

HIGHLY SMOOTH ZERO-ORDER METHODS FOR SOLVING OPTIMIZATION PROBLEMS UNDER THE PL CONDITION[†]

© 2023 A. V. Gasnikov^{a,b,c}, A. V. Lobanov^{a,c,*‡} and F. S. Stonyakin^{a,d}

^a 141700 Dolgoprudny, Institutskiy per., 9, Moscow Institute of Physics and Technology, Russia

^b 121205 Moscow, B. Boulevard 30, bld. 1, Skolkovo Institute of Science and Technology, Russia

^c 125047 Moscow, A. Solzhenitsyn st., 25, Institute for System Programming of the RAS, Russia

^d 295007 Simferopol, Prospekt Vernadskogo 4, V.I. Vernadsky Crimean Federal University, Russia

*e-mail: lobbsasha@mail.ru

Received 2023, revised 2023, accepted 2023.

Abstract – In this paper, we study the black box optimization problem under the Polyak–Lojasiewicz (PL) condition, assuming that the objective function is not just smooth, but has higher smoothness. By using "kernel-based" approximations instead of the exact gradient in the Stochastic Gradient Descent method, we improve the best-known results of convergence in the class of gradient-free algorithms solving problems under the PL condition. We generalize our results to the case where a zeroth-order oracle returns a function value at a point with some adversarial noise. We verify our theoretical results on the example of solving a system of nonlinear equations.

Keywords: Black-box optimization, gradient-free methods, kernel approximation, maximum noise level.

1. INTRODUCTION

The black box problem is a fundamental optimization problem when the objective function has only a zero-order oracle [1]. Recently, this problem has received significant attention in the setting of reinforcement learning [2; 3], federated learning [4; 5], and deep learning [6; 7]. Particularly in applied problems of multi-armed bandit [8; 9], online optimization [10–12], huge-scale optimization [13; 14], and hyperparameter tuning [15; 16]. In addition, the black box problem arises if the information about the derivatives is too expensive or not available [17].

To solve such problems in the convex case, there are various techniques for developing optimal gradient-free algorithms based on first-order optimization algorithms, optimal by three criteria at once: oracle complexity, iteration complexity, and maximum level of admissible noise still allowing to guarantee certain accuracy [18]. The main idea of which is to calculate the gradient approximation instead of an exact gradient [19]. For example, L_1 , and L_2 randomized approximations [8; 12; 20–22] and "kernel-based" approximation [11; 23–25] are usually used for the smooth case.

[†]The research was supported by Russian Science Foundation (project No. 21-71-30005), <https://rscf.ru/en/project/21-71-30005/>.

[‡]The main contribution to the article belongs to Aleksandr Lobanov <lobbsasha@mail.ru>. According to the rules of the journal, the authors of the article are arranged in **alphabetical** order.

Table 1: Comparison of convergence results of Zero-Order Mini-batch SGD method via Kernel-based approximation, "Kernel approx." (This work) with the existing counterpart via Gaussian approximation, "Gaus. approx." [36]. Notation: d = dimension, $\delta(x)$ = adversarial deterministic noise, ξ = adversarial stochastic noise, Δ = level noise, β = smoothness order parameter.

Case	Zero-order oracle	Error floor		Stochastic?	Adversarial noise	
		Gaussian approx.	Kernel approx.		ADN?	ASN?
deterministic	$\tilde{f}(x) = f(x) + \delta(x)$	$\mathcal{O}(d^2\Delta)$	$\mathcal{O}\left(d^2\Delta^{\frac{2(\beta-1)}{\beta}}\right)$	✗	✓	✗
two-point	$\tilde{f}(x, \xi) = f(x, \xi) + \delta(x)$	$\mathcal{O}(d^2\Delta)$	$\mathcal{O}\left(d^2\Delta^{\frac{2(\beta-1)}{\beta}}\right)$	✓	✓	✗
one-point	$\tilde{f}(x, \xi) = f(x) + \xi$	ε	ε	✓	✗	✓
	$\tilde{f}(x, \xi) = f(x) + \xi + \delta(x)$	$\mathcal{O}(d^2\Delta)$	$\mathcal{O}\left(d^2\Delta^{\frac{2(\beta-1)}{\beta}}\right)$	✓	✓	✓

Where the "kernel-based" approximation takes advantage of higher-order smoothness as opposed to the L_2 and L_1 randomized approximations. For the non-smooth case, there are smoothing techniques via L_1 and L_2 randomized approximation [5; 26].

In contrast to the convex case, where algorithms have been intensively appearing and being studied in theory and practice in recent years [e.g., 27, and see above], the analysis of the black box problem in the nonconvex case is just beginning to be actively studied [28–30]. One such problem that is frequently common in the literature lately is the smooth (nonconvex) optimization problem under the Polyak–Lojasiewicz condition [31–35]. In this formulation of the problem, [36] studied the biased Stochastic Gradient Descent (SGD) method and also proposed its gradient-free counterpart.

We focus on solving the smooth black-box optimization problem (in particular, the nonconvex optimization) under the Polyak–Lojasiewicz condition. We propose an optimal zero-order algorithm (see Algorithm 1) based on Mini-batch SGD (added batching to biased Stochastic Gradient Descent [36] to improve iteration complexity). Using kernel approximation as a technique for developing a gradient-free algorithm, we show in theory and in practice a significant improvement in convergence rate estimates. We provide an extended analysis of Algorithm 1, namely we consider a (close to reality) stochastic formulation of the problem with adversarial noise [37–39].

1.1. Contribution

- We present a Zero-Order Mini-batch SGD method that significantly improves existing estimates of the convergence results (in particular, estimates of the error floor, see the "deterministic" case of the zero-order oracle (first line) in Table 1 and Table 2) in a class of gradient-free algorithms for solving smooth nonconvex problems under the PL condition.
- We provide an extended theoretical analysis of the Zero-Order Mini-batch SGD algorithm by considering a stochastic setting using a gradient approximation with one-point and two-point feedback, when the zero-order oracle is corrupted by an bounded adversarial deterministic ("ADN") and stochastic ("ASN") noise. We show that also in the stochastic setting a Zero-Order Mini-batch Stochastic Gradient Descent method is superior to the existing counterpart [36] (see stochastic cases of the zero-order oracle $\tilde{f}(x, \xi)$ in Table 1 and Table 2).

Table 2: Comparison of iteration complexity $\#N$ oracle complexity, $\#T$ and maximum noise level $\#\Delta$ of the algorithm with the following approaches: Kernel and Gaussian approximations for zero-order oracle cases ① : $\tilde{f}(x) = f(x) + \delta(x)$, ② : $\tilde{f}(x, \xi) = f(x, \xi) + \delta(x)$, ③ : $\tilde{f}(x, \xi) = f(x) + \xi$, ④ : $\tilde{f}(x, \xi) = f(x) + \xi + \delta(x)$. In both approaches, the dependence of the estimates on the batch size B is presented. Notation: β = order of smoothness, ε = accuracy of the solution to the problem, d = dimension, μ = Polyak–Lojasiewicz condition constant.

Case	Batch Size	Gaussian approximation			Kernel approximation		
		$\# N$	$\# T$	$\# \Delta$	$\# N$	$\# T$	$\# \Delta$
①, ②, ④	$B \in [1, \beta^3 d]$	$\tilde{O}\left(\frac{d}{B}\mu^{-1}\right)$	$\tilde{O}(d\mu^{-1})$	$\Delta \leq \mu\varepsilon d^{-3/2}$	$\tilde{O}\left(\frac{d}{B}\mu^{-1}\right)$	$\tilde{O}(d\mu^{-1})$	$\Delta \leq \mu\varepsilon d^{-3/2}$
	$B > \beta^3 d$	$\tilde{O}(\mu^{-1})$	$\max\left\{\tilde{O}(\mu^{-1}B), \tilde{O}\left(\frac{d^4\Delta^2}{\varepsilon^2\mu^\beta}\right)\right\}$	$\Delta \leq \mu\varepsilon d^{-2}$	$\tilde{O}(\mu^{-1})$	$\max\left\{\tilde{O}(\mu^{-1}B), \tilde{O}\left(\frac{d^{2+\frac{2}{\beta-1}}\Delta^2}{\varepsilon^{\frac{2}{\beta-1}}\mu^{\frac{2\beta-1}{\beta-1}}}\right)\right\}$	$\Delta \leq \frac{(\mu\varepsilon)^{\frac{\beta}{2(\beta-1)}}}{d^{\frac{\beta}{\beta-1}}}$
③	$B \in [1, \beta^3 d]$	$\tilde{O}\left(\frac{d}{B}\mu^{-1}\right)$	$\tilde{O}(d\mu^{-1})$	$\Delta \leq \mu\varepsilon d^{-1}$	$\tilde{O}\left(\frac{d}{B}\mu^{-1}\right)$	$\tilde{O}(d\mu^{-1})$	$\Delta \leq \mu\varepsilon d^{-1}$
	$B > \beta^3 d$	$\tilde{O}(\mu^{-1})$	$\max\left\{\tilde{O}(\mu^{-1}B), \tilde{O}\left(\frac{d^4\Delta^2}{\varepsilon^2\mu^\beta}\right)\right\}$	$\Delta \leq \mu\varepsilon d^{-2}B^{1/2}$	$\tilde{O}(\mu^{-1})$	$\max\left\{\tilde{O}(\mu^{-1}B), \tilde{O}\left(\frac{d^{2+\frac{2}{\beta-1}}\Delta^2}{\varepsilon^{\frac{2}{\beta-1}}\mu^{\frac{2\beta-1}{\beta-1}}}\right)\right\}$	$\Delta \leq \frac{(\mu\varepsilon)^{\frac{\beta}{2(\beta-1)}}}{d^{\frac{\beta}{\beta-1}}}B^{1/2}$

- We empirically verify our theoretical results by comparing the Zero-Order Mini-batch Stochastic Gradient Descent method with the existing algorithm using a typical example for problems under the Polyak–Lojasiewicz condition: solving a system of nonlinear equations.
- We demonstrate the effectiveness of the randomized approximation in practice by comparing it to the Gaussian approximation from [36] and show the advantages of "kernel-based" approximation in experiments by additionally comparing it to the L_2 approximation.

1.2. Paper organization

This article has the following structure. In Section 2 we provide related works. We describe the statement of the problem in Section 3. We present the optimal gradient-free Algorithm 1 in Section 4. In Section 5, we extend our analysis of the Algorithm 1 to a stochastic setting with different models of adversarial noise. In Section 7, we discuss the results. While in Section 8 we analyze the effectiveness of our approach on a practical experiment. Section 9 concludes the paper. We give a detailed proof of the Theorems and Lemmas in the Appendix (Supplementary Materials).

2. RELATED WORK

Polyak–Lojasiewicz condition. The field of research on the nonconvex problem under the Polyak–Lojasiewicz condition can be traced back to [31], where it is shown that this condition is sufficient for the gradient descent to achieve a global linear convergence rate. Recently, problems with Polyak–Lojasiewicz condition have been actively investigated. For instance, [35] proposed a new formulation of the problem under the Polyak–Lojasiewicz condition, including a compromise and an early stopping rule to guarantee its achievement. Already in the new statement of the problem, [40] proposed a fully adaptive gradient algorithm (with respect to the Lipschitz constant of the gradient and the noise level in the gradient) for solving problems with this condition. However, in this work we consider the standard statement of the problem under the Polyak–Lojasiewicz condition [e.g. see 33; 41]. Also, it is shown in [41] that non-accelerated algorithms are optimal for smooth problems under the Polyak–Lojasiewicz condition. Therefore, in this paper,

we consider non-accelerated first-order optimization algorithms as a base for developing optimal gradient-free methods.

SGD type algorithms. Many works [42–47] investigated Stochastic Gradient Descent in various settings. For black-box problems, this algorithm and its modifications are the basis for developing gradient-free algorithms. For instance, [5] used the following first-order optimization algorithms as a basis for developing gradient-free algorithms in the federated learning setting: Minibatch Accelerated SGD and Single-Machine Accelerated SGD [from 48], and Federated Accelerated Stochastic Gradient Descent [from 49], which are accelerated modifications of the SGD, namely the AC-SA [50]. In a nonconvex optimization problem with a Polyak–Lojasiewicz condition, [36] studied a biased Stochastic Gradient Descent, where the oracle has access to biased and noisy gradient estimates, and showed that Stochastic Gradient Descent methods can generally converge only to the neighborhood of the solution. In this paper, we use the biased Stochastic Gradient Descent algorithm [36] as a basis for developing an optimal zero-order method for solving a smooth nonconvex optimization problem, assuming that the Polyak–Lojasiewicz condition is satisfied.

Kernel approximation. The works [11; 23–25] investigated and used kernel approximation as a technique for creating gradient-free algorithms. In the survey [18] showed that the significant difference of this approximation from others is taking into account the advantages of the high order of smoothness of the objective function, i.e. satisfying the Hölder condition. All these works [11; 24; 25] used the central finite difference (CFD) scheme instead of the forward finite difference scheme (FFD) in kernel approximation. It turns out there is an explanation, [51] showed that in a smooth case, one should use CFD, not FFD. In this paper, we use the kernel approximation since we assume that our function has higher smoothness, in contrast to [36], where only smoothness is assumed.

Stochastic optimization. Stochastic optimization problems have received special attention in the literature [46; 52–56]. For example, [57] investigated a stochastic problem with non-sub-Gaussian (heavy-tailed) noise in the convex case, and [56] investigated in the convex-concave case. For black-box problems, [24; 25] studied a one-point gradient approximation with additive stochastic noise. [58] studied a two-point gradient approximation corrupted by an adversarial deterministic noise, and [5] studied a one-point gradient approximation with the same adversarial deterministic noise. In this paper, we generalize our analysis to the stochastic optimization problem and consider gradient approximation with one-point and two-point feedback obtained via function realizations. We consider two cases of noise: noise as defined in [12; 24; 25] and noise as defined in [26; 58].

To the best of our knowledge for the moment there are lack of results around gradient-free methods for problems with Polyak–Lojasiewicz condition. The known results [59–61] are dominated by [36] as we consider to be state of the art results.

3. SETUP

We study black-box optimization problems of the form:

$$f^* := \min_{x \in Q \subset \mathbb{R}^d} f(x), \tag{1}$$

where $f : \mathbb{R}^d \rightarrow \mathbb{R}$ is a function that we want to minimize over a closed convex subset Q of \mathbb{R}^d . The problem (1) is a general statement problem in the field of optimization. In order to narrow down the class of problems to be solved, let us define the function class using the following assumptions.

3.1. Notation

We use $\langle x, y \rangle := \sum_{i=1}^d x_i y_i$ to denote standard inner product of $x, y \in \mathbb{R}^d$. We denote l_p -norms ($p \geq 1$) in \mathbb{R}^d as $\|x\|_p := \left(\sum_{i=1}^d |x_i|^p \right)^{1/p}$. Particularly, for l_2 -norm in \mathbb{R}^d it follows $\|x\|_2 := \sqrt{\langle x, x \rangle}$. We denote l_p -sphere as $S_p^d(r) := \{x \in \mathbb{R}^d : \|x\|_p = r\}$. Operator $\mathbb{E}[\cdot]$ denotes full mathematical expectation. We use the notation $\tilde{O}(\cdot)$ to hide logarithmic factors.

3.2. Assumptions on the objective function

For all our theoretical results we assume that f is not just smooth, but has a high order of smoothness.

Assumption 1 (Higher order smoothness) *Let l denote the maximal integer number strictly less than β . Let $\mathcal{F}_\beta(L)$ denote the set of all functions $f : \mathbb{R}^d \rightarrow \mathbb{R}$ which are differentiable l times and for all $x, z \in Q$ the Hölder-type condition holds:*

$$\left| f(z) - \sum_{0 \leq |n| \leq l} \frac{1}{n!} D^n f(x) (z - x)^n \right| \leq L_\beta \|z - x\|_2^\beta,$$

where $L_\beta > 0$, the sum is over multi-index $n = (n_1, \dots, n_d) \in \mathbb{N}^d$, we used the notation $n! = n_1! \cdots n_d!$, $|n| = n_1 + \cdots + n_d$, and $\forall v = (v_1, \dots, v_d) \in \mathbb{R}^d$ we defined

$$D^n f(x) v^n = \frac{\partial^{|n|} f(x)}{\partial^{n_1} x_1 \cdots \partial^{n_d} x_d} v_1^{n_1} \cdots v_d^{n_d}.$$

Also we assume the Polyak–Lojasiewicz condition (μ -PL) with parameter $\mu > 0$.

Assumption 2 (μ -PL) *The function is differentiable and there exists constant $\mu > 0$ s.t. $\forall x \in \mathbb{R}^d$*

$$\|\nabla f(x)\|_2^2 \geq 2\mu(f(x) - f^*). \quad (2)$$

Assumption 1 is quite common in the literature [e.g. 11; 24; 25]. We introduced this assumption in order to take advantage of the properties of higher smoothness. As far as we know Assumption 2 was introduced by [31]. Functions satisfying this assumption are called gradient-dominated functions [62].

