

Relationship between Batch Size and Number of Steps Needed for Nonconvex Optimization of Stochastic Gradient Descent using Armijo Line Search

Anonymous authors

Paper under double-blind review

Abstract

Stochastic gradient descent (SGD) is the simplest deep learning optimizer with which to train deep neural networks. While SGD can use various learning rates, such as constant or diminishing rates, the previous numerical results showed that SGD performs better than other deep learning optimizers using when it uses learning rates given by line search methods. In this paper, we perform a convergence analysis on SGD with a learning rate given by an Armijo line search for nonconvex optimization. The analysis indicates that the upper bound of the expectation of the squared norm of the full gradient becomes small when the number of steps and the batch size are large. Next, we show that, for SGD with the Armijo-line-search learning rate, the number of steps needed for nonconvex optimization is a monotone decreasing convex function of the batch size; that is, the number of steps needed for nonconvex optimization decreases as the batch size increases. Furthermore, we show that the stochastic first-order oracle (SFO) complexity, which is the stochastic gradient computation cost, is a convex function of the batch size; that is, there exists a critical batch size that minimizes the SFO complexity. Finally, we provide numerical results that support our theoretical results. The numerical results indicate that the number of steps needed for training deep neural networks decreases as the batch size increases and that there exist the critical batch sizes that can be estimated from the theoretical results.

1 Introduction

1.1 Background

Nonconvex optimization is useful for training deep neural networks, since the loss functions called the expected risk and empirical risk are nonconvex and they need only be minimized in order to find the model parameters. Deep-learning optimizers have been presented for minimizing the loss functions. The simplest one is stochastic gradient descent (SGD) (Robbins & Monro, 1951; Zinkevich, 2003; Nemirovski et al., 2009; Ghadimi & Lan, 2012; 2013) and there are numerous theoretical analyses on using SGD for nonconvex optimization (Jain et al., 2018; Vaswani et al., 2019; Fehrman et al., 2020; Chen et al., 2020; Scaman & Malherbe, 2020; Loizou et al., 2021). Variants have also been presented, such as momentum methods (Polyak, 1964; Nesterov, 1983) and adaptive methods including Adaptive Gradient (AdaGrad) (Duchi et al., 2011), Root Mean Square Propagation (RMSProp) (Tieleman & Hinton, 2012), Adaptive Moment Estimation (Adam) (Kingma & Ba, 2015), Adaptive Mean Square Gradient (AMSGrad) (Reddi et al., 2018), and Adam with decoupled weight decay (AdamW) (Loshchilov & Hutter, 2019). SGD and its variants are useful for training not only deep neural networks but also generative adversarial networks (Heusel et al., 2017; Naganuma & Iiduka, 2023; Sato & Iiduka, 2023).

The performance of deep-learning optimizers for nonconvex optimization depends on the batch size. The previous numerical results in (Shallue et al., 2019) and (Zhang et al., 2019) have shown that the number of steps K needed to train a deep neural network halves for each doubling of the batch size b and that there is a region of diminishing returns beyond the *critical batch size* b^* . This fact can be expressed as follows: there

is a positive number C such that $N := Kb \approx C$ for $b \leq b^*$ and $N := Kb \geq C$ for $b \geq b^*$. The deep neural network model uses b gradients of the loss functions per step. Hence, when K is the number of steps required to train a deep neural network, the model has a stochastic gradient computation cost of Kb . We will define the *stochastic first-order oracle (SFO) complexity* (Iiduka, 2022; Sato & Iiduka, 2023) of a deep-learning optimizer to be $N := Kb$. From the previous numerical results in (Shallue et al., 2019) and (Zhang et al., 2019), the SFO complexity is minimized at a critical batch size b^* and there are diminishing returns once the batch size exceeds b^* . Therefore, it is desirable to use the critical batch size when minimizing the SFO complexity of the deep-learning optimizer.

Not only a batch size but also a learning rate affects the performance of deep-learning optimizers for nonconvex optimization. A performance measure of a deep-learning optimizer generating a sequence $(\boldsymbol{\theta}_k)_{k \in \mathbb{N}}$ is the expectation of the squared norm of the gradient of a nonconvex loss function f , denoted by $\mathbb{E}[\|\nabla f(\boldsymbol{\theta}_k)\|^2]$. If this performance measure becomes small when the number of steps k is large, the deep-learning optimizer approximates a local minimizer of f . For example, let us consider the problem of minimizing a smooth function f (see Section 2.1 for the definition of smoothness). Here, SGD using a constant learning rate $\alpha = O(\frac{1}{L})$ satisfies $\min_{k \in [K]} \mathbb{E}[\|\nabla f(\boldsymbol{\theta}_k)\|^2] = O(\frac{1}{K} + \frac{\alpha}{b})$, where L is the Lipschitz constant of ∇f , b is the batch size, and $[K] := \{1, 2, \dots, K\}$ (see also Table 1). Moreover, SGD using a learning rate satisfying the Armijo condition was presented in (Vaswani et al., 2019). The *Armijo line search* (Nocedal & Wright, 2006, Chapter 3.1) is a standard method for finding an appropriate learning rate α_k giving a sufficient decrease in f , i.e., $f(\boldsymbol{\theta}_{k+1}) < f(\boldsymbol{\theta}_k)$ (see Section 2.3.1 for the definition of the Armijo condition).

1.2 Motivation

The numerical results in (Vaswani et al., 2019) indicated that the Armijo-line-search learning rate is superior to using a constant learning rate when using SGD to train deep neural networks in the sense of minimizing the training loss functions and improving test accuracy. Motivated by the useful numerical results in (Vaswani et al., 2019), we decided to perform convergence analyses on SGD with the Armijo-line-search learning rate for nonconvex optimization in deep neural networks.

Theorem 3 in (Vaswani et al., 2019) is a convergence analysis of SGD with the Armijo-line-search learning rate for nonconvex optimization under a strong growth condition that implies the interpolation property. Here, let $f: \mathbb{R}^d \rightarrow \mathbb{R}$ be an empirical risk defined by $f(\boldsymbol{\theta}) := \frac{1}{n} \sum_{i \in [n]} f_i(\boldsymbol{\theta})$, where n is the number of training data and $f_i: \mathbb{R}^d \rightarrow \mathbb{R}$ is a loss function corresponding to the i -th training data z_i . We say that f has the interpolation property if $\nabla f(\boldsymbol{\theta}) = \mathbf{0}$ implies $\nabla f_i(\boldsymbol{\theta}) = \mathbf{0}$ ($i \in [n]$). The interpolation property holds for optimization of a linear model with the squared hinge loss for binary classification on linearly separable data (Vaswani et al., 2019, Section 2). However, the interpolation condition would be unrealistic for deep neural networks, since their loss functions are nonconvex. The motivation behind this work is thus to show that SGD with the Armijo-line-search learning rate can solve nonconvex optimization problems in deep neural networks.

As indicated the second paragraph in Section 1.1, the batch size has a significant effect on the performance of SGD. Hence, in accordance with the first motivation stated above, we decided to investigate appropriate batch sizes for SGD with the Armijo-line-search learning rate. In particular, we were interested in verifying whether a critical batch size b^* minimizing the SFO complexity N exists for training deep neural networks with SGD using the Armijo condition in theory and in practice. This is because the previous studies in (Shallue et al., 2019; Zhang et al., 2019; Iiduka, 2022; Sato & Iiduka, 2023) showed the existence of critical batch sizes for training deep neural networks or generative adversarial networks with optimizers with constant or diminishing learning rates and without Armijo-line-search learning rates.

We are also interested in estimating critical batch sizes before implementing SGD with the Armijo-line-search learning rate. The previous results in (Iiduka, 2022; Sato & Iiduka, 2023) showed that, for optimizers using constant learning rates, the critical batch sizes determined from numerical results are close to the theoretically estimated sizes. Motivated by the results in (Iiduka, 2022; Sato & Iiduka, 2023), we sought to verify whether, for SGD with the Armijo-line-search learning rate, the measured critical batch sizes are close to the batch sizes estimated from theoretical results.

1.3 Contribution

1.3.1 Convergence analysis of SGD with Armijo-line-search learning rates

The first contribution of this paper is to present a convergence analysis of SGD with Armijo-line-search learning rates for general nonconvex optimization (Theorem 3.1); in particular, it is shown that SGD with this rate α_k satisfies that, for all $K \geq 1$,

$$\min_{k \in [0:K-1]} \mathbb{E} [\|\nabla f(\boldsymbol{\theta}_k)\|^2] \leq \underbrace{\frac{C_1}{\delta(1-c)(2-L_n\bar{\alpha})} \frac{1}{K}}_{B(\boldsymbol{\theta}_0, K)} + \underbrace{\frac{C_2}{2-L_n\bar{\alpha}} \frac{1}{b}}_{V(\sigma^2, b)}, \quad (1)$$

where the parameters are defined in Table 1 (see also Theorem 3.1). The inequality (1) indicates that the upper bound of the performance measure $\min_{k \in [0:K-1]} \mathbb{E} [\|\nabla f(\boldsymbol{\theta}_k)\|^2]$ that consists of a bias term $B(\boldsymbol{\theta}_0, K)$ and variance term $V(\sigma^2, b)$ becomes small when the number of steps K is large and the batch size b is large. Therefore, it is desirable to set K large and b large so that Algorithm 1 will approximate a local minimizer of f .

