
Escaping Saddle Points with Adaptive Gradient Methods

Matthew Staib^{1,2} Sashank Reddi³ Satyen Kale³ Sanjiv Kumar³ Suvrit Sra¹

Abstract

Adaptive methods such as Adam and RMSProp are widely used in deep learning but are not well understood. In this paper, we seek a crisp, clean and precise characterization of their behavior in nonconvex settings. To this end, we first provide a novel view of adaptive methods as preconditioned SGD, where the preconditioner is estimated in an online manner. By studying the preconditioner on its own, we elucidate its purpose: it rescales the stochastic gradient noise to be isotropic near stationary points, which helps escape saddle points. Furthermore, we show that adaptive methods can efficiently estimate the aforementioned preconditioner. By gluing together these two components, we provide the first (to our knowledge) second-order convergence result for any adaptive method. The key insight from our analysis is that, compared to SGD, adaptive methods escape saddle points faster, and can converge faster overall to second-order stationary points.

1. Introduction

Stochastic first-order methods are the algorithms of choice for training deep networks, or more generally optimization problems of the form $\operatorname{argmin}_x \mathbb{E}_z[f(x, z)]$. While vanilla stochastic gradient descent (SGD) is still the most popular such algorithm, there has been much recent interest in adaptive methods that adaptively change learning rates for each parameter. This is a very old idea, e.g. (Jacobs, 1988); modern variants such as Adagrad (Duchi et al., 2011; McMahan & Streeter, 2010) Adam (Kingma & Ba, 2014) and RMSProp (Tieleman & Hinton, 2012) are widely used in deep learning due to their good empirical performance.

Adagrad uses the square root of the sum of the outer product of the past gradients to achieve adaptivity. At time step t ,

¹MIT EECS ²Based on work performed at Google Research New York ³Google Research New York. Correspondence to: Matthew Staib <mstaib@mit.edu>.

Adagrad updates the parameters in the following manner:

$$x_{t+1} = x_t - G_t^{-1/2} g_t,$$

where g_t is a noisy stochastic gradient at x_t and $G_t = \sum_{i=1}^t g_i g_i^T$. More often, a diagonal version of Adagrad is used due to practical considerations, which effectively yields a per parameter learning rate. In the convex setting, Adagrad achieves provably good performance, especially when the gradients are sparse. Although Adagrad works well in sparse convex settings, its performance appears to deteriorate in (dense) nonconvex settings. This performance degradation is often attributed to the rapid decay of the learning rate in Adagrad over time, which is a consequence of rapid increase in eigenvalues of the matrix G_t .

To tackle this issue, variants of Adagrad such as Adam and RMSProp have been proposed, which replace the sum of the outer products with an exponential moving average i.e., $G_t = (1 - \beta) \sum_{i=1}^t \beta^{t-i} g_i g_i^T$ for some constant $\beta \in (0, 1)$. This connection with Adagrad is often used to justify the design of Adam and RMSProp (e.g. (Goodfellow et al., 2016)). Although this connection is simple and appealing, it is clearly superficial. For instance, while learning rates in Adagrad decrease monotonically, it is not necessarily the case with Adam or RMSProp as shown recently in Reddi et al. (2018b), leading to their non-convergence in even simple convex settings. Thus, a principled understanding of these adaptive methods is largely missing.

In this paper, we introduce a much simpler way of thinking about adaptive methods such as Adam and RMSProp. Roughly, adaptive methods try to precondition SGD by some matrix A , e.g. when A is diagonal, A_{ii} corresponds to the effective stepsize for coordinate i . For some choices of A the algorithms do not have oracle access to A , but instead form an estimate $\hat{A} \approx A$. We separate out these two steps, by 1) giving convergence guarantees for an idealized setting where we have access to A , then 2) proving bounds on the quality of the estimate \hat{A} . Our approach makes it possible to effectively intuit about the algorithms, prove convergence guarantees (including second-order convergence), and give insights about how to choose algorithm parameters. It also leads to a number of surprising results, including an understanding of why the Reddi et al. (2018b) counterexample is hard for adaptive methods, why adaptive methods tend to escape saddle points faster than SGD (observed empirically

in (Reddi et al., 2018a)), insights into how to tune Adam’s parameters, and (to our knowledge) the first *second-order* convergence proof for any adaptive method.

Contributions: In addition to the aforementioned novel viewpoint, we also make the following key contributions:

- We develop a new approach for analyzing convergence of adaptive methods leveraging the preconditioner viewpoint and by way of disentangling estimation from the behavior of the *idealized* preconditioner.
- We provide *second-order convergence* results for adaptive methods, and as a byproduct, first-order convergence results. To the best of our knowledge, ours is the first work to show second order convergence for any adaptive method.
- We provide theoretical insights on how adaptive methods escape saddle points quickly. In particular, we show that the preconditioner used in adaptive methods leads to isotropic noise near stationary points, which helps escape saddle points faster.
- Our analysis also provides practical suggestions for tuning the exponential moving average parameter β .

1.1. Related work

There is an immense amount of work studying nonconvex optimization for machine learning, which is too much to discuss here in detail. Thus, we only briefly discuss two lines of work that are most relevant to our paper here. First, the recent work e.g. (Chen et al., 2018; Reddi et al., 2018b; Zou et al., 2018) to understand and give theoretical guarantees for adaptive methods such as Adam and RMSProp. Second, the technical developments in using first-order algorithms to achieve nonconvex second-order convergence (see Definition 2.1) e.g. (Ge et al., 2015; Allen-Zhu & Li, 2018; Jin et al., 2017; Lee et al., 2016).

Nonconvex convergence of adaptive methods. Many recent works have investigated convergence properties of adaptive methods. However, to our knowledge, all these results either require convexity or show only first-order convergence to stationary points. Reddi et al. (2018b) showed non-convergence of Adam and RMSProp in simple convex settings and provided a variant of Adam, called AMSGrad, with guaranteed convergence in the convex setting; Zhou et al. (2018) generalized this to a nonconvex first-order convergence result. Zaheer et al. (2018) showed first-order convergence of Adam when the batch size grows over time. Chen et al. (2018) bound the nonconvex convergence rate for a large family of Adam-like algorithms, but they essentially need to assume the effective stepsize is well-behaved

(as in AMSGrad). Agarwal et al. (2018) give a convex convergence result for a full-matrix version of RMSProp, which they extend to the nonconvex case via iteratively optimizing convex functions. Their algorithm uses a fixed sliding window instead of an exponential moving average. Makkamala & Hein (2017) prove improved convergence bounds for Adagrad in the online strongly convex case; they prove similar results for RMSProp, but only in a regime where it is essentially the same as Adagrad. Ward et al. (2018) give a nonconvex convergence result for a variant of Adagrad which employs an adaptively decreasing single learning rate (not per-parameter). Zou et al. (2018) give sufficient conditions for first-order convergence of Adam.

Nonconvex second order convergence of first order methods.

