
Understanding the Generalization Benefit of Normalization Layers: Sharpness Reduction

Kaifeng Lyu Zhiyuan Li Sanjeev Arora
Department of Computer Science
Princeton University
{klyu,zhiyuanli,arora}@cs.princeton.edu

Abstract

Normalization layers (e.g., Batch Normalization, Layer Normalization) were introduced to help with optimization difficulties in very deep nets, but they clearly also help generalization, even in not-so-deep nets. Motivated by the long-held belief that flatter minima lead to better generalization, this paper gives mathematical analysis and supporting experiments suggesting that normalization (together with accompanying weight-decay) encourages GD to reduce the sharpness of loss surface. Here “sharpness” is carefully defined given that the loss is scale-invariant, a known consequence of normalization. Specifically, for a fairly broad class of neural nets with normalization, our theory explains how GD with a finite learning rate enters the so-called Edge of Stability (EoS) regime, and characterizes the trajectory of GD in this regime via a continuous sharpness-reduction flow.

1 Introduction

Training modern deep neural nets crucially relies on normalization layers to make the training process less sensitive to hyperparameters and initialization. The two of the most popular normalization layers are Batch Normalization (BN) [55] for vision tasks and Layer Normalization (LN) [9] for language tasks. Recent works also proposed other normalization layers aiming for better performance, most notably including Group Normalization (GN) [120], Weight Normalization (WN) [102], Scaled Weight Standardization (SWS) [97, 53, 14], etc. Most normalization layers amount to a reparametrization of the neural net so that the loss becomes invariant to the scale of most parameters, and with a minor change, to *all* parameters: $\mathcal{L}(c\mathbf{w}) = \mathcal{L}(\mathbf{w})$ for all scalings $c > 0$ [55, 7, 77]. The current paper assumes this scale-invariance for all parameters and analyzes the trajectory of gradient descent with *weight decay* (WD):

$$\mathbf{w}_{t+1} \leftarrow (1 - \hat{\eta}\lambda)\mathbf{w}_t - \hat{\eta}\nabla\mathcal{L}(\mathbf{w}_t). \quad (1)$$

The use of WD is a common practice that has been adopted in training state-of-the-art neural nets, such as ResNets [46, 47] and Transformers [29, 15]. Previous ablation studies showed that adding WD to normalized nets indeed leads to better generalization [126, 72, 125]. More notably, Liu et al. [83] conducted experiments of training ResNets initialized from global minima with poor test accuracy, and showed that SGD with WD escapes from those bad global minima and attains good test accuracy. In contrast, training with vanilla SGD yields significant generalization degradation.

In the traditional view, WD regularizes the model by penalizing the parameter norm, but this may appear nonsensical for scale-invariant loss because one can scale down the norm arbitrarily without changing the loss value. However, the scale of the parameter *does* matter in backward propagation, and thus WD can affect the training dynamics. In particular, simple calculus shows $\nabla\mathcal{L}(\mathbf{w}) = \frac{1}{\|\mathbf{w}\|_2}\nabla\mathcal{L}(\frac{\mathbf{w}}{\|\mathbf{w}\|_2}) \propto \frac{1}{\|\mathbf{w}\|_2}$ and $\nabla^2\mathcal{L}(\mathbf{w}) = \frac{1}{\|\mathbf{w}\|_2^2}\nabla^2\mathcal{L}(\frac{\mathbf{w}}{\|\mathbf{w}\|_2}) \propto \frac{1}{\|\mathbf{w}\|_2^2}$, so WD is in effect trying to enlarge the gradient and Hessian in training. This makes the training dynamics very different from unnormalized nets and requires revisiting classical convergence analyses [77, 78, 84, 80].

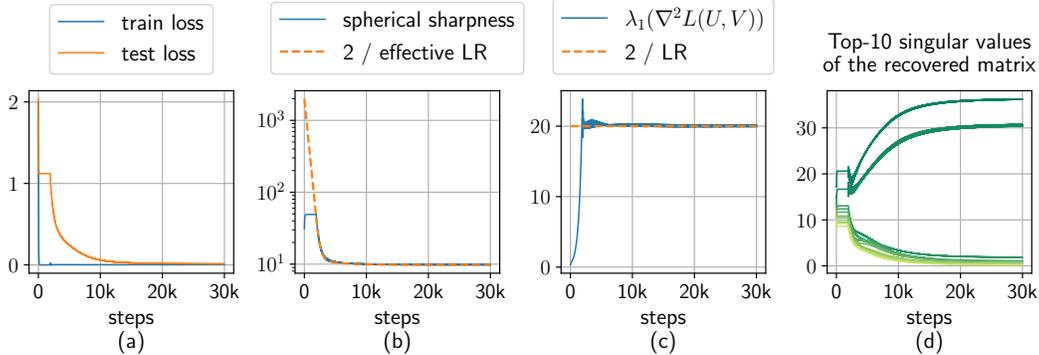


Figure 1: Experiment on overparameterized matrix completion with Batch Normalization. Given 800 (32%) entries Ω of a rank-2 matrix $M \in \mathbb{R}^{50 \times 50}$, use GD+WD to optimize the loss $\mathcal{L}(U, V) := \frac{1}{|\Omega|} \sum_{(i,j) \in \Omega} (\text{BN}([UV^\top]_{i,j}) - M_{i,j})^2$, where $U, V \in \mathbb{R}^{50 \times 50}$ (thus no explicit constraint on rank). Starting from step $\sim 2k$, spherical sharpness drops significantly (b), which encourages low-rank (d) and causes the test loss (MSE of all entries) to decrease from 1.12 to 0.013 (a). See also Appendix P.1.

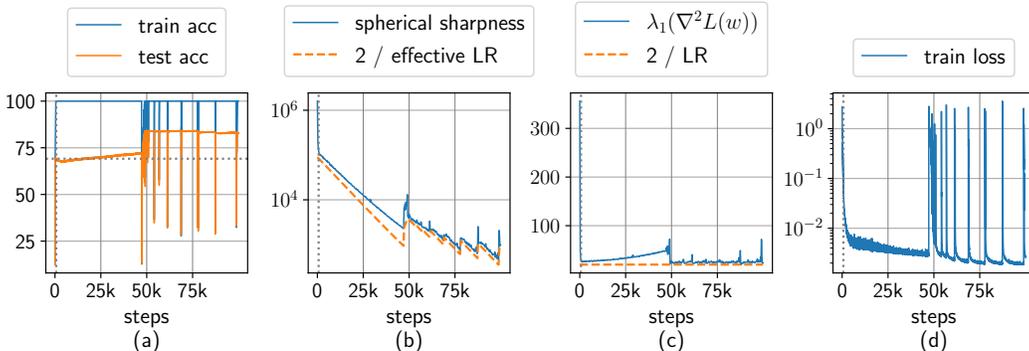


Figure 2: In training a smooth and scale-invariant VGG-11 on CIFAR-10 with (full-batch) GD+WD, the spherical sharpness keeps decreasing and the test accuracy keeps increasing. BN is added after every linear layer to ensure scale-invariance. 100% training accuracy is achieved after ~ 680 steps (dotted line), but as the training continues for $\sim 47k$ steps, the spherical sharpness keeps decreasing (b) and the test accuracy increases from 69.1% to 72.0% (a). Then the training exhibits destabilization but the test accuracy is further boosted to 84.3%. Removing either of BN or WD eliminates this phenomenon; see Appendices P.4 and P.5.

The current paper aims to improve mathematical understanding of how normalization improves generalization. While this may arise from many places, we focus on studying the dynamics of (full-batch) GD (1), which is a necessary first step towards understanding SGD. We show that the interplay between normalization and WD provably induces an implicit bias to persistently reduce the *sharpness* of the local loss landscape during the training process, which we call the *sharpness-reduction bias*.

It is long believed that flatter minima generalize better [50, 63, 95], but the notion of sharpness/flatness makes sense only if it is carefully defined in consideration of various symmetries in neural nets. One of the most straightforward measures of sharpness is the maximum eigenvalue of Hessian, namely $\lambda_1(\nabla^2 \mathcal{L}(w_t))$. But for normalized nets, this sharpness measure is vulnerable to weight rescaling, because one can scale the weight norm to make a minimizer arbitrarily flat [31]. Also, this sharpness measure may not decrease with the number of training steps: an empirical study by Cohen et al. [24] shows that for various neural nets (including normalized nets), GD has an overwhelming tendency to persistently increase $\lambda_1(\nabla^2 \mathcal{L}(w_t))$ until it reaches the *Edge of Stability (EoS) regime*, a regime where $\lambda_1(\nabla^2 \mathcal{L}(w_t))$ stays around $2/\hat{\eta}$ ($\hat{\eta}$ is the learning rate). See also Section 6 and Figure 2c.

1.1 Our Contributions

The sharpness measure we use in this paper takes care of the scale-invariance in normalized nets. We are motivated by our experiments on matrix completion (with BN) and CIFAR-10, where our sharpness measure decreases as the training proceeds, and the generalization improves accordingly; see Figures 1 and 2. We note that techniques from previous works [92, 95, 37] can be easily adopted here to establish a PAC-Bayes bound on the test error, where our sharpness measure appears as an additive term (see Appendix C).

Definition 1.1 (Spherical Sharpness). For a scale-invariant loss $\mathcal{L}(\mathbf{w})$ (i.e., $\mathcal{L}(c\mathbf{w}) = \mathcal{L}(\mathbf{w})$ for all $c > 0$), the spherical sharpness at $\mathbf{w} \in \mathbb{R}^D$ is defined by $\lambda_1(\nabla^2 \mathcal{L}(\frac{\mathbf{w}}{\|\mathbf{w}\|_2}))$, the maximum eigenvalue of the Hessian matrix after projecting \mathbf{w} onto the unit sphere.

Based on Definition 1.1, we study the aforementioned sharpness-reduction bias in training normalized nets with GD+WD (defined in (1)). For constant learning rate $\hat{\eta}$ and weight decay $\hat{\lambda}$, we can rewrite this rule equivalently as Projected Gradient Descent (PGD) on the unit sphere with *adaptive* learning rates, $\boldsymbol{\theta}_{t+1} \leftarrow \Pi(\boldsymbol{\theta}_t - \tilde{\eta}_t \nabla \mathcal{L}(\boldsymbol{\theta}_t))$, where $\boldsymbol{\theta}_t := \frac{\mathbf{w}_t}{\|\mathbf{w}_t\|_2}$ is the direction of \mathbf{w}_t , and $\tilde{\eta}_t$ is the “effective” learning rate at step t (see Lemma 3.1). We call $\tilde{\eta}_t$ adaptive because it can be shown to resemble the behaviors of adaptive gradient methods (e.g., RMSprop [49]): $\tilde{\eta}_t$ increases when gradient is small and decreases when gradient is large (Figure 3). Our main contributions are as follows:

1. After $\boldsymbol{\theta}_t$ reaches a point near the manifold of minimizers of \mathcal{L} , we theoretically show that the effective learning rate $\tilde{\eta}_t$ increases until GD enters a regime where $2/\tilde{\eta}_t$ roughly equals to the spherical sharpness (or equivalently $2/\hat{\eta} \approx \lambda_1(\nabla^2 \mathcal{L}(\mathbf{w}_t))$), namely the EoS regime (Section 4.1).
2. In the EoS regime, we show that for GD with a small (but finite) learning rate, $\boldsymbol{\theta}_t$ oscillates around the manifold and moves approximately along a sharpness-reduction flow, which is a gradient flow for minimizing spherical sharpness on the manifold (with gradient-dependent learning rate) (Section 4.2).
3. As an application of our theory, we show that for linear regression with BN, GD+WD finds the minimizer that corresponds to the linear model with minimum weight norm, which looks surprisingly the same as the conventional effect of WD but is achieved through the completely different sharpness-reduction mechanism (Section 5).
4. We experimentally verified the sharpness-reduction phenomenon predicted by our theorem and its benefits to generalization on CIFAR-10 with VGG-11 and ResNet-20, as well as matrix completion with BN (Appendix P).
5. We generalize our theoretical results of sharpness-reduction bias to a broader class of adaptive gradient methods, most notably a variant of RMSprop with scalar learning rate (Appendix B).

