# ON THE STABILITY OF NONLINEAR DYNAMICS IN GD AND SGD: BEYOND QUADRATIC POTENTIALS

## **Anonymous authors**

Paper under double-blind review

### **ABSTRACT**

The dynamical stability of the iterates during training plays a key role in determining the minima obtained by training algorithms. For example, stable solutions of gradient descent (GD) correspond to flat minima, and these have been associated with favorable features. While prior work often relies on linearization to determine stability, it remains unclear whether linearized dynamics faithfully capture the full nonlinear behavior. In this work, we explicitly study the effect of nonlinear terms. For GD, we show that linear analysis can be misleading. The iterates may stably oscillate near a linearly unstable minimum, and still converge once the step size decays. Here, we derive an exact condition for such stable oscillations, which depends on higher-order derivatives of the loss. Extending the analysis to stochastic gradient descent (SGD), we demonstrate that nonlinear dynamics can diverge in expectation if even a single batch is unstable. This implies that stability can be dictated by the worst-case batch, rather than an average effect, as linear analysis suggests. Finally, we prove that if all batches are linearly stable, then the nonlinear dynamics of SGD are stable in expectation.

## 1 Introduction

Understanding the nature of the minima reached by our training procedures is a central problem in machine learning and optimization (Neyshabur et al., 2014). A common way to investigate this issue is by analyzing the stability of the iterates as the algorithm approaches a minimum (Wu et al., 2018). For example, it has been shown that stable minimizers of gradient descent (GD) correspond to flat minima (Cohen et al., 2021), and those have been associated with flat predictor functions (Mulayoff et al., 2021; Nacson et al., 2023) and balanced networks (Mulayoff & Michaeli, 2020). These highlight the role of dynamical stability in shaping the properties of the obtained solutions.

In dynamical systems theory, stability analysis is often carried out via linearization. Once the iterates arrive at the vicinity of a fixed point, it is often sufficient to study the linearized system in order to determine whether convergence occurs (Thompson & Stewart, 2002). This technique has been widely applied to study GD (Cohen et al., 2021). Extending it to the stochastic regime, Wu et al. (2018) proposed using linearization to analyze the stability of stochastic gradient descent (SGD) in the mean-square sense.

This approach has inspired a large body of subsequent research. In particular, Ma & Ying (2021) demonstrated that the moments of the linearized dynamics evolve independently, and for the second moment (mean squared error), they provided an implicit expression for the exact stability criterion. Building on this result, Mulayoff & Michaeli (2024) derived an explicit form of the condition, yielding new insights into the linear stability of SGD. Importantly, the stability threshold of the step size depends on the curvature of all samples in the training set (see App. B). Despite this progress, however, it remains unclear whether—and under what conditions—the behavior of linearized iterates truly reflects the full nonlinear dynamics of SGD.

In this work, we aim to address this gap by studying the effect of nonlinear terms. We begin with the deterministic case of GD, as we observe that linearized dynamics can be misleading. GD's iterates may stably oscillate near a linearly unstable minimum and, after step size decay, eventually converge to it. This indicates that oscillations must be taken into account while considering minima stability. Such oscillations correspond to a flip (period doubling) bifurcation of the GD map, in which the

iterates oscillate along the sharpest direction of the minimum. The stability of this bifurcation is governed by the first Lyapunov coefficient in its normal form (see Sec. 5.1). Using this understanding, we derive a precise criterion for stable oscillations. This condition depends on the third- and fourth-order derivatives of the loss at the minimum (see Thm. 1).

We then extend our analysis to the stochastic setting of SGD. Following prior work, we focus on interpolating minima and assume the loss functions are analytic in a neighborhood of the minimum. In this setting, linearized dynamics, combined with mean-square analysis, suggest that the stability threshold of SGD depends on an average curvature over the different mini-batches. In contrast, we show that if the iterates are unstable even with respect to a single batch, the full nonlinear dynamics of SGD are unstable in expectation. This result suggests that stability is determined by the worst-case batch, contradicting prior assumptions about the averaging effect of stochasticity (see Thm. 2).

Finally, we provide a sufficient condition for the stability of SGD. Specifically, we prove that if the dynamics are linearly stable with respect to all possible batches, then there exists a neighborhood of the minimum from which the full nonlinear dynamics converge in expectation (see Thm. 3). Our analysis uses Koopman theory (Koopman, 1931), which allows us to formulate the finite-dimensional nonlinear dynamics as a linear dynamical system in an infinite-dimensional Hilbert space. This reformulation yields two notable benefits. First, nonlinear dynamics are reduced to linear ones, which are significantly more tractable. Second, the transformation provides a deterministic linear relation between the moments of the dynamics. Then, we use tools from functional analysis to derive the result. Considering our earlier findings, we see that this sufficient condition can also be necessary in certain cases, as we demonstrate in Sec. 2.2. Specifically, when unstable oscillations arise in batches with low stability thresholds. Since it takes only one such batch, and the number of possible batches is exponentially large, it is quite likely that SGD operates in this regime.

## 2 Warmup

In the following, we examine how nonlinear dynamics influence the solutions obtained by gradient-based methods. We begin with GD, showing that the iterates can stably oscillate near an unstable minimum. Combined with step size decay, this suggests that convergence to minima is often mediated by stable oscillations. We then turn to SGD, where we find that, contrary to linear predictions, stability in expectation can be dictated by the worst-case batch rather than an average effect.

#### 2.1 NORMAL FORM OF OSCILLATIONS IN GRADIENT DESCENT

A classical result states that GD with constant step size  $\eta$  converges to a minimizer in the general case only if it is linearly stable. Specifically, for  $\mathcal{L}:\mathbb{R}^d\to\mathbb{R}$  with a minimizer  $\boldsymbol{x}^*$ , the linear stability threshold is given by  $\eta_{\text{lin}}=2/\lambda_{\text{max}}(\nabla^2\mathcal{L}(\boldsymbol{x}^*))$ , and  $\boldsymbol{x}^*$  is linearly stable if and only if  $\eta<\eta_{\text{lin}}$ . Recently, it has been shown that GD typically operates at the edge of stability (EoS) when optimizing neural networks (Cohen et al., 2021). In this regime, the top eigenvalue of the Hessian hovers just above  $2/\eta$  as the parameters approach a minimum. This implies that GD often encounters linearly unstable minima. Although the algorithm cannot converge directly to such points, it can stably oscillate nearby and, after step size decay, eventually converge. Thus, the ability to endure oscillations near a minimum determines whether the algorithm will eventually converge to it.

In this section, we demonstrate that linear stability cannot predict whether GD will stably oscillate near a minimum. To illustrate this, we examine GD's iterates over two univariate functions, depicted in Fig. 1(a), that share the same curvature at the minimum but differ in higher-order terms:

$$f_{+}(x) = \frac{1}{2}x^{2} + \frac{1}{4}x^{4}$$
, and  $f_{-}(x) = \frac{1}{2}x^{2} - \frac{1}{4}x^{4}$ . (1)

Both have the same sharpness at the local minimizer  $x^* = 0$ , with  $f''_+(0) = f''_-(0) = 1$ , yielding a linear stability threshold of  $\eta_{\text{lin}} = 2$ . The iterates of GD are given by

$$x_{t+1} = -(\eta - 1)x_t \pm \eta x_t^3. (2)$$

Let us examine how the asymptotic value of the iterates depends on the step size  $\eta$ . Figures 1(b) and 1(c) plot the accumulation points of  $\{x_t\}$  for various values of  $\eta$  on  $f_+$  and  $f_-$ , where  $x_0$  is chosen at random from the interval (-1,1). For  $\eta < \eta_{\rm lin}$ , both dynamics converge to  $x^* = 0$ .

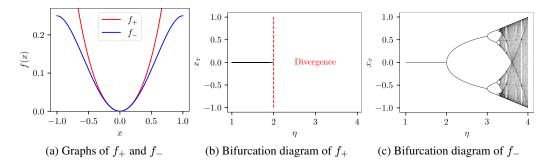


Figure 1: Stable vs. unstable oscillations near a minimum. We apply GD to  $f_+$  and  $f_-$  from (1) with various step sizes  $\eta \in (1,4)$ . The resulting dynamics (2) correspond to the normal form of a flip bifurcation. Once the step size exceeds the linear stability threshold  $\eta_{\rm lin}=2$ , stability is determined by the sign of the cubic term in the dynamics. Panel (a) shows  $f_+$  and  $f_-$ , whose minima share the same sharpness. Panel (b) visualizes GD's output on  $f_+$  with various step sizes. When the step size  $\eta$  crosses  $\eta_{\rm lin}$ , the minimum  $x^*=0$  loses stability, resulting in unstable oscillations, which lead to divergence. Panel (c) depicts GD's convergent points on  $f_-$  for various step sizes. At the threshold,  $\eta=\eta_{\rm lin}$ , the minimizer  $x^*=0$  loses stability and the iterates settle into a stable period-2 cycle, which then undergoes period doubling, chaos, and eventually divergence for  $\eta>4$ .

However, when  $\eta > \eta_{\rm lin}$ , the behavior of the dynamics differs. The iterates on  $f_+$  immediately diverge once the step size crosses  $\eta_{\rm lin}$ . In contrast, GD on  $f_-$  exhibits rich nonlinear dynamics, where it initially settles into stable cycles over a wide range of step sizes while featuring period doubling bifurcations, before transitioning into chaos, and finally diverging once  $\eta > 4$ . Importantly, when such stable oscillations occur, decaying the step size below  $\eta_{\rm lin}$  results in convergence to  $x^*$ .

This simple example demonstrates that exceeding the linear stability threshold does not necessarily imply that GD escapes the minimum. Interestingly, under mild assumptions, the behavior of any nonlinear dynamics along the critical manifold near a linearly unstable fixed point can be reduced to this simple one-dimensional map, called *normal form* (see Sec. 5.1). Then, as the example illustrates, the sign of the cubic term in this normal form can be used to determine whether stable oscillations arise. In Sec. 3 we extend this analysis to higher dimensions and derive a general condition for stable oscillations of GD at the edge of stability.

## 2.2 Worst case stability in Stochastic gradient descent

Linearized analyses of SGD in expectation predict stability by averaging curvature information across all samples (see App. B). In particular, under mean-square analysis, the distance of the iterates to a minimizer remains bounded as long as the step size  $\eta$  is below a threshold determined by an average sharpness of the loss. Here, we show that the full nonlinear dynamics behave differently. Stability in expectation may be governed by the worst-case batch rather than by an average.

To illustrate this discrepancy, we examine the dynamics of SGD on the following functions:

$$f_{+}(x) = \frac{1}{2}x^{2} + \frac{1}{4}x^{4}, \quad \text{and} \quad f_{a}(x) = \frac{a}{2}x^{2},$$
 (3)

where  $a \in (0,1)$  is a fixed parameter. More specifically, we consider the minimization of the average of  $f_+$  and  $f_a$ , where at each iteration, SGD takes a gradient step with respect to one of these functions, chosen at random. Here  $x^*=0$  is an interpolating minimizer, *i.e.*, it minimizes each function individually. The sharpness of these functions at  $x^*$ , given by their second derivative, is  $h_+=f_+''(0)=1$  and  $h_a=f_a''(0)=a$ . Consequently, the linear stability thresholds for optimizing each function separately are  $\eta_+=2/h_+=2$  and  $\eta_a=2/h_a=2/a$ . Under the linearized mean-square analysis of SGD, the combined stability threshold equals (see App. B)

$$\eta_{\text{lin}} = 2\frac{h_+ + h_a}{h_+^2 + h_a^2} = 2\frac{1+a}{1+a^2} > 2.$$
(4)

We now compare this prediction with the actual SGD dynamics. Proposition 1 shows that whenever  $\eta > 2$ , the nonlinear SGD iterates diverge in expectation (see proof in App. C).

**Proposition 1 (Worst case batch)** Let  $\{x_t\}$  be SGD's iterates on  $f_+$  and  $f_a$  from (3), s.t.  $x_0 \neq 0$ . If  $\eta > 2$  then  $\mathbb{E}[|x_t - x^*|] \xrightarrow[t \to \infty]{} \infty$ .

In other words, because one of the two losses  $(f_+)$  becomes unstable at  $\eta>2$ , the entire stochastic process diverges despite the linearized analysis predicting stability up to  $\eta_{\rm lin}>2$ . This simple example shows that nonlinear SGD can be governed by the least stable batch rather than by an average stability criterion. In Sec.4, we formalize this observation and provide general necessary and sufficient conditions for nonlinear stability of SGD.

## 3 OSCILLATIONS IN GRADIENT DESCENT

In this section, we present a general condition for stable oscillations of GD near a minimum. As noted earlier, GD typically exhibits the edge-of-stability (EoS) phenomenon when optimizing neural networks (Cohen et al., 2021). During the early stages of training, a phase called progressive sharpening, the landscape becomes sharper as the top eigenvalue of the Hessian increases until it reaches the linear stability threshold of  $2/\eta$  (Wang et al., 2022). Beyond this point, the sharpness remains slightly above  $2/\eta$  for the rest of the training. Consequently, as the iterates approach a minimum, GD often encounters minima whose sharpness marginally exceeds the linear stability threshold. While direct convergence to such minimizers is impossible, the iterates can stably oscillate in their vicinity. Then, once the step size decays, these oscillations vanish, allowing the method to settle into the minimum. Thus, understanding the behavior of GD at the edge of stability in the vicinity of minima is critical for determining to which minima it converges. Here we have the following result, which uses the kth order derivative in multilinear form, denoted as  $\mathcal{D}^k$ .

