

APPENDIX

Generalized Smoothness in Stochastic Convex Optimization: First- and Zero-Order Methods

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 $a, b \in \mathbb{R}^d$ ($d \geq 1$) the following equality holds:

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

$$\langle a, b \rangle \leq \|a\| \cdot \|b\|. \quad (6)$$

Squared norm of the sum For all $a_1, \dots, a_n \in \mathbb{R}^d$

$$\|a_1 + \dots + a_n\|^2 \leq n\|a_1\|^2 + \dots + n\|a_n\|^2. \quad (7)$$

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 $x, y \in \mathbb{R}^d$:

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

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

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

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

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 $\forall x \in \mathbb{R}^d$:

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

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

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

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

□

756 A.3 AN UPPER BOUND ON THE GRADIENT NORM
757

758 If Assumption 1.2 holds, then it also holds that $\forall x \in \mathbb{R}^d$:
759

$$760 \quad \|\nabla f(x)\| \leq L_0 + \frac{4L_1 + 1}{2}(f(x) - f^*). \quad (11)$$

761
762 *Proof.* We start with equation 10:
763

$$764 \quad \begin{aligned} \|\nabla f(x)\|^2 &\leq 2(L_0 + L_1 \|\nabla f(x)\|)(f(x) - f^*) \\ \iff \|\nabla f(x)\|^2 - 2L_1 \|\nabla f(x)\| (f(x) - f^*) - 2L_0 \|\nabla f(x)\| &\leq 0. \end{aligned}$$

765 We need to solve this quadratic inequality w.r.t. $\|\nabla f(x)\|$. The discriminant is
766

$$767 \quad 4L_1^2(f(x) - f^*) + 8L_0(f(x) - f^*) > 0,$$

768 i.e., it is positive. Since $\|\nabla f(x)\| \geq 0$, we should also satisfy
769

$$770 \quad \begin{aligned} \|\nabla f(x)\| &\leq \frac{2L_1(f(x) - f^*) + \sqrt{4L_1^2(f(x) - f^*) + 8L_0(f(x) - f^*)}}{2} \\ &\leq L_1(f(x) - f^*) + \sqrt{L_1^2(f(x) - f^*) + 2L_0(f(x) - f^*)} \\ &\leq 2L_1(f(x) - f^*) + L_0 + \frac{1}{2}(f(x) - f^*) \\ &= L_0 + \frac{4L_1 + 1}{2}(f(x) - f^*) \end{aligned}$$

771
772
773
774
775
776
777
778
779 \square

780 A.4 WIRTINGER-POINCARÉ INEQUALITY
781

782 Let f is differentiable, then for all $x \in \mathbb{R}^d$, $\gamma e \in S^d(\gamma)$:
783

$$784 \quad \mathbb{E} [f(x + \gamma e)^2] \leq \frac{\gamma^2}{d} \mathbb{E} [\|\nabla f(x + \gamma e)\|^2]. \quad (12)$$

785
786
787
788 B CLIPPED STOCHASTIC GRADIENT DESCENT (PROOF OF THE
789 THEOREM 3.1)
790

791 We start by using (L_0, L_1) -smoothness (see Assumption 1.2):
792

$$793 \quad \begin{aligned} f(x^{k+1}) - f(x^k) &\stackrel{\text{equation 8}}{\leq} \langle \nabla f(x^k), x^{k+1} - x^k \rangle + \frac{L_0 + L_1 \|\nabla f(x^k)\|}{2} \|x^{k+1} - x^k\|^2 \\ &= -\eta \langle \nabla f(x^k), \text{clip}_c(\nabla f(x^k, \xi^k)) \rangle \\ &\quad + \frac{\eta^2(L_0 + L_1 \|\nabla f(x^k)\|)}{2} \|\text{clip}_c(\nabla f(x^k, \xi^k))\|^2. \end{aligned} \quad (13)$$

794
795
796
797
798 Next, we consider three cases depending on the gradient norm: $\|\nabla f(x^k)\| \geq c$ – the full gradient is
800 clipped and $\frac{c}{2} \leq \|\nabla f(x^k)\| \leq c$ and $\|\nabla f(x^k)\| \leq \frac{c}{2}$ – the full gradient is not clipped.
801

802 B.1 FIRST CASE: $\|\nabla f(x^k)\| \geq c$
803

804 In this case $\alpha \nabla f(x^k) = \text{clip}_c(\nabla f(x^k))$ with $\alpha = \min \left\{ 1, \frac{c}{\|\nabla f(x^k)\|} \right\} = \frac{c}{\|\nabla f(x^k)\|}$, therefore we
805 have the following
806

$$807 \quad \begin{aligned} -\eta \langle \nabla f(x^k), \text{clip}_c(\nabla f(x^k, \xi^k)) \rangle &\stackrel{\text{equation 5}}{=} -\frac{\alpha\eta}{2} \|\nabla f(x^k)\|^2 - \frac{\eta}{2\alpha} \|\text{clip}_c(\nabla f(x^k, \xi^k))\|^2 \\ &\quad + \frac{\eta}{2\alpha} \|\text{clip}_c(\nabla f(x^k, \xi^k)) - \alpha \nabla f(x^k)\|^2 \end{aligned}$$

$$\begin{aligned}
&= -\frac{\alpha\eta}{2} \|\nabla f(x^k)\|^2 - \frac{\eta}{2\alpha} \|\text{clip}_c(\nabla f(x^k, \xi^k))\|^2 \\
&\quad + \frac{\eta}{2\alpha} \|\text{clip}_c(\nabla f(x^k, \xi^k)) - \text{clip}_c(\nabla f(x^k))\|^2 \\
&= -\frac{c\eta}{2} \|\nabla f(x^k)\| - \frac{\eta}{2\alpha} \|\text{clip}_c(\nabla f(x^k, \xi^k))\|^2 \\
&\quad + \frac{\eta}{2\alpha} \|\text{clip}_c(\nabla f(x^k, \xi^k)) - \text{clip}_c(\nabla f(x^k))\|^2.
\end{aligned}$$

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

$$\begin{aligned}
-\eta \langle \nabla f(x^k), \mathbb{E}[\text{clip}_c(\nabla f(x^k, \xi^k))] \rangle &\leq -\frac{c\eta}{2} \|\nabla f(x^k)\| - \frac{\eta}{2\alpha} \mathbb{E}[\|\text{clip}_c(\nabla f(x^k, \xi^k))\|^2] \\
&\quad + \frac{\eta}{2\alpha} \mathbb{E}[\|\nabla f(x^k, \xi^k) - \nabla f(x^k)\|^2] \\
&\leq -\frac{c\eta}{2} \|\nabla f(x^k)\| - \frac{\eta}{2\alpha} \mathbb{E}[\|\text{clip}_c(\nabla f(x^k, \xi^k))\|^2] \\
&\quad + \frac{\eta\sigma^2}{2\alpha B} \\
&= -\frac{c\eta}{2} \|\nabla f(x^k)\| - \frac{\eta}{2\alpha} \mathbb{E}[\|\text{clip}_c(\nabla f(x^k, \xi^k))\|^2] \\
&\quad + \frac{\eta \|\nabla f(x^k)\| \sigma^2}{2cB}. \tag{14}
\end{aligned}$$

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

In the case $c \geq \sqrt{2}\sigma$ We have in equation 14:

$$\begin{aligned}
-\eta \langle \nabla f(x^k), \mathbb{E}[\text{clip}_c(\nabla f(x^k, \xi^k))] \rangle &\stackrel{\text{equation 14}}{\leq} -\frac{c\eta}{2} \|\nabla f(x^k)\| - \frac{\eta}{2\alpha} \mathbb{E}[\|\text{clip}_c(\nabla f(x^k, \xi^k))\|^2] \\
&\quad + \frac{\eta \|\nabla f(x^k)\| \sigma^2}{2cB} \\
&= -\frac{\eta}{2\alpha} \mathbb{E}[\|\text{clip}_c(\nabla f(x^k, \xi^k))\|^2] \\
&\quad - \frac{c\eta}{2} \|\nabla f(x^k)\| \left(1 - \frac{\sigma^2}{c^2 B}\right) \\
&\leq -\frac{\eta}{2\alpha} \mathbb{E}[\|\text{clip}_c(\nabla f(x^k, \xi^k))\|^2] \\
&\quad - \frac{c\eta}{4} \|\nabla f(x^k)\| \\
&= -\frac{\eta \|\nabla f(x^k)\|}{2c} \mathbb{E}[\|\text{clip}_c(\nabla f(x^k, \xi^k))\|^2] \\
&\quad - \frac{c\eta}{4} \|\nabla f(x^k)\|.
\end{aligned}$$

Plugging this into equation 13 and choosing $\eta \leq \frac{1}{4(L_0 + L_1 c)}$ we have:

$$\begin{aligned}
\mathbb{E}[f(x^{k+1})] - f(x^k) &\stackrel{\text{equation 13}}{\leq} -\frac{\eta \|\nabla f(x^k)\|}{2c} \mathbb{E}[\|\text{clip}_c(\nabla f(x^k, \xi^k))\|^2] - \frac{c\eta}{4} \|\nabla f(x^k)\| \\
&\quad + \frac{\eta^2(L_0 + L_1 \|\nabla f(x^k)\|)}{2} \mathbb{E}[\|\text{clip}_c(\nabla f(x^k, \xi^k))\|^2] \\
&= -\frac{\eta \|\nabla f(x^k)\|}{2c} \mathbb{E}[\|\text{clip}_c(\nabla f(x^k, \xi^k))\|^2] (1 - \eta L_1 c) \\
&\quad - \frac{c\eta}{4} \|\nabla f(x^k)\| + \frac{\eta^2 L_0}{2} \mathbb{E}[\|\text{clip}_c(\nabla f(x^k, \xi^k))\|^2] \\
&\leq -\frac{c\eta}{4} \|\nabla f(x^k)\| - \frac{\eta}{2} \mathbb{E}[\|\text{clip}_c(\nabla f(x^k, \xi^k))\|^2] (1 - \eta(L_0 + L_1 c))
\end{aligned}$$

$$\leq -\frac{c\eta}{4} \|\nabla f(x^k)\|. \quad (15)$$

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

$$\begin{aligned} f(x^k) - f^* &\leq \langle \nabla f(x^k), x^k - x^* \rangle \\ &\stackrel{\text{equation 6}}{\leq} \|\nabla f(x^k)\| \|x^k - x^*\| \\ &\leq \|\nabla f(x^k)\| \underbrace{\max_{k \in [0, N-1]} \|x^k - x^*\|}_R. \end{aligned}$$

Hence we have:

$$\|\nabla f(x^k)\| \geq \frac{f(x^k) - f^*}{R}. \quad (16)$$

Then substituting equation 16 into equation 15 we obtain:

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

This inequality is equivalent to the trailing inequality:

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

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

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

In the case $c \leq \sqrt{2}\sigma$ We have in equation 14:

$$\begin{aligned} -\eta \langle \nabla f(x^k), \mathbb{E} [\text{clip}_c(f(x^k, \xi^k))] \rangle &\stackrel{\text{equation 14}}{\leq} -\frac{c\eta}{2} \|\nabla f(x^k)\| - \frac{\eta}{2\alpha} \mathbb{E} [\|\text{clip}_c(\nabla f(x^k, \xi^k))\|^2] \\ &\quad + \frac{\eta \|\nabla f(x^k)\| \sigma^2}{2cB} \\ &= -\frac{c\eta}{2} \|\nabla f(x^k)\| - \frac{\eta}{2\alpha} \mathbb{E} [\|\text{clip}_c(\nabla f(x^k, \xi^k))\|^2] \\ &\quad + \frac{\eta M \sigma^2}{2cB}. \end{aligned}$$

Plugging this into equation 13 and choosing $\eta \leq \frac{1}{4(L_0 + L_1 c)}$ we have:

$$\begin{aligned} \mathbb{E} [f(x^{k+1})] - f(x^k) &\stackrel{\text{equation 13}}{\leq} -\frac{c\eta}{2} \|\nabla f(x^k)\| - \frac{\eta \|\nabla f(x^k)\|}{2c} \mathbb{E} [\|\text{clip}_c(\nabla f(x^k, \xi^k))\|^2] \\ &\quad + \frac{\eta^2(L_0 + L_1 \|\nabla f(x^k)\|)}{2} \mathbb{E} [\|\text{clip}_c(\nabla f(x^k, \xi^k))\|^2] + \frac{\eta M \sigma^2}{2cB} \\ &= -\frac{c\eta}{2} \|\nabla f(x^k)\| - \frac{\eta \|\nabla f(x^k)\|}{2c} \mathbb{E} [\|\text{clip}_c(\nabla f(x^k, \xi^k))\|^2] (1 - \eta L_1 c) \\ &\quad + \frac{\eta^2 L_0}{2} \mathbb{E} [\|\text{clip}_c(\nabla f(x^k, \xi^k))\|^2] + \frac{\eta M \sigma^2}{2cB} \\ &\leq -\frac{c\eta}{2} \|\nabla f(x^k)\| - \frac{\eta}{2} \mathbb{E} [\|\text{clip}_c(\nabla f(x^k, \xi^k))\|^2] (1 - \eta(L_0 + L_1 c)) \\ &\quad + \frac{\eta M \sigma^2}{2cB} \\ &\leq -\frac{c\eta}{2} \|\nabla f(x^k)\| + \frac{\eta M \sigma^2}{2cB}. \end{aligned} \quad (17)$$

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

$$\begin{aligned} f(x^k) - f^* &\leq \langle \nabla f(x^k), x^k - x^* \rangle \\ &\stackrel{\text{equation 6}}{\leq} \|\nabla f(x^k)\| \|x^k - x^*\| \\ &\leq \|\nabla f(x^k)\| \underbrace{\max_{k \in [0, N-1]} \|x^k - x^*\|}_R. \end{aligned}$$

Hence we have:

$$\|\nabla f(x^k)\| \geq \frac{f(x^k) - f^*}{R}. \quad (18)$$

Then substituting equation 18 into equation 17 we obtain:

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

This inequality is equivalent to the trailing inequality:

$$\begin{aligned} \mathbb{E} [f(x^{k+1})] - f^* &\leq \left(1 - \frac{\eta c}{2R}\right) (f(x^k) - f^*) + \frac{\eta M \sigma^2}{2cB} \\ &\stackrel{\text{equation 11}}{\leq} \left(1 - \frac{\eta c}{2R} + \frac{\eta \sigma^2 (4L_1 + 1)}{4cB}\right) (f(x^k) - f^*) + \frac{\eta \sigma^2 L_0}{2cB} \\ &\stackrel{\textcircled{1}}{\leq} \left(1 - \frac{\eta}{4R}\right) (f(x^k) - f^*) + \frac{\eta \sigma^2 L_0}{2cB}, \end{aligned}$$

where in $\textcircled{1}$ we used $B = \mathcal{O}\left(\frac{\sigma^2}{\varepsilon} \left[\eta + \frac{L_0 R}{c^2} + \frac{R}{c} + \frac{(4L_1 + 1)R\varepsilon}{c}\right]\right)$.

