

The “Black-Box” Optimization Problem: Zero-Order Accelerated Stochastic Method via Kernel Approximation

Aleksandr Lobanov · Nail Bashirov ·
Alexander Gasnikov

Received: date / Accepted: date

Abstract In this paper, we study the standard formulation of an optimization problem when the computation of gradient is not available. Such a problem can be classified as a “black box” optimization problem, since the oracle returns only the value of the objective function at the requested point, possibly with some stochastic noise. Assuming convex, and higher-order of smoothness of the objective function, this paper provides a zero-order accelerated stochastic gradient descent (ZO-AccSGD) method for solving this problem, which exploits the higher-order of smoothness information via kernel approximation. As theoretical results, we show that the ZO-AccSGD algorithm proposed in this paper improves the convergence results of state-of-the-art (SOTA) algorithms, namely the estimate of iteration complexity. In addition, our theo-

Aleksandr Lobanov, Corresponding author
Moscow Institute of Physics and Technology
9 Institutskiy per., Dolgoprudny, 141701, Russian Federation
Skolkovo Institute of Science and Technology
30 Bolshoy Boulevard, bld. 1, Moscow, 121205, Russian Federation
ISP RAS Research Center for Trusted Artificial Intelligence
25 A. Solzhenitsyn st., Moscow, 125047, Russian Federation
lobbsasha@mail.ru

Nail Bashirov
Moscow Institute of Physics and Technology
9 Institutskiy per., Dolgoprudny, 141701, Russian Federation
Institute for Information Transmission Problems
19 B. Karetny per., Moscow, 127051, Russian Federation
bashirov.nr@phystech.edu

Alexander Gasnikov
Moscow Institute of Physics and Technology
9 Institutskiy per., Dolgoprudny, 141701, Russian Federation
Innopolis University
1 Universitetskaya Str., Innopolis, 420500, Russian Federation
ISP RAS Research Center for Trusted Artificial Intelligence
25 A. Solzhenitsyn st., Moscow, 125047, Russian Federation
gasnikov@yandex.ru

retical analysis provides an estimate of the maximum allowable noise level at which the desired accuracy can be achieved. Validation of our theoretical results is demonstrated both on the model function and on functions of interest in the field of machine learning. We also provide a discussion in which we explain the results obtained and the superiority of the proposed algorithm over SOTA algorithms for solving the original problem.

Keywords Black-box optimization · Gradient-free methods · Kernel approximation · Maximum noise level

Mathematics Subject Classification (2000) 65K05 · 90C15 · 90C25

1 Introduction

Black-box optimization problems [27] (or also known as derivative-free optimization problems [14, 44]) arise when the gradient computation process is unavailable for some reason, e.g., the objective function $f(x)$ is not smooth [40, 16, 25, 30] or the process of computing the gradient $\nabla f(x)$ is too “expensive” compared to computing the value of the objective function $f(x)$ (in this case, the functions can be either smooth [1, 2] or of higher order of smoothness [7, 4, 32]). Moreover, there are often situations in practice [11] when the oracle returns a noisy value of the objective function (i.e., the value of the function $f(x)$ with some bounded noise ξ) at the requested point x , where the noise directly affects the “cost” of calling the oracle: the more inaccurate the oracle returns the value of the objective function (i.e., the greater the noise), the cheaper the oracle call. Such an oracle has a common “charactonym” name in the literature, namely a gradient-free oracle or a zero-order oracle [45]. Since this class of optimization problem has significant interest in settings such as federated learning [29, 39], distributed learning [5, 54, 33], overparameterized models [31] (in particular, in application problems such as hyperparameter tuning [28, 24], multi-armed bandits [17, 8], and many others [37]...), it is important to know and understand what approaches exist to solve this class of problem.

Apparently, the main way to solve the black-box problem is to apply gradient-free algorithms/zero-order methods to this problem. Among such methods there are two classes (or two approaches to the creation of gradient-free algorithms): the first is the class of evolutionary algorithms [50, 6, 22], which can often show their efficiency only empirically; the second is the class based on the advantages of first-order algorithms [26, 19], whose efficiency is provided in the form of theoretical estimates. Evolutionary algorithms are often used in the class of non-convex multimodal problems, in which the main goal is to find not a local but a global optimum. However, in the class of convex problems, it makes sense to use theoretically based algorithms, which often guarantee faster convergence by taking advantage of first-order algorithms.

The basic idea of creating efficient gradient-free algorithms for solving the convex black-box optimization problem is to use instead of the true gradient in first-order optimization algorithms some estimate of the gradient or also known

as a gradient approximation [18]. It is this seemingly simple idea that allows gradient-free algorithms to utilize the power of efficient first-order methods to solve the black box problem. However, it is important to correctly choose the algorithm and gradient approximation based on the original problem. Often, accelerated batched methods for solving corresponding optimization problems are chosen as efficient first-order algorithms, but there are also exceptions, e.g., for the class of problems satisfying the Polyak–Lojasiewicz condition, unaccelerated algorithms are already considered efficient (see [55] for more details). Regarding the question of the choice of gradient approximation: in [46] it is shown that central finite difference is a more preferable scheme for constructing a gradient approximation than forward finite difference. Also in [32] in the Experiments section, the authors have shown on a model practical experiment the advantage of using randomized approximations, in particular among the randomized approximations, they highlight l_2 randomization. However, this approximation is the most preferable for solving a smooth black-box optimization problem, but not for a problem with higher-order of smoothness. In 1990, B. Polyak and A. Tsybakov managed to propose such a gradient approximation, which takes into account the advantage of higher-order of smoothness of the function [42]. This gradient approximation is called the Kernel approximation. What distinguishes this approximation from l_2 randomization is the presence of a kernel by which the information about the higher-order of smoothness of the function is taken into account. And it was this paper [42] that became the starting point for the study of the solution of the black box problem with the assumption that the function has a high order of smoothness.

This paper investigates the improvement of the iteration complexity of gradient-free algorithms for solving a class of convex black-box optimization problems, assuming that the objective function has higher-order of smoothness. To create an optimal zero-order optimization method in terms of iteration complexity, we use the accelerated batched stochastic gradient descent method [52] (Nesterov–accelerated) as a basis. Among the gradient approximations that take into account the advantage of higher-order of smoothness, we choose the Kernel approximation because it is the one that requires only two calls to the gradient-free oracle per iteration (which guarantees a better estimate of oracle complexity), unlike the higher-order finite-difference gradient approximation [9]. However, the Kernel approximation is not a biased gradient estimator, so it is important to understand how noise is accounted for in the first-order algorithm. To this end, we generalize the accelerated first-order algorithm [52] to the case with a biased gradient oracle. Thus, to create a gradient-free algorithm, we base on the accelerated first-order method with a biased gradient oracle (see Section 3), using the Kernel approximation with a one-point zero-order oracle instead of the true gradient. In addition, we explicitly derive the estimate at the maximum noise level at which the desired accuracy can be achieved. Finally, we demonstrate our theoretical results on a model example.

1.1 Our contribution

Our contribution to this paper can be summarized as follows:

- We generalize the convergence results of the accelerated batched stochastic gradient descent algorithm (Nesterov–accelerated) [52] to the case with a biased gradient oracle;
- We provide a novel gradient-free optimization algorithm for solving a convex black-box optimization problem under an higher-order of smoothness condition: the zero-order accelerated stochastic gradient descent method (ZO-AccSGD, see Algorithm 1). This algorithm improves existing estimates on iteration complexity by working out the batching technique: $N = \mathcal{O}(\varepsilon^{-1/2})$. Moreover, using the well-known analysis of bias and second moment (variance) estimation, we were able to get rid of the smoothness order dependence β of the objective function in oracle complexity:
$$T = \mathcal{O}\left(\frac{d^2}{\varepsilon^{2+\frac{2}{\beta-1}}}\right);$$
- We provide an analysis that includes an elaboration on the maximum allowable noise level. We show that the noise level at which the desired accuracy is still achieved depends directly on the batch size B ;
- We confirm our theoretical results in Section “Experiments” by considering both a model problem (solving a system of linear equations), and problems of interest in machine learning.

1.2 Related works

Gradient-free oracle. The gradient approximation is typically a finite difference zero-order oracle. Therefore, every work on gradient-free oracle utilizes one or another zero-order oracle concept. For example, the works in [1] propose an oracle concept that returns the exact value of the objective function at the requested point: $\tilde{f} = f(x)$. This concept is intuitive and widely used, especially in tutorials such as introductions to optimization [41], etc. The following concept of gradient-free oracle, introduced in [16, 29], is oriented towards practical problems and is presented as follows: the oracle returns the value of the objective function at the requested point with some bounded deterministic noise $\tilde{f} = f(x) + \delta(x)$, where $|\delta(x)| \leq \Delta$. This concept applies well to deterministic optimization problems, but we can modernize it for a stochastic optimization problem [19]: $\tilde{f} = f(x, \xi) + \delta(x)$. This stochastic variant of the concept of a gradient-free oracle with bounded deterministic noise allows us to construct a gradient approximation depending on the availability of feedback. For example, if we can call the gradient-free oracle on one realization of the function twice, then the Kernel approximation with two-point feedback takes the following form: $\mathbf{g} = \frac{d}{2h} (f(x + h\mathbf{r}\mathbf{e}, \xi) + \delta(x + h\mathbf{r}\mathbf{e}) - f(x - h\mathbf{r}\mathbf{e}, \xi) - \delta(x - h\mathbf{r}\mathbf{e})) K(r)\mathbf{e}$. However, if we only have access to one-point feedback, i.e., we can call the oracle on one realization of the function only once, then the kernel approximation

takes the following form: $\mathbf{g} = \frac{d}{h} (f(x + h\mathbf{r}\mathbf{e}, \xi) + \delta(x + h\mathbf{r}\mathbf{e})) K(r)\mathbf{e}$. In addition, there is another concept of the gradient-free oracle that is quite controversial in gradient approximation, namely the kernel approximation [42, 4, 38, 32, 3] is as follows: $\mathbf{g} = \frac{d}{2h} (f(x + h\mathbf{r}\mathbf{e}) + \tilde{\xi}_1 - f(x - h\mathbf{r}\mathbf{e}) - \tilde{\xi}_2) K(r)\mathbf{e}$. Such a gradient approximation can quite rightly be called a one-point feedback approximation, although at first glance it is not even obvious that a gradient-free oracle that returns the value of the objective function at the requested point with some bounded stochastic noise $\tilde{f} = f(x) + \tilde{\xi}$ is a stochastic gradient-free oracle. However, if we consider the stochastic noise $\tilde{\xi}_1$ and $\tilde{\xi}_2$ as realizations of the function, it will be clear enough that such a gradient approximation can rightfully be called a one-point approximation, since the function is computed though twice in one iteration, but on different realizations. In our work, we also use the concept of a gradient-free oracle with stochastic noise, which generates a gradient approximation with one-point feedback.

Bounded gradient noises. Currently, there are a series of works [53, 43, 23, 10, 49, 32, 47, 48] that assume different constraints on gradient noise. For example, [43, 23] uses a standard and common in earlier works constraint on gradient noise, namely, they estimate some constant: $\mathbb{E} [\|\nabla f(x, \xi)\|^2] \leq \sigma^2$. However, there is some disadvantage of such a constraint, namely the large number constraint, that is, if the norm of the gradient decreases with the number of iterations, the estimate of the second moment will remain as large. To address this problem, some works [10, 49, 32] impose a constraint that is considered more adaptive: $\mathbb{E} [\|\nabla f(x, \xi)\|^2] \leq \rho \|\nabla f(x)\|^2 + \sigma^2$. There are also papers [47] that assume the strong growth condition is satisfied: $\mathbb{E} [\|\nabla f(x, \xi)\|^2] \leq \rho \|\nabla f(x)\|^2$. In the case where the model is overparameterized [53, 48], it is proposed to estimate the gradient noise as follows: $\mathbb{E} [\|\nabla f(x^*, \xi)\|^2] \leq \sigma_*^2$. The essential difference from the previous constraints is that the gradient estimate is evaluated at the solution point and depends directly on the solution of the problem, i.e., if $f^* = \min_x f(x)$ tends to zero, then also $\sigma_*^2 \leq Lf^*$ decreases. In our work, we use the following constraint on the noise of the biased gradient oracle $\mathbb{E} [\|\mathbf{g}(x, \xi)\|^2] \leq \rho \|\nabla f(x)\|^2 + \sigma^2$, since this constraint is adaptive and our approach in creating a gradient-free algorithm is based on the approach of the paper [52], which also addresses this constraint. The gradient oracle $\mathbf{g}(x, \xi)$ will be introduced in Subsection 2.2.

Iteration complexity. The study of a class of convex black-box optimization problems under the condition of higher-order of smoothness began with a 1990 paper [42], where a method of gradient estimation through the kernel and lower bound estimates of a gradient-free algorithm was proposed. At present, there already exist “pretty” results in this direction [7, 4, 38, 32, 3], which can be considered as “state of the art”. For example, in [4], the authors proposed a Zero-Order Stochastic Projected Gradient algorithm that used a central finite

difference kernel approximation and required the following iteration N (as well as oracle T) complexity to achieve a given accuracy ε : $T = N = \tilde{O}\left(\frac{d^{2+\frac{2}{\beta-1}}}{\varepsilon^{2+\frac{2}{\beta-1}}}\right)$. In another paper [38], the authors managed to improve this dimensionality estimate by using some “trick” in analyzing the bias of the gradient-free oracle: $T = N = \tilde{O}\left(\frac{d^{2+\frac{1}{\beta-1}}}{\varepsilon^{2+\frac{2}{\beta-1}}}\right)$. In a recent paper [3], the authors have managed to propose an improved analysis for estimating the bias and second moment (variance) of the gradient approximation, getting rid of the smoothness order dependence in the degree of dimensionality (in the strongly convex case). Moreover, it is not difficult to show that with the help of this analysis one can get rid of the dependence in the convex case as well (see [38] for the transformation from the strongly convex case to the convex case). However, these works focus on one of the three optimality criteria, namely oracle complexity. In our work, we use an improved analysis from [3] for gradient approximation to improve existing estimates of oracle complexity in the convex case, namely getting rid of the dependence of dimensionality on smoothness order, and by using an accelerated version of stochastic gradient descent and working out the batching technique we improve iteration complexity estimation.

1.3 Paper organization

This paper has the following structure. Section 2 introduces the formulation of the problem considered in this paper, as well as the main idea of its solution. Section 3 presents generalized results for the case of a biased oracle to solve the problem. The main result, which provides a novel gradient-free algorithm, can be found in Section 4. A discussion of the results is given in Section 5. Section 6 presents the experiments. While Section 7 concludes the paper.

2 Problem Formulation

In this section, we introduce the notations, definitions and assumptions used in our analysis to formulate the optimization problem. We also describe the main idea of our approach to solve the black-box optimization problem.

