

Accelerated Zero-Order SGD Method for Solving the Black Box Optimization Problem under “Overparametrization” Condition^{*}

Aleksandr Lobanov^{1,2}[0000–0003–1620–9581] and Alexander Gasnikov^{1,2,3}[0000–0002–7386–039X]

¹ Moscow Institute of Physics and Technology, Dolgoprudny, Russia

² ISP RAS Research Center for Trusted Artificial Intelligence, Moscow, Russia

³ Institute for Information Transmission Problems RAS, Moscow, Russia
{lobanov.av,gasnikov.av}@mipt.ru

Abstract. This paper is devoted to solving a convex stochastic optimization problem in a overparameterization setup for the case where the original gradient computation is not available, but an objective function value can be computed. For this class of problems we provide a novel gradient-free algorithm, whose creation approach is based on applying a gradient approximation with l_2 randomization instead of a gradient oracle in the biased Accelerated SGD algorithm, which generalizes the convergence results of the AC-SA algorithm to the case where the gradient oracle returns a noisy (inexact) objective function value. We also perform a detailed analysis to find the maximum admissible level of adversarial noise at which we can guarantee to achieve the desired accuracy. We verify the theoretical results of convergence using a model example.

Keywords: Black-Box Optimization · Overparametrization · Accelerated Zero-Order SGD Method · Biased Gradient Oracle.

1 Introduction

The black-box optimization problem [1–3] usually arises when we have few information about the objective function. Such problems can appear when the objective function is non-smooth [4] (i.e. no gradient computation is available) or, for instance, when the function is smooth [5] (may even have a higher order of smoothness [6]), but the process of computing the derivatives is too expensive in contrast to computing the value of the objective function. Moreover, this problem is often solved under conditions of privacy, when some data cannot be disseminated due to confidentiality. Then a gradient-free oracle [7] (or in other words zero-order oracle) acts as some “black box”, which returns only objective function value $f(x)$ at requested point x with some bounded adversarial noise $\delta(x) \leq \Delta$ (where Δ is level noise). Where the latter means that the gradient-free

^{*} The research was supported by Russian Science Foundation (project No. 21-71-30005), <https://rscf.ru/en/project/21-71-30005/>.

oracle can give an inexact value of the objective function. As practice shows [8], the lower the level of adversarial noise, the more expensive it is to call a gradient-free oracle, so it is important to understand how large the level of adversarial noise can be, at which a “good” convergence rate is still guaranteed (by “good” convergence rate is meant the convergence of the algorithm when $\Delta = 0$). This concept of a gradient-free oracle is common in the literature and has the interpretation of a adversarial attack on the black-box model [9]. The black box problem is currently actively researched in the optimization [10; 11] and machine learning [12; 13] community, since this problem has applications in the following areas: federated learning [14; 15], distributed learning [16; 17] and deep learning [18]. A particular need for solving such problems arises in the following applications: model hyperparameters tuning [19; 20], reinforcement learning [21; 22], multi-armed bandits [23–25], and many others.

There are various approaches to solving the black-box problem, but the most common in a theoretically proof-of-concept sense is to create gradient-free optimization algorithms based on the state of the art first-order methods and using various randomized gradient approximations [26]. The most common first-order optimization methods in machine learning models are momentum SGD, Adam, and others. But note that these algorithms are variants of Stochastic Gradient Descent (SGD) [27; 28], which use stochastic gradient estimates. By adding a procedure for batching stochastic gradient estimates, it is possible to obtain algorithms that are easily parallelized on several computers. It is data distribution (parallelization) that significantly reduces the computational costs that certainly arise in a huge number of modern machine learning models. Also, with the addition of acceleration one can still achieve improvement in terms of estimates on the number of successive iterations. Thus, for many optimization problems, the state of the art first-order algorithms are accelerated batched variants of SGD.

In this paper, we focus on solving a stochastic convex black-box optimization problem in a smooth setting with an overparameterization condition. Where the latter means that the model has many more parameters than the data available. We can summarize our contributions as follows:

- We derive the convergence rate for a biased Accelerated Stochastic Gradient Descent, which covers smooth convex stochastic optimization problems under the overparameterization setup.
- We provide a novel Accelerated Zero-Order Stochastic Gradient Descent Method (AZO-SGD) for solving the black-box problem in a smooth setting under the overparameterization condition. We analyze the robustness of AZO-SGD algorithm to adversarial noise, providing an estimate for the maximum admissible level of adversarial noise at which the desired accuracy can still be achieved. We show that our algorithm is optimal on oracle calls in the class of gradient-free algorithms.
- We show the convergence of the Accelerated Zero-Order Stochastic Gradient Descent Method proposed in this paper using a model example of finite sums in which the number of summands is less than the number of variables (the overparameterization condition).

1.1 Related works

Adversarial noise. Finding the maximum admissible noise level at which one can still guarantee convergence to the desired accuracy ε is an important issue for the black-box optimization problem. Special attention to this question was allocated by the works [14; 29–31]. For example, in [29] the authors found the maximum admissible level of *adversarial deterministic noise* for a non-smooth convex black-box optimization problem $\sim \mathcal{O}(\varepsilon^2 d^{-1/2})$, moreover in [14] it is shown that this estimate will be the same for l_1 and l_2 randomization. In [30], the authors showed that by assuming a strong convexity, the estimate maximum level of *adversarial deterministic noise* can be improved to $\sim \mathcal{O}(\mu^{1/2} \varepsilon^{3/2} d^{-1/2})$ in non-smooth setting. And in [31] the authors were able to show that this estimation in a non-smooth one can be also improved to $\sim \mathcal{O}(\varepsilon d^{-1/2})$ by using *adversarial stochastic noise*, since this concept of *adversarial stochastic noise* does not accumulate in the bias, but accumulates only in the variance. In addition, this paper shows that if the function is smooth, the estimate of the maximum level of adversarial stochastic noise can be improved to $\sim \mathcal{O}(\varepsilon^{1/2} d^{-1/2})$. In this paper, we will use a gradient approximation via l_2 randomization to create a novel gradient-free algorithm, and we will find the maximum admissible level of deterministic noise in a smooth setting under overparameterization condition.

SGD type algorithms. Many works [32–40] study the Stochastic Gradient Descent and its variant in different setups. For example, in [32] the authors proposed an accelerated method of stochastic gradient descent, AC-SA. Later in [33] authors proposed optimal algorithms for federated learning architecture, which is based on AC-SA (Single-Machine Accelerated SGD and Mini-Batch Accelerated SGD) method. In [34] proposed an clipped accelerated SGD method for heavy-tailed optimization problems based on the accelerated SGD method: Stochastic Similar Triangle Method (SSTM) [35]. In [36], the authors studied the biased SGD method in the Polak-Lojasiewicz [37; 38] setup. It is worth noting that in [39] it was shown that for problems satisfying the Polak-Lojasiewicz condition the non-accelerated SGD algorithm will be optimal. In [40], the study of the AC-SA (accelerated SGD) algorithm was continued already in the overparameterization setup. We in this paper generalize the analysis of the AC-SA algorithm from [40], to create a biased accelerated SGD algorithm in the overparameterization setting. A biased first-order algorithm is necessary because using l_2 randomization produces a bias in the case when $\delta(x) > 0$. Therefore, based on the biased batched accelerated stochastic gradient descent and using l_2 randomization, we create a new gradient-free optimization algorithm to solve a convex stochastic black-box optimization problem under overparameterization setup.

Gradient noise assumptions. Recently there is a trend in works [40–47] of relaxed stochastic gradient variance restriction condition. Very many works (e.g. see [41; 42]) use the standard assumption: $\mathbb{E}[\|\nabla f(x, \xi)\|^2] \leq \sigma^2$. However, already in the works [43–45; 48] used in the analysis of the algorithm a more relaxed

assumption of weak growth: $\mathbb{E} \left[\|\nabla f(x, \xi)\|^2 \right] \leq M \|\nabla f(x)\|^2 + \sigma^2$. The following work [46] have set the constants so that the strong growth condition assumption is satisfied: $\mathbb{E} \left[\|\nabla f(x, \xi)\|^2 \right] \leq M \|\nabla f(x)\|^2$. In [40; 47] the condition satisfying the overparameterized set: $\mathbb{E} \left[\|\nabla f(x^*, \xi)\|^2 \right] \leq \sigma_*^2$. In this paper, we will also assume uniform smoothness of the function over ξ as well as the overparameterized condition, since our approach to creating gradient-free algorithms is based on the Accelerated Batched Stochastic Gradient Descent (AC-SA) proposed in [40].

1.2 Notations

We use $\langle x, y \rangle := \sum_{i=1}^d x_i y_i$ to denote standard inner product of $x, y \in \mathbb{R}^d$, 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 notation $\tilde{O}(\cdot)$ to hide logarithmic factors.

1.3 Paper Organization

This paper has the following structure. In Section 1 we introduce this paper and also provide related works. In Section 2 we formulate the problem statement. While in Section 3 we provide an accelerated SGD algorithm with a biased gradient oracle in the reparameterization setup. Section 4 presents the main result of this paper. We confirm the theoretical results via a model example in Section 5. While Section 6 concludes this paper. Detailed proofs are presented in the supplementary materials (Appendix)[†].

