

UNIVERSAL CONCAVITY-AWARE DESCENT RATE FOR OPTIMIZERS

Anonymous authors

Paper under double-blind review

ABSTRACT

Many machine learning problems involve a challenging task of calibrating parameters in a computational model to fit the training data; this task is especially challenging for non-convex problems. Many optimization algorithms have been proposed to assist in calibrating these parameters, each with its respective advantages in different scenarios, but it is often difficult to determine the scenarios for which an algorithm is best suited. To contend with this challenge, much work has been done on proving the rate at which these optimizers converge to their final solution, however the wide variety of such convergence rate bounds, each with their own different assumptions, convergence metrics, tightnesses, and parameters (which may or may not be known to the practitioner) make comparing these convergence rates difficult. To help with this problem, we present a minmax-optimal algorithm and, by comparison to it, give a single descent bound which is applicable to a very wide family of optimizers, tasks, and data (including all of the most prevalent ones), which also puts special emphasis on being tight even in parameter subspaces in which the cost function is concave.

1 INTRODUCTION

Many machine learning problems involve calibrating the parameters of a given model to match the data distribution of a phenomenon one wishes to model, e.g. the structure of folded proteins, processing images to automatically generate appropriate labels for them, or generating images and text to interactively chat with a human engagingly. This process involves:

1. Collecting many samples ("data points") from the desired data distribution.
2. Measuring how well the model fits the collected data points (the "data set") with a given performance analysis metrics (the "loss function", a.k.a. the "objective function"). By convention, lower values of the loss function imply better performance on the model's part.
3. Adjusting the model's parameters to improve the performance, as measured by the loss function ("model parameter optimization").
4. Repeat until desired performance achieved.

However, no single existing 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. As a result, one must compare among various different optimization algorithms (henceforth, "optimizers") to select the one most suited to the current scenario.

In an effort to help practitioners select the best optimizer for their setup and estimate the absolute computational resources that will be required to obtain a given performance, many experiments have been run comparing the performance of different optimizers on a variety of applications (Xu et al., 2017; Schmidt et al., 2021), and on the theoretical side - convergence rate bounds have been proven for various different optimizers. However, due to the wide variety of assumptions, convergence rate

054 metrics, bound parameters (which may be expensive - if not impossible - to compute ahead of time),
 055 and tightness of the bounds in all of these works, comparing among them remains a challenging
 056 task. Secondly, there is a lack of convergence rate bounds general enough to be easily applicable
 057 to newly proposed optimizers. Thirdly, many of these bounds fail to demonstrate the empirically-
 058 verified convergence rate superiority of the more sophisticated methods that make use of second-order
 059 curvature information instead of exclusively the gradient. Lastly, although convergence rate bounds
 060 exist for non-convex functions, many of them fail to properly address the opportunities that lay
 061 in linear subspaces of the parameter space in which the loss function is concave (meaning that a
 062 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
 063 that more attention should be given to these subspaces of the function in the context of neural network
 064 optimization; Alain et al. (2018) and Ghorbani et al. (2019) demonstrate experimentally that there is
 065 much to be gained by taking optimal steps in these subspaces, often even orders of magnitude greater
 066 than the potential gains in convex subspaces.

067
 068 **Our contributions** In an effort to help practitioners select the best optimizer for their use case,
 069 we develop a tool for estimating the value of second-order optimization algorithms; this will help
 070 decide if the additional computational burden of these algorithms is worthwhile. We develop a
 071 minmax-optimal algorithm, rate algorithms by similarity to it, and demonstrate in theory and in
 072 practice that in general, second-order algorithms work best on mechanistically simple problems. Our
 073 algorithm-optimality bound satisfies the following good properties:

- 074
 075 1. **Concave tightness** Our bound exploits the opportunity for greater descent in subspaces of
 076 the parameter space in which the loss function is concave.
- 077
 078 2. **Universality** We make only weak and commonly satisfied assumptions for our bound, to
 079 allow for its application to a wide and prevalent family of optimizers and loss functions.
- 080
 081 3. **Tightness for any level of iteration step-quality** instead of assuming a bound on the quality
 082 of steps given in each iteration as some previous works have done, our theoretical bounds
 083 are given as a continuous function of the quality of each iteration’s step.
- 084
 085 4. **Bound on loss function descent** Our main result bounds the rate at which the model’s
 086 performance increases (as measured by the loss function). This is in contrast to previous
 087 works, which instead bound various indicators of local minimality, such as gradient norm,
 088 local near-convexity, or proximity to a local minimum (in Euclidean distance). Although Xu
 089 et al. (2020) write that the latter convergence rate metrics is more relevant to the non-convex
 090 optimization setting, we feel that the former is more practically useful, since generally
 091 real-world applications with limited computational resources simply demand a minimal
 092 performance guarantee of their model, without regard to the theoretical capabilities of a
 093 given model or optimization algorithm.
- 094
 095 5. **Simplicity of cubic minimization problem** We approach the multidimensional cubic
 096 polynomial minimization problem posed by Nesterov & Polyak (2006) by decomposing it
 097 into n 1-dimensional problems via eigendecomposition of the Hessian, making our approach
 098 to the solution of this minimization problem far simpler conceptually.

099
 100 Our paper is organized as follows: In section 2, we review previous work and describe the notation
 101 we will use throughout the paper. In section 3, we develop the minmax-optimal `ELMO` algorithm and
 102 analyze its descent rate. In section 4, we make claims as to the benefits of optimizer similarity to `ELMO`
 103 (proven in appendix H). In section 5 we show the value of our novel Lipschitz parameter separation
 104 scheme by showing how much lower the Lipschitz parameters of most relevant eigenspaces can be,
 105 thus giving optimizers a more accurate minimizable model of the loss function in each neighborhood
 106 it finds itself in. Finally, in section 6, we present experiments validating one particular use case of
 107 our bound: we show that the advantage second-order optimizers hold over first-order optimizers is
 inversely proportional to the convex Lipschitz parameters. In other words, second-order optimizers
 present strong performance (thus may be worth their additional computational burden) in settings
 with small convex Lipschitz parameters, and weak performance (thus not worthwhile) in settings
 with large convex Lipschitz parameters.

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 $\theta_{t+1} - \theta_t \triangleq \Delta\theta_t$.
- 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 $(v_i(\theta_t), \lambda_i(\theta_t))_{j=0}^n$ 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

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

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

Notation 3. Let $A, B \in \mathbb{R}^{n \times n}$. We use the following notations (when applicable):

- We mark A 's transpose as A^\top .
- 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 $A \succ B$).
- Mark $\lambda_{\min}(A), \lambda_{\max}(A)$ the minimal/maximal eigenvalue of A , respectively, and their ratio $\kappa(A) = \frac{\lambda_{\max}(A)}{\lambda_{\min}(A)}$ the condition number of A .

Notation 4. 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 $\xi(U \cdot D \cdot U^\top) = U \cdot \text{diag}(\xi(d_1), \xi(d_2), \dots, \xi(d_n)) \cdot U^\top$.

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.

We note that the overwhelming majority of gradient-based optimizers may be expressed as quasi-Newton optimizers (some popular examples may be seen in Martens (2020)). As a result, 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

As discussed in item 4 of the contributions section, 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) split the process of optimization with a general optimizer into the initial "search phase", in which the optimizer searches for an approximately convex region in which the point it will eventually converge to resides, and the later "final phase", in which the optimizer converges to its final solution within this convex region.

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 (which is rarely satisfied by the loss functions in neural network optimization scenarios).

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 on convergence rate.

3 A MINMAX HESSIAN LIPSCHITZ-AWARE OPTIMIZATION ALGORITHM

Any deterministic optimization algorithm is comprised of two parts: first, we gather information about the loss function to enable us to implicitly construct a local model of the loss function, and secondly we step to the minimum of this model. Accordingly, gradient descent and Newton’s method use first- and second-order Taylor approximations of f respectively, and while these models 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 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, leading to looseness in convergence rate bounds and subpar 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:

$$\forall \theta \in \mathbb{R}^n \exists (\bar{L}^i)_{i=1}^n \in (\mathbb{R}^+)^n \forall \varphi \in \mathbb{R}^n : \left| \lambda_i - \tilde{\lambda}_i \right| \leq \bar{L}^i \cdot \left| (\theta - \varphi)^\top v_i \right|$$

$$\forall \theta \in \mathbb{R}^n \exists L_R \in \mathbb{R}^+ \forall \varphi \in \mathbb{R}^n : \left\| V - \tilde{V} \right\|_2 \leq L_R \cdot \|\theta - \varphi\|_2$$

$$\exists L_H \in \mathbb{R} \forall \theta \in \mathbb{R}^n \forall_{i \in [n]} : \max \{L_R, \bar{L}^i\} \leq L_H \wedge \bar{L}^i \geq L_H^{-1}$$

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 value of this finer assumption later, by demonstrating that these parameters vary widely. In particular, and taking into account that optimization primarily occurs in a very limited subspace of the domain space (Gur-Ari et al., 2018), we will demonstrate that the Lipschitz parameters relevant to these subspaces are often orders of magnitude smaller than the others.

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

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

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

3.2 EXPLOITING THESE BOUNDS FOR A MINMAX ALGORITHM

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

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

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

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

Finally, we are ready to present algorithm Eigenspace-Lipschitz Minmax Optimizer (ELMO). We mark EIGEN an eigendecomposition subroutine and LIPSCHITZ a Lipschitz parameter oracle.

An important observation to make about the algorithm above is its equal applicability to convex and concave regions of the domain space. In fact, when $\lambda_i < 0$ (implying a concave subspace), the step size (and, correspondingly, the amount of descent on our model of the loss function M_t^i) is actually greater than otherwise. This is due to ELMO’s ability to make use of concave regions of the loss function for greater descent.

3.3 ALGORITHM ELMO’S DESCENT RATE

An important factor in deciding how much computational power to put into optimizing a model is the ratio between the cost of computational resources and the improvement to the model’s quality. 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 3.2. *Worst case-optimal descent rate Let f be a function with Lipschitz-continuous Hessian. After t iterations, algorithm ELMO satisfies*

$$f(\theta_0) - f(\theta_t) = \mathcal{O}(\log t) \quad (8)$$

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.

Theorem 3.3. *Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ satisfying assumptions 1 and 2. Algorithm ELMO satisfies*

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

4 DESCENT RATE OF QUASI-NEWTON OPTIMIZATION ALGORITHMS

Although algorithm ELMO is 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 (since its leading coefficient is not assumed to be nonzero) and its model minimization step is unique, its greater computational burden of computing the Lipschitz parameters may cause it to be an ineffective optimization algorithm in practice. Since most prevalent practical optimizers today belong to the Quasi-Newton family, we satisfy ourselves with a quantification of their quality based on their similarity to this ideal algorithm.

THE MINMAX PRECONDITIONER

Since Quasi-Newton methods are characterized by their preconditioners, we must first develop algorithm ELMO’s characteristic preconditioner. We begin by defining a metric of distance between optimization algorithms by the difference between their characteristic steps, 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, it is a function of that algorithm's defining preconditioner: $\Delta\Delta^i\theta_t = \Delta\Delta^i\theta_t(\Phi_t)$

Lemma 4.1. *Minmax 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

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

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 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 $\forall_{i \in [n].s.t. \lambda_i \geq 0}$:

$$M_t^i \left(\frac{|\nabla f(\theta_t)_t v_i|}{\lambda_i} \right) \geq 0 \text{ for any step } t \text{ and eigenspace } i \text{ with } L_t^i \geq -3 \frac{\lambda_i^2}{|\nabla f(\theta_t)^\top v_i|}$$

4.1 PER-ITERATION DESCENT OF ARBITRARY STEP

Due to the computational difficulty of computing ELMO's iteration step precisely, practitioners may prefer computationally cheaper alternatives. To address this, we provide guarantees for the worst-case rate of loss function descent of an arbitrary optimization algorithm relative to algorithm ELMO's descent, as a function of the algorithm's similarity to ELMO. 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.