3.3. Assumptions on the biased gradient oracle

We now formulate the standard assumptions about the gradient oracle for the biased SGD method [36]. To do this, we start by introducing the following definition.

Definition 1 (Biased Gradient Oracle) A map $\mathbf{g} : \mathbb{R}^d \times \mathcal{D} \rightarrow \mathbb{R}^d$ s.t.

$$\mathbf{g}(x, \xi) = \nabla f(x) + \mathbf{b}(x) + \mathbf{n}(x, \xi) \quad (3)$$

for a bias $\mathbf{b} : \mathbb{R}^d \rightarrow \mathbb{R}^d$ and zero-mean noise $\mathbf{n} : \mathbb{R}^d \times \mathcal{D} \rightarrow \mathbb{R}^d$, that is $\mathbb{E}_\xi \mathbf{n}(x, \xi) = 0, \forall x \in \mathbb{R}^d$.

We assume that this gradient oracle has bias and noise, and they are bounded.

Assumption 3 ((M, σ^2)-bounded noise) There exists constants $M, \sigma^2 \geq 0$ such that $\forall x \in \mathbb{R}^d$

$$\mathbb{E}_\xi \|\mathbf{n}(x, \xi)\|_2^2 \leq M \|\nabla f(x) + \mathbf{b}(x)\|_2^2 + \sigma^2. \quad (4)$$

Assumption 4 (Bounded bias) There exists constants $0 \leq m < 1$, and $\zeta^2 \geq 0$ s.t. $\forall x \in \mathbb{R}^d$

$$\|\mathbf{b}(x)\|_2^2 \leq m \|\nabla f(x)\|_2^2 + \zeta^2. \quad (5)$$

Many prior work in the context stochastic optimization often assumed the bounded noise and bounded bias. For example, the Assumption 3 is similar as in [46]. In the case without bias, Assumption 3 is referred to as the strong growth condition [e.g. 63]. Also the Assumption 4 was used by [64].

The assumptions introduced in this subsection are necessary to use the following auxiliary lemma, since our approach to creating a gradient-free method is based on the SGD algorithm [36] (see Appendix C).

Lemma 1 Let $\{x_k\}_{k \geq 0}$ denote the iterates of Algorithm SGD [36] with batching, function f satisfy Assumptions 1-2 with $\beta = 2$, and the gradient oracle (3) satisfy Assumptions 3-4. Then there exists step size $\eta \leq \frac{1}{(M+1)L_2}$ such that it holds for all $N \geq 0$ and an arbitrary batch size B

$$\mathbb{E}[f(x_N)] - f^* \leq (1 - \eta\mu(1 - m))^N (f(x_0) - f^*) + \frac{\zeta^2}{2\mu(1 - m)} + \frac{\eta L_2 \sigma^2}{2B\mu(1 - m)}.$$

4. ZERO-ORDER MINI-BATCH SGD

In this section, we present our approach to solving problem (1) when the gradient oracle (3) has no information about the derivatives of the objective function. The main idea of our approach is to create an optimal (on oracle complexity, iteration complexity, and maximum level of noise) gradient-free algorithm based on the first-order method (namely, Stochastic Gradient Descent). To begin, we introduce an approximation of the gradient oracle (3) via a zero-order oracle (value of the objective function $f(x)$ with some adversarial deterministic noise $\delta(x)$ such that $|\delta(x)| \leq \Delta$ and $\Delta > 0$):

$$\tilde{f}(x) = f(x) + \delta(x), \quad (6)$$

which take advantage of the higher smoothness of the function [11; 24; 25]. Such approximation is referred to as "kernel-based" approximation of gradient and has the following form

$$\tilde{\mathbf{g}}(x, \mathbf{e}) = d \frac{\tilde{f}(x + \gamma r \mathbf{e}) - \tilde{f}(x - \gamma r \mathbf{e})}{2\gamma} K(r) \mathbf{e}, \quad (7)$$

where $\gamma > 0$ is smoothing parameter, \mathbf{e} is uniformly distributed in $S_2^d(1)$, r is uniformly distributed in $[-1, 1]$, \mathbf{e} and r are independent, $K : [-1, 1] \rightarrow \mathbb{R}$ is a kernel function that satisfies

$$\mathbb{E}[K(u)] = 0, \quad \mathbb{E}[uK(u)] = 1, \quad \mathbb{E}[u^j K(u)] = 0, \quad j = 2, \dots, l, \quad \mathbb{E}[|u|^\beta |K(u)|] < \infty.$$

Now we can present Algorithm 1. This algorithm is a modification of SGD method. Instead of calculating the first-order gradient oracle from Definition 1, we calculate an approximation of the gradient (7). Also, to achieve an optimal iteration complexity, we add a batch size B .

Algorithm 1 Zero-Order Mini-batch Stochastic Gradient Descent (ZO-MB-SGD)

Input: step size η , iteration number N , batch size B , Kernel $K : [-1, 1] \rightarrow \mathbb{R}$, smoothing parameter γ , $x_0 \in \mathbb{R}^d$

for $k = 0$ **to** $N - 1$ **do**

1. Sample vectors $\mathbf{e}_1, \mathbf{e}_2, \dots, \mathbf{e}_B$ uniformly distributed on the unit sphere $S_2^d(1)$ and scalars r_1, r_2, \dots, r_B uniformly distributed on the interval $[-1, 1]$
2. Define $\tilde{\mathbf{g}}(x_k, \mathbf{e}_i) = d \frac{\tilde{f}(x_k + \gamma r_i \mathbf{e}_i) - \tilde{f}(x_k - \gamma r_i \mathbf{e}_i)}{2\gamma} K(r_i) \mathbf{e}_i$ using (6)
3. Calculate $\mathbf{g}_k = \frac{1}{B} \sum_{i=1}^B \tilde{\mathbf{g}}(x_k, \mathbf{e}_i)$
4. $x_{k+1} \leftarrow x_k - \eta \mathbf{g}_k$

end for

Return: x_N

The next theorem presents the convergence result of ZO-MB-SGD method in terms of the expectation.

Theorem 1 *Let function $f(x)$ satisfy Assumptions 1 - 2 and the gradient approximation (7) satisfy Assumptions 3 - 4 then by choosing the step size $\eta \leq \frac{1}{(M+1)L_2}$, there exists parameters*

$$M \leq 6\beta^3 d, \quad \sigma^2 \leq \frac{3\beta^3 d^2 L_2^2 \gamma^2}{4} + \frac{2\beta^3 d^2 \Delta^2}{\gamma^2}, \quad m = 0, \quad \zeta^2 \leq \beta^2 d^2 \left(L_\beta^2 \gamma^{2(\beta-1)} + \frac{\Delta^2}{\gamma^2} \right)$$

such that Algorithm 1 with smoothing parameter $\gamma \leq \mathcal{O}(\Delta^{1/\beta})$ achieves the following error floor

$$\mathbb{E}f(x_N) - f^* = \max \left\{ \varepsilon, \mathcal{O} \left(d^2 \Delta^{\frac{2(\beta-1)}{\beta}} \right) \right\},$$

where f^* is the solution to problem (1), L_2 is the Lipschitz gradient constant with respect to 2-norm.

The convergence result of Theorem 1 shows that Algorithm 1 with gradient approximation (7) achieves the error floor $\mathcal{O} \left(d^2 \Delta^{\frac{2(\beta-1)}{\beta}} \right)$ with linear convergence rate. This asymptote arises due to the accumulation of adversarial noise in the bias $\mathbf{b}(x)$. To achieve ε -accuracy of Problem (1), the maximum level of noise can be defined as $\Delta \leq \mathcal{O} \left(\varepsilon^{\frac{\beta}{2(\beta-1)}} d^{\frac{-\beta}{\beta-1}} \right)$. This result significantly improves the approach described in [36]. Namely, using the same concept of the zero-order oracle (6) Algorithm 1 via Gaussian approximation achieves the error floor $\mathcal{O}(d^2 \Delta)$, respectively, with the following maximum level of noise $\Delta \leq \mathcal{O}(\varepsilon d^{-2})$. For a detailed proof of Theorem 1, see Appendix D. Also see Appendix H for details on estimates for Gaussian smoothing approximation.

Note that the results of Theorem 1 were obtained by using the deterministic setting of the zero-order oracle (6). However, if we reformulate problem (1) to a stochastic setting, specifying that $f(x) := \mathbb{E}_\xi f(x, \xi)$ and we use a gradient approximation (7) with two-point feedback, replacing in the zero-order oracle (6) the calculation of objective function value $f(x)$ with the calculation of objective function value on realizations $f(x, \xi)$, then it is not difficult to show that the convergence results of Theorem 1 also hold for Algorithm 1, which uses the gradient approximation with two-point feedback, where the term two-point feedback implies that it is possible to get the value of the objective function at two points on one realization. See Appendix E and I for a detailed proof of the case where Algorithm 1 uses a gradient approximation with two-point feedback, applying the "kernel-based" approximation and Gaussian approximation approaches, respectively. The concepts of a stochastic zero-order oracle when two-point feedback is not available are explored in the Section 5.

Remark 1 *It should be noted that Algorithm 1 has a linear convergence rate (see Theorem 1). Then, since the ZO-MB-SGD method uses a randomized scheme (randomization on the L_2 sphere), we can obtain exact estimates of the large deviation probabilities using Markov's inequality [65]:*

$$\mathcal{P}(f(x_{N(\varepsilon\omega)}) - f^* \geq \varepsilon) \leq \omega \frac{\mathbb{E}[f(x_{N(\varepsilon\omega)})] - f^*}{\varepsilon\omega} \leq \omega.$$

5. EXTENDED ANALYSIS OF CONVERGENCE VIA STOCHASTIC ZERO-ORDER ORACLE

In this section, we continue our study of gradient-free methods for solving Problem (1) with the μ -PL condition (see Assumption 2). We want to extend the analysis of Algorithm 1. Specifically, for cases where the zero-order oracle has a stochastic setting $\tilde{f}(x, \xi)$. To do this, we consider the concept of gradient approximation when only one-point feedback is available [see, e.g. 66] with two types of adversarial noise: adversarial stochastic noise, which is quite common in the following works [12; 24; 25], and adversarial deterministic noise, which is actively used in [5; 26; 58]. To introduce stochastic, let us rewrite the initial problem (1), where zero-order oracle returns an unbiased noisy stochastic function value $f(x, \xi)$:

$$\min_{x \in Q \subset \mathbb{R}^d} \{f(x) := \mathbb{E}_\xi f(x, \xi)\}. \quad (8)$$

5.1. Stochastic zero-order oracle with additive adversarial stochastic noise

If for some reason two-point feedback is not available, you can use Zero-Order Mini-batch Stochastic Gradient Descent method (see Algorithm 1) with one-point feedback. Then using the concept of additive stochastic noise, the "kernel-based" approximation of the gradient (7) takes the following form:

$$\tilde{\mathbf{g}}(x, \xi, \mathbf{e}) = d \frac{f(x + \gamma r \mathbf{e}) + \xi_1 - f(x - \gamma r \mathbf{e}) - \xi_2}{2\gamma} K(r) \mathbf{e}, \quad (9)$$

where the zero-order oracle in this case is defined in the following form:

$$\tilde{f}(x, \xi) = f(x) + \xi, \quad (10)$$

and $\xi_1 \neq \xi_2$ are adversarial stochastic noises such that $\mathbb{E}[\xi_1^2] \leq \tilde{\Delta}^2$ and $\mathbb{E}[\xi_2^2] \leq \tilde{\Delta}^2$, $\tilde{\Delta} \geq 0$, and the random variables ξ_1 and ξ_2 are independent from \mathbf{e} and r . Also, this concept does not require the assumption of zero-mean ξ_1 and ξ_2 . It is enough that $\mathbb{E}[\xi_1 \mathbf{e}] = 0$ and $\mathbb{E}[\xi_2 \mathbf{e}] = 0$. The gradient approximation (9) has some similarities with the two-point gradient approximation (see discussion following Theorem 1), but it is a one-point gradient approximation because it is impossible to obtain the value of the objective function on one realization twice. The following theorem provides the convergence results of ZO-MB-SGD method (see Algorithm 1) with the gradient approximation (9).

Theorem 2 *Let function $f(x)$ satisfy Assumptions 1 - 2 and the gradient approximation (9) satisfy Assumptions 3 - 4 then by choosing the step size $\eta \leq \frac{1}{(M+1)L_2}$ there exists parameters*

$$M \leq 18\beta^3 d, \quad \sigma^2 \leq \frac{3d^2 L_2^2 \gamma^2}{4} + \frac{2d^2 \tilde{\Delta}^2}{\gamma^2}, \quad m = 0, \quad \zeta^2 \leq 8\beta^2 d^2 L_\beta^2 \gamma^{2(\beta-1)}$$

such that Algorithm 1 with approximation of the gradient (9) has the following convergence rate

$$\mathbb{E}f(x_N) - f^* = \mathcal{O}((1 - \eta\mu(1 - m))^N (f(x_0) - f^*)),$$

where f^* is the solution to problem (8).

The results of Theorem 2 show that Algorithm 1 with gradient approximation (9) has a linear convergence rate. Also, in contrast to previous Theorem 1, it does not have a pronounced asymptote. This effect is observed because the concept of zero-order oracle (10) does not imply an accumulation of adversarial stochastic noise in the bias, and also reduces the variance by the large batch size B . It is worth noting that the same convergence results of ZO-MB-SGD method (see Algorithm 1) are obtained using the approach of Gaussian approximation with a zero-order oracle concept (10). The proof of Theorem 2 for the kernel approximation approach can be found in more detail in Appendix F. Also see Appendix J for a detailed proof of the convergence rate of Algorithm 1, using a Gaussian gradient approximation with one-point feedback via the zero-order oracle concept (10).

5.2. Stochastic zero-order oracle with mixed adversarial noise

Another concept with one-point feedback is combining the gradient approximation (9) (from Subsection 5.1, adversarial stochastic noise) with adversarial deterministic noise (e.g., from zero-order oracle (6)). Then the gradient approximation (7) with one-point feedback takes the following form:

$$\tilde{\mathbf{g}}(x, \xi, \mathbf{e}) = d \frac{\tilde{f}(x + \gamma r \mathbf{e}, \xi_1) - \tilde{f}(x - \gamma r \mathbf{e}, \xi_2)}{2\gamma} K(r) \mathbf{e}, \quad (11)$$

where the zero-order oracle

$$\tilde{f}(x, \xi_1) = f(x) + \xi_1 + \delta(x), \quad (12)$$

and $\xi_1 \neq \xi_2$ are adversarial stochastic noises such that $\mathbb{E}[\xi_1^2] \leq \tilde{\Delta}^2$ and $\mathbb{E}[\xi_2^2] \leq \tilde{\Delta}^2$, $\tilde{\Delta} \geq 0$, and $\delta(x)$ is adversarial deterministic noise, $|\delta(x)| \leq \Delta$, $\Delta \geq 0$ is a level of noise. The random variables ξ_1 and ξ_2 are independent from \mathbf{e} and r . Also, this concept does not require the assumption of zero-mean ξ_1 and ξ_2 . It is enough that $\mathbb{E}[\xi_1 \mathbf{e}] = 0$ and $\mathbb{E}[\xi_2 \mathbf{e}] = 0$. It is worth noting that the

approximation of gradient (11) is also a gradient approximation with one-point feedback as in Subsection 5.1. This approximation implies calculating the value of objective function on different realizations $\xi_1 \neq \xi_2$. In the case when $\xi_1 = \xi_2$ we can say that the gradient approximations considered in this subsection and in the discussion of Theorem 1 are identical. The following theorem presents the convergence results of ZO-MB-SGD method with the gradient approximation (11) via a zero-order oracle (12).

Theorem 3 *Let function $f(x)$ satisfy Assumptions 1 - 2 and the gradient approximation (11) satisfy Assumptions 3 - 4 then by choosing the step size $\eta \leq \frac{1}{(M+1)L_2}$ there exists parameters*

$$M \leq \beta^3 d, \quad \sigma^2 \leq \beta^3 d^2 L_2^2 \gamma^2 + \frac{\beta^3 d^2 (\Delta^2 + \tilde{\Delta}^2)}{\gamma^2}, \quad m = 0, \quad \zeta^2 \leq \beta^2 d^2 L_\beta^2 \gamma^{2(\beta-1)} + \frac{\beta^2 d^2 \Delta^2}{\gamma^2}$$

such that Algorithm 1 with gradient approximation (11) and $\gamma \leq \mathcal{O}(\Delta^{1/\beta})$ achieves the following error floor

$$\mathbb{E}f(x_N) - f^* = \max \left\{ \varepsilon, \mathcal{O} \left(d^2 \Delta^{\frac{2(\beta-1)}{\beta}} \right) \right\},$$

where f^* is the solution to problem (8).