**Theorem 1 (Stable oscillations)** Let  $\mathcal{L}: \mathbb{R}^d \to \mathbb{R}$  and  $\boldsymbol{x}^*$  be its local minimizer, such that  $\mathcal{L}$  is four times differentiable at  $\boldsymbol{x}^*$ . Assume  $\nabla^2 \mathcal{L}(\boldsymbol{x}^*)$  is strictly positive and let  $\boldsymbol{v}$  be a top eigenvector corresponding to the maximal eigenvalue. Suppose GD on  $\mathcal{L}$  with step size  $\eta$  operates at the edge of stability, i.e.,  $\lambda_{\max}\left(\nabla^2 \mathcal{L}(\boldsymbol{x}^*)\right) = 2/\eta$ , and that  $\lambda_{\max}$  has multiplicity one. Then a stable period-2 cycle exists at the vicinity of  $\boldsymbol{x}^*$  if and only if

$$\mathcal{D}^{3}\mathcal{L}(\boldsymbol{x}^{*})[\boldsymbol{v},\boldsymbol{v},\boldsymbol{q}] > \mathcal{D}^{4}\mathcal{L}(\boldsymbol{x}^{*})[\boldsymbol{v},\boldsymbol{v},\boldsymbol{v},\boldsymbol{v}], \tag{5}$$

where

$$\boldsymbol{q} \triangleq \left[ \nabla^2 \mathcal{L}(\boldsymbol{x}^*) \right]^{-1} \nabla_{\boldsymbol{v}} \mathcal{D}^3 \mathcal{L}(\boldsymbol{x}^*) [\boldsymbol{v}, \boldsymbol{v}, \boldsymbol{v}]. \tag{6}$$

This theorem states that GD can stably oscillate near a minimum if and only if the condition in (5) holds. This condition is composed of high-order derivatives of the loss. Intuitively, it suggests that when the third derivative dominates over the fourth across the sharpest direction, we have stable oscillations and vice versa. The expression for  $\boldsymbol{q}$  has a Newton-like structure, where the inverse Hessian is applied to a gradient. However, this gradient acts only on the cubic term in the Taylor expansion of the loss, not on the full objective. Obviously,  $\nabla_{\boldsymbol{v}}\mathcal{D}^3\mathcal{L}(\boldsymbol{x}^*)[\boldsymbol{v},\boldsymbol{v},\boldsymbol{v}]$  equals  $3\mathcal{D}^3\mathcal{L}(\boldsymbol{x}^*)[\boldsymbol{v},\boldsymbol{v}]$ , and thus  $\boldsymbol{v}$ 's scale and polarity do not affect the condition. When the condition is satisfied, step sizes slightly above  $2/\lambda_{\max}$  produce stable periodic oscillations, whose amplitude grows with  $\eta$ , while smaller step sizes converge to the minimum. Conversely, if the condition is not met, any step size larger than  $2/\lambda_{\max}$  leads the iterates to escape the small neighborhood of the minimum. The proof outline, along with additional information about bifurcations, are given in Sec. 5.

To illustrate Thm. 1, we consider a simple example, shown in Fig. 2. Let

$$f_{\alpha}(x) = \frac{1}{2}x^2 + \frac{\alpha}{6}x^3 + \frac{1}{8}x^4,\tag{7}$$

with minimizer at  $x^*=0$ . Figure 2(a) depicts  $f_{\alpha}$  near the minimum for a few values of  $\alpha$ . The linear stability threshold of GD around  $x^*$  is  $\eta_{\text{lin}}=2$ , at which the update rule becomes

$$x_{t+1} = -x_t - \alpha x_t^2 - x_t^3. (8)$$

In this case, the condition for stable oscillations (5) simplifies to  $|\alpha|>1$  (see App. E). Figure 2(b) presents the accumulation points of the iterates for a range of  $\alpha$ . When  $|\alpha|>1$ , GD's iterates converge to stable period-2 cycles, whereas for  $|\alpha|<1$ , the iterates diverge. This example demonstrates that Thm. 1 captures the precise phase transition from stable to unstable oscillations.

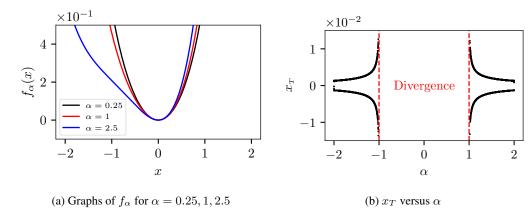


Figure 2: **Demonstration of Thm. 1.** Consider  $f_{\alpha}(x) = \frac{1}{2}x^2 + \frac{\alpha}{6}x^3 + \frac{1}{8}x^4$ , whose linear stability threshold under GD is  $\eta_{\text{lin}} = 2$ . At this step size, the update rule becomes  $x_{t+1} = -x_t - \alpha x_t^2 - x_t^3$ . According to Thm. 1, GD oscillates stably around the minimum  $x^* = 0$  if and only if  $|\alpha| > 1$  (see App. E). Panel (a) plots  $f_{\alpha}$  near  $x^*$  for three choices of  $\alpha$ , highlighting the asymmetry introduced by the cubic term. Panel (b) shows the long-term value  $x_T$  across a range of  $\alpha$ . When  $|\alpha| > 1$ , GD converges to a stable period-2 cycle, whereas for  $|\alpha| < 1$  the iterates diverge. This confirms that condition (5) precisely captures the transition from stability to instability.

## 4 STABILITY OF NONLINEAR DYNAMICS IN SGD

In this section, we present our results on the stability of nonlinear dynamics in SGD. Let  $f_i : \mathbb{R}^d \to \mathbb{R}$  be analytic for all  $i \in [n]$ . We define the loss function and its batch approximation as

$$\mathcal{L}(\boldsymbol{x}) = \frac{1}{n} \sum_{i=1}^{n} f_i(\boldsymbol{x}), \quad \text{and} \quad \hat{\mathcal{L}}_{\mathcal{B}}(\boldsymbol{x}) = \frac{1}{B} \sum_{i \in \mathcal{B}} f_i(\boldsymbol{x}),$$
 (9)

where  $\mathcal{B} \subseteq [n]$  is a batch (set) of size  $|\mathcal{B}| = B$ . The iterates of SGD are given by

$$\mathbf{x}_{t+1} = \mathbf{x}_t - \eta \nabla \hat{\mathcal{L}}_{\mathcal{B}_t}(\mathbf{x}_t). \tag{10}$$

Here,  $\mathcal{B}_t$  refers to a stochastic batch sampled at iteration t. We assume that the batches  $\{\mathcal{B}_t\}$  are drawn without replacement, independently across iterations. Namely, there are distinct samples within each batch and possible repetitions between different batches.

Our analysis focuses on the dynamics of SGD near interpolating minimizers. This setting has been extensively studied by prior work, particularly in the context of dynamical stability and overparameterized models (Wu et al., 2018; Ma & Ying, 2021; Mulayoff & Michaeli, 2024).

**Definition 1 (Interpolating minimizer)** We say that 
$$\mathbf{x}^* \in \mathbb{R}^d$$
 is an interpolating minimizer of  $\mathcal{L}$  if  $\forall i \in [n]$   $\nabla f_i(\mathbf{x}^*) = \mathbf{0}$  and  $\nabla^2 f_i(\mathbf{x}^*) \succ \mathbf{0}$  (positive definite). (11)

To gain intuition about the stability of SGD near interpolating minimizers, it is useful to examine the dynamics of the iterates across all possible batches. Concretely, consider running GD separately on every batch. For a given step size, some batches may converge, while others may exhibit stable oscillations or even diverge. To capture the average behavior of the algorithm, we adopt the notion of stability in expectation (Ma & Ying, 2021). A popular instance of this approach is the mean-square (Wu et al., 2018), whose stability threshold in the linear setting aggregates curvature information from all samples (Mulayoff & Michaeli, 2024). However, as shown in Sec. 2, nonlinear dynamics behave differently. Instability of even a single batch can be enough to cause the mean to diverge.

**Theorem 2 (Necessary condition)** Let  $x^*$  be an interpolating minimizer of  $\mathcal{L}$ ,  $x_0 \in \mathbb{R}^d$ , and  $\mathcal{B}_*$  be a batch of size B. Denote GD's iterates with step size  $\eta$  over  $\hat{\mathcal{L}}_{\mathcal{B}_*}$  by  $x_t^{(\mathcal{B}_*)}$ . If

$$\sqrt[t]{\left\|\boldsymbol{x}_{t}^{(\mathcal{B}_{*})}-\boldsymbol{x}^{*}\right\|} \underset{t\to\infty}{\longrightarrow} \infty, \tag{12}$$

then SGD's iterates  $\{x_t\}$  of (10) with step size  $\eta$  diverge in expectation, i.e.,  $\mathbb{E}[\|x_t - x^*\|] \xrightarrow[t \to \infty]{} \infty$ .

In simple terms, this theorem states that if GD on even a single batch diverges at a rate higher than linear, then SGD as a whole will also diverge in expectation (see proof in App. F). Section 2.1 provides a concrete example, where GD's iterates on the function  $f_+$  in (1) diverge superlinearly (see App. C). Consequently, if the finite-sum loss  $\mathcal{L}$  contains a batch loss  $\hat{\mathcal{L}}_B = f_+$ , then SGD will diverge in expectation. This is the underlying principle behind the observation in Sec. 2.2.

What can we learn from this result? Suppose the iterates reach a neighborhood of a minimizer  $x^*$ , and let  $\eta_{\mathcal{B}} = 2/\lambda_{\max}(\nabla^2 \hat{\mathcal{L}}_{\mathcal{B}}(x^*))$  denote the linear stability threshold of a batch loss  $\hat{\mathcal{L}}_{\mathcal{B}}$ . Clearly, for small enough neighbourhood, GD can diverge only if the step size satisfies  $\eta \geq \eta_{\mathcal{B}}$ . Notably, if the condition for stable oscillations in Thm. 1 is violated, then superlinear divergence may already occur at the threshold  $\eta = \eta_{\mathcal{B}}$ . In this case, the stability threshold of SGD is effectively capped by  $\eta_{\mathcal{B}}$ . This naturally raises the following question. Under what conditions can we guarantee the stability of SGD? The result below addresses this point (see proof in App. 5.2).

**Theorem 3 (Sufficient condition)** Let  $x^*$  be an interpolating minimizer of  $\mathcal{L}$ , and consider SGD's iterates (10) denoted by  $\{x_t\}$ . If

$$\eta < \min_{\mathcal{B}: |\mathcal{B}| = B} \frac{2}{\lambda_{\max} \left( \nabla^2 \hat{\mathcal{L}}_{\mathcal{B}}(\boldsymbol{x}^*) \right)},$$
(13)

then there exists a neighborhood  $\{x_0: \|x_0 - x^*\| < \rho\}$  s.t.  $\mathbb{E}[\|x_t - x^*\|_k^k] \rho^{-k} \underset{t \to \infty}{\longrightarrow} 0$  for all even k.

This result shows that if the step size is linearly stable with respect to all batches, then the full nonlinear dynamics of SGD are stable in expectation. As demonstrated in Sec. 2.2, this sufficient condition can also be necessary in certain cases. Specifically, when unstable oscillations leading to superlinear divergence arise in batches with low linear stability thresholds. Since it takes only one such batch, and the number of possible batches is exponentially large, it is quite likely that SGD operates in this regime.

## 5 DERIVATIONS

# 5.1 Gradient descent oscillations as a flip bifurcation

In this section, we give a brief review of bifurcations and formulate GD's dynamics in this framework. For a comprehensive overview of bifurcations, see Kuznetsov (1998). Consider the parameter-dependent nonlinear system  $x_{t+1} = \psi(x_t, \eta)$  with fixed point  $x^*$ , i.e.,  $\psi(x^*, \eta) = x^*$ . In general, bifurcations of fixed points occur when a parameter changes its value, while affecting the stability of the dynamics. The special case of flip bifurcation, also called period-doubling, happens when the fixed point  $x^*$  loses stability as the parameter  $\eta$  changes, and a period-2 cycle emerges. Mathematically, let us define the critical value of the parameter  $\eta_c$  such that the dominant eigenvalue of the Jacobian  $\mathcal{D}_x \psi(x^*, \eta_c)$  equals -1. Then a flip bifurcation takes place while this eigenvalue crosses minus one on the real line as  $\eta$  exceeds  $\eta_c$ . Here  $\eta_c$  is the linear stability threshold.

In this case, for  $\eta < \eta_c$ , the fixed point  $x^*$  is stable, and if the iterates happen to arrive close by, they will be attracted to it. If the Jacobian has full rank, the iterates will in fact converge to  $x^*$ . However, when  $\eta$  is slightly above  $\eta_c$ , the fixed point  $x^*$  is no longer stable, and a period-2 cycle appears as

$$\psi(x^{(1)}, \eta) = x^{(2)}, \quad \text{and} \quad \psi(x^{(2)}, \eta) = x^{(1)}.$$
 (14)

The stability of the resulting period-2 cycle is governed by the coefficient of the cubic term in the corresponding normal form of the bifurcation. This form provides a canonical (standard) dynamics to which any flip bifurcation can be reduced. Concretely, consider the dynamics along the one-dimensional critical manifold, tangent to the dominant eigenvector of the Jacobian. Then this dynamics can be transformed into (Kuznetsov, 1998, Sec 5.4)

$$\xi_{t+1} = -\xi_t + C_0 \xi_t^3 + O(\xi_t^4), \tag{15}$$

where  $C_0$  is the first Lyapunov coefficient. When  $C_0 > 0$ , the resulting cycle is stable (supercritical bifurcation), and the dynamics in the long run will alternate between  $x^{(1)}$  and  $x^{(2)}$ . Whereas for

<sup>&</sup>lt;sup>1</sup>Dominant eigenvalue is an eigenvalue that has maximal absolute value. Here we assume that it is unique.

 $C_0 < 0$ , the cycle is unstable (subcritical bifurcation) and the iterates will diverge from  $x^*$ . The expression for  $C_0$ , involving the second- and third-order derivatives of  $\psi$  at  $x^*$ , is given in App. D.

We now turn to apply this theory to prove Thm. 1. In the context of GD's iterates near a minimum, the dynamics evolve according to the update rule

$$\boldsymbol{x}_{t+1} = \boldsymbol{x}_t - \eta \nabla \mathcal{L}(\boldsymbol{x}_t) \triangleq \boldsymbol{\psi}_t(\boldsymbol{x}_t, \eta), \tag{16}$$

where  $\mathcal{L}$  is an objective function to be minimized, and  $\eta$  is the step size. Obviously, minimizers of  $\mathcal{L}$  are the fixed points of  $\psi$ , as the gradient vanishes at these points. The Jacobian of the GD map is

$$\mathcal{D}_{x}\psi(x,\eta) = I - \eta \nabla^{2}\mathcal{L}(x). \tag{17}$$

Note that the eigenvalues of the Jacobian are given by  $\{1 - \eta \lambda_i(\nabla^2 \mathcal{L})\}$ . Let  $x^*$  be a minimizer of  $\mathcal{L}$ , then the critical value of  $\eta$  is the well known linear stability threshold

$$\eta_{\rm lin} = \frac{2}{\lambda_{\rm max}(\nabla^2 \mathcal{L}(\boldsymbol{x}^*))}.$$
 (18)

Thus, as  $\eta$  exceeds  $\eta_{\rm lin}$ , the dominant eigenvalue of the Jacobian crosses -1 on the real axis, matching the scenario of the flip bifurcation. Assuming  $\nabla^2 \mathcal{L}(\boldsymbol{x}^*)$  is strictly positive and  $\lambda_{\rm max}(\nabla^2 \mathcal{L}(\boldsymbol{x}^*))$  has multiplicity one, the stability of oscillations in the small neighborhood of  $\boldsymbol{x}^*$  is governed by  $C_0$ . Overall, we see that oscillations near a minimum  $\boldsymbol{x}^*$  are stable if and only if  $C_0$  is positive. In App. D we show that a positive Lyapunov coefficient is equivalent to the condition in (5).