Theorem 4.2. *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

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

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

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

4.2 GENERALIZATION OF PREVIOUS QUASI-NEWTON PRECONDITIONER QUALITY METRICS

Notation 9. 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 concave subspaces of the loss function domain space and when the curvature shift is negligible, we examine

the special case of our metric when $\lambda_i > 0$ (when the loss function is concave 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

$$|\Delta\Delta^i\theta'_t| = \left| 1 - \frac{\nabla f(\theta_t)^\top}{\nabla f(\theta_t)^\top v_i} \cdot (\alpha_t \Phi_t \mathcal{H}(\theta_t)) \cdot v_i \cdot \frac{\sqrt{1 + 2L_t^i \cdot \frac{|\nabla f(\theta_t)^\top v_i|}{\lambda_i^2}} + 1}{2} \right| \quad (10)$$

Županski (1993) introduce the "Effective Hessian" (a.k.a. the "Preconditioned Hessian") as $\mathcal{I}_t = \alpha_t \Phi_t \mathcal{H}(\theta_t)$, with its 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 10.

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 10 simplifies to

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

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.

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

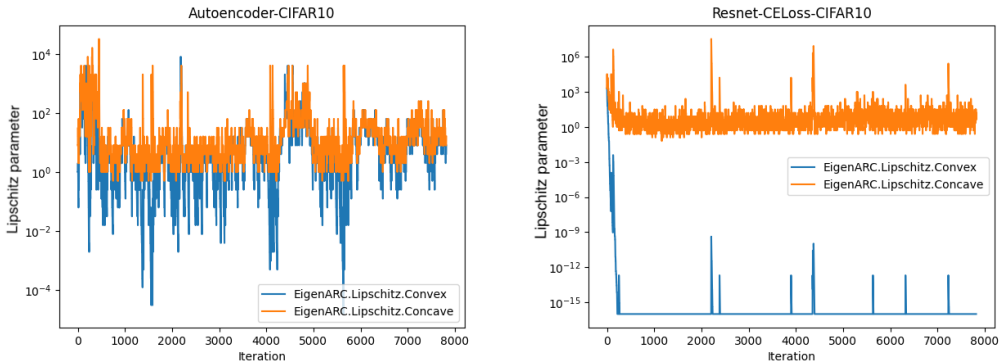
Another source of interest in these parameters’ distribution is in explaining the effectiveness of second-order quasi-Newton algorithms that implicitly assume the Lipschitz parameters are insignificant (i.e. very close to zero), since their model of the loss function is a quadratic Taylor polynomial (i.e. no curvature shift); this may be seen from equation 10 which shows optimality of Newton’s method only when $\lambda_i > 0$ and $L^+ = 0$. We will show that they are not generally small by any means, however we will show that the Lipschitz parameters of the subspaces in which they work (the convex subspaces - see the implementation of Sivan et al. (2024), for instance, which applies Newton’s method only on subspaces with significantly convex subspaces) are in fact small in certain settings.

5.1 EXPERIMENTS

The first source of evidence for our claim is from existing literature on the subject; we defer a discussion of this to appendix E. 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 $\frac{f(\theta_t) - f(\theta_{t+1})}{|\sum_{i=1}^n M_t^i|}$ 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 F as well as the full set of our experiment results, however we present two experiments in figure 1 for completeness. Our experiments show

432
433
434
435
436
437
438
439
440
441
442
443
444



445
446
447
448
449
450
451
452
453
454
455
456
457
458

Figure 1: Comparisons of convex-subspace Lipschitz parameters to concave-subspace Lipschitz parameters. *Logarithmic scale*

that as expected, $L^+ \ll L^-$, and the gap widens exponentially as training progresses in all cases except the autoencoders. Since we will see 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. Interestingly, the Lipschitz parameters seem to depend primarily on the task, and are much less affected by network structure or model output-target loss function.

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

459
460
461
462
463
464
465
466
467
468
469
470

- 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.

471

6 A QUALITY PREDICTOR FOR NEWTON’S METHOD

472
473
474
475
476
477
478

Expanding on the latter application in section 5, 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 10 by showing that the expected descent in a given eigenspace is an approximately monotonically decreasing function of the corresponding Lipschitz parameter.

479
480
481
482
483
484
485

Figure 2 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:

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

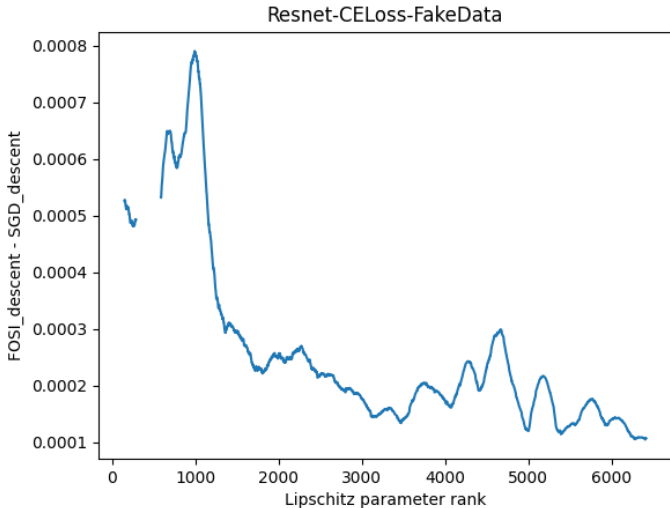


Figure 2: Inverse relation between a convex-subspace Lipschitz parameter and corresponding descent superiority of Quasi-Newton method

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 F, as well as further detailing on figure 2.

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 G. Alternatively, practitioners may choose to use ARC steps instead of first-order methods, when the Lipschitz parameter is significantly large. These findings may instead be used to construct a meta-optimizer, that periodically computes Lipschitz parameters and adaptively selects optimizers and optimization hyperparameters throughout the optimization process accordingly. We leave this direction to future research.

| 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

7 CONCLUSION

In this work we developed and analyzed a Hessian eigenspace Lipschitz-aware minmax optimization algorithm ELMO 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 ELMO. 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.

REFERENCES

- 540
541
542 Naman Agarwal, Brian Bullins, Xinyi Chen, Elad Hazan, Karan Singh, Cyril Zhang, and Yi Zhang.
543 Efficient full-matrix adaptive regularization. In Kamalika Chaudhuri and Ruslan Salakhutdinov
544 (eds.), *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of
545 *Proceedings of Machine Learning Research*, pp. 102–110. PMLR, 09–15 Jun 2019. URL <https://proceedings.mlr.press/v97/agarwall19b.html>.
546
- 547 Naman Agarwal, Rohan Anil, Elad Hazan, Tomer Koren, and Cyril Zhang. Learning rate grafting:
548 Transferability of optimizer tuning, 2022. URL https://openreview.net/forum?id=FpKgG31Z_i9.
549
- 550 Guillaume Alain, Nicolas Le Roux, and Pierre-Antoine Manzagol. Negative eigenvalues of the
551 hessian in deep neural networks. In *6th International Conference on Learning Representations,*
552 *ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Workshop Track Proceedings.*
553 OpenReview.net, 2018. URL <https://openreview.net/forum?id=S1iiddyDG>.
554
- 555 Shun-ichi Amari. Natural Gradient Works Efficiently in Learning. *Neural Computation*, 10(2):
556 251–276, 02 1998. ISSN 0899-7667. doi: 10.1162/089976698300017746. URL <https://doi.org/10.1162/089976698300017746>.
557
- 558 Rohan Anil, Vineet Gupta, Tomer Koren, Kevin Regan, and Yoram Singer. Second order optimization
559 made practical. *CoRR*, abs/2002.09018, 2020. URL <https://arxiv.org/abs/2002.09018>.
560
- 561 Somenath Bera and Vimal Shrivastava. Analysis of various optimizers on deep convolutional
562 neural network model in the application of hyperspectral remote sensing image classification.
563 *International Journal of Remote Sensing*, 41:2664–2683, 04 2020. doi: 10.1080/01431161.2019.
564 1694725.
565
- 566 Dimitri P. Bertsekas. Incremental least squares methods and the extended kalman filter. *SIAM J.*
567 *Optim.*, 6(3):807–822, 1996. doi: 10.1137/S1052623494268522. URL <https://doi.org/10.1137/S1052623494268522>.
568
- 569 Aleksandar Botev, Hippolyt Ritter, and David Barber. Practical Gauss-Newton optimisation for
570 deep learning. In Doina Precup and Yee Whye Teh (eds.), *Proceedings of the 34th International*
571 *Conference on Machine Learning*, volume 70 of *Proceedings of Machine Learning Research*,
572 pp. 557–565. PMLR, 06–11 Aug 2017. URL <https://proceedings.mlr.press/v70/botev17a.html>.
573
- 574 Léon Bottou. Stochastic learning. In Olivier Bousquet, Ulrike von Luxburg, and Gunnar Rätsch
575 (eds.), *Advanced Lectures on Machine Learning: ML Summer Schools 2003, Canberra, Australia,*
576 *February 2 - 14, 2003, Tübingen, Germany, August 4 - 16, 2003, Revised Lectures*, pp. 146–168.
577 Springer Berlin Heidelberg, Berlin, Heidelberg, 2004. ISBN 978-3-540-28650-9. doi: 10.1007/
578 978-3-540-28650-9_7. URL https://doi.org/10.1007/978-3-540-28650-9_7.
579
- 580 Léon Bottou and Yann Lecun. On-line learning for very large datasets. *J. Applied Stochastic Models*
581 *in Business and Industry*, 01 2004.
- 582 C. Cartis, N.I.M. Gould, and Ph.L. Toint. Complexity bounds for second-order optimality in
583 unconstrained optimization. *Journal of Complexity*, 28(1):93–108, 2012a. ISSN 0885-064X.
584 doi: <https://doi.org/10.1016/j.jco.2011.06.001>. URL <https://www.sciencedirect.com/science/article/pii/S0885064X11000537>.
585
- 586 Coralía Cartis, Nicholas I. M. Gould, and Philippe L. Toint. Adaptive cubic regularisation methods for
587 unconstrained optimization. part i: motivation, convergence and numerical results. *Mathematical*
588 *Programming*, 127(2):245–295, 2011a. doi: 10.1007/s10107-009-0286-5. URL <https://doi.org/10.1007/s10107-009-0286-5>.
589
- 590 Coralía Cartis, Nicholas I. M. Gould, and Philippe L. Toint. Adaptive cubic regularisation methods
591 for unconstrained optimization. part ii: worst-case function- and derivative-evaluation complexity.
592 *Mathematical Programming*, 130(2):295–319, 2011b. doi: 10.1007/s10107-009-0337-y. URL
593 <https://doi.org/10.1007/s10107-009-0337-y>.