The essential lemma to proving (1) is the guarantee of the existence of a lower bound on the learning rates satisfying the Armijo condition (Lemma 2.1). Although, in general, learning rates satisfying the Armijo condition do not have any lower bound (Lemma 2.1(i)), the corresponding learning rates computed by a backtracking line search (Algorithm 2) have a lower bound (Lemma 2.1(ii)). In addition, the descent lemma (i.e., $f(\mathbf{y}) \leq f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle + \frac{L_n}{2} \|\mathbf{y} - \mathbf{x}\|^2$ ($\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$)) holds from the smoothness condition on f . Thus, we can prove (1) by using the existence of a lower bound on the learning rates satisfying the Armijo condition and the descent lemma (see Appendix A.2 for details of the proof of Theorem 3.1).

Table 1: Relationship between batch size b and number of steps K to achieve an ϵ -approximation defined by $\min_{k \in [0:K-1]} \mathbb{E} [\|\nabla f(\boldsymbol{\theta}_k)\|^2] \leq \frac{C_1}{K} + \frac{C_2}{b} = \epsilon^2$ for SGD with a constant learning rate $\alpha \in (0, \frac{2}{L_n})$ and for SGD with the Armijo-line-search learning rate ($[0 : K - 1] := \{0, 1, \dots, K - 1\}$. $f := \frac{1}{n} \sum_{i \in [n]} f_i$ is bounded below by f_* , L_i is the Lipschitz constant of ∇f_i , $L_n := \frac{1}{n} \sum_{i \in [n]} L_i$, $C(\boldsymbol{\theta}_0, f_*)$ is a positive real number depending on the initial point $\boldsymbol{\theta}_0 \in \mathbb{R}^d$ and f_* , σ^2 is the upper bound of the variance of the stochastic gradient, and $\bar{\alpha}$ is the upper bound of the learning rates)

Learning Rate		Upper Bound $\frac{C_1}{K} + \frac{C_2}{b}$	Steps K	SFO N	Critical Batch b^*
Constant $\alpha \in \left(0, \frac{2}{L_n}\right)$	C_1	$\frac{2(f(\boldsymbol{\theta}_0) - f_*)}{(2 - L_n\alpha)\alpha}$	$K = \frac{C_1 b}{\epsilon^2 b - C_2}$	$N = \frac{C_1 b^2}{\epsilon^2 b - C_2}$	$b^* = \frac{2C_2}{\epsilon^2}$
	C_2	$\frac{L_n \alpha \sigma^2}{2 - L_n \alpha}$			
Armijo ($c, \delta \in (0, 1)$)	C_1	$\frac{LC(\boldsymbol{\theta}_0, f_*)}{\delta(1-c)(2-L_n\bar{\alpha})}$	$K = \frac{C_1 b}{\epsilon^2 b - C_2}$	$N = \frac{C_1 b^2}{\epsilon^2 b - C_2}$	$b^* = \frac{2C_2}{\epsilon^2}$
	C_2	$\frac{L_n \bar{\alpha} \sigma^2}{2 - L_n \bar{\alpha}}$			

To show the merit of SGD with the Armijo line search learning rate, we compare an implementation using this rate with one using a constant learning rate (see, e.g., (Scaman & Malherbe, 2020, Section 4) for convergence analyses of SGD with constant learning rates). In this case, we need to set a constant learning rate $\alpha \in (0, \frac{2}{L_n})$ depending on the Lipschitz constant L_n of ∇f (see also Table 1). However, computing L_n is NP-hard (Virmaux & Scaman, 2018), so it would be unrealistic to set a constant learning rate depending on L_n before implementing SGD. Meanwhile, we need to set $c, \delta \in (0, 1)$ in order to use SGD with the Armijo-line-search learning rate. We would like to emphasize here that we can choose *any* $c, \delta \in (0, 1)$ to implement SGD. That is, for any $c, \delta \in (0, 1)$, there exists a learning rate satisfying the Armijo condition (see (9) for the definition of the Armijo condition) and it can be found by conducting a simple backtracking

line search (Algorithm 2) instead of performing a complicated computation such as of the Lipschitz constant of ∇f (Lemma 2.1).

1.3.2 Steps needed for ϵ -approximation of SGD with Armijo line-search-learning rates

The previous results in (Shallue et al., 2019; Zhang et al., 2019; Iiduka, 2022; Sato & Iiduka, 2023) indicated that, for optimizers, the number of steps K needed to train a deep neural network or generative adversarial networks decreases as the batch size increases. The second contribution of this paper is to show that, for SGD with the Armijo-line-search learning rate, the number of steps K needed for nonconvex optimization decreases as the batch size increases. Let us consider the case in which the right-hand side of (1) is equal to ϵ^2 , where $\epsilon > 0$ is the precision. Then, K is a rational function defined for a batch size b by

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

where C_1 and C_2 are the positive constants defined in (1) (see also Table 1). We can easily show that K defined above is a monotone decreasing and convex function with respect to b (Theorem 3.2). Accordingly, the number of steps needed for nonconvex optimization decreases as the batch size increases.

1.3.3 Critical batch size minimizing SFO complexity of SGD with Armijo-line-search learning rates

Using K defined by (2) above, we can further define the SFO complexity N of SGD with Armijo-line-search learning rates (see also Table 1):

$$N = Kb = K(b)b = \frac{C_1 b^2}{\epsilon^2 b - C_2}. \quad (3)$$

We can easily show that N is convex with respect to b and that a global minimizer

$$b^* = \frac{2C_2}{\epsilon^2} = \frac{2L_n \bar{\alpha} \sigma^2}{(2 - L_n \bar{\alpha}) \epsilon^2} \quad (4)$$

exists for it (Theorem 3.3). Accordingly, there is a critical batch size b^* at which N is minimized.

1.3.4 Numerical results supporting our theoretical results

The numerical results in (Vaswani et al., 2019) showed that SGD with the Armijo-line-search learning rate performs better than other optimizers in training deep neural networks. Hence, we sought to verify whether the numerical results match our theoretical results (Sections 1.3.1, 1.3.2, and 1.3.3). We trained residual networks (ResNets) on the CIFAR-10 and CIFAR-100 datasets and a two-hidden-layer multi-layer perceptron (MLP) on the MNIST dataset. We numerically found that increasing the batch size b decreases the number of steps K needed to achieve high training accuracies and that there are critical batch sizes minimizing the SFO complexities. We also estimated batch sizes using (4) for the critical batch size b^* and compared them with ones determined from the numerical results. We found that the estimated batch sizes are close to the ones determined from the numerical results (Section 4).

2 Mathematical Preliminaries

2.1 Definitions

Let \mathbb{N} be the set of nonnegative integers, $[n] := \{1, 2, \dots, n\}$ for $n \geq 1$, and $[0 : n] := \{0, 1, \dots, n\}$ for $n \geq 0$. Let \mathbb{R}^d be a d -dimensional Euclidean space with inner product $\langle \cdot, \cdot \rangle$ inducing the norm $\|\cdot\|$.

Let $f: \mathbb{R}^d \rightarrow \mathbb{R}$ be continuously differentiable. We denote the gradient of f by $\nabla f: \mathbb{R}^d \rightarrow \mathbb{R}^d$. Let $L > 0$. $f: \mathbb{R}^d \rightarrow \mathbb{R}$ is said to be L -smooth if $\nabla f: \mathbb{R}^d \rightarrow \mathbb{R}^d$ is L -Lipschitz continuous, i.e., for all $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$, $\|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\| \leq L\|\mathbf{x} - \mathbf{y}\|$. When $f: \mathbb{R}^d \rightarrow \mathbb{R}$ is L -smooth, the following inequality, called the descent lemma (Beck, 2017, Lemma 5.7), holds: for all $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$, $f(\mathbf{y}) \leq f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle + \frac{L}{2}\|\mathbf{y} - \mathbf{x}\|^2$. Let $f_* \in \mathbb{R}$. $f: \mathbb{R}^d \rightarrow \mathbb{R}$ is said to be bounded below by f_* if, for all $\mathbf{x} \in \mathbb{R}^d$, $f(\mathbf{x}) \geq f_*$.

2.2 Assumptions and problem

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 by a differentiable nonconvex loss function $f(\boldsymbol{\theta}; z)$. We aim to minimize the empirical average loss defined for all $\boldsymbol{\theta} \in \mathbb{R}^d$ by

$$f(\boldsymbol{\theta}) = \frac{1}{n} \sum_{i \in [n]} f(\boldsymbol{\theta}; z_i) = \frac{1}{n} \sum_{i \in [n]} f_i(\boldsymbol{\theta}),$$

where $S = (z_1, z_2, \dots, z_n)$ denotes the training set and $f_i(\cdot) := f(\cdot; z_i)$ denotes the loss function corresponding to the i -th training data z_i .

This paper considers the following smooth nonconvex optimization problem.