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

$$\mathbb{E} [f(x^N)] - f^* \leq \left(1 - \frac{\eta c}{4R}\right)^N (f(x^0) - f^*) + \frac{L_0 R \sigma^2}{c^2 B}.$$

B.2 SECOND CASE: $\frac{c}{2} \leq \|\nabla f(x^k)\| \leq c$

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

$$\begin{aligned} -\eta \langle \nabla f(x^k), \text{clip}_c(\nabla f(x^k, \xi^k)) \rangle &\stackrel{\text{equation 5}}{=} -\frac{\alpha \eta}{2} \|\nabla f(x^k)\|^2 - \frac{\eta}{2\alpha} \|\text{clip}_c(\nabla f(x^k, \xi^k))\|^2 \\ &\quad + \frac{\eta}{2\alpha} \|\text{clip}_c(\nabla f(x^k, \xi^k)) - \alpha \nabla f(x^k)\|^2 \\ &= -\frac{\eta}{2} \|\nabla f(x^k)\|^2 - \frac{\eta}{2} \|\text{clip}_c(\nabla f(x^k, \xi^k))\|^2 \\ &\quad + \frac{\eta}{2} \|\text{clip}_c(\nabla f(x^k, \xi^k)) - \text{clip}_c(\nabla f(x^k))\|^2 \\ &\leq -\frac{c\eta}{4} \|\nabla f(x^k)\| - \frac{\eta}{2} \|\text{clip}_c(\nabla f(x^k, \xi^k))\|^2 \\ &\quad + \frac{\eta}{2} \|\text{clip}_c(\nabla f(x^k, \xi^k)) - \text{clip}_c(\nabla f(x^k))\|^2. \end{aligned}$$

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

$$\begin{aligned} -\eta \langle \nabla f(x^k), \mathbb{E} [\text{clip}_c(\nabla f(x^k, \xi^k))] \rangle &\leq -\frac{c\eta}{4} \|\nabla f(x^k)\| - \frac{\eta}{2} \mathbb{E} \left[\|\text{clip}_c(\nabla f(x^k, \xi^k))\|^2 \right] \\ &\quad + \frac{\eta}{2} \mathbb{E} \left[\|\nabla f(x^k, \xi^k) - \nabla f(x^k)\|^2 \right] \end{aligned}$$

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

Plugging this into equation 13 and choosing $\eta \leq \frac{1}{4(L_0+L_1c)}$ we have:

$$\begin{aligned}
\mathbb{E} [f(x^{k+1})] - f(x^k) &\stackrel{\text{equation 13}}{\leq} -\frac{c\eta}{4} \|\nabla f(x^k)\| - \frac{\eta}{2} \mathbb{E} \left[\|\text{clip}_c(\nabla f(x^k, \xi^k))\|^2 \right] \\
&\quad + \frac{\eta^2(L_0 + L_1 \|\nabla f(x^k)\|)}{2} \mathbb{E} \left[\|\text{clip}_c(\nabla f(x^k, \xi^k))\|^2 \right] + \frac{\eta\sigma^2}{2B} \\
&= -\frac{c\eta}{4} \|\nabla f(x^k)\| + \frac{\eta\sigma^2}{2B} \\
&\quad - \frac{\eta}{2} \mathbb{E} \left[\|\text{clip}_c(\nabla f(x^k, \xi^k))\|^2 \right] (1 - \eta(L_0 + L_1 \|\nabla f(x^k)\|)) \\
&\leq -\frac{c\eta}{4} \|\nabla f(x^k)\| + \frac{\eta\sigma^2}{2B}. \tag{19}
\end{aligned}$$

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

$$\begin{aligned}
f(x^k) - f^* &\leq \langle \nabla f(x^k), x^k - x^* \rangle \\
&\stackrel{\text{equation 6}}{\leq} \|\nabla f(x^k)\| \|x^k - x^*\| \\
&\leq \|\nabla f(x^k)\| \underbrace{\max_{k \in [0, N-1]} \|x^k - x^*\|}_R.
\end{aligned}$$

Hence we have:

$$\|\nabla f(x^k)\| \geq \frac{f(x^k) - f^*}{R}. \tag{20}$$

Then substituting equation 20 into equation 19 we obtain:

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

This inequality is equivalent to the trailing inequality:

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

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

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

B.3 THIRD CASE: $\|\nabla f(x^k)\| \leq \frac{c}{2}$

We introduce an indicative function:

$$\aleph_k = \mathbb{1} \{ \|\nabla f(x^k, \xi^k)\| > c \}. \tag{21}$$

Then the following is true:

$$\mathbb{E} [\aleph_k] = \mathbb{E} [\aleph_k^2] = \mathcal{P} [\|\nabla f(x^k, \xi^k)\| > c] \stackrel{\textcircled{1}}{\leq} \mathcal{P} [\|\nabla f(x^k, \xi^k) - \nabla f(x^k)\| > \frac{c}{2}] \stackrel{\textcircled{2}}{\leq} \frac{4\sigma^2}{c^2 B}, \tag{22}$$

where in $\textcircled{1}$ we used $\|\nabla f(x^k, \xi^k)\| \leq \|\nabla f(x^k, \xi^k) - \nabla f(x^k)\| + \|\nabla f(x^k)\| \leq \|\nabla f(x^k, \xi^k) - \nabla f(x^k)\| + \frac{c}{2}$, and in $\textcircled{2}$ we used Markov's inequality.

Let $r_{k+1} = \mathbb{E} [\|x^{k+1} - x^*\|]$ and $F_{k+1} = \mathbb{E} [f(x^{k+1}) - f^*]$, then given that

$$\begin{aligned} \text{clip}_c(\nabla f(x^k, \xi^k)) &= \nabla f(x^k, \xi^k)(1 - \aleph_k) + \frac{c}{\|\nabla f(x^k, \xi^k)\|} \nabla f(x^k, \xi^k) \aleph_k \\ &= \nabla f(x^k, \xi^k) + \left(\frac{c}{\|\nabla f(x^k, \xi^k)\|} - 1 \right) \nabla f(x^k, \xi^k) \aleph_k \end{aligned}$$

we get with $\eta \leq \frac{1}{4(L_0 + L_1 c)}$:

$$\begin{aligned} r_{k+1}^2 &= r_k^2 - 2\eta \langle \mathbb{E} [\text{clip}_c(\nabla f(x^k, \xi^k))], x^k - x^* \rangle + \eta^2 \mathbb{E} [\|\text{clip}_c(\nabla f(x^k, \xi^k))\|^2] \\ &= r_k^2 - 2\eta \langle \nabla f(x^k), x^k - x^* \rangle - 2\eta \left\langle \mathbb{E} \left[\left(\frac{c}{\|\nabla f(x^k, \xi^k)\|} - 1 \right) \nabla f(x^k, \xi^k) \aleph_k \right], x^k - x^* \right\rangle \\ &\quad + \eta^2 \mathbb{E} [\|\text{clip}_c(\nabla f(x^k, \xi^k))\|^2] \\ &\stackrel{\text{equation 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}{\|\nabla f(x^k, \xi^k)\|} - 1 \right) \nabla f(x^k, \xi^k) \aleph_k \right] \right\| \|x^k - x^*\| \\ &\quad + \eta^2 \mathbb{E} [\|\text{clip}_c(\nabla f(x^k, \xi^k))\|^2] \\ &\stackrel{\textcircled{1}}{\leq} r_k^2 - 2\eta F_k + 2\eta \left\| \mathbb{E} \left[\left(\frac{c}{\|\nabla f(x^k, \xi^k)\|} - 1 \right) \nabla f(x^k, \xi^k) \aleph_k \right] \right\| \|x^0 - x^*\| \\ &\quad + \eta^2 \mathbb{E} [\|\text{clip}_c(\nabla f(x^k, \xi^k))\|^2] \\ &\stackrel{\text{equation 7}}{\leq} r_k^2 - 2\eta F_k + 2\eta \left\| \mathbb{E} \left[\left(\frac{c}{\|\nabla f(x^k, \xi^k)\|} - 1 \right) \nabla f(x^k, \xi^k) \aleph_k \right] \right\| \|x^0 - x^*\| \\ &\quad + 2\eta^2 \mathbb{E} [\|\text{clip}_c(\nabla f(x^k, \xi^k)) - \nabla f(x^k)\|^2] + 2\eta^2 \|\nabla f(x^k)\|^2 \\ &= r_k^2 - 2\eta F_k + 2\eta \left\| \mathbb{E} \left[\left(\frac{c}{\|\nabla f(x^k, \xi^k)\|} - 1 \right) \nabla f(x^k, \xi^k) \aleph_k \right] \right\| R \\ &\quad + 2\eta^2 \mathbb{E} [\|\text{clip}_c(\nabla f(x^k, \xi^k)) - \text{clip}_c(\nabla f(x^k))\|^2] + 2\eta^2 \|\nabla f(x^k)\|^2 \\ &\stackrel{\textcircled{2}}{\leq} r_k^2 - 2\eta F_k + 2\eta \left\| \mathbb{E} \left[\left(\frac{c}{\|\nabla f(x^k, \xi^k)\|} - 1 \right) \nabla f(x^k, \xi^k) \aleph_k \right] \right\| R \\ &\quad + 2\eta^2 \mathbb{E} [\|\nabla f(x^k, \xi^k) - \nabla f(x^k)\|^2] + 2\eta^2 \|\nabla f(x^k)\|^2 \\ &\leq r_k^2 - 2\eta F_k + 2\eta \left\| \mathbb{E} \left[\left(\frac{c}{\|\nabla f(x^k, \xi^k)\|} - 1 \right) \nabla f(x^k, \xi^k) \aleph_k \right] \right\| R \\ &\quad + \frac{2\eta^2 \sigma^2}{B} + 2\eta^2 \|\nabla f(x^k)\|^2 \\ &\stackrel{\text{equation 10}}{\leq} r_k^2 - 2\eta F_k + 2\eta \left\| \mathbb{E} \left[\left(\frac{c}{\|\nabla f(x^k, \xi^k)\|} - 1 \right) \nabla f(x^k, \xi^k) \aleph_k \right] \right\| R \\ &\quad + \frac{2\eta^2 \sigma^2}{B} + 4\eta^2 (L_0 + L_1 \|\nabla f(x^k)\|) F_k \\ &\leq r_k^2 - 2\eta F_k + 2\eta \left\| \mathbb{E} \left[\left(\frac{c}{\|\nabla f(x^k, \xi^k)\|} - 1 \right) \nabla f(x^k, \xi^k) \aleph_k \right] \right\| R \\ &\quad + \frac{2\eta^2 \sigma^2}{B} + 4\eta^2 (L_0 + L_1 c) F_k \\ &= r_k^2 - 2\eta F_k (1 - 2\eta (L_0 + L_1 c)) + \frac{2\eta^2 \sigma^2}{B} \\ &\quad + 2\eta \left\| \mathbb{E} \left[\left(\frac{c}{\|\nabla f(x^k, \xi^k)\|} - 1 \right) \nabla f(x^k, \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, \xi^k)\|} - 1 \right) \nabla f(x^k, \xi^k) \aleph_k \right] \right\| R. \end{aligned} \tag{23}$$

1080
1081
1082
1083
1084
1085
1086
1087
1088
1089
1090
1091
1092
1093
1094
1095
1096
1097
1098

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

$$\begin{aligned}
& 2\eta R \left\| \mathbb{E} \left[\left(\frac{c}{\|\nabla f(x^k, \xi^k)\|} - 1 \right) \nabla f(x^k, \xi^k) \aleph_k \right] \right\| \\
& \stackrel{\text{equation 21}}{\leq} 2\eta R \mathbb{E} \left[\|\nabla f(x^k, \xi^k)\| \cdot \left(1 - \frac{c}{\|\nabla f(x^k, \xi^k)\|} \right) \aleph_k \right] \\
& \leq 2\eta R \mathbb{E} [\|\nabla f(x^k, \xi^k)\| \cdot \aleph_k] \\
& \leq 2\eta R (\mathbb{E} [\|\nabla f(x^k, \xi^k) - \nabla f(x^k)\| \cdot \aleph_k] + \|\nabla f(x^k)\| \mathbb{E} [\aleph_k]) \\
& \leq 2\eta R \left(\sqrt{\mathbb{E} [\|\nabla f(x^k, \xi^k) - \nabla f(x^k)\|^2]} \cdot \mathbb{E} [\aleph_k^2] + \|\nabla f(x^k)\| \mathbb{E} [\aleph_k] \right) \\
& \stackrel{\text{equation 22}}{\leq} 2\eta R \left(\frac{2\sigma^2}{cB} + \frac{c}{2} \cdot \frac{4\sigma^2}{c^2B} \right) \\
& = \frac{8\eta\sigma^2 R}{cB}. \tag{24}
\end{aligned}$$

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

1100
1101
1102
1103
1104
1105

$$\begin{aligned}
\eta F_k & \stackrel{\text{equation 23}}{\leq} r_k^2 - r_{k+1}^2 + \frac{2\eta^2\sigma^2}{B} + 2\eta \left\| \mathbb{E} \left[\left(\frac{c}{\|\nabla f(x^k, \xi^k)\|} - 1 \right) \nabla f(x^k, \xi^k) \aleph_k \right] \right\| R \\
& \stackrel{\text{equation 24}}{\leq} r_k^2 - r_{k+1}^2 + \frac{2\eta^2\sigma^2}{B} + \frac{8\eta\sigma^2 R}{cB}
\end{aligned}$$

Combining all cases we have:

1106
1107
1108
1109
1110

$$\begin{aligned}
\mathbb{E} [f(x^N)] - f^* & \leq F_N \cdot \mathbb{1}[\mathcal{T}_1] + F_N \cdot \mathbb{1}[\mathcal{T}_2] \\
& \leq \left(1 - \frac{\eta c}{4R} \right)^N F_0 + \frac{R^2}{\eta N} + \frac{\sigma^2 L_0 R}{c^2 B} + \frac{2\eta\sigma^2}{B} + \frac{8\sigma^2 R}{cB},
\end{aligned}$$

where \mathcal{T}_1 describes case $\|\nabla f(x^k)\| \geq \frac{c}{2}$, and \mathcal{T}_2 describes case $\|\nabla f(x^k)\| < \frac{c}{2}$.

1111
1112
1113
1114
1115

C NORMALIZED STOCHASTIC GRADIENT DESCENT (PROOF OF THE THEOREM 4.1)

1116
1117
1118
1119

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

1120
1121
1122
1123

$$\begin{aligned}
f(x^{k+1}) - f(x^k) & \stackrel{\text{equation 8}}{\leq} \langle \nabla f(x^k), x^{k+1} - x^k \rangle + \frac{L_0 + L_1 \|\nabla f(x^k)\|}{2} \|x^{k+1} - x^k\|^2 \\
& = -\eta \langle \nabla f(x^k), G(x^k, \xi^k) \rangle + \frac{\eta^2 (L_0 + L_1 \|\nabla f(x^k)\|)}{2} \|G(x^k, \xi^k)\|^2. \tag{25}
\end{aligned}$$

1124
1125
1126

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

1127
1128

C.1 FIRST CASE: $\|\nabla f(x^k)\| \geq \lambda$ AND $\|\nabla f(x^k, \xi^k)\| \geq \lambda$

1129
1130

Let us evaluate first summand of equation 25 with $\alpha = \|\nabla f(x^k)\|^{-1}$:

1131
1132
1133

$$\begin{aligned}
-\eta \langle \nabla f(x^k), G(x^k, \xi^k) \rangle & \stackrel{\text{equation 5}}{=} -\frac{\alpha\eta}{2} \|\nabla f(x^k)\|^2 - \frac{\eta}{2\alpha} \|G(x^k, \xi^k)\|^2 \\
& \quad + \frac{\eta}{2\alpha} \|G(x^k, \xi^k) - \alpha \nabla f(x^k)\|^2
\end{aligned}$$

$$\begin{aligned}
&= -\frac{\eta}{2} \|\nabla f(x^k)\| - \frac{\eta}{2\alpha} \|G(x^k, \xi^k)\|^2 \\
&\quad + \frac{\eta}{2\lambda^2\alpha} \|\lambda G(x^k, \xi^k) - \lambda\alpha \nabla f(x^k)\|^2 \\
&= -\frac{\eta}{2} \|\nabla f(x^k)\| - \frac{\eta}{2\alpha} \|G(x^k, \xi^k)\|^2 \\
&\quad + \frac{\eta}{2\lambda^2\alpha} \|\text{clip}_\lambda(\nabla f(x^k, \xi^k)) - \text{clip}_\lambda(\nabla f(x^k))\|^2
\end{aligned}$$

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

$$\begin{aligned}
-\eta \langle \nabla f(x^k), \mathbb{E}[G(x^k, \xi^k)] \rangle &\leq -\frac{\eta}{2} \|\nabla f(x^k)\| - \frac{\eta}{2\alpha} \mathbb{E}[\|G(x^k, \xi^k)\|^2] \\
&\quad + \frac{\eta}{2\lambda^2\alpha} \mathbb{E}[\|\nabla f(x^k, \xi^k) - \nabla f(x^k)\|^2]. \quad (26)
\end{aligned}$$