The results of Theorem 3 show that the Algorithm 1 with gradient approximation (11), which is corrupted by adversarial deterministic and stochastic noises, converges to the following asymptote $\mathcal{O} \left(d^2 \Delta^{\frac{2(\beta-1)}{\beta}} \right)$ with linear rate. This result can be restated: the maximum permissible level of adversarial deterministic noise to achieve ε -accuracy is $\Delta \leq \mathcal{O} \left(\varepsilon^{\frac{\beta}{2(\beta-1)}} d^{\frac{-\beta}{\beta-1}} \right)$. As a result, we indicate exactly adversarial deterministic noise, because the adversarial stochastic noise does not accumulate in the bias as well as in the variance (if the batch size is large enough). In contrast to stochastic noise, it is the adversarial deterministic noise that defines the level of asymptote (since it accumulates in the bias). Note that in this zero-order oracle concept, too, gradient approximation with "kernel-based" approach achieves a better error floor than Gaussian approach $\mathcal{O}(d^2 \Delta)$. See Appendix G and K for a proof of convergence results through kernel and Gaussian approximation approaches, respectively.

6. IMPROVED ESTIMATES OF TABLE 2

In this section, we show how the estimates obtained in this paper on the oracle T complexity, as well as on the maximum noise level Δ , change using an improved analysis (presented in [67]) of the bias and second moment estimates of the Kernel approximation of the gradient, see e.g., (7).

Chronologically speaking, this paper was submitted to arXiv in May 2023 and is the first paper in which a gradient-free algorithm (via Kernel approximation) is proposed for solving including non-convex optimization problems satisfying the Polyak–Lojasiewicz condition. However, we believe that it is not correct not to mention the work [67], which appeared on arXiv in June 2023, in which the authors proposed an improved analysis for the Kernel approximation and showed an improved estimate on the total number of oracle calls. Despite the superiority in terms of oracle complexity, the results of this paper are independently interesting and partially state of the art, at least in terms of iteration complexity (since we achieve optimal estimates in the case $\beta = 2$, as well as the

Table 3: Comparison of iteration complexity $\#N$ oracle complexity, $\#T$ and maximum noise level $\#\Delta$ of the algorithm with the following approaches: Kernel and Gaussian approximations for zero-order oracle cases ① : $\tilde{f}(x) = f(x) + \delta(x)$, ② : $\tilde{f}(x, \xi) = f(x, \xi) + \delta(x)$, ③ : $\tilde{f}(x, \xi) = f(x) + \xi$, ④ : $\tilde{f}(x, \xi) = f(x) + \xi + \delta(x)$. In both approaches, the dependence of the estimates on the batch size B is presented. Notation: β = order of smoothness, ε = accuracy of the solution to the problem, d = dimension, μ = Polyak-Lojasiewicz condition constant. **Updated Table 2 with new analysis from paper [67]**

Case	Batch Size	Gaussian approximation			Kernel approximation		
		$\#N$	$\#T$	$\#\Delta$	$\#N$	$\#T$	$\#\Delta$
①, ②, ④	$B \in [1, \beta^3 d]$	$\tilde{O}(\frac{d}{B}\mu^{-1})$	$\tilde{O}(d\mu^{-1})$	$\Delta \leq \mu\varepsilon d^{-1}$	$\tilde{O}(\frac{d}{B}\mu^{-1})$	$\tilde{O}(d\mu^{-1})$	$\Delta \leq \mu\varepsilon d^{-1}$
	$B > \beta^3 d$	$\tilde{O}(\mu^{-1})$	$\max\{\tilde{O}(\mu^{-1}B), \tilde{O}(\frac{d^2\Delta^2}{\varepsilon^2\mu^\beta})\}$	$\Delta \leq \mu\varepsilon d^{-1}$	$\tilde{O}(\mu^{-1})$	$\max\{\tilde{O}(\mu^{-1}B), \tilde{O}(\frac{d^2\Delta^2}{\varepsilon^{\frac{\beta}{\mu-1}}\mu^{\frac{2\beta-1}{\beta-1}}})\}$	$\Delta \leq \frac{(\mu\varepsilon)^{\frac{\beta}{d}}}{d}$
③	$B \in [1, \beta^3 d]$	$\tilde{O}(\frac{d}{B}\mu^{-1})$	$\tilde{O}(d\mu^{-1})$	$\Delta \leq \mu\varepsilon d^{-1/2}$	$\tilde{O}(\frac{d}{B}\mu^{-1})$	$\tilde{O}(d\mu^{-1})$	$\Delta \leq \mu\varepsilon d^{-1/2}$
	$B > \beta^3 d$	$\tilde{O}(\mu^{-1})$	$\max\{\tilde{O}(\mu^{-1}B), \tilde{O}(\frac{d^2\Delta^2}{\varepsilon^2\mu^\beta})\}$	$\Delta \leq \mu\varepsilon d^{-1}B^{1/2}$	$\tilde{O}(\mu^{-1})$	$\max\{\tilde{O}(\mu^{-1}B), \tilde{O}(\frac{d^2\Delta^2}{\varepsilon^{\frac{\beta}{\mu-1}}\mu^{\frac{2\beta-1}{\beta-1}}})\}$	$\Delta \leq \frac{(\mu\varepsilon)^{\frac{\beta}{d}}}{d}B^{1/2}$

best we know in the case $\beta > 2$), as well as in terms of maximum noise (since we explicitly derive $\Delta > 0$, at which we can still guarantee «good» convergence).

But applying the improved estimates on the bias

$$\|\mathbb{E}[\tilde{\mathbf{g}}(x, \xi, e)] - \nabla f(x)\|_2 \leq \kappa_\beta \frac{L}{(l-1)!} \cdot \frac{d}{d+\beta-1} \tau^{\beta-1},$$

and second moment

$$\mathbb{E}[\|\tilde{\mathbf{g}}(x, \xi, e)\|_2^2] \leq 4d\mathbb{E}[\|\nabla f(x)\|_2^2] + 4d\kappa L^2\tau^2 + \frac{\kappa d^2 \tilde{\Delta}^2}{\tau^2}.$$

from [67] (here the authors consider the case of a zero-order oracle (10) with stochastic noise $\mathbb{E}[\xi_1^2] \leq \tilde{\Delta}^2$, which is generalized to deterministic noise $|\delta(x)| \leq \Delta$ (see, for example, Appendix D)) to the analysis of the gradient-free algorithm proposed in this paper we obtain the following results (see Table 3, changes are highlighted). It is not difficult to see that the changes are mainly in the dimensionality of the problem d , namely the problem dimensionality in the obtained estimates does not depend on the order of smoothness β . Thus, at this point it is safe to say that Table 3 represents the state-of-the-art results.

7. DISCUSSION

To create optimal algorithms for three criteria: oracle complexity, iteration complexity, and maximum level of admissible noise still allowing to guarantee certain accuracy usually use accelerated batched algorithms as a base. This way, optimal iterative and oracle complexities are achieved. However, for a smooth optimization problem with a PL condition [41] showed that the non-accelerated algorithms are optimal. That is why we chose the stochastic gradient descent method (the batched version), which is frequently used in almost all fields of machine learning. By doing so, we guaranteed the optimal oracle and iteration complexity. Also using known techniques from [18], we obtained an estimate for the maximum allowable adversarial noise level, which is the best among the estimates we know. Thus, the proposed Algorithm 1 can highly likely be considered

optimal by three criteria at once: the iteration number [41], the oracle complexity [20], and the maximum adversarial noise level [37].

Our theoretical results show that in the presence of adversarial noise, Algorithm 1 or its analog (e.g., from [36]) converges only to an asymptote (error floor). Since adaptive to adversarial noise is an important property of zero-order algorithm [see, e.g., 14], we considered two models of adversarial noise: deterministic and stochastic. From Theorems 2 and 3, we can see that adversarial stochastic noise does not accumulate in the bias. This is an essential difference between the two noise models.

The theoretical results show the advantage of our approach over the approach that used approximation with forward finite difference via Gaussian smoothing from [36]. They studied a SGD with biased gradient in an optimization problem under PL condition. But there exists for the smooth case a randomized approximation [18], e.g., via L_2 -randomization, which was not considered in [36]. A natural question arises: Is Algorithm 1 superior to the gradient-free counterpart of Algorithm with approximation of the gradient via L_2 -randomization in optimization problems under the PL condition? After all, the only difference between these approaches is that our approach (kernel approximation) uses increased smoothness information. We explore this question in the experiments. And it turns out the answer to that question is positive (see more details Section 8). Thus, the Zero-Order Mini-batch Stochastic Gradient Descent method (see Algorithm 1) is robust for solving the problem under the PL condition.

8. EXPERIMENTS

In this section, we focus on verifying whether our theoretical bounds are aligned with the numerical performance of the Zero-Order Mini-batch SGD method. In particular, we compare Algorithm 1 with the gradient-free counterpart from [36] which uses Gaussian smoothing approximation instead of the exact gradient. In all tests we understand adversarial noise as a computational error (mantissa).

We consider a standard problem satisfying the Polyak–Lojasiewicz condition. Namely, the solution of a system of p nonlinear equations [40]. The optimization problem (1) have the following form:

$$\min_{x \in \mathbb{R}^d} f(x) := \|g(x)\|_2^2,$$

where $g(x) = 0$ is a system of p nonlinear equations such that $p \leq d$,

$$g(x) = C \sin(x) + D \cos(x) - b,$$

$x \in \mathbb{R}^d$, $C, D \in \mathbb{R}^{p \times d}$ and $b \in \mathbb{R}^p$. We use the weighted sums of the Legendre polynomials as the kernel $K(r)$. For example, we have the following values for $\beta = \{1, 2, 3, 4, 5, 6\}$ [11]:

$$\begin{aligned} K_\beta(r) &= 3r & \beta &= 1, 2; \\ K_\beta(r) &= \frac{15r}{4}(5 - 7r^3) & \beta &= 3, 4; \\ K_\beta(r) &= \frac{195r}{64}(99r^4 - 126r^2 + 35) & \beta &= 5, 6. \end{aligned}$$

In Figure 1, we compare Algorithm 1 ("Kernel" approximation) with the gradient-free counterpart of the algorithm from [36] ("Gaussian" approximation). We observe that Algorithm 1

significantly surpasses its counterpart in terms of convergence rate and error floor. We also see that when we increase batch size (e.g., from $B = 1$ to $B = 10$) convergence neighborhood of the ZO-MB-SGD method decreases.

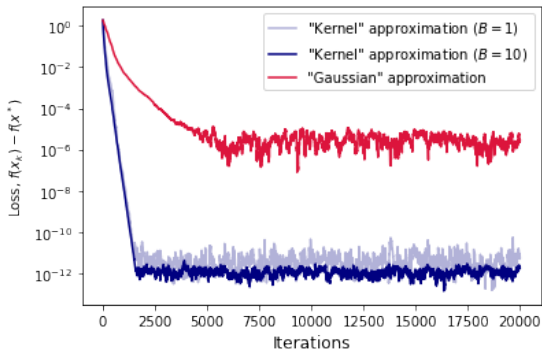


Figure 1: Comparing Algorithm 1 with gradient-free Algorithm from [36] and effect of the parameter B (batch size) on the convergence neighborhood of Zero-Order Mini-batch Stochastic Gradient Descent. Here we optimize $f(x)$ with the parameters: $d = 16$ (dimensional of problem), $p = 5$ (number of nonlinear equations), $\gamma = 0.01$ (smoothing parameter), $\eta = 0.01$ (fixed step size), $B = \{1, 10\}$ (batch size), $\beta = 3$ (order of smoothness).

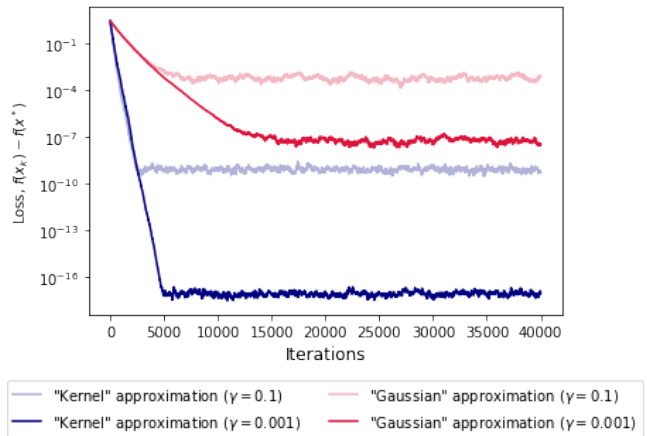


Figure 2: Effect of the parameter γ (smoothing parameter) on the error floor. Here we optimize $f(x)$ with parameters: $d = 128$ (dimensional of problem), $p = 16$ (number of nonlinear equations), $\gamma = \{0.1, 0.001\}$ (smoothing parameter), $\eta = 0.01$ (fixed step size), $B = 2$ (batch size), $\beta = 3$ (order of smoothness).

Figure 2 shows the effect of the smoothing parameter γ on the error floor. We can observe that the parameter γ directly affects the error floor. We also note that with different parameters γ the Zero-Order Mini-batch SGD (see Algorithm 1) retains its significant superiority over its counterpart.

To validate the robustness of our Algorithm (1), we now introduce the following gradient approximation via L_2 randomization which is well considered e.g. in [26]:

$$\tilde{\mathbf{g}}(x, \mathbf{e}) = d \frac{f(x + \gamma \mathbf{e}) - f(x - \gamma \mathbf{e})}{2\gamma} \mathbf{e}, \quad (13)$$

where \mathbf{e} is uniformly distributed on $S_2^d(1)$. Then Figure 3 demonstrates the effect of the smoothness order parameter on the convergence rate and the error floor, where " L_2 approximation" is a gradient-free counterpart of algorithm [36] with (13). We see that L_2 approximation has the same convergence rate as "Gaussian" approximation, but has a better error floor (i.e. shows efficiency of the randomization). We also observe that the advantages of the high order of smoothness improve the error floor of Algorithm 1. By comparing the L_2 approximation with Algorithm 1, we see the efficiency of the "Kernel" (in terms of the error floor), which uses information about highly smoothness. See Appendix A for the effect of the parameters d and p on the convergence rate and the error floor.

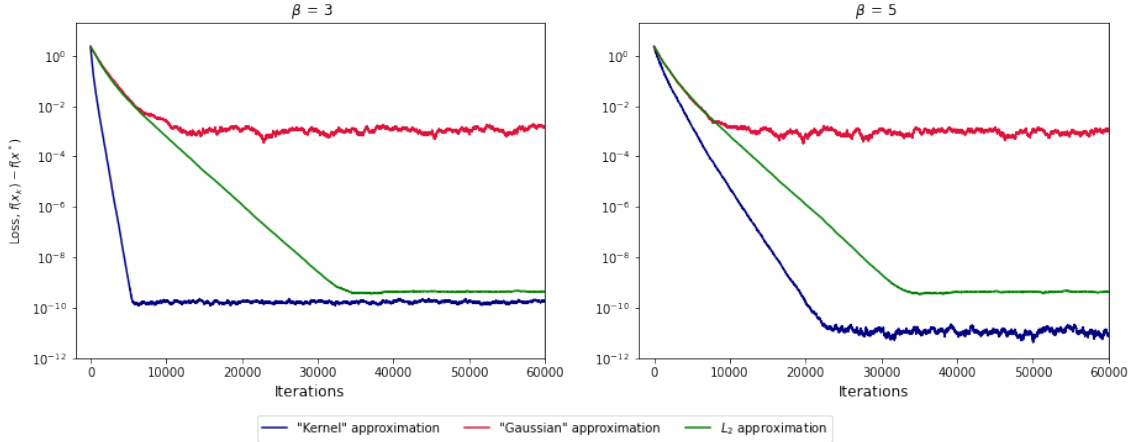


Figure 3: Effect of the parameter β on the convergence rate and the error floor. Here we optimize $f(x)$ with the parameters: $d = 256$ (dimensional of problem), $p = 32$ (number of nonlinear equations), $\gamma = 0.1$ (smoothing parameter), $\eta = 0.01$ (fixed step size), $B = 10$ (batch size), $\beta = \{3, 5\}$ (order of smoothness).

9. CONCLUSIONS

We proposed a Zero-Order Minibatch SGD method for solving smooth nonconvex optimization problems under the Polyak–Lojasiewicz condition. We generalized our results to stochastic problems under adversarial noise, considering the possible options: when two-point feedback is available and when only one-point feedback is available. Our algorithm showed efficiency on the standard (for problems under the Polyak–Lojasiewicz condition) example of solving a system of nonlinear equations.

CONFLICT OF INTEREST

The authors declare that they have no conflict if interest.

