# Iteration and Stochastic First-order Oracle Complexities of Stochastic Gradient Descent using Constant and Decaying Learning Rates

Kento Imaizumi[*1] and Hideaki Iiduka[*1]  [*]Equal contribution; [1]Department of Computer Science, Meiji University, Japan

**ABSTRACT**
The performance of stochastic gradient descent (SGD), which is the simplest first-order optimizer for training deep neural networks, depends on not only the learning rate but also the batch size. They both affect the number of iterations and the stochastic first-order oracle (SFO) complexity needed for training. In particular, the previous numerical results indicated that, for SGD using a constant learning rate, the number of iterations needed for training decreases when the batch size increases, and the SFO complexity needed for training is minimized at a critical batch size and that it increases once the batch size exceeds that size. Here, we study the relationship between batch size and the iteration and SFO complexities needed for nonconvex optimization in deep learning with SGD using constant or decaying learning rates and show that SGD using the critical batch size minimizes the SFO complexity. We also provide numerical comparisons of SGD with the existing first-order optimizers and show the usefulness of SGD using a critical batch size. Moreover, we show that measured critical batch sizes are close to the sizes estimated from our theoretical results.

**KEYWORDS**
SGD, batch size, iteration complexity, SFO complexity, nonconvex optimization

## 1. Introduction

### 1.1. Background

First-order optimizers can train deep neural networks by minimizing loss functions called the expected and empirical risks. They use stochastic first-order derivatives (stochastic gradients), which are estimated from the full gradient of the loss function. The simplest first-order optimizer is stochastic gradient descent (SGD) [1–5], which has a number of variants, including momentum variants[6,7] and numerous adaptive variants, such as adaptive gradient (AdaGrad) [8], root mean square propagation (RM-SProp) [9], adaptive moment estimation (Adam) [10], adaptive mean square gradient (AMSGrad) [11], and Adam with decoupled weight decay (AdamW) [12].

SGD can be applied to nonconvex optimization [13–22], where its performance strongly depends on the learning rate $\alpha_k$. For example, under the bounded variance assumption, SGD using a constant learning rate $\alpha_k = \alpha$ satisfies $\frac{1}{K} \sum_{k=0}^{K-1} \|\nabla f(\boldsymbol{\theta}_k)\|^2 =$

$O(\frac{1}{K}) + \sigma^2$ [18, Theorem 12] and SGD using a decaying learning rate (i.e., $\alpha_k \to 0$) satisfies that $\frac{1}{K} \sum_{k=0}^{K-1} \mathbb{E}[\|\nabla f(\boldsymbol{\theta}_k)\|^2] = O(\frac{1}{\sqrt{K}})$ [18, Theorem 11], where $(\boldsymbol{\theta}_k)_{k \in \mathbb{N}}$ is the sequence generated by SGD to find a local minimizer of $f$, $K$ is the number of iterations, and $\sigma^2$ is the upper bound of the variance.

The performance of SGD also depends on the batch size $b$. The convergence analyses reported in [14,17,21,23,24] indicated that SGD with a decaying learning rate and large batch size converges to a local minimizer of the loss function. In [25], it was numerically shown that using an enormous batch reduces both the number of parameter updates and model training time. Moreover, setting appropriate batch sizes for optimizers used in training generative adversarial networks were investigated in [26].

## 1.2. Motivation

The previous numerical results in [27] indicated that, for SGD using constant or linearly decaying learning rates, the number of iterations $K$ needed to train a deep neural network decreases as the batch size $b$ increases. Motivated by the numerical results in [27], we decided to clarify *the theoretical iteration complexity* of SGD with a constant or decaying learning rate in training a deep neural network. We used the performance measure of previous theoretical analyses of SGD, i.e., $\min_{k \in [0:K-1]} \mathbb{E}[\|\nabla f(\boldsymbol{\theta}_k)\|] \leq \epsilon$, where $\epsilon \ (> 0)$ is the precision and $[0 : K - 1] := \{0, 1, \ldots, K - 1\}$. We found that, if SGD is an $\epsilon$–approximation, i.e., $\min_{k \in [0:K-1]} \mathbb{E}[\|\nabla f(\boldsymbol{\theta}_k)\|] \leq \epsilon$, then it can train a deep neural network in $K$ iterations.

In addition, the numerical results in [27] indicated an interesting fact wherein diminishing returns exist beyond a critical batch size; i.e., the number of iterations needed to train a deep neural network does not strictly decrease beyond the critical batch size. Here, we define *the stochastic first-order oracle (SFO) complexity* as $N := Kb$, where $K$ is the number of iterations needed to train a deep neural network and $b$ is the batch size, as stated above. The deep neural network model uses $b$ gradients of the loss functions per iteration. The model has a stochastic gradient computation cost of $N = Kb$. From the numerical results in [27, Figures 4 and 5], we can conclude that the critical batch size $b^\star$ (if it exists) is useful for SGD, since the SFO complexity $N(b)$ is minimized at $b = b^\star$ and the SFO complexity increases once the batch size exceeds $b^\star$. Hence, on the basis of the first motivation stated above, we decided to clarify the SFO complexities needed for SGD using a constant or decaying learning rate to be an $\epsilon$–approximation.

## 1.3. Contribution

### 1.3.1. Upper bound of theoretical performance measure

To clarify the iteration and SFO complexities needed for SGD to be an $\epsilon$–approximation, we first give upper bounds of $\min_{k \in [0:K-1]} \mathbb{E}[\|\nabla f(\boldsymbol{\theta}_k)\|^2]$ for SGD to generate a sequence $(\boldsymbol{\theta}_k)_{k \in \mathbb{N}}$ with constant or decaying learning rates (see Theorem 3.1 for the definitions of $C_i$ and $D_i$). As our aim is to show that SGD is an $\epsilon$–approximation $\min_{k \in [0:K-1]} \mathbb{E}[\|\nabla f(\boldsymbol{\theta}_k)\|^2] \leq \epsilon^2$, it is desirable that the upper bounds of $\min_{k \in [0:K-1]} \mathbb{E}[\|\nabla f(\boldsymbol{\theta}_k)\|^2]$ be small. Table 1 indicates that the upper bounds become small when the number of iterations and batch size are large. The table also indicates that the convergence of SGD strongly depends on the batch size, since the variance terms (including $\sigma^2$ and $b$; see Theorem 3.1 for the definitions of $C_2$ and $D_2$)

in the upper bounds of $\min_{k \in [0:K-1]} \mathbb{E}[\|\nabla f(\boldsymbol{\theta}_k)\|^2]$ decrease as the batch size becomes larger.

**Table 1.** Upper bounds of $\min_{k \in [0:K-1]} \mathbb{E}[\|\nabla f(\boldsymbol{\theta}_k)\|^2]$ for SGD using a constant or decaying learning rate and the critical batch size to minimize the SFO complexities and achieve $\min_{k \in [0:K-1]} \mathbb{E}[\|\nabla f(\boldsymbol{\theta}_k)\|] \leq \epsilon$ ($C_i$ and $D_i$ are positive constants, $K$ is the number of iterations, $b$ is the batch size, $T \geq 1$, $\epsilon > 0$, and $L$ is the Lipschitz constant of $\nabla f$)

| Learning Rate | | Upper Bound | Critical Batch Size |
|---|---|---|---|
| Constant $\alpha \in (0, \frac{2}{L})$ | | $\dfrac{C_1}{K} + \dfrac{C_2}{b}$ | $\dfrac{2C_2}{\epsilon^2}$ |
| Decay $\alpha_k = \dfrac{1}{\left(\lfloor \frac{k}{T} \rfloor + 1\right)^a}$ | $a \in (0, \frac{1}{2})$ | $\dfrac{D_1}{K^a} + \dfrac{D_2}{(1-2a)K^a b}$ | $\dfrac{(1-a)D_2}{a(1-2a)D_1}$ |
| | $a = \frac{1}{2}$ | $\dfrac{D_1}{\sqrt{K}} + \left(\dfrac{1}{\sqrt{K}} + 1\right)\dfrac{D_2}{b}$ | $\approx \dfrac{D_2}{\epsilon^2}$ |
| | $a \in (\frac{1}{2}, 1)$ | $\dfrac{D_1}{K^{1-a}} + \dfrac{2aD_2}{(2a-1)K^{1-a}b}$ | $\dfrac{2a^2 D_2}{(1-a)(2a-1)D_1}$ |

### 1.3.2. Critical batch size to reduce SFO complexity

Section 1.3.1 showed that using large batches is appropriate for SGD in the sense of minimizing the upper bound of the performance measure. Here, we are interested in finding appropriate batch sizes from the viewpoint of the computation cost. This is because the SFO complexity increases with the batch size. As indicated in Section 1.2, the critical batch size $b^\star$ minimizes the SFO complexity, $N = Kb$. Hence, we will investigate the properties of the SFO complexity $N = Kb$ needed to achieve an $\epsilon$–approximation. Here, let us consider SGD using a constant learning rate. From the "Upper Bound" row in Table 1, we have

$$\min_{k \in [0:K-1]} \mathbb{E}[\|\nabla f(\boldsymbol{\theta}_k)\|^2] \leq \underbrace{\frac{C_1}{K} + \frac{C_2}{b} \leq \epsilon^2}_{\Leftrightarrow K \geq K(b) := \frac{C_1 b}{\epsilon^2 b - C_2} \ \left(b > \frac{C_2}{\epsilon^2}\right)} \quad .$$

We can check that the number of iterations, $K(b) := \frac{C_1 b}{\epsilon^2 b - C_2}$, needed to achieve an $\epsilon$–approximation is monotone decreasing and convex with respect to the batch size (Theorem 3.2). Accordingly, we have that $K(b) \geq \inf\{K\colon \min_{k \in [0:K-1]} \mathbb{E}[\|\nabla f(\boldsymbol{\theta}_k)\|] \leq \epsilon\}$, where SGD using the batch size $b$ generates $(\boldsymbol{\theta}_k)_{k=0}^{K-1}$. Moreover, we find that the SFO complexity is $N(b) = K(b)b = \frac{C_1 b^2}{\epsilon^2 b - C_2}$. The convexity of $N(b) = \frac{C_1 b^2}{\epsilon^2 b - C_2}$ (Theorem 3.3) ensures that a critical batch size $b^\star = \frac{2C_2}{\epsilon^2}$ whereby $N'(b^\star) = 0$ exists such that $N(b)$ is minimized at $b^\star$ (see the "Critical Batch Size" row in Table 1). A similar discussion guarantees the existence of a critical batch size for SGD using a decaying learning rate $\alpha_k = \frac{1}{\left(\lfloor \frac{k}{T} \rfloor + 1\right)^a}$, where $T \geq 1$, $a \in (0, \frac{1}{2}) \cup (\frac{1}{2}, 1)$, and $\lfloor \cdot \rfloor$ is the floor function (see the "Critical Batch Size" row in Table 1).