Notation. We use $\langle x, y \rangle := \sum_{i=1}^d x_i y_i$ to denote standard inner product of $x, y \in \mathbb{R}^d$, where x_i and y_i are the i -th component of x and y respectively. We denote Euclidean norm (l_2 -norm) in \mathbb{R}^d as $\|x\| = \|x\|_2 := \sqrt{\langle x, x \rangle}$. We use the following notation $B_2^d(r) := \{x \in \mathbb{R}^d : \|x\| \leq r\}$ to denote Euclidean ball (l_2 -ball) and $S_2^d(r) := \{x \in \mathbb{R}^d : \|x\| = r\}$ to denote Euclidean sphere. Operator $\mathbb{E}[\cdot]$ denotes full mathematical expectation.

We consider a standard optimization problem of the following form, which is commonly encountered in the literature, especially at the first acquaintance with optimization methods:

$$f^* = \min_{x \in Q \subseteq \mathbb{R}^d} f(x), \quad (1)$$

where $f : \mathbb{R}^d \rightarrow \mathbb{R}$ convex function that we want to minimize on the convex set Q . This general formulation is a broad class of optimization problems. To narrow down the class of optimization problems, we use a standard formulation of the problem and impose constraints on the function and the gradient oracle in the form of assumptions that will be used in our analysis throughout paper.

2.1 Assumptions on objective function

In our analysis presented in Section 3, we assume that function f is L -smooth.

Assumption 2.1 (L -smooth) *Function f is L -smooth if it holds*

$$f(y) \leq f(x) + \langle \nabla f(x), y - x \rangle + \frac{L}{2} \|y - x\|^2, \quad \forall x, y \in Q.$$

And in Section 4 our theoretical reasoning assumes that the objective function $f(x)$ is not just smooth, but has a higher order of smoothness.

Assumption 2.2 (Higher order smoothness) *Let l denote 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:*

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

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

The assumptions introduced in this subsection are standard and common in the literature in related works, e.g., see Assumption 2.1 in [34, 36], and Assumption 2.2 in the following works [42, 7, 3]. Moreover, it is not hard to see the connection between two assumptions, namely, in the case $\beta = 2 : L_2 = \frac{L}{2}$.

2.2 Assumptions on gradient oracle

Before presenting the assumptions on the gradient oracle, we introduce a formal definition, which is used extensively in the analysis for the convergence results of first-order algorithm in Section 3. Regarding convergence of zero-order algorithm, gradient-free oracle will be introduced in Subsection 4.1.

Definition 2.1 (Biased Gradient Oracle) A map $\mathbf{g} : \mathbb{R}^d \times \mathcal{D} \rightarrow \mathbb{R}^d$ s.t.

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

for a bias $\mathbf{b} : \mathbb{R}^d \rightarrow \mathbb{R}^d$ and unbiased stochastic gradient $\mathbb{E}[\nabla f(x, \xi)] = \nabla f(x)$.

We assume that the bias and gradient noise are bounded.

Assumption 2.3 (Bounded bias) *There exists constant $\delta \geq 0$ s.t. $\forall x \in \mathbb{R}^d$*

$$\|\mathbf{b}(x)\| = \|\mathbb{E}[\mathbf{g}(x, \xi)] - \nabla f(x)\| \leq \delta. \quad (3)$$

Assumption 2.4 (Bounded noise) *There exists constants $\rho, \sigma^2 \geq 0$ such that the more general condition of strong growth is satisfied $\forall x \in \mathbb{R}^d$*

$$\mathbb{E}[\|\mathbf{g}(x, \xi)\|^2] \leq \rho \|\nabla f(x)\|^2 + \sigma^2. \quad (4)$$

Assumptions 2.3 and 2.4 are not uncommon in works studying optimization algorithms with biased oracle (see Definition 2.1), such as [1, 32, 31].

2.3 The main idea of problem solving

The problem presented above does not strongly correspond to the black-box optimization problem, which is also stated in the title of the paper. This is done in order to present in Section 3 a first-order algorithm that has access to the noisy value of the gradient (see Definition 2.1). However, our approach to solving the optimization problem (1) when the gradient is still not available to the algorithm (black-box problems) is to create a gradient-free optimization algorithm based on and exploiting the power of the first-order method. Despite the fact that the original problem (1) is deterministic, we must rely on a first-order optimization algorithm that solves exactly the stochastic optimization problem $f(x) := \mathbb{E}[f(x, \xi)]$, since “stochasticity” is artificially created in the gradient approximation (see Subsection 4.1). Moreover, the Kernel approximation is a biased gradient estimator, so it is important to choose an algorithm that accounts for the imprecision in the gradient oracle. Thus, to summarize our approach, in Section 3 we generalize the SOTA results to the case with a biased gradient oracle (see Definition 2.1) that satisfies Assumptions 2.3 and 2.4, and use this first-order algorithm to create a gradient-free algorithm (see Section 4) for solving the black-box optimization problem under the condition of higher-order of smoothness of the objective function (see Assumption 2.2).

3 Generalization of Convergence Results for Accelerated SGD to the Biased Oracle

In this section, we provide the first-order algorithm on which the novel gradient-free method for solving the black-box optimization problem in Section 1 will

be based. Since this first-order algorithm must solve a stochastic optimization problem (due to the artificial “stochasticity” in the gradient approximation: $\mathbf{e} \in S_2^d(1)$, which will be introduced later), we reformulate the initial optimization problem as follows (1):

$$f^* = \min_{x \in Q \subseteq \mathbb{R}^d} \{f(x) := \mathbb{E}[f(x, \xi)]\}. \quad (5)$$

Next, before providing the convergence of the first-order algorithm with the biased gradient oracle, we present the known convergence results of accelerated stochastic gradient descent for solving problem (5).

3.1 Background

In 2019, the authors of [52] provided convergence results of Nesterov-accelerated Stochastic Gradient Descent [35] for the problem when the unbiased gradient oracle $\mathbf{g}(x, \xi) = \nabla f(x, \xi)$ (see Definition 2.1 with $\delta = 0$ in Assumption 2.3) satisfies the strong growth condition (see Assumption 2.4). In particular, this algorithm consists of the following update rules:

$$\begin{cases} x_{k+1} = y_k - \eta \mathbf{g}(y_k, \xi_k) \\ y_k = \alpha_k z_k + (1 - \alpha_k) x_k \\ z_{k+1} = \zeta_k z_k + (1 - \zeta_k) y_k - \gamma_k \eta \mathbf{g}(y_k, \xi_k). \end{cases}$$

And has the following convergence result in the case where $\sigma \neq 0$:

Lemma 3.1 ([52], **Theorem 1**) *Let the function f satisfy Assumption 2.1, and the unbiased gradient oracle $\mathbf{g}(x, \xi) = \nabla f(x, \xi)$ satisfies the strong growth condition (Assumption 2.4), then the accelerated Stochastic Gradient Descent by Nesterov with chosen parameters:*

$$\gamma_k = \frac{\tilde{\rho}^{-1} + \sqrt{\tilde{\rho}^{-2} + 4\gamma_{k-1}^2}}{2}; \quad a_{k+1} = \gamma_k \sqrt{\eta \tilde{\rho}}; \quad \alpha_k = \frac{\gamma_k \eta}{\gamma_k \eta + a_k^2}; \quad \eta = \frac{1}{\tilde{\rho} L},$$

where $\tilde{\rho} = \max\{1, \rho\}$, has the following rate of convergence:

$$\mathbb{E}[f(x_N)] - f^* \lesssim \frac{\tilde{\rho}^2 L R^2}{N^2} + \frac{N \sigma^2}{L \tilde{\rho}^2}.$$

This result is considered a state of the art for this problem formulation, however, as mentioned earlier, due to the Kernel approximation, which accumulates noise (i.e., has bias), we can't use this algorithm as a basis for creating a gradient-free optimization method. Therefore, in the next subsection, we extend the convergence results of the algorithm (Lemma 3.1) to the case where the gradient oracle $\mathbf{g}(x, \xi)$ (see Definition 2.1) can return a noisy gradient value $\nabla f(x, \xi) + \mathbf{b}(x)$, i.e., Assumption 2.3 is satisfied.

3.2 Accelerated SGD with biased gradients

In this subsection, we present the main result of Section 3, namely, we provide a first-order algorithm that we will base on in the next section to create a gradient-free algorithm. To achieve one of the main goals of our work, namely to improve and, if necessary, to obtain an optimal estimate of the iteration complexity N of the gradient-free optimization algorithm, we not only generalize the convergence result of Lemma 3.1 to the case with a biased oracle, but also get rid of the constant ρ from the first term, thus improving the estimate by the number of successive iterations of the accelerated Stochastic Gradient Descent. We improve the results of Lemma 3.1 in terms of iteration complexity by applying the batching technique. Thus, Accelerated Stochastic Gradient Descent (Nesterov acceleration) with a biased gradient oracle $\mathbf{g}(x, \xi)$ (see Definition 2.1) has the following convergence results presented in Theorem 3.1.

Theorem 3.1 (Biased AccSGD) *Let the function f satisfy Assumption 2.1, and the gradient oracle $\mathbf{g}(x, \xi)$ from Definition 2.1 satisfies Assumptions 2.3 and 2.4, then the accelerated Stochastic Gradient Descent with batching (B is a batch size) by Nesterov with $\rho_B = \max\{1, \frac{\rho}{B}\}$ and chosen parameters:*

$$\gamma_k = \frac{\rho_B^{-1} + \sqrt{\rho_B^{-2} + 4\gamma_{k-1}^2}}{2}; \quad a_{k+1} = \gamma_k \sqrt{\eta \rho_B}; \quad \alpha_k = \frac{\gamma_k \eta}{\gamma_k \eta + a_k^2}; \quad \eta = \frac{1}{\rho_B L}$$

has the following rate of convergence:

$$\mathbb{E}[f(x_N)] - f^* \lesssim \frac{\rho_B^2 L R^2}{N^2} + \frac{N \sigma^2}{\rho_B^2 L B} + \delta \tilde{R} + \frac{N}{L} \delta^2.$$

It is not hard to see that the convergence result of the accelerated batched algorithm presented in Theorem 3.1 is a generalization to the case when the gradient oracle $\mathbf{g}(x, \xi)$ (see Definition 2.1) returns a noisy gradient value. If we put $\delta = 0$ we get exactly the same convergence as in Lemma 3.1 up to the constants ρ from the condition of strongly growing (see Assumption 2.4) and B the size of the batch. The two terms accounting for noise accumulation are standard for the accelerated algorithm (see, e.g., [20, 15, 51]) and can be found, for example, by using the (δ, L) -oracle technique [19]. It is through the use of the batched technique in Theorem 3.1 that we will be able to achieve an optimal estimate on the iteration complexity $N = \mathcal{O}\left(\sqrt{\varepsilon^{-1} L R^2}\right)$ that will be obtained from the first term, since it dominates the second term for a sufficiently large value of the batch size $B \geq \rho$. This result allows us to use this accelerated batched algorithm to create a gradient-free method for solving the black-box optimization problem under the condition of higher-order of smoothness of the objective function f . A detailed proof of Theorem 3.1 can be found in Appendix B.

4 Main Results

In this section, we present the main result of our paper, namely a novel gradient-free method for solving the black-box optimization problem (1) with the condition that the objective function is not only smooth but also has a higher order of smoothness (i.e., the Assumption 2.2 is satisfied). Our approach to create a gradient-free algorithm is to choose and use a gradient estimate $\mathbf{g}(x, \mathbf{e})$ (an approximation of the gradient that will account for the higher-order of smoothness of the function) instead of the real gradient oracle $\mathbf{g}(x, \xi)$ see Definition 2.1 in the accelerated batched first-order method.

4.1 Gradient approximation

To solve a deterministic convex black-box optimization problem (1), where the “black box” plays the role of a gradient-free oracle \tilde{f} , which is formally defined as follows: we assume that the oracle \tilde{f} can only return the value of the objective function $f(x)$ at the requested point with some stochastic noise $\tilde{\xi}$:

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

where $\tilde{\xi}$ is stochastic, possibly adversarial, noise, $\mathbb{E}[\tilde{\xi}^2] \leq \Delta^2$. Then we use the so-called “Kernel-based approximation”, which was presented in 1990 in the paper [42] and was recognized years later in a number of papers [7, 4, 38, 32, 3], as an approximation of the gradient that takes into account the information about the higher-order of smoothness, has the following form:

$$\mathbf{g}(x, \mathbf{e}) = d \frac{f(x + h\mathbf{r}\mathbf{e}) + \tilde{\xi}_1 - f(x - h\mathbf{r}\mathbf{e}) - \tilde{\xi}_2}{2h} K(r)\mathbf{e}, \quad (7)$$

where $h > 0$ is a smoothing parameter, $\mathbf{e} \in S_2^d(1)$ is a vector uniformly distributed on the Euclidean unit sphere, r is a vector uniformly distributed on the interval $r \in [0, 1]$, $K : [-1, 1] \rightarrow \mathbb{R}$ is a kernel function that satisfies

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

This concept of noise is often found in the literature [4, 32], where the $\tilde{\xi}_1 \neq \tilde{\xi}_2$ such that $\mathbb{E}[\tilde{\xi}_1^2] \leq \Delta^2$ and $\mathbb{E}[\tilde{\xi}_2^2] \leq \Delta^2$, $\Delta \geq 0$ is level noise, and the random variables $\tilde{\xi}_1$ and $\tilde{\xi}_2$ are independent from \mathbf{e} and r . Also, this concept does not necessarily have to have a zero mean $\tilde{\xi}_1$ and $\tilde{\xi}_2$. It is enough that $\mathbb{E}[\tilde{\xi}_1 \mathbf{e}] = 0$ and $\mathbb{E}[\tilde{\xi}_2 \mathbf{e}] = 0$. Moreover, the gradient approximation (7) may at first glance appear to be an approximation with two-point feedback because of the structure of the central finite difference, but this is not entirely true. Since $\tilde{\xi}_1 \neq \tilde{\xi}_2$ and if we consider $\tilde{\xi}_1$ and $\tilde{\xi}_2$ as concrete realizations of the objective function $f(x)$, it is clear that the function cannot be computed on the same realization twice per iteration. Thus, the approximation with this concept of gradient-free oracle (6) is an approximation with one-point feedback.

4.2 Zero-order accelerated stochastic gradient descent

Now that we have chosen the gradient approximation and the accelerated batched first-order method, we can present a novel gradient-free algorithm Zero-Order Accelerated Stochastic Gradient Descent (ZO-AccSGD), which is obtained by replacing the real gradient with the gradient approximation (7).

Algorithm 1 Zero-Order Accelerated Stochastic Gradient Descent

Input: iteration number N , batch size B , Kernel $K : [-1, 1] \rightarrow \mathbb{R}$, step size η , smoothing parameter h , $x_0 = y_0 = z_0 \in \mathbb{R}^d$, $\alpha_0 = \gamma_0 = 0$.

