

On the Convergence of AdaGrad with Momentum for Training Deep Neural Networks

Fangyu Zou[†] and Li Shen[‡]

[†]Department of Mathematics, Stony Brook University, Stony Brook, USA

[‡]Tencent AI Lab, Shenzhen, China

fangyu.zou@stonybrook.edu, mathshenli@gmail.com

August 10, 2018

Abstract

Adaptive stochastic gradient descent methods, such as AdaGrad, Adam, AdaDelta, Nadam, AMSGrad, *etc.*, have been demonstrated efficacious in solving non-convex stochastic optimization, such as training deep neural networks. However, their convergence rates have not been touched under the non-convex stochastic circumstance except recent breakthrough results on AdaGrad [34] and perturbed AdaGrad [22]. In this paper, we propose two new adaptive stochastic gradient methods called AdaHB and AdaNAG which integrate coordinate-wise AdaGrad with heavy ball momentum and Nesterov accelerated gradient momentum, respectively. The $\mathcal{O}(\frac{\log T}{\sqrt{T}})$ non-asymptotic convergence rates of AdaHB and AdaNAG in non-convex stochastic setting are also jointly characterized by leveraging a newly developed unified formulation of these two momentum mechanisms. In particular, when momentum term vanishes we obtain convergence rate of coordinate-wise AdaGrad in non-convex stochastic setting as a byproduct.

1 Introduction

Let \mathbb{R}^d be the d -dimensional Euclidean space. We revisit the non-convex stochastic optimization:

$$\min_{x \in \mathbb{R}^d} f(x) = \mathbb{E}[\tilde{f}(x, \xi)], \quad (1)$$

where $\mathbb{E}[\tilde{f}(x, \xi)]$ denotes the expectation with respect to some random variable ξ , and the objective function $f : \mathbb{R}^d \rightarrow (-\infty, \infty)$ is a L -smooth function whose gradient is a L -Lipschitz continuous mapping and satisfies

$$\|\nabla f(x) - \nabla f(y)\| \leq L\|x - y\|, \quad \forall x, y \in \mathbb{R}^d. \quad (2)$$

In the whole paper, we assume that f is bounded from below, *i.e.*, $f^* = \min_{\mathbb{R}^d} f(x) > -\infty$.

One typical instance of the above non-convex stochastic problem is the empirical risk minimization:

$$\min_{x \in \mathbb{R}^d} f(x) = \frac{1}{n} \sum_{i=1}^n \tilde{f}_i(x, \xi_i), \quad (3)$$

where n is the number of samples $\{\xi_i\}$. The number n is usually extremely large so that the problem (3) is a sufficient approximation to the problem (1) by the law of large numbers. In general, gradient of f in problems (1) and (3) can only be accessed with noise, denoted as $g(x, \xi)$, due to the absence of true distribution of the random variable ξ and too many samples, respectively. For problem (3),

noise gradient can be estimated as $g(x, \xi) = \nabla \tilde{f}_i(x, \xi_i)$ ¹ for a random index $i \in \{1, \dots, n\}$. In addition, problem (3) arises from many statistical learning [24] and deep learning [14, 20] models with finite training samples, such as the logistic regression, AUC maximization [31], *etc.* When data perturbations/augmentation and dropout techniques [32] are introduced in the training process, the training model falls into problem (1) in which infinite datasets and noise gradient are used [2, 16].

Stochastic Gradient Decent (SGD) [4, 5] is the most popular method to solve the two non-convex problems (1) and (3). The SGD iteratively updates the solution

$$x_{t+1} = x_t - \eta_t g(x_t, \xi_t), \quad (4)$$

where η_t and $g(x_t, \xi_t)$ are the learning rate and the noise gradient estimation in the t -th iteration, respectively. The convergence and convergence rate of vanilla SGD have been established [3, 12, 40]. However, it is well known that the vanilla SGD suffers from poor convergence performance because it is too sensitive to the learning rate and is hard to tune a suitable one. To obtain an easy-to-tune learning rate and improve the robustness of SGD, many attempts have been tried, such as reducing the variance of the noise gradient [7, 17, 21, 38], using adaptive learning rate [10, 18, 37, 30] and using momentum acceleration mechanisms [29, 27]. The best promising variance reducing techniques are summarized as mini-batch [21], SVRG/SAGA [17, 7] and importance sampling [25, 38]. However, SVRG/SAGA technique requires to store/calculate the full batch gradient of f in each epoch. It is unbearable if the memory storage is limited, and even impossible for problem (1) since the number of training samples is infinite. Comparing to variance reducing, adaptive learning rate and momentum mechanisms are more suitable for problems (1) and (3), since they do not require large memory and merely require a little more computation in each iteration comparing with the vanilla SGD method. The both techniques have been widely used and demonstrated efficacy for training deep neural networks [19, 33, 18, 37].

Momentum mechanisms Momentum acceleration mechanisms, including **Heavy Ball** (HB) momentum [29] and **Nesterov Accelerated Gradient** (NAG) momentum [27], have been theoretically and numerically investigated for both convex and non-convex optimization problems [27, 26, 28, 11, 6] for constant global learning rate. Due to the conciseness and effectiveness of momentum accelerated mechanisms, the HB momentum and the NAG momentum have already been embedded in SGD method to speed up the training of deep neural networks [19, 33]. Given initial points $x_1 \in \mathbb{R}^d$, $m_0 = \mathbf{0}$ and $y_1 = x_1$, the stochastic gradient descent methods with HB (SHB) and NAG (SNAG) accelerated mechanisms are respectively defined as below [36]:

$$(\text{SHB}) : \begin{cases} m_t = \mu m_{t-1} - \eta_t g(x_t, \xi_t); \\ x_{t+1} = x_t + m_t, \end{cases} \quad (\text{SNAG}) : \begin{cases} y_{t+1} = x_t - \eta_t g(x_t, \xi_t); \\ x_{t+1} = y_{t+1} + \mu(y_{t+1} - y_t), \end{cases}$$

for $t = 1, 2, \dots$. The constant μ is called the momentum factor and $0 \leq \mu < 1$. Moreover, with a pre-fixed global learning rate $\eta_t = \mathcal{O}(\frac{1}{\sqrt{T}})$, the $\mathcal{O}(\frac{1}{\sqrt{T}})$ convergence rates of SHB and SNAG have been established in a unified analysis framework in non-convex stochastic setting in [36]. In addition, they have numerically illustrated that SHB and SNAG have faster training process and better generalization ability compared with vanilla SGD. Prior to this, convergence rate of SNAG has already been established in [13] individually.

Adaptive learning rate The first popular adaptive stochastic gradient descent method is the AdaGrad [23, 10]. Comparing with the vanilla SGD, the AdaGrad achieves significantly better performance especially when the gradients are sparse. However, when the loss functions are non-convex and the gradients are dense, the performance of AdaGrad deteriorates due to the rapid decay of the learning rate. To address this issue, several variants of AdaGrad, such as RMSProp [15], Adam/AdaMax [18], AdaDelta [37], Nadam [9], FTML [39], *etc.*, have been proposed by leveraging the exponential moving

¹One may also consider the mini-batch case $g(x, \xi) = \frac{1}{|\mathbb{B}|} \sum_{i \in \mathbb{B}} \nabla \tilde{f}_i(x, \xi_i)$ for a random index set $\mathbb{B} \subseteq \{1, 2, \dots, n\}$.

averages of second order moment of past noise gradient and/or momentum accelerated mechanisms to further improve their efficiency. Although these variants of adaptive stochastic methods have been widely used and have achieved great success in practice, especially in the non-convex stochastic setting for training deep neural networks, these algorithms are still short of convergence guarantees even in convex setting. Particularly, a recent paper [30] has pointed out that Adam can be divergent even for convex problems by a concrete counterexample.

On the other hand, several convergent adaptive stochastic gradient descent methods, such as AdaGrad [10, 23], AdaBatch [8], WNgrad [35] and AMSGrad [30], only guarantees convergence limited to the (strongly) convex setting so far. Very recently, Li and Orabona [22] proved for the first time the convergence of an adaptive method with perturbed AdaGrad adaptive learning rate in the non-convex stochastic setting². Ward *et al* [34] independently and almost simultaneously established $\mathcal{O}(\frac{\log T}{\sqrt{T}})$ convergence rate for the original AdaGrad. Later on, global convergences of full batch RMSProp and Adam [1] are established in the non-convex deterministic setting. To the best of our knowledge, there are still no theoretical guarantees for momentum accelerated adaptive stochastic gradient descent algorithms, like Adam, Nadam and AMSGrad, for non-convex stochastic problems (1) and (3).