Problem 2.1 Suppose that $f_i: \mathbb{R}^d \rightarrow \mathbb{R}$ ($i \in [n]$) is L_i -smooth and bounded below by $f_{i,*}$. Then,

$$\text{minimize } f(\boldsymbol{\theta}) := \frac{1}{n} \sum_{i \in [n]} f_i(\boldsymbol{\theta}) \text{ subject to } \boldsymbol{\theta} \in \mathbb{R}^d.$$

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 $G_\xi(\boldsymbol{\theta})$ of the function f , where a random variable ξ is supported on Ξ independently of $\boldsymbol{\theta}$. We make the following standard assumptions.

Assumption 2.1

(A1) Let $(\boldsymbol{\theta}_k)_{k \in \mathbb{N}} \subset \mathbb{R}^d$ be the sequence generated by SGD. For each iteration k ,

$$\mathbb{E}_{\xi_k} [G_{\xi_k}(\boldsymbol{\theta}_k)] = \nabla f(\boldsymbol{\theta}_k),$$

where ξ_0, ξ_1, \dots are independent samples and the random variable ξ_k is independent of $(\boldsymbol{\theta}_l)_{l=0}^k$. There exists a nonnegative constant σ^2 such that

$$\mathbb{E}_{\xi_k} [\|G_{\xi_k}(\boldsymbol{\theta}_k) - \nabla f(\boldsymbol{\theta}_k)\|^2] \leq \sigma^2. \quad (5)$$

(A2) For each iteration k , SGD samples a batch B_k of size b independently of k and estimates the full gradient ∇f as

$$\nabla f_{B_k}(\boldsymbol{\theta}_k) := \frac{1}{b} \sum_{i \in [b]} G_{\xi_{k,i}}(\boldsymbol{\theta}_k) = \frac{1}{b} \sum_{i \in [b]} \nabla f_{\xi_{k,i}}(\boldsymbol{\theta}_k),$$

where $\xi_{k,i}$ is a random variable generated by the i -th sampling in the k -th iteration.

2.3 Stochastic gradient descent using Armijo line search

2.3.1 Armijo condition

Suppose that $f: \mathbb{R}^d \rightarrow \mathbb{R}$ is continuously differentiable. We would like to find a stationary point $\boldsymbol{\theta}^* \in \mathbb{R}^d$ such that $\nabla f(\boldsymbol{\theta}^*) = \mathbf{0}$ by using an iterative method defined by

$$\boldsymbol{\theta}_{k+1} := \boldsymbol{\theta}_k + \alpha_k \mathbf{d}_k, \quad (6)$$

where $\alpha_k > 0$ is the step size (called a learning rate in machine learning field) and $\mathbf{d}_k \in \mathbb{R}^d$ is the search direction. Various methods can be used depending on the search direction \mathbf{d}_k . For example, the method (6) with $\mathbf{d}_k := -\nabla f(\boldsymbol{\theta}_k)$ is gradient descent, while the method (6) with $\mathbf{d}_k := -\nabla f(\boldsymbol{\theta}_k) + \beta_{k-1} \mathbf{d}_{k-1}$, where $\beta_k \geq 0$, is the conjugate gradient method. If we define \mathbf{d}_k (e.g., $\mathbf{d}_k := -\nabla f(\boldsymbol{\theta}_k)$), it is desirable to set α_k^* satisfying

$$f(\boldsymbol{\theta}_k + \alpha_k^* \mathbf{d}_k) = \min_{\alpha > 0} f(\boldsymbol{\theta}_k + \alpha \mathbf{d}_k). \quad (7)$$

The step size α_k^* defined by (7) can be easily computed when f is quadratic and convex. However, for a general nonconvex function f , it is difficult to compute the step size α_k^* in (7) exactly. Here, we can use the *Armijo condition* for finding an appropriate step size α_k : Let $c \in (0, 1)$. We would like to find $\alpha_k > 0$ such that

$$f(\boldsymbol{\theta}_k + \alpha_k \mathbf{d}_k) \leq f(\boldsymbol{\theta}_k) + c\alpha_k \langle \nabla f(\boldsymbol{\theta}_k), \mathbf{d}_k \rangle. \quad (8)$$

When \mathbf{d}_k satisfies the descent property defined by $\langle \nabla f(\boldsymbol{\theta}_k), \mathbf{d}_k \rangle < 0$ (e.g., gradient descent using $\mathbf{d}_k := -\nabla f(\boldsymbol{\theta}_k)$ has the property such that $\langle \nabla f(\boldsymbol{\theta}_k), \mathbf{d}_k \rangle = -\|\nabla f(\boldsymbol{\theta}_k)\|^2 < 0$), the Armijo condition ensures that $f(\boldsymbol{\theta}_{k+1}) = f(\boldsymbol{\theta}_k + \alpha_k \mathbf{d}_k) < f(\boldsymbol{\theta}_k)$. Accordingly, α_k satisfying the Armijo condition (8) is appropriate in the sense of minimizing f .

The existence of step sizes satisfying the Armijo condition (8) is guaranteed.

Proposition 2.1 (Nocedal & Wright, 2006, Lemma 3.1) *Let $f: \mathbb{R}^d \rightarrow \mathbb{R}$ be continuously differentiable. Let $\boldsymbol{\theta}_k \in \mathbb{R}^d$ and let $\mathbf{d}_k (\neq \mathbf{0})$ have the descent property defined by $\langle \nabla f(\boldsymbol{\theta}_k), \mathbf{d}_k \rangle < 0$. Let $c \in (0, 1)$. Then, there exists $\gamma_k > 0$ such that, for all $\alpha_k \in (0, \gamma_k]$, the Armijo condition (8) holds.*

2.3.2 Stochastic gradient descent under Armijo condition

The objective of this paper is to solve Problem 2.1 using mini-batch SGD under Assumption 2.1 defined by

$$\boldsymbol{\theta}_{k+1} = \boldsymbol{\theta}_k + \alpha_k \mathbf{d}_k = \boldsymbol{\theta}_k - \alpha_k \nabla f_{B_k}(\boldsymbol{\theta}_k) = \boldsymbol{\theta}_k - \frac{\alpha_k}{b} \sum_{i \in [b]} \mathbf{G}_{\xi_{k,i}}(\boldsymbol{\theta}_k),$$

where $b > 0$ is the batch size and $\alpha_k > 0$ is the learning rate. For each iteration k , we can use $\boldsymbol{\theta}_k$, f_{B_k} , and ∇f_{B_k} . Hence, the Armijo condition (Vaswani et al., 2019, (1)) at the k -th iteration for SGD can be obtained by replacing f in (8) with f_{B_k} and using $\mathbf{d}_k = -\nabla f_{B_k}(\boldsymbol{\theta}_k)$:

$$f_{B_k}(\boldsymbol{\theta}_k - \alpha_k \nabla f_{B_k}(\boldsymbol{\theta}_k)) \leq f_{B_k}(\boldsymbol{\theta}_k) - c\alpha_k \|\nabla f_{B_k}(\boldsymbol{\theta}_k)\|^2. \quad (9)$$

The Armijo condition (9) ensures that $f_{B_k}(\boldsymbol{\theta}_{k+1}) = f_{B_k}(\boldsymbol{\theta}_k - \alpha_k \nabla f_{B_k}(\boldsymbol{\theta}_k)) < f_{B_k}(\boldsymbol{\theta}_k)$; i.e., the Armijo condition (9) is appropriate in the sense of minimizing the estimated objective function f_{B_k} from the full objective function f . In fact, the numerical results in (Vaswani et al., 2019, Section 7) indicate that SGD using the Armijo condition (9) is superior to using other deep-learning optimizers to train deep neural networks.

Algorithm 1 is the SGD algorithm using the Armijo condition (9).

Algorithm 1 Stochastic gradient descent using Armijo line search

Require: $c \in (0, 1)$ (hyperparameter), $b > 0$ (batch size), $\boldsymbol{\theta}_0 \in \mathbb{R}^d$ (initial point), $K \geq 1$ (steps)

Ensure: $\boldsymbol{\theta}_K \in \mathbb{R}^d$

$k \leftarrow 0$

for $k = 0, 1, \dots, K - 1$ **do**

 Compute $\alpha_k > 0$ satisfying $f_{B_k}(\boldsymbol{\theta}_k - \alpha_k \nabla f_{B_k}(\boldsymbol{\theta}_k)) \leq f_{B_k}(\boldsymbol{\theta}_k) - c\alpha_k \|\nabla f_{B_k}(\boldsymbol{\theta}_k)\|^2 \triangleleft$ Algorithm 2

 Compute $\boldsymbol{\theta}_{k+1} = \boldsymbol{\theta}_k - \alpha_k \nabla f_{B_k}(\boldsymbol{\theta}_k)$

end for

The search direction of Algorithm 1 is $\mathbf{d}_k = -\nabla f_{B_k}(\boldsymbol{\theta}_k) (\neq \mathbf{0})$ which has the descent property defined by $\langle \nabla f_{B_k}(\boldsymbol{\theta}_k), \mathbf{d}_k \rangle = -\|\nabla f_{B_k}(\boldsymbol{\theta}_k)\|^2 < 0$. Hence, from Proposition 2.1, there exists a learning rate $\alpha_k \in (0, \gamma_k]$ satisfying the Armijo condition (9). Moreover, the proposition guarantees that the learning rate can be chosen to be sufficiently small, e.g., $\liminf_{k \rightarrow +\infty} \alpha_k = 0$.