Starting with Ge et al. (2015) there has been a resurgence in interest in giving first-order algorithms that find *second* order stationary points of nonconvex objectives, where the gradient is small and the Hessian is nearly positive semidefinite. Most other results in this space operate in the deterministic setting where we have exact gradients, with carefully injected isotropic noise to escape saddle points. Levy (2016) show improved results for normalized gradient descent. Some algorithms rely on Hessian-vector products instead of pure gradient information e.g. (Agarwal et al., 2017; Carmon et al., 2018); it is possible to reduce Hessian-vector based algorithms to gradient algorithms (Xu et al., 2018; Allen-Zhu & Li, 2018). Jin et al. (2017) improve the dependence on dimension to polylogarithmic. Mokhtari et al. (2018) work towards adapting these techniques for constrained optimization. Most relevant to our work is that of Daneshmand et al. (2018), who prove convergence of SGD with better rates than Ge et al. (2015). Concurrent with our paper, Fang et al. (2019) give even better rates for SGD. Our work differs in that we provide second-order results for *preconditioned* SGD.

2. Notation and definitions

The objective function is f , and the gradient and Hessian of f are ∇f and $H = \nabla^2 f$, respectively. Denote by $x_t \in \mathbb{R}^d$ the iterate at time t , by g_t an unbiased stochastic gradient at x_t and by ∇_t the expected gradient at t . The matrix G_t refers to $\mathbb{E}[g_t g_t^T]$. Denote by $\lambda_{\max}(G)$ and $\lambda_{\min}(G)$ the largest and smallest eigenvalues of G , and $\kappa(G)$ is the condition number $\lambda_{\max}(G)/\lambda_{\min}(G)$ of G . For a vector v , its elementwise p -th power is written v^p . The objective $f(x)$ has global minimizer x^* , and we write $f^* = f(x^*)$. The Euclidean norm of a vector v is written as $\|v\|$, while for a matrix M , $\|M\|$ refers to the operator norm of M . The matrix I is the identity matrix, whose dimension should be clear from context.

Definition 2.1 (Second-order stationary point). A (τ_g, τ_h) -stationary point of f is a point x so that $\|\nabla f(x)\| \leq \tau_g$ and

Algorithm 1 Preconditioned SGD

Input: initial x_0 , time T , stepsize η , preconditioner $A(x)$
for $t = 0, \dots, T$ **do**
 $g_t \leftarrow$ stochastic gradient at x_t
 $A_t \leftarrow A(x_t)$ \triangleright e.g. $A_t = \mathbb{E}[g_t g_t^T]^{-1/2}$
 $x_{t+1} \leftarrow x_t - \eta A_t g_t$
end for

Algorithm 2 Full-matrix RMSProp

Input: initial x_0 , time T , stepsize η , small number $\varepsilon > 0$
 for stability
for $t = 0, \dots, T$ **do**
 $g_t \leftarrow$ stochastic gradient
 $\hat{G}_t = \beta \hat{G}_{t-1} + (1 - \beta) g_t g_t^T$
 $A_t = (\hat{G}_t + \varepsilon I)^{-1/2}$
 $x_{t+1} \leftarrow x_t - \eta A_t g_t$
end for

$\lambda_{\min}(\nabla^2 f(x)) \geq -\tau_h$, where $\tau_g, \tau_h > 0$.

As is standard (e.g. Nesterov & Polyak (2006)), we will discuss only $(\tau, \sqrt{\rho\tau})$ -stationary points, where ρ is the Lipschitz constant of the Hessian.

3. The RMSProp Preconditioner

Recall that methods like Adam and RMSProp replace the running sum $\sum_{i=1}^t g_i g_i^T$ used in Adagrad with an exponential moving average (EMA) of the form $(1 - \beta) \sum_{i=1}^t \beta^{t-i} g_i g_i^T$, e.g. full-matrix RMSProp is described formally in Algorithm 2. One key observation is that $\hat{G}_t = (1 - \beta) \sum_{i=1}^t \beta^{t-i} g_i g_i^T \approx \mathbb{E}[g_t g_t^T] =: G_t$ if β is chosen appropriately; in other words, at time t , the accumulated \hat{G}_t can be seen as an approximation of the true second moment matrix $G_t = \mathbb{E}[g_t g_t^T]$ at the current iterate. Thus, RMSProp can be viewed as preconditioned SGD (Algorithm 1) with the preconditioner being $A_t = G_t^{-1/2}$. In practice, it is too expensive to compute G_t exactly since it requires summing over all training samples. Practical adaptive methods (see Algorithm 2) estimate this preconditioner (or a diagonal approximation) on-the-fly via an EMA.

Before developing our formal results, we will build intuition about the behavior of adaptive methods by studying an idealized adaptive method (IAM) with perfect access to G_t . In the rest of this section, we make use of idealized RMSProp to answer some simple questions about adaptive methods that we feel have not yet been addressed satisfactorily.

3.1. What is the purpose of the preconditioner?

Why should preconditioning by $A = \mathbb{E}[g g^T]^{-1/2}$ help optimization? The original Adam paper (Kingma & Ba, 2014) argues that Adam is an approximation to natural gradient descent, since if the objective f is a log-likelihood, $\mathbb{E}[g g^T]$

approximates the Fisher information matrix, which captures curvature information in the space of distributions. There are multiple issues with comparing adaptive methods to natural gradient descent, which we discuss in Appendix A. Instead, Balles & Hennig (2018) argue that the primary function of adaptive methods is to equalize the stochastic gradient noise in each direction. But it is still *not* clear why or how equalized noise should help optimization.

Our IAM abstraction makes it easy to explain precisely how rescaling the gradient noise helps. Specifically, we manipulate the update rule for idealized RMSProp:

$$x_{t+1} \leftarrow x_t - \eta A_t g_t \quad (1)$$

$$= x_t - \eta A_t \nabla_t - \underbrace{\eta A_t (g_t - \nabla_t)}_{=: \xi_t} \quad (2)$$

The $A_t \nabla_t$ term is deterministic; only ξ_t is stochastic, with mean $\mathbb{E}[A_t (g_t - \nabla_t)] = A_t \mathbb{E}[g_t - \nabla_t] = 0$. Take $\varepsilon = 0$ and assume $G_t = \mathbb{E}[g_t g_t^T]$ is invertible, so that $\xi_t = G_t^{-1/2} (g_t - \nabla_t)$. Now we can be more precise about how RMSProp rescales gradient noise. Specifically, we compute the covariance of the noise ξ_t :

$$\text{Cov}(\xi_t) = I - G_t^{-1/2} \nabla_t \nabla_t^T G_t^{-1/2}. \quad (3)$$

The key insight is: near stationary points, ∇_t will be small, so that the noise covariance $\text{Cov}(\xi_t)$ is approximately the identity matrix I . In other words, at stationary points, the gradient noise is approximately isotropic. This observation hints at why adaptive methods are so successful for non-convex problems, where one of the main challenges is to escape saddle points (Reddi et al., 2018a). Essentially all first-order approaches for escaping saddlepoints rely on adding carefully tuned isotropic noise, so that regardless of what the escape direction is, there is enough noise in that direction to escape with high probability.