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 $[-1, 1]$ independently
2. Define $\mathbf{g}(x_k, \mathbf{e}_i) = d \frac{\tilde{f}(x_k + hr_i \mathbf{e}_i) - \tilde{f}(x_k - hr_i \mathbf{e}_i)}{2h} K(r_i) \mathbf{e}_i$ via (6)
3. Calculate $\mathbf{g}_k = \frac{1}{B} \sum_{i=1}^B \mathbf{g}(x_k, \mathbf{e}_i)$
4. $x_{k+1} \leftarrow y_k - \eta \mathbf{g}_k$
5. $z_{k+1} \leftarrow z_k - \gamma_k \eta \mathbf{g}_k$
6. $y_{k+1} \leftarrow \alpha_{k+1} z_{k+1} + (1 - \alpha_{k+1}) x_{k+1}$

end for

Return: x_N

To obtain the convergence rate of the Algorithm 1, we need to first estimate the bias $\|\mathbb{E}[\mathbf{g}(x, \mathbf{e})] - \nabla f(x)\|$ and second moment (variance) $\mathbb{E}[\|\mathbf{g}(x, \mathbf{e})\|^2]$ of the gradient approximation $\mathbf{g}(x, \mathbf{e})$ of (7). Then, by substituting these estimates into the convergence result of the first-order algorithm we plan to rely on (in our case it is the Biased Accelerated Stochastic Gradient Descent, see Theorem 3.1), in particular instead of σ^2 (see Assumption 2.4) from the second term we need to substitute the obtained estimate on the second moment (variance) $\mathbb{E}[\|\mathbf{g}(x, \mathbf{e})\|^2]$, and instead of δ (see Assumption 2.3) of the third and fourth terms we need to substitute the obtained estimate for the bias $\|\mathbb{E}[\mathbf{g}(x, \mathbf{e})] - \nabla f(x)\|$, we get the convergence rate of the novel gradient-free algorithm. Then, using the bias and second moment estimates for the gradient approximation that takes into account information about the higher order of smoothness (Kernel approximation (7)) presented in [3] we have the following convergence results for Zero-Order Accelerated Stochastic Gradient Descent. Note that in Theorem 4.1 we present convergence results for an arbitrary kernel function satisfying the conditions presented in Subsection 4.1.

Theorem 4.1 (Convergence results) *Let the function f satisfy Assumption 2.2 and the gradient approximation $\mathbf{g}(x, \mathbf{e})$ of (7) satisfies Assumptions 2.3 and 2.4, then Zero-Order Accelerated Stochastic Gradient Descent (see Algorithm 1) with $\rho_B = \max\{1, \frac{4d\kappa}{B}\}$, and with the chosen algorithm parameters:*

$$\gamma_k = \frac{\rho_B^{-1} + \sqrt{\rho_B^{-2} + 4\gamma_{k-1}^2}}{2}; \quad a_{k+1} = \gamma_k \sqrt{\eta \rho_B}; \quad \alpha_k = \frac{\gamma_k \eta}{\gamma_k \eta + a_k^2}; \quad \eta = \frac{1}{\rho_B L}$$

converges to the desired ε accuracy, $\mathbb{E}[f(x_N)] - f^* \leq \varepsilon$

- in the case $B \in [1, 4d\kappa]$, $h \lesssim \varepsilon^{3/4}$ and $\beta \geq \frac{7}{3}$ after

$$N = \mathcal{O}\left(\sqrt{\frac{d^2 LR^2}{B^2 \varepsilon}}\right); \quad T = N \cdot B = \mathcal{O}\left(\sqrt{\frac{d^2 LR^2}{\varepsilon}}\right)$$

number of iterations and gradient-free oracle calls, respectively, at

$$\Delta \lesssim \frac{\varepsilon^{3/2}}{\sqrt{d}} \quad \text{maximum noise level;}$$

- in the case $B > 4d\kappa$ and $h \lesssim \varepsilon^{1/(\beta-1)}$ after

$$N = \mathcal{O}\left(\sqrt{\frac{LR^2}{\varepsilon}}\right); \quad T = N \cdot B = \max\left\{\mathcal{O}\left(\sqrt{\frac{d^2 LR^2}{\varepsilon}}\right), \mathcal{O}\left(\frac{d^2 \Delta^2}{\varepsilon^{2+\frac{2}{\beta-1}}}\right)\right\}$$

number of iterations and gradient-free oracle calls, respectively, at

$$\Delta \lesssim \frac{\varepsilon^{\frac{3\beta+1}{4(\beta-1)}}}{d} B^{1/2} \quad \text{maximum noise level;}$$

From the results of Theorem 4.1, it is not difficult to see that at the batch size $B = 1$ the dimensionality factor d comes out in the iteration complexity. The dimensionality can be eliminated by batching, i.e. the larger the batch size B , the better the iteration complexity N becomes, in particular, starting from $B = 4d\kappa$ the iteration complexity completely gets rid of dimensionality and reaches the optimal estimate for the accelerated algorithm. It is worth noting that the maximum noise level when the batch size is $1 \leq B \leq 4d\kappa$, in particular when $\beta \geq \frac{7}{3}$ has, perhaps close to the optimal value for the smooth case (i.e., the case where the Assumption 2.1 holds): $\Delta \lesssim d^{-1/2} \varepsilon^{3/2}$. However, this estimate is invariant regardless of the order of smoothness. But there is a way to improve the maximum noise level at which the algorithm is still guaranteed to achieve the desired ε accuracy. If we take the size of the batches larger than $B > 4d\kappa$, then the maximum allowable noise level Δ will improve and, in particular, will depend on the order of smoothness β . That is, the maximum noise level can be maximized in two ways: taking a larger batch size or using a function with higher-order of smoothness. However, improving the maximum noise level entails a deterioration of the oracle complexity, but has no effect on the iteration complexity $N \sim \varepsilon^{-1/2}$. It is not difficult to see that in any case considered, our results outperform all known results (see subsection Related works), in particular, we improve the iterative complexity estimate, as well as get rid of the dimensionality dependence of oracle complexity, and finally we present estimates of the maximum noise level as a function of batch size.

Proof First, we write the bias and the second moment of the gradient approximation (7) from the improved analysis of the paper [3].

Bias of gradient approximation Using the variational representation of the Euclidean norm, and definition of gradient approximation (7) we can write:

$$\begin{aligned}
\|\mathbf{b}(x)\| &= \|\mathbb{E}[\mathbf{g}(x_k, \xi, \mathbf{e})] - \nabla f(x_k)\| \\
&= \left\| \frac{d}{2h} \mathbb{E} \left[\left(\tilde{f}(x_k + h\mathbf{r}\mathbf{e}) - \tilde{f}(x_k - h\mathbf{r}\mathbf{e}) \right) K(r)\mathbf{e} \right] - \nabla f(x_k) \right\| \\
&\stackrel{\textcircled{1}}{=} \left\| \frac{d}{h} \mathbb{E} [f(x_k + h\mathbf{r}\mathbf{e})K(r)\mathbf{e}] - \nabla f(x_k) \right\| \\
&\stackrel{\textcircled{2}}{=} \|\mathbb{E}[\nabla f(x_k + h\mathbf{r}\mathbf{u})rK(r)] - \nabla f(x_k)\| \\
&= \sup_{z \in S_2^d(1)} \mathbb{E}[(\nabla_z f(x_k + h\mathbf{r}\mathbf{u}) - \nabla_z f(x_k))rK(r)] \\
&\stackrel{(19),(20)}{\leq} \kappa_\beta h^{\beta-1} \frac{L}{(l-1)!} \mathbb{E}[\|u\|^{\beta-1}] \\
&\leq \kappa_\beta h^{\beta-1} \frac{L}{(l-1)!} \frac{d}{d+\beta-1} \\
&\lesssim \kappa_\beta L h^{\beta-1}, \tag{8}
\end{aligned}$$

where $u \in B_2^d(1)$, $\textcircled{1}$ = the equality is obtained from the fact, namely, distribution of e is symmetric, $\textcircled{2}$ = the equality is obtained from a version of Stokes' theorem [56] (see Section 13.3.5, Exercise 14a).

Bounding second moment of gradient approximation By definition gradient approximation (7) and Wirtinger-Poincare inequality (18) we have

$$\begin{aligned}
\mathbb{E}[\|\mathbf{g}(x_k, \xi, \mathbf{e})\|^2] &= \frac{d^2}{4h^2} \mathbb{E} \left[\left\| \left(\tilde{f}(x_k + h\mathbf{r}\mathbf{e}) - \tilde{f}(x_k - h\mathbf{r}\mathbf{e}) \right) K(r)\mathbf{e} \right\|^2 \right] \\
&= \frac{d^2}{4h^2} \mathbb{E} \left[\left(f(x_k + h\mathbf{r}\mathbf{e}) - f(x_k - h\mathbf{r}\mathbf{e}) + (\tilde{\xi}_1 - \tilde{\xi}_2) \right)^2 K^2(r) \right] \\
&\stackrel{(15)}{\leq} \frac{\kappa d^2}{2h^2} \left(\mathbb{E} \left[(f(x_k + h\mathbf{r}\mathbf{e}) - f(x_k - h\mathbf{r}\mathbf{e}))^2 \right] + 2\Delta^2 \right) \\
&\stackrel{(18)}{\leq} \frac{\kappa d^2}{2h^2} \left(\frac{h^2}{d} \mathbb{E} \left[\|\nabla f(x_k + h\mathbf{r}\mathbf{e}) + \nabla f(x_k - h\mathbf{r}\mathbf{e})\|^2 \right] + 2\Delta^2 \right) \\
&= \frac{\kappa d^2}{2h^2} \left(\frac{h^2}{d} \mathbb{E} \left[\|\nabla f(x_k + h\mathbf{r}\mathbf{e}) + \nabla f(x_k - h\mathbf{r}\mathbf{e}) \pm 2\nabla f(x_k)\|^2 \right] + 2\Delta^2 \right) \\
&\stackrel{(17)}{\leq} \underbrace{4d\kappa}_{\rho} \|\nabla f(x_k)\|^2 + \underbrace{4d\kappa L^2 h^2 + \frac{\kappa d^2 \Delta^2}{h^2}}_{\sigma^2}. \tag{9}
\end{aligned}$$

We can now explicitly obtain the convergence rate of the novel gradient-free Algorithm 1: Zero-Order Accelerated Stochastic Gradient Descent by substituting the bias (8) and the second moment (9) of the gradient approximation (7) into the convergence rate of the first-order algorithm, which we use as the base for creating zero-order algorithm, namely Biased Accelerated Stochastic

Gradient Descent (see Theorem 3.1) with $\rho_B = \max\{1, \frac{4\kappa d}{B}\}$:

$$\mathbb{E}[f(x_N)] - f^* \lesssim \frac{\rho_B^2 LR^2}{N^2} + \frac{Nd\kappa L^2 h^2}{\rho_B^2 LB} + \frac{N\kappa d^2 \Delta^2}{h^2 \rho_B^2 LB} + \tilde{R}\kappa_\beta L_\beta h^{\beta-1} + \frac{N\kappa_\beta^2 L_\beta^2 h^{2(\beta-1)}}{L}.$$

To obtain estimates for the iteration number N , the total number of gradient-free oracle calls T , and the maximum noise level Δ , we consider 4 cases depending on the batch size B .

Case 1, when $B = 1$ we have the following convergence rate:

$$\mathbb{E}[f(x_N)] - f^* \lesssim \underbrace{\frac{\kappa^2 d^2 LR^2}{N^2}}_{\textcircled{1}} + \underbrace{\frac{Nd\kappa L^2 h^2}{\kappa^2 d^2 L}}_{\textcircled{2}} + \underbrace{\frac{N\kappa d^2 \Delta^2}{h^2 \kappa^2 d^2 L}}_{\textcircled{3}} + \underbrace{\tilde{R}\kappa_\beta L_\beta h^{\beta-1}}_{\textcircled{4}} + \underbrace{\frac{N\kappa_\beta^2 L_\beta^2 h^{2(\beta-1)}}{L}}_{\textcircled{5}}.$$

From term $\textcircled{1}$, we find iteration number N required for Algorithm 1 to achieve ε -accuracy:

$$\begin{aligned} \textcircled{1}: \quad \frac{\kappa^2 d^2 LR^2}{N^2} \leq \varepsilon &\Rightarrow N \geq \sqrt{\frac{\kappa^2 d^2 LR^2}{\varepsilon}}; \\ N &= \mathcal{O}\left(\sqrt{\frac{d^2 LR^2}{\varepsilon}}\right). \end{aligned} \quad (10)$$

From terms $\textcircled{2}$, $\textcircled{4}$ and $\textcircled{5}$ we find the smoothing parameter h :

$$\begin{aligned} \textcircled{2}: \quad \frac{Nd\kappa L^2 h^2}{\kappa^2 d^2 L} \leq \varepsilon &\Rightarrow h^2 \stackrel{(10)}{\lesssim} \frac{\kappa^2 d^2 \varepsilon^{3/2}}{\kappa^2 d^2} \Rightarrow h \lesssim \varepsilon^{3/4}; \\ \textcircled{4}: \quad \tilde{R}\kappa_\beta L_\beta h^{\beta-1} \leq \varepsilon &\Rightarrow h \lesssim \varepsilon^{1/(\beta-1)}; \\ \textcircled{5}: \quad \frac{N\kappa_\beta^2 L_\beta^2 h^{2(\beta-1)}}{L} \leq \varepsilon &\Rightarrow h^{2(\beta-1)} \stackrel{(10)}{\lesssim} d^{-1} \varepsilon^{3/2} \quad h \lesssim \frac{\varepsilon^{\frac{3}{4(\beta-1)}}}{d^{\frac{1}{2(\beta-1)}}}. \end{aligned}$$

- When $\beta \geq \frac{7}{3}$, we have that $h \lesssim \min\{\varepsilon^{3/4}, \varepsilon^{1/(\beta-1)}, \frac{\varepsilon^{\frac{3}{4(\beta-1)}}}{d^{\frac{1}{2(\beta-1)}}}\} = \varepsilon^{3/4}$.

From term $\textcircled{3}$, we find the maximum noise level Δ at which Algorithm 1 can still achieve the desired accuracy:

$$\textcircled{3}: \quad \frac{N\kappa d^2 \Delta^2}{h^2 \kappa^2 d^2 L} \leq \varepsilon \Rightarrow \Delta^2 \lesssim \frac{\varepsilon^3 d^2}{d^3} \Rightarrow \Delta \lesssim \frac{\varepsilon^{3/2}}{\sqrt{d}}.$$

- When $\beta < \frac{7}{3}$, we have that $h \lesssim \min\{\varepsilon^{3/4}, \varepsilon^{1/(\beta-1)}, \frac{\varepsilon^{\frac{3}{4(\beta-1)}}}{d^{\frac{1}{2(\beta-1)}}}\} = \varepsilon^{1/(\beta-1)}$.