The convergence analyses of Algorithm 1 use a lower bound of $\alpha_k \in (0, \gamma_k]$ satisfying the Armijo condition (9). To guarantee the existence of such a lower bound, we use the backtracking method ((Nocedal & Wright, 2006, Algorithm 3.1) and (Vaswani et al., 2019, Algorithm 2)) described in Algorithm 2.

The following lemma guarantees the existence of a lower bound on the learning rates computed by Algorithm 2. The proof is given in Appendix A.1.

Algorithm 2 Backtracking Armijo-line-search method (Nocedal & Wright, 2006, Algorithm 3.1)**Require:** $c, \delta, \frac{1}{\gamma} \in (0, 1)$ (hyperparameters), $\alpha = \gamma^{\frac{b}{n}} \alpha_{k-1}$ (initialization), $\theta_k \in \mathbb{R}^d$, $f_{B_k}: \mathbb{R}^d \rightarrow \mathbb{R}$ **Ensure:** α_k satisfying $f_{B_k}(\theta_k - \alpha_k \nabla f_{B_k}(\theta_k)) \leq f_{B_k}(\theta_k) - c\alpha_k \|\nabla f_{B_k}(\theta_k)\|^2$ **repeat** $\alpha \leftarrow \delta\alpha$ **until** $f_{B_k}(\theta_k - \alpha \nabla f_{B_k}(\theta_k)) \leq f_{B_k}(\theta_k) - c\alpha \|\nabla f_{B_k}(\theta_k)\|^2$

Lemma 2.1 Consider Algorithm 1 under Assumption 2.1 for solving Problem 2.1. Let α_k be a learning rate satisfying the Armijo condition (9) (whose existence is guaranteed by Proposition 2.1), let L_{B_k} be the Lipschitz constant of ∇f_{B_k} , and let L be the maximum value of the Lipschitz constant L_i of ∇f_i . Then, the following hold.

- (i) [Counter-example of (Vaswani et al., 2019, Lemma 1)] There exists Problem 2.1 such that α_k does not satisfy $\min\{\frac{2(1-c)}{L_{B_k}}, \bar{\alpha}\} \leq \alpha_k$, where $\bar{\alpha}$ is an upper bound of α_k .
- (ii) [Lower bound on learning rate determined by backtracking line search method] If α_k can be computed by Algorithm 2, then there exists a lower bound of α_k such that $0 < \underline{\alpha} := \frac{2\delta(1-c)}{L} \leq \frac{2\delta(1-c)}{L_{B_k}} \leq \alpha_k$.

3 Analysis of SGD using Armijo Line Search

3.1 Convergence analysis of Algorithm 1

Here, we present a convergence analysis of Algorithm 1. The proof of Theorem 3.1 is given in Appendix A.2.

Theorem 3.1 (Upper bound of the squared norm of the full gradient) Consider the sequence $(\theta_k)_{k \in \mathbb{N}}$ generated by Algorithm 1 under Assumption 2.1 for solving Problem 2.1 and suppose that the learning rate $\alpha_k \in [\underline{\alpha}, \bar{\alpha}]$ is computed by Algorithm 2. Then, for all $K \geq 1$,

$$\min_{k \in [0:K-1]} \mathbb{E} [\|\nabla f(\theta_k)\|^2] \leq \underbrace{\frac{LC(\theta_0, f_*)}{\delta(1-c)(2-L_n\bar{\alpha})}}_{B(\theta_0, K)} \frac{1}{K} + \underbrace{\frac{L_n\bar{\alpha}\sigma^2}{2-L_n\bar{\alpha}} \frac{1}{b}}_{V(\sigma^2, b)},$$

where $\delta, c \in (0, 1)$, $L := \max_{i \in [n]} L_i$, $L_n := \frac{1}{n} \sum_{i \in [n]} L_i$ ($\leq L$), $f_* := \frac{1}{n} \sum_{i \in [n]} f_{i,*}$, $\underline{\alpha} := \frac{2\delta(1-c)}{L}$ ($< \frac{2}{L_n}$), $\bar{\alpha} < \frac{2}{L_n}$, and $C(\theta_0, f_*)$ is defined by $\mathbb{E}[f(\theta_0) - f_*]$ or $\mathbb{E}[f(\theta_0) + f(\theta_{k_0+1}) - 2f_*]$ for some $k_0 \in [0:K-1]$.

Theorem 3.1 indicates that the upper bound of the minimum value of $\mathbb{E}[\|\nabla f(\theta_k)\|^2]$ consists of a bias term $B(\theta_0, K)$ and variance term $V(\sigma^2, b)$. When the number of steps K is large and the batch size b is large, $B(\theta_0, K)$ and $V(\sigma^2, b)$ become small. Therefore, we need to set K large and b large so that Algorithm 1 will approximate a local minimizer of f .

Here, we compare Theorem 3.1 with the convergence analysis of SGD using a constant learning rate. SGD using a constant learning rate $\alpha \in (0, \frac{2}{L_n})$ satisfies that

$$\min_{k \in [0:K-1]} \mathbb{E} [\|\nabla f(\theta_k)\|^2] \leq \frac{2(f(\theta_0) - f_*)}{(2 - L_n\alpha)\alpha} \frac{1}{K} + \frac{L_n\alpha\sigma^2}{2 - L_n\alpha} \frac{1}{b} \quad (10)$$

(The proof of (10) is given in Appendix A.5). We need to set a constant learning rate $\alpha \in (0, \frac{2}{L_n})$ depending on the Lipschitz constant L_n of ∇f . However, since computing L_n is NP-hard (Virmaux & Scaman, 2018), it is difficult to set $\alpha \in (0, \frac{2}{L_n})$. Meanwhile, it is sufficient to set $c, \delta \in (0, 1)$ in Algorithms 1 and 2 without computing the Lipschitz constant of ∇f .

We also compare Theorem 3.1 with Theorem 3 in (Vaswani et al., 2019). Theorem 3 in (Vaswani et al., 2019) indicates that, under a strong growth condition with a constant ρ (i.e., $\mathbb{E}_i[\|\nabla f_i(\theta)\|^2] \leq \rho \|\nabla f(\theta)\|^2$

($\theta \in \mathbb{R}^d$)) and the Armijo condition, SGD satisfies that

$$\min_{k \in [0:K-1]} \mathbb{E} [\|\nabla f(\theta_k)\|^2] \leq \frac{f(\theta_0) - f(\theta^*)}{\Delta K},$$

where $c > 1 - \frac{L}{\rho L_n}$, $\bar{\alpha} < \frac{2}{\rho L_n}$, $\Delta := (\bar{\alpha} + \frac{2(1-c)}{L}) - \rho(\bar{\alpha} - \frac{2(1-c)}{L} + L_n \bar{\alpha}^2)$, and θ^* is a local minimizer of f . Theorem 3.1 is a convergence analysis of Algorithm 1 without assuming the strong growth condition or limiting the hyperparameter c . Moreover, Theorem 3.1 shows that using large batch size is appropriate for SGD using the Armijo line search (Algorithm 1).

3.2 Steps needed for ϵ -approximation

To investigate the relationship between the number of steps K needed for nonconvex optimization and the batch size b , we consider an ϵ -approximation of Algorithm 1 defined as follows:

$$\min_{k \in [0:K-1]} \mathbb{E} [\|\nabla f(\theta_k)\|^2] \leq \epsilon^2, \quad (11)$$

where $\epsilon > 0$ is the precision.

Theorem 3.1 leads to the following theorem indicating the relationship between b and the values of K that achieves an ϵ -approximation. The proof of Theorem 3.2 is given in Appendix A.3.

Theorem 3.2 (Steps needed for nonconvex optimization of SGD using Armijo line search)

Suppose that the assumptions in Theorem 3.1 hold. Define $K: \mathbb{R} \rightarrow \mathbb{R}$ for all $b > \frac{C_2}{\epsilon^2}$ by

$$K(b) = \frac{C_1 b}{\epsilon^2 b - C_2}, \quad (12)$$

where the positive constants C_1 and C_2 are defined as in Theorem 3.1. Then, the following hold:

- (i) [Steps needed for nonconvex optimization] K defined by (12) achieves an ϵ -approximation (11).
- (ii) [Properties of the steps] K defined by (12) is monotone decreasing and convex for $b > \frac{C_2}{\epsilon^2}$.

Theorem 3.2 ensures that the number of steps K needed for SGD using the Armijo line search to be an ϵ -approximation is small when the batch size b is large. Therefore, it is useful to set a sufficiently large batch size in the sense of minimizing the steps needed for an ϵ -approximation of SGD using the Armijo line search.

3.3 Critical batch size minimizing SFO complexity

The following theorem shows the existence of a critical batch size for SGD using the Armijo line search. The proof of Theorem 3.3 is given in Appendix A.4.

