On the Convergence of Adaptive Gradient Methods for Nonconvex Optimization

Dongruo Zhou*† Yiqi Tang*‡ Ziyan Yang*§ Yuan Cao¶ Quanquan Gu $^{\parallel}$

Abstract

Adaptive gradient methods are workhorses in deep learning. However, the convergence guarantees of adaptive gradient methods for nonconvex optimization have not been sufficiently studied. In this paper, we provide a sharp analysis of a recently proposed adaptive gradient method namely partially adaptive momentum estimation method (Padam) (Chen and Gu, 2018), which admits many existing adaptive gradient methods such as RMSProp and AMSGrad as special cases. Our analysis shows that, for smooth nonconvex functions, Padam converges to a first-order stationary point at the rate of $O((\sum_{i=1}^{d} \|\mathbf{g}_{1:T,i}\|_2)^{1/2}/T^{3/4} + d/T)$, where T is the number of iterations, d is the dimension, $\mathbf{g}_1, \dots, \mathbf{g}_T$ are the stochastic gradients, and $\mathbf{g}_{1:T,i} = [g_{1,i}, g_{2,i}, \dots, g_{T,i}]^{\top}$. Our theoretical result also suggests that in order to achieve faster convergence rate, it is necessary to use Padam instead of AMSGrad. This is well-aligned with the empirical results of deep learning reported in Chen and Gu (2018).

1 Introduction

Stochastic gradient descent (SGD) (Robbins and Monro, 1951) and its variants have been widely used in training deep neural networks. Among those variants, adaptive gradient methods (AdaGrad) (Duchi et al., 2011; McMahan and Streeter, 2010), which scale each coordinate of the gradient by a function of past gradients, can achieve better performance than vanilla SGD in practice when the gradients are sparse. An intuitive explanation for the success of AdaGrad is that it automatically adjusts the learning rate for each feature based on the partial gradient, which accelerates the convergence. However, AdaGrad was later found to demonstrate degraded performance especially in cases where the loss function is nonconvex or the gradient is dense, due to rapid decay of learning rate. This problem is especially exacerbated in deep learning due to the huge number of optimization

^{*}Equal Contribution

[†]Department of Computer Science, University of California, Los Angeles, CA 90095, USA; e-mail: zhoudongruo@gmail.com

[‡]Department of Computer Science, University of Virginia, Charlottesville, VA 22904, USA; e-mail: yt6ze@virginia.edu

[§]Department of Computer Science, University of Virginia, Charlottesville, VA 22904, USA; e-mail: zy3cx@virginia.edu

[¶]The Program in Applied and Computational Mathematics, Princeton University, Princeton, NJ 08544, USA; e-mail: yuanc@princeton.edu

Department of Computer Science, University of California, Los Angeles, CA 90095, USA; e-mail: qgu@cs.ucla.edu

variables. To overcome this issue, RMSProp (Tieleman and Hinton, 2012) was proposed to use exponential moving average rather than the arithmetic average to scale the gradient, which mitigates the rapid decay of the learning rate. Kingma and Ba (2014) proposed an adaptive momentum estimation method (Adam), which incorporates the idea of momentum (Polyak, 1964; Sutskever et al., 2013) into RMSProp. Other related algorithms include AdaDelta (Zeiler, 2012) and Nadam (Dozat, 2016), which combine the idea of exponential moving average of the historical gradients, Polyak's heavy ball (Polyak, 1964) and Nesterov's accelerated gradient descent (Nesterov, 2013). Recently, by revisiting the original convergence analysis of Adam, Reddi et al. (2018) found that for some handcrafted simple convex optimization problem, Adam does not even converge to the global minimizer. In order to address this convergence issue of Adam, Reddi et al. (2018) proposed a new variant of the Adam algorithm namely AMSGrad, which has guaranteed convergence in the convex optimization setting. The update rule of AMSGrad is as follows¹:

$$\mathbf{x}_{t+1} = \mathbf{x}_t - \alpha_t \frac{\mathbf{m}_t}{\sqrt{\widehat{\mathbf{v}}_t}}, \text{ with } \widehat{\mathbf{v}}_t = \max(\widehat{\mathbf{v}}_{t-1}, \mathbf{v}_t),$$

where $\alpha_t > 0$ is the step size, $\mathbf{x}_t \in \mathbb{R}^d$ is the iterate in the t-th iteration, and $\mathbf{m}_t, \mathbf{v}_t \in \mathbb{R}^d$ are the exponential moving averages of the gradient and the squared gradient at the t-th iteration respectively. More specifically, \mathbf{m}_t and \mathbf{v}_t are defined as follows²:

$$\mathbf{m}_{t} = \beta_{1} \mathbf{m}_{t-1} + (1 - \beta_{1}) \mathbf{g}_{t}, \quad \mathbf{v}_{t} = \beta_{2} \mathbf{v}_{t-1} + (1 - \beta_{2}) \mathbf{g}_{t}^{2}$$

where $\beta_1 \in [0,1]$ and $\beta_2 \in [0,1]$ are hyperparameters of the algorithm, and \mathbf{g}_t is the stochastic gradient at the t-th iteration. However, Wilson et al. (2017) found that for over-parameterized neural networks, training with Adam or its variants typically generalizes worse than SGD, even when the training performance is better. In particular, they found that carefully-tuned SGD with momentum, weight decay and appropriate learning rate decay strategies can significantly outperform adaptive gradient algorithms in terms of test error. This problem is often referred to as the generalization gap for adaptive gradient methods. In order to close this generalization gap of Adam and AMSGrad, Chen and Gu (2018) proposed a partially adaptive momentum estimation method (Padam). Instead of scaling the gradient by $\widehat{\mathbf{v}}_t$, this method chooses to scale the gradient by $\widehat{\mathbf{v}}_t^{-p}$, where $p \in (0,1]$ is a hyper parameter. This gives rise to the following update formula³:

$$\mathbf{x}_{t+1} = \mathbf{x}_t - \alpha_t \frac{\mathbf{m}_t}{\widehat{\mathbf{v}}_t^p}, \text{ with } \widehat{\mathbf{v}}_t = \max(\widehat{\mathbf{v}}_{t-1}, \mathbf{v}_t).$$

Evidently, when p = 1/2, Padam reduces to AMSGrad. Padam also reduces to the corrected version of RMSProp (Reddi et al., 2018) when p = 1/2 and $\beta_1 = 0$.

Despite the successes of adaptive gradient methods for training deep neural networks, the convergence guarantees for these algorithms are mostly restricted to online convex optimization (Duchi et al., 2011; Kingma and Ba, 2014; Reddi et al., 2018; Chen and Gu, 2018). Therefore, there

¹With slight abuse of notation, here we denote by $\sqrt{\mathbf{v}_t}$ the element-wise square root of the vector \mathbf{v}_t , by $\mathbf{m}_t/\sqrt{\mathbf{v}_t}$ the element-wise division between \mathbf{m}_t and $\sqrt{\mathbf{v}_t}$, and by $\max(\hat{\mathbf{v}}_{t-1}, \mathbf{v}_t)$ the element-wise maximum between $\hat{\mathbf{v}}_{t-1}$ and \mathbf{v}_t .

²Here we denote by \mathbf{g}_t^2 the element-wise square of the vector \mathbf{g}_t .

³We denote by $\hat{\mathbf{v}}_t^{-p} = [v_{t,1}^{-p}, \cdots, v_{t,d}^{-p}]^{\top}$ the element-wise -p-th power of the vector $\hat{\mathbf{v}}_t$

is a huge gap between existing online convex optimization guarantees for adaptive gradient methods and the empirical successes of adaptive gradient methods in nonconvex optimization. In order to bridge this gap, there are a few recent attempts to prove the nonconvex optimization guarantees for adaptive gradient methods. More specifically, Basu et al. (2018) proved the convergence rate of RMSProp and Adam when using deterministic gradient rather than stochastic gradient. Li and Orabona (2018) achieves convergence rate of AdaGrad, assuming the gradient is *L*-Lipschitz continuous. Ward et al. (2018) proved the convergence rate of a simplified AdaGrad where the moving average of the norms of the gradient vectors is used to adjust the gradient vector in both deterministic and stochastic settings for smooth nonconvex functions. Nevertheless, the convergence guarantees in Basu et al. (2018); Ward et al. (2018) are still limited to simplified algorithms. Another attempt to obtain the convergence rate under stochastic setting is prompted recently by Zou and Shen (2018), in which they only focus on the condition when the momentum vanishes.

In this paper, we provide a sharp convergence analysis of the adaptive gradient methods. In particular, we analyze the state-of-the-art adaptive gradient method, i.e., Padam (Chen and Gu, 2018), and prove its convergence rate for smooth nonconvex objective functions in the stochastic optimization setting. Our results directly imply the convergence rates for AMSGrad (the corrected version of Adam) and the corrected version of RMSProp (Reddi et al., 2018). Our analyses can be extended to other adaptive gradient methods such as AdaGrad, AdaDelta (Zeiler, 2012) and Nadam (Dozat, 2016) mentioned above, but we omit these extensions in this paper for the sake of conciseness. It is worth noting that our convergence analysis emphasizes equally on the dependence of number of iterations T and dimension d in the convergence rate. This is motivated by the fact that modern machine learning methods, especially the training of deep neural networks, usually requires solving a very high-dimensional nonconvex optimization problem. The order of dimension d is usually comparable to or even larger than the total number of iterations T. Take training the latest convolutional neural network DenseNet-BC (Huang et al., 2017) with depth L=100 and growth rate k=12 on CIFAR-10 (Krizhevsky, 2009) as an example. According to Huang et al. (2017), the network is trained with in total 0.28 million iterations, however the number of parameters in the network is 0.8 million. This example shows that d can indeed be in the same order of T in practice. Therefore, we argue that it is very important to show the precise dependence on both T and d in the convergence analysis of adaptive gradient methods for modern machine learning.

When we were preparing this manuscript, we noticed that there was a paper (Chen et al., 2018) released on arXiv on August 8th, 2018, which analyzes the convergence of a class of Adamtype algorithms including AMSGrad and AdaGrad for nonconvex optimization. Our work is an independent work, and our derived convergence rate for AMSGrad is faster than theirs.

1.1 Our Contributions

The main contributions of our work are summarized as follows:

• We prove that the convergence rate of Padam to a stationary point for stochastic nonconvex optimization is

$$O\left(\frac{\left(\sum_{i=1}^{d} \|\mathbf{g}_{1:T,i}\|_{2}\right)^{1/2}}{T^{3/4}} + \frac{d}{T}\right),\tag{1.1}$$

where $\mathbf{g}_1, \dots, \mathbf{g}_T$ are the stochastic gradients and $\mathbf{g}_{1:T,i} = [g_{1,i}, g_{2,i}, \dots, g_{T,i}]^{\top}$. When the stochastic gradients are ℓ_{∞} -bounded, (1.1) matches the convergence rate of vanilla SGD in terms of the rate of T.