From term $\textcircled{3}$, we find the maximum noise level Δ at which Algorithm 1 can still achieve the desired accuracy:

$$\textcircled{3}: \quad \frac{N\kappa d^2 \Delta^2}{h^2 \kappa^2 d^2 L} \leq \varepsilon \Rightarrow \Delta^2 \lesssim \frac{\varepsilon^{\frac{3}{2} + \frac{2}{\beta-1}} d^2}{d^3} \Rightarrow \Delta \lesssim \frac{\varepsilon^{\frac{3\beta+1}{4(\beta-1)}}}{\sqrt{d}}.$$

The oracle complexity T in this case coincides with the iteration complexity N and has the following form:

$$T = N \cdot B = \mathcal{O} \left(\sqrt{\frac{d^2 LR^2}{\varepsilon}} \right).$$

Case 2, when $1 < B < 4\kappa d$ we have the following convergence rate:

$$\mathbb{E}[f(x_N)] - f^* \lesssim \underbrace{\frac{\kappa^2 d^2 LR^2}{N^2 B^2}}_{\textcircled{1}} + \underbrace{\frac{Nd\kappa L^2 h^2 B^2}{\kappa^2 d^2 LB}}_{\textcircled{2}} + \underbrace{\frac{N\kappa d^2 \Delta^2 B^2}{h^2 \kappa^2 d^2 LB}}_{\textcircled{3}} + \underbrace{\tilde{R}\kappa_\beta L_\beta h^{\beta-1}}_{\textcircled{4}} + \underbrace{\frac{N\kappa_\beta^2 L_\beta^2 h^{2(\beta-1)}}{L}}_{\textcircled{5}}.$$

From term ①, we find iteration number N required for Algorithm 1 to achieve ε -accuracy:

$$\begin{aligned} \textcircled{1}: \quad \frac{\kappa^2 d^2 LR^2}{N^2 B^2} \leq \varepsilon &\Rightarrow N \geq \sqrt{\frac{\kappa^2 d^2 LR^2}{\varepsilon B^2}}; \\ N &= \mathcal{O} \left(\sqrt{\frac{d^2 LR^2}{\varepsilon B^2}} \right). \end{aligned} \quad (11)$$

From terms ②, ④ and ⑤ we find the smoothing parameter h :

$$\begin{aligned} \textcircled{2}: \quad \frac{Nd\kappa L^2 h^2 B^2}{\kappa^2 d^2 LB} \leq \varepsilon &\Rightarrow h^2 \stackrel{(11)}{\lesssim} \frac{\kappa^2 d^2 B^2 \varepsilon^{3/2}}{\kappa^2 d^2 B^2} \Rightarrow h \lesssim \varepsilon^{3/4}; \\ \textcircled{4}: \quad \tilde{R}\kappa_\beta L_\beta h^{\beta-1} \leq \varepsilon &\Rightarrow h \lesssim \varepsilon^{1/(\beta-1)}; \\ \textcircled{5}: \quad \frac{N\kappa_\beta^2 L_\beta^2 h^{2(\beta-1)}}{L} \leq \varepsilon &\Rightarrow h^{2(\beta-1)} \stackrel{(11)}{\lesssim} d^{-1} \varepsilon^{3/2} \quad h \lesssim \frac{\varepsilon^{\frac{3}{4(\beta-1)}}}{d^{\frac{1}{2(\beta-1)}}}. \end{aligned}$$

- When $\beta \geq \frac{7}{3}$, we have that $h \lesssim \min\{\varepsilon^{3/4}, \varepsilon^{1/(\beta-1)}, \frac{\varepsilon^{\frac{3}{4(\beta-1)}}}{d^{\frac{1}{2(\beta-1)}}}\} = \varepsilon^{3/4}$.

From term ③, we find the maximum noise level Δ at which Algorithm 1 can still achieve the desired accuracy:

$$\textcircled{3}: \quad \frac{N\kappa d^2 B^2 \Delta^2}{h^2 \kappa^2 d^2 LB} \leq \varepsilon \Rightarrow \Delta^2 \lesssim \frac{\varepsilon^3 d^2 B^2}{d^3 B^2} \Rightarrow \Delta \lesssim \frac{\varepsilon^{3/2}}{\sqrt{d}}.$$

- When $\beta < \frac{7}{3}$, we have that $h \lesssim \min\{\varepsilon^{3/4}, \varepsilon^{1/(\beta-1)}, \frac{\varepsilon^{\frac{3}{4(\beta-1)}}}{d^{\frac{1}{2(\beta-1)}}}\} = \varepsilon^{1/(\beta-1)}$.

From term ③, we find the maximum noise level Δ at which Algorithm 1 can still achieve the desired accuracy:

$$\textcircled{3}: \quad \frac{N\kappa d^2 B^2 \Delta^2}{h^2 \kappa^2 d^2 LB} \leq \varepsilon \Rightarrow \Delta^2 \lesssim \frac{\varepsilon^{\frac{3}{2} + \frac{2}{\beta-1}} d^2 B^2}{d^3 B^2} \Rightarrow \Delta \lesssim \frac{\varepsilon^{\frac{3\beta+1}{4(\beta-1)}}}{\sqrt{d}}.$$

The oracle complexity T in this case has the following form:

$$T = N \cdot B = \mathcal{O} \left(\sqrt{\frac{d^2 LR^2}{\varepsilon}} \right).$$

Case 3, when $B = 4\kappa d$ we have the following convergence rate:

$$\mathbb{E}[f(x_N)] - f^* \lesssim \underbrace{\frac{LR^2}{N^2}}_{\textcircled{1}} + \underbrace{\frac{Nd\kappa L^2 h^2}{Ld\kappa}}_{\textcircled{2}} + \underbrace{\frac{N\kappa d^2 \Delta^2}{h^2 L\kappa d}}_{\textcircled{3}} + \underbrace{\tilde{R}\kappa_\beta L_\beta h^{\beta-1}}_{\textcircled{4}} + \underbrace{\frac{N\kappa_\beta^2 L_\beta^2 h^{2(\beta-1)}}{L}}_{\textcircled{5}}.$$

From term $\textcircled{1}$, we find iteration number N required for Algorithm 1 to achieve ε -accuracy:

$$\begin{aligned} \textcircled{1}: \quad \frac{LR^2}{N^2} \leq \varepsilon &\Rightarrow N \geq \sqrt{\frac{LR^2}{\varepsilon}}; \\ N &= \mathcal{O}\left(\sqrt{\frac{LR^2}{\varepsilon}}\right). \end{aligned} \quad (12)$$

From terms $\textcircled{2}$, $\textcircled{4}$ and $\textcircled{5}$ we find the smoothing parameter h :

$$\begin{aligned} \textcircled{2}: \quad \frac{Nd\kappa L^2 h^2}{Ld\kappa} \leq \varepsilon &\Rightarrow h^2 \stackrel{(12)}{\lesssim} \varepsilon^{3/2} \Rightarrow h \lesssim \varepsilon^{3/4}, \\ \textcircled{4}: \quad \tilde{R}\kappa_\beta L_\beta h^{\beta-1} \leq \varepsilon &\Rightarrow h \lesssim \varepsilon^{1/(\beta-1)}, \\ \textcircled{5}: \quad \frac{N\kappa_\beta^2 L_\beta^2 h^{2(\beta-1)}}{L} \leq \varepsilon &\Rightarrow h^{2(\beta-1)} \stackrel{(12)}{\lesssim} d^{-1} \varepsilon^{3/2} \Rightarrow h \lesssim \frac{\varepsilon^{\frac{3}{4(\beta-1)}}}{d^{\frac{1}{2(\beta-1)}}}. \end{aligned}$$

- When $\beta \geq \frac{7}{3}$, we have that $h \lesssim \min\{\varepsilon^{3/4}, \varepsilon^{1/(\beta-1)}, \frac{\varepsilon^{\frac{3}{4(\beta-1)}}}{d^{\frac{1}{2(\beta-1)}}}\} = \varepsilon^{3/4}$.

From term $\textcircled{3}$, we find the maximum noise level Δ at which Algorithm 1 can still achieve the desired accuracy:

$$\textcircled{3}: \quad \frac{N\kappa d^2 \Delta^2}{h^2 \kappa d L} \leq \varepsilon \Rightarrow \Delta^2 \lesssim \frac{\varepsilon^3 d}{d^2} \Rightarrow \Delta \lesssim \frac{\varepsilon^{3/2}}{\sqrt{d}}.$$

- When $\beta < \frac{7}{3}$, we have that $h \lesssim \min\{\varepsilon^{3/4}, \varepsilon^{1/(\beta-1)}, \frac{\varepsilon^{\frac{3}{4(\beta-1)}}}{d^{\frac{1}{2(\beta-1)}}}\} = \varepsilon^{1/(\beta-1)}$.

From term $\textcircled{3}$, we find the maximum noise level Δ at which Algorithm 1 can still achieve the desired accuracy:

$$\textcircled{3}: \quad \frac{N\kappa d^2 \Delta^2}{h^2 \kappa d L} \leq \varepsilon \Rightarrow \Delta^2 \lesssim \frac{\varepsilon^{\frac{3}{2} + \frac{2}{\beta-1}} d}{d^2} \Rightarrow \Delta \lesssim \frac{\varepsilon^{\frac{3\beta+1}{4(\beta-1)}}}{\sqrt{d}}.$$

The oracle complexity T in this case has the following form:

$$T = N \cdot B = \mathcal{O}\left(\sqrt{\frac{d^2 LR^2}{\varepsilon}}\right).$$

Case 4, when $B > 4\kappa d$ we have the following convergence rate:

$$\mathbb{E}[f(x_N)] - f^* \lesssim \underbrace{\frac{LR^2}{N^2}}_{\textcircled{1}} + \underbrace{\frac{Nd\kappa L^2 h^2}{LB}}_{\textcircled{2}} + \underbrace{\frac{N\kappa d^2 \Delta^2}{h^2 LB}}_{\textcircled{3}} + \underbrace{\tilde{R}\kappa_\beta L_\beta h^{\beta-1}}_{\textcircled{4}} + \underbrace{\frac{N\kappa_\beta^2 L_\beta^2 h^{2(\beta-1)}}{L}}_{\textcircled{5}}.$$

From term $\textcircled{1}$, we find iteration number N required for Algorithm 1 to achieve ε -accuracy:

$$\begin{aligned} \textcircled{1}: \quad \frac{LR^2}{N^2} \leq \varepsilon &\Rightarrow N \geq \sqrt{\frac{LR^2}{\varepsilon}}; \\ N &= \mathcal{O}\left(\sqrt{\frac{LR^2}{\varepsilon}}\right). \end{aligned} \quad (13)$$

From terms $\textcircled{2}$, $\textcircled{4}$ and $\textcircled{5}$ we find the smoothing parameter h :

$$\begin{aligned} \textcircled{2}: \quad \frac{Nd\kappa L^2 h^2}{LB} \leq \varepsilon &\Rightarrow h^2 \stackrel{(13)}{\lesssim} \frac{\varepsilon^{3/2} B}{d} \Rightarrow h \lesssim \frac{\varepsilon^{3/4} B^{1/2}}{d^{1/2}}; \\ \textcircled{4}: \quad \tilde{R}\kappa_\beta L_\beta h^{\beta-1} \leq \varepsilon &\Rightarrow h \lesssim \varepsilon^{1/(\beta-1)}; \\ \textcircled{5}: \quad \frac{N\kappa_\beta^2 L_\beta^2 h^{2(\beta-1)}}{L} \leq \varepsilon &\Rightarrow h^{2(\beta-1)} \stackrel{(13)}{\lesssim} d^{-1} \varepsilon^{3/2} \Rightarrow h \lesssim \frac{\varepsilon^{\frac{3}{4(\beta-1)}}}{d^{\frac{1}{2(\beta-1)}}}. \end{aligned}$$

The smoothing parameter can be estimated as $h \lesssim \min\{\varepsilon^{3/4}, \varepsilon^{1/(\beta-1)}, \frac{\varepsilon^{\frac{3}{4(\beta-1)}}}{d^{\frac{1}{2(\beta-1)}}}\} = \varepsilon^{1/(\beta-1)}$. **From term $\textcircled{3}$** , we find the maximum noise level Δ (via batch size B) at which Algorithm 1 can still achieve the desired accuracy:

$$\textcircled{3}: \quad \frac{N\kappa d^2 \Delta^2}{h^2 LB} \leq \varepsilon \Rightarrow \Delta^2 \lesssim \frac{\varepsilon^{\frac{3}{2} + \frac{2}{\beta-1}} B}{d^2} \Rightarrow \Delta \lesssim \frac{\varepsilon^{\frac{3\beta+1}{4(\beta-1)}} B^{1/2}}{d}$$

or let's represent the batch size B via the maximum noise level Δ :

$$\textcircled{3}: \quad \frac{N\kappa d^2 \Delta^2}{h^2 LB} \leq \varepsilon \Rightarrow B \stackrel{(13)}{\gtrsim} \frac{\kappa d^2 \Delta^2}{\varepsilon^{\frac{3}{2} + \frac{2}{\beta-1}}} \Rightarrow B = \mathcal{O}\left(\frac{d^2 \Delta^2}{\varepsilon^{\frac{3}{2} + \frac{2}{\beta-1}}}\right).$$

Then the oracle complexity T in this case has the following form:

$$T = N \cdot B = \max\left\{\mathcal{O}\left(\sqrt{\frac{d^2 LR^2}{\varepsilon}}\right), \mathcal{O}\left(\frac{d^2 \Delta^2}{\varepsilon^{\frac{3}{2} + \frac{2}{\beta-1}}}\right)\right\}.$$

□

5 Discussion and Further Work

Section 4 focuses on solving the convex deterministic black-box optimization problem (1), however, when constructing the gradient-free algorithm (see Subsection 4.2) is based on a first-order method that solves the convex stochastic optimization problem (5) due to the arising of artificial “stochasticity” in the gradient approximation (7). It is not difficult to show that the results of Theorem 4.1 will be robust if the original problem of Section 4 is replaced by a stochastic black-box optimization problem, since there will already be two stochasticities in the analysis that can be formally combined into one $(\tilde{\xi}, \mathbf{e})$.

If we pay attention to the results presented in the works of [7, 4, 38] and others, we can see that they “struggle” for oracle complexity T . However, in high dimensional problems, it is important to be able to distribute the computational power loads, thereby reducing the time taken to solve a particular problem. Therefore, with the help of a not tricky technique, namely with the help of batching technique and using the accelerated algorithm as a base (in particular, Accelerated Stochastic Gradient Descent with accelerated of Nesterov, see Theorem 3.1), we managed to improve the estimate of the number of consecutive iterations $N \sim \varepsilon^{-1/2}$ to achieve the desired accuracy ε of the solution of the original problem, without worsening the oracle complexity T , and moreover improving in terms of dimensionality d for a class of convex optimization problems. It is due to this fact, namely the ability to improve one optimality criterion without compromising the second one, that recently authors of works on gradient-free optimization algorithms have been evaluating the efficiency of their algorithms by three optimality criteria at once [18]: iteration complexity, total number of calls to the gradient-free oracle, and maximum noise level Δ at which it is still possible to achieve desired accuracy.

