

OPTIMIZER SELECTION BASED ON FUNCTION-PROXY PROXIMITY

Anonymous authors

Paper under double-blind review

ABSTRACT

Many machine learning problems involve the challenging task of training a model to fit the training data; this task is especially challenging for nonconvex problems. Many model training algorithms have been proposed, but it is often difficult to determine which algorithm is best suited to a given machine learning problem. To contend with this challenge, we study the effect of loss function curvature shift on optimizers' proxies of the loss function, and obtain a bound on the rate at which a very large family of optimizers (including all of the most prevalent ones) descend towards the loss function's minimum, while only making relatively weak assumptions on the loss function. Uniquely, our bound is tight even in parameter subspaces in which the loss function is concave, which have been shown to bear potential for fast descent while being neglected by existing convergence rate bounds. We demonstrate the applicability of our bound by developing a meta-algorithm for optimizer selection based on it, and validate our meta-algorithm experimentally.

1 INTRODUCTION

Machine learning problems involve training a model to fit a target distribution (such as protein structures, image labels, text, etc.), usually given as a set of samples from said distribution. However, no single existing training algorithm (a.k.a. "optimization algorithm", or "optimizer") is best suited to all machine learning problems - each has its unique strengths and weaknesses (see Vaswani et al. (2020); Sivan et al. (2024); Ruder (2016); Mustapha et al. (2021); Bera & Shrivastava (2020); Zeiler (2012); Duchi et al. (2011b); Xu et al. (2017); Wadia et al. (2021); Mittal et al. (2019); Zhou et al. (2020); Schmidt et al. (2021)), such as generalization capability, convergence rate, saddle-point and flat region evasion capability, robustness to hyperparameter choice, computational complexity per-iteration, memory complexity, etc., and different areas in which it empirically seems to work best. This diversity of optimizers raises the question of how to select an optimizer for a given problem; ideally, we would like to have a meta-algorithm for optimizer selection. Moreover, since the best optimizer may vary throughout the training process and among different model parameters, our ideal algorithm should combine optimizers by continually selecting the best optimizer for each individual parameter online, all throughout the training process.

One way to compare the computational overheads of different methods is experimentally; the literature is rich with such experiments on a variety of applications (Xu et al., 2017; Schmidt et al., 2021). On the theoretical side - convergence rate bounds have been proven for various optimizers. However, due to the wide variety of assumptions, convergence rate metrics, bound parameters (which may be expensive - if not impossible - to compute ahead of time), and tightness of the bounds in all of these works, comparing amongst them and applying these bounds for optimizer selection remains a challenging task. Existing bounds are also rarely applicable to novel optimizers, and are insufficiently tight, as may be seen by their failure to demonstrate the empirically-verified faster convergence of the more sophisticated methods that make use of second-order curvature information instead of exclusively gradients. Lastly, most broadly applicable bounds for non-convex functions fail to sufficiently address the rate of curvature change across the loss function surface. One particular area that suffers from neglect is linear subspaces of the parameter space in which the loss function is concave (meaning that a restriction $f|_{\mathbb{S}} : \mathbb{S} \rightarrow \mathbb{R}$ of the loss function f to a linear subspace \mathbb{S} is locally concave). We believe that more attention should be given to these subspaces in the context of neural network optimization; Alain et al. (2018) and Ghorbani et al. (2019) experimentally demonstrate that

there is much to be gained by taking optimal steps in these subspaces, with gains that are often orders of magnitude greater than the potential gains in convex subspaces.

Our contributions

1. We design a new and refined cubic approximation which gives rise to an easy to compute algorithmic baseline named `ELMO` (Alg. 1), and prove its minmax-optimality¹ up to the third order. We also show the equivalence between our approach and a special preconditioning matrix which results in a natural refinement of Newton’s method.
2. Based on this algorithmic baseline, we design a novel meta-algorithm named `AMOS` (Alg. 2) that enables us to adaptively select between first and second order training methods in order to boost performance.
3. We present an experimental study on NanoGPT that demonstrates `AMOS`’s applicability in practice.
4. We present an experimental study of the Hessian’s rate of change throughout training.

Our paper is organized as follows: In section 2, we review previous work and describe the notation we will use throughout the paper. In section 3, we develop the `ELMO` algorithm and use it as a baseline for an analysis of arbitrary quasi-Newton optimizer regret. In section 4, we use the bounds we developed to construct a meta-optimizer that selects the appropriate optimizer on-the-fly. Its performance is demonstrated experimentally. Finally, in section 5 we show the value of our novel Lipschitz parameter separation scheme by showing that the Hessian-Lipschitz parameters of certain eigenspaces are much smaller than others, giving optimizers a more accurate loss-function proxy.

2 BACKGROUND

Assumption 1. For a given optimization problem with loss function $f : \mathbb{R}^n \rightarrow \mathbb{R}$, we assume f is twice differentiable.

We note that this assumption is satisfied for all prevalent deep learning optimization problems for all but a zero-measure set of parameters.

2.1 NOTATIONS AND DEFINITIONS

Notation 1. Let $\theta_{t+1}, \theta_t \in \mathbb{R}^n$ the parameter vectors of a pair of consecutive iterations of a given optimization algorithm.

- For brevity of notation, we mark $\Delta\theta_t \triangleq \theta_{t+1} - \theta_t$.
- We mark $\nabla f(\theta_t)$ the gradient of f and $\mathcal{H}(\theta_t)$ the Hessian of f at θ_t .

Notation 2. Let $\theta_t \in \mathbb{R}^n$. We mark $(v_i(\theta_t), \lambda_i(\theta_t))_{i=0}^n$ an orthogonal eigendecomposition of $\mathcal{H}(\theta_t)$ (which exists due to the Hessian symmetry property). For brevity of notation, we will sometimes drop the (θ_t) and just write v_i, λ_i when the meaning is clear.

Since v_i and $-v_i$ are both equally viable eigenvectors, we eliminate ambiguity by assuming

$$\forall_{i \in [n]} : \nabla f(\theta_t)^\top v_i \leq 0 \quad (1)$$

with \top indicating transposition.

Definition 1. We say an algorithm is a *k-order algorithm* if it requires oracle access to the first k derivatives of f .

¹Given an optimization algorithm, there exists an "adversarial" loss function that minimizes the descent obtained by this optimizer, while matching all of the information gathered about the loss-function at a given iteration. We define "regret" as the difference between the actual descent obtained by a given algorithm in this adversarial setting, and the descent obtainable by the most adversary-robust algorithm (meaning that it maximizes loss function descent despite the adversary). An algorithm that computes the minimizer of this regret at each iteration is called "minmax-optimal".

Notation 3. For $\tau \in \mathbb{N}$ we mark $[\tau] = \{t \in \mathbb{N} : t \leq \tau\}$.

Definition 2. Let $U, D \in \mathbb{R}^{n \times n}$ s.t. $D = \text{diag}(d_1, d_2, \dots, d_n)$ is diagonal and U orthogonal, and let $\xi : \mathbb{R} \rightarrow \mathbb{R}$. We mark $\xi(U \cdot D \cdot U^\top) = U \cdot \text{diag}(\xi(d_1), \xi(d_2), \dots, \xi(d_n)) \cdot U^\top$.

Notation 4. Let $A, B \in \mathbb{R}^{n \times n}$. We write $A \succeq 0$ iff A is positive semi-definite, $A \succ 0$ if A is positive definite, $A \succeq B$ if $A - B \succeq 0$ (and likewise for $A \succ B$).

Definition 3. We say that an optimization algorithm is a *quasi-Newton optimization algorithm* if its characteristic update rule may be expressed as:

$$\theta_{t+1} = \theta_t - \alpha_t \Phi_t \nabla f(\theta_t)$$

for $\Phi_t \in \mathbb{R}^{n \times n}$, $\Phi_t \succeq 0$, $\Phi_t^\top = \Phi_t$, $\alpha_t \in \mathbb{R}^+$. We call Φ_t in such algorithms the "preconditioner matrix".

This approach is inspired by Newton's method in convex optimization (see Nocedal & Wright (2006, Chapter 3)) where $\Phi_t = (\mathcal{H}(\theta_t))^{-1}$; see appendix A for a discussion of the challenges and proposed solutions involved in these algorithms. The overwhelming majority of gradient-based optimizers may be expressed as quasi-Newton optimizers (see Martens (2020) for examples), so this paper will concern itself exclusively with this family of optimizers.

Notation 5. Throughout this paper, we will mark the point a convergent quasi-Newton algorithm converges to by θ^* .

2.2 RELATED WORK

The value of the loss function after t iterations is of particular importance to practitioners, due to its implications on the quality of model. One measure of optimizer quality relating to this value is the objective function sub-optimality gap (OFSOG), defined as $f(\theta_T) - f(\theta^*)$. The ARC algorithm is a second-order algorithm that uses a low-rank SVD approximation of the Hessian and estimates a single Hessian-Lipschitz parameter adaptively; Cartis et al. (2012b) prove that OFSOG-optimality (bounding the OFSOG to below ϵ) is achieved by a variant of the ARC algorithm after $\mathcal{O}(\epsilon^{-1})$ iterations in the convex regime, or $\mathcal{O}(\log(\epsilon^{-1}))$ iterations in the strongly convex regime. Garmanjani (2020) show similar bounds for the Nonlinear Stepsize Control algorithm family, and Toint (2013) demonstrate that this is a generalization of ARC and trust-region methods. Liu et al. (2024) prove OFSOG-optimality for the Sophia optimizer (a second-order algorithm that approximates the Hessian as a diagonal matrix, which is estimated with Hutchinson's estimator (Hutchinson, 1989)) after $\mathcal{O}(\epsilon^{-1})$ iterations in the convex regime.

Bottou (2004) partition the optimization process into the initial "search phase", in which the optimizer searches for an approximately convex region in which a local minimum resides, and the later "final phase", in which the optimizer converges to this local minimum.

In the machine learning literature, many common loss functions are "empirical risk functions" - that is, loss functions which can be written as a sum of terms, each of which is a function of only a single sample from the data distribution. When this sum ranges over a very large number of samples, a common approach to estimating it is to perform a Monte Carlo approximation, summing over only a small subset of the terms; this approach is known as the "minibatch approach". Amari (1998) then note that when using this approach, θ_t may be seen as a statistical estimator for θ^* . Working in the "final phase" (and thus assuming convexity), and adopting the estimator approach to θ_t taken by Amari (1998); Bottou & Lecun (2004) give a convergence rate bound for this estimator's variance parameterized by the first- and second-order derivatives at θ^* , assuming only that $\lim_{t \rightarrow \infty} \Phi_t = \mathcal{H}^{-1}(\theta^*)$. Martens (2020) takes these convergence rates and plugs them into a Taylor approximation of $f(\theta_t)$ to obtain the asymptotic OFSOG, given by $f(\theta_T) - f(\theta^*) = \frac{n}{2T} + o(\frac{1}{T})$.

Since the goal of optimization is to minimize a loss function, arguably the best metric for measuring an optimization algorithm's quality are the gains it makes as measured by the loss function values, i.e. its rate of loss function descent. Nevertheless, most algorithms' convergence rate bounds relate to their gradient norms; we note, however, that a bound on an algorithm's gradient norm may be a poor proxy for its descent rate in the early, nonconvex "search" phase, since convergence rate bounds may only imply proximity to a critical point of the gradient, which is neither guaranteed to be the point the algorithm will ultimately converge to nor even to have a small loss function value by any measure.

To the best of our knowledge, our bound is the first to directly address the problem of bounding the loss function value in the "search" phase without assuming convexity (despite neural network loss functions rarely being convex in practice).

We refer the reader to appendix B for discussion on previous attempts at universal convergence rate bounds, other convergence rate measures, and the effect of the preconditioner (and correspondingly, the effect of optimizer choice) on convergence rate.

3 A MINMAX HESSIAN LIPSCHITZ-AWARE OPTIMIZATION ALGORITHM

Optimization algorithm iterations are comprised of two parts: first, we gather information about the loss function and construct a local model of the loss function, and second - we step to the minimum of this proxy. Accordingly, gradient descent and Newton's method use first- and second-order Taylor approximations of f respectively, and while these proxies do give a direction of descent in every subspace of the domain space, they do not indicate optimal step sizes in concave subspaces of the domain space (that is, subspaces in which the loss function is locally concave), since concave first- and second-order polynomials have no minima. To obtain a unique step in all settings (so that our optimizer will be sufficiently general to apply to nonconvex and nonquadratic regions of neural network loss functions), we must therefore model f with a third-order Taylor polynomial.

3.1 GENERAL BOUNDS ON PER-ITERATION DESCENT

A recurring theme in the neural network optimization literature is that the greatest-magnitude eigenvalues of the Hessian are slow to change, as well as their eigenvectors; see, for instance, Sivan et al. (2024); Alain et al. (2018); Sagun et al. (2016); Ghorbani et al. (2019); Gur-Ari et al. (2018); Liu et al. (2024). It is common to formalize this as an assumption (see, e.g., O'Leary-Roseberry et al. (2019); Nesterov & Polyak (2006)) of Hessian-Lipschitz continuity with the matrix spectral norm:

$$\exists_{L_H \in \mathbb{R}} \forall_{\theta, \varphi \in \mathbb{R}^n} : \|\mathcal{H}(\theta) - \mathcal{H}(\varphi)\|_2 \leq L_H \cdot \|\theta - \varphi\|_2 \quad (2)$$

This assumption relies on a single scalar $L_H \in \mathbb{R}$ to describe the the entire Hessian's rate of change. With $\frac{n^2}{2}$ independent entries, however, the Hessian can shift in a far more subtle manner, leading this assumption to be overly conservative, requiring a very large L_H for the assumption to be satisfied, which leads to looseness in convergence rate bounds and suboptimal performance of algorithms that rely on this scalar. We instead make the following finer-grained assumption on the rate of change of the Hessian's eigendecomposition:

Assumption 2. *Hessian Lipschitz-Continuity in each Eigenspace*

For any $\theta, \varphi \in \mathbb{R}^n$, let (eigendecompositions) $\mathcal{H}(\theta) = V \cdot \Lambda \cdot V^\top$, $\mathcal{H}(\varphi) = \tilde{V} \cdot \tilde{\Lambda} \cdot \tilde{V}^\top$ with $V, \tilde{V} \in \mathbb{R}^{n \times n}$ orthogonal matrices and $\Lambda = \text{diag}(\lambda_i)_{i=1}^n$, $\tilde{\Lambda} = \text{diag}(\tilde{\lambda}_i)_{i=1}^n \in \mathbb{R}^{n \times n}$ diagonal matrices, sorted s.t. $\forall_{i \in [n-1]} : \lambda_i \leq \lambda_{i+1}$, $\tilde{\lambda}_i \leq \tilde{\lambda}_{i+1}$. Then the following are satisfied:

$$\begin{aligned} \forall_{\theta \in \mathbb{R}^n} \exists_{(\tilde{L}^i)_{i=1}^n \in (\mathbb{R}^+)^n} \forall_{\varphi \in \mathbb{R}^n} : & \left| \lambda_i - \tilde{\lambda}_i \right| \leq \tilde{L}^i \cdot \left| (\theta - \varphi)^\top v_i \right| \\ \forall_{\theta \in \mathbb{R}^n} \exists_{L_R \in \mathbb{R}^+} \forall_{\varphi \in \mathbb{R}^n} : & \left\| V - \tilde{V} \right\|_2 \leq L_R \cdot \|\theta - \varphi\|_2 \\ \exists_{L_H \in \mathbb{R}} \forall_{\theta \in \mathbb{R}^n} \forall_{i \in [n]} : & \max \{ L_R, \tilde{L}^i \} \leq L_H \wedge \tilde{L}^i \geq L_H^{-1} \end{aligned}$$

When θ is the t -th iterate θ_t of an optimization algorithm, we'll mark the corresponding Lipschitz parameters as L_t^i . We will experimentally demonstrate the significance of this refinement later, by demonstrating that these parameters vary by orders of magnitude.

The above assumption allows us to bound the loss function in each eigenspace of the Hessian; these bounds will then be applicable as tight (since the bounds satisfy assumptions 1 and 2) pessimistic and optimistic models of the loss function in the neighborhood of some iterate θ_t :

Notation 6. Let $\theta_t \in \mathbb{R}^n$, $v_i \in \mathbb{R}^n$ an eigenvector of $\mathcal{H}(\theta_t)$.

$$M_t^i(x) \triangleq \nabla f(\theta_t)^\top v_i \cdot x + \frac{v_i^\top \mathcal{H}(\theta_t) v_i}{2} \cdot x^2 + \frac{L_t^i}{6} \cdot |x|^3 \quad (3)$$

$$m_t^i(x) \triangleq \nabla f(\theta_t)^\top v_i \cdot x + \frac{v_i^\top \mathcal{H}(\theta_t) v_i}{2} \cdot x^2 - \frac{L_t^i}{6} \cdot |x|^3 \quad (4)$$

Lemma 3.1. *Eigenspace Descent Bounds*

Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be a function satisfying assumptions 1 and 2, and let $\theta_{t+1} \in \mathbb{R}^n$. Marking $\Delta\theta_t = \theta_{t+1} - \theta_t$, we have

$$\exists_{(L_t^i)_{i=1}^n \in (\mathbb{R}^+)^n} : f(\theta_{t+1}) - f(\theta_t) \leq \sum_{i=1}^n M_t^i(\Delta\theta_t^\top v_i) \quad (5)$$

$$\exists_{(L_t^i)_{i=1}^n \in (\mathbb{R}^+)^n} : f(\theta_{t+1}) - f(\theta_t) \geq \sum_{i=1}^n m_t^i(\Delta\theta_t^\top v_i) \quad (6)$$

We emphasize that unlike most previous works, our bound does not assume convexity and moreover it demonstrates the strong opportunity for descent in subspaces of the parameter space in which the loss function is concave, giving the first theoretical explanation for the experimental results of Alain et al. (2018) and Ghorbani et al. (2019), to the best of our knowledge (since the negative-definite Hessian corresponding to a locally concave loss function lowers the upper bound M_t^i of the per-iteration descent).

3.2 EXPLOITING THESE BOUNDS FOR A MINMAX ALGORITHM

Here we describe how to exploit the refined approximations that we obtain above to design a minimax optimal algorithm named ELMO which is depicted in Alg. 1.

To gain perspective on the upcoming algorithm as minmax-optimal, we restate a special case of the above lemma as follows: M_t^i is the pointwise maximal function satisfying assumptions 1 and 2:

$$M_t^i(x) = \max_{\substack{\tilde{f}: \mathbb{R} \rightarrow \mathbb{R} \\ \tilde{f}(\theta_t) = f(\theta_t)}} \tilde{f}(\theta_t + x \cdot v_i(\theta_t))$$

In other words, M_t is the loss function (up to domain- and range-space translation) with the least opportunity for descent while still ensuring that all information collected by a second-order optimizer (with Hessian-Lipschitz oracle access) remains unchanged - i.e. *the adversarial loss function detailed in regret's definition*.

Since each element of the sum is a 1-dimensional trinomial, the minmax step is now easily obtained (due to orthogonality of the eigenspaces, resulting from the Hessian's symmetry) by taking the positive root of each term's derivative:

$$\Delta\theta_t^{*\top} v_i \triangleq \arg \min_{\Delta\theta_t} \sum_{i=1}^n M_t^i(\Delta\theta_t^\top v_i) = \frac{\sqrt{\lambda_i^2 + 2L_t^i |\nabla f(\theta_t)^\top v_i|} - \lambda_i}{L_t^i} \quad (7)$$

The above naturally induces a second-order training algorithm that we name Eigenspace-Lipschitz Minmax Optimizer (ELMO); which we describe in Alg. 1. We present a theoretical investigation of ELMO in Appendix C, and we leave its experimental evaluation to future work to maintain our focus on optimizer selection. We mark EIGEN an eigendecomposition subroutine and LIPSCHITZ a Lipschitz parameter oracle.

Importantly, we observe equation 7's equal applicability to convex and concave regions of

the domain space (and ELMO correspondingly) due to the lack of convexity assumptions. Indeed, ELMO makes use of concavity by taking larger steps in concave subspaces than otherwise (since ELMO's second term's step size is proportional to $-\lambda_i$).

Algorithm 1 Algorithm ELMO

Require: $\epsilon \in \mathbb{R}^+$, $\theta_0 \in \mathbb{R}^n$, EIGEN, LIPSCHITZ
 $t \leftarrow 0$

while $f(\theta_t) - f(\theta^*) > \epsilon$ **do**

$(\lambda_i, v_i)_{i=1}^n \leftarrow \text{EIGEN}(\mathcal{H}(\theta_t))$

$(L_t^i)_{i=1}^n \leftarrow \text{LIPSCHITZ}(\theta_t, (v_i)_{i=1}^n)$

$(\Delta\theta_t^i)_{i=1}^n \leftarrow \frac{\sqrt{\lambda_i^2 + 2L_t^i |\nabla f(\theta_t)^\top v_i|} - \lambda_i}{L_t^i}$

$\theta_{t+1} \leftarrow \theta_t + \sum_{i=1}^n \Delta\theta_t^i \cdot v_i$

$t \leftarrow t + 1$

end while

3.3 THE MINMAX PRECONDITIONER

Here we find algorithm ELMO’s characteristic preconditioner. That is, we find Φ_t s.t. ELMO’s step $\Delta\theta_t^*$ can be expressed as in definition 3. We begin by defining a metric of distance between optimization algorithms by the difference between the steps they take, and find the preconditioner matrix whose corresponding quasi-Newton algorithm is equivalent to algorithm ELMO.

Notation 7. For a given algorithm with step $\Delta\theta_t$ at iteration t , mark $\Delta\Delta^i\theta_t = \Delta\theta_t^\top v_i - \Delta\theta_t^{*\top} v_i$ the step’s distance from ELMO’s step. Since $\Delta\Delta^i\theta_t$ is a function of the algorithm chosen (through $\Delta\theta_t$), it is a function of that algorithm’s defining preconditioner: $\Delta\Delta^i\theta_t = \Delta\Delta^i\theta_t(\Phi_t)$

This may be understood intuitively by noting that $\sum_{i=1}^n \Delta\Delta^i\theta_t$ is the L_1 norm of the difference between steps, as measured over the Hessian eigenbasis.

Lemma 3.2. *ELMO’s preconditioner*

Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ satisfying assumptions 1 and 2. The preconditioner of the quasi-Newton algorithm that is equivalent to ELMO (meaning $|\Delta\Delta^i\theta_t| = 0$) is

$$\Phi_t^{\text{ELMO}} = \arg \min_{\Phi_t \in \mathbb{R}^{n \times n}} |\Delta\Delta^i\theta_t(\Phi_t)| = \left(\frac{\mathcal{H}(\theta_t) + \sqrt{(\mathcal{H}(\theta_t))^2 + 2V \cdot \text{diag} \left(L_t^i \cdot \left| \nabla f(\theta_t)^\top v_i \right| \right)_{i=1}^n \cdot V^\top}}{2} \right)^{-1}$$

This preconditioner shows the mechanistic similarity of our algorithm to Newton’s method: while Newton’s method’s preconditioner is simply the inverse Hessian (which may not be positive definite), the matrix whose inverse is our algorithm’s preconditioner is an average between the Hessian and a positive definite, regularized version of the Hessian, whose every eigenvalue is no less than the corresponding Hessian eigenvalue’s magnitude. This ensures positive semi-definiteness of our preconditioner, with regularization dependent on the loss function’s rate of curvature shift.