• Our result implies the convergence rate for AMSGrad is

$$O\left(\sqrt{\frac{d}{T}} + \frac{d}{T}\right),$$

which has a better dependence on the dimension d and T than the convergence rate proved in Chen et al. (2018), i.e.,

$$O\left(\frac{\log T + d^2}{\sqrt{T}}\right).$$

1.2 Additional Related Work

Here we briefly review other related work on nonconvex stochastic optimization.

Ghadimi and Lan (2013) proposed a randomized stochastic gradient (RSG) method, and proved its $O(1/\sqrt{T})$ convergence rate to a stationary point. Ghadimi and Lan (2016) proposed an randomized stochastic accelerated gradient (RSAG) method, which achieves $O(1/T + \sigma^2/\sqrt{T})$ convergence rate, where σ^2 is an upper bound on the variance of the stochastic gradient. Motivated by the success of stochastic momentum methods in deep learning (Sutskever et al., 2013), Yang et al. (2016) provided a unified convergence analysis for both stochastic heavy-ball method and the stochastic variant of Nesterov's accelerated gradient method, and proved $O(1/\sqrt{T})$ convergence rate to a stationary point for smooth nonconvex functions. Reddi et al. (2016); Allen-Zhu and Hazan (2016) proposed variants of stochastic variance-reduced gradient (SVRG) method (Johnson and Zhang, 2013) that is provably faster than gradient descent in the nonconvex finite-sum setting. Lei et al. (2017) proposed a stochastically controlled stochastic gradient (SCSG), which further improves convergence rate of SVRG for finite-sum smooth nonconvex optimization. Very recently, Zhou et al. (2018) proposed a new algorithm called stochastic nested variance-reduced gradient (SNVRG), which achieves strictly better gradient complexity than both SVRG and SCSG for finite-sum and stochastic smooth nonconvex optimization.

There is another line of research in stochastic smooth nonconvex optimization, which makes use of the λ -nonconvexity of a nonconvex function f (i.e., $\nabla^2 f \succeq -\lambda \mathbf{I}$). More specifically, Natasha 1 (Allen-Zhu, 2017b) and Natasha 1.5 (Allen-Zhu, 2017a) have been proposed, which solve a modified regularized problem and achieve faster convergence rate to first-order stationary points than SVRG and SCSG in the finite-sum and stochastic settings respectively. In addition, Allen-Zhu (2018) proposed an SGD4 algorithm, which optimizes a series of regularized problems, and is able to achieve a faster convergence rate than SGD.

1.3 Organization and Notation

The remainder of this paper is organized as follows: We present the problem setup and review the algorithms in Section 2. We provide the convergence guarantee of Padam for stochastic smooth nonconvex optimization in Section 3. Finally, we conclude our paper in Section 4.

Notation. Scalars are denoted by lower case letters, vectors by lower case bold face letters, and matrices by upper case bold face letters. For a vector $\mathbf{x} = [x_i] \in \mathbb{R}^d$, we denote the ℓ_p norm $(p \ge 1)$ of \mathbf{x} by $\|\mathbf{x}\|_p = \left(\sum_{i=1}^d |x_i|^p\right)^{1/p}$, the ℓ_∞ norm of \mathbf{x} by $\|\mathbf{x}\|_\infty = \max_{i=1}^d |x_i|$. For a sequence of vectors $\{\mathbf{g}_j\}_{j=1}^t$, we denote by $g_{j,i}$ the i-th element in \mathbf{g}_j . We also denote $\mathbf{g}_{1:t,i} = [g_{1,i}, g_{2,i}, \dots, g_{t,i}]^\top$. With slightly abuse of notation, for any two vectors \mathbf{a} and \mathbf{b} , we denote \mathbf{a}^2 as the element-wise square, \mathbf{a}^p as the element-wise power operation, \mathbf{a}/\mathbf{b} as the element-wise division and $\max(\mathbf{a}, \mathbf{b})$ as the element-wise maximum. For a matrix $\mathbf{A} = [A_{ij}] \in \mathbb{R}^{d \times d}$, we define $\|\mathbf{A}\|_{1,1} = \sum_{i,j=1}^d |A_{ij}|$. Given two sequences $\{a_n\}$ and $\{b_n\}$, we write $a_n = O(b_n)$ if there exists a constant $0 < C < +\infty$ such that $a_n \le C b_n$. We use notation $\widetilde{O}(\cdot)$ to hide logarithmic factors.

2 Problem Setup and Algorithms

In this section, we first introduce the preliminary definitions used in this paper, followed by the problem setup of stochastic nonconvex optimization. Then we review the state-of-the-art adaptive gradient method, i.e., Padam (Chen and Gu, 2018), along with AMSGrad (the corrected version of Adam) (Reddi et al., 2018) and the corrected version of RMSProp (Tieleman and Hinton, 2012; Reddi et al., 2018).

2.1 Problem Setup

We study the following stochastic nonconvex optimization problem

$$\min_{\mathbf{x} \in \mathbb{R}^d} f(\mathbf{x}) := \mathbb{E}_{\xi} [f(\mathbf{x}; \xi)],$$

where ξ is a random variable satisfying certain distribution, $f(\mathbf{x}; \xi) : \mathbb{R}^d \to \mathbb{R}$ is a L-smooth nonconvex function. In the stochastic setting, one cannot directly access the full gradient of $f(\mathbf{x})$. Instead, one can only get unbiased estimators of the gradient of $f(\mathbf{x})$, which is $\nabla f(\mathbf{x}; \xi)$. This setting has been studied in Ghadimi and Lan (2013, 2016).

2.2 Algorithms

In this section we introduce the algorithms we study in this paper. We mainly consider three algorithms: Padam (Chen and Gu, 2018), AMSGrad (Reddi et al., 2018) and a corrected version of RMSProp (Tieleman and Hinton, 2012; Reddi et al., 2018).

The Padam algorithm is given in Algorithm 1. It is originally proposed by Chen and Gu (2018) to improve the generalization performance of adaptive gradient methods. As is shown in Algorithm 1, the learning rate of Padam is $\alpha_t \hat{\mathbf{V}}_t^{-p}$, where p is a partially adaptive parameter. With this parameter p, Padam unifies AMSGrad and SGD with momentum, and gives a general framework of algorithms with exponential moving average. Padam reduces to the AMSGrad algorithm when p = 1/2. If p = 1/2 and $\beta_1 = 0$, Padam reduces to a corrected version of the RMSProp algorithm given by Reddi et al. (2018). As important special cases of Padam, we show AMSGrad and the corrected version of RMSProp in Algorithms 2 and 3 respectively.

```
Algorithm 1 Partially adaptive momentum estimation method (Padam) (Chen and Gu, 2018)
```

```
Input: \mathbf{x}_1, step size \{\alpha_t\}_{t=1}^T, \beta_1, \beta_2, p.
   1: \mathbf{m}_0 \leftarrow 0, \hat{\mathbf{v}}_0 \leftarrow 0, \mathbf{v}_0 \leftarrow 0
   2: for t = 1 to T do
              \mathbf{g}_t = \nabla f(\mathbf{x}_t, \xi_t)
            \mathbf{m}_t = \beta_1 \mathbf{m}_{t-1} + (1 - \beta_1) \mathbf{g}_t
             \mathbf{v}_t = \beta_2 \mathbf{v}_{t-1} + (1 - \beta_2) \mathbf{g}_t^2
              \hat{\mathbf{v}}_t = \max(\hat{\mathbf{v}}_{t-1}, \mathbf{v}_t)
              \mathbf{x}_{t+1} = \mathbf{x}_t - \alpha_t \widehat{\mathbf{V}}_t^{-p} \mathbf{m}_t \text{ with } \widehat{\mathbf{V}}_t = \operatorname{diag}(\widehat{\mathbf{v}}_t)
   8: end for
Output: Choose \mathbf{x}_{\text{out}} from \{\mathbf{x}_t\}, 2 \leq t \leq T with probability \alpha_{t-1}/(\sum_{i=1}^{T-1} \alpha_i).
```

Algorithm 2 AMSGrad (Reddi et al., 2018)

```
Input: \mathbf{x}_1, step size \{\alpha_t\}_{t=1}^T, \beta_1, \beta_2.
   1: \mathbf{m}_0 \leftarrow 0, \widehat{\mathbf{v}}_0 \leftarrow 0, \mathbf{v}_0 \leftarrow 0
   2: for t = 1 to T do
              \mathbf{g}_t = \nabla f(\mathbf{x}_t, \xi_t)
            \mathbf{m}_t = \beta_1 \mathbf{m}_{t-1} + (1 - \beta_1) \mathbf{g}_t
            \mathbf{v}_t = \beta_2 \mathbf{v}_{t-1} + (1 - \beta_2) \mathbf{g}_t^2
            \widehat{\mathbf{v}}_t = \max(\widehat{\mathbf{v}}_{t-1}, \mathbf{v}_t)
            \mathbf{x}_{t+1} = \mathbf{x}_t - \alpha_t \hat{\mathbf{V}}_t^{-1/2} \mathbf{m}_t \text{ with } \hat{\mathbf{V}}_t = \operatorname{diag}(\hat{\mathbf{v}}_t)
   8: end for
Output: Choose \mathbf{x}_{\text{out}} from \{\mathbf{x}_t\}, 2 \leq t \leq T with probability \alpha_{t-1} / \sum_{i=1}^{T-1} \alpha_i.
```

Algorithm 3 RMSProp (Tieleman and Hinton, 2012) (corrected version by Reddi et al. (2018))

```
Input: \mathbf{x}_1, step size \{\alpha_t\}_{t=1}^T, \beta.
   1: \hat{\mathbf{v}}_0 \leftarrow 0, \, \mathbf{v}_0 \leftarrow 0
   2: for t = 1 to T do
   3: \mathbf{g}_t = \nabla f(\mathbf{x}_t, \xi_t)
   4: \mathbf{v}_t = \beta \mathbf{v}_{t-1} + (1 - \beta) \mathbf{g}_t^2
   5: \widehat{\mathbf{v}}_t = \max(\widehat{\mathbf{v}}_{t-1}, \mathbf{v}_t)
               \mathbf{x}_{t+1} = \mathbf{x}_t - \alpha_t \hat{\mathbf{V}}_t^{-1/2} \mathbf{g}_t \text{ with } \hat{\mathbf{V}}_t = \operatorname{diag}(\hat{\mathbf{v}}_t)
```

Output: Choose \mathbf{x}_{out} from $\{\mathbf{x}_t\}, 2 \leq t \leq T$ with probability $\alpha_{t-1} / \sum_{i=1}^{T-1} \alpha_i$.

3 Main Theory

In this section we present our main theoretical results. We first introduce the following assumptions.