Theorem 4.1 clearly demonstrates the advantage of Algorithm 1 (which uses enhanced smoothness information) over existing gradient-free algorithms that do not use this information. For example, when using the proposed algorithm (in the case, for example, when the smoothness of the function $\beta = 5$ and the batch size $B = 4d\kappa$) we will converge to the same error floor as the existing algorithms (solving smooth problems $\beta = 2$), only faster. However, we want our algorithm not only to converge faster, but also more accurately (to achieve a better error floor). This is exactly what we have achieved (which is really not trivial). We show that using $B > 4d\kappa$ the maximum noise level (or at a fixed noise level, error floor) accounts for smoothness order and batching costs, preserving optimality on iteration complexity.

We see the following directions as the development of this work: obtaining convergence results for a μ -strongly convex black-box optimization problem. For this formulation of the problem there are already some results presented in [3], we expect that using similar reasoning, namely generalizing the convergence results of the accelerated algorithm [52] for solving a strongly convex stochastic optimization problem to the case with a biased gradient oracle (see Definition 2.1) and using the kernel approximation (7), we will be able to improve the performance of [3] in terms of iteration complexity N , achieving

the same oracle complexity estimates T , and provide an explicit condition on the maximum noise level Δ . Another direction of development of our work is to improve oracle complexity T for convex and strongly convex black-box optimization problem. It is worth noting that in the class of convex functions we managed to improve oracle complexity, but this upper bound does not match the lower bound presented in [4,38]. Finally, the last direction that looks promising at the moment is the study of the maximum allowable noise level. In this paper, we have provided a maximum noise level at which convergence to the desired accuracy ε is guaranteed; however, we have not guaranteed the optimality of this estimate since the upper bound on the noise level is not yet known. Furthermore, we expect that this estimator can be improved by using a different concept of a gradient-free oracle, in particular when the oracle can output the objective function value with some bounded adversarial deterministic noise (see [16] for details). In this case, we can also expect an improvement in oracle complexity, since the gradient approximation with a central finite difference structure will already have access to two-point feedback.

6 Experiments

In this section, we verify the performance of the proposed algorithm in Section 4 on model functions as well as on functions of interest in machine learning. In all experiments as the Kernel $K(r)$ of gradient approximation (7) we use the already standard function, namely Legendre polynomials, for which it was proved in the paper [7] that the constants κ and κ_β do not depend on the dimensionality d , but only on the smoothness order of β . We have the following values for different β :

$$\begin{aligned} K(r) &= \frac{15r}{4}(5 - 7r^2) && \text{for } \beta = 3, 4; \\ K(r) &= \frac{195r}{16}(99r^4 - 126r^2 + 35) && \text{for } \beta = 5, 6. \end{aligned}$$

6.1 System of linear equations

We consider solving a system p of linear equations, where the problem (1) looks like $\min_{x \in \mathbb{R}^d} \|Ax - b\|^2$, where $x \in \mathbb{R}^d$, $A \in \mathbb{R}^{p \times d}$ and $b \in \mathbb{R}^p$. Figure 1a shows the advantage of the Kernel approximation for the case when the function has a higher order of smoothness. Here we optimize $f(x)$ with the parameters: $d = 64$ (dimension of the problem); $B = 50$ (batch size); $\Delta = 10^{-5}$ (noise level), $\eta = 0.02$ (step-size); $h = 0.5$ (smoothing parameter). Figure 1b shows the advantage of the algorithm proposed in Section 4 compared to existing accelerated algorithms: ARDFDS from [21], ZO-VARAG from [13]. We can see that ZO-VARAG hits the asymptote rather quickly, since the work of [13] did not consider the problem formulation with adversarial noise, i.e.,

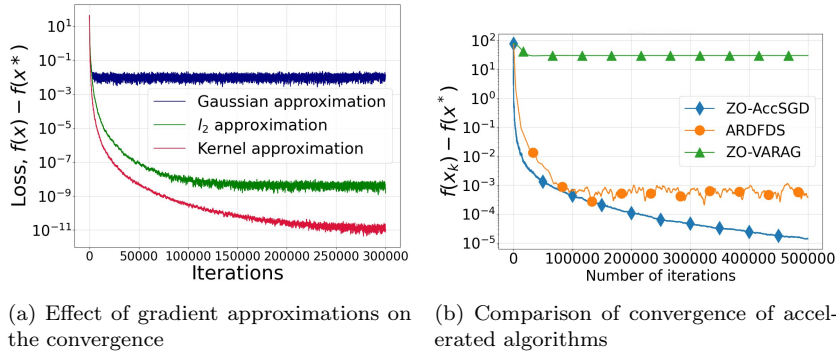
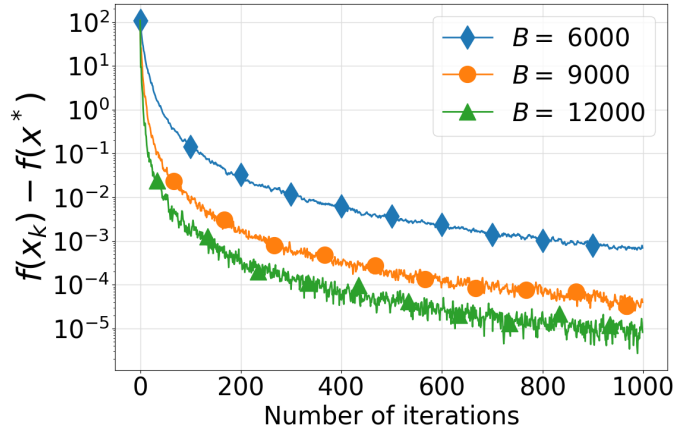


Fig. 1: Numerical experiments

the algorithm is not adaptive to noise. It is also worth noting that our algorithm outperforms ARDFDS from [21] in particular because ZO-AccSGD uses information about the high smoothness of the function.

From Figure 2 we can see the same convergence rates for the first 20-30 iterates with $B = 4d\kappa = 9000$ and $B = 12000$, due to the same values of ρ_B (optimal in terms of iteration complexity, see Theorem 4.1). However, going further through iterations, the error floor separation becomes more clear. This validates our theoretical results: the overbatching effect does improve the distance to the solution x^* of the proposed Algorithm 1.

Fig. 2: Demonstration of the overbatching effect on a model problem of solving a linear system of p equations

6.2 Logistic regression

This Subsection is devoted to numerical experiments on the optimization of a function of interest in machine learning, namely logistic regression:

$$\min_{x \in \mathbb{R}^d} f(x) = \frac{1}{M} \sum_{i=1}^M f_i(x), \quad (14)$$

where $f_i(x) = \log(1 + \exp(-y_i \cdot (Ax)_i))$ stands for the loss on the i -th data point, $A \in \mathbb{R}^{M \times d}$ is an instances matrix, $y \in \{-1, 1\}^M$ is a label vector and $x \in \mathbb{R}^d$ is a vector of weights. It is easy to show, that logistic regression function is L -smooth with $L = \frac{1}{4M} \sqrt{\lambda_{\max}(A^T A)}$, where $\lambda_{\max}(A^T A)$ denotes the largest eigenvalue of the matrix $A^T A$.

In our experiments we use data from LIBSVM[12] library, specifically such datasets as *phishing*, *diabetes* and *hearts*. The information about them can be found in the Table 1:

	<i>phishing</i>	<i>diabetes</i>	<i>hearts</i>
Size (M)	11055	768	270
Dimension (d)	68	8	13

Table 1: Summary of used datasets

In all tests we use standard solvers from scipy library to get an accurate approximation of the solution $f(x^*)$. Next, we set x_0 such that $f(x_0) - f(x^*) \sim 1$, and in all experiments smoothing parameter h is equal to 10^{-10} . We compare the results of our Algorithm 1 (ZO-AccSGD) with two methods presented in [20]: RDFDS and its accelerated version ARDFDS. To find better learning rates we grid searched it for all methods and then took the best parameters for each method. The results for different batching constant B are presented in Figure 3:

As an accelerated method, our Algorithm 1 is sensible to stochasticity, however compared to ARDFDS has a better robustness. Also, it is clear from graphs that bigger batchsize implies lower error floor and faster convergence rate. Even with gridsearched parameters for stepsize in all methods, ZO-AccSGD (see Algorithm 1) performance is better compared to RDFDS and ARDFDS ones.

Technical aspects. All experiments were written on Python and are available in anonymous repository at the following link: <https://github.com/MetalistForever/ZO-AccSGD>. They were run on the local machine with mobile CPU Ryzen 4800H. Average time for single run with number of iterations $N = 10^5$ and batchsize $B = 10$ on the *hearts* dataset is 9 min (≈ 205 it/s), nonetheless the are more time-consuming tests, for example to run optimizer on *phishing* dataset with $N = 10^5$ and $B = 100$ it took more than an hour.

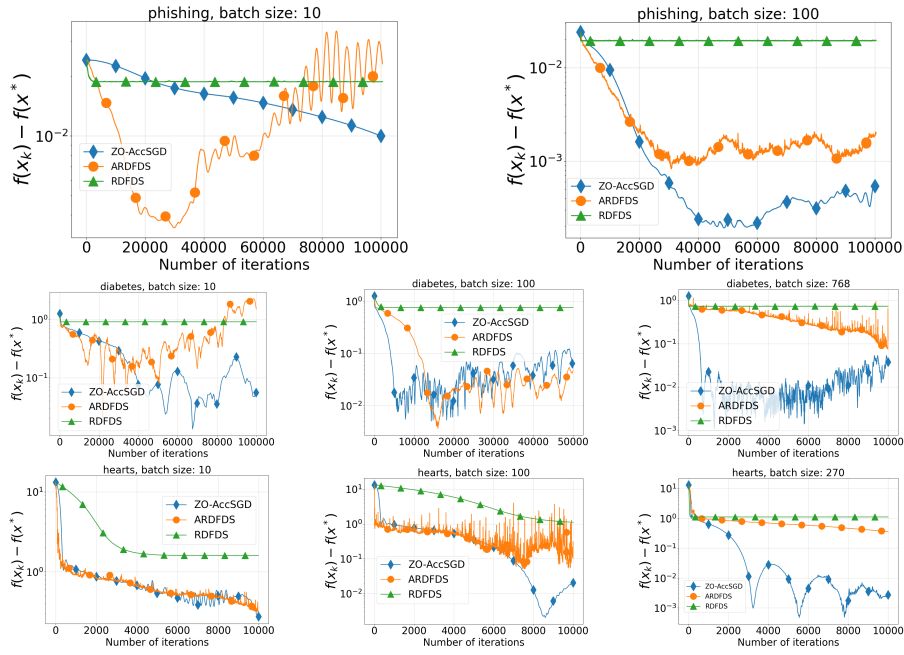


Fig. 3: Numerical results for solving logistic regression problem (14) for different datasets

7 Conclusion

In this paper, we proposed a novel gradient-free algorithm (see Algorithm 1) that improves the iteration N , oracle complexities T , and explicitly defines the maximum noise level Δ at which the algorithm can still be guaranteed to converge to the desired ε accuracy in a class of convex black-box optimization problem where the function is not just smooth but has a higher order of smoothness. Our approach for the gradient-free algorithm was based on the work of [52], however, due to the biased gradient approximation (Kernel approximation), we generalized the convergence results of this work to the gradient oracle with bias, and applied a batting technique to improve the first term in the convergence of [52] (this result may be of independent interest). In the experiments section, we confirmed our theoretical results obtained in this paper. In addition, we considered possible developments of this paper.

Acknowledgements The work of Alexander Gasnikov, Aleksandr Lobanov was supported by a grant for research centers in the field of artificial intelligence, provided by the Analytical Center for the Government of the Russian Federation in accordance with the subsidy agreement (agreement identifier 000000D730321P5Q0002) and the agreement with the Ivannikov Institute for System Programming of the Russian Academy of Sciences dated November 2, 2021 No. 70-2021-00142.

References

1. Ajalloeian, A., Stich, S.U.: On the convergence of sgd with biased gradients. arXiv preprint arXiv:2008.00051 (2020)
2. Akhavan, A., Chzhen, E., Pontil, M., Tsybakov, A.: A gradient estimator via l_1 -randomization for online zero-order optimization with two point feedback. *Advances in Neural Information Processing Systems* **35**, 7685–7696 (2022)
3. Akhavan, A., Chzhen, E., Pontil, M., Tsybakov, A.B.: Gradient-free optimization of highly smooth functions: improved analysis and a new algorithm. arXiv preprint arXiv:2306.02159 (2023)
4. Akhavan, A., Pontil, M., Tsybakov, A.: Exploiting higher order smoothness in derivative-free optimization and continuous bandits. *Advances in Neural Information Processing Systems* **33**, 9017–9027 (2020)
5. Akhavan, A., Pontil, M., Tsybakov, A.: Distributed zero-order optimization under adversarial noise. *Advances in Neural Information Processing Systems* **34**, 10,209–10,220 (2021)
6. Auger, A., Hansen, N.: A restart cma evolution strategy with increasing population size. In: 2005 IEEE congress on evolutionary computation, vol. 2, pp. 1769–1776. IEEE (2005)
7. Bach, F., Perchet, V.: Highly-smooth zero-th order online optimization. In: *Conference on Learning Theory*, pp. 257–283. PMLR (2016)
8. Bartlett, P., Dani, V., Hayes, T., Kakade, S., Rakhlin, A., Tewari, A.: High-probability regret bounds for bandit online linear optimization. In: *Proceedings of the 21st Annual Conference on Learning Theory-COLT 2008*, pp. 335–342. Omnipress (2008)
9. Berahas, A.S., Cao, L., Choromanski, K., Scheinberg, K.: A theoretical and empirical comparison of gradient approximations in derivative-free optimization. *Foundations of Computational Mathematics* **22**(2), 507–560 (2022)
10. Bertsekas, D., Tsitsiklis, J.N.: *Neuro-dynamic programming*. Athena Scientific (1996)
11. Bogolubsky, L., Dvurechenskii, P., Gasnikov, A., Gusev, G., Nesterov, Y., Raigorodskii, A.M., Tikhonov, A., Zhukovskii, M.: Learning supervised pagerank with gradient-based and gradient-free optimization methods. *Advances in neural information processing systems* **29** (2016)
12. Chang, C.C., Lin, C.J.: LIBSVM: A library for support vector machines. *ACM Transactions on Intelligent Systems and Technology* **2**, 27:1–27:27 (2011). Software available at <http://www.csie.ntu.edu.tw/~cjlin/libsvm>
13. Chen, Y., Orvieto, A., Lucchi, A.: An accelerated dfo algorithm for finite-sum convex functions. In: *International Conference on Machine Learning*, pp. 1681–1690. PMLR (2020)
14. Conn, A.R., Scheinberg, K., Vicente, L.N.: *Introduction to derivative-free optimization*. SIAM (2009)
15. Dvinskikh, D., Gasnikov, A.: Decentralized and parallel primal and dual accelerated methods for stochastic convex programming problems. *Journal of Inverse and Ill-posed Problems* **29**(3), 385–405 (2021)
16. Dvinskikh, D., Tominin, V., Tominin, I., Gasnikov, A.: Noisy zeroth-order optimization for non-smooth saddle point problems. In: *International Conference on Mathematical Optimization Theory and Operations Research*, pp. 18–33. Springer (2022)
17. Flaxman, A.D., Kalai, A.T., McMahan, H.B.: Online convex optimization in the bandit setting: gradient descent without a gradient. arXiv preprint cs/0408007 (2004)
18. Gasnikov, A., Dvinskikh, D., Dvurechensky, P., Gorbunov, E., Beznosikov, A., Lobanov, A.: Randomized gradient-free methods in convex optimization. arXiv preprint arXiv:2211.13566 (2022)
19. Gasnikov, A., Novitskii, A., Novitskii, V., Abdukhakimov, F., Kamzolov, D., Beznosikov, A., Takac, M., Dvurechensky, P., Gu, B.: The power of first-order smooth optimization for black-box non-smooth problems. In: *International Conference on Machine Learning*, pp. 7241–7265. PMLR (2022)
20. Gorbunov, E., Dvinskikh, D., Gasnikov, A.: Optimal decentralized distributed algorithms for stochastic convex optimization. arXiv preprint arXiv:1911.07363 (2019)