In the case: $0 \leq \sigma \leq \frac{\lambda}{\sqrt{2}}$. Using this in equation 26, we have the following with $\eta_k \leq$

$$\begin{aligned}
&\frac{\|\nabla f(x^k)\|}{2(L_0 + L_1 \|\nabla f(x^k)\|)}: \\
\mathbb{E}[f(x^{k+1})] - f(x^k) &\stackrel{\text{equation 25}}{\leq} -\eta \langle \nabla f(x^k), \mathbb{E}[G(x^k, \xi^k)] \rangle + \frac{\eta^2(L_0 + L_1 \|\nabla f(x^k)\|)}{2} \mathbb{E}[\|G(x^k, \xi^k)\|^2] \\
&\stackrel{\text{equation 26}}{\leq} -\frac{\eta}{2} \|\nabla f(x^k)\| - \frac{\eta}{2\alpha} \mathbb{E}[\|G(x^k, \xi^k)\|^2] + \frac{\eta}{2\lambda^2\alpha} \mathbb{E}[\|\nabla f(x^k, \xi^k) - \nabla f(x^k)\|^2] \\
&\quad + \frac{\eta^2(L_0 + L_1 \|\nabla f(x^k)\|)}{2} \mathbb{E}[\|G(x^k, \xi^k)\|^2] \\
&= -\frac{\eta}{2} \|\nabla f(x^k)\| + \frac{\eta}{2\lambda^2\alpha} \mathbb{E}[\|\nabla f(x^k, \xi^k) - \nabla f(x^k)\|^2] \\
&\quad - \frac{\eta}{2} \mathbb{E}[\|G(x^k, \xi^k)\|^2] \left(1 - \frac{\eta(L_0 + L_1 \|\nabla f(x^k)\|)}{\|\nabla f(x^k)\|}\right) \\
&\leq -\frac{\eta}{2} \|\nabla f(x^k)\| + \frac{\eta\sigma^2}{2\lambda^2\alpha} \\
&\leq -\frac{\eta}{2} \|\nabla f(x^k)\| + \frac{\eta}{4} \|\nabla f(x^k)\| \\
&= -\frac{\eta}{4} \|\nabla f(x^k)\|. \quad (27)
\end{aligned}$$

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

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

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

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

$$\begin{aligned}
f(x^k) - f^* &\leq \langle \nabla f(x^k), x^k - x^* \rangle \\
&\stackrel{\text{equation 6}}{\leq} \|\nabla f(x^k)\| \|x^k - x^*\| \\
&\leq \|\nabla f(x^k)\| \underbrace{\max_{k \in [0, N-1]} \|x^k - x^*\|}_R.
\end{aligned}$$

Hence we have:

$$\|\nabla f(x^k)\| \geq \frac{f(x^k) - f^*}{R}. \quad (28)$$

1188 Then substituting equation 28 into equation 27 we obtain:

$$1189 \mathbb{E} [f(x^{k+1})] - f(x^k) \leq -\frac{\eta}{4} \|\nabla f(x^k)\| \leq -\frac{\eta}{4R} (f(x^k) - f^*).$$

1192 This inequality is equivalent to the trailing inequality:

$$1193 \mathbb{E} [f(x^{k+1})] - f^* \leq \left(1 - \frac{\eta}{4R}\right) (f(x^k) - f^*).$$

1196 Then for $k = 0, 1, 2, \dots, N - 1$ iterations that satisfy the conditions $\|\nabla f(x^k, \xi^k)\| \geq \sqrt{2}\sigma$ and $\|\nabla f(x^k)\| \geq \sqrt{2}\sigma$ NSGD shows linear convergence:

$$1199 \mathbb{E} [f(x^N)] - f^* \leq \left(1 - \frac{\eta}{4R}\right)^N (f(x^0) - f^*).$$

1202 **In the case:** $\frac{\lambda}{\sqrt{2}} \leq \sigma$. Using this in equation 26, we have the following with $\eta_k \leq \frac{\|\nabla f(x^k)\|}{2(L_0 + L_1 \|\nabla f(x^k)\|)}$:

$$1205 \mathbb{E} [f(x^{k+1})] - f(x^k) \stackrel{\text{equation 25}}{\leq} -\eta \langle \nabla f(x^k), \mathbb{E} [G(x^k, \xi^k)] \rangle + \frac{\eta^2(L_0 + L_1 \|\nabla f(x^k)\|)}{2} \mathbb{E} [\|G(x^k, \xi^k)\|^2]$$

$$1207 \stackrel{\text{equation 26}}{\leq} -\frac{\eta}{2} \|\nabla f(x^k)\| - \frac{\eta}{2\alpha} \mathbb{E} [\|G(x^k, \xi^k)\|^2] + \frac{\eta}{2\lambda^2\alpha} \mathbb{E} [\|\nabla f(x^k, \xi^k) - \nabla f(x^k)\|^2]$$

$$1209 + \frac{\eta^2(L_0 + L_1 \|\nabla f(x^k)\|)}{2} \mathbb{E} [\|G(x^k, \xi^k)\|^2]$$

$$1211 = -\frac{\eta}{2} \|\nabla f(x^k)\| + \frac{\eta}{2\lambda^2\alpha} \mathbb{E} [\|\nabla f(x^k, \xi^k) - \nabla f(x^k)\|^2]$$

$$1213 - \frac{\eta}{2} \mathbb{E} [\|G(x^k, \xi^k)\|^2] \left(1 - \frac{\eta(L_0 + L_1 \|\nabla f(x^k)\|)}{\|\nabla f(x^k)\|}\right)$$

$$1215 \leq -\frac{\eta}{2} \|\nabla f(x^k)\| + \frac{\eta\sigma^2}{2\lambda^2\alpha B}$$

$$1217 \leq -\frac{\eta}{2} \|\nabla f(x^k)\| + \frac{\eta\sigma^2 M}{2\lambda^2 B}. \quad (29)$$

1221 The step size will be constant, depending on the hyperparameter λ :

$$1222 \frac{\|\nabla f(x^k)\|}{2(L_0 + L_1 \|\nabla f(x^k)\|)} = \frac{1}{2\left(L_0 \frac{1}{\|\nabla f(x^k)\|} + L_1\right)} = \frac{\lambda}{2\left(L_0 \frac{\lambda}{\|\nabla f(x^k)\|} + L_1\lambda\right)} \geq \frac{\lambda}{2(L_0 + L_1\lambda)}.$$

1225 Thus, $\eta_k = \eta \leq \frac{\lambda}{2(L_0 + L_1\lambda)}$.

1228 Using the convexity assumption of the function, we have the following:

$$1230 f(x^k) - f^* \leq \langle \nabla f(x^k), x^k - x^* \rangle$$

$$1231 \stackrel{\text{equation 6}}{\leq} \|\nabla f(x^k)\| \|x^k - x^*\|$$

$$1232 \leq \|\nabla f(x^k)\| \underbrace{\max_{k \in [0, N-1]} \|x^k - x^*\|}_R.$$

1236 Hence we have:

$$1237 \|\nabla f(x^k)\| \geq \frac{f(x^k) - f^*}{R}. \quad (30)$$

1239 Then substituting equation 30 into equation 29 we obtain:

$$1240 \mathbb{E} [f(x^{k+1})] - f(x^k) \leq -\frac{\eta}{2} \|\nabla f(x^k)\| + \frac{\eta\sigma^2 M}{2\lambda^2 B} \leq -\frac{\eta}{2R} (f(x^k) - f^*) + \frac{\eta\sigma^2 M}{2\lambda^2 B}.$$

This inequality is equivalent to the trailing inequality:

$$\begin{aligned} \mathbb{E} [f(x^{k+1})] - f^* &\leq \left(1 - \frac{\eta}{2R}\right) (f(x^k) - f^*) + \frac{\eta\sigma^2 M}{2\lambda^2 B} \\ &\stackrel{\text{equation 11}}{\leq} \left(1 - \frac{\eta}{2R} + \frac{\eta\sigma^2(4L_1 + 1)}{4\lambda^2 B}\right) (f(x^k) - f^*) + \frac{\eta\sigma^2 L_0}{2\lambda^2 B} \\ &\stackrel{\textcircled{1}}{\leq} \left(1 - \frac{\eta}{4R}\right) (f(x^k) - f^*) + \frac{\eta\sigma^2 L_0}{2\lambda^2 B}, \end{aligned}$$

where in $\textcircled{1}$ we used $B \geq \max\left\{\frac{\sigma^2 R(4L_1 + 1)}{\lambda^2}, \frac{\sigma^2 L_0 R}{\lambda^2 \varepsilon}\right\}$.

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

$$\mathbb{E} [f(x^N)] - f^* \leq \left(1 - \frac{\eta}{4R}\right)^N (f(x^0) - f^*) + \frac{\sigma^2 L_0 R}{\lambda^2 B}.$$

C.2 SECOND CASE: $\|\nabla f(x^k)\| \leq \lambda$ AND $\|\nabla f(x^k, \xi^k)\| \geq \lambda$

Let us evaluate first summand of equation 25 with $\alpha = \lambda^{-1}$:

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

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

$$\begin{aligned} -\eta \langle \nabla f(x^k), \mathbb{E} [G(x^k, \xi^k)] \rangle &\leq -\frac{\eta}{2} \|\nabla f(x^k)\| - \frac{\eta}{2\alpha} \mathbb{E} [\|G(x^k, \xi^k)\|^2] \\ &\quad + \frac{\eta}{2\lambda} \mathbb{E} [\|\nabla f(x^k, \xi^k) - \nabla f(x^k)\|^2]. \end{aligned} \quad (31)$$

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

$$\begin{aligned} \mathbb{E} [f(x^{k+1})] - f(x^k) &\stackrel{\text{equation 25}}{\leq} -\eta \langle \nabla f(x^k), \mathbb{E} [G(x^k, \xi^k)] \rangle + \frac{\eta^2(L_0 + L_1 \|\nabla f(x^k)\|)}{2} \mathbb{E} [\|G(x^k, \xi^k)\|^2] \\ &\stackrel{\text{equation 31}}{\leq} -\frac{\eta}{2} \|\nabla f(x^k)\| - \frac{\eta}{2\alpha} \mathbb{E} [\|G(x^k, \xi^k)\|^2] + \frac{\eta}{2\lambda} \mathbb{E} [\|\nabla f(x^k, \xi^k) - \nabla f(x^k)\|^2] \\ &\quad + \frac{\eta^2(L_0 + L_1 \|\nabla f(x^k)\|)}{2} \mathbb{E} [\|G(x^k, \xi^k)\|^2] \\ &= -\frac{\eta}{2} \|\nabla f(x^k)\| + \frac{\eta}{2\lambda} \mathbb{E} [\|\nabla f(x^k, \xi^k) - \nabla f(x^k)\|^2] \\ &\quad - \frac{\eta}{2\alpha} \mathbb{E} [\|G(x^k, \xi^k)\|^2] \left(1 - \frac{\eta(L_0 + L_1 \|\nabla f(x^k)\|)}{\|\nabla f(x^k)\|}\right) \\ &\leq -\frac{\eta}{2} \|\nabla f(x^k)\| + \frac{\eta\sigma^2}{2\lambda B} \end{aligned}$$

$$\leq -\frac{\eta}{2} \|\nabla f(x^k)\| + \frac{\eta\sigma^2}{2\lambda B}. \quad (32)$$

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

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

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

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

$$\begin{aligned} f(x^k) - f^* &\leq \langle \nabla f(x^k), x^k - x^* \rangle \\ &\stackrel{\text{equation 6}}{\leq} \|\nabla f(x^k)\| \|x^k - x^*\| \\ &\leq \|\nabla f(x^k)\| \underbrace{\max_{k \in [0, N-1]} \|x^k - x^*\|}_R. \end{aligned}$$

Hence we have:

$$\|\nabla f(x^k)\| \geq \frac{f(x^k) - f^*}{R}. \quad (33)$$

Then substituting equation 33 into equation 32 we obtain:

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

This inequality is equivalent to the trailing inequality:

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

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

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

C.3 THIRD CASE: $\|\nabla f(x^k)\| \leq \lambda$ AND $\|\nabla f(x^k, \xi^k)\| \leq \lambda$

Using this in equation 25, we have the following with $\eta_k \leq \frac{\|\nabla f(x^k)\|}{2(L_0 + L_1 \|\nabla f(x^k)\|)}$ and $\alpha = \|\nabla f(x^k)\|^{-1}$:

$$\begin{aligned} \mathbb{E}[f(x^{k+1})] - f(x^k) &\stackrel{\text{equation 25}}{\leq} -\eta \langle \nabla f(x^k), \mathbb{E}[G(x^k, \xi^k)] \rangle \\ &\quad + \frac{\eta^2(L_0 + L_1 \|\nabla f(x^k)\|)}{2} \mathbb{E}[\|G(x^k, \xi^k)\|^2] \\ &= -\frac{\eta\alpha}{2} \|\nabla f(x^k)\|^2 - \frac{\eta}{2\alpha} \mathbb{E}[\|G(x^k, \xi^k)\|^2] \\ &\quad + \frac{\eta}{2\alpha} \mathbb{E}[\|G(x^k, \xi^k) - \alpha \nabla f(x^k)\|^2] \\ &\quad + \frac{\eta^2(L_0 + L_1 \|\nabla f(x^k)\|)}{2} \mathbb{E}[\|G(x^k, \xi^k)\|^2] \\ &= -\frac{\eta}{2} \|\nabla f(x^k)\| + \frac{\eta}{2\alpha} \mathbb{E}[\|G(x^k, \xi^k) - \alpha \nabla f(x^k)\|^2] \end{aligned}$$

$$\begin{aligned}
& -\frac{\eta}{2}\mathbb{E}\left[\|G(x^k, \xi^k)\|^2\right]\left(1 - \frac{\eta(L_0 + L_1\|\nabla f(x^k)\|)}{\|\nabla f(x^k)\|}\right) \\
& \leq -\frac{\eta}{2}\|\nabla f(x^k)\| + \frac{\eta}{2\alpha}\mathbb{E}\left[\|G(x^k, \xi^k) - \alpha\nabla f(x^k)\|^2\right] \\
& \leq -\frac{\eta}{2}\|\nabla f(x^k)\| + \frac{\eta}{\alpha}\mathbb{E}\left[\|G(x^k, \xi^k)\|^2 + \|\alpha\nabla f(x^k)\|^2\right] \\
& = -\frac{\eta}{2}\|\nabla f(x^k)\| + \frac{\eta}{\alpha}\mathbb{E}\left[\left\|\frac{\nabla f(x^k, \xi^k)}{\|\nabla f(x^k, \xi^k)\|}\right\|^2 + \left\|\frac{\nabla f(x^k)}{\|\nabla f(x^k)\|}\right\|^2\right] \\
& = -\frac{\eta}{2}\|\nabla f(x^k)\| + \frac{2\eta\lambda\|\nabla f(x^k)\|}{\lambda} \\
& \leq -\frac{\eta}{2}\|\nabla f(x^k)\| + 2\eta\lambda. \tag{34}
\end{aligned}$$

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

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

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

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

$$\begin{aligned}
f(x^k) - f^* & \leq \langle \nabla f(x^k), x^k - x^* \rangle \\
& \stackrel{\text{equation 6}}{\leq} \|\nabla f(x^k)\| \|x^k - x^*\| \\
& \leq \|\nabla f(x^k)\| \underbrace{\max_{k \in [0, N-1]} \|x^k - x^*\|}_R.
\end{aligned}$$

Hence we have:

$$\|\nabla f(x^k)\| \geq \frac{f(x^k) - f^*}{R}. \tag{35}$$

Then substituting equation 35 into equation 34 we obtain:

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

This inequality is equivalent to the trailing inequality:

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

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

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

C.4 FOURTH CASE: $\|\nabla f(x^k)\| \geq \lambda$ AND $\|\nabla f(x^k, \xi^k)\| \leq \lambda$

