Maximum Noise Level as Third Optimality Criterion in Black-Box Optimization Problem

Anonymous Authors¹

Abstract

This paper is devoted to the study (common in many applications) of the black-box optimization problem, where the black-box represents a gradient-free oracle $f_p = f(x) + \xi_p$ providing the objective function value with some stochastic noise. Assuming that the objective function is μ -strongly convex, and also not just L-smooth, but has a higher order of smoothness ($\beta > 2$) we provide a novel optimization method: Zero-Order Accelerated Batched Stochastic Gradient Descent, whose theoretical analysis closes the question regarding the iteration complexity, achieving optimal estimates. Moreover, we provide a thorough analysis of the maximum noise level, and show under which condition the maximum noise level will take into account information about batch size B as well as information about the smoothness order of the function β . Finally, we show the importance of considering the maximum noise level Δ as a third optimality criterion along with the standard two on the example of a numerical experiment of interest to the machine learning community, where we compare with state-of-theart gradient-free algorithms.

1. Introduction

This paper focuses on solving a standard optimization problem:

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

where $f: Q \to \mathbb{R}$ is function that we want to minimize, f^* is the solution, which we want to find. It is known that if there are no obstacles to compute the gradient of the objective function f or to compute a higher order of the derivative

of the function, then optimal first- or higher-order optimizations algorithms (Nesterov, 2003) should be used to solve the original optimization problem (1). However, if computing the function gradient $\nabla f(x)$ is impossible for any reason, then perhaps the only way to solve the original problem is to use gradient-free (zero-order) optimization algorithms (Conn et al., 2009; Rios & Sahinidis, 2013). Among the situations in which information about the derivatives of the objective function is unavailable are the following:

- a) non-smoothness of the objective function. This situation is probably the most widespread among theoretical works (Gasnikov et al., 2022; Alashqar et al., 2023; Kornilov et al., 2024);
- b) the desire to save computational resources, i.e., computing the gradient $\nabla f(x)$ can sometimes be much "more expensive" than computing the objective function value f(x). This situation is quite popular and extremely understandable in the real world (Bogolubsky et al., 2016);
- c) *inaccessibility of the function gradient*. A vivid example of this situation is the problem of creating an ideal product for a particular person (Lobanov et al., 2024).

Like first-order optimization algorithms, gradient-free algorithms have the following optimality criteria: #N – the number of consecutive iterations required to achieve the desired accuracy of the solution ε and #T – the total number of calls (in this case) to the gradient-free oracle, where by gradient-free/derivative-free oracle we mean that we have access only to the objective function f(x) with some bounded stochastic noise ξ_p ($\mathbb{E}[\xi_p^2] \leq \Delta^2$). It should be noted that because the objective function is subject to noise, the gradient-free oracle plays the role of a black box. That is why there is a tendency in the literature when the initial problem formulation (1) with a gradient-free oracle is called a black-box optimization problem (Kimiaei & Neumaier, 2022). However, unlike higher-order algorithms, gradient-free algorithms have a third optimality criterion: the maximum noise level Δ at which the algorithm will still converge "good", where by "good convergence" we mean convergence as in the case when $\Delta = 0$. The existence of

¹Anonymous Institution, Anonymous City, Anonymous Region, Anonymous Country. Correspondence to: Anonymous Author <anon.email@domain.com>.

Preliminary work. Under review by the International Conference on Machine Learning (ICML). Do not distribute.



Figure 1. Motivation to find the maximum noise level Δ

067 such a seemingly unusual criterion can be explained by the 068 following motivational examples (see Figure 1*). Among 069 the motivations we can highlight the most demanded espe-070 cially by companies (and not only). Resource saving (Fig-071 ure 1a): The more accurately the objective function value is calculated, the more expensive this process to be performed. 073 Robustness to Attacks (Figure 1b): Improving the maximum 074 noise level makes the algorithm more robust to adversarial 075 attacks. Confidentiality (Figire 1c): Some companies, due 076 to secrecy, can't hand over all the information. Therefore, 077 it is important to be able to answer the following question: 078 How much can the objective function be noisy? 079

063 064

065 066

109

The basic idea to create algorithms with a gradient-free ora-081 cle that will be efficient according to the above three criteria 082 is to take advantage of first-order algorithms by substituting 083 a gradient approximation instead of the true gradient (Gas-084 nikov et al., 2023). The choice of the first-order optimization algorithm depends on the formulation of the original 086 problem (on the Assumptions on the function and the gra-087 dient oracle). But the choice of gradient approximation 088 depends on the smoothness of the function. For example, 089 if the function is non-smooth, a smoothing scheme with 090 l_1 randomization (Alashqar et al., 2023; Lobanov, 2023) 091 or with l_2 randomization (Dvinskikh et al., 2022; Lobanov 092 et al., 2023a;b) should be used to solve the original prob-093 lem. If the function is smooth, it is enough to use choose l_1 094 randomization (Akhavan et al., 2022) or l_2 randomization 095 (Gorbunov et al., 2018; Lobanov & Gasnikov, 2023). But 096 if the objective function is not just smooth but also has a 097 higher order of smoothness ($\beta \geq 2$), then the so-called Ker-098 nel approximation (Akhavan et al., 2023; Gasnikov et al., 099 2024b;a), which takes into account the information about 100 the increased smoothness of the function using two-point feedback, should be used as the gradient approximation.