- 594 Coralia Cartis, Nicholas I. M. Gould, and Philippe L. Toint. Evaluation complexity of adaptive cubic
595 regularization methods for convex unconstrained optimization. *Optimization Methods and Soft-*
596 *ware*, 27:197 – 219, 2012b. URL [https://api.semanticscholar.org/CorpusID:](https://api.semanticscholar.org/CorpusID:16647191)
597 16647191.
- 598 Andrew R. Conn, Nicholas I. M. Gould, and Philippe L. Toint. *Trust-Region Methods*. SIAM,
599 Philadelphia, PA, USA, 2000.
- 600
601 Yann N. Dauphin, Razvan Pascanu, Caglar Gulcehre, Kyunghyun Cho, Surya Ganguli, and Yoshua
602 Bengio. Identifying and attacking the saddle point problem in high-dimensional non-convex opti-
603 mization. In *Proceedings of the 27th International Conference on Neural Information Processing*
604 *Systems - Volume 2*, NIPS’14, pp. 2933–2941, Cambridge, MA, USA, 2014. MIT Press.
- 605
606 Alexandre D’efossez, Léon Bottou, Francis R. Bach, and Nicolas Usunier. A simple convergence
607 proof of adam and adagrad. *Trans. Mach. Learn. Res.*, 2022, 2020. URL [https://api.](https://api.semanticscholar.org/CorpusID:225213299)
608 [semanticscholar.org/CorpusID:225213299](https://api.semanticscholar.org/CorpusID:225213299).
- 609
610 Ron S. Dembo, Stanley C. Eisenstat, and Trond Steihaug. Inexact newton methods. *SIAM Journal on*
611 *Numerical Analysis*, 19(2):400–408, 1982. doi: 10.1137/0719025. URL [https://doi.org/](https://doi.org/10.1137/0719025)
612 10.1137/0719025.
- 613 J.E. Dennis and R.B. Schnabel. *Numerical Methods for Unconstrained Optimization and Non-*
614 *linear Equations*. Prentice-Hall Civil Engineering and Engineering Mechanics Se. Prentice-
615 Hall, 1983. ISBN 9780136272168. URL [https://books.google.co.il/books?id=](https://books.google.co.il/books?id=4fFQAAAAMAAJ)
616 4fFQAAAAMAAJ.
- 617 Tian Ding, Dawei Li, and Ruoyu Sun. Sub-optimal local minima exist for almost all over-
618 parameterized neural networks. *ArXiv*, abs/1911.01413, 2019. URL [https://api.](https://api.semanticscholar.org/CorpusID:207870322)
619 [semanticscholar.org/CorpusID:207870322](https://api.semanticscholar.org/CorpusID:207870322).
- 620
621 John Duchi, Elad Hazan, and Yoram Singer. Adaptive subgradient methods for online learning and
622 stochastic optimization. *Journal of Machine Learning Research*, 12(61):2121–2159, 2011a. URL
623 <http://jmlr.org/papers/v12/duchilla.html>.
- 624
625 John C. Duchi, Elad Hazan, and Yoram Singer. Adaptive subgradient methods for online learning
626 and stochastic optimization. *J. Mach. Learn. Res.*, 12:2121–2159, 2011b. URL [https://api.](https://api.semanticscholar.org/CorpusID:538820)
627 [semanticscholar.org/CorpusID:538820](https://api.semanticscholar.org/CorpusID:538820).
- 628
629 S C Eisenstat and H F Walker. Choosing the forcing terms in an inexact newton method. *SIAM*
630 *Journal on Scientific Computing*, 17(1), 1 1996. doi: 10.1137/0917003. URL [https://www.](https://www.osti.gov/biblio/218521)
631 [osti.gov/biblio/218521](https://www.osti.gov/biblio/218521).
- 632
633 William R. Esposito and Christodoulos A. Floudas. Gauss–newton method: Least squares, relation to
634 newton’s method. In Christodoulos A. Floudas and Panos M. Pardalos (eds.), *Encyclopedia of Optimization*, pp. 733–738. Springer
635 US, Boston, MA, 2001. ISBN 978-0-306-48332-5. doi: 10.1007/0-306-48332-7_151. URL
636 https://doi.org/10.1007/0-306-48332-7_151.
- 637
638 Reuben Feinman. Pytorch-minimize: a library for numerical optimization with autograd, 2021. URL
639 <https://github.com/rfeinman/pytorch-minimize>.
- 640
641 Dan Garber, Elad Hazan, Chi Jin, Kakade Sham, Cameron Musco, Praneeth Netrapalli, and Aaron
642 Sidford. Faster eigenvector computation via shift-and-invert preconditioning. In Maria Florina
643 Balcan and Kilian Q. Weinberger (eds.), *Proceedings of The 33rd International Conference on*
644 *Machine Learning*, volume 48 of *Proceedings of Machine Learning Research*, pp. 2626–2634,
645 New York, New York, USA, 20–22 Jun 2016. PMLR. URL [https://proceedings.mlr.](https://proceedings.mlr.press/v48/garber16.html)
646 [press/v48/garber16.html](https://proceedings.mlr.press/v48/garber16.html).
- 647
648 R. Garmanjani. A note on the worst-case complexity of nonlinear stepsize control methods for convex
649 smooth unconstrained optimization. *Optimization*, 71:1–11, 10 2020. doi: 10.1080/02331934.
650 2020.1830088.

- 648 Rong Ge, Furong Huang, Chi Jin, and Yang Yuan. Escaping from saddle points — online stochastic
649 gradient for tensor decomposition. In Peter Grünwald, Elad Hazan, and Satyen Kale (eds.),
650 *Proceedings of The 28th Conference on Learning Theory*, volume 40 of *Proceedings of Machine
651 Learning Research*, pp. 797–842, Paris, France, 03–06 Jul 2015. PMLR. URL [https://
652 proceedings.mlr.press/v40/Ge15.html](https://proceedings.mlr.press/v40/Ge15.html).
- 653 Behrooz Ghorbani, Shankar Krishnan, and Ying Xiao. An investigation into neural net opti-
654 mization via hessian eigenvalue density. In Kamalika Chaudhuri and Ruslan Salakhutdinov
655 (eds.), *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of
656 *Proceedings of Machine Learning Research*, pp. 2232–2241. PMLR, 09–15 Jun 2019. URL
657 <https://proceedings.mlr.press/v97/ghorbani19b.html>.
- 658 Donald Goldfarb, Yi Ren, and Achraf Bahamou. Practical quasi-newton methods for training deep
659 neural networks. In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin (eds.),
660 *Advances in Neural Information Processing Systems*, volume 33, pp. 2386–2396. Curran Asso-
661 ciates, Inc., 2020. URL [https://proceedings.neurips.cc/paper_files/paper/
662 2020/file/192fc044e74df9ac5dc9f3395-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2020/file/192fc044e74df9ac5dc9f3395-Paper.pdf).
- 663 I. J. Good. Rational decisions. *Journal of the Royal Statistical Society. Series B (Methodological)*, 14
664 (1):107–114, 1952. ISSN 00359246. URL <http://www.jstor.org/stable/2984087>.
- 665 Andreas Griewank. The modification of newton’s method for unconstrained optimization by bounding
666 cubic terms. Technical report, Department of Applied Mathematics and Theoretical Physics,
667 University of Cambridge, United Kingdom, 1981.
- 668 Vineet Gupta, Tomer Koren, and Yoram Singer. Shampoo: Preconditioned stochastic tensor opti-
669 mization. In Jennifer Dy and Andreas Krause (eds.), *Proceedings of the 35th International
670 Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*,
671 pp. 1842–1850. PMLR, 10–15 Jul 2018. URL [https://proceedings.mlr.press/v80/
672 gupta18a.html](https://proceedings.mlr.press/v80/gupta18a.html).
- 673 Guy Gur-Ari, Daniel A. Roberts, and Ethan Dyer. Gradient descent happens in a tiny subspace.
674 *ArXiv*, abs/1812.04754, 2018. URL [https://api.semanticscholar.org/CorpusID:
675 54480858](https://api.semanticscholar.org/CorpusID:54480858).
- 676 Kaiming He, X. Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition.
677 *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 770–778, 2015.
678 URL <https://api.semanticscholar.org/CorpusID:206594692>.
- 679 M.F. Hutchinson. A stochastic estimator of the trace of the influence matrix for laplacian smoothing
680 splines. *Communication in Statistics- Simulation and Computation*, 18:1059–1076, 01 1989. doi:
681 10.1080/03610919008812866.
- 682 Maarten Jansen and Gerda Claeskens. Cramér–rao inequality. In Miodrag Lovric (ed.), *International
683 Encyclopedia of Statistical Science*, pp. 322–323. Springer Berlin Heidelberg, Berlin, Heidelberg,
684 2011. ISBN 978-3-642-04898-2. doi: 10.1007/978-3-642-04898-2_197. URL [https://doi.
685 org/10.1007/978-3-642-04898-2_197](https://doi.org/10.1007/978-3-642-04898-2_197).
- 686 Chi Jin, Rong Ge, Praneeth Netrapalli, Sham M. Kakade, and Michael I. Jordan. How to escape
687 saddle points efficiently. In Doina Precup and Yee Whye Teh (eds.), *Proceedings of the 34th
688 International Conference on Machine Learning*, volume 70 of *Proceedings of Machine Learning
689 Research*, pp. 1724–1732. PMLR, 06–11 Aug 2017. URL [https://proceedings.mlr.
690 press/v70/jin17a.html](https://proceedings.mlr.press/v70/jin17a.html).
- 691 Kenji Kawaguchi. Deep learning without poor local minima. In D. Lee, M. Sugiyama, U. Luxburg,
692 I. Guyon, and R. Garnett (eds.), *Advances in Neural Information Processing Systems*, volume 29.
693 Curran Associates, Inc., 2016. URL [https://proceedings.neurips.cc/paper_
694 files/paper/2016/file/f2fc990265c712c49d51a18a32b39f0c-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2016/file/f2fc990265c712c49d51a18a32b39f0c-Paper.pdf).
- 695 Kenji Kawaguchi and Yoshua Bengio. Depth with nonlinearity creates no bad local minima in
696 resnets. *Neural Networks*, 118:167–174, 2019. ISSN 0893-6080. doi: [https://doi.org/10.1016/
697 j.neunet.2019.06.009](https://doi.org/10.1016/j.neunet.2019.06.009). URL [https://www.sciencedirect.com/science/article/
698 pii/S0893608019301820](https://www.sciencedirect.com/science/article/pii/S0893608019301820).