Assumption 3.1 (Bounded Gradient). $f(\mathbf{x}) = \mathbb{E}_{\xi} f(\mathbf{x}; \xi)$ has G_{∞} -bounded stochastic gradient. That is, for any ξ , we assume that

$$\|\nabla f(\mathbf{x};\xi)\|_{\infty} \leq G_{\infty}.$$

It is worth mentioning that Assumption 3.1 is slightly weaker than the ℓ_2 -boundedness assumption $\|\nabla f(\mathbf{x};\xi)\|_2 \leq G_2$ used in Reddi et al. (2016); Chen et al. (2018). Since $\|\nabla f(\mathbf{x};\xi)\|_{\infty} \leq \|\nabla f(\mathbf{x};\xi)\|_2 \leq \sqrt{d}\|\nabla f(\mathbf{x};\xi)\|_{\infty}$, the ℓ_2 -boundedness assumption implies Assumption 3.1 with $G_{\infty} = G_2$. Meanwhile, G_{∞} will be tighter than G_2 by a factor of \sqrt{d} when each coordinate of $\nabla f(\mathbf{x};\xi)$ almost equals to each other.

Assumption 3.2 (*L*-smooth). $f(\mathbf{x}) = \mathbb{E}_{\xi} f(\mathbf{x}; \xi)$ is *L*-smooth: for any $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$, we have

$$|f(\mathbf{x}) - f(\mathbf{y}) - \langle \nabla f(\mathbf{y}), \mathbf{x} - \mathbf{y} \rangle| \le \frac{L}{2} ||\mathbf{x} - \mathbf{y}||_2^2.$$

Assumption 3.2 is a standard assumption frequently used in analysis of gradient-based algorithms. It is equivalent to the *L*-gradient Lipschitz condition, which is often written as $\|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\|_2 \le L\|\mathbf{x} - \mathbf{y}\|_2$.

We are now ready to present our main result.

Theorem 3.3 (Padam). In Algorithm 1, suppose that $p \in [0, 1/2]$, $\beta_1 < \beta_2^{2p}$ and $\alpha_t = \alpha$ for t = 1, ..., T. Then under Assumptions 3.1 and 3.2, for any $q \in [\max\{0, 4p - 1\}, 1]$, the output \mathbf{x}_{out} of Algorithm 1 satisfies that

$$\mathbb{E}\Big[\|\nabla f(\mathbf{x}_{\text{out}})\|_{2}^{2} \Big] \le \frac{M_{1}}{T\alpha} + \frac{M_{2}d}{T} + \frac{M_{3}\alpha d^{q}}{T^{(1-q)/2}} \mathbb{E}\left(\sum_{i=1}^{d} \|\mathbf{g}_{1:T,i}\|_{2}\right)^{1-q}, \tag{3.1}$$

where

$$M_{1} = 2G_{\infty}^{2p} \Delta f, \quad M_{2} = \frac{4G_{\infty}^{2+2p} \mathbb{E} \|\widehat{\mathbf{v}}_{1}^{-p}\|_{1}}{d(1 - \beta_{1})} + 4G_{\infty}^{2},$$

$$M_{3} = \frac{4LG_{\infty}^{1+q-2p}}{(1 - \beta_{2})^{2p}} + \frac{8LG_{\infty}^{1+q-2p}(1 - \beta_{1})}{(1 - \beta_{2})^{2p}(1 - \beta_{1}/\beta_{2}^{2p})} \left(\frac{\beta_{1}}{1 - \beta_{1}}\right)^{2},$$

and $\Delta f = f(\mathbf{x}_1) - \inf_{\mathbf{x}} f(\mathbf{x}).$

Remark 3.4. From Theorem 3.3, we can see that M_1 and M_3 are independent of the number of iterations T and dimension d. In addition, if $\|\widehat{\mathbf{v}}_1^{-1}\|_{\infty} = O(1)$, it is easy to see that M_2 also has an upper bound that is independent of T and d.

The following corollary is a special case of Theorem 3.3 when $p \in [0, 1/4]$ and q = 0.

Corollary 3.5. Under the same conditions of Theorem 3.3, if $p \in [0, 1/4]$, then the output of Padam satisifies

$$\mathbb{E}\left[\left\|\nabla f(\mathbf{x}_{\text{out}})\right\|_{2}^{2}\right] \leq \frac{M_{1}}{T\alpha} + \frac{M_{2}d}{T} + \frac{M_{3}'\alpha}{\sqrt{T}}\mathbb{E}\left(\sum_{i=1}^{a} \|\mathbf{g}_{1:T,i}\|_{2}\right),\tag{3.2}$$

where M_1 and M_2 and Δf are the same as in Theorem 3.3, and M'_3 is defined as follows:

$$M_3' = \frac{4LG_{\infty}^{1-2p}}{(1-\beta_2)^{2p}} + \frac{8LG_{\infty}^{1-2p}(1-\beta_1)}{(1-\beta_2)^{2p}(1-\beta_1/\beta_2^{2p})} \left(\frac{\beta_1}{1-\beta_1}\right)^2.$$

Remark 3.6. Corollary 3.5 simplifies the result of Theorem 3.3 by choosing q=0 under the condition $p \in [0,1/4]$. We remark that this choice of q is optimal in an important special case studied in Duchi et al. (2011); Reddi et al. (2018): when the gradient vectors are sparse, we assume that $\sum_{i=1}^{d} \|\mathbf{g}_{1:T,i}\|_2 \ll \sqrt{dT}$. Then for q > 0, it follows that

$$\frac{\sum_{i=1}^{d} \|\mathbf{g}_{1:T,i}\|_{2}}{T} \ll \frac{d^{q} \left(\sum_{i=1}^{d} \|\mathbf{g}_{1:T,i}\|_{2}\right)^{1-q}}{T^{1-q/2}}.$$
(3.3)

(3.3) implies that the upper bound provided by (3.2) is strictly better than (3.1) with q > 0. Therefore when the gradient vectors are sparse, Padam achieves faster convergence when p is located in [0, 1/4].

Remark 3.7. We show the convergence rate under different choices of step size α . If

$$\alpha = \Theta\left(T^{1/4}\left(\sum_{i=1}^{d} \|\mathbf{g}_{1:T,i}\|_{2}\right)^{1/2}\right)^{-1},$$

then by (3.2), we have

$$\mathbb{E}\left[\left\|\nabla f(\mathbf{x}_{\text{out}})\right\|_{2}^{2}\right] = O\left(\frac{\left(\sum_{i=1}^{d} \|\mathbf{g}_{1:T,i}\|_{2}\right)^{1/2}}{T^{3/4}} + \frac{d}{T}\right). \tag{3.4}$$

Note that the convergence rate given by (3.4) is related to the sum of gradient norms $\sum_{i=1}^{d} \|\mathbf{g}_{1:T,i}\|_2$. As is mentioned in Remark 3.6, when the stochastic gradients $\mathbf{g}_{1:T,i}$, $i=1,\ldots,d$ are sparse, we follow the assumption given by Duchi et al. (2011) that $\sum_{i=1}^{d} \|\mathbf{g}_{1:T,i}\|_2 \ll \sqrt{dT}$. More specifically, suppose $\sum_{i=1}^{d} \|\mathbf{g}_{1:T,i}\|_2 = O(d^s\sqrt{T})$ for some $0 \le s \le 1/2$. We have

$$\mathbb{E}\left[\left\|\nabla f(\mathbf{x}_{\text{out}})\right\|_{2}^{2}\right] = O\left(\frac{d^{s/2}}{T^{1/2}} + \frac{d}{T}\right).$$

When s = 1/2, we have

$$\mathbb{E}\Big[\left\|\nabla f(\mathbf{x}_{\text{out}})\right\|_{2}^{2}\Big] = O\left(\frac{d^{1/4}}{\sqrt{T}} + \frac{d}{T}\right),$$

which matches the rate $O(1/\sqrt{T})$ achieved by nonconvex SGD (Ghadimi and Lan, 2016), considering the dependence of T.

Remark 3.8. If we set $\alpha = 1/\sqrt{T}$ which is not related to $\sum_{i=1}^{d} \|\mathbf{g}_{1:T,i}\|_2$, then (3.2) suggests that

$$\mathbb{E}\left[\left\|\nabla f(\mathbf{x}_{\text{out}})\right\|_{2}^{2}\right] = O\left(\frac{1}{\sqrt{T}} + \frac{d}{T} + \frac{1}{T}\sum_{i=1}^{d}\|\mathbf{g}_{1:T,i}\|_{2}\right). \tag{3.5}$$

When $\sum_{i=1}^{d} \|\mathbf{g}_{1:T,i}\|_{2} \ll \sqrt{dT}$ (Duchi et al., 2011; Reddi et al., 2018), by (3.5) we have

$$\mathbb{E}\left[\left\|\nabla f(\mathbf{x}_{\text{out}})\right\|_{2}^{2}\right] = O\left(\frac{1}{\sqrt{T}} + \frac{d}{T} + \sqrt{\frac{d}{T}}\right),\,$$

which matches the convergence result in nonconvex SGD (Ghadimi and Lan, 2016) considering the dependence of T.

Next we show the convergence analysis of two popular algorithms: AMSGrad and RMSProp. Since AMSGrad and RMSProp can be seen as two specific instances of Padam, we can apply Theorem 3.3 with specific parameter choice, and obtain the following two corollaries.

Corollary 3.9 (AMSGrad). Under the same conditions of Theorem 3.3, for AMSGrad in Algorithm 2, if $\alpha_t = \alpha = 1/\sqrt{dT}$ for t = 1, ..., T, then the output \mathbf{x}_{out} satisfies that

$$\mathbb{E}\left[\left\|\nabla f(\mathbf{x}_{\text{out}})\right\|_{2}^{2}\right] \leq \frac{M_{1}^{A}\sqrt{d}}{\sqrt{T}} + \frac{M_{2}^{A}d}{T} + \frac{M_{3}^{A}\sqrt{d}}{\sqrt{T}},\tag{3.6}$$

where $\{M_i^A\}_{i=1}^3$ are defined as follows:

$$M_1^A = 2G_{\infty}\Delta f, \quad M_2^A = \frac{4G_{\infty}^3 \mathbb{E} \|\widehat{\mathbf{v}}_1^{-1/2}\|_1}{d(1-\beta_1)} + 4G_{\infty}^2,$$
$$M_3^A = \frac{4LG_{\infty}}{(1-\beta_2)} + \frac{8LG_{\infty}(1-\beta_1)}{(1-\beta_2)(1-\beta_1/\beta_2)} \left(\frac{\beta_1}{1-\beta_1}\right)^2.$$

Remark 3.10. As what has been illustrated in Theorem 3.3, $\{M_i^A\}_{i=1}^3$ are independent of T and essentially independent of d. Thus, (3.6) implies that AMSGrad achieves

$$O\bigg(\sqrt{\frac{d}{T}} + \frac{d}{T}\bigg)$$

convergence rate, which matches the convergence rate of nonconvex SGD (Ghadimi and Lan, 2016). Chen et al. (2018) also provided similar bound for AMSGrad. They showed that

$$\mathbb{E}\Big[\left\|\nabla f(\mathbf{x}_{\text{out}})\right\|_{2}^{2}\Big] = O\left(\frac{\log T + d^{2}}{\sqrt{T}}\right).$$

It can be seen that the dependence of d in their bound is quadratic, which is worse than the linear dependence implied by (3.6). Moreover, by Corollary 3.5, Corollary 3.9 and (3.3), it is easy to see that Padam with $p \in [0, 1/4]$ is faster than AMSGrad where p = 1/2, which backups the experimental results in Chen and Gu (2018).