In this paper we try to close the gap by establishing the convergence rates of adaptive stochastic gradient descent methods with momentum accelerated mechanisms in non-convex stochastic setting. Based on the two stochastic momentum mechanisms, *i.e.*, SHB and SNAG, we obtain two new convergent adaptive stochastic gradient methods, named as AdaHB and AdaNAG. Moreover, we jointly characterize the convergence rates of AdaHB and AdaNAG by reshaping them to a unified iteration scheme. We emphasize that our unified formulation is totally different with the one in [36] which takes a more complicated form of three step recursions. In addition, to make AdaHB and AdaNAG more practical for large scale problem, coordinate-wise adaptive learning rate [10] with low computational cost is used. When momentum parameter $\mu = 0$, AdaHB and AdaNAG reduce to original coordinate-wise AdaGrad in [10], and its convergence rate is directly derived in the non-convex stochastic circumstance as byproduct. We emphasize that our result does not need to perturb the adaptive learning rate as [34]. To the end of this section, we summarize the contributions of this paper as follows:

- We integrate coordinate-wise AdaGrad with momentum mechanisms that yields two typical adaptive stochastic momentum methods called AdaHB and AdaNAG.
- We establish the $\mathcal{O}(\frac{\log T}{\sqrt{T}})$ convergence rates of AdaHB and AdaNAG methods in a unified scheme in non-convex stochastic setting. Our assumptions are natural and minimal.
- We derive convergence rate of the original coordinate-wise AdaGrad in non-convex stochastic setting as byproduct when the momentum factor vanishes.

2 Notations and Preliminaries

Notations We use T to denote the maximum iterations. Random variable in the t -th iteration is denoted by ξ_t . The noise gradient estimation of $f(x)$ at t -th iteration will be denoted by $g(x_t, \xi_t)$ where $\xi_t, t = 1, 2, \dots$ are random variables. We denote $\mathbb{E}[\cdot]$ as the expectation with respect to the underlying probability space and \mathbb{E}_t as the conditional expectation with respect to ξ_t when ξ_1, \dots, ξ_{t-1} are fixed.

In this paper we consider the coordinate-wise updating. The adaptive learning rate η_t is a vector in \mathbb{R}^d . Given a vector $v \in \mathbb{R}^d$ we denote its k -th coordinate by v_k . For a vector x_t in the t -th iteration, the k -th coordinate of x_t is denoted as $x_{t,k}$ by adding a subscript k . This works the same for g_t , m_t , η_t , etc. The k -th coordinate of the gradient $\nabla f(x)$ is denoted by $\nabla_k f(x_t)$. Given two vectors $v, w \in \mathbb{R}^d$, the inner product of v and w is denoted by $\langle v, w \rangle = \sum_{k=1}^d v_k w_k$, while we will also heavily use the coordinate-wise product of v and w and is simply denoted as vw which is another vector in

²The adaptive learning rate has to be perturbed due to the issue of convergence.

\mathbb{R}^d and whose k -th coordinate is given by $(vw)_k = v_k w_k$. Division by a vector and the square root of a vector are defined in this coordinate-wise sense similarly, namely, computations are conducted coordinate-wise. Given a vector $v \in \mathbb{R}^d$, we define the weighted norm

$$\|\nabla f(x)\|_v^2 := \langle \nabla f(x), v \nabla f(x) \rangle = \sum_{k=1}^d v_k |\nabla_k f(x)|^2. \quad (5)$$

The $\|\cdot\|$ without any subscript is Euclidean norm by default. The $\|\cdot\|_1$ is defined as $\|v\|_1 = \sum_{k=1}^d |v_k|$. We use bold letters to denote the constant vectors with the same coordinates. In particular, we denote $\mathbf{0} = (0, 0, \dots, 0) \in \mathbb{R}^d$ and $\boldsymbol{\epsilon} = (\epsilon, \epsilon, \dots, \epsilon) \in \mathbb{R}^d$.

Assumption In the stochastic setting, instead of observing a full gradient $\nabla f(x_t)$ in the t -th iteration, we observe a noise gradient $g(x_t, \xi_t)$ of f . In what follows, we will denote $g(x_t, \xi_t)$ by g_t for the sake of brevity for $t = 1, 2, \dots$. We assume that the random variables ξ_t , $t = 1, 2, \dots$ are independent of each other and also of x_t . We assume additionally that:

(A1) $\mathbb{E}_t[g_t] = \nabla f(x_t)$; that is, g_t is an unbiased estimator of $\nabla f(x_t)$;

(A2) $\mathbb{E}_t\|g_t\|^2 \leq \sigma^2$; that is, g_t has uniform bounded second moment.

Notice that $\mathbb{E}_t\|g_t\|^2 = \mathbb{E}_t\|g_t - \nabla f(x_t)\|^2 + \|\nabla f(x_t)\|^2$. Therefore, assumption (A2) indicates that both the variance of the noise gradient g_t of f and the full batch gradient $\nabla f(x_t)$ of f are bounded.

Unified stochastic momentum scheme In this paragraph, we propose a new unified formulation for SHB and SNAG. By introducing $m_t = y_{t+1} - y_t$ with $m_0 = 0$, SNAG can be equivalently written as

$$(\text{SNAG}) : \begin{cases} m_t = \mu m_{t-1} - \eta_t g_t \\ x_{t+1} = x_t + m_t + \mu(m_t - m_{t-1}) \end{cases}, \quad t = 1, 2, \dots$$

Here, m_t in both SHB and SNAG is called the momentum. Comparing SHB and above SNAG we can see that the difference between SHB and SNAG lies in that SNAG has more weight on the current momentum m_t . We reshape the two forms to the following unified stochastic momentum (USM) form:

$$(\text{USM}) \begin{cases} m_t = \mu m_{t-1} - \eta_t g_t \\ x_{t+1} = x_t + m_t + \lambda \mu(m_t - m_{t-1}) \end{cases}, \quad t = 1, 2, \dots,$$

where $\lambda \geq 0$ is a constant. When $\lambda = 0$, it is the SHB; when $\lambda = 1$, it is the NAG; when $0 < \lambda < 1$, USM can be regarded as an interpolation between SHB and SNAG. Indeed, our theoretical analysis shows that λ remains valid as long as $\lambda \leq 1/(1 - \mu)$. We call the parameter λ the interpolation factor.

Remark 1. In [36], a three steps version of unified stochastic momentum method has been provided

$$\begin{cases} y_{t+1} = x_t - \eta_t g_t \\ y_{t+1}^s = x_t - s \eta_t g_t \\ x_{t+1} = y_{t+1} + \mu(y_{t+1}^s - y_t^s) \end{cases}, \quad t = 1, 2, \dots, \quad (6)$$

where $y_0^s = x_0$. The convergence rate of above iteration has been established if $\eta_t = \mathcal{O}(\frac{1}{\sqrt{t}})$. Notably, above USM is slightly simpler than Eq. (6) and the learning rate η_t in USM is adaptively determined.

3 AdaGrad with Momentum

In above USM algorithm, if we take the adaptive learning rate η_t as [10], we get a class of new adaptive stochastic method with momentum accelerated mechanisms. In this paper we focus on the coordinate-wise adaptive learning rate in [10] which requires low-computational cost and is defined as:

$$\eta_{t,k} = \frac{\eta}{\sqrt{\epsilon + \sum_{i=1}^t g_{i,k}^2}} \quad \text{for } k = 1, 2, \dots, d, \quad (7)$$

where $\eta > 0$ is a fixed parameter and the constant $\epsilon > 0$ is the initial accumulator value added to make sure the denominator is strictly positive. In practice ϵ is usually chosen small. Below, we formally present AdaUSM algorithm which effectively integrates coordinate-wise AdaGrad with USM method.

Algorithm 1 AdaUSM: AdaGrad with unified stochastic momentum

Parameters: Choose $x_1 \in \mathbb{R}^d$, fixed parameter $\eta \geq 0$, momentum factor μ and initial accumulator value $\epsilon \geq 0$. Set $m_0 = \mathbf{0}$, $v_0 = \epsilon$, $0 \leq \mu < 1$ and $0 \leq \lambda \leq 1/(1 - \mu)$.

for $t = 1, 2, \dots, T$ **do**
 Sample a stochastic gradient g_t .
 for $k = 1, 2, \dots, d$ **do**
 $v_{t,k} = v_{t-1,k} + g_{t,k}^2$.
 $m_{t,k} = \mu m_{t-1,k} - \eta g_{t,k} / \sqrt{v_{t,k}}$.
 $x_{t+1,k} = x_{t,k} + m_{t,k} + \lambda \mu (m_{t,k} - m_{t-1,k})$.
 end for
end for

It is easy to see that $\eta_t = \frac{\eta}{\sqrt{v_t}}$ where the division operation is coordinate-wise as Eq. (7). When $\lambda = 0$, AdaUSM reduces to AdaGrad method with heavy ball momentum which we denote as AdaHB for short. When $\lambda = 1$, AdaUSM reduces to AdaGrad method with Nesterov acceleration gradient momentum which we denote as AdaNAG for short. Below, we specify AdaHB and AdaNAG for readability.