2 Technical Preliminaries

We study a standard stochastic convex optimization problem:

$$f^* = \min_{x \in \mathbb{R}^d} \{f(x) := \mathbb{E}[f(x, \xi)]\}, \quad (1)$$

where $f : \mathbb{R}^d \rightarrow \mathbb{R}$ is smooth convex function that we want to minimize over \mathbb{R}^d . This problem statement is a general smooth convex stochastic optimization problem, so to define the class of the problem considered in this paper we will introduce some assumptions on the objective function and on the gradient oracle.

[†] The full version of this article, which includes the Appendix can be found by the article title in the arXiv at the following link: <https://arxiv.org/abs/2307.12725>.

2.1 Assumptions on the Objective Function

In all proofs we assume convexity and smoothness of the function $f(x, \xi)$.

Assumption 1. *For almost every ξ , $f(x, \xi)$ is non-negative, convex w.r.t. x , i.e.*

$$\forall x, y, \xi \quad f(y, \xi) \geq f(x, \xi) + \langle \nabla f(x, \xi), y - x \rangle.$$

Assumption 2. *For almost every ξ , $f(x, \xi)$ is non-negative, L -smooth w.r.t. x , i.e.*

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

Assumptions 1 and 2 are common in the literature (see, e.g., [49; 50]). However, it is worth noting that these assumptions require uniform convexity and smoothness over $\xi \sim \mathcal{D}$ [51]. This is an essential difference from the standard assumptions, which define a narrower class of the objective function.

Assumption 3. *The function $f(x)$ is a convex and has the minimum value $f^* = \min_x f(x)$, which is attained at a point x^* with $\|x^*\| \leq R$.*

We explicitly introduce the problem solution f^* in Assumption 3, since our approach implies that the convergence rate depends on the solution f^* (see, e.g., [52]), i.e., our analysis will show an improvement in convergence at $f^* \rightarrow 0$.

2.2 Assumptions on the Gradient Oracle

In our analysis we consider the case when we obtain an inexact gradient value when calling the oracle. Therefore we first define a biased gradient oracle.

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

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

where $\mathbf{b} : \mathbb{R}^d \rightarrow \mathbb{R}^d$ such that $\forall x \in \mathbb{R}^d : \|\mathbf{b}(x)\|^2 \leq \zeta^2$.

Next, we assume that gradient noise is bounded as follows.

Assumption 4. *There exists $\sigma_*^2 \geq 0$ such that $\forall x \in \mathbb{R}^d$*

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

Assumption 4 is a common assumption for overparameterized optimization problems, since in this setup in many problems [53–55] f^* can be expected to be small. We introduced Assumptions 1–4 because in this paper we based on the results of [40], in which it was proved that the convergence rate of the AC-SA method (from [32]) has the following form in overparameterization setup of problem (1):

$$\mathbb{E} [f(x_N^{ag}) - f^*] \leq c \cdot \left(\frac{LR^2}{N^2} + \frac{LR^2}{BN} + \sqrt{\frac{LR^2 f^*}{BN}} \right), \quad (2)$$

where N is a iteration number, B is a batch size and $\sigma_*^2 \leq 2Lf^*$ (it's proven [40]).

3 Accelerated SGD with Biased Gradient

Our approach to create a gradient-free algorithm (for an overparameterized setup) implies the use of l_2 randomization (see Section 4), based on the AC-SA algorithm from [40]. However, the standard concept of a gradient-free oracle implies the presence of adversarial noise, which can accumulate in both variance and bias. Therefore, it is important to investigate the adversarial noise for the question of maximum admissible noise level, when the desired accuracy can be guaranteed. As can be seen, the result (2) obtained in [40] does not account for the bias in the gradient oracle (see Definition 1 in the case when $\|\mathbf{b}(x)\| = 0$). Consequently, in Theorem 1 we provide a novel biased AC-SA algorithm that is robust to the overparameterized setup and accounts for bias in gradient oracle.

Theorem 1 (Convergence of Biased AC-SA). *Let f satisfy Assumptions 1–3 and gradient oracle from Definition 1 satisfy Assumption 4, then Biased AC-SA algorithm guarantees the convergence with a universal constant c*

$$\mathbb{E}[f(x_N^{ag}) - f^*] \leq c \cdot \left(\frac{LR^2}{N^2} + \frac{LR^2}{BN} + \frac{\sigma_* R}{\sqrt{BN}} + \zeta R + \frac{\zeta^2}{2L} N \right).$$

The results of Theorem 1 show the convergence of the AC-SA algorithm, considering the bias $\mathbf{b}(x)$ in the gradient oracle. It is not difficult to see that if we consider the case without bias ($\zeta = 0$), the convergence result of Theorem 1 will fully correspond to the result (2). The last two terms in Theorem 1 are standard for the accelerated algorithm (see, for example, [56; 57]). There are several ways to obtain these results: using the (δ, L) -oracle technique [58], modifying Assumptions 1, 2, or performing sequential reasoning with current assumptions. The proof of Theorem 1 can be found in supplementary materials (Appendix B).

4 Main Result

In this section, we present the main result of this work, namely a gradient-free algorithm for solving a convex smooth stochastic black-box optimization problem in an overparameterized setup. We further narrow the problem class (1) considered in this section to the black box problem, that is, when the calculation of the gradient oracle is not available for some reason. Unfortunately, we cannot apply the AC-SA algorithm or even the biased AC-SA algorithm to solving this problem class. Therefore, there is a need to create an algorithm that only requires calculations of function values. Such algorithms are usually called *gradient-free*, since the efficiency of this class of algorithms is determined by three criteria: the maximum admissible level of adversarial noise Δ , iterative complexity N , and in particular the total number of calls to the *gradient-free oracle* T . Our approach in creating gradient-free algorithm based on the biased AC-SA algorithm. Instead of the gradient oracle (see Definition 1) we use the gradient approximation:

$$\mathbf{g}(x, \xi, e) = \frac{d}{2\tau} (f_\delta(x + \tau e, \xi) - f_\delta(x - \tau e, \xi)) e, \quad (3)$$

where e is a vector uniformly distributed on unit sphere $S_2^d(1)$, τ is a smoothing parameter and $f_\delta(x, \xi) = f(x, \xi) + \delta(x)$ ($|\delta(x)| \leq \Delta$) is a *gradient-free* oracle. Thus Algorithm 1 presents a novel gradient-free method, namely Accelerated Zero-Order Stochastic Gradient Descent (AZO-SGD) Method for solving the black-box optimization problem (1) under the overparameterization condition.

Algorithm 1 Accelerated Zero-Order Stochastic Gradient Descent (AZO-SGD)

Input: Start point $x_0^{ag} = x_0 \in \mathbb{R}^d$, maximum number of iterations $N \in \mathbb{Z}_+$.

Let stepsize $\gamma_k > 0$, parameters $\beta_k, \tau > 0$, batch size $B \in \mathbb{Z}_+$.

- 1: **for** $k = 0, \dots, N - 1$ **do**
- 2: $\beta_k = 1 + \frac{k}{6}$ and $\gamma_k = \gamma(k + 1)$ for $\gamma = \min \left\{ \frac{1}{12L}, \frac{B}{24L(N+1)}, \sqrt{\frac{BR^2}{Lf^*N^3}} \right\}$
- 3: $x_k^{md} = \beta_k^{-1}x_k + (1 - \beta_k^{-1})x_k^{ag}$
- 4: Sample $\{e_1, \dots, e_B\}$ and $\{\xi_1, \dots, \xi_B\}$ independently
- 5: Define $\mathbf{g}_k = \frac{1}{B} \sum_{i=1}^B \mathbf{g}(x_k^{md}, \xi_i, e_i)$ using (3)
- 6: $\tilde{x}_{k+1} = x_k - \gamma_k \mathbf{g}_k$
- 7: $x_{k+1} = \min \left\{ 1, \frac{R}{\|\tilde{x}_{k+1}\|} \right\} \tilde{x}_{k+1}$
- 8: $x_{k+1}^{ag} = \beta_k^{-1}x_{k+1} + (1 - \beta_k^{-1})x_k^{ag}$
- 9: **end for**

Output: x_N^{ag} .

Next, we provide Theorem 2, in which we show the convergence results for Accelerated Zero-Order Stochastic Gradient Descent (AZO-SGD) Method.

Theorem 2. *Let f satisfy Assumptions 1–3 and gradient approximation (3) with parameter $\tau \leq \frac{\varepsilon}{LR}$ satisfy Assumption 4, then Accelerated Zero-Order Stochastic Gradient Descent (AZO-SGD) Method (see Algorithm 1) achieves ε -accuracy: $\mathbb{E}[f(x_N^{ag}) - f^*] \leq \varepsilon$ after*

$$N = \mathcal{O} \left(\sqrt{\frac{LR^2}{\varepsilon}} \right), \quad T = \max \left\{ \mathcal{O} \left(\frac{LR^2}{\varepsilon} \right), \mathcal{O} \left(\frac{d\sigma_*^2 R^2}{\varepsilon^2} \right) \right\}$$

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

$$\Delta \leq \frac{\varepsilon^2}{dLR^2}$$

the maximum admissible level of adversarial noise.