- 702 Diederik Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint*
703 *arXiv:1412.6980*, 2014.
704
- 705 Yann LeCun. *PhD thesis: Modeles connexionnistes de l'apprentissage (connectionist learning*
706 *models)*. PhD thesis, Université Pierre et Marie Curie, Paris, France, 1987. URL <https://api.semanticscholar.org/CorpusID:151887454>.
707
- 708 Kfir Yehuda Levy. The power of normalization: Faster evasion of saddle points. *ArXiv*,
709 [abs/1611.04831](https://api.semanticscholar.org/CorpusID:16706102), 2016. URL <https://api.semanticscholar.org/CorpusID:16706102>.
710
711
- 712 Hao Li, Zheng Xu, Gavin Taylor, Christoph Studer, and Tom Goldstein. Visualizing the loss land-
713 scape of neural nets. In S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and
714 R. Garnett (eds.), *Advances in Neural Information Processing Systems*, volume 31. Curran Asso-
715 ciates, Inc., 2018. URL [https://proceedings.neurips.cc/paper_files/paper/](https://proceedings.neurips.cc/paper_files/paper/2018/file/a41b3bb3e6b050b6c9067c67f663b915-Paper.pdf)
716 [2018/file/a41b3bb3e6b050b6c9067c67f663b915-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2018/file/a41b3bb3e6b050b6c9067c67f663b915-Paper.pdf).
- 717 Hong Liu, Zhiyuan Li, David Leo Wright Hall, Percy Liang, and Tengyu Ma. Sophia: A scalable
718 stochastic second-order optimizer for language model pre-training. In *The Twelfth International*
719 *Conference on Learning Representations*, 2024. URL [https://openreview.net/forum?](https://openreview.net/forum?id=3xHDeA8Noi)
720 [id=3xHDeA8Noi](https://openreview.net/forum?id=3xHDeA8Noi).
- 721
722 Kenneth L. Manders and Leonard Adleman. Np-complete decision problems for binary quadrat-
723 ics. *Journal of Computer and System Sciences*, 16(2):168–184, 1978. ISSN 0022-0000.
724 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.](https://www.sciencedirect.com/science/article/pii/0022000078900442)
725 [com/science/article/pii/0022000078900442](https://www.sciencedirect.com/science/article/pii/0022000078900442).
- 726 James Martens. Deep learning via hessian-free optimization. In Johannes Fürnkranz and Thorsten
727 Joachims (eds.), *ICML*, pp. 735–742. Omnipress, 2010. URL [http://dblp.uni-trier.](http://dblp.uni-trier.de/db/conf/icml/icml2010.html#Martens10)
728 [de/db/conf/icml/icml2010.html#Martens10](http://db/conf/icml/icml2010.html#Martens10).
- 729
730 James Martens. New insights and perspectives on the natural gradient method. *J. Mach. Learn. Res.*,
731 21(1), jan 2020. ISSN 1532-4435.
- 732 James Martens and Roger Grosse. Optimizing neural networks with kronecker-factored approximate
733 curvature. In Francis Bach and David Blei (eds.), *Proceedings of the 32nd International Conference*
734 *on Machine Learning*, volume 37 of *Proceedings of Machine Learning Research*, pp. 2408–2417,
735 Lille, France, 07–09 Jul 2015. PMLR. URL [https://proceedings.mlr.press/v37/](https://proceedings.mlr.press/v37/martens15.html)
736 [martens15.html](https://proceedings.mlr.press/v37/martens15.html).
- 737 Harshal Mittal, Kartikey Pandey, and Yash Kant. Iclr reproducibility challenge report (padam :
738 Closing the generalization gap of adaptive gradient methods in training deep neural networks).
739 *ArXiv*, [abs/1901.09517](https://api.semanticscholar.org/CorpusID:249647677), 2019. URL [https://api.semanticscholar.org/CorpusID:](https://api.semanticscholar.org/CorpusID:249647677)
740 [249647677](https://api.semanticscholar.org/CorpusID:249647677).
- 741
742 Aatila Mustapha, Lachgar Mohamed, and Kartit Ali. Comparative study of optimization techniques
743 in deep learning: Application in the ophthalmology field. *Journal of Physics: Conference Series*,
744 1743, 2021. URL <https://api.semanticscholar.org/CorpusID:234179873>.
- 745 Yurii Nesterov and B. T. Polyak. Cubic regularization of newton method and its global performance.
746 *Mathematical Programming*, 108(1):177–205, 2006. doi: [10.1007/s10107-006-0706-8](https://doi.org/10.1007/s10107-006-0706-8). URL
747 <https://doi.org/10.1007/s10107-006-0706-8>.
748
- 749 Andrew Ng. Cs229 lecture notes - supervised learning. Available at
750 https://cs229.stanford.edu/lectures-spring2022/main_notes.pdf, 2012.
- 751 Quynh Nguyen, Mahesh Chandra Mukkamala, and Matthias Hein. On the loss landscape of a class
752 of deep neural networks with no bad local valleys. In *International Conference on Learning*
753 *Representations*, 2019. URL <https://openreview.net/forum?id=HJgXsjA5tQ>.
754
- 755 Jorge Nocedal and Stephen J. Wright. *Numerical Optimization*. Springer, New York, NY, USA, 2e
edition, 2006.

- 756 Thomas O’Leary-Roseberry, Nick Alger, and Omar Ghattas. Inexact newton methods for stochastic
757 non-convex optimization with applications to neural network training. *arXiv: Optimization and*
758 *Control*, 2019. URL <https://api.semanticscholar.org/CorpusID:155100112>.
- 759 P.J. Olver and C. Shakiban. *Applied Linear Algebra*. Prentice Hall, 2006. ISBN 9780131473829.
760 URL <https://books.google.co.il/books?id=D2tyQgAACAAJ>.
- 761
762 OpenAI, Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni
763 Aleman, Diogo Almeida, Janko Altschmidt, Sam Altman, Shyamal Anadkat, Red Avila, Igor
764 Babuschkin, Suchir Balaji, Valerie Balcom, Paul Baltescu, Haiming Bao, Mohammad Bavarian,
765 Jeff Belgum, Irwan Bello, Jake Berdine, Gabriel Bernadett-Shapiro, Christopher Berner, Lenny
766 Bogdonoff, Oleg Boiko, Madelaine Boyd, Anna-Luisa Brakman, Greg Brockman, Tim Brooks,
767 Miles Brundage, Kevin Button, Trevor Cai, Rosie Campbell, Andrew Cann, Brittany Carey, Chelsea
768 Carlson, Rory Carmichael, Brooke Chan, Che Chang, Fotis Chantzis, Derek Chen, Sully Chen,
769 Ruby Chen, Jason Chen, Mark Chen, Ben Chess, Chester Cho, Casey Chu, Hyung Won Chung,
770 Dave Cummings, Jeremiah Currier, Yunxing Dai, Cory Decareaux, Thomas Degry, Noah Deutsch,
771 Damien Deville, Arka Dhar, David Dohan, Steve Dowling, Sheila Dunning, Adrien Ecoffet, Atty
772 Eleti, Tyna Eloundou, David Farhi, Liam Fedus, Niko Felix, Simón Posada Fishman, Juston Forte,
773 Isabella Fulford, Leo Gao, Elie Georges, Christian Gibson, Vik Goel, Tarun Gogineni, Gabriel
774 Goh, Rapha Gontijo-Lopes, Jonathan Gordon, Morgan Grafstein, Scott Gray, Ryan Greene, Joshua
775 Gross, Shixiang Shane Gu, Yufei Guo, Chris Hallacy, Jesse Han, Jeff Harris, Yuchen He, Mike
776 Heaton, Johannes Heidecke, Chris Hesse, Alan Hickey, Wade Hickey, Peter Hoeschele, Brandon
777 Houghton, Kenny Hsu, Shengli Hu, Xin Hu, Joost Huizinga, Shantanu Jain, Shawn Jain, Joanne
778 Jang, Angela Jiang, Roger Jiang, Haozhun Jin, Denny Jin, Shino Jomoto, Billie Jonn, Heewoo
779 Jun, Tomer Kaftan, Łukasz Kaiser, Ali Kamali, Ingmar Kanitscheider, Nitish Shirish Keskar,
780 Tabarak Khan, Logan Kilpatrick, Jong Wook Kim, Christina Kim, Yongjik Kim, Jan Hendrik
781 Kirchner, Jamie Kiros, Matt Knight, Daniel Kokotajlo, Łukasz Kondraciuk, Andrew Kondrich,
782 Aris Konstantinidis, Kyle Kosic, Gretchen Krueger, Vishal Kuo, Michael Lampe, Ikai Lan, Teddy
783 Lee, Jan Leike, Jade Leung, Daniel Levy, Chak Ming Li, Rachel Lim, Molly Lin, Stephanie
784 Lin, Mateusz Litwin, Theresa Lopez, Ryan Lowe, Patricia Lue, Anna Makanju, Kim Malfacini,
785 Sam Manning, Todor Markov, Yaniv Markovski, Bianca Martin, Katie Mayer, Andrew Mayne,
786 Bob McGrew, Scott Mayer McKinney, Christine McLeavey, Paul McMillan, Jake McNeil, David
787 Medina, Aalok Mehta, Jacob Menick, Luke Metz, Andrey Mishchenko, Pamela Mishkin, Vinnie
788 Monaco, Evan Morikawa, Daniel Mossing, Tong Mu, Mira Murati, Oleg Murk, David Mély,
789 Ashvin Nair, Reiichiro Nakano, Rajeev Nayak, Arvind Neelakantan, Richard Ngo, Hyeonwoo
790 Noh, Long Ouyang, Cullen O’Keefe, Jakub Pachocki, Alex Paino, Joe Palermo, Ashley Pantuliano,
791 Giambattista Parascandolo, Joel Parish, Emy Parparita, Alex Passos, Mikhail Pavlov, Andrew Peng,
792 Adam Perelman, Filipe de Avila Belbute Peres, Michael Petrov, Henrique Ponde de Oliveira Pinto,
793 Michael, Pokorný, Michelle Pokrass, Vitchyr H. Pong, Tolly Powell, Alethea Power, Boris Power,
794 Elizabeth Proehl, Raul Puri, Alec Radford, Jack Rae, Aditya Ramesh, Cameron Raymond, Francis
795 Real, Kendra Rimbach, Carl Ross, Bob Rotsted, Henri Roussez, Nick Ryder, Mario Saltarelli, Ted
796 Sanders, Shibani Santurkar, Girish Sastry, Heather Schmidt, David Schnurr, John Schulman, Daniel
797 Selsam, Kyla Sheppard, Toki Sherbakov, Jessica Shieh, Sarah Shoker, Pranav Shyam, Szymon
798 Sidor, Eric Sigler, Maddie Simens, Jordan Sitkin, Katarina Slama, Ian Sohl, Benjamin Sokolowsky,
799 Yang Song, Natalie Staudacher, Felipe Petroski Such, Natalie Summers, Ilya Sutskever, Jie Tang,
800 Nikolas Tezak, Madeleine B. Thompson, Phil Tillet, Amin Tootoonchian, Elizabeth Tseng, Preston
801 Tuggle, Nick Turley, Jerry Tworek, Juan Felipe Cerón Uribe, Andrea Vallone, Arun Vijayvergiya,
802 Chelsea Voss, Carroll Wainwright, Justin Jay Wang, Alvin Wang, Ben Wang, Jonathan Ward, Jason
803 Wei, CJ Weinmann, Akila Welihinda, Peter Welinder, Jiayi Weng, Lilian Weng, Matt Wiethoff,
804 Dave Willner, Clemens Winter, Samuel Wolrich, Hannah Wong, Lauren Workman, Sherwin Wu,
805 Jeff Wu, Michael Wu, Kai Xiao, Tao Xu, Sarah Yoo, Kevin Yu, Qiming Yuan, Wojciech Zaremba,
806 Rowan Zellers, Chong Zhang, Marvin Zhang, Shengjia Zhao, Tianhao Zheng, Juntang Zhuang,
807 William Zhuk, and Barret Zoph. Gpt-4 technical report, 2024.
- 808 Panos M. Pardalos and Stephen A. Vavasis. Quadratic programming with one negative eigenvalue
809 is np-hard. *Journal of Global Optimization*, 1(1):15–22, 1991. doi: 10.1007/BF00120662. URL
<https://doi.org/10.1007/BF00120662>.
- 808 Dylan Patel and Gerald Wong. Gpt-4 architecture, infrastructure, training dataset, costs, vision, moe.
809 <https://www.semianalysis.com/p/gpt-4-architecture-infrastructure>,
2023. Accessed: 2024-05-01.