Algorithm 2 AdaHB: AdaGrad with HB

Parameters: Choose $x_1 \in \mathbb{R}^d$, fixed parameter $\eta \geq 0$, momentum factor μ and initial accumulator value $\epsilon \geq 0$. Set $m_0 = \mathbf{0}$, $v_0 = \epsilon$.

for $t = 1, 2, \dots, T$ **do**
 Sample a stochastic gradient g_t .
 for $k = 1, 2, \dots, d$ **do**
 $v_{t,k} = v_{t-1,k} + g_{t,k}^2$.
 $m_{t,k} = \mu m_{t-1,k} - \eta g_{t,k} / \sqrt{v_{t,k}}$.
 $x_{t+1,k} = x_{t,k} + m_{t,k}$.
 end for
end for

Algorithm 3 AdaNAG: AdaGrad with NAG

Parameters: Choose $x_1 \in \mathbb{R}^d$, fixed parameter $\eta \geq 0$, momentum factor μ and initial accumulator value $\epsilon \geq 0$. Set $m_0 = \mathbf{0}$, $v_0 = \epsilon$.

for $t = 1, 2, \dots, T$ **do**
 Sample a stochastic gradient g_t .
 for $k = 1, 2, \dots, d$ **do**
 $v_{t,k} = v_{t-1,k} + g_{t,k}^2$.
 $m_{t,k} = \mu m_{t-1,k} - \eta g_{t,k} / \sqrt{v_{t,k}}$.
 $x_{t+1,k} = x_{t,k} + m_{t,k} + \mu (m_{t,k} - m_{t-1,k})$.
 end for
end for

Remark 2. Let $\mu = 0$. AdaUSM, AdaHB and AdaNAG all reduce to the coordinate-wise AdaGrad [10]. Hence, its convergence rate in non-convex stochastic setting can be directly derived via AdaUSM.

Remark 3. The widely used exponential average moving (EMA) technique in Adam, AdaMax and AMSGrad for the gradient can be viewed as a momentum form. Let $m_0 = 0$. EMA method updates

$$(\mathbf{EMA}) : \begin{cases} m_t = \mu m_{t-1} + (1 - \mu) g_t, & t = 1, 2, \dots \\ x_{t+1} = x_t - \eta_t m_t \end{cases}$$

When the stepsize $\eta_t = \eta$ is constant, the EMA is actually the HB with a different stepsize $(1 - \mu)\eta$. To see this, we make $\hat{m}_t = -\eta m_t$ in EMA and obtain $\hat{m}_t = \mu\hat{m}_{t-1} - (1 - \mu)\eta g_t$ and $x_{t+1} = x_t + \hat{m}_t$, which corresponds to taking the stepsize in SHB as $(1 - \mu)\eta$. However, when the stepsize η_t is adaptive, the two forms (i.e., SHB and EMA) are not equivalent. Let $\hat{m}_t = -\eta_t m_t$ in EMA. We have

$$\hat{m}_t = -\eta_t m_{t-1} - (1 - \mu)\eta_t g_t = \hat{m}_{t-1} - (1 - \mu)\eta_t g_t + (\eta_t - \eta_{t-1})m_{t-1}.$$

There is an additional error term $(\eta_t - \eta_{t-1})m_{t-1}$ which vanishes if $\eta_t = \eta_{t-1}$ for all t . More precisely, if we solve out the momentum m_t in terms of the past noise gradient estimations g_1, \dots, g_t and eliminate the momentum term m_t , we have for SHB: $x_{t+1} = x_t - \sum_{i=1}^t \eta_i g_i \mu^{t-i}$, while for EMA: $x_{t+1} = x_t - (1 - \mu)\eta_t \sum_{i=1}^t g_i \mu^{t-i}$. One can clearly see that the key difference lies in whether we use the current stepsize or the past stepsize in the exponential moving average.

3.1 Main results

In this subsection, we theoretically prove the convergence rates of AdaUSM, AdaHB and AdaNAG in non-convex stochastic setting. All the proof procedures are presented in supplementary material.

Theorem 1. Suppose the objective function f is L -smooth and $f^* = \min_{x \in \mathbb{R}^d} f(x) > -\infty$. Let $\{x_t\} \subseteq \mathbb{R}^d$ be a sequence generated by AdaUSM. Assume that the noise gradient g_t in each iteration satisfies the assumptions (A1) and (A2). Then for any $\delta > 0$, the convergence rate of AdaUSM holds with probability at least $1 - \delta^{2/3}$ as below:

$$\min_{1 \leq t \leq T} \|\nabla f(x_t)\|^2 \leq \frac{C(T)}{\delta} = \mathcal{O}\left(\frac{\log T}{\sqrt{T}}\right),$$

$$\text{where } C(T) = \frac{\sqrt{\epsilon d + \sigma^2 T}}{T} \left[\frac{2}{(1 + \lambda\mu)\eta} (f(x_1) - f^*) + \left(\frac{2\eta(1 + 2\lambda)^2 L}{(1 + \lambda\mu)(1 - \mu)^3} + \frac{4\sigma}{1 - \mu} \right) d \log \left(1 + \frac{\sigma^2 T}{\epsilon} \right) \right].$$

When the momentum μ and interpolation factor λ vanish, AdaSUM reduces to the coordinate-wise AdaGrad. Hence, we have the following corollary for coordinate-wise AdaGrad.

Corollary 1. Assuming the same setting as Theorem 1. For any $\delta > 0$, the convergence rate of coordinate-wise AdaGrad holds with probability at least $1 - \delta^{2/3}$ as below:

$$\min_{1 \leq t \leq T} \|\nabla f(x_t)\|^2 \leq \frac{C(T)}{\delta} = \mathcal{O}\left(\frac{\log T}{\sqrt{T}}\right),$$

$$\text{where } C(T) = \frac{\sqrt{\epsilon d + \sigma^2 T}}{T} \left[\frac{2}{\eta} (f(x_1) - f^*) + (2\eta L + 4\sigma) d \log \left(1 + \frac{\sigma^2 T}{\epsilon} \right) \right] = \mathcal{O}\left(\frac{\log T}{\sqrt{T}}\right).$$

Remark 4. Above corollary provides the first convergence rate for native coordinate-wise AdaGrad in non-convex stochastic setting. Recently, convergence rate for AdaGrad with global adaptive learning rate $\eta_t = \eta / \sqrt{\epsilon + \sum_{i=1}^t \|g_i\|^2}$ has been established in [34]. Our argument also works for AdaUSM with the global adaptive learning rate with little modification. The proof is essentially the same as $d = 1$.

Remark 5. In [22], convergence rates of an AdaGrad-like and its coordinate-wise version have been established with the following perturbed adaptive learning rate:

$$\eta_t = \frac{\eta}{(\epsilon + \sum_{i=1}^t \|g_i\|^2)^{\frac{1}{2} + \beta}}, \quad \eta_{t,k} = \frac{\eta}{(\epsilon + \sum_{i=1}^t g_{i,k}^2)^{\frac{1}{2} + \beta}} \text{ for } k = 1, 2, \dots, d,$$

under slightly different assumptions. All that matters is that the perturbation term $\beta > 0$ is unavoidable to make their convergence rate valid. Notably, the perturbed adaptive learning rate decreases to zero faster than the native one in Eq. (7).

4 Conclusion

In this paper, we integrate coordinate-wise AdaGrad with heavy ball and Nesterov accelerate gradient momentum mechanisms which yields two new adaptive stochastic gradient decent methods called AdaHB and AdaNAG, respectively. Their convergence rates are jointly analyzed in a unified framework under the non-convex stochastic circumstance. Moreover, when the momentum vanishes we obtain the convergence rate of coordinate-wise AdaGrad [10] in non-convex stochastic setting as byproduct.

References

- [1] A. Basu, S. De, A. Mukherjee, and E. Ullah. Convergence guarantees for rmsprop and adam in non-convex optimization and their comparison to nesterov acceleration on autoencoders. *arXiv preprint arXiv:1807.06766*, 2018.
- [2] A. Bietti and J. Mairal. Stochastic optimization with variance reduction for infinite datasets with finite sum structure. In *Advances in Neural Information Processing Systems*, pages 1623–1633, 2017.
- [3] L. Bottou. Online learning and stochastic approximations. *On-line learning in neural networks*, 17(9):142, 1998.
- [4] L. Bottou. Large-scale machine learning with stochastic gradient descent. In *Proceedings of COMPSTAT’2010*, pages 177–186. Springer, 2010.
- [5] L. Bottou, F. E. Curtis, and J. Nocedal. Optimization methods for large-scale machine learning. *SIAM Review*, 60(2):223–311, 2018.
- [6] Y. Carmon, J. C. Duchi, O. Hinder, and A. Sidford. Accelerated methods for nonconvex optimization. *SIAM Journal on Optimization*, 28(2):1751–1772, 2018.
- [7] A. Defazio, F. Bach, and S. Lacoste-Julien. Saga: A fast incremental gradient method with support for non-strongly convex composite objectives. In *Advances in neural information processing systems*, pages 1646–1654, 2014.
- [8] A. Défossez and F. Bach. Adabatch: Efficient gradient aggregation rules for sequential and parallel stochastic gradient methods. *arXiv preprint arXiv:1711.01761*, 2017.
- [9] T. Dozat. Incorporating nesterov momentum into adam. In *International Conference on Learning Representations Workshop*, 2016.
- [10] J. Duchi, E. Hazan, and Y. Singer. Adaptive subgradient methods for online learning and stochastic optimization. *Journal of Machine Learning Research*, 12(Jul):2121–2159, 2011.
- [11] E. Ghadimi, H. R. Feyzmahdavian, and M. Johansson. Global convergence of the heavy-ball method for convex optimization. In *Control Conference (ECC), 2015 European*, pages 310–315. IEEE, 2015.
- [12] S. Ghadimi and G. Lan. Stochastic first-and zeroth-order methods for nonconvex stochastic programming. *SIAM Journal on Optimization*, 23(4):2341–2368, 2013.
- [13] S. Ghadimi and G. Lan. Accelerated gradient methods for nonconvex nonlinear and stochastic programming. *Mathematical Programming*, 156(1-2):59–99, 2016.
- [14] I. Goodfellow, Y. Bengio, A. Courville, and Y. Bengio. *Deep learning*, volume 1. MIT press Cambridge, 2016.
- [15] G. Hinton, N. Srivastava, and K. Swersky. Neural networks for machine learning lecture 6a overview of mini-batch gradient descent. *Cited on*, page 14, 2012.
- [16] T. Hofmann, A. Lucchi, S. Lacoste-Julien, and B. McWilliams. Variance reduced stochastic gradient descent with neighbors. In *Advances in Neural Information Processing Systems*, pages 2305–2313, 2015.
- [17] R. Johnson and T. Zhang. Accelerating stochastic gradient descent using predictive variance reduction. In *Advances in neural information processing systems*, pages 315–323, 2013.