In fact, Newton’s algorithm (and correspondingly, most existing quasi-Newton algorithms that approximate it) may even lead to a worst-case *decrease* in model quality, even when the associated loss function is convex, for sufficiently great curvature shift (measured by Lipschitz parameter). Plugging

Newton’s step into equation 3 and rearranging tells us that $\forall_{i \in [n].s.t. \lambda_i \geq 0} : M_t^i \left(\frac{|\nabla f(\theta_t)^\top v_i|}{\lambda_i} \right) \geq 0$

for any step t and eigenspace i with

$$L_t^i \geq 3 \frac{\lambda_i^2}{\left| \nabla f(\theta_t)^\top v_i \right|} \quad (8)$$

3.4 DESCENT REGRET OF QUASI-NEWTON OPTIMIZATION ALGORITHMS

Algorithm ELMO is regret-optimal among first- and second-order methods in the sense that its model of the loss function is a generalization of quasi-Newton methods’ and Gradient Descent’s models (Appendix F), so it is natural to use it as a baseline against which to measure other optimizers’ regret.

Different datasets imply different units for the parameters of neural networks, however (e.g. a neural network whose training dataset measures distances in miles will have differently scaled parameters than one whose distances are measured in kilometers, and different optimizer step sizes accordingly), so to maintain consistency across a variety of settings, we turn to dimensional analysis (Lin & Segel, 1988) to eliminate this effect on $\Delta\Delta^i\theta_t$. For simplicity, we restrict our discussion to the descent of the loss function’s restriction to a given eigenspace $\text{span}(v_i)$.

Notation 8. Mark $\Delta\Delta^i\theta'_t = \frac{\Delta\Delta^i\theta_t}{\Delta\theta_t^{*\top} v_i}$ the step’s distance from ELMO’s step relative to ELMO’s step.

In fact, this optimizer metric generalizes several existing metrics, see Appendix D for details.

Theorem 3.3. *Worst-case descent rate for arbitrary optimizers*

Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ a twice-differentiable function satisfying assumptions 1 and 2, and let $\Delta\theta_t$ satisfy $M_t^i(\Delta\theta_t^\top v_i) \leq 0$. Then

$$|M_t^i(\Delta\theta_t^\top v_i)| = \Theta\left(1 + |\Delta\Delta^i\theta'_t|^2\right) \cdot |M_t^i(\Delta\theta_t^{*\top} v_i)|$$

$$|m_t^i(\Delta\theta_t^\top v_i)| = \Theta\left(1 + |\Delta\Delta^i\theta'_t|^p\right) \cdot |m_t^i(\Delta\theta_t^{*\top} v_i)| \quad (9)$$

$$\text{with } p = \begin{cases} 2 & \lambda_i > 0 \wedge \frac{|\nabla f(\theta_t)^\top v_i|}{\lambda_i^2} = 0 \\ 1 & \text{else} \end{cases}$$

In other words, an arbitrary algorithm’s descent regret (the ratio between the algorithm’s descent in the adversarial setting and the optimal descent in the adversarial setting) is proportional to a polynomial of $\Delta\Delta^i\theta'_t$.

4 DYNAMIC OPTIMIZER SELECTION

Here we describe a Meta-Optimizer which we name AMOS (Alg. 2) that dynamically combines different optimizers throughout the training process in order to boost performance. AMOS leans on our ELMO minimax optimal baseline.

Algorithm 2 Algorithm AMOS

Require: $\epsilon_1, \epsilon_2, L_{thresh}, \alpha_L \in \mathbb{R}^+, \theta_0 \in \mathbb{R}^n, I_{hess} \in 2^{\mathbb{N}}, \text{OPT}^1, \text{OPT}^2$

- 1: $t \leftarrow 0$
- 2: **while** $f(\theta_t) - f(\theta_{t+1}) > \epsilon_1$ **do**
- 3: **if** $t \in I_{hess}$ **then**
- 4: $(\lambda_i^t)_{i=1}^n \leftarrow \text{EIGEN}(\mathcal{H}(\theta_t))$
- 5: $(L_t^i)_{i=1}^n \leftarrow \alpha_L \cdot L_{t-1}^i + (1 - \alpha_L) \cdot \frac{|\lambda_i^t - \lambda_i^{t-1}|}{\epsilon_2 + \|\theta_t - \theta_{t-1}\|_2}$
- 6: **for** $i = 1..n$ **do**
- 7: **if** $L_t^i > L_{thresh}$ **then**
- 8: $\text{OPTIM}_i \leftarrow \text{OPT}^1$
- 9: **else**
- 10: $\text{OPTIM}_i \leftarrow \text{OPT}^2$
- 11: **end if**
- 12: **end for**
- 13: **end if**
- 14: $(\Delta\theta_t)_i = \text{OPTIM}_i.\text{step}()$
- 15: $(\theta_{t+1})_i \leftarrow (\theta_t)_i + (\Delta\theta_t)_i$
- 16: $t \leftarrow t + 1$
- 17: **end while**

Concretely, at each Hessian computation, Lipschitz parameters are approximated (line 5) as a running average of the slope of the Hessian eigenvalues, $L_t^i \approx \alpha_L L_{t-1}^i + (1 - \alpha_L) \frac{|\lambda_i(x_t) - \lambda_i(x_{t-1})|}{\|x_t - x_{t-1}\|}$; model parameters² whose Lipschitz parameters are above a given threshold are assigned to the first-order optimizer until the next Hessian computation, while parameters whose Lipschitz parameters are beneath said threshold are assigned to the second-order optimizer (lines 7-10).

²Most popular second-order optimizers assume diagonal Hessian (e.g. Sophia) or block-diagonal Hessian (e.g. K-FAC) structure, the former of which associates each model parameter with its respective eigenvalue, allowing us to associate it with a Lipschitz parameter (and an analogous Lipschitz parameter may be associated with each block of model parameters in the latter case). For the remaining second-order optimizers, we partition model parameters into Hessian *eigenspaces*, and assign each eigenspace to its respective optimizer, which reduces overhead by restricting the second-order optimizer to fewer eigenspaces, similarly to Sivan et al. (2024). See appendix A.

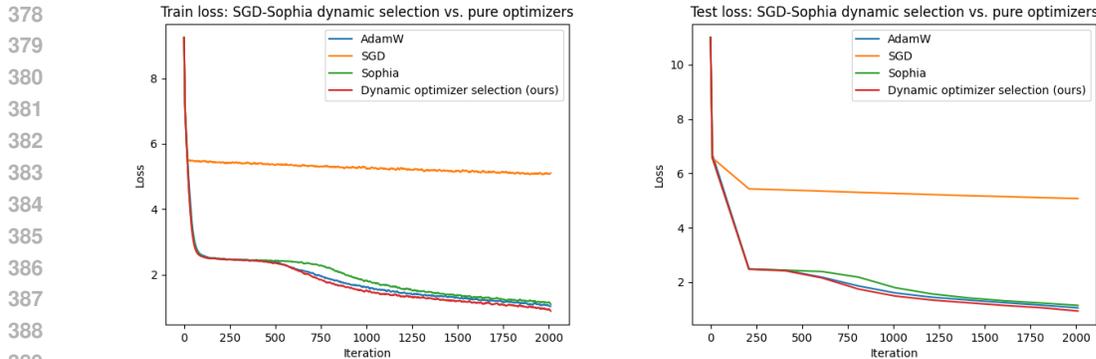


Figure 1: Training trajectory comparison. Our SGD-Sophia hybrid selection strategy outperforms all the rest while offloading responsibility for a significant portion of weights to SGD.

AMOS EXPERIMENTS

We tested AMOS³ in the NanoGPT language modeling framework by Karpathy (2022); further details on settings in appendix K.

- We trained a 12-layer transformer with 12-head multihead attention on the TinyShakespeare dataset (Karpathy, 2015) for 2010 training iterations. This task involves modeling the character-level distribution of 40,000 lines taken from Shakespeare plays.
- Our baseline algorithm is AdamW with the tuned hyperparameters from vanilla NanoGPT.
- We use AMOS to adaptively select between SGD and Sophia throughout training.

Figure 1 presents the training trajectories from our experiments. We see that while pure SGD and Sophia are inferior to AdamW, by using AMOS we are able to combine SGD and Sophia to surpass AdamW on NanoGPT while also offloading 9% of the parameters on average (throughout training) and up to 91% (in a single iteration) to the cheaper first-order optimizer. We ascribe this additional loss minimization to our optimizer avoiding destructive Newton steps (as discussed in section 3.3).

5 LIPSCHITZ DISTRIBUTION

Previous works using the Hessian Lipschitz continuity assumption (e.g. ARC (Nesterov & Polyak, 2006) and its variants, O’Leary-Roseberry et al. (2019)) assume a single Lipschitz parameter for all eigenspaces. Although a finite number n of eigenspaces ensures that such a Lipschitz parameter exists (the maximal Lipschitz parameter), they fail to account for the distribution of these Lipschitz parameters over the eigenspaces. We claim that these parameters vary widely both over the eigenspaces and over the course of training, so that a single constant value fails to capture this structure; in this section we provide evidence for this claim. A discussion of sources of interest in this distribution is deferred to Appendix G.

The first source of evidence for our claim is from existing literature on the subject; we defer a discussion of this to appendix I. To test our claim directly, we modify an ARC implementation (Simpson & Wang, 2023) to compute the steps called for by ELM0 at each point reached by a quasi-Newton algorithm, restricted to the subspace spanned by the eigenvectors corresponding to the single most positive and single most negative eigenvalues of each Hessian, and to use distinct Lipschitz parameters for each. Due to the computational difficulty of computing Lipschitz parameters precisely, we use these Lipschitz parameter values as an estimate for L^+ , L^- . We note the crudeness of these adaptive measurements, merely adapting to keep $\frac{f(\theta_t) - f(\theta_{t+1})}{|\sum_{i=1}^n M_i^t|}$ within a given range with a restriction to powers of 2; nevertheless, the point is made.

A detailing of our experiment settings is given in appendix J as well as the full set of our experiment results, however we present two experiments in figure 2 for completeness. Our experiments show that

³AMOS: Adaptive Meta-algorithm for Optimizer Selection

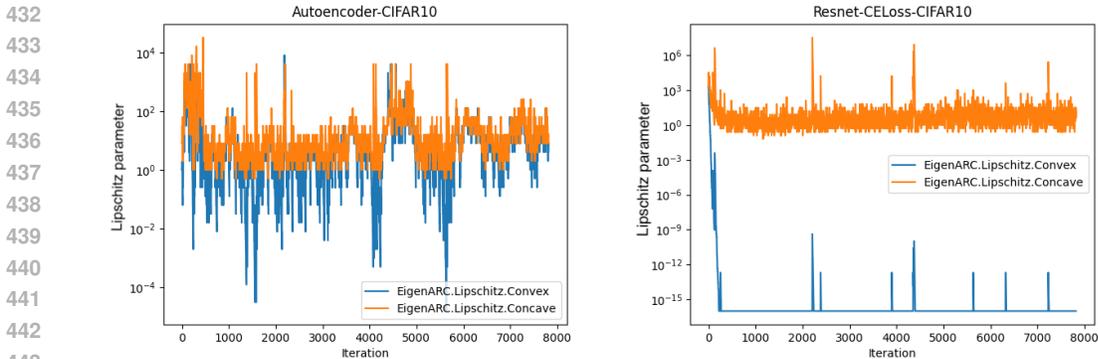


Figure 2: Comparisons of convex-subspace Lipschitz parameters to concave-subspace Lipschitz parameters. *Logarithmic scale*

as expected, $L^+ \ll L^-$, and the gap widens exponentially as training progresses in all cases except the autoencoders. Since we saw in section 3.3 that small convex Lipschitz parameters imply effective second-order optimization, this justifies common practice as noted by, e.g. O’Leary-Roseberry et al. (2019), of requiring the preconditioner to be an increasingly better approximation of the inverse Hessian (by increasing the strictness of the inverse Hessian approximation algorithm’s stopping condition) as training progresses. Our findings also support Gur-Ari et al. (2018)’s observation that neural network optimization succeeds despite its nonconvexity due to the gradient residing primarily in convex subspaces, by demonstrating the stability of these convex subspaces.

Several factors seem to impact the size (by orders of magnitude) of the convex Lipschitz parameters, and they appear to be correlated with an intuitive sense of the difficulty of the setting being trained.

- The convex Lipschitz parameters are many orders of magnitude greater in the autoencoder task than in the classification task. We ascribe this gap to the more difficult task of learning a generative representation of the data instead of merely a discriminative representation of it (see Ng (2012, Chapter 4) for a discussion on generative vs. discriminative models).
- The convex Lipschitz parameters are reduced approximately 100x in the image classification task by adding residual connections. It is well known that residual connections reduce training difficulty (Li et al., 2018).
- The convex Lipschitz parameters are approximately 100x smaller when training ResNet to perform classification of natural images instead of Gaussian noise with random labels. We ascribe this to greater difficulty involved in discriminating noise, which requires partial memorization of the training set.

The Lipschitz parameters’ indication of Newton’s applicability suggests that the above observations may be applicable as a rule of thumb for optimizer selection; this is validated in experiments in Appendix E.

6 CONCLUSION

In this work we developed and analyzed a Hessian eigenspace Lipschitz-aware minmax optimization algorithm ELM0 by taking an eigendecomposition-centric approach to locally modelling a loss function. We then proved a widely applicable worst-case relative descent rate bound for quasi-Newton optimizers by comparison to ELM0. We experimented with the Lipschitz distributions, discovering that they are correlated with task difficulty and that they are helpful for optimizer and optimization hyperparameters selection — specifically, integrating second-order information into optimizers at the cost of additional computational complexity is worthwhile in settings where the convex Lipschitz parameters are small, but not those where they are large. Finally, we used all of the above to design algorithm AMOS for dynamic optimizer selection, and demonstrated its performance in experiments.

REFERENCES

- 486
487
488 Naman Agarwal, Brian Bullins, Xinyi Chen, Elad Hazan, Karan Singh, Cyril Zhang, and Yi Zhang.
489 Efficient full-matrix adaptive regularization. In Kamalika Chaudhuri and Ruslan Salakhutdinov
490 (eds.), *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of
491 *Proceedings of Machine Learning Research*, pp. 102–110. PMLR, 09–15 Jun 2019. URL <https://proceedings.mlr.press/v97/agarwal19b.html>.
492
- 493 Naman Agarwal, Rohan Anil, Elad Hazan, Tomer Koren, and Cyril Zhang. Learning rate grafting:
494 Transferability of optimizer tuning, 2022. URL https://openreview.net/forum?id=FpKgG31Z_i9.
495
- 496 Guillaume Alain, Nicolas Le Roux, and Pierre-Antoine Manzagol. Negative eigenvalues of the
497 hessian in deep neural networks. In *6th International Conference on Learning Representations,*
498 *ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Workshop Track Proceedings.*
499 OpenReview.net, 2018. URL <https://openreview.net/forum?id=SliiddyDG>.
500
- 501 Zeyuan Allen-Zhu. Katyusha: the first direct acceleration of stochastic gradient methods. *Proceedings*
502 *of the 49th Annual ACM SIGACT Symposium on Theory of Computing*, 2016. URL <https://api.semanticscholar.org/CorpusID:51918825>.
503
- 504 Shun-ichi Amari. Natural Gradient Works Efficiently in Learning. *Neural Computation*, 10(2):
505 251–276, 02 1998. ISSN 0899-7667. doi: 10.1162/089976698300017746. URL <https://doi.org/10.1162/089976698300017746>.
506
507
- 508 Rohan Anil, Vineet Gupta, Tomer Koren, Kevin Regan, and Yoram Singer. Second order optimization
509 made practical. *CoRR*, abs/2002.09018, 2020. URL <https://arxiv.org/abs/2002.09018>.
510
- 511 Somenath Bera and Vimal Shrivastava. Analysis of various optimizers on deep convolutional
512 neural network model in the application of hyperspectral remote sensing image classification.
513 *International Journal of Remote Sensing*, 41:2664–2683, 04 2020. doi: 10.1080/01431161.2019.
514 1694725.
515
- 516 Dimitri P. Bertsekas. Incremental least squares methods and the extended kalman filter. *SIAM J.*
517 *Optim.*, 6(3):807–822, 1996. doi: 10.1137/S1052623494268522. URL <https://doi.org/10.1137/S1052623494268522>.
518
- 519 Aleksandar Botev, Hippolyt Ritter, and David Barber. Practical Gauss-Newton optimisation for
520 deep learning. In Doina Precup and Yee Whye Teh (eds.), *Proceedings of the 34th International*
521 *Conference on Machine Learning*, volume 70 of *Proceedings of Machine Learning Research*,
522 pp. 557–565. PMLR, 06–11 Aug 2017. URL <https://proceedings.mlr.press/v70/botev17a.html>.
523
524
- 525 Léon Bottou. Stochastic learning. In Olivier Bousquet, Ulrike von Luxburg, and Gunnar Rätsch
526 (eds.), *Advanced Lectures on Machine Learning: ML Summer Schools 2003, Canberra, Australia,*
527 *February 2 - 14, 2003, Tübingen, Germany, August 4 - 16, 2003, Revised Lectures*, pp. 146–168.
528 Springer Berlin Heidelberg, Berlin, Heidelberg, 2004. ISBN 978-3-540-28650-9. doi: 10.1007/
529 978-3-540-28650-9_7. URL https://doi.org/10.1007/978-3-540-28650-9_7.
- 530 Léon Bottou and Yann Lecun. On-line learning for very large datasets. *J. Applied Stochastic Models*
531 *in Business and Industry*, 01 2004.
532
- 533 C. Cartis, N.I.M. Gould, and Ph.L. Toint. Complexity bounds for second-order optimality in
534 unconstrained optimization. *Journal of Complexity*, 28(1):93–108, 2012a. ISSN 0885-064X.
535 doi: <https://doi.org/10.1016/j.jco.2011.06.001>. URL <https://www.sciencedirect.com/science/article/pii/S0885064X11000537>.
536
- 537 Coralía Cartis, Nicholas I. M. Gould, and Philippe L. Toint. Adaptive cubic regularisation methods for
538 unconstrained optimization. part i: motivation, convergence and numerical results. *Mathematical*
539 *Programming*, 127(2):245–295, 2011a. doi: 10.1007/s10107-009-0286-5. URL <https://doi.org/10.1007/s10107-009-0286-5>.

- 540 Coralia Cartis, Nicholas I. M. Gould, and Philippe L. Toint. Adaptive cubic regularisation methods
541 for unconstrained optimization. part ii: worst-case function- and derivative-evaluation complexity.
542 *Mathematical Programming*, 130(2):295–319, 2011b. doi: 10.1007/s10107-009-0337-y. URL
543 <https://doi.org/10.1007/s10107-009-0337-y>.
- 544 Coralia Cartis, Nicholas I. M. Gould, and Philippe L. Toint. Evaluation complexity of adaptive cubic
545 regularization methods for convex unconstrained optimization. *Optimization Methods and Soft-*
546 *ware*, 27:197 – 219, 2012b. URL [https://api.semanticscholar.org/CorpusID:](https://api.semanticscholar.org/CorpusID:16647191)
547 16647191.
- 549 Andrew R. Conn, Nicholas I. M. Gould, and Philippe L. Toint. *Trust-Region Methods*. SIAM,
550 Philadelphia, PA, USA, 2000.
- 551 Yann N. Dauphin, Razvan Pascanu, Caglar Gulcehre, Kyunghyun Cho, Surya Ganguli, and Yoshua
552 Bengio. Identifying and attacking the saddle point problem in high-dimensional non-convex opti-
553 mization. In *Proceedings of the 27th International Conference on Neural Information Processing*
554 *Systems - Volume 2*, NIPS’14, pp. 2933–2941, Cambridge, MA, USA, 2014. MIT Press.
- 556 Alexandre D’efossez, Léon Bottou, Francis R. Bach, and Nicolas Usunier. A simple convergence
557 proof of adam and adagrad. *Trans. Mach. Learn. Res.*, 2022, 2020. URL [https://api.](https://api.semanticscholar.org/CorpusID:225213299)
558 [semanticscholar.org/CorpusID:225213299](https://api.semanticscholar.org/CorpusID:225213299).
- 559 Ron S. Dembo, Stanley C. Eisenstat, and Trond Steihaug. Inexact newton methods. *SIAM Journal on*
560 *Numerical Analysis*, 19(2):400–408, 1982. doi: 10.1137/0719025. URL [https://doi.org/](https://doi.org/10.1137/0719025)
561 [10.1137/0719025](https://doi.org/10.1137/0719025).
- 562 J.E. Dennis and R.B. Schnabel. *Numerical Methods for Unconstrained Optimization and Non-*
563 *linear Equations*. Prentice-Hall Civil Engineering and Engineering Mechanics Se. Prentice-
564 Hall, 1983. ISBN 9780136272168. URL [https://books.google.co.il/books?id=](https://books.google.co.il/books?id=4fFQAAAAMAAJ)
565 [4fFQAAAAMAAJ](https://books.google.co.il/books?id=4fFQAAAAMAAJ).
- 567 Tian Ding, Dawei Li, and Ruoyu Sun. Sub-optimal local minima exist for almost all over-
568 parameterized neural networks. *ArXiv*, abs/1911.01413, 2019. URL [https://api.](https://api.semanticscholar.org/CorpusID:207870322)
569 [semanticscholar.org/CorpusID:207870322](https://api.semanticscholar.org/CorpusID:207870322).
- 571 John Duchi, Elad Hazan, and Yoram Singer. Adaptive subgradient methods for online learning and
572 stochastic optimization. *Journal of Machine Learning Research*, 12(61):2121–2159, 2011a. URL
573 <http://jmlr.org/papers/v12/duchi11a.html>.
- 574 John C. Duchi, Elad Hazan, and Yoram Singer. Adaptive subgradient methods for online learning
575 and stochastic optimization. *J. Mach. Learn. Res.*, 12:2121–2159, 2011b. URL [https://api.](https://api.semanticscholar.org/CorpusID:538820)
576 [semanticscholar.org/CorpusID:538820](https://api.semanticscholar.org/CorpusID:538820).
- 577 S C Eisenstat and H F Walker. Choosing the forcing terms in an inexact newton method. *SIAM*
578 *Journal on Scientific Computing*, 17(1), 1 1996. doi: 10.1137/0917003. URL [https://www.](https://www.osti.gov/biblio/218521)
579 [osti.gov/biblio/218521](https://www.osti.gov/biblio/218521).
- 581 William R. Esposito and Christodoulos A. Floudas. Gauss–newton method: Least squares, relation to
582 newton’s method. In Christodoulos A. Floudas and Panos M. Pardalos (eds.), *Encyclopedia of Optimization*, pp. 733–738. Springer
583 US, Boston, MA, 2001. ISBN 978-0-306-48332-5. doi: 10.1007/0-306-48332-7_151. URL
584 https://doi.org/10.1007/0-306-48332-7_151.
- 586 Reuben Feinman. Pytorch-minimize: a library for numerical optimization with autograd, 2021. URL
587 <https://github.com/rfeinman/pytorch-minimize>.
- 588 Dan Garber, Elad Hazan, Chi Jin, Kakade Sham, Cameron Musco, Praneeth Netrapalli, and Aaron
589 Sidford. Faster eigenvector computation via shift-and-invert preconditioning. In Maria Florina
590 Balcan and Kilian Q. Weinberger (eds.), *Proceedings of The 33rd International Conference on*
591 *Machine Learning*, volume 48 of *Proceedings of Machine Learning Research*, pp. 2626–2634,
592 New York, New York, USA, 20–22 Jun 2016. PMLR. URL [https://proceedings.mlr.](https://proceedings.mlr.press/v48/garber16.html)
593 [press/v48/garber16.html](https://proceedings.mlr.press/v48/garber16.html).