## 5.2 Sufficient condition for stability of SGD

In this section, we derive Thm. 3. SGD update rule with step size  $\eta$  is given by

$$\boldsymbol{x}_{t+1} = \boldsymbol{x}_t - \eta \nabla \hat{\mathcal{L}}_{\mathcal{B}_t}(\boldsymbol{x}_t) \triangleq \hat{\boldsymbol{\psi}}_{\mathcal{B}_t}(\boldsymbol{x}_t), \tag{19}$$

where  $\hat{\psi}_{\mathcal{B}_t}: \mathbb{R}^d \to \mathbb{R}^d$  is the SGD map. As  $x^*$  is an interpolating minimizer of  $\mathcal{L}$ , we have

$$\hat{\psi}_{\mathcal{B}_t}(\boldsymbol{x}^*) = \boldsymbol{x}^* \quad \text{w.p. } 1. \tag{20}$$

Since  $\{f_i\}$  are analytic, we can use the Taylor expansions of  $\hat{\psi}_{\mathcal{B}_t}$  at  $x^*$  to get

$$\hat{\psi}_{\mathcal{B}_t}(\boldsymbol{x}) = \hat{\psi}_{\mathcal{B}_t}(\boldsymbol{x}^*) + \sum_{k=1}^{\infty} \frac{1}{k!} \mathcal{D}^k \hat{\psi}_{\mathcal{B}_t}(\boldsymbol{x}^*) (\boldsymbol{x} - \boldsymbol{x}^*)^{\otimes k}, \tag{21}$$

where  $\mathcal{D}^k$  is the kth order derivative in matrix form (not to be confused with  $\mathcal{D}^k$ ). Let

$$\Delta \boldsymbol{x}_t^k \triangleq (\boldsymbol{x}_t - \boldsymbol{x}^*)^{\otimes k} \in \mathbb{R}^{d^k}$$
 (22)

be the kth Kronecker power of the distance to the minimum<sup>2</sup>. Then from the update rule in (19)

$$\mathbb{E}\left[\Delta \boldsymbol{x}_{t+1}\right] = \mathbb{E}\left[\sum_{k=1}^{\infty} \frac{1}{k!} \mathcal{D}^{k} \hat{\boldsymbol{\psi}}_{\mathcal{B}_{t}}(\boldsymbol{x}^{*}) \Delta \boldsymbol{x}_{t}^{k}\right] = \sum_{k=1}^{\infty} \frac{1}{k!} \mathbb{E}\left[\mathcal{D}^{k} \hat{\boldsymbol{\psi}}_{\mathcal{B}_{t}}(\boldsymbol{x}^{*})\right] \mathbb{E}\left[\Delta \boldsymbol{x}_{t}^{k}\right], \tag{23}$$

where we used the fact that  $\{\mathcal{D}^k\hat{\psi}_{\mathcal{B}_t}(\boldsymbol{x}^*)\}_{k=1}^{\infty}$  and  $\boldsymbol{x}_t$  are statistically independent. We see that the evolution of the first moment of the distance to the minimum,  $\mathbb{E}[\Delta \boldsymbol{x}]$ , depends *linearly* on all higher-order moments  $\{\mathbb{E}[\Delta \boldsymbol{x}^k]\}_{k=1}^{\infty}$ . Consequently, analyzing the stability of SGD in expectation requires studying the joint dynamics of all moments. In App. H, we show that the evolution of the kth moment over time is

$$\mathbb{E}\left[\Delta \boldsymbol{x}_{t+1}^{k}\right] = \mathbb{E}\left[\left(\Delta \boldsymbol{x}_{t+1}\right)^{\otimes k}\right] = \mathbb{E}\left[\left(\sum_{p=1}^{\infty} \frac{1}{p!} \mathcal{D}^{p} \hat{\boldsymbol{\psi}}_{\mathcal{B}_{t}}(\boldsymbol{x}^{*}) \Delta \boldsymbol{x}_{t}^{p}\right)^{\otimes k}\right] = \sum_{p=k}^{\infty} \boldsymbol{\Psi}_{k,p} \mathbb{E}\left[\Delta \boldsymbol{x}_{t}^{p}\right], \quad (24)$$

where explicit expression for  $\Psi_{k,p} \in \mathbb{R}^{d^k \times d^p}$  is given in App. H. Once again, we obtain a linear relation between the moments at successive times. This motivates us to express the mapping from  $\{\mathbb{E}[\Delta x_t^k]\}_{k=1}^{\infty}$  to  $\{\mathbb{E}[\Delta x_{t+1}^k]\}_{k=1}^{\infty}$  as a linear operator on the infinite-dimensional Hilbert space  $\ell_2$ .

The first power  $\Delta x_t^1$  is denoted simply  $\Delta x_t$ .

This formulation is meaningful only if the sequence has finite norm and the operator is bounded. To ensure this, we introduce a radius  $\rho > 0$  and analyze a scaled version of the moments. Let

$$\bar{\boldsymbol{\mu}}_{t}^{k} \triangleq \mathbb{E}\left[\left(\frac{\boldsymbol{x}_{t} - \boldsymbol{x}^{*}}{\rho}\right)^{\otimes k}\right] = \rho^{-k}\mathbb{E}\left[\Delta\boldsymbol{x}_{t}^{k}\right]. \tag{25}$$

Therefore,

378

379

380

381

382

384

386

387

388

389 390

391 392 393

394

396

397

398

399

400 401

402

403 404

405

406

407

408

409

411 412 413

414 415

416 417

418 419

420

421

422

423

424

425

426

427

428

429

430

431

$$\bar{\boldsymbol{\mu}}_{t+1}^{k} = \rho^{-k} \mathbb{E} \left[ \Delta \boldsymbol{x}_{t+1}^{k} \right] = \sum_{p=k}^{\infty} \rho^{-k} \boldsymbol{\Psi}_{k,p} \mathbb{E} \left[ \Delta \boldsymbol{x}_{t}^{p} \right] = \sum_{p=k}^{\infty} \rho^{p-k} \boldsymbol{\Psi}_{k,p} \mathbb{E} \left[ \rho^{-p} \Delta \boldsymbol{x}_{t}^{p} \right] = \sum_{p=k}^{\infty} \rho^{p-k} \boldsymbol{\Psi}_{k,p} \bar{\boldsymbol{\mu}}_{t}^{k}.$$
(26)

Define the linear operator  $\Psi_{\rho}$  in Hilbert space  $\ell_2$  and the moments vector  $\bar{\mu}_t$  as

$$\bar{\boldsymbol{\mu}}_{t} = \begin{bmatrix} \bar{\boldsymbol{\mu}}_{t}^{1} \\ \bar{\boldsymbol{\mu}}_{t}^{2} \\ \bar{\boldsymbol{\mu}}_{t}^{3} \\ \vdots \end{bmatrix} \quad \text{and} \quad \boldsymbol{\Psi}_{\rho} = \begin{bmatrix} \boldsymbol{\Psi}_{1,1} & \rho \boldsymbol{\Psi}_{1,2} & \rho^{2} \boldsymbol{\Psi}_{1,3} & \cdots \\ \mathbf{0} & \boldsymbol{\Psi}_{2,2} & \rho \boldsymbol{\Psi}_{2,3} & \cdots \\ \mathbf{0} & \mathbf{0} & \boldsymbol{\Psi}_{3,3} & \cdots \\ \vdots & \vdots & \vdots & \ddots \end{bmatrix}, \quad (27)$$

then

$$\bar{\boldsymbol{\mu}}_{t+1} = \boldsymbol{\Psi}_{o} \bar{\boldsymbol{\mu}}_{t}. \tag{28}$$

 $\bar{\mu}_{t+1} = \Psi_{\rho}\bar{\mu}_t. \tag{28}$  This relation is valid only when  $\Psi_{\rho}$  is bounded. Intuitively, taking smaller values of  $\rho$  can help bound the operator. Assuming the operator is bounded, (28) unfolds as  $\bar{\mu}_t = \Psi_\rho^t \bar{\mu}_0$ . To impose a condition on the initial point  $x_0$ , observe that  $\bar{\mu}_0 \in \ell_2$ , and thus must be square-summable. Hence,

$$\|\bar{\boldsymbol{\mu}}_0\|^2 = \sum_{k=1}^{\infty} \|\bar{\boldsymbol{\mu}}_0^k\|^2 = \sum_{k=1}^{\infty} \left\| \left( \frac{\boldsymbol{x}_0 - \boldsymbol{x}^*}{\rho} \right)^{\otimes k} \right\|^2 = \sum_{k=1}^{\infty} \left( \frac{\|\boldsymbol{x}_0 - \boldsymbol{x}^*\|}{\rho} \right)^{2k}. \tag{29}$$

The above expression is finite if and only if  $||x_0 - x^*|| < \rho$ , which defines the neighborhood around the minimum where our analysis applies. We see that choosing a smaller  $\rho$  to ensure boundedness of the operator correspondingly shrinks this neighborhood.

For stable dynamics in expectation, the linear system in (28) must be stable. Note that under  $\Psi_{\rho}$ , moment vectors map naturally to moment vectors. Thus, to get the exact stability threshold, we would need the response of  $\Psi_{\rho}$  to be smaller than one on this restricted set. Instead, we relax this constraint to obtain a sufficient condition, requiring stability for any vector in  $\ell_2$ . In App. I, we prove that under the condition (13), there exists a value of  $\rho > 0$  ensuring boundedness of the operator. Then, in App. K we show that once the operator is bounded, its spectral radius is strictly less than one. Therefore,  $\|\bar{\mu}_t\| \to 0$  as t tends to infinity (see App. G). Hence, each normalized moment also tends to zero elementwise, i.e.,  $\bar{\mu}_t^k \to 0$ . Since  $\bar{\mu}_t^k$  contains all degree-k monomials of the components of  $\Delta x_t$ , denoted  $\{\Delta x_{t,i}\}_{i=1}^d$ , summing over the subset of single-variable terms we get

$$\sum_{i=1}^{d} \mathbb{E}\left[ (\Delta \boldsymbol{x}_{t,i})^{k} \right] = \mathbb{E}\left[ \sum_{i=1}^{d} (\boldsymbol{x}_{t,i} - \boldsymbol{x}_{i}^{*})^{k} \right] \rho^{-k} \underset{t \to \infty}{\longrightarrow} 0.$$
 (30)

Restricting this to even order moments (even k), we get  $\mathbb{E}\left[\|\boldsymbol{x}_t - \boldsymbol{x}^*\|_k^k\right] \rho^{-k} \longrightarrow 0$ .

## RELATED WORK

**Bifurcation, oscillations and EoS in GD.** Cohen et al. (2021) examined the behavior of GD, and showed that it typically happens at the edge of stability. Wang et al. (2022) proved progressive sharpening for a two-layer network and analyzed the EOS dynamics through four phases, depending on the change in the sharpness value. Zhu et al. (2022) gave a simple example that exhibits EoS. Ma et al. (2022) analyzed EoS under the assumption of subquadratic growth of the loss. Ahn et al. (2022) illustrated that unstable convergence is possible in specific cases. Damian et al. (2023) showed how GD self-stabilizes. Specifically, they demonstrated 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  $2/\eta$ . Kreisler et al. (2023); Song & Yun (2023) proved that under EoS, different GD trajectories align on a specific bifurcation diagram independent of initialization. Ghosh et al. (2025) analyzed the dynamics of deep linear networks, focusing on 2-period cycle, while showing that oscillations occur within a small subspace, where the dimension of the subspace is controlled by the step size. Chen et al. (2024) studied GD dynamics on quadratic loss from stability up to the chaos phase.

**Stability of SGD.** Empirically, Keskar et al. (2016); Jastrzębski et al. (2017); Jastrzębski et al. (2019; 2020) have shown that SGD with a large step size or small batch size leads to flatter minima. Cohen et al. (2021, App. G) found that with large batches, the sharpness behaves similarly to full-batch gradient descent. Gilmer et al. (2022) studied how the curvature of the loss affects the training dynamics in multiple settings.

On the theoretical side, Wu et al. (2018) analyzed stability in the mean-square sense and provided an implicit sufficient condition. Granziol et al. (2022) used random matrix theory to characterize the maximal stable learning rate as a function of batch size, under certain assumptions on Hessian noise. Velikanov et al. (2023) studied SGD with momentum and derived an implicit upper bound on the learning rate using spectrally expressible approximations and a moment-generating function. Ma & Ying (2021) investigated higher-order moments of SGD and established an implicit necessary and sufficient stability condition. Wu et al. (2022) proposed a necessary condition based on an alignment property, though a general analytic bound for this property is missing. Ziyin et al. (2023) examined stability in probability rather than in mean square, showing that SGD can in theory converge with high probability to linearly unstable minima for GD, *i.e.*, where  $\eta \gg 2/\lambda_{\rm max}(\nabla^2 \mathcal{L})$ . However, this prediction was not observed empirically. Mulayoff et al. (2021) considered non-differentiable minima and derived a necessary condition for strong stability, meaning SGD remains within a ball around the minimum. Finally, Mulayoff & Michaeli (2024) provided the exact stability criterion explicitly in closed-form expression for the linearized dynamics.

Additionally, Liu et al. (2021) analyzed the covariance matrix of the stationary distribution of iterates near minima, and Ziyin et al. (2022) extended these results by deriving an implicit relation between this covariance and that of the gradient noise. However, both works leave open the question of when the dynamics actually converge to a stationary state. Recently, Lee & Jang (2023) examined the stability of SGD along its trajectory and established an explicit exact condition for objective decrease via a descent lemma in expectation.