- [18] D. P. Kingma and J. Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- [19] A. Krizhevsky, I. Sutskever, and G. E. Hinton. Imagenet classification with deep convolutional neural networks. In *Advances in neural information processing systems*, pages 1097–1105, 2012.
- [20] Y. LeCun, Y. Bengio, and G. Hinton. Deep learning. *nature*, 521(7553):436, 2015.
- [21] M. Li, T. Zhang, Y. Chen, and A. J. Smola. Efficient mini-batch training for stochastic optimization. In *Proceedings of the 20th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 661–670. ACM, 2014.
- [22] X. Li and F. Orabona. On the convergence of stochastic gradient descent with adaptive stepsizes. *arXiv preprint arXiv:1805.08114*, 2018.
- [23] H. B. McMahan and M. Streeter. Adaptive bound optimization for online convex optimization. *COLT 2010*, page 244, 2010.
- [24] N. M. Nasrabadi. Pattern recognition and machine learning. *Journal of electronic imaging*, 16(4):049901, 2007.
- [25] D. Needell, R. Ward, and N. Srebro. Stochastic gradient descent, weighted sampling, and the randomized kaczmarz algorithm. In *Advances in Neural Information Processing Systems*, pages 1017–1025, 2014.
- [26] Y. Nesterov. *Introductory lectures on convex optimization: A basic course*, volume 87. Springer Science & Business Media, 2013.
- [27] Y. E. Nesterov. A method for solving the convex programming problem with convergence rate $O(1/k^2)$. In *Dokl. Akad. Nauk SSSR*, volume 269, pages 543–547, 1983.
- [28] P. Ochs, T. Brox, and T. Pock. ipiasco: Inertial proximal algorithm for strongly convex optimization. *Journal of Mathematical Imaging and Vision*, 53(2):171–181, 2015.
- [29] B. T. Polyak. Some methods of speeding up the convergence of iteration methods. *USSR Computational Mathematics and Mathematical Physics*, 4(5):1–17, 1964.
- [30] S. J. Reddi, S. Kale, and S. Kumar. On the convergence of adam and beyond. In *International Conference on Learning Representations*, 2018.
- [31] D. Sculley. Large scale learning to rank. In *NIPS Workshop on Advances in Ranking*, pages 58–63, 2009.
- [32] N. Srivastava, G. Hinton, A. Krizhevsky, I. Sutskever, and R. Salakhutdinov. Dropout: a simple way to prevent neural networks from overfitting. *The Journal of Machine Learning Research*, 15(1):1929–1958, 2014.
- [33] I. Sutskever, J. Martens, G. Dahl, and G. Hinton. On the importance of initialization and momentum in deep learning. In *International conference on machine learning*, pages 1139–1147, 2013.
- [34] R. Ward, X. Wu, and L. Bottou. Adagrad stepsizes: Sharp convergence over nonconvex landscapes, from any initialization. *arXiv preprint arXiv:1806.01811*, 2018.
- [35] X. Wu, R. Ward, and L. Bottou. Wngrad: Learn the learning rate in gradient descent. *arXiv preprint arXiv:1803.02865*, 2018.
- [36] Y. Yan, T. Yang, Z. Li, Q. Lin, and Y. Yang. A unified analysis of stochastic momentum methods for deep learning. In *IJCAI*, pages 2955–2961, 2018.

- [37] M. D. Zeiler. Adadelta: an adaptive learning rate method. *arXiv preprint arXiv:1212.5701*, 2012.
- [38] P. Zhao and T. Zhang. Stochastic optimization with importance sampling for regularized loss minimization. In *international conference on machine learning*, pages 1–9, 2015.
- [39] S. Zheng and J. T. Kwok. Follow the moving leader in deep learning. In *International Conference on Machine Learning*, pages 4110–4119, 2017.
- [40] Z. Zhou, P. Mertikopoulos, N. Bambos, S. Boyd, and P. W. Glynn. Stochastic mirror descent in variationally coherent optimization problems. In *Advances in Neural Information Processing Systems*, pages 7040–7049, 2017.

Supplementary Material for “ On the Convergence of AdaGrad with Momentum for Training Deep Neural Networks ”

In this supplementary we give a complete proof of Theorem 1. The material arranges as follows. In Section A, we provide preliminary lemmas that will be used to establish the main result of this paper. In Section B, we give the detailed proof of Theorem 1.

A Preliminary lemmas

In this section we provide preliminary lemmas that will be used to prove our main theorem. The readers may skip this part for the first time and come back whenever the lemmas are needed.

Lemma 1. *Let $S_t = S_0 + \sum_{i=1}^t a_i$ where $\{a_t\}$ is a non-negative sequence and $S_0 > 0$. We have $\sum_{t=1}^T \frac{a_t}{S_t} \leq \log(S_T) - \log(S_0)$.*

Proof. The finite sum $\sum_{t=1}^T \frac{a_t}{S_t}$ can be interpreted as a Riemann sum as follows $\sum_{t=1}^T \frac{1}{S_t} (S_t - S_{t-1})$. Since $1/x$ is decreasing on the interval $(0, \infty)$, we have

$$\sum_{t=1}^T \frac{1}{S_t} (S_t - S_{t-1}) \leq \int_{S_0}^{S_T} \frac{1}{x} dx = \log(S_T) - \log(S_0).$$

The proof is finished. □

The following lemma is a direct result of the momentum updating rule.

Lemma 2. *Suppose $m_t = \mu m_{t-1} - \eta_t g_t$ with $m_0 = \mathbf{0}$ and $0 \leq \mu < 1$. We have the following estimate*

$$\sum_{t=1}^T \|m_t\|^2 \leq \frac{1}{(1-\mu)^2} \sum_{t=1}^T \|\eta_t g_t\|^2.$$

Proof. If $\mu = 0$, the equality directly holds due to $m_t = -\eta_t g_t$. Otherwise, $0 < \mu < 1$. Since $m_t = \mu m_{t-1} - \eta_t g_t$, for any $\theta > 0$ we have

$$\begin{aligned} \|m_t\|^2 &= \|\mu m_{t-1}\|^2 + \|\eta_t g_t\|^2 - 2\langle \mu m_{t-1}, \eta_t g_t \rangle \\ &\leq \|\mu m_{t-1}\|^2 + \|\eta_t g_t\|^2 + \theta \|\mu m_{t-1}\|^2 + 1/\theta \|\eta_t g_t\|^2 \\ &= (1+\theta)\mu^2 \|m_{t-1}\|^2 + (1+1/\theta) \|\eta_t g_t\|^2. \end{aligned}$$

In particular, we can pick θ in the above so that $(1+\theta)\mu = 1$, namely, $\theta = 1/\mu - 1$. Then $1+1/\theta = 1/(1-\mu)$. Hence, $\|m_t\|^2 \leq \mu \|m_{t-1}\|^2 + 1/(1-\mu) \|\eta_t g_t\|^2$. Dividing both sides by μ^t , we get

$$\frac{\|m_t\|^2}{\mu^t} \leq \frac{\|m_{t-1}\|^2}{\mu^{t-1}} + \frac{1}{1-\mu} \frac{\|\eta_t g_t\|^2}{\mu^t}.$$

Note that $m_0 = \mathbf{0}$. Hence, $\frac{\|m_t\|^2}{\mu^t} \leq \frac{1}{1-\mu} \sum_{i=1}^t \|\eta_i g_i\|^2 \mu^{-i}$. Then multiplying both sides by μ^t we obtain

$$\|m_t\|^2 \leq \frac{1}{1-\mu} \sum_{i=1}^t \|\eta_i g_i\|^2 \mu^{t-i}.$$

Take summation of above inequality over $t = 1, 2, \dots, T$, we have

$$\sum_{t=1}^T \|m_t\|^2 \leq \frac{1}{1-\mu} \sum_{t=1}^T \sum_{i=1}^t \|\eta_i g_i\|^2 \mu^{t-i} = \frac{1}{1-\mu} \sum_{i=1}^T \sum_{t=i}^T \|\eta_i g_i\|^2 \mu^{t-i} \leq \frac{1}{(1-\mu)^2} \sum_{i=1}^T \|\eta_i g_i\|^2.$$

The equality above is due to a double-sum trick. The second inequality is due to the following fact of geometric series

$$\sum_{t=0}^N \mu^t \leq \sum_{t=0}^{\infty} \mu^t = \frac{1}{1-\mu} \quad (\text{for } 0 < \mu < 1).$$

We finish the proof. \square

The following lemma is a result of USM formulation for any general adaptive learning rate.

