Power of Generalized Smoothness in Stochastic Convex Optimization: First- and Zero-Order Algorithms

Aleksandr Lobanov MIPT, Skoltech, HSE lobbsasha@mail.ru Alexander Gasnikov Innopolis, MIPT, ISP RAS gasnikov@yandex.ru

Abstract

This paper is devoted to the study of stochastic optimization problems under the generalized smoothness assumption. By considering the unbiased gradient oracle in *Stochastic Gradient Descent*, we provide strategies to achieve in bounds the summands describing linear rate. In particular, in the case $L_0 = 0$, we obtain in the **convex setup** the iteration complexity: $N = O\left(L_1R\log\frac{1}{\varepsilon} + \frac{L_1CR^2}{\varepsilon}\right)$ for *Clipped Stochastic Gradient Descent* and $N = O\left(L_1R\log\frac{1}{\varepsilon}\right)$ for *Normalized Stochastic Gradient Descent* and $N = O\left(L_1R\log\frac{1}{\varepsilon}\right)$ for *Normalized Stochastic Gradient Descent*. Furthermore, we generalize the convergence results to the case with a biased gradient oracle, and show that the power of (L_0, L_1) -smoothness extends to *zero-order algorithms*. Finally, we demonstrate the possibility of linear convergence in the convex setup through numerical experimentation, which has aroused some interest in the machine learning community.

1 Introduction

In many real-world scenarios, systems are often noisy and complex, making deterministic optimization infeasible. Therefore, this work focuses on a stochastic optimization problem:

$$\min_{x \in \mathbb{P}^d} \left\{ f(x) := \mathbb{E}_{\xi \sim \mathcal{D}} \left[f(x, \xi) \right] \right\},\tag{1}$$

where $f : \mathbb{R}^d \to \mathbb{R}$ is a convex function and where we assume that optimization algorithms only have access to the gradient oracle $\mathbf{g} : \mathbb{R}^d \times \mathcal{D} \to \mathbb{R}^d$ with stochastic gradient $\mathbb{E}\left[\nabla f(x,\xi)\right] = \nabla f(x)$ and bias $\mathbf{b}(x)$ terms:

$$\mathbf{g}(x,\xi) = \nabla f(x,\xi) + \mathbf{b}(x). \tag{2}$$

Frequently, to solve problem (1) one uses what is likely already a classic optimization algorithm, namely Stochastic Gradient Descent (SGD) [10] or its variations, which have demonstrated their effectiveness in different settings, for instance, federated learning [60, 33, 58], deep learning [15, 64, 18], reinforcement learning [7, 39] and others. Among the variants of SGD, it is worth noting the Normalized Stochastic Gradient Descent (NSGD) [27, 66] which has received widely attention from the community because it addresses challenges in optimization for machine learning [8]. And it's also worth noting the Clipped Stochastic Gradient Descent (ClipSGD) [24], which is commonly used to stabilize the training of deep learning models [48, 25].

Many standard literatures analyze stochastic optimization algorithms with unbiased gradient oracle (2). In particular, SGD [36, 11], NSGD [65, 30], ClipSGD [25, 34]. However, there are a number of applications where gradient oracle (2) is biased. For example, sparsified SGD [4], delayed SGD [53], etc. Zero-order algorithms [46, 16] occupy a special place in the class of stochastic methods with biased gradient oracle (2). They are motivated by various applications, including multi-armed bandit [52, 38], online optimization [1, 6, 3], hyperparameter tuning [28, 47].

Preprint. Under review.

In our work, we investigate the convergence of first-order algorithms: ClipSGD, NSGD, and zeroorder algorithms: ZO-ClipSGD, ZO-NSGD, assuming convexity and (L_0, L_1) -smoothness.

We emphasize the following points:

Algorithm step size. Zero-order algorithms do not have access to the exact (stochastic) gradient in particular, as well as every algorithm with a biased gradient oracle, so we focus on creating first-order methods whose step size does not depend on knowledge of the gradient at a given point. We use the developed first-order algorithms as a basis for creating zero-order methods [22].

Linear convergence. Historically [45], stochastic optimization first- and zero-order algorithms have achieved the desired accuracy with a linear rate of convergence only in strongly convex case and under assumption of standard smoothness. However, the work of [44] showed that if the generalized smoothness assumption is satisfied in a <u>deterministic convex</u> optimization problem, then gradient descent has two regimes: linear convergence rate as long as $\|\nabla f(x^k)\| \geq \frac{L_0}{L_1}$, and a sublinear convergence rate in the other case (see the example of the power of



in the other case (see the example of the power of Figure 1: Changing regimes demonstration norm function in Figure 1). Considering these points, our work answers the following question:

Can linear convergence rate in stochastic convex optimization be achieved for first- and zero-order algorithms with constant step size?

1.1 Main Contributions

More specifically, our contributions are the following:

- We provide strategies to obtain summands that describe the linear convergence rate. In particular, we show that using clipping or normalization techniques can achieve the desired results.
- We improve convergence results for ClipSGD and NSGD with unbiased gradient oracle (2) in the convex setting assuming (L_0, L_1) -smoothness (see Table 1). Moreover, we show that in the case $L_0 = 0$, NSGD can converge in the convex setup with a linear convergence rate to the desired accuracy, requiring $N = \tilde{\mathcal{O}} (L_1 R)$ iterations and $B = \mathcal{O} \left(\frac{\sigma^2 M R^3}{\varepsilon^3} \right)$ batch size.
- We generalize ClipSGD, and NSGD to the case of a biased gradient oracle, showing how the bias accumulates over iterations (this result may be of independent interest).
- We provide the first convergence results for the zero-order algorithms ZO-ClipSGD (Algorithm 3), and ZO-NSGD (Algorithm 4) in the convex and (L_0, L_1) -smooth setting. We show that the power of generalized smoothness extends to zero-order methods as well, achieving summands characterizing the linear convergence rate (see Table 1).
- We demonstrate on a numerical example of logistic regression (which is of particular interest to the machine learning community) that indeed, zero- and first-order stochastic algorithms can converge with linear rates in a convex setup.

1.2 Formal Setting and Assumptions

In this subsection, we introduce and discuss main assumptions and notations used throughout paper.

Notations. We use $\langle x, y \rangle := \sum_{i=1}^{d} x_i y_i$ to denote standard inner product of $x, y \in \mathbb{R}^d$. We denote Euclidean norm in \mathbb{R}^d as $\|x\| := \sqrt{\sum_{i=1}^{d} x_i^2}$. In particular, this norm $\|x\| := \sqrt{\langle x, x \rangle}$ is related to the inner product. We use $\mathcal{P}[\cdot]$ to define probability measure which is always known from the context, $\mathbb{E}[\cdot]$ denotes mathematical expectation. We use the following notation $B^d(r) := \{x \in \mathbb{R}^d : \|x\| \le r\}$ to denote Euclidean ball $(l_2$ -ball) and $S^d(r) := \{x \in \mathbb{R}^d : \|x\| = r\}$ to denote Euclidean sphere. We denote $0 \le M < \infty$ as the upper bound of the gradient norm $\|\nabla f(x^k)\|$. For simplicity, we denote $f^* := f(x^*)$ and $R = \|x^0 - x^*\|$. We use $\tilde{O}(\cdot)$ to hide the logarithmic coefficients.

Table 1: Comparison of convergence results of SGD variants to the most related work [20] in the convex and (L_0, L_1) -smooth setup. Notation: $\eta \leq (L_0 + L_1 c)^{-1}$ – step size; c > 0 – clipping radius; $\mathcal{R} = \left(\eta + \frac{MR}{c^2} + \frac{R}{c}\right)$; ε = desired accuracy; d = dimension; SLCR = summand with linear convergence rate.

| Algorithm | Number of Iterations $\#N$ | Batch Size #B | Maximum Noise Level $\#\Delta$ | SLCR? | Reference |
|------------|---|---|---|-------|--------------------|
| | $\mathcal{O}\left(\frac{L_0R^2}{\varepsilon} + \frac{\sigma^2R^2}{\varepsilon^2} + L_1^2R^2\right)$ | × | × | × | Gaash et al. [20] |
| ClipSGD | $\mathcal{O}\left(\frac{R}{\eta c}\log\frac{1}{\varepsilon}+\frac{R^2}{\eta \varepsilon}\right)$ | $\mathcal{O}\left(\frac{\sigma^2 \mathcal{R}}{\varepsilon}\right)$ | × | 1 | Theorem 3.1 (Ours) |
| NSGD | $\mathcal{O}\left(\left(L_1R + \frac{L_0R^2}{\varepsilon}\right)\log\frac{1}{\varepsilon}\right)$ | $\mathcal{O}\left(\frac{\sigma^2 M R^3}{\varepsilon^3}\right)$ | × | 1 | Theorem 4.1 (Ours) |
| ZO-ClipSGD | $\mathcal{O}\left(\frac{R}{\eta c}\log \frac{1}{\varepsilon} + \frac{R^2}{\eta \varepsilon}\right)$ | $O\left(\frac{dMR\tilde{\sigma}^2}{\varepsilon c^2}\right)$ | $\mathcal{O}\left(\frac{\varepsilon}{\sqrt{dR(L_0+L_1M)}}\min\left\{\tilde{\sigma},\frac{\varepsilon}{\sqrt{dR}}\right\}\right)$ | 1 | Theorem 5.2 (Ours) |
| ZO-NSGD | $\mathcal{O}\left(\left(L_1R + \frac{L_0R^2}{\varepsilon}\right)\log\frac{1}{\varepsilon}\right)$ | $\mathcal{O}\left(\frac{dMR^{3}\tilde{\sigma}^{2}}{\varepsilon^{3}}\right)$ | $-\mathcal{O}\left(\frac{\varepsilon^{3/2}}{\sqrt{dR^{3/2}(L_0+L_1M)}}\min\left\{\tilde{\sigma},\frac{\varepsilon^{3/2}}{\sqrt{dR^{3/2}}}\right\}\right)$ | 1 | Theorem 5.4 (Ours) |

Assumptions on objective function. Throughout this paper, we refer to the standard *L*-smoothness assumption, which is widely used in the literature [e.g. 51] and has the following form:

Assumption 1.1 (*L*-smoothness). Function f is *L*-smooth if for any $x, y \in \mathbb{R}^d$ is satisfied:

$$\left\|\nabla f(y) - \nabla f(x)\right\| \le L \left\|y - x\right\|.$$

Despite the widespread use of Assumption 1.1, our work focuses on the more general smoothness assumption, which has recently attracted increased interest. In particular, in [62] it was shown that norm of Hesse matrix correlates with norm of gradient function when training neural networks, and in [44] it was shown that using generalized smoothness it is possible to significantly improve the convergence of algorithms. (L_0, L_1) -smoothness [62, 63] has been proposed as a natural relaxation of standard smoothness assumption.

Assumption 1.2 ((L_0, L_1) -smoothness). A function $f : \mathbb{R}^d \to \mathbb{R}$ is (L_0, L_1) -smooth if the following inequality is satisfied for any $x, y \in \mathbb{R}^d$ with $||y - x|| \le \frac{1}{L_1}$:

$$\|\nabla f(y) - \nabla f(x)\| \le (L_0 + L_1 \|\nabla f(x)\|) \|y - x\|.$$

Assumption 1.2 in the case $L_1 = 0$ covers the standard Assumption 1.1. Moreover, (L_0, L_1) -smoothness is strictly more general than L-smoothness, see the examples in [62, 12, 34, 26].

Remark 1.3 (Clarification regarding $L_0 = 0$). In this paper we often emphasize the case $L_0 = 0$ in Assumption 1.2. It is worth noting that the class of functions that do not reach their infimum x^* (converge to an asymptote) satisfies this case. Explicit examples of functions with $L_0 = 0$ are the exponent of the inner product and the logistic function (see [26] for details).

Assumptions on gradient oracle. In our analysis, we consider cases with both unbiased and biased gradient oracle (2). Therefore, we assume that the bias and variance of gradient oracle (2) are bounded: Assumption 1.4 (Bounded bias). There exists constant $\zeta \ge 0$ such that the bias is bounded if $\forall x \in \mathbb{R}^d$:

$$\|\mathbf{b}(x)\| \le \zeta.$$

Assumption 1.5 (Bounded variance). There exists constant $\sigma^2 \ge 0$ such that the variance is bounded if $\forall x \in \mathbb{R}^d$:

$$\mathbb{E}\left[\left\|\mathbf{g}(x,\xi) - \mathbb{E}\left[\mathbf{g}(x,\xi)\right]\right\|^{2}\right] \leq \sigma^{2}.$$

Assumption 1.4 is organic [see, e.g. 43], and the case $\zeta = 0$ corresponds to the unbiased gradient oracle (2). Assumption 1.5 is often used by the community [e.g. 32, 37], and is sometimes called heavy-tailed noise [25].

1.3 Paper Organization

Next, our paper has the following structure. In Section 2, we discuss related work. In Section 3, We start to present the main results of our work, in particular, we provide the first strategy for obtaining a summand characterizing the linear rate in the convergence estimate. In Section 4, we analyze NSGD in the convex setting, showing in which regime linear convergence can be observed. We provide a first analysis of zero-order algorithms under (L_0, L_1) -smoothness in Section 5. In Section 6, we discuss the results obtained. While, in Section 7, we show experimentally about the possibility of linear convergence in the convex setting. Finally, Section 8 concludes our paper. All missing proofs of Lemmas and Theorems are provided in the supplementary materials (Appendix).

2 Related Works

In this section, we will discuss the most related works.

Algorithms under (L_0, L_1) -smoothness. Generalized smoothness was first introduced in [62], which analyzed ClipSGD in the non-convex setting. A number of works [56, 41, 19, 40, 29, 59, 57] followed that also focused on the non-convex setup, including ClipSGD [63, 61, 34], NSGD [65, 31]. After that, there was interest in research on algorithms in the convex deterministic setting: Clipped Gradient Descent [34], Normalized Gradient Descent [13], Gradient Descent with Polyak step size $\eta_k = \frac{f(x^k) - f^*}{\|\nabla f(x^k)\|^2}$ [54], and $\eta_k = \frac{1}{L_0 + L_1 \|\nabla f(x^k)\|}$ [26, 55]. Moreover, in [44], it was theoretically shown that it is possible to significantly improve the convergence of algorithms in the (strongly) convex setting by achieving linear convergence rate. However, much less attention has been paid to the stochastic convex setting. Perhaps the only results are [26, 20], which considers SGD and ClipSGD achieving only a sublinear convergence to the desired accuracy. In our work, we focus on the stochastic convex setup, showing that existing convergence results can be significantly improved.

Zero-order algorithms. The work of [21] showed that to achieve optimal estimates of iteration N and oracle T complexity in zero-order algorithms, one should base it on a first-order algorithm using a gradient approximation as the biased gradient oracle (2), which uses only information about the objective function f. Using this technique a number of works have achieved the best convergence results in various settings including distributed optimization [2], federated optimization [49], overparameterization [42], Polyak-Lojasiewicz condition [23], etc. However, all these works assumed standard smoothness (Assumption 1.1) and achieved only sublinear convergence rates. In our work, we present convergence results for zero-order algorithms under (L_0, L_1) -smoothness.

3 Clipped Stochastic Gradient Descent

In this section we begin to present the main results of our work. In particular, we analyze the convergence of SGD variants under convexity and (L_0, L_1) -smoothness with step size independent of the gradient norm. We assume that the gradient oracle (2) is unbiased $\zeta = 0$, i.e., Assumption 1.5 takes:

$$\mathbb{E}\left[\left\|\nabla f(x,\xi) - \nabla f(x)\right\|^{2}\right] \leq \sigma^{2}.$$

As a first strategy to obtain the summands that characterize the linear rate, we consider the clipping technique. Applying this technique we produce the ClipSGD, which has the following form:

Algorithm 1 Clipped Stochastic Gradient Descent Method (ClipSGD)

Input: initial point $x_0 \in \mathbb{R}^d$, iterations N, batch size B, step size $\eta_k > 0$ and clipping radius c > 0for k = 0 to N - 1 do 1. Draw fresh i.i.d. samples $\xi_1^k, ..., \xi_B^k$ 2. $\nabla f(x^k, \boldsymbol{\xi}^k) = \frac{1}{B} \sum_{i=1}^B \nabla f(x^k, \xi_i^k)$ 3. $\operatorname{clip}_c(\nabla f(x^k, \boldsymbol{\xi}^k)) = \min\left\{1, \frac{c}{\|\nabla f(x^k, \boldsymbol{\xi}^k)\|}\right\} \nabla f(x^k, \boldsymbol{\xi}^k)$ 4. $x^{k+1} \leftarrow x^k - \eta_k \cdot \operatorname{clip}_c(\nabla f(x^k, \boldsymbol{\xi}^k))$ end for Return: x^N

Algorithm 1 uses the clipped stochastic gradient $\operatorname{clip}_c(\nabla f(x, \boldsymbol{\xi}))$, which normalizes the gradient only if $\|\nabla f(x, \boldsymbol{\xi})\| > c$. Next theorem provides the convergence result for ClipSGD.

Theorem 3.1. Let function f satisfy Assumption 1.2 ((L_0, L_1) -smoothness) and unbiased gradient oracle (2) satisfy Assumption 1.5 (bounded variance), then Algorithm 1 with constant step size $\eta_k = \eta \leq [4(L_0 + L_1c)]^{-1}$ and arbitrary clipping radius c guarantees error:

$$\mathbb{E}\left[f(x^N)\right] - f^* \lesssim \left(1 - \frac{\eta c}{R}\right)^K (f(x^0) - f^*) + \frac{R^2}{\eta(N-K)} + \frac{\sigma^2}{B}\left(\eta + \frac{MR}{c^2} + \frac{R}{c}\right),$$

where $0 \le K < N$ is number of iterations for which $\|\nabla f(x^k)\| \le \frac{c}{2}$ is satisfied.

It should be noted that the results of Theorem 3.1 are given with a choice of step size independent of the gradient at the current point. This choice of step allows us to separate the constants L_0 and L_1 in the final estimates. The summand $\frac{R^2}{\eta(N-K)}$ is a typical ClipSGD characterizing the sublinear rate (see e.g. [25]). However, it is worth noting that by substituting $\eta = (L_0 + L_1 c)^{-1}$, then the summand with $L_1: \frac{L_1 cR^2}{N-K}$ already improves existing results both assuming standard smoothness [25] and generalized smoothness [20]. The first summand $(1 - \frac{\eta c}{R})^K (f(x^0) - f^*)$, which characterizes the linear rate, deserves special attention. To the best of our knowledge, this is the first result for ClipSGD showing such a summand in a convex setting. Moreover, by substituting $\eta = (L_0 + L_1 c)^{-1}$, it is not hard to see that in the regime $L_0 = 0$ (see Remark 1.3), Algorithm 1 at a batch size $B = \mathcal{O}\left(\frac{\sigma^2}{\varepsilon} \left(\eta + \frac{MR}{c^2} + \frac{R}{c}\right)\right)$ requires only $N = \mathcal{O}\left(L_1R\log\frac{1}{\varepsilon} + \frac{L_1cR^2}{\varepsilon}\right)$ iterations. This iteration complexity significantly outperforms standard results in the L-smoothness setting (Assumption 1.1), since [37] shows a lower bound consisting only of a sublinear convergence rate. Moreover, comparing to the closest work to the problem setting, then even when $\sigma = 0$ [20] does not guarantee convergence to the desired accuracy, offering an estimate of $N = \mathcal{O}\left(L_1^2R^2\right)$ that is independent of accuracy.

4 Normalized Stochastic Gradient Descent

In the previous section, we showed that it is possible to obtain in the final convergence estimate a summand characterizing the linear rate. In addition, we highlighted the regime $L_0 = 0$, in which ClipSGD has the following iteration complexity: $N = O\left(L_1R\log\frac{1}{\varepsilon} + \frac{L_1cR^2}{\varepsilon}\right)$. However, with this iteration complexity, it cannot be said that the algorithm can converge with linear rate to the desired accuracy. Such an estimate can characterize that the algorithm converges with linear rate as long as the gradient norm is large $\|\nabla f(x^k)\| \ge c$, and then slows down to the sublinear rate. However, it is worth noting that the summand responsible for the sublinear rate depends on the clipping radius: $\frac{L_1cR^2}{\varepsilon}$. That is, if we take c smaller, the ClipSGD will take longer to converge to the linear rate. Thus, noticing that the regime $\|\nabla f(x^k)\| \ge c$ is a gradient normalization, then considering NSGD, it seems that one can achieve a true linear convergence rate to the desired accuracy.

In this section, we consider a normalization technique to obtain the summand characterizing the linear rate. Applying this technique we produce the NSGD, which has the following form (see Algorithm 2):

Algorithm 2 Normalized Stochastic Gradient Descent Method (NSGD)

Input: initial point $x_0 \in \mathbb{R}^d$, iterations number N, batch size B and step size $\eta_k > 0$ for k = 0 to N - 1 do 1. Draw fresh i.i.d. samples $\xi_1^k, ..., \xi_B^k$ 2. $\nabla f(x^k, \boldsymbol{\xi}^k) = \frac{1}{B} \sum_{i=1}^{B} \nabla f(x^k, \xi_i^k)$ 3. $x^{k+1} \leftarrow x^k - \eta_k \cdot \frac{\nabla f(x^k, \boldsymbol{\xi}^k)}{\|\nabla f(x^k, \boldsymbol{\xi}^k)\|}$ end for Return: x^N

The following theorem provides a convergence result for Algorithm 2.

Theorem 4.1. Let function f satisfy Assumption 1.2 ((L_0, L_1) -smoothness) and unbiased gradient oracle (2) satisfy Assumption 1.5 (bounded variance), then Algorithm 2 with hyperparameter $\lambda > 0$ and constant step size $\eta_k = \eta \le \lambda / [2(L_0 + L_1\lambda)]$ guarantees:

$$\mathbb{E}\left[f(x^N)\right] - f^* \lesssim \left(1 - \frac{\eta}{R}\right)^N (f(x^0) - f^*) + \frac{\sigma^2 M R}{B\lambda^2} + \lambda R.$$

From Theorem 4.1 we can see that we have indeed got rid of the summand characterizing the sublinear rate from the deterministic part. Thus, we see that it is normalization that allows us to achieve the summand characterizing the linear rate $\left(1 - \frac{\eta}{R}\right)^N (f(x^0) - f^*)$. However, note that by substituting $\eta = \lambda / [2(L_0 + L_1\lambda)]$, we obtain a summand with $L_0 : \left(1 - \frac{\lambda}{RL_0}\right)^N (f(x^0) - f^*)$, which is in fact

sublinear since it depends on the hyperparameter λ (it follows from the third summand that $\lambda \sim \varepsilon/R$), and with $L_1: \left(1-\frac{1}{RL_1}\right)^N (f(x^0)-f^*)$, which is indeed linear since it does not depend on λ in any way. That is, Algorithm 2, which uses batch parallelization, requires $N = \mathcal{O}\left((L_1R + \frac{L_0R^2}{\epsilon})\log\frac{1}{\epsilon}\right)$ iterations. Similar to the reasoning in the previous section, it is worth highlighting the regime $L_0 = 0$ (see Remark 1.3). Then we obtain a very surprising result on iteration complexity, namely, to achieve the desired accuracy NSGD converge with a linear rate of $N = \mathcal{O}\left(L_1 R \log \frac{1}{\epsilon}\right)$ iterations. This estimate breaks all existing bounds on first-order algorithms [37], given the specificity of the problem formulation, namely convexity. However, to achieve this rate over iterations, NSGD requires a batch size $B = O\left(\frac{\sigma^2 M R^3}{\varepsilon^3}\right)$. We emphasize that the fact that NSGD requires a large batch size is not surprising (see, e.g., [14]), in contrast to the true linear convergence rate in the convex setup.

5 Zero-Order Algorithms

In this section, we consider another class of algorithms: optimization algorithms that have access only to an objective function value f possibly with some bounded adversarial noise $|\delta(x)| \leq \Delta$:

$$\tilde{f}(x,\xi) = f(x,\xi) + \delta(x).$$
(3)

In (3), Δ means the maximum possible allowable noise level at which the desired accuracy can still be achieved. In [5], the importance of considering Δ as a third optimality criterion for zero-order algorithms was shown. In particular, in some applications [9], the larger noise level Δ is, the cheaper the call to the inexact oracle \tilde{f} in (3).

Since this class of algorithms does not have access to the stochastic gradient $\nabla f(x,\xi)$, the gradient oracle (2) will be the gradient approximation with L_2 randomization [52, 46]:

$$\mathbf{g}(x,\{e,\xi\}) = \frac{d}{2\gamma} \left(\tilde{f}(x+\gamma e,\xi) - \tilde{f}(x-\gamma e,\xi) \right) e, \tag{4}$$

where $\gamma > 0$ is a smoothing parameter, e is a random vector uniformly distributed in $S^{d}(1)$.

Due to the fact that the gradient approximation is the biased gradient oracle (2), in order to create zero-order algorithms by basing on the results in Sections 3 and 4, it is necessary to first generalize the results of Theorems 3.1, 4.1 (note that in these regimes it is not necessary to know the (stochastic) norm of the gradient with step size) to the case of gradient oracle with bias.

Next, we present convergence results for the following two zero-order algorithms: ZO-ClipSGD (Algorithm 3) and ZO-NSGD (Algorithm 4).

5.1 ZO-ClipSGD Method

er

The first algorithm we consider in this section is ZO-ClipSGD. This algorithm is a modification of ClipSGD (Algorithm 1), which uses instead of the original $\|\nabla f(x,\xi)\|$, the stochastic gradient approximation (4), which is the biased gradient oracle (2). The ZO-ClipSGD has the following form:

Algorithm 3 Zero-Order Clipped Stochastic Gradient Descent Method (ZO-ClipSGD)

Input: initial point $x_0 \in \mathbb{R}^d$, iterations N, batch size B, step size $\eta_k > 0$ and clipping radius c > 0for k = 0 to N - 1 do 1. Draw fresh i.i.d. samples $\xi_1^k,...,\xi_B^k$ and $e_1^k,...,e_B^k$ 2. $\mathbf{g}(x^k, {\mathbf{e}^k, \boldsymbol{\xi}^k}) = \frac{1}{B} \sum_{i=1}^{B} \mathbf{g}(x^k, {e_i^k, \xi_i^k})$ via (4)

$$\begin{array}{l} 3. \ \operatorname{clip}_c(\mathbf{g}(x^k, \{\mathbf{e}^k, \boldsymbol{\xi}^k\})) = \min\left\{1, \frac{c}{\|\mathbf{g}(x^k, \{\mathbf{e}^k, \boldsymbol{\xi}^k\})\|}\right\} \mathbf{g}(x^k, \{\mathbf{e}^k, \boldsymbol{\xi}^k\}) \\ 4. \ x^{k+1} \leftarrow x^k - \eta_k \cdot \operatorname{clip}_c(\mathbf{g}(x^k, \{\mathbf{e}^k, \boldsymbol{\xi}^k\})) \\ \text{end for} \\ \text{Return: } x^N \end{array}$$

Before proceeding to present the convergence results of Algorithm 3, we note that the gradient approximation (4) is a biased gradient oracle, so we cannot directly use the results obtained in Theorem 3.1. Thus, in order to obtain estimates for the iteration complexity N, oracle complexity T and maximum noise Δ , we first need to generalize the results of Theorem 3.1 to the case with a biased gradient oracle (2).

Lemma 5.1. Let function f satisfy Assumption 1.2 and biased gradient oracle (2) ($\zeta > 0$) satisfy Assumption 1.5, then Algorithm 1 with step size $\eta_k = \eta \leq [4(L_0 + L_1c)]^{-1}$ guarantees error:

$$\mathbb{E}\left[f(x^N)\right] - f^* \lesssim \left(1 - \frac{\eta c}{R}\right)^K \left(f(x^0) - f^*\right) + \frac{R^2}{\eta(N - K)} + \mathcal{R} \cdot \left(\frac{\sigma^2}{B} + \zeta^2\right) + R\zeta,$$

where c is arbitrary clipping radius, $\mathcal{R} = \left(\eta + \frac{MR}{c^2} + \frac{R}{c}\right), 0 \le K < N$ is the number of iterations for which the condition $\|\nabla f(x^k)\| \le \frac{c}{3}$ is satisfied.

Lemma 5.1 shows how the bias accumulates over iterations, thus converging to the error floor. By estimating the second moment and the bias of the gradient approximation (4) and substituting them into Lemma 5.1, we find three optimality criteria for ZO-ClipSGD.

Theorem 5.2. Let function f satisfy Assumption 1.2, gradient approximation (4) satisfy Assumption 1.5, then Algorithm 3 with step size $\eta_k = \eta \leq [4(L_0 + L_1c)]^{-1}$ converges to desired ε accuracy after:

$$N = \mathcal{O}\left(\frac{R}{\eta c}\log\frac{1}{\varepsilon} + \frac{R^2}{\eta \varepsilon}\right); \qquad T = \mathcal{O}\left(\frac{d\tilde{\sigma}^2 M R^2}{\varepsilon c^2 \eta} \left(\frac{1}{c}\log\frac{1}{\varepsilon} + \frac{R}{\varepsilon}\right)\right)$$

number of iterations and zero-order oracle calls at

$$\Delta \lesssim \frac{\varepsilon}{\sqrt{d}R(L_0 + L_1M)} \min\left\{\tilde{\sigma}, \frac{\varepsilon}{\sqrt{d}R}\right\}$$

maximum noise level, where c > 0 is clipping radius, $\mathbb{E}\left[\left\|\nabla f(x,\xi)\right\|^2\right] \leq \tilde{\sigma}^2$.