Meanwhile, for a decaying learning rate $\alpha_k = \frac{1}{\sqrt{\lfloor \frac{k}{T} \rfloor + 1}}$, although $N(b)$ is convex with respect to $b$, we have that $N'(b) > 0$ for all $b > \frac{D_2}{\epsilon^2}$ (Theorem 3.3(iii)). Hence, for

3

this case, a critical batch size $b^\star$ defined by $N'(b^\star) = 0$ does not exist. However, since the critical batch size minimizes the SFO complexity $N$, we can define one as follows: $b^\star \approx \frac{D_2}{\epsilon^2}$. Accordingly, we have that $N(b^\star) \geq \inf\{Kb\colon \min_{k\in[0:K-1]} \mathbb{E}[\|\nabla f(\boldsymbol{\theta}_k)\|] \leq \epsilon\}$, where SGD using $b^\star$ generates $(\boldsymbol{\theta}_k)_{k=0}^{K-1}$.

### 1.3.3. Iteration and SFO complexities

Let $\mathcal{F}(n, \Delta_0, L)$ be an $L$–smooth function class with $f := \frac{1}{n}\sum_{i=1}^n f_i$ and $f(\boldsymbol{\theta}_0) - f_\star \leq \Delta_0$ (see (C1)) and let $\mathcal{O}(b, \sigma^2)$ be a stochastic first-order oracle class (see (C2) and (C3)). The iteration complexity $\mathcal{K}_\epsilon$ [21, (7)] and SFO complexity $\mathcal{N}_\epsilon$ needed for SGD to be an $\epsilon$–approximation are defined as

$$
\mathcal{K}_\epsilon(n, b, \alpha_k, \Delta_0, L, \sigma^2) := \sup_{\mathsf{O}\in\mathcal{O}(b,\sigma^2)} \sup_{f\in\mathcal{F}(n,\Delta_0,L)} \inf\left\{K\colon \min_{k\in[0:K-1]} \mathbb{E}[\|\nabla f(\boldsymbol{\theta}_k)\|] \leq \epsilon\right\},
$$
(1)

$$
\mathcal{N}_\epsilon(n, b, \alpha_k, \Delta_0, L, \sigma^2) := \mathcal{K}_\epsilon(n, b, \alpha_k, \Delta_0, L, \sigma^2)b.
$$
(2)

Table 2 summarizes the iteration and SFO complexities (see also Theorem 3.4). Corollaries 6 and 7 in [14] are the same as our results for SGD with a constant learning rate in Theorems 3.1 and 3.3, since the randomized stochastic projected gradient free algorithm in [14] which is a stochastic zeroth-order (SZO) method that coincides with SGD and it can be applied to the situations where only noisy function values are available. In particular, Corollary 6 in [14] gave the convergence rate of the SZO methods using a fixed batch size, and Corollary 7 indicated the SZO complexity of the SZO method is the same as the SFO complexity. Hence, Corollaries 6 and 7 in [14] lead to the finding that the iteration complexity of SGD using a constant learning rate is $O(1/\epsilon^2)$ and the SFO complexity of SGD using a constant learning rate is $O(1/\epsilon^4)$.

Since the positive constants $C_i$ and $D_i$ depend on the learning rate, we need to compare numerically the performance of SGD with a constant learning rate with that of SGD with a decaying learning rate. Moreover, we also need to compare SGD with the existing first-order optimizers in order to verify its usefulness. Section 4 presents numerical comparisons showing that SGD using the critical batch size outperforms the existing first-order optimizers. We also show that the measured critical batch sizes are close to the theoretical sizes.

## 2. Nonconvex Optimization and SGD

### 2.1. Nonconvex optimization in deep learning

Let $\mathbb{R}^d$ be a $d$-dimensional Euclidean space with inner product $\langle \boldsymbol{x}, \boldsymbol{y} \rangle := \boldsymbol{x}^\top \boldsymbol{y}$ inducing the norm $\|\boldsymbol{x}\|$ and $\mathbb{N}$ be the set of nonnegative integers. Define $[0:n] := \{0, 1, \ldots, n\}$ for $n \geq 1$. Let $(x_k)_{k\in\mathbb{N}}$ and $(y_k)_{k\in\mathbb{N}}$ be positive real sequences and let $x(\epsilon), y(\epsilon) > 0$, where $\epsilon > 0$. $O$ denotes Landau's symbol; i.e., $y_k = O(x_k)$ if there exist $c > 0$ and $k_0 \in \mathbb{N}$ such that $y_k \leq cx_k$ for all $k \geq k_0$, and $y(\epsilon) = O(x(\epsilon))$ if there exists $c > 0$ such that $y(\epsilon) \leq cx(\epsilon)$. Given a parameter $\boldsymbol{\theta} \in \mathbb{R}^d$ and a data point $z$ in a data domain $Z$, a machine-learning model provides a prediction whose quality can be measured in terms of a differentiable nonconvex loss function $\ell(\boldsymbol{\theta}; z)$. We aim to minimize the empirical loss defined for all $\boldsymbol{\theta} \in \mathbb{R}^d$ by $f(\boldsymbol{\theta}) = \frac{1}{n}\sum_{i=1}^n \ell(\boldsymbol{\theta}; z_i) = \frac{1}{n}\sum_{i=1}^n f_i(\boldsymbol{\theta})$, where $S = (z_1, z_2, \ldots, z_n)$ denotes the training set (We assume that the number of training

**Table 2.** Iteration and SFO complexities needed for SGD using a constant or decaying learning rate to be an $\epsilon$–approximation (The critical batch sizes are used to compute $\mathcal{K}_\epsilon$ and $\mathcal{N}_\epsilon$)

| Learning Rate | | Iteration Complexity $\mathcal{K}_\epsilon$ | SFO Complexity $\mathcal{N}_\epsilon(n, b, \alpha_k, \Delta_0, L, \sigma^2)$ |
|---|---|---|---|
| Constant $\alpha \in (0, \frac{2}{L})$ | | $O\left(\dfrac{1}{\epsilon^2}\right) = \sup_{f,\mathsf{O}} K(b^\star)$ | $O\left(\dfrac{1}{\epsilon^4}\right) = \sup_{f,\mathsf{O}} \dfrac{4C_1 C_2}{\epsilon^4}$ |
| Decay $\alpha_k = \frac{1}{\left(\lfloor \frac{k}{T} \rfloor + 1\right)^a}$ | $a \in (0, \frac{1}{2})$ | $O\left(\dfrac{1}{\epsilon^{\frac{2}{a}}}\right) = \sup_{f,\mathsf{O}} K(b^\star)$ | $O\left(\dfrac{1}{\epsilon^{\frac{2}{a}}}\right) = \sup_{f,\mathsf{O}} \dfrac{(1-a)^{1-\frac{1}{a}} D_2}{a(1-2a)D_1^{1-\frac{1}{a}} \epsilon^{\frac{2}{a}}}$ |
| | $a = \frac{1}{2}$ | $O\left(\dfrac{1}{\epsilon^4}\right) = \sup_{f,\mathsf{O}} K(b^\star)$ | $O\left(\dfrac{1}{\epsilon^6}\right) = \sup_{f,\mathsf{O}} \left(\dfrac{D_1(D_2+1)}{\epsilon^2} + D_2\right)^2 \dfrac{D_1+1}{\epsilon^2}$ |
| | $a \in (\frac{1}{2}, 1)$ | $O\left(\dfrac{1}{\epsilon^{\frac{2}{1-a}}}\right) = \sup_{f,\mathsf{O}} K(b^\star)$ | $O\left(\dfrac{1}{\epsilon^{\frac{2}{1-a}}}\right) = \sup_{f,\mathsf{O}} \dfrac{2a^{2-\frac{1}{1-a}}(1-a)^{-1} D_2}{(2a-1)D_1^{1-\frac{1}{1-a}} \epsilon^{\frac{2}{1-a}}}$ |

data $n$ is large) and $f_i(\cdot) := \ell(\cdot; z_i)$ denotes the loss function corresponding to the $i$-th training data $z_i$.

## 2.2. SGD

### 2.2.1. Conditions and algorithm

We assume that a stochastic first-order oracle (SFO) exists such that, for a given $\boldsymbol{\theta} \in \mathbb{R}^d$, it returns a stochastic gradient $\mathsf{G}_\xi(\boldsymbol{\theta})$ of the function $f$, where a random variable $\xi$ is independent of $\boldsymbol{\theta}$. Let $\mathbb{E}_\xi[\cdot]$ be the expectation taken with respect to $\xi$. The following are standard conditions.

(C1) $f := \frac{1}{n}\sum_{i=1}^n f_i \colon \mathbb{R}^d \to \mathbb{R}$ is $L$–smooth, i.e., $\nabla f \colon \mathbb{R}^d \to \mathbb{R}^d$ is $L$–Lipschitz continuous (i.e., $\|\nabla f(\boldsymbol{x}) - \nabla f(\boldsymbol{y})\| \leq L\|\boldsymbol{x} - \boldsymbol{y}\|$). $f$ is bounded below from $f_\star \in \mathbb{R}$. Let $\Delta_0 > 0$ satisfy $f(\boldsymbol{\theta}_0) - f_\star \leq \Delta_0$, where $\boldsymbol{\theta}_0$ is an initial point.

(C2) Let $(\boldsymbol{\theta}_k)_{k\in\mathbb{N}} \subset \mathbb{R}^d$ be the sequence generated by SGD. For each iteration $k$, $\mathbb{E}_{\xi_k}[\mathsf{G}_{\xi_k}(\boldsymbol{\theta}_k)] = \nabla f(\boldsymbol{\theta}_k)$, where $\xi_0, \xi_1, \ldots$ are independent samples and the random variable $\xi_k$ is independent of $(\boldsymbol{\theta}_l)_{l=0}^k$. There exists a nonnegative constant $\sigma^2$ such that $\mathbb{E}_{\xi_k}[\|\mathsf{G}_{\xi_k}(\boldsymbol{\theta}_k) - \nabla f(\boldsymbol{\theta}_k)\|^2] \leq \sigma^2$.