- 594 R. Garmanjani. A note on the worst-case complexity of nonlinear stepsize control methods for convex
595 smooth unconstrained optimization. *Optimization*, 71:1–11, 10 2020. doi: 10.1080/02331934.
596 2020.1830088.
- 597
- 598 Rong Ge, Furong Huang, Chi Jin, and Yang Yuan. Escaping from saddle points — online stochastic
599 gradient for tensor decomposition. In Peter Grünwald, Elad Hazan, and Satyen Kale (eds.),
600 *Proceedings of The 28th Conference on Learning Theory*, volume 40 of *Proceedings of Machine
601 Learning Research*, pp. 797–842, Paris, France, 03–06 Jul 2015. PMLR. URL [https://
602 proceedings.mlr.press/v40/Ge15.html](https://proceedings.mlr.press/v40/Ge15.html).
- 603 Behrooz Ghorbani, Shankar Krishnan, and Ying Xiao. An investigation into neural net opti-
604 mization via hessian eigenvalue density. In Kamalika Chaudhuri and Ruslan Salakhutdinov
605 (eds.), *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of
606 *Proceedings of Machine Learning Research*, pp. 2232–2241. PMLR, 09–15 Jun 2019. URL
607 <https://proceedings.mlr.press/v97/ghorbani19b.html>.
- 608
- 609 Donald Goldfarb, Yi Ren, and Achraf Bahamou. Practical quasi-newton methods for training deep
610 neural networks. In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin (eds.),
611 *Advances in Neural Information Processing Systems*, volume 33, pp. 2386–2396. Curran Asso-
612 ciates, Inc., 2020. URL [https://proceedings.neurips.cc/paper_files/paper/
613 2020/file/192fc044e74df9ac5dc9f3395-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2020/file/192fc044e74df9ac5dc9f3395-Paper.pdf).
- 614 I. J. Good. Rational decisions. *Journal of the Royal Statistical Society. Series B (Methodological)*, 14
615 (1):107–114, 1952. ISSN 00359246. URL <http://www.jstor.org/stable/2984087>.
- 616
- 617 Andreas Griewank. The modification of newton’s method for unconstrained optimization by bounding
618 cubic terms. Technical report, Department of Applied Mathematics and Theoretical Physics,
619 University of Cambridge, United Kingdom, 1981.
- 620 Vineet Gupta, Tomer Koren, and Yoram Singer. Shampoo: Preconditioned stochastic tensor op-
621 timization. In Jennifer Dy and Andreas Krause (eds.), *Proceedings of the 35th International
622 Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*,
623 pp. 1842–1850. PMLR, 10–15 Jul 2018. URL [https://proceedings.mlr.press/v80/
624 gupta18a.html](https://proceedings.mlr.press/v80/gupta18a.html).
- 625
- 626 Guy Gur-Ari, Daniel A. Roberts, and Ethan Dyer. Gradient descent happens in a tiny subspace.
627 *ArXiv*, abs/1812.04754, 2018. URL [https://api.semanticscholar.org/CorpusID:
628 54480858](https://api.semanticscholar.org/CorpusID:54480858).
- 629 Kaiming He, X. Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition.
630 *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 770–778, 2015.
631 URL <https://api.semanticscholar.org/CorpusID:206594692>.
- 632
- 633 M.F. Hutchinson. A stochastic estimator of the trace of the influence matrix for laplacian smoothing
634 splines. *Communication in Statistics- Simulation and Computation*, 18:1059–1076, 01 1989. doi:
635 10.1080/03610919008812866.
- 636 Maarten Jansen and Gerda Claeskens. Cramér–rao inequality. In Miodrag Lovric (ed.), *International
637 Encyclopedia of Statistical Science*, pp. 322–323. Springer Berlin Heidelberg, Berlin, Heidelberg,
638 2011. ISBN 978-3-642-04898-2. doi: 10.1007/978-3-642-04898-2_197. URL [https://doi.
639 org/10.1007/978-3-642-04898-2_197](https://doi.org/10.1007/978-3-642-04898-2_197).
- 640
- 641 Chi Jin, Rong Ge, Praneeth Netrapalli, Sham M. Kakade, and Michael I. Jordan. How to escape
642 saddle points efficiently. In Doina Precup and Yee Whye Teh (eds.), *Proceedings of the 34th
643 International Conference on Machine Learning*, volume 70 of *Proceedings of Machine Learning
644 Research*, pp. 1724–1732. PMLR, 06–11 Aug 2017. URL [https://proceedings.mlr.
645 press/v70/jin17a.html](https://proceedings.mlr.press/v70/jin17a.html).
- 646 Andrej Karpathy. char-rnn. <https://github.com/karpathy/char-rnn>, 2015.
- 647 Andrej Karpathy. NanoGPT. <https://github.com/karpathy/nanoGPT>, 2022.

- 648 Kenji Kawaguchi. Deep learning without poor local minima. In D. Lee, M. Sugiyama, U. Luxburg,
649 I. Guyon, and R. Garnett (eds.), *Advances in Neural Information Processing Systems*, volume 29.
650 Curran Associates, Inc., 2016. URL [https://proceedings.neurips.cc/paper_
651 files/paper/2016/file/f2fc990265c712c49d51a18a32b39f0c-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2016/file/f2fc990265c712c49d51a18a32b39f0c-Paper.pdf).
652
- 653 Kenji Kawaguchi and Yoshua Bengio. Depth with nonlinearity creates no bad local minima in
654 resnets. *Neural Networks*, 118:167–174, 2019. ISSN 0893-6080. doi: [https://doi.org/10.1016/
655 j.neunet.2019.06.009](https://doi.org/10.1016/j.neunet.2019.06.009). URL [https://www.sciencedirect.com/science/article/
656 pii/S0893608019301820](https://www.sciencedirect.com/science/article/pii/S0893608019301820).
- 657 Diederik Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint
658 arXiv:1412.6980*, 2014.
659
- 660 Yann LeCun. *PhD thesis: Modeles connexionnistes de l'apprentissage (connectionist learning
661 models)*. PhD thesis, Université Pierre et Marie Curie, Paris, France, 1987. URL [https:
662 //api.semanticscholar.org/CorpusID:151887454](https://api.semanticscholar.org/CorpusID:151887454).
- 663 Kfir Yehuda Levy. The power of normalization: Faster evasion of saddle points. *ArXiv*,
664 abs/1611.04831, 2016. URL [https://api.semanticscholar.org/CorpusID:
665 16706102](https://api.semanticscholar.org/CorpusID:16706102).
- 666 Hao Li, Zheng Xu, Gavin Taylor, Christoph Studer, and Tom Goldstein. Visualizing the loss land-
667 scape of neural nets. In S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and
668 R. Garnett (eds.), *Advances in Neural Information Processing Systems*, volume 31. Curran Asso-
669 ciates, Inc., 2018. URL [https://proceedings.neurips.cc/paper_files/paper/
670 2018/file/a41b3bb3e6b050b6c9067c67f663b915-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2018/file/a41b3bb3e6b050b6c9067c67f663b915-Paper.pdf).
671
- 672 Lin and Segel. *Mathematics Applied to Deterministic Problems in the Natural Sciences*. SIAM,
673 Philadelphia, 1988. ISBN 978-0-89871-229-2. doi: 10.1137/1.9781611971347. URL [https:
674 //doi.org/10.1137/1.9781611971347](https://doi.org/10.1137/1.9781611971347).
- 675 Hong Liu, Zhiyuan Li, David Leo Wright Hall, Percy Liang, and Tengyu Ma. Sophia: A scalable
676 stochastic second-order optimizer for language model pre-training. In *The Twelfth International
677 Conference on Learning Representations*, 2024. URL [https://openreview.net/forum?
678 id=3xHDeA8Noi](https://openreview.net/forum?id=3xHDeA8Noi).
- 679 Kenneth L. Manders and Leonard Adleman. Np-complete decision problems for binary quadrat-
680 ics. *Journal of Computer and System Sciences*, 16(2):168–184, 1978. ISSN 0022-0000.
681 doi: [https://doi.org/10.1016/0022-0000\(78\)90044-2](https://doi.org/10.1016/0022-0000(78)90044-2). URL [https://www.sciencedirect.
682 com/science/article/pii/0022000078900442](https://www.sciencedirect.com/science/article/pii/0022000078900442).
683
- 684 James Martens. Deep learning via hessian-free optimization. In Johannes Fürnkranz and Thorsten
685 Joachims (eds.), *ICML*, pp. 735–742. Omnipress, 2010. URL [http://dblp.uni-trier.
686 de/db/conf/icml/icml2010.html#Martens10](http://dblp.uni-trier.de/db/conf/icml/icml2010.html#Martens10).
- 687 James Martens. New insights and perspectives on the natural gradient method. *J. Mach. Learn. Res.*,
688 21(1), jan 2020. ISSN 1532-4435.
689
- 690 James Martens and Roger Grosse. Optimizing neural networks with kronecker-factored approximate
691 curvature. In Francis Bach and David Blei (eds.), *Proceedings of the 32nd International Conference
692 on Machine Learning*, volume 37 of *Proceedings of Machine Learning Research*, pp. 2408–2417,
693 Lille, France, 07–09 Jul 2015. PMLR. URL [https://proceedings.mlr.press/v37/
694 martens15.html](https://proceedings.mlr.press/v37/martens15.html).
- 695 Harshal Mittal, Kartikey Pandey, and Yash Kant. Iclr reproducibility challenge report (padam :
696 Closing the generalization gap of adaptive gradient methods in training deep neural networks).
697 *ArXiv*, abs/1901.09517, 2019. URL [https://api.semanticscholar.org/CorpusID:
698 249647677](https://api.semanticscholar.org/CorpusID:249647677).
699
- 700 Aatila Mustapha, Lachgar Mohamed, and Kartit Ali. Comparative study of optimization techniques
701 in deep learning: Application in the ophthalmology field. *Journal of Physics: Conference Series*,
1743, 2021. URL <https://api.semanticscholar.org/CorpusID:234179873>.

- 702 Yuri Nesterov and B. T. Polyak. Cubic regularization of newton method and its global performance.
703 *Mathematical Programming*, 108(1):177–205, 2006. doi: 10.1007/s10107-006-0706-8. URL
704 <https://doi.org/10.1007/s10107-006-0706-8>.
- 705 Andrew Ng. Cs229 lecture notes - supervised learning. Available at
706 https://cs229.stanford.edu/lectures-spring2022/main_notes.pdf, 2012.
- 707
- 708 Quynh Nguyen, Mahesh Chandra Mukkamala, and Matthias Hein. On the loss landscape of a class
709 of deep neural networks with no bad local valleys. In *International Conference on Learning*
710 *Representations*, 2019. URL <https://openreview.net/forum?id=HJgXsjA5tQ>.
- 711
- 712 Jorge Nocedal and Stephen J. Wright. *Numerical Optimization*. Springer, New York, NY, USA, 2e
713 edition, 2006.
- 714 Thomas O’Leary-Roseberry, Nick Alger, and Omar Ghattas. Inexact newton methods for stochastic
715 non-convex optimization with applications to neural network training. *arXiv: Optimization and*
716 *Control*, 2019. URL <https://api.semanticscholar.org/CorpusID:155100112>.
- 717
- 718 P.J. Olver and C. Shakiban. *Applied Linear Algebra*. Prentice Hall, 2006. ISBN 9780131473829.
719 URL <https://books.google.co.il/books?id=D2tyQgAACAAJ>.
- 720 OpenAI, Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni
721 Aleman, Diogo Almeida, Janko Altschmidt, Sam Altman, Shyamal Anadkat, Red Avila, Igor
722 Babuschkin, Suchir Balaji, Valerie Balcom, Paul Baltescu, Haiming Bao, Mohammad Bavarian,
723 Jeff Belgum, Irwan Bello, Jake Berdine, Gabriel Bernadett-Shapiro, Christopher Berner, Lenny
724 Bogdonoff, Oleg Boiko, Madelaine Boyd, Anna-Luisa Brakman, Greg Brockman, Tim Brooks,
725 Miles Brundage, Kevin Button, Trevor Cai, Rosie Campbell, Andrew Cann, Brittany Carey, Chelsea
726 Carlson, Rory Carmichael, Brooke Chan, Che Chang, Fotis Chantzis, Derek Chen, Sully Chen,
727 Ruby Chen, Jason Chen, Mark Chen, Ben Chess, Chester Cho, Casey Chu, Hyung Won Chung,
728 Dave Cummings, Jeremiah Currier, Yunxing Dai, Cory Decareaux, Thomas Degry, Noah Deutsch,
729 Damien Deville, Arka Dhar, David Dohan, Steve Dowling, Sheila Dunning, Adrien Ecoffet, Atty
730 Eleti, Tyna Eloundou, David Farhi, Liam Fedus, Niko Felix, Simón Posada Fishman, Juston Forte,
731 Isabella Fulford, Leo Gao, Elie Georges, Christian Gibson, Vik Goel, Tarun Gogineni, Gabriel
732 Goh, Rapha Gontijo-Lopes, Jonathan Gordon, Morgan Grafstein, Scott Gray, Ryan Greene, Joshua
733 Gross, Shixiang Shane Gu, Yufe i Guo, Chris Hallacy, Jesse Han, Jeff Harris, Yuchen He, Mike
734 Heaton, Johannes Heidecke, Chris Hesse, Alan Hickey, Wade Hickey, Peter Hoeschele, Brandon
735 Houghton, Kenny Hsu, Shengli Hu, Xin Hu, Joost Huizinga, Shantanu Jain, Shawn Jain, Joanne
736 Jang, Angela Jiang, Roger Jiang, Haozhun Jin, Denny Jin, Shino Jomoto, Billie Jonn, Heewoo
737 Jun, Tomer Kaftan, Łukasz Kaiser, Ali Kamali, Ingmar Kanitscheider, Nitish Shirish Keskar,
738 Tabarak Khan, Logan Kilpatrick, Jong Wook Kim, Christina Kim, Yongjik Kim, Jan Hendrik
739 Kirchner, Jamie Kiros, Matt Knight, Daniel Kokotajlo, Łukasz Kondraciuk, Andrew Kondrich,
740 Aris Konstantinidis, Kyle Kosic, Gretchen Krueger, Vishal Kuo, Michael Lampe, Ikai Lan, Teddy
741 Lee, Jan Leike, Jade Leung, Daniel Levy, Chak Ming Li, Rachel Lim, Molly Lin, Stephanie
742 Lin, Mateusz Litwin, Theresa Lopez, Ryan Lowe, Patricia Lue, Anna Makanju, Kim Malfacini,
743 Sam Manning, Todor Markov, Yaniv Markovski, Bianca Martin, Katie Mayer, Andrew Mayne,
744 Bob McGrew, Scott Mayer McKinney, Christine McLeavey, Paul McMillan, Jake McNeil, David
745 Medina, Aalok Mehta, Jacob Menick, Luke Metz, Andrey Mishchenko, Pamela Mishkin, Vinnie
746 Monaco, Evan Morikawa, Daniel Mossing, Tong Mu, Mira Murati, Oleg Murk, David Mély,
747 Ashvin Nair, Reiichiro Nakano, Rajeev Nayak, Arvind Neelakantan, Richard Ngo, Hyeonwoo
748 Noh, Long Ouyang, Cullen O’Keefe, Jakub Pachocki, Alex Paino, Joe Palermo, Ashley Pantuliano,
749 Giambattista Parascandolo, Joel Parish, Emy Parparita, Alex Passos, Mikhail Pavlov, Andrew Peng,
750 Adam Perelman, Filipe de Avila Belbute Peres, Michael Petrov, Henrique Ponde de Oliveira Pinto,
751 Michael, Pokorný, Michelle Pokrass, Vitchyr H. Pong, Tolly Powell, Alethea Power, Boris Power,
752 Elizabeth Proehl, Raul Puri, Alec Radford, Jack Rae, Aditya Ramesh, Cameron Raymond, Francis
753 Real, Kendra Rimbach, Carl Ross, Bob Rotsted, Henri Roussez, Nick Ryder, Mario Saltarelli, Ted
754 Sanders, Shibani Santurkar, Girish Sastry, Heather Schmidt, David Schnurr, John Schulman, Daniel
755 Selsam, Kyla Sheppard, Toki Sherbakov, Jessica Shieh, Sarah Shoker, Pranav Shyam, Szymon
Sidor, Eric Sigler, Maddie Simens, Jordan Sitkin, Katarina Slama, Ian Sohl, Benjamin Sokolowsky,
Yang Song, Natalie Staudacher, Felipe Petroski Such, Natalie Summers, Ilya Sutskever, Jie Tang,
Nikolas Tezak, Madeleine B. Thompson, Phil Tillet, Amin Tootoonchian, Elizabeth Tseng, Preston
Tuggle, Nick Turley, Jerry Tworek, Juan Felipe Cerón Uribe, Andrea Vallone, Arun Vijayvergiya,

- 756 Chelsea Voss, Carroll Wainwright, Justin Jay Wang, Alvin Wang, Ben Wang, Jonathan Ward, Jason
757 Wei, CJ Weinmann, Akila Welihinda, Peter Welinder, Jiayi Weng, Lilian Weng, Matt Wiethoff,
758 Dave Willner, Clemens Winter, Samuel Wolrich, Hannah Wong, Lauren Workman, Sherwin Wu,
759 Jeff Wu, Michael Wu, Kai Xiao, Tao Xu, Sarah Yoo, Kevin Yu, Qiming Yuan, Wojciech Zaremba,
760 Rowan Zellers, Chong Zhang, Marvin Zhang, Shengjia Zhao, Tianhao Zheng, Juntang Zhuang,
761 William Zhuk, and Barret Zoph. Gpt-4 technical report, 2024.
- 762 Panos M. Pardalos and Stephen A. Vavasis. Quadratic programming with one negative eigenvalue
763 is np-hard. *Journal of Global Optimization*, 1(1):15–22, 1991. doi: 10.1007/BF00120662. URL
764 <https://doi.org/10.1007/BF00120662>.
765
- 766 Dylan Patel and Gerald Wong. Gpt-4 architecture, infrastructure, training dataset, costs, vision, moe.
767 <https://www.semianalysis.com/p/gpt-4-architecture-infrastructure>,
768 2023. Accessed: 2024-05-01.
- 769 Barak A. Pearlmutter. Fast Exact Multiplication by the Hessian. *Neural Computation*, 6(1):147–160,
770 01 1994. ISSN 0899-7667. doi: 10.1162/neco.1994.6.1.147. URL [https://doi.org/10.](https://doi.org/10.1162/neco.1994.6.1.147)
771 [1162/neco.1994.6.1.147](https://doi.org/10.1162/neco.1994.6.1.147).
- 772 Karl Pearson. Note on Regression and Inheritance in the Case of Two Parents. *Proceedings of the*
773 *Royal Society of London Series I*, 58:240–242, January 1895.
- 774 Sebastian Ruder. An overview of gradient descent optimization algorithms. *ArXiv*, abs/1609.04747,
775 2016. URL <https://api.semanticscholar.org/CorpusID:17485266>.
776
- 777 Levent Sagun, Léon Bottou, and Yann LeCun. Eigenvalues of the hessian in deep learning: Singu-
778 larity and beyond. *arXiv: Learning*, 2016. URL [https://api.semanticscholar.org/](https://api.semanticscholar.org/CorpusID:35723845)
779 [CorpusID:35723845](https://api.semanticscholar.org/CorpusID:35723845).
- 780 Robin M Schmidt, Frank Schneider, and Philipp Hennig. Descending through a crowded valley -
781 benchmarking deep learning optimizers. In Marina Meila and Tong Zhang (eds.), *Proceedings of*
782 *the 38th International Conference on Machine Learning*, volume 139 of *Proceedings of Machine*
783 *Learning Research*, pp. 9367–9376. PMLR, 18–24 Jul 2021. URL [https://proceedings.](https://proceedings.mlr.press/v139/schmidt21a.html)
784 [mlr.press/v139/schmidt21a.html](https://proceedings.mlr.press/v139/schmidt21a.html).
785
- 786 Nicol Schraudolph. Fast curvature matrix-vector products for second-order gradient descent. *Neural*
787 *computation*, 14:1723–38, 08 2002. doi: 10.1162/08997660260028683.
- 788 Cooper Simpson and Jaden Wang. PyTorch-ARC. github.com/RS-Coop/PyTorch-ARC,
789 2023. Adaptive Regularization with Cubics (ARC) optimizer for PyTorch.
- 790 Hadar Sivan, Moshe Gabel, and Assaf Schuster. FOSI: Hybrid first and second order optimization.
791 In *The Twelfth International Conference on Learning Representations*, 2024. URL [https:](https://openreview.net/forum?id=NvbeD9Ttkx)
792 [//openreview.net/forum?id=NvbeD9Ttkx](https://openreview.net/forum?id=NvbeD9Ttkx).
- 793 Daniel Soudry and Yair Carmon. No bad local minima: Data independent training error guar-
794 antees for multilayer neural networks. *ArXiv*, abs/1605.08361, 2016. URL [https://api.](https://api.semanticscholar.org/CorpusID:3029264)
795 [semanticscholar.org/CorpusID:3029264](https://api.semanticscholar.org/CorpusID:3029264).
796
- 797 Tijmen Tieleman and Geoffrey Hinton. Lecture 6.5-rmsprop: Divide the gradient by a running
798 average of its recent magnitude. *COURSERA: Neural networks for machine learning*, 4(2):26–31,
799 2012.
- 800 Philippe L. Toint. Nonlinear stepsize control, trust regions and regularizations for unconstrained
801 optimization. *Optimization Methods and Software*, 28(1):82–95, 2013. doi: 10.1080/10556788.
802 [2011.610458](https://doi.org/10.1080/10556788.2011.610458). URL <https://doi.org/10.1080/10556788.2011.610458>.
- 803 J.F. Traub. *Iterative Methods for the Solution of Equations*. AMS Chelsea Publishing Series.
804 Chelsea, 1982. ISBN 9780828403122. URL [https://books.google.co.il/books?](https://books.google.co.il/books?id=se3YdgFgz4YC)
805 [id=se3YdgFgz4YC](https://books.google.co.il/books?id=se3YdgFgz4YC).
806
- 807 Sharan Vaswani, Reza Babanezhad, Jose Gallego, Aaron Mishkin, Simon Lacoste-Julien, and
808 Nicolas Le Roux. To each optimizer a norm, to each norm its generalization. *ArXiv*, abs/2006.06821,
809 2020. URL <https://api.semanticscholar.org/CorpusID:219636073>.

- 810 Neha Wadia, Daniel Duckworth, Samuel S Schoenholz, Ethan Dyer, and Jascha Sohl-Dickstein.
811 Whitening and second order optimization both make information in the dataset unusable during
812 training, and can reduce or prevent generalization. In Marina Meila and Tong Zhang (eds.),
813 *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of *Pro-*
814 *ceedings of Machine Learning Research*, pp. 10617–10629. PMLR, 18–24 Jul 2021. URL
815 <https://proceedings.mlr.press/v139/wadia21a.html>.
- 816 Xiao Wang, Shiqian Ma, Donald Goldfarb, and Wei Liu. Stochastic quasi-newton methods for
817 nonconvex stochastic optimization. *SIAM Journal on Optimization*, 27, 07 2016. doi: 10.1137/
818 15M1053141.
- 819 Rachel Ward, Xiaoxia Wu, and Leon Bottou. AdaGrad stepsizes: Sharp convergence over nonconvex
820 landscapes. In Kamalika Chaudhuri and Ruslan Salakhutdinov (eds.), *Proceedings of the 36th*
821 *International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning*
822 *Research*, pp. 6677–6686. PMLR, 09–15 Jun 2019. URL [https://proceedings.mlr.](https://proceedings.mlr.press/v97/ward19a.html)
823 [press/v97/ward19a.html](https://proceedings.mlr.press/v97/ward19a.html).
- 824 Peng Xu, Farbod Roosta-Khorasan, and Michael Mahoney. Second-order optimization for non-convex
825 machine learning: An empirical study. In *Proceedings of the 2020 SIAM International Conference*
826 *on Data Mining*, 08 2017.
- 827 Peng Xu, Fred Roosta, and Michael W. Mahoney. Newton-type methods for non-convex optimization
828 under inexact hessian information. *Mathematical Programming*, 184(1):35–70, 2020. doi: 10.1007/
829 s10107-019-01405-z. URL <https://doi.org/10.1007/s10107-019-01405-z>.
- 830 Robert M. Young. 75.9 euler’s constant. *The Mathematical Gazette*, 75(472):187–190, 1991. ISSN
831 00255572. URL <http://www.jstor.org/stable/3620251>.
- 832 Xiao-Hu Yu and Guo-An Chen. On the local minima free condition of backpropagation learning.
833 *IEEE Transactions on Neural Networks*, 6(5):1300–1303, 1995. doi: 10.1109/72.410380.
- 834 Matthew D. Zeiler. Adadelta: An adaptive learning rate method. *ArXiv*, abs/1212.5701, 2012. URL
835 <https://api.semanticscholar.org/CorpusID:7365802>.
- 836 Pan Zhou, Jiashi Feng, Chao Ma, Caiming Xiong, Steven Chu Hong Hoi, and Weinan E.
837 Towards theoretically understanding why sgd generalizes better than adam in deep learning.
838 In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin (eds.), *Ad-*
839 *vances in Neural Information Processing Systems*, volume 33, pp. 21285–21296. Curran Asso-
840 ciates, Inc., 2020. URL [https://proceedings.neurips.cc/paper_files/paper/](https://proceedings.neurips.cc/paper_files/paper/2020/file/f3f27a324736617f20abbf2fffd806f6d-Paper.pdf)
841 [2020/file/f3f27a324736617f20abbf2fffd806f6d-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2020/file/f3f27a324736617f20abbf2fffd806f6d-Paper.pdf).
- 842 Milija Županski. A preconditioning algorithm for large-scale minimization problems. *Tellus A:*
843 *Dynamic Meteorology and Oceanography*, Jan 1993. doi: 10.3402/tellusa.v45i5.15048.