## 7 CONCLUSION, LIMITATIONS AND FUTURE DIRECTIONS

In this paper, we investigated the nonlinear stability of gradient descent GD and SGD. For GD, we derived an explicit condition characterizing when stable oscillations arise at the edge of stability, namely, when the cubic term dominates the quartic term in the local Taylor expansion of the objective. For SGD, we showed that the instability of even a single batch can be sufficient to render the entire dynamics unstable in expectation, implying that stability is dictated by the worst-case batch rather than by an average effect. Finally, we proved that if the step size is stable with respect to all batches, then all moments of the full nonlinear SGD dynamics remain stable in a neighborhood of the minimizer. Together, these results reveal that nonlinear effects can fundamentally reshape the stability landscape compared with standard linear analyses.

**Limitations and future directions.** Our analysis of oscillations in GD focuses on isolated minima, where the Jacobian of the dynamical system has a single critical eigenvalue. In deep learning, however, minima often form low-dimensional manifolds with multiple near-critical directions. This can lead to richer local dynamics, including combinations of fold and flip behaviors, and a more complex stability picture. Extending our results to this setting, potentially via generalized fold-flip bifurcations (Kuznetsov et al., 2004), is an important direction for future work.

For our analysis of SGD we assume an interpolating setting in which all mini-batches share the same minimizer. This allows us to prove that stability of each batch implies stability of the full dynamics. In practice, however, batches may have distinct or only approximately aligned minima. In such cases, SGD cannot converge exactly to the minimizer; even if the dynamics are stable, the algorithm exhibits an inherent bias in expectation (Défossez & Bach, 2015). In this work we adopt a convergence-in-expectation perspective for SGD, but in non-interpolating settings the limiting point may be biased or correspond to a different minimum. Developing a principled notion of nonlinear stability that captures this behavior remains an important direction for future research.

# REFERENCES

- Kwangjun Ahn, Jingzhao Zhang, and Suvrit Sra. Understanding the unstable convergence of gradient descent. In *International conference on machine learning*, pp. 247–257. PMLR, 2022.
- Xuxing Chen, Krishnakumar Balasubramanian, Promit Ghosal, and Bhavya Agrawalla. From stability to chaos: Analyzing gradient descent dynamics in quadratic regression. *Transactions on machine learning research*, 2024.
- Jeremy Cohen, Simran Kaur, Yuanzhi Li, J Zico Kolter, and Ameet Talwalkar. Gradient descent on neural networks typically occurs at the edge of stability. In *International Conference on Learning Representations*, 2021.
  - Alex Damian, Eshaan Nichani, and Jason D. Lee. Self-stabilization: The implicit bias of gradient descent at the edge of stability. In *The Eleventh International Conference on Learning Represen*tations, 2023.
  - Alexandre Défossez and Francis Bach. Averaged least-mean-squares: Bias-variance trade-offs and optimal sampling distributions. In *Artificial Intelligence and Statistics*, pp. 205–213. PMLR, 2015.
  - Nelson Dunford and Jacob T Schwartz. *Linear operators II: spectral theory*. New York: Wiley-Interscience,, 1964.
  - Avrajit Ghosh, Soo Min Kwon, Rongrong Wang, Saiprasad Ravishankar, and Qing Qu. Learning dynamics of deep matrix factorization beyond the edge of stability. In *The Second Conference on Parsimony and Learning (Recent Spotlight Track)*, 2025.
  - Justin Gilmer, Behrooz Ghorbani, Ankush Garg, Sneha Kudugunta, Behnam Neyshabur, David Cardoze, George Edward Dahl, Zachary Nado, and Orhan Firat. A loss curvature perspective on training instabilities of deep learning models. In *International Conference on Learning Representations*, 2022.
  - Diego Granziol, Stefan Zohren, and Stephen Roberts. Learning rates as a function of batch size: A random matrix theory approach to neural network training. *J. Mach. Learn. Res*, 23:1–65, 2022.
  - Stanisław Jastrzębski, Zachary Kenton, Devansh Arpit, Nicolas Ballas, Asja Fischer, Yoshua Bengio, and Amos Storkey. Three factors influencing minima in SGD. *arXiv preprint arXiv:1711.04623*, 2017.
  - Stanisław Jastrzębski, Zachary Kenton, Nicolas Ballas, Asja Fischer, Yoshua Bengio, and Amos Storkey. On the relation between the sharpest directions of DNN loss and the SGD step length. In *International Conference on Learning Representations*, 2019.
  - Stanisław Jastrzębski, Maciej Szymczak, Stanislav Fort, Devansh Arpit, Jacek Tabor, Kyunghyun Cho\*, and Krzysztof Geras\*. The break-even point on optimization trajectories of deep neural networks. In *International Conference on Learning Representations*, 2020.
  - Nitish Shirish Keskar, Dheevatsa Mudigere, Jorge Nocedal, Mikhail Smelyanskiy, and Ping Tak Peter Tang. On large-batch training for deep learning: Generalization gap and sharp minima. *arXiv* preprint arXiv:1609.04836, 2016.
  - Bernard O Koopman. Hamiltonian systems and transformation in hilbert space. *Proceedings of the National Academy of Sciences*, 17(5):315–318, 1931.
  - Itai Kreisler, Mor Shpigel Nacson, Daniel Soudry, and Yair Carmon. Gradient descent monotonically decreases the sharpness of gradient flow solutions in scalar networks and beyond. In *International Conference on Machine Learning*, pp. 17684–17744. PMLR, 2023.
  - Yu A Kuznetsov, Hil GE Meijer, and Lennaert van Veen. The fold-flip bifurcation. *International Journal of Bifurcation and Chaos*, 14(07):2253–2282, 2004.
  - Yuri A Kuznetsov. Elements of applied bifurcation theory. Springer, 1998.

543

544

546

547 548

549

550

551

552 553

554

555

556

558 559

560

561 562

563

564

565

566

567 568

569

570

571

572 573

574

575

576

577

578

579 580

581

582

583

584

585

586

588

589

590

591 592

- 540 Sungyoon Lee and Cheongjae Jang. A new characterization of the edge of stability based on a sharpness measure aware of batch gradient distribution. In The Eleventh International Conference 542 on Learning Representations, 2023.
  - Kangqiao Liu, Liu Ziyin, and Masahito Ueda. Noise and fluctuation of finite learning rate stochastic gradient descent. In Marina Meila and Tong Zhang (eds.), Proceedings of the 38th International Conference on Machine Learning, volume 139 of Proceedings of Machine Learning Research, pp. 7045–7056. PMLR, 18–24 Jul 2021.
  - Chao Ma and Lexing Ying. On linear stability of SGD and input-smoothness of neural networks. In Thirty-Fifth Conference on Neural Information Processing Systems, 2021.
  - Chao Ma, Daniel Kunin, Lei Wu, and Lexing Ying. Beyond the quadratic approximation: The multiscale structure of neural network loss landscapes. arXiv preprint arXiv:2204.11326, 2022.
  - Rotem Mulayoff and Tomer Michaeli. Unique properties of flat minima in deep networks. In International Conference on Machine Learning, pp. 7108–7118. PMLR, 2020.
  - Rotem Mulayoff and Tomer Michaeli. Exact mean square linear stability analysis for SGD. In Proceedings of Thirty Seventh Conference on Learning Theory, volume 247 of Proceedings of Machine Learning Research, pp. 3915–3969. PMLR, July 2024.
  - Rotem Mulayoff, Tomer Michaeli, and Daniel Soudry. The implicit bias of minima stability: A view from function space. Advances in Neural Information Processing Systems, 34:17749–17761, 2021.
  - Mor Shpigel Nacson, Rotem Mulayoff, Greg Ongie, Tomer Michaeli, and Daniel Soudry. The implicit bias of minima stability in multivariate shallow ReLU networks. In The Eleventh International Conference on Learning Representations, 2023.
  - Behnam Neyshabur, Ryota Tomioka, and Nathan Srebro. In search of the real inductive bias: On the role of implicit regularization in deep learning. arXiv preprint arXiv:1412.6614, 2014.
  - Minhak Song and Chulhee Yun. Trajectory alignment: understanding the edge of stability phenomenon via bifurcation theory. arXiv preprint arXiv:2307.04204, 2023.
  - John Michael Tutill Thompson and H Bruce Stewart. *Nonlinear dynamics and chaos*. John Wiley & Sons, 2002.
  - Maksim Velikanov, Denis Kuznedelev, and Dmitry Yarotsky. A view of mini-batch SGD via generating functions: conditions of convergence, phase transitions, benefit from negative momenta. In The Eleventh International Conference on Learning Representations, 2023.
  - Zixuan Wang, Zhouzi Li, and Jian Li. Analyzing sharpness along GD trajectory: Progressive sharpening and edge of stability. Advances in Neural Information Processing Systems, 35:9983–9994, 2022.
  - Lei Wu, Chao Ma, and E Weinan. How SGD selects the global minima in over-parameterized learning: A dynamical stability perspective. In Advances in Neural Information Processing Systems, pp. 8279–8288, 2018.
  - Lei Wu, Mingze Wang, and Weijie Su. The alignment property of SGD noise and how it helps select flat minima: A stability analysis. Advances in Neural Information Processing Systems, 35: 4680–4693, 2022.
  - Xingyu Zhu, Zixuan Wang, Xiang Wang, Mo Zhou, and Rong Ge. Understanding edge-of-stability training dynamics with a minimalist example. arXiv preprint arXiv:2210.03294, 2022.
  - Liu Ziyin, Kangqiao Liu, Takashi Mori, and Masahito Ueda. Strength of minibatch noise in SGD. In International Conference on Learning Representations, 2022.
  - Liu Ziyin, Botao Li, Tomer Galanti, and Masahito Ueda. The probabilistic stability of stochastic gradient descent. arXiv preprint arXiv:2303.13093, 2023.

# **APPENDICES**

# A LARGE LANGUAGE MODEL (LLM) USAGE DISCLOSURE

In this paper we used LLMs only to polish the text. The LLM was given a text written by the authors, and it suggested alternative ways of writing. These suggestions, if accepted, were further refined by the authors and not been used as is. Specifically, LLMs were *not* used to generate any text from scratch, or to suggest research direction or to derive the results.

# B BACKGROUND ON LINEAR STABILITY OF SGD

Analyzing the full dynamics of SGD can be hard. Therefore, many works opt to study the linearized dynamic near minima (Wu et al., 2018; Ma & Ying, 2021; Mulayoff et al., 2021; Mulayoff & Michaeli, 2024), as it is common in the analysis of nonlinear systems. In our paper we focus on interpolating minimizers, defined in Def. 1. In this case, the linearized dynamics is defined below.

**Definition 2 (Linearized dynamics)** Let  $\mathcal{L}$  form (9), and  $x^*$  be its interpolating minimizer, s.t.  $\mathcal{L}$  is twice differentiable at  $x^*$ . Then the linearized dynamics of SGD near  $x^*$  are given by

$$\tilde{\boldsymbol{x}}_{t+1} = \tilde{\boldsymbol{x}}_t - \frac{\eta}{B} \sum_{i \in \mathcal{B}_t} \nabla^2 f_i(\boldsymbol{x}^*) (\tilde{\boldsymbol{x}}_t - \boldsymbol{x}^*).$$
(31)

The linearized dynamics can be viewed as SGD on the second-order approximation of  $\mathcal{L}$  at  $x^*$ ,

$$\tilde{\mathcal{L}}(\boldsymbol{x}) = \mathcal{L}(\boldsymbol{x}^*) + \frac{1}{2}(\boldsymbol{x} - \boldsymbol{x}^*)^{\mathrm{T}} \nabla^2 \mathcal{L}(\boldsymbol{x}^*)(\boldsymbol{x} - \boldsymbol{x}^*).$$
(32)

Therefore, linear dynamics analysis is exact only when  $\{f_i\}$  are all quadratic potentials.

There are few traditional ways to define the convergence of random processes, such as the iterates of SGD. One prominent choice is to use the mean square sense of convergence to define stability. For univariate optimization, the mean square linear stability threshold is as follows. Generalization to higher dimensions and non-interpolating minima can be found in Mulayoff & Michaeli (2024).

**Theorem 4 (Univariate linear stability threshold, Wu et al. (2018))** Let  $f_i : \mathbb{R} \to \mathbb{R}$  be twice differentiable functions and let  $x^*$  be an interpolating minimum of the loss, i.e.,

$$\forall 1 \le i \le n \qquad f_i'(x^*) = 0 \qquad \text{and} \qquad h_i \triangleq f_i''(x^*) > 0.$$
 (33)

Define

$$h = \frac{1}{n} \sum_{i=1}^{n} h_i, \qquad s^2 = \frac{1}{n} \sum_{i=1}^{n} (h_i - h)^2, \qquad and \qquad p = \frac{n - B}{B(n - 1)}.$$
 (34)

Consider the iterates of the linearized SGD  $\{\tilde{x}_t\}$  in (31). Then,  $\mathbb{E}[(\tilde{x}_t - x^*)^2]$  is bounded if and only if  $\eta \leq \eta_{\text{lin}}$ , where

$$\eta_{\rm lin} \triangleq \frac{2h}{h^2 + ps^2}.\tag{35}$$

From this result, we see that the linear stability threshold  $\eta_{\text{lin}}$  takes into account the sharpness of all functions  $\{f_i\}$ . When the batch size B equals one, we get p=1 and then

$$\eta_{\text{lin}} = \frac{2h}{h^2 + s^2} = 2 \frac{\sum_{i=1}^{n} h_i}{\sum_{i=1}^{n} h_i^2}.$$
 (36)

# C PROOF OF PROPOSITION 1

Let

$$f_{+}(x) = \frac{1}{2}x^{2} + \frac{1}{4}x^{4}. (37)$$

Here we show that when GD is applied on  $f_+$  with step size  $\eta > 2$ , its iterates diverge in rate higher than linear. In this case, Thm. 2 tell us that SGD dynamics also diverge.

