

The Number of Steps Needed for Nonconvex Optimization of a Deep Learning Optimizer is a Rational Function of Batch Size

HIDEAKI IIDUKA

Department of Computer Science, Meiji University, 1-1-1 Higashimita, Tama-ku, Kawasaki-shi, Kanagawa 214-8571 Japan. (iiduka@cs.meiji.ac.jp)

Abstract: Recently, convergence as well as convergence rate analyses of deep learning optimizers for nonconvex optimization have been widely studied. Meanwhile, numerical evaluations for the optimizers have precisely clarified the relationship between batch size and the number of steps needed for training deep neural networks. The main contribution of this paper is to show theoretically that the number of steps needed for nonconvex optimization of each of the optimizers can be expressed as a rational function of batch size. Having these rational functions leads to two particularly important facts, which were validated numerically in previous studies. The first fact is that there exists an optimal batch size such that the number of steps needed for nonconvex optimization is minimized. This implies that using larger batch sizes than the optimal batch size does not decrease the number of steps needed for nonconvex optimization. The second fact is that the optimal batch size depends on the optimizer. In particular, it is shown theoretically that momentum and Adam-type optimizers can exploit larger optimal batches and further reduce the minimum number of steps needed for nonconvex optimization than can the stochastic gradient descent optimizer.

1. INTRODUCTION

One way to train deep neural networks is to find the model parameters of the deep neural networks that minimize loss functions called the expected risk and empirical risk using first-order optimization methods [2, Section 4]. The simplest optimizer is stochastic gradient descent (SGD) [22, 33, 18, 8, 9]. There have been many deep learning optimizers to accelerate SGD, such as momentum methods [20, 19] and adaptive methods, e.g., Adaptive Gradient (AdaGrad) [5], Root Mean Square Propagation (RMSProp) [28], Adaptive Moment Estimation (Adam) [13], and Adaptive Mean Square Gradient (AMSGrad) [21] (Table 2 in [25] lists useful deep learning optimizers).

Convergence and convergence rate analyses of deep learning optimizers have been widely studied for convex optimization [34, 13, 21, 15, 17]. Meanwhile, theoretical investigation of deep learning optimizers for nonconvex optimization is needed so that these optimizers can put into practice for nonconvex optimization in deep learning [30, 1, 29].

Convergence analyses of SGD for nonconvex optimization were studied in [7, 3, 24, 14] (see [10, 14] for convergence analyses of SGD for two classes of nonconvex optimization problems, quasar-convex and Polyak–Lojasiewicz optimization problems). For example, Theorem 11 in [24] indicates that SGD with a diminishing learning rate $\alpha_k = 1/\sqrt{k}$ has $\mathcal{O}(1/\sqrt{K})$ convergence, where K denotes the number of steps. Convergence analyses of SGD depending on the batch size were presented in [3]. In particular, Theorem 3.2 in [3] indicates that running SGD with a diminishing learning rate $\alpha_k = 1/k$ and large batch size for sufficiently many steps leads to convergence to a local minimizer of a sum of loss functions.

Convergence analyses of adaptive methods for nonconvex optimization were studied in [6, 4, 32, 12]. In [4], it was shown that generalized Adam, which includes the Heavy-ball method, AdaGrad, RMSProp, AMSGrad, and AdaGrad with First Order Momentum (AdaFom), using a diminishing learning rate $\alpha_k = 1/\sqrt{k}$ has an $\mathcal{O}(\log K/\sqrt{K})$ convergence rate. AdaBelief (named for adapting stepsizes by the belief in observed gradients) using $\alpha_k = 1/\sqrt{k}$ has $\mathcal{O}(\log K/\sqrt{K})$ convergence [32]. In [12], a method was presented to unify useful adaptive methods such as AMSGrad and AdaBelief, and it was shown that the method with $\alpha_k = 1/\sqrt{k}$ has an $\mathcal{O}(1/\sqrt{K})$ convergence rate, which improves on the results in [4, 32]. A theoretical investigation of Stochastic Path-Integrated Differential Estimator (SPIDER) for ϵ -approximation in nonconvex optimization was reported in [6]. In particular, Theorem 2 in [6] clarified that SPIDER, which has a constant learning rate, for ϵ -approximation must use the full-batch gradient with the number of samples n or the stochastic gradient with batch size \sqrt{n} .

Meanwhile, in [26], it was studied how increasing the batch size affects the performances of SGD, SGD with momentum [20, 23], and Nesterov momentum [19, 27]. The relationships between batch size and performance for Adam and K-FAC (Kronecker-Factored Approximate Curvature [16]) were studied in [31]. In both studies, it was numerically shown that increasing batch size tends to decrease the number of steps K needed for training deep neural networks, but with diminishing returns [26, Figure 4], [31, Figure 8]. Moreover, it was shown that SGD with momentum and Nesterov momentum can exploit larger batches than SGD [26, Figure 4], and that K-FAC and Adam can exploit larger batches than SGD with momentum [31, Figure 5]. Thus, it was shown that momentum and adaptive methods can significantly reduce the number of steps K needed for training deep neural networks [26, Figure 4], [31, Figure 5].

1.1. Contribution. The contribution of this paper is to construct a theory guaranteeing the useful numerical results in [26, 31]. Table 1 (resp. Table 2) summarizes our results for SGD, Nesterov momentum (N-Momentum), and Adam-type optimizers with a constant learning rate rule (resp. diminishing learning rate rule), described in Theorem 3.1 (resp. Theorem 3.2). See Theorem A.2 in Appendix for other result for the optimizers with a diminishing learning rate rule. Figure 1 (resp. Figure 2) visualizes the relationships between the optimizers for the results shown in Table 1 (resp. Table 2) for an appropriately set momentum coefficient β .

The main contribution of this paper is to clarify that

- the number of steps $K = K_\epsilon$ needed for nonconvex optimization in the sense of ¹

$$(1) \quad \min_{k \in [K]} \mathbb{E} [\|\nabla f(\mathbf{x}_k)\|^2] \leq \epsilon^2,$$

where $\nabla f: \mathbb{R}^d \rightarrow \mathbb{R}^d$ denotes the gradient of a nonconvex loss function $f: \mathbb{R}^d \rightarrow \mathbb{R}$, $\epsilon > 0$ is a precision accuracy, and the sequence $(\mathbf{x}_k)_{k \in \mathbb{N}} \subset \mathbb{R}^d$ is generated by a particular optimizer, such as one of SGD, N-Momentum, and Adam-type optimizers, can be expressed as a rational function of batch size s (see the ‘‘Rational Function’’ columns of Tables 1 and 2).

The explicit forms of the rational functions imply the following two significant facts:

- (I) There exists an optimal batch size s^* such that $K_\epsilon(s)$ is minimized; specifically, $K_\epsilon(s)$ is monotone decreasing for $s \leq s^*$ and monotone increasing for $s \geq s^*$. This fact guarantees theoretically the existences of the diminishing returns shown in

¹Jensen’s inequality guarantees that (1) implies that $\min_{k \in [K]} \mathbb{E} [\|\nabla f(\mathbf{x}_k)\|] \leq \epsilon$.

TABLE 1. Relationship between batch size s and the number of steps K_ϵ needed for nonconvex optimization in the sense of (1) of optimizers with constant learning rates

	Constant Learning Rate Rule ($\alpha_k = \alpha \in (0, 1], \beta_k = \beta \in [0, b] \subset [0, 1)$)		
	Rational Function	Optimal Batch Size s^*	Minimum Steps $K_\epsilon(s^*)$
SGD	$K_\epsilon = \frac{A_\alpha s^2}{\epsilon^2 s - B_\alpha}$	$\frac{dDL^2 n^2 \alpha}{\epsilon^2}$	$\frac{(dDLn)^2}{\epsilon^4}$
N-Momentum	$K_\epsilon = \frac{A_\alpha s^2}{(\epsilon^2 - C_\beta)s - B_\alpha}$	$\frac{dDL^2 n^2 \alpha}{\tilde{b}\epsilon^2 - dDLn\beta}$	$\frac{(dDLn)^2}{(\tilde{b}\epsilon^2 - dDLn\beta)^2}$
Adam-type	$K_\epsilon = \frac{A_\alpha s^2}{(\epsilon^2 - C_\beta)s - B_\alpha}$	$\frac{dDL^2 n^2 \alpha}{\tilde{\gamma}^2(\tilde{b}\epsilon^2 - dDLn\beta)h_0^*}$	$\frac{(dDLn)^2 H}{\tilde{\gamma}^2(\tilde{b}\epsilon^2 - dDLn\beta)^2 h_0^*}$

Let $\epsilon > 0$, $\tilde{b} := 1 - b$, $\tilde{\gamma} := 1 - \gamma$ ($\gamma \in [0, 1)$), and $H \geq h_0^* > 0$. The number of samples is denoted by n , $\nabla f_i: \mathbb{R}^d \rightarrow \mathbb{R}^d$ ($i \in [n] := \{1, 2, \dots, n\}$) is Lipschitz continuous with Lipschitz constant L_i , and L denotes the maximum value of L_i . D is the upper bound of $(x_{k,i} - x_i)^2$ ($(x_i) \in \mathbb{R}^d$), where $(\mathbf{x}_k)_{k \in \mathbb{N}} = ((x_{k,i})_{i \in [n]})_{k \in \mathbb{N}}$ is generated by an optimizer. A_α and B_α are positive constants depending on a learning rate α and C_β is a positive constant depending on a momentum coefficient β (see Theorem 3.1 for detailed definitions of the constants).

TABLE 2. Relationship between batch size s and the number of steps K_ϵ needed for nonconvex optimization in the sense of (1) of optimizers with diminishing learning rates

	Diminishing Learning Rate Rule ($\alpha_k = \frac{\alpha}{\sqrt{k}}, \beta_k = \beta \in [0, b] \subset [0, 1)$)		
	Rational Function	Optimal Batch Size s^*	Minimum Steps $K_\epsilon(s^*)$
SGD	$K_\epsilon = \left\{ \frac{A_\alpha s^2 + B_\alpha}{\epsilon^2 s} \right\}^2$	$\sqrt{2}Ln\alpha$	$\frac{2(dDLn)^2}{\epsilon^4}$
N-Momentum	$K_\epsilon = \left\{ \frac{A_\alpha s^2 + B_\alpha}{(\epsilon^2 - C_\beta)s} \right\}^2$	$\sqrt{2}Ln\alpha$	$\frac{2(dDLn)^2}{(\tilde{b}\epsilon^2 - dDLn\beta)^2}$
Adam-type	$K_\epsilon = \left\{ \frac{A_\alpha s^2 + B_\alpha}{(\epsilon^2 - C_\beta)s} \right\}^2$	$\frac{\sqrt{2}Ln\alpha}{\tilde{\gamma}\sqrt{Hh_0^*}}$	$\frac{2(dDLn)^2 H}{\tilde{\gamma}^2(\tilde{b}\epsilon^2 - dDLn\beta)^2 h_0^*}$

See Table 1 and Theorem 3.2 for definitions of the constants.

[26, Figure 4], [31, Figure 8], which are such that increasing the batch size does not decrease the number of steps K_ϵ .