Corollary 3.11 (corrected version of RMSProp). Under the same conditions of Theorem 3.3, for RMSProp in Algorithm 3, if $\alpha_t = \alpha = 1/\sqrt{dT}$ for t = 1, ..., T, then the output \mathbf{x}_{out} satisfies that

$$\mathbb{E}\left[\left\|\nabla f(\mathbf{x}_{\text{out}})\right\|_{2}^{2}\right] \leq \frac{M_{1}^{R}\sqrt{d}}{\sqrt{T}} + \frac{M_{2}^{R}d}{T} + \frac{M_{3}^{R}\sqrt{d}}{\sqrt{T}},\tag{3.7}$$

where $\{M_i^R\}_{i=1}^3$ are defined in the following:

$$M_1^R = 2G_{\infty}\Delta f$$
, $M_2^R = 4G_{\infty}^3 \mathbb{E} \| \hat{\mathbf{v}}_1^{-1/2} \|_1 / d + 4G_{\infty}^2$, $M_3^R = \frac{4LG_{\infty}}{(1 - \beta_2)}$.

Remark 3.12. $\{M_i^R\}_{i=1}^3$ are independent of T and essentially independent of d. Thus, (3.7) implies that RMSProp achieves $O(\sqrt{d/T} + d/T)$ convergence rate, which matches the convergence rate of nonconvex SGD given by Ghadimi and Lan (2016).

4 Conclusions

In this paper, we provided a sharp analysis of the state-of-the-art adaptive gradient method Padam (Chen and Gu, 2018), and proved its convergence rate for smooth nonconvex optimization. Our results directly imply the convergence rates of AMSGrad and the corrected version of RMSProp for smooth nonconvex optimization. In terms of the number of iterations T, the derived convergence rates in this paper match the $O(1/\sqrt{T})$ rate achieved by SGD; in terms of dimension d, our results give better rate than existing work. Our results also offer some insights into the choice of the partially adaptive parameter p in the Padam algorithm: when the gradients are sparse, Padam with $p \in [0, 1/4]$ achieves the fastest convergence rate. This theoretically backups the experimental results in existing work (Chen and Gu, 2018).

Acknowledgement

We would like to thank Jinghui Chen for discussion on this work.

A Proof of the Main Theory

Here we provide the detailed proof of the main theorem.

A.1 Proof of Theorem 3.3

Let $\mathbf{x}_0 = \mathbf{x}_1$. To prove Theorem 3.3, we need the following lemmas:

Lemma A.1 (Restatement of Lemma). Let $\hat{\mathbf{v}}_t$ and \mathbf{m}_t be as defined in Algorithm 1. Then under Assumption 3.1, we have $\|\nabla f(\mathbf{x})\|_{\infty} \leq G_{\infty}$, $\|\hat{\mathbf{v}}_t\|_{\infty} \leq G_{\infty}^2$ and $\|\mathbf{m}_t\|_{\infty} \leq G_{\infty}$.

Lemma A.2. Suppose that f has G_{∞} -bounded stochastic gradient. Let β_1, β_2 be the weight parameters, $\alpha_t, t = 1, ..., T$ be the step sizes in Algorithm 1 and $q \in [\max\{4p-1, 0\}, 1]$. We denote

 $\gamma = \beta_1/\beta_2^{2p}$. Suppose that $\alpha_t = \alpha$ and $\gamma \leq 1$, then under Assumption 3.1, we have the following two results:

$$\sum_{t=1}^{T} \alpha_t^2 \mathbb{E} \Big[\| \widehat{\mathbf{V}}_t^{-p} \mathbf{m}_t \|_2^2 \Big] \le \frac{T^{(1+q)/2} d^q \alpha^2 (1-\beta_1) G_{\infty}^{(1+q-4p)}}{(1-\beta_2)^{2p} (1-\gamma)} \mathbb{E} \bigg(\sum_{i=1}^{d} \| \mathbf{g}_{1:T,i} \|_2 \bigg)^{1-q},$$

and

$$\sum_{t=1}^{T} \alpha_t^2 \mathbb{E} \Big[\big\| \widehat{\mathbf{V}}_t^{-p} \mathbf{g}_t \big\|_2^2 \Big] \leq \frac{T^{(1+q)/2} d^q \alpha^2 G_{\infty}^{(1+q-4p)}}{(1-\beta_2)^{2p}} \mathbb{E} \bigg(\sum_{i=1}^{d} \| \mathbf{g}_{1:T,i} \|_2 \bigg)^{1-q}.$$

To deal with stochastic momentum \mathbf{m}_t and stochastic weight $\widehat{\mathbf{V}}_t^{-p}$, following Yang et al. (2016), we define an auxiliary sequence \mathbf{z}_t as follows: let $\mathbf{x}_0 = \mathbf{x}_1$, and for each $t \geq 1$,

$$\mathbf{z}_{t} = \mathbf{x}_{t} + \frac{\beta_{1}}{1 - \beta_{1}} (\mathbf{x}_{t} - \mathbf{x}_{t-1}) = \frac{1}{1 - \beta_{1}} \mathbf{x}_{t} - \frac{\beta_{1}}{1 - \beta_{1}} \mathbf{x}_{t-1}.$$
(A.1)

Lemma A.3 shows that $\mathbf{z}_{t+1} - \mathbf{z}_t$ can be represented in two different ways.

Lemma A.3. Let \mathbf{z}_t be defined in (A.1). For $t \geq 2$, we have

$$\mathbf{z}_{t+1} - \mathbf{z}_t = \frac{\beta_1}{1 - \beta_1} \left[\mathbf{I} - \left(\alpha_t \widehat{\mathbf{V}}_t^{-p} \right) \left(\alpha_{t-1} \widehat{\mathbf{V}}_{t-1}^{-p} \right)^{-1} \right] (\mathbf{x}_{t-1} - \mathbf{x}_t) - \alpha_t \widehat{\mathbf{V}}_t^{-p} \mathbf{g}_t. \tag{A.2}$$

and

$$\mathbf{z}_{t+1} - \mathbf{z}_t = \frac{\beta_1}{1 - \beta_1} \left(\alpha_{t-1} \widehat{\mathbf{V}}_{t-1}^{-p} - \alpha_t \widehat{\mathbf{V}}_t^{-p} \right) \mathbf{m}_{t-1} - \alpha_t \widehat{\mathbf{V}}_t^{-p} \mathbf{g}_t.$$
(A.3)

For t = 1, we have

$$\mathbf{z}_2 - \mathbf{z}_1 = -\alpha_1 \widehat{\mathbf{V}}_1^{-p} \mathbf{g}_1. \tag{A.4}$$

By Lemma A.3, we connect $\mathbf{z}_{t+1} - \mathbf{z}_t$ with $\mathbf{x}_{t+1} - \mathbf{x}_t$ and $\alpha_t \widehat{\mathbf{V}}_t^{-p} \mathbf{g}_t$

The following two lemmas give bounds on $\|\mathbf{z}_{t+1} - \mathbf{z}_t\|_2$ and $\|\nabla f(\mathbf{z}_t) - \nabla f(\mathbf{x}_t)\|_2$, which play important roles in our proof.

Lemma A.4. Let \mathbf{z}_t be defined in (A.1). For $t \geq 2$, we have

$$\|\mathbf{z}_{t+1} - \mathbf{z}_t\|_2 \le \|\alpha \widehat{\mathbf{V}}_t^{-p} \mathbf{g}_t\|_2 + \frac{\beta_1}{1 - \beta_1} \|\mathbf{x}_{t-1} - \mathbf{x}_t\|_2.$$

Lemma A.5. Let \mathbf{z}_t be defined in (A.1). For $t \geq 2$, we have

$$\|\nabla f(\mathbf{z}_t) - \nabla f(\mathbf{x}_t)\|_2 \le L\left(\frac{\beta_1}{1-\beta_1}\right) \cdot \|\mathbf{x}_t - \mathbf{x}_{t-1}\|_2.$$

Now we are ready to prove Theorem 3.3.

Proof of Theorem 3.3. Since f is L-smooth, we have:

$$f(\mathbf{z}_{t+1}) \leq f(\mathbf{z}_t) + \nabla f(\mathbf{z}_t)^{\top} (\mathbf{z}_{t+1} - \mathbf{z}_t) + \underbrace{\frac{L}{2} \|\mathbf{z}_{t+1} - \mathbf{z}_t\|_{2}^{2}}_{= f(\mathbf{z}_t) + \underbrace{\nabla f(\mathbf{x}_t)^{\top} (\mathbf{z}_{t+1} - \mathbf{z}_t)}_{I_1} + \underbrace{\frac{(\nabla f(\mathbf{z}_t) - \nabla f(\mathbf{x}_t))^{\top} (\mathbf{z}_{t+1} - \mathbf{z}_t)}_{I_2} + \underbrace{\frac{L}{2} \|\mathbf{z}_{t+1} - \mathbf{z}_t\|_{2}^{2}}_{I_2}$$
(A.5)

In the following, we bound I_1 , I_2 and I_3 separately.