(C3) For each iteration $k$, SGD samples a batch $B_k$ of size $b$ independently of $k$ and estimates the full gradient $\nabla f$ as $\nabla f_{B_k}(\boldsymbol{\theta}_k) := \frac{1}{b}\sum_{i\in[b]} \mathsf{G}_{\xi_{k,i}}(\boldsymbol{\theta}_k)$, where $b \leq n$ and $\xi_{k,i}$ is a random variable generated by the $i$-th sampling in the $k$-th iteration.

Algorithm 1 is the SGD optimizer under (C1)–(C3).

### 2.2.2. Learning rates

We use the following learning rates:

**(Constant rate)** $\alpha_k$ does not depend on $k \in \mathbb{N}$, i.e., $\alpha_k = \alpha < \frac{2}{L}$ ($k \in \mathbb{N}$), where the upper bound $\frac{2}{L}$ of $\alpha$ is needed to analyze SGD (see Appendix A.2).

**(Decaying rate)** $(\alpha_k)_{k\in\mathbb{N}} \subset (0, +\infty)$ is monotone decreasing for $k$ (i.e., $\alpha_k \geq \alpha_{k+1}$) and converges to 0. In particular, we will use $\alpha_k = \frac{1}{\left(\lfloor \frac{k}{T} \rfloor + 1\right)^a}$, where

---

**Algorithm 1** SGD

---

**Require:** $\alpha_k \in (0, +\infty)$ (learning rate), $b \geq 1$ (batch size), $K \geq 1$ (iteration)
**Ensure:** $\boldsymbol{\theta}_K$
 1: $k \leftarrow 0$, $\boldsymbol{\theta}_0 \in \mathbb{R}^d$
 2: **loop**
 3: $\quad \nabla f_{B_k}(\boldsymbol{\theta}_k) := \frac{1}{b} \sum_{i \in [b]} \mathsf{G}_{\xi_{k,i}}(\boldsymbol{\theta}_k)$
 4: $\quad \boldsymbol{\theta}_{k+1} := \boldsymbol{\theta}_k - \alpha_k \nabla f_{B_k}(\boldsymbol{\theta}_k)$
 5: $\quad k \leftarrow k + 1$
 6: **end loop**

---

**(Decay 1)** $a \in (0, \frac{1}{2}) \vee$ **(Decay 2)** $a = \frac{1}{2} \vee$ **(Decay 3)** $a \in (\frac{1}{2}, 1)$. It is guaranteed that there exists $k_0 \in \mathbb{N}$ such that, for all $k \geq k_0$, $\alpha_k < \frac{2}{L}$. Furthermore, we assume that $k_0 = 0$, since we can replace $\alpha_k$ with $\frac{\alpha}{(\lfloor \frac{k}{T} \rfloor + 1)^a} \leq \alpha < \frac{2}{L}$ ($k \in \mathbb{N}$), where $\alpha \in (0, \frac{2}{L})$ is defined as in **(Constant)**.

## 3. Our Results

### 3.1. Upper bound of the squared norm of the full gradient

Here, we give an upper bound of $\min_{k \in [0:K-1]} \mathbb{E}[\|\nabla f(\boldsymbol{\theta}_k)\|^2]$, where $\mathbb{E}[\cdot]$ stands for the total expectation, for the sequence generated by SGD using each of the learning rates defined in Section 2.2.2.

**Theorem 3.1** (Upper bound of the squared norm of the full gradient). *The sequence* $(\boldsymbol{\theta}_k)_{k \in \mathbb{N}}$ *generated by Algorithm 1 under (C1)–(C3) satisfies that, for all $K \geq 1$,*

$$
\min_{k \in [0:K-1]} \mathbb{E}\left[\|\nabla f(\boldsymbol{\theta}_k)\|^2\right] \leq
\begin{cases}
\dfrac{C_1}{K} + \dfrac{C_2}{b} & \textbf{\textit{(Constant)}} \\[2ex]
\dfrac{D_1}{K^a} + \dfrac{D_2}{(1-2a)K^a b} & \textbf{\textit{(Decay 1)}} \\[2ex]
\dfrac{D_1}{\sqrt{K}} + \left(\dfrac{1}{\sqrt{K}} + 1\right)\dfrac{D_2}{b} & \textbf{\textit{(Decay 2)}} \\[2ex]
\dfrac{D_1}{K^{1-a}} + \dfrac{2aD_2}{(2a-1)K^{1-a}b} & \textbf{\textit{(Decay 3)}}
\end{cases}
$$

*where*

$$
C_1 := \frac{2(f(\boldsymbol{\theta}_0) - f_\star)}{(2 - L\alpha)\alpha}, \; C_2 := \frac{L\sigma^2 \alpha}{2 - L\alpha},
$$
$$
D_1 := \frac{2(f(\boldsymbol{\theta}_0) - f_\star)}{\alpha(2 - L\alpha)}, \; D_2 := \frac{T\alpha^2 L\sigma^2}{2 - L\alpha}.
$$

Theorem 3.1 indicates that the upper bound of $\min_{k \in [0:K-1]} \mathbb{E}[\|\nabla f(\boldsymbol{\theta}_k)\|^2]$ consists of a bias term including $f(\boldsymbol{\theta}_0) - f_\star$ and a variance term including $\sigma^2$ and that these terms become small when the number of iterations and the batch size are large. In particular, the bias term using **(Constant)** is $O(\frac{1}{K})$, which is a better rate than using **(Decay 1)**–**(Decay 3)**.

### 3.2. Number of iterations needed for SGD to be an $\epsilon$-approximation

Let us suppose that SGD is an $\epsilon$-approximation defined as follows:

$$\mathbb{E}\left[\|\nabla f(\boldsymbol{\theta}_{K^*})\|^2\right] := \min_{k \in [0:K-1]} \mathbb{E}\left[\|\nabla f(\boldsymbol{\theta}_k)\|^2\right] \leq \epsilon^2, \tag{3}$$

where $\epsilon > 0$ is the precision and $K^* \in [0 : K - 1]$. Condition (3) implies that $\mathbb{E}[\|\nabla f(\boldsymbol{\theta}_{K^*})\|] \leq \epsilon$. Theorem 3.1 below gives the number of iterations needed to be an $\epsilon$-approximation (3).

**Theorem 3.2** (Numbers of iterations needed for nonconvex optimization using SGD)**.** *Let $(\boldsymbol{\theta}_k)_{k \in \mathbb{N}}$ be the sequence generated by Algorithm 1 under (C1)–(C3) and let $K \colon \mathbb{R} \to \mathbb{R}$ be*

$$K(b) = \begin{cases} \dfrac{C_1 b}{\epsilon^2 b - C_2} & \textbf{\textit{(Constant)}} \\[3mm] \left\{ \dfrac{1}{\epsilon^2}\left(\dfrac{D_2}{(1-2a)b} + D_1\right)\right\}^{\frac{1}{a}} & \textbf{\textit{(Decay 1)}} \\[3mm] \left(\dfrac{D_1 b + D_2}{\epsilon^2 b - D_2}\right)^2 & \textbf{\textit{(Decay 2)}} \\[3mm] \left\{ \dfrac{1}{\epsilon^2}\left(\dfrac{2aD_2}{(2a-1)b} + D_1\right)\right\}^{\frac{1}{1-a}} & \textbf{\textit{(Decay 3)}} \end{cases}$$

*where $C_1$, $C_2$, $D_1$ $(> \epsilon^2)$, and $D_2$ are defined as in Theorem 3.1, the domain of $K$ in* **(Constant)** *is $b > \frac{C_2}{\epsilon^2}$, and the domain of $K$ in* **(Decay 2)** *is $b > \frac{D_2}{\epsilon^2}$. Then, we have the following:*

(i) *The above $K$ achieves an $\epsilon$-approximation (3).*
(ii) *The above $K$ is a monotone decreasing and convex function with respect to the batch size $b$.*

Theorem 3.2 indicates that the number of iterations needed for SGD using constant or decay learning rates to be an $\epsilon$-approximation is small when the batch size is large. Hence, it is appropriate to set a large batch size in order to minimize the iterations needed for an $\epsilon$-approximation (3). However, the SFO complexity, which is the cost of the stochastic gradient computation, grows larger with $b$. Hence, the appropriate batch size should also minimize the SFO complexity.

### 3.3. SFO complexity needed for SGD to be an $\epsilon$-approximation

Theorem 3.2 leads to the following theorem on the properties of the SFO complexity $N$ needed for SGD to be an $\epsilon$-approximation (3).

**Theorem 3.3** (SFO complexity needed for nonconvex optimization of SGD)**.** *Let $(\boldsymbol{\theta}_k)_{k \in \mathbb{N}}$ be the sequence generated by Algorithm 1 under (C1)–(C3) and define $N \colon \mathbb{R} \to$*

$\mathbb{R}$ *by*

$$
N(b) = K(b)b =
\begin{cases}
\dfrac{C_1 b^2}{\epsilon^2 b - C_2} & \textbf{\textit{(Constant)}} \\[2ex]
\left\{ \dfrac{1}{\epsilon^2} \left( \dfrac{D_2}{(1-2a)b} + D_1 \right) \right\}^{\frac{1}{a}} b & \textbf{\textit{(Decay 1)}} \\[2ex]
\left( \dfrac{D_1 b + D_2}{\epsilon^2 b - D_2} \right)^2 b & \textbf{\textit{(Decay 2)}} \\[2ex]
\left\{ \dfrac{1}{\epsilon^2} \left( \dfrac{2aD_2}{(2a-1)b} + D_1 \right) \right\}^{\frac{1}{1-a}} b & \textbf{\textit{(Decay 3)}}
\end{cases}
$$

*where $C_1$, $C_2$, $D_1$ and $D_2$ are as in Theorem 3.1, the domain of $N$ in* **(Constant)** *is $b > \frac{C_2}{\epsilon^2}$, and the domain of $N$ in* **(Decay 2)** *is $b > \frac{D_2}{\epsilon^2}$. Then, we have the following:*

(i) *The above $N$ is convex with respect to the batch size $b$.*

(ii) *There exists a critical batch size*

$$
b^\star =
\begin{cases}
\dfrac{2C_2}{\epsilon^2} & \textbf{\textit{(Constant)}} \\[2ex]
\dfrac{(1-a)D_2}{a(1-2a)D_1} & \textbf{\textit{(Decay 1)}} \\[2ex]
\dfrac{2a^2 D_2}{(1-a)(2a-1)D_1} & \textbf{\textit{(Decay 3)}}
\end{cases}
\tag{4}
$$

*satisfying $N'(b^\star) = 0$ such that $b^\star$ minimizes the SFO complexity $N$.*

(iii) *For* **(Decay 2)**, *$N'(b) > 0$ holds for all $b > \dfrac{D_2}{\epsilon^2}$.*

Theorem 3.3(ii) indicates that, if we can set a critical batch size (4) for each of **(Constant)**, **(Decay 1)**, and **(Decay 3)**, then the SFO complexity will be minimized. However, it would be difficult to set $b^\star$ in (4) before implementing SGD, since $b^\star$ in (4) involves unknown parameters, such as $L$ and $\sigma^2$ (computing $L$ is NP-hard [28]). Hence, we would like to estimate the critical batch sizes by using Theorem 3.3(ii) and (iii) (see Section 4.3). Theorem 3.3(ii) indicates that the smaller $\epsilon$ is, the larger the critical batch size $b^\star$ in **(Constant)** becomes. Theorem 3.3(iii) indicates that the critical batch size is close to $\frac{D_2}{\epsilon^2}$ when using **(Decay 2)** to minimize the SFO complexity $N$.