- 810 Barak A. Pearlmutter. Fast Exact Multiplication by the Hessian. *Neural Computation*, 6(1):147–160,
811 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)
812 [1162/neco.1994.6.1.147](https://doi.org/10.1162/neco.1994.6.1.147).
- 813
- 814 Karl Pearson. Note on Regression and Inheritance in the Case of Two Parents. *Proceedings of the*
815 *Royal Society of London Series I*, 58:240–242, January 1895.
- 816
- 817 Sebastian Ruder. An overview of gradient descent optimization algorithms. *ArXiv*, abs/1609.04747,
818 2016. URL <https://api.semanticscholar.org/CorpusID:17485266>.
- 819
- 820 Levent Sagun, Léon Bottou, and Yann LeCun. Eigenvalues of the hessian in deep learning: Singu-
821 larity and beyond. *arXiv: Learning*, 2016. URL [https://api.semanticscholar.org/](https://api.semanticscholar.org/CorpusID:35723845)
822 [CorpusID:35723845](https://api.semanticscholar.org/CorpusID:35723845).
- 823
- 824 Robin M Schmidt, Frank Schneider, and Philipp Hennig. Descending through a crowded valley -
825 benchmarking deep learning optimizers. In Marina Meila and Tong Zhang (eds.), *Proceedings of*
826 *the 38th International Conference on Machine Learning*, volume 139 of *Proceedings of Machine*
827 *Learning Research*, pp. 9367–9376. PMLR, 18–24 Jul 2021. URL [https://proceedings.](https://proceedings.mlr.press/v139/schmidt21a.html)
828 [mlr.press/v139/schmidt21a.html](https://proceedings.mlr.press/v139/schmidt21a.html).
- 829
- 830 Nicol Schraudolph. Fast curvature matrix-vector products for second-order gradient descent. *Neural*
831 *computation*, 14:1723–38, 08 2002. doi: 10.1162/08997660260028683.
- 832
- 833 Cooper Simpson and Jaden Wang. PyTorch-ARC. github.com/RS-Coop/PyTorch-ARC,
834 2023. Adaptive Regularization with Cubics (ARC) optimizer for PyTorch.
- 835
- 836 Hadar Sivan, Moshe Gabel, and Assaf Schuster. FOSI: Hybrid first and second order optimization.
837 In *The Twelfth International Conference on Learning Representations*, 2024. URL [https:](https://openreview.net/forum?id=NvbeD9Ttkx)
838 [//openreview.net/forum?id=NvbeD9Ttkx](https://openreview.net/forum?id=NvbeD9Ttkx).
- 839
- 840 Daniel Soudry and Yair Carmon. No bad local minima: Data independent training error guar-
841 antees for multilayer neural networks. *ArXiv*, abs/1605.08361, 2016. URL [https://api.](https://api.semanticscholar.org/CorpusID:3029264)
842 [semanticscholar.org/CorpusID:3029264](https://api.semanticscholar.org/CorpusID:3029264).
- 843
- 844 Tijmen Tieleman and Geoffrey Hinton. Lecture 6.5-rmsprop: Divide the gradient by a running
845 average of its recent magnitude. *COURSERA: Neural networks for machine learning*, 4(2):26–31,
846 2012.
- 847
- 848 Philippe L. Toint. Nonlinear stepsize control, trust regions and regularizations for unconstrained
849 optimization. *Optimization Methods and Software*, 28(1):82–95, 2013. doi: 10.1080/10556788.
850 2011.610458. URL <https://doi.org/10.1080/10556788.2011.610458>.
- 851
- 852 J.F. Traub. *Iterative Methods for the Solution of Equations*. AMS Chelsea Publishing Series.
853 Chelsea, 1982. ISBN 9780828403122. URL [https://books.google.co.il/books?](https://books.google.co.il/books?id=se3YdgFgz4YC)
854 [id=se3YdgFgz4YC](https://books.google.co.il/books?id=se3YdgFgz4YC).
- 855
- 856 Sharan Vaswani, Reza Babanezhad, Jose Gallego, Aaron Mishkin, Simon Lacoste-Julien, and
857 Nicolas Le Roux. To each optimizer a norm, to each norm its generalization. *ArXiv*, abs/2006.06821,
858 2020. URL <https://api.semanticscholar.org/CorpusID:219636073>.
- 859
- 860 Neha Wadia, Daniel Duckworth, Samuel S Schoenholz, Ethan Dyer, and Jascha Sohl-Dickstein.
861 Whitening and second order optimization both make information in the dataset unusable during
862 training, and can reduce or prevent generalization. In Marina Meila and Tong Zhang (eds.),
863 *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of *Pro-*
864 *ceedings of Machine Learning Research*, pp. 10617–10629. PMLR, 18–24 Jul 2021. URL
865 <https://proceedings.mlr.press/v139/wadia21a.html>.
- 866
- 867 Xiao Wang, Shiqian Ma, Donald Goldfarb, and Wei Liu. Stochastic quasi-newton methods for
868 nonconvex stochastic optimization. *SIAM Journal on Optimization*, 27, 07 2016. doi: 10.1137/
869 15M1053141.

- 864 Rachel Ward, Xiaoxia Wu, and Leon Bottou. AdaGrad stepsizes: Sharp convergence over nonconvex
865 landscapes. In Kamalika Chaudhuri and Ruslan Salakhutdinov (eds.), *Proceedings of the 36th*
866 *International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning*
867 *Research*, pp. 6677–6686. PMLR, 09–15 Jun 2019. URL [https://proceedings.mlr.](https://proceedings.mlr.press/v97/ward19a.html)
868 [press/v97/ward19a.html](https://proceedings.mlr.press/v97/ward19a.html).
- 869 Peng Xu, Farbod Roosta-Khorasan, and Michael Mahoney. Second-order optimization for non-convex
870 machine learning: An empirical study. In *Proceedings of the 2020 SIAM International Conference*
871 *on Data Mining*, 08 2017.
- 872 Peng Xu, Fred Roosta, and Michael W. Mahoney. Newton-type methods for non-convex optimization
873 under inexact hessian information. *Mathematical Programming*, 184(1):35–70, 2020. doi: 10.1007/
874 s10107-019-01405-z. URL <https://doi.org/10.1007/s10107-019-01405-z>.
- 875 Robert M. Young. 75.9 euler’s constant. *The Mathematical Gazette*, 75(472):187–190, 1991. ISSN
876 00255572. URL <http://www.jstor.org/stable/3620251>.
- 877 Xiao-Hu Yu and Guo-An Chen. On the local minima free condition of backpropagation learning.
878 *IEEE Transactions on Neural Networks*, 6(5):1300–1303, 1995. doi: 10.1109/72.410380.
- 879 Matthew D. Zeiler. Adadelta: An adaptive learning rate method. *ArXiv*, abs/1212.5701, 2012. URL
880 <https://api.semanticscholar.org/CorpusID:7365802>.
- 881 Pan Zhou, Jiashi Feng, Chao Ma, Caiming Xiong, Steven Chu Hong Hoi, and Weinan E.
882 Towards theoretically understanding why sgd generalizes better than adam in deep learn-
883 ing. In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin (eds.), *Ad-*
884 *vances in Neural Information Processing Systems*, volume 33, pp. 21285–21296. Curran Asso-
885 ciates, Inc., 2020. URL [https://proceedings.neurips.cc/paper_files/paper/](https://proceedings.neurips.cc/paper_files/paper/2020/file/f3f27a324736617f20abbf2ffd806f6d-Paper.pdf)
886 [2020/file/f3f27a324736617f20abbf2ffd806f6d-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2020/file/f3f27a324736617f20abbf2ffd806f6d-Paper.pdf).
- 887 Milija Županski. A preconditioning algorithm for large-scale minimization problems. *Tellus A:*
888 *Dynamic Meteorology and Oceanography*, Jan 1993. doi: 10.3402/tellusa.v45i5.15048.

894 A QUASI-NEWTON CHALLENGES AND PROPOSED SOLUTIONS

895 Some of the challenges involved in training neural networks include:

- 896 • Because the models and data often have very complex structures, obtaining precisely optimal
897 parameters is often computationally prohibitive. As a result, one must satisfy oneself with
898 a small degree of suboptimality in the model’s parameters, chosen to be small enough to
899 satisfy one’s needs while not exhausting the computational capacity at hand.
- 900 • Since many model architectures (e.g. artificial neural networks) have very complex struc-
901 tures, the loss function is generally non-convex as a function of the model’s parameters.
902 This makes finding the globally optimal choice of parameters an NP-hard problem (Pardalos
903 & Vavasis, 1991; Manders & Adleman, 1978). As a result, one must satisfy oneself with
904 merely a local minimum of the loss function (that is, a point at which the norm of the gradient
905 w.r.t. the model parameters is zero, and the function is locally convex, or equivalently, the
906 Hessian is positive semi-definite). This is often considered sufficient (see Soudry & Carmon
907 (2016); Kawaguchi & Bengio (2019); Kawaguchi (2016); Nguyen et al. (2019)), however
908 this does not apply to saddle points and local maxima (points at which the gradient norm is
909 zero but the Hessian is not positive definite). Although some work has been done on trying to
910 eliminate this problem by eliminating local- but not global- minima via overparameterization
911 (Yu & Chen, 1995), further work (Ding et al., 2019) has shown that this does not scale to
912 deep neural networks.
- 913 • As mentioned previously, there is no single universally optimal optimizer, even among
914 existing optimizers.

915 As a result, sophisticated optimizers are necessary to contend with different neural network training
916 scenarios. Restricting ourselves to quasi-newton optimization algorithms, scenarios with $n \gg 0$
917

(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

$$\Phi_t \succeq 0 \quad (11)$$

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 $(\tilde{\lambda}_j, \tilde{v}_j)$, $\tilde{v}_j = v_j$, $\tilde{\lambda}_j =$

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

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

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

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

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

B OTHER CONVERGENCE RATE MEASURES

Convergence rates to first-order criticality Most works on convergence rates in the non-convex regime bound the number of iterations necessary to achieve first-order criticality ($\|\nabla f(\theta_t)\|_2 = 0$) by means of finding an ϵ_g -stationary point (a point at which $\|\nabla f(\theta_t)\|_2 \leq \epsilon_g$). The seminal work Wang et al. (2016) provide a convergence rate bound for general optimizers (with very weak assumptions) in 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 $\nu \in (0.5, 1)$.

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

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

Convergence rate dependence on preconditioner quality One possible quality metric for Φ_t 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, Chapter 7.1) assume $\sup_t(\eta_t) < 1$ and prove that first-order criticality may be reached within $\mathcal{O}\left(\frac{\log \epsilon}{\log \frac{1 + \sup_t(\eta_t)}{2}}\right)$ 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 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.

C.1 COMPARISON TO CAUCHY’S METHOD

Cauchy’s method (Traub, 1982) is nearly identical to ELMO:

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

The difference between our minimization step and their step is merely the sign on the squareroot. The difference lies in 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.

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

$$\begin{aligned}
\lim_{L_H \rightarrow 0^+} \Delta\theta_t^{*\top} v_i &= \lim_{L_H \rightarrow 0^+} -\frac{2\nabla f(\theta_t)^\top v_i}{\sqrt{\lambda_i^2(\theta_t) - 2L_t^i \cdot \nabla f(\theta_t)^\top v_i} + \lambda_i(\theta_t)} \\
&= \begin{cases} -\frac{\nabla f(\theta_t)^\top v_i}{\lambda_i} & \lambda_i > 0 \\ \infty & \lambda_i < 0 \end{cases} \quad (12)
\end{aligned}$$

D CONVERGENCE RATE DEPENDENCE ON HESSIAN-LIPSCHITZ PARAMETER

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

D.1 LIPSCHITZ ROBUSTNESS

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

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

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

when $L_t^i \rightarrow \infty$

Proof.