- (II) The optimal batch size s^* and the minimum number of steps $K_\epsilon(s^*)$ depend on the optimizer. In particular, N-Momentum and Adam-type optimizers can exploit the same sized or larger batches (s^* in Tables 1 and 2 and Figures 1 and 2) than can SGD. Furthermore, the dependence of N-Momentum and Adam-type optimizers on β allows them to reduce the minimum number of steps ($K_\epsilon(s^*)$ in Tables 1 and 2 and Figures 1 and 2) more than can SGD (see Section 3 for details).

Comparisons of Optimal Batch Sizes for Different Learning Rate Rules. Tables 1 and 2 ensure that $K_\epsilon(s^*)$ for Algorithm 1 using constant learning rates is almost the

FIGURE

1. Relationships between the optimizers in terms of the results in Table 1 (relations $s_{C,SGD}^* \leq s_{C,NM}^* \leq s_{C,A}^*$ hold generally, but those of $K_\epsilon(s_{C,SGD}^*) \geq K_\epsilon(s_{C,NM}^*) \geq K_\epsilon(s_{C,A}^*)$ depend on momentum coefficient β)

FIGURE

2. Relationships between the optimizers in terms of the results in Table 2 (relations $s_{D,SGD}^* = s_{D,NM}^* \leq s_{D,A}^*$ hold generally, but those of $K_\epsilon(s_{D,SGD}^*) \geq K_\epsilon(s_{D,NM}^*) \geq K_\epsilon(s_{D,A}^*)$ depend on momentum coefficient β)

same as $K_\epsilon(s^*)$ for Algorithm 1 using diminishing learning rates. Meanwhile, we would like to emphasize that the optimal batch size s_C^* for Algorithm 1 using constant learning rates depends on ϵ and β , and the optimal batch size s_D^* for Algorithm 1 using diminishing learning rates does not depend on ϵ and β . For example, under the precision accuracy $\epsilon = 10^{-1}$, we can know the optimal batch sizes for N-Momentum with the frequently used parameter value $\alpha = 10^{-3}$ are respectively

$$s_{C,NM}^* = \frac{dDL^2n^2\epsilon^3}{b\epsilon^2 - dDLn\beta} \text{ and } s_{D,NM}^* = \sqrt{2}Ln\epsilon^3$$

before implementing N-Momentum.

1.2. Notation. \mathbb{N} denotes the set of nonnegative integers. Let $n \in \mathbb{N} \setminus \{0\}$. We define $[n] := \{1, 2, \dots, n\}$. \mathbb{R}^d denotes d -dimensional Euclidean space with inner product $\langle \cdot, \cdot \rangle$ inducing norm $\|\cdot\|$. Let \mathbb{S}_{++}^d be the set of $d \times d$ symmetric positive-definite matrices and let \mathbb{D}^d be the set of $d \times d$ diagonal matrices: $\mathbb{D}^d = \{M \in \mathbb{R}^{d \times d}: M = \text{diag}(x_i), x_i \in \mathbb{R} (i \in [d])\}$. For a random variable Z , we use $\mathbb{E}[Z]$ to indicate its expectation.

2. NONCONVEX OPTIMIZATION AND DEEP LEARNING OPTIMIZERS

2.1. Assumptions Regarding Loss Function and Gradient Estimation. This paper considers optimization problems under the following assumptions.

Assumption 2.1.

(A1) [Loss function] $f_i: \mathbb{R}^d \rightarrow \mathbb{R}$ ($i \in [n]$) is differentiable and $f: \mathbb{R}^d \rightarrow \mathbb{R}$ is defined for all $\mathbf{x} \in \mathbb{R}^d$ by

$$f(\mathbf{x}) := \frac{1}{n} \sum_{i=1}^n f_i(\mathbf{x}),$$

where n denotes the number of samples.

(A2) [Gradient estimation] For each iteration k , optimizers sample a batch $\mathcal{S}_k \subset [n]$ of size $s := |\mathcal{S}_k|$ independently of k and estimate the full gradient ∇f as

$$\nabla f_{\mathcal{S}_k} := \frac{1}{s} \sum_{i \in \mathcal{S}_k} \nabla f_i.$$

(A3) [Gradient boundedness] There exists a positive number G such that, for all $\mathbf{x} \in X$,

$$(2) \quad \mathbb{E} [\|\nabla f_{\mathcal{S}_k}(\mathbf{x})\|^2] \leq \frac{G^2}{s^2},$$

where X is a subset of \mathbb{R}^d .

Assumption (A1) is a standard one for nonconvex optimization in deep neural networks (see, e.g., [4, (2)] and [6, (1.2)]). Assumption (A2) is needed for the optimizers to work (see, e.g., [4, Section 2] and [6, Notation section]). Assumption (A3) is used to analyze the optimizers. Assumption (A3) holds if each of the following holds (see Proposition A.1 in Appendix for details):

- (G1) $X \subset \mathbb{R}^d$ is bounded, the gradient ∇f_i is Lipschitz continuous with Lipschitz constant L_i , and $S_i := \{\mathbf{x}^* \in \mathbb{R}^d: \nabla f_i(\mathbf{x}^*) = \mathbf{0}\} \neq \emptyset$ ($i \in [n]$), where $L := \max_{i \in [n]} L_i$. (If we define $G_{k,L} := \sup_{\mathbf{x} \in X} \sum_{i \in \mathcal{S}_k} \|\nabla f_i(\mathbf{x})\|$, then we can take $G := \sup_{k \in \mathbb{N}} G_{k,L}$.)
- (G2) $X \subset \mathbb{R}^d$ is bounded and closed. (If we define $G_k := \sup_{\mathbf{x} \in X} \sum_{i \in \mathcal{S}_k} \|\nabla f_i(\mathbf{x})\|$, then we can take $G := \sup_{k \in \mathbb{N}} G_k$.)

2.2. Nonconvex Optimization in Deep Learning. This paper considers the following problem [4, 32].

Problem 2.1. *Under Assumption 2.1, we would like to find a local minimizer \mathbf{x}^* of f over \mathbb{R}^d , i.e.,*

$$\mathbf{x}^* \in X^* := \{\mathbf{x} \in \mathbb{R}^d: \nabla f(\mathbf{x}) = \mathbf{0}\}.$$

If f is convex [13, 21], then the solution to Problem 2.1 is a global minimizer of f over \mathbb{R}^d . See the third and fourth paragraphs of Section 1 for the previous studies on Problem 2.1.

2.3. Deep Learning Optimizers. There are many deep learning optimizers [25, Table 2]. In this paper, we consider the following algorithm (Algorithm 1), which is a unified algorithm for useful optimizers, for example, N-Momentum [19, 27], AMSGrad [21, 4], AMSBound [15], and AdaBelief [32], listed in Table 3 in Appendix.

The useful optimizers, such as N-Momentum, AMSGrad, AMSBound, and AdaBelief (Table 3), all satisfy the following conditions:

Assumption 2.2. The sequence $(\mathbf{H}_k)_{k \in \mathbb{N}} \subset \mathbb{S}_{++}^d \cap \mathbb{D}^d$, with $\mathbf{H}_k := \text{diag}(h_{k,i})$, in Algorithm 1 satisfies the following conditions:

- (A4) $h_{k+1,i} \geq h_{k,i}$ almost surely for all $k \in \mathbb{N}$ and all $i \in [d]$;
- (A5) For all $i \in [d]$, a positive number H_i exists such that $\sup_{k \in \mathbb{N}} \mathbb{E}[h_{k,i}] \leq H_i$.

Algorithm 1 Deep learning optimizer for solving Problem 2.1

Require: $(\alpha_k)_{k \in \mathbb{N}} \subset (0, 1]$, $(\beta_k)_{k \in \mathbb{N}} \subset [0, b] \subset [0, 1]$, $\gamma \in [0, 1]$

- 1: $k \leftarrow 0$, $\mathbf{x}_0, \mathbf{m}_{-1} := \mathbf{0} \in \mathbb{R}^d$, $\mathbf{H}_0 \in \mathbb{S}_{++}^d \cap \mathbb{D}^d$, $\mathcal{S}_0 \subset [n]$
 - 2: **loop**
 - 3: $\mathbf{m}_k := \beta_k \mathbf{m}_{k-1} + (1 - \beta_k) \nabla f_{\mathcal{S}_k}(\mathbf{x}_k)$
 - 4: $\hat{\mathbf{m}}_k := \frac{\mathbf{m}_k}{1 - \gamma^{k+1}}$
 - 5: $\mathbf{H}_k \in \mathbb{S}_{++}^d \cap \mathbb{D}^d$ (see Table 3 for examples of \mathbf{H}_k)
 - 6: Find $\mathbf{d}_k \in \mathbb{R}^d$ that solves $\mathbf{H}_k \mathbf{d} = -\hat{\mathbf{m}}_k$
 - 7: $\mathbf{x}_{k+1} := \mathbf{x}_k + \alpha_k \mathbf{d}_k$
 - 8: $k \leftarrow k + 1$
 - 9: **end loop**
-

Moreover, the following condition holds:

- (A6) $D := \max_{i \in [d]} \sup_{k \in \mathbb{N}} (x_{k+1,i} - x_i)^2 < +\infty$, where $\mathbf{x} := (x_i) \in \mathbb{R}^d$ and $(\mathbf{x}_k)_{k \in \mathbb{N}} := ((x_{k,i}))_{k \in \mathbb{N}}$ is the sequence generated by Algorithm 1.

The previous results in [4, p.29], [32, p.18], and [12] show that $(\mathbf{H}_k)_{k \in \mathbb{N}}$ in Table 3 satisfies (A4) and (A5). Assumption (A6) is assumed in [18, p.1574], [13, Theorem 4.1], [21, p.2], and [32, Theorem 2.1]. If (A6) holds, then there exists a bounded set $X \subset \mathbb{R}^d$ such that $(\mathbf{x}_k)_{k \in \mathbb{N}} \subset X$. Accordingly, the Lipschitz continuity of ∇f_i , the nonemptiness of S_i , and (A6) imply that (G1) with $G := Ln\sqrt{dD}$ holds (see Proposition A.1 in Appendix for details). We define

$$h_0^* := \min_{i \in [d]} h_{0,i} \text{ and } H := \max_{i \in [d]} H_i,$$

where $h_{k,i}$ and H_i are defined as in Assumption 2.2.

3. MAIN RESULTS

3.1. Constant Learning Rate Rule. The following theorem gives the relationship between batch size s and the number of steps K_ϵ needed for (1) for Algorithm 1 with a constant learning rate $\alpha_k = \alpha$ (see Table 1 for the specific results in Theorem 3.1 with $G := Ln\sqrt{dD}$ (i.e., under condition (G1))).