### 3.4. Iteration and SFO complexities of SGD

Theorems 3.2 and 3.3 lead to the following theorem indicating the iteration and SFO complexities needed for SGD to be an $\epsilon$–approximation (see also Table 2).

**Theorem 3.4** (Iteration and SFO complexities of SGD)**.** *The iteration and SFO complexities such that Algorithm 1 under (C1)–(C3) is an $\epsilon$–approximation (3) are as*

*follows:*

$$(\mathcal{K}_\epsilon(n, b^\star, \alpha_k, \Delta_0, L, \sigma^2), \mathcal{N}_\epsilon(n, b^\star, \alpha_k, \Delta_0, L, \sigma^2)) = \begin{cases} \left( O\left(\dfrac{1}{\epsilon^2}\right), O\left(\dfrac{1}{\epsilon^4}\right) \right) & \textbf{(Constant)} \\ \left( O\left(\dfrac{1}{\epsilon^{\frac{2}{a}}}\right), O\left(\dfrac{1}{\epsilon^{\frac{2}{a}}}\right) \right) & \textbf{(Decay 1)} \\ \left( O\left(\dfrac{1}{\epsilon^4}\right), O\left(\dfrac{1}{\epsilon^6}\right) \right) & \textbf{(Decay 2)} \\ \left( O\left(\dfrac{1}{\epsilon^{\frac{2}{1-a}}}\right), O\left(\dfrac{1}{\epsilon^{\frac{2}{1-a}}}\right) \right) & \textbf{(Decay 3)} \end{cases}$$

*where $\mathcal{K}_\epsilon(n, b, \alpha_k, \Delta_0, L, \sigma^2)$ and $\mathcal{N}_\epsilon(n, b, \alpha_k, \Delta_0, L, \sigma^2)$ are defined as in (1), the critical batch sizes in Theorem 3.3 are used to compute $\mathcal{K}_\epsilon(n, b^\star, \alpha_k, \Delta_0, L, \sigma^2)$ and $\mathcal{N}_\epsilon(n, b^\star, \alpha_k, \Delta_0, L, \sigma^2)$. In* **(Decay 2)**, *we assume that $b^\star = \frac{D_2 + 1}{\epsilon^2}$. (see also (4)).*

Theorem 3.4 indicates that the iteration and SFO complexities for **(Constant)** are smaller than those for **(Decay 1)**–**(Decay 3)**.

## 4. Numerical Results

We numerically verified the number of iterations and SFO complexities needed to achieve high test accuracy for different batch sizes in training ResNet [29] and Wide-ResNet [30]. The parameter $\alpha$ used in **(Constant)** was determined by conducting a grid search of $\{0.001, 0.005, 0.01, 0.05, 0.1, 0.5\}$. The parameters $\alpha$ and $T$ used in the decaying learning rates **(Decay 1)**–**(Decay 3)** defined by $\alpha_k = \frac{\alpha}{\left(\lfloor \frac{k}{T} \rfloor + 1\right)^a}$ were determined by a grid search of $\alpha \in \{0.001, 0.1, 0.125, 0.25, 0.5, 1.0\}$ and $T \in \{5, 10, 20, 30, 40, 50\}$. The parameter $a$ was set to $a = \frac{1}{4}$ in **(Decay 1)** and $a = \frac{3}{4}$ in **(Decay 3)**. We compared SGD with SGD with momentum (momentum), Adam, AdamW, and RMSProp. The learning rates and hyperparameters of these four optimizers were determined on the basis of the previous results [9,10,12] (The weight decay used in the momentum was $5 \times 10^{-4}$). The experimental environment consisted of an NVIDIA DGX A100×8GPU and Dual AMD Rome7742 2.25-GHz, 128 Cores×2CPU. The software environment was Python 3.10.6, PyTorch 1.13.1, and CUDA 11.6. The code is available at `https://github.com/imakn0907/SGD_using_decaying`.

### 4.1. Training ResNet-18 on the CIFAR-10 and CIFAR-100 datasets

First, we trained ResNet-18 on the CIFAR-10 dataset. The stopping condition of the optimizers was 200 epochs. Figure 1 indicates that the number of iterations is monotone decreasing and convex with respect to batch size for SGDs using a constant learning rate or a decaying learning rate. Figure 2 indicates that, in each case of SGD with **(Constant)**–**(Decay 3)**, a critical batch size $b^\star = 2^4$ exists at which the SFO complexity is minimized.

Figures 3 and 4 indicate that the number of iterations and the SFO complexity for four different learning rates in achieving a test accuracy of 0.6 when training ResNet-18 on the CIFAR-100 dataset. The figures indicate that critical batch sizes existed when using **(Constant)**–**(Decay 3)**.

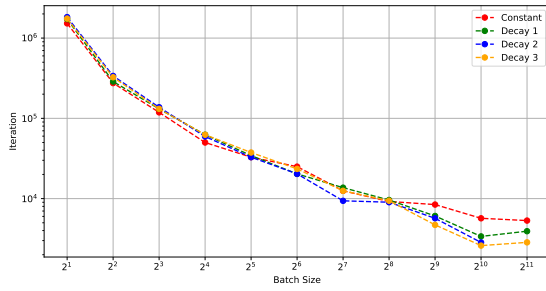Figures 5 and 6 compare SGD with **(Decay 1)** with the other optimizers in training

**Figure 1.** Number of iterations needed for SGD with (Constant), (Decay 1), (Decay 2), and (Decay 3) to achieve a test accuracy of 0.9 versus batch size (ResNet-18 on CIFAR-10)
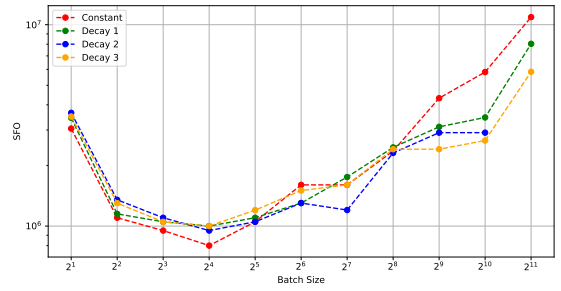


**Figure 2.** SFO complexity needed for SGD with (Constant), (Decay 1), (Decay 2), and (Decay 3) to achieve a test accuracy of 0.9 versus batch size (ResNet-18 on CIFAR-10)
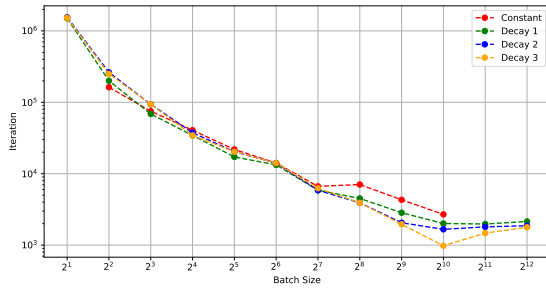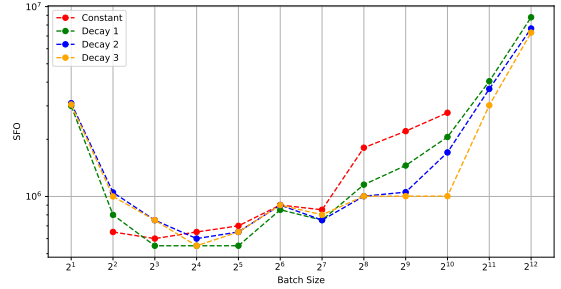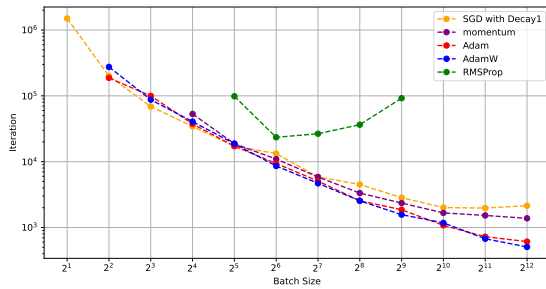


**Figure 3.** Number of iterations needed for SGD with (Constant), (Decay 1), (Decay 2), and (Decay 3) to achieve a test accuracy of 0.6 versus batch size (ResNet-18 on CIFAR-100)



**Figure 4.** SFO complexity needed for SGD with (Constant), (Decay 1), (Decay 2), and (Decay 3) to achieve a test accuracy of 0.6 versus batch size (ResNet-18 on CIFAR-100)



**Figure 5.** Number of iterations needed for SGD with (Decay 1), momentum, Adam, AdamW, and RMSProp to achieve a test accuracy of 0.6 versus batch size (ResNet-18 on CIFAR-100)
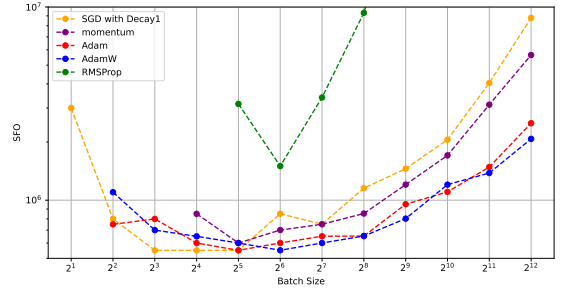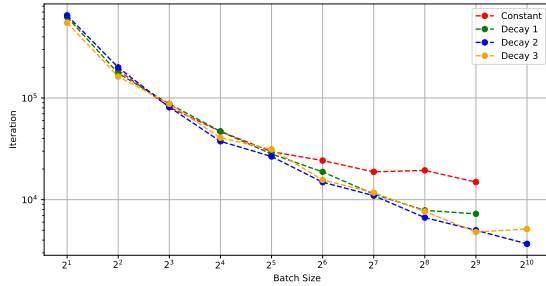


**Figure 6.** SFO complexity needed for SGD with (Decay 1), momentum, Adam, AdamW, and RMSProp to achieve a test accuracy of 0.6 versus batch size (ResNet-18 on CIFAR-100)

10

ResNet-18 on the CIFAR-100 dataset. These figures indicates that SGD with **(Decay 1)** and a critical batch size ($b = 2^4$) outperformed the other optimizers in the sense of minimizing the number of iterations and the SFO complexity. Figure 6 also indicates that the existing optimizers using constant learning rates had critical batch sizes minimizing the SFO complexities. In particular, AdamW using the critical batch size $b^\star = 2^5$ performed well.