It is not hard to see from Theorem 5.2 that in the generalized smoothness condition, the iteration complexity at ZO-ClipSGD is the same as the first-order method of Algorithm 1. This effect is similar in the standard smoothness condition as well. The oracle complexity is *d* times worse than its first-order counterpart due to the restriction to the oracle (Algorithm 3 uses only the zero-order oracle (3)). It is worth noting that the maximum noise level Δ outperforms [35], showing that generalized smoothness not only allows us to reach the summands characterizing the linear rate, but also improves the estimate on the maximum noise level (it is Δ that affects the error floor, in other words, the accuracy of the solution, to control the asymptote). See the proof in the Appendix D.

5.2 ZO-NSGD Method

Similar to the first-order algorithms, in this subsection we answer the question whether linear convergence can be achieved by the zero-order method in a convex setup. To answer this question, we consider the Zero-Order Normalized Stochastic Gradient Descent Method.

Algorithm 4 Zero-Order Normalized Stochastic Gradient Descent Method (ZO-NSGD)

Input: initial point $x_0 \in \mathbb{R}^d$, iterations number N, batch size B, step size $\eta_k > 0$ for k = 0 to N - 1 do 1. Draw fresh i.i.d. samples $\xi_1^k, ..., \xi_B^k$ and $e_1^k, ..., e_B^k$ 2. $\mathbf{g}(x^k, \{\mathbf{e}^k, \boldsymbol{\xi}^k\}) = \frac{1}{B} \sum_{i=1}^{B} \mathbf{g}(x^k, \{e_i^k, \xi_i^k\})$ via (4) 3. $x^{k+1} \leftarrow x^k - \eta_k \cdot \frac{\mathbf{g}(x^k, \{\mathbf{e}^k, \boldsymbol{\xi}^k\})}{\|\mathbf{g}(x^k, \{\mathbf{e}^k, \boldsymbol{\xi}^k\})\|}$ end for Return: x^N

In a similar way as in the previous subsection, we first generalize Theorem 4.1 to the case with a biased gradient oracle to substitute estimates on the bias and the second moment of the gradient approximation to find optimality criteria for the gradient-free algorithm ZO-NSGD.

Lemma 5.3. Let function f satisfy Assumption 1.2 $((L_0, L_1)$ -smoothness) and biased gradient oracle (2) $(\zeta > 0)$ satisfy Assumption 1.5 (bounded variance), then Algorithm 4 with hyperparameter $\lambda > 0$ and step size $\eta_k = \eta \leq \lambda / [2(L_0 + L_1\lambda)]$ guarantees:

$$\mathbb{E}\left[f(x^N)\right] - f^* \lesssim \left(1 - \frac{\eta}{R}\right)^N \left(f(x^0) - f^*\right) + \frac{MR}{\lambda^2} \left(\frac{\sigma^2}{B} + \zeta^2\right) + \lambda R.$$

From Lemma 5.3 we can see how the inaccuracy accumulates over the iteration. The summand with ζ^2 is unimprovable for first-order unaccelerated algorithms (see, e.g., [17, 23]). Now, having obtained the results for the biased NSGD we can use them to derive convergence results for Algorithm 4.

Theorem 5.4. Let function f satisfy Assumption 1.2 ((L_0, L_1) -smoothness) and gradient approximation (4) satisfy Assumption 1.5 (bounded variance), then ZO-NSGD with step size $\eta_k = \eta \leq \lambda / [2(L_0 + L_1\lambda)]$ converges to desired ε accuracy ($\mathbb{E}[f(x^N)] - f^* \leq \varepsilon$) after:

$$N = \mathcal{O}\left(\frac{R}{\eta}\log\frac{1}{\varepsilon}\right); \qquad T = \tilde{\mathcal{O}}\left(\frac{d\tilde{\sigma}^2 M R^4}{\varepsilon^3 \eta}\right)$$

number of iterations and zero-order oracle calls (3) at

$$\Delta \lesssim \frac{\varepsilon^{3/2}}{\sqrt{d}R^{3/2}(L_0 + L_1M)} \min\left\{\tilde{\sigma}, \frac{\varepsilon^{3/2}}{\sqrt{d}R^{3/2}}\right\}$$

maximum noise level, where $\lambda > 0$ is hyperparameter, $\mathbb{E}\left[\|\nabla f(x,\xi)\|^2 \right] \leq \tilde{\sigma}^2$.

It is not hard to see that given a restricted oracle (3), Theorem 5.4 shows that ZO-NSGD requires $N = \mathcal{O}\left(\frac{R}{\eta}\log\frac{1}{\varepsilon}\right)$ iterations to achieve the desired accuracy, which corresponds to a linear rate. However, it is worth noting that, as in Theorem 4.1, if we take the maximum step size $\eta = \lambda/\left[2(L_0 + L_1\lambda)\right]$, then the summand with L_0 still shows sublinear convergence $\tilde{\mathcal{O}}\left(\frac{L_0R^2}{\varepsilon}\right)$, despite the presence of the summand with L_1 . But in the $L_0 = 0$ regime, we can say unambiguously that the zero-order algorithm ZO-NSGD can converge with a linear rate to the desired accuracy in the convex setup if we take the batch size $B = \mathcal{O}\left(\frac{d\tilde{\sigma}^2 MR^3}{\varepsilon^3}\right)$. Theorem 5.4 shows that the power of generalized smoothness (together with Batch parallelization) extends to zero-order algorithms. Comparing the result of ZO-NSGD with Theorem 5.2, we can see that for a significant improvement in iteration complexity N we "pay" for a deterioration in both oracle complexity T and maximum noise level Δ , which seems quite natural. See the proof of Lemma 5.3 and Theorem 5.4 in Appendix E.

6 Discussion and Future Works

In Sections 3-5, we gave two strategies for obtaining summands that characterize the linear rate in convergence estimates of algorithms for the convex setting: clipping (see Section 3) and normalization (see Section 4) techniques. Although these summands are quite unexpected for the convex setup and improve the estimates on the iteration complexity, we cannot claim linear convergence in general, since the convergence is dominated by the summands characterizing the sublinear rate. However, as we noted in Theorems 4.1 and 5.4, in the regime $L_0 = 0$, the NSGD and ZO-NSGD methods break all existing bounds on iteration complexity, demonstrating that it is possible to converge with linear rate to the desired accuracy with the condition of using Batch parallization. This result is pleasantly surprising and opens up a number of directions for future research.

In this paper we have focused on iteration complexity, so we see a careful analysis of the optimality criterion in the aggregate as future work. In particular, it seems interesting to show that M is indeed bounded, e.g., using the technique from [40] and evaluating with respect to the smoothness constants L_0 and L_1 . The existence of the regime $L_0 = 0$ which allows one to achieve a linear convergence rate in the convex setting prompts the following question: can the iteration complexity be improved by assuming, for example, strong convexity, the PL condition, etc.? It is also interesting to see if similar effects are found in accelerated, adaptive algorithms, variational inequalities, distributed learning, nonsmooth (or increased smoothness) problems, overparameterization, online optimization, etc.

7 Numerical Experiments

In this section, we numerically analyze the algorithms presented in this paper and show that linear convergence in stochastic convex optimization is possible. For this illustration, we have chosen a problem that is of particular interest in the machine learning community: the logistic regression problem on w1a dataset [50]. We consider the following convex problem statement (1):

$$\min_{x \in \mathbb{R}^d} f(x) = \frac{1}{M} \sum_{i=1}^M f_i(x), \qquad f_i(x) = \log \left(1 + \exp(-y_i \cdot (Ax)_i)\right),$$

where $f_i(x)$ is the loss on the *i*-th data point, $A \in \mathbb{R}^{M \times d}$ is an instances matrix, $y \in \{-1, 1\}^M$ is a label vector and $x \in \mathbb{R}^d$ is a vector of weights. It is easy to show, that logistic regression function is L-smooth (see Assumption 1.1) with $L = \frac{1}{4M} \sqrt{\lambda_{\max}(A^T A)}$, where $\lambda_{\max}(A^T A)$ denotes the largest eigenvalue of the matrix $A^T A$. Moreover, such a problem statement is a special case of (1) with ξ being a random variable with the uniform distribution on $\{1, ..., M\}$.



Figure 2: Linear convergence demonstration.

Figure 3: Comparison of the considered algorithms.

In all tests we used the following parameters: M = 2477 - number of data, d = 300 - problem dimension, B = 10 - batch size, $h = 10^{-5}$ - smoothing parameter, $\Delta = 10^{-9}$ - noise level. Figure 2 shows a comparison of two step strategies of the NSGD algorithm. The blue line (see "Standard smoothness step"), which corresponds to the theoretical step size in the *L*-smoothness setting demonstrates slow (sublinear) convergence. In particular, NSGD with this choice of step advanced the function from 0.02014151259 to 0.01731953499 in 25000 iterations. The red line (see "Ours step size"), which corresponds to the next step size $\eta = \frac{1}{\|A\|_1}$, demonstrates linear convergence that significantly outperforms the strategy with the theoretical step size. Moreover, *Figure 2 demonstrates that indeed the first-order algorithm can converge with linear rate to the desired accuracy in a convex setup!*

Figure 3 demonstrates the convergence dynamics of all the algorithms considered in this paper. In particular, as in Figure 2, NSGD (see red line) shows a real linear convergence. ClipSGD (see green line), which used a step size $\eta = \frac{1}{c \cdot ||A||_1}$, where $c = 10^{-1}$ is the clipping radius, shows two modes of convergence: as long as $||\nabla f(x^k)|| \ge c$ the algorithm converges with a linear rate matching the NSGD, as soon as $||\nabla f(x^k)|| \le c$ the algorithm slows down to the sublinear rate. The dynamics are similar for zero-order algorithms: ZO-NSGD (see orange line), ZO-ClipSGD (see blue line). Expectedly, these algorithms converge slower on the first iterations than their first-order counterparts due to restricted access to the oracle (3). However, it is worth noting that ZO-NSGD also exhibits linear convergence, thereby outperforming the first-order ClipSGD algorithm after 55000 iterations.

8 Conclusion

In this paper, we considered a stochastic convex optimization problem under the generalized smoothness condition of the objective function. We are the first who have provided strategies to achieve summands characterizing linear rate, thereby improving the iteration complexity (see Sections 3 and 4). In Section 5, we showed that this effect of generalized smoothness extends to zero-order algorithms as well. Moreover, we highlight the regime $L_0 = 0$, under which we theoretically guarantee linear convergence for the NSGD and ZO-NSGD algorithms (subject to the use of batch parallization). This is the first result demonstrating linear convergence in such a problem setting, thus opening up a number of future works (see Section 6). Finally, in Section 7 we show using numerical experiments that linear convergence in such a problem formulation is also possible in practice.

References

- [1] Alekh Agarwal, Ofer Dekel, and Lin Xiao. Optimal algorithms for online convex optimization with multi-point bandit feedback. In *Colt*, pages 28–40. Citeseer, 2010.
- [2] Arya Akhavan, Massimiliano Pontil, and Alexandre Tsybakov. Distributed zero-order optimization under adversarial noise. Advances in Neural Information Processing Systems, 34:10209–10220, 2021.
- [3] Arya Akhavan, Evgenii Chzhen, Massimiliano Pontil, and Alexandre Tsybakov. A gradient estimator via 11-randomization for online zero-order optimization with two point feedback. *Advances in Neural Information Processing Systems*, 35:7685–7696, 2022.
- [4] Dan Alistarh, Torsten Hoefler, Mikael Johansson, Nikola Konstantinov, Sarit Khirirat, and Cédric Renggli. The convergence of sparsified gradient methods. *Advances in Neural Information Processing Systems*, 31, 2018.
- [5] Authors Anonymous. Maximum noise level as third optimality criterion in black-box optimization problem, 2025.
- [6] Francis Bach and Vianney Perchet. Highly-smooth zero-th order online optimization. In *Conference on Learning Theory*, pages 257–283. PMLR, 2016.
- [7] Irwan Bello, Barret Zoph, Vijay Vasudevan, and Quoc V Le. Neural optimizer search with reinforcement learning. In *International Conference on Machine Learning*, pages 459–468. PMLR, 2017.
- [8] Yoshua Bengio, Patrice Simard, and Paolo Frasconi. Learning long-term dependencies with gradient descent is difficult. *IEEE transactions on neural networks*, 5(2):157–166, 1994.
- [9] Lev Bogolubsky, Pavel Dvurechenskii, Alexander Gasnikov, Gleb Gusev, Yurii Nesterov, Andrei M Raigorodskii, Aleksey Tikhonov, and Maksim Zhukovskii. Learning supervised pagerank with gradientbased and gradient-free optimization methods. *Advances in neural information processing systems*, 29, 2016.
- [10] Léon Bottou. Online algorithms and stochastic approximations. Online learning in neural networks, 1998.
- [11] Léon Bottou, Frank E Curtis, and Jorge Nocedal. Optimization methods for large-scale machine learning. SIAM review, 60(2):223–311, 2018.
- [12] Ziyi Chen, Yi Zhou, Yingbin Liang, and Zhaosong Lu. Generalized-smooth nonconvex optimization is as efficient as smooth nonconvex optimization. In *International Conference on Machine Learning*, pages 5396–5427. PMLR, 2023.
- [13] Ziyi Chen, Yi Zhou, Yingbin Liang, and Zhaosong Lu. Generalized-smooth nonconvex optimization is as efficient as smooth nonconvex optimization. In *International Conference on Machine Learning*, pages 5396–5427. PMLR, 2023.
- [14] Ashok Cutkosky and Harsh Mehta. Momentum improves normalized sgd. In International conference on machine learning, pages 2260–2268. PMLR, 2020.
- [15] Jeffrey Dean, Greg Corrado, Rajat Monga, Kai Chen, Matthieu Devin, Mark Mao, Marc'aurelio Ranzato, Andrew Senior, Paul Tucker, Ke Yang, et al. Large scale distributed deep networks. Advances in neural information processing systems, 25, 2012.
- [16] Yury Demidovich, Grigory Malinovsky, Igor Sokolov, and Peter Richtárik. A guide through the zoo of biased sgd. *Advances in Neural Information Processing Systems*, 36:23158–23171, 2023.
- [17] Olivier Devolder. Exactness, inexactness and stochasticity in first-order methods for large-scale convex optimization. *Candidate's Dissertation (CORE UCLouvain Louvain-la-Neuve, Belgium)*, 2013.
- [18] Tolga Dimlioglu and Anna Choromanska. Grawa: Gradient-based weighted averaging for distributed training of deep learning models. In *International Conference on Artificial Intelligence and Statistics*, pages 2251–2259. PMLR, 2024.
- [19] Matthew Faw, Litu Rout, Constantine Caramanis, and Sanjay Shakkottai. Beyond uniform smoothness: A stopped analysis of adaptive sgd. In *The Thirty Sixth Annual Conference on Learning Theory*, pages 89–160. PMLR, 2023.
- [20] Ofir Gaash, Kfir Yehuda Levy, and Yair Carmon. Convergence of clipped sgd on convex (l_0, l_1)-smooth functions. arXiv preprint arXiv:2502.16492, 2025.

- [21] Alexander Gasnikov, Anton Novitskii, Vasilii Novitskii, Farshed Abdukhakimov, Dmitry Kamzolov, Aleksandr Beznosikov, Martin Takac, Pavel Dvurechensky, and Bin Gu. The power of first-order smooth optimization for black-box non-smooth problems. In *International Conference on Machine Learning*, pages 7241–7265. PMLR, 2022.
- [22] Alexander Gasnikov, Darina Dvinskikh, Pavel Dvurechensky, Eduard Gorbunov, Aleksandr Beznosikov, and Alexander Lobanov. Randomized gradient-free methods in convex optimization. In *Encyclopedia of Optimization*, pages 1–15. Springer, 2023.
- [23] AV Gasnikov, AV Lobanov, and FS Stonyakin. Highly smooth zeroth-order methods for solving optimization problems under the pl condition. *Computational Mathematics and Mathematical Physics*, 64(4): 739–770, 2024.
- [24] Ian Goodfellow. Deep learning, 2016.
- [25] Eduard Gorbunov, Marina Danilova, and Alexander Gasnikov. Stochastic optimization with heavy-tailed noise via accelerated gradient clipping. Advances in Neural Information Processing Systems, 33:15042– 15053, 2020.
- [26] Eduard Gorbunov, Nazarii Tupitsa, Sayantan Choudhury, Alen Aliev, Peter Richtárik, Samuel Horváth, and Martin Takáč. Methods for convex (l_0, l_1)-smooth optimization: Clipping, acceleration, and adaptivity. arXiv preprint arXiv:2409.14989, 2024.
- [27] Elad Hazan, Kfir Levy, and Shai Shalev-Shwartz. Beyond convexity: Stochastic quasi-convex optimization. Advances in neural information processing systems, 28, 2015.
- [28] José Miguel Hernández-Lobato, Matthew W Hoffman, and Zoubin Ghahramani. Predictive entropy search for efficient global optimization of black-box functions. *Advances in neural information processing systems*, 27, 2014.
- [29] Yusu Hong and Junhong Lin. On convergence of adam for stochastic optimization under relaxed assumptions. arXiv preprint arXiv:2402.03982, 2024.
- [30] Florian Hübler, Ilyas Fatkhullin, and Niao He. From gradient clipping to normalization for heavy tailed sgd. *arXiv preprint arXiv:2410.13849*, 2024.
- [31] Florian Hübler, Junchi Yang, Xiang Li, and Niao He. Parameter-agnostic optimization under relaxed smoothness. In *International Conference on Artificial Intelligence and Statistics*, pages 4861–4869. PMLR, 2024.
- [32] Anatoli B Juditsky and Arkadii S Nemirovski. First order methods for nonsmooth convex large-scale optimization, i: General purpose methods. *Optimization for Machine Learning*, pages 1–28, 2010.
- [33] Peter Kairouz, H Brendan McMahan, Brendan Avent, Aurélien Bellet, Mehdi Bennis, Arjun Nitin Bhagoji, Kallista Bonawitz, Zachary Charles, Graham Cormode, Rachel Cummings, et al. Advances and open problems in federated learning. *Foundations and trends*® *in machine learning*, 14(1–2):1–210, 2021.
- [34] Anastasia Koloskova, Hadrien Hendrikx, and Sebastian U Stich. Revisiting gradient clipping: Stochastic bias and tight convergence guarantees. In *International Conference on Machine Learning*, pages 17343– 17363. PMLR, 2023.
- [35] Nikita Kornilov, Ohad Shamir, Aleksandr Lobanov, Darina Dvinskikh, Alexander Gasnikov, Innokentiy Shibaev, Eduard Gorbunov, and Samuel Horváth. Accelerated zeroth-order method for non-smooth stochastic convex optimization problem with infinite variance. *Advances in Neural Information Processing Systems*, 36:64083–64102, 2023.
- [36] Simon Lacoste-Julien, Mark Schmidt, and Francis Bach. A simpler approach to obtaining an o (1/t) convergence rate for the projected stochastic subgradient method. arXiv preprint arXiv:1212.2002, 2012.
- [37] Guanghui Lan. An optimal method for stochastic composite optimization. *Mathematical Programming*, 133(1):365–397, 2012.
- [38] Tor Lattimore and Andras Gyorgy. Improved regret for zeroth-order stochastic convex bandits. In *Conference on Learning Theory*, pages 2938–2964. PMLR, 2021.
- [39] Hojoon Lee, Hanseul Cho, Hyunseung Kim, Daehoon Gwak, Joonkee Kim, Jaegul Choo, Se-Young Yun, and Chulhee Yun. Plastic: Improving input and label plasticity for sample efficient reinforcement learning. Advances in Neural Information Processing Systems, 36, 2024.

- [40] Haochuan Li, Jian Qian, Yi Tian, Alexander Rakhlin, and Ali Jadbabaie. Convex and non-convex optimization under generalized smoothness. *Advances in Neural Information Processing Systems*, 36: 40238–40271, 2023.
- [41] Haochuan Li, Alexander Rakhlin, and Ali Jadbabaie. Convergence of adam under relaxed assumptions. *Advances in Neural Information Processing Systems*, 36:52166–52196, 2023.
- [42] Aleksandr Lobanov and Alexander Gasnikov. Accelerated zero-order sgd method for solving the black box optimization problem under "overparametrization" condition. In *International Conference on Optimization* and Applications, pages 72–83. Springer, 2023.
- [43] Aleksandr Lobanov, Nail Bashirov, and Alexander Gasnikov. The "black-box" optimization problem: Zero-order accelerated stochastic method via kernel approximation. *Journal of Optimization Theory and Applications*, pages 1–36, 2024.
- [44] Aleksandr Lobanov, Alexander Gasnikov, Eduard Gorbunov, and Martin Takác. Linear convergence rate in convex setup is possible! gradient descent method variants under (l_0, l_1) -smoothness. *arXiv preprint arXiv:2412.17050*, 2024.
- [45] Yurii Nesterov. Lectures on convex optimization, volume 137. Springer, 2018.
- [46] Yurii Nesterov and Vladimir Spokoiny. Random gradient-free minimization of convex functions. Foundations of Computational Mathematics, 17(2):527–566, 2017.
- [47] Anthony Nguyen and Krishnakumar Balasubramanian. Stochastic zeroth-order functional constrained optimization: Oracle complexity and applications. *INFORMS Journal on Optimization*, 2022.
- [48] Razvan Pascanu, Tomas Mikolov, and Yoshua Bengio. On the difficulty of training recurrent neural networks. In Sanjoy Dasgupta and David McAllester, editors, *Proceedings of the 30th International Conference on Machine Learning*, volume 28 of *Proceedings of Machine Learning Research*, pages 1310–1318, Atlanta, Georgia, USA, 17–19 Jun 2013. PMLR.
- [49] Kumar Kshitij Patel, Aadirupa Saha, Lingxiao Wang, and Nathan Srebro. Distributed online and bandit convex optimization. In OPT 2022: Optimization for Machine Learning (NeurIPS 2022 Workshop), 2022.
- [50] John C Platt. Fast training of support vector machines using sequential minimal optimization. Advances in Kernel Methods Support Vector Learning, MIT Press, 1998.
- [51] Boris T Polyak. Introduction to optimization. Optimization Software, Inc. Publications Division, New York, 1987.
- [52] Ohad Shamir. An optimal algorithm for bandit and zero-order convex optimization with two-point feedback. *The Journal of Machine Learning Research*, 18(1):1703–1713, 2017.
- [53] Sebastian U Stich and Sai Praneeth Karimireddy. The error-feedback framework: Better rates for sgd with delayed gradients and compressed communication. *arXiv preprint arXiv:1909.05350*, 2019.
- [54] Yuki Takezawa, Han Bao, Ryoma Sato, Kenta Niwa, and Makoto Yamada. Parameter-free clipped gradient descent meets polyak. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024.
- [55] Daniil Vankov, Anton Rodomanov, Angelia Nedich, Lalitha Sankar, and Sebastian U Stich. Optimizing (*l*_0, *l*_1)-smooth functions by gradient methods. *arXiv preprint arXiv:2410.10800*, 2024.
- [56] Bohan Wang, Huishuai Zhang, Zhiming Ma, and Wei Chen. Convergence of adagrad for non-convex objectives: Simple proofs and relaxed assumptions. In *The Thirty Sixth Annual Conference on Learning Theory*, pages 161–190. PMLR, 2023.
- [57] Bohan Wang, Huishuai Zhang, Qi Meng, Ruoyu Sun, Zhi-Ming Ma, and Wei Chen. On the convergence of adam under non-uniform smoothness: Separability from sgdm and beyond. *arXiv preprint arXiv:2403.15146*, 2024.
- [58] Blake E Woodworth, Brian Bullins, Ohad Shamir, and Nathan Srebro. The min-max complexity of distributed stochastic convex optimization with intermittent communication. In *Conference on Learning Theory*, pages 4386–4437. PMLR, 2021.
- [59] Chenghan Xie, Chenxi Li, Chuwen Zhang, Qi Deng, Dongdong Ge, and Yinyu Ye. Trust region methods for nonconvex stochastic optimization beyond lipschitz smoothness. In *Proceedings of the AAAI Conference* on Artificial Intelligence, pages 16049–16057, 2024.

- [60] Honglin Yuan and Tengyu Ma. Federated accelerated stochastic gradient descent. Advances in Neural Information Processing Systems, 33:5332–5344, 2020.
- [61] Bohang Zhang, Jikai Jin, Cong Fang, and Liwei Wang. Improved analysis of clipping algorithms for non-convex optimization. Advances in Neural Information Processing Systems, 33:15511–15521, 2020.
- [62] Jingzhao Zhang, Tianxing He, Suvrit Sra, and Ali Jadbabaie. Why gradient clipping accelerates training: A theoretical justification for adaptivity. In *International Conference on Learning Representations*, 2019.
- [63] Jingzhao Zhang, Sai Praneeth Karimireddy, Andreas Veit, Seungyeon Kim, Sashank Reddi, Sanjiv Kumar, and Suvrit Sra. Why are adaptive methods good for attention models? *Advances in Neural Information Processing Systems*, 33:15383–15393, 2020.
- [64] Sixin Zhang, Anna E Choromanska, and Yann LeCun. Deep learning with elastic averaging sgd. Advances in neural information processing systems, 28, 2015.
- [65] Shen-Yi Zhao, Yin-Peng Xie, and Wu-Jun Li. On the convergence and improvement of stochastic normalized gradient descent. *Science China Information Sciences*, 64:1–13, 2021.
- [66] Shen-Yi Zhao, Chang-Wei Shi, Yin-Peng Xie, and Wu-Jun Li. Stochastic normalized gradient descent with momentum for large-batch training. *Science China Information Sciences*, 67(11):212101, 2024.
- [67] Vladimir Antonovich Zorich and Octavio Paniagua. Mathematical analysis II, volume 220. Springer, 2016.

APPENDIX

Power of Generalized Smoothness in Stochastic Convex Optimization: First- and Zero-Order Algorithms

A Auxiliary Results

In this section we provide auxiliary materials that are used in the proof of Theorems.

A.1 Basic inequalities and assumptions

Basic inequalities. For all $a, b \in \mathbb{R}^d$ $(d \ge 1)$ the following equality holds:

$$2\langle a,b\rangle - \|b\|^2 = \|a\|^2 - \|a-b\|^2,$$
(5)

$$\langle a, b \rangle \le \|a\| \cdot \|b\|. \tag{6}$$

Squared norm of the sum For all $a_1, ..., a_n \in \mathbb{R}^d$, where $n = \{2, 3\}$ $\|a_1 + ... + a_n\|_2^2 \le n\|a_1\|_2^2 + ... + n\|a_n\|_2^2.$ (7)

Generalized-Lipschitz-smoothness. Throughout this paper, we assume that the (L_0, L_1) -smoothness condition (Assumption 1.2) is satisfied. This inequality can be represented in the equivalent form for any $x, y \in \mathbb{R}^d$:

$$f(y) - f(x) \le \langle \nabla f(x), y - x \rangle + \frac{L_0 + L_1 \| \nabla f(x) \|}{2} \| y - x \|^2,$$
(8)

where $L_0, L_1 \ge 0$ for any $x \in \mathbb{R}^d$ and $||y - x|| \le \frac{1}{L_1}$.

Variance decomposition. If ξ is random vector in \mathbb{R}^d with bounded second moment, then

$$\mathbb{E}\left[\left\|\xi+a\right\|^{2}\right] = \mathbb{E}\left[\left\|\xi-\mathbb{E}\left[\xi\right]\right\|^{2}\right] + \mathbb{E}\left[\left\|\mathbb{E}\left[\xi\right]-a\right\|^{2}\right],\tag{9}$$

for any deterministic vector $a \in \mathbb{R}^d$.

A.2 Auxiliary Lemma about Generalized Smoothness

If Assumption 1.2 holds, then it also holds that $\forall x \in \mathbb{R}^d$:

$$\|\nabla f(x)\|^2 \le 2(L_0 + L_1 \|\nabla f(x)\|)(f(x) - f^*),$$
(10)

where $f^* = \inf_x f(x)$.

Proof. We start the proof by applying (8) for $y = x - \frac{1}{L_0 + L_1 \|\nabla f(x)\|} \nabla f(x)$, where $\|y - x\| = \frac{\|\nabla f(x)\|}{L_0 + L_1 \|\nabla f(x)\|} \le \frac{1}{L_1}$. Then we can obtain:

$$f^* \le f\left(x - \frac{1}{L_0 + L_1 \|\nabla f(x)\|} \nabla f(x)\right) \stackrel{(8)}{\le} f(x) - \frac{1}{2(L_0 + L_1 \|\nabla f(x)\|)} \|\nabla f(x)\|^2.$$

A.3 Wirtinger-Poincare inequality

Let f is differentiable, then for all $x \in \mathbb{R}^d$, $\gamma e \in S^d(\gamma)$:

$$\mathbb{E}\left[f(x+\gamma e)^2\right] \le \frac{\gamma^2}{d} \mathbb{E}\left[\left\|\nabla f(x+\gamma e)\right\|^2\right].$$
(11)

B Clipped Stochastic Gradient Descent (Proof of the Theorem 3.1)

We start by using (L_0, L_1) -smoothness (see Assumption 1.2):

$$f(x^{k+1}) - f(x^{k}) \stackrel{(8)}{\leq} \langle \nabla f(x^{k}), x^{k+1} - x^{k} \rangle + \frac{L_{0} + L_{1} \left\| \nabla f(x^{k}) \right\|}{2} \left\| x^{k+1} - x^{k} \right\|^{2} \\ = -\eta \left\langle \nabla f(x^{k}), \operatorname{clip}_{c} \left(\nabla f(x^{k}, \boldsymbol{\xi}^{k}) \right) \right\rangle \\ + \frac{\eta^{2} (L_{0} + L_{1} \left\| \nabla f(x^{k}) \right\|)}{2} \left\| \operatorname{clip}_{c} \left(\nabla f(x^{k}, \boldsymbol{\xi}^{k}) \right) \right\|^{2}.$$
(12)

Next, we consider three cases depending on the gradient norm: $\|\nabla f(x^k)\| \ge c$ – the full gradient is clipped and $\frac{c}{2} \le \|\nabla f(x^k)\| \le c$ and $\|\nabla f(x^k)\| \le \frac{c}{2}$ – the full gradient is not clipped.