Theorem 3.1. *Suppose that Assumptions 2.1 and 2.2 hold and let $s, \epsilon > 0$.*

- (i) *Consider Algorithm 1 with*

$$\alpha_k := \alpha \in (0, 1] \text{ and } \beta_k := \beta \in [0, b] \subset [0, 1].$$

Then, for all $K \geq 1$,

$$\min_{k \in [K]} \mathbb{E} [\|\nabla f(\mathbf{x}_k)\|^2] \leq \underbrace{\frac{dDH}{2(1-b)\alpha}}_{A_\alpha} \frac{s}{K} + \underbrace{\frac{G^2\alpha}{2(1-b)(1-\gamma)^2 h_0^*}}_{B_\alpha} \frac{1}{s} + \underbrace{\frac{\sqrt{dDG}}{1-b}}_{C_\beta} \beta.$$

- (ii) *Consider Algorithm 1 with*

$$\alpha_k := \alpha \in (0, 1] \text{ and } \beta_k := \beta < \min \left\{ \frac{1-b}{\sqrt{dDG}} \epsilon^2, b \right\}.$$

Then, the number of steps K_ϵ needed to achieve (1) is expressed as the following rational function of batch size s :

$$(3) \quad K_\epsilon(s) = \frac{A_\alpha s^2}{(\epsilon^2 - C_\beta)s - B_\alpha} \quad \left(s \in \left(\frac{B_\alpha}{\epsilon^2 - C_\beta}, +\infty \right) \right).$$

In particular, the minimum value of K_ϵ needed to achieve (1) is

$$K_\epsilon(s^*) = \frac{4A_\alpha B_\alpha}{(\epsilon^2 - C_\beta)^2} = \frac{dDG^2H}{(1-\gamma)^2\{(1-b)\epsilon^2 - \sqrt{dDG\beta}\}^2 h_0^*}$$

when

$$s^* = \frac{2B_\alpha}{\epsilon^2 - C_\beta} = \frac{G^2\alpha}{(1-\gamma)^2\{(1-b)\epsilon^2 - \sqrt{dDG\beta}\}h_0^*}.$$

3.1.1. *Discussion of Theorem 3.1.* Let us examine the results in Theorem 3.1 for SGD, N-Momentum, and Adam-type optimizers.

[Performance of Algorithm 1] SGD is Algorithm 1 with $\beta = b = \gamma = 0$ and $h_0^* = H = 1$, N-Momentum is Algorithm 1 with $\gamma = 0$ and $h_0^* = H = 1$, and the Adam-type optimizer is Algorithm 1 with $\gamma \in [0, 1)$ and $h_{k,i}$ defined by one of $\sqrt{\hat{v}_{k,i}}$, $\sqrt{\bar{v}_{k,i}}$, and $\sqrt{\hat{s}_{k,i}}$ (see Table 3). Theorem 3.1(i) indicates that, for all $K \geq 1$, all $\alpha \in (0, 1]$, all $\beta \in [0, b] \subset [0, 1)$, and all $s > 0$,

$$(4) \quad \min_{k \in [K]} \mathbb{E} [\|\nabla f(\mathbf{x}_k)\|^2] \leq \begin{cases} \frac{dD_{\text{SGD}}}{2\alpha} \frac{s}{K} + \frac{G^2\alpha}{2} \frac{1}{s} & \text{(SGD),} \\ \frac{dD_{\text{NM}}}{2(1-b)\alpha} \frac{s}{K} + \frac{G^2\alpha}{2(1-b)} \frac{1}{s} + \frac{\sqrt{dD_{\text{NM}}G}}{1-b} \beta & \text{(N-Momentum),} \\ \frac{dD_{\text{A}}H}{2(1-b)\alpha} \frac{s}{K} + \frac{G^2\alpha}{2(1-b)(1-\gamma)^2 h_0^*} \frac{1}{s} + \frac{\sqrt{dD_{\text{A}}G}}{1-b} \beta & \text{(Adam-type).} \end{cases}$$

Note that D depends on the optimizer, which we distinguish by the notation D_{SGD} , D_{NM} , and D_{A} . For fixed s , if α and β are sufficiently small, (4) indicates that SGD, N-Momentum, and Adam-type optimizers have approximately $\mathcal{O}(1/K)$ convergence. For fixed s and K , if α is sufficiently small, the second term on the right-hand side of (4) will be small, whereas the first term will be large. Hence, there is no evidence that Algorithm 1 with a sufficiently small learning rate α would perform arbitrarily well. For fixed α and K , if s is sufficiently large, again the second term of the right-hand side of (4) will be small and the first term will be large. Hence, (4) indicates that there is no evidence that Algorithm 1 with a large batch size s performs better than with a smaller batch size.

[Existence of optimal batch size] The function $K_\epsilon(s)$ defined by (3) satisfies the following:

$$\frac{dK_\epsilon(s)}{ds} \begin{cases} < 0 & \text{if } s \in \left(\frac{B_\alpha}{\epsilon^2 - C_\beta}, s^* \right), \\ = 0 & \text{if } s = s^* = \frac{2B_\alpha}{\epsilon^2 - C_\beta}, \\ > 0 & \text{if } s \in (s^*, +\infty). \end{cases}$$

The above shows that increasing the batch size initially decreases the number of steps K_ϵ needed to achieve (1). Then, there is an optimal batch size ($s = s^*$) minimizing $K_\epsilon(s)$; thus increasing the batch size does not always decrease the number of steps K_ϵ .

[Comparison of optimal batch sizes] We assume that SGD, N-Momentum, and Adam-type optimizers all use the same G . For example, under (G1), we have $G = Ln\sqrt{dD}$, where $D = \max\{D_{\text{SGD}}, D_{\text{NM}}, D_{\text{A}}\}$. From Theorem 3.1(ii), we find that

$$(5) \quad s_{\text{C,SGD}}^* = \frac{G^2\alpha}{\epsilon^2} \leq s_{\text{C,NM}}^* = \frac{G^2\alpha}{(1-b)\epsilon^2 - \sqrt{dD_{\text{NM}}G\beta}}.$$

This implies that N-Momentum exploits larger batches than SGD. Moreover, if²

$$(6) \quad (1 - \gamma)^2 \leq \frac{1}{h_0^*},$$

then we have that

$$(7) \quad s_{\mathcal{C},\text{SGD}}^* = \frac{G^2 \alpha}{\epsilon^2} \leq s_{\mathcal{C},\text{A}}^* = \frac{G^2 \alpha}{(1 - \gamma)^2 \{(1 - b)\epsilon^2 - \sqrt{dD_{\text{A}}G\beta}\} h_0^*}.$$

Therefore, N-Momentum and Adam-type optimizers exploit larger batches than SGD. Moreover, if (6) holds and if³

$$(8) \quad D_{\text{NM}} \leq D_{\text{A}},$$

then

$$s_{\mathcal{C},\text{NM}}^* \leq s_{\mathcal{C},\text{A}}^*.$$

[Comparison of minimum numbers of steps] Theorem 3.1(ii) guarantees that, if β satisfies the condition in Theorem 3.1(ii) and if

$$(9) \quad \beta \leq \frac{(1 - b)\sqrt{D_{\text{SGD}}} - \sqrt{D_{\text{NM}}}}{\sqrt{dD_{\text{SGD}}D_{\text{NM}}G}} \epsilon^2,$$

then

$$(10) \quad K_{\epsilon}(s_{\mathcal{C},\text{SGD}}^*) = \frac{dD_{\text{SGD}}G^2}{\epsilon^4} \geq K_{\epsilon}(s_{\mathcal{C},\text{NM}}^*) = \frac{dD_{\text{NM}}G^2}{\{(1 - b)\epsilon^2 - \sqrt{dD_{\text{NM}}G\beta}\}^2}.$$

Moreover, if β satisfies the condition in Theorem 3.1(ii) and if

$$(11) \quad \beta \leq \frac{(1 - b)(1 - \gamma)\sqrt{D_{\text{SGD}}h_0^*} - \sqrt{D_{\text{A}}H}}{(1 - \gamma)\sqrt{dD_{\text{SGD}}D_{\text{A}}h_0^*G}} \epsilon^2,$$

then

$$(12) \quad K_{\epsilon}(s_{\mathcal{C},\text{SGD}}^*) \geq K_{\epsilon}(s_{\mathcal{C},\text{A}}^*) = \frac{dD_{\text{A}}G^2H}{(1 - \gamma)^2 \{(1 - b)\epsilon^2 - \sqrt{dD_{\text{A}}G\beta}\}^2 h_0^*}.$$

Additionally, if β satisfies the condition in Theorem 3.1(ii) and if

$$(13) \quad \beta \leq \frac{(1 - b)\{(1 - \gamma)\sqrt{D_{\text{NM}}h_0^*} - \sqrt{D_{\text{A}}H}\}}{\{(1 - \gamma)\sqrt{dD_{\text{NM}}D_{\text{A}}h_0^*} - \sqrt{dD_{\text{NM}}D_{\text{A}}H}\}G} \epsilon^2,$$

then

$$(14) \quad K_{\epsilon}(s_{\mathcal{C},\text{NM}}^*) \geq K_{\epsilon}(s_{\mathcal{C},\text{A}}^*).$$

See (17) for more specific β satisfying (13).

² γ and h_0^* can be chosen before implementing optimizers. For example, let $\gamma = 0.9$, which is used in [13]. Then, for all $i \in [d]$, we can set $h_{0,i} \leq 100$ (e.g., $h_0^* = h_{0,i} = 1$) in order to satisfy (6).

³We may assume that $D_{\text{NM}} = D_{\text{A}}$ in place of (8).

3.2. Diminishing Learning Rate Rule. The following theorem gives the relationships between batch size s and the number of steps K_ϵ needed for (1) for Algorithm 1 with a diminishing learning rate $\alpha_k := \alpha/\sqrt{k}$ (see Table 2 for the specific results in Theorem 3.2 with $G := Ln\sqrt{dD}$ (i.e., under condition (G1)) and Theorem A.2 for other results of Algorithm 1 with diminishing learning rates).

Theorem 3.2. *Suppose that Assumptions 2.1 and 2.2 hold and also $s, \epsilon > 0$ and $\alpha \in (0, 1]$.*

(i) *Consider Algorithm 1 with*

$$\alpha_k := \frac{\alpha}{\sqrt{k}} \text{ and } \beta_k := \beta \in [0, b] \subset [0, 1).$$

Then, for all $K \geq 1$,

$$\min_{k \in [K]} \mathbb{E} [\|\nabla f(\mathbf{x}_k)\|^2] \leq \underbrace{\frac{dDH}{2(1-b)\alpha}}_{A_\alpha} \frac{s}{\sqrt{K}} + \underbrace{\frac{G^2\alpha}{(1-b)(1-\gamma)^2 h_0^*}}_{B_\alpha} \frac{1}{s\sqrt{K}} + \underbrace{\frac{\sqrt{dDG}}{1-b}}_{C_\beta} \beta.$$

(ii) *Consider Algorithm 1 with*

$$\alpha_k := \frac{\alpha}{\sqrt{k}} \text{ and } \beta_k := \beta < \min \left\{ \frac{1-b}{\sqrt{dDG}} \epsilon^2, b \right\}.$$

Then, the number of steps K_ϵ needed to achieve (1) is expressed as the following rational function of batch size s :

$$(15) \quad K_\epsilon(s) = \left\{ \frac{A_\alpha s^2 + B_\alpha}{(\epsilon^2 - C_\beta)s} \right\}^2.$$

In particular, the minimum value of K_ϵ needed to achieve (1) is

$$K_\epsilon(s^*) = \frac{4A_\alpha B_\alpha}{(\epsilon^2 - C_\beta)^2} = \frac{2dDG^2H}{(1-\gamma)^2 \{(1-b)\epsilon^2 - \sqrt{dDG}\beta\}^2 h_0^*}$$

when

$$s^* = \sqrt{\frac{B_\alpha}{A_\alpha}} = \frac{\sqrt{2G\alpha}}{(1-\gamma)\sqrt{dDH}h_0^*}.$$

3.2.1. Discussion of Theorem 3.2. Let us discuss the results in Theorem 3.2 and compare them with those in Theorem 3.1 for SGD, N-Momentum, and Adam-type optimizers.