849 A QUASI-NEWTON CHALLENGES AND PROPOSED SOLUTIONS

850 **Notation 9.** Let $A \in \mathbb{R}^{n \times n}$. We mark $\lambda_{\min}(A)$, $\lambda_{\max}(A)$ the minimal/maximal eigenvalue of A ,
851 respectively, and their ratio $\kappa(A) = \frac{\lambda_{\max}(A)}{\lambda_{\min}(A)}$ the condition number of A .
852

853 Some of the challenges involved in training neural networks include:
854

- 855 • Because the models and data often have very complex structures, obtaining precisely optimal
856 parameters is often computationally prohibitive. As a result, one must satisfy oneself with
857 a small degree of suboptimality in the model’s parameters, chosen to be small enough to
858 satisfy one’s needs while not exhausting the computational capacity at hand.
- 859 • Since many model architectures (e.g. artificial neural networks) have very complex struc-
860 tures, the loss function is generally non-convex as a function of the model’s parameters.
861 This makes finding the globally optimal choice of parameters an NP-hard problem (Pardalos
862 & Vavasis, 1991; Manders & Adleman, 1978). As a result, one must satisfy oneself with
863 merely a local minimum of the loss function (that is, a point at which the norm of the gradient

w.r.t. the model parameters is zero, and the function is locally convex, or equivalently, the Hessian is positive semi-definite). This is often considered sufficient (see Soudry & Carmon (2016); Kawaguchi & Bengio (2019); Kawaguchi (2016); Nguyen et al. (2019)), however this does not apply to saddle points and local maxima (points at which the gradient norm is zero but the Hessian is not positive definite). Although some work has been done on trying to eliminate this problem by eliminating local- but not global- minima via overparameterization (Yu & Chen, 1995), further work (Ding et al., 2019) has shown that this does not scale to deep neural networks.

- As mentioned previously, there is no single universally optimal optimizer, even among existing optimizers.

As a result, sophisticated optimizers are necessary to contend with different neural network training scenarios. Restricting ourselves to quasi-newton optimization algorithms, scenarios with $n \gg 0$ (a common theme in machine learning, where n may be in the millions, billions, or even trillions, as GPT4 (OpenAI et al., 2024) is rumored to have. See Patel & Wong (2023)) are that computing and inverting the Hessian (with respective complexities $\mathcal{O}(n^2)$, $\mathcal{O}(n^3)$) may be computationally prohibitive. Also, one must ensure that

$$\Phi_t \succeq 0 \tag{10}$$

to ensure that $\theta_{t+1} - \theta_t$ is a descent direction of f . This is because $-\nabla f(\theta_t)$ is a descent direction of f , which implies that for all $v \in \mathbb{R}^n$, $-\alpha \cdot \nabla f(\theta_t)^\top v \cdot v^\top$ is a descent direction for $\alpha > 0$ and an ascent direction for $\alpha < 0$. However, if (λ_i, v_i) is an eigenvalue-eigenvector pair of Φ_t with $\lambda_i < 0$ then $-v_i^\top \cdot \Phi_t \nabla f(\theta_t) \cdot v_i = -\lambda_i \nabla f(\theta_t)^\top \cdot v_i \cdot v_i$ which is an ascent direction, and then a better preconditioner could immediately be obtained by taking $\tilde{\Phi}_t$ with eigenpairs $(\tilde{\lambda}_j, \tilde{v}_j)$, $\tilde{v}_j = v_j$, $\tilde{\lambda}_j =$

$$\begin{cases} \lambda_j & j \neq i \\ 0 & j = i \end{cases} \text{ to prevent an ascent in the subspace (a.k.a. eigenspace) } \textit{span}(v_i).$$

Three common ways to contend with these challenges are:

- **The Hessian-Free approach** Making use of Pearlmutter (1994) to compute Hessian-vector products without explicit computation of the Hessian, one uses conjugate-gradient (Olver & Shakiban, 2006) iterations to compute progressively finer approximations to $(\mathcal{H}(\theta_t))^{-1} \cdot \nabla f(\theta_t)$, stopping when one reaches a dimension with negative curvature. See, for instance, Martens (2010).
- **The Lanczos eigendecomposition approach** Making use of Lanczos iterations (Olver & Shakiban, 2006), one decomposes the Hessian into its eigendecomposition, and explicitly edits its eigenvalues. See, for instance, Dauphin et al. (2014); Sivan et al. (2024).
- **The Gauss-Newton approach** Using the generalized Gauss-Newton approximation to the Hessian (Esposito & Floudas, 2001; Schraudolph, 2002), one can obtain a matrix which has the following good properties:
 - Well approximated by a Kronecker product (sparse representation), which allows one to represent it and multiply by it very cheaply
 - Positive semi-definite
 - Can be computed with only a first-order loss function gradient oracle
 - Well approximates the true loss Hessian, when the second derivative of the model or the residual loss $(f(\theta_t) - f(\theta^*))$ is insignificant next to the generalized Gauss Newton

Some examples of this approach include Agarwal et al. (2019); Botev et al. (2017); Gupta et al. (2018); Martens & Grosse (2015); Goldfarb et al. (2020); Anil et al. (2020). Of particular note are examples that make diagonal approximations to the Gauss-Newton, as noted by Martens (2020), that are most often viewed as first-order methods, such as Adagrad (Duchi et al., 2011a), RMSProp (Tieleman & Hinton, 2012), and Adam (Kingma & Ba, 2014). As noted by Martens (2020), due to the strong connection between the generalized Gauss-Newton and the Fischer Information matrix (when the loss function is cross-entropy loss (Good, 1952)), one can achieve certain theoretical benefits when using such methods, such as Fischer efficiency; see Amari (1998) for instance, which views θ_t as an unbiased estimator of θ^* of f , and uses the Cramer-Rao inequality (Jansen & Claeskens, 2011) to

918 lower-bound the minimal number of iterations required to minimize the variance of said
 919 estimator as a function of the Fischer Information due to the number of samples consumed
 920 by each iteration.
 921

922 See Nocedal & Wright (2006, Chapters 3.3,3.4) for further discussion of these approaches.

923 In order for a minimization problem to be well-defined, one must assume that f is lower-bounded.
 924 We can infer from this that any subset of the domain space in which f is concave must be a bounded
 925 set (because nonconstant concave functions with unbounded domains are not lower-bounded); this
 926 means that the second-order Taylor approximation of the function must have a bounded neighborhood
 927 in which it approximates the function well. Additionally, even in subsets of the domain space in
 928 which f is convex, the neighborhood in which the second-order Taylor approximation of the loss
 929 function well-approximates the true loss function may be bounded. To address this, two common
 930 approaches been proposed in the literature, namely:

- 931 • **The Trust Regions Approach**, which explicitly maintains a radius of the neighborhood
 932 in which the second-order Taylor polynomial is a good approximation of the function, and
 933 bounds the step size to that radius. See Conn et al. (2000), Nocedal & Wright (2006, Chapter
 934 4).
- 935 • **The Cubic Regularization Approach**, which assumes that Hessian is Lipschitz continuous
 936 (using the L2 vector-induced matrix norm to measure distances between Hessians), and as
 937 such can upper bound the distance between two points of the function using a third-order
 938 polynomial (discussed below, see Lemma 4.1.14 from Dennis & Schnabel (1983)). See
 939 Nesterov & Polyak (2006) for an algorithm based on this approach that adaptively estimates
 940 the Hessian-Lipschitz parameter.
 941

942 B OTHER CONVERGENCE RATE MEASURES

943 **Convergence rates to first-order criticality** Most works on convergence rates in the non-convex
 944 regime bound the number of iterations necessary to achieve first-order criticality ($\|\nabla f(\theta_t)\|_2 = 0$) by
 945 means of finding an ϵ_g -stationary point (a point at which $\|\nabla f(\theta_t)\|_2 \leq \epsilon_g$). The seminal work Wang
 946 et al. (2016) provide a convergence rate bound for general optimizers (with very weak assumptions) in
 947 the non-convex regime of $\mathcal{O}\left(\kappa^{\frac{2}{1-\nu}}(\Phi_t) \cdot \epsilon_g^{-\frac{1}{1-\nu}}\right)$ with learning rate $\alpha_t = \mathcal{O}(t^{-\nu})$ and $\nu \in (0.5, 1)$.
 948

949 However, this bound is minimized by setting Φ_t to the minimizer of $\kappa(\Phi_t)$, which is a scalar matrix;
 950 this is equivalent to gradient descent, a first-order method. Experiments (see Sivan et al. (2024), for
 951 instance) and theory show that higher-order methods can achieve faster rates of convergence in our
 952 setting, demonstrating looseness of this convergence rate bound. See also D’efossez et al. (2020) who
 953 give such convergence rate bounds (requiring t iterations, for t s.t. $\frac{\sqrt{t}}{\log(t)} = \Omega(\epsilon_g^{-1})$) for Adam and
 954 Adagrad, and Ward et al. (2019) who give such convergence rate bounds (at $\mathcal{O}(\epsilon_g^{-1})$) for gradient
 955 descent with Adagrad-grafted step-sizes (see Agarwal et al. (2022) for a discussion on learning rate
 956 grafting).
 957

958 **Convergence rates to second-order criticality** A few go further in bounding the number of steps
 959 required to achieve second-order criticality (a point satisfying $\|\nabla f(\theta_t)\|_2 < \epsilon_g, -\lambda_{\min}(\mathcal{H}(\theta_t)) <$
 960 ϵ_H). For instance, Nesterov & Polyak (2006); Cartis et al. (2011b); Xu et al. (2020) provide such
 961 bounds (at $\mathcal{O}(\max(\epsilon_g, \epsilon_H)^{-3})$) on variants of the ARC algorithm, and Levy (2016); Jin et al. (2017);
 962 Ge et al. (2015) provide such bounds for varieties of SGD. This is of great importance since as noted,
 963 local minima are generally considered sufficiently optimal while local maxima/saddle points are not,
 964 despite being impossible to distinguish with only first-order criticality information. To the best of
 965 our knowledge, however, no such bounds exist in the general setting, nor do they exist for the vast
 966 majority of existing optimization algorithms.
 967

968 **Convergence rate dependence on preconditioner quality** One possible quality metric for Φ_t
 969 is given by $\eta_t \triangleq \left\| (I - \mathcal{H}(\theta_t) \cdot \Phi_t) \cdot \frac{\nabla f(\theta_t)}{\|\nabla f(\theta_t)\|_2} \right\|_2$. In the convex regime, Nocedal & Wright (2006,
 970 Chapter 7.1) assume $\sup_t \eta_t < 1$ and prove that first-order criticality may be reached within
 971

$\mathcal{O}\left(\frac{\log \epsilon}{\log \frac{1+\sup_t(\eta_t)}{2}}\right)$ iterations. Adding an assumption of Lipschitz-continuity of the Hessian, they prove quadratic convergence to first-order criticality. O’Leary-Roseberry et al. (2019), in contrast, do not assume convexity but provide a bound on the parameter gap $\|\theta_t - \theta^*\|_2$ for η_t satisfying the Eisenstat-Walker (Eisenstat & Walker, 1996; Dembo et al., 1982) condition $\eta_t \leq \|\nabla f(\theta_t)\|_2$ on a Tikhonov-regularized Hessian. Like Wang et al. (2016), however, here too the constant in their bound is inversely proportional to $\zeta - \lambda_{\min}(\mathcal{H}(\theta_t))$ with ζ the Tikhonov regularization constant, thus is minimized by taking $\zeta \rightarrow \infty$, eliminating all second-order information and reverting to simple gradient descent. As before, this implies looseness due to the empirical success of making use of second-order methods.

C ALGORITHM ELMO’S DESCENT RATE

An important factor in deciding how much computational effort to invest each optimization iteration is the ratio between the additional computational burden and the corresponding model quality improvement (as measured by the loss function) obtained by that iteration. To that end, we demonstrate that an upper bound on algorithm ELMO’s performance has quickly diminishing rewards for additional iterations. Counter-intuitively, this is a good thing - it means that as long as the algorithm converges to an acceptable minimum point, just a few iterations are likely to be necessary in practice - since any more than that will not have much of an effect on the model’s quality anyway.

Theorem C.1. *Worst case-optimal descent rate* Let f be a function with Lipschitz-continuous Hessian. After t iterations, algorithm ELMO satisfies

$$f(\theta_0) - f(\theta_t) = \mathcal{O}(\log t) \quad (11)$$

Although the above theorem gives only an upper bound on the model’s performance, we demonstrate that it is actually within a constant multiplicative factor of the algorithm’s lower bound.

Lemma C.2. *Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ satisfying assumptions 1 and 2. Algorithm ELMO satisfies*

$$|m_t^i(\Delta\theta_t^{*\top} v_i)| \leq 5 |M_t^i(\Delta\theta_t^{*\top} v_i)|$$

D GENERALIZATION OF PREVIOUS QUASI-NEWTON PRECONDITIONER QUALITY METRICS

Notation 10. Taking $(\lambda_i)_{i=1}^n$ the eigenvalues of $\mathcal{H}(\theta)$ for some θ , note that since n is finite, there exist $L^+ \triangleq \max_i\{L^i : \lambda_i > 0\}$, $L^- \triangleq \max_i\{L^i : \lambda_i \leq 0\}$.

Since most prevalent quasi-Newton algorithms apply a principled approach only to the convex subspaces of the loss function domain space and even then only when the curvature shift is negligible, we examine the special case of our metric when $\lambda_i > 0$ (when the loss function is convex over the domain subspace under examination) and show that our quality metric for quasi-Newton algorithm steps generalizes previous metrics. When $\lambda_i > 0$, we have

$$|\Delta\Delta^i\theta'_t| = \left| 1 - \frac{\nabla f(\theta_t)^\top}{\nabla f(\theta_t)^\top v_i} \cdot (\alpha_t \Phi_t \mathcal{H}(\theta_t)) \cdot v_i \cdot \frac{\sqrt{1 + 2L_t^i \cdot \frac{|\nabla f(\theta_t)^\top v_i|}{\lambda_i^2}} + 1}{2} \right| \quad (12)$$

Županski (1993) introduce the "Effective Hessian" (a.k.a. the "Preconditioned Hessian") as $\mathcal{I}_t = \alpha_t \Phi_t \mathcal{H}(\theta_t)$, whose condition number used as a quality metric for preconditioners; ideally, $\kappa(\mathcal{I}_t) < \kappa(\mathcal{H}(\theta_t))$. The Effective Hessian may be plainly seen in equation 12.

Mark $r_t \triangleq (I - \mathcal{H}(\theta_t) \cdot \Phi_t) \cdot \frac{\nabla f(\theta_t)}{\nabla f(\theta_t)^\top v_i}$; this is the 1-dimensional version of the quality metric η_t for Φ_t used by Nocedal & Wright (2006, Chapter 7.1) and mentioned in appendix B (now redefined by projecting $\nabla f(\theta_t)$ onto the i -th eigenspace instead of taking its full norm). When $L^+ \approx 0$ (i.e. when the loss function curvature shift is negligible), equation 12 simplifies to

$$|\Delta\Delta^i\theta'_t| \approx |r_t^\top \cdot v_i|$$

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

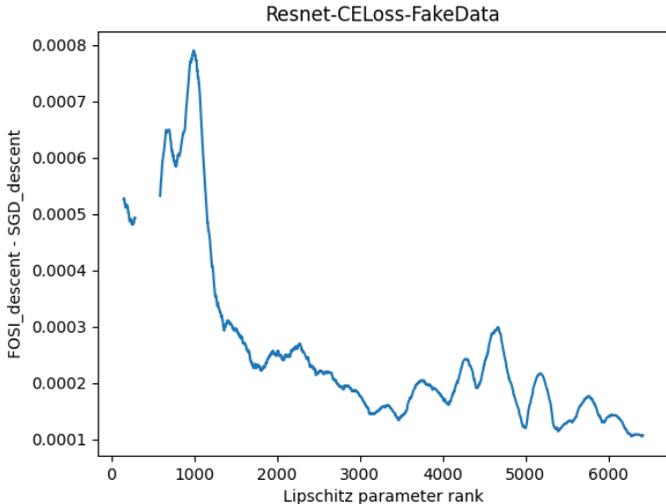


Figure 3: Inverse relation between a convex-subspace Lipschitz parameter and corresponding descent superiority of quasi-Newton method

E A QUALITY PREDICTOR FOR NEWTON’S METHOD

An important challenge is finding the best balance between per-iteration computational burden and expected loss function descent. We set out to provide such a metric due to equation 12 by showing that the expected descent in a given eigenspace is an approximately monotonically decreasing function of the corresponding Lipschitz parameter.

Figure 3 shows an example of this phenomenon by plotting the quasi-Newton superiority (how much better a quasi-Newton method will work than a first-order method, defined as $(f(\theta_t) - f(\theta_{t+1}^{Newton})) - (f(\theta_t) - f(\theta_{t+1}^{SGD}))$) against the convex Lipschitz parameter rank. Here too we represent the full spectrum of convex Lipschitz parameters with the single Lipschitz parameter representing the eigenspace with the greatest eigenvalue; nevertheless, a qualitative inverse correlation is clear. Pearson correlation coefficient values (Pearson, 1895) are shown in table 1, as well as p-values of a test of the null hypothesis that the distributions underlying the samples are uncorrelated and normally distributed. The Scipy manual writes:

The p-value roughly indicates the probability of an uncorrelated system producing datasets that have a Pearson correlation at least as extreme as the one computed from these datasets.

Here too, the detailing of our experiment settings is given in appendix J, as well as further detailing on figure 3.

[(0, 70.16), (0.7, 16.19), (0.8, 17.9), (0.9, 14.65), (0.95, 0.65), (0.99, 0.91), (0.995, 0.65), (0.999, 0.02)]

Dataset	Pearson r	p -value
CIFAR10	-0.245341	10^{-107}
FakeData	-0.026608	0.031120
MNIST	-0.368788	10^{-300}

Table 1: Pearson r inverse correlation between quasi-Newton superiority and Lipschitz parameter

Since the Lipschitz parameters are approximately locally constant throughout training as shown in the previous section, this reverse correlation may be used to help practitioners decide how much computational burden is worth putting into each iteration, given that even an exact Newton step may not be significantly superior to first-order methods when the curvature drift (as measured by Lipschitz parameters) is significantly large; hyperparameter optimization algorithm selection may then follow accordingly. We present experiments validating this selection method in appendix L. Alternatively, practitioners may

choose to use ARC steps instead of first-order methods, when the Lipschitz parameter is significantly large.

F COMPARISON OF ELMO TO SELECT RELATED METHODS

ELMO is strikingly similar to Cauchy’s method (not to be confused with Cauchy’s Steepest Descent method (Nocedal & Wright, 2006, Chapter 4.1)) and Newton’s method mentioned above. In this section, we note the similarity between them, and the sources of the differences between them.

F.1 COMPARISON TO CAUCHY’S METHOD

Cauchy’s method (Traub, 1982) is very reminiscent of ELMO:

$$\begin{aligned}
(\theta_{t+1} - \theta_t)_{cauchy}^\top \cdot v_i &\triangleq -\frac{2}{1 + \sqrt{1 - 2\frac{L_t^i \cdot \nabla f(\theta_t)^\top v_i}{\lambda_i^2(\theta_t)}}} \cdot \frac{\nabla f(\theta_t)^\top v_i}{\lambda_i(\theta_t)} \\
&= -\frac{2}{1 + \frac{1}{|\lambda_i(\theta_t)|} \cdot \sqrt{\lambda_i^2(\theta_t) - 2L_t^i \cdot \nabla f(\theta_t)^\top v_i}} \frac{\nabla f(\theta_t)^\top v_i}{\lambda_i(\theta_t)} \\
&= \frac{-2\nabla f(\theta_t)^\top v_i}{\lambda_i(\theta_t) + \frac{\lambda_i(\theta_t)}{|\lambda_i(\theta_t)|} \cdot \sqrt{\lambda_i^2(\theta_t) - 2L_t^i \cdot \nabla f(\theta_t)^\top v_i}} \\
&= \frac{1}{L_t^i} \cdot \frac{-2L_t^i \cdot \nabla f(\theta_t)^\top v_i}{\lambda_i(\theta_t) + \frac{\lambda_i(\theta_t)}{|\lambda_i(\theta_t)|} \cdot \sqrt{\lambda_i^2(\theta_t) - 2L_t^i \cdot \nabla f(\theta_t)^\top v_i}} \\
&= \frac{1}{L_t^i} \cdot \frac{(\lambda_i^2(\theta_t) - 2L_t^i \cdot \nabla f(\theta_t)^\top v_i) - \lambda_i^2(\theta_t)}{\lambda_i(\theta_t) + \frac{\lambda_i(\theta_t)}{|\lambda_i(\theta_t)|} \cdot \sqrt{\lambda_i^2(\theta_t) - 2L_t^i \cdot \nabla f(\theta_t)^\top v_i}} \\
&= \frac{\sqrt{\lambda_i^2(\theta_t) - 2L_t^i \cdot \nabla f(\theta_t)^\top v_i} - \sqrt{\lambda_i^2(\theta_t)}}{L_t^i} \cdot \frac{|\lambda_i(\theta_t)|}{\lambda_i(\theta_t)} \\
&= \frac{\sqrt{\lambda_i^2(\theta_t) - 2L_t^i \cdot \nabla f(\theta_t)^\top v_i} - \frac{|\lambda_i(\theta_t)|}{\lambda_i(\theta_t)} \cdot \lambda_i(\theta_t)}{L_t^i} \cdot \frac{|\lambda_i(\theta_t)|}{\lambda_i(\theta_t)} \\
&= \frac{\frac{|\lambda_i(\theta_t)|}{\lambda_i(\theta_t)} \cdot \sqrt{\lambda_i^2(\theta_t) - 2L_t^i \cdot \nabla f(\theta_t)^\top v_i} - \lambda_i(\theta_t)}{L_t^i}
\end{aligned}$$

The difference between our minimization step and their step is the sign on the squareroot. The difference stems from removing the absolute value in equation 3’s 3rd-order term and taking the negative root of its derivative, due to the difference in goals: we attempt to minimize the function, leading us to select the positive step. They attempt to find the function’s critical points, leading them to select the negative step.