Using this in equation 25, we have the following with $\eta_k \leq \frac{\|\nabla f(x^k)\|}{2(L_0 + L_1\|\nabla f(x^k)\|)}$ and $\alpha = \lambda^{-1}$:

$$\begin{aligned}
\mathbb{E}[f(x^{k+1})] - f(x^k) & \stackrel{\text{equation 25}}{\leq} -\eta \langle \nabla f(x^k), \mathbb{E}[G(x^k, \xi^k)] \rangle \\
& \quad + \frac{\eta^2(L_0 + L_1\|\nabla f(x^k)\|)}{2} \mathbb{E}\left[\|G(x^k, \xi^k)\|^2\right]
\end{aligned}$$

$$\begin{aligned}
&= -\frac{\eta\alpha}{2} \|\nabla f(x^k)\|^2 - \frac{\eta}{2\alpha} \|\mathbb{E}[G(x^k, \boldsymbol{\xi}^k)]\|^2 \\
&\quad + \frac{\eta}{2\alpha} \|\mathbb{E}[G(x^k, \boldsymbol{\xi}^k)] - \alpha\nabla f(x^k)\|^2 \\
&\quad + \frac{\eta^2(L_0 + L_1 \|\nabla f(x^k)\|)}{2} \mathbb{E}[\|G(x^k, \boldsymbol{\xi}^k)\|^2] \\
&= -\frac{\eta}{2\lambda} \|\nabla f(x^k)\|^2 + \frac{\eta}{2\lambda} \|\mathbb{E}[\lambda G(x^k, \boldsymbol{\xi}^k)] - \nabla f(x^k)\|^2 \\
&\quad + \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} \left\| \mathbb{E} \left[\frac{\lambda \nabla f(x^k, \boldsymbol{\xi}^k)}{\|\nabla f(x^k, \boldsymbol{\xi}^k)\|} - \nabla f(x^k, \boldsymbol{\xi}^k) \right] \right\|^2 \\
&\quad + \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} \left\| \mathbb{E} \left[\left(\frac{\lambda}{\|\nabla f(x^k, \boldsymbol{\xi}^k)\|} - 1 \right) \nabla f(x^k, \boldsymbol{\xi}^k) \right] \right\|^2 \\
&\quad + \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, \boldsymbol{\xi}^k)\|} - 1 \right)^2 \|\nabla f(x^k, \boldsymbol{\xi}^k)\|^2 \right] \\
&\quad + \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}{\|\nabla f(x^k, \boldsymbol{\xi}^k)\|^2} \|\nabla f(x^k, \boldsymbol{\xi}^k)\|^2 \right] \\
&\quad + \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\lambda}{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}{2} \|\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} \\
&\leq -\frac{\eta}{4} \|\nabla f(x^k)\| + \frac{\eta\lambda}{2}. \tag{36}
\end{aligned}$$

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

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

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

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

$$\begin{aligned}
f(x^k) - f^* &\leq \langle \nabla f(x^k), x^k - x^* \rangle \\
&\stackrel{\text{equation 6}}{\leq} \|\nabla f(x^k)\| \|x^k - x^*\| \\
&\leq \|\nabla f(x^k)\| \underbrace{\max_{k \in [0, N-1]} \|x^k - x^*\|}_R.
\end{aligned}$$

Hence we have:

$$\|\nabla f(x^k)\| \geq \frac{f(x^k) - f^*}{R}. \quad (37)$$

Then substituting equation 37 into equation 36 we obtain:

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

This inequality is equivalent to the trailing inequality:

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

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

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

Combining all the cases considered, we obtain the convergence rate Normalized Stochastic Gradient Descent with batch size $B \geq \max \left\{ \frac{\sigma^2 R(4L_1+1)}{\lambda^2}, \frac{\sigma^2 L_0 R}{\lambda^2 \varepsilon} \right\}$:

$$\mathbb{E} [f(x^N)] - f^* \lesssim \left(1 - \frac{\eta}{R}\right)^N (f(x^0) - f^*) + \frac{\sigma^2 L_0 R}{B\lambda^2} + \lambda R.$$

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 $\mathbf{g}(x^k, \xi^k) = \nabla f(x^k, \xi^k) + \mathbf{b}(x^k)$, where $\mathbf{b}(x^k)$ is biased bounded by $\zeta \geq 0 : \|\mathbf{b}(x^k)\| \leq \zeta$; 2) deriving convergence estimates of ZO-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):

$$\begin{aligned} f(x^{k+1}) - f(x^k) &\stackrel{\text{equation 8}}{\leq} \langle \nabla f(x^k), x^{k+1} - x^k \rangle + \frac{L_0 + L_1 \|\nabla f(x^k)\|}{2} \|x^{k+1} - x^k\|^2 \\ &= -\eta \langle \nabla f(x^k), \text{clip}_c(\mathbf{g}(x^k, \xi^k)) \rangle \\ &\quad + \frac{\eta^2 (L_0 + L_1 \|\nabla f(x^k)\|)}{2} \|\text{clip}_c(\mathbf{g}(x^k, \xi^k))\|^2. \end{aligned} \quad (38)$$

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

D.1.1 FIRST CASE: $\|\nabla f(x^k)\| \geq c$

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

$$\begin{aligned} -\eta \langle \nabla f(x^k), \text{clip}_c(\mathbf{g}(x^k, \xi^k)) \rangle &\stackrel{\text{equation 5}}{=} -\frac{\alpha\eta}{2} \|\nabla f(x^k)\|^2 - \frac{\eta}{2\alpha} \|\text{clip}_c(\mathbf{g}(x^k, \xi^k))\|^2 \\ &\quad + \frac{\eta}{2\alpha} \|\text{clip}_c(\mathbf{g}(x^k, \xi^k)) - \alpha \nabla f(x^k)\|^2 \\ &= -\frac{\alpha\eta}{2} \|\nabla f(x^k)\|^2 - \frac{\eta}{2\alpha} \|\text{clip}_c(\mathbf{g}(x^k, \xi^k))\|^2 \end{aligned}$$

$$\begin{aligned}
& + \frac{\eta}{2\alpha} \|\text{clip}_c(\mathbf{g}(x^k, \boldsymbol{\xi}^k)) - \text{clip}_c(\nabla f(x^k))\|^2 \\
& = -\frac{c\eta}{2} \|\nabla f(x^k)\| - \frac{\eta}{2\alpha} \|\text{clip}_c(\mathbf{g}(x^k, \boldsymbol{\xi}^k))\|^2 \\
& \quad + \frac{\eta}{2\alpha} \|\text{clip}_c(\mathbf{g}(x^k, \boldsymbol{\xi}^k)) - \text{clip}_c(\nabla f(x^k))\|^2.
\end{aligned}$$

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

$$\begin{aligned}
& -\eta \langle \nabla f(x^k), \mathbb{E}[\text{clip}_c(\mathbf{g}(x^k, \boldsymbol{\xi}^k))] \rangle \leq -\frac{c\eta}{2} \|\nabla f(x^k)\| - \frac{\eta}{2\alpha} \mathbb{E}[\|\text{clip}_c(\mathbf{g}(x^k, \boldsymbol{\xi}^k))\|^2] \\
& \quad + \frac{\eta}{2\alpha} \mathbb{E}[\|\mathbf{g}(x^k, \boldsymbol{\xi}^k) - \nabla f(x^k)\|^2] \\
& \stackrel{\text{equation 9}}{=} -\frac{c\eta}{2} \|\nabla f(x^k)\| - \frac{\eta}{2\alpha} \mathbb{E}[\|\text{clip}_c(\mathbf{g}(x^k, \boldsymbol{\xi}^k))\|^2] \\
& \quad + \frac{\eta}{2\alpha} \mathbb{E}[\|\mathbf{g}(x^k, \boldsymbol{\xi}^k) - \mathbb{E}[\mathbf{g}(x^k, \boldsymbol{\xi}^k)]\|^2] \\
& \quad + \frac{\eta}{2\alpha} \|\mathbf{b}(x^k)\|^2 \\
& \leq -\frac{c\eta}{2} \|\nabla f(x^k)\| - \frac{\eta}{2\alpha} \mathbb{E}[\|\text{clip}_c(\nabla f(x^k, \boldsymbol{\xi}^k))\|^2] \\
& \quad + \frac{\eta\sigma^2 M}{2cB} + \frac{\eta \|\nabla f(x^k)\| \zeta^2}{2c}. \tag{39}
\end{aligned}$$

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

In the case $c \geq \sqrt{2}\zeta$ We have in equation 39:

$$\begin{aligned}
& -\eta \langle \nabla f(x^k), \mathbb{E}[\text{clip}_c(\mathbf{g}(x^k, \boldsymbol{\xi}^k))] \rangle \stackrel{\text{equation 39}}{\leq} -\frac{c\eta}{2} \|\nabla f(x^k)\| - \frac{\eta}{2\alpha} \mathbb{E}[\|\text{clip}_c(\mathbf{g}(x^k, \boldsymbol{\xi}^k))\|^2] \\
& \quad + \frac{\eta\sigma^2 M}{2cB} + \frac{\eta \|\nabla f(x^k)\| \zeta^2}{2c} \\
& = -\frac{\eta}{2\alpha} \mathbb{E}[\|\text{clip}_c(\mathbf{g}(x^k, \boldsymbol{\xi}^k))\|^2] \\
& \quad - \frac{c\eta}{2} \|\nabla f(x^k)\| \left(1 - \frac{\zeta^2}{c^2}\right) + \frac{\eta\sigma^2 M}{2cB} \\
& \leq -\frac{\eta}{2\alpha} \mathbb{E}[\|\text{clip}_c(\mathbf{g}(x^k, \boldsymbol{\xi}^k))\|^2] - \frac{c\eta}{4} \|\nabla f(x^k)\| + \frac{\eta\sigma^2 M}{2cB} \\
& = -\frac{\eta \|\nabla f(x^k)\|}{2c} \mathbb{E}[\|\text{clip}_c(\nabla f(x^k, \boldsymbol{\xi}^k))\|^2] \\
& \quad - \frac{c\eta}{4} \|\nabla f(x^k)\| + \frac{\eta\sigma^2 M}{2cB}.
\end{aligned}$$

Plugging this into equation 38 and choosing $\eta \leq \frac{1}{4(L_0 + L_1 c)}$ we have:

$$\begin{aligned}
& \mathbb{E}[\nabla f(x^{k+1})] - f(x^k) \stackrel{\text{equation 13}}{\leq} -\frac{\eta \|\nabla f(x^k)\|}{2c} \mathbb{E}[\|\text{clip}_c(\mathbf{g}(x^k, \boldsymbol{\xi}^k))\|^2] - \frac{c\eta}{4} \|\nabla f(x^k)\| + \frac{\eta\sigma^2 M}{2cB} \\
& \quad + \frac{\eta^2(L_0 + L_1 \|\nabla f(x^k)\|)}{2} \mathbb{E}[\|\text{clip}_c(\mathbf{g}(x^k, \boldsymbol{\xi}^k))\|^2] \\
& = -\frac{\eta \|\nabla f(x^k)\|}{2c} \mathbb{E}[\|\text{clip}_c(\mathbf{g}(x^k, \boldsymbol{\xi}^k))\|^2] (1 - \eta L_1 c) - \frac{c\eta}{4} \|\nabla f(x^k)\| \\
& \quad + \frac{\eta^2 L_0}{2} \mathbb{E}[\|\text{clip}_c(\mathbf{g}(x^k, \boldsymbol{\xi}^k))\|^2] + \frac{\eta\sigma^2 M}{2cB} \\
& \leq -\frac{c\eta}{4} \|\nabla f(x^k)\| - \frac{\eta}{2} \mathbb{E}[\|\text{clip}_c(\mathbf{g}(x^k, \boldsymbol{\xi}^k))\|^2] (1 - \eta(L_0 + L_1 c)) \\
& \quad + \frac{\eta\sigma^2 M}{2cB}
\end{aligned}$$

$$\leq -\frac{c\eta}{4} \|\nabla f(x^k)\| + \frac{\eta\sigma^2 M}{2cB}. \quad (40)$$

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

$$f(x^k) - f^* \leq \langle \nabla f(x^k), x^k - x^* \rangle \stackrel{\text{equation 6}}{\leq} \|\nabla f(x^k)\| \|x^k - x^*\| \leq \|\nabla f(x^k)\| \underbrace{\max_{k \in [0, N-1]} \|x^k - x^*\|}_R.$$

Hence we have:

$$\|\nabla f(x^k)\| \geq \frac{f(x^k) - f^*}{R}. \quad (41)$$

Then substituting equation 41 into equation 40 we obtain:

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

This inequality is equivalent to the trailing inequality:

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

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

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

In the case $c \leq \sqrt{2}\zeta$ We have in equation 39:

$$\begin{aligned} -\eta \langle \nabla f(x^k), \mathbb{E}[\text{clip}_c(\mathbf{g}(x^k, \boldsymbol{\xi}^k))] \rangle &\stackrel{\text{equation 39}}{\leq} -\frac{c\eta}{2} \|\nabla f(x^k)\| - \frac{\eta}{2\alpha} \mathbb{E}[\|\text{clip}_c(\mathbf{g}(x^k, \boldsymbol{\xi}^k))\|^2] \\ &\quad + \frac{\eta\sigma^2 M}{2cB} + \frac{\eta \|\nabla f(x^k)\| \zeta^2}{2c} \\ &= -\frac{c\eta}{2} \|\nabla f(x^k)\| - \frac{\eta}{2\alpha} \mathbb{E}[\|\text{clip}_c(\mathbf{g}(x^k, \boldsymbol{\xi}^k))\|^2] \\ &\quad + \frac{\eta M}{2c} \left(\frac{\sigma^2}{B} + \zeta^2\right). \end{aligned}$$

Plugging this into equation 38 and choosing $\eta \leq \frac{1}{4(L_0 + L_1 c)}$ we have:

$$\begin{aligned} \mathbb{E}[f(x^{k+1})] - f(x^k) &\stackrel{\text{equation 38}}{\leq} -\frac{c\eta}{2} \|\nabla f(x^k)\| - \frac{\eta \|\nabla f(x^k)\|}{2c} \mathbb{E}[\|\text{clip}_c(\mathbf{g}(x^k, \boldsymbol{\xi}^k))\|^2] \\ &\quad + \frac{\eta^2(L_0 + L_1 \|\nabla f(x^k)\|)}{2} \mathbb{E}[\|\text{clip}_c(\mathbf{g}(x^k, \boldsymbol{\xi}^k))\|^2] + \frac{\eta M}{2c} \left(\frac{\sigma^2}{B} + \zeta^2\right) \\ &= -\frac{c\eta}{2} \|\nabla f(x^k)\| - \frac{\eta \|\nabla f(x^k)\|}{2c} \mathbb{E}[\|\text{clip}_c(\mathbf{g}(x^k, \boldsymbol{\xi}^k))\|^2] (1 - \eta L_1 c) \\ &\quad + \frac{\eta^2 L_0}{2} \mathbb{E}[\|\text{clip}_c(\mathbf{g}(x^k, \boldsymbol{\xi}^k))\|^2] + \frac{\eta M}{2c} \left(\frac{\sigma^2}{B} + \zeta^2\right) \\ &\leq -\frac{c\eta}{2} \|\nabla f(x^k)\| - \frac{\eta}{2} \mathbb{E}[\|\text{clip}_c(\mathbf{g}(x^k, \boldsymbol{\xi}^k))\|^2] (1 - \eta(L_0 + L_1 c)) \\ &\quad + \frac{\eta M}{2c} \left(\frac{\sigma^2}{B} + \zeta^2\right) \\ &\leq -\frac{c\eta}{2} \|\nabla f(x^k)\| + \frac{\eta M}{2c} \left(\frac{\sigma^2}{B} + \zeta^2\right). \quad (42) \end{aligned}$$

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

$$f(x^k) - f^* \leq \langle \nabla f(x^k), x^k - x^* \rangle \stackrel{\text{equation 6}}{\leq} \|\nabla f(x^k)\| \|x^k - x^*\| \leq \|\nabla f(x^k)\| \underbrace{\max_{k \in [0, N-1]} \|x^k - x^*\|}_R.$$

Hence we have:

$$\|\nabla f(x^k)\| \geq \frac{f(x^k) - f^*}{R}. \quad (43)$$

Then substituting equation 43 into equation 42 we obtain:

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

This inequality is equivalent to the trailing inequality:

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

Then for $k = 0, 1, 2, \dots, N-1$ iterations that satisfy the conditions $\|\nabla f(x^k)\| \geq c$ and $c \leq \sqrt{2}\zeta$, then ClipSGD