### 4.2. Training Wide-ResNet on the CIFAR-10 and CIFAR-100 datasets



**Figure 7.** Number of iterations needed for SGD with (Constant), (Decay 1), (Decay 2), and (Decay 3) to achieve a test accuracy of 0.9 versus batch size (Wide-ResNet-28-10 on CIFAR-10)
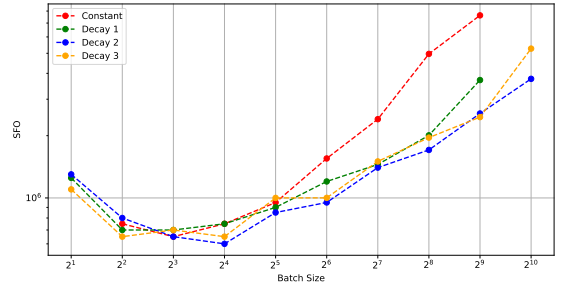
**Figure 8.** SFO complexity needed for SGD with (Constant), (Decay 1), (Decay 2), and (Decay 3) to achieve a test accuracy of 0.9 versus batch size (Wide-ResNet-28-10 on CIFAR-10)

Next, we trained Wide-ResNet-28 [30] on the CIFAR-10 and CIFAR-100 datasets. The stopping condition of the optimizers was 200 epochs. Figures 7 and 8 show that the number of iterations and the SFO complexity of SGD to achieve a test accuracy of 0.9 (CIFAR-10) versus batch size. Figures 7 and 8 indicate that the critical batch size was $b^\star = 2^4$ in each case of SGD using **(Constant)**–**(Decay 3)**.
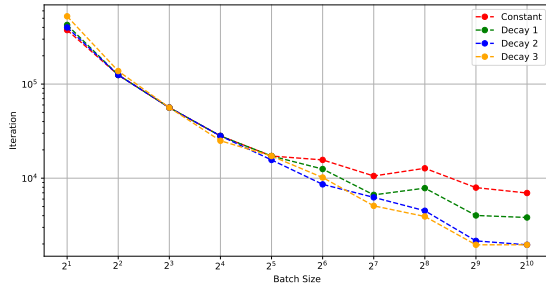


**Figure 9.** Number of iterations needed for SGD with (Constant), (Decay 1), (Decay 2), and (Decay 3) to achieve a test accuracy of 0.6 versus batch size (WideResNet-28-12 on CIFAR-100)

**Figure 10.** SFO complexity needed for SGD with (Constant), (Decay 1), (Decay 2), and (Decay 3) to achieve a test accuracy of 0.6 versus batch size (WideResNet-28-12 on CIFAR-100)
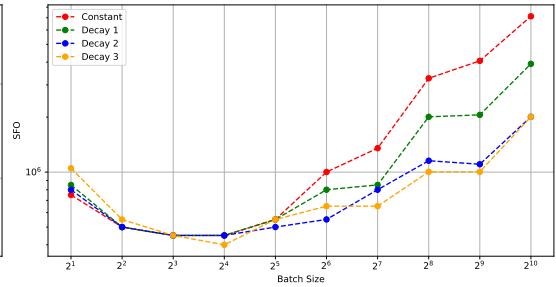
Figures 9 and 10 indicate that the number of iterations and the SFO complexity of SGD to achieve a test accuracy of 0.6 (CIFAR-100) versus batch size and show that a critical batch size existed for **(Constant)**–**(Decay 3)**.

As in Figures 5 and 6, Figures 11 and 12 indicate that SGD using **(Decay 3)** and the existing optimizers using constant learning rates had critical batch sizes minimizing the SFO complexities.
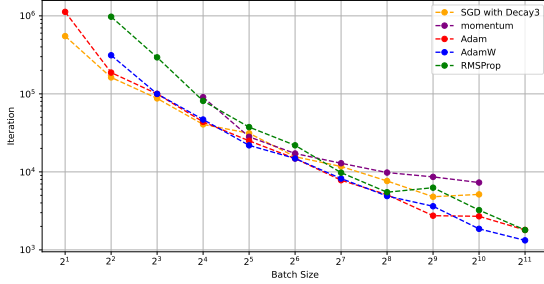
11

**Figure 11.** Number of iterations needed for SGD with (Decay 3), momentum, Adam, AdamW, and RMSProp to achieve a test accuracy of 0.9 versus batch size (WideResNet-28-10 on CIFAR-10)

**Figure 12.** SFO complexity needed for SGD with (Decay 3), momentum, Adam, AdamW, and RMSProp to achieve a test accuracy of 0.9 versus batch size (WideResNet-28-10 on CIFAR-10)
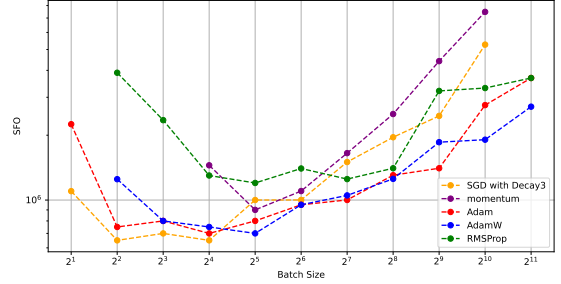
## *4.3. Estimation of critical batch sizes*

**Table 3.** Measured (left) and estimated (right; bold) critical batch sizes (D1 and D3 stand for (Decay 1) and (Decay 3))

|  | ResNet-18 | | Wide-ResNet-28 | |
|---|---|---|---|---|
|  | CIFAR10 | CIFAR100 | CIFAR10 | CIFAR100 |
| D1 | $2^4$ | $2^4$ | $2^2$ | $2^3, 2^4$ |
| D3 | $2^4$    **24** | $2^4$    **24** | $2^2, 2^4$    **6** | $2^4$    **12, 24** |

We estimated the critical batch sizes of **(Decay 3)** using Theorem 3.3 and measured the critical batch sizes of **(Decay 1)**. From (4) and $a = \frac{1}{4}$ (**(Decay 1)**), we have that, for training ResNet-18 on the CIFAR-10 dataset, $b^\star = 2^4 = \frac{(1-a)D_2}{a(1-2a)D_1}$, i.e., $\frac{D_2}{D_1} = \frac{8}{3}$. Then, the estimated critical batch size of SGD using **(Decay 3)** ($a = \frac{3}{4}$) for training ResNet-18 on the CIFAR-10 dataset is

$$b^\star = \frac{2a^2}{(1-a)(2a-1)} \frac{D_2}{D_1} = \frac{2a^2}{(1-a)(2a-1)} \frac{8}{3}$$
$$= 24 \in (2^4, 2^5),$$

which implies that the estimated critical batch size $b^\star = 24$ is close to the measured size $b = 2^4$. We also found that the estimated critical batch sizes are close to the measured critical batch sizes (see Table 3).

## 5. Conclusion and future work

This paper investigated the required number of iterations and SFO complexities for SGD using constant or decay learning rates to achieve an $\epsilon$–approximation. Our theoretical analyses indicated that the number of iterations needed for an $\epsilon$–approximation is monotone decreasing and convex with respect to the batch size and the SFO complexity needed for an $\epsilon$–approximation is convex with respect to the batch size. Moreover, we showed that SGD using a critical batch size reduces the SFO complexity.

12

The numerical results indicated that SGD using the critical batch size performs better than the existing optimizers in the sense of minimizing the SFO complexity. We also estimated critical batch sizes of SGD using our theoretical results and showed that they are close to the measured critical batch sizes.

The results in this paper can be only applied to SGD. This is a limitation of our work. Hence, in the future, we should investigate whether our results can be applied to variants of SGD, such as the momentum methods and adaptive methods.

(3669 words)

## References

[1] Robbins H, Monro H. A stochastic approximation method. The Annals of Mathematical Statistics. 1951;22:400–407.

[2] Zinkevich M. Online convex programming and generalized infinitesimal gradient ascent. In: Proceedings of the 20th International Conference on Machine Learning; 2003. p. 928–936.

[3] Nemirovski A, Juditsky A, Lan G, et al. Robust stochastic approximation approach to stochastic programming. SIAM Journal on Optimization. 2009;19:1574–1609.

[4] Ghadimi S, Lan G. Optimal stochastic approximation algorithms for strongly convex stochastic composite optimization I: A generic algorithmic framework. SIAM Journal on Optimization. 2012;22:1469–1492.

[5] Ghadimi S, Lan G. Optimal stochastic approximation algorithms for strongly convex stochastic composite optimization II: Shrinking procedures and optimal algorithms. SIAM Journal on Optimization. 2013;23:2061–2089.

[6] Polyak BT. Some methods of speeding up the convergence of iteration methods. USSR Computational Mathematics and Mathematical Physics. 1964;4:1–17.

[7] Nesterov Y. A method for unconstrained convex minimization problem with the rate of convergence $O(1/k^2)$. Doklady AN USSR. 1983;269:543–547.

[8] Duchi J, Hazan E, Singer Y. Adaptive subgradient methods for online learning and stochastic optimization. Journal of Machine Learning Research. 2011;12:2121–2159.

[9] Tieleman T, Hinton G. RMSProp: Divide the gradient by a running average of its recent magnitude. COURSERA: Neural networks for machine learning. 2012;4:26–31.

[10] Kingma DP, Ba J. Adam: A method for stochastic optimization. In: Proceedings of The International Conference on Learning Representations; 2015.

[11] Reddi SJ, Kale S, Kumar S. On the convergence of Adam and beyond. In: Proceedings of The International Conference on Learning Representations; 2018.

[12] Loshchilov I, Hutter F. Decoupled weight decay regularization. In: Proceedings of The International Conference on Learning Representations; 2019.

[13] Ghadimi S, Lan G. Stochastic first- and zeroth-order methods for nonconvex stochastic programming. SIAM Journal on Optimization. 2013;23(4):2341–2368. Available from: https://doi.org/10.1137/120880811.