The result of Theorem 2 shows the effective iterative complexity N , since an accelerated method was taken as the base (biased AC-SA, see Theorem 1). It is also worth noting that the batch size $B = \max \left\{ \mathcal{O} \left(\sqrt{\frac{LR^2}{\varepsilon}} \right), \mathcal{O} \left(\frac{d\sigma_*^2 R}{L^{1/2}\varepsilon^{3/2}} \right) \right\}$ can change with time, i.e., it directly depends on σ_*^2 , which leads to an optimal estimate of the number of gradient-free oracle calls T . Finally, one of the main

results of this paper is the estimation for the maximum admissible noise level Δ . This estimation is inferior to the other estimations $\sim \mathcal{O}(\varepsilon^2 d^{-1/2})$ in the non-smooth setting and $\sim \mathcal{O}(\varepsilon^{3/2} d^{-1/2})$ in smooth setting, but this is expected, since Algorithm 1 is working in a different setup, namely in overparameterization condition. It's also noted that we don't guarantee that this evaluation is unimproved, but only states that at $\Delta \leq \mathcal{O}(\varepsilon^2 d^{-1})$ there will be convergence to ε -accuracy. We present an open question for future research: improving the estimate to the maximum admissible noise level, as well as finding upper bounds on noise level beyond which convergence in the overparameterized setup cannot be guaranteed. The proof of Theorem 1 can be found in supplementary materials (Appendix C).

Remark 1 (General case). It is worth noting that this paper, and in particular Theorem 1 and Theorem 2, focus on the Euclidean case. However, using the work of [59] (namely, Algorithm 2) as a basis, and similarly generalizing the convergence results (Corollary 1, [59]) to the case with a biased gradient oracle (see Definition 1), we can obtain a gradient-free algorithm for a more general class of problems (L_p -norm and presence of constraints). In this case, the parameters (number of iterations N , the maximum admissible level of adversarial noise Δ , smoothing parameter τ) of the gradient-free algorithm will be the same as in Theorem 2, except for the total number of oracle calls: let given $1/p + 1/q = 1$

$$T = N \cdot B = \max \left\{ \mathcal{O} \left(\frac{LR^2}{\varepsilon} \right), \mathcal{O} \left(\frac{\min\{q, \ln d\} d^{2-\frac{2}{p}} \sigma_*^2 R^2}{\varepsilon^2} \right) \right\}.$$

This generalization allows one to solve problem (1) in a broader setting, for instance, by imposing a constraint. In particular, by solving the problem on a simplex ($p = 1$, $q = \infty$) we can achieve a reduction of the total number of calls to the zero-order oracle T by $\ln(d)$ compared to the Euclidean case.

5 Experiments

In this section, we will use a simple example to verify the theoretical results, namely to show the convergence of the proposed algorithm Accelerated Zero-Order Stochastic Gradient Descent Method (AZO-SGD, see Algorithm 1). The optimization problem (1) is as follows:

$$\min_{x \in \mathbb{R}^d} f(x) := \frac{1}{m} \sum_{i=1}^m (l_i(x))^2, \quad (4)$$

where $l(x) = Ax - b$ is system of m linear equations under overparameterization ($d > m$), $A \in \mathbb{R}^{m \times d}$, $x, b \in \mathbb{R}^d$. Problem (4) is a convex stochastic optimization problem (1), also known as the Empirical Risk Minimization problem, where $\xi = i$ is one of m linear equations.

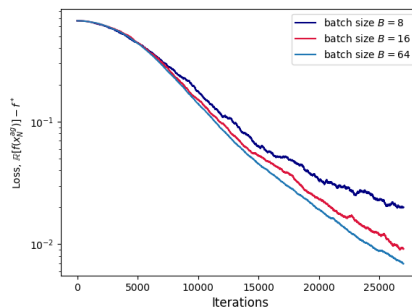


Fig. 1: Convergence of the Accelerated Zero-Order Stochastic Gradient Descent Method and the effect of parameter B (batch size) on the iteration complexity.

In Figure 1 we see the convergence of the proposed gradient-free algorithm. We can also conclude that as the batch size increases, the number of iterations required to achieve the desired accuracy decreases. This effect occurs especially in applications with huge-scale machine learning models. We optimize $f(x)$ (4) with parameters: $d = 256$ (dimensional of problem), $m = 128$ (number of linear equations), $\tau = 0.001$ (smoothing parameter), $\gamma = 0.0001$ (step size), $B = \{8, 16, 64\}$ (batch size). We also understand machine inaccuracy by noise level.

6 Conclusion

The optimization problem in the overparameterization setup has not yet been sufficiently studied. In this paper, we proposed a novel gradient-free algorithm: Accelerated Zero-Order Stochastic Gradient Descent Method for solving the smooth convex stochastic black-box optimization problem in the overparameterization setup. Our approach in creating the gradient-free Algorithm 1 was based on accelerated stochastic gradient descent. However, since there is an accumulation of adversarial noise in l_2 randomization, the result of [40] was generalized to the case of a biased gradient oracle. We also showed that the proposed gradient-free algorithm (AZO-SGD) is optimal in terms of iteration and oracle complexities. In addition, we obtained the first estimate, as far as we know, of the level of adversarial noise in the overparameterization setup, thereby opening up many potentially interesting future research questions in this setup.

The authors are grateful to Daniil Vostrikov.

References

1. *Audet C., Hare W.* Derivative-free and blackbox optimization. — 2017.
2. *Conn A. R., Scheinberg K., Vicente L. N.* Introduction to derivative-free optimization. — SIAM, 2009.
3. *Black Box Optimization, Machine Learning, and No-Free Lunch Theorems / P. M. Pardalos, V. Rasskazova, M. N. Vrahatis, [et al.].* — Springer, 2021.

4. The power of first-order smooth optimization for black-box non-smooth problems / A. Gasnikov [et al.] // International Conference on Machine Learning. — PMLR. 2022. — P. 7241–7265.
5. A gradient estimator via L1-randomization for online zero-order optimization with two point feedback / A. Akhavan [et al.] // Advances in Neural Information Processing Systems. — 2022. — Vol. 35. — P. 7685–7696.
6. *Bach F., Perchet V.* Highly-smooth zero-th order online optimization // Conference on Learning Theory. — PMLR. 2016. — P. 257–283.
7. *Rosenbrock H.* An automatic method for finding the greatest or least value of a function // The computer journal. — 1960. — Vol. 3, no. 3. — P. 175–184.
8. Learning supervised pagerank with gradient-based and gradient-free optimization methods / L. Bogolubsky [et al.] // Advances in neural information processing systems. — 2016. — Vol. 29.
9. Zoo: Zeroth order optimization based black-box attacks to deep neural networks without training substitute models / P.-Y. Chen [et al.] // Proceedings of the 10th ACM workshop on artificial intelligence and security. — 2017. — P. 15–26.
10. Optimal rates for zero-order convex optimization: The power of two function evaluations / J. C. Duchi [et al.] // IEEE Transactions on Information Theory. — 2015. — Vol. 61, no. 5. — P. 2788–2806.
11. *Shibaev I., Dvurechensky P., Gasnikov A.* Zeroth-order methods for noisy Hölder-gradient functions // Optimization Letters. — 2022. — Vol. 16, no. 7. — P. 2123–2143.
12. *Papernot N., McDaniel P., Goodfellow I.* Transferability in machine learning: from phenomena to black-box attacks using adversarial samples // arXiv preprint arXiv:1605.07277. — 2016.
13. Practical black-box attacks against machine learning / N. Papernot [et al.] // Proceedings of the 2017 ACM on Asia conference on computer and communications security. — 2017. — P. 506–519.
14. Gradient-Free Federated Learning Methods with l_1 and l_2 -Randomization for Non-Smooth Convex Stochastic Optimization Problems / A. Lobanov [et al.] // arXiv preprint arXiv:2211.10783. — 2022.
15. Distributed online and bandit convex optimization / K. K. Patel [et al.] // OPT 2022: Optimization for Machine Learning (NeurIPS 2022 Workshop). — 2022.
16. Optimal convergence rates for convex distributed optimization in networks / K. Scaman [et al.] // Journal of Machine Learning Research. — 2019. — Vol. 20. — P. 1–31.
17. Non-Smooth Setting of Stochastic Decentralized Convex Optimization Problem Over Time-Varying Graphs / A. Lobanov [et al.] // arXiv preprint arXiv:2307.00392. — 2023.
18. Black-box generation of adversarial text sequences to evade deep learning classifiers / J. Gao [et al.] // 2018 IEEE Security and Privacy Workshops (SPW). — IEEE. 2018. — P. 50–56.