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

D.1.2 SECOND CASE: $\frac{c}{3} \leq \|\nabla f(x^k)\| \leq c$

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

$$\begin{aligned} -\eta \langle \nabla f(x^k), \text{clip}_c(\mathbf{g}(x^k, \boldsymbol{\xi}^k)) \rangle &\stackrel{\text{equation 5}}{=} -\frac{\alpha\eta}{2} \|\nabla f(x^k)\|^2 - \frac{\eta}{2\alpha} \|\text{clip}_c(\mathbf{g}(x^k, \boldsymbol{\xi}^k))\|^2 \\ &\quad + \frac{\eta}{2\alpha} \|\text{clip}_c(\mathbf{g}(x^k, \boldsymbol{\xi}^k)) - \alpha \nabla f(x^k)\|^2 \\ &= -\frac{\eta}{2} \|\nabla f(x^k)\|^2 - \frac{\eta}{2} \|\text{clip}_c(\mathbf{g}(x^k, \boldsymbol{\xi}^k))\|^2 \\ &\quad + \frac{\eta}{2} \|\text{clip}_c(\mathbf{g}(x^k, \boldsymbol{\xi}^k)) - \text{clip}_c(\nabla f(x^k))\|^2 \\ &\leq -\frac{c\eta}{6} \|\nabla f(x^k)\| - \frac{\eta}{2} \|\text{clip}_c(\mathbf{g}(x^k, \boldsymbol{\xi}^k))\|^2 \\ &\quad + \frac{\eta}{2} \|\text{clip}_c(\mathbf{g}(x^k, \boldsymbol{\xi}^k)) - \text{clip}_c(\nabla f(x^k))\|^2. \end{aligned}$$

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

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

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

Plugging this into equation 38 and choosing $\eta \leq \frac{1}{4(L_0 + L_1 c)}$ we have:

$$\begin{aligned}
\mathbb{E} [f(x^{k+1})] - f(x^k) &\stackrel{\text{equation 38}}{\leq} -\frac{c\eta}{6} \|\nabla f(x^k)\| - \frac{\eta}{2} \mathbb{E} \left[\|\text{clip}_c(\mathbf{g}(x^k, \boldsymbol{\xi}^k))\|^2 \right] \\
&\quad + \frac{\eta^2(L_0 + L_1 \|\nabla f(x^k)\|)}{2} \mathbb{E} \left[\|\text{clip}_c(\mathbf{g}(x^k, \boldsymbol{\xi}^k))\|^2 \right] + \frac{\eta}{2} \left(\frac{\sigma^2}{B} + \zeta^2 \right) \\
&= -\frac{c\eta}{6} \|\nabla f(x^k)\| - \frac{\eta}{2} \mathbb{E} \left[\|\text{clip}_c(\mathbf{g}(x^k, \boldsymbol{\xi}^k))\|^2 \right] (1 - \eta(L_0 + L_1 \|\nabla f(x^k)\|)) \\
&\quad + \frac{\eta}{2} \left(\frac{\sigma^2}{B} + \zeta^2 \right) \\
&\leq -\frac{c\eta}{6} \|\nabla f(x^k)\| + \frac{\eta}{2} \left(\frac{\sigma^2}{B} + \zeta^2 \right). \tag{44}
\end{aligned}$$

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

$$f(x^k) - f^* \leq \langle \nabla f(x^k), x^k - x^* \rangle \stackrel{\text{equation 6}}{\leq} \|\nabla f(x^k)\| \|x^k - x^*\| \leq \|\nabla f(x^k)\| \underbrace{\max_{k \in [0, N-1]} \|x^k - x^*\|}_R.$$

Hence we have:

$$\|\nabla f(x^k)\| \geq \frac{f(x^k) - f^*}{R}. \tag{45}$$

Then substituting equation 45 into equation 44 we obtain:

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

This inequality is equivalent to the trailing inequality:

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

Then for $k = 0, 1, 2, \dots, N-1$ 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

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

D.1.3 THIRD CASE: $\|\nabla f(x^k)\| \leq \frac{c}{3}$

We introduce an indicative function:

$$\aleph_k = \mathbf{1} \{ \|\mathbf{g}(x^k, \boldsymbol{\xi}^k)\| > c \}. \tag{46}$$

Then the following is true:

$$\mathbb{E} [\aleph_k] = \mathbb{E} [\aleph_k^2] = \mathcal{P} \left[\|\mathbf{g}(x^k, \boldsymbol{\xi}^k)\| > c \right] \stackrel{\textcircled{1}}{\leq} \mathcal{P} \left[\|\mathbf{g}(x^k, \boldsymbol{\xi}^k) - \mathbb{E} [\mathbf{g}(x^k, \boldsymbol{\xi}^k)]\| > \frac{c}{3} \right] \stackrel{\textcircled{2}}{\leq} \frac{9\sigma^2}{c^2 B}, \tag{47}$$

where in $\textcircled{1}$ we used $\|\mathbf{g}(x^k, \boldsymbol{\xi}^k)\| \leq \|\mathbf{g}(x^k, \boldsymbol{\xi}^k) - \mathbb{E} [\mathbf{g}(x^k, \boldsymbol{\xi}^k)]\| + \|\mathbb{E} [\mathbf{g}(x^k, \boldsymbol{\xi}^k)]\| \leq \|\mathbf{g}(x^k, \boldsymbol{\xi}^k) - \mathbb{E} [\mathbf{g}(x^k, \boldsymbol{\xi}^k)]\| + \frac{c}{2}$, where assume that $\zeta \leq \frac{2c}{3}$: and in $\textcircled{2}$ we used Markov's inequality.

Let $r_{k+1} = \mathbb{E} [\|x^{k+1} - x^*\|]$ and $F_{k+1} = \mathbb{E} [f(x^{k+1}) - f^*]$, then given that

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

we get with $\eta \leq \frac{1}{4(L_0 + L_1 c)}$:

$$\begin{aligned} r_{k+1}^2 &= r_k^2 - 2\eta \langle \mathbb{E} [\text{clip}_c(\mathbf{g}(x^k, \boldsymbol{\xi}^k))], x^k - x^* \rangle + \eta^2 \mathbb{E} [\|\text{clip}_c(\mathbf{g}(x^k, \boldsymbol{\xi}^k))\|^2] \\ &= r_k^2 - 2\eta \langle \nabla f(x^k), x^k - x^* \rangle - 2\eta \left\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^* \right\rangle \\ &\quad + \eta^2 \mathbb{E} [\|\text{clip}_c(\mathbf{g}(x^k, \boldsymbol{\xi}^k))\|^2] + 2\eta \langle \mathbf{b}(x^k), x^k - x^* \rangle \\ &\stackrel{\text{equation 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^*\| \\ &\quad + \eta^2 \mathbb{E} [\|\text{clip}_c(\mathbf{g}(x^k, \boldsymbol{\xi}^k))\|^2] + 2\eta \|\mathbf{b}(x^k)\| \|x^k - x^*\| \\ &\stackrel{\textcircled{1}}{\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\| \|x^0 - x^*\| \\ &\quad + \eta^2 \mathbb{E} [\|\text{clip}_c(\mathbf{g}(x^k, \boldsymbol{\xi}^k))\|^2] + 2\eta \|\mathbf{b}(x^k)\| \|x^0 - x^*\| \\ &\stackrel{\text{equation 7}}{\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\| \|x^0 - x^*\| \\ &\quad + 2\eta^2 \mathbb{E} [\|\text{clip}_c(\mathbf{g}(x^k, \boldsymbol{\xi}^k)) - \nabla f(x^k)\|^2] + 2\eta^2 \|\nabla f(x^k)\|^2 + 2\eta \zeta \|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 \\ &\quad + 2\eta^2 \mathbb{E} [\|\text{clip}_c(\mathbf{g}(x^k, \boldsymbol{\xi}^k)) - \text{clip}_c(\nabla f(x^k))\|^2] + 2\eta^2 \|\nabla f(x^k)\|^2 + 2\eta \zeta R \\ &\stackrel{\textcircled{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 \\ &\quad + 2\eta^2 \mathbb{E} [\|\mathbf{g}(x^k, \boldsymbol{\xi}^k) - \nabla f(x^k)\|^2] + 2\eta^2 \|\nabla f(x^k)\|^2 + 2\eta \zeta R \\ &\stackrel{\text{equation 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 \\ &\quad + 2\eta^2 \mathbb{E} [\|\mathbf{g}(x^k, \boldsymbol{\xi}^k) - \mathbb{E} [\mathbf{g}(x^k, \boldsymbol{\xi}^k)]\|^2] + 2\eta^2 \|\nabla f(x^k)\|^2 + 2\eta \zeta R + 2\eta^2 \|\mathbf{b}(x^k)\|^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 \\ &\quad + \frac{2\eta^2 \sigma^2}{B} + 2\eta^2 \|\nabla f(x^k)\|^2 + 2\eta \zeta R + 2\eta^2 \zeta^2 \\ &\stackrel{\text{equation 10}}{\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 \\ &\quad + \frac{2\eta^2 \sigma^2}{B} + 4\eta^2 (L_0 + L_1 \|\nabla f(x^k)\|) 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 \\ &\quad + \frac{2\eta^2 \sigma^2}{B} + 4\eta^2 (L_0 + L_1 c) F_k + 2\eta \zeta R + 2\eta^2 \zeta^2 \\ &= r_k^2 - 2\eta F_k (1 - 2\eta (L_0 + L_1 c)) + \frac{2\eta^2 \sigma^2}{B} + 2\eta \zeta R + 2\eta^2 \zeta^2 \end{aligned}$$

$$\begin{aligned}
& + 2\eta \left\| \mathbb{E} \left[\left(\frac{c}{\|\mathbf{g}(x^k, \boldsymbol{\xi}^k)\|} - 1 \right) \mathbf{g}(x^k, \boldsymbol{\xi}^k) \aleph_k \right] \right\| R \\
& \leq r_k^2 - \eta F_k + \frac{2\eta^2 \sigma^2}{B} + 2\eta \zeta R + 2\eta^2 \zeta^2 \\
& + 2\eta \left\| \mathbb{E} \left[\left(\frac{c}{\|\mathbf{g}(x^k, \boldsymbol{\xi}^k)\|} - 1 \right) \mathbf{g}(x^k, \boldsymbol{\xi}^k) \aleph_k \right] \right\| R. \tag{48}
\end{aligned}$$

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

$$\begin{aligned}
& 2\eta R \left\| \mathbb{E} \left[\left(\frac{c}{\|\mathbf{g}(x^k, \boldsymbol{\xi}^k)\|} - 1 \right) \mathbf{g}(x^k, \boldsymbol{\xi}^k) \aleph_k \right] \right\| \\
& \stackrel{\text{equation 46}}{\leq} 2\eta R \mathbb{E} \left[\|\mathbf{g}(x^k, \boldsymbol{\xi}^k)\| \cdot \left(1 - \frac{c}{\|\mathbf{g}(x^k, \boldsymbol{\xi}^k)\|} \right) \aleph_k \right] \\
& \leq 2\eta R \mathbb{E} [\|\mathbf{g}(x^k, \boldsymbol{\xi}^k)\| \cdot \aleph_k] \\
& \leq 2\eta R (\mathbb{E} [\|\mathbf{g}(x^k, \boldsymbol{\xi}^k) - \mathbb{E}[\mathbf{g}(x^k, \boldsymbol{\xi}^k)]\|] \cdot \mathbb{E}[\aleph_k] + \|\nabla f(x^k)\| \mathbb{E}[\aleph_k] + \|\mathbf{b}(x^k)\| \mathbb{E}[\aleph_k]) \\
& \leq 2\eta R \left(\sqrt{\mathbb{E} [\|\mathbf{g}(x^k, \boldsymbol{\xi}^k) - \mathbb{E}[\mathbf{g}(x^k, \boldsymbol{\xi}^k)]\|^2]} \cdot \mathbb{E}[\aleph_k^2] + \frac{2c}{3} \mathbb{E}[\aleph_k] \right) \\
& \stackrel{\text{equation 47}}{\leq} 2\eta R \left(\frac{3\sigma^2}{cB} + \frac{2c}{3} \cdot \frac{9\sigma^2}{c^2 B} \right) \\
& = \frac{18\eta\sigma^2 R}{cB}. \tag{49}
\end{aligned}$$

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

$$\begin{aligned}
\eta F_k & \stackrel{\text{equation 48}}{\leq} r_k^2 - r_{k+1}^2 + \frac{2\eta^2 \sigma^2}{B} + 2\eta \zeta R + 2\eta^2 \zeta^2 + 2\eta \left\| \mathbb{E} \left[\left(\frac{c}{\|\nabla f(x^k, \boldsymbol{\xi}^k)\|} - 1 \right) \nabla f(x^k, \boldsymbol{\xi}^k) \aleph_k \right] \right\| R \\
& \stackrel{\text{equation 49}}{\leq} r_k^2 - r_{k+1}^2 + \frac{2\eta^2 \sigma^2}{B} + \frac{18\eta\sigma^2 R}{cB} + 2\eta \zeta R + 2\eta^2 \zeta^2
\end{aligned}$$

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

$$\begin{aligned}
\mathbb{E} [f(x^N)] - f^* & \leq F_N \cdot \mathbb{1}[\mathcal{T}_1] + F_N \cdot \mathbb{1}[\mathcal{T}_2] \\
& \lesssim \left(1 - \frac{\eta c}{R} \right)^K F_0 + \frac{R^2}{\eta(N-K)} + \left(\frac{MR}{c^2} + \frac{R}{c} + \eta \right) \cdot \left(\frac{\sigma^2}{B} + \zeta^2 \right) + R\zeta, \tag{50}
\end{aligned}$$

where \mathcal{T}_1 describes case $\|\nabla f(x^k)\| \geq \frac{c}{3}$, and \mathcal{T}_2 describes case $\|\nabla f(x^k)\| < \frac{c}{3}$.

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 equation 4.