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

□

D.2 BENEFIT OF LIPSCHITZ TIGHTNESS

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

Theorem D.2. *Benefit of Lipschitz tightness: concave subspaces*

1134 Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ satisfying assumptions 1 and 2. If $\lambda_i(\theta_t) \leq 0$ then

$$1137 \quad M_t^i(\Delta\theta_t^{*\top} v_i) = \Theta\left(-\frac{1}{L_t^{i,2}}\right)$$

1141 when $L_t^i \rightarrow 0^+$

1148 *Proof.*

$$\begin{aligned}
1151 & \lim_{L_t^i \rightarrow 0^+} \frac{M_t^i(\Delta\theta_t^{*\top} v_i)}{\frac{\frac{2}{3}\lambda_i^3(\theta_t)}{L_t^{i,2}}} \\
1152 & = \lim_{L_t^i \rightarrow 0^+} \left(3 \left(\frac{\sqrt{1 - 2L_t^i \cdot \frac{\nabla f(\theta_t)^\top v_i}{(-\lambda_i(\theta_t))^2}} + 1}{2} \right)^2 - 2 \left(\frac{\sqrt{1 - 2L_t^i \cdot \frac{\nabla f(\theta_t)^\top v_i}{(-\lambda_i(\theta_t))^2}} + 1}{2} \right)^3 \right) \\
1153 & \quad + L_t^i \cdot \frac{3}{2\lambda_i^3(\theta_t)} \cdot \left(\sqrt{\lambda_i^2(\theta_t) - 2L_t^i \cdot \nabla f(\theta_t)^\top v_i} - \lambda_i(\theta_t) \right) \cdot \nabla f(\theta_t)^\top v_i \\
1154 & = 1
\end{aligned}$$

1164 □

1170 **Theorem D.3.** *Benefit of Lipschitz tightness: convex subspaces*

1171 Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ satisfying assumptions 1 and 2. If $\lambda_i(\theta_t) > 0$ then

$$1174 \quad M_t^i(\Delta\theta_t^{*\top} v_i) - M_t^i\left(\lim_{L_t^i \rightarrow 0^+} \Delta\theta_t^{*\top} v_i\right) = \Theta(L_t^i)$$

1179 when $L_t^i \rightarrow 0^+$

1182 *Proof.* We begin by noting that by equation 12, $\lim_{L_t^i \rightarrow 0^+} \Delta\theta_t^{*\top} v_i = \frac{|\nabla f(\theta_t)^\top v_i|}{\lambda_i(\theta_t)}$. Plugging this into M_t^i :

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

□

E EVIDENCE FROM THE LITERATURE

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

F LIPSCHITZ PARAMETER EXPERIMENTS

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

- 1242 • $\sigma_0^+ = \sigma_0^- = 1$
- 1243 • $\eta_1 = 0.1, \eta_2 = 0.9$
- 1244 • $\gamma_1 = \gamma_2 = 2$
- 1245 • Maximum sub-problem failures = 11
- 1246 • Maximum sub-problem iterations = 50,
- 1247 • Sub-problem tolerance = 10^{-6}
- 1248 • Lanczos eigendecomposition (Garber et al., 2016)
- 1249 • BFGS trinomial sub-problem solver (Nocedal & Wright, 2006, Chapter 6.1), (Feinman, 2021)
- 1250 • trinomial sub-problem maximal failures=11
- 1251 • trinomial sub-problem maximal iterations=50

1252 The test scenarios¹ include combinations of:

- 1253 • Training ResNet18 artificial neural networks (He et al., 2015) for image classification on MNIST, CIFAR10, and FakeData (random noise in place of images) with Cross-Entropy Loss (CELoss), to evaluate the effect of changing data on the Lipschitz parameters
- 1254 • Training a CNN for image classification on CIFAR10 with CELoss, to evaluate the effect of changing neural network architecture on the Lipschitz parameters.
- 1255 • Training a CNN² autoencoder³ (LeCun, 1987) to compress MNIST

1256 The classification CNN architecture:

- 1257 1. A feature extractor consisting of 2 2D convolutional layers with 6 output channels for the first and 16 output channels for the second. Both had kernel sizes of 5 pixels. Each of these is followed by a ReLU nonlinearity and then 2x2 2D max pooling
- 1258 2. A 3-layer MLP classification head with hidden sizes (120, 84) and ReLU nonlinearities

1259 The autoencoder CNN architecture:

- 1260 • Encoder: 4 2D convolutional layers with respective output channel numbers (16,32,32,64) and kernel sizes of 3 pixels for the first two and 5 pixels for the last two. The first two have 1 pixel padding and the last two have 2 pixel paddings. After every layer we apply LeakyReLU nonlinearity and after every 2 layers we apply 2x2 2D max pooling.
- 1261 • Decoder: 4 composite layers consisting of
 - 1262 1. a 2D transpose-convolutional layer
 - 1263 2. LeakyReLU nonlinearity (only for first and third composite layers)
 - 1264 3. a 2d convolutional layer
 - 1265 4. LeakyReLU nonlinearity
- 1266 • Decoder hidden channel sizes: 32-32-16-16-16-16-3
- 1267 • Decoder kernel sizes: 2-5-5-5-2-3-5-3
- 1268 • All strides are of size 1, except for the first and third transpose-convolutional layers of the decoder, with stride of size 2
- 1269 • Decoder paddings: 0-2-2-2-0-1-2-1

1270 Each experiment took several hours to run. All experiments (including those from above) shown in figure 3.

1271 ¹Code available on Github at REDACTED

1272 ²convolutional neural network

1273 ³With hidden dimensions 128-64-36-18-9-18-36-64-128, ReLU nonlinearities, and sigmoid nonlinearity on the output

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

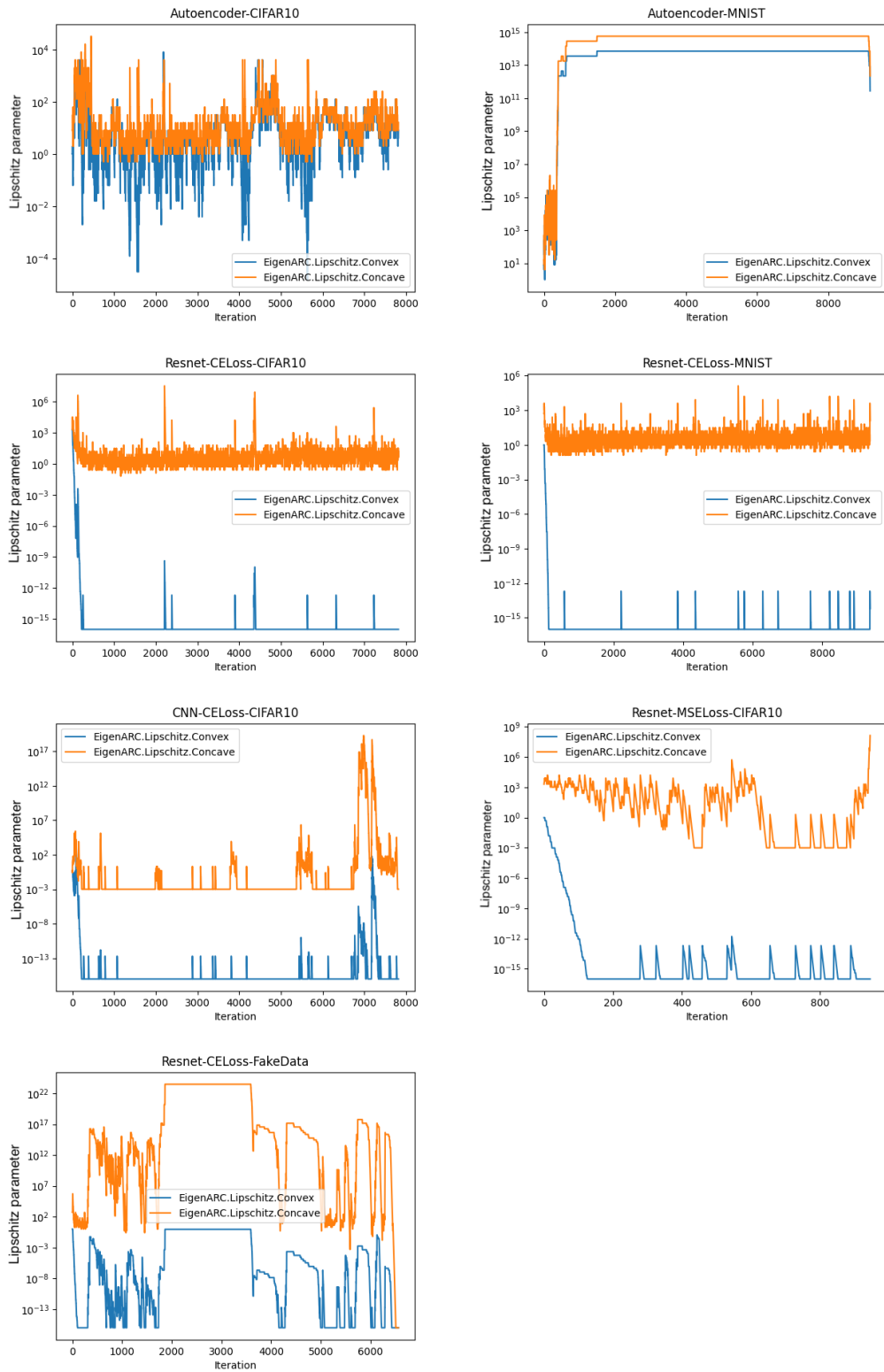


Figure 3: Comparisons of convex-subspace Lipschitz parameters to concave-subspace Lipschitz parameters. *Logarithmic scale*

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

For visual clarity, the quasi-Newton superiority measurements in 2 are presented after:

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

F.1 COMPUTATION OF LIPSCHITZ PARAMETERS

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

Algorithm 2 Algorithm EigenARC

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

```

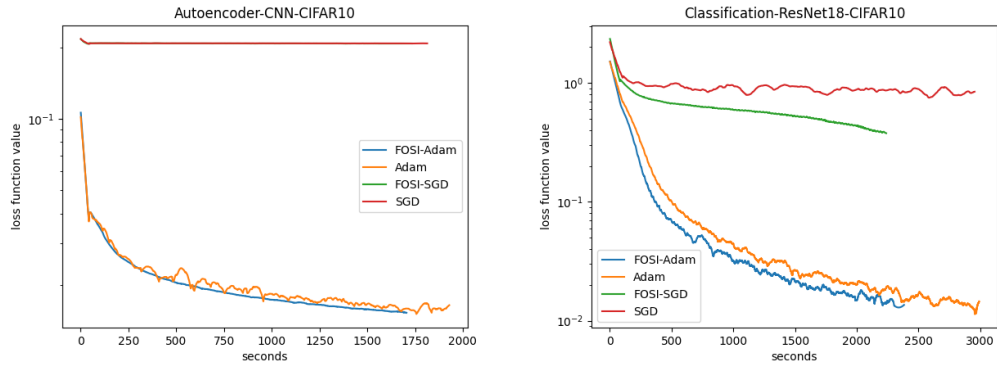
1:  $t \leftarrow 0$ 
2: while  $\|\nabla f(\theta_t)\|_2 > \epsilon$  do ▷ While BASE_OPT hasn't converged yet
3:    $(\lambda_i, v_i)_{i=1}^n \leftarrow \text{EIGEN}(\mathcal{H}(\theta_t))$ 
4:   if  $\text{ASSESS\_LIPSCHITZ}((L_t^i)_{i=1}^n) > \eta_2$  then ▷ Overly conservative  $L_t^i$ 
5:      $(L_t^i)_{i=1}^n \leftarrow \gamma_2 \cdot (L_t^i)_{i=1}^n$ 
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  $(L^i)_{i=1, \hat{i}=1}^{n, t}$ 

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

```

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

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



(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

Figure 4: Comparison of second-order optimizers against first-order optimizers in settings with different sized convex Lipschitz parameters. Second-order optimizers only hold an advantage over first-order optimizers (thus justifying their additional computational complexity) when the convex Lipschitz parameters are small.

G LIPSCHITZ-AIDED OPTIMIZER SELECTION

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.

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 4.

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.