REFERENCES

1. *Rosenbrock H.* An automatic method for finding the greatest or least value of a function // The computer journal. — 1960. — Vol. 3, no. 3. — P. 175–184.
2. Structured evolution with compact architectures for scalable policy optimization / K. Choromanski [et al.] // International Conference on Machine Learning. — PMLR. 2018. — P. 970–978.
3. *Mania H., Guy A., Recht B.* Simple random search of static linear policies is competitive for reinforcement learning // Advances in Neural Information Processing Systems. — 2018. — Vol. 31.
4. Distributed Online and Bandit Convex Optimization / K. K. Patel [et al.] // OPT 2022: Optimization for Machine Learning (NeurIPS 2022 Workshop). — 2022.
5. Gradient-Free Federated Learning Methods with l_1 and l_2 -Randomization for Non-Smooth Convex Stochastic Optimization Problems / A. Lobanov [et al.] // arXiv preprint arXiv:2211.10783. — 2022.
6. Zoo: Zeroth order optimization based black-box attacks to deep neural networks without training substitute models / P.-Y. Chen [et al.] // Proceedings of the 10th ACM workshop on artificial intelligence and security. — 2017. — P. 15–26.
7. Black-box generation of adversarial text sequences to evade deep learning classifiers / J. Gao [et al.] // 2018 IEEE Security and Privacy Workshops (SPW). — IEEE. 2018. — P. 50–56.
8. *Shamir O.* An optimal algorithm for bandit and zero-order convex optimization with two-point feedback // The Journal of Machine Learning Research. — 2017. — Vol. 18, no. 1. — P. 1703–1713.
9. *Lattimore T., Gyorgy A.* Improved regret for zeroth-order stochastic convex bandits // Conference on Learning Theory. — PMLR. 2021. — P. 2938–2964.
10. *Agarwal A., Dekel O., Xiao L.* Optimal Algorithms for Online Convex Optimization with Multi-Point Bandit Feedback. // *Colt*. — Citeseer. 2010. — P. 28–40.
11. *Bach F., Perchet V.* Highly-smooth zero-th order online optimization // Conference on Learning Theory. — PMLR. 2016. — P. 257–283.
12. A gradient estimator via L_1 -randomization for online zero-order optimization with two point feedback / A. Akhavan [et al.] // Advances in Neural Information Processing Systems. — 2022. — Vol. 35. — P. 7685–7696.
13. Convex optimization: Algorithms and complexity / S. Bubeck [et al.] // Foundations and Trends® in Machine Learning. — 2015. — Vol. 8, no. 3/4. — P. 231–357.
14. Learning supervised pagerank with gradient-based and gradient-free optimization methods / L. Bogolubsky [et al.] // Advances in neural information processing systems. — 2016. — Vol. 29.
15. *Hernández-Lobato J. M., Hoffman M. W., Ghahramani Z.* Predictive entropy search for efficient global optimization of black-box functions // Advances in neural information processing systems. — 2014. — Vol. 27.

16. *Nguyen A., Balasubramanian K.* Stochastic Zeroth-Order Functional Constrained Optimization: Oracle Complexity and Applications // INFORMS Journal on Optimization. — 2022.
17. *Conn A. R., Scheinberg K., Vicente L. N.* Introduction to derivative-free optimization. — SIAM, 2009.
18. Randomized gradient-free methods in convex optimization / A. Gasnikov [et al.] // arXiv preprint arXiv:2211.13566. — 2022.
19. *Wasan M. T.* Stochastic approximation. — Cambridge University Press, 2004.
20. Optimal rates for zero-order convex optimization: The power of two function evaluations / J. C. Duchi [et al.] // IEEE Transactions on Information Theory. — 2015. — Vol. 61, no. 5. — P. 2788–2806.
21. *Nesterov Y., Spokoiny V.* Random gradient-free minimization of convex functions // Foundations of Computational Mathematics. — 2017. — Vol. 17, no. 2. — P. 527–566.
22. Gradient-free proximal methods with inexact oracle for convex stochastic nonsmooth optimization problems on the simplex / A. V. Gasnikov [et al.] // Automation and Remote Control. — 2016. — Vol. 77. — P. 2018–2034.
23. *Polyak B. T., Tsybakov A. B.* Optimal order of accuracy of search algorithms in stochastic optimization // Problemy Peredachi Informatsii. — 1990. — Vol. 26, no. 2. — P. 45–53.
24. *Akhavan A., Pontil M., Tsybakov A.* Exploiting higher order smoothness in derivative-free optimization and continuous bandits // Advances in Neural Information Processing Systems. — 2020. — Vol. 33. — P. 9017–9027.
25. *Novitskii V., Gasnikov A.* Improved exploiting higher order smoothness in derivative-free optimization and continuous bandit // arXiv preprint arXiv:2101.03821. — 2021.
26. The power of first-order smooth optimization for black-box non-smooth problems / A. Gasnikov [et al.] // International Conference on Machine Learning. — PMLR. 2022. — P. 7241–7265.
27. *Karimireddy S. P. R., Stich S., Jaggi M.* Adaptive balancing of gradient and update computation times using global geometry and approximate subproblems // International Conference on Artificial Intelligence and Statistics. — PMLR. 2018. — P. 1204–1213.
28. *Ghadimi S., Lan G.* Stochastic first-and zeroth-order methods for nonconvex stochastic programming // SIAM Journal on Optimization. — 2013. — Vol. 23, no. 4. — P. 2341–2368.
29. *Hajinezhad D., Hong M., Garcia A.* Zeroth order nonconvex multi-agent optimization over networks // arXiv preprint arXiv:1710.09997. — 2017.
30. Zeroth-order stochastic variance reduction for nonconvex optimization / S. Liu [et al.] // Advances in Neural Information Processing Systems. — 2018. — Vol. 31.
31. *Polyak B. T.* Gradient methods for the minimisation of functionals // USSR Computational Mathematics and Mathematical Physics. — 1963. — Vol. 3, no. 4. — P. 864–878.
32. *Lojasiewicz S.* Une propriété topologique des sous-ensembles analytiques réels // Les équations aux dérivées partielles. — 1963. — Vol. 117. — P. 87–89.

33. *Karimi H., Nutini J., Schmidt M.* Linear convergence of gradient and proximal-gradient methods under the polyak-lojasiewicz condition // Joint European conference on machine learning and knowledge discovery in databases. — Springer. 2016. — P. 795–811.
34. *Belkin M.* Fit without fear: remarkable mathematical phenomena of deep learning through the prism of interpolation // Acta Numerica. — 2021. — Vol. 30. — P. 203–248.
35. *Polyak B. T., Kuruzov I. A., Stonyakin F. S.* Stopping Rules for Gradient Methods for Non-Convex Problems with Additive Noise in Gradient // arXiv preprint arXiv:2205.07544. — 2022.
36. *Ajalloeian A., Stich S. U.* On the convergence of SGD with biased gradients // arXiv preprint arXiv:2008.00051. — 2020.
37. *Risteski A., Li Y.* Algorithms and matching lower bounds for approximately-convex optimization // Advances in Neural Information Processing Systems. — 2016. — Vol. 29.
38. *Vasin A., Gasnikov A., Spokoiny V.* Stopping rules for accelerated gradient methods with additive noise in gradient // arXiv preprint.—2021.—<https://arxiv.org/abs/2102.02921>. — 2021.
39. One-point gradient-free methods for composite optimization with applications to distributed optimization / I. Stepanov [et al.] // arXiv preprint arXiv:2107.05951. — 2021.
40. *Kuruzov I. A., Stonyakin F. S., Alkousa M. S.* Gradient-Type Methods for Optimization Problems with Polyak-Lojasiewicz Condition: Early Stopping and Adaptivity to Inexactness Parameter // Advances in Optimization and Applications: 13th International Conference, OPTIMA 2022, Petrovac, Montenegro, September 26–30, 2022, Revised Selected Papers. — Springer. 2023. — P. 18–32.
41. *Yue P., Fang C., Lin Z.* On the Lower Bound of Minimizing Polyak-Lojasiewicz functions // arXiv preprint arXiv:2212.13551. — 2022.
42. *Zhang T.* Solving large scale linear prediction problems using stochastic gradient descent algorithms // Proceedings of the twenty-first international conference on Machine learning. — 2004. — P. 116.
43. *Bottou L.* Large-scale machine learning with stochastic gradient descent // Proceedings of COMPSTAT'2010. — Springer, 2010. — P. 177–186.
44. The convergence of sparsified gradient methods / D. Alistarh [et al.] // Advances in Neural Information Processing Systems. — 2018. — Vol. 31.
45. Gradient sparsification for communication-efficient distributed optimization / J. Wangni [et al.] // Advances in Neural Information Processing Systems. — 2018. — Vol. 31.
46. *Stich S. U.* Unified optimal analysis of the (stochastic) gradient method // arXiv preprint arXiv:1907.04232. — 2019.
47. *Varre A. V., Pillaud-Vivien L., Flammarion N.* Last iterate convergence of SGD for Least-Squares in the Interpolation regime. // Advances in Neural Information Processing Systems. — 2021. — Vol. 34. — P. 21581–21591.

48. The min-max complexity of distributed stochastic convex optimization with intermittent communication / B. E. Woodworth [et al.] // Conference on Learning Theory. — PMLR. 2021. — P. 4386–4437.
49. *Yuan H., Ma T.* Federated accelerated stochastic gradient descent // Advances in Neural Information Processing Systems. — 2020. — Vol. 33. — P. 5332–5344.
50. *Lan G.* An optimal method for stochastic composite optimization // Mathematical Programming. — 2012. — Vol. 133, no. 1. — P. 365–397.
51. *Scheinberg K.* Finite Difference Gradient Approximation: To Randomize or Not? // INFORMS Journal on Computing. — 2022. — Vol. 34, no. 5. — P. 2384–2388.
52. Linearly converging error compensated SGD / E. Gorbunov [et al.] // Advances in Neural Information Processing Systems. — 2020. — Vol. 33. — P. 20889–20900.
53. Is local SGD better than minibatch SGD? / B. Woodworth [et al.] // International Conference on Machine Learning. — PMLR. 2020. — P. 10334–10343.
54. *Mishchenko K., Khaled A., Richtárik P.* Random reshuffling: Simple analysis with vast improvements // Advances in Neural Information Processing Systems. — 2020. — Vol. 33. — P. 17309–17320.
55. Asynchronous Stochastic Optimization Robust to Arbitrary Delays / A. Cohen [et al.] // Advances in Neural Information Processing Systems. — 2021. — Vol. 34. — P. 9024–9035.
56. Clipped Stochastic Methods for Variational Inequalities with Heavy-Tailed Noise / E. Gorbunov [et al.] // arXiv preprint arXiv:2206.01095. — 2022.
57. Near-optimal high probability complexity bounds for non-smooth stochastic optimization with heavy-tailed noise / E. Gorbunov [et al.] // arXiv preprint arXiv:2106.05958. — 2021.
58. Noisy Zeroth-Order Optimization for Non-smooth Saddle Point Problems / D. Dvinskikh [et al.] // International Conference on Mathematical Optimization Theory and Operations Research. — Springer. 2022. — P. 18–33.
59. A geometric integration approach to smooth optimisation: Foundations of the discrete gradient method / M. J. Ehrhardt [et al.] // arXiv preprint arXiv:1805.06444. — 2018.
60. Derivative-free methods for policy optimization: Guarantees for linear quadratic systems / D. Malik [et al.] // The 22nd international conference on artificial intelligence and statistics. — PMLR. 2019. — P. 2916–2925.
61. *Luo X., Xu X.* Stochastic gradient-free descents // arXiv preprint arXiv:1912.13305. — 2019.
62. *Nesterov Y., Polyak B. T.* Cubic regularization of Newton method and its global performance // Mathematical Programming. — 2006. — Vol. 108, no. 1. — P. 177–205.
63. *Vaswani S., Bach F., Schmidt M.* Fast and faster convergence of sgd for over-parameterized models and an accelerated perceptron // The 22nd international conference on artificial intelligence and statistics. — PMLR. 2019. — P. 1195–1204.
64. *Hu B., Seiler P., Lessard L.* Analysis of biased stochastic gradient descent using sequential semidefinite programs // Mathematical Programming. — 2021. — Vol. 187, no. 1. — P. 383–408.

65. Modern efficient numerical approaches to regularized regression problems in application to traffic demands matrix calculation from link loads / A. Anikin [et al.] // Proceedings of International conference ITAS-2015. Russia, Sochi. — 2015.
66. Stochastic online optimization. Single-point and multi-point non-linear multi-armed bandits. Convex and strongly-convex case / A. V. Gasnikov [et al.] // Automation and remote control. — 2017. — Vol. 78, no. 2. — P. 224–234.
67. Gradient-free optimization of highly smooth functions: improved analysis and a new algorithm / A. Akhavan [et al.] // arXiv preprint arXiv:2306.02159. — 2023.

A. ANALYZING EFFECT OF PROBLEM PARAMETERS ON ALGORITHMS

In Figure 4 we show the effect of the parameters d and p on the convergence rate and error floor. We can observe the behavior of the algorithms via the following approaches of approximating the gradient oracle (see Definition 1): "Kernel" approximation, "Gaussian" approximation and L_2 approximation at different parameters of the dimensional of the problem and the number of nonlinear equations. When increasing the number of nonlinear equations, we see that all algorithms worsen the convergence rate, but do not change the error floor. And when increasing the dimensional of the problem, we observe that only the "Kernel" approximation and the L_2 approximation improve the error floor (i.e. shows the efficiency of the randomization). We note that Algorithm 1 surpasses its counterparts in all the above cases (at different parameters), thereby confirming that the Algorithm 1 is robust for solving the problem under the Polyak–Lojasiewicz condition in general.

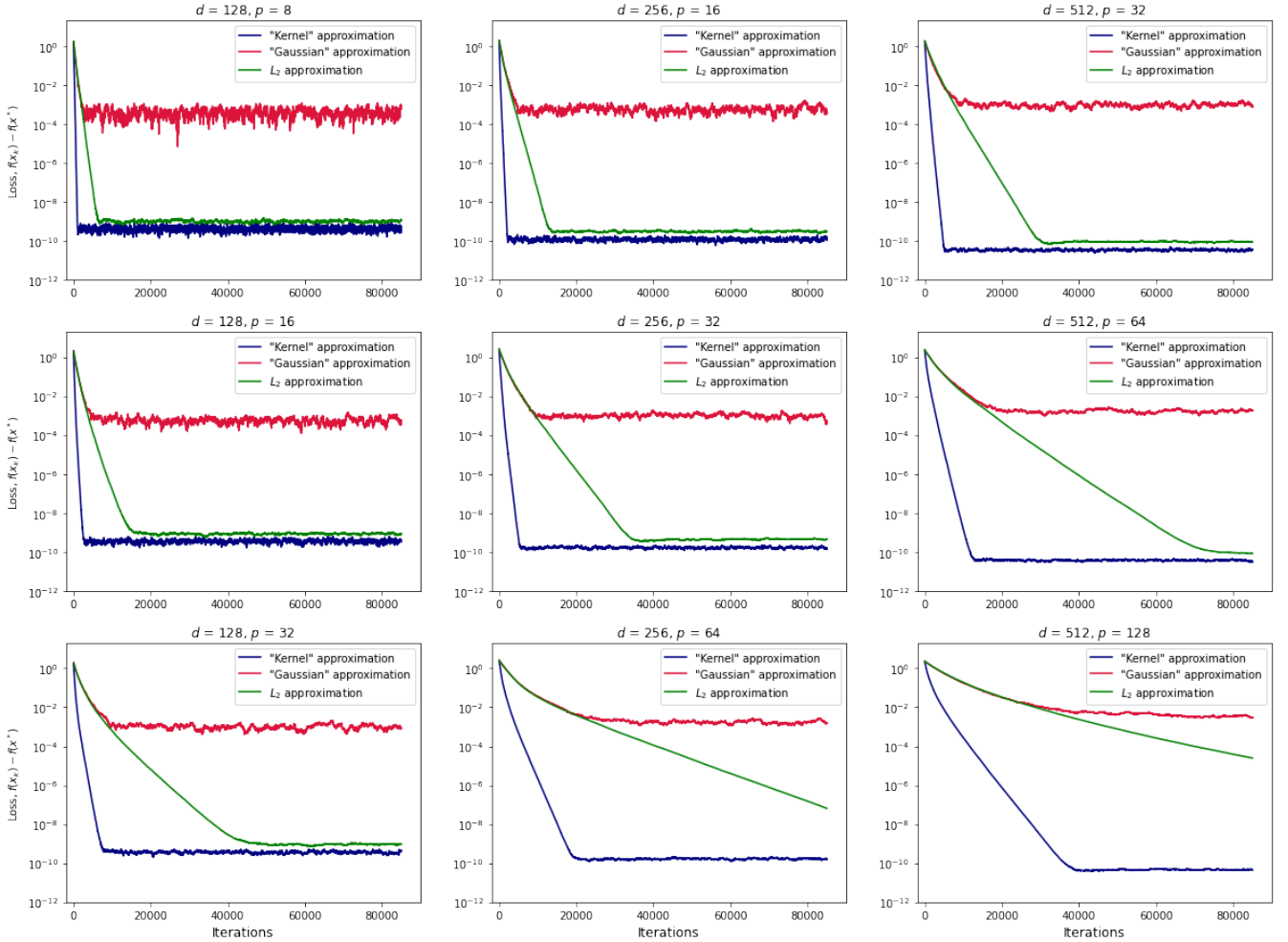


Figure 4: Effect of the parameters d (dimensional of problem) and p (number of nonlinear equations) on the convergence rate and the error floor. Here we optimize $f(x)$ with the following parameters: $d = \{128, 256, 512\}$ (dimensional of problem), $p = \{8, 16, 32, 64, 128\}$ (number of nonlinear equations), $\gamma = 0.1$ (smoothing parameter), $\eta = 0.01$ (fixed step size), $B = 10$ (batch size), $\beta = 3$ (order of smoothness).