[Performance of Algorithm 1] Theorem 3.2(i) indicates that Algorithm 1 satisfies that, for all $K \geq 1$, all $\alpha \in (0, 1]$, all $\beta \in [0, b]$, and all $s > 0$,

$$\min_{k \in [K]} \mathbb{E} [\|\nabla f(\mathbf{x}_k)\|^2] \leq \begin{cases} \frac{dD_{\text{SGD}}}{2\alpha} \frac{s}{\sqrt{K}} + \frac{G^2\alpha}{s\sqrt{K}} & \text{(SGD),} \\ \frac{dD_{\text{NM}}}{2(1-b)\alpha} \frac{s}{\sqrt{K}} + \frac{G^2\alpha}{(1-b)s\sqrt{K}} + \frac{\sqrt{dD_{\text{NM}}G}}{1-b} \beta & \text{(N-Momentum),} \\ \frac{dD_{\text{A}}H}{2(1-b)\alpha} \frac{s}{\sqrt{K}} + \frac{G^2\alpha}{(1-b)(1-\gamma)^2 h_0^*} \frac{1}{s\sqrt{K}} + \frac{\sqrt{dD_{\text{A}}G}}{1-b} \beta & \text{(Adam-type).} \end{cases}$$

By a similar argument to that in Section 3.1.1, SGD, N-Momentum, and Adam-type optimizers have approximately $\mathcal{O}(1/\sqrt{K})$ convergence (see also Theorem A.2, which indicates that Algorithm 1 with $\alpha_k = \alpha/\sqrt{k}$ and $\beta_k = \beta^k$ has an only $\mathcal{O}(1/\sqrt{K})$ convergence rate) and that there is no evidence that Algorithm 1 with a large batch size s performs better than with a smaller batch size.

[Existence of optimal batch size] K_ϵ defined by (15) guarantees that there exists s^* such that $dK_\epsilon(s^*)/ds = 0$, the same as seen in Section 3.1.1 for Theorem 3.1. This implies that there is an optimal batch size ($s = s^*$) such that $K_\epsilon(s)$ is minimized, i.e., that increasing the batch size does not always decrease the number of steps K_ϵ .

[Comparison of optimal batch sizes] For simplicity, let us consider the case where (G1) holds. Theorem 3.2(ii) with $G = Ln\sqrt{dD}$ ensures that the optimal batch sizes for SGD, N-Momentum, and Adam-type optimizers with $\alpha_k = \alpha/\sqrt{k}$ and $\beta_k = \beta$ satisfy that

$$s_{\text{D,SGD}}^* = \frac{\sqrt{2}G\alpha}{\sqrt{dD_{\text{SGD}}}} = \sqrt{2}Ln\alpha = \frac{\sqrt{2}G\alpha}{\sqrt{dD_{\text{NM}}}} = s_{\text{D,NM}}^*.$$

Furthermore, if⁴

$$(16) \quad (1 - \gamma)^2 \leq \frac{1}{Hh_0^*},$$

then

$$s_{\text{D,SGD}}^* = s_{\text{D,NM}}^* \leq s_{\text{D,A}}^* = \frac{\sqrt{2}Ln\alpha}{(1 - \gamma)\sqrt{Hh_0^*}}.$$

Therefore, N-Momentum and Adam-type optimizers exploit the same sized or larger batches than SGD. Here, we notice that $s_{\text{C,SGD}}^*$, $s_{\text{C,NM}}^*$, and $s_{\text{C,A}}^*$ defined as in (5) and (7) depend on ϵ and β , while $s_{\text{D,SGD}}^*$, $s_{\text{D,NM}}^*$, and $s_{\text{D,A}}^*$ do not depend on ϵ and β .

[Comparison of minimum numbers of steps] Again, by a similar argument to that in Section 3.1.1, the restrictions on β (9), (11), and (13) imply that (10), (12), and (14) hold, respectively, i.e., that

$$K_\epsilon(s_{\text{D,A}}^*) \leq K_\epsilon(s_{\text{D,NM}}^*) \leq K_\epsilon(s_{\text{D,SGD}}^*).$$

The previous studies [13, 21, 15] used $\beta = 0.9$ or 0.99 , which is close to 1, for adaptive methods. Meanwhile, a sufficient condition for $K_\epsilon(s_{\text{D,A}}^*) \leq K_\epsilon(s_{\text{D,NM}}^*)$ is (13) with $D = D_{\text{NM}} = D_{\text{A}}$ and $G = Ln\sqrt{dD}$, i.e.,

$$(17) \quad \beta \leq \frac{(1 - b)\{(1 - \gamma)\sqrt{D_{\text{NM}}h_0^*} - \sqrt{D_{\text{A}}H}\}}{\{(1 - \gamma)\sqrt{dD_{\text{NM}}D_{\text{A}}h_0^*} - \sqrt{dD_{\text{NM}}D_{\text{A}}H}\}G} \epsilon^2 = \frac{(1 - b)\epsilon^2}{\sqrt{dDG}} = \frac{(1 - b)\epsilon^2}{LndD},$$

which implies that adaptive methods using the above β (which is small when the number of samples n and the number of dimension d are both large and the precision accuracy ϵ is small) are good for training deep neural networks in the sense that $K_\epsilon(s_{\text{D,A}}^*) \leq K_\epsilon(s_{\text{D,NM}}^*)$.

4. CONCLUSION AND FUTURE WORK

The main contribution of this paper was to show that the number of steps $K_\epsilon(s)$ needed for nonconvex optimization, $\min_{k \in [K]} \mathbb{E}[\|\nabla f(\mathbf{x}_k)\|^2] \leq \epsilon^2$, of a deep learning optimizer is a rational function of batch size. We showed that there exists an optimal batch size s^* such that $K_\epsilon(s)$ is minimized. This means that the optimizer using the optimal batch size s^* converges to a local minimizer of the sum of loss functions in at most $K_\epsilon(s^*)$ steps and is most desirable for training deep neural networks. Hence, there is no guarantee that the optimizer with a sufficiently large batch size s ($> s^*$) would perform better than with a smaller batch size. We also showed that the optimal batch size depends on the optimizer. In particular, it was shown that momentum and adaptive methods can exploit larger optimal batches than can SGD and that, if we can set an appropriate momentum coefficient β , then momentum and adaptive methods reduce $K_\epsilon(s^*)$ more than can SGD.

The results in this paper support theoretically the detailed numerical validations in recent papers [26, 31]. The learning rate used in [26, p.15] decayed linearly, which is distinctly different from both constant and diminishing learning rates. In the future, we should check

⁴The definitions of H and h_0^* imply that $(1 - \gamma)\sqrt{Hh_0^*} \leq (1 - \gamma)H$. The condition $(1 - \gamma)H \leq 1$ is sufficient to guarantee (16).

numerically the existences of optimal batch sizes of optimizers with not only constant but also diminishing learning rates to fully support all of the results in this paper.

Acknowledgments. This work was supported by Japan Society for the Promotion of Science (JSPS) KAKENHI Grant Number 21K11773.

REFERENCES

- [1] Arjovsky, M., Chintala, S., Bottou, L.: Wasserstein GAN <https://arxiv.org/pdf/1701.07875.pdf> (2017)
- [2] Bottou, L., Curtis, F.E., Nocedal, J.: Optimization methods for large-scale machine learning. *SIAM Review* **60**, 223–311 (2018)
- [3] Chen, H., Zheng, L., AL Kontar, R., Raskutti, G.: Stochastic gradient descent in correlated settings: A study on Gaussian processes. In: *Advances in Neural Information Processing Systems*, vol. 33, pp. 2722–2733 (2020)
- [4] Chen, X., Liu, S., Sun, R., Hong, M.: On the convergence of a class of Adam-type algorithms for non-convex optimization. In: *International Conference on Learning Representations* (2019)
- [5] Duchi, J., Hazan, E., Singer, Y.: Adaptive subgradient methods for online learning and stochastic optimization. *Journal of Machine Learning Research* **12**, 2121–2159 (2011)
- [6] Fang, C., Li, C.J., Lin, Z., Zhang, T.: SPIDER: Near-optimal non-convex optimization via stochastic path-integrated differential estimator. In: *Advances in Neural Information Processing Systems*, vol. 31 (2018)
- [7] Fehrman, B., Gess, B., Jentzen, A.: Convergence rates for the stochastic gradient descent method for non-convex objective functions. *Journal of Machine Learning Research* **21**, 1–48 (2020)
- [8] Ghadimi, S., Lan, G.: Optimal stochastic approximation algorithms for strongly convex stochastic composite optimization I: A generic algorithmic framework. *SIAM Journal on Optimization* **22**, 1469–1492 (2012)
- [9] Ghadimi, S., Lan, G.: Optimal stochastic approximation algorithms for strongly convex stochastic composite optimization II: Shrinking procedures and optimal algorithms. *SIAM Journal on Optimization* **23**, 2061–2089 (2013)
- [10] Gower, R.M., Sebbouh, O., Loizou, N.: SGD for structured nonconvex functions: Learning rates, minibatching and interpolation. In: *Proceedings of the 24th International Conference on Artificial Intelligence and Statistics (AISTATS)*, vol. 130 (2021)
- [11] Horn, R.A., Johnson, C.R.: *Matrix Analysis*. Cambridge University Press, Cambridge (1985)
- [12] Iiduka, H.: Appropriate learning rates of adaptive learning rate optimization algorithms for training deep neural networks. *IEEE Transactions on Cybernetics* https://iiduka.net/_media/iiduka/cyb-e-2021-05-1174.pdf (2021)
- [13] Kingma, D.P., Ba, J.: Adam: A method for stochastic optimization. In: *Proceedings of The International Conference on Learning Representations* (2015)
- [14] Loizou, N., Vaswani, S., Laradji, I., Lacoste-Julien, S.: Stochastic polyak step-size for SGD: An adaptive learning rate for fast convergence. In: *Proceedings of the 24th International Conference on Artificial Intelligence and Statistics (AISTATS)*, vol. 130 (2021)
- [15] Luo, L., Xiong, Y., Liu, Y., Sun, X.: Adaptive gradient methods with dynamic bound of learning rate. In: *Proceedings of The International Conference on Learning Representations* (2019)