Theorem 3.3 (Existence of critical batch size for SGD using Armijo line search) *Suppose that the assumptions in Theorem 3.1 hold. Define SFO complexity $N: \mathbb{R} \rightarrow \mathbb{R}$ for the number of steps K , defined by (12), needed for an ϵ -approximation (11) and for a batch size $b > \frac{C_2}{\epsilon^2}$ by*

$$N(b) = K(b)b = \frac{C_1 b^2}{\epsilon^2 b - C_2}, \quad (13)$$

where the positive constants C_1 and C_2 are defined as in Theorem 3.1. Then, the following hold:

- (i) [SFO complexity] N defined by (13) is convex for $b > \frac{C_2}{\epsilon^2}$.
- (ii) [Critical batch size] *There exists a critical batch size*

$$b^* = \frac{2C_2}{\epsilon^2} = \frac{2L_n \bar{\alpha} \sigma^2}{(2 - L_n \bar{\alpha}) \epsilon^2} \quad (14)$$

such that b^ minimizes the SFO complexity (13).*

(iii) [Upper bound on critical batch size] *The critical batch size b^* defined by (14) satisfies*

$$b^* \leq \frac{2\delta(1-c)\bar{\alpha}\sigma^2}{\{\bar{\alpha} - \delta(1-c)\bar{\alpha}\}\epsilon^2}. \quad (15)$$

The previous results in (Shallue et al., 2019; Zhang et al., 2019; Iiduka, 2022) show that, for deep-learning optimizers, there are critical batch sizes at which the SFO complexities are minimized. We are interested in verifying whether a critical batch size exists for SGD using the Armijo line search. Theorem 3.3(iii) indicates that an upper bound on the critical batch size can be obtained from some hyperparameters. Accordingly, we would like to estimate the critical batch size using the upper bound (15). Therefore, the next section numerically examines the relationship between the batch size b and the number of steps K needed for nonconvex optimization and the relationship between b and the SFO complexity N to check if there is a critical batch size b^* minimizing N and if the critical batch size b^* can be estimated from our theoretical results.

4 Numerical Results

The numerical results in (Vaswani et al., 2019) showed that SGD with the Armijo-line-search learning rate performs better than other optimizers in training deep neural networks. Hence, here, we verify whether the numerical results match our theoretical results (Theorems 3.2 and 3.3), that is, the relationship between K and b and the relationship between N and b for Algorithm 1.

The metrics were the number of steps K and the SFO complexity $N = Kb$ indicating that the training accuracy is higher than a certain score. We used Algorithm 1 with the Armijo-line-search learning rate computed by Algorithm 2 with $\gamma = 2$, $\delta = 0.9$, $\bar{\alpha} = 10$ (see <https://github.com/IssamLaradji/sls> for the setting of parameters), and c ($= 0.05, 0.1, 0.15, 0.2, 0.25, 0.3$), where n is the number of training samples. The stopping condition was 200 epochs. The experimental environment consisted of eight NVIDIA DGX A100 GPUs and two Dual AMD Rome7742 2.25-GHz, 128-Core CPUs. The software environment was Python 3.8.2, Pytorch 1.6.0, and CUDA 11.6. The code is available at https://anonymous.4open.science/r/armijo_linesearch-C1C3.

4.1 Training ResNet and MLP on the CIFAR-10, CIFAR-100, and MNIST datasets

We trained ResNet-34 on the CIFAR-10 dataset ($n = 50000$). Figure 1 plots the number of steps needed for the training accuracy to be more than 0.99 for Algorithm 1 versus batch size. It can be seen that Algorithm 1 decreases the number of steps as the batch size increases. Figure 2 plots the SFO complexities of Algorithm 1 versus the batch size. It indicates that there are critical batch sizes that minimize the SFO complexities.

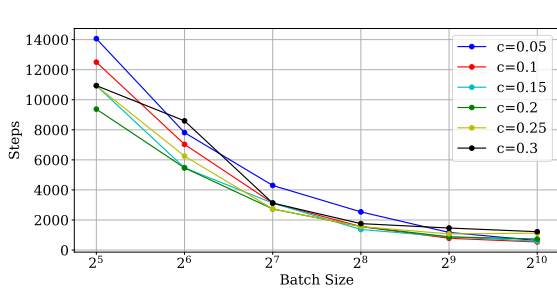


Figure 1: Number of steps for Algorithm 1 versus batch size needed to train ResNet-34 on CIFAR-10

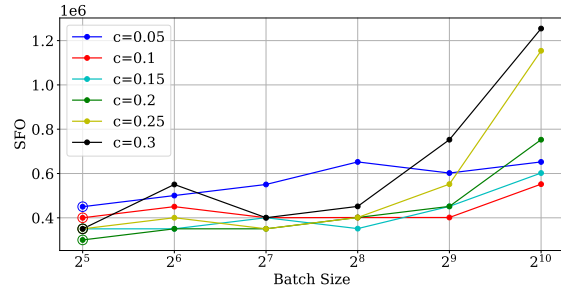


Figure 2: SFO complexity for Algorithm 1 versus batch size needed to train ResNet-34 on CIFAR-10 (The double-circle symbol denotes the measured critical batch size)

We also considered the case of training ResNet-34 on the CIFAR-100 dataset ($n = 50000$). Figure 3 plots the number of steps needed for the training accuracy to be more than 0.99 for Algorithm 1 versus the batch

size, and Figure 4 plots the SFO complexities of Algorithm 1 versus the batch size. As in Figures 1 and 2, these figures show that Algorithm 1 decreases the number of steps as the batch size increases and there are critical batch sizes that minimize the SFO complexities.

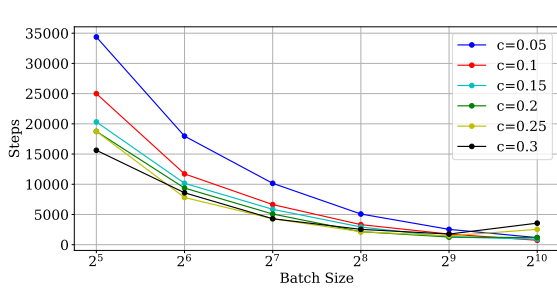


Figure 3: Number of steps for Algorithm 1 versus batch size needed to train ResNet-34 on CIFAR-100

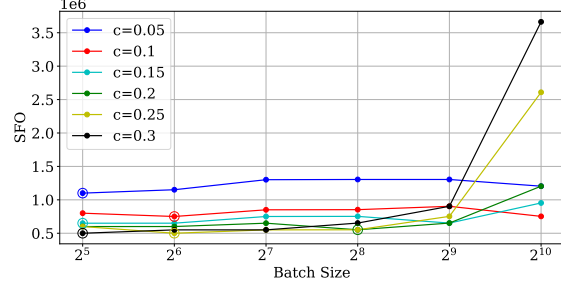


Figure 4: SFO complexity for Algorithm 1 versus batch size needed to train ResNet-34 on CIFAR-100 (The double-circle symbol denotes the measured critical batch size)

We trained a two-hidden-layer MLP with widths of 512 and 256 on the MNIST dataset ($n = 60000$). Figure 5 plots the number of steps needed for the training accuracy to be more than 0.97 for Algorithm 1 versus the batch size, and Figure 6 plots the SFO complexities of Algorithm 1 versus the batch size. As in Figures 1–4, these figures show that Algorithm 1 decreases the number of steps as the batch size increases and there are critical batch sizes that minimize the SFO complexities.

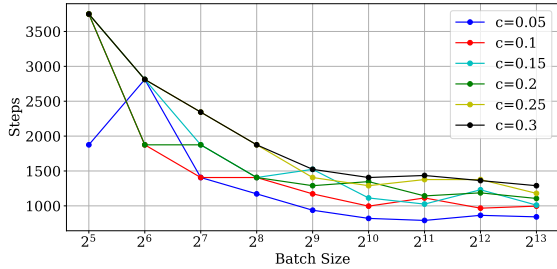


Figure 5: Number of steps for Algorithm 1 versus batch size needed to train MLP on MNIST

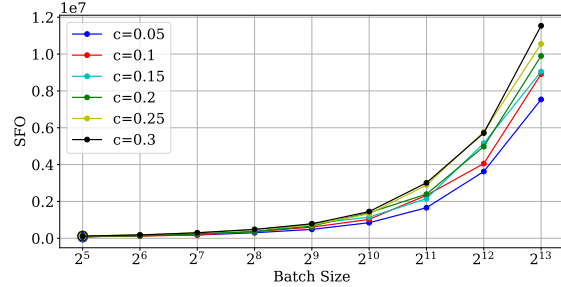


Figure 6: SFO complexity for Algorithm 1 versus batch size needed to train MLP on MNIST (The double-circle symbol denotes the measured critical batch size)