F.2 COMPARISON TO NEWTON’S METHOD

Unlike Cauchy’s method, Newton’s method (in optimization) makes a second-order approximation to the function’s gradient. This is equivalent to the Hessian being constant, which is equivalent to

1134 $L_H = 0$. Indeed, taking the limit of equation 7 when $L_H \rightarrow 0^+$, we recover Newton’s method:
 1135
 1136

$$\begin{aligned}
 \lim_{L_H \rightarrow 0^+} \Delta \theta_t^{*\top} v_i &= \lim_{L_H \rightarrow 0^+} - \frac{2 \nabla f(\theta_t)^\top v_i}{\sqrt{\lambda_i^2(\theta_t) - 2L_t^i \cdot \nabla f(\theta_t)^\top v_i + \lambda_i(\theta_t)}} \\
 &= \begin{cases} -\frac{\nabla f(\theta_t)^\top v_i}{\lambda_i} & \lambda_i > 0 \\ \infty & \lambda_i < 0 \end{cases}
 \end{aligned} \tag{13}$$

1145 G RELEVANCE OF THE HESSIAN-LIPSCHITZ DISTRIBUTION

1148 One source of interest in the Hessian-Lipschitz distribution is for optimization algorithms (e.g. ARC)
 1149 that make use of these parameters for the loss function modelling stage of each iteration. This may
 1150 reduce computational complexity by reducing the number of parameters one must compute at each
 1151 iteration, however appendix H shows that poorly estimating the Lipschitz parameters can have a
 1152 detrimental effect on an algorithm’s descent rate (thereby increasing the number of iterations the
 1153 algorithm will require to converge).

1154 Another source of interest in these parameters’ distribution is in explaining the effectiveness of second-
 1155 order quasi-Newton algorithms that implicitly assume the Lipschitz parameters are insignificant (i.e.
 1156 very close to zero), since their model of the loss function is a quadratic Taylor polynomial (i.e. no
 1157 curvature shift); this may be seen from equation 12 which shows optimality of Newton’s method only
 1158 when $\lambda_i > 0$ and $L^+ = 0$. As seen in section 5, the Lipschitz parameters are not generally small
 1159 by any means, however the Lipschitz parameters of the subspaces in which they work (the convex
 1160 subspaces - see the implementation of Sivan et al. (2024), for instance, which applies Newton’s
 1161 method only on subspaces with significantly convex subspaces) are in fact small in certain settings.

1164 H CONVERGENCE RATE DEPENDENCE ON HESSIAN-LIPSCHITZ PARAMETER

1167 As noted by Griewank (1981), the Hessian-Lipschitz parameter (in our case, the respective constants
 1168 of each eigenspace) may be computationally difficult to obtain precisely, leading some optimization
 1169 algorithms to estimate it approximately instead of computing it precisely (e.g. ARC). In order to
 1170 balance the computational burden of computing it to a high degree of exactitude with the degradation
 1171 of an algorithm’s convergence rate that comes with poor estimations, we study the effects of the
 1172 Hessian-Lipschitz parameter on $M_t^i (\Delta \theta_t^\top v_i)$.

1175 H.1 LIPSCHITZ ROBUSTNESS

1177 To address the convergence rate’s robustness to overly conservative L_t^i , we consider the case when
 1178 $L_t^i \rightarrow \infty$.

1180 **Theorem H.1.** *Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ satisfying assumptions 1 and 2. Then*

$$M_t^i (\Delta \theta_t^{*\top} v_i) = \Theta \left(\frac{1}{\sqrt{L_t^i}} \right)$$

1187 when $L_t^i \rightarrow \infty$

1188 *Proof.*

$$\begin{aligned}
1189 & \lim_{L_t^i \rightarrow \infty} \frac{M_t^i (\Delta \theta_t^{*\top} v_i)}{-12(-\nabla f(\theta_t)^\top v_i)^{1.5} - 36(-\nabla f(\theta_t)^\top v_i) + \sqrt{2}\sqrt{-\nabla f(\theta_t)^\top v_i} + 3\sqrt{2}}{\frac{3}{-\nabla f(\theta_t)^\top v_i} - 18\sqrt{2}} \cdot \frac{1}{\sqrt{L_t^i}} \\
1190 & = \lim_{L_t^i \rightarrow \infty} -\frac{1}{\sqrt{\frac{\lambda_i^2(\theta_t)}{L_t^i(-2 \cdot \nabla f(\theta_t)^\top v_i)} + 1} + \frac{\lambda_i(\theta_t)}{\sqrt{L_t^i}\sqrt{-2 \cdot \nabla f(\theta_t)^\top v_i}}} \\
1191 & \quad \cdot \frac{6\sqrt{-\nabla f(\theta_t)^\top v_i}}{6\sqrt{-\nabla f(\theta_t)^\top v_i} - 2 \cdot \nabla f(\theta_t)^\top v_i} \\
1192 & \quad \cdot \frac{2 \cdot \nabla f(\theta_t)^\top v_i}{6\sqrt{-\nabla f(\theta_t)^\top v_i} - 2 \cdot \nabla f(\theta_t)^\top v_i} \\
1193 & \quad \cdot \left(\sqrt{\frac{\lambda_i^2(\theta_t)}{L_t^i(-2 \cdot \nabla f(\theta_t)^\top v_i)} + 1} - \frac{\lambda_i(\theta_t)}{\sqrt{L_t^i}\sqrt{-2 \cdot \nabla f(\theta_t)^\top v_i}} \right)^3 \\
1194 & \quad + \frac{1}{\sqrt{L_t^i}} \cdot \frac{1 + 6\sqrt{2}\nabla f(\theta_t)^\top v_i}{2(6\sqrt{-\nabla f(\theta_t)^\top v_i} - 2 \cdot \nabla f(\theta_t)^\top v_i)} \left(\frac{\lambda_i(\theta_t)}{\sqrt{\frac{\lambda_i^2(\theta_t)}{L_t^i} - 2 \cdot \nabla f(\theta_t)^\top v_i + \frac{\lambda_i(\theta_t)}{\sqrt{L_t^i}}}} \right)^2 \\
1195 & \quad \cdot \left(\frac{1}{6} + \sqrt{2}\nabla f(\theta_t)^\top v_i \right) \sqrt{-2 \cdot \nabla f(\theta_t)^\top v_i} \\
1196 & = 1
\end{aligned}$$

1214 \square

1217 H.2 BENEFIT OF LIPSCHITZ TIGHTNESS

1218 To see how minimizing L_t^i as much as possible benefits the bound, we consider the case when $L_t^i \rightarrow 0^+$.

1219 **Theorem H.2.** *Benefit of Lipschitz tightness: concave subspaces*

1220 *Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ satisfying assumptions 1 and 2. If $\lambda_i(\theta_t) \leq 0$ then*

$$1221 M_t^i (\Delta \theta_t^{*\top} v_i) = \Theta \left(-\frac{1}{L_t^{i,2}} \right)$$

1222 when $L_t^i \rightarrow 0^+$

1223 *Proof.*

$$\begin{aligned}
1224 & \lim_{L_t^i \rightarrow 0^+} \frac{M_t^i (\Delta \theta_t^{*\top} v_i)}{\frac{\frac{2}{3}\lambda_i^3(\theta_t)}{L_t^{i,2}}} \\
1225 & = \lim_{L_t^i \rightarrow 0^+} \left(3 \left(\frac{\sqrt{1 - 2L_t^i \cdot \frac{\nabla f(\theta_t)^\top v_i}{(-\lambda_i(\theta_t))^2}} + 1}{2} \right)^2 - 2 \left(\frac{\sqrt{1 - 2L_t^i \cdot \frac{\nabla f(\theta_t)^\top v_i}{(-\lambda_i(\theta_t))^2}} + 1}{2} \right)^3 \right) \\
1226 & \quad + L_t^i \cdot \frac{3}{2\lambda_i^3(\theta_t)} \cdot \left(\sqrt{\lambda_i^2(\theta_t) - 2L_t^i \cdot \nabla f(\theta_t)^\top v_i} - \lambda_i(\theta_t) \right) \cdot \nabla f(\theta_t)^\top v_i \\
1227 & = 1
\end{aligned}$$

1240 \square

Theorem H.3. *Benefit of Lipschitz tightness: convex subspaces*

Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ satisfying assumptions 1 and 2. If $\lambda_i(\theta_t) > 0$ then

$$M_t^i(\Delta\theta_t^{*\top} v_i) - M_t^i\left(\lim_{L_t^i \rightarrow 0^+} \Delta\theta_t^{*\top} v_i\right) = \Theta(L_t^i)$$

when $L_t^i \rightarrow 0^+$

Proof. We begin by noting that by equation 13, $\lim_{L_t^i \rightarrow 0^+} \Delta\theta_t^{*\top} v_i = \frac{|\nabla f(\theta_t)^\top v_i|}{\lambda_i(\theta_t)}$. Plugging this into M_t^i :

$$\begin{aligned} & \lim_{L_t^i \rightarrow 0^+} \frac{M_t^i(\Delta\theta_t^{*\top} v_i) - M_t^i\left(\lim_{L_t^i \rightarrow 0^+} \Delta\theta_t^{*\top} v_i\right)}{\frac{L_t^i}{\frac{1}{2}\lambda_i(\theta_t) \cdot \frac{(\nabla f(\theta_t)^\top v_i)^3}{\lambda_i^4(\theta_t)} - \frac{(\nabla f(\theta_t)^\top v_i)^3}{2\lambda_i^2(\theta_t)}} \\ &= \lim_{L_t^i \rightarrow 0^+} \frac{(\nabla f(\theta_t)^\top v_i)^3}{2\lambda_i^2(\theta_t)} \cdot \left(\frac{4}{\left(\sqrt{1 - 2L_t^i \cdot \frac{\nabla f(\theta_t)^\top v_i}{\lambda_i(\theta_t)}} + 1\right)^2} \right) \\ &+ \frac{1}{2}\lambda_i(\theta_t) \cdot \frac{(\nabla f(\theta_t)^\top v_i)^3}{\lambda_i^4(\theta_t)} \cdot \left(\frac{2\lambda_i(\theta_t)}{\sqrt{\lambda_i^2(\theta_t) - 2L_t^i \cdot \nabla f(\theta_t)^\top v_i + \lambda_i(\theta_t)}} + 1 \right) \\ &\cdot \left(\frac{2\lambda_i^2(\theta_t)}{\lambda_i(\theta_t) \sqrt{\lambda_i^2(\theta_t) - 2L_t^i \cdot \nabla f(\theta_t)^\top v_i + \lambda_i(\theta_t)}} \right) \left(\frac{2}{1 + \sqrt{1 - 2L_t^i \cdot \frac{\nabla f(\theta_t)^\top v_i}{\lambda_i(\theta_t)}}} \right) \\ &- \frac{L_t^i}{4} \cdot \frac{(\nabla f(\theta_t)^\top v_i)^4}{\lambda_i^5(\theta_t)} \left(\frac{2}{\sqrt{1 - 2L_t^i \cdot \frac{\nabla f(\theta_t)^\top v_i}{\lambda_i(\theta_t)}} + 1} \right) \left(\frac{2}{\sqrt{1 - 2L_t^i \cdot \frac{\nabla f(\theta_t)^\top v_i}{\lambda_i(\theta_t)}} + 1} \right) \\ &\cdot \left(\frac{1 + \frac{2}{\sqrt{1 - 2L_t^i \cdot \frac{\nabla f(\theta_t)^\top v_i}{\lambda_i(\theta_t)}} + 1} + \frac{4}{\left(\sqrt{1 - 2L_t^i \cdot \frac{\nabla f(\theta_t)^\top v_i}{\lambda_i(\theta_t)}} + 1\right)^2} \right) \\ &= 1 \end{aligned}$$

□

I QUALITATIVE EVIDENCE FROM PREVIOUS WORK ON LIPSCHITZ PARAMETERS

Experiments by Alain et al. (2018); Sagun et al. (2016); Ghorbani et al. (2019); Gur-Ari et al. (2018) show that the positive eigenvalues of the Hessian remain relatively stable throughout training, while the negative eigenvalues shrink rapidly. Alain et al. (2018) and Sagun et al. (2016) also show that the negative eigenvalues shift chaotically. Gur-Ari et al. (2018) show that when training a network on a classification task with k classes, then at least the eigenspace spanned by the k eigenvectors corresponding to the top k eigenvalues remains very stable. Sivan et al. (2024); Liu et al. (2024) also show that when training a neural network on a variety of tasks, the top k eigenvalues and their corresponding eigenvectors change very slowly. Alain et al. (2018) also show explicitly that the second-order Taylor approximation is a poor approximation of the loss function in the eigenspace corresponding to the negative eigenvalues (the concave eigenspace), but an excellent approximation

1296 in the eigenspace corresponding to the positive eigenvalues (the convex eigenspace); indeed, they
 1297 show that the optimal step in the convex eigenspace is well estimated by the Newton step, while there
 1298 is no correlation between the Hessian and the optimal step in the concave eigenspace. Using the
 1299 Lipschitz parameter as a measure of the rate of change of the Hessian in a given subspace (hence a
 1300 measure of the quality of a second-order Taylor approximation to a function and its corresponding
 1301 Newton step), this supports the claim that $L^+ \ll L^-$.

1302

1303 J LIPSCHITZ PARAMETER EXPERIMENTS

1304

1305 We tested our algorithm in 7 scenarios with the PyTorch 1.13.0 framework, each on a single NVIDIA
 1306 GeForce RTX 3090 GPU with the standard hyperparameters and settings for ARC:

1307

- 1308 • $\sigma_0^+ = \sigma_0^- = 1$
- 1309 • $\eta_1 = 0.1, \eta_2 = 0.9$
- 1310 • $\gamma_1 = \gamma_2 = 2$
- 1311 • Maximum sub-problem failures = 11
- 1312 • Maximum sub-problem iterations = 50,
- 1313 • Sub-problem tolerance = 10^{-6}
- 1314 • Lanczos eigendecomposition (Garber et al., 2016)
- 1315 • BFGS trinomial sub-problem solver (Nocedal & Wright, 2006, Chapter 6.1), (Feinman,
 1316 2021)
- 1317 • trinomial sub-problem maximal failures=11
- 1318 • trinomial sub-problem maximal iterations=50

1319

1320 The test scenarios⁴ include combinations of:

1321

- 1322 • Training ResNet18 artificial neural networks (He et al., 2015) for image classification on
 1323 MNIST, CIFAR10, and FakeData (random noise in place of images) with Cross-Entropy
 1324 Loss (CELoss), to evaluate the effect of changing data on the Lipschitz parameters
- 1325 • Training a CNN for image classification on CIFAR10 with CELoss, to evaluate the effect of
 1326 changing neural network architecture on the Lipschitz parameters.
- 1327 • Training a CNN⁵ autoencoder⁶ (LeCun, 1987) to compress MNIST

1328

1329 The classification CNN architecture:

1330

- 1331 1. A feature extractor consisting of 2 2D convolutional layers with 6 output channels for the
 1332 first and 16 output channels for the second. Both had kernel sizes of 5 pixels. Each of these
 1333 is followed by a ReLU nonlinearity and then 2x2 2D max pooling
- 1334 2. A 3-layer MLP classification head with hidden sizes (120, 84) and ReLU nonlinearities

1335

1336 The autoencoder CNN architecture:

1337

- 1338 • Encoder: 4 2D convolutional layers with respective output channel numbers (16,32,32,64)
 1339 and kernel sizes of 3 pixels for the first two and 5 pixels for the last two. The first two
 1340 have 1 pixel padding and the last two have 2 pixel paddings. After every layer we apply
 1341 LeakyReLU nonlinearity and after every 2 layers we apply 2x2 2D max pooling.
- 1342 • Decoder: 4 composite layers consisting of
- 1343 1. a 2D transpose-convolutional layer
- 1344 2. LeakyReLU nonlinearity (only for first and third composite layers)

1345

1346 ⁴Code available on Github at REDACTED

1347

1348 ⁵convolutional neural network

1349 ⁶With hidden dimensions 128-64-36-18-9-18-36-64-128, ReLU nonlinearities, and sigmoid nonlinearity on
 the output

- 1350 3. a 2d convolutional layer
 1351 4. LeakyReLU nonlinearity
 1352
 1353 • Decoder hidden channel sizes: 32-32-16-16-16-16-3
 1354 • Decoder kernel sizes: 2-5-5-5-2-3-5-3
 1355 • All strides are of size 1, except for the first and third transpose-convolutional layers of the
 1356 decoder, with stride of size 2
 1357 • Decoder paddings: 0-2-2-2-0-1-2-1
 1358

1359 Each experiment took several hours to run. All experiments (including those from above) shown in
 1360 figure 4.
 1361

1362 One caveat is that due to computational constraints, we use stochastic minibatch training for the
 1363 neural networks instead of using the full batch to compute the gradient and Hessian-vector products
 1364 at each iteration (see Bertsekas (1996) for an introduction to minibatch Monte-Carlo estimation of
 1365 a sum). However, Cartis et al. (2011a), notes that the adaptive Lipschitz parameter estimates may
 1366 account for this variance by being greater than the actual Lipschitz parameters. Thus, our claims of
 1367 $L^+ \ll 1$ are not affected (since our experiments effectively provide an upper bound on L^+) while
 1368 our claims of $L^- \gg 0$ are weakened. Since there is no reason to expect the variance on L^- to be
 1369 significantly greater than the variance on L^+ , however, our experiments remain valid.
 1370

1370 For visual clarity, the quasi-Newton superiority measurements in 3 are presented after:

- 1371 1. Clipping extreme values to the 10% - 90% quantile range
 1372 2. Gaussian smoothing, consisting of a rolling window of size 300 and standard deviation of
 1373 100
 1374

1375 J.1 COMPUTATION OF LIPSCHITZ PARAMETERS

1376
 1377 We modified the standard ARC algorithm to compute distinct Lipschitz parameters for the eigenspaces
 1378 corresponding to the minimal and maximal eigenvalues. Pseudocode for this algorithm is given
 1379 below.
 1380

1381 While lines 3-18 of EigenARC may technically be usable as the LIPSCHITZ subroutine of al-
 1382 gorithm ELMO above, each iteration requires $\Omega(n)$ evaluations of the loss function, which will be
 1383 computationally expensive if $n \gg 0$ and if the loss function is computationally heavy. This may be
 1384 ameliorated by performing these calculations only for a small subset of the eigenspaces like Sivan
 1385 et al. (2024), however we leave this to future work.
 1386

1387 K AMOS EXPERIMENT DETAILS

1388 In section 4 we present tests demonstrating the application of AMOS to training NanoGP on the
 1389 TinyShakespeare dataset. Below are the hyperparameters we selected for the experiment.
 1390

- 1391 • Our AdamW reference run takes the tuned hyperparameters from vanilla NanoGPT, namely
 1392 - learning rate = $6 \cdot 10^{-4}$, weight decay = 10^{-1} , $\beta_1 = 0.9$, $\beta_2 = 0.95$.
 1393 • Our first-order optimizer is SGD with learning rate set to $1.25 \cdot 10^{-5}$, selected by grid
 1394 hyperparameter search. We also tried adding a Cosine Annealing scheduler for SGD, but
 1395 this harmed SGD’s performance and had no qualitative effect on the training trajectory
 1396 comparison graph, so it is not included in our results.
 1397 • Our second-order optimizer is Sophia, with $\rho = 0.01$ and otherwise Karpathy (2022)’s
 1398 standard AdamW hyperparameters.
 1399 • Our Sophia and SGD reference optimization trajectories use optimizers with the same
 1400 hyperparameters but without dynamic optimizer selection (instead using pure SGD all the
 1401 way through or pure Sophia).
 1402 • The momentum coefficient for our meta-optimizer $\alpha_L = 0.99$, and the Lipschitz threshold
 1403 $L_{thresh} = 2 \cdot 10^{-10}$, $\epsilon_2 = 10^{-12}$.

1404
 1405
 1406
 1407
 1408
 1409
 1410
 1411
 1412
 1413
 1414
 1415
 1416
 1417
 1418
 1419
 1420
 1421
 1422
 1423
 1424
 1425
 1426
 1427
 1428
 1429
 1430
 1431
 1432
 1433
 1434
 1435
 1436
 1437
 1438
 1439
 1440
 1441
 1442
 1443
 1444
 1445
 1446
 1447
 1448
 1449
 1450
 1451
 1452
 1453
 1454
 1455
 1456
 1457

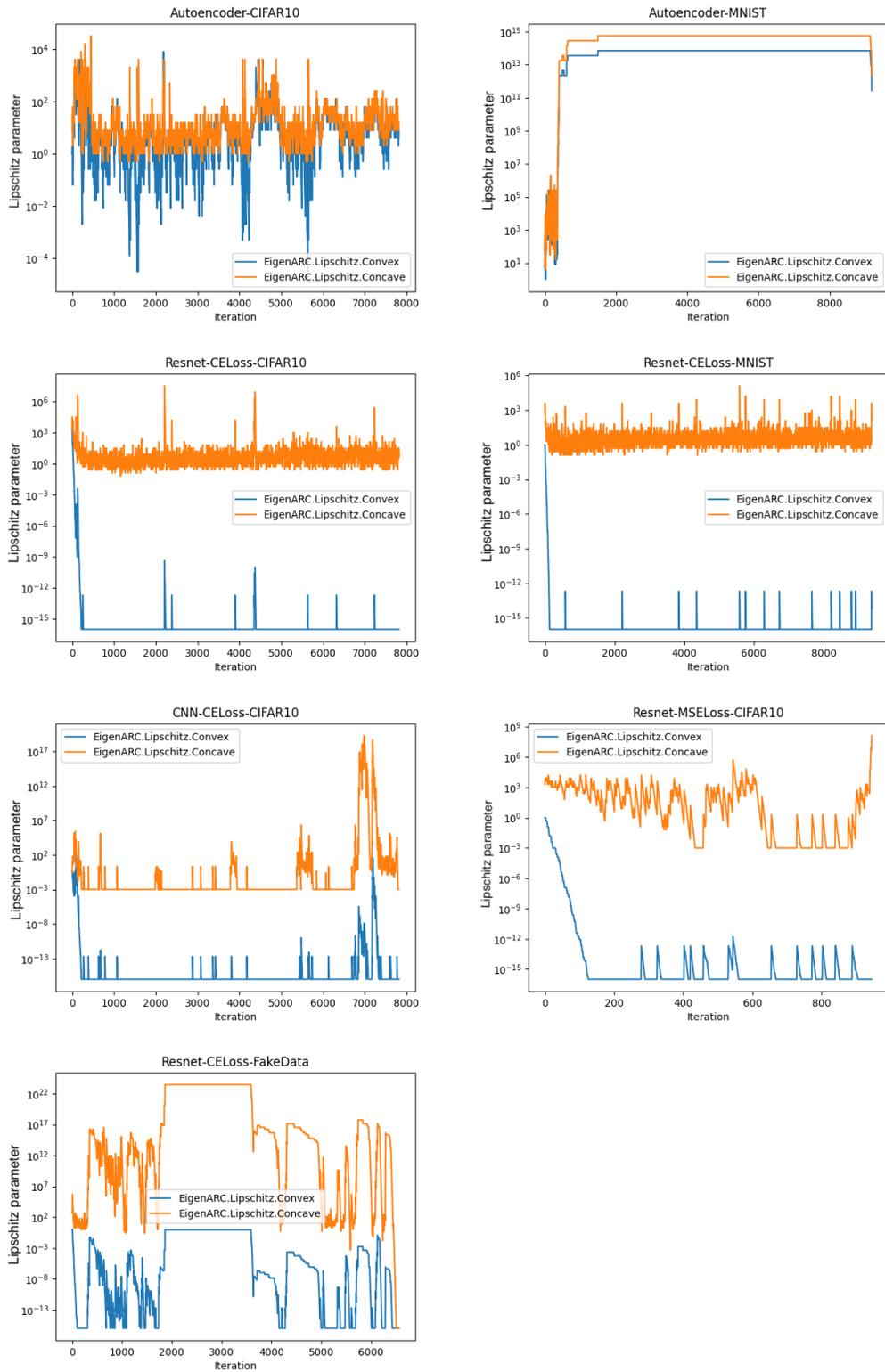


Figure 4: Comparisons of convex-subspace Lipschitz parameters to concave-subspace Lipschitz parameters. *Logarithmic scale*