Bounding term I_1 : When t = 1, we have

$$\nabla f(\mathbf{x}_1)^{\top}(\mathbf{z}_2 - \mathbf{z}_1) = -\nabla f(\mathbf{x}_1)^{\top} \alpha_1 \widehat{\mathbf{V}}_t^{-p} \mathbf{g}_1. \tag{A.6}$$

For $t \geq 2$, we have

$$\nabla f(\mathbf{x}_{t})^{\top} (\mathbf{z}_{t+1} - \mathbf{z}_{t})$$

$$= \nabla f(\mathbf{x}_{t})^{\top} \left[\frac{\beta_{1}}{1 - \beta_{1}} (\alpha_{t-1} \widehat{\mathbf{V}}_{t-1}^{-p} - \alpha_{t} \widehat{\mathbf{V}}_{t}^{-p}) \mathbf{m}_{t-1} - \alpha_{t} \widehat{\mathbf{V}}_{t}^{-p} \mathbf{g}_{t} \right]$$

$$= \frac{\beta_{1}}{1 - \beta_{1}} \nabla f(\mathbf{x}_{t})^{\top} (\alpha_{t-1} \widehat{\mathbf{V}}_{t-1}^{-p} - \alpha_{t} \widehat{\mathbf{V}}_{t}^{-p}) \mathbf{m}_{t-1} - \nabla f(\mathbf{x}_{t})^{\top} \alpha_{t} \widehat{\mathbf{V}}_{t}^{-p} \mathbf{g}_{t}, \tag{A.7}$$

where the first equality holds due to (A.3) in Lemma A.3. For $\nabla f(\mathbf{x}_t)^{\top} (\alpha_{t-1} \widehat{\mathbf{V}}_{t-1}^{-p} - \alpha_t \widehat{\mathbf{V}}_t^{-p}) \mathbf{m}_{t-1}$ in (A.7), we have

$$\nabla f(\mathbf{x}_{t})^{\top} (\alpha_{t-1} \widehat{\mathbf{V}}_{t-1}^{-p} - \alpha_{t} \widehat{\mathbf{V}}_{t}^{-p}) \mathbf{m}_{t-1} \leq \|\nabla f(\mathbf{x}_{t})\|_{\infty} \cdot \|\alpha_{t-1} \widehat{\mathbf{V}}_{t-1}^{-p} - \alpha_{t} \widehat{\mathbf{V}}_{t}^{-p}\|_{1,1} \cdot \|\mathbf{m}_{t-1}\|_{\infty}$$

$$\leq G_{\infty}^{2} \left[\|\alpha_{t-1} \widehat{\mathbf{V}}_{t-1}^{-p}\|_{1,1} - \|\alpha_{t} \widehat{\mathbf{V}}_{t}^{-p}\|_{1,1} \right]$$

$$= G_{\infty}^{2} \left[\|\alpha_{t-1} \widehat{\mathbf{V}}_{t-1}^{-p}\|_{1} - \|\alpha_{t} \widehat{\mathbf{V}}_{t}^{-p}\|_{1} \right]. \tag{A.8}$$

The first inequality holds because for a positive diagonal matrix \mathbf{A} , we have $\mathbf{x}^{\top}\mathbf{A}\mathbf{y} \leq \|\mathbf{x}\|_{\infty} \cdot \|\mathbf{A}\|_{1,1} \cdot \|\mathbf{y}\|_{\infty}$. The second inequality holds due to $\alpha_{t-1}\widehat{\mathbf{V}}_{t-1}^{-p} \succeq \alpha_t\widehat{\mathbf{V}}_t^{-p} \succeq 0$. Next we bound $-\nabla f(\mathbf{x}_t)^{\top}\alpha_t\widehat{\mathbf{V}}_t^{-p}\mathbf{g}_t$. We have

$$-\nabla f(\mathbf{x}_{t})^{\top} \alpha_{t} \widehat{\mathbf{V}}_{t}^{-p} \mathbf{g}_{t}$$

$$= -\nabla f(\mathbf{x}_{t})^{\top} \alpha_{t-1} \widehat{\mathbf{V}}_{t-1}^{-p} \mathbf{g}_{t} - \nabla f(\mathbf{x}_{t})^{\top} \left(\alpha_{t} \widehat{\mathbf{V}}_{t}^{-p} - \alpha_{t-1} \widehat{\mathbf{V}}_{t-1}^{-p}\right) \mathbf{g}_{t}$$

$$\leq -\nabla f(\mathbf{x}_{t})^{\top} \alpha_{t-1} \widehat{\mathbf{V}}_{t-1}^{-p} \mathbf{g}_{t} + \|\nabla f(\mathbf{x}_{t})\|_{\infty} \cdot \|\alpha_{t} \widehat{\mathbf{V}}_{t}^{-p} - \alpha_{t-1} \widehat{\mathbf{V}}_{t-1}^{-p}\|_{1,1} \cdot \|\mathbf{g}_{t}\|_{\infty}$$

$$\leq -\nabla f(\mathbf{x}_{t})^{\top} \alpha_{t-1} \widehat{\mathbf{V}}_{t-1}^{-p} \mathbf{g}_{t} + G_{\infty}^{2} \left(\|\alpha_{t-1} \widehat{\mathbf{V}}_{t-1}^{-p}\|_{1,1} - \|\alpha_{t} \widehat{\mathbf{V}}_{t}^{-p}\|_{1,1} \right)$$

$$= -\nabla f(\mathbf{x}_{t})^{\top} \alpha_{t-1} \widehat{\mathbf{V}}_{t-1}^{-p} \mathbf{g}_{t} + G_{\infty}^{2} \left(\|\alpha_{t-1} \widehat{\mathbf{V}}_{t-1}^{-p}\|_{1} - \|\alpha_{t} \widehat{\mathbf{V}}_{t}^{-p}\|_{1} \right). \tag{A.9}$$

The first inequality holds because for a positive diagonal matrix \mathbf{A} , we have $\mathbf{x}^{\top} \mathbf{A} \mathbf{y} \leq \|\mathbf{x}\|_{\infty} \cdot \|\mathbf{A}\|_{1,1} \cdot \|\mathbf{y}\|_{\infty}$. The second inequality holds due to $\alpha_{t-1} \hat{\mathbf{V}}_{t-1}^{-p} \succeq \alpha_t \hat{\mathbf{V}}_t^{-p} \succeq 0$. Substituting (A.8) and (A.9)

into (A.7), we have

$$\nabla f(\mathbf{x}_t)^{\top} (\mathbf{z}_{t+1} - \mathbf{z}_t) \le -\nabla f(\mathbf{x}_t)^{\top} \alpha_{t-1} \widehat{\mathbf{V}}_{t-1}^{-p} \mathbf{g}_t + \frac{1}{1 - \beta_1} G_{\infty}^2 \left(\left\| \alpha_{t-1} \widehat{\mathbf{v}}_{t-1}^{-p} \right\|_1 - \left\| \alpha_t \widehat{\mathbf{v}}_t^{-p} \right\|_1 \right). \tag{A.10}$$

Bounding term I_2 : For $t \ge 1$, we have

$$(\nabla f(\mathbf{z}_{t}) - \nabla f(\mathbf{x}_{t}))^{\top} (\mathbf{z}_{t+1} - \mathbf{z}_{t})$$

$$\leq \|\nabla f(\mathbf{z}_{t}) - \nabla f(\mathbf{x}_{t})\|_{2} \cdot \|\mathbf{z}_{t+1} - \mathbf{z}_{t}\|_{2}$$

$$\leq (\|\alpha_{t} \widehat{\mathbf{V}}_{t}^{-p} \mathbf{g}_{t}\|_{2} + \frac{\beta_{1}}{1 - \beta_{1}} \|\mathbf{x}_{t-1} - \mathbf{x}_{t}\|_{2}) \cdot \frac{\beta_{1}}{1 - \beta_{1}} \cdot L \|\mathbf{x}_{t} - \mathbf{x}_{t-1}\|_{2}$$

$$= L \frac{\beta_{1}}{1 - \beta_{1}} \|\alpha_{t} \widehat{\mathbf{V}}_{t}^{-p} \mathbf{g}_{t}\|_{2} \cdot \|\mathbf{x}_{t} - \mathbf{x}_{t-1}\|_{2} + L \left(\frac{\beta_{1}}{1 - \beta_{1}}\right)^{2} \|\mathbf{x}_{t} - \mathbf{x}_{t-1}\|_{2}^{2}$$

$$\leq L \|\alpha_{t} \widehat{\mathbf{V}}_{t}^{-p} \mathbf{g}_{t}\|_{2}^{2} + 2L \left(\frac{\beta_{1}}{1 - \beta_{1}}\right)^{2} \|\mathbf{x}_{t} - \mathbf{x}_{t-1}\|_{2}^{2}, \tag{A.11}$$

where the second inequality holds because of Lemma A.3 and Lemma A.4, the last inequality holds due to Young's inequality.

Bounding term I_3 : For $t \ge 1$, we have

$$\frac{L}{2} \|\mathbf{z}_{t+1} - \mathbf{z}_{t}\|_{2}^{2} \leq \frac{L}{2} \left[\|\alpha_{t} \widehat{\mathbf{V}}_{t}^{-p} \mathbf{g}_{t}\|_{2} + \frac{\beta_{1}}{1 - \beta_{1}} \|\mathbf{x}_{t-1} - \mathbf{x}_{t}\|_{2} \right]^{2} \\
\leq L \|\alpha_{t} \widehat{\mathbf{V}}_{t}^{-p} \mathbf{g}_{t}\|_{2}^{2} + 2L \left(\frac{\beta_{1}}{1 - \beta_{1}} \right)^{2} \|\mathbf{x}_{t-1} - \mathbf{x}_{t}\|_{2}^{2}. \tag{A.12}$$

The first inequality is obtained by introducing Lemma A.3.

For t = 1, substituting (A.6), (A.11) and (A.12) into (A.5), taking expectation and rearranging terms, we have

$$\mathbb{E}[f(\mathbf{z}_{2}) - f(\mathbf{z}_{1})] \\
\leq \mathbb{E}\left[-\nabla f(\mathbf{x}_{1})^{\top} \alpha_{1} \widehat{\mathbf{V}}_{1}^{-p} \mathbf{g}_{1} + 2L \|\alpha_{1} \widehat{\mathbf{V}}_{1}^{-p} \mathbf{g}_{1}\|_{2}^{2} + 4L \left(\frac{\beta_{1}}{1-\beta_{1}}\right)^{2} \|\mathbf{x}_{1} - \mathbf{x}_{0}\|_{2}^{2}\right] \\
= \mathbb{E}[-\nabla f(\mathbf{x}_{1})^{\top} \alpha_{1} \widehat{\mathbf{V}}_{1}^{-p} \mathbf{g}_{1} + 2L \|\alpha_{1} \widehat{\mathbf{V}}_{1}^{-p} \mathbf{g}_{1}\|_{2}^{2}] \\
\leq \mathbb{E}[d\alpha_{1} G_{\infty}^{2-2p} + 2L \|\alpha_{1} \widehat{\mathbf{V}}_{1}^{-p} \mathbf{g}_{1}\|_{2}^{2}], \tag{A.13}$$

where the last inequality holds because

$$-\nabla f(\mathbf{x}_1)^{\top} \widehat{\mathbf{V}}_1^{-p} \mathbf{g}_1 \leq d \cdot \|\nabla f(\mathbf{x}_1)\|_{\infty} \cdot \|\widehat{\mathbf{V}}_1^{-p} \mathbf{g}_1\|_{\infty} \leq d \cdot G_{\infty} \cdot G_{\infty}^{1-2p} = dG_{\infty}^{2-2p}.$$