[14] Ghadimi S, Lan G, Zhang H. Mini-batch stochastic approximation methods for nonconvex stochastic composite optimization. Mathematical Programming. 2016;155(1):267–305. Available from: https://doi.org/10.1007/s10107-014-0846-1.

[15] Vaswani S, Mishkin A, Laradji I, et al. Painless stochastic gradient: Interpolation, line-search, and convergence rates. In: Advances in Neural Information Processing Systems; Vol. 32; 2019.

[16] Fehrman B, Gess B, Jentzen A. Convergence rates for the stochastic gradient descent method for non-convex objective functions. Journal of Machine Learning Research. 2020; 21:1–48.

[17] Chen H, Zheng L, AL Kontar R, et al. Stochastic gradient descent in correlated settings: A study on Gaussian processes. In: Advances in Neural Information Processing Systems;

Vol. 33; 2020.

[18] 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.

[19] Loizou N, Vaswani S, Laradji I, et al. 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; Vol. 130; 2021.

[20] Wang X, Magnússon S, Johansson M. On the convergence of step decay step-size for stochastic optimization. In: Beygelzimer A, Dauphin Y, Liang P, et al., editors. Advances in Neural Information Processing Systems; 2021. Available from: `https://openreview.net/forum?id=M-W0asp3fD`.

[21] Arjevani Y, Carmon Y, Duchi JC, et al. Lower bounds for non-convex stochastic optimization. Mathematical Programming. 2023;199(1):165–214.

[22] Khaled A, Richtárik P. Better theory for SGD in the nonconvex world. Transactions on Machine Learning Research. 2023;.

[23] Jain P, Kakade SM, Kidambi R, et al. Parallelizing stochastic gradient descent for least squares regression: Mini-batching, averaging, and model misspecification. Journal of Machine Learning Research. 2018;18(223):1–42.

[24] Cotter A, Shamir O, Srebro N, et al. Better mini-batch algorithms via accelerated gradient methods. In: Advances in Neural Information Processing Systems; Vol. 24; 2011.

[25] Smith SL, Kindermans PJ, Le QV. Don't decay the learning rate, increase the batch size. In: International Conference on Learning Representations; 2018.

[26] Sato N, Iiduka H. Existence and estimation of critical batch size for training generative adversarial networks with two time-scale update rule. In: Proceedings of the 40th International Conference on Machine Learning; (Proceedings of Machine Learning Research; Vol. 202). PMLR; 2023. p. 30080–30104.

[27] Shallue CJ, Lee J, Antognini J, et al. Measuring the effects of data parallelism on neural network training. Journal of Machine Learning Research. 2019;20:1–49.

[28] Virmaux A, Scaman K. Lipschitz regularity of deep neural networks: analysis and efficient estimation. In: Advances in Neural Information Processing Systems; Vol. 31; 2018.

[29] He K, Zhang X, Ren S, et al. Deep residual learning for image recognition. In: Proceedings of the IEEE conference on computer vision and pattern recognition; 2016. p. 770–778.

[30] Zagoruyko S, Komodakis N. Wide residual networks. arXiv preprint arXiv:160507146. 2016;.

## Appendix A. Appendix

### A.1. Lemma

First, we will prove the following lemma.

**Lemma A.1.** *The sequence $(\boldsymbol{\theta}_k)_{k \in \mathbb{N}}$ generated by Algorithm 1 under (C1)–(C3) satisfies that, for all $K \geq 1$,*

$$\sum_{k=0}^{K-1} \alpha_k \left(1 - \frac{L\alpha_k}{2}\right) \mathbb{E}\left[\|\nabla f(\boldsymbol{\theta}_k)\|^2\right] \leq \mathbb{E}\left[f(\boldsymbol{\theta}_0) - f_\star\right] + \frac{L\sigma^2}{2b} \sum_{k=0}^{K-1} \alpha_k^2,$$

*where $\mathbb{E}$ stands for the total expectation.*

*Proof:* Condition (C1) (*L*-smoothness of $f$) implies that the descent lemma holds, i.e., for all $k \in \mathbb{N}$,

$$f(\boldsymbol{\theta}_{k+1}) \leq f(\boldsymbol{\theta}_k) + \langle \nabla f(\boldsymbol{\theta}_k), \boldsymbol{\theta}_{k+1} - \boldsymbol{\theta}_k \rangle + \frac{L}{2}\|\boldsymbol{\theta}_{k+1} - \boldsymbol{\theta}_k\|^2,$$

which, together with $\boldsymbol{\theta}_{k+1} := \boldsymbol{\theta}_k - \alpha_k \nabla f_{B_k}(\boldsymbol{\theta}_k)$, implies that

$$f(\boldsymbol{\theta}_{k+1}) \leq f(\boldsymbol{\theta}_k) - \alpha_k \langle \nabla f(\boldsymbol{\theta}_k), \nabla f_{B_k}(\boldsymbol{\theta}_k) \rangle + \frac{L\alpha_k^2}{2}\|\nabla f_{B_k}(\boldsymbol{\theta}_k)\|^2. \qquad \text{(A1)}$$

Condition (C2) guarantees that

$$\mathbb{E}_{\xi_k}\left[\nabla f_{B_k}(\boldsymbol{\theta}_k)|\boldsymbol{\theta}_k\right] = \nabla f(\boldsymbol{\theta}_k) \text{ and } \mathbb{E}_{\xi_k}\left[\|\nabla f_{B_k}(\boldsymbol{\theta}_k) - \nabla f(\boldsymbol{\theta}_k)\|^2|\boldsymbol{\theta}_k\right] \leq \frac{\sigma^2}{b}. \qquad \text{(A2)}$$

Hence, we have

$$\begin{aligned}
\mathbb{E}_{\xi_k}\left[\|\nabla f_{B_k}(\boldsymbol{\theta}_k)\|^2|\boldsymbol{\theta}_k\right] &= \mathbb{E}_{\xi_k}\left[\|\nabla f_{B_k}(\boldsymbol{\theta}_k) - \nabla f(\boldsymbol{\theta}_k) + \nabla f(\boldsymbol{\theta}_k)\|^2|\boldsymbol{\theta}_k\right] \\
&= \mathbb{E}_{\xi_k}\left[\|\nabla f_{B_k}(\boldsymbol{\theta}_k) - \nabla f(\boldsymbol{\theta}_k)\|^2|\boldsymbol{\theta}_k\right] + 2\mathbb{E}_{\xi_k}\left[\langle \nabla f_{B_k}(\boldsymbol{\theta}_k) - \nabla f(\boldsymbol{\theta}_k), \nabla f(\boldsymbol{\theta}_k) \rangle|\boldsymbol{\theta}_k\right] \\
&\quad + \mathbb{E}_{\xi_k}\left[\|\nabla f(\boldsymbol{\theta}_k)\|^2|\boldsymbol{\theta}_k\right] \\
&\leq \frac{\sigma^2}{b} + \mathbb{E}_{\xi_k}\left[\|\nabla f(\boldsymbol{\theta}_k)\|^2\right]. \qquad \text{(A3)}
\end{aligned}$$

Taking the expectation conditioned on $\boldsymbol{\theta}_k$ on both sides of (A1), together with (A2) and (A3), guarantees that, for all $k \in \mathbb{N}$,

$$\begin{aligned}
\mathbb{E}_{\xi_k}\left[f(\boldsymbol{\theta}_{k+1})|\boldsymbol{\theta}_k\right] &\leq f(\boldsymbol{\theta}_k) - \alpha_k \mathbb{E}_{\xi_k}\left[\langle \nabla f(\boldsymbol{\theta}_k), \nabla f_{B_k}(\boldsymbol{\theta}_k) \rangle|\boldsymbol{\theta}_k\right] + \frac{L\alpha_k^2}{2}\mathbb{E}_{\xi_k}\left[\|\nabla f_{B_k}(\boldsymbol{\theta}_k)\|^2|\boldsymbol{\theta}_k\right] \\
&\leq f(\boldsymbol{\theta}_k) - \alpha_k\|\nabla f(\boldsymbol{\theta}_k)\|^2 + \frac{L\alpha_k^2}{2}\left(\frac{\sigma^2}{b} + \|\nabla f(\boldsymbol{\theta}_k)\|^2\right).
\end{aligned}$$

Hence, taking the total expectation on both sides of the above inequality ensures that, for all $k \in \mathbb{N}$,

$$\alpha_k \left( 1 - \frac{L\alpha_k}{2} \right) \mathbb{E} \left[ \|\nabla f(\boldsymbol{\theta}_k)\|^2 \right] \leq \mathbb{E} \left[ f(\boldsymbol{\theta}_k) - f(\boldsymbol{\theta}_{k+1}) \right] + \frac{L\sigma^2 \alpha_k^2}{2b}.$$

Let $K \geq 1$. Summing the above inequality from $k = 0$ to $k = K - 1$ ensures that

$$\sum_{k=0}^{K-1} \alpha_k \left( 1 - \frac{L\alpha_k}{2} \right) \mathbb{E} \left[ \|\nabla f(\boldsymbol{\theta}_k)\|^2 \right] \leq \mathbb{E} \left[ f(\boldsymbol{\theta}_0) - f(\boldsymbol{\theta}_K) \right] + \frac{L\sigma^2}{2b} \sum_{k=0}^{K-1} \alpha_k^2,$$

which, together with (C1) (the lower bound $f_\star$ of $f$), implies that the assertion in Lemma A.1 holds. $\qquad\square$

### A.2. Proof of Theorem 3.1

**(Constant):** Lemma A.1 with $\alpha_k = \alpha$ implies that

$$\alpha \left( 1 - \frac{L\alpha}{2} \right) \sum_{k=0}^{K-1} \mathbb{E} \left[ \|\nabla f(\boldsymbol{\theta}_k)\|^2 \right] \leq \mathbb{E} \left[ f(\boldsymbol{\theta}_0) - f_\star \right] + \frac{L\sigma^2 \alpha^2 K}{2b}.$$

Since $\alpha < \frac{2}{L}$, we have that

$$\min_{k \in [0:K-1]} \mathbb{E} \left[ \|\nabla f(\boldsymbol{\theta}_k)\|^2 \right] \leq \frac{1}{K} \sum_{k=0}^{K-1} \mathbb{E} \left[ \|\nabla f(\boldsymbol{\theta}_k)\|^2 \right] \leq \underbrace{\frac{2(f(\boldsymbol{\theta}_0) - f_\star)}{(2 - L\alpha)\alpha}}_{C_1} \frac{1}{K} + \underbrace{\frac{L\sigma^2 \alpha}{2 - L\alpha}}_{C_2} \frac{1}{b}.$$