1458
1459
1460
1461
1462
1463
1464
1465
1466
1467
1468
1469
1470
1471
1472
1473
1474
1475
1476
1477
1478
1479
1480
1481
1482
1483
1484
1485
1486
1487
1488
1489
1490
1491
1492
1493
1494
1495
1496
1497
1498
1499
1500
1501
1502
1503
1504
1505
1506
1507
1508
1509
1510
1511

Algorithm 3 Algorithm EigenARC

Require: $\epsilon \in \mathbb{R}^+, \theta_0 \in \mathbb{R}^n, \gamma_1 > 1 > \gamma_2 > 0, \eta_2 \geq \eta_1 > 0, (L_0^i)_{i=1}^n > 0, \text{EIGEN, BASE_OPT}$

- 1: $t \leftarrow 0$
- 2: **while** $\|\nabla f(\theta_t)\|_2 > \epsilon$ **do** ▷ While BASE_OPT hasn't converged yet
- 3: $(\lambda_i, v_i)_{i=1}^n \leftarrow \text{EIGEN}(\mathcal{H}(\theta_t))$
- 4: **if** $\text{ASSESS_LIPSCHITZ}((L_t^i)_{i=1}^n) > \eta_2$ **then** ▷ Overly conservative L_t^i
- 5: $(L_t^i)_{i=1}^n \leftarrow \gamma_2 \cdot (L_t^i)_{i=1}^n$
- 6: **else**
- 7: **if** $\text{ASSESS_LIPSCHITZ}((L_t^i)_{i=1}^n) < \eta_1$ **then** ▷ Overly liberal L_t^i
- 8: **while** $\text{ASSESS_LIPSCHITZ}((L_t^i)_{i=1}^n) > \eta_1$ **do** ▷ Raise all L_t^i assessment is passed
- 9: $(L_t^i)_{i=1}^n \leftarrow \gamma_1 \cdot (L_t^i)_{i=1}^n$
- 10: **end while**
- 11: **for** $i=1, \dots, n$ **do** ▷ Reduce the L_t^i that can be reduced without violating assessment
- 12: **while** $\text{ASSESS_LIPSCHITZ}((L_t^i)_{i=1}^n) > \eta_1$ **do**
- 13: $L_t^i \leftarrow \frac{L_t^i}{\gamma_1}$
- 14: **end while**
- 15: $L_t^i \leftarrow \gamma_1 \cdot L_t^i$
- 16: **end for**
- 17: **end if**
- 18: **end if**
- 19: $\theta_{t+1} \leftarrow \text{BASE_OPT}(\theta_t)$
- 20: $t \leftarrow t + 1$
- 21: **end while**

return $(L_t^i)_{i=1, \hat{t}=1}^{n, t}$

procedure $\text{ASSESS_LIPSCHITZ}((\hat{L}_t^i)_{i=1}^n)$

return
$$\frac{f(\theta_t) - f(\theta_t + \sum_{i=1}^n \Delta \theta_t^{*T} v_i (\hat{L}_t^i) \cdot v_i)}{-\sum_{i=1}^n M_t^i(\Delta \theta_t^{*T} v_i (\hat{L}_t^i))}$$

end procedure

1512
 1513
 1514
 1515
 1516
 1517
 1518
 1519
 1520
 1521
 1522
 1523
 1524
 1525
 1526
 1527
 1528
 1529
 1530
 1531
 1532
 1533
 1534
 1535
 1536
 1537
 1538
 1539
 1540
 1541
 1542
 1543
 1544
 1545
 1546
 1547
 1548
 1549
 1550
 1551
 1552
 1553
 1554
 1555
 1556
 1557
 1558
 1559
 1560
 1561
 1562
 1563
 1564
 1565

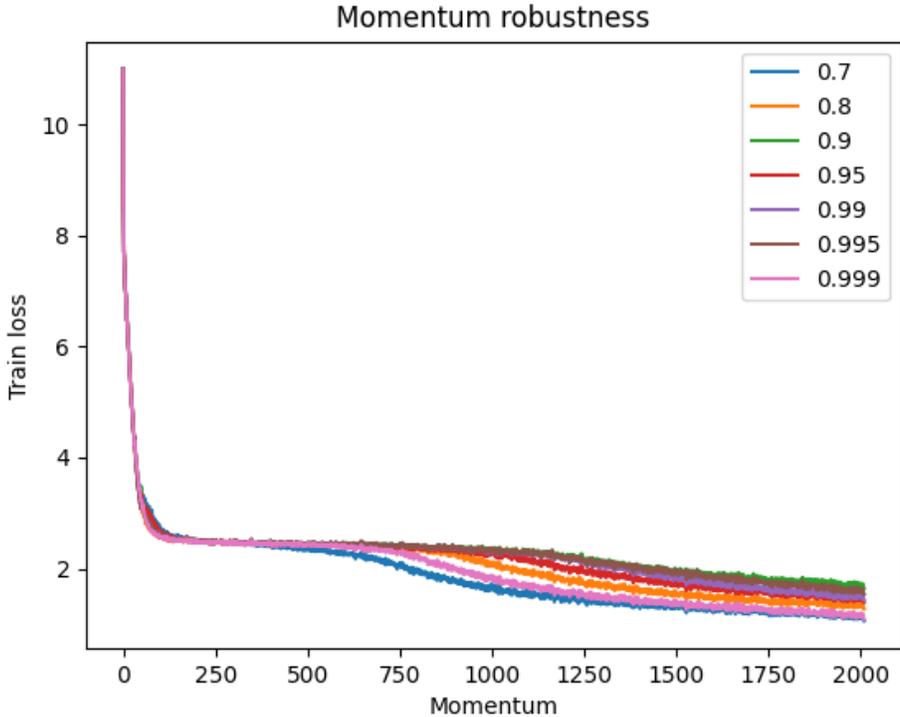


Figure 5: Training trajectories of AMOS with various different momentum settings. There is little variation in the long run.

Lipschitz momentum (α_L)	Offloaded parameters (%)
0.7	16.19
0.8	17.90
0.9	14.65
0.95	0.65
0.99	0.91
0.995	0.65
0.999	0.02

Table 2: Percent of parameters offloaded to SGD

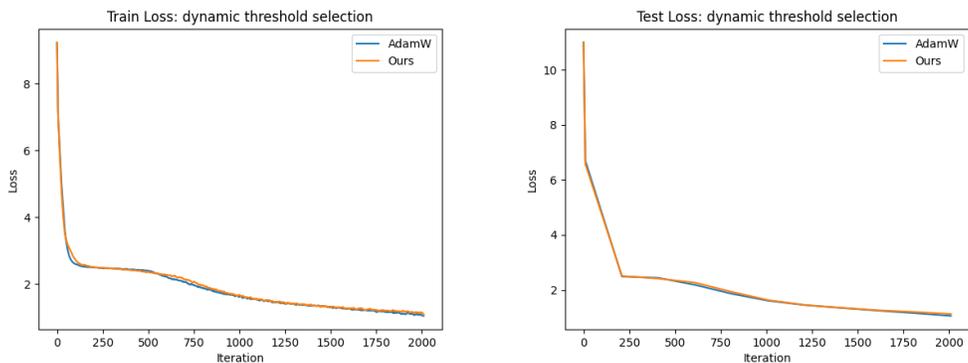
- For hardware efficiency reasons, we transfer parameters to SGD at a per-tensor granularity instead of per-scalar (although we suggest subdividing tensors into constant sizes in future work). The mean Lipschitz parameter of the entire tensor is compared against the preset threshold to determine if it is transferred to the first-order optimizer or not.

To demonstrate our AMOS’s robustness to hyperparameter selection, we rerun the above test with threshold L_{thresh} selected by equation 8 (thus requiring no tuning). Training trajectories for various Lipschitz-estimation momentums α_L shown in figure 5; we see that the momentum has little effect on the training trajectory. In particular, we compare the training- and test-loss of the experiment with $\alpha_L = 0.7$ against the aforementioned standard AdamW baseline in figure 6; it is evident that there is little difference between the training trajectories.

In each of the above runs, we report the percent of parameters offloaded to SGD in table 2.

Caveat: our Lipschitz initialization scheme (initialization to 0) biases the Lipschitz parameters towards 0, with greater bias maintained with greater momentum (hence the reverse correlation of parameter offload with Lipschitz-momentum). We leave initialization schemes to future work.

1566
1567
1568
1569
1570
1571
1572
1573
1574
1575
1576
1577
1578



1579 Figure 6: Comparison of AMOS-combined SGD and Sophia vs. AdamW. We see little difference
1580 between the two despite offloading 16% of parameters to cheaper SGD optimizer.
1581

1582
1583

1584 The settings for this experiment are identical to those above, with the following exceptions:

1585
1586
1587
1588
1589
1590
1591
1592

- SGD’s learning rate reduced to $3.125 \cdot 10^{-6}$
- $\epsilon_2 = 10^{-15}$
- Each tensor’s Lipschitz parameter is chosen with maximum aggregation across the tensor instead of mean aggregation

1593
1594
1595
1596

L LIPSCHITZ-AIDED OPTIMIZER SELECTION

1597
1598
1599
1600
1601
1602
1603
1604

In this section, we demonstrate the use of convex Lipschitz parameters to select the best optimizer for our use case. Due to the relative constancy of Lipschitz parameters throughout the training process (after an initial warmup phase) in different settings, we can select optimizers for each setting based on the following rule: **quasi-Newton optimizers hold an advantage over first-order optimizers when the convex Lipschitz parameters are small**. As discussed in section 5, the convex Lipschitz parameters in the image autoencoder training setting are far larger than those in the image classification setting, so we compare a quasi-Newton optimizer against first-order methods in these settings to validate our rule.

1605
1606
1607
1608
1609
1610
1611
1612
1613
1614
1615
1616

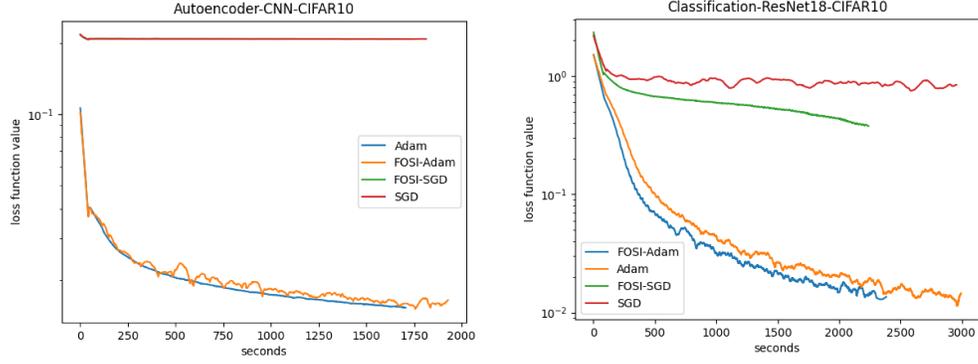
FOSI Sivan et al. (2024) is a variant of Saddle-Free Newton Dauphin et al. (2014) which applies Newton iterations in the domain space subspaces spanned by the dominant eigenvectors of the Hessian, and a first-order "base optimizer" in the remaining subspaces. We use FOSI as our representative second-order optimizer due to its computational effectiveness, capability to adjust the computational complexity of each iteration by adjusting the number of "dominant" eigenvectors to compute (fewer eigenvectors comes at the cost of a poorer Hessian approximation by approximating the Hessian with a lower-rank matrix, although this is somewhat mitigated by applying the base optimizer in these subspaces), and fairness of comparison (since its integration of first-order optimizers allows us to compare the effect of second-order optimization in the dominant eigenspaces against first-order optimization in these spaces, while all else is held equal - the remaining subspaces are both treated by the same first-order optimizers).

Experiment results may be seen in figure 7.

1617
1618
1619

The experiments are run with the same settings as before, with FOSI augmenting SGD and Adam respectively and Savitzky-Golay order-2 filtering with a window size of 5000 for clarity of visualization. It may be clearly seen that FOSI second-order augmentation is beneficial only in the classification setting, due to the small convex Lipschitz parameters.

1620
1621
1622
1623
1624
1625
1626
1627
1628
1629
1630
1631
1632



(a) The setting in which we train a CNN on an au- (b) The setting in which we train ResNet18 on a clas-
toencoder task has large convex Lipschitz parameters sification task has small convex Lipschitz parameters
throughout training throughout training

1633 Figure 7: Comparison of second-order optimizers against first-order optimizers in settings with
1634 different sized convex Lipschitz parameters. Second-order optimizers only hold an advantage over
1635 first-order optimizers (thus justifying their additional computational complexity) when the convex
1636 Lipschitz parameters are small.
1637
1638
1639
1640

1641 M PROOFS

1642 M.1 PRELIMINARY LEMMAS

1643 Before we can get started, we prove a few basic lemmas.

1644 **Lemma M.1.**

$$1645 \forall_{x \geq -1} : \sqrt{1+x} \leq 1 + \frac{x}{2}$$

1646 *Proof.* Mark $g : \mathbb{R} \rightarrow \mathbb{R}, g(x) = 1 + \frac{x}{2} - \sqrt{1+x}$. We have

$$1647 g'(x) = \frac{1}{2} \left(1 - \frac{1}{\sqrt{1+x}} \right)$$

$$1648 g''(x) = \frac{1}{4} \left(\frac{1}{(1+x)^{\frac{3}{2}}} \right)$$

1649 g is convex due to its second derivative being positive for all $x > -1$. Therefore, its sole critical
1650 point $x = 0$ obtained from the derivative is a minimum, and $\forall_{x \geq -1} : g(x) \geq g(0) = 0$ \square

1651 **Corollary M.1.1.**

$$1652 \forall_{x \in \mathbb{R}^+ \forall_{y \geq -x} : \sqrt{x+y} \leq \sqrt{x} + \frac{y}{2\sqrt{x}}$$

1653 *Proof.*

$$1654 \sqrt{x+y} = \sqrt{x} \sqrt{1 + \frac{y}{x}} \leq \sqrt{x} \left(1 + \frac{y}{2x} \right) = \sqrt{x} + \frac{y}{2\sqrt{x}}$$

1655 \square

1656 **Lemma.** Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ satisfy assumptions 1 and 2. Then

$$1657 m_t^i (\Delta \theta_t^{*\top} v_i) \leq M_t^i (\Delta \theta_t^{*\top} v_i) \leq 0 \quad (14)$$

1658 *Proof.* The first inequality stems from the trivial fact that $m_t^i \leq M_t^i$.

1659 The second inequality follows from the fact that (by design), $\Delta \theta_t^{*\top} v_i$ is a minimizer of $M_t^i (\Delta \theta_t^\top v_i)$,
1660 but

$$1661 M_t^i(0) = 0$$

1662
1663
1664
1665
1666
1667
1668
1669
1670
1671
1672
1673

1674

1675

1676

Lemma. 3.2 Minmax preconditioner

1677

1678

1679

1680

1681

$$\arg \min_{\Phi_t \in \mathbb{R}^{n \times n}} |\Delta \Delta^i \theta_t(\Phi_t)| = \left(\frac{\mathcal{H}(\theta_t) + \sqrt{(\mathcal{H}(\theta_t))^2 + 2V \cdot \text{diag} \left(L_t^i \cdot |\nabla f(\theta_t)^\top v_i| \right)_{i=1}^n \cdot V^\top}}{2} \right)^{-1}$$

1682

Proof.

1683

1684

1685

1686

1687

1688

1689

1690

1691

$$\begin{aligned} |\Delta \Delta^i \theta_t| &= \left| \nabla f(\theta_t)^\top \left(\Phi_t - \frac{2}{\lambda_i + \sqrt{\lambda_i^2 - 2L_t^i \cdot \nabla f(\theta_t)^\top v_i}} I \right) v_i \right| \\ &= \left| \nabla f(\theta_t)^\top \left(\Phi_t - \left(\frac{\mathcal{H}(\theta_t) + \sqrt{(\mathcal{H}(\theta_t))^2 + 2L_t^i \cdot |\nabla f(\theta_t)^\top v_i|} I}{2} \right)^{-1} \right) v_i \right| \end{aligned}$$

1692

and the result follows from developing the second parenthesized term for all n dimensions of the domain space. \square

1693

1694

1695

M.2 LEMMA 3.1: EIGENSPACE DESCENT

1696

1697

1698

Working with assumptions 1 and equation 2, Dennis & Schnabel (1983, Lemma 4.1.14) prove the following:

1699

Lemma M.2.

1700

1701

$$f(\theta_{t+1}) - f(\theta_t) \leq \nabla f(\theta_t)^\top \cdot \Delta \theta_t + \frac{1}{2} \Delta \theta_t^\top \mathcal{H}(\theta_t) \Delta \theta_t + \frac{1}{6} L_H \|\Delta \theta_t\|_2^3 \quad (15)$$

1702

1703

1704

1705

1706

1707

1708

1709

1710

1711

1712

1713

1714

1715

Lemma. 3.1 Eigenspace Descent Bounds

1716

1717

1718

Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be a function satisfying assumptions 1 and 2, and let $\theta_{t+1} \in \mathbb{R}^n$. Marking $\Delta \theta_t = \theta_{t+1} - \theta_t$, we have

1719

1720

1721

$$\exists_{(L_t^i)_{i=1}^n \in (\mathbb{R}^+)^n} : f(\theta_{t+1}) - f(\theta_t) \leq \sum_{i=1}^n M_t^i (\Delta \theta_t^\top v_i) \quad (16)$$

1722

1723

1724

1725

1726

1727

$$\exists_{(L_t^i)_{i=1}^n \in (\mathbb{R}^+)^n} : f(\theta_{t+1}) - f(\theta_t) \geq \sum_{i=1}^n m_t^i (\Delta \theta_t^\top v_i) \quad (17)$$

We give 2 proofs of the above lemma. The first proof is far simpler and relies on the standard spectral norm-Lipschitz continuous Hessian assumption given by equation 2 instead of assumption 2:

1728 *Proof.* Beginning with Nesterov & Polyak (2006, Lemma 1) for the first inequality,
 1729
 1730

$$\begin{aligned}
 & \left| f(\theta_{t+1}) - f(\theta_t) - \left(\nabla f(\theta_t)^\top (\theta_{t+1} - \theta_t) + (\theta_{t+1} - \theta_t)^\top \mathcal{H}(\theta_t) (\theta_{t+1} - \theta_t) \right) \right| \\
 & \leq L_H \|\theta_{t+1} - \theta_t\|_2^3 \\
 & \left| f(\theta_{t+1}) - f(\theta_t) - \left(\sum_{i=1}^n \left(\nabla f(\theta_t)^\top v_i \cdot (\theta_{t+1} - \theta_t)^\top v_i + \lambda_i \left((\theta_{t+1} - \theta_t)^\top v_i \right)^2 \right) \right) \right| \\
 & \leq L_H \left\| \sum_{i=1}^n (\theta_{t+1} - \theta_t)^\top v_i v_i^\top \right\|_2^3 \\
 & \leq L_H \left(\sum_{i=1}^n \left\| (\theta_{t+1} - \theta_t)^\top v_i v_i^\top \right\|_2 \right)^3 \\
 & = L_H \cdot n^3 \left(\sum_{i=1}^n \frac{1}{n} \left| (\theta_{t+1} - \theta_t)^\top v_i \right| \right)^3 \\
 & \leq L_H \cdot n^3 \left(\frac{1}{n} \sum_{i=1}^n \left| (\theta_{t+1} - \theta_t)^\top v_i \right|^3 \right) \\
 & = L_H \cdot n^2 \sum_{i=1}^n \left| (\theta_{t+1} - \theta_t)^\top v_i \right|^3 \\
 & = \sum_{i=1}^n \tilde{L}_H \cdot \left| (\theta_{t+1} - \theta_t)^\top v_i \right|^3
 \end{aligned}$$

1756 with

- 1757 1. the second inequality being a representation of $\theta_{t+1} - \theta_t$ over the (orthogonal) Hessian
 1758 eigenbasis
- 1759 2. the third inequality due to the triangle inequality
- 1760 3. the fourth inequality due to Jensen's inequality

1771 \square

1776 Our first proof of lemma 3.1 is simple, but leaves something to be desired due to its lack of per-
 1777 eigenspace Lipschitz parameters and due to the presence of n^2 in the bound, which can be very
 1778 large, as noted in section A. The first proof's assumption of equation 2 is also easily seen to be
 1779 no weaker than assumption 2 (meaning that assuming equation 2 implies assumption 2) by taking
 1780 $L_R \triangleq L_t^i \triangleq L_H$ for all $t \in [T], i \in [n]$. To address these concerns, we make use of assumption 2,
 1781 and give a second (though more complicated) proof for lemma 3.1. We'll prove only the upper bound,
 as a proof for the lower bound is similar.

1782 *Proof.*

$$\begin{aligned}
1783 & f(\theta_{t+1}) - f(\theta_t) \\
1784 & = \int_0^1 \nabla f(\theta_t + y(\theta_{t+1} - \theta_t))^\top (\theta_{t+1} - \theta_t) dy \\
1785 & = \nabla f(\theta_t)^\top (\theta_{t+1} - \theta_t) + \int_0^1 (\nabla f(\theta_t + y(\theta_{t+1} - \theta_t)) - \nabla f(\theta_t))^\top (\theta_{t+1} - \theta_t) dy \\
1786 & = \nabla f(\theta_t)^\top (\theta_{t+1} - \theta_t) + \int_0^1 \left(\int_0^1 y \mathcal{H}(\theta_t + yz(\theta_{t+1} - \theta_t)) (\theta_{t+1} - \theta_t) dz \right)^\top (\theta_{t+1} - \theta_t) dy \\
1787 & = \nabla f(\theta_t)^\top (\theta_{t+1} - \theta_t) + \int_0^1 \int_0^1 y (\theta_{t+1} - \theta_t)^\top \mathcal{H}(\theta_t + yz(\theta_{t+1} - \theta_t)) (\theta_{t+1} - \theta_t) dydz \\
1788 & = \nabla f(\theta_t)^\top (\theta_{t+1} - \theta_t) + (\theta_{t+1} - \theta_t)^\top \mathcal{H}(\theta_t) (\theta_{t+1} - \theta_t) \\
1789 & + \underbrace{\int_0^1 \int_0^1 y (\theta_{t+1} - \theta_t)^\top (\mathcal{H}(\theta_t + yz(\theta_{t+1} - \theta_t)) - \mathcal{H}(\theta_t)) (\theta_{t+1} - \theta_t) dydz}_Z \\
1790 &
\end{aligned}$$

1801 with the first and third equalities due to the fundamental theorem of calculus.