- [16] Martens, J., Grosse, R.: Optimizing neural networks with Kronecker-factored approximate curvature. In: Proceedings of Machine Learning Research, vol. 37, pp. 2408–2417 (2015)
- [17] Mendler-Dünner, C., Perdomo, J.C., Zrnic, T., Hardt, M.: Stochastic optimization for performative prediction. In: Advances in Neural Information Processing Systems, vol. 33 (2020)
- [18] Nemirovski, A., Juditsky, A., Lan, G., Shapiro, A.: Robust stochastic approximation approach to stochastic programming. *SIAM Journal on Optimization* **19**, 1574–1609 (2009)
- [19] Nesterov, Y.: A method for unconstrained convex minimization problem with the rate of convergence $O(1/k^2)$. *Doklady AN USSR* **269**, 543–547 (1983)
- [20] Polyak, B.T.: Some methods of speeding up the convergence of iteration methods. *USSR Computational Mathematics and Mathematical Physics* **4**, 1–17 (1964)
- [21] Reddi, S.J., Kale, S., Kumar, S.: On the convergence of Adam and beyond. In: Proceedings of The International Conference on Learning Representations (2018)
- [22] Robbins, H., Monro, H.: A stochastic approximation method. *The Annals of Mathematical Statistics* **22**, 400–407 (1951)
- [23] Rumelhart, D.E., Hinton, G.E., Williams, R.J.: Learning representations by back-propagating errors. *Nature* **323**, 533–536 (1986)
- [24] Scaman, K., Malherbe, C.: Robustness analysis of non-convex stochastic gradient descent using biased expectations. In: Advances in Neural Information Processing Systems, vol. 33 (2020)
- [25] Schmidt, R.M., Schneider, F., Hennig, P.: Descending through a crowded valley—Benchmarking deep learning optimizers. arXiv, <https://arxiv.org/pdf/2007.01547.pdf> (2021)
- [26] Shallue, C.J., Lee, J., Antognini, J., Sohl-Dickstein, J., Frostig, R., Dahl, G.E.: Measuring the effects of data parallelism on neural network training. *Journal of Machine Learning Research* **20**, 1–49 (2019)
- [27] Sutskever, I., Martens, J., Dahl, G., Hinton, G.: On the importance of initialization and momentum in deep learning. In: Proceedings of the 30th International Conference on Machine Learning, pp. 1139–1147 (2013)
- [28] Tieleman, T., Hinton, G.: RMSProp: Divide the gradient by a running average of its recent magnitude. *COURSERA: Neural networks for machine learning* **4(2)**, 26–31 (2012)
- [29] Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, L., Polosukhin, I.: Attention is All you Need. In: Advances in Neural Information Processing Systems, vol. 30, pp. 5998–6008 (2017)
- [30] Xu, K., Ba, J., Kiros, R., Cho, K., Courville, A., Salakhudinov, R., Zemel, R., Bengio, Y.: Show, attend and tell: Neural image caption generation with visual attention. In: Proceedings of the 32nd International Conference on Machine Learning, vol. 37, pp. 2048–2057 (2015)
- [31] Zhang, G., Li, L., Nado, Z., Martens, J., Sachdeva, S., Dahl, G.E., Shallue, C.J., Grosse, R.: Which algorithmic choices matter at which batch sizes? Insights from a noisy quadratic model. In: Advances in Neural Information Processing Systems (2019)
- [32] Zhuang, J., Tang, T., Ding, Y., Tatikonda, S., Dvornek, N., Papademetris, X., Duncan, J.S.: AdaBelief optimizer: Adapting stepsizes by the belief in observed gradients. In: Advances in Neural Information Processing Systems (2020)

- [33] Zinkevich, M.: Online convex programming and generalized infinitesimal gradient ascent. In: Proceedings of the Twentieth International Conference on Machine Learning (ICML-2003), pp. 928–936 (2003)
- [34] Zinkevich, M., Weimer, M., Li, L., Smola, A.: Parallelized stochastic gradient descent. In: Advances in Neural Information Processing Systems, vol. 23 (2010)

APPENDIX A. APPENDIX

Unless stated otherwise, all relations between random variables are supported to hold almost surely. Let $S \in \mathbb{S}_{++}^d$. The S -inner product of \mathbb{R}^d is defined for all $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$ by $\langle \mathbf{x}, \mathbf{y} \rangle_S := \langle \mathbf{x}, S\mathbf{y} \rangle$ and the S -norm is defined by $\|\mathbf{x}\|_S := \sqrt{\langle \mathbf{x}, S\mathbf{x} \rangle}$. The history of process ξ_0, ξ_1, \dots to time step k is denoted by $\xi_{[k]} = (\xi_0, \xi_1, \dots, \xi_k)$.

A.1. Sufficient Conditions for Assumption (A3).

Proposition A.1. *Assumption (A3) holds if each of the following holds:*

- (G1) $X \subset \mathbb{R}^d$ is bounded, the gradient ∇f_i is Lipschitz continuous with Lipschitz constant L_i , $S_i := \{\mathbf{x}^* \in \mathbb{R}^d : \nabla f_i(\mathbf{x}^*) = \mathbf{0}\} \neq \emptyset$ ($i \in [n]$), where $L := \max_{i \in [n]} L_i$. (If we define $G_{k,L} := \sup_{\mathbf{x} \in X} \sum_{i \in \mathcal{S}_k} \|\nabla f_i(\mathbf{x})\|$, then we can take $G := \sup_{k \in \mathbb{N}} G_{k,L}$.)
- (G2) $X \subset \mathbb{R}^d$ is bounded and closed. (If we define $G_k := \sup_{\mathbf{x} \in X} \sum_{i \in \mathcal{S}_k} \|\nabla f_i(\mathbf{x})\|$, then we can take $G := \sup_{k \in \mathbb{N}} G_k$.)

Under (A6), G in (G1) and (G2) are respectively $G = Ln\sqrt{dD}$ and $G = n\tilde{G}$, where $\tilde{G} := \max_{i \in [n]} \sup_{\mathbf{x} \in X} \|\nabla f_i(\mathbf{x})\|$.

Proof: The definition of $\nabla f_{\mathcal{S}_k}$ and the triangle inequality imply that, for all $\mathbf{x} \in \mathbb{R}^d$ and all $k \in \mathbb{N}$,

$$(18) \quad \|\nabla f_{\mathcal{S}_k}(\mathbf{x})\|^2 = \left\| \frac{1}{s} \sum_{i \in \mathcal{S}_k} \nabla f_i(\mathbf{x}) \right\|^2 \leq \frac{1}{s^2} \left(\sum_{i \in \mathcal{S}_k} \|\nabla f_i(\mathbf{x})\| \right)^2.$$

Suppose that (G1) holds. Let $\mathbf{x}^* \in S_i$ ($i \in [n]$). The Cauchy–Schwarz inequality and the Lipschitz continuity of ∇f_i , together with the definition of L , ensure that, for all $\mathbf{x} \in \mathbb{R}^d$ and all $i \in [n]$,

$$\|\nabla f_i(\mathbf{x})\| \leq \|\nabla f_i(\mathbf{x}^*)\| + L_i \|\mathbf{x} - \mathbf{x}^*\| \leq L \|\mathbf{x} - \mathbf{x}^*\|.$$

Accordingly, we have that, for all $\mathbf{x} \in X$ and all $k \in \mathbb{N}$,

$$\sum_{i \in \mathcal{S}_k} \|\nabla f_i(\mathbf{x})\| \leq Ls \|\mathbf{x} - \mathbf{x}^*\| \leq Ln \|\mathbf{x} - \mathbf{x}^*\|.$$

Hence, $G_{k,L} \leq Ln \sup_{\mathbf{x} \in X} \|\mathbf{x} - \mathbf{x}^*\| < +\infty$. Taking the expectation of (18) thus implies (A3). Assumption (A6) implies that there exists a bounded set $X \subset \mathbb{R}^d$ such that $(\mathbf{x}_k)_{k \in \mathbb{N}} \subset X$. From $\|\mathbf{x}_k - \mathbf{x}^*\|^2 = \sum_{i \in [d]} (x_{k,i} - x_i)^2 \leq dD$, we have that, for all $k \in \mathbb{N}$,

$$G_{k,L} \leq Ln\sqrt{dD} =: G.$$

Suppose that (G2) holds. Since ∇f_i is continuous and X is compact, we have that $G = \sup_{k \in \mathbb{N}} G_k < +\infty$. Taking the expectation of (18) thus implies (A3). Assumption (A6) ensures that there exists a bounded, closed set $X \subset \mathbb{R}^d$ such that $(\mathbf{x}_k)_{k \in \mathbb{N}} \subset X$. Define $\tilde{G}_i := \sup_{\mathbf{x} \in X} \|\nabla f_i(\mathbf{x})\| < +\infty$ and $\tilde{G} := \max_{i \in [n]} \tilde{G}_i$. Then, we have that, for all $\mathbf{x} \in X$,

$$\sum_{i \in \mathcal{S}_k} \|\nabla f_i(\mathbf{x})\| \leq s\tilde{G} \leq n\tilde{G} =: G.$$

This completes the proof. \square

A.2. **Examples of Algorithm 1.** We list some examples of $\mathbf{H}_k \in \mathbb{S}_{++}^d \cap \mathbb{D}^d$ (step 5) in Algorithm 1.

TABLE 3. Examples of $\mathbf{H}_k \in \mathbb{S}_{++}^d \cap \mathbb{D}^d$ (step 5) in Algorithm 1 ($\delta, \zeta \in [0, 1)$)

	\mathbf{H}_k
SGD ($\beta_k = \gamma = 0$)	\mathbf{H}_k is the identity matrix.
N-Momentum [19] ($\gamma = 0$)	\mathbf{H}_k is the identity matrix.
AMSGrad [4] ($\gamma = 0$)	$\mathbf{v}_k = \delta \mathbf{v}_{k-1} + (1 - \delta) \nabla f_{\mathcal{S}_k}(\mathbf{x}_k) \odot \nabla f_{\mathcal{S}_k}(\mathbf{x}_k)$ $\hat{\mathbf{v}}_k = (\max\{\hat{v}_{k-1,i}, v_{k,i}\})_{i=1}^d$ $\mathbf{H}_k = \text{diag}(\sqrt{\hat{v}_{k,i}})$
AMSBound [15] ($\gamma = 0$)	$\mathbf{v}_k = \delta \mathbf{v}_{k-1} + (1 - \delta) \nabla f_{\mathcal{S}_k}(\mathbf{x}_k) \odot \nabla f_{\mathcal{S}_k}(\mathbf{x}_k)$ $\hat{\mathbf{v}}_k = (\max\{\hat{v}_{k-1,i}, v_{k,i}\})_{i=1}^d$ $\tilde{\mathbf{v}}_k = \left(\text{Clip} \left(\frac{1}{\sqrt{\hat{v}_{k,i}}}, l_k, u_k \right)^{-1} \right)_{i=1}^d$ $\mathbf{H}_k = \text{diag}(\sqrt{\tilde{v}_{k,i}})$
AdaBelief [32] ($s_{k,i} \leq s_{k+1,i}$ is needed)	$\tilde{\mathbf{s}}_k = (\nabla f_{\mathcal{S}_k}(\mathbf{x}_k) - \mathbf{m}_k) \odot (\nabla f_{\mathcal{S}_k}(\mathbf{x}_k) - \mathbf{m}_k)$ $\mathbf{s}_k = \delta \mathbf{v}_{k-1} + (1 - \delta) \tilde{\mathbf{s}}_k$ $\hat{\mathbf{s}}_k = \frac{\mathbf{s}_k}{1 - \zeta^k}$ $\mathbf{H}_k = \text{diag}(\sqrt{\hat{s}_{k,i}})$

We define $\mathbf{x} \odot \mathbf{x}$ for $\mathbf{x} := (x_i)_{i=1}^d \in \mathbb{R}^d$ by $\mathbf{x} \odot \mathbf{x} := (x_i^2)_{i=1}^d \in \mathbb{R}^d$. $\text{Clip}(\cdot, l, u): \mathbb{R} \rightarrow \mathbb{R}$ in AMSBound ($l, u \in \mathbb{R}$ with $l \leq u$ are given) is defined for all $x \in \mathbb{R}$ by

$$\text{Clip}(x, l, u) := \begin{cases} l & \text{if } x < l, \\ x & \text{if } l \leq x \leq u, \\ u & \text{if } x > u. \end{cases}$$

A.3. **Lemmas and Theorem.** The following are the key lemmas to prove the main theorems in this paper.