Lemma 3. *Suppose the objective function f is L -smooth. Let $\{x_t\}$ and $\{m_t\}$ be sequences generated by the following general SGD with USM momentum: starting from initial values x_1 and $m_0 = 0$, update through*

$$\begin{cases} m_t = \mu m_{t-1} - \eta_t g_t \\ x_{t+1} = x_t + m_t + \lambda \mu (m_t - m_{t-1}), \end{cases}$$

where the momentum factor μ and the interpolation factor λ satisfy $0 \leq \mu < 1$ and $0 \leq \lambda \leq 1/(1-\mu)$, respectively. Assume that the noise gradients g_t satisfy the assumptions (A1) and (A2). For any $t \geq 2$, we have the following estimate

$$\langle \nabla f(x_t), m_t \rangle \leq \mu \langle \nabla f(x_{t-1}), m_{t-1} \rangle + (1 + \frac{3}{2} \lambda \mu) \mu L \|m_{t-1}\|^2 + \frac{1}{2} \lambda \mu^2 L \|m_{t-2}\|^2 - \langle \nabla f(x_t), \eta_t g_t \rangle. \quad (8)$$

In particular, the following estimate holds

$$\langle \nabla f(x_t), m_t \rangle \leq (1 + 2\lambda) L \sum_{i=1}^{t-1} \|m_i\|^2 \mu^{t-i} - \sum_{i=1}^t \langle \nabla f(x_i), \eta_i g_i \rangle \mu^{t-i}. \quad (9)$$

Proof. Since $m_t = \mu m_{t-1} - \eta_t g_t$, we have

$$\begin{aligned} \langle \nabla f(x_t), m_t \rangle &= \mu \langle \nabla f(x_t), m_{t-1} \rangle - \langle \nabla f(x_t), \eta_t g_t \rangle \\ &= \mu \langle \nabla f(x_{t-1}), m_{t-1} \rangle + \mu \langle \nabla f(x_t) - \nabla f(x_{t-1}), m_{t-1} \rangle - \langle \nabla f(x_t), \eta_t g_t \rangle \end{aligned} \quad (10)$$

Note that by L -smoothness of the function f , we have

$$\begin{aligned} \|\nabla f(x_t) - \nabla f(x_{t-1})\| &\leq L \|x_t - x_{t-1}\| = L \|m_{t-1} + \lambda \mu (m_{t-1} - m_{t-2})\| \\ &\leq (1 + \lambda \mu) L \|m_{t-1}\| + \lambda \mu L \|m_{t-2}\|. \end{aligned} \quad (11)$$

Hence, by Cauchy-Schwartz inequality and Eq. (11), we have

$$\begin{aligned} \langle \nabla f(x_t) - \nabla f(x_{t-1}), m_{t-1} \rangle &\leq \|\nabla f(x_t) - \nabla f(x_{t-1})\| \|m_{t-1}\| \\ &\leq (1 + \lambda \mu) L \|m_{t-1}\|^2 + \lambda \mu L \|m_{t-2}\| \|m_{t-1}\| \\ &\leq (1 + \frac{3}{2} \lambda \mu) L \|m_{t-1}\|^2 + \frac{1}{2} \lambda \mu L \|m_{t-2}\|^2. \end{aligned} \quad (12)$$

Combining Eq. (10) and Eq. (12) we get the desired inequality Eq. (8).

To get the second estimate, let $B_t = \langle \nabla f(x_t), m_t \rangle$. If $\mu = 0$, the equality holds trivially. Otherwise $0 < \mu < 1$. We divide μ^t from both sides of Eq. (8) and obtain

$$\frac{B_t}{\mu^t} \leq \frac{B_{t-1}}{\mu^{t-1}} + (1 + \frac{3}{2} \lambda \mu) L \frac{\|m_{t-1}\|^2}{\mu^{t-1}} + \frac{1}{2} \lambda L \frac{\|m_{t-2}\|^2}{\mu^{t-2}} - \langle \nabla f(x_t), \eta_t g_t \rangle \mu^{-t}. \quad (13)$$

Note that $m_0 = \mathbf{0}$, and $B_1 = -\langle \nabla f(x_1), \eta_1 g_1 \rangle$. Therefore,

$$\begin{aligned} \frac{B_t}{\mu^t} &\leq \frac{B_1}{\mu} + (1 + \frac{3}{2}\lambda\mu)L \sum_{i=2}^t \|m_{i-1}\|^2 \mu^{-(i-1)} + \frac{1}{2}\lambda L \sum_{i=2}^t \|m_{i-2}\|^2 \mu^{-(i-2)} - \sum_{i=2}^t \langle \nabla f(x_i), \eta_i g_i \rangle \mu^{-i} \\ &\leq (1 + 2\lambda)L \sum_{i=1}^{t-1} \|m_i\|^2 \mu^{-i} - \sum_{i=1}^t \langle \nabla f(x_i), \eta_i g_i \rangle \mu^{-i}. \end{aligned} \quad (14)$$

Multiplying both sides of Eq. (14) by μ^t , we obtain the desired estimate Eq. (9). This finishes the proof. \square

The following two lemmas, which are first introduced in [34, Theorem 10], are particular results of the AdaGrad adaptive learning rate. Here we adjust their proof to the coordinate-wise version and represent it here for readers' convenience.

Lemma 4 (Proof of Theorem 10 in [34]). *Let $\sigma_t = \sqrt{\mathbb{E}_t g_t^2}$ and let $\hat{\eta}_t = \eta / \sqrt{\epsilon + \sum_{i=1}^{t-1} g_i^2 + \sigma_t^2}$. Assume that the noise gradients g_t satisfy the assumptions (A1) and (A2), then we have the following estimate*

$$-\mathbb{E}_t \langle \nabla f(x_t), \eta_t g_t \rangle \leq -\frac{1}{2} \|\nabla f(x_t)\|_{\hat{\eta}_t}^2 + \frac{2\sigma}{\eta} \mathbb{E}_t \|\eta_t g_t\|^2, \quad (15)$$

where $\|\nabla f(x_t)\|_{\hat{\eta}_t}^2 = \sum_{k=1}^d \hat{\eta}_{t,k} |\nabla_k f(x_t)|^2$ defined as Eq. (5), and the constant σ is the one in assumption (A2).

Proof. First, we have

$$-\langle \nabla f(x_t), \eta_t g_t \rangle = -\langle \nabla f(x_t), \hat{\eta}_t g_t \rangle + \langle \nabla f(x_t), (\hat{\eta}_t - \eta_t) g_t \rangle. \quad (16)$$

Note that $\hat{\eta}_t$ is independent of g_t , and $\mathbb{E}_t g_t = \nabla f(x_t)$ by assumption (A1). Hence,

$$\mathbb{E}_t \langle \nabla f(x_t), \hat{\eta}_t g_t \rangle = \langle \nabla f(x_t), \hat{\eta}_t \nabla f(x_t) \rangle = \|\nabla f(x_t)\|_{\hat{\eta}_t}^2. \quad (17)$$

Taking the conditional expectation of Eq. (16) with respect to ξ_t while ξ_1, \dots, ξ_{t-1} are fixed, we have

$$-\mathbb{E}_t \langle \nabla f(x_t), \eta_t g_t \rangle = -\|\nabla f(x_t)\|_{\hat{\eta}_t}^2 + \mathbb{E}_t \langle \nabla f(x_t), (\hat{\eta}_t - \eta_t) g_t \rangle. \quad (18)$$

To estimate the second term of Eq. (18), we first have

$$\langle \nabla f(x_t), (\hat{\eta}_t - \eta_t) g_t \rangle \leq \sum_{k=1}^d |\hat{\eta}_{t,k} - \eta_{t,k}| |\nabla_k f(x_t)| |g_{t,k}|. \quad (19)$$

Let $V_t = \epsilon + \sum_{i=1}^t g_i^2$, then $\eta_t = \eta / \sqrt{V_t}$ and $\hat{\eta}_t = \eta / \sqrt{V_{t-1} + \sigma_t^2}$. We then get

$$|\hat{\eta}_t - \eta_t| = \frac{\eta |g_t^2 - \sigma_t^2|}{\sqrt{V_{t-1} + \sigma_t^2} \sqrt{V_t} (\sqrt{V_{t-1} + \sigma_t^2} + \sqrt{V_t})} \leq \frac{\eta (|g_t| + \sigma_t)}{\sqrt{V_{t-1} + \sigma_t^2} \sqrt{V_t}} = \hat{\eta}_t \eta_t (|g_t| + \sigma_t) / \eta. \quad (20)$$