21. Gorbunov, E., Dvurechensky, P., Gasnikov, A.: An accelerated method for derivative-free smooth stochastic convex optimization. *SIAM Journal on Optimization* **32**(2), 1210–1238 (2022)
22. Hansen, N.: The cma evolution strategy: a comparing review. *Towards a new evolutionary computation: Advances in the estimation of distribution algorithms* pp. 75–102 (2006)
23. Hazan, E., Kale, S.: Beyond the regret minimization barrier: optimal algorithms for stochastic strongly-convex optimization. *The Journal of Machine Learning Research* **15**(1), 2489–2512 (2014)
24. Hazan, E., Klivans, A., Yuan, Y.: Hyperparameter optimization: A spectral approach. *arXiv preprint arXiv:1706.00764* (2017)
25. Huang, Y., Lin, Q.: Single-loop switching subgradient methods for non-smooth weakly convex optimization with non-smooth convex constraints. *arXiv preprint* (2023)
26. Kiefer, J., Wolfowitz, J.: Stochastic estimation of the maximum of a regression function. *The Annals of Mathematical Statistics* pp. 462–466 (1952)
27. Kimiaei, M., Neumaier, A.: Efficient unconstrained black box optimization. *Mathematical Programming Computation* **14**(2), 365–414 (2022)
28. Li, L., Jamieson, K., DeSalvo, G., Rostamizadeh, A., Talwalkar, A.: Hyperband: A novel bandit-based approach to hyperparameter optimization. *The journal of machine learning research* **18**(1), 6765–6816 (2017)
29. Lobanov, A., Alashqar, B., Dvinskikh, D., Gasnikov, A.: Gradient-free federated learning methods with l_1 and l_2 -randomization for non-smooth convex stochastic optimization problems. *arXiv preprint arXiv:2211.10783* (2022)
30. Lobanov, A., Anikin, A., Gasnikov, A., Gornov, A., Chukanov, S.: Zero-order stochastic conditional gradient sliding method for non-smooth convex optimization. *arXiv preprint arXiv:2303.02778* (2023)
31. Lobanov, A., Gasnikov, A.: Accelerated zero-order sgd method for solving the black box optimization problem under “overparametrization” condition. *arXiv preprint arXiv:2307.12725* (2023)
32. Lobanov, A., Gasnikov, A., Stonyakin, F.: Highly smoothness zero-order methods for solving optimization problems under pl condition. *arXiv preprint arXiv:2305.15828* (2023)
33. Lobanov, A., Konin, G., Gasnikov, A., Kovalev, D.: Non-smooth setting of stochastic decentralized convex optimization problem over time-varying graphs. *arXiv preprint arXiv:2307.00392* (2023)
34. Nemirovski, A., Juditsky, A., Lan, G., Shapiro, A.: Robust stochastic approximation approach to stochastic programming. *SIAM Journal on optimization* **19**(4), 1574–1609 (2009)
35. Nesterov, Y.: Efficiency of coordinate descent methods on huge-scale optimization problems. *SIAM Journal on Optimization* **22**(2), 341–362 (2012)
36. Nesterov, Y., et al.: *Lectures on convex optimization*, vol. 137. Springer (2018)
37. Nguyen, A., Balasubramanian, K.: Stochastic zeroth-order functional constrained optimization: Oracle complexity and applications. *INFORMS Journal on Optimization* **5**(3), 256–272 (2023)
38. Novitskii, V., Gasnikov, A.: Improved exploiting higher order smoothness in derivative-free optimization and continuous bandit. *arXiv preprint arXiv:2101.03821* (2021)
39. Patel, K.K., Saha, A., Wang, L., Srebro, N.: Distributed online and bandit convex optimization. In: *OPT 2022: Optimization for Machine Learning (NeurIPS 2022 Workshop)* (2022)
40. Polyak, B.T.: Minimization of unsmooth functionals. *USSR Computational Mathematics and Mathematical Physics* **9**(3), 14–29 (1969)
41. Polyak, B.T.: *Introduction to optimization* (1987)
42. Polyak, B.T., Tsybakov, A.B.: Optimal order of accuracy of search algorithms in stochastic optimization. *Problemy Peredachi Informatsii* **26**(2), 45–53 (1990)
43. Rakhlin, A., Shamir, O., Sridharan, K.: Making gradient descent optimal for strongly convex stochastic optimization. In: *Proceedings of the 29th International Conference on Machine Learning*, pp. 1571–1578 (2012)

44. Rios, L.M., Sahinidis, N.V.: Derivative-free optimization: a review of algorithms and comparison of software implementations. *Journal of Global Optimization* **56**, 1247–1293 (2013)
45. Rosenbrock, H.: An automatic method for finding the greatest or least value of a function. *The computer journal* **3**(3), 175–184 (1960)
46. Scheinberg, K.: Finite difference gradient approximation: To randomize or not? *INFORMS Journal on Computing* **34**(5), 2384–2388 (2022)
47. Schmidt, M., Roux, N.L.: Fast convergence of stochastic gradient descent under a strong growth condition. *arXiv preprint arXiv:1308.6370* (2013)
48. Srebro, N., Sridharan, K., Tewari, A.: Optimistic rates for learning with a smooth loss. *arXiv preprint arXiv:1009.3896* (2010)
49. Stich, S.U.: Unified optimal analysis of the (stochastic) gradient method. *arXiv preprint arXiv:1907.04232* (2019)
50. Storn, R., Price, K.: Differential evolution—a simple and efficient heuristic for global optimization over continuous spaces. *Journal of global optimization* **11**, 341–359 (1997)
51. Vasin, A., Gasnikov, A., Dvurechensky, P., Spokoiny, V.: Accelerated gradient methods with absolute and relative noise in the gradient. *Optimization Methods and Software* pp. 1–50 (2023)
52. Vaswani, S., Bach, F., Schmidt, M.: Fast and faster convergence of sgd for overparameterized models and an accelerated perceptron. In: *The 22nd international conference on artificial intelligence and statistics*, pp. 1195–1204. PMLR (2019)
53. Woodworth, B.E., Srebro, N.: An even more optimal stochastic optimization algorithm: minibatching and interpolation learning. *Advances in Neural Information Processing Systems* **34**, 7333–7345 (2021)
54. Yu, Z., Ho, D.W., Yuan, D.: Distributed randomized gradient-free mirror descent algorithm for constrained optimization. *IEEE Transactions on Automatic Control* **67**(2), 957–964 (2021)
55. Yue, P., Fang, C., Lin, Z.: On the lower bound of minimizing polyak-łojasiewicz functions. *arXiv preprint arXiv:2212.13551* (2022)
56. Zorich, V.A., Paniagua, O.: *Mathematical analysis II*, vol. 220. Springer (2016)

APPENDIX

A Auxiliary Facts and Results

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

A.1 Squared norm of the sum

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

$$\|a_1 + \dots + a_n\|^2 \leq n \|a_1\|^2 + \dots + n \|a_n\|^2. \quad (15)$$

A.2 Fenchel-Young inequality

For all $a, b \in \mathbb{R}^d$ and $\lambda > 0$

$$\langle a, b \rangle \leq \frac{\|a\|_2^2}{2\lambda} + \frac{\lambda \|b\|_2^2}{2}. \quad (16)$$

A.3 L smoothness function

Function f is called L -smooth on \mathbb{R}^d with $L > 0$ when it is differentiable and its gradient is L -Lipschitz continuous on \mathbb{R}^d , i.e.

$$\|\nabla f(x) - \nabla f(y)\| \leq L \|x - y\|, \quad \forall x, y \in \mathbb{R}^d. \quad (17)$$

It is well-known that L -smoothness implies (see e.g., Assumption 2.1)

$$f(y) \leq f(x) + \langle \nabla f(x), y - x \rangle + \frac{L}{2} \|y - x\|^2 \quad \forall x, y \in \mathbb{R}^d,$$

and if f is additionally convex, then

$$\|\nabla f(x) - \nabla f(y)\|^2 \leq 2L (f(x) - f(y) - \langle \nabla f(y), x - y \rangle) \quad \forall x, y \in \mathbb{R}^d.$$

A.4 Wirtinger-Poincare inequality

Let f is differentiable, then for all $x \in \mathbb{R}^d$, $h\mathbf{e} \in S_2^d(h)$:

$$\mathbb{E} [f(x + h\mathbf{e})^2] \leq \frac{h^2}{d} \mathbb{E} [\|\nabla f(x + h\mathbf{e})\|^2]. \quad (18)$$

A.5 Taylor expansion

Using the Taylor expansion we have

$$\nabla_z f(x + h\mathbf{r}\mathbf{u}) = \nabla_z f(x) + \sum_{1 \leq |n| \leq l-1} \frac{(h\mathbf{r})^{|n|}}{n!} D^{(n)} \nabla_z f(x) \mathbf{u}^n + R(h\mathbf{r}\mathbf{u}), \quad (19)$$

where by assumption

$$|R(h\mathbf{r}\mathbf{u})| \leq \frac{L}{(l-1)!} \|h\mathbf{r}\mathbf{u}\|^{\beta-1} = \frac{L}{(l-1)!} |r|^{\beta-1} h^{\beta-1} \|\mathbf{u}\|^{\beta-1}. \quad (20)$$

A.6 Kernel property

If \mathbf{e} is uniformly distributed on $S_2^d(1)$ we have $\mathbb{E}[\mathbf{e}\mathbf{e}^\top] = (1/d)I_{d \times d}$, where $I_{d \times d}$ is the identity matrix. Therefore, using the facts $\mathbb{E}[rK(r)] = 1$ and $\mathbb{E}[r^{|n|}K(r)] = 0$ for $2 \leq |n| \leq l$ we have

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

A.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, [7]) that κ_β and κ do not depend on d , they depend only on β , such that for $\beta \geq 1$:

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

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

B Proof of Theorem 3.1

In this section, we present a detailed description of the derivation of the results of Theorem 3.1, which generalize the result of Lemma 3.1 to the case of a biased gradient oracle (see Definition 2.1). Therefore, our analysis will rely on the proof of Theorem 1 of [52], working through the summands responsible for the accumulation of noise in the gradient oracle. Before starting the proof, we recall that the update equations for SGD with Nesterov acceleration have the following general form:

$$x_{k+1} = y_k - \eta \mathbf{g}(y_k, \xi_k) \quad (24)$$

$$y_k = \alpha_k z_k + (1 - \alpha_k) x_k \quad (25)$$

$$z_{k+1} = \zeta_k z_k + (1 - \zeta_k) y_k - \gamma_k \eta \mathbf{g}(y_k, \xi_k) \quad (26)$$

where $\mathbf{g}(y_k, \xi_k)$ is a biased gradient oracle (see Definition 2.1) and updates of the parameters:

$$\gamma_k = \frac{1}{\rho} \left(1 + \frac{\zeta_k(1 - \alpha_k)}{\alpha_k} \right), \quad (27)$$

$$\zeta_k \geq 1 - \gamma_k \mu \eta, \quad (28)$$

$$a_{k+1} = \gamma_k \sqrt{\eta \rho} b_{k+1}, \quad (29)$$

$$b_{k+1} = \frac{b_k}{\sqrt{\zeta_k}}, \quad (30)$$

$$\alpha_k = \frac{\gamma_k \zeta_k b_{k+1}^2 \eta}{\gamma_k \zeta_k b_{k+1}^2 \eta + a_k^2}. \quad (31)$$

We now prove the following lemma assuming that the function $f(\cdot)$ is convex and smooth.