3.2. (Reddi et al., 2018b) counterexample resolution

Recently, Reddi et al. (2018b) provided a simple *convex* stochastic counterexample on which RMSProp and Adam do not converge. Their reasoning is that RMSProp and Adam too quickly forget about large gradients from the past, in favor of small (but poor) gradients at the present. In contrast, for RMSProp with the idealized preconditioner (Algorithm 1 with $A = \mathbb{E}[g g^T]^{-1/2}$), there is no issue, but the preconditioner A cannot be computed in practice. Rather, for this example, the exponential moving average estimation scheme fails to adequately estimate the preconditioner.

The counterexample is an optimization problem of the form

$$\min_{x \in [-1, 1]} F(x) = p f_1(x) + (1 - p) f_2(x), \quad (4)$$

where the stochastic gradient oracle returns ∇f_1 with probability p and ∇f_2 otherwise. Let $\zeta > 0$ be ‘‘small,’’ and $C >$

0 be “large.” Reddi et al. (2018b) set $p = (1 + \zeta)/(C + 1)$, $f_1(x) = Cx$, and $f_2(x) = -x$. Overall, then, $F(x) = \zeta x$ which is minimized at $x = -1$, however Reddi et al. (2018b) show that RMSProp has $\mathbb{E}[F(x_t)] \geq 0$ and so incurs suboptimality gap at least ζ . In contrast, the idealized preconditioner is a function of

$$\mathbb{E}[g^2] = p \left(\frac{\partial f_1}{\partial x} \right)^2 + (1 - p) \left(\frac{\partial f_2}{\partial x} \right)^2 = C(1 + \zeta) - \zeta$$

which is a constant independent of x . Hence the preconditioner is constant, and, up to the choice of stepsize, idealized RMSProp on this problem is the same as SGD, which of course will converge.

The difficulty for practical adaptive methods (which estimate $\mathbb{E}[g^2]$ via an EMA) is that as C grows, the variance of the estimate of $\mathbb{E}[g^2]$ grows too. Thus Reddi et al. (2018b) break Adam by making estimation of $\mathbb{E}[g^2]$ harder.

4. Main Results: Gluing Estimation and Optimization

The key enabling insight of this paper is to separately study the preconditioner and its estimation via EMA, then combine these to give proofs for practical adaptive methods. We will prove a formal guarantee that the EMA estimate \hat{G}_t is close to the true G_t . By combining our estimation results with the underlying behavior of the preconditioner, we will be able to give convergence proofs for practical adaptive methods that are constructed in a novel, modular way.

Separating these two components enables more general results: we actually analyze preconditioned SGD (Algorithm 1) with oracle access to an arbitrary preconditioner $A(x)$. Idealized RMSProp is but one particular instance. Our convergence results depend only on specific properties of the preconditioner $A(x)$, with which we can recover convergence results for many RMSProp variants simply by bounding the appropriate constants. For example, $A = (\mathbb{E}[gg^T]^{1/2} + \varepsilon I)^{-1}$ corresponds to full-matrix Adam with $\beta_1 = 0$ or RMSProp as commonly implemented. For cleaner presentation, we instead focus on the variant $A = (\mathbb{E}[gg^T] + \varepsilon I)^{-1/2}$, but our proof technique can handle either case or its diagonal approximation.

4.1. Estimating from Moving Sequences

The above discussion about IAM is helpful for intuition, and as a base algorithm for analyzing convergence. But it remains to understand how well the estimation procedure works, both for intuition’s sake and for later use in a convergence proof. In this section we introduce an abstraction we name “estimation from moving sequences.” This abstraction will allow us to guarantee high quality estimates of the preconditioner, or, for that matter, any similarly constructed

preconditioner. Our results will moreover make apparent how to choose the β parameter in the exponential moving average: β should increase with the stepsize η . Increasing β over time has been supported both empirically (Shazeer & Stern, 2018) as well as theoretically (Mukkamala & Hein, 2017; Zou et al., 2018; Reddi et al., 2018b), though to our knowledge, the precise pinning of β to the stepsize η is new.

Suppose there is a sequence of states $x_1, x_2, \dots, x_T \in \mathcal{X}$, e.g. the parameters of our model at each time step. We have access to the states x_t , but more importantly we know the states are not changing too fast: $\|x_t - x_{t-1}\|$ is bounded for all t . There is a Lipschitz function $G : \mathcal{X} \rightarrow \mathbb{R}^{d \times d}$, which in our case is the second moment matrix of the stochastic gradients, but could be more general. We would like to estimate $G(x)$ for each $x = x_t$, but we have only a noisy oracle $Y(x)$ for $G(x)$, which we assume is unbiased and has bounded variance. Our goal is, given noisy reads Y_1, \dots, Y_T of $G(x_1), \dots, G(x_T)$, to estimate $G(x_T)$ at the current point x_T as well as possible.

We consider estimators of the form $\sum_{t=1}^T w_t Y_t$. For example, setting $w_T = 1$ and all others to zero would yield an unbiased (but high variance) estimate of $G(x_T)$. We could assign more mass to older samples Y_t , but this will introduce bias into the estimate. By optimizing this bias-variance tradeoff, we can get a good estimator. In particular, taking w to be an exponential moving average (EMA) of $\{Y_t\}_{t=1}^T$ will prioritize more recent and relevant estimates, while placing enough weight on old estimates to reduce the variance. The tradeoff is controlled by the EMA parameter β ; e.g. if the sequence x_t moves slowly (the stepsize is small), we will want large β because older iterates are still very relevant.

In adaptive methods, the underlying function $G(x)$ we want to estimate is $\mathbb{E}[gg^T]$, and every stochastic gradient g gives us an unbiased estimate $Y = gg^T$. The diagonal approximation also fits this formulation: we can set $Y = \text{diag}(g^2)$ as an unbiased estimate of the matrix $\mathbb{E}[\text{diag}(g^2)]$. With these applications in mind, we formalize our results in terms of matrix estimation. By combining standard matrix concentration inequalities (e.g. Tropp (2015)) with bounds on how fast the sequence moves, we obtain the following result:

Theorem 4.1. *Assume $\|x_t - x_{t+1}\| \leq \eta M$. The function $G : \mathbb{R}^d \rightarrow \mathbb{R}^{d \times d}$ is matrix-valued and L -Lipschitz. For each t , Y_t is a random matrix with $\mathbb{E}[Y_t] = G(x_t)$, $\|G_t - \mathbb{E}[G_t]\| \leq R$, and $\|\mathbb{E}[(G_t - \mathbb{E}[G_t])^2]\| \leq \sigma_{\max}^2$. Set $w_t \propto \beta^{T-t}$ with $\sum_{t=1}^T w_t = 1$ and assume $T > 4/(1 - \beta)$. Then with probability $1 - \delta$, the estimation error $\Phi = \left\| \sum_{t=1}^T w_t Y_t - G(x_T) \right\|$ is bounded by*

$$\Phi \leq O(\sigma_{\max} \sqrt{1 - \beta} \sqrt{\log(d/\delta)} + ML\eta/(1 - \beta)).$$

This is optimized by $\beta = 1 - C\eta^{2/3}$, for which the bound is $O((\eta M \sigma_{\max}^2 (\log(d/\delta)) L)^{1/3})$ as long as $T > C'\eta^{-2/3}$.