Note that the above inequality is coordinate-wise. By Eq. (19) and Eq. (20), we have

$$\langle \nabla f(x_t), (\hat{\eta}_t - \eta_t) g_t \rangle \leq \sum_{k=1}^d \frac{1}{\eta} \hat{\eta}_{t,k} \eta_{t,k} |\nabla_k f(x_t)| |g_{t,k}|^2 + \sum_{k=1}^d \frac{\sigma_{t,k}}{\eta} \hat{\eta}_{t,k} \eta_{t,k} |\nabla_k f(x_t)| |g_{t,k}| \quad (21)$$

We claim that

$$\mathbb{E}_t \sum_{k=1}^d \frac{1}{\eta} \hat{\eta}_{t,k} \eta_{t,k} |\nabla_k f(x_t)| |g_{t,k}|^2 \leq \frac{1}{4} \|\nabla f(x_t)\|_{\hat{\eta}_t}^2 + \frac{\sigma}{\eta} \mathbb{E}_t \|\eta_t g_t\|^2. \quad (22)$$

To see this, first, if $\sigma_{t,k} > 0$, we apply the arithmetic inequality $2ab \leq a^2 + b^2$ with

$$a = \frac{1}{2\sigma_{t,k}} \sqrt{\hat{\eta}_{t,k}} |\nabla_k f(x_t)| |g_{t,k}| \quad \text{and} \quad b = \frac{\sigma_{t,k}}{\eta} \sqrt{\hat{\eta}_{t,k} \eta_{t,k}} |g_{t,k}|$$

to the left hand side of Eq. (22), we have

$$\frac{1}{\eta} \hat{\eta}_{t,k} \eta_{t,k} |\nabla_k f(x_t)| |g_{t,k}|^2 \leq \frac{|g_{t,k}|^2}{4\sigma_{t,k}^2} \hat{\eta}_{t,k} |\nabla_k f(x_t)|^2 + \frac{\sigma_{t,k}^2}{\eta^2} \hat{\eta}_{t,k} |\eta_{t,k} g_{t,k}|^2 \quad (23)$$

Note that $\mathbb{E}_t |g_{t,k}|^2 = \sigma_{t,k}$, $\hat{\eta}_{t,k} = \eta / \sqrt{V_{t-1,k} + \sigma_{t,k}^2} \leq \eta / \sigma_{t,k}$, and $\sigma_{t,k} \leq \sqrt{\mathbb{E}_t \|g_t\|^2} \leq \sigma$ by assumption (A2). Therefore,

$$\begin{aligned} \mathbb{E}_t \left[\frac{1}{\eta} \hat{\eta}_{t,k} \eta_{t,k} |\nabla_k f(x_t)| |g_{t,k}|^2 \right] &\leq \frac{1}{4} \hat{\eta}_{t,k} |\nabla_k f(x_t)|^2 + \frac{1}{\eta} \mathbb{E}_t (\sigma_{t,k} |\eta_{t,k} g_{t,k}|^2) \\ &\leq \frac{1}{4} \hat{\eta}_{t,k} |\nabla_k f(x_t)|^2 + \frac{\sigma}{\eta} \mathbb{E}_t |\eta_{t,k} g_{t,k}|^2. \end{aligned} \quad (24)$$

On the other hand, if $\sigma_{t,k} = 0$, then $g_{t,k} = 0$, Eq. (24) holds automatically. By taking sum of the components for $k = 1, 2, \dots, d$, we then get the desired claim Eq. (22).

Similarly, we apply the arithmetic inequality with

$$a = \frac{1}{2} \sqrt{\hat{\eta}_{t,k}} |\nabla_k f(x_t)| \quad \text{and} \quad b = \frac{\sigma_{t,k}}{\eta} \sqrt{\hat{\eta}_{t,k} \eta_{t,k}} |g_{t,k}|$$

to the second term of Eq. (21), we have

$$\begin{aligned} \frac{\sigma_{t,k}}{\eta} \hat{\eta}_{t,k} \eta_{t,k} |\nabla_k f(x_t)| |g_{t,k}| &\leq \frac{1}{4} \hat{\eta}_{t,k} |\nabla_k f(x_t)|^2 + \frac{\sigma_{t,k}^2}{\eta^2} \hat{\eta}_{t,k} |\eta_{t,k} g_{t,k}|^2 \\ &\leq \frac{1}{4} \hat{\eta}_{t,k} |\nabla_k f(x_t)|^2 + \frac{\sigma}{\eta} |\eta_{t,k} g_{t,k}|^2. \end{aligned} \quad (25)$$

Hence,

$$\mathbb{E}_t \sum_{k=1}^d \frac{\sigma}{\eta} \hat{\eta}_{t,k} |\nabla_k f(x_t)| |g_{t,k}| \leq \frac{1}{4} \|\nabla f(x_t)\|_{\hat{\eta}_t} + \frac{\sigma}{\eta} \mathbb{E}_t \|\eta_t g_t\|^2 \quad (26)$$

Combining Eq. (21), Eq. (24) and Eq. (26), we obtain the following estimate

$$\mathbb{E}_t \langle \nabla f(x_t), (\hat{\eta}_t - \eta_t) g_t \rangle \leq \frac{1}{2} \|\nabla f(x_t)\|_{\hat{\eta}_t}^2 + \frac{2\sigma}{\eta} \mathbb{E}_t \|\eta_t g_t\|^2. \quad (27)$$

The proof is finished by taking the estimate Eq. (27) into Eq. (18). \square

Lemma 5 (Proof of Theorem 10 in [34]). *Let $\hat{\eta}_t$ be defined as in Lemma 4. Assume that the noise gradients g_t satisfy the assumptions (A1) and (A2). We have the following estimate*

$$\mathbb{E} \left[\min_{1 \leq t \leq T} \|f(x_t)\|^{4/3} \right]^{3/2} \leq \frac{\sqrt{\epsilon d + \sigma^2 T}}{\eta T} \mathbb{E} \sum_{t=1}^T \|\nabla f(x_t)\|_{\hat{\eta}_t}^2. \quad (28)$$

Proof. We cite the proof form [34]. Let $\hat{V}_t = \epsilon + \sum_{i=1}^{t-1} g_i^2 + \sigma_t^2$ where $\sigma_t = \sqrt{\mathbb{E}_t g_t^2}$, then $\hat{\eta}_t = \eta / \sqrt{\hat{V}_t}$. For any $0 < p, q < 1$ with $1/p + 1/q = 1$, by Hölder's inequality we have $\mathbb{E} |XY| \leq (\mathbb{E} |X|^p)^{1/p} (\mathbb{E} |Y|^q)^{1/q}$. Taking $p = 3/2$ and $q = 3$, and

$$X = \left(\frac{\|\nabla f(x_t)\|^2}{\sqrt{\hat{V}_t}} \right)^{2/3}, \quad Y = \left(\|\hat{V}_t\|_1 \right)^{1/3},$$

we have

$$\mathbb{E}\|\nabla f(x_t)\|^{4/3} \leq \left(\mathbb{E} \frac{\|\nabla f(x_t)\|^2}{\sqrt{\|\hat{V}_t\|_1}} \right)^{2/3} \left(\mathbb{E}\|\hat{V}_t\|_1 \right)^{1/3}.$$

Namely,

$$\left(\mathbb{E}\|\nabla f(x_t)\|^{4/3} \right)^{3/2} \leq \left(\mathbb{E} \frac{\|\nabla f(x_t)\|^2}{\sqrt{\|\hat{V}_t\|_1}} \right) \left(\mathbb{E}\|\hat{V}_t\|_1 \right)^{1/2}. \quad (29)$$

Note that

$$\frac{\|\nabla f(x_t)\|^2}{\sqrt{\|\hat{V}_t\|_1}} = \sum_{k=1}^d \frac{|\nabla_k f(x_t)|^2}{\sqrt{\|\hat{V}_t\|_1}} \leq \frac{1}{\eta} \sum_{k=1}^d \frac{\eta |\nabla_k f(x_t)|^2}{\sqrt{\hat{V}_{t,k}}} = \frac{1}{\eta} \sum_{k=1}^d \hat{\eta}_{t,k} |\nabla_k f(x_t)|^2 = \frac{1}{\eta} \|\nabla f(x_t)\|_{\hat{\eta}_t}^2. \quad (30)$$

On the other hand, for any $t \leq T$ we have

$$\mathbb{E}\|\hat{V}_t\|_1 = \epsilon d + \mathbb{E} \sum_{k=1}^d \left(\sum_{i=1}^{t-1} g_{i,k}^2 + \sigma_{t,k}^2 \right) = \epsilon d + \sum_{i=1}^t \mathbb{E}\|g_i\|^2 \leq \epsilon d + \sigma^2 t \leq \epsilon d + \sigma^2 T. \quad (31)$$

Hence, by Eq. (29), Eq. (30) and Eq. (31), we have

$$\left(\mathbb{E}\|\nabla f(x_t)\|^{4/3} \right)^{3/2} \leq \frac{\sqrt{\epsilon d + \sigma^2 T}}{\eta} \mathbb{E}\|\nabla f(x_t)\|_{\hat{\eta}_t}^2, \quad \forall t \leq T \quad (32)$$

It follows that

$$\mathbb{E} \left[\min_{1 \leq t \leq T} \|\nabla f(x_t)\|^{4/3} \right]^{3/2} \leq \frac{1}{T} \sum_{t=1}^T \left(\mathbb{E}\|\nabla f(x_t)\|^{4/3} \right)^{3/2} \leq \frac{\sqrt{\epsilon d + \sigma^2 T}}{\eta T} \mathbb{E} \sum_{t=1}^T \|\nabla f(x_t)\|_{\hat{\eta}_t}^2. \quad (33)$$

The proof is finished. \square

B Proof of Theorem 1

In this section we prove our main theorem 1. We restate the theorem for readers' convenience.

