9} from the minimum, in the direction of the Hessian eigenvector corresponding to the largest eigenvalue. The vector \mathbf{v}_1 denotes the eigenvector of the Hessian corresponding to the largest eigenvalue, while \mathbf{v}_N denotes the eigenvector corresponding to the smallest eigenvalue. The value of $\lambda_{\max}(\nabla^2 \mathcal{L}(\mathbf{w}^*))$ is computed using the closed-form expression derived in Theorem 4.

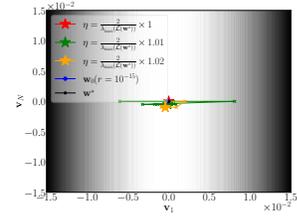
972
973
974
975
976
977
978
979



(a) Normalized training error across iterations, i.e., $\mathcal{L}(\mathbf{w}_k)/\|\mathbf{M}\|_F^2$.

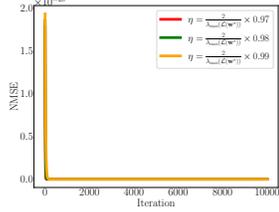


(b) Normalized ℓ^2 distance of \mathbf{w}_k from the minimum \mathbf{w}^* , i.e., $\|\mathbf{w}_k - \mathbf{w}^*\|_2^2/\|\mathbf{w}^*\|_2^2$.

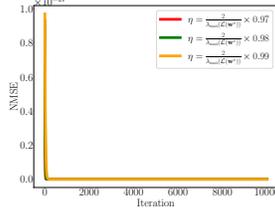


(c) Trajectories of GD on the contour map of the loss landscape around the minimum.

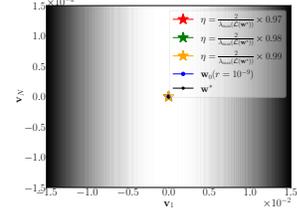
984
985
986
987
988
989
990



(d) Normalized training error across iterations, i.e., $\mathcal{L}(\mathbf{w}_k)/\|\mathbf{M}\|_F^2$.



(e) Normalized ℓ^2 distance of \mathbf{w}_k from the minimum \mathbf{w}^* , i.e., $\|\mathbf{w}_k - \mathbf{w}^*\|_2^2/\|\mathbf{w}^*\|_2^2$.

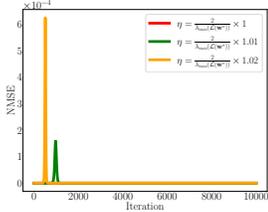


(f) Trajectories of GD on the contour map of the loss landscape around the minimum.

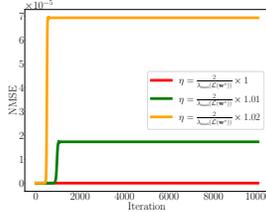
995
996
997
998

Figure 4: GD dynamics with different step sizes indicated by different colors for general matrix factorization, $\mathbf{M} = \mathbf{L}\mathbf{R}^\top$, of a random Gaussian matrix, where $\mathbf{L} \in \mathbb{R}^{10 \times 30}$ and $\mathbf{R} \in \mathbb{R}^{20 \times 30}$. The value of $\lambda_{\max}(\nabla^2 \mathcal{L}(\mathbf{w}^*))$ is computed using the closed-form expression derived in Corollary 6.

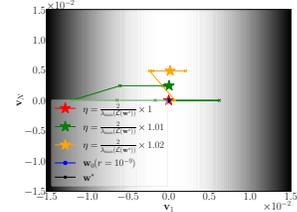
999
1000
1001
1002
1003
1004
1005



(a) Normalized training error across iterations, i.e., $\mathcal{L}(\mathbf{w}_k)/\|\mathbf{M}\|_F^2$.

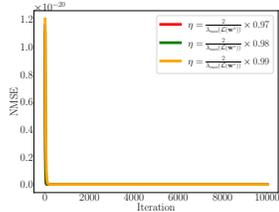


(b) Normalized ℓ^2 distance of \mathbf{w}_k from the minimum \mathbf{w}^* , i.e., $\|\mathbf{w}_k - \mathbf{w}^*\|_2^2/\|\mathbf{w}^*\|_2^2$.