4.2 Estimation of critical batch size

We estimated the critical batch size by using Theorem 3.3(iii) and the ideas presented in (Iiduka, 2022) and (Sato & Iiduka, 2023). We used Algorithm 1 with $c = 0.05$ for training ResNet-34 on the CIFAR-100 dataset (Figures 3 and 4). Theorem 3.3(iii) indicates that the upper bound of the critical batch size involves the unknown value σ^2 . We checked that the Armijo-line-search learning rates for Algorithm 1 with $c = 0.05$ are about 10 (see also (Vaswani et al., 2019, Figure 5 (Left))). Hence, we used $\underline{\alpha} \approx \bar{\alpha} \approx 10$. We estimated the unknown value $X = \frac{\sigma^2}{\epsilon^2}$ in the upper bound (15) of the critical batch size by using $\delta = 0.9$, $b^* = 2^5$ (see Figure 4), and $\underline{\alpha} \approx \bar{\alpha} \approx 10$ as follows:

$$b^* \approx \frac{2\delta(1-c)\bar{\alpha}\sigma^2}{\{\underline{\alpha} - \delta(1-c)\bar{\alpha}\}\epsilon^2} \approx 2\delta(1-c)X, \text{ i.e., } X \approx 18.71.$$

Let us estimate the critical batch size using $X \approx 18.71$ and Theorem 3.3(iii). For example, when using Algorithm 1 with $c = 0.25$ for training ResNet-34 on the CIFAR-100 dataset, the upper bound of the critical batch size is

$$\frac{2\delta(1-c)\bar{\alpha}}{\bar{\alpha} - \delta(1-c)\bar{\alpha}}X \approx 77.71 \approx 2^6 = b^*,$$

which implies that the estimated critical batch size 77.71 is close to the measured critical batch size $b^* = 2^6 = 64$ in Figure 4.

5 Conclusion

This paper presented a convergence analysis of SGD using the Armijo line search for nonconvex optimization. The analysis indicates that the upper bound of the expectation of the squared norm of the full gradient becomes smaller as the number of steps and the batch size grow. Moreover, we investigated the relationship between the number of steps and the batch size needed for nonconvex optimization of SGD using the Armijo line search. We showed that the number of steps needed for nonconvex optimization is monotone decreasing and convex with respect to the batch size; i.e., the steps decrease in number as the batch size increases. We also showed that the SFO complexity needed for nonconvex optimization is convex with respect to the batch size and that there exists a critical batch size at which the SFO complexity is minimized. In addition, we gave an upper bound on the critical batch size and showed that it can be estimated by using some parameters. Finally, we provided numerical results that support our theoretical findings. We trained ResNet-34 on the CIFAR-10 and CIFAR-100 datasets and MLP on the MNIST dataset and found that SGD using the Armijo line search decreases the number of steps as the batch size increases. Moreover, we showed that the batch sizes estimated from the upper bound of the critical batch size are close to those of the numerical results.

References

- Amir Beck. *First-Order Methods in Optimization*. Society for Industrial and Applied Mathematics, Philadelphia, PA, 2017.
- Hao Chen, Lili Zheng, Raed AL Kontar, and Garvesh Raskutti. Stochastic gradient descent in correlated settings: A study on Gaussian processes. In *Advances in Neural Information Processing Systems*, volume 33, 2020.
- John Duchi, Elad Hazan, and Yoram Singer. Adaptive subgradient methods for online learning and stochastic optimization. *Journal of Machine Learning Research*, 12:2121–2159, 2011.
- Benjamin Fehrman, Benjamin Gess, and Arnulf Jentzen. Convergence rates for the stochastic gradient descent method for non-convex objective functions. *Journal of Machine Learning Research*, 21:1–48, 2020.
- Saeed Ghadimi and Guanghui Lan. Optimal stochastic approximation algorithms for strongly convex stochastic composite optimization I: A generic algorithmic framework. *SIAM Journal on Optimization*, 22:1469–1492, 2012.
- Saeed Ghadimi and Guanghui Lan. Optimal stochastic approximation algorithms for strongly convex stochastic composite optimization II: Shrinking procedures and optimal algorithms. *SIAM Journal on Optimization*, 23:2061–2089, 2013.
- Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. GANs trained by a two time-scale update rule converge to a local nash equilibrium. In *Advances in Neural Information Processing Systems*, volume 30, 2017.
- Hideaki Iiduka. Critical batch size minimizes stochastic first-order oracle complexity of deep learning optimizer using hyperparameters close to one. arXiv: 2208.09814, 2022.

- Prateek Jain, Sham M. Kakade, Rahul Kidambi, Praneeth Netrapalli, and Aaron Sidford. Parallelizing stochastic gradient descent for least squares regression: Mini-batching, averaging, and model misspecification. *Journal of Machine Learning Research*, 18(223):1–42, 2018.
- Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *Proceedings of The International Conference on Learning Representations*, 2015.
- Nicolas Loizou, Sharan Vaswani, Issam Laradji, and Simon Lacoste-Julien. 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*, volume 130, 2021.
- Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *Proceedings of The International Conference on Learning Representations*, 2019.
- Hiroki Naganuma and Hideaki Iiduka. Conjugate gradient method for generative adversarial networks. In *Proceedings of The 26th International Conference on Artificial Intelligence and Statistics*, volume 206 of *Proceedings of Machine Learning Research*, pp. 4381–4408. PMLR, 2023.
- Arkadi Nemirovski, Anatoli Juditsky, Guanghui Lan, and Alexander Shapiro. Robust stochastic approximation approach to stochastic programming. *SIAM Journal on Optimization*, 19:1574–1609, 2009.
- Yurii Nesterov. A method for unconstrained convex minimization problem with the rate of convergence $O(1/k^2)$. *Doklady AN USSR*, 269:543–547, 1983.
- J. Nocedal and S. J. Wright. *Numerical Optimization*. Springer Series in Operations Research and Financial Engineering. Springer, New York, 2nd edition, 2006.
- Boris T. Polyak. Some methods of speeding up the convergence of iteration methods. *USSR Computational Mathematics and Mathematical Physics*, 4:1–17, 1964.
- Sashank J. Reddi, Satyen Kale, and Sanjiv Kumar. On the convergence of Adam and beyond. In *Proceedings of The International Conference on Learning Representations*, 2018.
- Herbert Robbins and Herbert Monro. A stochastic approximation method. *The Annals of Mathematical Statistics*, 22:400–407, 1951.
- Naoki Sato and Hideaki Iiduka. 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*, volume 202 of *Proceedings of Machine Learning Research*, pp. 30080–30104. PMLR, 23–29 Jul 2023.
- Kevin Scaman and Cédric Malherbe. Robustness analysis of non-convex stochastic gradient descent using biased expectations. In *Advances in Neural Information Processing Systems*, volume 33, 2020.
- Christopher J. Shallue, Jaehoon Lee, Joseph Antognini, Jascha Sohl-Dickstein, Roy Frostig, and George E. Dahl. Measuring the effects of data parallelism on neural network training. *Journal of Machine Learning Research*, 20:1–49, 2019.
- Tijmen Tieleman and Geoffrey Hinton. RMSProp: Divide the gradient by a running average of its recent magnitude. *COURSERA: Neural networks for machine learning*, 4:26–31, 2012.
- Sharan Vaswani, Aaron Mishkin, Issam Laradji, Mark Schmidt, Gauthier Gidel, and Simon Lacoste-Julien. Painless stochastic gradient: Interpolation, line-search, and convergence rates. In *Advances in Neural Information Processing Systems*, volume 32, 2019.
- Aladin Virmaux and Kevin Scaman. Lipschitz regularity of deep neural networks: analysis and efficient estimation. In *Advances in Neural Information Processing Systems*, volume 31, 2018.
- Guodong Zhang, Lala Li, Zachary Nado, James Martens, Sushant Sachdeva, George E. Dahl, Christopher J. Shallue, and Roger Grosse. Which algorithmic choices matter at which batch sizes? Insights from a noisy quadratic model. In *Advances in Neural Information Processing Systems*, volume 32, 2019.

Martin Zinkevich. Online convex programming and generalized infinitesimal gradient ascent. In *Proceedings of the 20th International Conference on Machine Learning*, pp. 928–936, 2003.