Theorem. Suppose the objective function f is L -smooth and $f^* = \min_{x \in \mathbb{R}^d} f(x) > -\infty$. Let $\{x_t\} \subseteq \mathbb{R}^d$ be a sequence generated by AdaUSM. Assume that the noise gradient g_t in each iteration satisfies the assumptions (A1) and (A2). Then for any $\delta > 0$, the convergence rate of AdaUSM holds with probability at least $1 - \delta^{2/3}$ as below:

$$\min_{1 \leq t \leq T} \|\nabla f(x_t)\|^2 \leq \frac{C(T)}{\delta} = \mathcal{O} \left(\frac{\log T}{\sqrt{T}} \right),$$

$$\text{where } C(T) = \frac{\sqrt{\epsilon d + \sigma^2 T}}{T} \left[\frac{2}{(1 + \lambda\mu)\eta} (f(x_1) - f^*) + \left(\frac{2\eta(1+2\lambda)^2 L}{(1 + \lambda\mu)(1 - \mu)^3} + \frac{4\sigma}{1 - \mu} \right) d \log \left(1 + \frac{\sigma^2 T}{\epsilon} \right) \right].$$

The key point of the Theorem is the estimate in the following proposition.

Proposition. Suppose the objective function f is L -smooth and $f^* = \min_{x \in \mathbb{R}^d} f(x) > -\infty$. Let $\{x_t\} \subseteq \mathbb{R}^d$ be a sequence generated by AdaUSM. Assume that the stochastic gradient g_t in each iteration satisfies the assumptions (A1) and (A2), then we have the following estimate

$$\left(\mathbb{E} \left[\min_{1 \leq t \leq T} \|\nabla f(x_t)\|^{4/3} \right] \right)^{3/2} \leq C(T) = \mathcal{O} \left(\frac{\log T}{\sqrt{T}} \right) \quad (34)$$

where

$$C(T) = \frac{\sqrt{\epsilon d + \sigma^2 T}}{T} \left[\frac{2}{(1 + \lambda\mu)\eta} (f(x_1) - f^*) + \left(\frac{2\eta(1 + 2\lambda)^2 L}{(1 + \lambda\mu)(1 - \mu)^3} + \frac{4\sigma}{1 - \mu} \right) d \log \left(1 + \frac{\sigma^2 T}{\epsilon} \right) \right].$$

The Theorem is an immediate result of the Proposition. Let us first prove the Theorem assuming the proposition.

Proof of the Theorem assuming the Proposition. Let $\xi = \min_{1 \leq t \leq T} \|\nabla f(x_t)\|^2$. By Proposition we have $\mathbb{E}|\xi|^{2/3} \leq C(T)^{2/3}$. Let \mathcal{P} denote the probability measure. By Chebyshev's inequality, we have

$$\mathcal{P}\left(|\xi|^{2/3} > \frac{C(T)^{2/3}}{\delta^{2/3}}\right) \leq \frac{\mathbb{E}|\xi|^{2/3}}{\frac{C(T)^{2/3}}{\delta^{2/3}}} \leq \delta^{2/3}. \quad (35)$$

Namely, $\mathcal{P}\left(|\xi| > \frac{C(T)}{\delta}\right) \leq \delta^{2/3}$. Therefore, we have $\mathcal{P}\left(|\xi| \leq \frac{C(T)}{\delta}\right) \geq 1 - \delta^{2/3}$. This finishes the proof. \square

Now we are ready to prove our main estimate in the Proposition. For clarity we first introduce the following lemmas which assemble the main estimate of the Proposition.

Lemma 6. *Assume the same setting as the Proposition, we have the following estimate*

$$\mathbb{E} \sum_{t=1}^T \|\nabla f(x_t)\|_{\hat{\eta}_t}^2 \leq \frac{2}{1+\lambda\mu}(f(x_1) - f^*) + \left(\frac{2(1+2\lambda)^2 L}{(1+\lambda\mu)(1-\mu)^3} + \frac{4\sigma}{\eta(1-\mu)} \right) \mathbb{E} \sum_{t=1}^T \|\eta_t g_t\|^2 \quad (36)$$

where $\hat{\eta}_t$ is defined as in Lemma 4.

The following lemma estimates the right hand side of Lemma 6.

Lemma 7. *Let $\eta_t = \sqrt{\epsilon + \sum_{i=1}^t g_i^2}$ be the AdaGrad adaptive learning rate. Assume the noise gradient g_t in each iteration satisfies the assumptions (A1) and (A2). We have the following estimate*

$$\mathbb{E} \sum_{t=1}^T \|\eta_t g_t\|^2 \leq \eta^2 d \log \left(1 + \frac{\sigma^2 T}{\epsilon} \right). \quad (37)$$

We first prove the Proposition assuming we already have the lemmas and leave the proof of the lemmas in the end.

Proof of the Proposition. The proof simply assembles Lemma 5, Lemma 6 and Lemma 7. By Lemma 5, we have

$$\left(\mathbb{E} \left[\min_{1 \leq t \leq T} \|\nabla f(x_t)\|^{4/3} \right] \right)^{3/2} \leq \frac{\sqrt{\epsilon d + \sigma^2 T}}{\eta T} \mathbb{E} \sum_{t=1}^T \|\nabla f(x_t)\|_{\hat{\eta}_t}^2 \quad (38)$$

On the other hand, by Lemma 6 and 7, we have the following estimate for the right hand side of Eq. (38)

$$\mathbb{E} \sum_{t=1}^T \|\nabla f(x_t)\|_{\hat{\eta}_t}^2 \leq \frac{2}{1+\lambda\mu}(f(x_1) - f^*) + \left(\frac{2(1+2\lambda)^2 L}{(1+\lambda\mu)(1-\mu)^3} + \frac{4\sigma}{\eta(1-\mu)} \right) \eta^2 d \log \left(1 + \frac{\sigma^2 T}{\epsilon} \right) \quad (39)$$

Combining Eq. (38) and Eq. (39) we get the estimate in the proposition. This finishes the proof. \square

Finally, we are left to prove Lemma 6 and Lemma 7.

Proof of Lemma 6. Since $x_{t+1} = x_t + m_t + \lambda\mu(m_t - m_{t-1})$, it follows by Lipschitz continuity of the gradient of f and descent lemma in [26] that

$$f(x_{t+1}) \leq f(x_t) + \langle \nabla f(x_t), m_t + \lambda\mu(m_t - m_{t-1}) \rangle + \frac{L}{2} \|m_t + \lambda\mu(m_t - m_{t-1})\|^2. \quad (40)$$

Since $m_t = \mu m_{t-1} - \eta_t g_t$, it follows

$$\begin{aligned} f(x_{t+1}) &\leq f(x_t) + \langle \nabla f(x_t), (1 + \lambda\mu - \lambda)m_t - \lambda\eta_t g_t \rangle + \frac{L}{2} \|m_t + \lambda\mu(m_t - m_{t-1})\|^2 \\ &\leq f(x_t) + (1 + \lambda\mu - \lambda) \langle \nabla f(x_t), m_t \rangle - \lambda \langle \nabla f(x_t), \eta_t g_t \rangle + \frac{L}{2} \|m_t + \lambda\mu(m_t - m_{t-1})\|^2. \end{aligned} \quad (41)$$

By Lemma 3, we have

$$\langle \nabla f(x_t), m_t \rangle \leq (1 + 2\lambda)L \sum_{i=1}^{t-1} \|m_i\|^2 \mu^{t-i} - \sum_{i=1}^t \langle \nabla f(x_i), \eta_i g_i \rangle \mu^{t-i}. \quad (42)$$