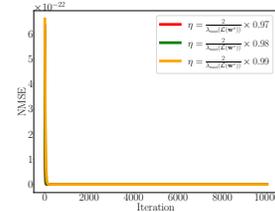


(c) Trajectories of GD on the contour map of the loss landscape around the minimum.

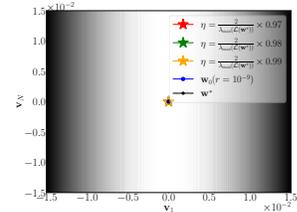
1006
1007
1008
1009
1010
1011
1012
1013
1014
1015
1016
1017



(d) Normalized training error across iterations, i.e., $\mathcal{L}(\mathbf{w}_k)/\|\mathbf{M}\|_F^2$.



(e) Normalized ℓ^2 distance of \mathbf{w}_k from the minimum \mathbf{w}^* , i.e., $\|\mathbf{w}_k - \mathbf{w}^*\|_2^2/\|\mathbf{w}^*\|_2^2$.

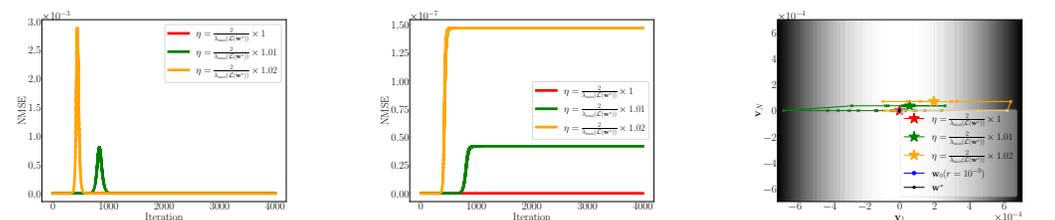


(f) Trajectories of GD on the contour map of the loss landscape around the minimum.

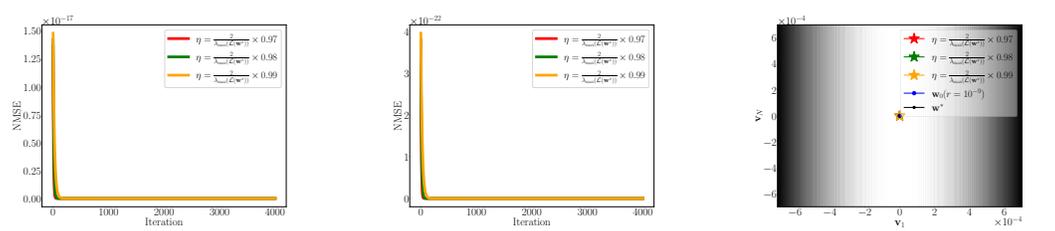
1021
1022
1023
1024
1025

Figure 5: GD dynamics with different step sizes indicated by different colors for general matrix factorization, $\mathbf{M} = \mathbf{L}\mathbf{R}^\top$, of a random Gaussian matrix, where $\mathbf{L} \in \mathbb{R}^{25 \times 30}$ and $\mathbf{R} \in \mathbb{R}^{20 \times 30}$. The value of $\lambda_{\max}(\nabla^2 \mathcal{L}(\mathbf{w}^*))$ is computed using the closed-form expression derived in Corollary 6.

1026
 1027
 1028
 1029
 1030
 1031
 1032
 1033
 1034
 1035
 1036
 1037
 1038
 1039
 1040
 1041
 1042
 1043
 1044
 1045
 1046
 1047
 1048
 1049
 1050
 1051
 1052
 1053
 1054
 1055
 1056
 1057
 1058
 1059
 1060
 1061
 1062
 1063
 1064
 1065
 1066
 1067
 1068
 1069
 1070
 1071
 1072
 1073
 1074
 1075
 1076
 1077
 1078
 1079