Technical Contribution. Our proof technique is novel and may have independent interest to the ML community. The main challenge is that we need to analyze the implicit bias of GD in the EoS regime which crucially relies on step size being finite — this is in sharp contrast to many previous works on implicit bias of GD [107, 106, 87, 59, 43, 42, 76, 100, 4, 22, 79, 88, 101, 108, 38] where the same bias exists at infinitesimal LR. Our analysis is inspired by a previous line of works [13, 25, 81] showing that label noise can drive SGD to move on the minimizer manifold along the direction of minimizing the trace of Hessian. We borrow a few lemmas from those analyses, but the overall proof strategy is very different because our setting does not even have any stochastic gradient noise. Instead, we connect the dynamics in the EoS regime to power methods and show that GD oscillates around the minimizer manifold. This oscillation then becomes a driving power that pushes the parameter to move on the manifold. Finally, we analyze the speed of this movement by modeling two key parameters of the dynamics as a 1-dimensional Hamiltonian system (Figure 6). To the best of our knowledge, we are the first to provide theoretical proof for a sharpness measure to decrease during the standard GD training, without any additional regularization (e.g., label noise [13, 25, 81]) and without involving uncommon variants of GD (e.g., normalized GD or non-smooth wrappings on the loss function [8]).

2 Related Works

Sharpness and Generalization. It has been long believed that flat minima generalize better [50]. Several empirical studies [63, 74, 117, 57] verified the positive correlation between flatness and generalization. Neyshabur et al. [95] justified this via PAC-Bayes theory [92]. Several other theoretical papers explored the generalization properties of flat minima specifically for two-layer nets [13, 94, 44, 81, 30] and deep linear nets [93]. Jiang et al. [60] conducted extensive experiments for all existing generalization measures to evaluate their correlation and causal relationships with generalization error, concluding that sharpness-based measures perform the best overall. In light of this, Foret et al. [37] proposed SAM algorithm to improve the generalization by minimizing the sharpness. Despite so many positive results on sharpness-based measures, a common issue of many works is that the measures may suffer from sensitivity to rescaling of parameters in deep nets [31]. Another issue is that the minima could lie in asymmetric valleys that are flat on one side and sharp on the other [45].

Understanding Normalization Layers. The benefits of normalization layers can be shown in various aspects. A series of works studied the forward propagation of deep nets at random initialization, showing that normalization layers stabilize the growth of intermediate layer outputs with depth [14, 10, 28], provably avoid rank collapse [26] and orthogonalize representations [27]. Although these works mainly focused on BN [55], Lubana et al. [85], Labatie et al. [67] provided thorough discussions on the applicability of these arguments to other normalization layers. It is also believed that BN has a unique regularization effect through the noise in batch statistics [86, 111, 104]. Several other works argued that normalization layers lead to a smoothening or preconditioning effect of the loss landscape [103, 12, 39, 61, 82, 68], which may help optimization. By analyzing the training dynamics, Arora et al. [7] rigorously proved that normalization yields an auto-tuning effect of the effective learning rate $\tilde{\eta}_t$, which makes the asymptotic speed of optimization much less sensitive to the learning rate and initialization. In linear regression settings, Cai et al. [16], Kohler et al. [65] showed that training with BN leads to a faster convergence rate; Wu et al. [119] studied the implicit regularization effect of WN [102]. For two-layer nets with normalization, Ma and Ying [90] derived a mean-field formulation of the training dynamics; Dukler et al. [33] proved a convergence rate via NTK-based analysis. The current paper focuses on the interplay between normalization and WD during training, whereas all the above works either do not analyze the dynamics or assume no WD.

Interplay Between Normalization and WD. A common feature of normalization layers (including but not limited to BN, WN, LN, GN, SWS) is that they make the loss invariant to the scale of layer weights. In presence of both scale-invariance and WD, training dynamics can go out of the scope of the classical optimization theory, e.g., one can train the net to small loss even with learning rates exponentially increasing [77]. A series of works investigated into the interplay between normalization and WD and argued that the training dynamic with SGD eventually reaches an “equilibrium” state, where the parameter norm [78, 113, 21] and the size of angular update [114] become stable. Li et al. [78], Wang and Wang [115] provided empirical and theoretical evidence that the function represented by the net also equilibrates to a stationary distribution that is independent of initialization. This could be related to Liu et al. [83]’s experiments on the ability of SGD with WD to escape from bad initialization, but it remains unclear why the generalization should be good at the equilibrium state. In this paper, we focus on (full-batch) GD, which is the most basic and important special case of SGD.

3 Preliminaries

Let $\mathbb{S}^{D-1} := \{\boldsymbol{\theta} \in \mathbb{R}^D : \|\boldsymbol{\theta}\|_2 = 1\}$ be the unit sphere equipped with subspace topology. We say a loss function $\mathcal{L}(\boldsymbol{w})$ defined on $\mathbb{R}^D \setminus \{\mathbf{0}\}$ is *scale-invariant* if $\mathcal{L}(c\boldsymbol{w}) = \mathcal{L}(\boldsymbol{w})$ for all $c > 0$. In other words, the loss value does not change with the parameter norm. For a differentiable scale-invariant function $\mathcal{L}(\boldsymbol{w})$, the gradient is (-1) -homogeneous and it is always perpendicular to \boldsymbol{w} , i.e., $\nabla\mathcal{L}(c\boldsymbol{w}) = c^{-1}\nabla\mathcal{L}(\boldsymbol{w})$ for all $c > 0$ and $\langle \nabla\mathcal{L}(\boldsymbol{w}), \boldsymbol{w} \rangle = 0$ (see Lemma D.1).

The focus of this paper is the dynamics of GD+WD on scale-invariant loss. (1) gives the update rule for learning rate (LR) $\hat{\eta}$ and weight decay (WD) $\hat{\lambda}$. We use $\boldsymbol{\theta}_t := \frac{\boldsymbol{w}_t}{\|\boldsymbol{w}_t\|_2}$ to denote the projection of \boldsymbol{w}_t onto \mathbb{S}^{D-1} at step t . We write GD+WD on scale-invariant loss as a specific kind of Projected Gradient Descent (PGD) and define the *effective learning rate* to be the LR $\tilde{\eta}_t := \frac{\hat{\eta}}{(1-\hat{\eta}\hat{\lambda})\|\boldsymbol{w}_t\|_2^2}$ that appears in the update rule of PGD. This notion is slightly different from the effective learning rate $\frac{\hat{\eta}}{\|\boldsymbol{w}_t\|_2^2}$ defined in previous works [113, 52, 7], but ours is more convenient for our analysis.

Lemma 3.1. *When the parameters \boldsymbol{w}_t are updated as (1), $\boldsymbol{\theta}_t$ satisfies the following equation:*

$$\boldsymbol{\theta}_{t+1} = \Pi(\boldsymbol{\theta}_t - \tilde{\eta}_t \nabla\mathcal{L}(\boldsymbol{\theta}_t)), \quad (2)$$

where $\tilde{\eta}_t := \frac{\hat{\eta}}{(1-\hat{\eta}\hat{\lambda})\|\boldsymbol{w}_t\|_2^2}$ is called the *effective learning rate* at step t , and $\Pi : \boldsymbol{w} \mapsto \frac{\boldsymbol{w}}{\|\boldsymbol{w}\|_2}$ is the projection operator that projects any vector onto the unit sphere.

4 GD+WD on Scale-Invariant Loss Functions

This section analyzes GD+WD (1) on a scale-invariant loss $\mathcal{L}(\boldsymbol{w})$, in particular what happens after approaching a manifold of local minimizers. Section 4.1 analyzes the dynamics in the stable regime, where loss is guaranteed to decrease monotonically, and Theorem 4.2 suggests \boldsymbol{w}_t can get close to a local minimizer at some time t_0 . We show that the effective LR keeps increasing after t_0 , causing

GD+WD to eventually leave this stable regime and enter a new regime which we call the Edge of Stability (EoS). In Section 4.2, we establish our main theorem, which connects the dynamics of w_t in the EoS regime to a sharpness-reduction flow.

4.1 GD+WD Eventually Leaves the Stable Regime

A standard step of analyzing optimization methods is to do Taylor expansion locally for the loss function, and show that how the optimization method decreases the loss using a *descent lemma*. In our case of scale-invariant loss functions, we use $\mathbf{H}(\mathbf{w}) := \nabla^2 \mathcal{L}(\mathbf{w}) \in \mathbb{R}^{D \times D}$ to denote the Hessian matrix of \mathcal{L} at $\mathbf{w} \in \mathbb{R}^D$, and $\lambda_1^{\mathbf{H}}(\mathbf{w}) := \lambda_1(\mathbf{H}(\mathbf{w}))$ to denote the top eigenvalue of $\mathbf{H}(\mathbf{w})$.

Lemma 4.1 (Descent Lemma). *For scale-invariant loss $\mathcal{L}(\mathbf{w})$, at step t of GD+WD we have*

$$\mathcal{L}(\boldsymbol{\theta}_{t+1}) \leq \mathcal{L}(\boldsymbol{\theta}_t) - \tilde{\eta}_t(1 - \tilde{\eta}_t \lambda_{\max}^{(t)}/2) \|\nabla \mathcal{L}(\boldsymbol{\theta}_t)\|_2^2.$$

where $\lambda_{\max}^{(t)} := \sup_{\alpha \in [0, \tilde{\eta}_t]} \{\lambda_1^{\mathbf{H}}(\boldsymbol{\theta}_t - \alpha \nabla \mathcal{L}(\boldsymbol{\theta}_t))\}$ is an upper bound of spherical sharpness locally.

This descent lemma shows that the training loss $\mathcal{L}(\boldsymbol{\theta}_t)$ keeps decreasing as long as the effective LR $\tilde{\eta}_t$ is smaller than $2/\lambda_{\max}^{(t)}$. We call the regime of $\tilde{\eta}_t < 2/\lambda_{\max}^{(t)}$ as the *stable regime* of GD+WD. If $\tilde{\eta}_t \approx 2/\lambda_{\max}^{(t)}$ with a small difference, then we call it as the *Edge of Stability (EoS) regime*. We remark that this condition for EoS regime is essentially the same as $\hat{\eta} \approx 2/\lambda_1^{\mathbf{H}}(\mathbf{w})$ in Cohen et al. [24]’s definition because $\tilde{\eta}_t \cdot \lambda_{\max}^{(t)} \approx \hat{\eta} \cdot \lambda_1^{\mathbf{H}}(\mathbf{w})$; see Appendix G.3.

Fix an initial point $\mathbf{w}_0 \in \mathbb{R}^D \setminus \{\mathbf{0}\}$. Now we aim to characterize the dynamics of GD+WD when LR $\hat{\eta}$ and WD $\hat{\lambda}$ are small enough. The convergence rate of GD+WD has been analyzed by Li et al. [80]. Here we present a variant of their theorem that bounds both the gradient and effective LR.