Note that $1 + \lambda\mu - \lambda \geq 0$ since $\lambda \leq 1/(1 - \mu)$. Combining Eq. (41) and Eq. (42), we have

$$\begin{aligned} f(x_{t+1}) &\leq f(x_t) + (1 + \lambda\mu - \lambda)(1 + 2\lambda)L \sum_{i=1}^{t-1} \|m_i\|^2 \mu^{t-i} + \frac{L}{2} \|m_t + \lambda\mu(m_t - m_{t-1})\|^2 \\ &\quad - (1 + \lambda\mu - \lambda) \sum_{i=1}^t \langle \nabla f(x_i), \eta_i g_i \rangle \mu^{t-i} - \lambda \langle \nabla f(x_t), \eta_t g_t \rangle \end{aligned} \quad (43)$$

On one hand, by arithmetic inequality, we have

$$\|m_t + \lambda\mu(m_t - m_{t-1})\|^2 \leq 2(1 + \lambda\mu)^2 \|m_t\|^2 + 2(\lambda\mu)^2 \|m_{t-1}\|^2. \quad (44)$$

Hence,

$$\begin{aligned} &(1 + \lambda\mu - \lambda)(1 + 2\lambda)L \sum_{i=1}^{t-1} \|m_i\|^2 \mu^{t-i} + \frac{L}{2} \|m_t + \lambda\mu(m_t - m_{t-1})\|^2 \\ &\leq (1 + \lambda\mu)(1 + 2\lambda)L \sum_{i=1}^{t-1} \|m_i\|^2 \mu^{t-i} - \lambda(1 + 2\lambda)\mu L \|m_{t-1}\|^2 + (1 + \lambda\mu)^2 L \|m_t\|^2 + (\lambda\mu)^2 L \|m_{t-1}\|^2 \\ &\leq (1 + 2\lambda)^2 L \sum_{i=1}^t \|m_i\|^2 \mu^{t-i}. \end{aligned} \quad (45)$$

Summarize Eq. (43) and Eq. (45), we have the following cleaner inequality

$$\begin{aligned} f(x_{t+1}) &\leq f(x_t) + (1 + 2\lambda)^2 L \sum_{i=1}^t \|m_i\|^2 \mu^{t-i} - (1 + \lambda\mu - \lambda) \sum_{i=1}^t \langle \nabla f(x_i), \eta_i g_i \rangle \mu^{t-i} \\ &\quad - \lambda \langle \nabla f(x_t), \eta_t g_t \rangle. \end{aligned} \quad (46)$$

On the other hand, by Lemma 4, we have

$$- \mathbb{E}_i \langle f(x_i), \eta_i g_i \rangle \leq -\frac{1}{2} \|\nabla f(x_i)\|_{\eta_i}^2 + \frac{2\sigma}{\eta} \mathbb{E}_i \|\eta_i g_i\|^2, \quad \forall i. \quad (47)$$

Combining Eq. (46) and Eq. (47), taking sum from 1 to T and taking expectation, followed by moving

the gradient square terms to the left hand side, we get

$$\begin{aligned}
& (1 + \lambda\mu - \lambda)\mathbb{E} \sum_{t=1}^T \sum_{i=1}^t \frac{1}{2} \|\nabla f(x_i)\|_{\hat{\eta}_i}^2 \mu^{t-i} + \lambda \mathbb{E} \sum_{t=1}^T \frac{1}{2} \|\nabla f(x_t)\|_{\hat{\eta}_t}^2 \\
& \leq f(x_1) - f^* + (1 + 2\lambda)^2 L \mathbb{E} \sum_{t=1}^T \sum_{i=1}^t \|m_i\|^2 \mu^{t-i} \\
& \quad + \frac{2\sigma}{\eta} \left[(1 + \lambda\mu - \lambda) \mathbb{E} \sum_{t=1}^T \sum_{i=1}^t \|\eta_i g_i\|^2 \mu^{t-i} + \lambda \mathbb{E} \sum_{t=1}^T \|\eta_t g_t\|^2 \right] \\
& \leq f(x_1) - f^* + (1 + 2\lambda)^2 L \mathbb{E} \sum_{t=1}^T \sum_{i=1}^t \|m_i\|^2 \mu^{t-i} + \frac{2\sigma(1 + \lambda\mu)}{\eta} \mathbb{E} \sum_{t=1}^T \sum_{i=1}^t \|\eta_i g_i\|^2 \mu^{t-i}.
\end{aligned} \tag{48}$$

The last inequality is due to $\lambda \sum_{t=1}^T \|\eta_t g_t\|^2 \leq \lambda \sum_{t=1}^T \sum_{i=1}^t \|\eta_i g_i\|^2 \mu^{t-i}$. Similarly, for the left hand side, note that $1 + \lambda\mu - \lambda \geq 0$ since $\lambda \leq 1/(1 - \mu)$, we have

$$\begin{aligned}
& (1 + \lambda\mu - \lambda) \sum_{t=1}^T \sum_{i=1}^t \|\nabla f(x_i)\|_{\hat{\eta}_i}^2 \mu^{t-i} + \lambda \sum_{t=1}^T \|\nabla f(x_t)\|_{\hat{\eta}_t}^2 \\
& \geq (1 + \lambda\mu - \lambda) \sum_{t=1}^T \|\nabla f(x_t)\|_{\hat{\eta}_t}^2 + \lambda \sum_{t=1}^T \|\nabla f(x_t)\|_{\hat{\eta}_t}^2 \\
& = (1 + \lambda\mu) \sum_{t=1}^T \|\nabla f(x_t)\|_{\hat{\eta}_t}^2.
\end{aligned} \tag{49}$$

Now we are left to estimate the third term and the last term in the right hand side of Eq. (48). We apply the double-sum trick

$$\sum_{t=1}^T \sum_{i=1}^t \|m_i\|^2 \mu^{t-i} = \sum_{i=1}^T \sum_{t=i}^T \|m_i\|^2 \mu^{t-i} \leq \frac{1}{1 - \mu} \sum_{i=1}^T \|m_i\|^2. \tag{50}$$

$$\sum_{t=1}^T \sum_{i=1}^t \|\eta_i g_i\|^2 \mu^{t-i} = \sum_{i=1}^T \sum_{t=i}^T \|\eta_i g_i\|^2 \mu^{t-i} \leq \frac{1}{1 - \mu} \sum_{i=1}^T \|\eta_i g_i\|^2. \tag{51}$$

Combining Eq. (48), Eq. (49), Eq. (50) and Eq. (51), we have

$$\frac{1 + \lambda\mu}{2} \mathbb{E} \sum_{t=1}^T \|\nabla f(x_t)\|_{\hat{\eta}_t}^2 \leq f(x_1) - f^* + \frac{(1 + 2\lambda)^2 L}{1 - \mu} \mathbb{E} \sum_{t=1}^T \|m_t\|^2 + \frac{2\sigma(1 + \lambda\mu)}{\eta(1 - \mu)} \mathbb{E} \sum_{t=1}^T \|\eta_t g_t\|^2. \tag{52}$$

Finally, By Lemma 2 we have

$$\mathbb{E} \sum_{t=1}^T \|m_t\|^2 \leq \frac{1}{(1 - \mu)^2} \mathbb{E} \sum_{t=1}^T \|\eta_t g_t\|^2. \tag{53}$$

Combining Eq. (52) and Eq. (53), we get

$$\mathbb{E} \sum_{t=1}^T \|\nabla f(x_t)\|_{\hat{\eta}_t}^2 \leq \frac{2}{1 + \lambda\mu} (f(x_1) - f^*) + \left(\frac{2(1 + 2\lambda)^2 L}{(1 + \lambda\mu)(1 - \mu)^3} + \frac{4\sigma}{\eta(1 - \mu)} \right) \mathbb{E} \sum_{t=1}^T \|\eta_t g_t\|^2. \tag{54}$$

The proof is finished. \square

Proof of Lemma 7. Let $V_t = \epsilon + \sum_{i=1}^t g_i^2$. Then $\eta_t = \eta/\sqrt{V_t}$. We have

$$\sum_{t=1}^T \|\eta_t g_t\|^2 = \eta^2 \sum_{t=1}^T \sum_{k=1}^d \frac{g_{t,k}^2}{V_{t,k}} = \eta^2 \sum_{k=1}^d \sum_{t=1}^T \frac{g_{t,k}^2}{V_{t,k}}. \quad (55)$$

Note that $V_{t,k} = \epsilon + \sum_{i=1}^t g_{i,k}^2$. By Lemma 1, we have $\sum_{t=1}^T \frac{g_{t,k}^2}{V_{t,k}} \leq \log(V_{T,k}) - \log(\epsilon)$. On the other hand, since $\log(x)$ is concave, we have $\mathbb{E}[\log(V_{T,k})] \leq \log(\mathbb{E}[V_{T,k}]) \leq \log(\epsilon + \sigma^2 T)$. Hence,

$$\begin{aligned} \mathbb{E} \sum_{t=1}^T \|\eta_t g_t\|^2 &\leq \eta^2 \mathbb{E} \sum_{k=1}^d (\log(V_{T,k}) - \log(\epsilon)) \\ &\leq \eta^2 \sum_{k=1}^d (\log(\mathbb{E}[V_{T,k}]) - \log(\epsilon)) \leq \eta^2 d \log\left(1 + \frac{\sigma^2 T}{\epsilon}\right). \end{aligned} \quad (56)$$

The proof is finished. \square