A Appendix

A.1 Proof of lemma 2.1

(i) Let $k \in \mathbb{N}$ and let L_{B_k} be the Lipschitz constant of ∇f_{B_k} . Lemma 1 in (Vaswani et al., 2019) is as follows:

$$\begin{aligned} & \forall f_{B_k} : \mathbb{R}^d \rightarrow \mathbb{R} \quad \forall c \in (0, 1) \quad \forall \theta_k \in \mathbb{R}^d \quad \forall \bar{\alpha} > 0 \\ & \exists \alpha_k \in (0, \bar{\alpha}] \quad (f_{B_k}(\theta_k - \alpha_k \nabla f_{B_k}(\theta_k)) \leq f_{B_k}(\theta_k) - c\alpha_k \|\nabla f_{B_k}(\theta_k)\|^2) \\ & \Rightarrow \min \left\{ \frac{2(1-c)}{L_{B_k}}, \bar{\alpha} \right\} \leq \alpha_k. \end{aligned} \quad (16)$$

The negative proposition of (16) is as follows:

$$\begin{aligned} & \exists f_{B_k} : \mathbb{R}^d \rightarrow \mathbb{R} \quad \exists c \in (0, 1) \quad \exists \theta_k \in \mathbb{R}^d \quad \exists \bar{\alpha} > 0 \\ & \exists \alpha_k \in (0, \bar{\alpha}] \quad (f_{B_k}(\theta_k - \alpha_k \nabla f_{B_k}(\theta_k)) \leq f_{B_k}(\theta_k) - c\alpha_k \|\nabla f_{B_k}(\theta_k)\|^2) \\ & \wedge \min \left\{ \frac{2(1-c)}{L_{B_k}}, \bar{\alpha} \right\} > \alpha_k. \end{aligned} \quad (17)$$

We will prove that (17) holds. Let $n = b = 1$, $d = 1$, $c = 0.1$, $\bar{\alpha} = 1$, and $f(\theta) = f_{B_k}(\theta) = \theta^2$. From $\nabla f(\theta) = 2\theta$, we have that $L_{B_k} = 2$. Since $\theta^* = 0$ is the global minimizer of f , we set $\theta_k \in \mathbb{R}$ such that $\theta_k \neq \theta^*$. The Armijo condition in this case is such that $(\theta_k - 2\alpha_k \theta_k)^2 \leq \theta_k^2 - c\alpha_k (2\theta_k)^2$, which is equivalent to $\alpha_k \leq 1 - c = 0.9$. Hence,

$$\begin{aligned} & \exists \alpha_k \in (0, 1] \quad (\alpha_k \leq 0.9) \wedge (\min \{0.9, 1\} > \alpha_k) \\ & \Leftrightarrow \exists \alpha_k \in (0, \bar{\alpha}] \quad (\alpha_k \leq 1 - c) \wedge \left(\min \left\{ \frac{2(1-c)}{L_{B_k}}, \bar{\alpha} \right\} > \alpha_k \right) \\ & \Leftrightarrow \exists \alpha_k \in (0, \bar{\alpha}] \quad (f_{B_k}(\theta_k - \alpha_k \nabla f_{B_k}(\theta_k)) \leq f_{B_k}(\theta_k) - c\alpha_k \|\nabla f_{B_k}(\theta_k)\|^2) \\ & \wedge \min \left\{ \frac{2(1-c)}{L_{B_k}}, \bar{\alpha} \right\} > \alpha_k, \end{aligned}$$

which implies that (17) holds.

(ii) Since $\frac{\alpha_k}{\delta}$ does not satisfy the Armijo condition (9), we have that

$$f_{B_k} \left(\theta_k - \frac{\alpha_k}{\delta} \nabla f_{B_k}(\theta_k) \right) > f_{B_k}(\theta_k) - c \frac{\alpha_k}{\delta} \|\nabla f_{B_k}(\theta_k)\|^2. \quad (18)$$

The L_{B_k} -smoothness of f_{B_k} ensures that the descent lemma is true, i.e.,

$$\begin{aligned} & f_{B_k} \left(\theta_k - \frac{\alpha_k}{\delta} \nabla f_{B_k}(\theta_k) \right) \\ & \leq f_{B_k}(\theta_k) + \left\langle \nabla f_{B_k}(\theta_k), \left(\theta_k - \frac{\alpha_k}{\delta} \nabla f_{B_k}(\theta_k) \right) - \theta_k \right\rangle + \frac{L_{B_k}}{2} \left\| \left(\theta_k - \frac{\alpha_k}{\delta} \nabla f_{B_k}(\theta_k) \right) - \theta_k \right\|^2, \end{aligned}$$

which implies that

$$f_{B_k} \left(\theta_k - \frac{\alpha_k}{\delta} \nabla f_{B_k}(\theta_k) \right) \leq f_{B_k}(\theta_k) + \frac{\alpha_k}{\delta} \left(\frac{L_{B_k} \alpha_k}{2\delta} - 1 \right) \|\nabla f_{B_k}(\theta_k)\|^2. \quad (19)$$

Hence, (18) and (19) imply that

$$-c \frac{\alpha_k}{\delta} \|\nabla f_{B_k}(\theta_k)\|^2 \leq \frac{\alpha_k}{\delta} \left(\frac{L_{B_k} \alpha_k}{2\delta} - 1 \right) \|\nabla f_{B_k}(\theta_k)\|^2,$$

which in turn implies that

$$\frac{\alpha_k}{\delta} \left(\frac{L_{B_k} \alpha_k}{2\delta} - (1-c) \right) \|\nabla f_{B_k}(\boldsymbol{\theta}_k)\|^2 \geq 0.$$

Accordingly,

$$\frac{L_{B_k} \alpha_k}{2\delta} - (1-c) \geq 0, \text{ i.e., } \alpha_k \geq \frac{2\delta(1-c)}{L_{B_k}}.$$

From $L_{B_k} = \frac{1}{b} \sum_{i \in [b]} L_{\xi_{k,i}} \leq L := \max_{i \in [n]} L_i$ ($k \in \mathbb{N}$), we also have that $\alpha_k \geq \frac{2\delta(1-c)}{L_{B_k}} \geq \frac{2\delta(1-c)}{L}$.

A.2 Proof of Theorem 3.1

The definition of $f(\boldsymbol{\theta}) := \frac{1}{n} \sum_{i \in [n]} f_i(\boldsymbol{\theta})$ and the L_i -smoothness of f_i ($i \in [n]$) imply that, for all $\boldsymbol{\theta}_1, \boldsymbol{\theta}_2 \in \mathbb{R}^d$,

$$\|\nabla f(\boldsymbol{\theta}_1) - \nabla f(\boldsymbol{\theta}_2)\| \leq \frac{1}{n} \sum_{i \in [n]} \|\nabla f_i(\boldsymbol{\theta}_1) - \nabla f_i(\boldsymbol{\theta}_2)\| \leq \frac{\sum_{i \in [n]} L_i}{n} \|\boldsymbol{\theta}_1 - \boldsymbol{\theta}_2\|, \quad (20)$$

which in turn implies that ∇f is Lipschitz continuous with Lipschitz constant $L_n := \frac{1}{n} \sum_{i \in [n]} L_i$. Hence, the descent lemma ensures that, for all $k \in \mathbb{N}$,

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

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

$$f(\boldsymbol{\theta}_{k+1}) \leq f(\boldsymbol{\theta}_k) - \alpha_k \langle \nabla f(\boldsymbol{\theta}_k), \nabla f_{B_k}(\boldsymbol{\theta}_k) \rangle + \frac{L_n \alpha_k^2}{2} \|\nabla f_{B_k}(\boldsymbol{\theta}_k)\|^2. \quad (21)$$

From $0 < \alpha_k \leq \bar{\alpha}$, we have that, for all $k \in \mathbb{N}$,

$$\langle \nabla f(\boldsymbol{\theta}_k), \nabla f_{B_k}(\boldsymbol{\theta}_k) \rangle \leq \frac{f(\boldsymbol{\theta}_k) - f(\boldsymbol{\theta}_{k+1})}{\alpha_k} + \frac{L_n \bar{\alpha}}{2} \|\nabla f_{B_k}(\boldsymbol{\theta}_k)\|^2. \quad (22)$$

Assumption 2.1 guarantees that

$$\mathbb{E} [\nabla f_{B_k}(\boldsymbol{\theta}_k) | \boldsymbol{\theta}_k] = \nabla f(\boldsymbol{\theta}_k) \text{ and } \mathbb{E} [\|\nabla f_{B_k}(\boldsymbol{\theta}_k) - \nabla f(\boldsymbol{\theta}_k)\|^2 | \boldsymbol{\theta}_k] \leq \frac{\sigma^2}{b}. \quad (23)$$

Hence, we have

$$\begin{aligned} & \mathbb{E} [\|\nabla f_{B_k}(\boldsymbol{\theta}_k)\|^2 | \boldsymbol{\theta}_k] \\ &= \mathbb{E} [\|\nabla f_{B_k}(\boldsymbol{\theta}_k) - \nabla f(\boldsymbol{\theta}_k) + \nabla f(\boldsymbol{\theta}_k)\|^2 | \boldsymbol{\theta}_k] \\ &= \mathbb{E} [\|\nabla f_{B_k}(\boldsymbol{\theta}_k) - \nabla f(\boldsymbol{\theta}_k)\|^2 | \boldsymbol{\theta}_k] + 2\mathbb{E} [\langle \nabla f_{B_k}(\boldsymbol{\theta}_k) - \nabla f(\boldsymbol{\theta}_k), \nabla f(\boldsymbol{\theta}_k) \rangle | \boldsymbol{\theta}_k] + \mathbb{E} [\|\nabla f(\boldsymbol{\theta}_k)\|^2 | \boldsymbol{\theta}_k] \\ &\leq \frac{\sigma^2}{b} + \mathbb{E} [\|\nabla f(\boldsymbol{\theta}_k)\|^2 | \boldsymbol{\theta}_k]. \end{aligned} \quad (24)$$