The GD map on  $f_+$  is

$$\psi(x) = x - \eta f'_{+}(x) = (1 - \eta)x - \eta x^{3}.$$
(38)

Define

$$\psi^t = \underbrace{\psi \circ \dots \circ \psi \circ \psi}_{t \text{ times}}. \tag{39}$$

Assume that  $\eta > 2$ , and note that for all  $x \in \mathbb{R}$ 

$$|\psi(x)| = |(1 - \eta)x - \eta x^3| = (\eta - 1)|x| + \eta|x|^3 = |\psi(|x|)|. \tag{40}$$

Let  $\tilde{\psi}(x) = |\psi|(|x|)$ , then

$$|\psi^t(x_0)| = \tilde{\psi}^t(|x_0|). \tag{41}$$

Since  $\tilde{\psi}: \mathbb{R}_+ \to \mathbb{R}_+$  is monotonically increasing on  $\mathbb{R}_+$ , we have that its composition  $\tilde{\psi}^t$  of any order is also monotonically increasing. Furthermore, we can bound  $\tilde{\psi}$  from below with

$$\tilde{\psi}(x) = (\eta - 1)|x| + \eta|x|^3 \ge \max\{(\eta - 1)|x|, \eta|x|^3\} \triangleq \varphi(x). \tag{42}$$

Thus,

$$|x_t| = |\psi^t(x_0)| \ge \varphi^t(x_0) \triangleq \chi_t. \tag{43}$$

Again,  $\varphi: \mathbb{R}_+ \to \mathbb{R}_+$  is monotonically increasing, and therefore any composition with itself is also monotonically increasing on  $\mathbb{R}_+$ . Obviously,

$$\varphi(x) \ge (\eta - 1)|x|$$
 and  $\varphi(x) \ge \eta |x|^3$ . (44)

Then we can bound  $\chi_t$  by

$$\chi_t \ge (\eta - 1)^t |x_0|. \tag{45}$$

Since  $\eta > 2$ , there exists  $T \in \mathbb{N}$  such that

$$\chi_T \ge (\eta - 1)^T |x_0| > 2. \tag{46}$$

Now, for all t > T we have

$$\chi_t = \varphi^{t-T}(\chi_T). \tag{47}$$

Here, we will use the second bound, i.e.,  $\varphi(x) \ge \eta |x|^3$ , t-T times as

$$\chi_{t} = \underbrace{\eta |\eta| \cdots \eta |\chi_{T}|^{3} \cdots |^{3}|^{3}}_{t-T \text{ times}}$$

$$= \eta^{(3^{t-T}-1)/2} |\chi_{T}|^{3^{t-T}}$$

$$\geq 2^{(3^{t-T}-1)/2} 2^{3^{t-T}}$$

$$= 2^{(3^{t-T+1}-1)/2}.$$
(48)

Therefore,  $\chi_t$  diverges with superlinear rate and so does  $|x_t|$ .

## D CONDITION FOR STABLE OSCILLATIONS IN GD

In Sec. 5.1 we formulate the oscillations of GD as a flip bifurcation. We saw that the first Lyapunov coefficient  $C_0$  controls the stability of the oscillations. For a general nonlinear map  $\psi(x, \eta)$ , this coefficient is given by (Kuznetsov, 1998, Sec. 5.4)

$$C_0 = \frac{1}{6} \left\langle \boldsymbol{u}, \mathcal{D}_{\boldsymbol{x}}^3 \boldsymbol{\psi}(\boldsymbol{x}^*, \eta) [\boldsymbol{v}, \boldsymbol{v}, \boldsymbol{v}] \right\rangle - \frac{1}{2} \left\langle \boldsymbol{u}, \mathcal{D}_{\boldsymbol{x}}^2 \boldsymbol{\psi}(\boldsymbol{x}^*, \eta) [\boldsymbol{v}, \boldsymbol{p}] \right\rangle, \tag{49}$$

where u and v are normalized left and right eigenvectors of the Jacobian  $\mathcal{D}_x \psi(x^*, \eta)$  corresponding to the eigenvalue -1, such that  $\langle u, v \rangle = 1$ , and

$$\boldsymbol{p} = \left[ \mathcal{D}_{\boldsymbol{x}} \boldsymbol{\psi}(\boldsymbol{x}^*, \eta) - \boldsymbol{I} \right]^{-1} \mathcal{D}_{\boldsymbol{x}}^2 \boldsymbol{\psi}(\boldsymbol{x}^*, \eta) [\boldsymbol{v}, \boldsymbol{v}]. \tag{50}$$

We would like to write these expressions in terms of the loss function  $\mathcal{L}$  for the GD dynamics. In this case,  $\psi(x, \eta) = x - \eta \nabla \mathcal{L}(x)$ , and we have that the Jacobian is

$$\mathcal{D}_{\boldsymbol{x}}\psi(\boldsymbol{x}^*,\eta) = \boldsymbol{I} - \eta \nabla^2 \mathcal{L}(\boldsymbol{x}^*). \tag{51}$$

Since this Jacobian is systematic, v equals u. Moreover, it is easy to see that v is the top eigenvector of the Hessian, corresponding to  $\lambda_{\max}(\nabla^2 \mathcal{L}(x^*))$ . Additionally,

$$\mathcal{D}_{\boldsymbol{x}}^{2}\boldsymbol{\psi}(\boldsymbol{x}^{*},\eta)[\boldsymbol{v},\boldsymbol{v}] = -\eta\mathcal{D}^{3}\mathcal{L}(\boldsymbol{x}^{*})[\boldsymbol{v},\boldsymbol{v}] = -\frac{\eta}{3}\nabla_{\boldsymbol{v}}\mathcal{D}^{3}\mathcal{L}(\boldsymbol{x}^{*})[\boldsymbol{v},\boldsymbol{v},\boldsymbol{v}]. \tag{52}$$

Thus,

$$p = \left[ -\eta \nabla^{2} \mathcal{L}(\boldsymbol{x}^{*}) \right]^{-1} \left( -\frac{\eta}{3} \nabla_{\boldsymbol{v}} \mathcal{D}^{3} \mathcal{L}(\boldsymbol{x}^{*}) [\boldsymbol{v}, \boldsymbol{v}, \boldsymbol{v}] \right)$$

$$= \frac{1}{3} \left[ \nabla^{2} \mathcal{L}(\boldsymbol{x}^{*}) \right]^{-1} \nabla_{\boldsymbol{v}} \mathcal{D}^{3} \mathcal{L}(\boldsymbol{x}^{*}) [\boldsymbol{v}, \boldsymbol{v}, \boldsymbol{v}]$$

$$= \frac{1}{3} \boldsymbol{q},$$
(53)

where

$$q = \left[\nabla^2 \mathcal{L}(\boldsymbol{x}^*)\right]^{-1} \nabla_{\boldsymbol{v}} \mathcal{D}^3 \mathcal{L}(\boldsymbol{x}^*) [\boldsymbol{v}, \boldsymbol{v}, \boldsymbol{v}]. \tag{54}$$

Next we have,

$$\langle \boldsymbol{u}, \mathcal{D}_{\boldsymbol{x}}^2 \boldsymbol{\psi}(\boldsymbol{x}^*, \eta) [\boldsymbol{v}, \boldsymbol{p}] \rangle = \frac{1}{3} \langle \boldsymbol{v}, \mathcal{D}_{\boldsymbol{x}}^2 \boldsymbol{\psi}(\boldsymbol{x}^*, \eta) [\boldsymbol{v}, \boldsymbol{q}] \rangle = -\frac{\eta}{3} \mathcal{D}^3 \mathcal{L}(\boldsymbol{x}^*) [\boldsymbol{v}, \boldsymbol{v}, \boldsymbol{q}].$$
 (55)

And the first term in  $C_0$  is

$$\langle \boldsymbol{u}, \mathcal{D}_{\boldsymbol{x}}^{3} \boldsymbol{\psi}(\boldsymbol{x}^{*}, \eta) [\boldsymbol{v}, \boldsymbol{v}, \boldsymbol{v}] \rangle = \langle \boldsymbol{v}, -\eta \mathcal{D}^{4} \mathcal{L}(\boldsymbol{x}^{*}) [\boldsymbol{v}, \boldsymbol{v}, \boldsymbol{v}] \rangle = -\eta \mathcal{D}^{4} \mathcal{L}(\boldsymbol{x}^{*}) [\boldsymbol{v}, \boldsymbol{v}, \boldsymbol{v}, \boldsymbol{v}].$$
 (56)

Overall.

$$C_0 = -\frac{\eta}{6} \mathcal{D}^4 \mathcal{L}(\boldsymbol{x}^*) [\boldsymbol{v}, \boldsymbol{v}, \boldsymbol{v}, \boldsymbol{v}] + \frac{\eta}{6} \mathcal{D}^3 \mathcal{L}(\boldsymbol{x}^*) [\boldsymbol{v}, \boldsymbol{v}, \boldsymbol{q}]. \tag{57}$$

A period-2 cycle near  $x^*$  is stable if and only if  $C_0 > 0$  (Kuznetsov, 1998), which results in the condition

$$\mathcal{D}^{3}\mathcal{L}(\boldsymbol{x}^{*})[\boldsymbol{v},\boldsymbol{v},\boldsymbol{q}] > \mathcal{D}^{4}\mathcal{L}(\boldsymbol{x}^{*})[\boldsymbol{v},\boldsymbol{v},\boldsymbol{v},\boldsymbol{v}]. \tag{58}$$

Originally, the scale of v was important for the magnitude of  $C_0$ . However, the scale of v has no effect on the sign of  $C_0$ , and thus has no impact on this condition.

## E ANALYTIC EXAMPLE OF THEOREM 1

In this section we consider the oscillations of GD at the edge of stability on the function

$$f_{\alpha}(x) = \frac{1}{2}x^2 + \frac{\alpha}{6}x^3 + \frac{1}{8}x^4.$$
 (59)

To apply Thm. 1, we denote  $\mathcal{L} = f_{\alpha}$ , and therefore

$$v = 1, \quad f_{\alpha}''(0) = 1, \quad f_{\alpha}^{(3)}(0) = \alpha, \quad f_{\alpha}^{(4)}(0) = 3.$$
 (60)

Then,

$$q = \left[ f_{\alpha}''(0) \right]^{-1} \frac{d}{dv} \left( f_{\alpha}^{(3)}(0)v^3 \right) \Big|_{v=1} = 3\alpha v^2 \Big|_{v=1} = 3\alpha.$$
 (61)

Overall,

$$\mathcal{D}^{3}\mathcal{L}(x^{*})[v, v, q] = f_{\alpha}^{(3)}(0)v^{2}q = \alpha \cdot 3\alpha = 3\alpha^{2}.$$
 (62)

$$\mathcal{D}^4 \mathcal{L}(x^*)[v, v, v, v] = f_{\alpha}^{(4)}(0)v^4 = 3.$$
(63)

Thus, the stability condition for oscillations is

$$\mathcal{D}^{3}\mathcal{L}(x^{*})[v,v,q] > \mathcal{D}^{4}\mathcal{L}(x^{*})[v,v,v,v] \quad \Longleftrightarrow \quad 3\alpha^{2} > 3 \quad \Longleftrightarrow \quad |\alpha| > 1. \tag{64}$$

## F PROOF OF THEOREM 2

Let  $\{\mathcal{B}_i\}_{i=1}^N$  be all possible different batches of size B form the dataset  $\{f_i\}_{i=1}^n$ , where  $N = \binom{n}{B}$ . Recall that  $\hat{\psi}_{\mathcal{B}}$  denotes the of GD transform on batch  $\mathcal{B}$ , *i.e.*, taking a single gradient step with respect to  $\hat{\mathcal{L}}_{\mathcal{B}}$  (see (19)). Moreover, let  $\hat{\psi}_{\mathcal{B}}^t$  denote the application of  $\hat{\psi}_{\mathcal{B}}$  for t times. Namely,

$$\hat{\psi}_{\mathcal{B}}^{t} = \underbrace{\hat{\psi}_{\mathcal{B}} \circ \cdots \circ \hat{\psi}_{\mathcal{B}} \circ \hat{\psi}_{\mathcal{B}}}_{t \text{ times}}.$$
(65)

For a stochastic batch  $\mathcal{B}_t$ ,  $\hat{\psi}_{\mathcal{B}_t}$  is distributed uniformly over  $\{\hat{\psi}_{\mathcal{B}_i}\}$ , *i.e.*, for any  $x \in \mathbb{R}^d$ 

$$\hat{\psi}_{\mathcal{B}_i}(\boldsymbol{x}) \sim \mathcal{U}\left(\left\{\hat{\psi}_{\mathcal{B}_i}(\boldsymbol{x})\right\}_{i=1}^N\right).$$
 (66)

Given an initial point  $x_0 \in \mathbb{R}^d$ , assume that for some batch  $\mathcal{B}_{i^*}$  (with index  $i^*$ ) GD's iterates, denoted by  $\{x_t^{(\mathcal{B}_{i^*})}\}$ , diverge with superlinear rate. That is

$$\|\boldsymbol{x}_{t}^{(\mathcal{B}_{i^{*}})} - \boldsymbol{x}^{*}\|^{\frac{1}{t}} \underset{t \to \infty}{\longrightarrow} \infty,$$
 (67)

Using our notation, we have  $x_t^{(\mathcal{B}_{i^*})} = \hat{\psi}_{\mathcal{B}_{i^*}}^t(x_0)$ . Let us look at the expectation of the distance between SGD iterates  $\{x_t\}$  form (10) and the minimizer  $x^*$ 

$$\mathbb{E}\left[\left\|\boldsymbol{x}_{t}-\boldsymbol{x}^{*}\right\|\right] = \frac{1}{N^{t}} \sum_{(i_{1},i_{2},...,i_{t})\in\{1,...,N\}^{t}} \left\|\hat{\psi}_{\mathcal{B}_{i_{t}}}\circ\cdots\circ\hat{\psi}_{\mathcal{B}_{i_{2}}}\circ\hat{\psi}_{\mathcal{B}_{i_{1}}}(\boldsymbol{x}_{0}) - \boldsymbol{x}^{*}\right\| \\
\geq \frac{1}{N^{t}} \left\|\hat{\psi}_{\mathcal{B}_{i_{t}}}\circ\cdots\circ\hat{\psi}_{\mathcal{B}_{i_{2}}}\circ\hat{\psi}_{\mathcal{B}_{i_{1}}}(\boldsymbol{x}_{0}) - \boldsymbol{x}^{*}\right\|_{i_{1}=i_{2}=\cdots=i_{t}=i^{*}} \\
= \frac{1}{N^{t}} \left\|\hat{\psi}_{\mathcal{B}_{i^{*}}}^{t}(\boldsymbol{x}_{0}) - \boldsymbol{x}^{*}\right\| \\
= \exp\left\{\log\left(\left\|\hat{\psi}_{\mathcal{B}_{i^{*}}}^{t}(\boldsymbol{x}_{0}) - \boldsymbol{x}^{*}\right\|\right) - t\log(N)\right\} \\
= \exp\left\{t\left[\frac{1}{t}\log\left(\left\|\hat{\psi}_{\mathcal{B}_{i^{*}}}^{t}(\boldsymbol{x}_{0}) - \boldsymbol{x}^{*}\right\|\right) - \log(N)\right]\right\} \xrightarrow{t\to\infty} \infty. \tag{68}$$

# G INTRODUCTION TO SPECTRAL ANALYSIS OF LINEAR OPERATORS

Let us start with the following definitions.