1802 Mark the Hessian eigendecompositions as follows:

$$\begin{aligned}
1803 & \mathcal{H}(\theta_t + yz(\theta_{t+1} - \theta_t)) = \tilde{V} \tilde{\Lambda} \tilde{V}^\top \\
1804 & \mathcal{H}(\theta_t) = V \Lambda V^\top
\end{aligned}$$

1805 with diagonal $\Lambda = \text{diag}(\lambda_i)_{i=1}^n$, $\tilde{\Lambda} = \text{diag}(\tilde{\lambda}_i)_{i=1}^n$ and orthogonal (due to the Hermitian nature of
1806 Hessian matrices) matrices V, \tilde{V} .

$$\begin{aligned}
1807 & Z = \int_0^1 \int_0^1 y (\theta_{t+1} - \theta_t)^\top (\tilde{V} \tilde{\Lambda} \tilde{V}^\top - V \Lambda V^\top) (\theta_{t+1} - \theta_t) dydz \\
1808 & = \int_0^1 \int_0^1 y (\theta_{t+1} - \theta_t)^\top (V \tilde{\Lambda} V^\top - V \Lambda V^\top) (\theta_{t+1} - \theta_t) dydz \\
1809 & + \int_0^1 \int_0^1 y (\theta_{t+1} - \theta_t)^\top (\tilde{V} \tilde{\Lambda} \tilde{V}^\top - V \tilde{\Lambda} V^\top) (\theta_{t+1} - \theta_t) dydz \\
1810 &
\end{aligned}$$

1811 Focusing on the first term,

$$\begin{aligned}
1812 & = \int_0^1 \int_0^1 y (\theta_{t+1} - \theta_t)^\top V (\tilde{\Lambda} - \Lambda) V^\top (\theta_{t+1} - \theta_t) dydz \\
1813 & = \int_0^1 \int_0^1 y \sum_{j=1}^n \sum_{i=1}^n (\theta_{t+1} - \theta_t)^\top v_i \cdot (\theta_{t+1} - \theta_t)^\top v_j \cdot v_j^\top V (\tilde{\Lambda} - \Lambda) V^\top v_i dydz \\
1814 & = \int_0^1 \int_0^1 y \cdot \sum_{i=1}^n \left((\theta_{t+1} - \theta_t)^\top v_i \right)^2 \cdot (\tilde{\lambda}_i - \lambda_i) dydz \\
1815 & \leq \sum_{i=1}^n L_t^i \cdot \left| (\theta_{t+1} - \theta_t)^\top v_i \right|^3 \cdot \int_0^1 \int_0^1 y \cdot yz \cdot dydz \\
1816 & = \sum_{i=1}^n \frac{L_t^i}{6} \cdot \left| (\theta_{t+1} - \theta_t)^\top v_i \right|^3 \\
1817 &
\end{aligned}$$

As for the second term,

$$\begin{aligned}
& \int_0^1 \int_0^1 y (\theta_{t+1} - \theta_t)^\top \left(\tilde{V} \tilde{\Lambda} \tilde{V}^\top - V \tilde{\Lambda} V^\top \right) (\theta_{t+1} - \theta_t) dy dz \\
&= \int_0^1 \int_0^1 y (\theta_{t+1} - \theta_t)^\top \left(\tilde{V} - V \right) \tilde{\Lambda} \left(\tilde{V} + V \right)^\top (\theta_{t+1} - \theta_t) dy dz \\
&+ \int_0^1 \int_0^1 y (\theta_{t+1} - \theta_t)^\top \left(\left(\tilde{V} \tilde{\Lambda} V^\top \right)^\top - \tilde{V} \tilde{\Lambda} V^\top \right) (\theta_{t+1} - \theta_t) dy dz \\
&= \int_0^1 \int_0^1 y (\theta_{t+1} - \theta_t)^\top \left(\tilde{V} - V \right) \tilde{\Lambda} \left(\tilde{V} + V \right)^\top (\theta_{t+1} - \theta_t) dy dz \\
&\leq \int_0^1 \int_0^1 y \cdot \|\tilde{V} - V\|_2 \cdot \|\tilde{\Lambda}\|_2 \cdot \|\tilde{V} + V\|_2 \cdot \|\theta_{t+1} - \theta_t\|_2^2 \cdot dy dz \\
&\leq \frac{1}{3} \cdot L_R \|\tilde{\Lambda}\|_2 \cdot \|\theta_{t+1} - \theta_t\|_2^3 \\
&= \frac{1}{3} \cdot L_R \|\tilde{\Lambda}\|_2 \cdot \|(\theta_{t+1} - \theta_t) V^\top\|_2^3 \\
&\leq \frac{\sqrt{n}}{3} \cdot L_R \|\tilde{\Lambda}\|_2 \cdot \|(\theta_{t+1} - \theta_t) V^\top\|_3^3 \\
&= \frac{\sqrt{n}}{3} \cdot L_R \|\tilde{\Lambda}\|_2 \cdot \sum_{i=1}^n |(\theta_{t+1} - \theta_t) V^\top \cdot e_i|^3 \\
&= \sum_{i=1}^n \frac{\sqrt{n}}{3} \cdot L_R \|\tilde{\Lambda}\|_2 \cdot |(\theta_{t+1} - \theta_t) v_i|^3
\end{aligned}$$

with

- the first 4 transfers similar to those in the proof of lemma M.4
- the third equality due to orthonormality
- the third inequality due to the L_p norms inequality
- e_i indicating the 1-hot vector with a 1 in the i -th entry

Putting it all together (and representing $\theta_{t+1} - \theta_t$ by its coordinate vector over the eigenbasis of $\mathcal{H}(\theta_t)$):

$$\begin{aligned}
f(\theta_{t+1}) - f(\theta_t) &\leq \sum_{i=1}^n \nabla f(\theta_t)^\top v_i \cdot (\theta_{t+1} - \theta_t)^\top v_i + \lambda_i(\theta_t) \cdot \left((\theta_{t+1} - \theta_t)^\top v_i \right)^2 \\
&\quad + \left(\frac{L_t^i}{6} + \frac{\sqrt{n}}{3} \cdot L_R \cdot L_H \right) \cdot |(\theta_{t+1} - \theta_t) v_i|^3
\end{aligned}$$

□

To understand the relationship between our assumption 2 and the more standard equation 2, we further prove that the combination of assumption 2 and a bounded spectrum assumption will be no weaker than equation 2:

Lemma M.3. *Let $A \in \mathbb{R}^{n \times n}$, $v \in \mathbb{R}^n$. Then $v^\top (A^\top - A) v = 0$.*

Proof.

$$v^\top (A^\top - A) v = (v^\top (A^\top - A) v)^\top = v^\top (A - A^\top) v = -v^\top (A^\top - A) v$$

□

Theorem M.4. Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be a function satisfying assumptions 1 and 2, and assume $\exists \lambda_{\text{sup}} \in \mathbb{R} \forall \theta \in \mathbb{R}^n \forall i \in [n] : \lambda_i(\theta) \leq \lambda_{\text{sup}}$. Then equation 2 is satisfied.

Proof.

$$\begin{aligned}
|p^\top (\mathcal{H}(\theta) - \mathcal{H}(\varphi)) p| &= \left| p^\top (V \Lambda V^\top - \tilde{V} \tilde{\Lambda} \tilde{V}^\top) p \right| \\
&\leq \left| p^\top (V \Lambda V^\top - \tilde{V} \Lambda \tilde{V}^\top) p \right| + \left| p^\top \tilde{V} (\Lambda - \tilde{\Lambda}) \tilde{V}^\top p \right| \\
&= \left| p^\top (V - \tilde{V}) \Lambda (V + \tilde{V})^\top p \right| + \left| p^\top \tilde{V} (\Lambda - \tilde{\Lambda}) \tilde{V}^\top p \right| \\
&\leq \|V - \tilde{V}\|_2 \cdot \|\Lambda\|_2 \cdot \|V + \tilde{V}\|_2 \cdot \|p\|_2^2 + \|\tilde{V}^\top p\|_2^2 \cdot \|\Lambda - \tilde{\Lambda}\|_2 \\
&\leq \left(2L_R \cdot \sup_{\theta' \in \mathbb{R}^n, i \in \mathbb{R}} \lambda_i(\theta') \right) \cdot \|\theta - \varphi\|_2 + \max_i L^i \cdot \|\theta - \varphi\|_2 \\
&\leq (2L_H \cdot \lambda_{\text{sup}} + L_H) \cdot \|\theta - \varphi\|_2
\end{aligned}$$

with

- the first inequality due to the triangle inequality
- the second equality due to lemma M.3
- the second inequality due to the Cauchy-Schwartz inequality
- the third inequality due to the triangle inequality, and the fact that all of an orthonormal matrix's eigenvalues equal one of $\{-1, 1\}$

□

M.3 THEOREM C.1: WORST CASE-OPTIMAL DESCENT RATE

Before we can prove theorem C.1, we need to upper bound equation 7.

Lemma M.5. *Minmax stepsize bound*

If $\lambda_i \geq 0$ then

$$\Delta \theta_t^{*\top} v_i = \mathcal{O} \left(\sqrt{|\nabla f(\theta_t)^\top v_i|} \right)$$

If $\lambda_i < 0$ then

$$\Delta \theta_t^{*\top} v_i = \mathcal{O}(|\lambda_i|)$$

1944 *Proof.* For i s.t. $0 \leq \lambda_i \leq \sqrt{L_t^i \cdot |\nabla f(\theta_t)^\top v_i|}$ we use corollary M.1.1 with $x = 2L_t^i \cdot |\nabla f(\theta_t)^\top v_i|$
 1945 to obtain
 1946

$$\begin{aligned}
 1947 \Delta \theta_t^{*\top} v_i &= \frac{\sqrt{\lambda_i^2 + 2L_t^i \cdot |\nabla f(\theta_t)^\top v_i|} - \lambda_i}{L_t^i} \\
 1948 &\leq \sqrt{2} \cdot \sqrt{\frac{|\nabla f(\theta_t)^\top v_i|}{L_t^i}} + \frac{\frac{\lambda_i^2}{2} \cdot \frac{1}{\sqrt{2L_t^i \cdot |\nabla f(\theta_t)^\top v_i|}} - \lambda_i}{L_t^i} \\
 1949 &\leq \sqrt{2} \cdot \sqrt{\frac{|\nabla f(\theta_t)^\top v_i|}{L_t^i}} + \frac{1}{2\sqrt{2}} \cdot \sqrt{\frac{|\nabla f(\theta_t)^\top v_i|}{L_t^i}} \\
 1950 &= \frac{5}{2\sqrt{2}} \cdot \sqrt{\frac{|\nabla f(\theta_t)^\top v_i|}{L_t^i}}
 \end{aligned}$$

1951 For i s.t. $\lambda_i > \sqrt{L_t^i \cdot |\nabla f(\theta_t)^\top v_i|}$ we use corollary M.1.1 with $x = \lambda_i^2$ to obtain
 1952

$$\begin{aligned}
 1953 \Delta \theta_t^{*\top} v_i &\leq \frac{|\nabla f(\theta_t)^\top v_i|}{\lambda_i} \leq \frac{|\nabla f(\theta_t)^\top v_i|}{\sqrt{L_t^i \cdot |\nabla f(\theta_t)^\top v_i|}} = \sqrt{\frac{|\nabla f(\theta_t)^\top v_i|}{L_t^i}} \\
 1954 &
 \end{aligned}$$

1955 For i s.t. $\lambda_i < 0$, we again use corollary M.1.1 with $x = \lambda_i^2$ to obtain
 1956

$$\begin{aligned}
 1957 \Delta \theta_t^{*\top} v_i &= \frac{2L_t^i \cdot |\nabla f(\theta_t)^\top v_i|}{L_t^i \left(\sqrt{\lambda_i^2 + 2L_t^i \cdot |\nabla f(\theta_t)^\top v_i|} - |\lambda_i| \right)} \leq \frac{2|\lambda_i|}{L_t^i} = \mathcal{O}(|\lambda_i|) \\
 1958 &
 \end{aligned}$$

□

1959 We are now ready to prove theorem C.1.

1960 **Theorem.** *Worst case-optimal descent rate* Let f be a function with Lipschitz-continuous Hessian.
 1961 After t iterations, algorithm ELMO satisfies

$$\begin{aligned}
 1962 f(\theta_0) - f(\theta_t) &= \mathcal{O}(\log t) \tag{18} \\
 1963 &
 \end{aligned}$$

1964 *Proof.* Cartis et al. (2012a) give $|\nabla f(\theta_t)^\top v_i| = \mathcal{O}\left(\frac{1}{t^{\frac{2}{3}}}\right)$ and $\forall_{i:\lambda_i < 0} : |\lambda_i| = \mathcal{O}\left(\frac{1}{\sqrt[3]{t}}\right)$ for the
 1965 ARC optimization algorithm, of which algorithm ELMO is a special case (the case where ARC
 1966 perfectly estimates the Hessian Lipschitz parameter).
 1967

1968 Making use of lemma M.5 and noting that $m_t^i (\Delta \theta_t^{*\top} v_i) \leq 0$ by equation 14:

$$\begin{aligned}
 1969 |m_t^i (\Delta \theta_t^{*\top} v_i)| &= |\nabla f(\theta_t)^\top v_i| \cdot |\Delta \theta_t^{*\top} v_i| + \frac{-\lambda_i}{2} \cdot (\Delta \theta_t^{*\top} v_i)^2 + \frac{L_t^i}{6} \cdot (\Delta \theta_t^{*\top} v_i)^3 \\
 1970 &
 \end{aligned}$$

1971 For i s.t. $\lambda_i \geq 0$:

$$\begin{aligned}
 1972 &\leq |\nabla f(\theta_t)^\top v_i| \cdot \mathcal{O}\left(\sqrt{\frac{|\nabla f(\theta_t)^\top v_i|}{L_t^i}}\right) + \mathcal{O}\left(\sqrt{\frac{|\nabla f(\theta_t)^\top v_i|}{L_t^i}}\right)^3 \\
 1973 &= \mathcal{O}\left(|\nabla f(\theta_t)^\top v_i|^{1.5}\right) = \mathcal{O}\left(\left(\frac{1}{t^{\frac{2}{3}}}\right)^{1.5}\right) = \mathcal{O}\left(\frac{1}{t}\right) \\
 1974 &
 \end{aligned}$$

1975

For i s.t. $\lambda_i < 0$:

$$\begin{aligned} &\leq \left| \nabla f(\theta_t)^\top v_i \right| \cdot \mathcal{O}(|\lambda_i|) + |\lambda_i| \cdot (\mathcal{O}(|\lambda_i|))^2 + \mathcal{O}(|\lambda_i|)^3 = \mathcal{O}\left(\left| \nabla f(\theta_t)^\top v_i \right| \cdot |\lambda_i| + |\lambda_i|^3\right) \\ &= \mathcal{O}\left(\frac{1}{t^{\frac{2}{3}}} \cdot \frac{1}{\sqrt[3]{t}} + \frac{1}{t}\right) = \mathcal{O}\left(\frac{1}{t}\right) \end{aligned}$$

Finally, we have

$$\begin{aligned} f(\theta_0) - f(\theta_T) &= \sum_{t=0}^{T-1} f(\theta_t) - f(\theta_{t+1}) \\ &\leq \sum_{t=0}^{T-1} |m_t^i(\Delta\theta_t^{*\top} v_i)| \\ &\leq \sum_{t=1}^{T-1} \mathcal{O}\left(\frac{1}{t}\right) \\ &\leq \mathcal{O}\left(\log t + \gamma + \frac{1}{2t}\right) \\ &= \mathcal{O}(\log t) \end{aligned}$$

with $\gamma \approx 0.57721$ as the Euler-Mascheroni constant and Young (1991) for the last inequality. \square

M.4 THEOREM C.2

Theorem.

$$|m_t^i(\Delta\theta_t^{*\top} v_i)| \leq 5 |M_t^i(\Delta\theta_t^{*\top} v_i)|$$

Proof. Due to equation 14, we have $\frac{m_t^i(\Delta\theta_t^{*\top} v_i)}{M_t^i(\Delta\theta_t^{*\top} v_i)} = \left| \frac{m_t^i(\Delta\theta_t^{*\top} v_i)}{M_t^i(\Delta\theta_t^{*\top} v_i)} \right|$. Now:

$$\begin{aligned} &\frac{m_t^i(\Delta\theta_t^{*\top} v_i)}{M_t^i(\Delta\theta_t^{*\top} v_i)} \\ &= \frac{\nabla f(\theta_t)^\top v_i \cdot \Delta\theta_t^{*\top} v_i + \frac{\lambda_i}{2} (\Delta\theta_t^{*\top} v_i)^2 - \frac{L_t^i}{6} \cdot |\Delta\theta_t^{*\top} v_i|^3}{\nabla f(\theta_t)^\top v_i \cdot \Delta\theta_t^{*\top} v_i + \frac{\lambda_i}{2} (\Delta\theta_t^{*\top} v_i)^2 + \frac{L_t^i}{6} \cdot |\Delta\theta_t^{*\top} v_i|^3} \\ &= \frac{-\left| \nabla f(\theta_t)^\top v_i \right| + \frac{\lambda_i}{2} \frac{\sqrt{\lambda_i^2 + 2L_t^i |\nabla f(\theta_t)^\top v_i|} - \lambda_i}{L_t^i} - \frac{L_t^i}{6} \cdot \left(\frac{\sqrt{\lambda_i^2 + 2L_t^i |\nabla f(\theta_t)^\top v_i|} - \lambda_i}{L_t^i} \right)^2}{-\left| \nabla f(\theta_t)^\top v_i \right| + \frac{\lambda_i}{2} \frac{\sqrt{\lambda_i^2 + 2L_t^i |\nabla f(\theta_t)^\top v_i|} - \lambda_i}{L_t^i} + \frac{L_t^i}{6} \cdot \left(\frac{\sqrt{\lambda_i^2 + 2L_t^i |\nabla f(\theta_t)^\top v_i|} - \lambda_i}{L_t^i} \right)^2} \\ &= \frac{5 \left(\lambda_i \sqrt{\lambda_i^2 + 2L_t^i |\nabla f(\theta_t)^\top v_i|} - \lambda_i^2 \right) - 8L_t^i |\nabla f(\theta_t)^\top v_i|}{\left(\lambda_i \sqrt{\lambda_i^2 + 2L_t^i |\nabla f(\theta_t)^\top v_i|} - \lambda_i^2 \right) - 4L_t^i |\nabla f(\theta_t)^\top v_i|} \end{aligned}$$

If $\lambda_i = 0$, then

$$\frac{m_t^i(\Delta\theta_t^{*\top} v_i)}{M_t^i(\Delta\theta_t^{*\top} v_i)} = 2$$

2052 If $\lambda_i > 0$, then

$$2053$$

$$2054$$

$$2055$$

$$2056$$

$$2057 \frac{m_t^i(\Delta\theta_t^{*\top} v_i)}{M_t^i(\Delta\theta_t^{*\top} v_i)} = \frac{5 \sqrt{1+2 \frac{L_t^i |\nabla f(\theta_t)^\top v_i|}{\lambda_i^2}} - 1}{\frac{L_t^i |\nabla f(\theta_t)^\top v_i|}{\lambda_i^2}} - 8 \leq \lim_{x \rightarrow \infty} \frac{5 \frac{\sqrt{1+2x-1}}{x} - 8}{\frac{\sqrt{1+2x-1}}{x} - 4} = 2$$

$$2058$$

$$2059$$

$$2060$$

$$2061$$

$$2062$$

$$2063$$

$$2064$$

2065 due to the monotonic increasing nature of $\psi_5 : \mathbb{R}^+ \rightarrow \mathbb{R}$, $\psi_5(x) = \frac{5 \frac{\sqrt{1+2x-1}}{x} - 8}{\frac{\sqrt{1+2x-1}}{x} - 4}$.

2067 If $\lambda_i < 0$, then

$$2068$$

$$2069$$

$$2070$$

$$2071$$

$$2072 \frac{m_t^i(\Delta\theta_t^{*\top} v_i)}{M_t^i(\Delta\theta_t^{*\top} v_i)} = \frac{5 \sqrt{1+2 \frac{L_t^i |\nabla f(\theta_t)^\top v_i|}{|\lambda_i|^2}} + 1}{\frac{L_t^i |\nabla f(\theta_t)^\top v_i|}{|\lambda_i|^2}} + 8 \leq \lim_{x \rightarrow 0^+} \frac{5 \frac{\sqrt{1+2x+1}}{x} + 8}{\frac{\sqrt{1+2x+1}}{x} + 4} = 5$$

$$2073$$

$$2074$$

$$2075$$

$$2076$$

$$2077$$

$$2078$$

$$2079$$

$$2080$$

2081 due to the monotonic decreasing nature of $\psi_6 : \mathbb{R}^+ \rightarrow \mathbb{R}$, $\psi_6(x) = \frac{5 \frac{\sqrt{1+2x+1}}{x} + 8}{\frac{\sqrt{1+2x+1}}{x} + 4}$. □

2086 M.5 THEOREM 3.3: WORST-CASE DESCENT RATE FOR ARBITRARY OPTIMIZERS

2087

2088 **Theorem.** *Relative Descent Rate*

2089 •

$$2090$$

$$2091$$

$$2092$$

$$2093$$

$$2094$$

$$2095 \left| \frac{M_t^i(\Delta\theta_t^\top v_i) - M_t^i(\Delta\theta_t^{*\top} v_i)}{M_t^i(\Delta\theta_t^{*\top} v_i)} \right| = \Theta(|\Delta\Delta^i \theta_t'|^2)$$

$$2096$$

$$2097$$

$$2098$$

2099 •

$$2100$$

$$2101$$

$$2102$$

$$2103$$

$$2104 \left| \frac{m_t^i(\Delta\theta_t^\top v_i) - m_t^i(\Delta\theta_t^{*\top} v_i)}{m_t^i(\Delta\theta_t^{*\top} v_i)} \right| = \Theta(|\Delta\Delta^i \theta_t'|)$$