Theorem 4.2 (Variant of Theorem D.2, Li et al. [80]). *Let $\mathcal{L}(\mathbf{w})$ be a scale-invariant loss function and $\rho_2 := \sup\{\|\nabla^2 \mathcal{L}(\mathbf{w})\|_2 : \mathbf{w} \in \mathbb{S}^{D-1}\}$ be the smoothness constant of \mathcal{L} restricted on the unit sphere. For GD+WD (1) with $\hat{\eta}\hat{\lambda} \leq 1/2$ and $\tilde{\eta}_0 \leq \frac{1}{\pi^2 \rho_2 (1 - \hat{\eta}\hat{\lambda})}$, let $T_0 := \left\lceil \frac{1}{2\hat{\eta}\hat{\lambda}} \ln \frac{\|\mathbf{w}_0\|_2^2}{\rho_2 \pi^2 \hat{\eta}} \right\rceil$ steps, there must exist $0 \leq t \leq T_0$ such that $\|\nabla \mathcal{L}(\boldsymbol{\theta}_t)\|_2^2 \leq 8\pi^4 \rho_2^2 \hat{\lambda} \hat{\eta}$ and $\tilde{\eta}_t \leq \frac{2}{\pi^2 \rho_2 (1 - \hat{\eta}\hat{\lambda})}$.*

Theorem 4.2 shows that for some $t_0 \leq T_0$, $\|\nabla \mathcal{L}(\boldsymbol{\theta}_{t_0})\|_2^2 \leq O(\hat{\lambda}\hat{\eta})$ and $\tilde{\eta}_{t_0} \leq \frac{1}{\pi^2 \rho_2} < \frac{2}{\rho_2}$, which means $\boldsymbol{\theta}_{t_0}$ is an approximate first-order stationary point of \mathcal{L} on the unit sphere. This does not guarantee that $\boldsymbol{\theta}_{t_0}$ is close to any global minimizer, but in practice the training loss rarely gets stuck at a non-optimal value when the model is overparameterized [70, 96, 71, 125]. We are thus motivated to study the case where $\boldsymbol{\theta}_{t_0}$ not only has small gradient $\|\nabla \mathcal{L}(\boldsymbol{\theta}_{t_0})\|_2^2 \leq O(\hat{\lambda}\hat{\eta})$ but also is close to a local minimizer $\boldsymbol{\theta}^* \in \mathbb{S}^{D-1}$ of \mathcal{L} in the sense that $\|\boldsymbol{\theta}_{t_0} - \boldsymbol{\theta}^*\|_2 \leq O((\hat{\lambda}\hat{\eta})^{1/2})$ (assuming smoothness, the latter implies the former).

As the gradient is small near the local minimizer $\boldsymbol{\theta}^*$, starting from step t_0 , the norm of \mathbf{w}_t decreases due to the effect of WD. See Figure 3a. Since the effective LR is inversely proportional to $\|\mathbf{w}_t\|_2^2$, this leads to the effective LR to increase. Then Theorem 4.4 will show that the GD+WD dynamic eventually leaves the stable regime at some time $t_1 > t_0$, and enters the EoS regime where $\tilde{\eta}_t \approx 2/\lambda_{\max}^{(t)}$.

To establish Theorem 4.4, we need to assume that \mathcal{L} satisfies Polyak-Łojasiewicz (PL) condition locally, which is a standard regularity condition in the optimization literature to ease theoretical analysis around a minimizer. Intuitively, PL condition guarantees that the gradient grows faster than a quadratic function as we move a parameter $\boldsymbol{\theta}$ away from $\boldsymbol{\theta}^*$. Note that PL condition is strictly weaker than convexity as the function can still be non-convex under PL condition (see, e.g., [62]).

Definition 4.3 (Polyak-Łojasiewicz Condition). *For a scale-invariant loss $\mathcal{L}(\mathbf{w})$ and $\mu > 0$, we say that \mathcal{L} satisfies μ -Polyak-Łojasiewicz condition (or μ -PL) locally around a local minimizer $\boldsymbol{\theta}^*$ on \mathbb{S}^{D-1} if for some neighborhood $U \subseteq \mathbb{S}^{D-1}$ of $\boldsymbol{\theta}^*$, $\forall \boldsymbol{\theta} \in U : \frac{1}{2} \|\nabla \mathcal{L}(\boldsymbol{\theta})\|_2^2 \geq \mu \cdot (\mathcal{L}(\boldsymbol{\theta}) - \mathcal{L}(\boldsymbol{\theta}^*))$.*

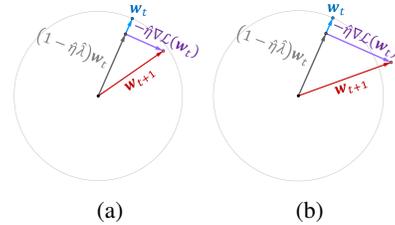


Figure 3: The norm of \mathbf{w}_t decreases when gradient is small and increases when gradient is large.

Theorem 4.4. Let $\mathcal{L}(\mathbf{w})$ be a \mathcal{C}^2 -smooth scale-invariant loss that satisfies μ -PL around a local minimizer θ^* on the unit sphere, and $\rho_2 := \sup\{\|\nabla^2\mathcal{L}(\mathbf{w})\|_2 : \mathbf{w} \in \mathbb{S}^{D-1}\}$. For GD+WD on $\mathcal{L}(\mathbf{w})$ with learning rate $\hat{\eta}$ and weight decay $\hat{\lambda}$, if at some step t_0 , $\|\theta_{t_0} - \theta^*\|_2 \leq O((\hat{\lambda}\hat{\eta})^{1/2})$ and $\tilde{\eta}_{t_0} \leq \frac{2}{\rho_2} < \frac{2}{\lambda_1^H(\theta^*)}$, and if $\hat{\lambda}\hat{\eta}$ is small enough, then there exists a time $t_1 > t_0$ such that $\|\theta_{t_1} - \theta^*\|_2 = O((\hat{\lambda}\hat{\eta})^{1/2})$ and $\tilde{\eta}_{t_1} = \frac{2}{\lambda_1^H(\theta^*)} + O((\hat{\lambda}\hat{\eta})^{1/2})$.

4.2 Dynamics at the Edge of Stability

From the analysis in the previous subsection, we know that θ_t can get close to a local minimizer θ^* and enter the EoS regime at some step t_1 . But what happens after t_1 ?

Figure 4 gives a warm-up example on a 3D scale-invariant loss $\mathcal{L} : \mathbb{R}^3 \setminus \{0\} \rightarrow \mathbb{R}$, where the black line is a manifold Γ consisting of all the minimizers. In training with GD+WD, θ_t first goes close to a local minimizer ζ_0 , then Theorem 4.4 suggests that WD causes the effective LR to steadily increase until the dynamic enters the EoS regime. Now something interesting happens — θ_t moves a bit away from ζ_0 and starts to oscillate around the manifold Γ . This oscillation is not completely perpendicular to Γ but actually forms a small angle that pushes θ_t to move downward persistently until θ_t approaches the minimizer ζ_* denoted in the plot.

For a general scale-invariant loss $\mathcal{L} : \mathbb{R}^D \setminus \{0\} \rightarrow \mathbb{R}$, which minimizer does θ_t move towards? In this work, we consider the setting where there is a manifold Γ consisting only of local minimizers (but not necessarily all of them). We show that θ_t always oscillates around the manifold once it approaches the manifold and enters the EoS regime, and meanwhile θ_t keeps moving in a direction of reducing spherical sharpness.

4.2.1 Assumptions

Now we formally introduce our main assumption on the local minimizer manifold Γ .

Assumption 4.5. The loss function $\mathcal{L} : \mathbb{R}^D \setminus \{0\} \rightarrow \mathbb{R}$ is \mathcal{C}^4 -smooth and scale-invariant. Γ is a \mathcal{C}^2 -smooth, $(D_\Gamma - 1)$ -dimensional submanifold of \mathbb{S}^{D-1} for some $0 \leq D_\Gamma < D$, where every $\theta \in \Gamma$ is a local minimizer of \mathcal{L} on \mathbb{S}^{D-1} and $\text{rank}(\mathbf{H}(\theta)) = D - D_\Gamma$.

Scale-invariance has become a standard assumption in studying neural nets with normalization layers [77, 78, 84]. For VGG and ResNet, the scale-invariance can be ensured after making minor changes to the architectures (see Appendix Q.1). The training loss \mathcal{L} may not be smooth if the activation is ReLU, but lately it has become clear that differentiable activations such as Swish [98], GeLU [48] can perform equally well. Swish is indeed used in our VGG-11 experiments (Figure 2), but ResNet with ReLU activation also exhibits a sharpness-reduction bias empirically (see Appendix P.2).

For any local minimizer $\theta \in \Gamma$, the eigenvalues $\lambda_k^H(\theta)$ must be non-negative. And $\lambda_k^H(\theta) = 0$ for all $D - D_\Gamma < k \leq D$, since Γ is of dimension $D_\Gamma - 1$. The condition $\text{rank}(\mathbf{H}(\theta)) = D - D_\Gamma$ ensures that the Hessian is maximally non-degenerate on Γ , which also appears as a key assumption in previous works [81, 8, 35]. This condition simplifies the calculus on Γ in our analysis as it ensures that the null space of the matrix $\mathbf{H}(\theta)$ equals to the tangent space of Γ at $\theta \in \Gamma$. It is also closely related to PL condition (Definition 4.3) as Assumption 4.5 implies that $\mathcal{L}(\theta)$ satisfies μ -PL (for some $\mu > 0$) locally around every $\theta \in \Gamma$ on the unit sphere (Arora et al. [8], Lemma B.3).

To ease our analysis, we also need the following regularity condition to ensure that the largest eigenvalue is unique. In our experiments, sharpness reduction happens even when the multiplicity of the top eigenvalue is more than 1, but we leave the analysis of that case to future work.

Assumption 4.6. For all $\theta \in \Gamma$, $\lambda_1^H(\theta) > \lambda_2^H(\theta)$. That is, the top eigenvalue of $\mathbf{H}(\theta)$ is unique.

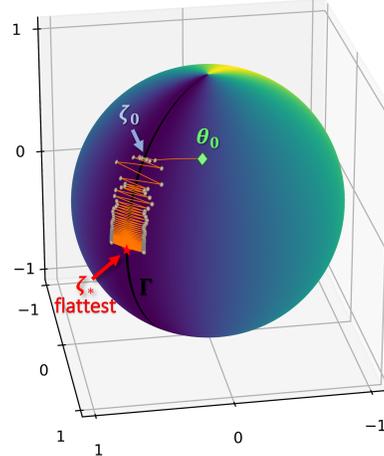


Figure 4: The trajectory of θ_t on a 3D scale-invariant loss function. Darker color means lower loss on the unit sphere, and points in the black line are minimizers (see Appendix F). In the end, θ_t approaches the flattest one (red star).