B.1 First case: $\left\|\nabla f(x^k)\right\| \ge c$

In this case $\alpha \nabla f(x^k) = \operatorname{clip}_c \left(\nabla f(x^k) \right)$ with $\alpha = \min \left\{ 1, \frac{c}{\|\nabla f(x^k)\|} \right\} = \frac{c}{\|\nabla f(x^k)\|}$, therefore we have the following

$$\begin{split} -\eta \left\langle \nabla f(x^{k}), \operatorname{clip}_{c}\left(f(x^{k}, \boldsymbol{\xi}^{k})\right)\right\rangle &\stackrel{(5)}{=} -\frac{\alpha\eta}{2} \left\|\nabla f(x^{k})\right\|^{2} - \frac{\eta}{2\alpha} \left\|\operatorname{clip}_{c}\left(\nabla f(x^{k}, \boldsymbol{\xi}^{k})\right)\right\|^{2} \\ &\quad + \frac{\eta}{2\alpha} \left\|\operatorname{clip}_{c}\left(\nabla f(x^{k}, \boldsymbol{\xi}^{k})\right) - \alpha \nabla f(x^{k})\right\|^{2} \\ &\quad = -\frac{\alpha\eta}{2} \left\|\nabla f(x^{k})\right\|^{2} - \frac{\eta}{2\alpha} \left\|\operatorname{clip}_{c}\left(\nabla f(x^{k}, \boldsymbol{\xi}^{k})\right)\right\|^{2} \\ &\quad + \frac{\eta}{2\alpha} \left\|\operatorname{clip}_{c}\left(\nabla f(x^{k}, \boldsymbol{\xi}^{k})\right) - \operatorname{clip}_{c}\left(\nabla f(x^{k})\right)\right\|^{2} \\ &\quad = -\frac{c\eta}{2} \left\|\nabla f(x^{k})\right\| - \frac{\eta}{2\alpha} \left\|\operatorname{clip}_{c}\left(\nabla f(x^{k}, \boldsymbol{\xi}^{k})\right)\right\|^{2} \\ &\quad + \frac{\eta}{2\alpha} \left\|\operatorname{clip}_{c}\left(\nabla f(x^{k}, \boldsymbol{\xi}^{k})\right) - \operatorname{clip}_{c}\left(\nabla f(x^{k})\right)\right\|^{2} \end{split}$$

Using that clipping is a projection on onto a convex set, namely ball with radius *c*, and thus is Lipshitz operator with Lipshitz constant 1, we can obtain:

$$-\eta \left\langle \nabla f(x^{k}), \mathbb{E}\left[\operatorname{clip}_{c}\left(f(x^{k},\boldsymbol{\xi}^{k})\right)\right]\right\rangle \leq -\frac{c\eta}{2} \left\|\nabla f(x^{k})\right\| - \frac{\eta}{2\alpha} \mathbb{E}\left[\left\|\operatorname{clip}_{c}\left(\nabla f(x^{k},\boldsymbol{\xi}^{k})\right)\right\|^{2}\right] \\ + \frac{\eta}{2\alpha} \mathbb{E}\left[\left\|\nabla f(x^{k},\boldsymbol{\xi}^{k}) - \nabla f(x^{k})\right\|^{2}\right] \\ \leq -\frac{c\eta}{2} \left\|\nabla f(x^{k})\right\| - \frac{\eta}{2\alpha} \mathbb{E}\left[\left\|\operatorname{clip}_{c}\left(\nabla f(x^{k},\boldsymbol{\xi}^{k})\right)\right\|^{2}\right] \\ + \frac{\eta\sigma^{2}}{2\alpha B} \\ = -\frac{c\eta}{2} \left\|\nabla f(x^{k})\right\| - \frac{\eta}{2\alpha} \mathbb{E}\left[\left\|\operatorname{clip}_{c}\left(\nabla f(x^{k},\boldsymbol{\xi}^{k})\right)\right\|^{2}\right] \\ + \frac{\eta \left\|\nabla f(x^{k})\right\|\sigma^{2}}{2cB}.$$
(13)

We now consider the cases depending on the relation between c and σ :

In the case
$$c \ge \sqrt{2}\sigma$$
 We have in (13):
 $-\eta \langle \nabla f(x^k), \mathbb{E} \left[\operatorname{clip}_c \left(\nabla f(x^k, \boldsymbol{\xi}^k) \right) \right] \rangle \stackrel{(13)}{\le} -\frac{c\eta}{2} \| \nabla f(x^k) \| - \frac{\eta}{2\alpha} \mathbb{E} \left[\| \operatorname{clip}_c \left(\nabla f(x^k, \boldsymbol{\xi}^k) \right) \|^2 \right] + \frac{\eta \| \nabla f(x^k) \| \sigma^2}{2cB} = -\frac{\eta}{2\alpha} \mathbb{E} \left[\| \operatorname{clip}_c \left(\nabla f(x^k, \boldsymbol{\xi}^k) \right) \|^2 \right]$

$$\begin{split} &-\frac{c\eta}{2} \left\| \nabla f(x^k) \right\| \left(1 - \frac{\sigma^2}{c^2 B} \right) \\ \leq &-\frac{\eta}{2\alpha} \mathbb{E} \left[\left\| \operatorname{clip}_c \left(\nabla f(x^k, \boldsymbol{\xi}^k) \right) \right\|^2 \right] \\ &- \frac{c\eta}{4} \left\| \nabla f(x^k) \right\| \\ = &-\frac{\eta \left\| \nabla f(x^k) \right\|}{2c} \mathbb{E} \left[\left\| \operatorname{clip}_c \left(\nabla f(x^k, \boldsymbol{\xi}^k) \right) \right\|^2 \right] \\ &- \frac{c\eta}{4} \left\| \nabla f(x^k) \right\|. \end{split}$$

Plugging this into (12) and choosing $\eta \leq \frac{1}{4(L_0+L_1c)}$ we have:

$$\mathbb{E}\left[\nabla f(x^{k+1})\right] - f(x^{k}) \stackrel{(12)}{\leq} -\frac{\eta \left\|\nabla f(x^{k})\right\|}{2c} \mathbb{E}\left[\left\|\operatorname{clip}_{c}\left(\nabla f(x^{k},\boldsymbol{\xi}^{k})\right)\right\|^{2}\right] - \frac{c\eta}{4} \left\|\nabla f(x^{k})\right\| \\ + \frac{\eta^{2}(L_{0} + L_{1} \left\|\nabla f(x^{k})\right\|)}{2} \mathbb{E}\left[\left\|\operatorname{clip}_{c}\left(\nabla f(x^{k},\boldsymbol{\xi}^{k})\right)\right\|^{2}\right] \\ = -\frac{\eta \left\|\nabla f(x^{k})\right\|}{2c} \mathbb{E}\left[\left\|\operatorname{clip}_{c}\left(\nabla f(x^{k},\boldsymbol{\xi}^{k})\right)\right\|^{2}\right] (1 - \eta L_{1}c) \\ - \frac{c\eta}{4} \left\|\nabla f(x^{k})\right\| + \frac{\eta^{2}L_{0}}{2} \mathbb{E}\left[\left\|\operatorname{clip}_{c}\left(\nabla f(x^{k},\boldsymbol{\xi}^{k})\right)\right\|^{2}\right] \\ \leq -\frac{c\eta}{4} \left\|\nabla f(x^{k})\right\| - \frac{\eta}{2} \mathbb{E}\left[\left\|\operatorname{clip}_{c}\left(\nabla f(x^{k},\boldsymbol{\xi}^{k})\right)\right\|^{2}\right] (1 - \eta(L_{0} + L_{1}c)) \\ \leq -\frac{c\eta}{4} \left\|\nabla f(x^{k})\right\|. \tag{14}$$

Using the convexity assumption of the function, we have the following:

$$f(x^{k}) - f^{*} \leq \left\langle \nabla f(x^{k}), x^{k} - x^{*} \right\rangle$$

$$\stackrel{(6)}{\leq} \left\| \nabla f(x^{k}) \right\| \left\| x^{k} - x^{*} \right\|$$

$$\leq \left\| \nabla f(x^{k}) \right\| \underbrace{\left\| x^{0} - x^{*} \right\|}_{R}.$$

Hence we have:

_

$$\left\|\nabla f(x^k)\right\| \ge \frac{f(x^k) - f^*}{R}.$$
(15)

Then substituting (15) into (14) we obtain:

$$\mathbb{E}\left[f(x^{k+1})\right] - f(x^k) \le -\frac{\eta c}{4} \left\|\nabla f(x^k)\right\| \le -\frac{\eta c}{4R} (f(x^k) - f^*).$$

This inequality is equivalent to the trailing inequality:

$$\mathbb{E}\left[f(x^{k+1})\right] - f^* \le \left(1 - \frac{\eta c}{4R}\right) \left(f(x^k) - f^*\right).$$

Then for k = 0, 1, 2, ..., N - 1 iterations that satisfy the conditions $\|\nabla f(x^k)\| \ge c \ge \sqrt{2}\sigma$, then ClipSGD has linear convergence

$$\mathbb{E}\left[f(x^N)\right] - f^* \le \left(1 - \frac{\eta}{2R}\right)^N \left(f(x^0) - f^*\right).$$

In the case $c \le \sqrt{2}\sigma$ We have in (13):

$$\begin{aligned} -\eta \left\langle \nabla f(x^k), \mathbb{E} \left[\mathsf{clip}_c \left(f(x^k, \boldsymbol{\xi}^k) \right) \right] \right\rangle &\stackrel{(13)}{\leq} -\frac{c\eta}{2} \left\| \nabla f(x^k) \right\| - \frac{\eta}{2\alpha} \mathbb{E} \left[\left\| \mathsf{clip}_c \left(\nabla f(x^k, \boldsymbol{\xi}^k) \right) \right\|^2 \right] \\ &+ \frac{\eta \left\| \nabla f(x^k) \right\| \sigma^2}{2cB} \end{aligned}$$

$$\begin{split} &= -\frac{c\eta}{2} \left\| \nabla f(x^k) \right\| - \frac{\eta}{2\alpha} \mathbb{E} \left[\left\| \mathsf{clip}_c \left(\nabla f(x^k, \pmb{\xi}^k) \right) \right\|^2 \right] \\ &+ \frac{\eta M \sigma^2}{2cB}. \end{split}$$

Plugging this into (12) and choosing $\eta \leq \frac{1}{4(L_0+L_1c)}$ we have:

$$\mathbb{E}\left[f(x^{k+1})\right] - f(x^{k}) \stackrel{(12)}{\leq} -\frac{c\eta}{2} \left\|\nabla f(x^{k})\right\| - \frac{\eta \left\|\nabla f(x^{k})\right\|}{2c} \mathbb{E}\left[\left\|\operatorname{clip}_{c}\left(\nabla f(x^{k},\boldsymbol{\xi}^{k})\right)\right\|^{2}\right] + \frac{\eta M \sigma^{2}}{2cB} \\
+ \frac{\eta^{2}(L_{0} + L_{1} \left\|\nabla f(x^{k})\right\|)}{2} \mathbb{E}\left[\left\|\operatorname{clip}_{c}\left(\nabla f(x^{k},\boldsymbol{\xi}^{k})\right)\right\|^{2}\right] + \frac{\eta M \sigma^{2}}{2cB} \\
= -\frac{c\eta}{2} \left\|\nabla f(x^{k})\right\| - \frac{\eta \left\|\nabla f(x^{k})\right\|}{2c} \mathbb{E}\left[\left\|\operatorname{clip}_{c}\left(\nabla f(x^{k},\boldsymbol{\xi}^{k})\right)\right\|^{2}\right] (1 - \eta L_{1}c) \\
+ \frac{\eta^{2}L_{0}}{2} \mathbb{E}\left[\left\|\operatorname{clip}_{c}\left(\nabla f(x^{k},\boldsymbol{\xi}^{k})\right)\right\|^{2}\right] + \frac{\eta M \sigma^{2}}{2cB} \\
\leq -\frac{c\eta}{2} \left\|\nabla f(x^{k})\right\| - \frac{\eta}{2} \mathbb{E}\left[\left\|\operatorname{clip}_{c}\left(\nabla f(x^{k},\boldsymbol{\xi}^{k})\right)\right\|^{2}\right] (1 - \eta (L_{0} + L_{1}c)) \\
+ \frac{\eta M \sigma^{2}}{2cB} \\
\leq -\frac{c\eta}{2} \left\|\nabla f(x^{k})\right\| + \frac{\eta M \sigma^{2}}{2cB}.$$
(16)

Using the convexity assumption of the function, we have the following:

$$f(x^{k}) - f^{*} \leq \left\langle \nabla f(x^{k}), x^{k} - x^{*} \right\rangle$$

$$\stackrel{(6)}{\leq} \left\| \nabla f(x^{k}) \right\| \left\| x^{k} - x^{*} \right\|$$

$$\leq \left\| \nabla f(x^{k}) \right\| \underbrace{\left\| x^{0} - x^{*} \right\|}_{R}.$$

Hence we have:

$$\left\|\nabla f(x^k)\right\| \ge \frac{f(x^k) - f^*}{R}.$$
(17)

Then substituting (17) into (16) we obtain:

$$\mathbb{E}\left[f(x^{k+1})\right] - f(x^k) \le -\frac{\eta c}{2} \left\|\nabla f(x^k)\right\| + \frac{\eta M \sigma^2}{2cB} \le -\frac{\eta c}{2R}(f(x^k) - f^*) + \frac{\eta M \sigma^2}{2cB}.$$

This inequality is equivalent to the trailing inequality:

$$\mathbb{E}\left[f(x^{k+1})\right] - f^* \le \left(1 - \frac{\eta c}{2R}\right) \left(f(x^k) - f^*\right) + \frac{\eta M \sigma^2}{2cB}$$

Then for k = 0, 1, 2, ..., N - 1 iterations that satisfy the conditions $\|\nabla f(x^k)\| \ge c$ and $c \le \sqrt{2}\sigma$, then ClipSGD has linear convergence

$$\mathbb{E}\left[f(x^{N})\right] - f^{*} \leq \left(1 - \frac{\eta c}{2R}\right)^{N} \left(f(x^{0}) - f^{*}\right) + \frac{MR\sigma^{2}}{c^{2}B}.$$

B.2 Second case: $\frac{c}{2} \le \left\| \nabla f(x^k) \right\| \le c$

In this case $\nabla f(x^k) = \operatorname{clip}_c \left(\nabla f(x^k) \right)$ with $\alpha = \min \left\{ 1, \frac{c}{\|\nabla f(x^k)\|} \right\} = 1$, therefore we have the following

$$-\eta \left\langle \nabla f(x^{k}), \operatorname{clip}_{c}\left(\nabla f(x^{k}, \boldsymbol{\xi}^{k})\right)\right\rangle \stackrel{(5)}{=} -\frac{\alpha\eta}{2} \left\|\nabla f(x^{k})\right\|^{2} - \frac{\eta}{2\alpha} \left\|\operatorname{clip}_{c}\left(\nabla f(x^{k}, \boldsymbol{\xi}^{k})\right)\right\|^{2} + \frac{\eta}{2\alpha} \left\|\operatorname{clip}_{c}\left(\nabla f(x^{k}, \boldsymbol{\xi}^{k})\right) - \alpha \nabla f(x^{k})\right\|^{2}$$

$$\begin{split} &= -\frac{\eta}{2} \left\| \nabla f(x^k) \right\|^2 - \frac{\eta}{2} \left\| \operatorname{clip}_c \left(\nabla f(x^k, \boldsymbol{\xi}^k) \right) \right\|^2 \\ &\quad + \frac{\eta}{2} \left\| \operatorname{clip}_c \left(\nabla f(x^k, \boldsymbol{\xi}^k) \right) - \operatorname{clip}_c \left(\nabla f(x^k) \right) \right\|^2 \\ &\leq -\frac{c\eta}{4} \left\| \nabla f(x^k) \right\| - \frac{\eta}{2} \left\| \operatorname{clip}_c \left(\nabla f(x^k, \boldsymbol{\xi}^k) \right) \right\|^2 \\ &\quad + \frac{\eta}{2} \left\| \operatorname{clip}_c \left(\nabla f(x^k, \boldsymbol{\xi}^k) \right) - \operatorname{clip}_c \left(\nabla f(x^k) \right) \right\|^2. \end{split}$$

Using that clipping is a projection on onto a convex set, namely ball with radius c, and thus is Lipshitz operator with Lipshitz constant 1, we can obtain:

$$\begin{split} -\eta \left\langle \nabla f(x^k), \mathbb{E} \left[\operatorname{clip}_c \left(\nabla f(x^k, \boldsymbol{\xi}^k) \right) \right] \right\rangle &\leq -\frac{c\eta}{4} \left\| \nabla f(x^k) \right\| - \frac{\eta}{2} \mathbb{E} \left[\left\| \operatorname{clip}_c \left(\nabla f(x^k, \boldsymbol{\xi}^k) \right) \right\|^2 \right] \\ &+ \frac{\eta}{2} \mathbb{E} \left[\left\| \nabla f(x^k, \boldsymbol{\xi}^k) - \nabla f(x^k) \right\|^2 \right] \\ &\leq -\frac{c\eta}{4} \left\| \nabla f(x^k) \right\| - \frac{\eta}{2} \mathbb{E} \left[\left\| \operatorname{clip}_c \left(\nabla f(x^k, \boldsymbol{\xi}^k) \right) \right\|^2 \right] + \frac{\eta \sigma^2}{2B} \\ &= -\frac{c\eta}{4} \left\| \nabla f(x^k) \right\| - \frac{\eta}{2} \mathbb{E} \left[\left\| \operatorname{clip}_c \left(\nabla f(x^k, \boldsymbol{\xi}^k) \right) \right\|^2 \right] + \frac{\eta \sigma^2}{2B} \end{split}$$

Plugging this into (12) and choosing $\eta \leq \frac{1}{4(L_0+L_1c)}$ we have:

$$\mathbb{E}\left[f(x^{k+1})\right] - f(x^{k}) \stackrel{(12)}{\leq} -\frac{c\eta}{4} \left\|\nabla f(x^{k})\right\| - \frac{\eta}{2} \mathbb{E}\left[\left\|\operatorname{clip}_{c}\left(\nabla f(x^{k},\boldsymbol{\xi}^{k})\right)\right\|^{2}\right] + \frac{\eta^{\sigma^{2}}}{2B} + \frac{\eta^{2}(L_{0} + L_{1} \left\|\nabla f(x^{k})\right\|)}{2} \mathbb{E}\left[\left\|\operatorname{clip}_{c}\left(\nabla f(x^{k},\boldsymbol{\xi}^{k})\right)\right\|^{2}\right] + \frac{\eta^{\sigma^{2}}}{2B} - \frac{c\eta}{4} \left\|\nabla f(x^{k})\right\| + \frac{\eta^{\sigma^{2}}}{2B} - \frac{\eta}{2} \mathbb{E}\left[\left\|\operatorname{clip}_{c}\left(\nabla f(x^{k},\boldsymbol{\xi}^{k})\right)\right\|^{2}\right] \left(1 - \eta(L_{0} + L_{1} \left\|\nabla f(x^{k})\right\|)\right) \le -\frac{c\eta}{4} \left\|\nabla f(x^{k})\right\| + \frac{\eta^{\sigma^{2}}}{2B}.$$
(18)

Using the convexity assumption of the function, we have the following:

$$f(x^k) - f^* \leq \langle \nabla f(x^k), x^k - x^* \rangle$$

$$\stackrel{(6)}{\leq} \|\nabla f(x^k)\| \| x^k - x^* \|$$

$$\leq \|\nabla f(x^k)\| \underbrace{\|x^0 - x^*\|}_R.$$

Hence we have:

$$\left\|\nabla f(x^k)\right\| \ge \frac{f(x^k) - f^*}{R}.$$
(19)

Then substituting (19) into (18) we obtain:

$$\mathbb{E}\left[f(x^{k+1})\right] - f(x^k) \le -\frac{\eta c}{4} \left\|\nabla f(x^k)\right\| + \frac{\eta \sigma^2}{2B} \le -\frac{\eta c}{4R}(f(x^k) - f^*) + \frac{\eta \sigma^2}{2B}$$

This inequality is equivalent to the trailing inequality:

$$\mathbb{E}\left[f(x^{k+1})\right] - f^* \le \left(1 - \frac{\eta c}{4R}\right) \left(f(x^k) - f^*\right) + \frac{\eta \sigma^2}{2B}.$$

Then for k = 0, 1, 2, ..., N - 1 iterations that satisfy the conditions $\frac{c}{2} \leq ||\nabla f(x^k)|| \leq c$, then ClipSGD has linear convergence

$$\mathbb{E}\left[f(x^N)\right] - f^* \le \left(1 - \frac{\eta c}{4R}\right)^N \left(f(x^0) - f^*\right) + \frac{2\sigma^2 R}{cB}.$$

Let $\mathcal{T}_1 = \left\{ m_0^{\mathcal{T}_1}, m_1^{\mathcal{T}_1}, m_2^{\mathcal{T}_1}, ..., m_{K-1}^{\mathcal{T}_1} \right\} = \left\{ k \in \{0, 1, 2, ..., N-1\} | \left\| \nabla f(x^k, \xi^k) \right\| \ge \frac{c}{2} \right\}$, where $K = |\mathcal{T}_1|$. Then for $k \in \mathcal{T}_1$ ClipSGD shows linear convergence:

$$\begin{split} F_N \cdot 1\!\!1 \left[\mathcal{T}_1\right] &\leq \left(1 - \frac{1}{4L_1 R}\right) F_N \leq \left(1 - \frac{1}{4L_1 R}\right) F_{m_{K-1}^{\tau_1}} + \frac{\eta M \sigma^2}{2cB} + \frac{\eta \sigma^2}{2B} \\ &\leq \ldots \leq \left(1 - \frac{1}{4L_1 R}\right)^K F_{m_0^{\tau_1}} + \frac{M R \sigma^2}{c^2 B} + \frac{2\sigma^2 R}{cB} \\ &\leq \left(1 - \frac{1}{4L_1 R}\right)^K F_0 + \frac{M R \sigma^2}{c^2 B} + \frac{2\sigma^2 R}{cB}, \end{split}$$

where $F_k = \mathbb{E}\left[f(x^k)\right] - f^*$, and we used that $F_k \leq F_{k-1}$.

B.3 Third case: $\left\|\nabla f(x^k)\right\| \leq \frac{c}{2}$

We introduce an indicative function:

$$\aleph_k = \mathbb{1}\left\{ \left\| \nabla f(x^k, \boldsymbol{\xi}^k) \right\| > c \right\}.$$
(20)

Then the following is true:

$$\mathbb{E}\left[\aleph_{k}\right] = \mathbb{E}\left[\aleph_{k}^{2}\right] = \mathcal{P}\left[\left\|\nabla f(x^{k},\boldsymbol{\xi}^{k})\right\| > c\right] \stackrel{\otimes}{\leq} \mathcal{P}\left[\left\|\nabla f(x^{k},\boldsymbol{\xi}^{k}) - \nabla f(x^{k})\right\| > \frac{c}{2}\right] \stackrel{\otimes}{\leq} \frac{4\sigma^{2}}{c^{2}B}, \quad (21)$$

where in ① we used $\|\nabla f(x^k, \boldsymbol{\xi}^k)\| \leq \|\nabla f(x^k, \boldsymbol{\xi}^k) - \nabla f(x^k)\| + \|\nabla f(x^k)\| \leq \|\nabla f(x^k, \boldsymbol{\xi}^k) - \nabla f(x^k)\| + \|\nabla f(x^k)\| \leq \|\nabla f(x^k, \boldsymbol{\xi}^k) - \nabla f(x^k)\| + \frac{c}{2}$, and in ② we used Markov's inequality. Let $r_{k+1} = \mathbb{E}\left[\|x^{k+1} - x^*\|\right]$ and $F_{k+1} = \mathbb{E}\left[f(x^{k+1}) - f^*\right]$, then given that $\operatorname{clip}_c\left(\nabla f(x^k, \boldsymbol{\xi}^k)\right) = \nabla f(x^k, \boldsymbol{\xi}^k)(1 - \aleph_k) + \frac{c}{\|\nabla f(x^k, \boldsymbol{\xi}^k)\|} \nabla f(x^k, \boldsymbol{\xi}^k)\aleph_k$ $= \nabla f(x^k, \boldsymbol{\xi}^k) + \left(\frac{c}{\|\nabla f(x^k, \boldsymbol{\xi}^k)\|} - 1\right) \nabla f(x^k, \boldsymbol{\xi}^k)\aleph_k$

we get with $\eta \leq \frac{1}{4(L_0+L_1c)}$:

$$\overset{\otimes}{\leq} r_{k}^{2} - 2\eta F_{k} + 2\eta \left\| \mathbb{E} \left[\left(\frac{c}{\|\nabla f(x^{k}, \boldsymbol{\xi}^{k})\|} - 1 \right) \nabla f(x^{k}, \boldsymbol{\xi}^{k}) \aleph_{k} \right] \right\| R + 2\eta^{2} \mathbb{E} \left[\left\| \nabla f(x^{k}, \boldsymbol{\xi}^{k}) - \nabla f(x^{k}) \right\|^{2} \right] + 2\eta^{2} \left\| \nabla f(x^{k}) \right\|^{2} \leq r_{k}^{2} - 2\eta F_{k} + 2\eta \left\| \mathbb{E} \left[\left(\frac{c}{\|\nabla f(x^{k}, \boldsymbol{\xi}^{k})\|} - 1 \right) \nabla f(x^{k}, \boldsymbol{\xi}^{k}) \aleph_{k} \right] \right\| R + \frac{2\eta^{2}\sigma^{2}}{B} + 2\eta^{2} \left\| \nabla f(x^{k}) \right\|^{2} \overset{(10)}{\leq} r_{k}^{2} - 2\eta F_{k} + 2\eta \left\| \mathbb{E} \left[\left(\frac{c}{\|\nabla f(x^{k}, \boldsymbol{\xi}^{k})\|} - 1 \right) \nabla f(x^{k}, \boldsymbol{\xi}^{k}) \aleph_{k} \right] \right\| R + \frac{2\eta^{2}\sigma^{2}}{B} + 4\eta^{2} \left(L_{0} + L_{1} \left\| \nabla f(x^{k}) \right\| \right) F_{k} \leq r_{k}^{2} - 2\eta F_{k} + 2\eta \left\| \mathbb{E} \left[\left(\frac{c}{\|\nabla f(x^{k}, \boldsymbol{\xi}^{k})\|} - 1 \right) \nabla f(x^{k}, \boldsymbol{\xi}^{k}) \aleph_{k} \right] \right\| R + \frac{2\eta^{2}\sigma^{2}}{B} + 4\eta^{2} \left(L_{0} + L_{1}c \right) F_{k} \\ = r_{k}^{2} - 2\eta F_{k} \left(1 - 2\eta \left(L_{0} + L_{1}c \right) \right) + \frac{2\eta^{2}\sigma^{2}}{B} + 2\eta \left\| \mathbb{E} \left[\left(\frac{c}{\|\nabla f(x^{k}, \boldsymbol{\xi}^{k})\|} - 1 \right) \nabla f(x^{k}, \boldsymbol{\xi}^{k}) \aleph_{k} \right] \right\| R \\ \leq r_{k}^{2} - \eta F_{k} + \frac{2\eta^{2}\sigma^{2}}{B} + 2\eta \left\| \mathbb{E} \left[\left(\frac{c}{\|\nabla f(x^{k}, \boldsymbol{\xi}^{k})\|} - 1 \right) \nabla f(x^{k}, \boldsymbol{\xi}^{k}) \aleph_{k} \right] \right\| R.$$
 (22)

Let's find the upper bound of the last summand:

$$2\eta R \left\| \mathbb{E} \left[\left(\frac{c}{\|\nabla f(x^{k}, \boldsymbol{\xi}^{k})\|} - 1 \right) \nabla f(x^{k}, \boldsymbol{\xi}^{k}) \aleph_{k} \right] \right\|$$

$$\stackrel{(20)}{\leq} 2\eta R \mathbb{E} \left[\|\nabla f(x^{k}, \boldsymbol{\xi}^{k})\| \cdot \left(1 - \frac{c}{\|\nabla f(x^{k}, \boldsymbol{\xi}^{k})\|} \right) \aleph_{k} \right]$$

$$\leq 2\eta R \mathbb{E} \left[\|\nabla f(x^{k}, \boldsymbol{\xi}^{k})\| \cdot \aleph_{k} \right]$$

$$\leq 2\eta R \left(\mathbb{E} \left[\|\nabla f(x^{k}, \boldsymbol{\xi}^{k}) - \nabla f(x^{k})\| \cdot \aleph_{k} \right] + \|\nabla f(x^{k})\| \mathbb{E} [\aleph_{k}] \right)$$

$$\leq 2\eta R \left(\sqrt{\mathbb{E} \left[\|\nabla f(x^{k}, \boldsymbol{\xi}^{k}) - \nabla f(x^{k})\|^{2} \right] \cdot \mathbb{E} [\aleph_{k}^{2}]} + \|\nabla f(x^{k})\| \mathbb{E} [\aleph_{k}] \right)$$

$$\stackrel{(21)}{\leq} 2\eta R \left(\frac{2\sigma^{2}}{cB} + \frac{c}{2} \cdot \frac{4\sigma^{2}}{c^{2}B} \right)$$

$$= \frac{8\eta\sigma^{2}R}{cB}.$$
(23)

Substituting into the initial formula and rearrange the summands, we obtain

$$\eta F_{k} \stackrel{(22)}{\leq} r_{k}^{2} - r_{k+1}^{2} + \frac{2\eta^{2}\sigma^{2}}{B} + 2\eta \left\| \mathbb{E} \left[\left(\frac{c}{\|\nabla f(x^{k}, \boldsymbol{\xi}^{k})\|} - 1 \right) \nabla f(x^{k}, \boldsymbol{\xi}^{k}) \aleph_{k} \right] \right\| R$$

$$\stackrel{(23)}{\leq} r_{k}^{2} - r_{k+1}^{2} + \frac{2\eta^{2}\sigma^{2}}{B} + \frac{8\eta\sigma^{2}R}{cB}$$

Let $\mathcal{T}_2 = \left\{ m_0^{\mathcal{T}_2}, m_1^{\mathcal{T}_2}, m_2^{\mathcal{T}_2}, ..., m_{N-K}^{\mathcal{T}_2} \right\} = \left\{ k \in \{0, 1, 2, ..., N-1\} | \| \nabla f(x^k) \| < \frac{c}{2} \right\}$, where $|\mathcal{T}_2| = N - K$. Then rearranging and summing over all $k \in \mathcal{T}_2$ we obtain

$$F_N \cdot \mathbb{1}[\mathcal{T}_2] \le \frac{1}{N-K} \sum_{k \in \mathcal{T}_2} F_k \le \frac{1}{\eta(N-K)} \sum_{k \in \mathcal{T}_2} \left(r_k^2 - r_{k+1}^2 \right) + \frac{1}{N-K} \sum_{k \in \mathcal{T}_2} \left(\frac{2\eta\sigma^2}{B} + \frac{8\sigma^2 R}{cB} \right)$$

$$= \frac{r_0^2 - r_N^2}{\eta(N - K)} + \frac{2\eta\sigma^2}{B} + \frac{8\sigma^2 R}{cB} \le \frac{r_0^2}{\eta(N - K)} + \frac{2\eta\sigma^2}{B} + \frac{8\sigma^2 R}{cB}$$

Hence we obtain:

$$F_N \cdot \mathbb{1}\left[\mathcal{T}_2\right] \le \frac{R^2}{\eta(N-K)} + \frac{2\eta\sigma^2}{B} + \frac{8\sigma^2 R}{cB}.$$

Combining all cases we have:

$$\mathbb{E}\left[f(x^{N})\right] - f^{*} \leq F_{N} \cdot \mathbb{1}\left[\mathcal{T}_{1}\right] + F_{N} \cdot \mathbb{1}\left[\mathcal{T}_{2}\right]$$
$$\leq \left(1 - \frac{\eta c}{2R}\right)^{K} F_{0} + \frac{R^{2}}{\eta(N-K)} + \frac{\sigma^{2}MR}{c^{2}B} + \frac{2\eta\sigma^{2}}{B} + \frac{8\sigma^{2}R}{cB}.$$

C Normalized Stochastic Gradient Descent (Proof of the Theorem 4.1)

Let's introduce the notation $G(x^k, \boldsymbol{\xi}^k) = \frac{\nabla f(x^k, \boldsymbol{\xi}^k)}{\|\nabla f(x^k, \boldsymbol{\xi}^k)\|}$, then using (L_0, L_1) -smoothness (see Assumption 1.2):

$$f(x^{k+1}) - f(x^k) \stackrel{(8)}{\leq} \left\langle \nabla f(x^k), x^{k+1} - x^k \right\rangle + \frac{L_0 + L_1 \left\| \nabla f(x^k) \right\|}{2} \left\| x^{k+1} - x^k \right\|^2$$
$$= -\eta \left\langle \nabla f(x^k), G(x^k, \boldsymbol{\xi}^k) \right\rangle + \frac{\eta^2 (L_0 + L_1 \left\| \nabla f(x^k) \right\|)}{2} \left\| G(x^k, \boldsymbol{\xi}^k) \right\|^2.$$
(24)

Next, we consider 4 cases of the relation $\|\nabla f(x^k)\|$ and $\|\nabla f(x^k, \boldsymbol{\xi}^k)\|$ with respect to the hyperparameter λ .