The proof is given in Appendix G. As long as T is sufficiently large, we can get a high quality estimate of $G_t = \mathbb{E}[g_t g_t^T]$. For this, it suffices to start off the underlying optimization algorithm with $W = O(\eta^{-2/3})$ burn-in iterations where our estimate is updated but the algorithm is not started. This burn-in period will not affect asymptotic runtime as long as $W = O(\eta^{-2/3}) = O(T)$. In our non-convex convergence results we will require $T = O(\tau^{-4})$ and $\eta = O(\tau^2)$, so that $W = O(\tau^{-4/3})$ which is much smaller than T . In practice, one can get away with much shorter (or no) burn-in period.

If β is properly tuned, while running an adaptive method like RMSProp, we will get good estimates of $G = \mathbb{E}[gg^T]$ from samples gg^T . However, we actually require a good estimate of $A = \mathbb{E}[gg^T]^{-1/2}$ and variants. To treat estimation in a unified way, we introduce estimable matrix sequences:

Definition 4.1. A $(W, T, \eta, \Delta, \delta)$ -estimable matrix sequence is a sequence of matrices $\{A(x_t)\}_{t=1}^{W+T}$ generated from $\{x_t\}_t$ with $\|x_t - x_{t-1}\| \leq \eta$ so that with probability $1 - \delta$, after a burn-in of time W , we can achieve an estimate sequence $\{\hat{A}_t\}$ so that $\|\hat{A}_t - A_t\| \leq \Delta$ simultaneously for all times $t = W + 1, \dots, W + T$.

Applying Theorem 4.1 and union bounding over all time $t = W + 1, \dots, W + T$, we may state a concise result in terms of Definition 4.1:

Proposition 4.1. Suppose $G = \mathbb{E}[g_t g_t^T]$ is L -Lipschitz as a function of x . When applied to a generator sequence $\{x_t\}$ with $\|x_t - x_{t-1}\| \leq \eta M$ and samples $Y_t = g_t g_t^T$, the matrix sequence $G_t = \mathbb{E}[g_t g_t^T]$ is $(W, T, \eta M, \Delta, \delta)$ -estimable with $W = O(\eta^{-2/3})$, $T = \Omega(W)$, and $\Delta = O(\eta^{1/3} \sigma_{\max}^{2/3} (\log(2Td/\delta))^{1/3} M^{1/3} L^{1/3})$.

We are hence guaranteed a good estimate of G . What we actually want, though, is a good estimate of the preconditioner $A = (G + \varepsilon I)^{-1/2}$. In Appendix H we show how to bound the quality of an estimate of A . One simple result is:

Proposition 4.2. Suppose $G = \mathbb{E}[gg^T]$ is L -Lipschitz as a function of x . Further suppose a uniform bound $\lambda_{\min}(G)I \preceq G(x)$ for all x , with $\lambda_{\min}(G) > 0$. When applied to a generator sequence $\{x_t\}$ with $\|x_t - x_{t-1}\| \leq \eta M$ and samples $Y_t = g_t g_t^T$, the matrix sequence $A_t = (G_t + \varepsilon I)^{-1/2}$ is $(W, T, \eta M, \Delta, \delta)$ -estimable with $W = O(\eta^{-2/3})$, $T = \Omega(W)$, and $\Delta = O((\eta \sigma_{\max}^2 \log(2Td/\delta) ML)^{1/3} (\varepsilon + \lambda_{\min}(G))^{-3/2})$.

4.2. Convergence Results

We saw in the last two sections that it is simple to reason about adaptive methods via IAM, and that it is possible to compute a good estimate of the preconditioner. But we still need to glue the two together in order to get a convergence proof for practical adaptive methods.

In this section we will give non-convex convergence results, first for IAM and then for practical realizations thereof. We start with first-order convergence as a warm-up, and then move on to second-order convergence. In each case we give a bound for IAM, study it, and then give the corresponding bound for practical adaptive methods.

4.2.1. ASSUMPTIONS AND NOTATION

We want results for a wide variety of preconditioners A , e.g. $A = I$, the RMSProp preconditioner $A = (G + \varepsilon I)^{-1/2}$, and the diagonal version thereof, $A = (\text{diag}(G) + \varepsilon I)^{-1/2}$. To facilitate this and the future extension of our approach to other preconditioners, we give guarantees that hold for general preconditioners A . Our bounds depend on A via the following properties:

Definition 4.2. We say $A(x)$ is a $(\Lambda_1, \Lambda_2, \Gamma, \nu, \lambda_-)$ -preconditioner if, for all x , the following bounds hold. First, $\|A \nabla f\|^2 \leq \Lambda_1 \|A^{1/2} \nabla f\|^2$. Second, if $\tilde{f}(x)$ is the quadratic approximation of f at some point x_0 , we assume $\|A(\nabla f - \nabla \tilde{f})\| \leq \Lambda_2 \|\nabla f - \nabla \tilde{f}\|$. Third, $\Gamma \geq \mathbb{E}[\|Ag\|^2]$. Fourth, $\nu \leq \lambda_{\min}(A \mathbb{E}[gg^T] A^T)$. Finally, $\lambda_- \leq \lambda_{\min}(A)$.

Note that we could bound $\Lambda_1 = \Lambda_2 = \lambda_{\max}(A)$. but in practice Λ_1 and Λ_2 may be smaller, since they depend on the behavior of A only in specific directions. In particular, if the preconditioner A is well-aligned with the Hessian, as may be the case if the natural gradient approximation is valid, then Λ_1 would be very small. If f is exactly quadratic, Λ_2 can be taken as a constant. The constant Γ controls the magnitude of (rescaled) gradient noise, which affects stability at a local minimum. Finally, ν gives a lower bound on the amount of gradient noise in any direction; when ν is larger it is easier to escape saddle points¹. For shorthand, a $(\cdot, \cdot, \Gamma, \cdot, \lambda_-)$ -preconditioner needs to satisfy only the corresponding inequalities.

In Appendix C we provide bounds on these constants for variants of the second moment preconditioner. We highlight the two most relevant cases, for SGD and RMSProp:

Proposition 4.3. The preconditioner $A = I$ is a $(\Lambda_1, \Lambda_2, \Gamma, \nu, \lambda_-)$ -preconditioner, with $\Lambda_1 = \Lambda_2 = 1$, $\Gamma \leq \mathbb{E}[\|g\|^2] \leq d \cdot \text{tr}(G)$, $\nu \leq \lambda_{\min}(G)$, and $\lambda_- = 1$.