19. *Hazan E., Klivans A., Yuan Y.* Hyperparameter optimization: A spectral approach // arXiv preprint arXiv:1706.00764. — 2017.
20. *Elsken T., Metzen J. H., Hutter F.* Neural architecture search: A survey // The Journal of Machine Learning Research. — 2019. — Vol. 20, no. 1. — P. 1997–2017.
21. Structured evolution with compact architectures for scalable policy optimization / K. Choromanski [et al.] // International Conference on Machine Learning. — PMLR. 2018. — P. 970–978.
22. *Mania H., Guy A., Recht B.* Simple random search of static linear policies is competitive for reinforcement learning // Advances in Neural Information Processing Systems. — 2018. — Vol. 31.
23. High-probability regret bounds for bandit online linear optimization / P. Bartlett [et al.] // Proceedings of the 21st Annual Conference on Learning Theory-COLT 2008. — Omnipress. 2008. — P. 335–342.
24. Regret analysis of stochastic and nonstochastic multi-armed bandit problems / S. Bubeck, N. Cesa-Bianchi, [et al.] // Foundations and Trends® in Machine Learning. — 2012. — Vol. 5, no. 1. — P. 1–122.
25. *Shamir O.* An optimal algorithm for bandit and zero-order convex optimization with two-point feedback // The Journal of Machine Learning Research. — 2017. — Vol. 18, no. 1. — P. 1703–1713.
26. Randomized gradient-free methods in convex optimization / A. Gasnikov [et al.] // arXiv preprint arXiv:2211.13566. — 2022.
27. *Robbins H., Monro S.* A stochastic approximation method // The annals of mathematical statistics. — 1951. — P. 400–407.
28. *Bottou L., Curtis F. E., Nocedal J.* Optimization methods for large-scale machine learning // SIAM review. — 2018. — Vol. 60, no. 2. — P. 223–311.
29. Noisy zeroth-order optimization for non-smooth saddle point problems / D. Dvinskikh [et al.] // International Conference on Mathematical Optimization Theory and Operations Research. — Springer. 2022. — P. 18–33.
30. Gradient Free Methods for Non-Smooth Convex Optimization with Heavy Tails on Convex Compact / N. Kornilov [et al.] // arXiv preprint arXiv:2304.02442. — 2023.
31. *Lobanov A.* Stochastic Adversarial Noise in the “Black Box” Optimization Problem // arXiv preprint arXiv:2304.07861. — 2023.
32. *Lan G.* An optimal method for stochastic composite optimization // Mathematical Programming. — 2012. — Vol. 133, no. 1/2. — P. 365–397.
33. The min-max complexity of distributed stochastic convex optimization with intermittent communication / B. E. Woodworth [et al.] // Conference on Learning Theory. — PMLR. 2021. — P. 4386–4437.
34. *Gorbunov E., Danilova M., Gasnikov A.* Stochastic optimization with heavy-tailed noise via accelerated gradient clipping // Advances in Neural Information Processing Systems. — 2020. — Vol. 33. — P. 15042–15053.

35. *Gasnikov A., Nesterov Y.* Universal fast gradient method for stochastic composit optimization problems // arXiv preprint arXiv:1604.05275. — 2016.
36. *Ajalloeian A., Stich S. U.* On the convergence of SGD with biased gradients // arXiv preprint arXiv:2008.00051. — 2020.
37. *Polyak B. T.* Gradient methods for the minimisation of functionals // USSR Computational Mathematics and Mathematical Physics. — 1963. — Vol. 3, no. 4. — P. 864–878.
38. *Lojasiewicz S.* Une propriété topologique des sous-ensembles analytiques réels // Les équations aux dérivées partielles. — 1963. — Vol. 117. — P. 87–89.
39. *Yue P., Fang C., Lin Z.* On the Lower Bound of Minimizing Polyak-Lojasiewicz functions // arXiv preprint arXiv:2212.13551. — 2022.
40. *Woodworth B. E., Srebro N.* An even more optimal stochastic optimization algorithm: minibatching and interpolation learning // Advances in Neural Information Processing Systems. — 2021. — Vol. 34. — P. 7333–7345.
41. *Rakhlin A., Shamir O., Sridharan K.* Making gradient descent optimal for strongly convex stochastic optimization // Proceedings of the 29th International Conference on Machine Learning. — 2012. — P. 1571–1578.
42. *Hazan E., Kale S.* Beyond the regret minimization barrier: optimal algorithms for stochastic strongly-convex optimization // The Journal of Machine Learning Research. — 2014. — Vol. 15, no. 1. — P. 2489–2512.
43. *Bertsekas D., Tsitsiklis J. N.* Neuro-dynamic programming. — Athena Scientific, 1996.
44. *Stich S. U.* Unified optimal analysis of the (stochastic) gradient method // arXiv preprint arXiv:1907.04232. — 2019.
45. *Lobanov A., Gasnikov A., Stonyakin F.* Highly Smoothness Zero-Order Methods for Solving Optimization Problems under PL Condition // arXiv preprint arXiv:2305.15828. — 2023.
46. *Schmidt M., Roux N. L.* Fast convergence of stochastic gradient descent under a strong growth condition // arXiv preprint arXiv:1308.6370. — 2013.
47. *Srebro N., Sridharan K., Tewari A.* Optimistic rates for learning with a smooth loss // arXiv preprint arXiv:1009.3896. — 2010.
48. Adam can converge without any modification on update rules / Y. Zhang [et al.] // Advances in Neural Information Processing Systems. — 2022. — Vol. 35. — P. 28386–28399.
49. *Vaswani S., Bach F., Schmidt M.* Fast and faster convergence of sgd for over-parameterized models and an accelerated perceptron // The 22nd international conference on artificial intelligence and statistics. — PMLR. 2019. — P. 1195–1204.
50. Sharp analysis of stochastic optimization under global Kurdyka-Lojasiewicz inequality / I. Fatkhullin [et al.] // Advances in Neural Information Processing Systems. — 2022. — Vol. 35. — P. 15836–15848.

51. *Tran T. H., Scheinberg K., Nguyen L. M.* Nesterov accelerated shuffling gradient method for convex optimization // International Conference on Machine Learning. — PMLR. 2022. — P. 21703–21732.
52. Better mini-batch algorithms via accelerated gradient methods / A. Cotter [et al.] // Advances in neural information processing systems. — 2011. — Vol. 24.
53. *Jacot A., Gabriel F., Hongler C.* Neural tangent kernel: Convergence and generalization in neural networks // Advances in neural information processing systems. — 2018. — Vol. 31.
54. *Allen-Zhu Z., Li Y., Liang Y.* Learning and generalization in overparameterized neural networks, going beyond two layers // Advances in neural information processing systems. — 2019. — Vol. 32.
55. Reconciling modern machine-learning practice and the classical bias–variance trade-off / M. Belkin [et al.] // Proceedings of the National Academy of Sciences. — 2019. — Vol. 116, no. 32. — P. 15849–15854.
56. *Dvinskikh D., Gasnikov A.* Decentralized and parallel primal and dual accelerated methods for stochastic convex programming problems // Journal of Inverse and Ill-posed Problems. — 2021. — Vol. 29, no. 3. — P. 385–405.
57. Accelerated gradient methods with absolute and relative noise in the gradient / A. Vasin [et al.] // Optimization Methods and Software. — 2023. — P. 1–50.
58. *Devolder O.* Exactness, inexactness and stochasticity in first-order methods for large-scale convex optimization : PhD thesis / Devolder Olivier. — CORE UCLouvain Louvain-la-Neuve, Belgium, 2013.
59. Accelerated stochastic approximation with state-dependent noise / S. Ilandarideva [et al.] // arXiv preprint arXiv:2307.01497. — 2023.
60. *Nesterov Y.* Introductory lectures on convex optimization: A basic course. Vol. 87. — Springer Science & Business Media, 2003.
61. Gradient-free optimization of highly smooth functions: improved analysis and a new algorithm / A. Akhavan [et al.] // arXiv preprint arXiv:2306.02159. — 2023.
62. *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 $n = \{2, 3\}$

$$\|a_1 + \dots + a_n\|^2 \leq n \|a_1\|^2 + \dots + n \|a_n\|^2. \quad (5)$$

A.2 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.

$$\|\nabla f(x) - \nabla f(y)\| \leq L \|x - y\|, \quad \forall x, y \in \mathbb{R}^d. \quad (6)$$

It is well-known that L -smoothness implies (see e.g., [60])

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

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.3 Wirtinger-Poincare inequality

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

$$\mathbb{E} [f(x + \tau e)^2] \leq \frac{\tau^2}{d} \mathbb{E} [\|\nabla f(x + \tau e)\|^2]. \quad (7)$$

B Proof Theorem 1

In this section, our reasoning will be based on the proof from [40]. Initially, let us formally define a batched biased gradient oracle (see Definition 1):

$$\mathbf{g}^B(x) := \frac{1}{B} \sum_{i=1}^B \mathbf{g}(x, \xi_i), \quad \text{for i.i.d. } \xi_1, \xi_2, \dots, \xi_B \sim \mathcal{D}. \quad (8)$$

Then Algorithm 2 presents a Biased Accelerated Mini-batch Stochastic Gradient Descent (Biased AC-SA method) under the overparameterization condition.