C.1 First case: $\left\|\nabla f(x^k)\right\| \ge \lambda$ and $\left\|\nabla f(x^k, \boldsymbol{\xi}^k)\right\| \ge \lambda$

Let us evaluate first summand of (24) with $\alpha = \left\| \nabla f(x^k) \right\|^{-1}$:

$$\begin{aligned} -\eta \left\langle \nabla f(x^{k}), G(x^{k}, \boldsymbol{\xi}^{k}) \right\rangle &\stackrel{(5)}{=} -\frac{\alpha \eta}{2} \left\| \nabla f(x^{k}) \right\|^{2} - \frac{\eta}{2\alpha} \left\| G(x^{k}, \boldsymbol{\xi}^{k}) \right\|^{2} \\ &+ \frac{\eta}{2\alpha} \left\| G(x^{k}, \boldsymbol{\xi}^{k}) - \alpha \nabla f(x^{k}) \right\|^{2} \\ &= -\frac{\eta}{2} \left\| \nabla f(x^{k}) \right\| - \frac{\eta}{2\alpha} \left\| G(x^{k}, \boldsymbol{\xi}^{k}) \right\|^{2} \\ &+ \frac{\eta}{2\lambda^{2}\alpha} \left\| \lambda G(x^{k}, \boldsymbol{\xi}^{k}) - \lambda \alpha \nabla f(x^{k}) \right\|^{2} \\ &= -\frac{\eta}{2} \left\| \nabla f(x^{k}) \right\| - \frac{\eta}{2\alpha} \left\| G(x^{k}, \boldsymbol{\xi}^{k}) \right\|^{2} \\ &+ \frac{\eta}{2\lambda^{2}\alpha} \left\| \operatorname{clip}_{\lambda} \left(\nabla f(x^{k}, \boldsymbol{\xi}^{k}) \right) - \operatorname{clip}_{\lambda} \left(\nabla f(x^{k}) \right) \right\|^{2} \end{aligned}$$

Using that clipping is a projection on onto a convex set, namely ball with radius λ , and thus is Lipshitz operator with Lipshitz constant 1, we can obtain:

$$-\eta \left\langle \nabla f(x^{k}), \mathbb{E}\left[G(x^{k}, \boldsymbol{\xi}^{k})\right]\right\rangle \leq -\frac{\eta}{2} \left\|\nabla f(x^{k})\right\| - \frac{\eta}{2\alpha} \mathbb{E}\left[\left\|G(x^{k}, \boldsymbol{\xi}^{k})\right\|^{2}\right] + \frac{\eta}{2\lambda^{2}\alpha} \mathbb{E}\left[\left\|\nabla f(x^{k}, \boldsymbol{\xi}^{k}) - \nabla f(x^{k})\right\|^{2}\right].$$
(25)

In the case:
$$0 \le \sigma \le \frac{\lambda}{\sqrt{2}}$$
. Using this in (25), we have the following with $\eta_k \le \frac{\|\nabla f(x^k)\|}{2(L_0+L_1\|\nabla f(x^k)\|)}$:

$$\mathbb{E}\left[f(x^{k+1})\right] - f(x^k) \stackrel{(24)}{\le} -\eta \left\langle \nabla f(x^k), \mathbb{E}\left[G(x^k, \boldsymbol{\xi}^k)\right]\right\rangle + \frac{\eta^2(L_0+L_1\|\nabla f(x^k)\|)}{2} \mathbb{E}\left[\left\|G(x^k, \boldsymbol{\xi}^k)\right\|^2\right]$$

$$\stackrel{(25)}{\le} -\frac{\eta}{2} \left\|\nabla f(x^k)\right\| - \frac{\eta}{2\alpha} \mathbb{E}\left[\left\|G(x^k, \boldsymbol{\xi}^k)\right\|^2\right] + \frac{\eta}{2\lambda^2\alpha} \mathbb{E}\left[\left\|\nabla f(x^k, \boldsymbol{\xi}) - \nabla f(x^k)\right\|^2\right]$$

$$+ \frac{\eta^2(L_0+L_1\|\nabla f(x^k)\|)}{2} \mathbb{E}\left[\left\|G(x^k, \boldsymbol{\xi}^k)\right\|^2\right]$$

$$= -\frac{\eta}{2} \left\| \nabla f(x^{k}) \right\| + \frac{\eta}{2\lambda^{2}\alpha} \mathbb{E} \left[\left\| \nabla f(x^{k}, \boldsymbol{\xi}^{k}) - \nabla f(x^{k}) \right\|^{2} \right] - \frac{\eta}{2} \mathbb{E} \left[\left\| G(x^{k}, \boldsymbol{\xi}^{k}) \right\|^{2} \right] \left(1 - \frac{\eta (L_{0} + L_{1} \left\| \nabla f(x^{k}) \right\|)}{\left\| \nabla f(x^{k}) \right\|} \right) \leq -\frac{\eta}{2} \left\| \nabla f(x^{k}) \right\| + \frac{\eta \sigma^{2}}{2\lambda^{2}\alpha} \leq -\frac{\eta}{2} \left\| \nabla f(x^{k}) \right\| + \frac{\eta}{4} \left\| \nabla f(x^{k}) \right\| = -\frac{\eta}{4} \left\| \nabla f(x^{k}) \right\|.$$
(26)

$$\frac{\left\|\nabla f(x^k)\right\|}{2\left(L_0 + L_1 \left\|\nabla f(x^k)\right\|\right)} = \frac{1}{2\left(L_0 \frac{1}{\left\|\nabla f(x^k)\right\|} + L_1\right)} = \frac{\lambda}{2\left(L_0 \frac{\lambda}{\left\|\nabla f(x^k)\right\|} + L_1\lambda\right)} \ge \frac{\lambda}{2\left(L_0 + L_1\lambda\right)}.$$

Thus, $\eta_k = \eta \le \frac{\lambda}{2\left(L_0 + L_1\lambda\right)}.$

Using the convexity assumption of the function, we have the following:

$$f(x^{k}) - f^{*} \leq \left\langle \nabla f(x^{k}), x^{k} - x^{*} \right\rangle$$

$$\stackrel{(6)}{\leq} \left\| \nabla f(x^{k}) \right\| \left\| x^{k} - x^{*} \right\|$$

$$\leq \left\| \nabla f(x^{k}) \right\| \underbrace{\left\| x^{0} - x^{*} \right\|}_{R}$$

Hence we have:

$$\left\|\nabla f(x^k)\right\| \ge \frac{f(x^k) - f^*}{R}.$$
(27)

Then substituting (27) into (26) we obtain:

$$\mathbb{E}\left[f(x^{k+1})\right] - f(x^k) \le -\frac{\eta}{4} \left\|\nabla f(x^k)\right\| \le -\frac{\eta}{4R} (f(x^k) - f^*).$$

This inequality is equivalent to the trailing inequality:

$$\mathbb{E}\left[f(x^{k+1})\right] - f^* \le \left(1 - \frac{\eta}{4R}\right) \left(f(x^k) - f^*\right).$$

Then for k = 0, 1, 2, ..., N - 1 iterations that satisfy the conditions $\|\nabla f(x^k, \boldsymbol{\xi}^k)\| \ge \sqrt{2}\sigma$ and $\|\nabla f(x^k)\| \ge \sqrt{2}\sigma$ NSGD shows linear convergence:

$$\mathbb{E}\left[f(x^N)\right] - f^* \le \left(1 - \frac{\eta}{4R}\right)^N (f(x^0) - f^*).$$

 $\begin{aligned} \text{In the case: } \frac{\lambda}{\sqrt{2}} &\leq \sigma. \quad \text{Using this in (25), we have the following with } \eta_k \leq \frac{\|\nabla f(x^k)\|}{2(L_0 + L_1 \|\nabla f(x^k)\|)}:\\ \mathbb{E}\left[f(x^{k+1})\right] - f(x^k) \stackrel{(24)}{\leq} &-\eta \left\langle \nabla f(x^k), \mathbb{E}\left[G(x^k, \boldsymbol{\xi}^k)\right]\right\rangle + \frac{\eta^2(L_0 + L_1 \|\nabla f(x^k)\|)}{2} \mathbb{E}\left[\left\|G(x^k, \boldsymbol{\xi}^k)\right\|^2\right]\\ \stackrel{(25)}{\leq} &-\frac{\eta}{2} \|\nabla f(x^k)\| - \frac{\eta}{2\alpha} \mathbb{E}\left[\left\|G(x^k, \boldsymbol{\xi}^k)\right\|^2\right] + \frac{\eta}{2\lambda^2\alpha} \mathbb{E}\left[\left\|\nabla f(x^k, \boldsymbol{\xi}^k) - \nabla f(x^k)\right\|^2\right]\\ &+ \frac{\eta^2(L_0 + L_1 \|\nabla f(x^k)\|)}{2} \mathbb{E}\left[\left\|G(x^k, \boldsymbol{\xi}^k) - \nabla f(x^k)\right\|^2\right]\\ &= -\frac{\eta}{2} \left\|\nabla f(x^k)\right\| + \frac{\eta}{2\lambda^2\alpha} \mathbb{E}\left[\left\|\nabla f(x^k, \boldsymbol{\xi}^k) - \nabla f(x^k)\right\|^2\right]\end{aligned}$

$$-\frac{\eta}{2}\mathbb{E}\left[\left\|G(x^{k},\boldsymbol{\xi}^{k})\right\|^{2}\right]\left(1-\frac{\eta(L_{0}+L_{1}\left\|\nabla f(x^{k})\right\|)}{\left\|\nabla f(x^{k})\right\|}\right)$$

$$\leq -\frac{\eta}{2}\left\|\nabla f(x^{k})\right\|+\frac{\eta\sigma^{2}}{2\lambda^{2}\alpha B}$$

$$\leq -\frac{\eta}{2}\left\|\nabla f(x^{k})\right\|+\frac{\eta\sigma^{2}M}{2\lambda^{2}B}.$$
(28)

$$\frac{\left\|\nabla f(x^{k})\right\|}{2\left(L_{0}+L_{1}\left\|\nabla f(x^{k})\right\|\right)} = \frac{1}{2\left(L_{0}\frac{1}{\left\|\nabla f(x^{k})\right\|}+L_{1}\right)} = \frac{\lambda}{2\left(L_{0}\frac{\lambda}{\left\|\nabla f(x^{k})\right\|}+L_{1}\lambda\right)} \ge \frac{\lambda}{2\left(L_{0}+L_{1}\lambda\right)}.$$

Thus, $\eta_k = \eta \leq \frac{\lambda}{2(L_0 + L_1\lambda)}$.

Using the convexity assumption of the function, we have the following:

$$f(x^{k}) - f^{*} \leq \left\langle \nabla f(x^{k}), x^{k} - x^{*} \right\rangle$$

$$\stackrel{(6)}{\leq} \left\| \nabla f(x^{k}) \right\| \left\| x^{k} - x^{*} \right\|$$

$$\leq \left\| \nabla f(x^{k}) \right\| \underbrace{\left\| x^{0} - x^{*} \right\|}_{R}$$

Hence we have:

$$\left\|\nabla f(x^k)\right\| \ge \frac{f(x^k) - f^*}{R}.$$
(29)

Then substituting (29) into (28) we obtain:

$$\mathbb{E}\left[f(x^{k+1})\right] - f(x^k) \le -\frac{\eta}{2} \left\|\nabla f(x^k)\right\| + \frac{\eta \sigma^2 M}{2\lambda^2 B} \le -\frac{\eta}{2R}(f(x^k) - f^*) + \frac{\eta \sigma^2 M}{2\lambda^2 B}.$$

This inequality is equivalent to the trailing inequality:

$$\mathbb{E}\left[f(x^{k+1})\right] - f^* \le \left(1 - \frac{\eta}{2R}\right) \left(f(x^k) - f^*\right) + \frac{\eta \sigma^2 M}{2\lambda^2 B}.$$

Then for k = 0, 1, 2, ..., N - 1 iterations that satisfy the conditions $\|\nabla f(x^k, \boldsymbol{\xi}^k)\| \geq \lambda$ and $\|\nabla f(x^k)\| \geq \lambda$ and $\sigma \geq \sqrt{2}\lambda$ NSGD shows linear convergence:

$$\mathbb{E}\left[f(x^N)\right] - f^* \le \left(1 - \frac{\eta}{2R}\right)^N \left(f(x^0) - f^*\right) + \frac{\sigma^2 MR}{\lambda^2 B}.$$

 $\textbf{C.2} \quad \textbf{Second case: } \left\|\nabla f(x^k)\right\| \leq \lambda \text{ and } \left\|\nabla f(x^k, \pmb{\xi}^k)\right\| \geq \lambda$

Let us evaluate first summand of (24) with $\alpha = \lambda^{-1}$:

$$\begin{split} -\eta \left\langle \nabla f(x^{k}), G(x^{k}, \boldsymbol{\xi}^{k}) \right\rangle &\stackrel{(5)}{=} -\frac{\alpha \eta}{2} \left\| \nabla f(x^{k}) \right\|^{2} - \frac{\eta}{2\alpha} \left\| G(x^{k}, \boldsymbol{\xi}^{k}) \right\|^{2} \\ &+ \frac{\eta}{2\alpha} \left\| G(x^{k}, \boldsymbol{\xi}^{k}) - \alpha \nabla f(x^{k}) \right\|^{2} \\ &\leq -\frac{\eta}{2} \left\| \nabla f(x^{k}) \right\| - \frac{\eta}{2\alpha} \left\| G(x^{k}, \boldsymbol{\xi}^{k}) \right\|^{2} \\ &+ \frac{\eta}{2\lambda} \left\| \lambda G(x^{k}, \boldsymbol{\xi}^{k}) - \nabla f(x^{k}) \right\|^{2} \\ &= -\frac{\eta}{2} \left\| \nabla f(x^{k}) \right\| - \frac{\eta}{2\alpha} \left\| G(x^{k}, \boldsymbol{\xi}^{k}) \right\|^{2} \\ &+ \frac{\eta}{2\lambda} \left\| \operatorname{clip}_{\lambda} \left(\nabla f(x^{k}, \boldsymbol{\xi}^{k}) \right) - \operatorname{clip}_{\lambda} \left(\nabla f(x^{k}) \right) \right\|^{2} \end{split}$$

Using that clipping is a projection on onto a convex set, namely ball with radius λ , and thus is Lipshitz operator with Lipshitz constant 1, we can obtain:

$$-\eta \left\langle \nabla f(x^{k}), \mathbb{E}\left[G(x^{k}, \boldsymbol{\xi}^{k})\right]\right\rangle \leq -\frac{\eta}{2} \left\|\nabla f(x^{k})\right\| - \frac{\eta}{2\alpha} \mathbb{E}\left[\left\|G(x^{k}, \boldsymbol{\xi}^{k})\right\|^{2}\right] + \frac{\eta}{2\lambda} \mathbb{E}\left[\left\|\nabla f(x^{k}, \boldsymbol{\xi}^{k}) - \nabla f(x^{k})\right\|^{2}\right].$$
(30)

Using this, we have the following with $\eta_k \leq \frac{\|\nabla f(x^k)\|}{2(L_0+L_1\|\nabla f(x^k)\|)}$:

$$\mathbb{E}\left[f(x^{k+1})\right] - f(x^{k}) \stackrel{(24)}{\leq} -\eta \left\langle \nabla f(x^{k}), \mathbb{E}\left[G(x^{k}, \boldsymbol{\xi}^{k})\right]\right\rangle + \frac{\eta^{2}(L_{0} + L_{1} \left\|\nabla f(x^{k})\right\|)}{2} \mathbb{E}\left[\left\|G(x^{k}, \boldsymbol{\xi}^{k})\right\|^{2}\right] \\
\stackrel{(30)}{\leq} -\frac{\eta}{2} \left\|\nabla f(x^{k})\right\| - \frac{\eta}{2\alpha} \mathbb{E}\left[\left\|G(x^{k}, \boldsymbol{\xi}^{k})\right\|^{2}\right] + \frac{\eta}{2\lambda} \mathbb{E}\left[\left\|\nabla f(x^{k}, \boldsymbol{\xi}^{k}) - \nabla f(x^{k})\right\|^{2}\right] \\
\quad + \frac{\eta^{2}(L_{0} + L_{1} \left\|\nabla f(x^{k})\right\|)}{2} \mathbb{E}\left[\left\|G(x^{k}, \boldsymbol{\xi}^{k})\right\|^{2}\right] \\
= -\frac{\eta}{2} \left\|\nabla f(x^{k})\right\| + \frac{\eta}{2\lambda} \mathbb{E}\left[\left\|\nabla f(x^{k}, \boldsymbol{\xi}^{k}) - \nabla f(x^{k})\right\|^{2}\right] \\
\quad - \frac{\eta}{2} \mathbb{E}\left[\left\|G(x^{k}, \boldsymbol{\xi}^{k})\right\|^{2}\right] \left(1 - \frac{\eta(L_{0} + L_{1} \left\|\nabla f(x^{k})\right\|)}{\left\|\nabla f(x^{k})\right\|}\right) \\
\leq -\frac{\eta}{2} \left\|\nabla f(x^{k})\right\| + \frac{\eta\sigma^{2}}{2\lambda B} \\
\leq -\frac{\eta}{2} \left\|\nabla f(x^{k})\right\| + \frac{\eta\sigma^{2}}{2\lambda B}.$$
(31)

The step size will be constant, depending on the hyperparameter λ :

$$\frac{\left\|\nabla f(x^k)\right\|}{2\left(L_0+L_1\left\|\nabla f(x^k)\right\|\right)} = \frac{1}{2\left(L_0\frac{1}{\left\|\nabla f(x^k)\right\|}+L_1\right)} = \frac{\lambda}{2\left(L_0\frac{\lambda}{\left\|\nabla f(x^k)\right\|}+L_1\lambda\right)} \ge \frac{\lambda}{2\left(L_0+L_1\lambda\right)}$$

Thus, $\eta_k = \eta \le \frac{\lambda}{2\left(L_0+L_1\lambda\right)}$.

Using the convexity assumption of the function, we have the following:

$$f(x^{k}) - f^{*} \leq \left\langle \nabla f(x^{k}), x^{k} - x^{*} \right\rangle$$

$$\stackrel{(6)}{\leq} \left\| \nabla f(x^{k}) \right\| \left\| x^{k} - x^{*} \right\|$$

$$\leq \left\| \nabla f(x^{k}) \right\| \underbrace{\left\| x^{0} - x^{*} \right\|}_{R}.$$

Hence we have:

$$\left\|\nabla f(x^k)\right\| \ge \frac{f(x^k) - f^*}{R}.$$
(32)

Then substituting (32) into (31) we obtain:

$$\mathbb{E}\left[f(x^{k+1})\right] - f(x^k) \le -\frac{\eta}{2} \left\|\nabla f(x^k)\right\| + \frac{\eta\sigma^2}{2\lambda B} \le -\frac{\eta}{2R}(f(x^k) - f^*) + \frac{\eta\sigma^2}{2\lambda B}.$$

This inequality is equivalent to the trailing inequality:

$$\mathbb{E}\left[f(x^{k+1})\right] - f^* \le \left(1 - \frac{\eta}{2R}\right) \left(f(x^k) - f^*\right) + \frac{\eta\sigma^2}{2\lambda B}.$$

Then for k = 0, 1, 2, ..., N - 1 iterations that satisfy the conditions $\|\nabla f(x^k)\| \leq \lambda$ and $\|\nabla f(x^k, \boldsymbol{\xi}^k)\| \geq \lambda$ NSGD shows linear convergence:

$$\mathbb{E}\left[f(x^N)\right] - f^* \le \left(1 - \frac{\eta}{2R}\right)^N \left(f(x^0) - f^*\right) + \frac{\sigma^2 R}{\lambda B}.$$

$\textbf{C.3} \quad \textbf{Third case: } \left\| \nabla f(x^k) \right\| \leq \lambda \textbf{ and } \left\| \nabla f(x^k, \pmb{\xi}^k) \right\| \leq \lambda$

Using this in (24), we have the following with $\eta_k \leq \frac{\|\nabla f(x^k)\|}{2(L_0+L_1\|\nabla f(x^k)\|)}$ and $\alpha = \|\nabla f(x^k)\|^{-1}$:

$$\mathbb{E}\left[f(x^{k+1})\right] - f(x^{k}) \stackrel{(24)}{\leq} -\eta \left\langle \nabla f(x^{k}), \mathbb{E}\left[G(x^{k}, \boldsymbol{\xi}^{k})\right]\right\rangle \\
+ \frac{\eta^{2}(L_{0} + L_{1} \|\nabla f(x^{k})\|)}{2} \mathbb{E}\left[\left\|G(x^{k}, \boldsymbol{\xi}^{k})\right\|^{2}\right] \\
= -\frac{\eta\alpha}{2} \|\nabla f(x^{k})\|^{2} - \frac{\eta}{2\alpha} \mathbb{E}\left[\left\|G(x^{k}, \boldsymbol{\xi}^{k})\right\|^{2}\right] \\
+ \frac{\eta}{2\alpha} \mathbb{E}\left[\left\|G(x^{k}, \boldsymbol{\xi}^{k}) - \alpha \nabla f(x^{k})\right\|^{2}\right] \\
+ \frac{\eta^{2}(L_{0} + L_{1} \|\nabla f(x^{k})\|)}{2} \mathbb{E}\left[\left\|G(x^{k}, \boldsymbol{\xi}^{k}) - \alpha \nabla f(x^{k})\right\|^{2}\right] \\
= -\frac{\eta}{2} \|\nabla f(x^{k})\| + \frac{\eta}{2\alpha} \mathbb{E}\left[\left\|G(x^{k}, \boldsymbol{\xi}^{k}) - \alpha \nabla f(x^{k})\right\|^{2}\right] \\
- \frac{\eta}{2} \mathbb{E}\left[\left\|G(x^{k}, \boldsymbol{\xi}^{k})\right\|^{2}\right] \left(1 - \frac{\eta(L_{0} + L_{1} \|\nabla f(x^{k})\|)}{\|\nabla f(x^{k})\|}\right) \\
\leq -\frac{\eta}{2} \|\nabla f(x^{k})\| + \frac{\eta}{2\alpha} \mathbb{E}\left[\left\|G(x^{k}, \boldsymbol{\xi}^{k}) - \alpha \nabla f(x^{k})\right\|^{2}\right] \\
= -\frac{\eta}{2} \|\nabla f(x^{k})\| + \frac{\eta}{\alpha} \mathbb{E}\left[\left\|G(x^{k}, \boldsymbol{\xi}^{k}) - \alpha \nabla f(x^{k})\right\|^{2}\right] \\
= -\frac{\eta}{2} \|\nabla f(x^{k})\| + \frac{\eta}{\alpha} \mathbb{E}\left[\left\|G(x^{k}, \boldsymbol{\xi}^{k})\right\|^{2} + \left\|\frac{\nabla f(x^{k})}{\|\nabla f(x^{k})\|}\right\|^{2}\right] \\
= -\frac{\eta}{2} \|\nabla f(x^{k})\| + \frac{2\eta\lambda}{\alpha} \mathbb{E}\left[\left\|\frac{\nabla f(x^{k}, \boldsymbol{\xi}^{k})}{\|\nabla f(x^{k}, \boldsymbol{\xi}^{k})\|}\right\|^{2} + \left\|\frac{\nabla f(x^{k})}{\|\nabla f(x^{k})\|}\right\|^{2}\right] \\
= -\frac{\eta}{2} \|\nabla f(x^{k})\| + 2\eta\lambda. \tag{33}$$

The step size will be constant, depending on the hyperparameter λ :

$$\frac{\left\|\nabla f(x^k)\right\|}{2\left(L_0+L_1\left\|\nabla f(x^k)\right\|\right)} = \frac{1}{2\left(L_0\frac{1}{\left\|\nabla f(x^k)\right\|}+L_1\right)} = \frac{\lambda}{2\left(L_0\frac{\lambda}{\left\|\nabla f(x^k)\right\|}+L_1\lambda\right)} \ge \frac{\lambda}{2\left(L_0+L_1\lambda\right)}.$$

Thus, $\eta_k = \eta \le \frac{\lambda}{2\left(L_0+L_1\lambda\right)}.$

Using the convexity assumption of the function, we have the following:

$$f(x^{k}) - f^{*} \leq \left\langle \nabla f(x^{k}), x^{k} - x^{*} \right\rangle$$

$$\stackrel{(6)}{\leq} \left\| \nabla f(x^{k}) \right\| \left\| x^{k} - x^{*} \right\|$$

$$\leq \left\| \nabla f(x^{k}) \right\| \underbrace{\left\| x^{0} - x^{*} \right\|}_{R}$$

Hence we have:

$$\left\|\nabla f(x^k)\right\| \ge \frac{f(x^k) - f^*}{R}.$$
(34)

Then substituting (34) into (33) we obtain:

$$\mathbb{E}\left[f(x^{k+1})\right] - f(x^k) \le -\frac{\eta}{2} \left\|\nabla f(x^k)\right\| + 2\eta\lambda \le -\frac{\eta}{2R}(f(x^k) - f^*) + 2\eta\lambda.$$

This inequality is equivalent to the trailing inequality:

$$\mathbb{E}\left[f(x^{k+1})\right] - f^* \le \left(1 - \frac{\eta}{2R}\right)\left(f(x^k) - f^*\right) + 2\eta\lambda.$$

Then for k = 0, 1, 2, ..., N - 1 iterations that satisfy the conditions $\|\nabla f(x^k)\| \le \lambda$ NSGD shows linear convergence:

$$\mathbb{E}\left[f(x^N)\right] - f^* \le \left(1 - \frac{\eta}{2R}\right)^N \left(f(x^0) - f^*\right) + \lambda R.$$