For $t \geq 2$, substituting (A.10), (A.11) and (A.12) into (A.5), taking expectation and rearranging

terms, we have

$$\mathbb{E}\left[f(\mathbf{z}_{t+1}) + \frac{G_{\infty}^{2} \|\alpha_{t}\widehat{\mathbf{v}}_{t}^{-p}\|_{1}}{1 - \beta_{1}} - \left(f(\mathbf{z}_{t}) + \frac{G_{\infty}^{2} \|\alpha_{t-1}\widehat{\mathbf{v}}_{t-1}^{-p}\|_{1}}{1 - \beta_{1}}\right)\right] \\
\leq \mathbb{E}\left[-\nabla f(\mathbf{x}_{t})^{\top}\alpha_{t-1}\widehat{\mathbf{V}}_{t-1}^{-p}\mathbf{g}_{t} + 2L\|\alpha_{t}\widehat{\mathbf{V}}_{t}^{-p}\mathbf{g}_{t}\|_{2}^{2} + 4L\left(\frac{\beta_{1}}{1 - \beta_{1}}\right)^{2}\|\mathbf{x}_{t} - \mathbf{x}_{t-1}\|_{2}^{2}\right] \\
= \mathbb{E}\left[-\nabla f(\mathbf{x}_{t})^{\top}\alpha_{t-1}\widehat{\mathbf{V}}_{t-1}^{-p}\nabla f(\mathbf{x}_{t}) + 2L\|\alpha_{t}\widehat{\mathbf{V}}_{t}^{-p}\mathbf{g}_{t}\|_{2}^{2} + 4L\left(\frac{\beta_{1}}{1 - \beta_{1}}\right)^{2}\|\alpha_{t-1}\widehat{\mathbf{V}}_{t-1}^{-p}\mathbf{m}_{t-1}\|_{2}^{2}\right] \\
\leq \mathbb{E}\left[-\alpha_{t-1}\|\nabla f(\mathbf{x}_{t})\|_{2}^{2}(G_{\infty}^{2p})^{-1} + 2L\|\alpha_{t}\widehat{\mathbf{V}}_{t}^{-p}\mathbf{g}_{t}\|_{2}^{2} + 4L\left(\frac{\beta_{1}}{1 - \beta_{1}}\right)^{2}\|\alpha_{t-1}\widehat{\mathbf{V}}_{t-1}^{-p}\mathbf{m}_{t-1}\|_{2}^{2}\right], \quad (A.14)$$

where the equality holds because $\mathbb{E}[\mathbf{g}_t] = \nabla f(\mathbf{x}_t)$ conditioned on $\nabla f(\mathbf{x}_t)$ and $\hat{\mathbf{V}}_{t-1}^{-p}$, the second inequality holds because of Lemma A.1. Telescoping (A.14) for t = 2 to T and adding with (A.13), we have

$$(G_{\infty}^{2p})^{-1} \sum_{t=2}^{T} \alpha_{t-1} \mathbb{E} \| \nabla f(\mathbf{x}_{t}) \|_{2}^{2}$$

$$\leq \mathbb{E} \left[f(\mathbf{z}_{1}) + \frac{G_{\infty}^{2} \| \alpha_{1} \widehat{\mathbf{v}}_{1}^{-p} \|_{1}}{1 - \beta_{1}} + d\alpha_{1} G_{\infty}^{2-2p} - \left(f(\mathbf{z}_{T+1}) + \frac{G_{\infty}^{2} \| \alpha_{T} \widehat{\mathbf{v}}_{T}^{-p} \|_{1}}{1 - \beta_{1}} \right) \right]$$

$$+ 2L \sum_{t=1}^{T} \mathbb{E} \| \alpha_{t} \widehat{\mathbf{V}}_{t}^{-p} \mathbf{g}_{t} \|_{2}^{2} + 4L \left(\frac{\beta_{1}}{1 - \beta_{1}} \right)^{2} \sum_{t=2}^{T} \mathbb{E} \left[\| \alpha_{t-1} \widehat{\mathbf{V}}_{t-1}^{-p} \mathbf{m}_{t-1} \|_{2}^{2} \right]$$

$$\leq \mathbb{E} \left[\Delta f + \frac{G_{\infty}^{2} \| \alpha_{1} \widehat{\mathbf{v}}_{1}^{-p} \|_{1}}{1 - \beta_{1}} + d\alpha_{1} G_{\infty}^{2-2p} \right] + 2L \sum_{t=1}^{T} \mathbb{E} \| \alpha_{t} \widehat{\mathbf{V}}_{t}^{-p} \mathbf{g}_{t} \|_{2}^{2}$$

$$+ 4L \left(\frac{\beta_{1}}{1 - \beta_{1}} \right)^{2} \sum_{t=1}^{T} \mathbb{E} \left[\| \alpha_{t} \widehat{\mathbf{V}}_{t}^{-p} \mathbf{m}_{t} \|_{2}^{2} \right]. \tag{A.15}$$

By Lemma A.2, we have

$$\sum_{t=1}^{T} \alpha_t^2 \mathbb{E} \left[\| \widehat{\mathbf{V}}_t^{-p} \mathbf{m}_t \|_2^2 \right] \le \frac{T^{(1+q)/2} d^q \alpha^2 (1-\beta_1) G_{\infty}^{(1+q-4p)}}{(1-\beta_2)^{2p} (1-\gamma)} \mathbb{E} \left(\sum_{i=1}^{d} \| \mathbf{g}_{1:T,i} \|_2 \right)^{1-q}, \tag{A.16}$$

where $\gamma = \beta_1/\beta_2^{2p}$. We also have

$$\sum_{t=1}^{T} \alpha_t^2 \mathbb{E} \left[\| \widehat{\mathbf{V}}_t^{-p} \mathbf{g}_t \|_2^2 \right] \le \frac{T^{(1+q)/2} d^q \alpha^2 G_{\infty}^{(1+q-4p)}}{(1-\beta_2)^{2p}} \mathbb{E} \left(\sum_{i=1}^{d} \| \mathbf{g}_{1:T,i} \|_2 \right)^{1-q}. \tag{A.17}$$

Substituting (A.16) and (A.17) into (A.15), and rearranging (A.15), we have

$$\mathbb{E}\|\nabla f(\mathbf{x}_{\text{out}})\|_{2}^{2} = \frac{1}{\sum_{t=2}^{T} \alpha_{t-1}} \sum_{t=2}^{T} \alpha_{t-1} \mathbb{E}\|\nabla f(\mathbf{x}_{t})\|_{2}^{2} \\
\leq \frac{G_{\infty}^{2p}}{\sum_{t=2}^{T} \alpha_{t-1}} \mathbb{E}\left[\Delta f + \frac{G_{\infty}^{2} \|\alpha_{1} \hat{\mathbf{v}}_{1}^{-p}\|_{1}}{1 - \beta_{1}} + d\alpha_{1} G_{\infty}^{2-2p}\right] \\
+ \frac{2LG_{\infty}^{2p}}{\sum_{t=2}^{T} \alpha_{t-1}} \frac{T^{(1+q)/2} d^{q} \alpha^{2} G_{\infty}^{(1+q-4p)}}{(1 - \beta_{2})^{2p}} \mathbb{E}\left(\sum_{i=1}^{d} \|\mathbf{g}_{1:T,i}\|_{2}\right)^{1-q} \\
+ \frac{4LG_{\infty}^{2p}}{\sum_{t=2}^{T} \alpha_{t-1}} \left(\frac{\beta_{1}}{1 - \beta_{1}}\right)^{2} \frac{T^{(1+q)/2} d^{q} \alpha^{2} (1 - \beta_{1}) G_{\infty}^{(1+q-4p)}}{(1 - \beta_{2})^{2p} (1 - \gamma)} \mathbb{E}\left(\sum_{i=1}^{d} \|\mathbf{g}_{1:T,i}\|_{2}\right)^{1-q} \\
\leq \frac{1}{T\alpha} 2G_{\infty}^{2p} \Delta f + \frac{4}{T} \left(\frac{G_{\infty}^{2+2p} \mathbb{E}\|\hat{\mathbf{v}}_{1}^{-p}\|_{1}}{1 - \beta_{1}} + dG_{\infty}^{2}\right) \\
+ \frac{d^{q} \alpha}{T^{(1-q)/2}} \mathbb{E}\left(\sum_{i=1}^{d} \|\mathbf{g}_{1:T,i}\|_{2}\right)^{1-q} \left(\frac{4LG_{\infty}^{1+q-2p}}{(1 - \beta_{2})^{2p}} + \frac{8LG_{\infty}^{1+q-2p} (1 - \beta_{1})}{(1 - \beta_{2})^{2p} (1 - \gamma)} \left(\frac{\beta_{1}}{1 - \beta_{1}}\right)^{2}\right), \quad (A.18)$$

where the second inequality holds because $\alpha_t = \alpha$. Rearranging (A.18), we obtain

$$\mathbb{E}\|\nabla f(\mathbf{x}_{\text{out}})\|_{2}^{2} \leq \frac{M_{1}}{T\alpha} + \frac{M_{2}d}{T} + \frac{\alpha d^{q}M_{3}}{T^{(1-q)/2}}\mathbb{E}\left(\sum_{i=1}^{d} \|\mathbf{g}_{1:T,i}\|_{2}\right)^{1-q},$$

where $\{M_i\}_{i=1}^3$ are defined in Theorem 3.3. This completes the proof.

A.2 Proof of Corollary 3.5

Proof of Corollary 3.5. From Theorem 3.3, let $p \in [0, 1/4]$, we have $q \in [0, 1]$. Setting q = 0, we have

$$\mathbb{E}\left[\left\|\nabla f(\mathbf{x}_{\text{out}})\right\|_{2}^{2}\right] \leq \frac{M_{1}}{T\alpha} + \frac{M_{2} \cdot d}{T} + \frac{M_{3}'\alpha}{\sqrt{T}}\mathbb{E}\left(\sum_{i=1}^{d} \|\mathbf{g}_{1:T,i}\|_{2}\right),$$

where M_1 and M_2 are defined in Theorem 3.3 and M_3' is defined in Corollary 3.5. This completes the proof.

A.3 Proof of Corollary 3.9

Proof of Corollary 3.9. From Theorem 3.3, we get the conclusion by setting p = 1/2 and q = 1.

A.4 Proof of Corollary 3.11

Proof of Corollary 3.11. From Corollary 3.9, we get the conclusion by further setting $\beta_1 = 0$.