H PROOFS

H.1 PRELIMINARY LEMMAS

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

1458 **Lemma H.1.**

$$1459 \quad \forall_{x \geq -1} : \sqrt{1+x} \leq 1 + \frac{x}{2}$$

1460 *Proof.* Mark $g : \mathbb{R} \rightarrow \mathbb{R}, g(x) = 1 + \frac{x}{2} - \sqrt{1+x}$. We have

$$1464 \quad g'(x) = \frac{1}{2} \left(1 - \frac{1}{\sqrt{1+x}} \right)$$

$$1466 \quad g''(x) = \frac{1}{4} \left(\frac{1}{(1+x)^{\frac{3}{2}}} \right)$$

1468 g is convex due to its second derivative being positive for all $x > -1$. Therefore, its sole critical
1469 point $x = 0$ obtained from the derivative is a minimum, and $\forall_{x \geq -1} : g(x) \geq g(0) = 0$ \square

1471 **Corollary H.1.1.**

$$1472 \quad \forall_{x \in \mathbb{R}^+} \forall_{y \geq -x} : \sqrt{x+y} \leq \sqrt{x} + \frac{y}{2\sqrt{x}}$$

1475 *Proof.*

$$1477 \quad \sqrt{x+y} = \sqrt{x} \sqrt{1 + \frac{y}{x}} \leq \sqrt{x} \left(1 + \frac{y}{2x} \right) = \sqrt{x} + \frac{y}{2\sqrt{x}}$$

1480 \square

1481 **Lemma.** Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ satisfy assumptions 1 and 2. Then

$$1483 \quad m_t^i (\Delta \theta_t^{*\top} v_i) \leq M_t^i (\Delta \theta_t^{*\top} v_i) \leq 0 \quad (13)$$

1485 *Proof.* The first inequality stems from the trivial fact that $m_t^i \leq M_t^i$.

1487 The second inequality follows from the fact that (by design), $\Delta \theta_t^{*\top} v_i$ is a minimizer of $M_t^i (\Delta \theta_t^\top v_i)$,
1488 but

$$1489 \quad M_t^i (0) = 0$$

1491 \square

1492 **Lemma.** 4.1 Minmax preconditioner

$$1493 \quad \arg \min_{\Phi_t \in \mathbb{R}^{n \times n}} |\Delta \Delta^i \theta_t (\Phi_t)| = \left(\frac{\mathcal{H}(\theta_t) + \sqrt{(\mathcal{H}(\theta_t))^2 + 2V \cdot \text{diag} \left(L_t^i \cdot |\nabla f(\theta_t)^\top v_i| \right)_{i=1}^n \cdot V^\top}}{2} \right)^{-1}$$

1500 *Proof.*

$$1502 \quad |\Delta \Delta^i \theta_t| = \left| \nabla f(\theta_t)^\top \left(\Phi_t - \frac{2}{\lambda_i + \sqrt{\lambda_i^2 - 2L_t^i \cdot \nabla f(\theta_t)^\top v_i}} I \right) v_i \right|$$

$$1503 \quad = \left| \nabla f(\theta_t)^\top \left(\Phi_t - \left(\frac{\mathcal{H}(\theta_t) + \sqrt{(\mathcal{H}(\theta_t))^2 + 2L_t^i \cdot |\nabla f(\theta_t)^\top v_i|} I}{2} \right)^{-1} \right) v_i \right|$$

1509 and the result follows from developing the second parenthesized term for all n dimensions of the
1511 domain space. \square

1512 H.2 LEMMA 3.1: EIGENSPACE DESCENT
 1513
 1514

1515 Working with assumptions 1 and equation 2, Dennis & Schnabel (1983, Lemma 4.1.14) prove the
 1516 following:
 1517

1518 **Lemma H.2.**
 1519

$$1520 f(\theta_{t+1}) - f(\theta_t) \leq \nabla f(\theta_t)^T \cdot \Delta\theta_t + \frac{1}{2} \Delta\theta_t^T \mathcal{H}(\theta_t) \Delta\theta_t + \frac{1}{6} L_H \|\Delta\theta_t\|_2^3 \quad (14)$$

1521
 1522 Much like algorithm ELMO, minimizing equation 14 would maximize [a bound on] the descent given
 1523 by iteration t . However, previous works such as Nesterov & Polyak (2006) note the difficulty of
 1524 minimizing this 3rd-order n -dimensional polynomial, even when L_H is known. Indeed, Cartis et al.
 1525 (2011a) propose minimizing it iteratively over a growing subspace, with each iteration’s minimization
 1526 subspace a superset of the previous iterations’ (in practice, they use the Hessian’s Krylov subspaces,
 1527 initialized with the gradient). In our theoretical analysis however, we have the freedom to simply
 1528 take the most natural decomposition of the space into subspaces, the eigenspaces of the Hessian.
 1529 This does not limit the practicality of our approach, however, since Lanczos methods allow one
 1530 to obtain elements of this decomposition. In fact, Sivan et al. (2024) demonstrate experimentally
 1531 that decomposing the parameter space into multiple eigenspaces and optimizing each separately can
 1532 significantly speed up optimization wall time, despite the additional computational burden of the
 1533 Lanczos iterations, because of the regularizing effect this has on the function in each of the subspaces
 1534 (by reducing the variance of the Hessian eigenvalues). Ghorbani et al. (2019) also show the benefits
 1535 of reducing this variance.
 1536
 1537
 1538
 1539
 1540
 1541

1542 **Lemma. 3.1 Eigenspace Descent Bounds**
 1543

1544 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
 1545 $\Delta\theta_t = \theta_{t+1} - \theta_t$, we have
 1546
 1547
 1548

$$1549 \exists (L_t^i)_{i=1}^n \in (\mathbb{R}^+)^n : f(\theta_{t+1}) - f(\theta_t) \leq \sum_{i=1}^n M_t^i (\Delta\theta_t^\top v_i) \quad (15)$$

$$1550 \exists (L_t^i)_{i=1}^n \in (\mathbb{R}^+)^n : f(\theta_{t+1}) - f(\theta_t) \geq \sum_{i=1}^n m_t^i (\Delta\theta_t^\top v_i) \quad (16)$$

1551
 1552
 1553
 1554
 1555
 1556
 1557
 1558
 1559 We give 2 proofs of the above lemma. The first proof is far simpler and relies on the standard spectral
 1560 norm-Lipschitz continuous Hessian assumption given by equation 2 instead of assumption 2:
 1561
 1562
 1563
 1564
 1565

1566 *Proof.* Beginning with Nesterov & Polyak (2006, Lemma 1) for the first inequality,
 1567

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

1594 with

- 1595 1. the second inequality being a representation of $\theta_{t+1} - \theta_t$ over the (orthogonal) Hessian
 1596 eigenbasis
- 1597 2. the third inequality due to the triangle inequality
- 1598 3. the fourth inequality due to Jensen's inequality

1609 □

1614 Our first proof of lemma 3.1 is simple, but leaves something to be desired due to its lack of per-
 1615 eigenspace Lipschitz parameters and due to the presence of n^2 in the bound, which can be very
 1616 large, as noted in section A. The first proof's assumption of equation 2 is also easily seen to be
 1617 no weaker than assumption 2 (meaning that assuming equation 2 implies assumption 2) by taking
 1618 $L_R \triangleq L_t^i \triangleq L_H$ for all $t \in [T], i \in [n]$. To address these concerns, we make use of assumption 2,
 1619 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.

1620 *Proof.*

$$\begin{aligned}
1621 & f(\theta_{t+1}) - f(\theta_t) \\
1622 & = \int_0^1 \nabla f(\theta_t + y(\theta_{t+1} - \theta_t))^\top (\theta_{t+1} - \theta_t) dy \\
1623 & = \nabla f(\theta_t)^\top (\theta_{t+1} - \theta_t) + \int_0^1 (\nabla f(\theta_t + y(\theta_{t+1} - \theta_t)) - \nabla f(\theta_t))^\top (\theta_{t+1} - \theta_t) dy \\
1624 & = \nabla f(\theta_t)^\top (\theta_{t+1} - \theta_t) + \int_0^1 \left(\int_0^1 y \mathcal{H}(\theta_t + yz(\theta_{t+1} - \theta_t)) (\theta_{t+1} - \theta_t) dz \right)^\top (\theta_{t+1} - \theta_t) dy \\
1625 & = \nabla f(\theta_t)^\top (\theta_{t+1} - \theta_t) + \int_0^1 \int_0^1 y (\theta_{t+1} - \theta_t)^\top \mathcal{H}(\theta_t + yz(\theta_{t+1} - \theta_t)) (\theta_{t+1} - \theta_t) dydz \\
1626 & = \nabla f(\theta_t)^\top (\theta_{t+1} - \theta_t) + (\theta_{t+1} - \theta_t)^\top \mathcal{H}(\theta_t) (\theta_{t+1} - \theta_t) \\
1627 & + \underbrace{\int_0^1 \int_0^1 y (\theta_{t+1} - \theta_t)^\top (\mathcal{H}(\theta_t + yz(\theta_{t+1} - \theta_t)) - \mathcal{H}(\theta_t)) (\theta_{t+1} - \theta_t) dydz}_Z
\end{aligned}$$

1628 with the first and third equalities due to the fundamental theorem of calculus.

1629 Mark the Hessian eigendecompositions as follows:

$$\begin{aligned}
1630 & \mathcal{H}(\theta_t + yz(\theta_{t+1} - \theta_t)) = \tilde{V} \tilde{\Lambda} \tilde{V}^\top \\
1631 & \mathcal{H}(\theta_t) = V \Lambda V^\top
\end{aligned}$$

1632 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 Hessian matrices) matrices V, \tilde{V} .

$$\begin{aligned}
1633 & Z = \int_0^1 \int_0^1 y (\theta_{t+1} - \theta_t)^\top (\tilde{V} \tilde{\Lambda} \tilde{V}^\top - V \Lambda V^\top) (\theta_{t+1} - \theta_t) dydz \\
1634 & = \int_0^1 \int_0^1 y (\theta_{t+1} - \theta_t)^\top (V \tilde{\Lambda} V^\top - V \Lambda V^\top) (\theta_{t+1} - \theta_t) dydz \\
1635 & + \int_0^1 \int_0^1 y (\theta_{t+1} - \theta_t)^\top (\tilde{V} \tilde{\Lambda} \tilde{V}^\top - V \tilde{\Lambda} V^\top) (\theta_{t+1} - \theta_t) dydz
\end{aligned}$$

1636 Focusing on the first term,

$$\begin{aligned}
1637 & = \int_0^1 \int_0^1 y (\theta_{t+1} - \theta_t)^\top V (\tilde{\Lambda} - \Lambda) V^\top (\theta_{t+1} - \theta_t) dydz \\
1638 & = \int_0^1 \int_0^1 y \sum_{j=1}^n \sum_{i=1}^n (\theta_{t+1} - \theta_t)^\top v_i \cdot (\theta_{t+1} - \theta_t)^\top v_j \cdot v_j^\top V (\tilde{\Lambda} - \Lambda) V^\top v_i dydz \\
1639 & = \int_0^1 \int_0^1 y \cdot \sum_{i=1}^n \left((\theta_{t+1} - \theta_t)^\top v_i \right)^2 \cdot (\tilde{\lambda}_i - \lambda_i) dydz \\
1640 & \leq \sum_{i=1}^n L_t^i \cdot \left| (\theta_{t+1} - \theta_t)^\top v_i \right|^3 \cdot \int_0^1 \int_0^1 y \cdot yz \cdot dydz \\
1641 & = \sum_{i=1}^n \frac{L_t^i}{6} \cdot \left| (\theta_{t+1} - \theta_t)^\top v_i \right|^3
\end{aligned}$$

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 H.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 H.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 H.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.

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

with

- the first inequality due to the triangle inequality
- the second equality due to lemma H.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 $\{-1, 1\}$

□

H.3 THEOREM 3.2: WORST CASE-OPTIMAL DESCENT RATE

Before we can prove theorem 3.2, we need to upper bound equation 7.