Lemma A.1. *Suppose that (A1) and (A2) hold and consider Algorithm 1. Then, for all $\mathbf{x} \in \mathbb{R}^d$ and all $k \in \mathbb{N}$,*

$$\mathbb{E} \left[\|\mathbf{x}_{k+1} - \mathbf{x}\|_{\mathbf{H}_k}^2 \right] \leq \mathbb{E} \left[\|\mathbf{x}_k - \mathbf{x}\|_{\mathbf{H}_k}^2 \right] + \alpha_k^2 \mathbb{E} \left[\|\mathbf{d}_k\|_{\mathbf{H}_k}^2 \right] + 2\alpha_k \left\{ \frac{\tilde{\beta}_k}{s\tilde{\gamma}_k} \mathbb{E} [\langle \mathbf{x} - \mathbf{x}_k, \nabla f(\mathbf{x}_k) \rangle] + \frac{\beta_k}{\tilde{\gamma}_k} \mathbb{E} [\langle \mathbf{x} - \mathbf{x}_k, \mathbf{m}_{k-1} \rangle] \right\},$$

where $\tilde{\beta}_k := 1 - \beta_k$ and $\tilde{\gamma}_k := 1 - \gamma^{k+1}$.

Proof: Let $\mathbf{x} \in \mathbb{R}^d$ and $k \in \mathbb{N}$. The definition of \mathbf{x}_{k+1} implies that

$$\|\mathbf{x}_{k+1} - \mathbf{x}\|_{\mathbf{H}_k}^2 \leq \|\mathbf{x}_k - \mathbf{x}\|_{\mathbf{H}_k}^2 + 2\alpha_k \langle \mathbf{x}_k - \mathbf{x}, \mathbf{d}_k \rangle_{\mathbf{H}_k} + \alpha_k^2 \|\mathbf{d}_k\|_{\mathbf{H}_k}^2.$$

Moreover, the definitions of \mathbf{d}_k , \mathbf{m}_k , and $\hat{\mathbf{m}}_k$ ensure that

$$\langle \mathbf{x}_k - \mathbf{x}, \mathbf{d}_k \rangle_{\mathbf{H}_k} = \frac{1}{\tilde{\gamma}_k} \langle \mathbf{x} - \mathbf{x}_k, \mathbf{m}_k \rangle = \frac{\beta_k}{\tilde{\gamma}_k} \langle \mathbf{x} - \mathbf{x}_k, \mathbf{m}_{k-1} \rangle + \frac{\tilde{\beta}_k}{\tilde{\gamma}_k} \langle \mathbf{x} - \mathbf{x}_k, \nabla f_{\mathcal{S}_k}(\mathbf{x}_k) \rangle,$$

where $\tilde{\beta}_k := 1 - \beta_k$ and $\tilde{\gamma}_k := 1 - \gamma^{k+1}$. Hence,

$$(19) \quad \begin{aligned} \|\mathbf{x}_{k+1} - \mathbf{x}\|_{\mathbf{H}_k}^2 &\leq \|\mathbf{x}_k - \mathbf{x}\|_{\mathbf{H}_k}^2 + 2\alpha_k \left\{ \frac{\beta_k}{\tilde{\gamma}_k} \langle \mathbf{x} - \mathbf{x}_k, \mathbf{m}_{k-1} \rangle + \frac{\tilde{\beta}_k}{\tilde{\gamma}_k} \langle \mathbf{x} - \mathbf{x}_k, \nabla f_{\mathcal{S}_k}(\mathbf{x}_k) \rangle \right\} \\ &\quad + \alpha_k^2 \|\mathbf{d}_k\|_{\mathbf{H}_k}^2. \end{aligned}$$

Meanwhile, the relationship between the expectation of the stochastic gradient vector $\nabla f_{\mathcal{S}_k}(\mathbf{x})$ and the full gradient vector $\nabla f(\mathbf{x})$ is as follows: For all $\mathbf{x} \in \mathbb{R}^d$,

$$(20) \quad \mathbb{E} [\nabla f_{\mathcal{S}_k}(\mathbf{x})] = \mathbb{E} \left[\frac{1}{s} \sum_{i \in \mathcal{S}_k} \nabla f_i(\mathbf{x}) \right] = \frac{1}{s} \mathbb{E} [\nabla f_i(\mathbf{x})] = \frac{1}{s} \nabla f(\mathbf{x}),$$

where the first equation comes from (A2), the second equation comes from the existence of T such that $[n] = \cup_{k=1}^T \mathcal{S}_k$, and the third equation comes from (A1). Condition (20) guarantees that

$$\begin{aligned} \mathbb{E} [\langle \mathbf{x} - \mathbf{x}_k, \nabla f_{\mathcal{S}_k}(\mathbf{x}_k) \rangle] &= \mathbb{E} \left[\mathbb{E} [\langle \mathbf{x} - \mathbf{x}_k, \nabla f_{\mathcal{S}_k}(\mathbf{x}_k) \rangle | \xi_{[k-1]}] \right] \\ &= \mathbb{E} [\langle \mathbf{x} - \mathbf{x}_k, \mathbb{E} [\nabla f_{\mathcal{S}_k}(\mathbf{x}_k) | \xi_{[k-1]}] \rangle] \\ &= \frac{1}{s} \mathbb{E} [\langle \mathbf{x} - \mathbf{x}_k, \nabla f(\mathbf{x}_k) \rangle]. \end{aligned}$$

Therefore, the lemma follows by taking the expectation of (19). \square

Lemma A.2. *Algorithm 1 satisfies that, under (A3), for all $k \in \mathbb{N}$,*

$$\mathbb{E} [\|\mathbf{m}_k\|^2] \leq \frac{G^2}{s^2}.$$

Under (A3) and (A4), for all $k \in \mathbb{N}$,

$$\mathbb{E} [\|\mathbf{d}_k\|_{\mathbf{H}_k}^2] \leq \frac{G^2}{(1-\gamma)^2 h_0^* s^2},$$

where $h_0^* := \min_{i \in [d]} h_{0,i}$.

Proof: The convexity of $\|\cdot\|^2$, together with the definition of \mathbf{m}_k and (A3), guarantees that, for all $k \in \mathbb{N}$,

$$\begin{aligned} \mathbb{E} [\|\mathbf{m}_k\|^2] &\leq \beta_k \mathbb{E} [\|\mathbf{m}_{k-1}\|^2] + (1 - \beta_k) \mathbb{E} [\|\nabla f_{\mathcal{S}_k}(\mathbf{x}_k)\|^2] \\ &\leq \beta_k \mathbb{E} [\|\mathbf{m}_{k-1}\|^2] + (1 - \beta_k) \frac{G^2}{s^2}. \end{aligned}$$

Induction thus ensures that, for all $k \in \mathbb{N}$,

$$(21) \quad \mathbb{E} [\|\mathbf{m}_k\|^2] \leq \max \left\{ \|\mathbf{m}_{-1}\|^2, \frac{G^2}{s^2} \right\} = \frac{G^2}{s^2},$$

where $\mathbf{m}_{-1} = \mathbf{0}$ is used. For $k \in \mathbb{N}$, $\mathbf{H}_k \in \mathbb{S}_{++}^d$ guarantees the existence of a unique matrix $\bar{\mathbf{H}}_k \in \mathbb{S}_{++}^d$ such that $\mathbf{H}_k = \bar{\mathbf{H}}_k^2$ [11, Theorem 7.2.6]. We have that, for all $\mathbf{x} \in \mathbb{R}^d$, $\|\mathbf{x}\|_{\mathbf{H}_k}^2 = \|\bar{\mathbf{H}}_k \mathbf{x}\|^2$. Accordingly, the definitions of \mathbf{d}_k and $\hat{\mathbf{m}}_k$ imply that, for all $k \in \mathbb{N}$,

$$\mathbb{E} [\|\mathbf{d}_k\|_{\mathbf{H}_k}^2] = \mathbb{E} \left[\left\| \bar{\mathbf{H}}_k^{-1} \mathbf{H}_k \mathbf{d}_k \right\|^2 \right] \leq \frac{1}{\tilde{\gamma}_k^2} \mathbb{E} \left[\left\| \bar{\mathbf{H}}_k^{-1} \right\|^2 \|\mathbf{m}_k\|^2 \right] \leq \frac{1}{(1-\gamma)^2} \mathbb{E} \left[\left\| \bar{\mathbf{H}}_k^{-1} \right\|^2 \|\mathbf{m}_k\|^2 \right],$$

where

$$\left\| \bar{\mathbf{H}}_k^{-1} \right\| = \left\| \text{diag} \left(h_{k,i}^{-\frac{1}{2}} \right) \right\| = \max_{i \in [d]} h_{k,i}^{-\frac{1}{2}}$$

and $\tilde{\gamma}_k := 1 - \gamma^{k+1} \geq 1 - \gamma$. Moreover, (A4) ensures that, for all $k \in \mathbb{N}$,

$$h_{k,i} \geq h_{0,i} \geq h_0^* := \min_{i \in [d]} h_{0,i}.$$

Hence, (21) implies that, for all $k \in \mathbb{N}$,

$$\mathbb{E} [\|\mathbf{d}_k\|_{\mathbf{H}_k}^2] \leq \frac{G^2}{(1-\gamma)^2 h_0^* s^2},$$

completing the proof. \square

We are in the position to prove the following theorem, which leads to Theorems 3.1, 3.2, and A.2.