B Proof of Technical Lemmas

B.1 Proof of Lemma A.1

Proof of Lemma A.1. Since f has G_{∞} -bounded stochastic gradient, for any \mathbf{x} and ξ , $\|\nabla f(\mathbf{x};\xi)\|_{\infty} \leq G_{\infty}$. Thus, we have

$$\|\nabla f(\mathbf{x})\|_{\infty} = \|\mathbb{E}_{\xi} \nabla f(\mathbf{x}; \xi)\|_{\infty} \le \mathbb{E}_{\xi} \|\nabla f(\mathbf{x}; \xi)\|_{\infty} \le G_{\infty}.$$

Next we bound $\|\mathbf{m}_t\|_{\infty}$. We have $\|\mathbf{m}_0\|_{\infty} = 0 \le G_{\infty}$. Suppose that $\|\mathbf{m}_t\|_{\infty} \le G_{\infty}$, then for \mathbf{m}_{t+1} , we have

$$\|\mathbf{m}_{t+1}\|_{\infty} = \|\beta_1 \mathbf{m}_t + (1 - \beta_1) \mathbf{g}_{t+1}\|_{\infty}$$

$$\leq \beta_1 \|\mathbf{m}_t\|_{\infty} + (1 - \beta_1) \|\mathbf{g}_{t+1}\|_{\infty}$$

$$\leq \beta_1 G_{\infty} + (1 - \beta_1) G_{\infty}$$

$$= G_{\infty}.$$

Thus, for any $t \geq 0$, we have $\|\mathbf{m}_t\|_{\infty} \leq G_{\infty}$. Finally we bound $\|\widehat{\mathbf{v}}_t\|_{\infty}$. First we have $\|\mathbf{v}_0\|_{\infty} = \|\widehat{\mathbf{v}}_0\|_{\infty} = 0 \leq G_{\infty}^2$. Suppose that $\|\widehat{\mathbf{v}}_t\|_{\infty} \leq G_{\infty}^2$ and $\|\mathbf{v}_t\|_{\infty} \leq G_{\infty}^2$. Note that we have

$$\|\mathbf{v}_{t+1}\|_{\infty} = \|\beta_2 \mathbf{v}_t + (1 - \beta_2) \mathbf{g}_{t+1}^2\|_{\infty}$$

$$\leq \beta_2 \|\mathbf{v}_t\|_{\infty} + (1 - \beta_2) \|\mathbf{g}_{t+1}^2\|_{\infty}$$

$$\leq \beta_2 G_{\infty}^2 + (1 - \beta_2) G_{\infty}^2$$

$$= G_{\infty}^2,$$

and by definition, we have $\|\widehat{\mathbf{v}}_{t+1}\|_{\infty} = \max\{\|\widehat{\mathbf{v}}_t\|_{\infty}, \|\mathbf{v}_{t+1}\|_{\infty}\} \leq G_{\infty}^2$. Thus, for any $t \geq 0$, we have $\|\widehat{\mathbf{v}}_t\|_{\infty} \leq G_{\infty}^2$.

B.2 Proof of Lemma A.2

Proof. Recall that $\hat{v}_{t,j}, m_{t,j}, g_{t,j}$ denote the j-th coordinate of $\hat{\mathbf{v}}_t, \mathbf{m}_t$ and \mathbf{g}_t . We have

$$\alpha_{t}^{2} \mathbb{E} \left[\| \widehat{\mathbf{V}}_{t}^{-p} \mathbf{m}_{t} \|_{2}^{2} \right] = \alpha_{t}^{2} \mathbb{E} \left[\sum_{i=1}^{d} \frac{m_{t,i}^{2}}{\widehat{v}_{t,i}^{2p}} \right]$$

$$\leq \alpha_{t}^{2} \mathbb{E} \left[\sum_{i=1}^{d} \frac{m_{t,i}^{2}}{v_{t,i}^{2p}} \right]$$

$$= \alpha_{t}^{2} \mathbb{E} \left[\sum_{i=1}^{d} \frac{(\sum_{j=1}^{t} (1 - \beta_{1}) \beta_{1}^{t-j} g_{j,i})^{2}}{(\sum_{j=1}^{t} (1 - \beta_{2}) \beta_{2}^{t-j} g_{j,i}^{2})^{2p}} \right],$$
(B.1)

where the first inequality holds because $\hat{v}_{t,i} \geq v_{t,i}$. Next we have

$$\alpha_{t}^{2} \mathbb{E} \left[\sum_{i=1}^{d} \frac{\left(\sum_{j=1}^{t} (1 - \beta_{1}) \beta_{1}^{t-j} g_{j,i} \right)^{2}}{\left(\sum_{j=1}^{t} (1 - \beta_{2}) \beta_{2}^{t-j} g_{j,i}^{2} \right)^{2p}} \right] \\
\leq \frac{\alpha_{t}^{2} (1 - \beta_{1})^{2}}{(1 - \beta_{2})^{2p}} \mathbb{E} \left[\sum_{i=1}^{d} \frac{\left(\sum_{j=1}^{t} \beta_{1}^{t-j} |g_{j,i}|^{(1+q-4p)} \right) \left(\sum_{j=1}^{t} \beta_{1}^{t-j} |g_{j,i}|^{(1-q+4p)} \right)}{\left(\sum_{j=1}^{t} \beta_{2}^{t-j} g_{j,i}^{2} \right)^{2p}} \right] \\
\leq \frac{\alpha_{t}^{2} (1 - \beta_{1})^{2}}{(1 - \beta_{2})^{2p}} \mathbb{E} \left[\sum_{i=1}^{d} \frac{\left(\sum_{j=1}^{t} \beta_{1}^{t-j} G_{\infty}^{(1+q-4p)} \right) \left(\sum_{j=1}^{t} \beta_{1}^{t-j} |g_{j,i}|^{(1-q+4p)} \right)}{\left(\sum_{j=1}^{t} \beta_{2}^{t-j} g_{j,i}^{2} \right)^{2p}} \right] \\
\leq \frac{\alpha_{t}^{2} (1 - \beta_{1}) G_{\infty}^{(1+q-4p)}}{(1 - \beta_{2})^{2p}} \mathbb{E} \left[\sum_{i=1}^{d} \frac{\sum_{j=1}^{t} \beta_{1}^{t-j} |g_{j,i}|^{(1-q+4p)}}{\left(\sum_{j=1}^{t} \beta_{2}^{t-j} g_{j,i}^{2} \right)^{2p}} \right], \tag{B.2}$$

where the first inequality holds due to Cauchy inequality, the second inequality holds because $|g_{j,i}| \leq G_{\infty}$, the last inequality holds because $\sum_{j=1}^{t} \beta_1^{t-j} \leq (1-\beta_1)^{-1}$. Note that

$$\sum_{i=1}^{d} \frac{\sum_{j=1}^{t} \beta_{1}^{t-j} |g_{j,i}|^{(1-q+4p)}}{(\sum_{j=1}^{t} \beta_{2}^{t-j} g_{j,i}^{2})^{2p}} \leq \sum_{i=1}^{d} \sum_{j=1}^{t} \frac{\beta_{1}^{t-j} |g_{j,i}|^{(1-q+4p)}}{(\beta_{2}^{t-j} g_{j,i}^{2})^{2p}} = \sum_{i=1}^{d} \sum_{j=1}^{t} \gamma^{t-j} |g_{j,i}|^{1-q},$$
(B.3)

where the equality holds due to the definition of γ . Substituting (B.2) and (B.3) into (B.1), we have

$$\alpha_t^2 \mathbb{E}\left[\|\widehat{\mathbf{V}}_t^{-p} \mathbf{m}_t\|_2^2\right] \le \frac{\alpha_t^2 (1-\beta_1) G_{\infty}^{(1+q-4p)}}{(1-\beta_2)^{2p}} \mathbb{E}\left[\sum_{i=1}^d \sum_{j=1}^t \gamma^{t-j} |g_{j,i}|^{1-q}\right].$$
(B.4)

Telescoping (B.4) for t = 1 to T, we have

$$\sum_{t=1}^{T} \alpha_{t}^{2} \mathbb{E} \left[\| \hat{\mathbf{V}}_{t}^{-p} \mathbf{m}_{t} \|_{2}^{2} \right] \leq \frac{\alpha^{2} (1 - \beta_{1}) G_{\infty}^{(1+q-4p)}}{(1 - \beta_{2})^{2p}} \mathbb{E} \left[\sum_{t=1}^{T} \sum_{i=1}^{d} \sum_{j=1}^{t} \gamma^{t-j} |g_{j,i}|^{1-q} \right]$$

$$= \frac{\alpha^{2} (1 - \beta_{1}) G_{\infty}^{(1+q-4p)}}{(1 - \beta_{2})^{2p}} \mathbb{E} \left[\sum_{i=1}^{d} \sum_{j=1}^{T} |g_{j,i}|^{1-q} \sum_{t=j}^{T} \gamma^{t-j} \right]$$

$$\leq \frac{\alpha^{2} (1 - \beta_{1}) G_{\infty}^{(1+q-4p)}}{(1 - \beta_{2})^{2p} (1 - \gamma)} \mathbb{E} \left[\sum_{i=1}^{d} \sum_{j=1}^{T} |g_{j,i}|^{1-q} \right].$$
(B.5)

Finally, we have

$$\sum_{i=1}^{d} \sum_{j=1}^{T} |g_{j,i}|^{1-q} \leq \sum_{i=1}^{d} \left(\sum_{j=1}^{T} g_{j,i}^{2} \right)^{(1-q)/2} \cdot T^{(1+q)/2}
= T^{(1+q)/2} \sum_{i=1}^{d} \|\mathbf{g}_{1:T,i}\|_{2}^{1-q}
\leq T^{(1+q)/2} d^{q} \left(\sum_{i=1}^{d} \|\mathbf{g}_{1:T,i}\|_{2} \right)^{1-q},$$
(B.6)

where the first and second inequalities hold due to Hölder's inequality. Substituting (B.6) into (B.5), we have

$$\sum_{t=1}^{T} \alpha_t^2 \mathbb{E} \Big[\| \widehat{\mathbf{V}}_t^{-p} \mathbf{m}_t \|_2^2 \Big] \leq \frac{T^{(1+q)/2} d^q \alpha^2 (1-\beta_1) G_{\infty}^{(1+q-4p)}}{(1-\beta_2)^{2p} (1-\gamma)} \mathbb{E} \bigg(\sum_{i=1}^{d} \| \mathbf{g}_{1:T,i} \|_2 \bigg)^{1-q}.$$

Specifically, taking $\beta_1 = 0$, we have $\mathbf{m}_t = \mathbf{g}_t$, then

$$\sum_{t=1}^T \alpha_t^2 \mathbb{E} \Big[\| \widehat{\mathbf{V}}_t^{-p} \mathbf{g}_t \|_2^2 \Big] \leq \frac{T^{(1+q)/2} d^q \alpha^2 G_{\infty}^{(1+q-4p)}}{(1-\beta_2)^{2p}} \mathbb{E} \bigg(\sum_{i=1}^d \| \mathbf{g}_{1:T,i} \|_2 \bigg)^{1-q}.$$