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

$$\begin{aligned}
\|\mathbb{E}[\mathbf{g}(x, \{\xi, e\})] - \nabla f(x)\| & = \left\| \mathbb{E} \left[\frac{d}{2\gamma} \left(\tilde{f}(x + \gamma e, \xi) - \tilde{f}(x - \gamma e, \xi) \right) e \right] - \nabla f(x) \right\| \\
& \stackrel{\textcircled{1}}{=} \left\| \mathbb{E} \left[\frac{d}{\gamma} \left(f(x + \gamma e, \xi) + \delta(x + \gamma e) \right) e \right] - \nabla f(x) \right\| \\
& \stackrel{\textcircled{2}}{\leq} \left\| \mathbb{E} \left[\frac{d}{\gamma} f(x + \gamma e, \xi) e \right] - \nabla f(x) \right\| + \frac{d\Delta}{\gamma} \\
& \stackrel{\textcircled{3}}{=} \|\mathbb{E}[\nabla f(x + \gamma u, \xi)] - \nabla f(x)\| + \frac{d\Delta}{\gamma}
\end{aligned}$$

$$\begin{aligned}
1836 & \\
1837 & = \sup_{z \in S^d(1)} \mathbb{E} [\|\nabla_z f(x + \gamma u, \xi) - \nabla_z f(x)\|] + \frac{d\Delta}{\gamma} \\
1838 & \\
1839 & \stackrel{\text{equation 8}}{\leq} (L_0 + L_1 \|\nabla f(x^k)\|) \gamma \mathbb{E} [\|u\|] + \frac{d\Delta}{\gamma} \\
1840 & \\
1841 & \leq (L_0 + L_1 M) \gamma + \frac{d\Delta}{\gamma}, \tag{51} \\
1842 & \\
1843 &
\end{aligned}$$

1843 where $u \in B^d(1)$, ① = the equality is obtained from the fact, namely, distribution of e is symmetric,
1844 ② = the inequality is obtain from bounded noise $|\delta(x)| \leq \Delta$, ③ = the equality is obtained from a
1845 version of Stokes' theorem (see Section 13.3.5, Exercise 14a, Zorich & Paniagua, 2016).
1846

1847 **Bounding second moment (variance) of gradient approximation** By definition gradient approxi-
1848 mation equation 4 and Wirtinger-Poincare inequality equation 12 we have
1849

$$\begin{aligned}
1850 & \mathbb{E} [\|\mathbf{g}(x, \{\xi, e\}) - \mathbb{E}[\mathbf{g}(x, \{\xi, e\})]\|^2] \\
1851 & \leq \mathbb{E} [\|\mathbf{g}(x, \{\xi, e\})\|^2] \\
1852 & = \frac{d^2}{4\gamma^2} \mathbb{E} [\|\tilde{f}(x + \gamma e, \xi) - \tilde{f}(x - \gamma e, \xi)\|^2] \\
1853 & = \frac{d^2}{4\gamma^2} \mathbb{E} [(f(x + \gamma e, \xi) - f(x - \gamma e, \xi) + \delta(x + \gamma e) - \delta(x - \gamma e))^2] \\
1854 & \stackrel{\text{equation 7}}{\leq} \frac{d^2}{2\gamma^2} (\mathbb{E} [(f(x + \gamma e, \xi) - f(x - \gamma e, \xi))^2] + 2\Delta^2) \\
1855 & \stackrel{\text{equation 12}}{\leq} \frac{d^2}{2\gamma^2} \left(\frac{\gamma^2}{d} \mathbb{E} [\|\nabla f(x + \gamma e, \xi) + \nabla f(x - \gamma e, \xi)\|^2] + 2\Delta^2 \right) \\
1856 & = \frac{d^2}{2\gamma^2} \left(\frac{\gamma^2}{d} \mathbb{E} [\|\nabla f(x + \gamma e, \xi) + \nabla f(x - \gamma e, \xi) \pm 2\nabla f(x, \xi)\|^2] + 2\Delta^2 \right) \\
1857 & \stackrel{\text{equation 8}}{\leq} 4d\mathbb{E} [\|\nabla f(x, \xi)\|^2] + 4dL^2\gamma^2\mathbb{E} [\|e\|^2] + \frac{d^2\Delta^2}{\gamma^2} \\
1858 & \stackrel{\text{①}}{\leq} 4d\tilde{\sigma}^2 + 4d(L_0 + L_1 \|\nabla f(x^k)\|)^2 \gamma^2 \mathbb{E} [\|e\|^2] + \frac{d^2\Delta^2}{\gamma^2} \\
1859 & \leq 4d\tilde{\sigma}^2 + 4d(L_0 + L_1 M)^2 \gamma^2 + \frac{d^2\Delta^2}{\gamma^2}, \tag{52} \\
1860 & \\
1861 & \\
1862 & \\
1863 & \\
1864 & \\
1865 & \\
1866 & \\
1867 & \\
1868 & \\
1869 & \\
1870 & \\
1871 &
\end{aligned}$$

1872 where ① = the inequality is obtain from $\mathbb{E} [\|\nabla f(x, \xi)\|^2] \leq \tilde{\sigma}^2$.
1873

1874 D.2.1 PROOF OF THEOREM 5.2

1875
1876 In order to obtain the convergence rate of ZO-ClipSGD in the convex setting, we need to substitute
1877 the obtained estimates equation 51 and equation 52 into the convergence rate of ClipSGD equation 50
1878 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
1879 ZO-ClipSGD in the convex setup is as follows:

$$\begin{aligned}
1880 & \\
1881 & \mathbb{E} [f(x^N)] - f^* \lesssim \underbrace{\left(1 - \frac{\eta}{R}\right)^K (f(x^0) - f^*)}_{\text{①}} + \underbrace{\frac{R^2}{\eta(N-K)}}_{\text{②}} + \underbrace{\frac{dMR\tilde{\sigma}^2}{c^2B}}_{\text{③}} + \underbrace{\frac{dMR(L_0 + L_1M)^2\gamma^2}{c^2B}}_{\text{④}} \\
1882 & \\
1883 & \\
1884 & \quad + \underbrace{\frac{d^2MR\Delta^2}{c^2B\gamma^2}}_{\text{⑤}} + \underbrace{\frac{MR(L_0 + L_1M)^2\gamma^2}{c^2}}_{\text{⑥}} + \underbrace{\frac{d^2MR\Delta^2}{c^2\gamma^2}}_{\text{⑦}} \\
1885 & \\
1886 & \quad + \underbrace{(L_0 + L_1M)\gamma R}_{\text{⑧}} + \underbrace{\frac{d\Delta R}{\gamma}}_{\text{⑨}}. \\
1887 & \\
1888 & \\
1889 &
\end{aligned}$$

1890 **From term ①**, we find the K :

$$1891 \quad \textcircled{1} : \left(1 - \frac{\eta c}{R}\right)^K (f(x^0) - f^*) \leq \varepsilon \Rightarrow K \geq \frac{R}{\eta c} \log \frac{f(x^0) - f^*}{\varepsilon}. \quad (53)$$

1894 **From term ②**, we find the number of iterations N required for Algorithm 3 in convex setup to
1895 achieve ε -accuracy:

$$1896 \quad \textcircled{2} : \frac{R^2}{\eta(N-K)} \leq \varepsilon \Rightarrow N \stackrel{\text{equation 53}}{\geq} \frac{R^2}{\eta\varepsilon} + \frac{R}{\eta c} \log \frac{f(x^0) - f^*}{\varepsilon};$$

$$1899 \quad N = \mathcal{O}\left(\frac{R^2}{\eta\varepsilon} + \frac{R}{\eta c} \log \frac{1}{\varepsilon}\right). \quad (54)$$

1902 **From terms ③**, we find the batch size B :

$$1903 \quad \textcircled{3} : \frac{dMR\tilde{\sigma}^2}{c^2B} \leq \varepsilon \Rightarrow B \geq \frac{dMR\tilde{\sigma}^2}{\varepsilon c^2};$$

$$1906 \quad B = \mathcal{O}\left(\frac{dMR\tilde{\sigma}^2}{\varepsilon c^2}\right). \quad (55)$$

1908 **From terms ④, ⑥ and ⑧** we find the smoothing parameter γ :

$$1910 \quad \textcircled{4} : \frac{dMR(L_0 + L_1M)^2 \gamma^2}{c^2B} \leq \varepsilon \Rightarrow \gamma \leq \sqrt{\frac{\varepsilon c^2 B}{dMR(L_0 + L_1M)^2}} \stackrel{\text{equation 55}}{=} \frac{\tilde{\sigma}}{(L_0 + L_1M)};$$

$$1913 \quad \textcircled{6} : \frac{MR(L_0 + L_1M)^2 \gamma^2}{c^2} \leq \varepsilon \Rightarrow \gamma \leq \frac{\sqrt{\varepsilon c}}{\sqrt{MR}(L_0 + L_1M)};$$

$$1915 \quad \textcircled{8} : (L_0 + L_1M)R\gamma \leq \varepsilon \Rightarrow \gamma \leq \frac{\varepsilon}{R(L_0 + L_1M)};$$

$$1917 \quad \gamma \leq \frac{1}{(L_0 + L_1M)} \min\left\{\tilde{\sigma}, \frac{\sqrt{\varepsilon c}}{\sqrt{MR}}, \frac{\varepsilon}{R}\right\} = \frac{\varepsilon}{R(L_0 + L_1M)}. \quad (56)$$

1920 **From the remaining terms ⑤, ⑦ and ⑨**, we find the maximum allowable level of adversarial noise
1921 Δ that still guarantees the convergence of the ZO-ClipSGD to desired accuracy ε in convex setup:

$$1922 \quad \textcircled{5} : \frac{d^2MR\Delta^2}{c^2B\gamma^2} \leq \varepsilon \Rightarrow \Delta \leq \frac{\sqrt{\varepsilon c} \gamma \sqrt{B}}{d\sqrt{MR}} \stackrel{\text{equation 55, equation 56}}{=} \frac{\varepsilon \tilde{\sigma}}{\sqrt{d}(L_0 + L_1M)R};$$

$$1925 \quad \textcircled{7} : \frac{d^2MR\Delta^2}{\gamma^2 c^2} \leq \varepsilon \Rightarrow \Delta \leq \sqrt{\frac{\gamma^2 c^2 \varepsilon}{d^2 MR}} \stackrel{\text{equation 56}}{=} \frac{\varepsilon^{3/2} c}{d(L_0 + L_1M) \sqrt{MR}^{3/2}};$$

$$1928 \quad \textcircled{9} : \frac{d\Delta R}{\gamma} \leq \varepsilon \Rightarrow \Delta \leq \sqrt{\frac{\gamma \varepsilon}{dR}} \stackrel{\text{equation 56}}{=} \frac{\varepsilon^2}{d(L_0 + L_1M)R^2};$$

$$1930 \quad \Delta \leq \frac{\varepsilon}{\sqrt{d}(L_0 + L_1M)R} \min\left\{\tilde{\sigma}, \frac{\sqrt{\varepsilon c}}{\sqrt{d}\sqrt{MR}}, \frac{\varepsilon}{\sqrt{d}R}\right\}$$

$$1932 \quad = \frac{\varepsilon}{\sqrt{d}(L_0 + L_1M)R} \min\left\{\tilde{\sigma}, \frac{\varepsilon}{\sqrt{d}R}\right\}. \quad (57)$$

1935 In this way, the ZO-ClipSGD achieves ε -accuracy: $\mathbb{E}[f(x^N) - f^*] \leq \varepsilon$ in convex setup after

$$1936 \quad N \stackrel{\text{equation 54}}{=} \mathcal{O}\left(\frac{R^2}{\eta\varepsilon} + \frac{R}{\eta c} \log \frac{1}{\varepsilon}\right), \quad T = N \cdot B \stackrel{\text{equation 54, equation 55}}{=} \mathcal{O}\left(\frac{d\tilde{\sigma}^2 MR^2}{\varepsilon c^2 \eta} \left(\frac{1}{c} \log \frac{1}{\varepsilon} + \frac{R}{\varepsilon}\right)\right)$$

1939 number of iterations, total number of zero-order oracle calls and at

$$1940 \quad \Delta \stackrel{\text{equation 57}}{\lesssim} \frac{\varepsilon}{\sqrt{d}(L_0 + L_1M)R} \min\left\{\tilde{\sigma}, \frac{\varepsilon}{\sqrt{d}R}\right\}$$

1941 the maximum level of noise with smoothing parameter $\frac{\varepsilon}{(L_0 + L_1M)R}$ equation 56.

1944 E ZERO-ORDER NORMALIZED STOCHASTIC GRADIENT DESCENT METHOD

1945
1946 This section consists of two parts: 1) a generalization of the convergence result of NSGD (Algorithm 2)
1947 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
1948 $\zeta \geq 0 : \|\mathbf{b}(x^k)\| \leq \zeta$; 2) deriving convergence estimates of ZO-NSGD directly.
1949

1950 E.1 BIASED NORMALIZED STOCHASTIC GRADIENT DESCENT METHOD (PROOF OF THE

1951 LEMMA 5.3)

1952
1953 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 Assump-
1954 tion 1.2):

$$1955 \begin{aligned} 1956 f(x^{k+1}) - f(x^k) &\stackrel{\text{equation 8}}{\leq} \langle \nabla f(x^k), x^{k+1} - x^k \rangle + \frac{L_0 + L_1 \|\nabla f(x^k)\|}{2} \|x^{k+1} - x^k\|^2 \\ 1957 &= -\eta \langle \nabla f(x^k), G(x^k, \boldsymbol{\xi}^k) \rangle + \frac{\eta^2 (L_0 + L_1 \|\nabla f(x^k)\|)}{2} \|G(x^k, \boldsymbol{\xi}^k)\|^2. \end{aligned} \quad (58)$$

1960 Next, we consider 4 cases of the relation $\|\nabla f(x^k)\|$ and $\|\mathbf{g}(x^k, \boldsymbol{\xi}^k)\|$ with respect to the hyperpa-
1961 rameter λ .
1962

1963 E.1.1 FIRST CASE: $\|\nabla f(x^k)\| \geq \lambda$ AND $\|\mathbf{g}(x^k, \boldsymbol{\xi}^k)\| \geq \lambda$

1964
1965 Let us evaluate first summand of equation 58 with $\alpha = \|\nabla f(x^k)\|^{-1}$:

$$1966 \begin{aligned} 1967 -\eta \langle \nabla f(x^k), G(x^k, \boldsymbol{\xi}^k) \rangle &\stackrel{\text{equation 5}}{=} -\frac{\alpha\eta}{2} \|\nabla f(x^k)\|^2 - \frac{\eta}{2\alpha} \|G(x^k, \boldsymbol{\xi}^k)\|^2 \\ 1968 &\quad + \frac{\eta}{2\alpha} \|G(x^k, \boldsymbol{\xi}^k) - \alpha \nabla f(x^k)\|^2 \\ 1969 &= -\frac{\eta}{2} \|\nabla f(x^k)\| - \frac{\eta}{2\alpha} \|G(x^k, \boldsymbol{\xi}^k)\|^2 \\ 1970 &\quad + \frac{\eta}{2\lambda^2\alpha} \|\lambda G(x^k, \boldsymbol{\xi}^k) - \lambda \alpha \nabla f(x^k)\|^2 \\ 1971 &= -\frac{\eta}{2} \|\nabla f(x^k)\| - \frac{\eta}{2\alpha} \|G(x^k, \boldsymbol{\xi}^k)\|^2 \\ 1972 &\quad + \frac{\eta}{2\lambda^2\alpha} \|\text{clip}_\lambda(\mathbf{g}(x^k, \boldsymbol{\xi}^k)) - \text{clip}_\lambda(\nabla f(x^k))\|^2 \end{aligned}$$

1973
1974 Using that clipping is a projection on onto a convex set, namely ball with radius λ , and thus is Lipschitz
1975 operator with Lipschitz constant 1, we can obtain:
1976

$$1977 \begin{aligned} 1978 -\eta \langle \nabla f(x^k), \mathbb{E}[G(x^k, \boldsymbol{\xi}^k)] \rangle &\leq -\frac{\eta}{2} \|\nabla f(x^k)\| - \frac{\eta}{2\alpha} \mathbb{E}[\|G(x^k, \boldsymbol{\xi}^k)\|^2] \\ 1979 &\quad + \frac{\eta}{2\lambda^2\alpha} \mathbb{E}[\|\mathbf{g}(x^k, \boldsymbol{\xi}^k) - \nabla f(x^k)\|^2]. \end{aligned} \quad (59)$$

1980
1981 **In the case:** $0 \leq \zeta \leq \frac{\lambda}{\sqrt{2}}$. Using this in equation 59, we have the following with $\eta_k \leq$

$$1982 \frac{\|\nabla f(x^k)\|}{2(L_0 + L_1 \|\nabla f(x^k)\|)}:$$

$$1983 \begin{aligned} 1984 \mathbb{E}[f(x^{k+1})] - f(x^k) &\stackrel{\text{equation 58}}{\leq} -\eta \langle \nabla f(x^k), \mathbb{E}[G(x^k, \boldsymbol{\xi}^k)] \rangle + \frac{\eta^2 (L_0 + L_1 \|\nabla f(x^k)\|)}{2} \mathbb{E}[\|G(x^k, \boldsymbol{\xi}^k)\|^2] \\ 1985 &\stackrel{\text{equation 59}}{\leq} -\frac{\eta}{2} \|\nabla f(x^k)\| - \frac{\eta}{2\alpha} \mathbb{E}[\|G(x^k, \boldsymbol{\xi}^k)\|^2] + \frac{\eta}{2\lambda^2\alpha} \mathbb{E}[\|\mathbf{g}(x^k, \boldsymbol{\xi}^k) - \nabla f(x^k)\|^2] \\ 1986 &\quad + \frac{\eta^2 (L_0 + L_1 \|\nabla f(x^k)\|)}{2} \mathbb{E}[\|G(x^k, \boldsymbol{\xi}^k)\|^2] \\ 1987 &= -\frac{\eta}{2} \|\nabla f(x^k)\| + \frac{\eta}{2\lambda^2\alpha} \mathbb{E}[\|\mathbf{g}(x^k, \boldsymbol{\xi}^k) - \nabla f(x^k)\|^2] \\ 1988 &\quad - \frac{\eta}{2} \mathbb{E}[\|G(x^k, \boldsymbol{\xi}^k)\|^2] \left(1 - \frac{\eta(L_0 + L_1 \|\nabla f(x^k)\|)}{\|\nabla f(x^k)\|}\right) \end{aligned}$$

$$\begin{aligned}
&\leq -\frac{\eta}{2} \|\nabla f(x^k)\| + \frac{\eta}{2\lambda^2\alpha} \mathbb{E} \left[\|\mathbf{g}(x^k, \boldsymbol{\xi}^k) - \nabla f(x^k)\|^2 \right] \\
&\stackrel{\text{equation 9}}{=} -\frac{\eta}{2} \|\nabla f(x^k)\| + \frac{\eta}{2\lambda^2\alpha} \mathbb{E} \left[\|\mathbf{g}(x^k, \boldsymbol{\xi}^k) - \mathbb{E}[\mathbf{g}(x^k, \boldsymbol{\xi}^k)]\| \right] + \frac{\eta}{2\lambda^2\alpha} \|\mathbf{b}(x^k)\|^2 \\
&\leq -\frac{\eta}{2} \|\nabla f(x^k)\| + \frac{\eta\sigma^2}{2\lambda^2\alpha B} + \frac{\eta\zeta^2}{2\lambda^2\alpha} \\
&\leq -\frac{\eta}{2} \|\nabla f(x^k)\| + \frac{\eta}{4} \|\nabla f(x^k)\| + \frac{\eta\sigma^2 M}{2\lambda^2 B} \\
&= -\frac{\eta}{4} \|\nabla f(x^k)\| + \frac{\eta\sigma^2 M}{2\lambda^2 B}. \tag{60}
\end{aligned}$$