Taking the total expectation on both sides of (22) thus ensures that, for all $k \in \mathbb{N}$,

$$\mathbb{E} [\|\nabla f(\boldsymbol{\theta}_k)\|^2] \leq \mathbb{E} \left[\frac{f(\boldsymbol{\theta}_k) - f(\boldsymbol{\theta}_{k+1})}{\alpha_k} \right] + \frac{L_n \bar{\alpha}}{2} \left(\mathbb{E} [\|\nabla f(\boldsymbol{\theta}_k)\|^2] + \frac{\sigma^2}{b} \right). \quad (25)$$

Let $L := \max_{i \in [n]} L_i$. A discussion similar to the one showing (20) leads to the finding that the Lipschitz constant of ∇f_{B_k} is $L_{B_k} = \frac{1}{b} \sum_{i \in [b]} L_{\xi_{k,i}} \leq L$ ($k \in \mathbb{N}$). Moreover, Lemma 2.1(ii) implies that there exists

$\underline{\alpha} > 0$ such that, for all $k \in \mathbb{N}$, $\underline{\alpha} := \frac{2\delta(1-c)}{L} \leq \alpha_k$. Let $K \geq 1$. Consider the case where, for all $k \in [0 : K-1]$, $f(\boldsymbol{\theta}_k) - f(\boldsymbol{\theta}_{k+1}) \geq 0$. Then, (25) guarantees that, for all $k \in [0 : K-1]$,

$$\mathbb{E} [\|\nabla f(\boldsymbol{\theta}_k)\|^2] \leq \frac{L}{2\delta(1-c)} \mathbb{E} [f(\boldsymbol{\theta}_k) - f(\boldsymbol{\theta}_{k+1})] + \frac{L_n \bar{\alpha}}{2} \left(\mathbb{E} [\|\nabla f(\boldsymbol{\theta}_k)\|^2] + \frac{\sigma^2}{b} \right),$$

which implies that

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

Summing the above inequalities from $k = 0$ to $k = K-1$ ensures that

$$\left(1 - \frac{L_n \bar{\alpha}}{2}\right) \sum_{k=0}^{K-1} \mathbb{E} [\|\nabla f(\boldsymbol{\theta}_k)\|^2] \leq \frac{L}{2\delta(1-c)} \mathbb{E} [f(\boldsymbol{\theta}_0) - f(\boldsymbol{\theta}_K)] + \frac{L_n \bar{\alpha} \sigma^2 K}{2b}. \quad (26)$$

Since f is bounded below by $f_* := \frac{1}{n} \sum_{i \in [n]} f_{i,*}$ and $\bar{\alpha} < \frac{2}{L_n}$, we have

$$\frac{1}{K} \sum_{k=0}^{K-1} \mathbb{E} [\|\nabla f(\boldsymbol{\theta}_k)\|^2] \leq \underbrace{\frac{L \mathbb{E} [f(\boldsymbol{\theta}_0) - f_*]}{\delta(1-c)(2 - L_n \bar{\alpha})K}}_{A_K} + \frac{L_n \bar{\alpha} \sigma^2}{(2 - L_n \bar{\alpha})b}.$$

Consider the case where there exists $k_0 \in [0 : K-1]$ such that $f(\boldsymbol{\theta}_{k_0}) - f(\boldsymbol{\theta}_{k_0+1}) < 0$. If $f(\boldsymbol{\theta}_k) - f(\boldsymbol{\theta}_{k+1}) < 0$ holds for all $k \in [0 : K-1]$, (25) ensures that

$$\mathbb{E} [\|\nabla f(\boldsymbol{\theta}_k)\|^2] < \frac{L_n \bar{\alpha}}{2} \left(\mathbb{E} [\|\nabla f(\boldsymbol{\theta}_k)\|^2] + \frac{\sigma^2}{b} \right),$$

which, together with $\bar{\alpha} < \frac{2}{L_n}$ and $A_K \geq 0$, implies that

$$\frac{1}{K} \sum_{k=0}^{K-1} \mathbb{E} [\|\nabla f(\boldsymbol{\theta}_k)\|^2] < \frac{L_n \bar{\alpha} \sigma^2}{(2 - L_n \bar{\alpha})b} \leq A_K + \frac{L_n \bar{\alpha} \sigma^2}{(2 - L_n \bar{\alpha})b}.$$

Hence, we may assume without loss of generality that, for all $k \neq k_0$, $f(\boldsymbol{\theta}_k) - f(\boldsymbol{\theta}_{k+1}) \geq 0$. A discussion similar to the one showing (26) implies that

$$\left(1 - \frac{L_n \bar{\alpha}}{2}\right) \sum_{k=0}^{K-1} \mathbb{E} [\|\nabla f(\boldsymbol{\theta}_k)\|^2] \leq \frac{L}{2\delta(1-c)} \mathbb{E} [f(\boldsymbol{\theta}_0) - f(\boldsymbol{\theta}_{k_0}) + f(\boldsymbol{\theta}_{k_0+1}) - f(\boldsymbol{\theta}_K)] + \frac{L_n \bar{\alpha} \sigma^2 K}{2b}.$$

Since f is bounded below by $f_* := \frac{1}{n} \sum_{i \in [n]} f_{i,*}$ and $\bar{\alpha} < \frac{2}{L_n}$, we have

$$\frac{1}{K} \sum_{k=0}^{K-1} \mathbb{E} [\|\nabla f(\boldsymbol{\theta}_k)\|^2] \leq \frac{L \mathbb{E} [f(\boldsymbol{\theta}_0) + f(\boldsymbol{\theta}_{k_0+1}) - 2f_*]}{\delta(1-c)(2 - L_n \bar{\alpha})K} + \frac{L_n \bar{\alpha} \sigma^2}{(2 - L_n \bar{\alpha})b}.$$

Moreover, since we have

$$\min_{k \in [0:K-1]} \mathbb{E} [\|\nabla f(\boldsymbol{\theta}_k)\|^2] \leq \frac{1}{K} \sum_{k=0}^{K-1} \mathbb{E} [\|\nabla f(\boldsymbol{\theta}_k)\|^2],$$

the assertion in Theorem 3.1 holds.

A.3 Proof of Theorem 3.2

(i) We have

$$\frac{C_1}{K} + \frac{C_2}{b} = \epsilon^2$$

is equivalent to

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

Hence, Theorem 3.1 leads to an ϵ -approximation.

(ii) We have

$$\frac{dK(b)}{db} = \frac{-C_1 C_2}{(\epsilon^2 b - C_2)^2} \leq 0 \text{ and } \frac{d^2 K(b)}{db^2} = \frac{2C_1 C_2 \epsilon^2}{(\epsilon^2 b - C_2)^3} \geq 0,$$

which implies that K is monotone decreasing and convex with respect to b .

A.4 Proof of Theorem 3.3

(i) From

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

we have

$$\frac{dN(b)}{db} = \frac{C_1 b(\epsilon^2 b - 2C_2)}{(\epsilon^2 b - C_2)^2} \text{ and } \frac{d^2 N(b)}{db^2} = \frac{2C_1 C_2^2}{(\epsilon^2 b - C_2)^3} \geq 0,$$

which implies that N is convex with respect to b .

(ii) We have

$$\frac{dN(b)}{db} \begin{cases} < 0 & \text{if } b < b^*, \\ = 0 & \text{if } b = b^* = \frac{2C_2}{\epsilon^2}, \\ > 0 & \text{if } b > b^*. \end{cases}$$

Hence, the point b^* minimizes N .

(iii) Lemma 2.1(iii) ensures that

$$L_n := \frac{1}{n} \sum_{i \in [n]} L_i \leq L = \frac{2\delta(1-c)}{\underline{\alpha}}.$$

Hence,

$$b^* = \frac{2C_2}{\epsilon^2} = \frac{2L_n \bar{\alpha} \sigma^2}{(2 - L_n \bar{\alpha}) \epsilon^2} \leq \frac{4\delta(1-c) \bar{\alpha} \sigma^2}{\underline{\alpha}} \frac{\underline{\alpha}}{2\{\underline{\alpha} - \delta(1-c) \bar{\alpha}\} \epsilon^2} = \frac{2\delta(1-c) \bar{\alpha} \sigma^2}{\{\underline{\alpha} - \delta(1-c) \bar{\alpha}\} \epsilon^2}.$$

A.5 Proof of (10)

Let $K \geq 1$. From (21) and $\alpha_k := \alpha > 0$, we have that, for all $k \in \mathbb{N}$,

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

Hence, (23) and (24) ensure that, for all $k \in \mathbb{N}$,

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

which implies that, for all $k \in \mathbb{N}$,

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

Summing the above inequalities from $k = 0$ to $k = K - 1$ ensures that

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

Since f is bounded below by $f_* := \frac{1}{n} \sum_{i \in [n]} f_{i,*}$, we have

$$\min_{k \in [0:K-1]} \mathbb{E} [\|\nabla f(\boldsymbol{\theta}_k)\|^2] \leq \frac{1}{K} \sum_{k=0}^{K-1} \mathbb{E} [\|\nabla f(\boldsymbol{\theta}_k)\|^2] \leq \frac{2\mathbb{E} [f(\boldsymbol{\theta}_0) - f_*]}{\alpha(2 - L_n \alpha)K} + \frac{L_n \alpha \sigma^2}{(2 - L_n \alpha)b}.$$