Theorem A.1. *Suppose that Assumptions 2.1 and 2.2 hold and consider Algorithm 1. Let $(\delta_k)_{k \in \mathbb{N}} \subset (0, +\infty)$ be the sequence defined by $\delta_k := \alpha_k \tilde{\beta}_k / \tilde{\gamma}_k$ and $V_k(\mathbf{x}) := \mathbb{E}[\langle \mathbf{x}_k - \mathbf{x}, \nabla f(\mathbf{x}_k) \rangle]$ for all $\mathbf{x} \in \mathbb{R}^d$ and all $k \in \mathbb{N}$. Assume that $(\delta_k)_{k \in \mathbb{N}}$ is monotone decreasing. Then, for all $\mathbf{x} \in \mathbb{R}^d$ and all $K \geq 1$,*

$$\sum_{k=1}^K V_k(\mathbf{x}) \leq \frac{dsDH}{2\tilde{b}\alpha_K} + \frac{G^2}{2\tilde{b}\tilde{\gamma}^2 h_0^* s} \sum_{k=1}^K \alpha_k + \frac{\sqrt{dDG}}{\tilde{b}} \sum_{k=1}^K \beta_k,$$

where $\tilde{b} := 1 - b$, $\tilde{\gamma} := 1 - \gamma$, D and H_i are defined as in Assumption 2.2, and $H := \max_{i \in [d]} H_i$.

Proof: Let $\mathbf{x} \in \mathbb{R}^d$. Lemma A.1 guarantees that, for all $k \in \mathbb{N}$,

$$\begin{aligned} V_k(\mathbf{x}) &\leq \frac{s}{2\delta_k} \left\{ \mathbb{E} [\|\mathbf{x}_k - \mathbf{x}\|_{\mathbf{H}_k}^2] - \mathbb{E} [\|\mathbf{x}_{k+1} - \mathbf{x}\|_{\mathbf{H}_k}^2] \right\} + \frac{s\alpha_k \tilde{\gamma}_k}{2\tilde{\beta}_k} \mathbb{E} [\|\mathbf{d}_k\|_{\mathbf{H}_k}^2] \\ &\quad + \frac{s\beta_k}{\tilde{\beta}_k} \mathbb{E} [\langle \mathbf{x} - \mathbf{x}_k, \mathbf{m}_{k-1} \rangle]. \end{aligned}$$

Summing the above inequality from $k = 1$ to $K \geq 1$ implies that

$$\begin{aligned} \sum_{k=1}^K V_k(\mathbf{x}) &\leq \frac{1}{2} \underbrace{\sum_{k=1}^K \frac{s}{\delta_k} \left\{ \mathbb{E} [\|\mathbf{x}_k - \mathbf{x}\|_{\mathbf{H}_k}^2] - \mathbb{E} [\|\mathbf{x}_{k+1} - \mathbf{x}\|_{\mathbf{H}_k}^2] \right\}}_{\Delta_K} \\ (22) \quad &\quad + \frac{1}{2} \underbrace{\sum_{k=1}^K \frac{s\alpha_k \tilde{\gamma}_k}{\tilde{\beta}_k} \mathbb{E} [\|\mathbf{d}_k\|_{\mathbf{H}_k}^2]}_{A_K} + \underbrace{\sum_{k=1}^K \frac{s\beta_k}{\tilde{\beta}_k} \mathbb{E} [\langle \mathbf{x} - \mathbf{x}_k, \mathbf{m}_{k-1} \rangle]}_{B_K}. \end{aligned}$$

Let us define $\theta_k := \delta_k/s$. From the definition of Δ_K and $\mathbb{E}[\|\mathbf{x}_{k+1} - \mathbf{x}\|_{\mathbf{H}_k}^2]/\theta_k \geq 0$,

$$(23) \quad \Delta_K \leq \frac{\mathbb{E} [\|\mathbf{x}_1 - \mathbf{x}\|_{\mathbf{H}_1}^2]}{\theta_1} + \underbrace{\sum_{k=2}^K \left\{ \frac{\mathbb{E} [\|\mathbf{x}_k - \mathbf{x}\|_{\mathbf{H}_k}^2]}{\theta_k} - \frac{\mathbb{E} [\|\mathbf{x}_k - \mathbf{x}\|_{\mathbf{H}_{k-1}}^2]}{\theta_{k-1}} \right\}}_{\tilde{\Delta}_K}.$$

Since $\bar{\mathbf{H}}_k \in \mathbb{S}_{++}^d$ exists such that $\mathbf{H}_k = \bar{\mathbf{H}}_k^2$, we have $\|\mathbf{x}\|_{\mathbf{H}_k}^2 = \|\bar{\mathbf{H}}_k \mathbf{x}\|^2$ for all $\mathbf{x} \in \mathbb{R}^d$. Accordingly, we have

$$\tilde{\Delta}_K = \mathbb{E} \left[\sum_{k=2}^K \left\{ \frac{\|\bar{\mathbf{H}}_k(\mathbf{x}_k - \mathbf{x})\|^2}{\theta_k} - \frac{\|\bar{\mathbf{H}}_{k-1}(\mathbf{x}_k - \mathbf{x})\|^2}{\theta_{k-1}} \right\} \right].$$

From $\bar{\mathbf{H}}_k = \text{diag}(\sqrt{h_{k,i}})$, we have that, for all $\mathbf{x} = (x_i)_{i=1}^d \in \mathbb{R}^d$, $\|\bar{\mathbf{H}}_k \mathbf{x}\|^2 = \sum_{i=1}^d h_{k,i} x_i^2$. Hence, for all $K \geq 2$,

$$\tilde{\Delta}_K = \mathbb{E} \left[\sum_{k=2}^K \sum_{i=1}^d \left(\frac{h_{k,i}}{\theta_k} - \frac{h_{k-1,i}}{\theta_{k-1}} \right) (x_{k,i} - x_i)^2 \right].$$

Accordingly, from (A4) and the monotone decrease of $(\delta_k)_{k \in \mathbb{N}}$, we have that, for all $k \geq 1$ and all $i \in [d]$,

$$\frac{h_{k,i}}{\theta_k} - \frac{h_{k-1,i}}{\theta_{k-1}} \geq 0.$$

Moreover, from (A6), $D := \max_{i \in [d]} \sup_{k \in \mathbb{N}} (x_{k,i} - x_i)^2 < +\infty$. Accordingly, for all $K \geq 2$,

$$\tilde{\Delta}_K \leq D \mathbb{E} \left[\sum_{k=2}^K \sum_{i=1}^d \left(\frac{h_{k,i}}{\theta_k} - \frac{h_{k-1,i}}{\theta_{k-1}} \right) \right] = D \mathbb{E} \left[\sum_{i=1}^d \left(\frac{h_{K,i}}{\theta_K} - \frac{h_{1,i}}{\theta_1} \right) \right].$$

Therefore, (23), $\mathbb{E}[\|\mathbf{x}_1 - \mathbf{x}\|_{\mathbb{H}_1}^2] / \theta_1 \leq D \mathbb{E}[\sum_{i=1}^d h_{1,i} / \theta_1]$, and (A5) imply, for all $K \in \mathbb{N}$,

$$\Delta_K \leq D \mathbb{E} \left[\sum_{i=1}^d \frac{h_{1,i}}{\theta_1} \right] + D \mathbb{E} \left[\sum_{i=1}^d \left(\frac{h_{K,i}}{\theta_K} - \frac{h_{1,i}}{\theta_1} \right) \right] = \frac{D}{\theta_K} \mathbb{E} \left[\sum_{i=1}^d h_{K,i} \right] \leq \frac{D}{\theta_K} \sum_{i=1}^d H_i,$$

which, together with $\theta_K := \alpha_K(1 - \beta_K)/(s(1 - \gamma^{K+1})) \geq \tilde{b}\alpha_K/s$ and $H = \max_{i \in [d]} H_i$, implies

$$(24) \quad \frac{1}{2} \Delta_K \leq \frac{dsDH}{2\tilde{b}\alpha_K}.$$

Lemma A.2 implies that, for all $K \in \mathbb{N}$,

$$A_K := \sum_{k=1}^K \frac{s\alpha_k \tilde{\gamma}_k}{\tilde{\beta}_k} \mathbb{E} \left[\|\mathbf{d}_k\|_{\mathbb{H}_k}^2 \right] \leq \sum_{k=1}^K \frac{s\alpha_k \tilde{\gamma}_k}{\tilde{\beta}_k} \frac{G^2}{\tilde{\gamma}^2 h_0^* s^2},$$

which, together with $\tilde{\gamma}_k \leq 1$ and $\beta_k \leq b$, implies that

$$(25) \quad \frac{1}{2} A_K \leq \frac{G^2}{2\tilde{b}\tilde{\gamma}^2 h_0^* s} \sum_{k=1}^K \alpha_k.$$

Lemma A.2 and Jensen's inequality ensure that, for all $k \in \mathbb{N}$,

$$\mathbb{E}[\|\mathbf{m}_k\|] \leq \frac{G}{s}.$$

The Cauchy-Schwarz inequality and (A6) guarantee that, for all $K \in \mathbb{N}$,

$$(26) \quad B_K := \sum_{k=1}^K \frac{s\beta_k}{\tilde{\beta}_k} \mathbb{E}[\langle \mathbf{x} - \mathbf{x}_k, \mathbf{m}_{k-1} \rangle] \leq \sum_{k=1}^K \frac{s\sqrt{dD}\beta_k}{\tilde{b}} \mathbb{E}[\|\mathbf{m}_{k-1}\|] \leq \frac{\sqrt{dD}G}{\tilde{b}} \sum_{k=1}^K \beta_k.$$

Therefore, (22), (24), (25), and (26) lead to the assertion in Theorem A.1. This completes the proof. \square

A.4. Proof of Theorem 3.1. (i) Theorem A.1, together with $\alpha_k = \alpha$ and $\beta_k = \beta$, guarantees that , for all $K \geq 1$ and all $\mathbf{x} \in \mathbb{R}^d$,

$$(27) \quad \frac{1}{K} \sum_{k=1}^K V_k(\mathbf{x}) \leq \frac{dDH}{2\tilde{b}\alpha} \frac{s}{K} + \frac{G^2\alpha}{2\tilde{b}\tilde{\gamma}^2 h_0^*} \frac{1}{s} + \frac{\sqrt{dDG}}{\tilde{b}} \beta.$$

Moreover, there exists $m \in [K]$ such that, for all $\mathbf{x} \in \mathbb{R}^d$,

$$(28) \quad \mathbb{E}[\langle \mathbf{x}_m - \mathbf{x}, \nabla f(\mathbf{x}_m) \rangle] = V_m(\mathbf{x}) = \min_{k \in [K]} V_k(\mathbf{x}) \leq \frac{1}{K} \sum_{k=1}^K V_k(\mathbf{x}).$$

Setting $\mathbf{x} = \mathbf{x}_m - \nabla f(\mathbf{x}_m)$, together with (27) and (28), guarantees that

$$(29) \quad \min_{k \in [K]} \mathbb{E}[\|\nabla f(\mathbf{x}_k)\|^2] \leq \mathbb{E}[\|\nabla f(\mathbf{x}_m)\|^2] \leq \underbrace{\frac{dDH}{2\tilde{b}\alpha}}_{A_\alpha} \frac{s}{K} + \underbrace{\frac{G^2\alpha}{2\tilde{b}\tilde{\gamma}^2 h_0^*}}_{B_\alpha} \frac{1}{s} + \underbrace{\frac{\sqrt{dDG}}{\tilde{b}}}_{C_\beta} \beta.$$