4.2.2 Main Theorem

First, we define $\eta_{\text{in}} := \hat{\eta}\hat{\lambda}$ as the intrinsic learning rate (name from Li et al. [78]) for convenience. As suggested in Theorems 4.2 and 4.4, θ_t can get close to a local minimizer and be in the EoS regime at some step t_1 : if ζ_0 is the local minimizer, then $\|\theta_{t_1} - \zeta_0\|_2 = O(\eta_{\text{in}}^{1/2})$ and $\tilde{\eta}_{t_1} = \frac{2}{\lambda_1^{\text{H}}(\zeta_0)} + O(\eta_{\text{in}}^{1/2})$. In our main theorem, we start our analysis from step t_1 while setting $t_1 = 0$ WLOG (otherwise we can shift the step numbers). We connect GD+WD in the EoS regime to the following gradient flow (3) on the manifold Γ minimizing spherical sharpness (with gradient-dependent learning rate), and show that one step of GD+WD tracks a time interval of length η_{in} in the gradient flow.

$$\zeta(0) = \zeta_0 \in \Gamma, \quad \frac{d}{d\tau} \zeta(\tau) = -\frac{2\nabla_{\Gamma} \log \lambda_1^{\text{H}}(\zeta(\tau))}{4 + \|\nabla_{\Gamma} \log \lambda_1^{\text{H}}(\zeta(\tau))\|_2^2}. \quad (3)$$

Here we use the notation $\nabla_{\Gamma} R(\theta)$ for any $R : \mathbb{R}^D \rightarrow \mathbb{R}$ to denote the projection of $\nabla R(\theta)$ onto the tangent space $\mathbb{T}_{\theta}(\Gamma)$ at $\theta \in \Gamma$. $\zeta(\tau)$ reduces sharpness as it moves in direction of the negative gradient of $\log \lambda_1^{\text{H}}(\zeta(\tau))$ on Γ . A simple chain rule shows how fast the spherical sharpness decreases:

$$\frac{d}{dt} \log \lambda_1^{\text{H}}(\zeta(\tau)) = -\frac{2\|\nabla_{\Gamma} \log \lambda_1^{\text{H}}(\zeta(\tau))\|_2^2}{4 + \|\nabla_{\Gamma} \log \lambda_1^{\text{H}}(\zeta(\tau))\|_2^2} \approx \begin{cases} -\frac{1}{2}\|\nabla_{\Gamma} \log \lambda_1^{\text{H}}(\zeta(\tau))\|_2^2 & \text{for small gradient;} \\ -2 & \text{for large gradient.} \end{cases}$$

Note that it is not enough to just assume that θ_0 is close to ζ_0 . If $\theta_0 = \zeta_0$ holds exactly, then the subsequent dynamic of w_t is described by $w_t = (1 - \hat{\eta}\hat{\lambda})^t w_0$ with direction unchanged. There are also some other bad initial directions of w_0 that may not lead to the sharpness-reduction bias. This motivates us to do a smoothed analysis for the initial direction: the initial direction is ζ with tiny random perturbation, where the perturbation scale is allowed to vary from $\exp(-\eta_{\text{in}}^{-o(1)})$ to $\eta_{\text{in}}^{1/2-o(1)}$, and we show that a good initial direction is met with high probability as $\eta_{\text{in}} \rightarrow 0$.¹ Alternatively, one can regard it as a modeling of the tiny random noise in GD+WD due to the precision errors in floating-point operations. See Figure 5b; the training loss can never be exactly zero in practice.

Initialization Scheme. Given a local minimizer $\zeta_0 \in \Gamma$, we initialize $w_0 \in \mathbb{R}^D \setminus \{0\}$ as follows: draw $\xi \sim \mathcal{N}(0, \sigma_0^2 I/D)$ from Gaussian and set the direction of w_0 to $\frac{\zeta_0 + \xi}{\|\zeta_0 + \xi\|_2}$, where σ_0 can take any value in $[\exp(-\eta_{\text{in}}^{-o(1)}), \eta_{\text{in}}^{1/2-o(1)}]$; then set the parameter norm $\|w_0\|_2$ to be any value that satisfies $|\tilde{\eta}_0 - \frac{2}{\lambda_1^{\text{H}}(\zeta_0)}| \leq \eta_{\text{in}}^{1/2-o(1)}$, where $\tilde{\eta}_0 := \frac{\hat{\eta}}{(1-\hat{\eta}\hat{\lambda})\|w_0\|_2}$ is the effective LR for the first step.

Theorem 4.7. *Under Assumptions 4.5 and 4.6, for GD+WD (1) with sufficiently small intrinsic learning rate $\eta_{\text{in}} := \hat{\eta}\hat{\lambda}$, if we follow the above initialization scheme for some $\zeta_0 \in \Gamma$, then with probability $1 - O(\eta_{\text{in}}^{1/2-o(1)})$, the trajectory of $\theta_t := \frac{w_t}{\|w_t\|_2}$ approximately tracks a sharpness-reduction flow $\zeta : [0, T] \rightarrow \Gamma$ that starts from ζ_0 and evolves as the ODE (3) up to time T (if solution exists), in the sense that $\|\theta_t - \zeta(t\eta_{\text{in}})\|_2 = O(\eta_{\text{in}}^{1/4-o(1)})$ for all $0 \leq t \leq T/\eta_{\text{in}}$.*

Remark 4.8 (Magnitude of Oscillation). As suggested by Figure 4, θ_t actually oscillates around the manifold. But according to our analysis, the magnitude of oscillation is as small as $O(\eta_{\text{in}}^{1/2-o(1)})$, so it is absorbed into our final bound $O(\eta_{\text{in}}^{1/4-o(1)})$ for the distance between θ_t and $\zeta(t\eta_{\text{in}})$.

4.2.3 Proof Idea

Throughout our proof, we view GD+WD for w_t as a PGD for θ_t with effective LR $\tilde{\eta}_t$ (Lemma 3.1). To track θ_t with $\zeta(t\eta_{\text{in}})$, for each step t , we construct a local minimizer $\phi_t \in \Gamma$ that serves as the ‘‘projection’’ of θ_t onto the manifold Γ , in the sense that the displacement $x_t := \theta_t - \phi_t$ is approximately perpendicular to the tangent space of Γ at ϕ_t . Our entire proof works through induction. According to the initial conditions, the dynamic is initially in the EoS regime: $\|x_t\|_2 \leq \eta_{\text{in}}^{1/2-o(1)}$ and $|\tilde{\eta}_t - 2/\lambda_1^{\text{H}}(\phi_t)| \leq \eta_{\text{in}}^{1/2-o(1)}$ at $t = 0$. In our induction, we maintain the induction hypothesis that these two EoS conditions continue to hold for all $t \geq 0$.

¹Here $\eta_{\text{in}}^{-o(1)}$ can be constant, $O(\log(1/\eta_{\text{in}}))$, or $O(\text{polylog}(1/\eta_{\text{in}}))$, but not $\eta_{\text{in}}^{-\epsilon}$ if $\epsilon > 0$ is a constant. As mentioned later, this need for random initialization is very similar to the one needed in power method for computing eigenvalues.

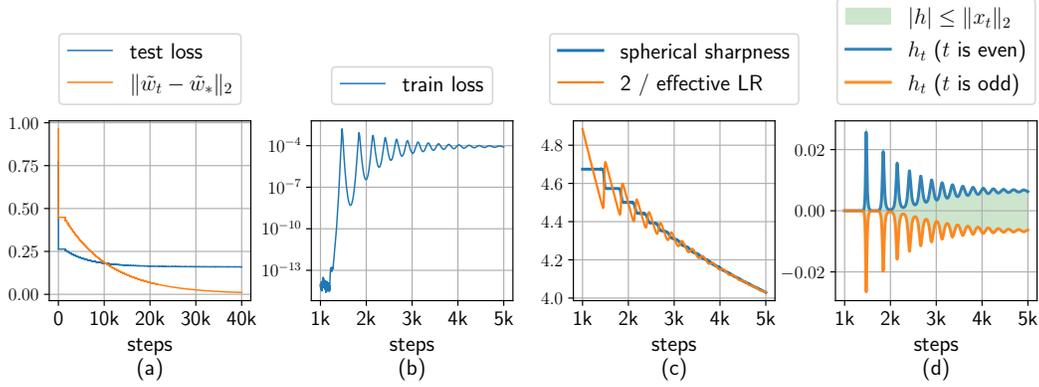


Figure 5: Illustration of the oscillation and periodic behaviors of GD+WD on linear regression with BN (see Sections 4.2.3 and 5). The training loss decreases to $\approx 10^{-14}$ in the first 1k steps and achieves test loss 0.26. Starting from step $\sim 1k$, the dynamic enters the EoS regime. (a). The test loss decreases to 0.16 as a distance measure to the flattest solution (M) decreases towards 0; (b). The training loss oscillates around $\sim 10^{-4}$ in the EoS regime; (c). $2/\tilde{\eta}_t$ switches back and forth between being smaller and larger than $\lambda_1^H(\phi_t)$; (d). The parameter oscillates around the minimizer manifold along the top eigenvector direction, and the magnitude of oscillation $|h_t|$ rises and falls periodically.

Period-Two Oscillation. A key insight in our proof is that after a few initial steps, θ_t is oscillating around ϕ_t along the $\pm v_1^H(\theta)$ directions, where $v_1^H(\theta)$ is a unit top eigenvector of $H(\theta)$ and is chosen in a way that $v_1^H(\theta)$ is continuous on Γ . More specifically, $x_t = h_t v_1^H(\phi_t) + O(\|x_t\|_2^2)$ for $h_t := \langle x_t, v_1^H(\phi_t) \rangle$. The oscillation is of period 2: $h_t > 0$ when t is even and $h_t < 0$ when t is odd. See Figure 5d for an example.

This oscillation can be connected to a power method for the matrix $I - \tilde{\eta}_t H(\phi_t)$. In the EoS regime, we can approximate θ_{t+1} (when x_t is small) as $\theta_{t+1} = \Pi(\theta_t - \tilde{\eta}_t \nabla \mathcal{L}(\theta_t)) \approx \Pi(\theta_t - \tilde{\eta}_t H(\phi_t) x_t) \approx \theta_t - \tilde{\eta}_t H(\phi_t) x_t$ by Taylor expansions of $\nabla \mathcal{L}$ and $\Pi : \mathbb{R}^D \setminus \{0\} \rightarrow \mathbb{S}^{D-1}$. We can further show that $\phi_{t+1} \approx \phi_t$ due to our choice of projections. Then the connection to power method is shown below:

$$x_{t+1} \approx \theta_{t+1} - \phi_t \approx (I - \tilde{\eta}_t H(\phi_t)) x_t.$$

By simple linear algebra, $v_1^H(\phi_t)$ is an eigenvector of $I - \tilde{\eta}_t H(\phi_t)$, associated with eigenvalue $1 - \tilde{\eta}_t \lambda_1^H(\phi_t) \approx -1$. The remaining eigenvalues are $\{1 - \tilde{\eta}_t \lambda_i^H(\phi_t)\}_{i=2}^D$, where $\lambda_i^H(\phi_t)$ is the i -th largest eigenvalue of $H(\theta_t)$, and they lie in the range $(-1, 1]$ since $\lambda_i^H(\phi_t) \in [0, \lambda_1^H(\phi_t)]$. Using a similar analysis to power method, we show that x_t quickly aligns to the direction of $\pm v_1^H(\phi_t)$ after a few initial steps, as the corresponding eigenvalue has approximately the largest absolute value.²

To formally establish the above result, we need a tiny initial alignment between x_0 and $v_1^H(\phi_0)$, just as the initial condition in power method. This is where we need the initial random perturbation.