B. AUXILIARY FACTS AND RESULTS

In this section we list the auxiliary facts and results that we use several times in our proofs.

B.1. Squared norm of the sum

For all $a_1, \dots, a_n \in \mathbb{R}^d$, where $n = \{2, 3\}$

$$\|a_1 + \dots + a_n\|_2^2 \leq n\|a_1\|_2^2 + \dots + n\|a_n\|_2^2. \quad (14)$$

B.2. Fact from concentration of the measure

Let \mathbf{e} is uniformly distributed on the Euclidean unit sphere, then, for $d \geq 8, \forall s \in \mathbb{R}^d$

$$\mathbb{E}_{\mathbf{e}} (\langle s, \mathbf{e} \rangle^2) \leq \frac{\|s\|_2^2}{d}. \quad (15)$$

B.3. Gaussian smoothing. Upper bounds for the moments

(Proved in Lemma 1, [21]). Let $\mathbf{u} \sim \mathcal{N}(0, 1)$ is a random Gaussian vector, then we have

$$\text{for } p \in [0, 2] : \quad \mathbb{E}_{\mathbf{u}} [\|\mathbf{u}\|_2^p] \leq d^{p/2}; \quad (16)$$

$$\text{for } p \geq 2 : \quad \mathbb{E}_{\mathbf{u}} [\|\mathbf{u}\|_2^p] \leq (p + d)^{p/2}. \quad (17)$$

B.4. Fact from Gaussian approximation

(Proved in Theorem 3, [21]). Let f is differentiable at $x \in \mathbb{R}^d$ and $\mathbf{u} \sim \mathcal{N}(0, 1)$ is normally distributed random Gaussian vector, then we have

$$\mathbb{E}_{\mathbf{u}} [\langle \nabla f(x), \mathbf{u} \rangle^2 \|\mathbf{u}\|_2^2] \leq (d + 4) \|\nabla f(x)\|_2^2. \quad (18)$$

B.5. Taylor expansion

Using the Taylor expansion we have

$$f(x + \gamma r e) = f(x) + \langle \nabla f(x), \gamma r e \rangle + \sum_{2 \leq |n| \leq l} \frac{(r\gamma)^{|n|}}{n!} D^{(n)} f(x) e^n + R(\gamma r e), \quad (19)$$

whereby assumption

$$|R(\gamma r e)| \leq L \|\gamma r e\|_2^\beta = L |r|^\beta \gamma^\beta. \quad (20)$$

B.6. Kernel property

If e is uniformly distributed on unit sphere we have $\mathbb{E}[e e^T] = (1/d) I_{d \times d}$, where $I_{d \times d}$ is the identity matrix. Therefore, using the facts $\mathbb{E}[r K(r)] = 1$ and $\mathbb{E}[r^{|n|} K(r)] = 0$ for $2 \leq |n| \leq l$ we have

$$\mathbb{E} \left[\frac{d}{\gamma} \left(\langle \nabla f(x), \gamma r e \rangle + \sum_{2 \leq |n| \leq l} \frac{(r\gamma)^{|n|}}{n!} D^{(n)} f(x) e^n \right) K(r) e \right] = \nabla f(x). \quad (21)$$

B.7. Bounds of the Weighted Sum of Legendre Polynomials

Let $\kappa_\beta = \int |u|^\beta |K(u)| du$ and set $\kappa = \int K^2(u) du$. Then if K be a weighted sum of Legendre polynomials, then it is proved in (see Appendix A.3, [11]) that κ_β and κ do not depend on d , they depend only on β , such that for $\beta \geq 1$:

$$\kappa_\beta \leq 2\sqrt{2}(\beta - 1), \quad (22)$$

$$\kappa \leq 3\beta^3. \quad (23)$$

C. THE CONVERGENCE RATE OF SGD WITH BIASED GRADIENT

Lemma 2 (Lemma 1) *Let $\{x_k\}_{k \geq 0}$ denote the iterates of Algorithm Mini-batch SGD, function f satisfy Assumptions 1-2 with $\beta = 2$, and the gradient oracle (3) satisfy Assumptions 3-4. Then there exists step size $\eta \leq \frac{1}{(M+1)L_2}$ such that it hold $\forall N \geq 0$ and an arbitrary batch size B*

$$\mathbb{E}[f(x_N)] - f^* \leq (1 - \eta\mu(1 - m))^N (f(x_0) - f^*) + \frac{\zeta^2}{2\mu(1 - m)} + \frac{\eta L_2 \sigma^2}{2B\mu(1 - m)}. \quad (24)$$

where L_2 is constant of the Lipschitz gradient such that $\|\nabla f(x) - \nabla f(y)\|_2 \leq L_2 \|x - y\|_2$.

P r o o f By the L_2 -smoothness of f and choice of the step size $\eta \leq \frac{1}{(M+1)L_2}$ we have

$$\begin{aligned} \mathbb{E}[f(x_{k+1})] &\leq f(x_k) + \langle \nabla f(x_k), x_{k+1} - x_k \rangle + \frac{L_2}{2} \|x_{k+1} - x_k\|_2^2 \\ &\leq f(x_k) - \eta \langle \nabla f(x_k), \mathbb{E}[\mathbf{g}_k] \rangle + \frac{\eta^2 L_2}{2} (\mathbb{E} [\|\mathbf{g}_k - \mathbb{E}[\mathbf{g}_k]\|_2^2] + \mathbb{E} [\|\mathbb{E}[\mathbf{g}_k]\|_2^2]) \\ &\stackrel{(3)}{=} f(x_k) - \eta \langle \nabla f(x_k), \nabla f(x_k) + \mathbf{b}(x_k) \rangle + \frac{\eta^2 L_2}{2} (\mathbb{E} [\|\mathbf{n}(x_k, \xi)\|_2^2] + \mathbb{E} [\|\nabla f(x_k) + \mathbf{b}(x_k)\|_2^2]) \\ &\stackrel{(4)}{\leq} f(x_k) - \eta \langle \nabla f(x_k), \nabla f(x_k) + \mathbf{b}(x_k) \rangle + \frac{\eta^2 L_2}{2} ((M + 1) \mathbb{E} [\|\nabla f(x_k) + \mathbf{b}(x_k)\|_2^2] + \sigma^2) \\ &= f(x_k) + \frac{\eta}{2} (\pm \|\nabla f(x_k)\|_2^2 - 2 \langle \nabla f(x_k), \nabla f(x_k) + \mathbf{b}(x_k) \rangle + \|\nabla f(x_k) + \mathbf{b}(x_k)\|_2^2) + \frac{\eta^2 L_2}{2} \sigma^2 \\ &= f(x_k) + \frac{\eta}{2} (-\|\nabla f(x_k)\|_2^2 + \|\mathbf{b}(x_k)\|_2^2) + \frac{\eta^2 L_2}{2} \sigma^2 \\ &\stackrel{(2),(5)}{\leq} (1 - \eta\mu(1 - m))(f(x_k) - f^*) + \frac{\eta \zeta^2}{2} + \frac{\eta^2 L_2}{2} \sigma^2 + f^*. \end{aligned} \quad (25)$$

Applying recursion to (25) and adding batching (with batch size B) we obtain

$$\mathbb{E}[f(x_N)] - f^* \leq (1 - \eta\mu(1 - m))^N (f(x_0) - f^*) + \frac{\zeta^2}{2\mu(1 - m)} + \frac{\eta L_2 \sigma^2}{2B\mu(1 - m)}.$$

Further structure of the appendix has close to block form. In particular, the first block (Appendices D-G) describes in detail the obtaining results for the approach proposed in this article to develop a gradient-free algorithm. Namely, the approach that implied using the "kernel-based"

approximation (Kernel approximation) instead of the gradient oracle (3). While the second block (Appendices H-K) provides a detailed description of obtaining results for the competing approach of developing gradient-free algorithms. In particular, the approach that uses Gaussian approximation instead of the gradient oracle (3). Each block considers different cases of a zero-order oracle. For instance, Appendix D and H consider zero-order oracle with adversarial deterministic noise described in Section 4. The stochastic case with the same concept of adversarial noise in a zero-order oracle generating a gradient approximation with two-point feedback, which is briefly described in the discussion of the results of Theorem 1, is discussed in Appendix E and I. Also in Appendix F and J, the case of a zero-order oracle with adversarial stochastic noise, described in Subsection 5.1, is considered in detail. The final case in each block (Appendix G and K) is a zero-order oracle case, combining adversarial deterministic and stochastic noise considered in Subsection 5.2. Note that the last two considered cases of zero-order oracle generate a gradient approximation with one-point feedback (see Section 5).

D. PROOF OF THEOREM 1

D.1. Kernel approximation

The "kernel-based" approximation of gradient has the following form (7):

$$\tilde{\mathbf{g}}(x, \mathbf{e}) = d \frac{\tilde{f}(x + \gamma r \mathbf{e}) - \tilde{f}(x - \gamma r \mathbf{e})}{2\gamma} K(r) \mathbf{e},$$

where $\tilde{f}(x)$ is defined in (6).

D.2. Bias square

By definition (7) we have

$$\begin{aligned} \|\mathbb{E}[\tilde{\mathbf{g}}(x, \mathbf{e})] - \nabla f(x)\|_2 &= \left\| \frac{d}{2\gamma} \mathbb{E} \left[\left(\tilde{f}(x + \gamma r \mathbf{e}) - \tilde{f}(x - \gamma r \mathbf{e}) \right) \mathbf{e} K(r) \right] - \nabla f(x) \right\|_2 \\ &\stackrel{(6)}{=} \left\| \frac{d}{2\gamma} \mathbb{E} \left[(f(x + \gamma r \mathbf{e}) - f(x - \gamma r \mathbf{e}) + \delta(x + \gamma r \mathbf{e}) - \delta(x - \gamma r \mathbf{e})) \mathbf{e} K(r) \right] - \nabla f(x) \right\|_2 \\ &\stackrel{(19)}{=} \left\| \frac{d}{\gamma} \mathbb{E} \left[\langle \nabla f(x), \gamma r \mathbf{e} \rangle + \sum_{2 \leq |n| \leq l \text{ odd}} \frac{(r\gamma)^{|n|}}{n!} D^{(n)} f(x) \mathbf{e}^n + \frac{R(\gamma r \mathbf{e}) - R(-\gamma r \mathbf{e})}{2} \right. \right. \\ &\quad \left. \left. + \frac{\delta(x + \gamma r \mathbf{e}) - \delta(x - \gamma r \mathbf{e})}{2} \right) \mathbf{e} K(r) \right] - \nabla f(x) \right\|_2 \\ &\stackrel{(21)}{=} \left\| \frac{d}{2\gamma} \mathbb{E} \left[(R(\gamma r \mathbf{e}) - R(-\gamma r \mathbf{e}) + \delta(x + \gamma r \mathbf{e}) - \delta(x - \gamma r \mathbf{e})) \mathbf{e} K(r) \right] \right\|_2 \\ &\leq \frac{d}{2\gamma} \mathbb{E} \left[\|R(\gamma r \mathbf{e}) - R(-\gamma r \mathbf{e}) + \delta(x + \gamma r \mathbf{e}) - \delta(x - \gamma r \mathbf{e})\| \|K(r)\| \right] \\ &\stackrel{(20)}{\leq} \kappa_\beta d \left(L_\beta \gamma^{\beta-1} + \frac{\Delta}{\gamma} \right). \quad / * \text{ the distribution of } \mathbf{e} \text{ is symmetric } * / \end{aligned}$$

Then, we can find bias square

$$\|\mathbf{b}(x)\|_2^2 = \|\mathbb{E}[\tilde{\mathbf{g}}(x, \mathbf{e})] - \nabla f(x)\|_2^2 \stackrel{(14)}{\leq} 2\kappa_\beta^2 d^2 L_\beta^2 \gamma^{2(\beta-1)} + 2\frac{\kappa_\beta^2 d^2 \Delta^2}{\gamma^2}. \quad (26)$$

D.3. Second moment

By definition (7) we have

$$\begin{aligned} \mathbb{E} [\|\tilde{\mathbf{g}}(x, e)\|_2^2] &= \frac{d^2}{4\gamma^2} \mathbb{E} [\|(f(x + \gamma re) - f(x - \gamma re) + \delta(x + \gamma re) - \delta(x - \gamma re)) K(r)e\|_2^2] \\ &= \frac{d^2}{4\gamma^2} \mathbb{E} [(f(x + \gamma re) - f(x - \gamma re) + 2\delta(x + \gamma re))^2 K^2(r) \|e\|_2^2] \\ &\stackrel{(14)}{\leq} \frac{d^2}{2\gamma^2} (\mathbb{E} [(f(x + \gamma re) - f(x - \gamma re))^2 K^2(r)] + 4\kappa\Delta^2). \end{aligned} \quad (27)$$

Using fact $\sqrt{\mathbb{E} [\|e\|_q^4]} \leq \min\{q-1, 16 \ln d - 8\} d^{2/q-1}$ to the first multiplier (27) we get

$$\begin{aligned} \frac{d^2}{2\gamma^2} \mathbb{E}_e [(f(x + \gamma re) - f(x - \gamma re))^2] &= \frac{d^2}{2\gamma^2} \mathbb{E}_e [((f(x + \gamma re) - f(x - \gamma re) \pm f(x) \pm 2\langle \nabla f(x), \gamma re \rangle)^2] \\ &\stackrel{(14)}{\leq} \frac{3d^2}{2\gamma^2} \mathbb{E}_e [(f(x + \gamma re) - f(x) - \langle \nabla f(x), \gamma re \rangle)^2 \\ &\quad + (f(x - \gamma re) - f(x) - \langle \nabla f(x), -\gamma re \rangle)^2 + 4\langle \nabla f(x), \gamma re \rangle^2] \end{aligned} \quad (28)$$

$$\begin{aligned} &\leq \frac{3d^2}{2\gamma^2} \mathbb{E}_e \left[\frac{L_2^2}{2} \|\gamma re\|_2^4 + 4\langle \nabla f(x), \gamma re \rangle^2 \right] \\ &\stackrel{(15)}{\leq} \frac{3d^2}{2\gamma^2} \left(\frac{L_2^2 \gamma^4}{2} \mathbb{E}_e [\|e\|_2^4] + \frac{4\gamma^2 \|\nabla f(x)\|_2^2}{d} \right) \\ &\leq 6d \|\nabla f(x)\|_2^2 + \frac{3d^2 L_2^2 \gamma^2}{4}, \end{aligned} \quad (29)$$

where (28) we obtained by applying the property of the Lipschitz continuous gradient.

By substituting (29) into (27) and using independence of e and r we obtain

$$\mathbb{E} [\|\tilde{\mathbf{g}}(x, e)\|_2^2] \leq \kappa \left(6d \|\nabla f(x)\|_2^2 + \frac{3d^2 L_2^2 \gamma^2}{4} + \frac{2d^2 \Delta^2}{\gamma^2} \right). \quad (30)$$

D.4. Convergence rate

From the inequalities (26) and (30) we can conclude that Assumption 3 holds with the choice

$$\sigma^2 \stackrel{(23)}{=} \mathcal{O}\left(\beta^3 d^2 L_2^2 \gamma^2 + \frac{\beta^3 d^2 \Delta^2}{\gamma^2}\right), \quad M = \mathcal{O}(\beta^3 d). \quad (31)$$

and that Assumption 4 also holds with the choice

$$m = 0, \zeta^2 \stackrel{(22)}{=} \mathcal{O}\left(\beta^2 d^2 L_\beta^2 \gamma^{2(\beta-1)} + \frac{\beta^2 d^2 \Delta^2}{\gamma^2}\right). \quad (32)$$

Now to find the asymptote to which the Algorithm 1 converges with the gradient approximation (7), substitute the parameters (31), (32) into the second and third terms (24) with $\eta = \mathcal{O}(1/M)$:

$$\mathbb{E}[f(x_N)] - f^* \leq \frac{\zeta^2}{2\mu(1-m)} + \frac{\eta L_2 \sigma^2}{2B\mu(1-m)} = \mathcal{O}\left(\beta^2 d^2 L_\beta^2 \gamma^{2(\beta-1)} + \frac{\beta^2 d^2 \Delta^2}{\gamma^2} + \frac{dL_2^2 \gamma^2}{B} + \frac{d\Delta^2}{B\gamma^2}\right).$$

Since B can be taken as large, the first two terms are responsible for the asymptote. We find the optimal smoothing parameter γ that minimizes the first two terms:

$$\mathbb{E}[f(x_N)] - f^* \leq \beta^2 d^2 \gamma^{2(\beta-1)} + \frac{\beta^2 d^2 \Delta^2}{\gamma^2} = \mathcal{O}\left(d^2 \Delta^{\frac{2(\beta-1)}{\beta}}\right), \quad (33)$$

where $\gamma = \Delta^{1/\beta}$ is optimal smoothing parameter. Then from (33) we can find the maximum noise level, assuming that $d^2 \Delta \leq \varepsilon$, for $\varepsilon > 0$ then we have

$$\Delta = \mathcal{O}\left(\varepsilon^{\frac{\beta}{2(\beta-1)}} d^{\frac{-\beta}{\beta-1}}\right).$$

E. PROOF OF CONVERGENCE OF ALGORITHM 1 WITH TWO-POINT FEEDBACK

E.1. Kernel approximation

The "kernel-based" approximation of gradient has the following form:

$$\tilde{\mathbf{g}}(x, \xi, \mathbf{e}) = d \frac{\tilde{f}(x + \gamma r \mathbf{e}, \xi) - \tilde{f}(x - \gamma r \mathbf{e}, \xi)}{2\gamma} K(r) \mathbf{e}, \quad (34)$$

where the zero-order oracle $\tilde{f}(x, \xi)$ is defined as follows:

$$\tilde{f}(x, \xi) = f(x, \xi) + \delta(x), \quad (35)$$

$\delta(x)$ is adversarial deterministic noise, $|\delta(x)| \leq \Delta$ is a level of noise.