Proposition 4.4. The preconditioner $A = (G + \varepsilon I)^{-1/2}$ is a $(\Lambda_1, \Lambda_2, \Gamma, \nu, \lambda_-)$ -preconditioner, with

$$\Lambda_1 = \Lambda_2 = \frac{1}{(\lambda_{\min}(G) + \varepsilon)^{1/2}}, \quad \Gamma = \frac{d \lambda_{\max}(G)}{\varepsilon + \lambda_{\max}(G)},$$

$$\nu = \frac{\lambda_{\min}(G)}{\lambda_{\min}(G) + \varepsilon}, \quad \text{and} \quad \lambda_- = (\lambda_{\max}(G) + \varepsilon)^{-1/2}.$$

¹In cases where $G = \mathbb{E}[gg^T]$ is rank deficient, e.g. in high-dimensional finite sum problems, lower bounds on $\lambda_{\min}(G)$ should be understood as lower bounds on $\mathbb{E}[(v^T g)^2]$ for escape directions v from saddle points, analogous to the ‘‘CNC condition’’ from (Daneshmand et al., 2018).

4.2.2. FIRST-ORDER CONVERGENCE

Proofs are given in Appendix F. For all first-order results, we assume that A is a $(\cdot, \cdot, \Gamma, \cdot, \lambda_-)$ -preconditioner. The proof technique is essentially standard, with minor changes in order to accommodate general preconditioners. First, suppose we have exact oracle access to the preconditioner:

Theorem 4.2. *Run preconditioned SGD with preconditioner A and stepsize $\eta = \tau^2 \lambda_- / (L\Gamma)$. For small enough τ , after $T = 2(f(x_0) - f^*)L\Gamma / (\tau^4 \lambda_-^2)$ iterations,*

$$\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E} [\|\nabla f(x_t)\|^2] \leq \tau^2. \quad (5)$$

Now we consider an alternate version where instead of the preconditioner A_t , we precondition by an noisy version \hat{A}_t that is close to A_t , i.e. $\|\hat{A}_t - A_t\| \leq \Delta$.

Theorem 4.3. *Suppose we have access to an inexact preconditioner \hat{A} , which satisfies $\|\hat{A} - A\| \leq \Delta$ for $\Delta < \lambda_- / 2$. Run preconditioned SGD with preconditioner \hat{A} and stepsize $\eta = \tau^2 \lambda_- / (4\sqrt{2}L\Gamma)$. For small enough τ , after $T = 32(f(x_0) - f^*)L\Gamma / (\tau^4 \lambda_-^2)$ iterations, we will have*

$$\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E} [\|\nabla f(x_t)\|^2] \leq \tau^2. \quad (6)$$

The results are the same up to constants. In other words, as long as we can achieve less than $\lambda_- / 2$ error, we will converge at essentially the same rate as if we had the exact preconditioner. In light of this, for the second-order convergence results, we treat only the noisy version.

Theorem 4.3 gives a convergence bound assuming a good estimate of the preconditioner, and our estimation results guarantee a good estimate. By gluing together Theorem 4.3 with our estimation results for the RMSProp preconditioner, i.e. Proposition 4.2, we can give a convergence result for bona fide RMSProp:

Corollary 4.1. *Consider RMSProp with burn-in, as in Algorithm 3, where we estimate $A = (G + \varepsilon I)^{-1/2}$. Retain the same choice of $\eta = O(\tau^2)$ and $T = O(\tau^{-4})$ as in Theorem 4.3. For small enough τ , such a choice of η will yield $\Delta < \lambda_- / 2$. Choose all other parameters e.g. β in accordance with Proposition 4.2. In particular, choose $W = \Theta(\eta^{-2/3}) = \Theta(\tau^{-4/3}) = O(T)$ for the burn-in parameter. Then with probability $1 - \delta$, in overall time $O(W + T) = O(\tau^{-4})$, we achieve*

$$\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E} [\|\nabla f(x_t)\|^2] \leq \tau^2. \quad (7)$$

Algorithm 3 RMSProp with burn-in

Input: initial x_0 , time T , stepsize η , burn-in length W
 $\hat{G}_0 \leftarrow \text{BURNIN}(W, \beta)$ ▷ Appendix B
for $t = 0, \dots, T$ **do**
 $g_t \leftarrow$ stochastic gradient
 $\hat{G}_t \leftarrow \beta \hat{G}_{t-1} + (1 - \beta) g_t g_t^T$
 $\hat{A}_t \leftarrow \hat{G}_t^{-1/2}$
 $x_{t+1} \leftarrow x_t - \eta \hat{A}_t g_t$
end for

4.2.3. SECOND-ORDER CONVERGENCE

Now we leverage the power of our high level approach to prove nonconvex second-order convergence for adaptive methods. Like the first-order results, we start by proving convergence bounds for a generic, possibly inexact preconditioner A . Our proof is based on that of Daneshmand et al. (2018) for SGD, and therefore we achieve the same $O(\tau^{-5})$ rate. It may be possible to improve our result using the technique of Fang et al. (2019), which is concurrent work to ours. However, our focus is on the preconditioner, and our study of it is wholly new. Accordingly, we study the convergence of Algorithm 4, which is the same as Algorithm 1 (generic preconditioned SGD) except that once in a while we take a large stepsize so we may escape saddlepoints. The proof is given completely in Appendix E. At a high level, we show the algorithm makes progress when the gradient is large and when we are at a saddle point, and does not escape from local minima. Our analysis uses all the constants specified in Definition 4.2, e.g. the speed of escape from saddle points depends on ν , the lower bound on stochastic gradient noise.

Then, as before, we simply fuse our convergence guarantees with our estimation guarantees. The end result is, to our knowledge, the first nonconvex second-order convergence result for any adaptive method.

Definitions for second-order results. Assume further that the Hessian H is ρ -Lipschitz and the preconditioner $A(x)$ is α -Lipschitz. The dependence on these constants is made precise in the proof, in Appendix E. The usual stepsize is η , while r is the occasional large stepsize that happens every t_{thresh} iterations. We tolerate a small probability of failure δ . For all results, we assume A is a $(\Lambda_1, \Lambda_2, \Gamma, \nu, \lambda_-)$ -preconditioner. For simplicity, we assume the noisy estimate \hat{A} also satisfies the Λ_1 inequality. We will also assume a uniform bound on $\|Ag\| \leq M = O(\sqrt{\Gamma})$.

The proofs rely on a few other quantities that we optimally determine as a function of the problem parameters: f_{thresh} is a threshold on the function value progress, and $g_{\text{thresh}} = f_{\text{thresh}} / t_{\text{thresh}}$ is the time-amortized average of f_{thresh} . We specify the precise values of all quantities in the proof.