Oscillation Drives ϕ_t to Move. This period-two oscillation is the driving power to push ϕ_t to move on the manifold. The main idea here is to realize that the oscillation direction deviates slightly from the direction of $\pm v_1^H(\phi_t)$ by using a higher-order approximation. We specifically use the Taylor approximation to show that this deviation leads ϕ_t to move slightly on Γ : after each cycle of oscillation, $\phi_{t+2} \approx \phi_t - 4h_t^2 \nabla_{\Gamma} \log \lambda_1^H(\phi_t) + O(\eta_{\text{in}}^{1.5-o(1)})$, which resembles two steps of gradient descent on Γ to minimize the logarithm of spherical sharpness with learning rate $2h_t^2$.

Periodic Behavior of h_t and $\tilde{\eta}_t$. It remains to analyze the dynamics of h_t so that we can know how fast the sharpness reduction is. Our analysis is inspired by an empirical study from Lobacheva et al. [84], which reveals a periodic behavior of gradients and effective learning rates in training normalized nets with weight decay. In our theoretical setting, we capture this periodic behavior by showing that h_t and $\tilde{\eta}_t$ do evolve periodically. See Figures 5c and 5d for an example.

The key is that $\tilde{\eta}_t$ changes as an adaptive gradient method: $\tilde{\eta}_t$ increases when gradient is small and decreases when gradient is large (due to the effect of WD; see Figures 3a and 3b), and in our case the gradient norm scales as $|h_t|$ since $\nabla \mathcal{L}(\theta_t) \approx h_t \lambda_1^H(\phi_t) v_1^H(\phi_t)$. By the power method approximation,

²Our construction of ϕ_t ensures that x_t only has a small overlap with the 1-eigenspace of $I - \tilde{\eta}_t H(\phi_t)$, so x_t can only align to $\pm v_1^H(\phi_t)$.

$h_{t+2} \approx (1 - \tilde{\eta}_t \lambda_1^H(\phi_t))^2 h_t$, so $|h_t|$ decreases when $\tilde{\eta}_t < 2/\lambda_1^H(\phi_t)$. But $|h_t|$ cannot decrease forever, since $\tilde{\eta}_t$ increases when $|h_t|$ is sufficiently small. When $\tilde{\eta}_t$ rises to over $2/\lambda_1^H(\phi_t)$, $|h_t|$ changes from decreasing to increasing according to our approximation. But h_t cannot increase indefinitely either, since $\tilde{\eta}_t$ decreases when $|h_t|$ is sufficiently large. A period finishes when $\tilde{\eta}_t < 2/\lambda_1^H(\phi_t)$ holds again.

In our theoretical analysis, we connect this periodic behavior with a 1-dimensional Hamiltonian system (see Appendix H.2), and show that $2h_t^2$ in each step can be approximated by its average value in the period without incurring a large error. Further calculations show that this average value is approximately $\frac{2\eta_{\text{in}}}{4 + \|\nabla_{\Gamma} \log \lambda_1^H(\zeta(t\eta_{\text{in}}))\|_2}$, the learning rate in the flow (3) multiplied with η_{in} . We can therefore conclude that each step of ϕ_t (or θ_t) tracks a time interval of η_{in} in the flow.

Extensions. We note that this periodic behavior is not limited to GD+WD on scale-invariant loss, since the above intuitive argument holds as long as the effective LR changes adaptively with respect to gradient change. Based on this intuition, an important notion called *Quasi-RMSprop scheduler* is proposed. For a PGD method, a learning rate scheduler is a rule for changing the effective LR in each step, and Quasi-RMSprop is a specific class of schedulers we define, including the way that the effective LR changes in GD+WD on scale-invariant loss (if viewed as PGD). Our proof is done in a unified way that works as long as the effective LR changes in each step according to a Quasi-RMSprop scheduler. As a by-product, a similar theorem can be proved for GD (without projection) on non-scale-invariant loss if the LR changes as a Quasi-RMSprop in each step. For example, we can extend our analysis to RMSprop with a scalar learning rate. See Appendix B.

5 Case Study: Linear Regression with Batch Normalization

In this section, we analyze the GD+WD dynamics on linear regression with Batch Normalization (BN), as a simple application of our theory. Let $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$ be a dataset, where $\mathbf{x}_i \in \mathbb{R}^d$ and $y_i \in \mathbb{R}$ are inputs and regression targets. We study the over-parameterized case where $d \gg n$, and we assume that the regression targets are generated by an unknown linear model.

A classic linear model is parameterized by $(\mathbf{w}, b) \in \mathbb{R}^d \times \mathbb{R}$ and outputs $\mathbf{w}^\top \mathbf{x} + b$ given input \mathbf{x} , but now we add a BN to the output. More specifically, we consider a batch-normalized linear model $\Phi(\mathbf{x}; \mathbf{w}, \gamma, \beta) := \gamma \cdot \frac{\mathbf{w}^\top \mathbf{x} - \mu_1}{\sigma_1} + \beta$, where μ_1, σ_1 are the mean and standard deviation of $\{\mathbf{w}^\top \mathbf{x}_i\}_{i=1}^n$ over the whole dataset³, and the bias term b is cancelled out due to BN. Note that $\Phi(\mathbf{x}; \mathbf{w}, \gamma, \beta)$ is still a linear function with respect to \mathbf{x} . Let $\boldsymbol{\mu}_x \in \mathbb{R}^d$ and $\boldsymbol{\Sigma}_x \in \mathbb{R}^{d \times d}$ be the mean and covariance of the input data $\{\mathbf{x}_i\}_{i=1}^n$. Then $\Phi(\mathbf{x}; \mathbf{w}, \gamma, \beta)$ can be rewritten as:

$$\Phi(\mathbf{x}; \mathbf{w}, \gamma, \beta) = \tilde{\mathbf{w}}^\top \mathbf{x} + \tilde{b}, \quad \text{where } \tilde{\mathbf{w}} := \gamma \mathbf{w} / \|\mathbf{w}\|_{\boldsymbol{\Sigma}_x}, \quad \tilde{b} := \beta - \tilde{\mathbf{w}}^\top \boldsymbol{\mu}_x. \quad (4)$$

No matter how \mathbf{w} is set, the output mean and variance of Φ are always β and γ^2 . To simplify our analysis, we fix β, γ to be non-trainable constants so that the mean and variance of Φ 's output match with those of $\{y_i\}_{i=1}^n$, that is, we set $\beta = \mu_y$ and $\gamma = \sigma_y$ to be the mean and standard deviation of y_i over the whole dataset. Then the training loss is $\mathcal{L}(\mathbf{w}) := \frac{1}{n} \sum_{i \in [n]} (\Phi(\mathbf{x}_i; \mathbf{w}, \gamma, \beta) - y_i)^2$.

Theorem 5.1. *In our setting of linear regression with BN, the sharpness-reduction flow ζ defined in (3) converges to the solution $\mathbf{w}^* \in \mathbb{S}^{d-1}$ that minimizes sharpness $\lambda_1^H(\mathbf{w}^*)$ on Γ , regardless of the initialization. Moreover, the coefficients $(\tilde{\mathbf{w}}, \tilde{b})$ associated with \mathbf{w}^* (defined in (4)) are the optimal solution of the following constrained optimization problem (M):*

$$\min \quad \|\mathbf{w}\|_2^2 \quad \text{s.t.} \quad \mathbf{w}^\top \mathbf{x}_i + b = y_i, \quad \forall i \in [n]. \quad (\text{M})$$

At first sight the result may appear trivial because the intent of WD is to regularize L^2 -norm. But this is deceptive because in scale-invariant nets the regularization effect of WD is not explicit. This result also challenges conventional view of optimization. GD is usually viewed as a discretization of its continuous counterpart, gradient flow (GF), and theoretical insight for the discrete update including convergence rate and implicit bias is achieved by analyzing the continuous counterpart (See Appendix A for a list). However, GF does not have the same sharpness-reduction bias as GD. As discussed in [77], adding WD only performs a time-rescaling on the GF trajectory on scale-invariant loss, but does not change the point that GF converge to if we project the trajectory onto the unit sphere. One can easily show that GF may converge to any zero-loss solution, but no matter how small

³Note that the batch size is n here as we are running full-batch GD

LR is, GD exhibits the sharpness-reduction bias towards the optimal solution of (M). To our best knowledge, this result is the first concrete example where even with arbitrarily small LR, GD can still generalize better than GF under natural settings.

6 Discussion

Experimental Verification of Sharpness Reduction. Besides Figures 1 and 2, Appendix P.1 provides additional matrix completion experiments with different data size, and Appendix P.2 provides CIFAR-10 experiments with ResNet-20. In all these experiments, we observed that GD continues to improve the test accuracy even after fitting the training set, and this phenomenon is correlated with the decreasing trend of spherical sharpness. See also Appendix P.3 for the validation for the periodic behavior we analyze in theory.

Ablation Studies on Normalization and Weight Decay. Our theoretical analysis crucially relies on the interplay between normalization and WD to establish the sharpness-reduction flow. We also conducted ablation studies on normalization and WD to highlight the importance of this interplay. First, if normalization is removed, the spherical sharpness becomes undefined, and we do not know if GD implicitly minimizes any sharpness measure. But even if a similar measure does exist, it cannot be strongly related to generalization, because we can verify that the test accuracy becomes very bad without normalization (56.8% on CIFAR-10, Figure 14), and continuing training after fitting the training set no longer improves test accuracy. Second, if WD is removed, the analysis in Arora et al. [7] guarantees convergence in the stable regime, and we can verify that the spherical sharpness and test accuracy stop changing when the loss is small. The final test accuracy is stuck at 66.4% (Figure 15), whereas training with WD leads to 84.3%.

Explaining the Progressive Sharpening and EoS Phenomena. Cohen et al. [24] conducted extensive empirical studies on the dynamics of GD in deep learning (without weight decay), formally $w_{t+1} \leftarrow w_t - \hat{\eta} \nabla \tilde{\mathcal{L}}(w_t)$. They observed the *progressive sharpening* phenomenon: $\lambda_1(\nabla^2 \tilde{\mathcal{L}}(w_t))$ tends to increase so long as it is less than $2/\hat{\eta}$. Then they observed that the training typically enters the EoS regime, which they define as a regime that (1) $\lambda_1(\nabla^2 \tilde{\mathcal{L}}(w_t))$ hovers right at, or just above $2/\hat{\eta}$; and (2) the training loss $\tilde{\mathcal{L}}(w_t)$ goes up and down over short timescales, yet still decreases in the long-term run. A recent research trend focuses on explaining the progressive sharpening and EoS phenomena [1, 91, 8, 18]. Our work corresponds to an important special case where $\tilde{\mathcal{L}}(w)$ is a scale-invariant loss with L^2 -regularization, namely $\mathcal{L}(w) + \frac{\lambda}{2} \|w\|_2^2$. By analyzing the interplay between normalization and WD, the first part of our results (Section 4.1) attributes progressive sharpening to norm change, and the second part (Section 4.2) justifies in theory that the training can make progress in the EoS regime. See Appendix G.3 for more discussion.

7 Conclusions and Future Work

We exhibited settings where gradient descent has an implicit bias to reduce spherical sharpness in training neural nets with normalization layers and weight decay, and we verified experimentally this sharpness-reduction bias predicted by our theorem as well as its generalization benefit on CIFAR-10.

Our theoretical analysis applies to dynamics around a minimizer manifold and requires a small (but finite) learning rate so that we can show that the parameter oscillates locally and approximately tracks a sharpness-reduction flow. We note that in practice a decrease in spherical sharpness is observed even with moderate LR and even before getting close to a minimizer manifold. Explaining these phenomena is left for future work. Now we list some other future directions. The first is to generalize our results to SGD, where the sharpness measure may not be the spherical sharpness and could depend on the structure of gradient noise. Second, to understand the benefit of reducing spherical sharpness on specific tasks, e.g., why does reducing spherical sharpness encourage low-rank on matrix completion with BN (Figure 1)? Third, to study sharpness-reduction bias for neural net architectures that are not scale-invariant on all parameters (e.g., with certain unnormalized layers).