Algorithm 2 Biased AC-SA

Input: Start point $x_0^{ag} = x_0 \in \mathbb{R}^d$, maximum number of iterations $N \in \mathbb{Z}_+$.

Let stepsize $\gamma_k > 0$, parameters $\beta_k, \gamma > 0$, batch size $B \in \mathbb{Z}_+$.

- 1: **for** $k = 0, \dots, N - 1$ **do**
- 2: $\beta_k = 1 + \frac{k}{6}$ and $\gamma_k = \gamma(k + 1)$ for $\gamma = \min \left\{ \frac{1}{12L}, \frac{B}{24L(N+1)}, \sqrt{\frac{BR^2}{Lf^*N^3}} \right\}$
- 3: $x_k^{md} = \beta_k^{-1} x_k + (1 - \beta_k^{-1}) x_k^{ag}$
- 4: $\tilde{x}_{k+1} = x_k - \gamma_k \mathbf{g}_k^B(x_k^{md})$, where $\mathbf{g}_k^B(x_k^{md})$ is defined from (8)
- 5: $x_{k+1} = \min \left\{ 1, \frac{R}{\|\tilde{x}_{k+1}\|} \right\} \tilde{x}_{k+1}$
- 6: $x_{k+1}^{ag} = \beta_k^{-1} x_{k+1} + (1 - \beta_k^{-1}) x_k^{ag}$
- 7: **end for**

Output: x_N^{ag} .

Then it is not hard to show that the following Lemma is also correct for the batched biased gradient oracle (8). Therefore, to avoid repetition, we formulate this lemma without proof, by referring to the original proof.

Lemma 1 (see Lemma 1, [40]). *Let x_{k+1}, x_k and x_k^{md} be updated as in Algorithm 2. Then for any $x \in \{x : \|x\| \leq R\}$*

$$\begin{aligned} \gamma_k \langle \mathbf{g}^B(x_k^{md}), x_{k+1} - x_k^{md} \rangle &\leq \gamma_k \langle \mathbf{g}^B(x_k^{md}), x - x_k^{md} \rangle + \frac{1}{2} \|x - x_k\|^2 \\ &\quad - \frac{1}{2} \|x - x_{k+1}\|^2 - \frac{1}{2} \|x_{k+1} - x_k\|^2. \end{aligned}$$

Next, we provide some auxiliary lemma before presenting proof of Theorem 1.

Lemma 2. *Let function $f(x, \xi)$ satisfy Assumptions 1-2 and the gradient oracle (see Definition 1) satisfy Assumption 4, and let $\mathbf{g}^B(x_k^{md})$ be defined in (8). Then*

$$\mathbb{E} \left[\|\mathbf{g}^B(x_k^{md}) - \nabla f(x_k^{md})\|^2 \right] \leq \frac{8L^2 R^2}{B\beta_k^2} + \frac{8L}{B} \mathbb{E}[f(x_k^{ag}) - f^*] + \frac{4\sigma_*^2}{B} + \|\mathbf{b}(x_k^{md})\|^2$$

Proof.

$$\begin{aligned}
& \mathbb{E} \|\mathbf{g}^B(x_k^{md}) - \nabla f(x_k^{md})\|^2 \\
&= \mathbb{E} \left[\left\| \frac{1}{B} \sum_{i=1}^B \nabla f(x_k^{md}, \xi_i) - \nabla f(x_k^{md}) \right\|^2 \right] \\
&\quad + \mathbb{E} \left[\left\| \mathbf{g}^B(x_k^{md}) - \frac{1}{B} \sum_{i=1}^B \nabla f(x_k^{md}, \xi_i) \right\|^2 \right] \\
&\stackrel{(8)}{=} \mathbb{E} \left[\left\| \frac{1}{B} \sum_{i=1}^B \nabla f(x_k^{md}, \xi_i) - \nabla f(x_k^{md}) \right\|^2 \right] + \mathbb{E} \left[\left\| \frac{1}{B} \sum_{i=1}^B \mathbf{b}(x_k^{md}) \right\|^2 \right] \\
&= \frac{1}{B^2} \sum_{i=1}^B \mathbb{E} \left[\|\nabla f(x_k^{md}, \xi_i) - \nabla f(x_k^{md})\|^2 \right] + \|\mathbf{b}(x_k^{md})\|^2 \\
&\leq \frac{1}{B} \mathbb{E} \left[\|\nabla f(x_k^{md}, \xi_1)\|^2 \right] + \|\mathbf{b}(x_k^{md})\|^2 \\
&\stackrel{(5)}{\leq} \frac{2}{B} \mathbb{E} \left[\|\nabla f(x_k^{md}, \xi_1) - \nabla f(x_k^{ag}, \xi_1)\|^2 \right] + \frac{2}{B} \mathbb{E} \left[\|\nabla f(x_k^{ag}, \xi_1)\|^2 \right] + \|\mathbf{b}(x_k^{md})\|^2 \\
&\stackrel{(6)}{\leq} \frac{2L^2}{B} \mathbb{E} \left[\|x_k^{md} - x_k^{ag}\|^2 \right] + \frac{4}{B} \mathbb{E} \left[\|\nabla f(x_k^{ag}, \xi_1) - \nabla f(x^*, \xi_1)\|^2 \right] \\
&\quad + \frac{4}{B} \mathbb{E} \left[\|\nabla f(x^*, \xi_1)\|^2 \right] + \|\mathbf{b}(x_k^{md})\|^2. \tag{9}
\end{aligned}$$

For the first term on the right hand side:

$$x_k^{md} = \beta_k^{-1} x_k + (1 - \beta_k^{-1}) x_k^{ag} \Rightarrow \|x_k^{md} - x_k^{ag}\| = \beta_k^{-1} \|x_k - x_k^{ag}\| \leq 2R\beta_k^{-1}. \tag{10}$$

For second term, we apply Theorem 2.1.5 from [60]:

$$\begin{aligned}
& \mathbb{E} \|\nabla f(x_k^{ag}, \xi_1) - \nabla f(x^*, \xi_1)\|^2 \\
&\leq 2L \mathbb{E} [f(x_k^{ag}, \xi_1) - f(x^*, \xi_1) - \langle \nabla f(x^*, \xi_1), x_k^{ag} - x^* \rangle] \\
&= 2L \mathbb{E} [f(x_k^{ag}) - f^*]. \tag{11}
\end{aligned}$$

For third term, we apply Assumption 4:

$$\mathbb{E} \left[\|\nabla f(x^*, \xi_1)\|^2 \right] = \mathbb{E} \left[\|\nabla f(x^*, \xi_1) - \nabla f(x^*)\|^2 \right] \leq \sigma_*^2. \tag{12}$$

Substituting (10)-(12) into (9) we complete the proof of the Lemma. \square

We can now proceed to prove the main theorem of Section 3.

Proof of the Theorem 1.

Using the convexity and L -smoothness of the function f we can obtain the following upper bound:

$$\begin{aligned}
\beta_k \gamma_k f(x_{k+1}^{ag}) &\leq \beta_k \gamma_k \left[f(x_k^{md}) + \langle \nabla f(x_k^{md}), x_{k+1}^{ag} - x_k^{md} \rangle + \frac{L}{2} \|x_{k+1}^{ag} - x_k^{md}\|^2 \right] \\
&= \beta_k \gamma_k \left[f(x_k^{md}) + \langle \nabla f(x_k^{md}), x_{k+1}^{ag} - x_k^{md} \rangle \right] + \frac{L\gamma_k}{2\beta_k} \|x_{k+1} - x_k\|^2 \\
&= \beta_k \gamma_k \left[f(x_k^{md}) + \langle \nabla f(x_k^{md}), \beta_k^{-1} x_{k+1} + (1 - \beta^{-1}) x_k^{ag} - x_k^{md} \rangle \right] \\
&\quad + \frac{L\gamma_k}{2\beta_k} \|x_{k+1} - x_k\|^2 \\
&= (\beta_k - 1) \gamma_k \left[f(x_k^{md}) + \langle \nabla f(x_k^{md}), x_k^{ag} - x_k^{md} \rangle \right] \\
&\quad + \gamma_k \left[f(x_k^{md}) + \langle \nabla f(x_k^{md}), x_{k+1} - x_k^{md} \rangle \right] + \frac{L\gamma_k}{2\beta_k} \|x_{k+1} - x_k\|^2 \\
&\leq (\beta_k - 1) \gamma_k f(x_k^{ag}) + \gamma_k \left[f(x_k^{md}) + \langle \mathbf{g}^B(x_k^{md}), x_{k+1} - x_k^{md} \rangle \right] \\
&\quad - \gamma_k \langle \mathbf{g}^B(x_k^{md}) - \nabla f(x_k^{md}), x_{k+1} - x_k^{md} \rangle + \frac{L\gamma_k}{2\beta_k} \|x_{k+1} - x_k\|^2
\end{aligned}$$