$$2105 \tag{19}$$

Proof. For the first part of the lemma,

$$\begin{aligned}
& \left| \frac{M_t^i (\Delta \theta_t^\top v_i) - M_t^i (\Delta \theta_t^{*\top} v_i)}{M_t^i (\Delta \theta_t^{*\top} v_i)} \right| \\
&= \frac{\nabla f(\theta_t)^\top \cdot v_i \cdot \left((\theta_{t+1} - \theta_t)^\top v_i - \Delta \theta_t^{*\top} v_i \right)}{\Delta \theta_t^{*\top} v_i \cdot \nabla f(\theta_t)^\top v_i + \frac{1}{2} \lambda_i (\Delta \theta_t^{*\top} v_i)^2 + \frac{L_t^i}{6} (\Delta \theta_t^{*\top} v_i)^3} \\
&+ \frac{\frac{1}{2} \lambda_i \cdot \left(\left((\theta_{t+1} - \theta_t)^\top v_i \right)^2 - (\Delta \theta_t^{*\top} v_i)^2 \right)}{\Delta \theta_t^{*\top} v_i \cdot \nabla f(\theta_t)^\top v_i + \frac{1}{2} \lambda_i (\Delta \theta_t^{*\top} v_i)^2 + \frac{L_t^i}{6} (\Delta \theta_t^{*\top} v_i)^3} \\
&+ \frac{\frac{L_t^i}{6} \cdot \left(\left((\theta_{t+1} - \theta_t)^\top \cdot v_i \right)^3 - (\Delta \theta_t^{*\top} v_i)^3 \right)}{\Delta \theta_t^{*\top} v_i \cdot \nabla f(\theta_t)^\top v_i + \frac{1}{2} \lambda_i (\Delta \theta_t^{*\top} v_i)^2 + \frac{L_t^i}{6} (\Delta \theta_t^{*\top} v_i)^3} \\
&= \Delta \Delta^i \theta_t' \left(\frac{\frac{L_t^i}{6} \cdot (\Delta \Delta^i \theta_t'^2 + 2 \Delta \Delta^i \theta_t' + 3) \cdot (\Delta \theta_t^{*\top} v_i)^2}{\frac{L_t^i}{6} (\Delta \theta_t^{*\top} v_i)^2 + \frac{1}{2} \lambda_i \Delta \theta_t^{*\top} v_i - |\nabla f(\theta_t)^\top v_i|} \right. \\
&+ \frac{\lambda_i \cdot \left(\frac{1}{2} \Delta \Delta^i \theta_t' + 1 \right) \cdot \Delta \theta_t^{*\top} v_i}{\frac{L_t^i}{6} (\Delta \theta_t^{*\top} v_i)^2 + \frac{1}{2} \lambda_i \Delta \theta_t^{*\top} v_i - |\nabla f(\theta_t)^\top v_i|} \\
&\left. - \frac{|\nabla f(\theta_t)^\top \cdot v_i|}{\frac{L_t^i}{6} (\Delta \theta_t^{*\top} v_i)^2 + \frac{1}{2} \lambda_i \Delta \theta_t^{*\top} v_i - |\nabla f(\theta_t)^\top v_i|} \right) \\
&= \Delta \Delta^i \theta_t' \left(\frac{\frac{1}{6} \cdot (\Delta \Delta^i \theta_t'^2 + 2 \Delta \Delta^i \theta_t' + 3) \cdot \left(\sqrt{\lambda_i^2 + 2L_t^i \cdot |\nabla f(\theta_t)^\top v_i|} - \lambda_i \right)^2}{L_t^i} \right. \\
&\left. \frac{\left(\sqrt{\lambda_i^2 + 2L_t^i \cdot |\nabla f(\theta_t)^\top v_i|} - \lambda_i \right)^2}{\frac{1}{6} \frac{\left(\sqrt{\lambda_i^2 + 2L_t^i \cdot |\nabla f(\theta_t)^\top v_i|} - \lambda_i \right)^2}{L_t^i} + \frac{1}{2} \lambda_i \frac{\sqrt{\lambda_i^2 + 2L_t^i \cdot |\nabla f(\theta_t)^\top v_i|} - \lambda_i}{L_t^i} - |\nabla f(\theta_t)^\top v_i|} \right. \\
&+ \frac{\lambda_i \cdot \left(\frac{1}{2} \Delta \Delta^i \theta_t' + 1 \right) \cdot \left(\frac{\sqrt{\lambda_i^2 + 2L_t^i \cdot |\nabla f(\theta_t)^\top v_i|} - \lambda_i}{L_t^i} \right)}{\frac{1}{6} \frac{\left(\sqrt{\lambda_i^2 + 2L_t^i \cdot |\nabla f(\theta_t)^\top v_i|} - \lambda_i \right)^2}{L_t^i} + \frac{1}{2} \lambda_i \frac{\sqrt{\lambda_i^2 + 2L_t^i \cdot |\nabla f(\theta_t)^\top v_i|} - \lambda_i}{L_t^i} - |\nabla f(\theta_t)^\top v_i|} \\
&\left. - \frac{|\nabla f(\theta_t)^\top \cdot v_i|}{\frac{1}{6} \frac{\left(\sqrt{\lambda_i^2 + 2L_t^i \cdot |\nabla f(\theta_t)^\top v_i|} - \lambda_i \right)^2}{L_t^i} + \frac{1}{2} \lambda_i \frac{\sqrt{\lambda_i^2 + 2L_t^i \cdot |\nabla f(\theta_t)^\top v_i|} - \lambda_i}{L_t^i} - |\nabla f(\theta_t)^\top v_i|} \right) \\
&= \Delta \Delta^i \theta_t' \cdot \left(\frac{\Delta \Delta^i \theta_t' + 2 \Delta \Delta^i \theta_t'^2}{\lambda_i \sqrt{\lambda_i^2 + 2L_t^i \cdot |\nabla f(\theta_t)^\top v_i|} - \lambda_i^2 - 4L_t^i |\nabla f(\theta_t)^\top v_i|} \cdot \lambda_i^2 \right. \\
&+ \frac{2 (\Delta \Delta^i \theta_t'^2 + 2 \Delta \Delta^i \theta_t')}{\lambda_i \sqrt{\lambda_i^2 + 2L_t^i \cdot |\nabla f(\theta_t)^\top v_i|} - \lambda_i^2 - 4L_t^i |\nabla f(\theta_t)^\top v_i|} \cdot L_t^i \cdot |\nabla f(\theta_t)^\top v_i| \\
&\left. - \frac{\Delta \Delta^i \theta_t' + 2 \Delta \Delta^i \theta_t'^2}{\lambda_i \sqrt{\lambda_i^2 + 2L_t^i \cdot |\nabla f(\theta_t)^\top v_i|} - \lambda_i^2 - 4L_t^i |\nabla f(\theta_t)^\top v_i|} \cdot \lambda_i \sqrt{\lambda_i^2 + 2L_t^i \cdot |\nabla f(\theta_t)^\top v_i|} \right)
\end{aligned}$$

2160 If $\lambda_i = 0$:

2161

2162

2163

2164

$$= -\Delta\Delta^i\theta_t'^2 \cdot \left(1 + \frac{1}{2}\Delta\Delta^i\theta_t'^2\right)$$

2165

2166

If $\lambda_i > 0$:

2167

2168

2169

2170

2171

2172

2173

2174

2175

2176

2177

2178

2179

2180

2181

2182

2183

2184

2185

2186

2187

2188

2189

2190

2191

2192

2193

2194

2195

2196

2197

2198

Proving that

2199

2200

2201

2202

2203

2204

would conclude the proof for this case. This is easily proven, by noting that

2205

2206

2207

2208

is monotonic and satisfies

2209

2210

2211

2212

2213

$$\lim_{x \rightarrow 0^+} \psi_1(x) = \frac{1}{3}$$

$$\lim_{x \rightarrow \infty} \psi_1(x) = 0$$

2214 If, on the other hand, $\lambda_i < 0$:

2215

2216

2217

2218

2219

2220

2221

2222

2223

2224

2225

2226

2227

2228

2229

2230

2231

2232 Proving that

2233

2234

2235

2236

2237

2238

2239

2240

2241

2242

2243

2244

would conclude the proof for this case as well. This is easily proven, by noting that

2245

2246

2247

2248

2249

2250

2251

2252

2253

2254

is monotonic and satisfies

2255

2256

2257

2258

2259

2260

2261

2262

2263

2264

2265

2266

2267

$$\begin{aligned}
&= -\Delta \Delta^i \theta_t'^2 \cdot \left((1 + 2\Delta \Delta^i \theta_t') - \frac{3}{2} \Delta \Delta^i \theta_t' \cdot \frac{\frac{L_t^i |\nabla f(\theta_t)^\top v_i|}{|\lambda_i|^2}}{\frac{L_t^i |\nabla f(\theta_t)^\top v_i|}{|\lambda_i|^2} + \frac{1}{4} \sqrt{1 + 2 \frac{L_t^i \cdot |\nabla f(\theta_t)^\top v_i|}{|\lambda_i|^2}} + \frac{1}{4}} \right) \\
&= \left(\frac{3}{2} \cdot \frac{\frac{L_t^i |\nabla f(\theta_t)^\top v_i|}{|\lambda_i|^2}}{\frac{L_t^i |\nabla f(\theta_t)^\top v_i|}{|\lambda_i|^2} + \frac{1}{4} \sqrt{1 + 2 \frac{L_t^i \cdot |\nabla f(\theta_t)^\top v_i|}{|\lambda_i|^2}} + \frac{1}{4}} - 2 \right) \Delta \Delta^i \theta_t'^3 - \Delta \Delta^i \theta_t'^2
\end{aligned}$$

$$\frac{\frac{L_t^i |\nabla f(\theta_t)^\top v_i|}{|\lambda_i|^2}}{\frac{L_t^i |\nabla f(\theta_t)^\top v_i|}{|\lambda_i|^2} + \frac{1}{4} \sqrt{1 + 2 \frac{L_t^i \cdot |\nabla f(\theta_t)^\top v_i|}{|\lambda_i|^2}} + \frac{1}{4}} \in [0, 1]$$

$$\psi_2 : \mathbb{R}^+ \rightarrow \mathbb{R}, \psi_2(x) = \frac{x}{x + \frac{1}{4} \sqrt{1 + 2x} + \frac{1}{4}}$$

$$\lim_{x \rightarrow 0^+} \psi_2(x) = 0$$

$$\lim_{x \rightarrow \infty} \psi_2(x) = 1$$

For the second part of the lemma,

$$\begin{aligned}
& \left| \frac{m_t^i (\Delta \theta_t^\top v_i) - m_t^i (\Delta \theta_t^{*\top} v_i)}{m_t^i (\Delta \theta_t^{*\top} v_i)} \right| \\
&= \frac{\nabla f(\theta_t)^\top \cdot v_i \cdot \left((\theta_{t+1} - \theta_t)^\top v_i - \Delta \theta_t^{*\top} v_i \right)}{\Delta \theta_t^{*\top} v_i \cdot \nabla f(\theta_t)^\top v_i + \frac{1}{2} \lambda_i (\Delta \theta_t^{*\top} v_i)^2 - \frac{L_t^i}{6} (\Delta \theta_t^{*\top} v_i)^3} \\
&+ \frac{\frac{1}{2} \lambda_i \cdot \left(\left((\theta_{t+1} - \theta_t)^\top v_i \right)^2 - (\Delta \theta_t^{*\top} v_i)^2 \right)}{\Delta \theta_t^{*\top} v_i \cdot \nabla f(\theta_t)^\top v_i + \frac{1}{2} \lambda_i (\Delta \theta_t^{*\top} v_i)^2 - \frac{L_t^i}{6} (\Delta \theta_t^{*\top} v_i)^3} \\
&- \frac{\frac{L_t^i}{6} \cdot \left(\left((\theta_{t+1} - \theta_t)^\top \cdot v_i \right)^3 - (\Delta \theta_t^{*\top} v_i)^3 \right)}{\Delta \theta_t^{*\top} v_i \cdot \nabla f(\theta_t)^\top v_i + \frac{1}{2} \lambda_i (\Delta \theta_t^{*\top} v_i)^2 - \frac{L_t^i}{6} (\Delta \theta_t^{*\top} v_i)^3} \\
&= \Delta \Delta^i \theta_t' \left(\frac{\frac{-L_t^i}{6} \cdot (\Delta \Delta^i \theta_t'^2 + 2 \Delta \Delta^i \theta_t' + 3)}{-\frac{L_t^i}{6} (\Delta \theta_t^{*\top} v_i)^2 + \frac{1}{2} \lambda_i \Delta \theta_t^{*\top} v_i - |\nabla f(\theta_t)^\top v_i|} \cdot (\Delta \theta_t^{*\top} v_i)^2 \right. \\
&+ \frac{\lambda_i \cdot \left(\frac{1}{2} \Delta \Delta^i \theta_t' + 1 \right)}{-\frac{L_t^i}{6} (\Delta \theta_t^{*\top} v_i)^2 + \frac{1}{2} \lambda_i \Delta \theta_t^{*\top} v_i - |\nabla f(\theta_t)^\top v_i|} \cdot \Delta \theta_t^{*\top} v_i \\
&\left. - \frac{|\nabla f(\theta_t)^\top \cdot v_i|}{-\frac{L_t^i}{6} (\Delta \theta_t^{*\top} v_i)^2 + \frac{1}{2} \lambda_i \Delta \theta_t^{*\top} v_i - |\nabla f(\theta_t)^\top v_i|} \right) \\
&= \Delta \Delta^i \theta_t' \left(\frac{-\frac{1}{6} \cdot (\Delta \Delta^i \theta_t'^2 + 2 \Delta \Delta^i \theta_t' + 3) \cdot \left(\sqrt{\lambda_i^2 + 2L_t^i \cdot |\nabla f(\theta_t)^\top v_i|} - \lambda_i \right)^2}{L_t^i} \right. \\
&- \frac{\left(\sqrt{\lambda_i^2 + 2L_t^i \cdot |\nabla f(\theta_t)^\top v_i|} - \lambda_i \right)^2}{L_t^i} + \frac{1}{2} \lambda_i \frac{\sqrt{\lambda_i^2 + 2L_t^i \cdot |\nabla f(\theta_t)^\top v_i|} - \lambda_i}{L_t^i} - |\nabla f(\theta_t)^\top v_i| \\
&+ \frac{\lambda_i \cdot \left(\frac{1}{2} \Delta \Delta^i \theta_t' + 1 \right) \cdot \left(\frac{\sqrt{\lambda_i^2 + 2L_t^i \cdot |\nabla f(\theta_t)^\top v_i|} - \lambda_i}{L_t^i} \right)}{\left(\sqrt{\lambda_i^2 + 2L_t^i \cdot |\nabla f(\theta_t)^\top v_i|} - \lambda_i \right)^2} \\
&- \frac{1}{6} \frac{\left(\sqrt{\lambda_i^2 + 2L_t^i \cdot |\nabla f(\theta_t)^\top v_i|} - \lambda_i \right)^2}{L_t^i} + \frac{1}{2} \lambda_i \frac{\sqrt{\lambda_i^2 + 2L_t^i \cdot |\nabla f(\theta_t)^\top v_i|} - \lambda_i}{L_t^i} - |\nabla f(\theta_t)^\top v_i| \\
&\left. - \frac{|\nabla f(\theta_t)^\top \cdot v_i|}{-\frac{1}{6} \frac{\left(\sqrt{\lambda_i^2 + 2L_t^i \cdot |\nabla f(\theta_t)^\top v_i|} - \lambda_i \right)^2}{L_t^i} + \frac{1}{2} \lambda_i \frac{\sqrt{\lambda_i^2 + 2L_t^i \cdot |\nabla f(\theta_t)^\top v_i|} - \lambda_i}{L_t^i} - |\nabla f(\theta_t)^\top v_i|} \right)
\end{aligned}$$

$$\begin{aligned}
&= \Delta \Delta^i \theta'_t \left(\frac{12\lambda_i \cdot \sqrt{\lambda_i^2 + 2L_t^i \cdot |\nabla f(\theta_t)^\top v_i|} - 12\lambda_i^2 - 12L_t^i |\nabla f(\theta_t)^\top v_i|}{5\lambda_i \sqrt{\lambda_i^2 + 2L_t^i \cdot |\nabla f(\theta_t)^\top v_i|} - 5\lambda_i^2 - 8L_t^i \cdot |\nabla f(\theta_t)^\top v_i|} \right. \\
&\quad + \frac{7\lambda_i \cdot \left(\sqrt{\lambda_i^2 + 2L_t^i \cdot |\nabla f(\theta_t)^\top v_i|} - \lambda_i \right) - 4L_t^i \cdot |\nabla f(\theta_t)^\top v_i|}{5\lambda_i \sqrt{\lambda_i^2 + 2L_t^i \cdot |\nabla f(\theta_t)^\top v_i|} - 5\lambda_i^2 - 8L_t^i \cdot |\nabla f(\theta_t)^\top v_i|} \cdot \Delta \Delta^i \theta'_t \\
&\quad \left. + \frac{2\lambda_i \sqrt{\lambda_i^2 + 2L_t^i \cdot |\nabla f(\theta_t)^\top v_i|} - 2\lambda_i^2 - 2L_t^i \cdot |\nabla f(\theta_t)^\top v_i|}{5\lambda_i \sqrt{\lambda_i^2 + 2L_t^i \cdot |\nabla f(\theta_t)^\top v_i|} - 5\lambda_i^2 - 8L_t^i \cdot |\nabla f(\theta_t)^\top v_i|} \cdot \Delta \Delta^i \theta_t'^2 \right)
\end{aligned}$$

Noting the common structure of each of the coefficients of $\Delta \Delta^i \theta_t'^1$, $\Delta \Delta^i \theta_t'^2$, $\Delta \Delta^i \theta_t'^3$, we prove the following to bound all three via appropriate settings of $a, b \in \{2, 4, 7, 12\}$:

If $\lambda_i > 0$:

$$\begin{aligned}
&\lim_{\frac{L_t^i \cdot |\nabla f(\theta_t)^\top v_i|}{\lambda_i^2} \rightarrow \mathcal{L}} \frac{a\lambda_i \cdot \left(\sqrt{\lambda_i^2 + 2L_t^i \cdot |\nabla f(\theta_t)^\top v_i|} - \lambda_i \right) - bL_t^i \cdot |\nabla f(\theta_t)^\top v_i|}{5\lambda_i \sqrt{\lambda_i^2 + 2L_t^i \cdot |\nabla f(\theta_t)^\top v_i|} - 5\lambda_i^2 - 8L_t^i \cdot |\nabla f(\theta_t)^\top v_i|} \\
&= \lim_{\frac{L_t^i \cdot |\nabla f(\theta_t)^\top v_i|}{\lambda_i^2} \rightarrow \mathcal{L}} \frac{a \frac{2}{\sqrt{1+2\frac{L_t^i \cdot |\nabla f(\theta_t)^\top v_i|}{\lambda_i^2} + 1}} - b}{5 \frac{2}{\sqrt{1+2\frac{L_t^i \cdot |\nabla f(\theta_t)^\top v_i|}{\lambda_i^2} + 1}} - 8} \\
&= \begin{cases} \frac{b-a}{3} & \mathcal{L} = 0^+ \\ \frac{b}{8} & \mathcal{L} = \infty \end{cases}
\end{aligned}$$

If $\lambda_i \leq 0$:

$$\begin{aligned}
&\lim_{\frac{L_t^i \cdot |\nabla f(\theta_t)^\top v_i|}{\lambda_i^2} \rightarrow \mathcal{L}} \frac{a|\lambda_i| \cdot \left(\sqrt{|\lambda_i|^2 + 2L_t^i \cdot |\nabla f(\theta_t)^\top v_i|} + |\lambda_i| \right) + bL_t^i \cdot |\nabla f(\theta_t)^\top v_i|}{5|\lambda_i| \left(\sqrt{|\lambda_i|^2 + 2L_t^i \cdot |\nabla f(\theta_t)^\top v_i|} + |\lambda_i| \right) + 8L_t^i \cdot |\nabla f(\theta_t)^\top v_i|} \\
&= \lim_{\frac{L_t^i \cdot |\nabla f(\theta_t)^\top v_i|}{\lambda_i^2} \rightarrow \mathcal{L}} \frac{a \left(\sqrt{1+2\frac{L_t^i \cdot |\nabla f(\theta_t)^\top v_i|}{|\lambda_i|^2}} + 1 \right) + b \frac{L_t^i \cdot |\nabla f(\theta_t)^\top v_i|}{|\lambda_i|^2}}{5 \left(\sqrt{1+2\frac{L_t^i \cdot |\nabla f(\theta_t)^\top v_i|}{|\lambda_i|^2}} + 1 \right) + 8 \frac{L_t^i \cdot |\nabla f(\theta_t)^\top v_i|}{|\lambda_i|^2}} \\
&= \begin{cases} \frac{a}{5} & \mathcal{L} = 0^+ \\ \frac{b}{8} & \mathcal{L} = \infty \end{cases}
\end{aligned}$$

Analogously to the first case, and due to the monotonic natures (for all $a, b \in \mathbb{R}$) of

$$\psi_3 : \mathbb{R}^+ \rightarrow \mathbb{R}, \psi_3(x) = \frac{a \frac{2}{\sqrt{1+2x+1}} - b}{5 \frac{2}{\sqrt{1+2x+1}} - 8}$$

and

$$\psi_4 : \mathbb{R}^+ \rightarrow \mathbb{R}, \psi_4(x) = \frac{a(\sqrt{1+2x+1}) + bx}{5(\sqrt{1+2x+1}) + 8x}$$

the term in the parentheses is bounded, thus we may conclude our proof of the lemma. \square

Remark. Note that when $\lambda_i > 0$, $\frac{L_i^i \cdot |\nabla f(\theta_t)^\top v_i|}{\lambda_i^2} \rightarrow 0^+$, the coefficients of $\Delta\Delta^i\theta_t^3$, $\Delta\Delta^i\theta_t^1$ shrink to 0 (since $a = b$ for those cases), so that $\left| \frac{m_t^i(\Delta\theta_t^\top v_i) - m_t^i(\Delta\theta_t^{*\top} v_i)}{m_t^i(\Delta\theta_t^{*\top} v_i)} \right| = \Theta\left(\Delta\Delta^i\theta_t^2\right)$

We are now ready to prove theorem 3.3.

Theorem. Worst-case descent rate for arbitrary optimizers

Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ a twice-differentiable function satisfying assumptions 1 and 2, and let $\Delta\theta_t$ satisfy $M_t^i(\Delta\theta_t^\top v_i) \leq 0$. Then

•

$$\left| \frac{M_t^i(\Delta\theta_t^\top v_i)}{M_t^i(\Delta\theta_t^{*\top} v_i)} \right| = \Theta\left(1 + |\Delta\Delta^i\theta_t|^2\right)$$

•

$$\left| \frac{m_t^i(\Delta\theta_t^\top v_i)}{m_t^i(\Delta\theta_t^{*\top} v_i)} \right| = \Theta\left(1 + |\Delta\Delta^i\theta_t|^p\right) \quad (20)$$

$$\text{with } p = \begin{cases} 2 & \lambda_i > 0 \wedge \frac{|\nabla f(\theta_t)^\top v_i|}{\lambda_i^2} = 0 \\ 1 & \text{else} \end{cases}.$$

Proof. Proof is immediate from lemma M.5, because we have

$$\frac{M_t^i(\Delta\theta_t^\top v_i)}{M_t^i(\Delta\theta_t^{*\top} v_i)} = 1 + \frac{M_t^i(\Delta\theta_t^\top v_i) - M_t^i(\Delta\theta_t^{*\top} v_i)}{M_t^i(\Delta\theta_t^{*\top} v_i)}$$

and similarly for $m_t^i(\Delta\theta_t^\top v_i)$. □

N LIMITATIONS AND FUTURE WORK

One interesting direction for future research is in putting the estimated Lipschitz parameters to work throughout the optimization process to increase the descent rate in hopes of matching and even surpassing ARC’s strong performance (Xu et al., 2017), i.e. applying algorithm ELMO in practice. We provide some basic methods for Lipschitz estimation in this paper, but future work may propose more accurate/low-overhead estimators. In this paper we focused primarily on the optimizer selection challenge, thus no experiments in this direction were attempted.

A second direction for future research is a bias correction for our simple Lipschitz parameter estimation in AMOS, much like Kingma & Ba (2014) did for Duchi et al. (2011a) and Tieleman & Hinton (2012). We also recall that the AMOS experiments in this work were performed on a per-tensor basis, but other partitions (e.g. constant chunk-size irrespective of tensor borders) may work better - no experiments were attempted in this direction.

One limitation of our Newton’s method performance predictor is the additional computational burden of computing the Lipschitz parameters. While some algorithms (e.g. Sophia) explicitly compute all the Hessian eigenvalues, others do not, which poses challenges in computing Lipschitz parameters with minimal overhead. However since the Hessian is usually only estimated at large intervals, the additional cost of Lipschitz estimation may be amortized across these iterations.

A second limitation is the need to store a copy of the model parameters at the position at which the previous Hessian was estimated; this adds pressure to the memory, which may pose a bottleneck in settings with large models or activations. Many modern optimizers already hold several buffers the same size as the model, however (e.g. Adam, which stores the first- and second gradient moments), such that this copy does not affect the order of magnitude of storage required. Furthermore, certain training models already store previous copies of the model parameters (e.g. Allen-Zhu (2016)), such that this may not add any additional overhead.

2430 A third limitation of our work is its inability to provide any indication of the number of iterations left
2431 to achieve convergence (unlike previous works which bound the gradient norm as a function of the
2432 number of iterations, which is a measure of convergence). We see this as an acceptable limitation
2433 however, since in practice a model is only required to achieve a certain level of performance on the
2434 data decided ahead of time, without regard to how much further it could be optimized. As noted in
2435 the introduction, performance is measured by the loss function, so our descent rate bound satisfies
2436 this practical requirement.

2437 A final limitation of our bound is its reliance on $\Delta\Delta^i\theta_t$ as a measure of algorithm optimality which
2438 is a function of $\Delta\theta_t^{*\top}v_i$, despite the fact that most optimizers do not compute that during training.
2439 This bound is therefore primarily of theoretical interest, as illustrated by its motivation of the very
2440 practical metric discussed in section 4.

2441
2442
2443
2444
2445
2446
2447
2448
2449
2450
2451
2452
2453
2454
2455
2456
2457
2458
2459
2460
2461
2462
2463
2464
2465
2466
2467
2468
2469
2470
2471
2472
2473
2474
2475
2476
2477
2478
2479
2480
2481
2482
2483