Acknowledgements

This work is funded by NSF, ONR, Simons Foundation, DARPA and SRC. ZL is also supported by Microsoft Research PhD Fellowship.

References

- [1] Kwangjun Ahn, Jingzhao Zhang, and Suvrit Sra. Understanding the unstable convergence of gradient descent. In Kamalika Chaudhuri, Stefanie Jegelka, Le Song, Csaba Szepesvari, Gang Niu, and Sivan Sabato, editors, *Proceedings of the 39th International Conference on Machine Learning*, volume 162 of *Proceedings of Machine Learning Research*, pages 247–257. PMLR, 17–23 Jul 2022.
- [2] Zeyuan Allen-Zhu, Yuanzhi Li, and Yingyu Liang. Learning and generalization in over-parameterized neural networks, going beyond two layers. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d’Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019.
- [3] Zeyuan Allen-Zhu, Yuanzhi Li, and Zhao Song. A convergence theory for deep learning via over-parameterization. In Kamalika Chaudhuri and Ruslan Salakhutdinov, editors, *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, pages 242–252. PMLR, 09–15 Jun 2019.
- [4] Sanjeev Arora, Nadav Cohen, Wei Hu, and Yuping Luo. Implicit regularization in deep matrix factorization. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d’ Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems 32*, pages 7411–7422. Curran Associates, Inc., 2019.
- [5] Sanjeev Arora, Simon Du, Wei Hu, Zhiyuan Li, and Ruosong Wang. Fine-grained analysis of optimization and generalization for overparameterized two-layer neural networks. In *International Conference on Machine Learning*, pages 322–332. PMLR, 2019.
- [6] Sanjeev Arora, Simon S Du, Wei Hu, Zhiyuan Li, Russ R Salakhutdinov, and Ruosong Wang. On exact computation with an infinitely wide neural net. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d’ Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems 32*, pages 8139–8148. Curran Associates, Inc., 2019.
- [7] Sanjeev Arora, Zhiyuan Li, and Kaifeng Lyu. Theoretical analysis of auto rate-tuning by batch normalization. In *International Conference on Learning Representations*, 2019.
- [8] Sanjeev Arora, Zhiyuan Li, and Abhishek Panigrahi. Understanding gradient descent on the edge of stability in deep learning. In Kamalika Chaudhuri, Stefanie Jegelka, Le Song, Csaba Szepesvari, Gang Niu, and Sivan Sabato, editors, *Proceedings of the 39th International Conference on Machine Learning*, volume 162 of *Proceedings of Machine Learning Research*, pages 948–1024. PMLR, 17–23 Jul 2022.
- [9] Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E Hinton. Layer normalization. *arXiv preprint arXiv:1607.06450*, 2016.
- [10] David Balduzzi, Marcus Frean, Lennox Leary, J. P. Lewis, Kurt Wan-Duo Ma, and Brian McWilliams. The shattered gradients problem: If resnets are the answer, then what is the question? In Doina Precup and Yee Whye Teh, editors, *Proceedings of the 34th International Conference on Machine Learning*, volume 70 of *Proceedings of Machine Learning Research*, pages 342–350. PMLR, 06–11 Aug 2017.
- [11] David Barrett and Benoit Dherin. Implicit gradient regularization. In *International Conference on Learning Representations*, 2021.
- [12] Johan Bjorck, Carla Gomes, and Bart Selman. Understanding batch normalization. *arXiv preprint arXiv:1806.02375*, 2018.
- [13] Guy Blanc, Neha Gupta, Gregory Valiant, and Paul Valiant. Implicit regularization for deep neural networks driven by an ornstein-uhlenbeck like process. In Jacob Abernethy and Shivani Agarwal, editors, *Proceedings of Thirty Third Conference on Learning Theory*, volume 125 of *Proceedings of Machine Learning Research*, pages 483–513. PMLR, 09–12 Jul 2020.
- [14] Andrew Brock, Soham De, and Samuel L Smith. Characterizing signal propagation to close the performance gap in unnormalized resnets. In *International Conference on Learning Representations*, 2021.

- [15] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. Language models are few-shot learners. In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33, pages 1877–1901. Curran Associates, Inc., 2020.
- [16] Yongqiang Cai, Qianxiao Li, and Zuowei Shen. A quantitative analysis of the effect of batch normalization on gradient descent. In Kamalika Chaudhuri and Ruslan Salakhutdinov, editors, *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, pages 882–890. PMLR, 09–15 Jun 2019.
- [17] Yuan Cao and Quanquan Gu. Generalization bounds of stochastic gradient descent for wide and deep neural networks. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019.
- [18] Lei Chen and Joan Bruna. On gradient descent convergence beyond the edge of stability. *arXiv preprint arXiv:2206.04172*, 2022.
- [19] Zixiang Chen, Yuan Cao, Difan Zou, and Quanquan Gu. How much over-parameterization is sufficient to learn deep ReLU networks? In *International Conference on Learning Representations*, 2021.
- [20] Yuejie Chi, Yue M. Lu, and Yuxin Chen. Nonconvex optimization meets low-rank matrix factorization: An overview. *IEEE Transactions on Signal Processing*, 67(20):5239–5269, 2019. doi: 10.1109/TSP.2019.2937282.
- [21] Vitaliy Chiley, Ilya Sharapov, Atli Kosson, Urs Koster, Ryan Reece, Sofia Samaniego de la Fuente, Vishal Subbiah, and Michael James. Online normalization for training neural networks. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019.
- [22] Léniaic Chizat and Francis Bach. Implicit bias of gradient descent for wide two-layer neural networks trained with the logistic loss. In *Conference on Learning Theory*, pages 1305–1338. PMLR, 2020.
- [23] Léniaic Chizat, Edouard Oyallon, and Francis Bach. On lazy training in differentiable programming. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems 32*, pages 2937–2947. Curran Associates, Inc., 2019.
- [24] 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.
- [25] Alex Damian, Tengyu Ma, and Jason D Lee. Label noise SGD provably prefers flat global minimizers. In M. Ranzato, A. Beygelzimer, Y. Dauphin, P.S. Liang, and J. Wortman Vaughan, editors, *Advances in Neural Information Processing Systems*, volume 34, pages 27449–27461. Curran Associates, Inc., 2021.
- [26] Hadi Daneshmand, Jonas Kohler, Francis Bach, Thomas Hofmann, and Aurelien Lucchi. Batch normalization provably avoids ranks collapse for randomly initialised deep networks. In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33, pages 18387–18398. Curran Associates, Inc., 2020.
- [27] Hadi Daneshmand, Amir Joudaki, and Francis Bach. Batch normalization orthogonalizes representations in deep random networks. In A. Beygelzimer, Y. Dauphin, P. Liang, and J. Wortman Vaughan, editors, *Advances in Neural Information Processing Systems*, 2021.

- [28] Soham De and Sam Smith. Batch normalization biases residual blocks towards the identity function in deep networks. In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33, pages 19964–19975. Curran Associates, Inc., 2020.
- [29] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota, June 2019. Association for Computational Linguistics. doi: 10.18653/v1/N19-1423.
- [30] Lijun Ding, Dmitriy Drusvyatskiy, and Maryam Fazel. Flat minima generalize for low-rank matrix recovery. *arXiv preprint arXiv:2203.03756*, 2022.
- [31] Laurent Dinh, Razvan Pascanu, Samy Bengio, and Yoshua Bengio. Sharp minima can generalize for deep nets. In Doina Precup and Yee Whye Teh, editors, *Proceedings of the 34th International Conference on Machine Learning*, volume 70 of *Proceedings of Machine Learning Research*, pages 1019–1028. PMLR, 06–11 Aug 2017.
- [32] Simon Du, Jason Lee, Haochuan Li, Liwei Wang, and Xiyu Zhai. Gradient descent finds global minima of deep neural networks. In *International Conference on Machine Learning*, pages 1675–1685. PMLR, 2019.
- [33] Yonatan Dukler, Quanquan Gu, and Guido Montufar. Optimization theory for ReLU neural networks trained with normalization layers. In Hal Daumé III and Aarti Singh, editors, *Proceedings of the 37th International Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning Research*, pages 2751–2760. PMLR, 13–18 Jul 2020.
- [34] K. J. Falconer. Differentiation of the Limit Mapping in a Dynamical System. *Journal of the London Mathematical Society*, s2-27(2):356–372, 04 1983. ISSN 0024-6107. doi: 10.1112/jlms/s2-27.2.356.
- [35] Benjamin Fehrman, Benjamin Gess, and Arnulf Jentzen. Convergence rates for the stochastic gradient descent method for non-convex objective functions. *Journal of Machine Learning Research*, 21(136):1–48, 2020.
- [36] Robert L. Foote. Shorter notes: Regularity of the distance function. *Proceedings of the American Mathematical Society*, 92(1):153–155, 1984. ISSN 00029939, 10886826.
- [37] Pierre Foret, Ariel Kleiner, Hossein Mobahi, and Behnam Neyshabur. Sharpness-aware minimization for efficiently improving generalization. In *International Conference on Learning Representations*, 2021.
- [38] Rong Ge, Yunwei Ren, Xiang Wang, and Mo Zhou. Understanding deflation process in over-parametrized tensor decomposition. *Advances in Neural Information Processing Systems*, 34, 2021.
- [39] Behrooz Ghorbani, Shankar Krishnan, and Ying Xiao. An investigation into neural net optimization via hessian eigenvalue density. In Kamalika Chaudhuri and Ruslan Salakhutdinov, editors, *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, pages 2232–2241. PMLR, 09–15 Jun 2019.
- [40] 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.
- [41] Suriya Gunasekar, Blake E Woodworth, Srinadh Bhojanapalli, Behnam Neyshabur, and Nati Srebro. Implicit regularization in matrix factorization. In I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing Systems 30*, pages 6151–6159. Curran Associates, Inc., 2017.