E.2. Bias square

By definition (34) we have

$$\begin{aligned}
\|\mathbb{E}[\tilde{\mathbf{g}}(x, \xi, \mathbf{e})] - \nabla f(x)\|_2 &= \left\| \frac{d}{2\gamma} \mathbb{E} \left[\left(\tilde{f}(x + \gamma r \mathbf{e}, \xi) - \tilde{f}(x - \gamma r \mathbf{e}, \xi) \right) \mathbf{e} K(r) \right] - \nabla f(x) \right\|_2 \\
&\stackrel{(8),(35)}{=} \left\| \frac{d}{2\gamma} \mathbb{E} \left[(f(x + \gamma r \mathbf{e}) - f(x - \gamma r \mathbf{e}) + \delta(x + \gamma r \mathbf{e}) - \delta(x - \gamma r \mathbf{e})) \mathbf{e} K(r) \right] - \nabla f(x) \right\|_2 \\
&\stackrel{(19)}{=} \left\| \frac{d}{\gamma} \mathbb{E} \left[\langle \nabla f(x), \gamma r \mathbf{e} \rangle + \sum_{2 \leq |n| \leq l \text{ odd}} \frac{(r\gamma)^{|n|}}{n!} D^{(n)} f(x) \mathbf{e}^n + \frac{R(\gamma r \mathbf{e}) - R(-\gamma r \mathbf{e})}{2} \right. \right. \\
&\quad \left. \left. + \frac{\delta(x + \gamma r \mathbf{e}) - \delta(x - \gamma r \mathbf{e})}{2} \right) \mathbf{e} K(r) \right] - \nabla f(x) \right\|_2 \\
&\stackrel{(21)}{=} \left\| \frac{d}{2\gamma} \mathbb{E} \left[(R(\gamma r \mathbf{e}) - R(-\gamma r \mathbf{e}) + \delta(x + \gamma r \mathbf{e}) - \delta(x - \gamma r \mathbf{e})) \mathbf{e} K(r) \right] \right\|_2 \\
&\leq \frac{d}{2\gamma} \mathbb{E} \left[\|R(\gamma r \mathbf{e}) - R(-\gamma r \mathbf{e}) + \delta(x + \gamma r \mathbf{e}) - \delta(x - \gamma r \mathbf{e})\| |K(r)| \right] \\
&\stackrel{(20)}{\leq} \kappa_\beta d \left(L_\beta \gamma^{\beta-1} + \frac{\Delta}{\gamma} \right). \quad / * \text{ the distribution of } \mathbf{e} \text{ is symmetric} * /
\end{aligned}$$

Then, we can find bias square

$$\|\mathbf{b}(x)\|_2^2 = \|\mathbb{E}[\tilde{\mathbf{g}}(x, \xi, \mathbf{e})] - \nabla f(x)\|_2^2 \leq \left(\kappa_\beta d L_\beta \gamma^{\beta-1} + \frac{\kappa_\beta d \Delta}{\gamma} \right)^2 \stackrel{(14)}{\leq} \frac{2\kappa_\beta^2 d^2}{\gamma^2} (L_\beta^2 \gamma^{2\beta} + \Delta^2). \quad (36)$$

E.3. Second moment

By definition (34) we have

$$\begin{aligned}
\mathbb{E} \left[\|\tilde{\mathbf{g}}(x, \xi, \mathbf{e})\|_2^2 \right] &= \frac{d^2}{4\gamma^2} \mathbb{E} \left[\left\| (f(x + \gamma r \mathbf{e}, \xi) - f(x - \gamma r \mathbf{e}, \xi) + \delta(x + \gamma r \mathbf{e}) - \delta(x - \gamma r \mathbf{e})) K(r) \mathbf{e} \right\|_2^2 \right] \\
&= \frac{d^2}{4\gamma^2} \mathbb{E} \left[(f(x + \gamma r \mathbf{e}, \xi) - f(x - \gamma r \mathbf{e}, \xi) + 2\delta(x + \gamma r \mathbf{e}))^2 K^2(r) \|\mathbf{e}\|_2^2 \right] \\
&\stackrel{(14)}{\leq} \frac{d^2}{2\gamma^2} \left(\mathbb{E} \left[(f(x + \gamma r \mathbf{e}, \xi) - f(x - \gamma r \mathbf{e}, \xi))^2 K^2(r) \right] + 4\kappa \Delta^2 \right). \quad (37)
\end{aligned}$$

Using fact $\sqrt{\mathbb{E}[\|e\|_q^4]} \leq \min\{q-1, 16 \ln d - 8\} d^{2/q-1}$ to the first multiplier (37) we get

$$\begin{aligned}
& \frac{d^2}{2\gamma^2} \mathbb{E}_e [(f(x + \gamma re, \xi) - f(x - \gamma re, \xi))^2] \\
&= \frac{d^2}{2\gamma^2} \mathbb{E}_e [((f(x + \gamma re, \xi) - f(x - \gamma re, \xi) \pm f(x, \xi) \pm 2\langle \nabla f(x, \xi), \gamma re \rangle)^2)] \\
&\stackrel{(14)}{\leq} \frac{3d^2}{2\gamma^2} \mathbb{E}_e [(f(x + \gamma re, \xi) - f(x, \xi) - \langle \nabla f(x, \xi), \gamma re \rangle)^2 \\
&\quad + (f(x - \gamma re, \xi) - f(x, \xi) - \langle \nabla f(x, \xi), -\gamma re \rangle)^2 + 4\langle \nabla f(x, \xi), \gamma re \rangle^2] \\
&\leq \frac{3d^2}{2\gamma^2} \mathbb{E}_e \left[\frac{L_2^2}{2} \|\gamma re\|_2^4 + 4\langle \nabla f(x, \xi), \gamma re \rangle^2 \right] \tag{38}
\end{aligned}$$

$$\begin{aligned}
&\stackrel{(15)}{\leq} \frac{3d^2}{2\gamma^2} \left(\frac{L_2^2 \gamma^4}{2} \mathbb{E}_e [\|e\|_2^4] + \frac{4\gamma^2 \|\nabla f(x, \xi)\|_2^2}{d} \right) \\
&\leq 6d \|\nabla f(x, \xi)\|_2^2 + \frac{3d^2 L_2^2 \gamma^2}{4}, \tag{39}
\end{aligned}$$

where (38) we obtained by applying the property of the Lipschitz continuous gradient.

By substituting (39) into (37) and using independence of e and r we obtain

$$\mathbb{E} [\|\tilde{\mathbf{g}}(x, \xi, e)\|_2^2] \leq \kappa \left(6d \|\nabla f(x, \xi)\|_2^2 + \frac{3d^2 L_2^2 \gamma^2}{4} + \frac{2d^2 \Delta^2}{\gamma^2} \right). \tag{40}$$

E.4. Convergence rate

From the inequalities (36) and (40) we can conclude that Assumption 3 holds with the choice

$$\sigma^2 \stackrel{(23)}{=} \mathcal{O} \left(\beta^3 d^2 L_2^2 \gamma^2 + \frac{\beta^3 d^2 \Delta^2}{\gamma^2} \right), \quad M = \mathcal{O}(\beta^3 d). \tag{41}$$

and that Assumption 4 also holds with the choice

$$m = 0, \zeta^2 \stackrel{(22)}{=} \mathcal{O} \left(\beta^2 d^2 L_\beta^2 \gamma^{2(\beta-1)} + \frac{\beta^2 d^2 \Delta^2}{\gamma^2} \right). \tag{42}$$

Now to find the asymptote to which the Algorithm 1 converges with the gradient approximation (34), substitute the parameters (41), (42) into the second and third terms (24) with $\eta = \mathcal{O}(1/M)$:

$$\mathbb{E}[f(x_N)] - f^* \leq \frac{\zeta^2}{2\mu(1-m)} + \frac{\eta L_2 \sigma^2}{2B\mu(1-m)} = \mathcal{O} \left(\beta^2 d^2 L_\beta^2 \gamma^{2(\beta-1)} + \frac{\beta^2 d^2 \Delta^2}{\gamma^2} + \frac{dL_2^2 \gamma^2}{B} + \frac{d\Delta^2}{B\gamma^2} \right).$$

Since B can be taken as large, the first two terms are responsible for the asymptote. We find the optimal smoothing parameter γ that minimizes the first two terms:

$$\mathbb{E}[f(x_N)] - f^* \leq \beta^2 d^2 \gamma^{2(\beta-1)} + \frac{\beta^2 d^2 \Delta^2}{\gamma^2} = \mathcal{O} \left(d^2 \Delta^{\frac{2(\beta-1)}{\beta}} \right), \tag{43}$$

where $\gamma = \Delta^{1/\beta}$ is optimal smoothing parameter. Then from (43) we can find the maximum noise level, assuming that $d^2\Delta \leq \varepsilon$, for $\varepsilon > 0$ then we have

$$\Delta = \mathcal{O}\left(\varepsilon^{\frac{\beta}{2(\beta-1)}} d^{\frac{-\beta}{\beta-1}}\right).$$

F. PROOF OF THEOREM 2

F.1. Kernel approximation

The "kernel-based" approximation of gradient has the following form (9):

$$\tilde{\mathbf{g}}(x, \xi, \mathbf{e}) = d \frac{\tilde{f}(x + \gamma r \mathbf{e}, \xi_1) - \tilde{f}(x - \gamma r \mathbf{e}, \xi_2)}{2\gamma} K(r) \mathbf{e},$$

where $\tilde{f}(x, \xi)$ is defined in (10).

F.2. Bias square

By definition (9) we have

$$\begin{aligned} \|\mathbb{E}[\tilde{\mathbf{g}}(x, \xi, \mathbf{e})] - \nabla f(x)\|_2 &= \left\| \frac{d}{2\gamma} \mathbb{E} \left[\left(\tilde{f}(x + \gamma r \mathbf{e}, \xi) - \tilde{f}(x - \gamma r \mathbf{e}, \xi) \right) \mathbf{e} K(r) \right] - \nabla f(x) \right\|_2 \\ &\stackrel{(10)}{=} \left\| \frac{d}{2\gamma} \mathbb{E} \left[(f(x + \gamma r \mathbf{e}) - f(x - \gamma r \mathbf{e}) + \xi_1 - \xi_2) \mathbf{e} K(r) \right] - \nabla f(x) \right\|_2 \\ &\stackrel{(19)}{=} \left\| \frac{d}{\gamma} \mathbb{E} \left[\langle \nabla f(x), \gamma r \mathbf{e} \rangle + \sum_{2 \leq |n| \leq l \text{ odd}} \frac{(r\gamma)^{|n|}}{n!} D^{(n)} f(x) \mathbf{e}^n \right. \right. \\ &\quad \left. \left. + \frac{R(\gamma r \mathbf{e}) - R(-\gamma r \mathbf{e})}{2} \mathbf{e} K(r) \right] - \nabla f(x) \right\|_2 \\ &\stackrel{(21)}{=} \left\| \frac{d}{2\gamma} \mathbb{E} \left[(R(\gamma r \mathbf{e}) - R(-\gamma r \mathbf{e})) \mathbf{e} K(r) \right] \right\|_2 \\ &\leq \frac{d}{2\gamma} \mathbb{E} \left[|R(\gamma r \mathbf{e}) - R(-\gamma r \mathbf{e})| |K(r)| \right] \\ &\stackrel{(20)}{\leq} \kappa_\beta d L_\beta \gamma^{\beta-1}. \end{aligned}$$

Then, we can find bias square

$$\|\mathbf{b}(x)\|_2^2 = \|\mathbb{E}[\tilde{\mathbf{g}}(x, \xi, \mathbf{e})] - \nabla f(x)\|_2^2 \leq (\kappa_\beta d L_\beta \gamma^{\beta-1})^2 = \kappa_\beta^2 d^2 L_\beta^2 \gamma^{2(\beta-1)}. \quad (44)$$

F.3. Second moment

By definition (9) we have

$$\begin{aligned}
\mathbb{E} [\|\tilde{\mathbf{g}}(x, \xi, e)\|_2^2] &= \frac{d^2}{4\gamma^2} \mathbb{E} [\|(f(x + \gamma re) - f(x - \gamma re) + \xi_1 - \xi_2) K(r)e\|_2^2] \\
&= \frac{d^2}{4\gamma^2} \mathbb{E} [(f(x + \gamma re) - f(x - \gamma re) + \xi_1 - \xi_2)^2 K^2(r) \|e\|_2^2] \\
&\stackrel{(14)}{\leq} \frac{d^2}{2\gamma^2} \left(\mathbb{E} [(f(x + \gamma re) - f(x - \gamma re))^2 K^2(r)] + 4\kappa\tilde{\Delta}^2 \right). \tag{45}
\end{aligned}$$

Using fact $\sqrt{\mathbb{E} [\|e\|_q^4]} \leq \min\{q-1, 16 \ln d - 8\} d^{2/q-1}$ to the first multiplier (37) we get

$$\begin{aligned}
&\frac{d^2}{2\gamma^2} \mathbb{E}_e [(f(x + \gamma re) - f(x - \gamma re))^2] \\
&= \frac{d^2}{2\gamma^2} \mathbb{E}_e [((f(x + \gamma re) - f(x - \gamma re) \pm f(x) \pm 2\langle \nabla f(x), \gamma re \rangle)^2] \\
&\stackrel{(14)}{\leq} \frac{3d^2}{2\gamma^2} \mathbb{E}_e [(f(x + \gamma re) - f(x) - \langle \nabla f(x), \gamma re \rangle)^2 \\
&\quad + (f(x - \gamma re) - f(x) - \langle \nabla f(x), -\gamma re \rangle)^2 + 4\langle \nabla f(x), \gamma re \rangle^2] \\
&\leq \frac{3d^2}{2\gamma^2} \mathbb{E}_e \left[\frac{L_2^2}{2} \|\gamma re\|_2^4 + 4\langle \nabla f(x), \gamma re \rangle^2 \right] \tag{46}
\end{aligned}$$

$$\begin{aligned}
&\stackrel{(15)}{\leq} \frac{3d^2}{2\gamma^2} \left(\frac{L_2^2 \gamma^4}{2} \mathbb{E}_e [\|e\|_2^4] + \frac{4\gamma^2 \|\nabla f(x)\|_2^2}{d} \right) \\
&\leq 6d \|\nabla f(x)\|_2^2 + \frac{3d^2 L_2^2 \gamma^2}{4}, \tag{47}
\end{aligned}$$

where (46) we obtained by applying the property of the Lipschitz continuous gradient.

By substituting (47) into (45) and using independence of e and r we obtain

$$\mathbb{E} [\|\tilde{\mathbf{g}}(x, \xi, e)\|_2^2] \leq \kappa \left(6d \|\nabla f(x)\|_2^2 + \frac{3d^2 L_2^2 \gamma^2}{4} + \frac{2d^2 \tilde{\Delta}^2}{\gamma^2} \right). \tag{48}$$

F.4. Convergence rate

From the inequalities (44) and (48) we can conclude that Assumption 3 holds with the choice

$$\sigma^2 \stackrel{(23)}{=} \mathcal{O} \left(\beta^3 d^2 L_2^2 \gamma^2 + \frac{\beta^3 d^2 \tilde{\Delta}^2}{\gamma^2} \right), \quad M = \mathcal{O}(\beta^3 d). \tag{49}$$

and that Assumption 4 also holds with the choice

$$m = 0, \zeta^2 \stackrel{(22)}{=} \mathcal{O}(\beta^2 d^2 L_\beta^2 \gamma^{2(\beta-1)}). \tag{50}$$

Now to find the asymptote to which the Algorithm 1 converges with the gradient approximation (9), substitute the parameters (49), (50) into the second and third terms (24) with $\eta = \mathcal{O}(1/M)$:

$$\mathbb{E}[f(x_N)] - f^* \leq \frac{\zeta^2}{2\mu(1-m)} + \frac{\eta L_2 \sigma^2}{2B\mu(1-m)} = \mathcal{O}\left(\beta^2 d^2 L_\beta^2 \gamma^{2(\beta-1)} + \frac{dL_2^2 \gamma^2}{B} + \frac{d\tilde{\Delta}^2}{B\gamma^2}\right).$$

Since the batch size B can be taken large enough, there are no restrictions on the smoothing parameter γ . Therefore, using the concept of the zero-order oracle (10) Zero-order Mini-batch SGD (see Algorithm 1) with gradient approximation (9) can achieve the desired accuracy.