Lemma H.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

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

1782 *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 H.1.1 with $x = 2L_t^i \cdot |\nabla f(\theta_t)^\top v_i|$
 1783 to obtain
 1784

$$\begin{aligned}
 1785 \Delta \theta_t^{*\top} v_i &= \frac{\sqrt{\lambda_i^2 + 2L_t^i \cdot |\nabla f(\theta_t)^\top v_i|} - \lambda_i}{L_t^i} \\
 1786 &\leq \sqrt{2} \cdot \sqrt{\frac{|\nabla f(\theta_t)^\top v_i|}{L_t^i}} + \frac{\frac{\lambda_i^2}{2} \cdot \frac{1}{\sqrt{2L_t^i \cdot |\nabla f(\theta_t)^\top v_i|}} - \lambda_i}{L_t^i} \\
 1787 &\leq \sqrt{2} \cdot \sqrt{\frac{|\nabla f(\theta_t)^\top v_i|}{L_t^i}} + \frac{1}{2\sqrt{2}} \cdot \sqrt{\frac{|\nabla f(\theta_t)^\top v_i|}{L_t^i}} \\
 1788 &= \frac{5}{2\sqrt{2}} \cdot \sqrt{\frac{|\nabla f(\theta_t)^\top v_i|}{L_t^i}}
 \end{aligned}$$

1789 For i s.t. $\lambda_i > \sqrt{L_t^i \cdot |\nabla f(\theta_t)^\top v_i|}$ we use corollary H.1.1 with $x = \lambda_i^2$ to obtain
 1800

$$\begin{aligned}
 1801 \Delta \theta_t^{*\top} v_i &\leq \frac{|\nabla f(\theta_t)^\top v_i|}{\lambda_i} \leq \frac{|\nabla f(\theta_t)^\top v_i|}{\sqrt{L_t^i \cdot |\nabla f(\theta_t)^\top v_i|}} = \sqrt{\frac{|\nabla f(\theta_t)^\top v_i|}{L_t^i}} \\
 1802 &
 \end{aligned}$$

1803 For i s.t. $\lambda_i < 0$, we again use corollary H.1.1 with $x = \lambda_i^2$ to obtain
 1804

$$\begin{aligned}
 1805 \Delta \theta_t^{*\top} v_i &= \frac{2L_t^i \cdot |\nabla f(\theta_t)^\top v_i|}{L_t^i \left(\sqrt{\lambda_i^2 + 2L_t^i \cdot |\nabla f(\theta_t)^\top v_i|} - |\lambda_i| \right)} \leq \frac{2|\lambda_i|}{L_t^i} = \mathcal{O}(|\lambda_i|) \\
 1806 &
 \end{aligned}$$

□

1807 We are now ready to prove theorem 3.2.
 1808

1809 **Theorem.** *Worst case-optimal descent rate* Let f be a function with Lipschitz-continuous Hessian.
 1810 After t iterations, algorithm ELMO satisfies
 1811

$$1812 f(\theta_0) - f(\theta_t) = \mathcal{O}(\log t) \tag{17}$$

1813 *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
 1814 ARC optimization algorithm, of which algorithm ELMO is a special case (the case where ARC
 1815 perfectly estimates the Hessian Lipschitz parameter).
 1816

1817 Making use of lemma H.5 and noting that $m_t^i (\Delta \theta_t^{*\top} v_i) \leq 0$ by equation 13:
 1818

$$1819 |m_t^i (\Delta \theta_t^{*\top} v_i)| = |\nabla f(\theta_t)^\top v_i| \cdot |\Delta \theta_t^{*\top} v_i| + \frac{-\lambda_i}{2} \cdot (\Delta \theta_t^{*\top} v_i)^2 + \frac{L_t^i}{6} \cdot (\Delta \theta_t^{*\top} v_i)^3$$

1820 For i s.t. $\lambda_i \geq 0$:
 1821

$$\begin{aligned}
 1822 &\leq |\nabla f(\theta_t)^\top v_i| \cdot \mathcal{O}\left(\sqrt{\frac{|\nabla f(\theta_t)^\top v_i|}{L_t^i}}\right) + \mathcal{O}\left(\sqrt{\frac{|\nabla f(\theta_t)^\top v_i|}{L_t^i}}\right)^3 \\
 1823 &= \mathcal{O}\left(|\nabla f(\theta_t)^\top v_i|^{1.5}\right) = \mathcal{O}\left(\left(\frac{1}{t^{\frac{2}{3}}}\right)^{1.5}\right) = \mathcal{O}\left(\frac{1}{t}\right)
 \end{aligned}$$

1835

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

H.4 THEOREM 3.3

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 13, 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

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

1890 If $\lambda_i > 0$, then

$$1891$$

$$1892$$

$$1893$$

$$1894$$

$$1895 \frac{m_t^i(\Delta\theta_t^{*\top} v_i)}{M_t^i(\Delta\theta_t^{*\top} v_i)} = \frac{5 \sqrt{1+2 \frac{L_t^i |\nabla f(\theta_t)^\top v_i|}{\lambda_i^2}} - 1}{\frac{L_t^i |\nabla f(\theta_t)^\top v_i|}{\lambda_i^2}} - 8 \leq \lim_{x \rightarrow \infty} \frac{5 \frac{\sqrt{1+2x-1}}{x} - 8}{\frac{\sqrt{1+2x-1}}{x} - 4} = 2$$

$$1896$$

$$1897$$

$$1898$$

$$1899$$

$$1900$$

1903 due to the monotonic increasing nature of $\psi_5 : \mathbb{R}^+ \rightarrow \mathbb{R}$, $\psi_5(x) = \frac{5 \frac{\sqrt{1+2x-1}}{x} - 8}{\frac{\sqrt{1+2x-1}}{x} - 4}$.

1905 If $\lambda_i < 0$, then

$$1906$$

$$1907$$

$$1908$$

$$1909$$

$$1910 \frac{m_t^i(\Delta\theta_t^{*\top} v_i)}{M_t^i(\Delta\theta_t^{*\top} v_i)} = \frac{5 \sqrt{1+2 \frac{L_t^i |\nabla f(\theta_t)^\top v_i|}{|\lambda_i|^2}} + 1}{\frac{L_t^i |\nabla f(\theta_t)^\top v_i|}{|\lambda_i|^2}} + 8 \leq \lim_{x \rightarrow 0^+} \frac{5 \frac{\sqrt{1+2x+1}}{x} + 8}{\frac{\sqrt{1+2x+1}}{x} + 4} = 5$$

$$1911$$

$$1912$$

$$1913$$

$$1914$$

$$1915$$

$$1916$$

1918 due to the monotonic decreasing nature of $\psi_6 : \mathbb{R}^+ \rightarrow \mathbb{R}$, $\psi_6(x) = \frac{5 \frac{\sqrt{1+2x+1}}{x} + 8}{\frac{\sqrt{1+2x+1}}{x} + 4}$. □

1924 H.5 THEOREM 4.2: WORST-CASE DESCENT RATE FOR ARBITRARY OPTIMIZERS

1926 **Theorem.** *Relative Descent*

1927 •

$$1932$$

$$1933 \left| \frac{M_t^i(\Delta\theta_t^\top v_i) - M_t^i(\Delta\theta_t^{*\top} v_i)}{M_t^i(\Delta\theta_t^{*\top} v_i)} \right| = \Theta(|\Delta\Delta^i \theta_t'|^2)$$

$$1934$$

$$1935$$

1938 •

$$1941$$

$$1942 \left| \frac{m_t^i(\Delta\theta_t^\top v_i) - m_t^i(\Delta\theta_t^{*\top} v_i)}{m_t^i(\Delta\theta_t^{*\top} v_i)} \right| = \Theta(|\Delta\Delta^i \theta_t'|) \quad (18)$$

$$1943$$

1944 *Proof.* For the first part of the lemma,

$$\begin{aligned}
1945 & \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| \\
1946 & \\
1947 & \\
1948 & = \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} \\
1949 & \\
1950 & \\
1951 & \quad + \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} \\
1952 & \\
1953 & \\
1954 & \quad + \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} \\
1955 & \\
1956 & \\
1957 & \\
1958 & \\
1959 & = \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. \\
1960 & \\
1961 & \\
1962 & \quad + \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|} \\
1963 & \\
1964 & \quad \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) \\
1965 & \\
1966 & \\
1967 & \\
1968 & \\
1969 & = \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(\frac{\sqrt{\lambda_i^2 + 2L_t^i \cdot |\nabla f(\theta_t)^\top v_i|} - \lambda_i}{L_t^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. \\
1970 & \\
1971 & \\
1972 & \quad + \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|} \\
1973 & \\
1974 & \\
1975 & \quad \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) \\
1976 & \\
1977 & \\
1978 & \\
1979 & \\
1980 & \\
1981 & \\
1982 & \\
1983 & \\
1984 & \\
1985 & \\
1986 & = \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. \\
1987 & \\
1988 & \\
1989 & \quad + \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| \\
1990 & \\
1991 & \\
1992 & \quad \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) \\
1993 & \\
1994 & \\
1995 & \\
1996 & \\
1997 &
\end{aligned}$$

1998 If $\lambda_i = 0$:

1999

2000

2001

2002

2003

2004

2005

2006

2007

2008

2009

2010

2011

2012

2013

2014

2015

2016

2017

2018

2019

2020

2021

2022

2023

2024

2025

2026

2027

2028

2029

2030

2031

2032

2033

2034

2035

2036

2037

2038

2039

2040

2041

2042

2043

2044

2045

2046

2047

2048

2049

2050

2051

If $\lambda_i > 0$:

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

Proving that

2036

2037

2038

2039

2040

2041

2042

2043

2044

2045

2046

2047

2048

2049

2050

2051

$$\frac{1}{1 + 2\sqrt{1 + 2\frac{L_i^i \cdot |\nabla f(\theta_t)^\top v_i|}{\lambda_i^2}}} \in \left(0, \frac{1}{3} \right]$$

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

2042

2043

2044

2045

2046

2047

2048

2049

2050

2051

$$\psi_1 : \mathbb{R}^+ \rightarrow \mathbb{R}, \psi_1(x) = (1 + 2\sqrt{1 + 2x})^{-1}$$

is monotonic and satisfies

2046

2047

2048

2049

2050

2051

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

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

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

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

Proving that

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

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

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

is monotonic and satisfies

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

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

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}{-\frac{1}{6} \cdot \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} \cdot \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} \cdot \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}$$

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

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 4.2.

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

•

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

•

$$\left| \frac{m_t^i(\Delta \theta_t^\top v_i)}{m_t^i(\Delta \theta_t^{*\top} v_i)} \right| = \Theta\left(1 + |\Delta \Delta^i \theta_t|^p\right) \quad (19)$$

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

Proof. Proof is immediate from lemma H.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 $m_t^i(\Delta \theta_t^\top v_i)$. □

I 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). Although the code attached to this paper is capable of estimating these parameters, it does so too slowly to be practically useful in computing all of an algorithm's steps, under most settings. We suggest future work could improve this algorithm's computational complexity.

A limitation of our Newton's method performance predictor is the additional computational burden of computing the Lipschitz parameters. We provide code for doing so in the attached code on Github, but we recommend performing these computations sparingly, since the Lipschitz parameters are approximately locally stable anyway.

A second limitation of our work is its inability to provide any indication of the number of iterations left to achieve convergence. We see this as an acceptable limitation however, since in practice a model is only required to achieve a certain level of performance on the data decided ahead of time, without regard to how much further it could be optimized. As noted in the introduction, performance is measured by the loss function, so our descent rate bound satisfies this practical requirement.

A final limitation of our bound is its reliance on $\Delta \Delta^i \theta_t$ as a measure of algorithm optimality which is a function of $\Delta \theta_t^{*\top} v_i$, despite the fact that most optimizers do not compute that during training. This bound is therefore primarily of theoretical interest, as illustrated by its motivation of the very practical metric discussed in section 6