Using Lemma 1 with $x = x^* \in \arg \min_{x: \|x\| \leq R} f(x)$ for second term we obtain:

$$\begin{aligned}
&\gamma_k f(x_k^{md}) + \gamma_k \langle \mathbf{g}^B(x_k^{md}), x_{k+1} - x_k^{md} \rangle \\
&= \gamma_k f(x_k^{md}) + \gamma_k \langle \mathbf{g}^B(x_k^{md}), x^* - x_k^{md} \rangle \\
&\quad + \frac{1}{2} \|x^* - x_k\|^2 - \frac{1}{2} \|x^* - x_{k+1}\|^2 - \frac{1}{2} \|x_{k+1} - x_k\|^2 \\
&= \gamma_k f(x_k^{md}) + \gamma_k \langle \nabla f(x_k^{md}), x^* - x_k^{md} \rangle + \gamma_k \langle \mathbf{g}^B(x_k^{md}) - \nabla f(x_k^{md}), x^* - x_k^{md} \rangle \\
&\quad + \frac{1}{2} \|x^* - x_k\|^2 - \frac{1}{2} \|x^* - x_{k+1}\|^2 - \frac{1}{2} \|x_{k+1} - x_k\|^2 \\
&\leq \gamma_k f^* + \gamma_k \langle \mathbf{g}^B(x_k^{md}) - \nabla f(x_k^{md}), x^* - x_k^{md} \rangle \\
&\quad + \frac{1}{2} \|x^* - x_k\|^2 - \frac{1}{2} \|x^* - x_{k+1}\|^2 - \frac{1}{2} \|x_{k+1} - x_k\|^2.
\end{aligned}$$

Substituting the obtained upper bound we obtain:

$$\begin{aligned}
\beta_k \gamma_k f(x_{k+1}^{ag}) &\leq (\beta_k - 1) \gamma_k f(x_k^{ag}) + \gamma_k f^* + \frac{L\gamma_k}{2\beta_k} \|x_{k+1} - x_k\|^2 \\
&\quad + \gamma_k \langle \mathbf{g}^B(x_k^{md}) - \nabla f(x_k^{md}), x^* - x_{k+1} \rangle \\
&\quad + \frac{1}{2} \left(\|x^* - x_k\|^2 - \|x^* - x_{k+1}\|^2 - \|x_{k+1} - x_k\|^2 \right).
\end{aligned}$$

Adding $\beta_k \gamma_k f^*$ to both sides we can obtain:

$$\begin{aligned}
\beta_k \gamma_k [f(x_{k+1}^{ag}) - f^*] &\leq (\beta_k - 1) \gamma_k [f(x_k^{ag}) - f^*] + \frac{1}{2} \|x_k - x^*\|^2 - \frac{1}{2} \|x_{k+1} - x^*\|^2 \\
&\quad + \gamma_k \langle \mathbf{g}^B(x_k^{md}) - \nabla f(x_k^{md}), x^* - x_{k+1} \rangle \\
&\quad + \frac{L\gamma_k - \beta_k}{2\beta_k} \|x_k - x_{k+1}\|^2
\end{aligned}$$

$$\begin{aligned}
&= (\beta_k - 1)\gamma_k [f(x_k^{ag}) - f^*] + \frac{1}{2} \|x_k - x^*\|^2 - \frac{1}{2} \|x_{k+1} - x^*\|^2 \\
&\quad + \gamma_k \langle \mathbf{g}^B(x_k^{md}) - \nabla f(x_k^{md}), x^* - x_k \rangle \\
&\quad + \gamma_k \langle \mathbf{g}^B(x_k^{md}) - \nabla f(x_k^{md}), x_k - x_{k+1} \rangle \\
&\quad + \frac{L\gamma_k - \beta_k}{2\beta_k} \|x_k - x_{k+1}\|^2 \\
&\leq (\beta_k - 1)\gamma_k [f(x_k^{ag}) - f^*] + \frac{1}{2} \|x_k - x^*\|^2 - \frac{1}{2} \|x_{k+1} - x^*\|^2 \\
&\quad + \gamma_k \langle \mathbf{g}^B(x_k^{md}) - \nabla f(x_k^{md}), x^* - x_k \rangle \\
&\quad + \gamma_k \|\mathbf{g}^B(x_k^{md}) - \nabla f(x_k^{md})\| \|x_k - x_{k+1}\| \\
&\quad + \frac{L\gamma_k - \beta_k}{2\beta_k} \|x_k - x_{k+1}\|^2.
\end{aligned}$$

Since $\beta_k = 1 + \frac{k}{6} > \frac{1+k}{6} \geq 2L\gamma_k$, then

$$\gamma_k \|\mathbf{g}^B(x_k^{md}) - \nabla f(x_k^{md})\| \|x_k - x_{k+1}\| + \frac{L\gamma_k - \beta_k}{2\beta_k} \|x_k - x_{k+1}\|^2$$

is a quadratic polynomial of the form: $-\frac{a}{2}y^2 + by$ (where $y = \|x_k - x_{k+1}\|$), which can be upper bounded by $-\frac{a}{2}y^2 + by \leq \max_y \{-\frac{a}{2}y^2 + by\} = \frac{b^2}{2a}$. We get

$$\begin{aligned}
\beta_k \gamma_k [f(x_{k+1}^{ag}) - f^*] &\leq (\beta_k - 1)\gamma_k [f(x_k^{ag}) - f^*] + \frac{1}{2} \|x_k - x^*\|^2 - \frac{1}{2} \|x_{k+1} - x^*\|^2 \\
&\quad + \gamma_k \langle \mathbf{g}^B(x_k^{md}) - \nabla f(x_k^{md}), x^* - x_k \rangle \\
&\quad + \frac{\beta_k \gamma_k^2}{2(\beta_k - L\gamma_k)} \|\mathbf{g}^B(x_k^{md}) - \nabla f(x_k^{md})\|^2 \\
&\leq (\beta_k - 1)\gamma_k [f(x_k^{ag}) - f^*] + \frac{1}{2} \|x_k - x^*\|^2 - \frac{1}{2} \|x_{k+1} - x^*\|^2 \\
&\quad + \gamma_k \langle \mathbf{g}^B(x_k^{md}) - \nabla f(x_k^{md}), x^* - x_k \rangle \\
&\quad + \gamma_k^2 \|\mathbf{g}^B(x_k^{md}) - \nabla f(x_k^{md})\|^2 \\
&\leq (\beta_k - 1)\gamma_k [f(x_k^{ag}) - f^*] + \frac{1}{2} \|x_k - x^*\|^2 - \frac{1}{2} \|x_{k+1} - x^*\|^2 \\
&\quad + \gamma_k \left\langle \frac{1}{B} \sum_{i=1}^B \nabla f(x_k^{md}, \xi_i) - \nabla f(x_k^{md}), x^* - x_k \right\rangle \\
&\quad + \gamma_k \left\langle \mathbf{g}^B(x_k^{md}) - \frac{1}{B} \sum_{i=1}^B \nabla f(x_k^{md}, \xi_i), x^* - x_k \right\rangle \\
&\quad + \gamma_k^2 \|\mathbf{g}^B(x_k^{md}) - \nabla f(x_k^{md})\|^2 \\
&\leq (\beta_k - 1)\gamma_k [f(x_k^{ag}) - f^*] + \frac{1}{2} \|x_k - x^*\|^2 - \frac{1}{2} \|x_{k+1} - x^*\|^2 \\
&\quad + \gamma_k \left\langle \frac{1}{B} \sum_{i=1}^B \nabla f(x_k^{md}, \xi_i) - \nabla f(x_k^{md}), x^* - x_k \right\rangle
\end{aligned}$$

$$+ \gamma_k^2 \|\mathbf{g}^B(x_k^{md}) - \nabla f(x_k^{md})\|^2 + \gamma_k \langle \mathbf{b}(x_k^{md}), x^* - x_k \rangle.$$

Taking the expectation of both sides we have:

$$\begin{aligned} \beta_k \gamma_k \mathbb{E} [f(x_{k+1}^{ag}) - f^*] &\leq (\beta_k - 1) \gamma_k \mathbb{E} [f(x_k^{ag}) - f^*] + \frac{1}{2} \mathbb{E} [\|x_k - x^*\|^2] \\ &\quad - \frac{1}{2} \mathbb{E} [\|x_{k+1} - x^*\|^2] + \gamma_k^2 \mathbb{E} [\|\mathbf{g}^B(x_k^{md}) - \nabla f(x_k^{md})\|^2] \\ &\quad + \gamma_k \langle \mathbf{b}(x_k^{md}), x^* - x_k \rangle. \end{aligned}$$