G. PROOF OF THEOREM 3

G.1. Kernel approximation

The "kernel-based" approximation of gradient has the following form (11):

$$\tilde{\mathbf{g}}(x, \xi, \mathbf{e}) = d \frac{\tilde{f}(x + \gamma r \mathbf{e}, \xi_1) - \tilde{f}(x - \gamma r \mathbf{e}, \xi_2)}{2\gamma} K(r) \mathbf{e},$$

where $\tilde{f}(x, \xi)$ is defined in (12).

G.2. Bias square

By definition (11) we have

$$\begin{aligned} \|\mathbb{E}[\tilde{\mathbf{g}}(x, \xi, \mathbf{e})] - \nabla f(x)\|_2 &= \left\| \frac{d}{2\gamma} \mathbb{E} \left[\left(\tilde{f}(x + \gamma r \mathbf{e}, \xi_1) - \tilde{f}(x - \gamma r \mathbf{e}, \xi_2) \right) \mathbf{e} K(r) \right] - \nabla f(x) \right\|_2 \\ &\stackrel{(8),(12)}{=} \left\| \frac{d}{2\gamma} \mathbb{E} \left[(f(x + \gamma r \mathbf{e}) - f(x - \gamma r \mathbf{e}) + \delta(x + \gamma r \mathbf{e}) - \delta(x - \gamma r \mathbf{e})) \mathbf{e} K(r) \right] - \nabla f(x) \right\|_2 \\ &\stackrel{(19)}{=} \left\| \frac{d}{\gamma} \mathbb{E} \left[\langle \nabla f(x), \gamma r \mathbf{e} \rangle + \sum_{2 \leq |n| \leq l \text{ odd}} \frac{(r\gamma)^{|n|}}{n!} D^{(n)} f(x) \mathbf{e}^n + \frac{R(\gamma r \mathbf{e}) - R(-\gamma r \mathbf{e})}{2} \right. \right. \\ &\quad \left. \left. + \frac{\delta(x + \gamma r \mathbf{e}) - \delta(x - \gamma r \mathbf{e})}{2} \right) \mathbf{e} K(r) \right] - \nabla f(x) \right\|_2 \\ &\stackrel{(21)}{=} \left\| \frac{d}{2\gamma} \mathbb{E} \left[(R(\gamma r \mathbf{e}) - R(-\gamma r \mathbf{e}) + \delta(x + \gamma r \mathbf{e}) - \delta(x - \gamma r \mathbf{e})) \mathbf{e} K(r) \right] \right\|_2 \\ &\leq \frac{d}{2\gamma} \mathbb{E} \left[|R(\gamma r \mathbf{e}) - R(-\gamma r \mathbf{e}) + \delta(x + \gamma r \mathbf{e}) - \delta(x - \gamma r \mathbf{e})| |K(r)| \right] \\ &\stackrel{(20)}{\leq} \kappa_\beta d \left(L_\beta \gamma^{\beta-1} + \frac{\Delta}{\gamma} \right). \quad / * \text{ the distribution of } \mathbf{e} \text{ is symmetric } * / \end{aligned}$$

Then, we can find bias square

$$\|\mathbf{b}(x)\|_2^2 = \|\mathbb{E}[\tilde{\mathbf{g}}(x, \xi, \mathbf{e})] - \nabla f(x)\|_2^2 \stackrel{(14)}{\leq} 2\kappa_\beta^2 d^2 L_\beta^2 \gamma^{2(\beta-1)} + 2 \frac{\kappa_\beta^2 d^2 \Delta^2}{\gamma^2}. \quad (51)$$

G.3. Second moment

By definition (11) we have

$$\begin{aligned}
\mathbb{E} [\|\tilde{\mathbf{g}}(x, \xi, e)\|_2^2] &= \frac{d^2}{4\gamma^2} \mathbb{E} [\|(f(x + \gamma re) - f(x - \gamma re) + \delta(x + \gamma re) - \delta(x - \gamma re) + \xi_1 - \xi_2) K(r)e\|_2^2] \\
&= \frac{d^2}{4\gamma^2} \mathbb{E} [(f(x + \gamma re) - f(x - \gamma re) + 2\delta(x + \gamma re) + \xi_1 - \xi_2)^2 K^2(r) \|e\|_2^2] \\
&\stackrel{(14)}{\leq} \frac{d^2}{2\gamma^2} \left(\mathbb{E} [(f(x + \gamma re) - f(x - \gamma re))^2 K^2(r)] + 4\kappa\Delta^2 + 4\kappa\tilde{\Delta}^2 \right). \tag{52}
\end{aligned}$$

Using fact $\sqrt{\mathbb{E} [\|e\|_q^4]} \leq \min \{q - 1, 16 \ln d - 8\} d^{2/q-1}$ to the first multiplier (52) we get

$$\begin{aligned}
&\frac{d^2}{2\gamma^2} \mathbb{E}_e [(f(x + \gamma re) - f(x - \gamma re))^2] \\
&= \frac{d^2}{2\gamma^2} \mathbb{E}_e [((f(x + \gamma re) - f(x - \gamma re) \pm f(x) \pm 2\langle \nabla f(x), \gamma re \rangle)^2] \\
&\stackrel{(14)}{\leq} \frac{3d^2}{2\gamma^2} \mathbb{E}_e [(f(x + \gamma re) - f(x) - \langle \nabla f(x), \gamma re \rangle)^2 \\
&\quad + (f(x - \gamma re) - f(x) - \langle \nabla f(x), -\gamma re \rangle)^2 + 4\langle \nabla f(x), \gamma re \rangle^2] \\
&\leq \frac{3d^2}{2\gamma^2} \mathbb{E}_e \left[\frac{L_2^2}{2} \|\gamma re\|_2^4 + 4\langle \nabla f(x), \gamma re \rangle^2 \right] \tag{53} \\
&\stackrel{(15)}{\leq} \frac{3d^2}{2\gamma^2} \left(\frac{L_2^2 \gamma^4}{2} \mathbb{E}_e [\|e\|_2^4] + \frac{4\gamma^2 \|\nabla f(x)\|_2^2}{d} \right)
\end{aligned}$$

$$\leq 6d \|\nabla f(x)\|_2^2 + \frac{3d^2 L_2^2 \gamma^2}{4}, \tag{54}$$

where (53) we obtained by applying the property of the Lipschitz continuous gradient.

By substituting (54) into (52) and using independence of e and r we obtain

$$\mathbb{E} [\|\tilde{\mathbf{g}}(x, \xi, e)\|_2^2] \leq \kappa \left(6d \|\nabla f(x, \xi)\|_2^2 + \frac{3d^2 L_2^2 \gamma^2}{4} + \frac{2d^2 (\Delta^2 + \tilde{\Delta}^2)}{\gamma^2} \right). \tag{55}$$

G.4. Convergence rate

From the inequalities (51) and (55) we can conclude that Assumption 3 holds with the choice

$$\sigma^2 \stackrel{(23)}{=} \mathcal{O} \left(\beta^3 d^2 L_2^2 \gamma^2 + \frac{\beta^3 d^2 (\Delta^2 + \tilde{\Delta}^2)}{\gamma^2} \right), \quad M = \mathcal{O}(\beta^3 d). \tag{56}$$

and that Assumption 4 also holds with the choice

$$m = 0, \zeta^2 \stackrel{(22)}{=} \mathcal{O} \left(\beta^2 d^2 L_\beta^2 \gamma^{2(\beta-1)} + \frac{\beta^2 d^2 \Delta^2}{\gamma^2} \right). \tag{57}$$

Now to find the asymptote to which the Algorithm 1 converges with the gradient approximation (11), substitute the parameters (56), (57) into the second and third terms (24) with $\eta = \mathcal{O}(1/M)$:

$$\mathbb{E}[f(x_N)] - f^* \leq \frac{\zeta^2}{2\mu(1-m)} + \frac{\eta L_2 \sigma^2}{2B\mu(1-m)} = \mathcal{O} \left(\beta^2 d^2 L_\beta^2 \gamma^{2(\beta-1)} + \frac{\beta^2 d^2 \Delta^2}{\gamma^2} + \frac{d L_2^2 \gamma^2}{B} + \frac{d(\Delta^2 + \tilde{\Delta}^2)}{B\gamma^2} \right).$$

Since B can be taken as large, the first two terms are responsible for the asymptote. We find the optimal smoothing parameter γ that minimizes the first two terms:

$$\mathbb{E}[f(x_N)] - f^* \leq \beta^2 d^2 \gamma^{2(\beta-1)} + \frac{\beta^2 d^2 \Delta^2}{\gamma^2} = \mathcal{O} \left(d^2 \Delta^{\frac{2(\beta-1)}{\beta}} \right), \quad (58)$$

where $\gamma = \Delta^{1/\beta}$ is optimal smoothing parameter. Then from (58) we can find the maximum noise level, assuming that $d^2 \Delta \leq \varepsilon$, for $\varepsilon > 0$ then we have

$$\Delta = \mathcal{O} \left(\varepsilon^{\frac{\beta}{2(\beta-1)}} d^{\frac{-\beta}{\beta-1}} \right).$$

H. GAUSSIAN SMOOTHING APPROACH FOR THEOREM 1

H.1. Definition Gaussian approximation

Let $\mathbf{u} \sim \mathcal{N}(0, 1)$ is a random Gaussian vector, $\gamma > 0$ is smoothing parameter then Gaussian smoothing approximation has the following form:

$$\tilde{\mathbf{g}}(x, \mathbf{u}) = \frac{\tilde{f}(x + \gamma \mathbf{u}) - \tilde{f}(x)}{\gamma} \mathbf{u}, \quad (59)$$

where \tilde{f} is defined in (6).

H.2. Bias square

By definition Gaussian approximation (59) we have

$$\begin{aligned} \|\mathbb{E}[\tilde{\mathbf{g}}(x, \mathbf{u})] - \nabla f(x)\|_2 &= \left\| \frac{1}{\gamma} \mathbb{E} \left[\left(\tilde{f}(x + \gamma \mathbf{u}) - \tilde{f}(x) \right) \mathbf{u} \right] - \nabla f(x) \right\|_2 \\ &\stackrel{(6)}{=} \left\| \frac{1}{\gamma} \mathbb{E} \left[(f(x + \gamma \mathbf{u}) - f(x) + \delta(x + \gamma \mathbf{u}) - \delta(x)) \mathbf{u} \right] - \nabla f(x) \right\|_2 \\ &= \left\| \frac{1}{\gamma} \mathbb{E}_{\mathbf{u}} \left[(f(x + \gamma \mathbf{u}) - f(x) - \gamma \langle \nabla f(x), \mathbf{u} \rangle + \delta(x + \gamma \mathbf{u}) - \delta(x)) \mathbf{u} \right] \right\|_2 \\ &\leq L_2 \gamma \mathbb{E}_{\mathbf{u}} \|\mathbf{u}\|_2^3 + \frac{\Delta}{\gamma} \mathbb{E}_{\mathbf{u}} \|\mathbf{u}\|_2 \quad / * \text{ the distribution of } \mathbf{u} \text{ is symmetric} * / \\ &\stackrel{(16),(17)}{\leq} L_2 \gamma d^{3/2} + \frac{\Delta}{\gamma} d^{1/2}. \end{aligned}$$

Then we can find the bias square:

$$\|\mathbf{b}(x)\|_2^2 = \|\mathbb{E}[\tilde{\mathbf{g}}(x, \mathbf{u})] - \nabla f(x)\|_2^2 \leq \left(L_2 \gamma d^{3/2} + \frac{\Delta}{\gamma} d^{1/2} \right)^2 \stackrel{(14)}{\leq} 2L_2^2 \gamma^2 d^3 + 2 \frac{\Delta^2}{\gamma^2} d. \quad (60)$$

H.3. Second moment

By definition (59) we have

$$\begin{aligned}
\mathbb{E}\|\tilde{\mathbf{g}}(x, \mathbf{u})\|_2^2 &= \mathbb{E}\left[\left\|\frac{1}{\gamma}\left(\tilde{f}(x + \gamma\mathbf{u}) - \tilde{f}(x)\right)\mathbf{u}\right\|_2^2\right] \\
&\stackrel{(6)}{\leq} \frac{1}{\gamma^2}\mathbb{E}\left[\left(f(x + \gamma\mathbf{u}) - f(x) + \delta(x + \gamma\mathbf{u}) - \delta(x)\right)^2\|\mathbf{u}\|_2^2\right] \\
&\leq \frac{1}{\gamma^2}\mathbb{E}\left[\left(f(x + \gamma\mathbf{u}) - f(x) \pm \gamma\langle\nabla f(x), \mathbf{u}\rangle + \delta(x + \gamma\mathbf{u}) - \delta(x)\right)^2\|\mathbf{u}\|_2^2\right] \\
&\stackrel{(14)}{\leq} \frac{3}{\gamma^2}L_2^2\gamma^4\mathbb{E}_{\mathbf{u}}[\|\mathbf{u}\|_2^6] + 3\mathbb{E}_{\mathbf{u}}[\langle\nabla f(x), \mathbf{u}\rangle^2\|\mathbf{u}\|_2^2] + \frac{3\Delta^2}{\gamma^2}\mathbb{E}_{\mathbf{u}}[\|\mathbf{u}\|_2^2] \\
&\stackrel{(16),(17),(18)}{\leq} 3L_2^2\gamma^2d^3 + 3d\|\nabla f(x)\|_2^2 + \frac{3d\Delta^2}{\gamma^2}. \tag{61}
\end{aligned}$$

H.4. Convergence rate

From the inequalities (60) and (61) we can conclude that Assumption 3 holds with the choice

$$\sigma^2 = \mathcal{O}\left(\gamma^2d^3 + \frac{d\Delta^2}{\gamma^2}\right), \quad M = \mathcal{O}(d). \tag{62}$$

and that Assumption 4 also holds with the choice

$$m = 0, \zeta^2 = \mathcal{O}\left(\gamma^2d^3 + \frac{\Delta^2}{\gamma^2}d\right). \tag{63}$$

Now to find the asymptote to which the counterpart of Algorithm 1 converges with the approximation (59), substitute the parameters (62), (63) into the second and third terms (24) with $\eta = \mathcal{O}(1/M)$:

$$\mathbb{E}[f(x_N)] - f^* \leq \frac{\zeta^2}{2\mu(1-m)} + \frac{\eta L_2 \sigma^2}{2B\mu(1-m)} = \mathcal{O}\left(\gamma^2d^3 + \frac{\Delta^2}{\gamma^2}d + \frac{\gamma^2d^2}{B} + \frac{\Delta^2}{B\gamma^2}\right).$$

Since B can be taken as large, the first two terms are responsible for the asymptote. We find the optimal smoothing parameter γ that minimizes the first two terms:

$$\mathbb{E}[f(x_N)] - f^* \leq \gamma^2d^3 + \frac{\Delta^2}{\gamma^2}d = \mathcal{O}(\Delta d^2), \tag{64}$$

where $\gamma = \Delta^{1/2}d^{-1/2}$ is optimal smoothing parameter.

Then from (64) we can find the maximum noise level, assuming that $\Delta d^2 \leq \varepsilon$, for $\varepsilon > 0$ then we have

$$\Delta = \mathcal{O}\left(\frac{\varepsilon}{d^2}\right).$$

I. GAUSSIAN SMOOTHING APPROACH WITH TWO-POINT FEEDBACK

I.1. Definition Gaussian approximation

Let $\mathbf{u} \sim \mathcal{N}(0, 1)$ is a random Gaussian vector, $\gamma > 0$ is smoothing parameter then Gaussian smoothing approximation has the following form:

$$\tilde{\mathbf{g}}(x, \xi, \mathbf{u}) = \frac{\tilde{f}(x + \gamma\mathbf{u}, \xi) - \tilde{f}(x, \xi)}{\gamma} \mathbf{u}, \quad (65)$$

where \tilde{f} is defined in (35).