**(Decay):** Since $(\alpha_k)_{k \in \mathbb{N}}$ converges to 0, there exists $k_0 \in \mathbb{N}$ such that, for all $k \geq k_0$, $\alpha_k < \frac{2}{L}$. We assume that $k_0 = 0$ (see Section 2.2.2). Lemma A.1 ensures that, for all $K \geq 1$,

$$\sum_{k=0}^{K-1} \alpha_k \left( 1 - \frac{L\alpha_k}{2} \right) \mathbb{E} \left[ \|\nabla f(\boldsymbol{\theta}_k)\|^2 \right] \leq \mathbb{E} \left[ f(\boldsymbol{\theta}_0) - f_\star \right] + \frac{L\sigma^2}{2b} \sum_{k=0}^{K-1} \alpha_k^2,$$

which, together with $\alpha_{k+1} \leq \alpha_k < \frac{2}{L}$ ($k \in \mathbb{N}$), implies that

$$\alpha_{K-1} \left( 1 - \frac{L\alpha_0}{2} \right) \sum_{k=0}^{K-1} \mathbb{E} \left[ \|\nabla f(\boldsymbol{\theta}_k)\|^2 \right] \leq \mathbb{E} \left[ f(\boldsymbol{\theta}_0) - f_\star \right] + \frac{L\sigma^2}{2b} \sum_{k=0}^{K-1} \alpha_k^2.$$

Hence, we have that

$$\sum_{k=0}^{K-1} \mathbb{E} \left[ \|\nabla f(\boldsymbol{\theta}_k)\|^2 \right] \leq \frac{2(f(\boldsymbol{\theta}_0) - f_\star)}{(2 - L\alpha_0)\alpha_{K-1}} + \frac{L\sigma^2}{b(2 - L\alpha_0)\alpha_{K-1}} \sum_{k=0}^{K-1} \alpha_k^2,$$

which implies that

$$\min_{k\in[0:K-1]} \mathbb{E}\left[\|\nabla f(\boldsymbol{\theta}_k)\|^2\right] \leq \frac{2(f(\boldsymbol{\theta}_0) - f_\star)}{2 - L\alpha_0} \frac{1}{K\alpha_{K-1}} + \frac{1}{b}\frac{L\sigma^2}{2 - L\alpha_0}\frac{1}{K\alpha_{K-1}}\sum_{k=0}^{K-1}\alpha_k^2.$$

Meanwhile, we have that

$$\sum_{k=0}^{K-1}\alpha_k^2 \leq \sum_{k=0}^{K-1}\frac{T\alpha^2}{(k+1)^{2a}} \leq T\alpha^2\left(1 + \int_0^{K-1}\frac{\mathrm{d}t}{(t+1)^{2a}}\right) \leq \begin{cases} \dfrac{T\alpha^2}{1-2a}K^{1-2a} & \textbf{(Decay 1)} \\ T\alpha^2(1 + \log K) & \textbf{(Decay 2)} \\ \dfrac{2aT\alpha^2}{2a-1} & \textbf{(Decay 3)} \end{cases}$$

and

$$\alpha_{K-1} = \frac{\alpha}{\left(\lfloor\frac{K-1}{T}\rfloor + 1\right)^a} \geq \frac{\alpha}{\left(\frac{K-1}{T} + 1\right)^a} \geq \frac{\alpha}{K^a}.$$

Here, we define

$$D_1 := \frac{2\left(f(\boldsymbol{\theta}_0) - f_\star\right)}{\alpha(2 - L\alpha)} \text{ and } D_2 := \frac{T\alpha^2L\sigma^2}{2 - L\alpha}.$$

Accordingly, we have that

$$\min_{k\in[0:K-1]} \mathbb{E}\left[\|\nabla f(\boldsymbol{\theta}_k)\|^2\right] \leq \begin{cases} \dfrac{D_1}{K^{1-a}} + \dfrac{D_2}{(1-2a)K^ab} & \textbf{(Decay 1)} \\ \dfrac{D_1}{\sqrt{K}} + \dfrac{D_2(1+\log K)}{\sqrt{K}b} & \textbf{(Decay 2)} \\ \dfrac{D_1}{K^{1-a}} + \dfrac{2aD_2}{(2a-1)K^{1-a}b} & \textbf{(Decay 3)} \end{cases}$$

which, together with $\log K < \sqrt{K}$ and the condition on $a$, implies that

$$\min_{k\in[0:K-1]} \mathbb{E}\left[\|\nabla f(\boldsymbol{\theta}_k)\|^2\right] \leq \begin{cases} \dfrac{D_1}{K^a} + \dfrac{D_2}{(1-2a)K^ab} & \textbf{(Decay 1)} \\ \dfrac{D_1}{\sqrt{K}} + \left(\dfrac{1}{\sqrt{K}} + 1\right)\dfrac{D_2}{b} & \textbf{(Decay 2)} \\ \dfrac{D_1}{K^{1-a}} + \dfrac{2aD_2}{(2a-1)K^{1-a}b} & \textbf{(Decay 3)}. \end{cases}$$

$\square$

### A.3. Proof of Theorem 3.2

(i) Let us consider the case of **(Constant)**. We consider that the upper bound $\frac{C_1}{K} + \frac{C_2}{b}$ in Theorem 3.1 is equal to $\epsilon^2$. This implies that $K = \frac{C_1b}{\epsilon^2b - C_2}$ achieves an $\epsilon$–approximation. A discussion similar to the one showing that $K = \frac{C_1b}{\epsilon^2b - C_2}$ is an $\epsilon$–approximation ensures that the assertion in Theorem 3.2(i) is true.

(ii) It is sufficient to prove that $K' = K'(b) < 0$ and $K'' = K''(b) > 0$ hold.

**(Constant):** Let $K = \frac{C_1 b}{\epsilon^2 b - C_2}$. Then, we have that

$$K' = \frac{C_1(\epsilon^2 b - C_2) - \epsilon^2 C_1 b}{(\epsilon^2 b - C_2)^2} = -\frac{C_1 C_2}{(\epsilon^2 b - C_2)^2} < 0,$$

$$K'' = \frac{2\epsilon^2 C_1 C_2(\epsilon^2 b - C_2)}{(\epsilon^2 b - C_2)^4} = \frac{2\epsilon^2 C_1 C_2((\frac{C_1}{K} + \frac{C_2}{b})b - C_2)}{(\epsilon^2 b - C_2)^4} = \frac{2\epsilon^2 C_1^2 C_2^2}{K(\epsilon^2 b - C_2)^4} > 0.$$

**(Decay 1):** Let $K = (\frac{1}{\epsilon^2}(D_1 + \frac{D_2}{(1-2a)b}))^{\frac{1}{a}}$. Then, we have that

$$K' = \frac{1}{a}\left\{ \frac{1}{\epsilon^2}\left( D_1 + \frac{D_2}{(1-2a)b} \right) \right\}^{\frac{1}{a}-1} \left( -\frac{D_2}{\epsilon^2(1-2a)b^2} \right)$$

$$= -\frac{D_2}{a\epsilon^2(1-2a)b^2}\left\{ \frac{1}{\epsilon^2}\left( D_1 + \frac{D_2}{(1-2a)b} \right) \right\}^{\frac{1-a}{a}} < 0,$$

$$K'' = \frac{2D_2}{a\epsilon^2(1-2a)b^3}\left\{ \frac{1}{\epsilon^2}\left( D_1 + \frac{D_2}{(1-2a)b} \right) \right\}^{\frac{1-a}{a}}$$

$$+ \frac{2(1-a)D_2}{a^2\epsilon^2(1-2a)b^3}\left\{ \frac{1}{\epsilon^2}\left( D_1 + \frac{D_2}{(1-2a)b} \right) \right\}^{\frac{1-a}{a}-1} \frac{D_2}{\epsilon^2(1-2a)b^2} > 0.$$

**(Decay 2):** Let $K = (\frac{bD_1 + D_2}{b\epsilon^2 - D_2})^2$. Then, we have that

$$K' = \frac{2D_1(bD_1 + D_2)(b\epsilon^2 - D_2)^2 - 2\epsilon^2(b\epsilon^2 - D_2)(bD_1 + D_2)^2}{(b\epsilon^2 - D_2)^4},$$

which, together with $b\epsilon^2 - D_2 > 0$ and $D_1 > \epsilon^2$, implies that

$$(b\epsilon^2 - D_2)^3 K' = 2D_1(bD_1 + D_2)(b\epsilon^2 - D_2) - 2\epsilon^2(bD_1 + D_2)^2$$

$$= 2(bD_1 + D_2)\left\{ D_1(b\epsilon^2 - D_2) - \epsilon^2(bD_1 + D_2) \right\}$$

$$= -2(bD_1 + D_2)(D_1 D_2 - \epsilon^2 D_2)$$

$$= -2D_2(bD_1 + D_2)(D_1 - \epsilon^2) < 0.$$

Moreover,

$$K'' = \frac{-2D_1 D_2(D_1 - \epsilon^2)(b\epsilon^2 - D_2)^3 + 6D_2\epsilon^2(b\epsilon^2 - D_2)^2(bD_1 + D_2)(D_1 - \epsilon^2)}{(b\epsilon^2 - D_2)^6},$$

which implies that

$$(b\epsilon^2 - D_2)^4 K'' = -2D_1 D_2(D_1 - \epsilon^2)(b\epsilon^2 - D_2) + 6D_2\epsilon^2(bD_1 + D_2)(D_1 - \epsilon^2)$$

$$= 2D_2(D_1 - \epsilon^2)\left\{ -D_1(b\epsilon^2 - D_2) + 3\epsilon^2(bD_1 + D_2) \right\}$$

$$= 2D_2(D_1 - \epsilon^2)(2D_1\epsilon^2 b + D_1 D_2 + 3D_2\epsilon^2) > 0.$$

18

**(Decay 3):** Let $K = (\frac{1}{\epsilon^2}(D_1 + \frac{2aD_2}{(2a-1)b}))^{\frac{1}{1-a}}$. Then, we have that

$$
K' = \frac{1}{1-a}\left\{\frac{1}{\epsilon^2}\left(D_1 + \frac{2aD_2}{(2a-1)b}\right)\right\}^{\frac{1}{1-a}-1}\left(-\frac{2aD_2}{\epsilon^2(2a-1)b^2}\right)
$$

$$
= -\frac{2aD_2}{\epsilon^2(1-a)(2a-1)b^2}\left\{\frac{1}{\epsilon^2}\left(D_1 + \frac{2aD_2}{(2a-1)b}\right)\right\}^{\frac{a}{1-a}} < 0,
$$

$$
K'' = \frac{4aD_2}{\epsilon^2(1-a)(2a-1)b^3}\left\{\frac{1}{\epsilon^2}\left(D_1 + \frac{2aD_2}{(2a-1)b}\right)\right\}^{\frac{a}{1-a}}
$$

$$
+ \frac{2a^2D_2}{\epsilon^2(1-a)^2(2a-1)b^2}\left\{\frac{1}{\epsilon^2}\left(D_1 + \frac{2aD_2}{(2a-1)b}\right)\right\}^{\frac{a}{1-a}-1}\frac{2aD_2}{\epsilon^2(2a-1)b^2} > 0.
$$