**Definition 3 (Operator norm)** Let A be a linear operator over a vector space V, then its operator norm is given by

$$\|A\| = \inf\{c > 0 : \|Av\| < c\|v\| \text{ for all } v \in V\}.$$
 (69)

**Definition 4 (Spectrum)** Let A be a linear operator over a Banach space V, then its spectrum is given by

$$\sigma(\mathbf{A}) = \{ \lambda \in \mathbb{C} : \mathbf{A} - \lambda \mathbf{I} \text{ is not bijective} \}, \tag{70}$$

where I is the identity operator.

**Definition 5 (Spectral radius)** The spectral radius of an operator **A** is given by

$$r(\mathbf{A}) = \sup_{\lambda \in \sigma(\mathbf{A})} |\lambda|. \tag{71}$$

Consider the following linear system

$$\mu_{t+1} = A\mu_t,\tag{72}$$

where A is a bounded linear operator. We want to have some condition such that the iterates are bounded or converging. Here, we can unfold the equation to get an explicit formula for any  $\mu_t$  as

$$\mu_t = A^t \mu_0. \tag{73}$$

 A naive way to ensure convergence, i.e.,  $\mu_t \to 0$  as  $t \to \infty$ , is by taking the operator norm of A to be less than one, i.e., ||A|| < 1. Then

$$\|\boldsymbol{\mu}_t\| = \|\boldsymbol{A}^t \boldsymbol{\mu}_0\| \le \|\boldsymbol{A}\|^t \|\boldsymbol{\mu}_0\| \to 0.$$
 (74)

However, this is quite restrictive and will give us a loss condition. Note that we only like to know if  $\|A^t\|$  is bounded or shrinks to zero. In special cases, we can easily compute  $A^t$ . For example, in the finite-dimensional case, where  $A = PDP^{-1}$  is diagonalizable (e.g., symmetric or normal). Then,

$$\mu_t = A^t \mu_0 = (PDP^{-1})^t \mu_0 = PD^t P^{-1} \mu_0.$$
 (75)

Here, the system is stable if and only if the spectral radius of A is less or equal to one. If it is strictly less than one, then  $\mu_t \to 0$ .

In the general case, we can use Gelfand's formula for bounded linear operators on Banach spaces. Let r(A) denote the spectral radius of A, then Gelfand's formula is

$$r(\boldsymbol{A}) = \lim_{t \to \infty} \|\boldsymbol{A}^t\|^{\frac{1}{t}} = \inf_{t \in \mathbb{N}} \|\boldsymbol{A}^t\|^{\frac{1}{t}}.$$
 (76)

From this formula, we can see that if r(A) < 1, then  $\mu_t \to 0$ .

## H COMPUTATION OF THE OPERATOR BLOCKS

In this section we give the missing steps from (24). To this end, denote

$$\mathbf{Y}_{t,k} = \frac{1}{k!} \mathcal{D}^k \hat{\mathbf{\psi}}_{\mathcal{B}_t}(\mathbf{x}^*) \in \mathbb{R}^{d \times d^k}. \tag{77}$$

Then, the evolution over time of the kth moment is

$$\mathbb{E}\left[\left(\sum_{p=1}^{\infty} \frac{1}{p!} \mathcal{D}^{p} \hat{\psi}_{\mathcal{B}_{t}}(\boldsymbol{x}^{*}) \Delta \boldsymbol{x}_{t}^{p}\right)^{\otimes k}\right]$$

$$= \mathbb{E}\left[\left(\sum_{p=1}^{\infty} \boldsymbol{Y}_{t,p} \Delta \boldsymbol{x}_{t}^{p}\right)^{\otimes k}\right]$$

$$= \mathbb{E}\left[\sum_{p=k}^{\infty} \sum_{\substack{1 \leq \kappa_{1}, \dots, \kappa_{k} \leq p-k+1 \\ \kappa_{1}+\kappa_{2}+\dots+\kappa_{k}=p}} (\boldsymbol{Y}_{t,\kappa_{1}} \Delta \boldsymbol{x}_{t}^{\kappa_{1}}) \otimes (\boldsymbol{Y}_{t,\kappa_{2}} \Delta \boldsymbol{x}_{t}^{\kappa_{2}}) \otimes \dots \otimes (\boldsymbol{Y}_{t,\kappa_{k}} \Delta \boldsymbol{x}_{t}^{\kappa_{k}})\right]$$

$$= \mathbb{E}\left[\sum_{p=k}^{\infty} \sum_{\substack{1 \leq \kappa_{1}, \dots, \kappa_{k} \leq p-k+1 \\ \kappa_{1}+\kappa_{2}+\dots+\kappa_{k}=p}} (\boldsymbol{Y}_{t,\kappa_{1}} \otimes \dots \otimes \boldsymbol{Y}_{t,\kappa_{k}}) (\Delta \boldsymbol{x}_{t}^{\kappa_{1}} \otimes \dots \otimes \Delta \boldsymbol{x}_{t}^{\kappa_{k}})\right]$$

$$= \mathbb{E}\left[\sum_{p=k}^{\infty} \sum_{\substack{1 \leq \kappa_{1}, \dots, \kappa_{k} \leq p-k+1 \\ \kappa_{1}+\kappa_{2}+\dots+\kappa_{k}=p}} (\boldsymbol{Y}_{t,\kappa_{1}} \otimes \boldsymbol{Y}_{t,\kappa_{2}} \otimes \dots \otimes \boldsymbol{Y}_{t,\kappa_{k}}) (\Delta \boldsymbol{x}_{t}^{\sum_{i=1}^{k} \kappa_{i}})\right]$$

$$= \mathbb{E}\left[\sum_{p=k}^{\infty} \sum_{\substack{1 \leq \kappa_{1}, \dots, \kappa_{k} \leq p-k+1 \\ \kappa_{1}+\kappa_{2}+\dots+\kappa_{k}=p}} (\boldsymbol{Y}_{t,\kappa_{1}} \otimes \boldsymbol{Y}_{t,\kappa_{2}} \otimes \dots \otimes \boldsymbol{Y}_{t,\kappa_{k}}) \Delta \boldsymbol{x}_{t}^{p}\right]$$

$$= \mathbb{E}\left[\sum_{p=k}^{\infty} \left(\sum_{\substack{1 \leq \kappa_{1}, \dots, \kappa_{k} \leq p-k+1 \\ \kappa_{1}+\kappa_{2}+\dots+\kappa_{k}=p}}} \boldsymbol{Y}_{t,\kappa_{1}} \otimes \boldsymbol{Y}_{t,\kappa_{2}} \otimes \dots \otimes \boldsymbol{Y}_{t,\kappa_{k}}\right) \Delta \boldsymbol{x}_{t}^{p}\right]$$

$$= \sum_{p=k}^{\infty} \mathbb{E} \left[ \sum_{\substack{1 \leq \kappa_{1}, \dots, \kappa_{k} \leq p-k+1 \\ \kappa_{1}+\kappa_{2}+\dots+\kappa_{k}=p}} \mathbf{Y}_{t,\kappa_{1}} \otimes \mathbf{Y}_{t,\kappa_{2}} \otimes \dots \otimes \mathbf{Y}_{t,\kappa_{k}} \right] \mathbb{E} \left[ \Delta \mathbf{x}_{t}^{p} \right]$$

$$= \sum_{p=k}^{\infty} \mathbf{\Psi}_{k,p} \mathbb{E} \left[ \Delta \mathbf{x}_{t}^{p} \right], \tag{78}$$

where

$$\Psi_{k,p} = \mathbb{E}\left[\sum_{\substack{1 \leq \kappa_1, \cdots, \kappa_k \leq p - k + 1 \\ \kappa_1 + \kappa_2 + \cdots + \kappa_k = p}} Y_{t,\kappa_1} \otimes Y_{t,\kappa_2} \otimes \cdots \otimes Y_{t,\kappa_k}\right] \in \mathbb{R}^{d^k \times d^p}, \tag{79}$$

# I BOUNDING THE OPERATOR

In this section, we assume that the condition of Thm. 3 holds. Then, we show that there exists a  $\rho > 0$  such that  $\Psi_{\rho}$  is bounded. For this, we use the following result (see proof in App. L).

**Theorem 5** Let T be an operator defined on  $\ell_2$  space. Denote by  $\{T_{i,j}\}$  a division of T into blocks, such that  $\forall i, j \ T_{i,j} \in \mathbb{R}^{d_i \times d_j}$  where  $\{d_i\}_{i=1}^{\infty}$  is a some sequence. Assume that

$$\forall j \in \mathbb{N}$$
  $\sum_{i=1}^{\infty} ||T_{i,j}|| \le \alpha$  and  $\forall i \in \mathbb{N}$   $\sum_{j=1}^{\infty} ||T_{i,j}|| \le \beta$ . (80)

Then T is a bounded linear operator and

$$||T||_2 \le \sqrt{\alpha \beta}. \tag{81}$$

Let us apply Thm. 5 to  $\Psi_{\rho}$ . Using the definition of the blocks  $\{\Psi_{k,p}\}$  given in (79) we have

$$\|\mathbf{\Psi}_{k,p}\| = \left\| \mathbb{E} \left[ \sum_{\substack{1 \leq \kappa_{1}, \cdots, \kappa_{k} \leq p-k+1 \\ \kappa_{1}+\kappa_{2}+\cdots+\kappa_{k}=p}} \mathbf{Y}_{t,\kappa_{1}} \otimes \mathbf{Y}_{t,\kappa_{2}} \otimes \cdots \otimes \mathbf{Y}_{t,\kappa_{k}} \right] \right\|$$

$$\leq \mathbb{E} \left[ \sum_{\substack{1 \leq \kappa_{1}, \cdots, \kappa_{k} \leq p-k+1 \\ \kappa_{1}+\kappa_{2}+\cdots+\kappa_{k}=p}} \|\mathbf{Y}_{t,\kappa_{1}} \otimes \mathbf{Y}_{t,\kappa_{2}} \otimes \cdots \otimes \mathbf{Y}_{t,\kappa_{k}} \| \right]$$

$$= \mathbb{E} \left[ \sum_{\substack{1 \leq \kappa_{1}, \cdots, \kappa_{k} \leq p-k+1 \\ \kappa_{1}+\kappa_{2}+\cdots+\kappa_{k}=p}} \|\mathbf{Y}_{t,\kappa_{1}}\| \|\mathbf{Y}_{t,\kappa_{2}}\| \cdots \|\mathbf{Y}_{t,\kappa_{k}}\| \right], \tag{82}$$

where  $Y_{t,p}$  is given in (77). Let  $\{\mathcal{B}_m\}_{m=1}^N$  be all possible different batches of size B form the dataset  $\{f_m\}_{m=1}^n$ , where  $N=\binom{n}{B}$ . Since  $\{f_m\}$  are analytic and  $\{\hat{\mathcal{L}}_{\mathcal{B}_m}\}$  are finite sum losses, then also  $\{\hat{\psi}_{\mathcal{B}_m}\}$  are analytic. Then, using Gevrey class theory, for each batch  $\mathcal{B}_m$  there exists  $C_m>0$  such that

$$\max_{i,j} \left| \left[ \mathcal{D}^p \hat{\psi}_{\mathcal{B}_m}(\boldsymbol{x}^*) \right]_{i,j} \right| \le C_m^{p+1} p! \qquad \forall p \ge 1,$$
(83)

where  $[\mathcal{D}^p \hat{\psi}_{\mathcal{B}_m}(\boldsymbol{x}^*)]_{i,j}$  are the elements of in the matrix  $\mathcal{D}^p \hat{\psi}_{\mathcal{B}_m}(\boldsymbol{x}^*)$ , which are all the (mixed) partial derivatives of degree p. Setting

$$C = \max_{m \in [N]} C_m, \tag{84}$$

we get for a random batch  $\mathcal{B}_t$ 

$$\max_{i,j} \left| \left[ \mathcal{D}^p \hat{\psi}_{\mathcal{B}_t}(\boldsymbol{x}^*) \right]_{i,j} \right| \le C^{p+1} p! \quad \text{w.p. } 1 \qquad \forall p \ge 1.$$
 (85)

Now that we have a uniform bound on all elements in the matrix, we can bound its norm. Specifically, it is well known that for a matrix  $A \in \mathbb{R}^{m \times n}$  with elements  $|A_{i,j}| \leq M$ , we have  $||A|| \leq M \sqrt{mn}$  (a simple application of Thm. 5 can give this result as well). Using this, and the fact that  $\mathcal{D}^p \hat{\psi}_{\mathcal{B}_*}(\mathbf{x}^*) \in \mathbb{R}^{d \times d^p}$  we get

$$\|\mathbf{Y}_{t,p}\| = \frac{1}{p!} \|\mathcal{D}^p \hat{\psi}_{\mathcal{B}_t}(\mathbf{x}^*)\| \le C^{p+1} d^{\frac{p+1}{2}} \quad \text{w.p. } 1.$$
 (86)

Define

$$Q_{t,p} = \begin{cases} \|\mathbf{Y}_{t,1}\|, & p = 1, \\ C^{p+1} d^{\frac{p+1}{2}}, & \text{otherwise.} \end{cases}$$
 (87)

Then for all  $p \geq 1$  and  $t \in \mathbb{N}$ 

$$\|Y_{t,p}\| \le Q_{t,p}$$
 w.p. 1, (88)

and

$$\|\Psi_{k,p}\| \le \mathbb{E} \left[ \sum_{\substack{1 \le \kappa_1, \dots, \kappa_k \le p-k+1\\ \kappa_1 + \kappa_2 + \dots + \kappa_k = p}} Q_{t,\kappa_1} Q_{t,\kappa_2} \cdots Q_{t,\kappa_k} \right].$$
 (89)