Theorem 4.4. *Consider Algorithm 4 with inexact preconditioner \hat{A}_t and exact preconditioner A_t satisfying the*

Algorithm 4 Preconditioned SGD with increasing stepsize

Input: initial x_0 , time T , stepsizes η, r , threshold t_{thresh} , matrix error Δ

for $t = 0, \dots, T$ **do**

$A_t \leftarrow A(x_t)$ \triangleright preconditioner at x_t

$\hat{A}_t \leftarrow$ any matrix with $\|\hat{A}_t - A_t\| \leq \Delta$

$g_t \leftarrow$ stochastic gradient at x_t

if $t \bmod t_{\text{thresh}} = 0$ **then**

$x_{t+1} \leftarrow x_t - r \hat{A}_t g_t$

else

$x_{t+1} \leftarrow x_t - \eta \hat{A}_t g_t$

end if

end for

preceding requirements. Suppose that for all t , we have $\|\hat{A}_t - A_t\| = O(\tau^{1/2})$. Then for small τ , with probability $1 - \delta$, we reach an $(\tau, \sqrt{\rho\tau})$ -stationary point in time

$$T = \tilde{O} \left(\frac{\Lambda_1^4 \Lambda_2^4 \Gamma^4}{\lambda_-^{10} \nu^4} \cdot \frac{L^3}{\rho \delta^3} \cdot \tau^{-5} \right). \quad (8)$$

The big- O suppresses other constants given in the proof.

To prove a result for bona fide RMSProp, we need to combine Theorem 4.4 with an algorithm that maintains a good estimate of $G = \mathbb{E}[gg^T]$ (and consequently $A = (G + \varepsilon I)^{-1/2}$). This is more delicate than the first-order case because now the stepsize varies. Whenever we take a large stepsize, the estimation algorithm will need to hallucinate S number of smaller steps in order to keep the estimate accurate. Our overall scheme is formalized in Appendix B, for which the following convergence result holds:

Corollary 4.2. Consider the RMSProp version of Algorithm 4 that is described in Appendix B. Retain the same choice of $\eta = O(\tau^{5/2})$, $r = O(\tau)$, and $T = O(\tau^{-5})$ as in Theorem 4.4. For small enough τ , such a choice of η will yield $\Delta < \lambda_-/2$. Choose $W = \Theta(\eta^{-2/3}) = \Theta(\tau^{-5/3}) = O(T)$ for the burn-in parameter. Choose $S = O(\tau^{-3/2})$, so that as far as the estimation scheme is concerned, the stepsize is bounded by $\max\{\eta, r/S\} = O(\tau^{5/2}) = O(\eta)$. Then as before, with probability $1 - \delta$, we can reach an $(\tau, \sqrt{\rho\tau})$ -stationary point in total time

$$W + T = \tilde{O} \left(\frac{\Lambda_1^4 \Lambda_2^4 \Gamma^4}{\lambda_-^{10} \nu^4} \cdot \frac{L^3}{\rho \delta^3} \cdot \tau^{-5} \right), \quad (9)$$

where $\Lambda_1, \Lambda_2, \Gamma, \nu, \lambda_-$ are the constants describing $A = (G + \varepsilon I)^{-1/2}$.

Again, as in the first order results, one could substitute in any other estimable preconditioner. In particular, in Appendix D we discuss the more common diagonal version of RMSProp.

5. Discussion

Separating the estimation step from the preconditioning enables evaluation of different choices for the preconditioner.

5.1. How to set the regularization parameter ε

In the adaptive methods literature, it is still a mystery how to properly set the regularization parameter ε that ensures invertibility of $G + \varepsilon I$. When the optimality tolerance τ is small enough, estimating the preconditioner is not the bottleneck. Thus, focusing only on the idealized case, one could just choose ε to minimize the bound. Our first-order results depend on ε only through the following term:

$$\frac{\Gamma}{\lambda_{\min}(A)} \leq \frac{d\lambda_{\min}(G)}{\varepsilon + \lambda_{\min}(G)} \cdot (\lambda_{\max}(G) + \varepsilon), \quad (10)$$

where we have used the preconditioner bounds from Proposition 4.4. This is minimized by taking $\varepsilon \rightarrow \infty$, which suggests using identity preconditioner, or SGD. In contrast, for second-order convergence, the bound is

$$\frac{\Lambda_1^4 \Lambda_2^4 \Gamma^4}{\lambda_-^{10} \nu^4} \leq d^4 \kappa(G)^4 (\lambda_{\max}(G) + \varepsilon), \quad (11)$$

which is instead minimized with $\varepsilon = 0$. So for the best second-order convergence rate, it is desirable to set ε as small as possible. Note that since our bounds hold only for small enough convergence tolerance τ , it is possible that the optimal ε should depend in some way on τ .

5.2. Comparison to SGD

Another important question we make progress towards is: when are adaptive methods better than SGD? Our second-order result depends on the preconditioner only through $\Lambda_1^4 \Lambda_2^4 \Gamma^4 / (\lambda_-^{10} \nu^4)$. Plugging in Proposition 4.3 for SGD, we may bound

$$\frac{\Lambda_1^4 \Lambda_2^4 \Gamma^4}{\lambda_-^{10} \nu^4} \leq \frac{\mathbb{E}[\|g\|^2]^4}{\lambda_{\min}(G)^4} \leq d^4 \kappa(G)^4, \quad (12)$$

while for full-matrix RMSProp, we have

$$\frac{\Lambda_1^4 \Lambda_2^4 \Gamma^4}{\lambda_-^{10} \nu^4} \leq d^4 \kappa(G)^4 (\lambda_{\max}(G) + \varepsilon). \quad (13)$$

Setting $\varepsilon = 0$ for simplicity, we conclude that full-matrix RMSProp converges faster if $\lambda_{\max}(G) \leq 1$.

Now suppose that for a given optimization problem, the preconditioner A is well-aligned with the Hessian so that $\Lambda_1 = O(1)$ (e.g. if the natural gradient approximation holds) and that near saddle points the objective is essentially quadratic so that $\Lambda_2 = O(1)$. In this regime, the preconditioner dependence of idealized full matrix RMSProp is $d^4 \lambda_{\max}(G)^5$, which yields a better result than SGD when $\lambda_{\max}(G) \leq \lambda_{\min}(G)^{-4}$. This will happen whenever $\lambda_{\min}(G)$ is relatively small. Thus, when there is not much noise in the escape direction, and the Hessian and $G^{-1/2}$ are not poorly aligned, RMSProp will converge faster overall. Similar phenomenon can be shown for the diagonal case when the approximation is good, per the results in Appendix C and D.

5.3. Alternative preconditioners

Our analysis inspires the design of other preconditioners: e.g., if at each iteration we sample two independent stochastic gradients g_1 and g_2 , we have unbiased sample access to $(g_1 - g_2)(g_1 - g_2)^T$, which in expectation yields the covariance $\Sigma = \text{Cov}(g)$ instead of the second moment matrix of g . It immediately follows that we can prove second-order convergence results for an algorithm that constructs an exponential moving average estimate of Σ and preconditions by $\Sigma^{-1/2}$, as advocated by [Ida et al. \(2017\)](#).