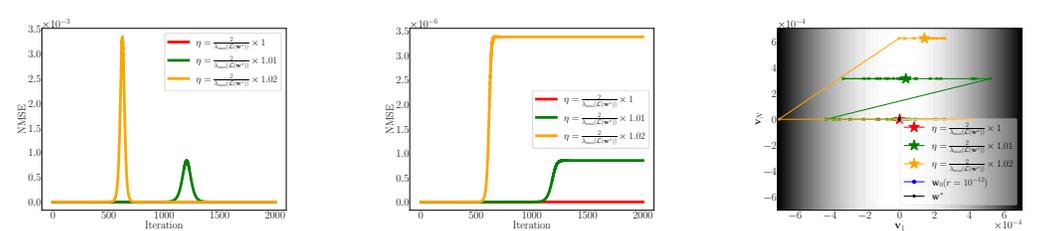


(a) Normalized training error across iterations, i.e., $\mathcal{L}(\mathbf{w}_k)/m^2$.
 (b) Normalized ℓ^2 distance of \mathbf{w}_k from the minimum \mathbf{w}^* , i.e., $\|\mathbf{w}_k - \mathbf{w}^*\|_2^2/\|\mathbf{w}^*\|_2^2$.
 (c) Trajectories of GD on the contour map of the loss landscape around the minimum.

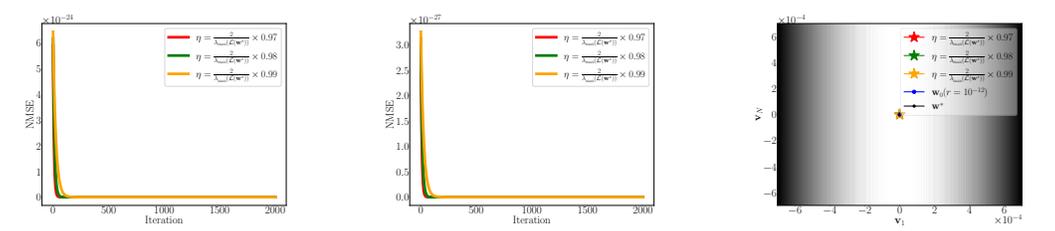


(d) Normalized training error across iterations, i.e., $\mathcal{L}(\mathbf{w}_k)/m^2$.
 (e) Normalized ℓ^2 distance of \mathbf{w}_k from the minimum \mathbf{w}^* , i.e., $\|\mathbf{w}_k - \mathbf{w}^*\|_2^2/\|\mathbf{w}^*\|_2^2$.
 (f) Trajectories of GD on the contour map of the loss landscape around the minimum.

Figure 6: GD dynamics with different step sizes indicated by different colors are initialized within a radius of 10^{-9} from the minimum in the direction of the Hessian eigenvector corresponding to the largest eigenvalue, for a 20-layer overparameterized scalar factorization of a random scalar. The value of $\lambda_{\max}(\nabla^2\mathcal{L}(\mathbf{w}^*))$ is computed using the closed-form expression derived in Theorem 4.



(a) Normalized training error across iterations, i.e., $\mathcal{L}(\mathbf{w}_k)/m^2$.
 (b) Normalized ℓ^2 distance of \mathbf{w}_k from the minimum \mathbf{w}^* , i.e., $\|\mathbf{w}_k - \mathbf{w}^*\|_2^2/\|\mathbf{w}^*\|_2^2$.
 (c) Trajectories of GD on the contour map of the loss landscape around the minimum.



(d) Normalized training error across iterations, i.e., $\mathcal{L}(\mathbf{w}_k)/m^2$.
 (e) Normalized ℓ^2 distance of \mathbf{w}_k from the minimum \mathbf{w}^* , i.e., $\|\mathbf{w}_k - \mathbf{w}^*\|_2^2/\|\mathbf{w}^*\|_2^2$.
 (f) Trajectories of GD on the contour map of the loss landscape around the minimum.

Figure 7: GD dynamics with different step sizes indicated by different colors are initialized within a radius of 10^{-12} from the minimum in the direction of the Hessian eigenvector corresponding to the largest eigenvalue, for a 5-layer overparameterized scalar factorization of a random scalar. The value of $\lambda_{\max}(\nabla^2\mathcal{L}(\mathbf{w}^*))$ is computed using the closed-form expression derived in Theorem 4.