Lemma B.1 *Let convex function satisfy Assumption 2.1 and the gradient oracle $\mathbf{g}(y_k, \xi_k)$ of Definition 2.1 satisfies Assumptions 2.3 and 2.4. Then, using the updates in Equation (24)-(26) and setting the parameters according to Equations 27-30, if $\eta \leq \frac{1}{\rho L}$ and $\Phi_k := \mathbb{E}[f(x_k)] - f^*$, then the following relation holds:*

$$\begin{aligned} b_N^2 \gamma_{N-1}^2 \Phi_N &\leq \frac{a_0}{\rho \eta} \Phi_0 + \frac{b_0^2}{\rho \eta} R_0^2 + \sum_{k=0}^{N-1} \frac{b_{k+1}^2 \gamma_k}{\rho} \tilde{R} \|\mathbb{E}[\mathbf{g}_k] - \nabla f(y_k)\| \\ &+ \sum_{k=0}^{N-1} \frac{a_{k+1}^2 \sigma^2}{2\rho^2} + \sum_{k=0}^{N-1} \frac{1}{2\rho} \|\mathbb{E}[\mathbf{g}_k] - \nabla f(y_k)\|^2. \end{aligned}$$

Proof Let $\mathbf{g}_k = \mathbf{g}(y_k, \xi_k)$ and $R_{k+1} = \|z_{k+1} - x^*\|$, then from equation (26):

$$\begin{aligned} R_{k+1}^2 &= \|\zeta_k z_k + (1 - \zeta_k)y_k - x^* - \gamma_k \eta \mathbf{g}_k\|^2 \\ &= \|\zeta_k z_k + (1 - \zeta_k)y_k - x^*\|^2 + \gamma_k^2 \eta^2 \|\mathbf{g}_k\|^2 + 2\gamma_k \eta \langle x^* - \zeta_k z_k - (1 - \zeta_k)y_k, \mathbf{g}_k \rangle. \end{aligned}$$

Taking expectation wrt to ξ_k :

$$\begin{aligned} \mathbb{E}[R_{k+1}^2] &= \mathbb{E}\left[\|\zeta_k z_k + (1 - \zeta_k)y_k - x^*\|^2\right] + \gamma_k^2 \eta^2 \mathbb{E}\left[\|\mathbf{g}_k\|^2\right] \\ &\quad + 2\gamma_k \eta \mathbb{E}\langle x^* - \zeta_k z_k - (1 - \zeta_k)y_k, \mathbf{g}_k \rangle \\ &\stackrel{(4)}{\leq} \|\zeta_k z_k + (1 - \zeta_k)y_k - x^*\|^2 + \gamma_k^2 \eta^2 \rho \|\nabla f(y_k)\|^2 \\ &\quad + 2\gamma_k \eta \langle x^* - \zeta_k z_k - (1 - \zeta_k)y_k, \mathbb{E}[\mathbf{g}_k] \rangle + \gamma_k^2 \eta^2 \sigma^2 \\ &= \|\zeta_k(z_k - x^*) + (1 - \zeta_k)(y_k - x^*)\|^2 + \gamma_k^2 \eta^2 \rho \|\nabla f(y_k)\|^2 \\ &\quad + 2\gamma_k \eta \langle x^* - \zeta_k z_k - (1 - \zeta_k)y_k, \mathbb{E}[\mathbf{g}_k] \rangle + \gamma_k^2 \eta^2 \sigma^2 \\ &\stackrel{\textcircled{1}}{\leq} \zeta_k \|z_k - x^*\|^2 + (1 - \zeta_k) \|y_k - x^*\|^2 + \gamma_k^2 \eta^2 \rho \|\nabla f(y_k)\|^2 \\ &\quad + 2\gamma_k \eta \langle x^* - \zeta_k z_k - (1 - \zeta_k)y_k, \mathbb{E}[\mathbf{g}_k] \rangle + \gamma_k^2 \eta^2 \sigma^2 \\ &= \zeta_k R_k^2 + (1 - \zeta_k) \|y_k - x^*\|^2 + \gamma_k^2 \eta^2 \rho \|\nabla f(y_k)\|^2 \\ &\quad + 2\gamma_k \eta \langle x^* - \zeta_k z_k - (1 - \zeta_k)y_k, \mathbb{E}[\mathbf{g}_k] \rangle + \gamma_k^2 \eta^2 \sigma^2 \\ &= \zeta_k R_k^2 + (1 - \zeta_k) \|y_k - x^*\|^2 + \gamma_k^2 \eta^2 \rho \|\nabla f(y_k)\|^2 \\ &\quad + 2\gamma_k \eta \langle \zeta_k(y_k - z_k) + x^* - y_k, \mathbb{E}[\mathbf{g}_k] \rangle + \gamma_k^2 \eta^2 \sigma^2 \\ &\stackrel{(25)}{=} \zeta_k R_k^2 + (1 - \zeta_k) \|y_k - x^*\|^2 + \gamma_k^2 \eta^2 \rho \|\nabla f(y_k)\|^2 + \gamma_k^2 \eta^2 \sigma^2 \\ &\quad + 2\gamma_k \eta \left\langle \frac{\zeta_k(1 - \alpha_k)}{\alpha_k} (x_k - y_k) + x^* - y_k, \mathbb{E}[\mathbf{g}_k] \right\rangle \\ &= \zeta_k R_k^2 + (1 - \zeta_k) \|y_k - x^*\|^2 + \gamma_k^2 \eta^2 \rho \|\nabla f(y_k)\|^2 + \gamma_k^2 \eta^2 \sigma^2 \\ &\quad + 2\gamma_k \eta \left(\frac{\zeta_k(1 - \alpha_k)}{\alpha_k} \langle \mathbb{E}[\mathbf{g}_k], (x_k - y_k) \rangle + \langle \mathbb{E}[\mathbf{g}_k], x^* - y_k \rangle \right) \\ &\stackrel{\textcircled{2}}{\leq} \zeta_k R_k^2 + (1 - \zeta_k) \|y_k - x^*\|^2 + \gamma_k^2 \eta^2 \rho \|\nabla f(y_k)\|^2 + \gamma_k^2 \eta^2 \sigma^2 \\ &\quad + 2\gamma_k \eta \left(\frac{\zeta_k(1 - \alpha_k)}{\alpha_k} [f(x_k) - f(y_k)] + f(x^*) - f(y_k) \right) \\ &\quad + 2\gamma_k \eta \left(\frac{\zeta_k(1 - \alpha_k)}{\alpha_k} [\langle \mathbb{E}[\mathbf{g}_k] - \nabla f(y_k), (x_k - y_k) \rangle] \right) \\ &\quad + 2\gamma_k \eta \langle \mathbb{E}[\mathbf{g}_k] - \nabla f(y_k), (x^* - y_k) \rangle, \tag{32} \end{aligned}$$

where in inequality $\textcircled{1}$ we used convexity of $\|\cdot\|^2$ and in inequality $\textcircled{2}$ we used convexity of $f(\cdot)$:

$$\begin{aligned} \langle \mathbb{E}[\mathbf{g}_k], x_k - y_k \rangle &= \langle \mathbb{E}[\mathbf{g}_k] - \nabla f(y_k), x_k - y_k \rangle + \langle \nabla f(y_k), x_k - y_k \rangle \\ &\leq \langle \mathbb{E}[\mathbf{g}_k] - \nabla f(y_k), x_k - y_k \rangle + f(x_k) - f(y_k). \end{aligned}$$

By Lipschitz continuity of the gradient (see Assumption 2.1):

$$f(x_{k+1}) - f(y_k) \leq \langle \nabla f(y_k), x_{k+1} - y_k \rangle + \frac{L}{2} \|x_{k+1} - y_k\|^2 \leq -\eta \langle \nabla f(y_k), \mathbf{g}_k \rangle + \frac{L\eta^2}{2} \|\mathbf{g}_k\|^2.$$

Taking expectation wrt ξ_k and choosing the parameter $\eta \leq \frac{1}{2\rho L}$ we have:

$$\begin{aligned} \mathbb{E}[f(x_{k+1}) - f(y_k)] &\leq -\eta \langle \nabla f(y_k), \mathbb{E}[\mathbf{g}_k] \rangle + \frac{L\eta^2}{2} \mathbb{E}\left[\|\mathbf{g}_k\|^2\right] \\ &\stackrel{(4)}{\leq} -\eta \langle \nabla f(y_k), \mathbb{E}[\mathbf{g}_k] \rangle + \frac{L\eta^2 \rho}{2} \|\nabla f(y_k)\|^2 + \frac{L\eta^2 \sigma^2}{2} \end{aligned}$$

$$\begin{aligned}
&= -\eta \|\nabla f(y_k)\|^2 + \eta \langle \nabla f(y_k), \nabla f(y_k) - \mathbb{E}[\mathbf{g}_k] \rangle \\
&\quad + \frac{L\eta^2\rho}{2} \|\nabla f(y_k)\|^2 + \frac{L\eta^2\sigma^2}{2} \\
&\stackrel{(16)}{\leq} -\frac{\eta}{2} \|\nabla f(y_k)\|^2 + \frac{\eta}{2} \|\mathbb{E}[\mathbf{g}_k] - \nabla f(y_k)\|^2 \\
&\quad + \frac{L\eta^2\rho}{2} \|\nabla f(y_k)\|^2 + \frac{L\eta^2\sigma^2}{2} \\
&= \left(-\frac{\eta}{2} + \frac{L\eta^2\rho}{2}\right) \|\nabla f(y_k)\|^2 + \frac{\eta}{2} \|\mathbb{E}[\mathbf{g}_k] - \nabla f(y_k)\|^2 + \frac{L\eta^2\sigma^2}{2} \\
&\leq -\frac{\eta}{4} \|\nabla f(y_k)\|^2 + \frac{\eta}{2} \|\mathbb{E}[\mathbf{g}_k] - \nabla f(y_k)\|^2 + \frac{L\eta^2\sigma^2}{2}.
\end{aligned}$$

Consequently, we can obtain upper bound:

$$\|\nabla f(y_k)\|^2 \leq \frac{4}{\eta} \mathbb{E}[f(y_k) - f(x_{k+1})] + 2 \|\mathbb{E}[\mathbf{g}_k] - \nabla f(y_k)\|^2 + 2L\eta\sigma^2. \quad (33)$$

Substituting the upper bound gradient norm (33) in the initial inequality (32) we have

$$\begin{aligned}
\mathbb{E}[R_{k+1}^2] &\leq \zeta_k R_k^2 + (1 - \zeta_k) \|y_k - x^*\|^2 + 4\gamma_k^2 \eta \rho \mathbb{E}[f(y_k) - f(x_{k+1})] \\
&\quad + 2\gamma_k \eta \left(\frac{\zeta_k(1 - \alpha_k)}{\alpha_k} [f(x_k) - f(y_k)] + f(x^*) - f(y_k) \right) \\
&\quad + 2\gamma_k \eta \frac{\zeta_k(1 - \alpha_k)}{\alpha_k} \langle \mathbb{E}[\mathbf{g}_k] - \nabla f(y_k), (x_k - y_k) \rangle \\
&\quad + \gamma_k^2 \eta^2 \sigma^2 + 2L\gamma_k^2 \eta^3 \rho \sigma^2 + 2\gamma_k^2 \eta^2 \rho \|\mathbb{E}[\mathbf{g}_k] - \nabla f(y_k)\|^2 \\
&\quad + 2\gamma_k \eta \langle \mathbb{E}[\mathbf{g}_k] - \nabla f(y_k), (x^* - y_k) \rangle \\
&\stackrel{\textcircled{1}}{\leq} \zeta_k R_k^2 + (1 - \zeta_k) \|y_k - x^*\|^2 + 4\gamma_k^2 \eta \rho \mathbb{E}[f(y_k) - f(x_{k+1})] \\
&\quad + 2\gamma_k \eta \left(\frac{\zeta_k(1 - \alpha_k)}{\alpha_k} [f(x_k) - f(y_k)] + f(x^*) - f(y_k) \right) \\
&\quad + 2\gamma_k \eta \frac{\zeta_k(1 - \alpha_k)}{\alpha_k} \langle \mathbb{E}[\mathbf{g}_k] - \nabla f(y_k), (x_k - y_k) \rangle \\
&\quad + 2\gamma_k^2 \eta^2 \sigma^2 + 2\gamma_k^2 \eta^2 \rho \|\mathbb{E}[\mathbf{g}_k] - \nabla f(y_k)\|^2 \\
&\quad + 2\gamma_k \eta \langle \mathbb{E}[\mathbf{g}_k] - \nabla f(y_k), (x^* - y_k) \rangle \\
&= \zeta_k R_k^2 + (1 - \zeta_k) \|y_k - x^*\|^2 + 2\gamma_k^2 \eta^2 \sigma^2 \\
&\quad + f(y_k) \left(4\gamma_k^2 \eta \rho - 2\gamma_k \eta \frac{\zeta_k(1 - \alpha_k)}{\alpha_k} - 2\gamma_k \eta \right) \\
&\quad - 4\gamma_k^2 \eta \rho \mathbb{E}[f(x_{k+1})] + 2\gamma_k \eta f(x^*) + \left(2\gamma_k \eta \frac{\zeta_k(1 - \alpha_k)}{\alpha_k} \right) f(x_k) \\
&\quad + 2\gamma_k^2 \eta^2 \rho \|\mathbb{E}[\mathbf{g}_k] - \nabla f(y_k)\|^2 + 2\gamma_k \eta \langle \mathbb{E}[\mathbf{g}_k] - \nabla f(y_k), (x^* - y_k) \rangle \\
&\quad + 2\gamma_k \eta \frac{\zeta_k(1 - \alpha_k)}{\alpha_k} \langle \mathbb{E}[\mathbf{g}_k] - \nabla f(y_k), (x_k - y_k) \rangle \\
&= \zeta_k R_k^2 + (1 - \zeta_k) \|y_k - x^*\|^2 + 2\gamma_k^2 \eta^2 \sigma^2 \\
&\quad + f(y_k) \left(4\gamma_k^2 \eta \rho - 2\gamma_k \eta \frac{\zeta_k(1 - \alpha_k)}{\alpha_k} - 2\gamma_k \eta \right) \\
&\quad - 4\gamma_k^2 \eta \rho \mathbb{E}[f(x_{k+1})] + 2\gamma_k \eta f(x^*) + \left(2\gamma_k \eta \frac{\zeta_k(1 - \alpha_k)}{\alpha_k} \right) f(x_k) \\
&\quad + 2\gamma_k^2 \eta^2 \rho \|\mathbb{E}[\mathbf{g}_k] - \nabla f(y_k)\|^2
\end{aligned}$$

$$+ 2\gamma_k \eta \left\langle \frac{\zeta_k(1-\alpha_k)}{\alpha_k} (x_k - y_k) + x^* - y_k, \mathbb{E}[\mathbf{g}_k] - \nabla f(y_k) \right\rangle,$$

where in inequality ① we used the fact that $\eta \leq \frac{1}{2\rho L}$.