Using the Lemma 2 we obtain:

$$\begin{aligned} \beta_k \gamma_k \mathbb{E} [f(x_{k+1}^{ag}) - f^*] &\leq (\beta_k - 1) \gamma_k \mathbb{E} [f(x_k^{ag}) - f^*] + \frac{1}{2} \mathbb{E} [\|x_k - x^*\|^2] \\ &\quad - \frac{1}{2} \mathbb{E} [\|x_{k+1} - x^*\|^2] + \frac{8L^2 R^2 \gamma_k^2}{B \beta_k^2} \\ &\quad + \frac{8L \gamma_k^2}{B} \mathbb{E} [f(x_k^{ag}) - f^*] + \frac{4\sigma_*^2 \gamma_k^2}{B} \\ &\quad + \gamma_k \langle \mathbf{b}(x_k^{md}), x^* - x_k \rangle + \gamma_k^2 \|\mathbf{b}(x_k^{md})\|^2 \\ &\leq \left(\beta_k - 1 + \frac{8L \gamma_k}{B} \right) \gamma_k \mathbb{E} [f(x_k^{ag}) - f^*] \\ &\quad + \frac{1}{2} \mathbb{E} [\|x_k - x^*\|^2] - \frac{1}{2} \mathbb{E} [\|x_{k+1} - x^*\|^2] \\ &\quad + \frac{8L^2 R^2 \gamma_k^2}{B \beta_k^2} + \frac{4\sigma_*^2 \gamma_k^2}{B} \\ &\quad + \gamma_k \langle \mathbf{b}(x_k^{md}), x^* - x_k \rangle + \gamma_k^2 \|\mathbf{b}(x_k^{md})\|^2. \end{aligned}$$

We now remind that

$$\begin{aligned} \beta_k &= 1 + \frac{k}{6}; \\ \gamma_k &= \gamma(k+1); \\ \gamma &\leq \min \left\{ \frac{1}{12L}, \frac{B}{24L(N+1)} \right\}. \end{aligned}$$

This ensure that $\forall k : \beta_k \geq 1$ and $2L\gamma_k \leq \beta_k$. Moreover, for $k \in [0; N-1]$:

$$\begin{aligned} &\left(\beta_{k+1} - 1 + \frac{8L \gamma_{k+1}}{B} \right) \gamma_{k+1} - \beta_k \gamma_k \\ &= \left(\beta_k - \frac{5}{6} + \frac{8L \gamma_{k+1}}{B} \right) \gamma(k+2) - \beta_k \gamma(k+1) \\ &= \gamma \left(1 + \frac{k}{6} - \frac{5(k+2)}{6} + \frac{8L \gamma(k+2)^2}{B} \right) \\ &= \gamma \left(-\frac{2}{3} - \frac{2k}{3} + \frac{(k+2)}{3} \cdot \frac{24L(k+2)\gamma}{B} \right) \end{aligned}$$

$$\leq \gamma \left(-\frac{k}{3} \right) \leq 0.$$

Thus we have shown that for $k \in [0; N-1]$: $\left(\beta_{k+1} - 1 + \frac{8L\gamma_{k+1}}{B} \right) \gamma_{k+1} \leq \beta_k \gamma_k$.
Therefore, we can conclude the following:

$$\begin{aligned} \beta_k \gamma_k \mathbb{E} [f(x_{k+1}^{ag}) - f^*] &\leq \left(\beta_k - 1 + \frac{8L\gamma_k}{B} \right) \gamma_k \mathbb{E} [f(x_k^{ag}) - f^*] \\ &\quad + \frac{1}{2} \mathbb{E} [\|x_k - x^*\|^2] - \frac{1}{2} \mathbb{E} [\|x_{k+1} - x^*\|^2] \\ &\quad + \frac{8L^2 R^2 \gamma_k^2}{B \beta_k^2} + \frac{4\sigma_*^2 \gamma_k^2}{B} \\ &\quad + \gamma_k \langle \mathbf{b}(x_k^{md}), x^* - x_k \rangle + \gamma_k^2 \|\mathbf{b}(x_k^{md})\|^2 \\ &\leq \beta_{k-1} \gamma_{k-1} \mathbb{E} [f(x_k^{ag}) - f^*] \\ &\quad + \frac{1}{2} \mathbb{E} [\|x_k - x^*\|^2] - \frac{1}{2} \mathbb{E} [\|x_{k+1} - x^*\|^2] \\ &\quad + \frac{8L^2 R^2 \gamma_k^2}{B \beta_k^2} + \frac{4\sigma_*^2 \gamma_k^2}{B} \\ &\quad + \gamma_k \langle \mathbf{b}(x_k^{md}), x^* - x_k \rangle + \gamma_k^2 \|\mathbf{b}(x_k^{md})\|^2. \end{aligned}$$

Summing the both sides over k we obtain:

$$\begin{aligned} \sum_{k=0}^{N-1} \beta_k \gamma_k \mathbb{E} [f(x_{k+1}^{ag}) - f^*] &\leq \sum_{k=0}^{N-1} \beta_{k-1} \gamma_{k-1} \mathbb{E} [f(x_k^{ag}) - f^*] \\ &\quad + \sum_{k=0}^{N-1} \left(\frac{1}{2} \mathbb{E} [\|x_k - x^*\|^2] - \frac{1}{2} \mathbb{E} [\|x_{k+1} - x^*\|^2] \right) \\ &\quad + \sum_{k=0}^{N-1} \frac{8L^2 R^2 \gamma_k^2}{B \beta_k^2} + \sum_{k=0}^{N-1} \frac{4\sigma_*^2 \gamma_k^2}{B} \\ &\quad + \sum_{k=0}^{N-1} \gamma_k \langle \mathbf{b}(x_k^{md}), x^* - x_k \rangle + \sum_{k=0}^{N-1} \gamma_k^2 \|\mathbf{b}(x_k^{md})\|^2. \end{aligned}$$

Simplifying the expression we have

$$\begin{aligned} \beta_{N-1} \gamma_{N-1} \mathbb{E} [f(x_N^{ag}) - f^*] &\leq \frac{1}{2} \mathbb{E} [\|x_0 - x^*\|^2] + \sum_{k=0}^{N-1} \frac{288L^2 R^2 \gamma^2 (k+1)^2}{B(k+6)^2} + \sum_{k=0}^{N-1} \frac{4\sigma_*^2 \gamma^2 (k+1)^2}{B} \\ &\quad + \sum_{k=0}^{N-1} \gamma(k+1) \langle \mathbf{b}(x_k^{md}), x^* - x_k \rangle + \sum_{k=0}^{N-1} \gamma^2 (k+1)^2 \|\mathbf{b}(x_k^{md})\|^2 \\ &\leq \frac{R^2}{2} + \frac{288L^2 R^2 \gamma^2 N}{B} + \frac{12\sigma_*^2 \gamma^2 N^3}{B} + 2\gamma N^2 R \zeta + \gamma^2 N^3 \zeta^2, \end{aligned}$$

where $\|\mathbf{b}(x_k^{md})\|^2 \leq \zeta^2$. Divide the left and right side by $\beta_{N-1}\gamma_{N-1} \simeq \gamma N^2$:

$$\mathbb{E}[f(x_N^{ag}) - f^*] \leq \frac{R^2}{2\gamma N^2} + \frac{288L^2R^2\gamma}{BN} + \frac{12\sigma_*^2\gamma N}{B} + 2R\zeta + \gamma N\zeta^2.$$

With our choice of $\gamma = \min\left\{\frac{1}{12L}, \frac{B}{24L(N+1)}, \sqrt{\frac{BR^2}{\sigma_*^2 N^3}}\right\}$ we obtain:

$$\mathbb{E}[f(x_N^{ag}) - f^*] \lesssim \frac{LR^2}{N^2} + \frac{LR^2}{BN} + \frac{\sigma_* R}{\sqrt{BN}} + \zeta R + \frac{\zeta^2 N}{2L}.$$

□

C Proof Theorem on the convergence of AZO-SGD

In this section, we present a detailed proof of the results of Theorem 2. First, we find the bias and the second moment of the gradient approximation (3) based on the improved analysis of the paper [61].

Bias of gradient approximation Using the variational representation of the Euclidean norm, and definition of gradient approximation (3) we can write:

$$\begin{aligned} \|\mathbb{E}[\mathbf{g}(x_k, \xi, e)] - \nabla f(x_k)\| &= \left\| \mathbb{E} \left[\frac{d}{2\tau} (f_\delta(x_k + \tau e, \xi) - f_\delta(x_k - \tau e, \xi)) e \right] - \nabla f(x_k) \right\| \\ &\stackrel{\textcircled{1}}{=} \left\| \mathbb{E} \left[\frac{d}{\tau} (f(x_k + \tau e, \xi) + \delta(x_k + \tau e)) e \right] - \nabla f(x_k) \right\| \\ &\stackrel{\textcircled{2}}{\leq} \left\| \mathbb{E} \left[\frac{d}{\tau} f(x_k + \tau e, \xi) e \right] - \nabla f(x_k) \right\| + \frac{d\Delta}{\tau} \\ &\stackrel{\textcircled{3}}{=} \|\mathbb{E}[\nabla f(x_k + \tau u, \xi)] - \nabla f(x_k)\| + \frac{d\Delta}{\tau} \\ &= \sup_{z \in S_2^d(1)} \mathbb{E}[\|\nabla_z f(x_k + \tau u, \xi) - \nabla_z f(x_k)\|] + \frac{d\Delta}{\tau} \\ &\stackrel{\textcircled{6}}{\leq} L\tau \mathbb{E}[\|u\|] + \frac{d\Delta}{\tau} \\ &\leq L\tau + \frac{d\Delta}{\tau}, \end{aligned} \tag{13}$$

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 inequality is obtain from bounded noise $|\delta(x)| \leq \Delta$, $\textcircled{3}$ = the equality is obtained from a version of Stokes' theorem [62].