(ii) A sufficient condition for (1), i.e.,

$$\min_{k \in [K]} \mathbb{E}[\|\nabla f(\mathbf{x}_k)\|^2] \leq \epsilon^2$$

is that the right-hand side of (29) is equal to ϵ^2 , i.e.,

$$A_\alpha s^2 + B_\alpha K + (C_\beta - \epsilon^2)sK = 0,$$

which implies that

$$K(s) = \frac{A_\alpha s^2}{(\epsilon^2 - C_\beta)s - B_\alpha} \quad \left(s \in \left(\frac{B_\alpha}{\epsilon^2 - C_\beta}, +\infty \right) \right),$$

where $\epsilon^2 - C_\beta > 0$ is guaranteed from $\beta < \tilde{b}\epsilon^2/\sqrt{dDG}$. We have that

$$\frac{dK(s)}{ds} = \frac{A_\alpha s}{\{(C_\beta - \epsilon^2)s + B_\alpha\}^2} \left\{ \begin{array}{l} < 0 & \text{if } s \in \left(\frac{B_\alpha}{\epsilon^2 - C_\beta}, s^* \right), \\ = 0 & \text{if } s = s^* = \frac{2B_\alpha}{\epsilon^2 - C_\beta}, \\ > 0 & \text{if } s \in (s^*, +\infty). \end{array} \right.$$

Hence, $K(s)$ attains the minimum $K(s^*)$ when $s = s^*$. \square

A.5. Proof of Theorem 3.2. (i) Theorem A.1, together with $\alpha_k = \alpha/\sqrt{k}$ and $\beta_k = \beta$, guarantees that, for all $K \geq 1$ and all $\mathbf{x} \in \mathbb{R}^d$,

$$(30) \quad \begin{aligned} \frac{1}{K} \sum_{k=1}^K V_k(\mathbf{x}) &\leq \frac{dDH}{2\tilde{b}} \frac{s}{\alpha_K K} + \frac{G^2}{2\tilde{b}\tilde{\gamma}^2 h_0^*} \frac{1}{sK} \sum_{k=1}^K \alpha_k + \frac{\sqrt{dDG}}{\tilde{b}} \beta \\ &\leq \frac{dDH}{2\tilde{b}\alpha} \frac{s}{\sqrt{K}} + \frac{G^2\alpha}{\tilde{b}\tilde{\gamma}^2 h_0^*} \frac{1}{s\sqrt{K}} + \frac{\sqrt{dDG}}{\tilde{b}} \beta, \end{aligned}$$

where we use

$$\frac{1}{K} \sum_{k=1}^K \frac{1}{\sqrt{k}} \leq \frac{1}{K} \left(1 + \int_1^K \frac{dt}{\sqrt{t}} \right) = \frac{1}{K} (2\sqrt{K} - 1) \leq \frac{2}{\sqrt{K}}.$$

An argument similar to the one for showing (28) and (29) ensures that (30) implies that

$$(31) \quad \min_{k \in [K]} \mathbb{E} [\|\nabla f(\mathbf{x}_k)\|^2] \leq \underbrace{\frac{dDH}{2\tilde{b}\alpha}}_{A_\alpha} \frac{s}{\sqrt{K}} + \underbrace{\frac{G^2\alpha}{\tilde{b}\tilde{\gamma}^2 h_0^*}}_{B_\alpha} \frac{1}{s\sqrt{K}} + \underbrace{\frac{\sqrt{dDG}}{\tilde{b}}}_{C_\beta} \beta.$$

(ii) A sufficient condition for (1), i.e.,

$$\min_{k \in [K]} \mathbb{E} [\|\nabla f(\mathbf{x}_k)\|^2] \leq \epsilon^2$$

is that the right-hand side of (31) is equal to ϵ^2 , i.e.,

$$A_\alpha s^2 + (C_\beta - \epsilon^2)s\sqrt{K} + B_\alpha = 0,$$

which implies that

$$K(s) = \left\{ \frac{A_\alpha s^2 + B_\alpha}{(\epsilon^2 - C_\beta)s} \right\}^2,$$

where $\epsilon^2 - C_\beta > 0$ is guaranteed from $\beta < \tilde{b}\epsilon^2/\sqrt{dDG}$. We have that

$$\frac{dK(s)}{ds} = \frac{2(A_\alpha s^2 + B_\alpha)}{(\epsilon^2 - C_\beta)^2 s^3} (A_\alpha s^2 - B_\alpha) \begin{cases} < 0 & \text{if } s \in (0, s^*), \\ = 0 & \text{if } s = s^* = \sqrt{\frac{B_\alpha}{A_\alpha}}, \\ > 0 & \text{if } s \in (s^*, +\infty), \end{cases}$$

which implies that $K(s)$ attains the minimum $K(s^*)$ when $s = s^*$. \square

A.6. Relationship between s and $K_\epsilon(s)$ for Algorithm 1 with diminishing learning rates. The following is a result for Algorithm 1 with diminishing sequences α_k and β_k .

Theorem A.2. *Suppose that Assumptions 2.1 and 2.2 hold and let $s, \epsilon > 0$, $\alpha \in (0, 1]$, and $\beta \in [0, b] \subset [0, 1]$.*

(i) *Consider Algorithm 1 with*

$$\alpha_k := \frac{\alpha}{\sqrt{k}} \text{ and } \beta_k := \beta^k.$$

Then, for all $K \geq 1$,

$$\min_{k \in [K]} \mathbb{E} [\|\nabla f(\mathbf{x}_k)\|^2] \leq \underbrace{\frac{dDH}{2(1-b)\alpha}}_{A_\alpha} \frac{s}{\sqrt{K}} + \underbrace{\frac{G^2\alpha}{(1-b)(1-\gamma)^2 h_0^*}}_{B_\alpha} \frac{1}{s\sqrt{K}} + \underbrace{\frac{\beta\sqrt{dDG}}{(1-b)(1-\beta)}}_{C_\beta} \frac{1}{K}.$$

(ii) *The number of steps K_ϵ needed to achieve (1) is expressed as the following rational function of batch size s :*

$$K_\epsilon(s) = \left\{ \frac{(A_\alpha s^2 + B_\alpha) + \sqrt{(A_\alpha s^2 + B_\alpha)^2 + 4\epsilon^2 C_\beta s^2}}{2\epsilon^2 s} \right\}^2.$$

In particular, the minimum value of K_ϵ needed to achieve (1) is

$$\begin{aligned} K_\epsilon(s^*) &= \left\{ \frac{\sqrt{A_\alpha B_\alpha} + \sqrt{A_\alpha B_\alpha + \epsilon^2 C_\beta}}{\epsilon^2} \right\}^2 \\ &= \frac{\left\{ \sqrt{(1-\beta)dDHG} + \sqrt{\left((1-\beta)dDG + 2(1-b)\beta(1-\gamma)^2 \epsilon^2 \sqrt{dD} h_0^* \right) G} \right\}^2}{2(1-b)^2(1-\beta)(1-\gamma)^2 \epsilon^4 h_0^*} \end{aligned}$$

when

$$s^* = \sqrt{\frac{B_\alpha}{A_\alpha}} = \frac{G\alpha}{(1-\gamma)\sqrt{dDH}h_0^*}.$$

Proof: (i) Theorem A.1, together with $\alpha_k = 1/\sqrt{k}$ and $\beta_k = \beta^k$, guarantees that, for all $K \geq 1$ and all $\mathbf{x} \in \mathbb{R}^d$,

$$(32) \quad \begin{aligned} \frac{1}{K} \sum_{k=1}^K V_k(\mathbf{x}) &\leq \frac{dDH}{2\tilde{b}} \frac{s}{\alpha_K K} + \frac{G^2}{2\tilde{b}\tilde{\gamma}^2 h_0^*} \frac{1}{sK} \sum_{k=1}^K \alpha_k + \frac{\sqrt{dDG}}{\tilde{b}} \frac{1}{K} \sum_{k=1}^K \beta^k \\ &\leq \frac{dDH}{2\tilde{b}\alpha} \frac{s}{\sqrt{K}} + \frac{G^2\alpha}{\tilde{b}\tilde{\gamma}^2 h_0^*} \frac{1}{s\sqrt{K}} + \frac{\beta\sqrt{dDG}}{\tilde{b}\tilde{\beta}} \frac{1}{K}, \end{aligned}$$

where we use $\tilde{\beta} := 1 - \beta$,

$$\begin{aligned} \frac{1}{K} \sum_{k=1}^K \frac{1}{\sqrt{k}} &\leq \frac{1}{K} \left(1 + \int_1^K \frac{dt}{\sqrt{t}} \right) = \frac{1}{K} (2\sqrt{K} - 1) \leq \frac{2}{\sqrt{K}}, \\ \frac{1}{K} \sum_{k=1}^K \beta^k &\leq \frac{1}{K} \sum_{k=1}^{+\infty} \beta^k = \frac{\beta}{\tilde{\beta}K}. \end{aligned}$$

An argument similar to the one for showing (28) and (29) ensures that (32) implies that

$$(33) \quad \min_{k \in [K]} \mathbb{E} [\|\nabla f(\mathbf{x}_k)\|^2] \leq \underbrace{\frac{dDH}{2\tilde{b}\alpha}}_{A_\alpha} \frac{s}{\sqrt{K}} + \underbrace{\frac{G^2\alpha}{\tilde{b}\tilde{\gamma}^2 h_0^*}}_{B_\alpha} \frac{1}{s\sqrt{K}} + \underbrace{\frac{\beta\sqrt{dDG}}{\tilde{b}\tilde{\beta}}}_{C_\beta} \frac{1}{K}.$$

(ii) A sufficient condition for (1), i.e.,

$$\min_{k \in [K]} \mathbb{E} [\|\nabla f(\mathbf{x}_k)\|^2] \leq \epsilon^2$$

is that the right-hand side of (33) is equal to ϵ^2 , i.e.,

$$A_\alpha s^2 \sqrt{K} + B_\alpha \sqrt{K} + C_\beta s - \epsilon^2 s K = 0,$$

which implies that

$$K(s) = \left\{ \frac{(A_\alpha s^2 + B_\alpha) + \sqrt{(A_\alpha s^2 + B_\alpha)^2 + 4\epsilon^2 C_\beta s^2}}{2\epsilon^2 s} \right\}^2.$$

We have that

$$\frac{d\sqrt{K(s)}}{ds} = \frac{(A_\alpha s^2 + B_\alpha) + \sqrt{(A_\alpha s^2 + B_\alpha)^2 + 4\epsilon^2 C_\beta s^2}}{2\epsilon^2 s^2 \sqrt{(A_\alpha s^2 + B_\alpha)^2 + 4\epsilon^2 C_\beta s^2}} (A_\alpha s^2 - B_\alpha) \begin{cases} < 0 & \text{if } s \in (0, s^*), \\ = 0 & \text{if } s = s^* = \sqrt{\frac{B_\alpha}{A_\alpha}}, \\ > 0 & \text{if } s \in (s^*, +\infty), \end{cases}$$

which implies that $K(s)$ attains the minimum $K(s^*)$ when $s = s^*$. \square