Since $\zeta_k \geq 1$ and $\gamma_k = \frac{1}{2\rho} \left(1 + \frac{\zeta_k(1-\alpha_k)}{\alpha_k}\right)$, we have

$$\begin{aligned} \mathbb{E}[R_{k+1}^2] &\leq \zeta_k R_k^2 + 2\gamma_k^2 \eta^2 \sigma^2 - 4\gamma_k^2 \eta \rho \mathbb{E}[f(x_{k+1})] + 2\gamma_k \eta f(x^*) \\ &\quad + 2\gamma_k \eta \left\langle \frac{\zeta_k(1-\alpha_k)}{\alpha_k} (x_k - y_k) + x^* - y_k, \mathbb{E}[\mathbf{g}_k] - \nabla f(y_k) \right\rangle \\ &\quad + \left(2\gamma_k \eta \frac{\zeta_k(1-\alpha_k)}{\alpha_k}\right) f(x_k) + 2\gamma_k^2 \eta^2 \rho \|\mathbb{E}[\mathbf{g}_k] - \nabla f(y_k)\|^2 \\ &\stackrel{(25)}{=} \zeta_k R_k^2 + 2\gamma_k^2 \eta^2 \sigma^2 - 4\gamma_k^2 \eta \rho \mathbb{E}[f(x_{k+1})] + 2\gamma_k \eta f(x^*) \\ &\quad + 2\gamma_k \eta \langle \zeta_k (y_k - z_k) + x^* - y_k, \mathbb{E}[\mathbf{g}_k] - \nabla f(y_k) \rangle \\ &\quad + \left(2\gamma_k \eta \frac{\zeta_k(1-\alpha_k)}{\alpha_k}\right) f(x_k) + 2\gamma_k^2 \eta^2 \rho \|\mathbb{E}[\mathbf{g}_k] - \nabla f(y_k)\|^2 \\ &= \zeta_k R_k^2 + 2\gamma_k^2 \eta^2 \sigma^2 - 4\gamma_k^2 \eta \rho \mathbb{E}[f(x_{k+1})] + 2\gamma_k \eta f(x^*) \\ &\quad + 2\gamma_k \eta \langle (\zeta_k - 1)(y_k - z_k) + x^* - z_k, \mathbb{E}[\mathbf{g}_k] - \nabla f(y_k) \rangle \\ &\quad + \left(2\gamma_k \eta \frac{\zeta_k(1-\alpha_k)}{\alpha_k}\right) f(x_k) + 2\gamma_k^2 \eta^2 \rho \|\mathbb{E}[\mathbf{g}_k] - \nabla f(y_k)\|^2 \\ &\leq \zeta_k R_k^2 + 2\gamma_k^2 \eta^2 \sigma^2 - 4\gamma_k^2 \eta \rho \mathbb{E}[f(x_{k+1})] + 2\gamma_k \eta f(x^*) \\ &\quad + 2\gamma_k \eta \langle (\zeta_k - 1)(y_k - z_k) + x^* - z_k, \mathbb{E}[\mathbf{g}_k] - \nabla f(y_k) \rangle \\ &\quad + \left(2\gamma_k \eta \frac{\zeta_k(1-\alpha_k)}{\alpha_k}\right) f(x_k) + 2\gamma_k^2 \eta^2 \rho \|\mathbb{E}[\mathbf{g}_k] - \nabla f(y_k)\|^2 \\ &\leq \zeta_k R_k^2 + 2\gamma_k^2 \eta^2 \sigma^2 - 4\gamma_k^2 \eta \rho \mathbb{E}[f(x_{k+1})] + 2\gamma_k \eta f(x^*) \\ &\quad + 2\gamma_k \eta (\|z_k - y_k\| + \|z_k - x^*\|) \|\mathbb{E}[\mathbf{g}_k] - \nabla f(y_k)\| \\ &\quad + \left(2\gamma_k \eta \frac{\zeta_k(1-\alpha_k)}{\alpha_k}\right) f(x_k) + 2\gamma_k^2 \eta^2 \rho \|\mathbb{E}[\mathbf{g}_k] - \nabla f(y_k)\|^2 \\ &\leq \zeta_k R_k^2 + 2\gamma_k^2 \eta^2 \sigma^2 - 4\gamma_k^2 \eta \rho \mathbb{E}[f(x_{k+1})] \\ &\quad + 2\gamma_k \eta f(x^*) + \left(2\gamma_k \eta \frac{\zeta_k(1-\alpha_k)}{\alpha_k}\right) f(x_k) \\ &\quad + 2\gamma_k \eta \tilde{R} \|\mathbb{E}[\mathbf{g}_k] - \nabla f(y_k)\| + 2\gamma_k^2 \eta^2 \rho \|\mathbb{E}[\mathbf{g}_k] - \nabla f(y_k)\|^2, \end{aligned}$$

where $\tilde{R} = \max_k \{\|z_k - x^*\|, \|z_k - y_k\|\}$.

Multiplying by b_{k+1}^2 :

$$\begin{aligned} b_{k+1}^2 \mathbb{E}[R_{k+1}^2] &\leq b_{k+1}^2 \zeta_k R_k^2 + 2b_{k+1}^2 \gamma_k^2 \eta^2 \sigma^2 - 4b_{k+1}^2 \gamma_k^2 \eta \rho \mathbb{E}[f(x_{k+1})] \\ &\quad + 2b_{k+1}^2 \gamma_k \eta f(x^*) + \left(2b_{k+1}^2 \gamma_k \eta \frac{\zeta_k(1-\alpha_k)}{\alpha_k}\right) f(x_k) \\ &\quad + 2b_{k+1}^2 \gamma_k \eta \tilde{R} \|\mathbb{E}[\mathbf{g}_k] - \nabla f(y_k)\| + 2b_{k+1}^2 \gamma_k^2 \eta^2 \rho \|\mathbb{E}[\mathbf{g}_k] - \nabla f(y_k)\|^2. \end{aligned}$$

Since

$$b_{k+1}^2 \zeta_k \leq b_k^2; \quad b_{k+1}^2 \gamma_k^2 \eta \rho = a_{k+1}^2; \quad \frac{b_{k+1}^2 \gamma_k \eta \zeta_k (1-\alpha_k)}{\alpha_k} = a_k^2$$

and $b_{k+1}^2 \gamma_k \eta - a_{k+1}^2 + a_k^2 = 0$ we have:

$$\begin{aligned} b_{k+1}^2 \mathbb{E}[R_{k+1}^2] &\leq b_k^2 R_k^2 + \frac{a_{k+1}^2 \eta \sigma^2}{\rho} - 2a_{k+1}^2 \mathbb{E}[f(x_{k+1})] + 2b_{k+1}^2 \gamma_k \eta f(x^*) + 2a_k^2 f(x_k) \\ &\quad + 2b_{k+1}^2 \gamma_k \eta \tilde{R} \|\mathbb{E}[\mathbf{g}_k] - \nabla f(y_k)\| + a_{k+1}^2 \eta \|\mathbb{E}[\mathbf{g}_k] - \nabla f(y_k)\|^2 \end{aligned}$$

$$\begin{aligned}
&= b_k^2 R_k^2 + \frac{a_{k+1}^2 \eta \sigma^2}{\rho} - 2a_{k+1}^2 [\mathbb{E}[f(x_{k+1})] - f^*] + 2a_k^2 [f(x_k) - f^*] \\
&\quad + 2 [b_{k+1}^2 \gamma_k \eta - a_{k+1}^2 + a_k^2] f(x^*) \\
&\quad + 2b_{k+1}^2 \gamma_k \eta \tilde{R} \|\mathbb{E}[\mathbf{g}_k] - \nabla f(y_k)\| + a_{k+1}^2 \eta \|\mathbb{E}[\mathbf{g}_k] - \nabla f(y_k)\|^2 \\
&= b_k^2 R_k^2 + \frac{a_{k+1}^2 \eta \sigma^2}{\rho} - 2a_{k+1}^2 [\mathbb{E}[f(x_{k+1})] - f^*] + 2a_k^2 [f(x_k) - f^*] \\
&\quad + 2b_{k+1}^2 \gamma_k \eta \tilde{R} \|\mathbb{E}[\mathbf{g}_k] - \nabla f(y_k)\| + a_{k+1}^2 \eta \|\mathbb{E}[\mathbf{g}_k] - \nabla f(y_k)\|^2.
\end{aligned}$$

By rearranging the terms and denoting $\Phi_k := \mathbb{E}[f(x_k)] - f^*$ we get:

$$\begin{aligned}
2a_{k+1}^2 \Phi_{k+1} - 2a_k^2 \Phi_k &\leq b_k^2 R_k^2 - b_{k+1}^2 \mathbb{E}[R_{k+1}^2] + \frac{a_{k+1}^2 \eta \sigma^2}{\rho} \\
&\quad + 2b_{k+1}^2 \gamma_k \eta \tilde{R} \|\mathbb{E}[\mathbf{g}_k] - \nabla f(y_k)\| + a_{k+1}^2 \eta \|\mathbb{E}[\mathbf{g}_k] - \nabla f(y_k)\|^2.
\end{aligned}$$

By summing over k we obtain:

$$\begin{aligned}
\sum_{k=0}^{N-1} 2(a_{k+1} \Phi_{k+1} - a_k \Phi_k) &\leq \sum_{k=0}^{N-1} b_k^2 R_k^2 - b_{k+1}^2 \mathbb{E}[R_{k+1}^2] + \sum_{k=0}^{N-1} \frac{a_{k+1}^2 \eta \sigma^2}{\rho} \\
&\quad + \sum_{k=0}^{N-1} 2b_{k+1}^2 \gamma_k \eta \tilde{R} \|\mathbb{E}[\mathbf{g}_k] - \nabla f(y_k)\| \\
&\quad + \sum_{k=0}^{N-1} a_{k+1}^2 \eta \|\mathbb{E}[\mathbf{g}_k] - \nabla f(y_k)\|^2.
\end{aligned}$$

Let's substitute a_{k+1}^2 :

$$\begin{aligned}
2b_N^2 \gamma_{N-1}^2 \rho \eta \Phi_N &\leq 2a_0 \Phi_0 + b_0^2 R_0^2 + \sum_{k=0}^{N-1} \frac{a_{k+1}^2 \eta \sigma^2}{\rho} + \sum_{k=0}^{N-1} 2b_{k+1}^2 \gamma_k \eta \tilde{R} \|\mathbb{E}[\mathbf{g}_k] - \nabla f(y_k)\| \\
&\quad + \sum_{k=0}^{N-1} a_{k+1}^2 \eta \|\mathbb{E}[\mathbf{g}_k] - \nabla f(y_k)\|^2.
\end{aligned}$$

Now dividing by $2\rho\eta$ we have:

$$\begin{aligned}
b_N^2 \gamma_{N-1}^2 \Phi_N &\leq \frac{a_0}{\rho \eta} \Phi_0 + \frac{b_0^2}{\rho \eta} R_0^2 + \sum_{k=0}^{N-1} \frac{a_{k+1}^2 \sigma^2}{2\rho^2} + \sum_{k=0}^{N-1} \frac{b_{k+1}^2 \gamma_k}{\rho} \tilde{R} \|\mathbb{E}[\mathbf{g}_k] - \nabla f(y_k)\| \\
&\quad + \sum_{k=0}^{N-1} \frac{a_{k+1}^2}{2\rho} \|\mathbb{E}[\mathbf{g}_k] - \nabla f(y_k)\|^2.
\end{aligned}$$

□

Lemma B.2 Under the parameter setting according to Equations 27–30, the following relation is true:

$$\gamma_k^2 - \gamma_k \frac{1}{\rho} = \gamma_{k-1}^2$$

Proof

$$\begin{aligned}
\gamma_k &= \frac{1}{2\rho} \left[1 + \frac{\zeta_k(1 - \alpha_k)}{\alpha_k} \right] \\
\gamma_k^2 - \frac{\gamma_k}{2\rho} &= \frac{\gamma_k \zeta_k (1 - \alpha_k)}{2\rho \alpha_k}
\end{aligned}$$

$$\begin{aligned}
&= \frac{1}{2\eta\rho} \frac{a_k^2}{b_{k+1}^2} \\
&= \frac{\zeta_k}{2\eta\rho} \frac{a_k^2}{b_k^2} \\
&= \frac{1}{2\eta\rho} \frac{a_k^2}{b_k^2} \\
&= \frac{1}{\eta\rho} (\gamma_{k-1}\sqrt{\eta\rho})^2 \\
&= \gamma_{k-1}^2.
\end{aligned}$$

□

Let us write the result of Lemma B.2 as:

$$\gamma_k^2 - \frac{\gamma_k}{\rho} - \gamma_{k-1}^2 = 0.$$

Next we can find the parameter γ_k :

$$\gamma_k = \frac{\frac{1}{\rho} + \sqrt{\frac{1}{\rho^2} + 4\gamma_{k-1}^2}}{2}.$$

Let $\gamma_0 = 0$, then for all k we have:

$$\begin{aligned}
\zeta_k &= 1 \\
b_{k+1} &= b_k = b_0 = 1 \\
a_{k+1} &= \gamma_k \sqrt{\eta\rho} b_0 \Rightarrow a_{k+1} = \gamma_k \sqrt{\eta\rho}.
\end{aligned}$$

The above equation implies that $a_0 = 0$. Then using the result of Lemma B.1 and by induction $\gamma_k \geq \frac{k}{2\rho}$ we have:

$$\begin{aligned}
\frac{N^2}{4\rho^2} \Phi_N &\leq \frac{a_0}{\rho\eta} \Phi_0 + \frac{b_0^2}{\rho\eta} R_0^2 + \sum_{k=0}^{N-1} \frac{a_{k+1}^2 \sigma^2}{2\rho^2} + \sum_{k=0}^{N-1} \frac{b_{k+1}^2 \gamma_k}{\rho} \tilde{R} \|\mathbb{E}[\mathbf{g}_k] - \nabla f(y_k)\| \\
&\quad + \sum_{k=0}^{N-1} \frac{a_{k+1}^2}{2\rho} \|\mathbb{E}[\mathbf{g}_k] - \nabla f(y_k)\|^2 \\
&\leq \frac{1}{\rho\eta} R_0^2 + \sum_{k=0}^{N-1} \frac{\gamma_k^2 \eta \sigma^2}{2\rho} + \sum_{k=0}^{N-1} \frac{\gamma_k}{\rho} \tilde{R} \|\mathbb{E}[\mathbf{g}_k] - \nabla f(y_k)\| \\
&\quad + \sum_{k=0}^{N-1} \frac{\gamma_k^2 \eta}{2} \|\mathbb{E}[\mathbf{g}_k] - \nabla f(y_k)\|^2 \\
&\leq \frac{1}{\rho\eta} R_0^2 + \sum_{k=0}^{N-1} \frac{k^2 \eta \sigma^2}{8\rho^3} + \sum_{k=0}^{N-1} \frac{k}{2\rho^2} \tilde{R} \|\mathbb{E}[\mathbf{g}_k] - \nabla f(y_k)\| \\
&\quad + \sum_{k=0}^{N-1} \frac{k^2}{4L\rho^2} \|\mathbb{E}[\mathbf{g}_k] - \nabla f(y_k)\|^2 \\
&\leq \frac{1}{\rho\eta} R_0^2 + \frac{k^3 \eta \sigma^2}{24\rho^3} + \frac{N^2}{4\rho^2} \tilde{R} \delta + \frac{N^3}{12L\rho^2} \delta^2.
\end{aligned}$$

Now we can get the convergence rate

$$\Phi_N \leq \frac{4\rho}{N^2\eta} R_0^2 + \frac{N\eta\sigma^2}{6\rho} + \tilde{R}\delta + \frac{N}{3L}\delta^2.$$

By adding batching, given that $\rho_B = \max\{1, \frac{\rho}{B}\}$, $\sigma_B^2 = \frac{\sigma^2}{B}$ and $R = R_0$ we have the convergence rate for accelerated SGD with biased gradient oracle and parameter $\eta \lesssim \frac{1}{\rho_B L}$:

$$\mathbb{E}[f(x_N)] - f^* \lesssim \frac{\rho_B^2 L R^2}{N^2} + \frac{N \sigma_B^2}{\rho_B^2 L} + \tilde{R} \delta + \frac{N}{L} \delta^2.$$