Bounding second moment of gradient approximation By definition gradient approximation (3) and Wirtinger-Poincare inequality (7) we have

$$\mathbb{E}[\|\mathbf{g}(x^*, \xi, e)\|^2] = \frac{d^2}{4\tau^2} \mathbb{E}[\|(f_\delta(x^* + \tau e, \xi) - f_\delta(x^* - \tau e, \xi)) e\|^2]$$

$$\begin{aligned}
&= \frac{d^2}{4\tau^2} \mathbb{E} \left[(f(x^* + \tau e, \xi) - f(x^* - \tau e, \xi) + \delta(x^* + \tau e) - \delta(x^* - \tau e))^2 \right] \\
&\stackrel{(5)}{\leq} \frac{d^2}{2\tau^2} \left(\mathbb{E} \left[(f(x^* + \tau e, \xi) - f(x^* - \tau e, \xi))^2 \right] + 2\Delta^2 \right) \\
&\stackrel{(7)}{\leq} \frac{d^2}{2\tau^2} \left(\frac{\tau^2}{d} \mathbb{E} \left[\|\nabla f(x^* + \tau e, \xi) + \nabla f(x^* - \tau e, \xi)\|^2 \right] + 2\Delta^2 \right) \\
&= \frac{d^2}{2\tau^2} \left(\frac{\tau^2}{d} \mathbb{E} \left[\|\nabla f(x^* + \tau e, \xi) + \nabla f(x^* - \tau e, \xi) \pm 2\nabla f(x^*, \xi)\|^2 \right] + 2\Delta^2 \right) \\
&\stackrel{(6)}{\leq} 4d \|\nabla f(x^*, \xi)\|^2 + 4dL^2\tau^2 \mathbb{E} \left[\|e\|^2 \right] + \frac{d^2\Delta^2}{\tau^2} \\
&\stackrel{\textcircled{1}}{\leq} 4d\sigma_*^2 + 4dL^2\tau^2 \mathbb{E} \left[\|e\|^2 \right] + \frac{d^2\Delta^2}{\tau^2}, \tag{14}
\end{aligned}$$

where $\textcircled{1}$ = the inequality is obtain from Assumption 4.

We can now explicitly write down the convergence of the gradient-free AZO-SGD method (see Section 4, Algorithm 1) by substituting upper bounds on the bias (13) and second moment (14) for the gradient approximation (3) in the convergence of the first-order method: Biased AC-SA (see Theorem 1):

$$\begin{aligned}
\mathbb{E} [f(x_N^{ag}) - f^*] &\lesssim \underbrace{\frac{LR^2}{N^2}}_{\textcircled{1}} + \underbrace{\frac{LR^2}{BN}}_{\textcircled{2}} + \underbrace{\frac{\sqrt{d}\sigma_*R}{\sqrt{BN}}}_{\textcircled{3}} + \underbrace{\frac{\sqrt{d}L\tau R}{\sqrt{BN}}}_{\textcircled{4}} + \underbrace{\frac{d\Delta R}{\tau\sqrt{BN}}}_{\textcircled{5}} \\
&\quad + \underbrace{L\tau R}_{\textcircled{6}} + \underbrace{\frac{d\Delta R}{\tau}}_{\textcircled{7}} + \underbrace{L\tau^2 N}_{\textcircled{8}} + \underbrace{\frac{d^2\Delta^2 N}{\tau^2 L}}_{\textcircled{9}}.
\end{aligned}$$

Proof of the Theorem 2.

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

$$\begin{aligned}
\textcircled{1} : \quad \frac{LR^2}{N^2} \leq \varepsilon &\Rightarrow N \geq \sqrt{\frac{LR^2}{\varepsilon}}; \\
N &= \mathcal{O} \left(\sqrt{\frac{LR^2}{\varepsilon}} \right). \tag{15}
\end{aligned}$$

From terms $\textcircled{2}$ and $\textcircled{3}$, we find the batch size B required to achieve optimality in iteration complexity N :

$$\begin{aligned}
\textcircled{2} : \quad \frac{LR^2}{BN} \leq \varepsilon &\Rightarrow B \geq \frac{LR^2}{\varepsilon N} \stackrel{(15)}{=} \mathcal{O} \left(\sqrt{\frac{LR^2}{\varepsilon}} \right); \\
\textcircled{3} : \quad \frac{\sqrt{d}\sigma_*R}{\sqrt{BN}} &\Rightarrow B \geq \frac{d\sigma_*^2 R^2}{\varepsilon^2 N} \stackrel{(15)}{=} \mathcal{O} \left(\frac{d\sigma_*^2 R}{\varepsilon^{3/2} L^{1/2}} \right); \\
B &= \max \left\{ \mathcal{O} \left(\sqrt{\frac{LR^2}{\varepsilon}} \right), \mathcal{O} \left(\frac{d\sigma_*^2 R}{\varepsilon^{3/2} L^{1/2}} \right) \right\}. \tag{16}
\end{aligned}$$

From terms ④, ⑥ and ⑧ we find the smoothing parameter τ :

$$\begin{aligned}
\text{④: } \quad & \frac{\sqrt{dL}\tau R}{\sqrt{BN}} \leq \varepsilon \Rightarrow \tau \leq \frac{\varepsilon\sqrt{BN}}{\sqrt{dLR}} \stackrel{(15),(16)}{=} \max\left\{\sqrt{\frac{\varepsilon}{dL}}, \frac{\sigma_*}{L}\right\}; \\
\text{⑥: } \quad & L\tau R \leq \varepsilon \Rightarrow \tau \leq \frac{\varepsilon}{LR}; \\
\text{⑧: } \quad & L\tau^2 N \leq \varepsilon \Rightarrow \tau \leq \sqrt{\frac{\varepsilon}{LN}} \stackrel{(15)}{=} \frac{\varepsilon^{3/4}}{L^{3/4}R^{1/2}}; \\
& \tau \leq \min\left\{\max\left\{\sqrt{\frac{\varepsilon}{dL}}, \frac{\sigma_*}{L}\right\}, \frac{\varepsilon}{LR}, \frac{\varepsilon^{3/4}}{L^{3/4}R^{1/2}}\right\} = \frac{\varepsilon}{LR}. \tag{17}
\end{aligned}$$

From the remaining terms ⑤, ⑦, and ⑨, we find the maximum allowable level of adversarial noise Δ that still guarantees the convergence of the Accelerated Zero-Order Stochastic Gradient Descent Method to desired accuracy ε :

$$\begin{aligned}
\text{⑤: } \quad & \frac{d\Delta R}{\tau\sqrt{BN}} \leq \varepsilon \Rightarrow \Delta \leq \frac{\varepsilon\tau\sqrt{BN}}{dR} \stackrel{(15),(16),(17)}{=} \max\left\{\frac{\varepsilon^{3/2}}{d\sqrt{LR}}, \frac{\varepsilon\sigma_*}{dLR}\right\}; \\
\text{⑦: } \quad & \frac{d\Delta R}{\tau} \leq \varepsilon \Rightarrow \Delta \leq \frac{\varepsilon\tau}{dR} \stackrel{(17)}{=} \frac{\varepsilon^2}{dLR}; \\
\text{⑨: } \quad & \frac{d^2\Delta^2 N}{\tau^2 L} \leq \varepsilon \Rightarrow \Delta \leq \sqrt{\frac{\varepsilon\tau^2 L}{d^2 N}} \stackrel{(15),(17)}{=} \frac{\varepsilon^{7/4}}{dL^{3/4}R^{3/2}}; \\
& \Delta \leq \min\left\{\max\left\{\frac{\varepsilon^{3/2}}{d\sqrt{LR}}, \frac{\varepsilon\sigma_*}{dLR}\right\}, \frac{\varepsilon^2}{dLR}, \frac{\varepsilon^{7/4}}{dL^{3/4}R^{3/2}}\right\} = \frac{\varepsilon^2}{dLR}. \tag{18}
\end{aligned}$$

In this way, the Accelerated Zero-Order Stochastic Gradient Descent (AZO-SGD) Method (see Algorithm 1) achieves ε -accuracy: $\mathbb{E}[f(x_N^{ag}) - f^*] \leq \varepsilon$ after

$$N = \mathcal{O}\left(\sqrt{\frac{LR^2}{\varepsilon}}\right), \quad T = N \cdot B = \max\left\{\mathcal{O}\left(\frac{LR^2}{\varepsilon}\right), \mathcal{O}\left(\frac{d\sigma_*^2 R^2}{\varepsilon^2}\right)\right\}$$

number of iterations (15), total number of gradient-free oracle calls (16) and at

$$\Delta \leq \frac{\varepsilon^2}{dLR^2}$$

the maximum level of noise (18) with smoothing parameter $\tau = \frac{\varepsilon}{LR}$ (17). \square