$\square$

## A.4. Proof of Theorem 3.3

**(Constant):** Let $N = \frac{C_1 b^2}{\epsilon^2 b - C_2}$. Then, we have that

$$
N' = \frac{2C_1 b(\epsilon^2 b - C_2) - \epsilon^2 C_1 b^2}{(\epsilon^2 b - C_2)^2} = \frac{C_1 b(\epsilon^2 b - 2C_2)}{(\epsilon^2 b - C_2)^2}.
$$

If $N' = 0$, we have that $\epsilon^2 b - 2C_2 = 0$, i.e., $b = \frac{2C_2}{\epsilon^2}$. Moreover,

$$
N'' = \frac{(2\epsilon^2 C_1 b - 2C_1 C_2)(\epsilon^2 b - C_2)^2 - 2\epsilon^2(\epsilon^2 b - C_2)(\epsilon^2 C_1 b^2 - 2C_1 C_2 b)}{(\epsilon^2 b - C_2)^4}
$$

$$
(\epsilon^2 b - C_2)^3 N'' = (2\epsilon^2 C_1 b - 2C_1 C_2)(\epsilon^2 b - C_2) - 2\epsilon^2(\epsilon^2 C_1 b^2 - 2C_1 C_2 b)
$$

$$
= 2C_1 C_2^2 > 0,
$$

which implies that $N$ is convex. Hence, there is a critical batch size $b^\star = \frac{2C_2}{\epsilon^2} > 0$ at which $N$ is minimized.

**(Decay 1):** Let $N = Kb$. Then, we have that

$$
N' = K + bK'
$$

$$
= \left\{\frac{1}{\epsilon^2}\left(D_1 + \frac{D_2}{(1-2a)b}\right)\right\}^{\frac{1}{a}} + \frac{1}{a}\left\{\frac{1}{\epsilon^2}\left(D_1 + \frac{D_2}{(1-2a)b}\right)\right\}^{\frac{1}{a}-1}\left(-\frac{D_2}{\epsilon^2(1-2a)b^2}\right)b
$$

$$
= \left\{\frac{1}{\epsilon^2}\left(D_1 + \frac{D_2}{(1-2a)b}\right)\right\}^{\frac{1}{a}-1}\left\{\frac{1}{\epsilon^2}\left(D_1 + \frac{D_2}{(1-2a)b}\right) - \frac{D_2}{a\epsilon^2(1-2a)b}\right\}.
$$

If $N' = 0$, we have that

$$
\frac{1}{\epsilon^2}\left(D_1 + \frac{D_2}{(1-2a)b}\right) - \frac{D_2}{a\epsilon^2(1-2a)b} = 0, \text{ i.e., } b = \frac{D_2(a-1)}{aD_1(2a-1)}.
$$

Moreover,

$$N'' = K' + (K' + bK'') = 2K' + bK''$$

$$= -\frac{2D_2}{a\epsilon^2(1-2a)b^2}\left\{\frac{1}{\epsilon^2}\left(D_1 + \frac{D_2}{(1-2a)b}\right)\right\}^{\frac{1-a}{a}} + \frac{2D_2}{a\epsilon^2(1-2a)b^2}\left\{\frac{1}{\epsilon^2}\left(D_1 + \frac{D_2}{(1-2a)b}\right)\right\}^{\frac{1-a}{a}}$$

$$+ \frac{2(1-a)D_2}{a^2\epsilon^2(1-2a)b^2}\left\{\frac{1}{\epsilon^2}\left(D_1 + \frac{D_2}{(1-2a)b}\right)\right\}^{\frac{1-a}{a}-1}\frac{D_2}{\epsilon^2(1-2a)b^2}$$

$$= \frac{2(1-a)D_2}{a^2\epsilon^2(1-2a)b^2}\left\{\frac{1}{\epsilon^2}\left(D_1 + \frac{D_2}{(1-2a)b}\right)\right\}^{\frac{1-a}{a}-1}\frac{D_2}{\epsilon^2(1-2a)b^2} > 0,$$

which implies that $N$ is convex. Hence, there is a critical batch size $b^\star = \frac{D_2(a-1)}{aD_1(2a-1)} > 0$.

**(Decay 2):** Let $N = bK$. Then, we have that

$$N' = K + bK'$$

$$= \frac{(bD_1 - D_2)^2}{(b\epsilon^2 - D_2)^2} - \frac{2D_2b(bD_1 + D_2)(D_1 - \epsilon^2)}{(b\epsilon^2 - D_2)^3}$$

$$= \frac{bD_1 + D_2}{(b\epsilon^2 - D_2)^3}\{(bD_1 + D_2)(b\epsilon^2 - D_2) - 2D_2b(D_1 - \epsilon^2)\}$$

$$= \frac{bD_1 + D_2}{(b\epsilon^2 - D_2)^3}\{D_1\epsilon^2b^2 + 3D_2(\epsilon^2 - D_1)b - D_2^2\}.$$

If $N' = 0$, we have that $D_1b + D_2 = 0$, i.e., $b = -\frac{D_2}{D_1} < 0$. Moreover,

$$N'' = 2K' + bK''$$

$$= -\frac{4D_2(bD_1 + D_2)(D_1 - \epsilon^2)}{(b\epsilon^2 - D_2)^3} + \frac{2D_2b(D_1 - \epsilon^2)(2D_1\epsilon^2b + D_1D_2 + 3D_2\epsilon^2)}{(b\epsilon^2 - D_2)^4}$$

$$= \frac{2D_2(D_1 - \epsilon^2)}{(b\epsilon^2 - D_2)^4}\{-2(bD_1 + D_2)(b\epsilon^2 - D_2) + b(2D_1\epsilon^2b + D_1D_2 + 3D_2\epsilon^2)\}$$

$$= \frac{2D_2(D_1 - \epsilon^2)}{(b\epsilon^2 - D_2)^4}(3D_1D_2b + D_2\epsilon^2b + 2D_2^2) > 0,$$

which implies that $N$ is convex. We can check that $N'(b) > 0$ for all $b > \frac{D_2}{\epsilon^2}$.

**(Decay 3):** Let $N = bK$. Then, we have that

$$N' = K + bK'$$

$$= \left\{\frac{1}{\epsilon^2}\left(D_1 + \frac{2aD_2}{(2a-1)b}\right)\right\}^{\frac{1}{1-a}} - \frac{2aD_2}{\epsilon^2(1-a)(2a-1)b}\left\{\frac{1}{\epsilon^2}\left(D_1 + \frac{2aD_2}{(2a-1)b}\right)\right\}^{\frac{1}{1-a}-1}$$

$$= \left\{\frac{1}{\epsilon^2}\left(D_1 + \frac{2aD_2}{(2a-1)b}\right)\right\}^{\frac{1}{1-a}-1}\left\{\frac{1}{\epsilon^2}\left(D_1 + \frac{2aD_2}{(2a-1)b}\right) - \frac{2aD_2}{\epsilon^2(1-a)(2a-1)b}\right\}.$$

If $N' = 0$, we have that

$$\frac{1}{\epsilon^2}\left(D_1 + \frac{2aD_2}{(2a-1)b}\right) - \frac{2aD_2}{\epsilon^2(1-a)(2a-1)b} = 0, \text{ i.e., } b = \frac{2a^2D_2}{(2a-1)(1-a)D_1}.$$

20

Moreover,

$$N'' = 2K' + bK''$$

$$= -\frac{2aD_2}{\epsilon^2(1-a)(2a-1)b^2}\left\{\frac{1}{\epsilon^2}\left(D_1 + \frac{2aD_2}{(2a-1)b}\right)\right\}^{\frac{a}{1-a}}$$

$$+ \frac{2aD_2}{\epsilon^2(1-a)(2a-1)b^2}\left\{\frac{1}{\epsilon^2}\left(D_1 + \frac{2aD_2}{(2a-1)b}\right)\right\}^{\frac{a}{1-a}}$$

$$+ \frac{2a^2D_2}{\epsilon^2(1-a)^2(2a-1)b}\left\{\frac{1}{\epsilon^2}\left(D_1 + \frac{2aD_2}{(2a-1)b}\right)\right\}^{\frac{2a-1}{1-a}}\frac{2aD_2}{\epsilon^2(2a-1)b^2}$$

$$= \frac{2a^2D_2}{\epsilon^2(1-a)^2(2a-1)b}\left\{\frac{1}{\epsilon^2(D_1 + \frac{2aD_2}{(2a-1)b})}\right\}^{\frac{2a-1}{1-a}}\frac{2aD_2}{\epsilon^2(2a-1)b^2} > 0,$$

which implies that $N$ is convex. Hence, there is a critical batch size $b^\star = \frac{2a^2D_2}{(2a-1)(1-a)D_1} > 0$. □

### *A.5. Proof of Theorem 3.4*

Using $K$ defined in Theorem 3.2 leads to the iteration complexity. For example, SGD using **(Constant)** satisfies $N(b) = \frac{C_1 b^2}{\epsilon^2 b - C_2}$ (Theorem 3.3). Using the critical batch size $b^\star = \frac{2C_2}{\epsilon^2}$ in (4) leads to

$$\inf\left\{N\colon \min_{k\in[0:K-1]}\mathbb{E}[\|\nabla f(\boldsymbol{\theta}_k)\|] \leq \epsilon\right\} \leq N(b^\star) = \frac{4C_1 C_2}{\epsilon^4},\ \text{i.e.,}\ \mathcal{N}_\epsilon = O\left(\frac{1}{\epsilon^4}\right).$$

A similar discussion, together with using $N$ defined in Theorem 3.3 and the critical batch size $b^\star$ in (4), leads to the SFO complexities of **(Decay 1)** and **(Decay 3)**. Using $N$ defined in Theorem 3.3 and a batch size $b = \frac{D_2+1}{\epsilon^2}$ leads to the SFO complexity of **(Decay 2)**. □