B.3 Proof of Lemma A.3

Proof. By definition, we have

$$\mathbf{z}_{t+1} = \mathbf{x}_{t+1} + \frac{\beta_1}{1 - \beta_1} (\mathbf{x}_{t+1} - \mathbf{x}_t) = \frac{1}{1 - \beta_1} \mathbf{x}_{t+1} - \frac{\beta_1}{1 - \beta_1} \mathbf{x}_t.$$

Then we have

$$\mathbf{z}_{t+1} - \mathbf{z}_t = \frac{1}{1 - \beta_1} (\mathbf{x}_{t+1} - \mathbf{x}_t) - \frac{\beta_1}{1 - \beta_1} (\mathbf{x}_t - \mathbf{x}_{t-1})$$
$$= \frac{1}{1 - \beta_1} (-\alpha_t \widehat{\mathbf{V}}_t^{-p} \mathbf{m}_t) + \frac{\beta_1}{1 - \beta_1} \alpha_{t-1} \widehat{\mathbf{V}}_{t-1}^{-p} \mathbf{m}_{t-1}.$$

The equities above are based on definition. Then we have

$$\mathbf{z}_{t+1} - \mathbf{z}_{t} = \frac{-\alpha_{t} \widehat{\mathbf{V}}_{t}^{-p}}{1 - \beta_{1}} \Big[\beta_{1} \mathbf{m}_{t-1} + (1 - \beta_{1}) \mathbf{g}_{t} \Big] + \frac{\beta_{1}}{1 - \beta_{1}} \alpha_{t-1} \widehat{\mathbf{V}}_{t-1}^{-p} \mathbf{m}_{t-1}$$

$$= \frac{\beta_{1}}{1 - \beta_{1}} \mathbf{m}_{t-1} (\alpha_{t-1} \widehat{\mathbf{V}}_{t-1}^{-p} - \alpha_{t} \widehat{\mathbf{V}}_{t}^{-p}) - \alpha_{t} \widehat{\mathbf{V}}_{t}^{-p} \mathbf{g}_{t}$$

$$= \frac{\beta_{1}}{1 - \beta_{1}} \alpha_{t-1} \widehat{\mathbf{V}}_{t-1}^{-p} \mathbf{m}_{t-1} \Big[\mathbf{I} - (\alpha_{t} \widehat{\mathbf{V}}_{t}^{-p}) (\alpha_{t-1} \widehat{\mathbf{V}}_{t-1}^{-p})^{-1} \Big] - \alpha_{t} \widehat{\mathbf{V}}_{t}^{-p} \mathbf{g}_{t}$$

$$= \frac{\beta_{1}}{1 - \beta_{1}} \Big[\mathbf{I} - (\alpha_{t} \widehat{\mathbf{V}}_{t}^{-p}) (\alpha_{t-1} \widehat{\mathbf{V}}_{t-1}^{-p})^{-1} \Big] (\mathbf{x}_{t-1} - \mathbf{x}_{t}) - \alpha_{t} \widehat{\mathbf{V}}_{t}^{-p} \mathbf{g}_{t}.$$

The equalities above follow by combining the like terms.

B.4 Proof of Lemma A.4

Proof. By Lemma A.3, we have

$$\|\mathbf{z}_{t+1} - \mathbf{z}_{t}\|_{2} = \left\| \frac{\beta_{1}}{1 - \beta_{1}} \left[\mathbf{I} - (\alpha_{t} \widehat{\mathbf{V}}_{t}^{-p}) (\alpha_{t-1} \widehat{\mathbf{V}}_{t-1}^{-p})^{-1} \right] (\mathbf{x}_{t-1} - \mathbf{x}_{t}) - \alpha_{t} \widehat{\mathbf{V}}_{t}^{-p} \mathbf{g}_{t} \right\|_{2}$$

$$\leq \frac{\beta_{1}}{1 - \beta_{1}} \left\| \mathbf{I} - (\alpha_{t} \widehat{\mathbf{V}}_{t}^{-p}) (\alpha_{t-1} \widehat{\mathbf{V}}_{t-1}^{-p})^{-1} \right\|_{\infty,\infty} \cdot \|\mathbf{x}_{t-1} - \mathbf{x}_{t}\|_{2} + \|\alpha \widehat{\mathbf{V}}_{t}^{-p} \mathbf{g}_{t}\|_{2},$$

where the inequality holds because the term $\beta_1/(1-\beta_1)$ is positive, and triangle inequality. Considering that $\alpha_t \hat{\mathbf{v}}_{t,j}^{-p} \leq \alpha_{t-1} \hat{\mathbf{v}}_{t-1,j}^{-p}$, when p > 0, we have $\left\| \mathbf{I} - (\alpha_t \hat{\mathbf{V}}_t^{-p})(\alpha_{t-1} \hat{\mathbf{V}}_{t-1}^{-p})^{-1} \right\|_{\infty,\infty} \leq 1$. With that fact, the term above can be bound as:

$$\|\mathbf{z}_{t+1} - \mathbf{z}_t\|_2 \le \|\alpha \widehat{\mathbf{V}}_t^{-p} \mathbf{g}_t\|_2 + \frac{\beta_1}{1 - \beta_1} \|\mathbf{x}_{t-1} - \mathbf{x}_t\|_2.$$

This completes the proof.

B.5 Proof of Lemma A.5

Proof. For term $\|\nabla f(\mathbf{z}_t) - \nabla f(\mathbf{x}_t)\|_2$, we have:

$$\begin{split} \|\nabla f(\mathbf{z}_t) - \nabla f(\mathbf{x}_t)\|_2 &\leq L \|\mathbf{z}_t - \mathbf{x}_t\|_2 \\ &\leq L \left\| \frac{\beta_1}{1 - \beta_1} (\mathbf{x}_t - \mathbf{x}_{t-1}) \right\|_2 \\ &\leq L \left(\frac{\beta_1}{1 - \beta_1} \right) \cdot \|\mathbf{x}_t - \mathbf{x}_{t-1}\|_2, \end{split}$$

where the last inequality holds because the term $\beta_1/(1-\beta_1)$ is positive.

References

- Allen-Zhu, Z. (2017a). Natasha 2: Faster non-convex optimization than sgd. arXiv preprint arXiv:1708.08694.
- Allen-Zhu, Z. (2017b). Natasha: Faster non-convex stochastic optimization via strongly non-convex parameter. In *International Conference on Machine Learning*.
- Allen-Zhu, Z. (2018). How to make the gradients small stochastically. arXiv preprint arXiv:1801.02982.
- ALLEN-ZHU, Z. and HAZAN, E. (2016). Variance reduction for faster non-convex optimization. In *International Conference on Machine Learning*.
- Basu, A., De, S., Mukherjee, A. and Ullah, E. (2018). Convergence guarantees for rmsprop and adam in non-convex optimization and their comparison to nesterov acceleration on autoencoders. arXiv preprint arXiv:1807.06766.
- CHEN, J. and Gu, Q. (2018). Closing the generalization gap of adaptive gradient methods in training deep neural networks. arXiv preprint arXiv:1806.06763.
- Chen, X., Liu, S., Sun, R. and Hong, M. (2018). On the convergence of a class of adam-type algorithms for nonconvex optimization. arXiv preprint arXiv:1808.02941.
- DOZAT, T. (2016). Incorporating nesterov momentum into adam.
- Duchi, J., Hazan, E. and Singer, Y. (2011). Adaptive subgradient methods for online learning and stochastic optimization. *Journal of Machine Learning Research* 12 2121–2159.
- GHADIMI, S. and LAN, G. (2013). Stochastic first-and zeroth-order methods for nonconvex stochastic programming. SIAM Journal on Optimization 23 2341–2368.
- GHADIMI, S. and LAN, G. (2016). Accelerated gradient methods for nonconvex nonlinear and stochastic programming. *Mathematical Programming* **156** 59–99.
- Huang, G., Liu, Z., Van Der Maaten, L. and Weinberger, K. Q. (2017). Densely connected convolutional networks. In *CVPR*, vol. 1.
- Johnson, R. and Zhang, T. (2013). Accelerating stochastic gradient descent using predictive variance reduction. In *Advances in neural information processing systems*.
- KINGMA, D. P. and BA, J. (2014). Adam: A method for stochastic optimization. arXiv preprint arXiv:1412.6980.
- Krizhevsky, A. (2009). Learning multiple layers of features from tiny images.
- Lei, L., Ju, C., Chen, J. and Jordan, M. I. (2017). Non-convex finite-sum optimization via scsg methods. In *Advances in Neural Information Processing Systems*.

- LI, X. and Orabona, F. (2018). On the convergence of stochastic gradient descent with adaptive stepsizes. arXiv preprint arXiv:1805.08114.
- MCMAHAN, H. B. and STREETER, M. (2010). Adaptive bound optimization for online convex optimization. arXiv preprint arXiv:1002.4908.
- NESTEROV, Y. (2013). Introductory lectures on convex optimization: A basic course, vol. 87. Springer Science & Business Media.
- POLYAK, B. T. (1964). Some methods of speeding up the convergence of iteration methods. *USSR Computational Mathematics and Mathematical Physics* 4 1–17.
- REDDI, S. J., HEFNY, A., SRA, S., POCZOS, B. and SMOLA, A. (2016). Stochastic variance reduction for nonconvex optimization 314–323.
- REDDI, S. J., KALE, S. and KUMAR, S. (2018). On the convergence of adam and beyond. In *International Conference on Learning Representations*.
- ROBBINS, H. and MONRO, S. (1951). A stochastic approximation method. The Annals of Mathematical Statistics 22 400–407.
- SUTSKEVER, I., MARTENS, J., DAHL, G. and HINTON, G. (2013). On the importance of initialization and momentum in deep learning. In *International conference on machine learning*.
- TIELEMAN, T. and HINTON, G. (2012). Lecture 6.5—RmsProp: Divide the gradient by a running average of its recent magnitude. COURSERA: Neural Networks for Machine Learning.
- WARD, R., Wu, X. and Bottou, L. (2018). Adagrad stepsizes: Sharp convergence over nonconvex landscapes, from any initialization. arXiv preprint arXiv:1806.01811.
- WILSON, A. C., ROELOFS, R., STERN, M., SREBRO, N. and RECHT, B. (2017). The marginal value of adaptive gradient methods in machine learning. In *Advances in Neural Information Processing Systems*.
- Yang, T., Lin, Q. and Li, Z. (2016). Unified convergence analysis of stochastic momentum methods for convex and non-convex optimization. arXiv preprint arXiv:1604.03257.
- ZEILER, M. D. (2012). Adadelta: an adaptive learning rate method. arXiv preprint arXiv:1212.5701
- Zhou, D., Xu, P. and Gu, Q. (2018). Stochastic nested variance reduction for nonconvex optimization. arXiv preprint arXiv:1806.07811.
- ZOU, F. and Shen, L. (2018). On the convergence of adagrad with momentum for training deep neural networks. arXiv preprint arXiv:1808.03408.