5.4. Tuning the EMA parameter β

Another mystery of adaptive methods is how to set the exponential moving average (EMA) parameter β . In practice β is typically set to a constant, e.g. 0.99, while other parameters such as the stepsize η are tuned more carefully and may vary over time. While our estimation guarantee [Theorem 4.1](#), suggests setting $\beta = 1 - O(\eta^{2/3})$, the specific formula depends on constants that may be unknown, e.g. Lipschitz constants and gradient norms. Instead, one could set $\beta = 1 - C\eta^{2/3}$, and search for a good choice of the hyperparameter C . For example, the common initial choice of $\eta = 0.001$ and $\beta = 0.99$ corresponds to $C = 1$.

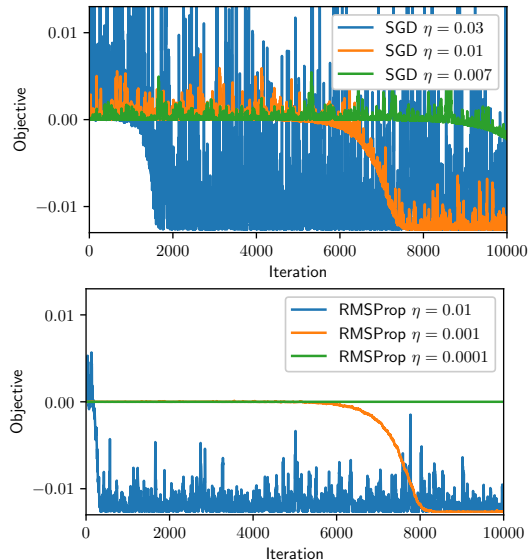
6. Experiments

We experimentally test our claims about adaptive methods escaping saddle points, and our suggestion for setting β .

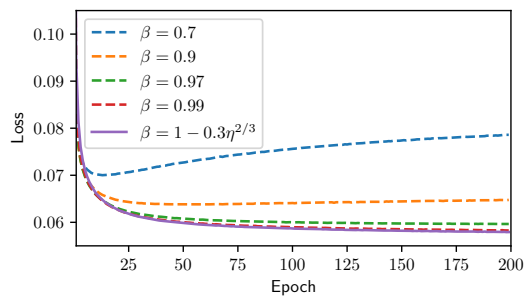
Escaping saddle points. First, we test our claim that when the gradient noise is ill-conditioned, adaptive methods escape saddle points faster than SGD, and often converge faster to (approximate) local minima. We construct a two dimensional² non-convex problem $f(x) = \frac{1}{n} \sum_{i=1}^n f_i(x)$ where $f_i(x) = \frac{1}{2}x^T Hx + b_i^T x + \|x\|_{10}^4$. Here, $H = \text{diag}([1, -0.1])$, so f has a saddle point at the origin with objective value zero. The vectors b_i are chosen so that sampling b uniformly from $\{b_i\}_{i=1}^n$ yields $\mathbb{E}[b] = 0$ and $\text{Cov}(b) = \text{diag}([1, 0.01])$. Hence at the origin there is an escape direction but little gradient noise in that direction.

We initialize SGD and (diagonal) RMSProp (with $\beta = 1 - \eta^{2/3}$) at the saddle point and test several stepsizes η for each. Results for the first 10^4 iterations are shown in [Figure 1](#). In order to escape the saddle point as fast as RMSProp, SGD requires a substantially larger stepsize, e.g. SGD needs $\eta = 0.01$ to escape as fast as RMSProp does with $\eta = 0.001$. But with such a large stepsize, SGD cannot converge to a small neighborhood of the local minimum, and instead bounces around due to gradient noise. Since RMSProp can escape with a small stepsize, it can converge

²The same phenomenon still holds in higher dimensions but the presentation is simpler with $d = 2$.



[Figure 1](#). SGD (top) vs RMSProp (bottom) performance escaping a saddle point with poorly conditioned gradient noise. Compared to RMSProp, SGD requires a larger stepsize to escape as quickly, which negatively impacts convergence to the local minimum.



[Figure 2](#). Performance on MNIST logistic regression of RMSProp with different choices of β and decreasing stepsize.

to a much smaller neighborhood of the local minimum. Overall, for any fixed final convergence criterion, RMSProp escapes faster and converges faster overall.

Setting the EMA parameter β . Next, we test our recommendations regarding setting the EMA parameter β . We consider logistic regression on MNIST. We use (diagonal) RMSProp with batch size 100, decreasing stepsize $\eta_t = 0.001/\sqrt{t}$ and $\varepsilon = 10^{-8}$, and compare different schedules for β . Specifically we test $\beta \in \{0.7, 0.9, 0.97, 0.99\}$ (so that $1 - \beta$ is spaced roughly logarithmically) as well as our recommendation of $\beta_t = 1 - C\eta_t^{2/3}$ for $C \in \{0.1, 0.3, 1\}$. As shown in [Figure 2](#), all options for β have similar performance initially, but as η_t decreases, large β yields substantially better performance. In particular, our decreasing β schedule achieved the best performance, and moreover was insensitive to how C was set.

Acknowledgements

SS acknowledges support from the DARPA Lagrange grant, an Amazon Research Award, and the NSF-CAREER Award (ID 1846088). We thank Nicolas Le Roux for helpful conversations.