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

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

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

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

$$f(x^k) - f^* \leq \langle \nabla f(x^k), x^k - x^* \rangle \stackrel{\text{equation 6}}{\leq} \|\nabla f(x^k)\| \|x^k - x^*\| \leq \|\nabla f(x^k)\| \underbrace{\max_{k \in [0, N-1]} \|x^k - x^*\|}_R.$$

Hence we have:

$$\|\nabla f(x^k)\| \geq \frac{f(x^k) - f^*}{R}. \tag{61}$$

Then substituting equation 61 into equation 60 we obtain:

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

This inequality is equivalent to the trailing inequality:

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

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

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

In the case: $\frac{\lambda}{\sqrt{2}} \leq \zeta$. Using this in equation 59, we have the following with $\eta_k \leq$

$$\begin{aligned}
&\frac{\|\nabla f(x^k)\|}{2(L_0 + L_1 \|\nabla f(x^k)\|)} \\
&\mathbb{E}[f(x^{k+1})] - f(x^k) \stackrel{\text{equation 58}}{\leq} -\eta \langle \nabla f(x^k), \mathbb{E}[G(x^k, \boldsymbol{\xi}^k)] \rangle + \frac{\eta^2(L_0 + L_1 \|\nabla f(x^k)\|)}{2} \mathbb{E} \left[\|G(x^k, \boldsymbol{\xi}^k)\|^2 \right] \\
&\stackrel{\text{equation 59}}{\leq} -\frac{\eta}{2} \|\nabla f(x^k)\| - \frac{\eta}{2\alpha} \mathbb{E} \left[\|G(x^k, \boldsymbol{\xi}^k)\|^2 \right] + \frac{\eta}{2\lambda^2\alpha} \mathbb{E} \left[\|\mathbf{g}(x^k, \boldsymbol{\xi}^k) - \nabla f(x^k)\|^2 \right] \\
&\quad + \frac{\eta^2(L_0 + L_1 \|\nabla f(x^k)\|)}{2} \mathbb{E} \left[\|G(x^k, \boldsymbol{\xi}^k)\|^2 \right] \\
&= -\frac{\eta}{2} \|\nabla f(x^k)\| + \frac{\eta}{2\lambda^2\alpha} \mathbb{E} \left[\|\mathbf{g}(x^k, \boldsymbol{\xi}^k) - \nabla f(x^k)\|^2 \right] \\
&\quad - \frac{\eta}{2} \mathbb{E} \left[\|G(x^k, \boldsymbol{\xi}^k)\|^2 \right] \left(1 - \frac{\eta(L_0 + L_1 \|\nabla f(x^k)\|)}{\|\nabla f(x^k)\|} \right)
\end{aligned}$$

$$\begin{aligned}
&\leq -\frac{\eta}{2} \|\nabla f(x^k)\| + \frac{\eta}{2\lambda^2\alpha} \mathbb{E} \left[\|\mathbf{g}(x^k, \boldsymbol{\xi}^k) - \nabla f(x^k)\|^2 \right] \\
&\stackrel{\text{equation 9}}{=} -\frac{\eta}{2} \|\nabla f(x^k)\| + \frac{\eta}{2\lambda^2\alpha} \mathbb{E} \left[\|\mathbf{g}(x^k, \boldsymbol{\xi}^k) - \mathbb{E}[\mathbf{g}(x^k, \boldsymbol{\xi}^k)]\| \right] + \frac{\eta}{2\lambda^2\alpha} \|\mathbf{b}(x^k)\|^2 \\
&\leq -\frac{\eta}{2} \|\nabla f(x^k)\| + \frac{\eta\sigma^2}{2\lambda^2\alpha B} + \frac{\eta\zeta^2}{2\lambda^2\alpha} \\
&\leq -\frac{\eta}{2} \|\nabla f(x^k)\| + \frac{\eta\sigma^2 M}{2\lambda^2 B} + \frac{\eta\zeta^2 M}{2\lambda^2} \\
&= -\frac{\eta}{2} \|\nabla f(x^k)\| + \frac{\eta\sigma^2 M}{2\lambda^2 B} + \frac{\eta\zeta^2 M}{2\lambda^2}. \tag{62}
\end{aligned}$$

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

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

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

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

$$f(x^k) - f^* \leq \langle \nabla f(x^k), x^k - x^* \rangle \stackrel{\text{equation 6}}{\leq} \|\nabla f(x^k)\| \|x^k - x^*\| \leq \|\nabla f(x^k)\| \underbrace{\max_{k \in [0, N-1]} \|x^k - x^*\|}_R.$$

Hence we have:

$$\|\nabla f(x^k)\| \geq \frac{f(x^k) - f^*}{R}. \tag{63}$$

Then substituting equation 63 into equation 62 we obtain:

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

This inequality is equivalent to the trailing inequality:

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

Then for $k = 0, 1, 2, \dots, N-1$ 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:

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

E.1.2 SECOND CASE: $\|\nabla f(x^k)\| \leq \lambda$ AND $\|\mathbf{g}(x^k, \boldsymbol{\xi}^k)\| \geq \lambda$

Let us evaluate first summand of equation 58 with $\alpha = \lambda^{-1}$:

$$\begin{aligned}
&-\eta \langle \nabla f(x^k), G(x^k, \boldsymbol{\xi}^k) \rangle \stackrel{\text{equation 5}}{=} -\frac{\alpha\eta}{2} \|\nabla f(x^k)\|^2 - \frac{\eta}{2\alpha} \|G(x^k, \boldsymbol{\xi}^k)\|^2 \\
&\quad + \frac{\eta}{2\alpha} \|G(x^k, \boldsymbol{\xi}^k) - \alpha \nabla f(x^k)\|^2 \\
&\leq -\frac{\eta}{2} \|\nabla f(x^k)\| - \frac{\eta}{2\alpha} \|G(x^k, \boldsymbol{\xi}^k)\|^2 \\
&\quad + \frac{\eta}{2\lambda} \|\lambda G(x^k, \boldsymbol{\xi}^k) - \nabla f(x^k)\|^2 \\
&= -\frac{\eta}{2} \|\nabla f(x^k)\| - \frac{\eta}{2\alpha} \|G(x^k, \boldsymbol{\xi}^k)\|^2
\end{aligned}$$

$$+ \frac{\eta}{2\lambda} \|\text{clip}_\lambda(\mathbf{g}(x^k, \boldsymbol{\xi}^k)) - \text{clip}_\lambda(\nabla f(x^k))\|^2$$

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

$$\begin{aligned} -\eta \langle \nabla f(x^k), \mathbb{E}[G(x^k, \boldsymbol{\xi}^k)] \rangle &\leq -\frac{\eta}{2} \|\nabla f(x^k)\| - \frac{\eta}{2\alpha} \mathbb{E}[\|G(x^k, \boldsymbol{\xi}^k)\|^2] \\ &\quad + \frac{\eta}{2\lambda} \mathbb{E}[\|\mathbf{g}(x^k, \boldsymbol{\xi}^k) - \nabla f(x^k)\|^2]. \end{aligned} \quad (64)$$

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

$$\begin{aligned} \mathbb{E}[f(x^{k+1})] - f(x^k) &\stackrel{\text{equation 58}}{\leq} -\eta \langle \nabla f(x^k), \mathbb{E}[G(x^k, \boldsymbol{\xi}^k)] \rangle + \frac{\eta^2(L_0 + L_1 \|\nabla f(x^k)\|)}{2} \mathbb{E}[\|G(x^k, \boldsymbol{\xi}^k)\|^2] \\ &\stackrel{\text{equation 64}}{\leq} -\frac{\eta}{2} \|\nabla f(x^k)\| - \frac{\eta}{2\alpha} \mathbb{E}[\|G(x^k, \boldsymbol{\xi}^k)\|^2] + \frac{\eta}{2\lambda} \mathbb{E}[\|\mathbf{g}(x^k, \boldsymbol{\xi}^k) - \nabla f(x^k)\|^2] \\ &\quad + \frac{\eta^2(L_0 + L_1 \|\nabla f(x^k)\|)}{2} \mathbb{E}[\|G(x^k, \boldsymbol{\xi}^k)\|^2] \\ &\stackrel{\text{equation 9}}{=} -\frac{\eta}{2} \|\nabla f(x^k)\| + \frac{\eta}{2\lambda} \mathbb{E}[\|\mathbf{g}(x^k, \boldsymbol{\xi}^k) - \mathbb{E}[\mathbf{g}(x^k, \boldsymbol{\xi}^k)]\|^2] + \frac{\eta}{2\lambda} \|\mathbf{b}(x^k)\|^2 \\ &\quad - \frac{\eta}{2} \mathbb{E}[\|G(x^k, \boldsymbol{\xi}^k)\|^2] \left(1 - \frac{\eta(L_0 + L_1 \|\nabla f(x^k)\|)}{\|\nabla f(x^k)\|}\right) \\ &\leq -\frac{\eta}{2} \|\nabla f(x^k)\| + \frac{\eta\sigma^2}{2\lambda B} + \frac{\eta\zeta^2}{2\lambda}. \end{aligned} \quad (65)$$

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

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

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

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

$$f(x^k) - f^* \leq \langle \nabla f(x^k), x^k - x^* \rangle \stackrel{\text{equation 6}}{\leq} \|\nabla f(x^k)\| \|x^k - x^*\| \leq \|\nabla f(x^k)\| \underbrace{\max_{k \in [0, N-1]} \|x^k - x^*\|}_R.$$

Hence we have:

$$\|\nabla f(x^k)\| \geq \frac{f(x^k) - f^*}{R}. \quad (66)$$

Then substituting equation 66 into equation 65 we obtain:

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

This inequality is equivalent to the trailing inequality:

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

Then for $k = 0, 1, 2, \dots, N-1$ 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:

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

2160 E.1.3 THIRD CASE: $\|\nabla f(x^k)\| \leq \lambda$ AND $\|\mathbf{g}(x^k, \boldsymbol{\xi}^k)\| \leq \lambda$

2162 Using this in equation 58, we have the following with $\eta_k \leq \frac{\|\nabla f(x^k)\|}{2(L_0+L_1\|\nabla f(x^k)\|)}$ and $\alpha =$
 2163 $\|\nabla f(x^k)\|^{-1}$:

$$\begin{aligned}
 2164 & \mathbb{E} [f(x^{k+1})] - f(x^k) \stackrel{\text{equation 58}}{\leq} -\eta \langle \nabla f(x^k), \mathbb{E} [G(x^k, \boldsymbol{\xi}^k)] \rangle + \frac{\eta^2(L_0 + L_1 \|\nabla f(x^k)\|)}{2} \mathbb{E} [\|G(x^k, \boldsymbol{\xi}^k)\|^2] \\
 2165 & \stackrel{\text{equation 5}}{=} -\frac{\eta\alpha}{2} \|\nabla f(x^k)\|^2 - \frac{\eta}{2\alpha} \mathbb{E} [\|G(x^k, \boldsymbol{\xi}^k)\|^2] + \frac{\eta}{2\alpha} \mathbb{E} [\|G(x^k, \boldsymbol{\xi}^k) - \alpha \nabla f(x^k)\|^2] \\
 2166 & \quad + \frac{\eta^2(L_0 + L_1 \|\nabla f(x^k)\|)}{2} \mathbb{E} [\|G(x^k, \boldsymbol{\xi}^k)\|^2] \\
 2167 & = -\frac{\eta}{2} \|\nabla f(x^k)\| + \frac{\eta}{2\alpha} \mathbb{E} [\|G(x^k, \boldsymbol{\xi}^k) - \alpha \nabla f(x^k)\|^2] \\
 2168 & \quad - \frac{\eta}{2} \mathbb{E} [\|G(x^k, \boldsymbol{\xi}^k)\|^2] \left(1 - \frac{\eta(L_0 + L_1 \|\nabla f(x^k)\|)}{\|\nabla f(x^k)\|}\right) \\
 2169 & \leq -\frac{\eta}{2} \|\nabla f(x^k)\| + \frac{\eta}{2\alpha} \mathbb{E} [\|G(x^k, \boldsymbol{\xi}^k) - \alpha \nabla f(x^k)\|^2] \\
 2170 & \leq -\frac{\eta}{2} \|\nabla f(x^k)\| + \frac{\eta}{\alpha} \mathbb{E} [\|G(x^k, \boldsymbol{\xi}^k)\|^2 + \|\alpha \nabla f(x^k)\|^2] \\
 2171 & = -\frac{\eta}{2} \|\nabla f(x^k)\| + \frac{\eta}{\alpha} \mathbb{E} \left[\left\| \frac{\mathbf{g}(x^k, \boldsymbol{\xi}^k)}{\|\mathbf{g}(x^k, \boldsymbol{\xi}^k)\|} \right\|^2 + \left\| \frac{\nabla f(x^k)}{\|\nabla f(x^k)\|} \right\|^2 \right] \\
 2172 & = -\frac{\eta}{2} \|\nabla f(x^k)\| + \frac{2\eta\lambda \|\nabla f(x^k)\|}{\lambda} \\
 2173 & \leq -\frac{\eta}{2} \|\nabla f(x^k)\| + 2\eta\lambda. \tag{67}
 \end{aligned}$$

2187 The step size will be constant, depending on the hyperparameter λ :

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

2188 Thus, $\eta_k = \eta \leq \frac{\lambda}{2(L_0+L_1\lambda)}$.

2189 Using the convexity assumption of the function, we have the following:

$$2190 f(x^k) - f^* \leq \langle \nabla f(x^k), x^k - x^* \rangle \stackrel{\text{equation 6}}{\leq} \|\nabla f(x^k)\| \|x^k - x^*\| \leq \|\nabla f(x^k)\| \underbrace{\max_{k \in [0, N-1]} \|x^k - x^*\|}_R.$$

2191 Hence we have:

$$2192 \|\nabla f(x^k)\| \geq \frac{f(x^k) - f^*}{R}. \tag{68}$$

2193 Then substituting equation 68 into equation 67 we obtain:

$$2194 \mathbb{E} [f(x^{k+1})] - f(x^k) \leq -\frac{\eta}{2} \|\nabla f(x^k)\| + 2\eta\lambda \leq -\frac{\eta}{2R} (f(x^k) - f^*) + 2\eta\lambda.$$

2195 This inequality is equivalent to the trailing inequality:

$$2196 \mathbb{E} [f(x^{k+1})] - f^* \leq \left(1 - \frac{\eta}{2R}\right) (f(x^k) - f^*) + 2\eta\lambda.$$