I.2. Bias square

By definition Gaussian approximation (65) we have

$$\begin{aligned} \|\mathbf{b}(x)\|_2 &= \|\mathbb{E}[\tilde{\mathbf{g}}(x, \xi, \mathbf{u})] - \nabla f(x)\|_2 \\ &= \left\| \frac{1}{\gamma} \mathbb{E} \left[\left(\tilde{f}(x + \gamma\mathbf{u}, \xi) - \tilde{f}(x, \xi) \right) \mathbf{u} \right] - \nabla f(x) \right\|_2 \\ &\stackrel{(35)}{=} \left\| \frac{1}{\gamma} \mathbb{E}_{\mathbf{u}} \left[\mathbb{E}_{\xi} [(f(x + \gamma\mathbf{u}, \xi)) - \mathbb{E}_{\xi} [f(x, \xi)] + \delta(x + \gamma\mathbf{u}) - \delta(x)] \mathbf{u} \right] - \nabla f(x) \right\|_2 \\ &\stackrel{(8)}{=} \left\| \frac{1}{\gamma} \mathbb{E}_{\mathbf{u}} [(f(x + \gamma\mathbf{u}) - f(x) + \delta(x + \gamma\mathbf{u}) - \delta(x)) \mathbf{u}] - \nabla f(x) \right\|_2 \\ &= \left\| \frac{1}{\gamma} \mathbb{E}_{\mathbf{u}} [(f(x + \gamma\mathbf{u}) - f(x) - \gamma \langle \nabla f(x), \mathbf{u} \rangle + \delta(x + \gamma\mathbf{u}) - \delta(x)) \mathbf{u}] \right\|_2 \\ &\leq L_2 \gamma \mathbb{E}_{\mathbf{u}} \|\mathbf{u}\|_2^3 + \frac{\Delta}{\gamma} \mathbb{E}_{\mathbf{u}} \|\mathbf{u}\|_2 \quad / * \text{ the distribution of } \mathbf{u} \text{ is symmetric} * / \\ &\stackrel{(16),(17)}{\leq} L_2 \gamma d^{3/2} + \frac{\Delta}{\gamma} d^{1/2}. \end{aligned}$$

Then we can find the bias square:

$$\|\mathbf{b}(x)\|_2^2 = \|\mathbb{E}[\tilde{\mathbf{g}}(x, \xi, \mathbf{u})] - \nabla f(x)\|_2^2 \leq \left(L_2 \gamma d^{3/2} + \frac{\Delta}{\gamma} d^{1/2} \right)^2 \stackrel{(14)}{\leq} 2L_2^2 \gamma^2 d^3 + 2 \frac{\Delta^2}{\gamma^2} d. \quad (66)$$

I.3. Second moment

By definition (65) we have

$$\begin{aligned}
\mathbb{E}\|\tilde{\mathbf{g}}(x, \xi, \mathbf{u})\|_2^2 &= \mathbb{E}\left[\left\|\frac{1}{\gamma}\left(\tilde{f}(x + \gamma\mathbf{u}, \xi) - \tilde{f}(x, \xi)\right)\mathbf{u}\right\|_2^2\right] \\
&\stackrel{(35)}{\leq} \frac{1}{\gamma^2}\mathbb{E}\left[(f(x + \gamma\mathbf{u}, \xi) - f(x, \xi) + \delta(x + \gamma\mathbf{u}) - \delta(x))^2\|\mathbf{u}\|_2^2\right] \\
&\leq \frac{1}{\gamma^2}\mathbb{E}\left[(f(x + \gamma\mathbf{u}, \xi) - f(x, \xi) \pm \gamma\langle\nabla f(x, \xi), \mathbf{u}\rangle + \delta(x + \gamma\mathbf{u}) - \delta(x))^2\|\mathbf{u}\|_2^2\right] \\
&\stackrel{(14)}{\leq} \frac{3}{\gamma^2}L_2^2\gamma^4\mathbb{E}_{\mathbf{u}}[\|\mathbf{u}\|_2^6] + 3\mathbb{E}_{\xi}\left[\mathbb{E}_{\mathbf{u}}[\langle\nabla f(x, \xi), \mathbf{u}\rangle^2\|\mathbf{u}\|_2^2]\right] + \frac{3\Delta^2}{\gamma^2}\mathbb{E}_{\mathbf{u}}[\|\mathbf{u}\|_2^2] \\
&\stackrel{(16),(17),(18)}{\leq} 3L_2^2\gamma^2d^3 + 3d\|\nabla f(x)\|_2^2 + \frac{3d\Delta^2}{\gamma^2}. \tag{67}
\end{aligned}$$

I.4. Convergence rate

From the inequalities (66) and (67) we can conclude that Assumption 3 holds with the choice

$$\sigma^2 = \mathcal{O}\left(\gamma^2d^3 + \frac{d\Delta^2}{\gamma^2}\right), \quad M = \mathcal{O}(d), \tag{68}$$

and that Assumption 4 also holds with the choice

$$m = 0, \zeta^2 = \mathcal{O}\left(\gamma^2d^3 + \frac{\Delta^2}{\gamma^2}d\right). \tag{69}$$

Now to find the asymptote to which the counterpart of Algorithm 1 converges with the approximation (65), substitute the parameters (68), (69) into the second and third terms (24) with $\eta = \mathcal{O}(1/M)$:

$$\mathbb{E}[f(x_N)] - f^* \leq \frac{\zeta^2}{2\mu(1-m)} + \frac{\eta L_2 \sigma^2}{2B\mu(1-m)} = \mathcal{O}\left(\gamma^2d^3 + \frac{\Delta^2}{\gamma^2}d + \frac{\gamma^2d^2}{B} + \frac{\Delta^2}{B\gamma^2}\right).$$

Since B can be taken as large, the first two terms are responsible for the asymptote. We find the optimal smoothing parameter γ that minimizes the first two terms:

$$\mathbb{E}[f(x_N)] - f^* \leq \gamma^2d^3 + \frac{\Delta^2}{\gamma^2}d = \mathcal{O}(\Delta d^2), \tag{70}$$

where $\gamma = \Delta^{1/2}d^{-1/2}$ is optimal smoothing parameter.

Then from (70) we can find the maximum noise level, assuming that $\Delta d^2 \leq \varepsilon$, for $\varepsilon > 0$ then we have

$$\Delta = \mathcal{O}\left(\frac{\varepsilon}{d^2}\right).$$

J. GAUSSIAN SMOOTHING APPROACH FOR THEOREM 2

J.1. Definition Gaussian approximation

Let $\mathbf{u} \sim \mathcal{N}(0, 1)$ is a random Gaussian vector, $\gamma > 0$ is smoothing parameter then Gaussian smoothing approximation has the following form:

$$\tilde{\mathbf{g}}(x, \xi, \mathbf{u}) = \frac{\tilde{f}(x + \gamma\mathbf{u}, \xi_1) - \tilde{f}(x, \xi_2)}{\gamma} \mathbf{u}, \quad (71)$$

where \tilde{f} is defined in (10).

J.2. Bias square

By definition Gaussian approximation (71) we have

$$\begin{aligned} \|\mathbb{E}[\tilde{\mathbf{g}}(x, \xi, \mathbf{u})] - \nabla f(x)\|_2 &= \left\| \frac{1}{\gamma} \mathbb{E} \left[\left(\tilde{f}(x + \gamma\mathbf{u}, \xi_1) - \tilde{f}(x, \xi_2) \right) \mathbf{u} \right] - \nabla f(x) \right\|_2 \\ &\stackrel{(10)}{=} \left\| \frac{1}{\gamma} \mathbb{E} [(f(x + \gamma\mathbf{u}) - f(x) + \xi_1 - \xi_2) \mathbf{u}] - \nabla f(x) \right\|_2 \\ &= \left\| \frac{1}{\gamma} \mathbb{E}_{\mathbf{u}} [(f(x + \gamma\mathbf{u}) - f(x) - \gamma \langle \nabla f(x), \mathbf{u} \rangle + \xi_1 - \xi_2) \mathbf{u}] \right\|_2 \\ &\leq L_2 \gamma \mathbb{E}_{\mathbf{u}} \|\mathbf{u}\|_2^3 \\ &\stackrel{(16),(17)}{\leq} L_2 \gamma d^{3/2}. \end{aligned} \quad (72)$$

where we receive (72) using that $\mathbb{E}[\xi_1 \mathbf{u}] = 0$ and $\mathbb{E}[\xi_2 \mathbf{u}] = 0$.

Then we can find the bias square:

$$\|\mathbf{b}(x)\|_2^2 = \|\mathbb{E}[\tilde{\mathbf{g}}(x, \xi, \mathbf{u})] - \nabla f(x)\|_2^2 \leq (L_2 \gamma d^{3/2})^2 = L_2^2 \gamma^2 d^3. \quad (73)$$

J.3. Second moment

By definition (71) we have

$$\begin{aligned} \mathbb{E} \|\tilde{\mathbf{g}}(x, \xi, \mathbf{u})\|_2^2 &= \mathbb{E} \left[\left\| \frac{1}{\gamma} \left(\tilde{f}(x + \gamma\mathbf{u}, \xi_1) - \tilde{f}(x, \xi_2) \right) \mathbf{u} \right\|_2^2 \right] \\ &\stackrel{(10)}{\leq} \frac{1}{\gamma^2} \mathbb{E} [(f(x + \gamma\mathbf{u}) - f(x) + \xi_1 - \xi_2)^2 \|\mathbf{u}\|_2^2] \\ &\leq \frac{1}{\gamma^2} \mathbb{E} [(f(x + \gamma\mathbf{u}) - f(x) \pm \gamma \langle \nabla f(x), \mathbf{u} \rangle + \xi_1 - \xi_2)^2 \|\mathbf{u}\|_2^2] \\ &\stackrel{(14)}{\leq} \frac{3}{\gamma^2} L_2^2 \gamma^4 \mathbb{E}_{\mathbf{u}} [\|\mathbf{u}\|_2^6] + 3 \mathbb{E}_{\mathbf{u}} [\langle \nabla f(x), \mathbf{u} \rangle^2 \|\mathbf{u}\|_2^2] + \frac{3\tilde{\Delta}^2}{\gamma^2} \mathbb{E} [\|\mathbf{u}\|_2^2] \\ &\stackrel{(16),(17),(18)}{\leq} 3L_2^2 \gamma^2 d^3 + 3d \|\nabla f(x)\|_2^2 + \frac{3d\tilde{\Delta}^2}{\gamma^2}. \end{aligned} \quad (74)$$

J.4. Convergence rate

From the inequalities (73) and (74) we can conclude that Assumption 3 holds with the choice

$$\sigma^2 = \mathcal{O}\left(\gamma^2 d^3 + \frac{d\Delta^2}{\gamma^2}\right), \quad M = \mathcal{O}(d). \quad (75)$$

and that Assumption 4 also holds with the choice

$$m = 0, \zeta^2 = \mathcal{O}(\gamma^2 d^3). \quad (76)$$

Now to find the asymptote to which the counterpart of Algorithm 1 converges with the approximation (71), substitute the parameters (75), (76) into the second and third terms (24) with $\eta = \mathcal{O}(1/M)$:

$$\mathbb{E}[f(x_N)] - f^* \leq \frac{\zeta^2}{2\mu(1-m)} + \frac{\eta L_2 \sigma^2}{2B\mu(1-m)} = \mathcal{O}\left(\gamma^2 d^3 + \frac{\gamma^2 d^2}{B} + \frac{\Delta^2}{B\gamma^2}\right).$$

Since the batch size B can be taken large enough, there are no restrictions on the smoothing parameter γ . Therefore, using the concept of the zero-order oracle (10) the gradient-free counterpart of ZO-MB-SGD (see Algorithm 1) with gradient approximation (71) can achieve the desired accuracy.

K. GAUSSIAN SMOOTHING APPROACH FOR THEOREM 3

K.1. Definition Gaussian approximation

Let $\mathbf{u} \sim \mathcal{N}(0, 1)$ is a random Gaussian vector, $\gamma > 0$ is smoothing parameter then Gaussian smoothing approximation has the following form:

$$\tilde{\mathbf{g}}(x, \xi, \mathbf{u}) = \frac{\tilde{f}(x + \gamma\mathbf{u}, \xi_1) - \tilde{f}(x, \xi_2)}{\gamma} \mathbf{u}, \quad (77)$$

where \tilde{f} is defined in (12).

K.2. Bias square

By definition Gaussian approximation (77) we have

$$\begin{aligned}
\|\mathbf{b}(x)\|_2 &= \|\mathbb{E}[\tilde{\mathbf{g}}(x, \xi, \mathbf{u})] - \nabla f(x)\|_2 \\
&= \left\| \frac{1}{\gamma} \mathbb{E} \left[\left(\tilde{f}(x + \gamma \mathbf{u}, \xi) - \tilde{f}(x, \xi) \right) \mathbf{u} \right] - \nabla f(x) \right\|_2 \\
&\stackrel{(12)}{=} \left\| \frac{1}{\gamma} \mathbb{E}_{\mathbf{u}} [(f(x + \gamma \mathbf{u}) - f(x) + \delta(x + \gamma \mathbf{u}) - \delta(x) + \xi_1 - \xi_2) \mathbf{u}] - \nabla f(x) \right\|_2 \\
&= \left\| \frac{1}{\gamma} \mathbb{E}_{\mathbf{u}} [(f(x + \gamma \mathbf{u}) - f(x) - \gamma \langle \nabla f(x), \mathbf{u} \rangle + \delta(x + \gamma \mathbf{u}) - \delta(x)) \mathbf{u}] \right\|_2 \\
&\leq L_2 \gamma \mathbb{E}_{\mathbf{u}} \|\mathbf{u}\|_2^3 + \frac{\Delta}{\gamma} \mathbb{E}_{\mathbf{u}} \|\mathbf{u}\|_2 \quad / * \text{ the distribution of } \mathbf{u} \text{ is symmetric } */ \\
&\stackrel{(16),(17)}{\leq} L_2 \gamma d^{3/2} + \frac{\Delta}{\gamma} d^{1/2}.
\end{aligned}$$

Then we can find the bias square:

$$\|\mathbf{b}(x)\|_2^2 = \|\mathbb{E}[\tilde{\mathbf{g}}(x, \xi, \mathbf{u})] - \nabla f(x)\|_2^2 \leq \left(L_2 \gamma d^{3/2} + \frac{\Delta}{\gamma} d^{1/2} \right)^2 \stackrel{(14)}{\leq} 2L_2^2 \gamma^2 d^3 + 2 \frac{\Delta^2}{\gamma^2} d. \quad (78)$$

K.3. Second moment

By definition (77) we have

$$\begin{aligned}
&\mathbb{E} \|\tilde{\mathbf{g}}(x, \xi, \mathbf{u})\|_2^2 \\
&= \mathbb{E} \left[\left\| \frac{1}{\gamma} \left(\tilde{f}(x + \gamma \mathbf{u}, \xi) - \tilde{f}(x, \xi) \right) \mathbf{u} \right\|_2^2 \right] \\
&\stackrel{(35)}{\leq} \frac{1}{\gamma^2} \mathbb{E} [(f(x + \gamma \mathbf{u}) - f(x) + \delta(x + \gamma \mathbf{u}) - \delta(x) + \xi_1 - \xi_2)^2 \|\mathbf{u}\|_2^2] \\
&\leq \frac{1}{\gamma^2} \mathbb{E} [(f(x + \gamma \mathbf{u}) - f(x) \pm \gamma \langle \nabla f(x), \mathbf{u} \rangle + \delta(x + \gamma \mathbf{u}) - \delta(x) + \xi_1 - \xi_2)^2 \|\mathbf{u}\|_2^2] \\
&\stackrel{(14)}{\leq} \frac{3}{\gamma^2} L_2^2 \gamma^4 \mathbb{E}_{\mathbf{u}} [\|\mathbf{u}\|_2^6] + 3 \mathbb{E}_{\mathbf{u}} [\langle \nabla f(x, \xi), \mathbf{u} \rangle^2 \|\mathbf{u}\|_2^2] + \frac{3(\Delta^2 + \tilde{\Delta}^2)}{\gamma^2} \mathbb{E}_{\mathbf{u}} [\|\mathbf{u}\|_2^2] \\
&\stackrel{(16),(17),(18)}{\leq} 3L_2^2 \gamma^2 d^3 + 3d \|\nabla f(x)\|_2^2 + \frac{3d(\Delta^2 + \tilde{\Delta}^2)}{\gamma^2}. \quad (79)
\end{aligned}$$

K.4. Convergence rate

From the inequalities (78) and (79) we can conclude that Assumption 3 holds with the choice

$$\sigma^2 = \mathcal{O} \left(\gamma^2 d^3 + \frac{d(\Delta^2 + \tilde{\Delta}^2)}{\gamma^2} \right), \quad M = \mathcal{O}(d). \quad (80)$$

and that Assumption 4 also holds with the choice

$$m = 0, \zeta^2 = \mathcal{O}\left(\gamma^2 d^3 + \frac{\Delta^2}{\gamma^2} d\right). \quad (81)$$

Now to find the asymptote to which the counterpart of Algorithm 1 converges with the approximation (77), substitute the parameters (80), (81) into the second and third terms (24) with $\eta = \mathcal{O}(1/M)$:

$$\mathbb{E}[f(x_N)] - f^* \leq \frac{\zeta^2}{2\mu(1-m)} + \frac{\eta L_2 \sigma^2}{2B\mu(1-m)} = \mathcal{O}\left(\gamma^2 d^3 + d \frac{\Delta^2}{\gamma^2} + \frac{\gamma^2 d^2}{B} + \frac{\Delta^2 + \tilde{\Delta}^2}{B\gamma^2}\right).$$

Since B can be taken as large, the first two terms are responsible for the asymptote. We find the optimal smoothing parameter γ that minimizes the first two terms:

$$\mathbb{E}[f(x_N)] - f^* \leq \gamma^2 d^3 + d \frac{\Delta^2}{\gamma^2} = \mathcal{O}(d^2 \Delta), \quad (82)$$

where $\gamma = \Delta^{1/2} d^{-1/2}$ is optimal smoothing parameter.

Then from (82) we can find the maximum noise level, assuming that $\Delta d^2 \leq \varepsilon$, for $\varepsilon > 0$ then we have

$$\Delta = \mathcal{O}\left(\frac{\varepsilon}{d^2}\right).$$