References

- Agarwal, N., Allen-Zhu, Z., Bullins, B., Hazan, E., and Ma, T. Finding approximate local minima faster than gradient descent. In *Proceedings of the 49th Annual ACM SIGACT Symposium on Theory of Computing, STOC 2017*, pp. 1195–1199, New York, NY, USA, 2017. ACM. ISBN 978-1-4503-4528-6. doi: 10.1145/3055399.3055464.
- Agarwal, N., Bullins, B., Chen, X., Hazan, E., Singh, K., Zhang, C., and Zhang, Y. The case for full-matrix adaptive regularization. *arXiv preprint arXiv:1806.02958*, 2018.
- Allen-Zhu, Z. and Li, Y. Neon2: Finding local minima via first-order oracles. In Bengio, S., Wallach, H., Larochelle, H., Grauman, K., Cesa-Bianchi, N., and Garnett, R. (eds.), *Advances in Neural Information Processing Systems 31*, pp. 3720–3730. Curran Associates, Inc., 2018.
- Balles, L. and Hennig, P. Dissecting Adam: The sign, magnitude and variance of stochastic gradients. In Dy, J. and Krause, A. (eds.), *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pp. 404–413, Stockholmsmssan, Stockholm Sweden, 10–15 Jul 2018. PMLR.
- Carmon, Y., Duchi, J., Hinder, O., and Sidford, A. Accelerated methods for nonconvex optimization. *SIAM Journal on Optimization*, 28(2):1751–1772, 2018. doi: 10.1137/17M1114296.
- Chen, X., Liu, S., Sun, R., and Hong, M. On the convergence of a class of adam-type algorithms for non-convex optimization. *arXiv preprint arXiv:1808.02941*, 2018.
- Daneshmand, H., Kohler, J., Lucchi, A., and Hofmann, T. Escaping saddles with stochastic gradients. In Dy, J. and Krause, A. (eds.), *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pp. 1155–1164, Stockholmsmssan, Stockholm Sweden, 10–15 Jul 2018. PMLR.
- Duchi, J., Hazan, E., and Singer, Y. Adaptive subgradient methods for online learning and stochastic optimization. *J. Mach. Learn. Res.*, 12:2121–2159, July 2011. ISSN 1532-4435.
- Fang, C., Lin, Z., and Zhang, T. Sharp analysis for non-convex sgd escaping from saddle points. *arXiv preprint arXiv:1902.00247*, 2019.
- Ge, R., Huang, F., Jin, C., and Yuan, Y. Escaping from saddle points — online stochastic gradient for tensor decomposition. In Grnwald, P., Hazan, E., and Kale, S. (eds.), *Proceedings of The 28th Conference on Learning Theory*, volume 40 of *Proceedings of Machine Learning Research*, pp. 797–842, Paris, France, 03–06 Jul 2015. PMLR.
- Goodfellow, I., Bengio, Y., and Courville, A. *Deep Learning*. MIT Press, 2016. <http://www.deeplearningbook.org>.
- Ida, Y., Fujiwara, Y., and Iwamura, S. Adaptive learning rate via covariance matrix based preconditioning for deep neural networks. In *Proceedings of the Twenty-Sixth International Joint Conference on Artificial Intelligence, IJCAI-17*, pp. 1923–1929, 2017. doi: 10.24963/ijcai.2017/267.
- Jacobs, R. A. Increased rates of convergence through learning rate adaptation. *Neural networks*, 1(4):295–307, 1988.
- Jin, C., Ge, R., Netrapalli, P., Kakade, S. M., and Jordan, M. I. How to escape saddle points efficiently. In Precup, D. and Teh, Y. W. (eds.), *Proceedings of the 34th International Conference on Machine Learning*, volume 70 of *Proceedings of Machine Learning Research*, pp. 1724–1732, International Convention Centre, Sydney, Australia, 06–11 Aug 2017. PMLR.
- Kingma, D. P. and Ba, J. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- Lee, J. D., Simchowitz, M., Jordan, M. I., and Recht, B. Gradient descent only converges to minimizers. In Feldman, V., Rakhlin, A., and Shamir, O. (eds.), *29th Annual Conference on Learning Theory*, volume 49 of *Proceedings of Machine Learning Research*, pp. 1246–1257, Columbia University, New York, New York, USA, 23–26 Jun 2016. PMLR.
- Levy, K. Y. The power of normalization: Faster evasion of saddle points. *arXiv preprint arXiv:1611.04831*, 2016.
- McMahan, H. B. and Streeter, M. J. Adaptive bound optimization for online convex optimization. In *COLT 2010 - The 23rd Conference on Learning Theory, Haifa, Israel, June 27-29, 2010*, pp. 244–256, 2010.
- Mokhtari, A., Ozdaglar, A., and Jadbabaie, A. Escaping saddle points in constrained optimization. In Bengio, S., Wallach, H., Larochelle, H., Grauman, K., Cesa-Bianchi,

- N., and Garnett, R. (eds.), *Advances in Neural Information Processing Systems 31*, pp. 3633–3643. Curran Associates, Inc., 2018.
- Mukkamala, M. C. and Hein, M. Variants of RMSProp and Adagrad with logarithmic regret bounds. In Precup, D. and Teh, Y. W. (eds.), *Proceedings of the 34th International Conference on Machine Learning*, volume 70 of *Proceedings of Machine Learning Research*, pp. 2545–2553, International Convention Centre, Sydney, Australia, 06–11 Aug 2017. PMLR.
- Nesterov. *Introductory lectures on convex optimization : a basic course*. Kluwer Academic Publishers, Boston, 2004. ISBN 1402075537.
- Nesterov, Y. and Polyak, B. T. Cubic regularization of Newton method and its global performance. *Mathematical Programming*, 108(1):177–205, 2006.
- Reddi, S., Zaheer, M., Sra, S., Póczos, B., Bach, F., Salakhutdinov, R., and Smola, A. A generic approach for escaping saddle points. In Storkey, A. and Perez-Cruz, F. (eds.), *Proceedings of the Twenty-First International Conference on Artificial Intelligence and Statistics*, volume 84 of *Proceedings of Machine Learning Research*, pp. 1233–1242, Playa Blanca, Lanzarote, Canary Islands, 09–11 Apr 2018a. PMLR.
- Reddi, S. J., Kale, S., and Kumar, S. On the convergence of Adam and beyond. In *International Conference on Learning Representations*, 2018b.
- Ruder, S. An overview of gradient descent optimization algorithms. *arXiv preprint arXiv:1609.04747*, 2016.
- Shazeer, N. and Stern, M. Adafactor: Adaptive learning rates with sublinear memory cost. In Dy, J. and Krause, A. (eds.), *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pp. 4596–4604, Stockholmsmssan, Stockholm Sweden, 10–15 Jul 2018. PMLR.
- Tieleman, T. and Hinton, G. Lecture 6.5-RMSProp: Divide the gradient by a running average of its recent magnitude. *COURSERA: Neural networks for machine learning*, 4(2):26–31, 2012.
- Tropp, J. A. An introduction to matrix concentration inequalities. *Foundations and Trends in Machine Learning*, 8(1-2):1–230, 2015. ISSN 1935-8237. doi: 10.1561/22000000048.
- Ward, R., Wu, X., and Bottou, L. Adagrad stepsizes: Sharp convergence over nonconvex landscapes, from any initialization. *arXiv preprint arXiv:1806.01811*, 2018.
- Xu, Y., Rong, J., and Yang, T. First-order stochastic algorithms for escaping from saddle points in almost linear time. In Bengio, S., Wallach, H., Larochelle, H., Grauman, K., Cesa-Bianchi, N., and Garnett, R. (eds.), *Advances in Neural Information Processing Systems 31*, pp. 5531–5541. Curran Associates, Inc., 2018.
- Zaheer, M., Reddi, S., Sachan, D., Kale, S., and Kumar, S. Adaptive methods for nonconvex optimization. In *NIPS*. 2018.
- Zhou, D., Tang, Y., Yang, Z., Cao, Y., and Gu, Q. On the convergence of adaptive gradient methods for nonconvex optimization. *arXiv preprint arXiv:1808.05671*, 2018.
- Zou, F., Shen, L., Jie, Z., Zhang, W., and Liu, W. A sufficient condition for convergences of Adam and RMSProp. *arXiv preprint arXiv:1811.09358*, 2018.