 $\textbf{C.4} \quad \textbf{Fourth case: } \left\|\nabla f(x^k)\right\| \geq \lambda \textbf{ and } \left\|\nabla f(x^k, \pmb{\xi}^k)\right\| \leq \lambda$

Using this in (24), we have the following with $\eta_k \leq \frac{\|\nabla f(x^k)\|}{2(L_0+L_1\|\nabla f(x^k)\|)}$ and $\alpha = \lambda^{-1}$:

$$\begin{split} \mathbb{E}\left[f(x^{k+1})\right] - f(x^{k}) \stackrel{(24)}{\leq} &-\eta \left\langle \nabla f(x^{k}), \mathbb{E}\left[G(x^{k}, \xi^{k})\right] \right\rangle \\ &+ \frac{\eta^{2}(L_{0} + L_{1} \|\nabla f(x^{k})\|)}{2} \mathbb{E}\left[\left\|G(x^{k}, \xi^{k})\right\|^{2}\right] \\ &= -\frac{\eta\alpha}{2\lambda} \|\nabla f(x^{k})\|^{2} - \frac{\eta}{2\alpha} \|\mathbb{E}\left[G(x^{k}, \xi^{k})\right]\|^{2} \\ &+ \frac{\eta^{2}(L_{0} + L_{1} \|\nabla f(x^{k})\|)}{2} \mathbb{E}\left[\left\|G(x^{k}, \xi^{k})\right\|^{2}\right] \\ &= -\frac{\eta}{2\lambda} \|\nabla f(x^{k})\|^{2} + \frac{\eta}{2\lambda} \|\mathbb{E}\left[\lambda G(x^{k}, \xi^{k})\right] - \nabla f(x^{k})\|^{2} \\ &+ \frac{\eta^{2}(L_{0} + L_{1} \|\nabla f(x^{k})\|)}{2} \\ &= -\frac{\eta}{2\lambda} \|\nabla f(x^{k})\|^{2} + \frac{\eta}{2\lambda} \|\mathbb{E}\left[\left(\frac{\lambda \nabla f(x^{k}, \xi^{k})}{\|\nabla f(x^{k}, \xi^{k})\|} - \nabla f(x^{k}, \xi^{k})\right]\right]\right|^{2} \\ &+ \frac{\eta^{2}(L_{0} + L_{1} \|\nabla f(x^{k})\|)}{2} \\ &= -\frac{\eta}{2\lambda} \|\nabla f(x^{k})\|^{2} + \frac{\eta}{2\lambda} \|\mathbb{E}\left[\left(\frac{\lambda}{\|\nabla f(x^{k}, \xi^{k})\|} - 1\right) \nabla f(x^{k}, \xi^{k})\right]\right\|^{2} \\ &+ \frac{\eta^{2}(L_{0} + L_{1} \|\nabla f(x^{k})\|)}{2} \\ &\leq -\frac{\eta}{2\lambda} \|\nabla f(x^{k})\|^{2} + \frac{\eta}{2\lambda} \mathbb{E}\left[\left(\frac{\lambda}{\|\nabla f(x^{k}, \xi^{k})\|} - 1\right)^{2} \|\nabla f(x^{k}, \xi^{k})\|^{2}\right] \\ &+ \frac{\eta^{2}(L_{0} + L_{1} \|\nabla f(x^{k})\|)}{2} \\ &\leq -\frac{\eta}{2\lambda} \|\nabla f(x^{k})\|^{2} + \frac{\eta}{2\lambda} \mathbb{E}\left[\left(\frac{\lambda}{\|\nabla f(x^{k}, \xi^{k})\|} - 1\right)^{2} \|\nabla f(x^{k}, \xi^{k})\|^{2}\right] \\ &+ \frac{\eta^{2}(L_{0} + L_{1} \|\nabla f(x^{k})\|)}{2} \\ &\leq -\frac{\eta}{2\lambda} \|\nabla f(x^{k})\|^{2} + \frac{\eta^{2}(L_{0} + L_{1} \|\nabla f(x^{k})\|)}{2} \\ &= -\frac{\eta}{2\lambda} \|\nabla f(x^{k})\|^{2} + \frac{\eta^{2}(L_{0} + L_{1} \|\nabla f(x^{k})\|)}{2} \\ &= -\frac{\eta}{2} \|\nabla f(x^{k})\| + \frac{\eta^{2}(L_{0} + L_{1} \|\nabla f(x^{k})\|)}{2} + \frac{\eta\lambda}{2} \\ &\leq -\frac{\eta}{2} \|\nabla f(x^{k})\| + \frac{\eta^{2}(L_{0} + L_{1} \|\nabla f(x^{k})\|)}{\|\nabla f(x^{k})\|} + \frac{\eta\lambda}{2} \\ &\leq -\frac{\eta}{4} \|\nabla f(x^{k})\| + \frac{\eta\lambda}{2}. \end{aligned}$$

(35)

$$\frac{\left\|\nabla f(x^k)\right\|}{2\left(L_0 + L_1 \left\|\nabla f(x^k)\right\|\right)} = \frac{1}{2\left(L_0 \frac{1}{\left\|\nabla f(x^k)\right\|} + L_1\right)} = \frac{\lambda}{2\left(L_0 \frac{\lambda}{\left\|\nabla f(x^k)\right\|} + L_1\lambda\right)} \ge \frac{\lambda}{2\left(L_0 + L_1\lambda\right)}.$$
Thus $m = n \le \frac{\lambda}{2}$

Thus, $\eta_k = \eta \leq \frac{\lambda}{2(L_0 + L_1\lambda)}$.

Using the convexity assumption of the function, we have the following:

$$f(x^{k}) - f^{*} \leq \left\langle \nabla f(x^{k}), x^{k} - x^{*} \right\rangle$$

$$\stackrel{(6)}{\leq} \left\| \nabla f(x^{k}) \right\| \left\| x^{k} - x^{*} \right\|$$

$$\leq \left\| \nabla f(x^{k}) \right\| \underbrace{\left\| x^{0} - x^{*} \right\|}_{R}$$

Hence we have:

$$\left\|\nabla f(x^k)\right\| \ge \frac{f(x^k) - f^*}{R}.$$
(36)

Then substituting (36) into (35) we obtain:

$$\mathbb{E}\left[f(x^{k+1})\right] - f(x^k) \le -\frac{\eta}{4} \left\|\nabla f(x^k)\right\| + \frac{\eta\lambda}{2} \le -\frac{\eta}{4R}(f(x^k) - f^*) + \frac{\eta\lambda}{2}$$

This inequality is equivalent to the trailing inequality:

$$\mathbb{E}\left[f(x^{k+1})\right] - f^* \le \left(1 - \frac{\eta}{4R}\right) \left(f(x^k) - f^*\right) + \frac{\eta\lambda}{2}.$$

Then for k = 0, 1, 2, ..., N - 1 iterations that satisfy the conditions $\|\nabla f(x^k)\| \ge \lambda$ and $\|\nabla f(x^k, \boldsymbol{\xi}^k)\| \le \lambda$ NSGD shows linear convergence:

$$\mathbb{E}\left[f(x^N)\right] - f^* \le \left(1 - \frac{\eta}{4R}\right)^N \left(f(x^0) - f^*\right) + 2\lambda R.$$

Combining all the cases considered, we obtain the convergence rate of NSGD:

$$\mathbb{E}\left[f(x^{N})\right] - f^{*} \lesssim \left(1 - \frac{\eta}{R}\right)^{N} \left(f(x^{0}) - f^{*}\right) + \frac{\sigma^{2}MR}{B\lambda^{2}} + \lambda R.$$

D Zero-Order Clipped Stochastic Gradient Descent Method

This section consists of two parts: 1) a generalization of the convergence result of ClipSGD (Algorithm 1) to the biased gradient oracle $\mathbf{g}(x^k, \boldsymbol{\xi}^k) = \nabla f(x^k, \boldsymbol{\xi}^k) + \mathbf{b}(x^k)$, where $\mathbf{b}(x^k)$ is biased bounded by $\zeta \ge 0$: $\|\mathbf{b}(x^k)\| \le \zeta$; 2) deriving convergence estimates of ZO-ClipSGD directly.

D.1 Biased Clipped Stochastic Gradient Descent Method (Proof of the Lemma 5.1)

We start by using (L_0, L_1) -smoothness (see Assumption 1.2):

$$f(x^{k+1}) - f(x^{k}) \stackrel{(8)}{\leq} \langle \nabla f(x^{k}), x^{k+1} - x^{k} \rangle + \frac{L_{0} + L_{1} \left\| \nabla f(x^{k}) \right\|}{2} \left\| x^{k+1} - x^{k} \right\|^{2} \\ = -\eta \left\langle \nabla f(x^{k}), \operatorname{clip}_{c} \left(\mathbf{g}(x^{k}, \boldsymbol{\xi}^{k}) \right) \right\rangle \\ + \frac{\eta^{2} (L_{0} + L_{1} \left\| \nabla f(x^{k}) \right\|)}{2} \left\| \operatorname{clip}_{c} \left(\mathbf{g}(x^{k}, \boldsymbol{\xi}^{k}) \right) \right\|^{2}.$$
(37)

Next, we consider three cases depending on the gradient norm: $\|\nabla f(x^k)\| \ge c$ – the full gradient is clipped and $\frac{c}{3} \le \|\nabla f(x^k)\| \le c$ and $\|\nabla f(x^k)\| \le \frac{c}{3}$ – the full gradient is not clipped.

D.1.1 First case: $\left\|\nabla f(x^k)\right\| \ge c$

In this case $\alpha \nabla f(x^k) = \operatorname{clip}_c \left(\nabla f(x^k) \right)$ with $\alpha = \min \left\{ 1, \frac{c}{\|\nabla f(x^k)\|} \right\} = \frac{c}{\|\nabla f(x^k)\|}$, therefore we have the following

$$\begin{split} -\eta \left\langle \nabla f(x^{k}), \operatorname{clip}_{c}\left(\mathbf{g}(x^{k}, \boldsymbol{\xi}^{k})\right)\right\rangle &\stackrel{(5)}{=} -\frac{\alpha\eta}{2} \left\| \nabla f(x^{k}) \right\|^{2} - \frac{\eta}{2\alpha} \left\| \operatorname{clip}_{c}\left(\mathbf{g}(x^{k}, \boldsymbol{\xi}^{k})\right) \right\|^{2} \\ &+ \frac{\eta}{2\alpha} \left\| \operatorname{clip}_{c}\left(\mathbf{g}(x^{k}, \boldsymbol{\xi}^{k})\right) - \alpha \nabla f(x^{k}) \right\|^{2} \\ &= -\frac{\alpha\eta}{2} \left\| \nabla f(x^{k}) \right\|^{2} - \frac{\eta}{2\alpha} \left\| \operatorname{clip}_{c}\left(\mathbf{g}(x^{k}, \boldsymbol{\xi}^{k})\right) \right\|^{2} \\ &+ \frac{\eta}{2\alpha} \left\| \operatorname{clip}_{c}\left(\mathbf{g}(x^{k}, \boldsymbol{\xi}^{k})\right) - \operatorname{clip}_{c}\left(\nabla f(x^{k})\right) \right\|^{2} \\ &= -\frac{c\eta}{2} \left\| \nabla f(x^{k}) \right\| - \frac{\eta}{2\alpha} \left\| \operatorname{clip}_{c}\left(\mathbf{g}(x^{k}, \boldsymbol{\xi}^{k})\right) \right\|^{2} \\ &+ \frac{\eta}{2\alpha} \left\| \operatorname{clip}_{c}\left(\mathbf{g}(x^{k}, \boldsymbol{\xi}^{k})\right) - \operatorname{clip}_{c}\left(\nabla f(x^{k})\right) \right\|^{2}. \end{split}$$

Using that clipping is a projection on onto a convex set, namely ball with radius c, and thus is Lipshitz operator with Lipshitz constant 1, we can obtain:

$$-\eta \left\langle \nabla f(x^{k}), \mathbb{E} \left[\operatorname{clip}_{c} \left(\mathbf{g}(x^{k}, \boldsymbol{\xi}^{k}) \right) \right] \right\rangle \leq -\frac{c\eta}{2} \left\| \nabla f(x^{k}) \right\| - \frac{\eta}{2\alpha} \mathbb{E} \left[\left\| \operatorname{clip}_{c} \left(\mathbf{g}(x^{k}, \boldsymbol{\xi}^{k}) \right) \right\|^{2} \right] \\ + \frac{\eta}{2\alpha} \mathbb{E} \left[\left\| \mathbf{g}(x^{k}, \boldsymbol{\xi}^{k}) - \nabla f(x^{k}) \right\|^{2} \right] \\ \stackrel{(9)}{=} -\frac{c\eta}{2} \left\| \nabla f(x^{k}) \right\| - \frac{\eta}{2\alpha} \mathbb{E} \left[\left\| \operatorname{clip}_{c} \left(\mathbf{g}(x^{k}, \boldsymbol{\xi}^{k}) \right) \right\|^{2} \right] \\ + \frac{\eta}{2\alpha} \mathbb{E} \left[\left\| \mathbf{g}(x^{k}, \boldsymbol{\xi}^{k}) - \mathbb{E} \left[\mathbf{g}(x^{k}, \boldsymbol{\xi}^{k}) \right] \right\|^{2} \right] \\ + \frac{\eta}{2\alpha} \left\| \mathbf{b}(x^{k}) \right\| \\ \leq -\frac{c\eta}{2} \left\| \nabla f(x^{k}) \right\| - \frac{\eta}{2\alpha} \mathbb{E} \left[\left\| \operatorname{clip}_{c} \left(\nabla f(x^{k}, \boldsymbol{\xi}^{k}) \right) \right\|^{2} \right] \\ + \frac{\eta \sigma^{2} M}{2cB} + \frac{\eta \left\| \nabla f(x^{k}) \right\| \zeta^{2}}{2c}. \tag{38}$$

We now consider the cases depending on the relation between c and ζ :

$$\begin{split} \boxed{\mathbf{In the case } c \geq \sqrt{2}\zeta} \quad & \text{We have in (38):} \\ & -\eta \left\langle \nabla f(x^k), \mathbb{E} \left[\text{clip}_c \left(\mathbf{g}(x^k, \boldsymbol{\xi}^k) \right) \right] \right\rangle \stackrel{(38)}{\leq} -\frac{c\eta}{2} \left\| \nabla f(x^k) \right\| - \frac{\eta}{2\alpha} \mathbb{E} \left[\left\| \text{clip}_c \left(\mathbf{g}(x^k, \boldsymbol{\xi}^k) \right) \right\|^2 \right] \\ & \quad + \frac{\eta \sigma^2 M}{2cB} + \frac{\eta \left\| \nabla f(x^k) \right\| \left\langle 2}{2c} \\ & = -\frac{\eta}{2\alpha} \mathbb{E} \left[\left\| \text{clip}_c \left(\mathbf{g}(x^k, \boldsymbol{\xi}^k) \right) \right\|^2 \right] \\ & \quad - \frac{c\eta}{2} \left\| \nabla f(x^k) \right\| \left(1 - \frac{\zeta^2}{c^2} \right) + \frac{\eta \sigma^2 M}{2cB} \\ & \leq -\frac{\eta}{2\alpha} \mathbb{E} \left[\left\| \text{clip}_c \left(\mathbf{g}(x^k, \boldsymbol{\xi}^k) \right) \right\|^2 \right] - \frac{c\eta}{4} \left\| \nabla f(x^k) \right\| + \frac{\eta \sigma^2 M}{2cB} \\ & = -\frac{\eta \left\| \nabla f(x^k) \right\|}{2c} \mathbb{E} \left[\left\| \text{clip}_c \left(\nabla f(x^k, \boldsymbol{\xi}^k) \right) \right\|^2 \right] \\ & \quad - \frac{c\eta}{4} \left\| \nabla f(x^k) \right\| + \frac{\eta \sigma^2 M}{2cB}. \end{split}$$

Plugging this into (37) and choosing $\eta \leq \frac{1}{4(L_0+L_1c)}$ we have:

$$\mathbb{E}\left[\nabla f(x^{k+1})\right] - f(x^k) \stackrel{(12)}{\leq} -\frac{\eta \left\|\nabla f(x^k)\right\|}{2c} \mathbb{E}\left[\left\|\operatorname{clip}_c\left(\mathbf{g}(x^k, \boldsymbol{\xi}^k)\right)\right\|^2\right] - \frac{c\eta}{4} \left\|\nabla f(x^k)\right\| + \frac{\eta \sigma^2 M}{2cB}$$

$$+ \frac{\eta^{2}(L_{0} + L_{1} \|\nabla f(x^{k})\|)}{2} \mathbb{E} \left[\left\| \operatorname{clip}_{c} \left(\mathbf{g}(x^{k}, \boldsymbol{\xi}^{k}) \right) \right\|^{2} \right]$$

$$= -\frac{\eta \|\nabla f(x^{k})\|}{2c} \mathbb{E} \left[\left\| \operatorname{clip}_{c} \left(\mathbf{g}(x^{k}, \boldsymbol{\xi}^{k}) \right) \right\|^{2} \right] (1 - \eta L_{1}c) - \frac{c\eta}{4} \|\nabla f(x^{k})\|$$

$$+ \frac{\eta^{2}L_{0}}{2} \mathbb{E} \left[\left\| \operatorname{clip}_{c} \left(\mathbf{g}(x^{k}, \boldsymbol{\xi}^{k}) \right) \right\|^{2} \right] + \frac{\eta\sigma^{2}M}{2cB}$$

$$\leq -\frac{c\eta}{4} \|\nabla f(x^{k})\| - \frac{\eta}{2} \mathbb{E} \left[\left\| \operatorname{clip}_{c} \left(\mathbf{g}(x^{k}, \boldsymbol{\xi}^{k}) \right) \right\|^{2} \right] (1 - \eta(L_{0} + L_{1}c))$$

$$+ \frac{\eta\sigma^{2}M}{2cB}$$

$$\leq -\frac{c\eta}{4} \|\nabla f(x^{k})\| + \frac{\eta\sigma^{2}M}{2cB}.$$
(39)

Using the convexity assumption of the function, we have the following:

$$f(x^{k}) - f^{*} \leq \left\langle \nabla f(x^{k}), x^{k} - x^{*} \right\rangle \stackrel{(6)}{\leq} \left\| \nabla f(x^{k}) \right\| \left\| x^{k} - x^{*} \right\| \leq \left\| \nabla f(x^{k}) \right\| \underbrace{\left\| x^{0} - x^{*} \right\|}_{R}.$$

Hence we have:

$$\left\|\nabla f(x^k)\right\| \ge \frac{f(x^k) - f^*}{R}.$$
(40)

Then substituting (40) into (39) we obtain:

$$\mathbb{E}\left[f(x^{k+1})\right] - f(x^k) \le -\frac{\eta c}{4} \left\|\nabla f(x^k)\right\| + \frac{\eta \sigma^2 M}{2cB} \le -\frac{\eta c}{4R}(f(x^k) - f^*) + \frac{\eta \sigma^2 M}{2cB}.$$

This inequality is equivalent to the trailing inequality:

$$\mathbb{E}\left[f(x^{k+1})\right] - f^* \le \left(1 - \frac{\eta c}{4R}\right) \left(f(x^k) - f^*\right) + \frac{\eta \sigma^2 M}{2cB}.$$

Then for k = 0, 1, 2, ..., N - 1 iterations that satisfy the conditions $\|\nabla f(x^k)\| \ge c \ge \sqrt{2}\zeta$, then ClipSGD with biased gradient oracle has linear convergence

$$\mathbb{E}\left[f(x^N)\right] - f^* \le \left(1 - \frac{\eta}{2R}\right)^N \left(f(x^0) - f^*\right) + \frac{\sigma^2 MR}{cB}.$$

In the case $c \le \sqrt{2}\zeta$ We have in (38):

$$\begin{split} -\eta \left\langle \nabla f(x^k), \mathbb{E}\left[\operatorname{clip}_c\left(\mathbf{g}(x^k, \boldsymbol{\xi}^k)\right)\right]\right\rangle \stackrel{(38)}{\leq} &-\frac{c\eta}{2} \left\| \nabla f(x^k) \right\| - \frac{\eta}{2\alpha} \mathbb{E}\left[\left\| \operatorname{clip}_c\left(\mathbf{g}(x^k, \boldsymbol{\xi}^k)\right) \right\|^2 \right] \\ &+ \frac{\eta \sigma^2 M}{2cB} + \frac{\eta \left\| \nabla f(x^k) \right\| \zeta^2}{2c} \\ &= -\frac{c\eta}{2} \left\| \nabla f(x^k) \right\| - \frac{\eta}{2\alpha} \mathbb{E}\left[\left\| \operatorname{clip}_c\left(\mathbf{g}(x^k, \boldsymbol{\xi}^k)\right) \right\|^2 \right] \\ &+ \frac{\eta M}{2c} \left(\frac{\sigma^2}{B} + \zeta^2 \right). \end{split}$$

Plugging this into (37) and choosing $\eta \leq \frac{1}{4(L_0+L_1c)}$ we have:

$$\begin{split} \mathbb{E}\left[f(x^{k+1})\right] - f(x^k) \stackrel{(37)}{\leq} &-\frac{c\eta}{2} \left\|\nabla f(x^k)\right\| - \frac{\eta \left\|\nabla f(x^k)\right\|}{2c} \mathbb{E}\left[\left\|\operatorname{clip}_c\left(\mathbf{g}(x^k, \boldsymbol{\xi}^k)\right)\right\|^2\right] \\ &+ \frac{\eta^2 (L_0 + L_1 \left\|\nabla f(x^k)\right\|)}{2} \mathbb{E}\left[\left\|\operatorname{clip}_c\left(\mathbf{g}(x^k, \boldsymbol{\xi}^k)\right)\right\|^2\right] + \frac{\eta M}{2c} \left(\frac{\sigma^2}{B} + \zeta^2\right) \\ &= -\frac{c\eta}{2} \left\|\nabla f(x^k)\right\| - \frac{\eta \left\|\nabla f(x^k)\right\|}{2c} \mathbb{E}\left[\left\|\operatorname{clip}_c\left(\mathbf{g}(x^k, \boldsymbol{\xi}^k)\right)\right\|^2\right] (1 - \eta L_1 c) \end{split}$$

$$+ \frac{\eta^{2}L_{0}}{2}\mathbb{E}\left[\left\|\operatorname{clip}_{c}\left(\mathbf{g}(x^{k},\boldsymbol{\xi}^{k})\right)\right\|^{2}\right] + \frac{\eta M}{2c}\left(\frac{\sigma^{2}}{B} + \zeta^{2}\right)$$

$$\leq -\frac{c\eta}{2}\left\|\nabla f(x^{k})\right\| - \frac{\eta}{2}\mathbb{E}\left[\left\|\operatorname{clip}_{c}\left(\mathbf{g}(x^{k},\boldsymbol{\xi}^{k})\right)\right\|^{2}\right]\left(1 - \eta(L_{0} + L_{1}c)\right)$$

$$+ \frac{\eta M}{2c}\left(\frac{\sigma^{2}}{B} + \zeta^{2}\right)$$

$$\leq -\frac{c\eta}{2}\left\|\nabla f(x^{k})\right\| + \frac{\eta M}{2c}\left(\frac{\sigma^{2}}{B} + \zeta^{2}\right). \tag{41}$$

Using the convexity assumption of the function, we have the following:

$$f(x^{k}) - f^{*} \leq \left\langle \nabla f(x^{k}), x^{k} - x^{*} \right\rangle^{(6)} \leq \left\| \nabla f(x^{k}) \right\| \left\| x^{k} - x^{*} \right\| \leq \left\| \nabla f(x^{k}) \right\| \underbrace{\left\| x^{0} - x^{*} \right\|}_{R}.$$

Hence we have:

$$\left\|\nabla f(x^k)\right\| \ge \frac{f(x^k) - f^*}{R}.$$
(42)

Then substituting (42) into (41) we obtain:

$$\mathbb{E}\left[f(x^{k+1})\right] - f(x^k) \le -\frac{\eta c}{2} \left\|\nabla f(x^k)\right\| + \frac{\eta M}{2c} \left(\frac{\sigma^2}{B} + \zeta^2\right) \le -\frac{\eta c}{2R} (f(x^k) - f^*) + \frac{\eta M}{2c} \left(\frac{\sigma^2}{B} + \zeta^2\right).$$

This inequality is equivalent to the trailing inequality:

$$\mathbb{E}\left[f(x^{k+1})\right] - f^* \le \left(1 - \frac{\eta c}{2R}\right) \left(f(x^k) - f^*\right) + \frac{\eta M}{2c} \left(\frac{\sigma^2}{B} + \zeta^2\right).$$

Then for k = 0, 1, 2, ..., N - 1 iterations that satisfy the conditions $\|\nabla f(x^k)\| \ge c$ and $c \le \sqrt{2}\zeta$, then ClipSGD has linear convergence

$$\mathbb{E}\left[f(x^N)\right] - f^* \le \left(1 - \frac{\eta c}{2R}\right)^N \left(f(x^0) - f^*\right) + \frac{MR}{c^2} \left(\frac{\sigma^2}{B} + \zeta^2\right).$$

D.1.2 Second case: $\frac{c}{3} \leq \left\| \nabla f(x^k) \right\| \leq c$

In this case $\nabla f(x^k) = \operatorname{clip}_c \left(\nabla f(x^k) \right)$ with $\alpha = \min \left\{ 1, \frac{c}{\|\nabla f(x^k)\|} \right\} = 1$, therefore we have the following

$$\begin{split} -\eta \left\langle \nabla f(x^{k}), \operatorname{clip}_{c}\left(\mathbf{g}(x^{k}, \boldsymbol{\xi}^{k})\right)\right\rangle &\stackrel{(5)}{=} -\frac{\alpha\eta}{2} \left\| \nabla f(x^{k}) \right\|^{2} - \frac{\eta}{2\alpha} \left\| \operatorname{clip}_{c}\left(\mathbf{g}(x^{k}, \boldsymbol{\xi}^{k})\right) - \alpha \nabla f(x^{k}) \right\|^{2} \\ &+ \frac{\eta}{2\alpha} \left\| \operatorname{clip}_{c}\left(\mathbf{g}(x^{k}, \boldsymbol{\xi}^{k})\right) - \alpha \nabla f(x^{k}) \right\|^{2} \\ &= -\frac{\eta}{2} \left\| \nabla f(x^{k}) \right\|^{2} - \frac{\eta}{2} \left\| \operatorname{clip}_{c}\left(\mathbf{g}(x^{k}, \boldsymbol{\xi}^{k})\right) \right\|^{2} \\ &+ \frac{\eta}{2} \left\| \operatorname{clip}_{c}\left(\mathbf{g}(x^{k}, \boldsymbol{\xi}^{k})\right) - \operatorname{clip}_{c}\left(\nabla f(x^{k})\right) \right\|^{2} \\ &\leq -\frac{c\eta}{6} \left\| \nabla f(x^{k}) \right\| - \frac{\eta}{2} \left\| \operatorname{clip}_{c}\left(\mathbf{g}(x^{k}, \boldsymbol{\xi}^{k})\right) \right\|^{2} \\ &+ \frac{\eta}{2} \left\| \operatorname{clip}_{c}\left(\mathbf{g}(x^{k}, \boldsymbol{\xi}^{k})\right) - \operatorname{clip}_{c}\left(\nabla f(x^{k})\right) \right\|^{2}. \end{split}$$

Using that clipping is a projection on onto a convex set, namely ball with radius c, and thus is Lipshitz operator with Lipshitz constant 1, we can obtain:

$$\begin{aligned} -\eta \left\langle \nabla f(x^{k}), \mathbb{E}\left[\operatorname{clip}_{c}\left(\mathbf{g}(x^{k},\boldsymbol{\xi}^{k})\right)\right]\right\rangle &\leq -\frac{c\eta}{6} \left\|\nabla f(x^{k})\right\| - \frac{\eta}{2} \mathbb{E}\left[\left\|\operatorname{clip}_{c}\left(\mathbf{g}(x^{k},\boldsymbol{\xi}^{k})\right)\right\|^{2}\right] \\ &+ \frac{\eta}{2} \mathbb{E}\left[\left\|\mathbf{g}(x^{k},\boldsymbol{\xi}^{k}) - \nabla f(x^{k})\right\|^{2}\right] \\ &\stackrel{(9)}{=} -\frac{c\eta}{6} \left\|\nabla f(x^{k})\right\| - \frac{\eta}{2} \mathbb{E}\left[\left\|\operatorname{clip}_{c}\left(\mathbf{g}(x^{k},\boldsymbol{\xi}^{k})\right)\right\|^{2}\right] \end{aligned}$$

$$\begin{split} &+ \frac{\eta}{2} \mathbb{E} \left[\left\| \mathbf{g}(x^{k}, \boldsymbol{\xi}^{k}) - \mathbb{E} \left[\mathbf{g}(x^{k}, \boldsymbol{\xi}^{k}) \right] \right\|^{2} \right] \\ &+ \frac{\eta}{2} \left\| \mathbf{b}(x^{k}) \right\|^{2} \\ &\leq - \frac{c\eta}{6} \left\| \nabla f(x^{k}) \right\| - \frac{\eta}{2} \mathbb{E} \left[\left\| \operatorname{clip}_{c} \left(\mathbf{g}(x^{k}, \boldsymbol{\xi}^{k}) \right) \right\|^{2} \right] \\ &+ \frac{\eta}{2} \left(\frac{\sigma^{2}}{B} + \zeta^{2} \right) \\ &= - \frac{c\eta}{6} \left\| \nabla f(x^{k}) \right\| - \frac{\eta}{2} \mathbb{E} \left[\left\| \operatorname{clip}_{c} \left(\mathbf{g}(x^{k}, \boldsymbol{\xi}^{k}) \right) \right\|^{2} \right] \\ &+ \frac{\eta}{2} \left(\frac{\sigma^{2}}{B} + \zeta^{2} \right). \end{split}$$