2197 Then for $k = 0, 1, 2, \dots, N-1$ iterations that satisfy the conditions $\|\nabla f(x^k)\| \leq \lambda$ NSGD with
 2198 biased gradient oracle shows linear convergence:

$$2199 \mathbb{E} [f(x^N)] - f^* \leq \left(1 - \frac{\eta}{2R}\right)^N (f(x^0) - f^*) + \lambda R.$$

E.1.4 FOURTH CASE: $\|\nabla f(x^k)\| \geq \lambda$ AND $\|\mathbf{g}(x^k, \boldsymbol{\xi}^k)\| \leq \lambda$

Using this in equation 58, 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{aligned}
\mathbb{E} [f(x^{k+1})] - f(x^k) &\stackrel{\text{equation 58}}{\leq} -\eta \langle \nabla f(x^k), \mathbb{E} [G(x^k, \boldsymbol{\xi}^k)] \rangle \\
&\quad + \frac{\eta^2(L_0 + L_1 \|\nabla f(x^k)\|)}{2} \mathbb{E} [\|G(x^k, \boldsymbol{\xi}^k)\|^2] \\
&\stackrel{\text{equation 5}}{=} -\frac{\eta\alpha}{2} \|\nabla f(x^k)\|^2 - \frac{\eta}{2\alpha} \|\mathbb{E} [G(x^k, \boldsymbol{\xi}^k)]\|^2 \\
&\quad + \frac{\eta}{2\alpha} \|\mathbb{E} [G(x^k, \boldsymbol{\xi}^k)] - \alpha \nabla f(x^k)\|^2 \\
&\quad + \frac{\eta^2(L_0 + L_1 \|\nabla f(x^k)\|)}{2} \mathbb{E} [\|G(x^k, \boldsymbol{\xi}^k)\|^2] \\
&= -\frac{\eta}{2\lambda} \|\nabla f(x^k)\|^2 + \frac{\eta}{2\lambda} \|\mathbb{E} [\lambda G(x^k, \boldsymbol{\xi}^k)] - \nabla f(x^k)\|^2 \\
&\quad + \frac{\eta^2(L_0 + L_1 \|\nabla f(x^k)\|)}{2} \\
&= -\frac{\eta}{2\lambda} \|\nabla f(x^k)\|^2 + \frac{\eta}{\lambda} \left\| \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] \right\|^2 \\
&\quad + \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} \left\| \mathbb{E} \left[\left(\frac{\lambda}{\|\mathbf{g}(x^k, \boldsymbol{\xi}^k)\|} - 1 \right) \mathbf{g}(x^k, \boldsymbol{\xi}^k) \right] \right\|^2 \\
&\quad + \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] \\
&\quad + \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] \\
&\quad + \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(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)\| \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}. \tag{69}
\end{aligned}$$

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

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

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

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

$$f(x^k) - f^* \leq \langle \nabla f(x^k), x^k - x^* \rangle \stackrel{\text{equation 6}}{\leq} \|\nabla f(x^k)\| \|x^k - x^*\| \leq \|\nabla f(x^k)\| \underbrace{\max_{k \in [0, N-1]} \|x^k - x^*\|}_R.$$

Hence we have:

$$\|\nabla f(x^k)\| \geq \frac{f(x^k) - f^*}{R}. \quad (70)$$

Then substituting equation 70 into equation 69 we obtain:

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

This inequality is equivalent to the trailing inequality:

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

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

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

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

$$\mathbb{E} [f(x^N)] - f^* \lesssim \left(1 - \frac{\eta}{R}\right)^N (f(x^0) - f^*) + \frac{MR}{\lambda^2} \left(\frac{\sigma^2}{B} + \zeta^2\right) + \lambda R. \quad (71)$$

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 equation 51 and equation 52 into the convergence rate of NSGD equation 71 instead of ζ and σ^2 , respectively. Then the convergence of ZO-NSGD in the convex setup is as follows:

$$\begin{aligned} \mathbb{E} [f(x^N)] - f^* \lesssim & \underbrace{\left(1 - \frac{\eta}{R}\right)^N (f(x^0) - f^*)}_{\textcircled{1}} + \underbrace{\frac{dMR\tilde{\sigma}^2}{\lambda^2 B}}_{\textcircled{2}} + \underbrace{\frac{dMR(L_0 + L_1 M)^2 \gamma^2}{\lambda^2 B}}_{\textcircled{3}} + \underbrace{\frac{d^2 MR \Delta^2}{\lambda^2 B \gamma^2}}_{\textcircled{4}} \\ & + \underbrace{\frac{MR(L_0 + L_1 M)^2 \gamma^2}{\lambda^2}}_{\textcircled{5}} + \underbrace{\frac{d^2 MR \Delta^2}{\lambda^2 \gamma^2}}_{\textcircled{6}} + \underbrace{\lambda R}_{\textcircled{7}}. \end{aligned}$$

From term $\textcircled{7}$, we find the hyperparameter λ :

$$\textcircled{1}: \quad \lambda R \leq \varepsilon \quad \Rightarrow \quad \lambda \leq \frac{\varepsilon}{R}. \quad (72)$$

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

$$\begin{aligned} \textcircled{1}: \quad \left(1 - \frac{\eta}{R}\right)^N (f(x^0) - f^*) \leq \varepsilon \quad \Rightarrow \quad N \geq \frac{R}{\eta} \log \frac{(f(x^0) - f^*)}{\varepsilon}; \\ N = \tilde{\mathcal{O}} \left(\frac{R}{\eta} \right). \end{aligned} \quad (73)$$

From terms $\textcircled{2}$, we find the batch size B :

$$\textcircled{2}: \quad \frac{dMR\tilde{\sigma}^2}{\lambda^2 B} \leq \varepsilon \quad \Rightarrow \quad B \stackrel{\text{equation 72}}{\geq} \frac{dMR^3 \tilde{\sigma}^2}{\varepsilon^3};$$

2322
2323
2324

$$B = \mathcal{O}\left(\frac{dMR^3\tilde{\sigma}^2}{\varepsilon^3}\right). \quad (74)$$

2325
2326

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

2327
2328
2329
2330
2331
2332
2333
2334

$$\begin{aligned} \textcircled{3} : \quad & \frac{dMR(L_0 + L_1M)^2 \gamma^2}{\lambda^2 B} \leq \varepsilon \Rightarrow \gamma \leq \sqrt{\frac{\varepsilon \lambda^2 B}{dMR(L_0 + L_1M)^2}} \stackrel{\text{equation 74, equation 72}}{=} \frac{\tilde{\sigma}}{(L_0 + L_1M)}; \\ \textcircled{5} : \quad & \frac{MR(L_0 + L_1M)^2 \gamma^2}{\lambda^2} \leq \varepsilon \Rightarrow \gamma \leq \frac{\sqrt{\varepsilon^3}}{\sqrt{MR^{3/2}}(L_0 + L_1M)}; \\ & \gamma \leq \frac{1}{(L_0 + L_1M)} \min\left\{\tilde{\sigma}, \frac{\varepsilon^{3/2}}{\sqrt{MR^{3/2}}}\right\} = \frac{\varepsilon^{3/2}}{(L_0 + L_1M)\sqrt{MR^{3/2}}}. \end{aligned} \quad (75)$$

2335
2336
2337

From the remaining terms ④ 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:

2338
2339
2340
2341
2342
2343
2344
2345

$$\begin{aligned} \textcircled{4} : \quad & \frac{d^2MR\Delta^2}{\lambda^2 B \gamma^2} \leq \varepsilon \Rightarrow \Delta \leq \frac{\sqrt{\varepsilon} \lambda \gamma \sqrt{B}}{d\sqrt{MR}} \stackrel{\text{equation 74, equation 75, equation 72}}{=} \frac{\varepsilon^{3/2} \tilde{\sigma}}{\sqrt{d}(L_0 + L_1M)R^{3/2}}; \\ \textcircled{6} : \quad & \frac{d^2MR\Delta^2}{\gamma^2 \lambda^2} \leq \varepsilon \Rightarrow \Delta \leq \sqrt{\frac{\gamma^2 \lambda^2 \varepsilon}{d^2 MR}} \stackrel{\text{equation 72, equation 75}}{=} \frac{\varepsilon^3}{d(L_0 + L_1M)R^3}; \\ & \Delta \leq \frac{\varepsilon^{3/2}}{\sqrt{d}(L_0 + L_1M)R^{3/2}} \min\left\{\tilde{\sigma}, \frac{\varepsilon^{3/2}}{\sqrt{d}R^{3/2}}\right\}. \end{aligned} \quad (76)$$

2346
2347

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

2348
2349

$$N \stackrel{\text{equation 73}}{=} \tilde{\mathcal{O}}\left(\frac{R}{\eta}\right), \quad T = N \cdot B \stackrel{\text{equation 73, equation 74}}{=} \mathcal{O}\left(\frac{d\tilde{\sigma}^2 MR^4}{\varepsilon^3 \eta}\right)$$

2350
2351

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

2352
2353
2354

$$\Delta \stackrel{\text{equation 76}}{\lesssim} \frac{\varepsilon^{3/2}}{\sqrt{d}(L_0 + L_1M)R^{3/2}} \min\left\{\tilde{\sigma}, \frac{\varepsilon^{3/2}}{\sqrt{d}R^{3/2}}\right\}$$

2355
2356

the maximum level of noise with smoothing parameter $\frac{\varepsilon^{3/2}}{(L_0 + L_1M)\sqrt{MR^{3/2}}}$ equation 75.

2357
2358

F ADDITIONAL CLARIFICATION

2359
2360
2361
2362

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 $\|\nabla f(x^k, \xi^k)\| \geq \sqrt{2}\sigma$ and $\|\nabla f(x^k)\| \geq \sqrt{2}\sigma$), shows that it is possible to achieve the desired accuracy ε in a finite number of iterations.

2363
2364
2365

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

2366
2367
2368
2369
2370

$$\begin{aligned} f(x^{k+1}) - f(x^k) & \stackrel{\text{equation 8}}{\leq} \langle \nabla f(x^k), x^{k+1} - x^k \rangle + \frac{L_0 + L_1 \|\nabla f(x^k)\|}{2} \|x^{k+1} - x^k\|^2 \\ & = -\eta \langle \nabla f(x^k), G(x^k, \xi^k) \rangle + \frac{\eta^2 (L_0 + L_1 \|\nabla f(x^k)\|)}{2} \|G(x^k, \xi^k)\|^2. \end{aligned} \quad (77)$$

2371
2372

Let us evaluate first summand of equation 77 with $\alpha = \|\nabla f(x^k)\|^{-1}$:

2373
2374
2375

$$\begin{aligned} -\eta \langle \nabla f(x^k), G(x^k, \xi^k) \rangle & \stackrel{\text{equation 5}}{=} -\frac{\alpha \eta}{2} \|\nabla f(x^k)\|^2 - \frac{\eta}{2\alpha} \|G(x^k, \xi^k)\|^2 \\ & \quad + \frac{\eta}{2\alpha} \|G(x^k, \xi^k) - \alpha \nabla f(x^k)\|^2 \end{aligned}$$

$$\begin{aligned}
&= -\frac{\eta}{2} \|\nabla f(x^k)\| - \frac{\eta}{2\alpha} \|G(x^k, \xi^k)\|^2 \\
&\quad + \frac{\eta}{2\lambda^2\alpha} \|\lambda G(x^k, \xi^k) - \lambda\alpha \nabla f(x^k)\|^2 \\
&= -\frac{\eta}{2} \|\nabla f(x^k)\| - \frac{\eta}{2\alpha} \|G(x^k, \xi^k)\|^2 \\
&\quad + \frac{\eta}{2\lambda^2\alpha} \|\text{clip}_\lambda(\nabla f(x^k, \xi^k)) - \text{clip}_\lambda(\nabla f(x^k))\|^2
\end{aligned}$$

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

$$\begin{aligned}
-\eta \langle \nabla f(x^k), \mathbb{E}[G(x^k, \xi^k)] \rangle &\leq -\frac{\eta}{2} \|\nabla f(x^k)\| - \frac{\eta}{2\alpha} \mathbb{E}[\|G(x^k, \xi^k)\|^2] \\
&\quad + \frac{\eta}{2\lambda^2\alpha} \mathbb{E}[\|\nabla f(x^k, \xi^k) - \nabla f(x^k)\|^2]. \tag{78}
\end{aligned}$$

Using this in equation 78, we have the following with $\eta_k \leq \frac{\|\nabla f(x^k)\|}{2(L_0 + L_1 \|\nabla f(x^k)\|)}$:

$$\begin{aligned}
\mathbb{E}[f(x^{k+1})] - f(x^k) &\stackrel{\text{equation 77}}{\leq} -\eta \langle \nabla f(x^k), \mathbb{E}[G(x^k, \xi^k)] \rangle + \frac{\eta^2(L_0 + L_1 \|\nabla f(x^k)\|)}{2} \mathbb{E}[\|G(x^k, \xi^k)\|^2] \\
&\stackrel{\text{equation 78}}{\leq} -\frac{\eta}{2} \|\nabla f(x^k)\| - \frac{\eta}{2\alpha} \mathbb{E}[\|G(x^k, \xi^k)\|^2] + \frac{\eta}{2\lambda^2\alpha} \mathbb{E}[\|\nabla f(x^k, \xi^k) - \nabla f(x^k)\|^2] \\
&\quad + \frac{\eta^2(L_0 + L_1 \|\nabla f(x^k)\|)}{2} \mathbb{E}[\|G(x^k, \xi^k)\|^2] \\
&= -\frac{\eta}{2} \|\nabla f(x^k)\| + \frac{\eta}{2\lambda^2\alpha} \mathbb{E}[\|\nabla f(x^k, \xi^k) - \nabla f(x^k)\|^2] \\
&\quad - \frac{\eta}{2} \mathbb{E}[\|G(x^k, \xi^k)\|^2] \left(1 - \frac{\eta(L_0 + L_1 \|\nabla f(x^k)\|)}{\|\nabla f(x^k)\|}\right) \\
&\leq -\frac{\eta}{2} \|\nabla f(x^k)\| + \frac{\eta\sigma^2}{2\lambda^2\alpha} \\
&\leq -\frac{\eta}{2} \|\nabla f(x^k)\| + \frac{\eta}{4} \|\nabla f(x^k)\| \\
&= -\frac{\eta}{4} \|\nabla f(x^k)\|. \tag{79}
\end{aligned}$$

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

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

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:

$$\begin{aligned}
f(x^k) - f(s) &\leq \langle \nabla f(x^k), x^k - s \rangle \\
&\stackrel{\text{equation 6}}{\leq} \|\nabla f(x^k)\| \|x^k - s\| \\
&\leq \|\nabla f(x^k)\| \underbrace{\|x^0 - s\|}_{R_s}.
\end{aligned}$$

Hence we have:

$$\|\nabla f(x^k)\| \geq \frac{f(x^k) - f(s)}{R_s}. \tag{80}$$

2430 Then substituting equation 80 into equation 79 we obtain:

$$2431 \mathbb{E} [f(x^{k+1})] - f(x^k) \leq -\frac{\eta}{4} \|\nabla f(x^k)\| \leq -\frac{\eta}{4R_s} (f(x^k) - f(s)).$$

2434 This inequality is equivalent to the trailing inequality:

$$2435 \mathbb{E} [f(x^{k+1})] - f^* \leq \left(1 - \frac{\eta}{4R_s}\right) (f(x^k) - f^*) + \frac{\eta}{4R_s} (f(s) - f^*).$$

2439 Then for $k = 0, 1, 2, \dots, N - 1$ iterations that satisfy the conditions $\|\nabla f(x^k, \xi^k)\| \geq \sqrt{2}\sigma$ and
 2440 $\|\nabla f(x^k)\| \geq \sqrt{2}\sigma$ NSGD shows linear convergence:

$$2442 f(x^N) - f^* \leq \left(1 - \frac{\eta}{4R_s}\right)^N (f(x^0) - f^*) + f(s) - f^*.$$

2445 Thus, we have shown that it is indeed possible to converge to a linear rate of convergence on logistic
 2446 regression using the hyperparameter R_s .

2447
2448
2449
2450
2451
2452
2453
2454
2455
2456
2457
2458
2459
2460
2461
2462
2463
2464
2465
2466
2467
2468
2469
2470
2471
2472
2473
2474
2475
2476
2477
2478
2479
2480
2481
2482
2483