- [42] Suriya Gunasekar, Jason D Lee, Daniel Soudry, and Nati Srebro. Implicit bias of gradient descent on linear convolutional networks. In S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett, editors, *Advances in Neural Information Processing Systems 31*, pages 9482–9491. Curran Associates, Inc., 2018.
- [43] Suriya Gunasekar, Jason D Lee, Daniel Soudry, and Nati Srebro. Implicit bias of gradient descent on linear convolutional networks. *Advances in Neural Information Processing Systems*, 31, 2018.
- [44] Jeff Z. HaoChen, Colin Wei, Jason Lee, and Tengyu Ma. Shape matters: Understanding the implicit bias of the noise covariance. In Mikhail Belkin and Samory Kpotufe, editors, *Proceedings of Thirty Fourth Conference on Learning Theory*, volume 134 of *Proceedings of Machine Learning Research*, pages 2315–2357. PMLR, 15–19 Aug 2021.
- [45] Haowei He, Gao Huang, and Yang Yuan. Asymmetric valleys: Beyond sharp and flat local minima. *Advances in neural information processing systems*, 32, 2019.
- [46] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.
- [47] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Identity mappings in deep residual networks. In Bastian Leibe, Jiri Matas, Nicu Sebe, and Max Welling, editors, *Computer Vision – ECCV 2016*, pages 630–645, Cham, 2016. Springer International Publishing. ISBN 978-3-319-46493-0.
- [48] Dan Hendrycks and Kevin Gimpel. Gaussian error linear units (GELUs). *arXiv preprint arXiv:1606.08415*, 2016.
- [49] Geoffrey Hinton, Nitish Srivastava, and Kevin Swersky. Neural networks for machine learning lecture 6a: Overview of mini-batch gradient descent. Technical report, 2012. URL https://www.cs.toronto.edu/~tijmen/csc321/slides/lecture_slides_lec6.pdf.
- [50] Sepp Hochreiter and Jürgen Schmidhuber. Flat minima. *Neural computation*, 9(1):1–42, 1997.
- [51] Elad Hoffer, Itay Hubara, and Daniel Soudry. Train longer, generalize better: closing the generalization gap in large batch training of neural networks. In I. Guyon, U. Von Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc., 2017.
- [52] Elad Hoffer, Ron Banner, Itay Golan, and Daniel Soudry. Norm matters: efficient and accurate normalization schemes in deep networks. In S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 31. Curran Associates, Inc., 2018.
- [53] Lei Huang, Xianglong Liu, Yang Liu, Bo Lang, and Dacheng Tao. Centered weight normalization in accelerating training of deep neural networks. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, Oct 2017.
- [54] Hikaru Ibayashi and Masaaki Imaizumi. Exponential escape efficiency of SGD from sharp minima in non-stationary regime. *arXiv preprint arXiv:2111.04004*, 2021.
- [55] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In Francis Bach and David Blei, editors, *Proceedings of the 32nd International Conference on Machine Learning*, volume 37 of *Proceedings of Machine Learning Research*, pages 448–456, Lille, France, 07–09 Jul 2015. PMLR.
- [56] Arthur Jacot, Franck Gabriel, and Clement Hongler. Neural tangent kernel: Convergence and generalization in neural networks. In S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett, editors, *Advances in Neural Information Processing Systems 31*, pages 8571–8580. Curran Associates, Inc., 2018.

- [57] 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.
- [58] Stanislaw Jastrzebski, 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.
- [59] Ziwei Ji and Matus Telgarsky. Directional convergence and alignment in deep learning. In H. Larochelle, M. Ranzato, R. Hadsell, M. F. Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33, pages 17176–17186. Curran Associates, Inc., 2020.
- [60] Yiding Jiang, Behnam Neyshabur, Hossein Mobahi, Dilip Krishnan, and Samy Bengio. Fantastic generalization measures and where to find them. In *International Conference on Learning Representations*, 2020.
- [61] Ryo Karakida, Shotaro Akaho, and Shun-ichi Amari. The normalization method for alleviating pathological sharpness in wide neural networks. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019.
- [62] Hamed Karimi, Julie Nutini, and Mark Schmidt. Linear convergence of gradient and proximal-gradient methods under the polyak-łojasiewicz condition. In *European Conference on Machine Learning and Knowledge Discovery in Databases - Volume 9851, ECML PKDD 2016*, pages 795–811, Berlin, Heidelberg, 2016. Springer-Verlag. ISBN 9783319461274. doi: 10.1007/978-3-319-46128-1_50.
- [63] 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. In *International Conference on Learning Representations*, 2017.
- [64] Bobby Kleinberg, Yuanzhi Li, and Yang Yuan. An alternative view: When does SGD escape local minima? In Jennifer Dy and Andreas Krause, editors, *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pages 2698–2707. PMLR, 10–15 Jul 2018.
- [65] Jonas Kohler, Hadi Daneshmand, Aurelien Lucchi, Thomas Hofmann, Ming Zhou, and Klaus Neymeyr. Exponential convergence rates for batch normalization: The power of length-direction decoupling in non-convex optimization. In Kamalika Chaudhuri and Masashi Sugiyama, editors, *Proceedings of the Twenty-Second International Conference on Artificial Intelligence and Statistics*, volume 89 of *Proceedings of Machine Learning Research*, pages 806–815. PMLR, 16–18 Apr 2019.
- [66] Ling kai Kong and Molei Tao. Stochasticity of deterministic gradient descent: Large learning rate for multiscale objective function. In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33, pages 2625–2638. Curran Associates, Inc., 2020.
- [67] Antoine Labatie, Dominic Masters, Zach Eaton-Rosen, and Carlo Luschi. Proxy-normalizing activations to match batch normalization while removing batch dependence. In M. Ranzato, A. Beygelzimer, Y. Dauphin, P.S. Liang, and J. Wortman Vaughan, editors, *Advances in Neural Information Processing Systems*, volume 34, pages 16990–17006. Curran Associates, Inc., 2021.
- [68] Susanna Lange, Kyle Helfrich, and Qiang Ye. Batch normalization preconditioning for neural network training. *Journal of Machine Learning Research*, 23(72):1–41, 2022.
- [69] Beatrice Laurent and Pascal Massart. Adaptive estimation of a quadratic functional by model selection. *Annals of Statistics*, pages 1302–1338, 2000.

- [70] Jason D. Lee, Max Simchowitz, Michael I. Jordan, and Benjamin Recht. Gradient descent only converges to minimizers. In Vitaly Feldman, Alexander Rakhlin, and Ohad Shamir, editors, *29th Annual Conference on Learning Theory*, volume 49 of *Proceedings of Machine Learning Research*, pages 1246–1257, Columbia University, New York, New York, USA, 23–26 Jun 2016. PMLR.
- [71] Jason D Lee, Ioannis Panageas, Georgios Piliouras, Max Simchowitz, Michael I Jordan, and Benjamin Recht. First-order methods almost always avoid saddle points. *arXiv preprint arXiv:1710.07406*, 2017.
- [72] Aitor Lewkowycz and Guy Gur-Ari. On the training dynamics of deep networks with L₂ regularization. In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33, pages 4790–4799. Curran Associates, Inc., 2020.
- [73] Aitor Lewkowycz, Yasaman Bahri, Ethan Dyer, Jascha Sohl-Dickstein, and Guy Gur-Ari. The large learning rate phase of deep learning: the catapult mechanism. *arXiv preprint arXiv:2003.02218*, 2020.
- [74] Hao Li, Zheng Xu, Gavin Taylor, Christoph Studer, and Tom Goldstein. Visualizing the loss landscape of neural nets. *Advances in neural information processing systems*, 31, 2018.
- [75] Yuanzhi Li and Yingyu Liang. Learning overparameterized neural networks via stochastic gradient descent on structured data. In S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 31. Curran Associates, Inc., 2018.
- [76] Yuanzhi Li, Tengyu Ma, and Hongyang Zhang. Algorithmic regularization in overparameterized matrix sensing and neural networks with quadratic activations. In Sébastien Bubeck, Vianney Perchet, and Philippe Rigollet, editors, *Proceedings of the 31st Conference On Learning Theory*, volume 75 of *Proceedings of Machine Learning Research*, pages 2–47. PMLR, 06–09 Jul 2018.
- [77] Zhiyuan Li and Sanjeev Arora. An exponential learning rate schedule for deep learning. In *International Conference on Learning Representations*, 2020.
- [78] Zhiyuan Li, Kaifeng Lyu, and Sanjeev Arora. Reconciling modern deep learning with traditional optimization analyses: The intrinsic learning rate. In H. Larochelle, M. Ranzato, R. Hadsell, M. F. Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33, pages 14544–14555. Curran Associates, Inc., 2020.
- [79] Zhiyuan Li, Yuping Luo, and Kaifeng Lyu. Towards resolving the implicit bias of gradient descent for matrix factorization: Greedy low-rank learning. In *International Conference on Learning Representations*, 2021.
- [80] Zhiyuan Li, Srinadh Bhojanapalli, Manzil Zaheer, Sashank Reddi, and Sanjiv Kumar. Robust training of neural networks using scale invariant architectures. In Kamalika Chaudhuri, Stefanie Jegelka, Le Song, Csaba Szepesvari, Gang Niu, and Sivan Sabato, editors, *Proceedings of the 39th International Conference on Machine Learning*, volume 162 of *Proceedings of Machine Learning Research*, pages 12656–12684. PMLR, 17–23 Jul 2022.
- [81] Zhiyuan Li, Tianhao Wang, and Sanjeev Arora. What happens after SGD reaches zero loss? –a mathematical framework. In *International Conference on Learning Representations*, 2022.
- [82] Zinan Lin, Vyas Sekar, and Giulia Fanti. Why spectral normalization stabilizes GANs: Analysis and improvements. In M. Ranzato, A. Beygelzimer, Y. Dauphin, P.S. Liang, and J. Wortman Vaughan, editors, *Advances in Neural Information Processing Systems*, volume 34, pages 9625–9638. Curran Associates, Inc., 2021.
- [83] Shengchao Liu, Dimitris Papailiopoulos, and Dimitris Achlioptas. Bad global minima exist and sgd can reach them. In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33, pages 8543–8552. Curran Associates, Inc., 2020.

- [84] Ekaterina Lobacheva, Maxim Kodryan, Nadezhda Chirkova, Andrey Malinin, and Dmitry P Vetrov. On the periodic behavior of neural network training with batch normalization and weight decay. In M. Ranzato, A. Beygelzimer, Y. Dauphin, P.S. Liang, and J. Wortman Vaughan, editors, *Advances in Neural Information Processing Systems*, volume 34, pages 21545–21556. Curran Associates, Inc., 2021.
- [85] Ekdeep S Lubana, Robert Dick, and Hidenori Tanaka. Beyond batchnorm: Towards a unified understanding of normalization in deep learning. In M. Ranzato, A. Beygelzimer, Y. Dauphin, P.S. Liang, and J. Wortman Vaughan, editors, *Advances in Neural Information Processing Systems*, volume 34, pages 4778–4791. Curran Associates, Inc., 2021.
- [86] Ping Luo, Xinjiang Wang, Wenqi Shao, and Zhanglin Peng. Towards understanding regularization in batch normalization. In *International Conference on Learning Representations*, 2019.
- [87] Kaifeng Lyu and Jian Li. Gradient descent maximizes the margin of homogeneous neural networks. In *International Conference on Learning Representations*, 2020.
- [88] Kaifeng Lyu, Zhiyuan Li, Runzhe Wang, and Sanjeev Arora. Gradient descent on two-layer nets: Margin maximization and simplicity bias. *Advances in Neural Information Processing Systems*, 34, 2021.
- [89] Chao Ma and Lexing Ying. On linear stability of SGD and input-smoothness of neural networks. In A. Beygelzimer, Y. Dauphin, P. Liang, and J. Wortman Vaughan, editors, *Advances in Neural Information Processing Systems*, 2021.
- [90] Chao Ma and Lexing Ying. A Riemannian mean field formulation for two-layer neural networks with batch normalization. *Research in the Mathematical Sciences*, 9(3):47, July 2022. ISSN 2197-9847.
- [91] Chao Ma, Daniel Kunin, Lei Wu, and Lexing Ying. Beyond the quadratic approximation: The multiscale structure of neural network loss landscapes. *Journal of Machine Learning*, 1(3): 247–267, 2022. ISSN 2790-2048.
- [92] David McAllester. Simplified PAC-Bayesian margin bounds. In *Learning theory and Kernel machines*, pages 203–215. Springer, 2003.
- [93] Rotem Mulayoff and Tomer Michaeli. Unique properties of flat minima in deep networks. In Hal Daumé III and Aarti Singh, editors, *Proceedings of the 37th International Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning Research*, pages 7108–7118. PMLR, 13–18 Jul 2020.
- [94] Rotem Mulayoff, Tomer Michaeli, and Daniel Soudry. The implicit bias of minima stability: A view from function space. In M. Ranzato, A. Beygelzimer, Y. Dauphin, P.S. Liang, and J. Wortman Vaughan, editors, *Advances in Neural Information Processing Systems*, volume 34, pages 17749–17761. Curran Associates, Inc., 2021.
- [95] Behnam Neyshabur, Srinadh Bhojanapalli, David McAllester, and Nati Srebro. Exploring generalization in deep learning. *Advances in neural information processing systems*, 30, 2017.
- [96] Ioannis Panageas and Georgios Piliouras. Gradient Descent Only Converges to Minimizers: Non-Isolated Critical Points and Invariant Regions. In Christos H. Papadimitriou, editor, *8th Innovations in Theoretical Computer Science Conference (ITCS 2017)*, volume 67 of *Leibniz International Proceedings in Informatics (LIPIcs)*, pages 2:1–2:12, Dagstuhl, Germany, 2017. Schloss Dagstuhl–Leibniz-Zentrum fuer Informatik. ISBN 978-3-95977-029-3. doi: 10.4230/LIPIcs.ITCS.2017.2.
- [97] Siyuan Qiao, Huiyu Wang, Chenxi Liu, Wei Shen, and Alan Yuille. Micro-batch training with batch-channel normalization and weight standardization. *arXiv preprint arXiv:1903.10520*, 2019.
- [98] Prajit Ramachandran, Barret Zoph, and Quoc V Le. Searching for activation functions. *arXiv preprint arXiv:1710.05941*, 2017.