Plugging this into (37) and choosing $\eta \leq \frac{1}{4(L_0+L_1c)}$ we have:

$$\mathbb{E}\left[f(x^{k+1})\right] - f(x^k) \stackrel{(37)}{\leq} -\frac{c\eta}{6} \left\|\nabla f(x^k)\right\| - \frac{\eta}{2} \mathbb{E}\left[\left\|\operatorname{clip}_c\left(\mathbf{g}(x^k, \boldsymbol{\xi}^k)\right)\right\|^2\right] \\
+ \frac{\eta^2(L_0 + L_1 \left\|\nabla f(x^k)\right\|\right)}{2} \mathbb{E}\left[\left\|\operatorname{clip}_c\left(\mathbf{g}(x^k, \boldsymbol{\xi}^k)\right)\right\|^2\right] + \frac{\eta}{2} \left(\frac{\sigma^2}{B} + \zeta^2\right) \\
= -\frac{c\eta}{6} \left\|\nabla f(x^k)\right\| - \frac{\eta}{2} \mathbb{E}\left[\left\|\operatorname{clip}_c\left(\mathbf{g}(x^k, \boldsymbol{\xi}^k)\right)\right\|^2\right] \left(1 - \eta(L_0 + L_1 \left\|\nabla f(x^k)\right\|)\right) \\
+ \frac{\eta}{2} \left(\frac{\sigma^2}{B} + \zeta^2\right) \\
\leq -\frac{c\eta}{6} \left\|\nabla f(x^k)\right\| + \frac{\eta}{2} \left(\frac{\sigma^2}{B} + \zeta^2\right).$$
(43)

Using the convexity assumption of the function, we have the following:

$$f(x^k) - f^* \le \left\langle \nabla f(x^k), x^k - x^* \right\rangle \stackrel{\text{(6)}}{\le} \left\| \nabla f(x^k) \right\| \left\| x^k - x^* \right\| \le \left\| \nabla f(x^k) \right\| \underbrace{\left\| x^0 - x^* \right\|}_R.$$

Hence we have:

$$\left\|\nabla f(x^k)\right\| \ge \frac{f(x^k) - f^*}{R}.$$
(44)

Then substituting (44) into (43) we obtain:

$$\mathbb{E}\left[f(x^{k+1})\right] - f(x^k) \le -\frac{\eta c}{6} \left\|\nabla f(x^k)\right\| + \frac{\eta}{2} \left(\frac{\sigma^2}{B} + \zeta^2\right) \le -\frac{\eta c}{6R} (f(x^k) - f^*) + \frac{\eta}{2} \left(\frac{\sigma^2}{B} + \zeta^2\right).$$

This inequality is equivalent to the trailing inequality:

$$\mathbb{E}\left[f(x^{k+1})\right] - f^* \le \left(1 - \frac{\eta c}{6R}\right) \left(f(x^k) - f^*\right) + \frac{\eta}{2} \left(\frac{\sigma^2}{B} + \zeta^2\right).$$

Then for k = 0, 1, 2, ..., N - 1 iterations that satisfy the conditions $\frac{c}{2} \leq ||\mathbf{g}(x^k, \boldsymbol{\xi}^k)|| \leq c$, then ClipSGD with biased gradient oracle has linear convergence

$$\mathbb{E}\left[f(x^N)\right] - f^* \le \left(1 - \frac{\eta c}{6R}\right)^N \left(f(x^0) - f^*\right) + \frac{3R}{c} \left(\frac{\sigma^2}{B} + \zeta^2\right).$$

Let $\mathcal{T}_1 = \left\{ m_0^{\mathcal{T}_1}, m_1^{\mathcal{T}_1}, m_2^{\mathcal{T}_1}, ..., m_{K-1}^{\mathcal{T}_1} \right\} = \left\{ k \in \{0, 1, 2, ..., N-1\} | \left\| \nabla f(x^k, \xi^k) \right\| \ge \frac{c}{3} \right\}$, where $K = |\mathcal{T}_1|$. Then for $k \in \mathcal{T}_1$ ClipSGD with biased gradient oracle shows linear convergence:

$$F_N \cdot \mathbb{1}\left[\mathcal{T}_1\right] \lesssim \left(1 - \frac{\eta c}{R}\right) F_N \le \left(1 - \frac{\eta c}{6R}\right) F_{m_{K-1}^{\mathcal{T}_1}} + \left(\frac{\eta M}{c} + \eta\right) \cdot \left(\frac{\sigma^2}{B} + \zeta^2\right) \le \dots \le$$

$$\leq \left(1 - \frac{\eta c}{R}\right)^{K} F_{m_{0}^{\tau_{1}}} + \left(\frac{MR}{c^{2}} + \frac{R}{c}\right) \cdot \left(\frac{\sigma^{2}}{B} + \zeta^{2}\right)$$
$$\leq \left(1 - \frac{\eta c}{R}\right)^{K} F_{0} + \left(\frac{MR}{c^{2}} + \frac{R}{c}\right) \cdot \left(\frac{\sigma^{2}}{B} + \zeta^{2}\right),$$

where $F_k = \mathbb{E}\left[f(x^k)\right] - f^*$, and we used that $F_k \leq F_{k-1}$.

D.1.3 Third case: $\left\|\nabla f(x^k)\right\| \leq \frac{c}{3}$

We introduce an indicative function:

$$\aleph_k = \mathbb{1}\left\{ \left\| \mathbf{g}(x^k, \boldsymbol{\xi}^k) \right\| > c \right\}.$$
(45)

Then the following is true:

$$\mathbb{E}\left[\aleph_{k}\right] = \mathbb{E}\left[\aleph_{k}^{2}\right] = \mathcal{P}\left[\left\|\mathbf{g}(x^{k},\boldsymbol{\xi}^{k})\right\| > c\right] \stackrel{\scriptscriptstyle{(0)}}{\leq} \mathcal{P}\left[\left\|\mathbf{g}(x^{k},\boldsymbol{\xi}^{k}) - \mathbb{E}\left[\mathbf{g}(x^{k},\boldsymbol{\xi}^{k})\right]\right\| > \frac{c}{3}\right] \stackrel{\scriptscriptstyle{(0)}}{\leq} \frac{9\sigma^{2}}{c^{2}B}, \quad (46)$$
where in \mathbb{O} we used $\left\|\mathbf{g}(x^{k},\boldsymbol{\xi}^{k})\right\| \leq \left\|\mathbf{g}(x^{k},\boldsymbol{\xi}^{k}) - \mathbb{E}\left[\mathbf{g}(x^{k},\boldsymbol{\xi}^{k})\right]\right\| + \mathbb{E}\left[\mathbf{g}(x^{k},\boldsymbol{\xi}^{k})\right] \leq \left\|\mathbf{g}(x^{k},\boldsymbol{\xi}^{k}) - \mathbb{E}\left[\mathbf{g}(x^{k},\boldsymbol{\xi}^{k})\right]\right\| + \mathbb{E}\left[\mathbf{g}(x^{k},\boldsymbol{\xi}^{k})\right] \leq \left\|\mathbf{g}(x^{k},\boldsymbol{\xi}^{k}) - \mathbb{E}\left[\mathbf{g}(x^{k},\boldsymbol{\xi}^{k})\right]\right\| + \mathbb{E}\left[\mathbf{g}(x^{k},\boldsymbol{\xi}^{k})\right] \leq \left\|\mathbf{g}(x^{k},\boldsymbol{\xi}^{k}) - \mathbb{E}\left[\mathbf{g}(x^{k},\boldsymbol{\xi}^{k})\right]\right\| + \frac{c}{2}, \quad \text{where assume that } \zeta \leq \frac{c}{3}: \text{ and in } \mathbb{O} \text{ we used Markov's inequality.}$

Let $r_{k+1} = \mathbb{E}\left[\left\|x^{k+1} - x^*\right\|\right]$ and $F_{k+1} = \mathbb{E}\left[f(x^{k+1}) - f^*\right]$, then given that $\operatorname{clip}_c\left(\mathbf{g}(x^k, \boldsymbol{\xi}^k)\right) = \mathbf{g}(x^k, \boldsymbol{\xi}^k)(1 - \aleph_k) + \frac{c}{\|\mathbf{g}(x^k, \boldsymbol{\xi}^k)\|}\mathbf{g}(x^k, \boldsymbol{\xi}^k)$

$$\begin{aligned} \left(\mathbf{g}(x^k, \boldsymbol{\xi}^k) \right) &= \mathbf{g}(x^k, \boldsymbol{\xi}^k) (1 - \aleph_k) + \frac{c}{\|\mathbf{g}(x^k, \boldsymbol{\xi}^k)\|} \mathbf{g}(x^k, \boldsymbol{\xi}^k) \aleph_k \\ &= \mathbf{g}(x^k, \boldsymbol{\xi}^k) + \left(\frac{c}{\|\mathbf{g}(x^k, \boldsymbol{\xi}^k)\|} - 1 \right) \mathbf{g}(x^k, \boldsymbol{\xi}^k) \aleph_k \end{aligned}$$

we get with
$$\eta \leq \frac{1}{4(L_0+L_1c)}$$
:
 $r_{k+1}^2 = r_k^2 - 2\eta \langle \mathbb{E} \left[\operatorname{clip}_c \left(\mathbf{g}(x^k, \boldsymbol{\xi}^k) \right) \right], x^k - x^* \rangle + \eta^2 \mathbb{E} \left[\left\| \operatorname{clip}_c \left(\mathbf{g}(x^k, \boldsymbol{\xi}^k) \right) \right\|^2 \right]$
 $= r_k^2 - 2\eta \langle \nabla f(x^k), x^k - x^* \rangle - 2\eta \langle \mathbb{E} \left[\left(\frac{c}{\|\mathbf{g}(x^k, \boldsymbol{\xi}^k)\|} - 1 \right) \mathbf{g}(x^k, \boldsymbol{\xi}^k) \aleph_k \right], x^k - x^* \rangle$
 $+ \eta^2 \mathbb{E} \left[\left\| \operatorname{clip}_c \left(\mathbf{g}(x^k, \boldsymbol{\xi}^k) \right) \right\|^2 \right] + 2\eta \langle \mathbf{b}(x^k), x^k - x^* \rangle$
 $\stackrel{(6)}{\leq} r_k^2 - 2\eta \langle \nabla f(x^k), x^k - x^* \rangle + 2\eta \left\| \mathbb{E} \left[\left(\frac{c}{\|\mathbf{g}(x^k, \boldsymbol{\xi}^k)\|} - 1 \right) \mathbf{g}(x^k, \boldsymbol{\xi}^k) \aleph_k \right] \right\| \|x^k - x^* \|$
 $+ \eta^2 \mathbb{E} \left[\left\| \operatorname{clip}_c \left(\mathbf{g}(x^k, \boldsymbol{\xi}^k) \right) \right\|^2 \right] + 2\eta \left\| \mathbf{b}(x^k) \right\| \|x^k - x^* \|$
 $\stackrel{(9)}{=} r_k^2 - 2\eta F_k + 2\eta \left\| \mathbb{E} \left[\left(\frac{c}{\|\mathbf{g}(x^k, \boldsymbol{\xi}^k)\|} - 1 \right) \mathbf{g}(x^k, \boldsymbol{\xi}^k) \aleph_k \right] \right\| \|x^0 - x^* \|$
 $+ \eta^2 \mathbb{E} \left[\left\| \operatorname{clip}_c \left(\mathbf{g}(x^k, \boldsymbol{\xi}^k) \right) \right\|^2 \right] + 2\eta \left\| \mathbf{b}(x^k) \right\| \|x^0 - x^* \|$
 $+ \eta^2 \mathbb{E} \left[\left\| \operatorname{clip}_c \left(\mathbf{g}(x^k, \boldsymbol{\xi}^k) \right) \right\|^2 \right] + 2\eta \left\| \mathbf{b}(x^k) \right\| \|x^0 - x^* \|$
 $= r_k^2 - 2\eta F_k + 2\eta \left\| \mathbb{E} \left[\left(\frac{c}{\|\mathbf{g}(x^k, \boldsymbol{\xi}^k)\|} - 1 \right) \mathbf{g}(x^k, \boldsymbol{\xi}^k) \aleph_k \right] \right\| \|x^0 - x^* \|$
 $= r_k^2 - 2\eta F_k + 2\eta \left\| \mathbb{E} \left[\left(\frac{c}{\|\mathbf{g}(x^k, \boldsymbol{\xi}^k)\|} - 1 \right) \mathbf{g}(x^k, \boldsymbol{\xi}^k) \aleph_k \right] \right\| R$
 $+ 2\eta^2 \mathbb{E} \left[\left\| \operatorname{clip}_c \left(\mathbf{g}(x^k, \boldsymbol{\xi}^k) \right) - \nabla f(x^k) \right\|^2 \right] + 2\eta^2 \left\| \nabla f(x^k) \right\|^2 + 2\eta \zeta \|x^0 - x^* \|$
 $\stackrel{(2)}{=} r_k^2 - 2\eta F_k + 2\eta \left\| \mathbb{E} \left[\left(\frac{c}{\|\mathbf{g}(x^k, \boldsymbol{\xi}^k)\|} - 1 \right) \mathbf{g}(x^k, \boldsymbol{\xi}^k) \aleph_k \right] \right\| R$
 $+ 2\eta^2 \mathbb{E} \left[\left\| \operatorname{clip}_c \left(\mathbf{g}(x^k, \boldsymbol{\xi}^k) \right\| - \operatorname{clip}_c \left(\nabla f(x^k) \right) \right\|^2 \right] + 2\eta^2 \left\| \nabla f(x^k) \right\|^2 + 2\eta \zeta R$
 $\stackrel{(2)}{\leq} r_k^2 - 2\eta F_k + 2\eta \left\| \mathbb{E} \left[\left(\frac{c}{\|\mathbf{g}(x^k, \boldsymbol{\xi}^k)\|} - 1 \right) \mathbf{g}(x^k, \boldsymbol{\xi}^k) \aleph_k \right] \right\| R$
 $+ 2\eta^2 \mathbb{E} \left[\left\| \mathbf{g}(x^k, \boldsymbol{\xi}^k) - \nabla f(x^k) \right\|^2 \right] + 2\eta^2 \left\| \nabla f(x^k) \right\|^2 + 2\eta \zeta R$
 $\stackrel{(9)}{=} r_k^2 - 2\eta F_k + 2\eta \left\| \mathbb{E} \left[\left(\frac{c}{\|\mathbf{g}(x^k, \boldsymbol{\xi}^k)\|} - 1 \right) \mathbf{g}(x^k, \boldsymbol{\xi}^k) \aleph_k \right] \right\| R$

$$+ 2\eta^{2}\mathbb{E}\left[\left\|\mathbf{g}(x^{k},\boldsymbol{\xi}^{k}) - \mathbb{E}\left[\mathbf{g}(x^{k},\boldsymbol{\xi}^{k})\right]\right\|^{2}\right] + 2\eta^{2}\left\|\nabla f(x^{k})\right\|^{2} + 2\eta\zeta R + 2\eta^{2}\left\|\mathbf{b}(x^{k})\right\|$$

$$\leq r_{k}^{2} - 2\eta F_{k} + 2\eta\left\|\mathbb{E}\left[\left(\frac{c}{\|\mathbf{g}(x^{k},\boldsymbol{\xi}^{k})\|} - 1\right)\mathbf{g}(x^{k},\boldsymbol{\xi}^{k})\aleph_{k}\right]\right\|R$$

$$+ \frac{2\eta^{2}\sigma^{2}}{B} + 2\eta^{2}\left\|\nabla f(x^{k})\right\|^{2} + 2\eta\zeta R + 2\eta^{2}\zeta^{2}$$

$$\left(\stackrel{(10)}{=} r_{k}^{2} - 2\eta F_{k} + 2\eta\left\|\mathbb{E}\left[\left(\frac{c}{\|\mathbf{g}(x^{k},\boldsymbol{\xi}^{k})\|} - 1\right)\mathbf{g}(x^{k},\boldsymbol{\xi}^{k})\aleph_{k}\right]\right\|R$$

$$+ \frac{2\eta^{2}\sigma^{2}}{B} + 4\eta^{2}\left(L_{0} + L_{1}\left\|\nabla f(x^{k})\right\|\right)F_{k} + 2\eta\zeta R + 2\eta^{2}\zeta^{2}$$

$$\leq r_{k}^{2} - 2\eta F_{k} + 2\eta\left\|\mathbb{E}\left[\left(\frac{c}{\|\mathbf{g}(x^{k},\boldsymbol{\xi}^{k})\|} - 1\right)\mathbf{g}(x^{k},\boldsymbol{\xi}^{k})\aleph_{k}\right]\right\|R$$

$$+ \frac{2\eta^{2}\sigma^{2}}{B} + 4\eta^{2}\left(L_{0} + L_{1}c\right)F_{k} + 2\eta\zeta R + 2\eta^{2}\zeta^{2}$$

$$= r_{k}^{2} - 2\eta F_{k}\left(1 - 2\eta\left(L_{0} + L_{1}c\right)\right) + \frac{2\eta^{2}\sigma^{2}}{B} + 2\eta\zeta R + 2\eta^{2}\zeta^{2}$$

$$= r_{k}^{2} - 2\eta F_{k}\left(1 - 2\eta\left(L_{0} + L_{1}c\right)\right) + \frac{2\eta^{2}\sigma^{2}}{B} + 2\eta\zeta R + 2\eta^{2}\zeta^{2}$$

$$+ 2\eta\left\|\mathbb{E}\left[\left(\frac{c}{\|\mathbf{g}(x^{k},\boldsymbol{\xi}^{k})\|} - 1\right)\mathbf{g}(x^{k},\boldsymbol{\xi}^{k})\aleph_{k}\right]\right\|R$$

$$\leq r_{k}^{2} - \eta F_{k} + \frac{2\eta^{2}\sigma^{2}}{B} + 2\eta\zeta R + 2\eta^{2}\zeta^{2}$$

$$+ 2\eta\left\|\mathbb{E}\left[\left(\frac{c}{\|\mathbf{g}(x^{k},\boldsymbol{\xi}^{k})\|} - 1\right)\mathbf{g}(x^{k},\boldsymbol{\xi}^{k})\aleph_{k}\right]\right\|R.$$

$$(47)$$

Let's find the upper bound of the last summand:

$$2\eta R \left\| \mathbb{E} \left[\left(\frac{c}{\|\mathbf{g}(x^{k}, \boldsymbol{\xi}^{k})\|} - 1 \right) \mathbf{g}(x^{k}, \boldsymbol{\xi}^{k}) \aleph_{k} \right] \right\|$$

$$\stackrel{(45)}{\leq} 2\eta R \mathbb{E} \left[\|\mathbf{g}(x^{k}, \boldsymbol{\xi}^{k})\| \cdot \left(1 - \frac{c}{\|\mathbf{g}(x^{k}, \boldsymbol{\xi}^{k})\|} \right) \aleph_{k} \right]$$

$$\leq 2\eta R \mathbb{E} \left[\|\mathbf{g}(x^{k}, \boldsymbol{\xi}^{k})\| \cdot \aleph_{k} \right]$$

$$\leq 2\eta R \left(\mathbb{E} \left[\|\mathbf{g}(x^{k}, \boldsymbol{\xi}^{k}) - \mathbb{E} \left[\mathbf{g}(x^{k}, \boldsymbol{\xi}^{k}) \right] \| \cdot \aleph_{k} \right] + \|\nabla f(x^{k})\| \mathbb{E} \left[\aleph_{k} \right] + \|\mathbf{b}(x^{k})\| \mathbb{E} \left[\aleph_{k} \right] \right)$$

$$\leq 2\eta R \left(\sqrt{\mathbb{E} \left[\|\mathbf{g}(x^{k}, \boldsymbol{\xi}^{k}) - \mathbb{E} \left[\mathbf{g}(x^{k}, \boldsymbol{\xi}^{k}) \right] \|^{2} \right] \cdot \mathbb{E} \left[\aleph_{k}^{2} \right]} + \frac{2c}{3} \mathbb{E} \left[\aleph_{k} \right] \right)$$

$$\stackrel{(46)}{\leq} 2\eta R \left(\frac{3\sigma^{2}}{cB} + \frac{2c}{3} \cdot \frac{9\sigma^{2}}{c^{2}B} \right)$$

$$= \frac{18\eta\sigma^{2}R}{cB}.$$
(48)

Substituting into the initial formula and rearrange the summands, we obtain

$$\eta F_{k} \stackrel{(47)}{\leq} r_{k}^{2} - r_{k+1}^{2} + \frac{2\eta^{2}\sigma^{2}}{B} + 2\eta\zeta R + 2\eta^{2}\zeta^{2} + 2\eta \left\| \mathbb{E}\left[\left(\frac{c}{\|\nabla f(x^{k}, \boldsymbol{\xi}^{k})\|} - 1 \right) \nabla f(x^{k}, \boldsymbol{\xi}^{k}) \aleph_{k} \right] \right\| R$$

$$\stackrel{(48)}{\leq} r_{k}^{2} - r_{k+1}^{2} + \frac{2\eta^{2}\sigma^{2}}{B} + \frac{18\eta\sigma^{2}R}{cB} + 2\eta\zeta R + 2\eta^{2}\zeta^{2}$$

Let $\mathcal{T}_2 = \left\{ m_0^{\mathcal{T}_2}, m_1^{\mathcal{T}_2}, m_2^{\mathcal{T}_2}, ..., m_{N-K}^{\mathcal{T}_2} \right\} = \left\{ k \in \{0, 1, 2, ..., N-1\} | \| \nabla f(x^k) \| < \frac{c}{3} \right\}$, where $|\mathcal{T}_2| = N - K$. Then rearranging and summing over all $k \in \mathcal{T}_2$ we obtain

$$F_N \cdot \mathbb{1}[\mathcal{T}_2] \le \frac{1}{N-K} \sum_{k \in \mathcal{T}_2} F_k \le \frac{1}{\eta(N-K)} \sum_{k \in \mathcal{T}_2} \left(r_k^2 - r_{k+1}^2 \right)$$

$$+ \frac{1}{N-K} \sum_{k \in \mathcal{T}_2} \left[\eta \left(\frac{\sigma^2}{B} + \zeta^2 \right) + 2R \left(\frac{9\sigma^2}{cB} + \zeta \right) \right]$$

$$= \frac{r_0^2 - r_N^2}{\eta (N-K)} \eta \left(\frac{\sigma^2}{B} + \zeta^2 \right) + 2R \left(\frac{9\sigma^2}{cB} + \zeta \right) \le \frac{r_0^2}{\eta (N-K)} \eta \left(\frac{\sigma^2}{B} + \zeta^2 \right) + 2R \left(\frac{9\sigma^2}{cB} + \zeta \right).$$

Hence we obtain:

$$F_N \cdot \mathbb{1}\left[\mathcal{T}_2\right] \le \frac{R^2}{\eta(N-K)} \eta\left(\frac{\sigma^2}{B} + \zeta^2\right) + 2R\left(\frac{9\sigma^2}{cB} + \zeta\right).$$

Combining all the cases considered, we obtain the convergence rate of ClipSGD with biased gradient oracle:

$$\mathbb{E}\left[f(x^{N})\right] - f^{*} \leq F_{N} \cdot \mathbb{1}\left[\mathcal{T}_{1}\right] + F_{N} \cdot \mathbb{1}\left[\mathcal{T}_{2}\right]$$

$$\lesssim \left(1 - \frac{\eta c}{R}\right)^{K} F_{0} + \frac{R^{2}}{\eta(N-K)} + \left(\frac{MR}{c^{2}} + \frac{R}{c} + \eta\right) \cdot \left(\frac{\sigma^{2}}{B} + \zeta^{2}\right) + R\zeta.$$
(49)

D.2 Convergence Results for ZO-ClipSGD

In order to obtain convergence results for ZO-ClipSGD it is necessary to estimate the bias and variance of the gradient approximation (4).

Bias of gradient approximation Using the variational representation of the Euclidean norm, and definition of gradient approximation (4) we can write:

$$\|\mathbb{E}\left[\mathbf{g}(x,\{\xi,e\})\right] - \nabla f(x)\| = \left\|\mathbb{E}\left[\frac{d}{2\gamma}\left(\tilde{f}(x+\gamma e,\xi) - \tilde{f}(x-\gamma e,\xi)\right)e\right] - \nabla f(x)\right\|$$

$$\stackrel{\text{@}}{=} \left\|\mathbb{E}\left[\frac{d}{\gamma}\left(f(x+\gamma e,\xi) + \delta(x+\gamma e)\right)e\right] - \nabla f(x)\right\|$$

$$\stackrel{\text{@}}{\leq} \left\|\mathbb{E}\left[\frac{d}{\gamma}f(x+\gamma e,\xi)e\right] - \nabla f(x)\right\| + \frac{d\Delta}{\gamma}$$

$$\stackrel{\text{@}}{=} \left\|\mathbb{E}\left[\nabla f(x+\gamma u,\xi)\right] - \nabla f(x)\right\| + \frac{d\Delta}{\gamma}$$

$$= \sup_{z\in S^{d}(1)}\mathbb{E}\left[\left\|\nabla_{z}f(x+\gamma u,\xi) - \nabla_{z}f(x)\right\|\right] + \frac{d\Delta}{\gamma}$$

$$\stackrel{\text{(8)}}{\leq} \left(L_{0} + L_{1}\left\|\nabla f(x^{k})\right\|\right)\gamma\mathbb{E}\left[\left\|u\right\|\right] + \frac{d\Delta}{\gamma}$$

$$\leq \left(L_{0} + L_{1}M\right)\gamma + \frac{d\Delta}{\gamma},$$
(50)

where $u \in B^d(1)$, (1) = the equality is obtained from the fact, namely, distribution of e is symmetric, (2) = the inequality is obtain from bounded noise $|\delta(x)| \leq \Delta$, (3) = the equality is obtained from a version of Stokes' theorem [see Section 13.3.5, Exercise 14a, 67].