Let us apply Thm. 5 on  $\Psi_{\rho}$ , while assuming that  $\rho < \frac{1}{C\sqrt{d}}$ . For the sum of the row block we have

$$\sum_{p=k}^{\infty} \rho^{p-k} \| \Psi_{k,p} \| \leq \sum_{p=k}^{\infty} \rho^{p-k} \mathbb{E} \left[ \sum_{\substack{1 \leq \kappa_1, \cdots, \kappa_k \leq p-k+1 \\ \kappa_1 + \kappa_2 + \cdots + \kappa_k = p}} Q_{t,\kappa_1} Q_{t,\kappa_2} \cdots Q_{t,\kappa_k} \right]$$

$$= \rho^{-k} \mathbb{E} \left[ \sum_{p=k}^{\infty} \left( \sum_{\substack{1 \leq \kappa_1, \cdots, \kappa_k \leq p-k+1 \\ \kappa_1 + \kappa_2 + \cdots + \kappa_k = p}} Q_{t,\kappa_1} Q_{t,\kappa_2} \cdots Q_{t,\kappa_k} \right) \rho^p \right]$$

$$= \rho^{-k} \mathbb{E} \left[ \left( \sum_{p=1}^{\infty} Q_{t,p} \rho^p \right)^k \right]$$

$$= \rho^{-k} \mathbb{E} \left[ \left( \| Y_{t,1} \| \rho + \sum_{p=2}^{\infty} C^{p+1} d^{\frac{p+1}{2}} \rho^p \right)^k \right]$$

$$= \rho^{-k} \mathbb{E} \left[ \left( \| Y_{t,1} \| \rho + C \sqrt{d} \sum_{p=2}^{\infty} \left( C \sqrt{d} \rho \right)^p \right)^k \right]$$

$$= \rho^{-k} \mathbb{E} \left[ \left( \| Y_{t,1} \| \rho + C \sqrt{d} \frac{C^2 d \rho^2}{1 - C \sqrt{d} \rho} \right)^k \right]$$

$$= \mathbb{E} \left[ \left( \| Y_{t,1} \| + \frac{C^3 d^{3/2} \rho}{1 - C \sqrt{d} \rho} \right)^k \right], \tag{90}$$

where in the sixth step we used

$$\sum_{p=2}^{\infty} q^p = \frac{q^2}{1-q},\tag{91}$$

for  $0 < q = C\sqrt{d}\rho < 1$ . Here we assume the condition in (13) holds and that  $\{\nabla^2 \hat{\mathcal{L}}_{\mathcal{B}_i}(\boldsymbol{x}^*)\}$  have full rank. This means that for every batch  $\mathcal{B}_i$ 

$$0 < \eta \lambda_{\min} \left( \nabla^2 \hat{\mathcal{L}}_{\mathcal{B}_i}(\boldsymbol{x}^*) \right) < \eta \lambda_{\max} \left( \nabla^2 \hat{\mathcal{L}}_{\mathcal{B}_i}(\boldsymbol{x}^*) \right) < 2.$$
 (92)

Recall that  $\mathcal{D}\hat{\psi}_{\mathcal{B}_i}(\boldsymbol{x}^*) = \boldsymbol{I} - \eta \nabla^2 \hat{\mathcal{L}}_{\mathcal{B}_i}(\boldsymbol{x}^*)$ . Then it is easy to show that there exists  $\varepsilon \in (0,1)$  such that

$$\max_{i \in [N]} \left\| \mathcal{D} \hat{\psi}_{\mathcal{B}_{i}}(\boldsymbol{x}^{*}) \right\| = \max_{i \in [N]} \left\| \boldsymbol{I} - \eta \nabla^{2} \hat{\mathcal{L}}_{\mathcal{B}_{i}}(\boldsymbol{x}^{*}) \right\|$$

$$= \max_{i \in [N]} \left\{ \max \left\{ 1 - \eta \lambda_{\min} \left( \nabla^{2} \hat{\mathcal{L}}_{\mathcal{B}_{i}}(\boldsymbol{x}^{*}) \right), \eta \lambda_{\max} \left( \nabla^{2} \hat{\mathcal{L}}_{\mathcal{B}_{i}}(\boldsymbol{x}^{*}) \right) - 1 \right\} \right\}$$

$$= 1 - \varepsilon. \tag{93}$$

This meaning that

$$\|\mathbf{Y}_{t,1}\| \le \max_{i \in [N]} \left\| \mathcal{D}\hat{\psi}_{\mathcal{B}_i}(\mathbf{x}^*) \right\| = 1 - \varepsilon \quad \text{w.p. } 1.$$
 (94)

Note that  $\frac{C^3 d^{3/2} \rho}{1 - C \sqrt{d} \rho} \ge 0$ , then we can further bound (90) by

$$\mathbb{E}\left[\left(\|\boldsymbol{Y}_{t,1}\| + \frac{C^3 d^{3/2} \rho}{1 - C\sqrt{d}\rho}\right)^k\right] \le \left(1 - \varepsilon + \frac{C^3 d^{3/2} \rho}{1 - C\sqrt{d}\rho}\right)^k. \tag{95}$$

In order for this to be bounded for any  $k \in \mathbb{N}$ , we will require

$$1 - \varepsilon + \frac{C^3 d^{3/2} \rho}{1 - C\sqrt{d}\rho} < 1$$

$$\Leftrightarrow \frac{C^3 d^{3/2} \rho}{1 - C\sqrt{d}\rho} < \varepsilon$$

$$\Leftrightarrow C^3 d^{3/2} \rho + \varepsilon C\sqrt{d}\rho < \varepsilon$$

$$\Leftrightarrow \rho < \frac{\varepsilon}{C^3 d^{3/2} + \varepsilon C\sqrt{d}} \triangleq \rho^*. \tag{96}$$

Therefore, under the condition of  $\rho < \rho^*$  there exists  $\gamma \in (0,1)$  (for example,  $\gamma = 1 - \varepsilon + \frac{C^3 d^{3/2} \rho}{1 - C \sqrt{d} \rho}$ ) such that

$$\sum_{p=k}^{\infty} \rho^{p-k} \| \mathbf{\Psi}_{k,p} \| \le \gamma^k. \tag{97}$$

This means that the rows sum is uniformly bounded. Now, for the column sums, under the same assumptions we get

$$\sup_{p} \sum_{k=1}^{p} \rho^{p-k} \| \Psi_{k,p} \| \leq \sum_{p=1}^{\infty} \sum_{k=1}^{p} \rho^{p-k} \| \Psi_{k,p} \| 
= \sum_{k=1}^{\infty} \sum_{p=k}^{\infty} \rho^{p-k} \| \Psi_{k,p} \| 
\leq \sum_{k=1}^{\infty} \gamma^{k} 
= \frac{\gamma}{1-\gamma},$$
(98)

where the change of summation order in second step is justified since all elements are non-negative. Hence, under the same assumptions, the absolute columns sum is also uniformly bounded.

Overall, we see that the conditions of Thm. 3 are sufficient to find a neighborhood around the minimum,  $\{x_0 : \|x_0 - x^*\| < \rho\}$ , such that the operator  $\Psi_\rho$  is bounded. For completeness, in App. J, we show that the condition in (13) is also necessary. Namely, if this condition is violated,  $\Psi_\rho$  is not bounded.

## J NECESSARY CONDITION FOR BOUNDNESS

In this section we show that the condition in (13) is also necessary to bound the operator  $\Psi_{\rho}$ . We bring this only to give a complete theoretical understanding, yet we do not use this derivation to prove our results.

For  $\Psi_{\rho}$  to be bounded, all of its submatrices must be bounded. Note that the diagonal blocks of this operator  $\{\Psi_{k,k}\}_{k=1}^{\infty}$  are independent of  $\rho$ . Therefore, we should have a condition, independent of  $\rho$ , for these submatrices to be bounded. These blocks are given by

$$\mathbf{\Psi}_{k,k} = \mathbb{E}\left[\left(\mathcal{D}\hat{\psi}_{\mathcal{B}_t}(\mathbf{x}^*)\right)^{\otimes k}\right]. \tag{99}$$

Note that

$$\mathcal{D}\hat{\psi}_{\mathcal{B}_t}(\boldsymbol{x}^*) = \mathcal{D}\left(\boldsymbol{x} - \eta \nabla \hat{\mathcal{L}}_{\mathcal{B}_t}(\boldsymbol{x})\right)\Big|_{\boldsymbol{x} = \boldsymbol{x}^*} = \boldsymbol{I} - \eta \nabla^2 \hat{\mathcal{L}}_{\mathcal{B}_t}(\boldsymbol{x}^*). \tag{100}$$

For ease of reading and better interpretability, let

$$\boldsymbol{H}_{\mathcal{B}} \triangleq \nabla^2 \hat{\mathcal{L}}_{\mathcal{B}}(\boldsymbol{x}^*) \tag{101}$$

denote the Hessian of the batch  $\mathcal{B}$ . Then

$$\Psi_{k,k} = \mathbb{E}\left[ \left( \mathbf{I} - \eta \mathbf{H}_{\mathcal{B}_t} \right)^{\otimes k} \right]. \tag{102}$$

Moreover, denote by  $H_{\rm max}$  the batch that has the largest maximal eigenvalue, that is

$$\boldsymbol{H}_{\max} = \underset{\boldsymbol{\mathcal{B}}: |\boldsymbol{\mathcal{B}}| = B}{\operatorname{arg} \max} \left\{ \lambda_{\max} (\boldsymbol{H}_{\mathcal{B}}) \right\}, \tag{103}$$

and by  $v_{\max}$  its corresponding eigenvector (normalized). Note that  $\Psi_{k,k}$  is symmetric for all  $k \in \mathbb{N}$ , therefore

$$\|\mathbf{\Psi}_{k,k}\| = \max_{\mathbf{u} \in \mathbb{R}^{d^k} : \|\mathbf{u}\| = 1} |\mathbf{u}^{\mathrm{T}} \mathbf{\Psi}_{k,k} \mathbf{u}|.$$
 (104)

Since  $\|\boldsymbol{v}_{\text{max}}^{\otimes k}\| = \|\boldsymbol{v}_{\text{max}}\|^k = 1$ , we have that

$$\|\boldsymbol{\Psi}_{k,k}\| \geq \left| \left(\boldsymbol{v}_{\max}^{\otimes k}\right)^{\mathrm{T}} \boldsymbol{\Psi}_{k,k} \boldsymbol{v}_{\max}^{\otimes k} \right|$$

$$= \left| \left(\boldsymbol{v}_{\max}^{\otimes k}\right)^{\mathrm{T}} \mathbb{E} \left[ \left(\boldsymbol{I} - \eta \boldsymbol{H}_{\mathcal{B}_{t}}\right)^{\otimes k} \right] \boldsymbol{v}_{\max}^{\otimes k} \right|$$

$$= \left| \mathbb{E} \left[ \left(\boldsymbol{v}_{\max}^{\otimes k}\right)^{\mathrm{T}} \left(\boldsymbol{I} - \eta \boldsymbol{H}_{\mathcal{B}_{t}}\right)^{\otimes k} \boldsymbol{v}_{\max}^{\otimes k} \right] \right|$$

$$= \left| \mathbb{E} \left[ \left(1 - \eta \boldsymbol{v}_{\max}^{\mathrm{T}} \boldsymbol{H}_{\mathcal{B}_{t}} \boldsymbol{v}_{\max}\right)^{k} \right] \right|. \tag{105}$$

Assume that

$$\eta > \frac{2}{\lambda_{\text{max}}(\boldsymbol{H}_{\text{max}})}.$$
(106)

Since  $\lambda_{\max}(H_{\max}) = v_{\max}^{\mathrm{T}} H_{\max} v_{\max}$ , under the assumption above, we have that

$$\mathbb{P}\left(\eta \boldsymbol{v}_{\max}^{\mathrm{T}} \boldsymbol{H}_{\mathcal{B}_{t}} \boldsymbol{v}_{\max} > 2\right) > 0. \tag{107}$$

Therefore, continuing from (105)

$$\|\boldsymbol{\Psi}_{k,k}\| \geq \left| \mathbb{E} \left[ \left( 1 - \eta \boldsymbol{v}_{\text{max}}^{\text{T}} \boldsymbol{H}_{\mathcal{B}_{t}} \boldsymbol{v}_{\text{max}} \right)^{k} \right] \right|$$

$$= \left| \mathbb{P} \left( \eta \boldsymbol{v}_{\text{max}}^{\text{T}} \boldsymbol{H}_{\mathcal{B}_{t}} \boldsymbol{v}_{\text{max}} > 2 \right) \mathbb{E} \left[ \left( 1 - \eta \boldsymbol{v}_{\text{max}}^{\text{T}} \boldsymbol{H}_{\mathcal{B}_{t}} \boldsymbol{v}_{\text{max}} \right)^{k} \middle| \eta \boldsymbol{v}_{\text{max}}^{\text{T}} \boldsymbol{H}_{\mathcal{B}_{t}} \boldsymbol{v}_{\text{max}} > 2 \right] \right.$$

$$+ \mathbb{P} \left( \eta \boldsymbol{v}_{\text{max}}^{\text{T}} \boldsymbol{H}_{\mathcal{B}_{t}} \boldsymbol{v}_{\text{max}} \leq 2 \right) \mathbb{E} \left[ \left( 1 - \eta \boldsymbol{v}_{\text{max}}^{\text{T}} \boldsymbol{H}_{\mathcal{B}_{t}} \boldsymbol{v}_{\text{max}} \right)^{k} \middle| \eta \boldsymbol{v}_{\text{max}}^{\text{T}} \boldsymbol{H}_{\mathcal{B}_{t}} \boldsymbol{v}_{\text{max}} \leq 2 \right] \right|$$

$$\geq \mathbb{P} \left( \eta \boldsymbol{v}_{\text{max}}^{\text{T}} \boldsymbol{H}_{\mathcal{B}_{t}} \boldsymbol{v}_{\text{max}} > 2 \right) \mathbb{E} \left[ \left( \eta \boldsymbol{v}_{\text{max}}^{\text{T}} \boldsymbol{H}_{\mathcal{B}_{t}} \boldsymbol{v}_{\text{max}} - 1 \right)^{k} \middle| \eta \boldsymbol{v}_{\text{max}}^{\text{T}} \boldsymbol{H}_{\mathcal{B}_{t}} \boldsymbol{v}_{\text{max}} > 2 \right]$$

$$- \mathbb{P} \left( \eta \boldsymbol{v}_{\text{max}}^{\text{T}} \boldsymbol{H}_{\mathcal{B}_{t}} \boldsymbol{v}_{\text{max}} \leq 2 \right) \middle| \mathbb{E} \left[ \left( 1 - \eta \boldsymbol{v}_{\text{max}}^{\text{T}} \boldsymbol{H}_{\mathcal{B}_{t}} \boldsymbol{v}_{\text{max}} \right)^{k} \middle| \eta \boldsymbol{v}_{\text{max}}^{\text{T}} \boldsymbol{H}_{\mathcal{B}_{t}} \boldsymbol{v}_{\text{max}} \leq 2 \right] \middle|,$$

$$(108)$$

where in the second step we used the law of total expectation, and in last step we used the triangle inequality. Since  $x^*$  is an interpolating minimum, then  $H_{\mathcal{B}_t}$  is PSD w.p. one, and  $0 \leq v_{\max}^T H_{\mathcal{B}_t} v_{\max}$ . Thus,

$$\left| \mathbb{E} \left[ \left( 1 - \eta \boldsymbol{v}_{\text{max}}^{\text{T}} \boldsymbol{H}_{\mathcal{B}_{t}} \boldsymbol{v}_{\text{max}} \right)^{k} \middle| \eta \boldsymbol{v}_{\text{max}}^{\text{T}} \boldsymbol{H}_{\mathcal{B}_{t}} \boldsymbol{v}_{\text{max}} \leq 2 \right] \right| \leq 1.$$
 (109)