- [99] Akshay Rangamani, Nam H. Nguyen, Abhishek Kumar, Dzung Phan, Sang Peter Chin, and Trac D. Tran. A scale invariant measure of flatness for deep network minima. In *ICASSP 2021 - 2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 1680–1684, 2021.
- [100] Noam Razin and Nadav Cohen. Implicit regularization in deep learning may not be explainable by norms. In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33, pages 21174–21187. Curran Associates, Inc., 2020.
- [101] Noam Razin, Asaf Maman, and Nadav Cohen. Implicit regularization in hierarchical tensor factorization and deep convolutional neural networks. In Kamalika Chaudhuri, Stefanie Jegelka, Le Song, Csaba Szepesvari, Gang Niu, and Sivan Sabato, editors, *Proceedings of the 39th International Conference on Machine Learning*, volume 162 of *Proceedings of Machine Learning Research*, pages 18422–18462. PMLR, 17–23 Jul 2022.
- [102] Tim Salimans and Durk P Kingma. Weight normalization: A simple reparameterization to accelerate training of deep neural networks. In D. Lee, M. Sugiyama, U. Luxburg, I. Guyon, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 29. Curran Associates, Inc., 2016.
- [103] Shibani Santurkar, Dimitris Tsipras, Andrew Ilyas, and Aleksander Madry. How does batch normalization help optimization? In S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett, editors, *Advances in Neural Information Processing Systems 31*, pages 2483–2493. Curran Associates, Inc., 2018.
- [104] Alexander Shekhovtsov and Boris Flach. Stochastic normalizations as bayesian learning. In C. V. Jawahar, Hongdong Li, Greg Mori, and Konrad Schindler, editors, *Computer Vision – ACCV 2018*, pages 463–479, Cham, 2019. Springer International Publishing. ISBN 978-3-030-20890-5.
- [105] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. In *International Conference on Learning Representations*, 2015.
- [106] Daniel Soudry, Elad Hoffer, Mor Shpigel Nacson, Suriya Gunasekar, and Nathan Srebro. The implicit bias of gradient descent on separable data. *Journal of Machine Learning Research*, 19(70):1–57, 2018.
- [107] Daniel Soudry, Elad Hoffer, and Nathan Srebro. The implicit bias of gradient descent on separable data. In *International Conference on Learning Representations*, 2018.
- [108] Dominik Stöger and Mahdi Soltanolkotabi. Small random initialization is akin to spectral learning: Optimization and generalization guarantees for overparameterized low-rank matrix reconstruction. *Advances in Neural Information Processing Systems*, 34, 2021.
- [109] Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe, Jon Shlens, and Zbigniew Wojna. Rethinking the inception architecture for computer vision. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2016.
- [110] Hidenori Tanaka and Daniel Kunin. Noether’s learning dynamics: Role of symmetry breaking in neural networks. In M. Ranzato, A. Beygelzimer, Y. Dauphin, P.S. Liang, and J. Wortman Vaughan, editors, *Advances in Neural Information Processing Systems*, volume 34, pages 25646–25660. Curran Associates, Inc., 2021.
- [111] Mattias Teye, Hossein Azizpour, and Kevin Smith. Bayesian uncertainty estimation for batch normalized deep networks. In Jennifer Dy and Andreas Krause, editors, *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pages 4907–4916. PMLR, 10–15 Jul 2018.
- [112] Yusuke Tsuzuku, Issei Sato, and Masashi Sugiyama. Normalized flat minima: Exploring scale invariant definition of flat minima for neural networks using PAC-Bayesian analysis. In Hal Daumé III and Aarti Singh, editors, *Proceedings of the 37th International Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning Research*, pages 9636–9647. PMLR, 13–18 Jul 2020.

- [113] Twan van Laarhoven. L2 regularization versus batch and weight normalization. *arXiv preprint arXiv:1706.05350*, 2017.
- [114] Ruosi Wan, Zhanxing Zhu, Xiangyu Zhang, and Jian Sun. Spherical motion dynamics: Learning dynamics of normalized neural network using sgd and weight decay. In M. Ranzato, A. Beygelzimer, Y. Dauphin, P.S. Liang, and J. Wortman Vaughan, editors, *Advances in Neural Information Processing Systems*, volume 34, pages 6380–6391. Curran Associates, Inc., 2021.
- [115] Yi Wang and Zhiren Wang. Three-stage evolution and fast equilibrium for SGD with non-degenerate critical points. In Kamalika Chaudhuri, Stefanie Jegelka, Le Song, Csaba Szepesvari, Gang Niu, and Sivan Sabato, editors, *Proceedings of the 39th International Conference on Machine Learning*, volume 162 of *Proceedings of Machine Learning Research*, pages 23092–23113. PMLR, 17–23 Jul 2022.
- [116] Yuqing Wang, Minshuo Chen, Tuo Zhao, and Molei Tao. Large learning rate tames homogeneity: Convergence and balancing effect. In *International Conference on Learning Representations*, 2022.
- [117] Lei Wu, Zhanxing Zhu, et al. Towards understanding generalization of deep learning: Perspective of loss landscapes. *arXiv preprint arXiv:1706.10239*, 2017.
- [118] Lei Wu, Chao Ma, and Weinan E. How sgd selects the global minima in over-parameterized learning: A dynamical stability perspective. In S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 31. Curran Associates, Inc., 2018.
- [119] Xiaoxia Wu, Edgar Dobriban, Tongzheng Ren, Shanshan Wu, Zhiyuan Li, Suriya Gunasekar, Rachel Ward, and Qiang Liu. Implicit regularization and convergence for weight normalization. In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33, pages 2835–2847. Curran Associates, Inc., 2020.
- [120] Yuxin Wu and Kaiming He. Group normalization. In *Proceedings of the European Conference on Computer Vision (ECCV)*, September 2018.
- [121] Zeke Xie, Issei Sato, and Masashi Sugiyama. A diffusion theory for deep learning dynamics: Stochastic gradient descent exponentially favors flat minima. In *International Conference on Learning Representations*, 2021.
- [122] Greg Yang. Scaling limits of wide neural networks with weight sharing: Gaussian process behavior, gradient independence, and neural tangent kernel derivation. *arXiv preprint arXiv:1902.04760*, 2019.
- [123] Mingyang Yi, Qi Meng, Wei Chen, Zhi-ming Ma, and Tie-Yan Liu. Positively scale-invariant flatness of ReLU neural networks. *arXiv preprint arXiv:1903.02237*, 2019.
- [124] Mingyang Yi, Huishuai Zhang, Wei Chen, Zhi-Ming Ma, and Tie-Yan Liu. Bn-invariant sharpness regularizes the training model to better generalization. In *Proceedings of the Twenty-Eighth International Joint Conference on Artificial Intelligence, IJCAI-19*, pages 4164–4170. International Joint Conferences on Artificial Intelligence Organization, 7 2019.
- [125] Chiyuan Zhang, Samy Bengio, Moritz Hardt, Benjamin Recht, and Oriol Vinyals. Understanding deep learning requires rethinking generalization. In *International Conference on Learning Representations*, 2017.
- [126] Guodong Zhang, Chaoqi Wang, Bowen Xu, and Roger Grosse. Three mechanisms of weight decay regularization. In *International Conference on Learning Representations*, 2019.
- [127] Zhanxing Zhu, Jingfeng Wu, Bing Yu, Lei Wu, and Jinwen Ma. The anisotropic noise in stochastic gradient descent: Its behavior of escaping from sharp minima and regularization effects. In Kamalika Chaudhuri and Ruslan Salakhutdinov, editors, *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, pages 7654–7663. PMLR, 09–15 Jun 2019.
- [128] Difan Zou, Yuan Cao, Dongruo Zhou, and Quanquan Gu. Stochastic gradient descent optimizes over-parameterized deep ReLU networks. *arXiv preprint arXiv:1811.08888*, 2018.

Checklist

1. For all authors...
 - (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s contributions and scope? [Yes]
 - (b) Did you describe the limitations of your work? [Yes] See Sections 4.2.1 and 7.
 - (c) Did you discuss any potential negative societal impacts of your work? [N/A] We are basically a theoretical work studying the generalization mystery in deep learning. We do not see any negative societal impact.
 - (d) Have you read the ethics review guidelines and ensured that your paper conforms to them? [Yes]
2. If you are including theoretical results...
 - (a) Did you state the full set of assumptions of all theoretical results? [Yes] The assumptions of our main theorem for GD+WD on scale-invariant loss are stated in Section 4.2.1. For GD/PGD with Quasi-RMSprop schedulers in general, see Appendix B.3.
 - (b) Did you include complete proofs of all theoretical results? [Yes] See Appendices G and H for the proofs for our main theorems, Appendix O for the proof for the linear example.
3. If you ran experiments...
 - (a) Did you include the code, data, and instructions needed to reproduce the main experimental results (either in the supplemental material or as a URL)? [Yes] See our supplementary material.
 - (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were chosen)? [Yes] See Appendix Q.
 - (c) Did you report error bars (e.g., with respect to the random seed after running experiments multiple times)? [No] Most of our experiments have only run once due to computational constraints, but we verified the sharpness-reduction bias across various settings.
 - (d) Did you include the total amount of compute and the type of resources used (e.g., type of GPUs, internal cluster, or cloud provider)? [Yes] See Appendix Q.
4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
 - (a) If your work uses existing assets, did you cite the creators? [N/A]
 - (b) Did you mention the license of the assets? [N/A]
 - (c) Did you include any new assets either in the supplemental material or as a URL? [N/A]
 - (d) Did you discuss whether and how consent was obtained from people whose data you’re using/curating? [N/A]
 - (e) Did you discuss whether the data you are using/curating contains personally identifiable information or offensive content? [N/A]
5. If you used crowdsourcing or conducted research with human subjects...
 - (a) Did you include the full text of instructions given to participants and screenshots, if applicable? [N/A]
 - (b) Did you describe any potential participant risks, with links to Institutional Review Board (IRB) approvals, if applicable? [N/A]
 - (c) Did you include the estimated hourly wage paid to participants and the total amount spent on participant compensation? [N/A]