Bounding second moment (variance) of gradient approximation By definition gradient approximation (4) and Wirtinger-Poincare inequality (11) we have

$$\mathbb{E}\left[\left\|\mathbf{g}(x,\{\xi,e\}) - \mathbb{E}\left[\mathbf{g}(x,\{\xi,e\})\right]\right\|^{2}\right]$$

$$\leq \mathbb{E}\left[\left\|\mathbf{g}(x,\{\xi,e\})\right\|^{2}\right]$$

$$= \frac{d^{2}}{4\gamma^{2}}\mathbb{E}\left[\left\|\left(\tilde{f}(x+\gamma e,\xi) - \tilde{f}(x-\gamma e,\xi)\right)e\right\|^{2}\right]$$

$$= \frac{d^{2}}{4\gamma^{2}}\mathbb{E}\left[\left(f(x+\gamma e,\xi) - f(x-\gamma e,\xi) + \delta(x+\gamma e) - \delta(x-\gamma e)\right)^{2}\right]$$

$$\overset{(7)}{\leq} \frac{d^2}{2\gamma^2} \left(\mathbb{E} \left[\left(f(x + \gamma e, \xi) - f(x - \gamma e, \xi) \right)^2 \right] + 2\Delta^2 \right)$$

$$\overset{(11)}{\leq} \frac{d^2}{2\gamma^2} \left(\frac{\gamma^2}{d} \mathbb{E} \left[\left\| \nabla f(x + \gamma e, \xi) + \nabla f(x - \gamma e, \xi) \right\|^2 \right] + 2\Delta^2 \right)$$

$$= \frac{d^2}{2\gamma^2} \left(\frac{\gamma^2}{d} \mathbb{E} \left[\left\| \nabla f(x + \gamma e, \xi) + \nabla f(x - \gamma e, \xi) \pm 2\nabla f(x, \xi) \right\|^2 \right] + 2\Delta^2 \right)$$

$$\overset{(8)}{\leq} 4d\mathbb{E} \left[\left\| \nabla f(x, \xi) \right\|^2 \right] + 4dL^2\gamma^2 \mathbb{E} \left[\left\| e \right\|^2 \right] + \frac{d^2\Delta^2}{\gamma^2}$$

$$\overset{@}{\leq} 4d\tilde{\sigma}^2 + 4d \left(L_0 + L_1 \left\| \nabla f(x^k) \right\| \right)^2 \gamma^2 \mathbb{E} \left[\left\| e \right\|^2 \right] + \frac{d^2\Delta^2}{\gamma^2}$$

$$\le 4d\tilde{\sigma}^2 + 4d \left(L_0 + L_1 M \right)^2 \gamma^2 + \frac{d^2\Delta^2}{\gamma^2},$$

$$\tag{51}$$

D.2.1 Proof of Theorem 5.2

In order to obtain the convergence rate of ZO-ClipSGD in the convex setting, we need to substitute the obtained estimates (50) and (51) into the convergence rate of ClipSGD (49) instead of ζ and σ^2 , respectively. Given that $\frac{MR}{c^2} + \frac{R}{c} + \eta \lesssim \frac{MR}{c^2}$ at small c, then the convergence of ZO-ClipSGD in the convex setup is as follows:

$$\begin{split} \mathbb{E}\left[f(x^{N})\right] - f^{*} \lesssim \underbrace{\left(1 - \frac{\eta}{R}\right)^{K} \left(f(x^{0}) - f^{*}\right)}_{\tiny \textcircled{0}} + \underbrace{\frac{R^{2}}{\eta(N-K)}}_{\tiny \textcircled{0}} + \underbrace{\frac{dMR\tilde{\sigma}^{2}}{c^{2}B}}_{\tiny \textcircled{0}} + \underbrace{\frac{dMR\left(L_{0} + L_{1}M\right)^{2}\gamma^{2}}{c^{2}B}}_{\tiny \textcircled{0}} + \underbrace{\frac{d^{2}MR\Delta^{2}}{c^{2}B\gamma^{2}}}_{\tiny \textcircled{0}} + \underbrace{\frac{MR\left(L_{0} + L_{1}M\right)^{2}\gamma^{2}}{c^{2}}}_{\tiny \textcircled{0}} + \underbrace{\frac{d^{2}MR\Delta^{2}}{c^{2}\gamma^{2}}}_{\tiny \textcircled{0}} + \underbrace{\frac{(L_{0} + L_{1}M)\gamma R}{e}}_{\tiny \textcircled{0}} + \underbrace{\frac{d\Delta R}{\gamma}}_{\tiny \textcircled{0}}. \end{split}$$

From term ①, we find the *K*:

$$①: \quad \left(1 - \frac{\eta c}{R}\right)^{K} \left(f(x^{0}) - f^{*}\right) \le \varepsilon \quad \Rightarrow \quad K \ge \frac{R}{\eta c} \log \frac{f(x^{0}) - f^{*}}{\varepsilon}.
 \tag{52}$$

From term (2), we find the number of iterations N required for Algorithm 3 in convex setup to achieve ε -accuracy:

From terms ③, we find the batch size B:

$$\Im: \quad \frac{dMR\tilde{\sigma}^2}{c^2B} \le \varepsilon \quad \Rightarrow \quad B \ge \frac{dMR\tilde{\sigma}^2}{\varepsilon c^2}; \\ B = \mathcal{O}\left(\frac{dMR\tilde{\sigma}^2}{\varepsilon c^2}\right).$$
 (54)

From terms (4), (6) and (8) we find the smoothing parameter γ :

$$\circledast: \quad \frac{dMR\left(L_0+L_1M\right)^2\gamma^2}{c^2B} \leq \varepsilon \quad \Rightarrow \quad \gamma \leq \sqrt{\frac{\varepsilon c^2B}{dMR\left(L_0+L_1M\right)^2}} \stackrel{(54)}{=} \frac{\tilde{\sigma}}{\left(L_0+L_1M\right)};$$

$$\begin{aligned}
& (\texttt{S}: \quad \frac{MR \left(L_0 + L_1 M\right)^2 \gamma^2}{c^2} \leq \varepsilon \quad \Rightarrow \quad \gamma \leq \frac{\sqrt{\varepsilon}c}{\sqrt{MR} \left(L_0 + L_1 M\right)};\\ & (\texttt{S}: \quad \left(L_0 + L_1 M\right) R\gamma \leq \varepsilon \quad \Rightarrow \quad \gamma \leq \frac{\varepsilon}{R \left(L_0 + L_1 M\right)};\\ & \gamma \leq \frac{1}{\left(L_0 + L_1 M\right)} \min\left\{\tilde{\sigma}, \frac{\sqrt{\varepsilon}c}{\sqrt{MR}}, \frac{\varepsilon}{R}\right\} = \frac{\varepsilon}{R \left(L_0 + L_1 M\right)}.\end{aligned}$$

$$(55)$$

From the remaining terms (5), (7) and (9), we find the maximum allowable level of adversarial noise Δ that still guarantees the convergence of the ZO-ClipSGD to desired accuracy ε in convex setup:

$$\begin{split} & (\texttt{S}: \quad \frac{d^2 M R \Delta^2}{c^2 B \gamma^2} \leq \varepsilon \quad \Rightarrow \quad \Delta \leq \frac{\sqrt{\varepsilon} c \gamma \sqrt{B}}{d \sqrt{MR}} \stackrel{(\texttt{54}),(\texttt{55})}{=} \frac{\varepsilon \tilde{\sigma}}{\sqrt{d} (L_0 + L_1 M) R}; \\ & (\texttt{T}: \quad \frac{d^2 M R \Delta^2}{\gamma^2 c^2} \leq \varepsilon \quad \Rightarrow \quad \Delta \leq \sqrt{\frac{\gamma^2 c^2 \varepsilon}{d^2 M R}} \stackrel{(\texttt{55})}{=} \frac{\varepsilon^{3/2} c}{d (L_0 + L_1 M) \sqrt{M} R^{3/2}}; \\ & (\texttt{S}: \quad \frac{d \Delta R}{\gamma} \leq \varepsilon \quad \Rightarrow \quad \Delta \leq \sqrt{\frac{\gamma \varepsilon}{d R}} \stackrel{(\texttt{55})}{=} \frac{\varepsilon^2}{d (L_0 + L_1 M) R^2}; \\ & \Delta \leq \frac{\varepsilon}{\sqrt{d} (L_0 + L_1 M) R} \min \left\{ \tilde{\sigma}, \frac{\sqrt{\varepsilon} c}{\sqrt{d} \sqrt{MR}}, \frac{\varepsilon}{\sqrt{d} R} \right\} \\ & = \frac{\varepsilon}{\sqrt{d} (L_0 + L_1 M) R} \min \left\{ \tilde{\sigma}, \frac{\varepsilon}{\sqrt{d} R} \right\}. \end{split}$$

In this way, the ZO-ClipSGD achieves ε -accuracy: $\mathbb{E}\left[f(x^N) - f^*\right] \le \varepsilon$ in convex setup after

$$N \stackrel{(53)}{=} \mathcal{O}\left(\frac{R^2}{\eta\varepsilon} + \frac{R}{\eta c}\log\frac{1}{\varepsilon}\right), \quad T = N \cdot B \stackrel{(53),(54)}{=} \mathcal{O}\left(\frac{d\tilde{\sigma}^2 M R^2}{\varepsilon c^2 \eta} \left(\frac{1}{c}\log\frac{1}{\varepsilon} + \frac{R}{\varepsilon}\right)\right)$$

number of iterations, total number of zero-order oracle calls and at

$$\Delta \stackrel{(56)}{\lesssim} \frac{\varepsilon}{\sqrt{d} \left(L_0 + L_1 M \right) R} \min \left\{ \tilde{\sigma}, \frac{\varepsilon}{\sqrt{d}R} \right\}$$

the maximum level of noise with smoothing parameter $\frac{\varepsilon}{(L_0+L_1M)R}$ (55).

E Zero-Order Normalized Stochastic Gradient Descent Method

This section consists of two parts: 1) a generalization of the convergence result of NSGD (Algorithm 2) to the biased gradient oracle $\mathbf{g}(x^k, \boldsymbol{\xi}^k) = \nabla f(x^k, \boldsymbol{\xi}^k) + \mathbf{b}(x^k)$, where $\mathbf{b}(x^k)$ is biased bounded by $\zeta \geq 0$: $\|\mathbf{b}(x^k)\| \leq \zeta$; 2) deriving convergence estimates of ZO-NSGD directly.

E.1 Biased Normalized Stochastic Gradient Descent Method (Proof of the Lemma 5.3)

Let's introduce the notation $G(x^k, \boldsymbol{\xi}^k) = \frac{\mathbf{g}(x^k, \boldsymbol{\xi}^k)}{\|\mathbf{g}(x^k, \boldsymbol{\xi}^k)\|}$, then using (L_0, L_1) -smoothness (see Assumption 1.2):

$$f(x^{k+1}) - f(x^k) \stackrel{(8)}{\leq} \left\langle \nabla f(x^k), x^{k+1} - x^k \right\rangle + \frac{L_0 + L_1 \left\| \nabla f(x^k) \right\|}{2} \left\| x^{k+1} - x^k \right\|^2$$
$$= -\eta \left\langle \nabla f(x^k), G(x^k, \boldsymbol{\xi}^k) \right\rangle + \frac{\eta^2 (L_0 + L_1 \left\| \nabla f(x^k) \right\|)}{2} \left\| G(x^k, \boldsymbol{\xi}^k) \right\|^2.$$
(57)

Next, we consider 4 cases of the relation $\|\nabla f(x^k)\|$ and $\|\mathbf{g}(x^k, \boldsymbol{\xi}^k)\|$ with respect to the hyperparameter λ .

E.1.1 First case: $\left\| \nabla f(x^k) \right\| \ge \lambda$ and $\left\| \mathbf{g}(x^k, \boldsymbol{\xi}^k) \right\| \ge \lambda$

Let us evaluate first summand of (57) with $\alpha = \left\| \nabla f(x^k) \right\|^{-1}$:

$$\begin{split} -\eta \left\langle \nabla f(x^{k}), G(x^{k}, \boldsymbol{\xi}^{k}) \right\rangle &\stackrel{(5)}{=} -\frac{\alpha \eta}{2} \left\| \nabla f(x^{k}) \right\|^{2} - \frac{\eta}{2\alpha} \left\| G(x^{k}, \boldsymbol{\xi}^{k}) \right\|^{2} \\ &+ \frac{\eta}{2\alpha} \left\| G(x^{k}, \boldsymbol{\xi}^{k}) - \alpha \nabla f(x^{k}) \right\|^{2} \\ &= -\frac{\eta}{2} \left\| \nabla f(x^{k}) \right\| - \frac{\eta}{2\alpha} \left\| G(x^{k}, \boldsymbol{\xi}^{k}) \right\|^{2} \\ &+ \frac{\eta}{2\lambda^{2}\alpha} \left\| \lambda G(x^{k}, \boldsymbol{\xi}^{k}) - \lambda \alpha \nabla f(x^{k}) \right\|^{2} \\ &= -\frac{\eta}{2} \left\| \nabla f(x^{k}) \right\| - \frac{\eta}{2\alpha} \left\| G(x^{k}, \boldsymbol{\xi}^{k}) \right\|^{2} \\ &+ \frac{\eta}{2\lambda^{2}\alpha} \left\| \operatorname{clip}_{\lambda} \left(\mathbf{g}(x^{k}, \boldsymbol{\xi}^{k}) \right) - \operatorname{clip}_{\lambda} \left(\nabla f(x^{k}) \right) \right\|^{2} \end{split}$$

Using that clipping is a projection on onto a convex set, namely ball with radius λ , and thus is Lipshitz operator with Lipshitz constant 1, we can obtain:

$$-\eta \left\langle \nabla f(x^{k}), \mathbb{E}\left[G(x^{k}, \boldsymbol{\xi}^{k})\right]\right\rangle \leq -\frac{\eta}{2} \left\|\nabla f(x^{k})\right\| - \frac{\eta}{2\alpha} \mathbb{E}\left[\left\|G(x^{k}, \boldsymbol{\xi}^{k})\right\|^{2}\right] + \frac{\eta}{2\lambda^{2}\alpha} \mathbb{E}\left[\left\|\mathbf{g}(x^{k}, \boldsymbol{\xi}^{k}) - \nabla f(x^{k})\right\|^{2}\right].$$
(58)

In the case:
$$0 \leq \zeta \leq \frac{\lambda}{\sqrt{2}}$$
. Using this in (58), we have the following with $\eta_k \leq \frac{\|\nabla f(x^k)\|}{2(L_0+L_1\|\nabla f(x^k)\|)}$:

$$\mathbb{E}\left[f(x^{k+1})\right] - f(x^k) \stackrel{(57)}{\leq} -\eta \langle \nabla f(x^k), \mathbb{E}\left[G(x^k, \boldsymbol{\xi}^k)\right] \rangle + \frac{\eta^2(L_0 + L_1 \|\nabla f(x^k)\|)}{2} \mathbb{E}\left[\left\|G(x^k, \boldsymbol{\xi}^k) - \nabla f(x^k)\right\|^2\right]$$

$$\stackrel{(58)}{\leq} -\frac{\eta}{2} \|\nabla f(x^k)\| - \frac{\eta}{2\alpha} \mathbb{E}\left[\left\|G(x^k, \boldsymbol{\xi}^k)\right\|^2\right] + \frac{\eta}{2\lambda^2\alpha} \mathbb{E}\left[\left\|\mathbf{g}(x^k, \boldsymbol{\xi}^k) - \nabla f(x^k)\right\|^2\right]$$

$$+ \frac{\eta^2(L_0 + L_1 \|\nabla f(x^k)\|)}{2} \mathbb{E}\left[\left\|\mathbf{g}(x^k, \boldsymbol{\xi}^k) - \nabla f(x^k)\right\|^2\right]$$

$$= -\frac{\eta}{2} \|\nabla f(x^k)\| + \frac{\eta}{2\lambda^2\alpha} \mathbb{E}\left[\left\|\mathbf{g}(x^k, \boldsymbol{\xi}^k) - \nabla f(x^k)\right\|^2\right]$$

$$= -\frac{\eta}{2} \mathbb{E}\left[\left\|G(x^k, \boldsymbol{\xi}^k)\right\|^2\right] \left(1 - \frac{\eta(L_0 + L_1 \|\nabla f(x^k)\|)}{\|\nabla f(x^k)\|}\right)$$

$$\leq -\frac{\eta}{2} \|\nabla f(x^k)\| + \frac{\eta}{2\lambda^2\alpha} \mathbb{E}\left[\left\|\mathbf{g}(x^k, \boldsymbol{\xi}^k) - \nabla f(x^k)\right\|^2\right]$$

$$\stackrel{(9)}{=} -\frac{\eta}{2} \|\nabla f(x^k)\| + \frac{\eta}{2\lambda^2\alpha} \mathbb{E}\left[\left\|\mathbf{g}(x^k, \boldsymbol{\xi}^k) - \mathbb{E}\left[\mathbf{g}(x^k, \boldsymbol{\xi}^k)\right]\right\|\right] + \frac{\eta}{2\lambda^2\alpha} \left\|\mathbf{b}(x^k)\right\|^2$$

$$\leq -\frac{\eta}{2} \|\nabla f(x^k)\| + \frac{\eta\sigma^2}{2\lambda^2\alpha B} + \frac{\eta\zeta^2}{2\lambda^2\alpha}$$

$$\leq -\frac{\eta}{2} \|\nabla f(x^k)\| + \frac{\eta}{4} \|\nabla f(x^k)\| + \frac{\eta\sigma^2 M}{2\lambda^2 B}$$

$$= -\frac{\eta}{4} \|\nabla f(x^k)\| + \frac{\eta\sigma^2 M}{2\lambda^2 B}.$$
(59)

The step size will be constant, depending on the hyperparameter λ :

$$\frac{\left\|\nabla f(x^k)\right\|}{2\left(L_0+L_1\left\|\nabla f(x^k)\right\|\right)} = \frac{1}{2\left(L_0\frac{1}{\left\|\nabla f(x^k)\right\|} + L_1\right)} = \frac{\lambda}{2\left(L_0\frac{\lambda}{\left\|\nabla f(x^k)\right\|} + L_1\lambda\right)} \ge \frac{\lambda}{2\left(L_0+L_1\lambda\right)}.$$

Thus, $\eta_k = \eta \le \frac{\lambda}{2\left(L_0+L_1\lambda\right)}.$

Using the convexity assumption of the function, we have the following:

$$f(x^{k}) - f^{*} \leq \left\langle \nabla f(x^{k}), x^{k} - x^{*} \right\rangle \stackrel{(6)}{\leq} \left\| \nabla f(x^{k}) \right\| \left\| x^{k} - x^{*} \right\| \leq \left\| \nabla f(x^{k}) \right\| \underbrace{\left\| x^{0} - x^{*} \right\|}_{R}.$$

Hence we have:

$$\left\|\nabla f(x^k)\right\| \ge \frac{f(x^k) - f^*}{R}.$$
(60)

Then substituting (60) into (59) we obtain:

$$\mathbb{E}\left[f(x^{k+1})\right] - f(x^k) \le -\frac{\eta}{4} \left\|\nabla f(x^k)\right\| + \frac{\eta\sigma^2 M}{2\lambda^2 B} \le -\frac{\eta}{4R}(f(x^k) - f^*) + \frac{\eta\sigma^2 M}{2\lambda^2 B}.$$

This inequality is equivalent to the trailing inequality:

$$\mathbb{E}\left[f(x^{k+1})\right] - f^* \le \left(1 - \frac{\eta}{4R}\right) \left(f(x^k) - f^*\right) + \frac{\eta \sigma^2 M}{2\lambda^2 B}.$$

Then for k = 0, 1, 2, ..., N - 1 iterations that satisfy the conditions $\|\mathbf{g}(x^k, \boldsymbol{\xi}^k)\| \ge \sqrt{2}\zeta$ and $\|\nabla f(x^k)\| \ge \sqrt{2}\zeta$ NSGD with biased gradient oracle shows linear convergence:

$$\mathbb{E}\left[f(x^N)\right] - f^* \le \left(1 - \frac{\eta}{4R}\right)^N \left(f(x^0) - f^*\right) + \frac{2\sigma^2 MR}{\lambda^2 B}.$$

In the case: $\frac{\lambda}{\sqrt{2}} \leq \zeta$. Using this in (58), we have the following with $\eta_k \leq \frac{\|\nabla f(x^k)\|}{2(L_0 + L_1 \|\nabla f(x^k)\|)}$:

$$\mathbb{E}\left[f(x^{k+1})\right] - f(x^{k}) \stackrel{(57)}{\leq} -\eta \left\langle \nabla f(x^{k}), \mathbb{E}\left[G(x^{k}, \boldsymbol{\xi}^{k})\right] \right\rangle + \frac{\eta^{2}(L_{0} + L_{1} \left\| \nabla f(x^{k}) \right\| \right)}{2} \mathbb{E}\left[\left\| G(x^{k}, \boldsymbol{\xi}^{k}) \right\|^{2}\right] \\
\stackrel{(58)}{\leq} -\frac{\eta}{2} \left\| \nabla f(x^{k}) \right\| - \frac{\eta}{2\alpha} \mathbb{E}\left[\left\| G(x^{k}, \boldsymbol{\xi}^{k}) \right\|^{2}\right] + \frac{\eta}{2\lambda^{2}\alpha} \mathbb{E}\left[\left\| g(x^{k}, \boldsymbol{\xi}^{k}) - \nabla f(x^{k}) \right\|^{2}\right] \\
\quad + \frac{\eta^{2}(L_{0} + L_{1} \left\| \nabla f(x^{k}) \right\| \right)}{2} \mathbb{E}\left[\left\| g(x^{k}, \boldsymbol{\xi}^{k}) - \nabla f(x^{k}) \right\|^{2}\right] \\
= -\frac{\eta}{2} \left\| \nabla f(x^{k}) \right\| + \frac{\eta}{2\lambda^{2}\alpha} \mathbb{E}\left[\left\| g(x^{k}, \boldsymbol{\xi}^{k}) - \nabla f(x^{k}) \right\|^{2}\right] \\
\quad - \frac{\eta}{2} \mathbb{E}\left[\left\| G(x^{k}, \boldsymbol{\xi}^{k}) \right\|^{2}\right] \left(1 - \frac{\eta(L_{0} + L_{1} \left\| \nabla f(x^{k}) \right\| \right)}{\left\| \nabla f(x^{k}) \right\|} \right) \\
\leq -\frac{\eta}{2} \left\| \nabla f(x^{k}) \right\| + \frac{\eta}{2\lambda^{2}\alpha} \mathbb{E}\left[\left\| g(x^{k}, \boldsymbol{\xi}^{k}) - \nabla f(x^{k}) \right\|^{2}\right] \\
\stackrel{(9)}{=} -\frac{\eta}{2} \left\| \nabla f(x^{k}) \right\| + \frac{\eta}{2\lambda^{2}\alpha} \mathbb{E}\left[\left\| g(x^{k}, \boldsymbol{\xi}^{k}) - \mathbb{E}\left[g(x^{k}, \boldsymbol{\xi}^{k}) \right] \right\|\right] + \frac{\eta}{2\lambda^{2}\alpha} \left\| \mathbf{b}(x^{k}) \right\|^{2} \\
\leq -\frac{\eta}{2} \left\| \nabla f(x^{k}) \right\| + \frac{\eta\sigma^{2}}{2\lambda^{2}\alpha B} + \frac{\eta\zeta^{2}}{2\lambda^{2}\alpha} \\
\leq -\frac{\eta}{2} \left\| \nabla f(x^{k}) \right\| + \frac{\eta\sigma^{2}M}{2\lambda^{2}B} + \frac{\eta\zeta^{2}M}{2\lambda^{2}} \\
= -\frac{\eta}{2} \left\| \nabla f(x^{k}) \right\| + \frac{\eta\sigma^{2}M}{2\lambda^{2}B} + \frac{\eta\zeta^{2}M}{2\lambda^{2}}.$$
(61)

The step size will be constant, depending on the hyperparameter λ :

$$\begin{split} \frac{\left\|\nabla f(x^k)\right\|}{2\left(L_0+L_1\left\|\nabla f(x^k)\right\|\right)} &= \frac{1}{2\left(L_0\frac{1}{\left\|\nabla f(x^k)\right\|}+L_1\right)} = \frac{\lambda}{2\left(L_0\frac{\lambda}{\left\|\nabla f(x^k)\right\|}+L_1\lambda\right)} \geq \frac{\lambda}{2\left(L_0+L_1\lambda\right)} \end{split}$$
 Thus, $\eta_k = \eta \leq \frac{\lambda}{2\left(L_0+L_1\lambda\right)}.$

Using the convexity assumption of the function, we have the following:

$$f(x^{k}) - f^{*} \leq \left\langle \nabla f(x^{k}), x^{k} - x^{*} \right\rangle \stackrel{(6)}{\leq} \left\| \nabla f(x^{k}) \right\| \left\| x^{k} - x^{*} \right\| \leq \left\| \nabla f(x^{k}) \right\| \underbrace{\left\| x^{0} - x^{*} \right\|}_{R}.$$

Hence we have:

$$\left\|\nabla f(x^k)\right\| \ge \frac{f(x^k) - f^*}{R}.$$
(62)

Then substituting (62) into (61) we obtain:

$$\mathbb{E}\left[f(x^{k+1})\right] - f(x^k) \le -\frac{\eta}{2} \left\|\nabla f(x^k)\right\| + \frac{\eta\sigma^2 M}{2\lambda^2 B} + \frac{\eta\zeta^2 M}{2\lambda^2} \le -\frac{\eta}{2R}(f(x^k) - f^*) + \frac{\eta\sigma^2 M}{2\lambda^2 B} + \frac{\eta\zeta^2 M}{2\lambda^2}.$$

This inequality is equivalent to the trailing inequality:

$$\mathbb{E}\left[f(x^{k+1})\right] - f^* \le \left(1 - \frac{\eta}{2R}\right) \left(f(x^k) - f^*\right) + \frac{\eta \sigma^2 M}{2\lambda^2 B} + \frac{\eta \zeta^2 M}{2\lambda^2}.$$

Then for k = 0, 1, 2, ..., N - 1 iterations that satisfy the conditions $\|\mathbf{g}(x^k, \boldsymbol{\xi}^k)\| \geq \lambda$ and $\|\nabla f(x^k)\| \geq \lambda$ and $\zeta \geq \sqrt{2}\lambda$ NSGD with biased gradient oracle shows linear convergence:

$$\mathbb{E}\left[f(x^N)\right] - f^* \le \left(1 - \frac{\eta}{2R}\right)^N \left(f(x^0) - f^*\right) + \frac{\sigma^2 MR}{\lambda^2 B} + \frac{\zeta^2 MR}{\lambda^2}.$$

E.1.2 Second case: $\left\| \nabla f(x^k) \right\| \leq \lambda$ and $\left\| \mathbf{g}(x^k, \boldsymbol{\xi}^k) \right\| \geq \lambda$

Let us evaluate first summand of (57) with $\alpha = \lambda^{-1}$:

$$\begin{aligned} -\eta \left\langle \nabla f(x^{k}), G(x^{k}, \boldsymbol{\xi}^{k}) \right\rangle &\stackrel{(5)}{=} -\frac{\alpha \eta}{2} \left\| \nabla f(x^{k}) \right\|^{2} - \frac{\eta}{2\alpha} \left\| G(x^{k}, \boldsymbol{\xi}^{k}) \right\|^{2} \\ &+ \frac{\eta}{2\alpha} \left\| G(x^{k}, \boldsymbol{\xi}^{k}) - \alpha \nabla f(x^{k}) \right\|^{2} \\ &\leq -\frac{\eta}{2} \left\| \nabla f(x^{k}) \right\| - \frac{\eta}{2\alpha} \left\| G(x^{k}, \boldsymbol{\xi}^{k}) \right\|^{2} \\ &+ \frac{\eta}{2\lambda} \left\| \lambda G(x^{k}, \boldsymbol{\xi}^{k}) - \nabla f(x^{k}) \right\|^{2} \\ &= -\frac{\eta}{2} \left\| \nabla f(x^{k}) \right\| - \frac{\eta}{2\alpha} \left\| G(x^{k}, \boldsymbol{\xi}^{k}) \right\|^{2} \\ &+ \frac{\eta}{2\lambda} \left\| \operatorname{clip}_{\lambda} \left(\mathbf{g}(x^{k}, \boldsymbol{\xi}^{k}) \right) - \operatorname{clip}_{\lambda} \left(\nabla f(x^{k}) \right) \right\|^{2} \end{aligned}$$

Using that clipping is a projection on onto a convex set, namely ball with radius λ , and thus is Lipshitz operator with Lipshitz constant 1, we can obtain:

$$-\eta \left\langle \nabla f(x^{k}), \mathbb{E}\left[G(x^{k}, \boldsymbol{\xi}^{k})\right]\right\rangle \leq -\frac{\eta}{2} \left\|\nabla f(x^{k})\right\| - \frac{\eta}{2\alpha} \mathbb{E}\left[\left\|G(x^{k}, \boldsymbol{\xi}^{k})\right\|^{2}\right] + \frac{\eta}{2\lambda} \mathbb{E}\left[\left\|\mathbf{g}(x^{k}, \boldsymbol{\xi}^{k}) - \nabla f(x^{k})\right\|^{2}\right].$$
(63)

Using this, we have the following with $\eta_k \leq \frac{\|\nabla f(x^k)\|}{2(L_0+L_1\|\nabla f(x^k)\|)}$:

$$\mathbb{E}\left[f(x^{k+1})\right] - f(x^k) \stackrel{(57)}{\leq} -\eta \left\langle \nabla f(x^k), \mathbb{E}\left[G(x^k, \boldsymbol{\xi}^k)\right] \right\rangle + \frac{\eta^2 (L_0 + L_1 \left\| \nabla f(x^k) \right\|)}{2} \mathbb{E}\left[\left\| G(x^k, \boldsymbol{\xi}^k) \right\|^2\right] \\ \stackrel{(63)}{\leq} -\frac{\eta}{2} \left\| \nabla f(x^k) \right\| - \frac{\eta}{2\alpha} \mathbb{E}\left[\left\| G(x^k, \boldsymbol{\xi}^k) \right\|^2\right] + \frac{\eta}{2\lambda} \mathbb{E}\left[\left\| \mathbf{g}(x^k, \boldsymbol{\xi}^k) - \nabla f(x^k) \right\|^2\right] \\ + \frac{\eta^2 (L_0 + L_1 \left\| \nabla f(x^k) \right\|)}{2} \mathbb{E}\left[\left\| G(x^k, \boldsymbol{\xi}^k) \right\|^2\right] \\ \stackrel{(9)}{=} -\frac{\eta}{2} \left\| \nabla f(x^k) \right\| + \frac{\eta}{2\lambda} \mathbb{E}\left[\left\| \mathbf{g}(x^k, \boldsymbol{\xi}^k) - \mathbb{E}\left[\mathbf{g}(x^k, \boldsymbol{\xi}^k)\right]\right\|^2\right] + \frac{\eta}{2\lambda} \left\| \mathbf{b}(x^k) \right\|^2$$

$$-\frac{\eta}{2}\mathbb{E}\left[\left\|G(x^{k},\boldsymbol{\xi}^{k})\right\|^{2}\right]\left(1-\frac{\eta(L_{0}+L_{1}\left\|\nabla f(x^{k})\right\|)}{\left\|\nabla f(x^{k})\right\|}\right)$$
$$\leq -\frac{\eta}{2}\left\|\nabla f(x^{k})\right\|+\frac{\eta\sigma^{2}}{2\lambda B}+\frac{\eta\zeta^{2}}{2\lambda}.$$
(64)

$$\frac{\left\|\nabla f(x^{k})\right\|}{2\left(L_{0}+L_{1}\left\|\nabla f(x^{k})\right\|\right)} = \frac{1}{2\left(L_{0}\frac{1}{\left\|\nabla f(x^{k})\right\|}+L_{1}\right)} = \frac{\lambda}{2\left(L_{0}\frac{\lambda}{\left\|\nabla f(x^{k})\right\|}+L_{1}\lambda\right)} \ge \frac{\lambda}{2\left(L_{0}+L_{1}\lambda\right)}.$$
Thus, $\eta_{k} = \eta < \frac{\lambda}{2\left(L_{0}+L_{1}\lambda\right)}.$

Thus, $\eta_k = \eta \leq \frac{\lambda}{2(L_0 + L_1\lambda)}$

Using the convexity assumption of the function, we have the following:

$$f(x^{k}) - f^{*} \leq \left\langle \nabla f(x^{k}), x^{k} - x^{*} \right\rangle \stackrel{(6)}{\leq} \left\| \nabla f(x^{k}) \right\| \left\| x^{k} - x^{*} \right\| \leq \left\| \nabla f(x^{k}) \right\| \underbrace{\left\| x^{0} - x^{*} \right\|}_{R}.$$

Hence we have:

$$\left\|\nabla f(x^k)\right\| \ge \frac{f(x^k) - f^*}{R}.$$
(65)

Then substituting (65) into (64) we obtain:

$$\mathbb{E}\left[f(x^{k+1})\right] - f(x^k) \le -\frac{\eta}{2} \left\|\nabla f(x^k)\right\| + \frac{\eta\sigma^2}{2\lambda B} + \frac{\eta\zeta^2}{2\lambda} \le -\frac{\eta}{2R}(f(x^k) - f^*) + \frac{\eta\sigma^2}{2\lambda B} + \frac{\eta\zeta^2}{2\lambda}.$$