However,

$$\mathbb{E}\left[\left(\eta \boldsymbol{v}_{\max}^{\mathrm{T}} \boldsymbol{H}_{\mathcal{B}_{t}} \boldsymbol{v}_{\max} - 1\right)^{k} \middle| \eta \boldsymbol{v}_{\max}^{\mathrm{T}} \boldsymbol{H}_{\mathcal{B}_{t}} \boldsymbol{v}_{\max} > 2\right] \underset{k \to \infty}{\longrightarrow} \infty. \tag{110}$$

This means that under the condition in (106), we have that  $\{\Psi_{k,k}\}$  are unbounded. Therefore, a necessary condition for boundness is

$$\eta \le \frac{2}{\lambda_{\max}(\boldsymbol{H}_{\max})}.\tag{111}$$

## K SPECTRAL ANALYSIS

In App.I we proved that, under the condition in (13) of Thm. 3, we can find a neighborhood  $\|x_0 - x^*\| < \rho$  such that the operator  $\Psi_\rho$  is bounded. In this section we show that under the same condition, the spectral radius of  $\Psi_\rho$ , denoted by  $r(\Psi_\rho)$ , is less than one. To do this, we first show that the operator is compact, which means that all the non-zero elements in its spectrum are eigenvalues (point spectrum). For this end, we define the following sequence of finite rank approximations (truncations)  $\{\Psi_\rho^k\}$ , comprised of the first  $k \times k$  blocks of  $\Psi_\rho$ . Namely,

$$\Psi_{\rho}^{k} = \begin{bmatrix}
\Psi_{1,1} & \rho \Psi_{1,2} & \rho^{2} \Psi_{1,3} & \cdots & \rho^{k-1} \Psi_{1,k} & \mathbf{0} & \cdots \\
\mathbf{0} & \Psi_{2,2} & \rho \Psi_{2,3} & \cdots & \rho^{k-2} \Psi_{2,k} & \mathbf{0} & \cdots \\
\mathbf{0} & \mathbf{0} & \Psi_{3,3} & \cdots & \rho^{k-3} \Psi_{3,k} & \mathbf{0} & \cdots \\
\vdots & \vdots & \vdots & \ddots & \vdots & \mathbf{0} & \cdots \\
\mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \Psi_{k,k} & \mathbf{0} & \cdots \\
\mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \cdots \\
\vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \ddots
\end{bmatrix} .$$
(112)

Furthermore, define  $\tilde{\Psi}_{i,j}$  as the embedding of the block  $\Psi_{i,j}$  to the full space, *i.e.* 

$$\tilde{\Psi}_{i,j} = \begin{bmatrix}
0 & \cdots & 0 & 0 & 0 & \cdots \\
\vdots & \ddots & \vdots & \vdots & \vdots & \cdots \\
0 & \cdots & 0 & 0 & 0 & \cdots \\
0 & \cdots & 0 & \Psi_{i,j} & 0 & \cdots \\
0 & \cdots & 0 & 0 & 0 & \cdots \\
\vdots & \vdots & \vdots & \vdots & \vdots & \ddots
\end{bmatrix},$$
(113)

such that

$$\Psi_{\rho} = \sum_{i=1}^{\infty} \sum_{j=i}^{\infty} \rho^{j-i} \tilde{\Psi}_{i,j} \quad \text{and} \quad \Psi_{\rho}^{k} = \sum_{i=1}^{k} \sum_{j=i}^{k} \rho^{j-i} \tilde{\Psi}_{i,j}. \quad (114)$$

Then,

$$\begin{aligned} \|\boldsymbol{\Psi}_{\rho} - \boldsymbol{\Psi}_{\rho}^{k}\| &= \left\| \sum_{i=1}^{\infty} \sum_{j=i}^{\infty} \rho^{j-i} \tilde{\boldsymbol{\Psi}}_{i,j} - \sum_{i=1}^{k} \sum_{j=i}^{k} \rho^{j-i} \tilde{\boldsymbol{\Psi}}_{i,j} \right\| \\ &= \left\| \sum_{i=k+1}^{\infty} \sum_{j=i}^{\infty} \rho^{j-i} \tilde{\boldsymbol{\Psi}}_{i,j} + \sum_{i=1}^{k} \sum_{j=k+1}^{\infty} \rho^{j-i} \tilde{\boldsymbol{\Psi}}_{i,j} \right\| \\ &\leq \sum_{i=k+1}^{\infty} \sum_{j=i}^{\infty} \rho^{j-i} \left\| \tilde{\boldsymbol{\Psi}}_{i,j} \right\| + \sum_{i=1}^{k} \sum_{j=k+1}^{\infty} \rho^{j-i} \left\| \tilde{\boldsymbol{\Psi}}_{i,j} \right\| \end{aligned}$$

$$= \sum_{i=k+1}^{\infty} \sum_{j=i}^{\infty} \rho^{j-i} \| \mathbf{\Psi}_{i,j} \| + \sum_{i=1}^{k} \sum_{j=k+1}^{\infty} \rho^{j-i} \| \mathbf{\Psi}_{i,j} \|$$

$$= \sum_{i=1}^{\infty} \sum_{j=i}^{\infty} \rho^{j-i} \| \mathbf{\Psi}_{i,j} \| - \sum_{i=1}^{k} \sum_{j=1}^{k} \rho^{j-i} \| \mathbf{\Psi}_{i,j} \| \xrightarrow[k \to \infty]{} 0,$$
(115)

where in the third step we used the fact that

$$\sum_{i=1}^{\infty} \sum_{j=i}^{\infty} \rho^{j-i} \left\| \tilde{\Psi}_{i,j} \right\| = \sum_{i=1}^{\infty} \sum_{j=i}^{\infty} \rho^{j-i} \left\| \Psi_{i,j} \right\|$$
(116)

is bounded (see (98)). Therefore,  $\Psi_{\rho}^{k} \xrightarrow{\|\cdot\|} \Psi_{\rho}$  as  $k \to \infty$ , and thus  $\Psi_{\rho}$  is compact. This means that the non-zero elements in the spectrum of  $\Psi_{\rho}$  are comprised of its eigenvalues only (point spectrum).

In the following we use a known result about the convergence of the spectrum of finite rank approximations.

**Lemma 1 (Dunford & Schwartz (1964, Cp. XI.9 Lemma 5))** Let  $\{T_k\}$  and T be compact operators, such that  $T_k \stackrel{\|\cdot\|}{\longrightarrow} T$ . Let  $\lambda_m(T)$  be an enumeration of the non-zero eigenvalues of T, each repeated according to its multiplicity. Then there exist enumerations  $\lambda_m(T_k)$  of the non-zero eigenvalues of  $\{T_k\}$ , with the repetitions according to multiplicity, such that

$$\lim_{k \to \infty} \lambda_m(T_k) = \lambda_m(T), \qquad m \ge 1, \tag{117}$$

where the limit is uniform in m.

Let  $\sigma(\cdot)$  denote the spectrum of an operator. Here, each  $\Psi_{\rho}^{k}$ , when restricted to its square support, is a block upper triangular matrix. Hence, its spectrum<sup>3</sup> is given by the union of the eigenvalues of the blocks in the diagonal. Namely,

$$\sigma\left(\mathbf{\Psi}_{\rho}^{k}\right) = \bigcup_{j=1}^{k} \sigma\left(\mathbf{\Psi}_{j,j}\right). \tag{118}$$

Thus, according to Lemma 1 we have that the non-zero spectrum of  $\Psi_{
ho}$  is

$$\sigma\left(\mathbf{\Psi}_{\rho}\right)\backslash\{0\} = \lim_{k\to\infty}\sigma\left(\mathbf{\Psi}_{\rho}^{k}\right)\backslash\{0\} = \bigcup_{k=1}^{\infty}\sigma\left(\mathbf{\Psi}_{k,k}\right)\backslash\{0\}. \tag{119}$$

Now that we have the spectrum of  $\Psi_{\rho}$  we turn to show that under the condition of Thm. 3 in (13), the spectral radius  $r(\Psi_{\rho})$  is less than one. Due to (119), it is sufficient to show that for all  $k \in \mathbb{N}$  we have  $r(\Psi_{k,k}) \leq c < 1$ , for some constant  $c \in (0,1)$ . Recall that

$$\Psi_{k,k} = \mathbb{E}\left[\left(\mathcal{D}\hat{\psi}_{\mathcal{B}_t}(\boldsymbol{x}^*)\right)^{\otimes k}\right],\tag{120}$$

where  $\mathcal{D}\hat{\psi}_{\mathcal{B}_t}(\boldsymbol{x}^*) = \boldsymbol{I} - \eta \nabla^2 \hat{\mathcal{L}}_{\mathcal{B}_t}(\boldsymbol{x}^*)$  is symmetric. Therefore,  $\boldsymbol{\Psi}_{k,k}$  is also symmetric, and we have that  $r(\boldsymbol{\Psi}_{k,k}) = \|\boldsymbol{\Psi}_{k,k}\|$ . Thus, using Jensen's inequality

$$r\left(\mathbf{\Psi}_{k,k}\right) = \|\mathbf{\Psi}_{k,k}\|$$

$$= \left\|\mathbb{E}\left[\left(\mathcal{D}\hat{\psi}_{\mathcal{B}_{t}}(\mathbf{x}^{*})\right)^{\otimes k}\right]\right\|$$

$$\leq \mathbb{E}\left[\left\|\left(\mathcal{D}\hat{\psi}_{\mathcal{B}_{t}}(\mathbf{x}^{*})\right)^{\otimes k}\right\|\right]$$

$$= \mathbb{E}\left[\left\|\mathcal{D}\hat{\psi}_{\mathcal{B}_{t}}(\mathbf{x}^{*})\right\|^{k}\right]. \tag{121}$$

<sup>&</sup>lt;sup>3</sup>Without the zero eigenvalue.

Note that under the conditions of Thm. 3 we have  $\|\mathcal{D}\hat{\psi}_{\mathcal{B}_t}(\boldsymbol{x}^*)\| \leq 1 - \varepsilon$  w.p. 1 for some  $\varepsilon \in (0, 1)$  (see (94), and the discussion above it). Therefore,

$$r(\mathbf{\Psi}_{k,k}) \leq \mathbb{E}\left[\left\|\mathcal{D}\hat{\psi}_{\mathcal{B}_t}(\mathbf{x}^*)\right\|^k\right] \leq (1-\varepsilon)^k.$$
 (122)

Overall,

$$r\left(\mathbf{\Psi}_{\rho}\right) = \sup_{k \in \mathbb{N}} r\left(\mathbf{\Psi}_{k,k}\right) \le \sup_{k \in \mathbb{N}} (1 - \varepsilon)^{k} = 1 - \varepsilon < 1. \tag{123}$$

## L Proof of Theorem 5

Let T be an operator defined on  $\ell_2$  space. Assume that T consists of blocks  $\{T_{i,j}\}$ , such that  $\forall i, j \ T_{i,j} \in \mathbb{R}^{d_i \times d_j}$  where  $\{d_i\}_{i=1}^{\infty}$  is a given sequence. Additionally, assume that

$$\forall j \in \mathbb{N}$$
  $\sum_{i=1}^{\infty} ||T_{i,j}|| \le \alpha$  and  $\forall i \in \mathbb{N}$   $\sum_{j=1}^{\infty} ||T_{i,j}|| \le \beta$ . (124)

Furthermore, for any  $u \in \ell_2$ , denote by  $u_i$  its *i*th segment, such that  $u_i \in \mathbb{R}^{d_i}$ . Then we have

$$\|T\boldsymbol{u}\|^{2} = \sum_{i=1}^{\infty} \left\| \sum_{j=1}^{\infty} T_{i,j} \boldsymbol{u}_{j} \right\|^{2}$$

$$\leq \sum_{i=1}^{\infty} \left( \sum_{j=1}^{\infty} \|T_{i,j}\| \|\boldsymbol{u}_{j}\| \right)^{2}$$

$$= \sum_{i=1}^{\infty} \left( \sum_{j=1}^{\infty} \sqrt{\|T_{i,j}\|} \sqrt{\|T_{i,j}\|} \|\boldsymbol{u}_{j}\| \right)^{2}$$

$$\leq \sum_{i=1}^{\infty} \left( \sum_{j=1}^{\infty} \|T_{i,j}\| \right) \left( \sum_{j=1}^{\infty} \|T_{i,j}\| \|\boldsymbol{u}_{j}\|^{2} \right)$$

$$\leq \sum_{i=1}^{\infty} \beta \left( \sum_{j=1}^{\infty} \|T_{i,j}\| \|\boldsymbol{u}_{j}\|^{2} \right)$$

$$= \beta \sum_{j=1}^{\infty} \|\boldsymbol{u}_{j}\|^{2} \sum_{i=1}^{\infty} \|T_{i,j}\|$$

$$\leq \beta \sum_{j=1}^{\infty} \|\boldsymbol{u}_{j}\|^{2} \alpha$$

$$= \alpha \beta \|\boldsymbol{u}\|^{2}. \tag{125}$$

where in the second step we used the triangle inequality, the fourth step is due to Cauchy-Schwarz inequality, and in the sixth step we used the fact that all summands are non-negative, and therefore we can change summation order.