This inequality is equivalent to the trailing inequality:

$$\mathbb{E}\left[f(x^{k+1})\right] - f^* \le \left(1 - \frac{\eta}{2R}\right) \left(f(x^k) - f^*\right) + \frac{\eta}{2\lambda} \left(\frac{\sigma^2}{B} + \zeta^2\right).$$

Then for k = 0, 1, 2, ..., N - 1 iterations that satisfy the conditions $\|\nabla f(x^k)\| \leq \lambda$ and $\|\mathbf{g}(x^k, \boldsymbol{\xi}^k)\| \geq \lambda$ NSGD with biased gradient oracle shows linear convergence:

$$\mathbb{E}\left[f(x^N)\right] - f^* \le \left(1 - \frac{\eta}{2R}\right)^N \left(f(x^0) - f^*\right) + \frac{R}{\lambda} \left(\frac{\sigma^2}{B} + \zeta^2\right).$$

E.1.3 Third case: $\|\nabla f(x^k)\| \leq \lambda$ and $\|\mathbf{g}(x^k, \boldsymbol{\xi}^k)\| \leq \lambda$

Using this in (57), we have the following with $\eta_k \leq \frac{\|\nabla f(x^k)\|}{2(L_0+L_1\|\nabla f(x^k)\|)}$ and $\alpha = \|\nabla f(x^k)\|^{-1}$:

$$\mathbb{E}\left[f(x^{k+1})\right] - f(x^k) \stackrel{(57)}{\leq} -\eta \left\langle \nabla f(x^k), \mathbb{E}\left[G(x^k, \boldsymbol{\xi}^k)\right] \right\rangle + \frac{\eta^2 (L_0 + L_1 \left\| \nabla f(x^k) \right\|)}{2} \mathbb{E}\left[\left\|G(x^k, \boldsymbol{\xi}^k)\right\|^2\right] \\ \stackrel{(5)}{=} -\frac{\eta \alpha}{2} \left\| \nabla f(x^k) \right\|^2 - \frac{\eta}{2\alpha} \mathbb{E}\left[\left\|G(x^k, \boldsymbol{\xi}^k)\right\|^2\right] + \frac{\eta}{2\alpha} \mathbb{E}\left[\left\|G(x^k, \boldsymbol{\xi}^k) - \alpha \nabla f(x^k)\right\|^2\right] \\ + \frac{\eta^2 (L_0 + L_1 \left\| \nabla f(x^k) \right\|)}{2} \mathbb{E}\left[\left\|G(x^k, \boldsymbol{\xi}^k) - \alpha \nabla f(x^k)\right\|^2\right] \\ = -\frac{\eta}{2} \left\| \nabla f(x^k) \right\| + \frac{\eta}{2\alpha} \mathbb{E}\left[\left\|G(x^k, \boldsymbol{\xi}^k) - \alpha \nabla f(x^k)\right\|^2\right] \\ - \frac{\eta}{2} \mathbb{E}\left[\left\|G(x^k, \boldsymbol{\xi}^k)\right\|^2\right] \left(1 - \frac{\eta (L_0 + L_1 \left\| \nabla f(x^k) \right\|)}{\left\| \nabla f(x^k) \right\|} \right) \\ \leq -\frac{\eta}{2} \left\| \nabla f(x^k) \right\| + \frac{\eta}{2\alpha} \mathbb{E}\left[\left\|G(x^k, \boldsymbol{\xi}^k) - \alpha \nabla f(x^k)\right\|^2\right] \\ \leq -\frac{\eta}{2} \left\| \nabla f(x^k) \right\| + \frac{\eta}{\alpha} \mathbb{E}\left[\left\|G(x^k, \boldsymbol{\xi}^k) - \alpha \nabla f(x^k)\right\|^2\right]$$

$$= -\frac{\eta}{2} \left\| \nabla f(x^{k}) \right\| + \frac{\eta}{\alpha} \mathbb{E} \left[\left\| \frac{\mathbf{g}(x^{k}, \boldsymbol{\xi}^{k})}{\|\mathbf{g}(x^{k}, \boldsymbol{\xi}^{k})\|} \right\|^{2} + \left\| \frac{\nabla f(x^{k})}{\|\nabla f(x^{k})\|} \right\|^{2} \right]$$
$$= -\frac{\eta}{2} \left\| \nabla f(x^{k}) \right\| + \frac{2\eta\lambda \left\| \nabla f(x^{k}) \right\|}{\lambda}$$
$$\leq -\frac{\eta}{2} \left\| \nabla f(x^{k}) \right\| + 2\eta\lambda.$$
(66)

$$\frac{\left\|\nabla f(x^{k})\right\|}{2\left(L_{0}+L_{1}\left\|\nabla f(x^{k})\right\|\right)} = \frac{1}{2\left(L_{0}\frac{1}{\left\|\nabla f(x^{k})\right\|}+L_{1}\right)} = \frac{\lambda}{2\left(L_{0}\frac{\lambda}{\left\|\nabla f(x^{k})\right\|}+L_{1}\lambda\right)} \ge \frac{\lambda}{2\left(L_{0}+L_{1}\lambda\right)}.$$

Thus, $\eta_k = \eta \leq \frac{\lambda}{2(L_0 + L_1\lambda)}$.

Using the convexity assumption of the function, we have the following:

$$f(x^{k}) - f^{*} \leq \left\langle \nabla f(x^{k}), x^{k} - x^{*} \right\rangle^{(6)} \leq \left\| \nabla f(x^{k}) \right\| \left\| x^{k} - x^{*} \right\| \leq \left\| \nabla f(x^{k}) \right\| \underbrace{\left\| x^{0} - x^{*} \right\|}_{R}$$

Hence we have:

$$\left\|\nabla f(x^k)\right\| \ge \frac{f(x^k) - f^*}{R}.$$
(67)

Then substituting (67) into (66) we obtain:

$$\mathbb{E}\left[f(x^{k+1})\right] - f(x^k) \le -\frac{\eta}{2} \left\|\nabla f(x^k)\right\| + 2\eta\lambda \le -\frac{\eta}{2R}(f(x^k) - f^*) + 2\eta\lambda.$$

This inequality is equivalent to the trailing inequality:

$$\mathbb{E}\left[f(x^{k+1})\right] - f^* \le \left(1 - \frac{\eta}{2R}\right)\left(f(x^k) - f^*\right) + 2\eta\lambda.$$

Then for k = 0, 1, 2, ..., N - 1 iterations that satisfy the conditions $\|\nabla f(x^k)\| \le \lambda$ NSGD with biased gradient oracle shows linear convergence:

$$\mathbb{E}\left[f(x^N)\right] - f^* \le \left(1 - \frac{\eta}{2R}\right)^N \left(f(x^0) - f^*\right) + \lambda R.$$

E.1.4 Fourth case: $\left\| \nabla f(x^k) \right\| \ge \lambda$ and $\left\| \mathbf{g}(x^k, \boldsymbol{\xi}^k) \right\| \le \lambda$

Using this in (57), we have the following with $\eta_k \leq \frac{\|\nabla f(x^k)\|}{2(L_0+L_1\|\nabla f(x^k)\|)}$ and $\alpha = \lambda^{-1}$:

$$\mathbb{E}\left[f(x^{k+1})\right] - f(x^k) \stackrel{(5)}{\leq} -\eta \left\langle \nabla f(x^k), \mathbb{E}\left[G(x^k, \boldsymbol{\xi}^k)\right] \right\rangle \\ + \frac{\eta^2(L_0 + L_1 \left\| \nabla f(x^k) \right\|}{2} \mathbb{E}\left[\left\| G(x^k, \boldsymbol{\xi}^k) \right\|^2\right] \\ \stackrel{(5)}{\equiv} -\frac{\eta \alpha}{2} \left\| \nabla f(x^k) \right\|^2 - \frac{\eta}{2\alpha} \left\| \mathbb{E}\left[G(x^k, \boldsymbol{\xi}^k)\right] \right\|^2 \\ + \frac{\eta}{2\alpha} \left\| \mathbb{E}\left[G(x^k, \boldsymbol{\xi}^k)\right] - \alpha \nabla f(x^k) \right\|^2 \\ + \frac{\eta^2(L_0 + L_1 \left\| \nabla f(x^k) \right\|}{2} \mathbb{E}\left[\left\| G(x^k, \boldsymbol{\xi}^k) \right\|^2\right] \\ = -\frac{\eta}{2\lambda} \left\| \nabla f(x^k) \right\|^2 + \frac{\eta}{2\lambda} \left\| \mathbb{E}\left[\lambda G(x^k, \boldsymbol{\xi}^k)\right] - \nabla f(x^k) \right\|^2 \\ + \frac{\eta^2(L_0 + L_1 \left\| \nabla f(x^k) \right\|}{2} \right]$$

$$= -\frac{\eta}{2\lambda} \|\nabla f(x^{k})\|^{2} + \frac{\eta}{\lambda} \|\mathbb{E} \left[\frac{\lambda \mathbf{g}(x^{k}, \boldsymbol{\xi}^{k})}{\|\mathbf{g}(x^{k}, \boldsymbol{\xi}^{k})\|} - \mathbf{g}(x^{k}, \boldsymbol{\xi}^{k}) \right] \|^{2} \\ + \frac{\eta}{\lambda} \|\mathbf{b}(x^{k})\|^{2} + \frac{\eta^{2}(L_{0} + L_{1} \|\nabla f(x^{k})\|)}{2} \\ = -\frac{\eta}{2\lambda} \|\nabla f(x^{k})\|^{2} + \frac{\eta}{2\lambda} \|\mathbb{E} \left[\left(\frac{\lambda}{\|\mathbf{g}(x^{k}, \boldsymbol{\xi}^{k})\|} - 1 \right) \mathbf{g}(x^{k}, \boldsymbol{\xi}^{k}) \right] \|^{2} \\ + \frac{\eta}{\lambda} \|\mathbf{b}(x^{k})\|^{2} + \frac{\eta^{2}(L_{0} + L_{1} \|\nabla f(x^{k})\|)}{2} \\ \leq -\frac{\eta}{2\lambda} \|\nabla f(x^{k})\|^{2} + \frac{\eta}{2\lambda} \mathbb{E} \left[\left(\frac{\lambda}{\|\mathbf{g}(x^{k}, \boldsymbol{\xi}^{k})\|} - 1 \right)^{2} \|\mathbf{g}(x^{k}, \boldsymbol{\xi}^{k})\|^{2} \right] \\ + \frac{\eta}{\lambda} \|\mathbf{b}(x^{k})\|^{2} + \frac{\eta^{2}(L_{0} + L_{1} \|\nabla f(x^{k})\|)}{2} \\ \leq -\frac{\eta}{2\lambda} \|\nabla f(x^{k})\|^{2} + \frac{\eta}{2\lambda} \mathbb{E} \left[\frac{\lambda^{2}}{\|\mathbf{g}(x^{k}, \boldsymbol{\xi}^{k})\|^{2}} \|\mathbf{g}(x^{k}, \boldsymbol{\xi}^{k})\|^{2} \right] \\ + \frac{\eta}{\lambda} \|\mathbf{b}(x^{k})\|^{2} + \frac{\eta^{2}(L_{0} + L_{1} \|\nabla f(x^{k})\|)}{2} \\ \leq -\frac{\eta}{2\lambda} \|\nabla f(x^{k})\|^{2} + \frac{\eta^{2}(L_{0} + L_{1} \|\nabla f(x^{k})\|)}{2} + \frac{\eta\lambda}{2} + \frac{\eta}{\lambda} \|\mathbf{b}(x^{k})\|^{2} \\ \leq -\frac{\eta}{2} \|\nabla f(x^{k})\| + \frac{\eta^{2}(L_{0} + L_{1} \|\nabla f(x^{k})\|)}{2} + \frac{\eta\lambda}{2} + \frac{\eta}{\lambda} \|\mathbf{b}(x^{k})\|^{2} \\ = -\frac{\eta}{2} \|\nabla f(x^{k})\| + \frac{\eta^{2}(L_{0} + L_{1} \|\nabla f(x^{k})\|)}{\|\nabla f(x^{k})\|} + \frac{\eta\lambda}{2} + \frac{\eta}{\lambda} \|\mathbf{b}(x^{k})\|^{2} \\ \leq -\frac{\eta}{4} \|\nabla f(x^{k})\| \left(1 - \frac{\eta(L_{0} + L_{1} \|\nabla f(x^{k})\|)}{\|\nabla f(x^{k})\|} \right) + \frac{\eta\lambda}{2} + \frac{\eta}{\lambda} \|\mathbf{b}(x^{k})\|^{2} \\ \leq -\frac{\eta}{4} \|\nabla f(x^{k})\| + \frac{\eta\lambda}{2} + \frac{\eta\zeta^{2}}{\lambda}.$$
(68)

$$\frac{\left\|\nabla f(x^{k})\right\|}{2\left(L_{0}+L_{1}\left\|\nabla f(x^{k})\right\|\right)} = \frac{1}{2\left(L_{0}\frac{1}{\left\|\nabla f(x^{k})\right\|}+L_{1}\right)} = \frac{\lambda}{2\left(L_{0}\frac{\lambda}{\left\|\nabla f(x^{k})\right\|}+L_{1}\lambda\right)} \ge \frac{\lambda}{2\left(L_{0}+L_{1}\lambda\right)}.$$

Thus, $\eta_k = \eta \leq \frac{\lambda}{2(L_0 + L_1\lambda)}$.

Using the convexity assumption of the function, we have the following:

$$f(x^{k}) - f^{*} \leq \left\langle \nabla f(x^{k}), x^{k} - x^{*} \right\rangle \stackrel{(6)}{\leq} \left\| \nabla f(x^{k}) \right\| \left\| x^{k} - x^{*} \right\| \leq \left\| \nabla f(x^{k}) \right\| \underbrace{\left\| x^{0} - x^{*} \right\|}_{R}.$$

Hence we have:

$$\left\|\nabla f(x^k)\right\| \ge \frac{f(x^k) - f^*}{R}.$$
(69)

Then substituting (69) into (68) we obtain:

$$\mathbb{E}\left[f(x^{k+1})\right] - f(x^k) \le -\frac{\eta}{4} \left\|\nabla f(x^k)\right\| + \frac{\eta\lambda}{2} + \frac{\eta\zeta^2}{\lambda} \le -\frac{\eta}{4R}(f(x^k) - f^*) + \frac{\eta\lambda}{2} + \frac{\eta\zeta^2}{\lambda}.$$

This inequality is equivalent to the trailing inequality:

$$\mathbb{E}\left[f(x^{k+1})\right] - f^* \le \left(1 - \frac{\eta}{4R}\right)\left(f(x^k) - f^*\right) + \frac{\eta\lambda}{2} + \frac{\eta\zeta^2}{\lambda}.$$

Then for k = 0, 1, 2, ..., N - 1 iterations that satisfy the conditions $\|\nabla f(x^k)\| \ge \lambda$ and $\|\mathbf{g}(x^k, \boldsymbol{\xi}^k)\| \le \lambda$ NSGD with biased gradient oracle shows linear convergence:

$$\mathbb{E}\left[f(x^N)\right] - f^* \le \left(1 - \frac{\eta}{4R}\right)^N \left(f(x^0) - f^*\right) + 2\lambda R + \frac{2\zeta^2 R}{\lambda}.$$

Combining all the cases considered, we obtain the convergence rate of NSGD with biased gradient oracle:

$$\mathbb{E}\left[f(x^{N})\right] - f^* \lesssim \left(1 - \frac{\eta}{R}\right)^{N} \left(f(x^{0}) - f^*\right) + \frac{MR}{\lambda^2} \left(\frac{\sigma^2}{B} + \zeta^2\right) + \lambda R.$$
(70)

E.2 Convergence Results for ZO-NSGD (Proof of the Theorem 5.4)

In order to obtain the convergence rate of ZO-NSGD in the convex setting, we need to substitute the obtained estimates (50) and (51) into the convergence rate of NSGD (70) instead of ζ and σ^2 , respectively. Then the convergence of ZO-NSGD in the convex setup is as follows:

$$\begin{split} \mathbb{E}\left[f(x^{N})\right] - f^{*} \lesssim \underbrace{\left(1 - \frac{\eta}{R}\right)^{N} \left(f(x^{0}) - f^{*}\right)}_{\oplus} + \underbrace{\frac{dMR\tilde{\sigma}^{2}}{\lambda^{2}B}}_{\circledast} + \underbrace{\frac{dMR\left(L_{0} + L_{1}M\right)^{2}\gamma^{2}}{\lambda^{2}B}}_{\circledast} + \underbrace{\frac{dMR\left(L_{0} + L_{1}M\right)^{2}\gamma^{2}}{\lambda^{2}B\gamma^{2}}}_{\circledast} + \underbrace{\frac{dMR\left(L_{0} + L_{1}M\right)^{2}\gamma^{2}}{\varepsilon}}_{\circledast} + \underbrace{\frac{d^{2}MR\Delta^{2}}{\lambda^{2}\gamma^{2}}}_{\circledast} + \underbrace{\frac{\lambda R}{\varepsilon}}_{\circledast}. \end{split}$$

From term $\overline{\mathbb{O}}$, we find the hyperparameter λ :

0

From term ①, we find the number of iterations N required for Algorithm 4 in convex setup to achieve ε -accuracy:

$$①: \quad \left(1 - \frac{\eta}{R}\right)^{N} \left(f(x^{0}) - f^{*}\right) \leq \varepsilon \quad \Rightarrow \quad N \geq \frac{R}{\eta} \log \frac{\left(f(x^{0}) - f^{*}\right)}{\varepsilon}; \\
 N = \tilde{\mathcal{O}}\left(\frac{R}{\eta}\right).$$
(72)

From terms ②, we find the batch size B:

From terms (3) and (5) we find the smoothing parameter γ :

$$\begin{aligned}
\textcircled{3}: \quad \frac{dMR\left(L_{0}+L_{1}M\right)^{2}\gamma^{2}}{\lambda^{2}B} &\leq \varepsilon \quad \Rightarrow \quad \gamma \leq \sqrt{\frac{\varepsilon\lambda^{2}B}{dMR\left(L_{0}+L_{1}M\right)^{2}}} \stackrel{(73):(71)}{=} \frac{\tilde{\sigma}}{(L_{0}+L_{1}M)}; \\
\textcircled{3}: \quad \frac{MR\left(L_{0}+L_{1}M\right)^{2}\gamma^{2}}{\lambda^{2}} \leq \varepsilon \quad \Rightarrow \quad \gamma \leq \frac{\sqrt{\varepsilon^{3}}}{\sqrt{M}R^{3/2}\left(L_{0}+L_{1}M\right)}; \\
\gamma \leq \frac{1}{(L_{0}+L_{1}M)}\min\left\{\tilde{\sigma}, \frac{\varepsilon^{3/2}}{\sqrt{M}R^{3/2}}\right\} = \frac{\varepsilon^{3/2}}{(L_{0}+L_{1}M)\sqrt{M}R^{3/2}}.
\end{aligned}$$

$$(74)$$

From the remaining terms (1) and (6), we find the maximum allowable level of adversarial noise Δ that still guarantees the convergence of the ZO-NSGD to desired accuracy ε in convex setup:

$$\textcircled{ : } \quad \underbrace{ \frac{d^2 M R \Delta^2}{\lambda^2 B \gamma^2} \leq \varepsilon }_{\lambda^2 B \gamma^2} \leq \varepsilon \quad \Rightarrow \quad \Delta \leq \frac{\sqrt{\varepsilon} \lambda \gamma \sqrt{B}}{d \sqrt{M R}} \stackrel{(73),(74),(71)}{=} \frac{\varepsilon^{3/2} \tilde{\sigma}}{\sqrt{d} \left(L_0 + L_1 M \right) R^{3/2}};$$

In this way, the ZO-NSGD achieves ε -accuracy: $\mathbb{E}\left[f(x^N) - f^*\right] \leq \varepsilon$ in convex setup after

$$N \stackrel{(72)}{=} \tilde{\mathcal{O}}\left(\frac{R}{\eta}\right), \quad T = N \cdot B \stackrel{(72),(73)}{=} \mathcal{O}\left(\frac{d\tilde{\sigma}^2 M R^4}{\varepsilon^3 \eta}\right)$$

number of iterations, total number of zero-order oracle calls and at

$$\Delta \stackrel{(75)}{\lesssim} \frac{\varepsilon^{3/2}}{\sqrt{d} \left(L_0 + L_1 M\right) R^{3/2}} \min\left\{\tilde{\sigma}, \frac{\varepsilon^{3/2}}{\sqrt{d} R^{3/2}}\right\}$$

the maximum level of noise with smoothing parameter $\frac{\varepsilon^{3/2}}{(L_0+L_1M)\sqrt{MR^{3/2}}}$ (74).

F Additional Clarification

In this section, we would like to clarify the convergence in the case $L_0 = 0$ (Remark 1.3). In this case the problem does not reach a minimum (hence $R = \arg \inf f(x) = +\infty$). Therefore, we exemplify the special case of NSGD (when $\|\nabla f(x^k, \boldsymbol{\xi}^k)\| \ge \sqrt{2}\sigma$ and $\|\nabla f(x^k)\| \ge \sqrt{2}\sigma$), shows that it is possible to achieve the desired accuracy ε in a finite number of iterations.

Let's introduce the notation $G(x^k, \boldsymbol{\xi}^k) = \frac{\nabla f(x^k, \boldsymbol{\xi}^k)}{\|\nabla f(x^k, \boldsymbol{\xi}^k)\|}$, then using (L_0, L_1) -smoothness (see Assumption 1.2):

$$f(x^{k+1}) - f(x^k) \stackrel{(8)}{\leq} \left\langle \nabla f(x^k), x^{k+1} - x^k \right\rangle + \frac{L_0 + L_1 \left\| \nabla f(x^k) \right\|}{2} \left\| x^{k+1} - x^k \right\|^2$$
$$= -\eta \left\langle \nabla f(x^k), G(x^k, \boldsymbol{\xi}^k) \right\rangle + \frac{\eta^2 (L_0 + L_1 \left\| \nabla f(x^k) \right\|)}{2} \left\| G(x^k, \boldsymbol{\xi}^k) \right\|^2.$$
(76)

Let us evaluate first summand of (76) with $\alpha = \left\| \nabla f(x^k) \right\|^{-1}$:

$$\begin{split} -\eta \left\langle \nabla f(x^k), G(x^k, \boldsymbol{\xi}^k) \right\rangle &\stackrel{(5)}{=} -\frac{\alpha \eta}{2} \left\| \nabla f(x^k) \right\|^2 - \frac{\eta}{2\alpha} \left\| G(x^k, \boldsymbol{\xi}^k) \right\|^2 \\ &\quad + \frac{\eta}{2\alpha} \left\| G(x^k, \boldsymbol{\xi}^k) - \alpha \nabla f(x^k) \right\|^2 \\ &\quad = -\frac{\eta}{2} \left\| \nabla f(x^k) \right\| - \frac{\eta}{2\alpha} \left\| G(x^k, \boldsymbol{\xi}^k) \right\|^2 \\ &\quad + \frac{\eta}{2\lambda^2 \alpha} \left\| \lambda G(x^k, \boldsymbol{\xi}^k) - \lambda \alpha \nabla f(x^k) \right\|^2 \\ &\quad = -\frac{\eta}{2} \left\| \nabla f(x^k) \right\| - \frac{\eta}{2\alpha} \left\| G(x^k, \boldsymbol{\xi}^k) \right\|^2 \\ &\quad + \frac{\eta}{2\lambda^2 \alpha} \left\| \operatorname{clip}_{\lambda} \left(\nabla f(x^k, \boldsymbol{\xi}^k) \right) - \operatorname{clip}_{\lambda} \left(\nabla f(x^k) \right) \right\|^2 \end{split}$$

Using that clipping is a projection on onto a convex set, namely ball with radius λ , and thus is Lipshitz operator with Lipshitz constant 1, we can obtain:

$$-\eta \left\langle \nabla f(x^{k}), \mathbb{E}\left[G(x^{k}, \boldsymbol{\xi}^{k})\right] \right\rangle \leq -\frac{\eta}{2} \left\| \nabla f(x^{k}) \right\| - \frac{\eta}{2\alpha} \mathbb{E}\left[\left\| G(x^{k}, \boldsymbol{\xi}^{k}) \right\|^{2} \right] + \frac{\eta}{2\lambda^{2}\alpha} \mathbb{E}\left[\left\| \nabla f(x^{k}, \boldsymbol{\xi}^{k}) - \nabla f(x^{k}) \right\|^{2} \right].$$
(77)

Using this in (77), we have the following with $\eta_k \leq \frac{\|\nabla f(x^k)\|}{2(L_0+L_1\|\nabla f(x^k)\|)}$:

$$\mathbb{E}\left[f(x^{k+1})\right] - f(x^k) \stackrel{(76)}{\leq} -\eta \left\langle \nabla f(x^k), \mathbb{E}\left[G(x^k, \boldsymbol{\xi}^k)\right] \right\rangle + \frac{\eta^2 (L_0 + L_1 \left\| \nabla f(x^k) \right\|)}{2} \mathbb{E}\left[\left\|G(x^k, \boldsymbol{\xi}^k)\right\|^2\right]$$

$$\stackrel{(77)}{\leq} -\frac{\eta}{2} \left\| \nabla f(x^{k}) \right\| - \frac{\eta}{2\alpha} \mathbb{E} \left[\left\| G(x^{k}, \boldsymbol{\xi}^{k}) \right\|^{2} \right] + \frac{\eta}{2\lambda^{2}\alpha} \mathbb{E} \left[\left\| \nabla f(x^{k}, \boldsymbol{\xi}) - \nabla f(x^{k}) \right\|^{2} \right]$$

$$+ \frac{\eta^{2} (L_{0} + L_{1} \left\| \nabla f(x^{k}) \right\|)}{2} \mathbb{E} \left[\left\| G(x^{k}, \boldsymbol{\xi}^{k}) \right\|^{2} \right]$$

$$= -\frac{\eta}{2} \left\| \nabla f(x^{k}) \right\| + \frac{\eta}{2\lambda^{2}\alpha} \mathbb{E} \left[\left\| \nabla f(x^{k}, \boldsymbol{\xi}^{k}) - \nabla f(x^{k}) \right\|^{2} \right]$$

$$- \frac{\eta}{2} \mathbb{E} \left[\left\| G(x^{k}, \boldsymbol{\xi}^{k}) \right\|^{2} \right] \left(1 - \frac{\eta (L_{0} + L_{1} \left\| \nabla f(x^{k}) \right\|)}{\left\| \nabla f(x^{k}) \right\|} \right)$$

$$\leq -\frac{\eta}{2} \left\| \nabla f(x^{k}) \right\| + \frac{\eta \sigma^{2}}{2\lambda^{2}\alpha}$$

$$\leq -\frac{\eta}{2} \left\| \nabla f(x^{k}) \right\| + \frac{\eta}{4} \left\| \nabla f(x^{k}) \right\|$$

$$= -\frac{\eta}{4} \left\| \nabla f(x^{k}) \right\|.$$

$$(78)$$

$$\frac{\left\|\nabla f(x^k)\right\|}{2\left(L_0 + L_1 \left\|\nabla f(x^k)\right\|\right)} = \frac{1}{2\left(L_0 \frac{1}{\left\|\nabla f(x^k)\right\|} + L_1\right)} = \frac{\lambda}{2\left(L_0 \frac{\lambda}{\left\|\nabla f(x^k)\right\|} + L_1\lambda\right)} \ge \frac{\lambda}{2\left(L_0 + L_1\lambda\right)}.$$
Thus $n_i = n \le \frac{\lambda}{2}$

Thus, $\eta_k = \eta \leq \frac{\lambda}{2(L_0 + L_1\lambda)}$.

We introduce the hyperparameter of the algorithm $R_s = ||x^0 - s||$. Then using the convexity assumption of the function, we have the following:

$$f(x^{k}) - f(s) \leq \left\langle \nabla f(x^{k}), x^{k} - s \right\rangle$$

$$\stackrel{(6)}{\leq} \left\| \nabla f(x^{k}) \right\| \left\| x^{k} - s \right\|$$

$$\leq \left\| \nabla f(x^{k}) \right\| \underbrace{\left\| x^{0} - s \right\|}_{R_{s}}.$$

Hence we have:

$$\left\|\nabla f(x^k)\right\| \ge \frac{f(x^k) - f(s)}{R_s}.$$
(79)

Then substituting (79) into (78) we obtain:

$$\mathbb{E}\left[f(x^{k+1})\right] - f(x^k) \le -\frac{\eta}{4} \left\|\nabla f(x^k)\right\| \le -\frac{\eta}{4R_s}(f(x^k) - f(s)).$$

This inequality is equivalent to the trailing inequality:

$$\mathbb{E}\left[f(x^{k+1})\right] - f^* \le \left(1 - \frac{\eta}{4R_s}\right) \left(f(x^k) - f^*\right) + \frac{\eta}{4R_s} (f(s) - f^*).$$

Then for k = 0, 1, 2, ..., N - 1 iterations that satisfy the conditions $\|\nabla f(x^k, \boldsymbol{\xi}^k)\| \ge \sqrt{2}\sigma$ and $\|\nabla f(x^k)\| \ge \sqrt{2}\sigma$ NSGD shows linear convergence:

$$f(x^N) - f^* \le \left(1 - \frac{\eta}{4R_s}\right)^N (f(x^0) - f^*) + f(s) - f^*.$$

Thus, we have shown that it is indeed possible to converge to a linear rate of convergence on logistic regression using the hyperparameter R_s .