In this paper, we consider the black-box optimization problem (1), assuming strong convexity as well as increased smoothness of the objective function. We choose accelerated stochastic gradient descent (Vaswani et al., 2019) as the basis for a gradient-free algorithm. Since the Kernel approx-

*The pictures are taken from the following resource

imation (which accounts for the advantages of increased smoothness) is biased, we generalize the result of (Vaswani et al., 2019) to the biased gradient oracle. We use the resulting accelerated stochastic gradient descent with a biased gradient oracle to create a gradient-free algorithm. Finally, we explicitly derive estimates on the three optimality criteria of the gradient-free algorithm.

1.1. Main Assumptions and Notations

Since the original problem (1) is general, in this subsection we further define the problem by imposing constraints on the objective function as well as the zero-order oracle. In particular, we assume that the function f is not just *L*-smooth, but has increased smoothness, and is also μ -strongly convex.

Assumption 1.1 (Higher order smoothness). Let l denote maximal integer number strictly less than β . Let $\mathcal{F}_{\beta}(L)$ denote the set of all functions $f : \mathbb{R}^d \to \mathbb{R}$ which are differentiable l times and $\forall x, z \in Q$ the Hölder-type condition:

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

where $l < \beta$ (β is smoothness order), $L_{\beta} > 0$, the sum is over multi-index $n = (n_1, ..., n_d) \in \mathbb{N}^d$, we used the notation $n! = n_1! \cdots n_d!$, $|n| = n_1 + \cdots + n_d$, $\forall v = (v_1, ..., v_d) \in \mathbb{R}^d$, and we defined $D^n f(x) v^n = \frac{\partial^{|n|} f(x)}{\partial^{n_1} x_1 \cdots \partial^{n_d} x_d} v_1^{n_1} \cdots v_d^{n_d}$.

Assumption 1.2 (Strongly convex). Function $f : \mathbb{R}^d \to \mathbb{R}$ is μ -strongly convex with some constant $\mu > 0$ if for any $x, y \in \mathbb{R}^d$ it holds that

$$f(y) \ge f(x) + \langle \nabla f(x), y - x \rangle + \frac{\mu}{2} \left\| y - x \right\|^2.$$

Assumption 1.1 is commonly appeared in papers (Bach & Perchet, 2016; Akhavan et al., 2023), which consider the case when the objective function has smoothness order $\beta \ge 2$. It is worth noting that the smoothness constant L_{β} in the case when $\beta = 2$ has the following relation with the standard Lipschitz gradient constant $L = 2 \cdot L_2$. In

Maximum Noise Level as Third Optimality Criterion

J, J	8	·····	
References	Iteration Complexity $(\#N)$	Oracle Complexity $(\#T)$	Maximum Noise Level (Δ)
(Bach & Perchet, 2016)	$\mathcal{O}\left(rac{d^{2+rac{2}{eta-1}}\Delta^2}{\muarepsilon^{rac{eta+1}{eta-1}}} ight)$	$\mathcal{O}\left(rac{d^{2+rac{2}{eta-1}}\Delta^2}{\muarepsilon^{rac{eta+1}{eta-1}}} ight)$	×
(Akhavan et al., 2020)	$ ilde{\mathcal{O}}\left(rac{d^{2+rac{2}{eta=1}}\Delta^{2}}{(\muarepsilon)^{rac{eta}{eta=1}}} ight)$	$ ilde{\mathcal{O}}\left(rac{d^{2+rac{2}{eta=1}}\Delta^2}{(\muarepsilon)^{rac{eta}{eta=1}}} ight)$	×
(Novitskii & Gasnikov, 2021)	$ ilde{\mathcal{O}}\left(rac{d^{2+rac{1}{eta-1}}\Delta^{2}}{(\muarepsilon)^{rac{eta}{eta-1}}} ight)$	$ ilde{\mathcal{O}}\left(rac{d^{2+rac{1}{eta-1}}\Delta^2}{(\muarepsilon)^{eta-1}} ight)$	×
(Akhavan et al., 2023)	$ ilde{\mathcal{O}}\left(rac{d^2\Delta^2}{(\muarepsilon)^{rac{eta}{eta-1}}} ight)$	$ ilde{\mathcal{O}}\left(rac{d^2\Delta^2}{(\muarepsilon)^{rac{eta}{eta-1}}} ight)$	×
Theorem 3.1 (Our work)	$\mathcal{O}\left(\sqrt{\frac{L}{\mu}}\log \frac{1}{\varepsilon} ight)$	$\max\left\{\tilde{\mathcal{O}}\left(\sqrt{\frac{d^2L}{\mu}}\right),\tilde{\mathcal{O}}\left(\frac{d^2\Delta^2}{(\mu\varepsilon)^{\frac{\beta}{\beta-1}}}\right)\right\}$	1

Table 1. Overview of convergence results of previous works. Notations: d = dimensionality of the problem (1); $\beta =$ smoothness order of the objective function f; $\mu =$ strong convexity constant; $\varepsilon =$ desired accuracy of the problem solution by function.

addition, Assumption 1.2 is standard among optimization works (Nesterov, 2003; Stich, 2019).

In this paper, we assume that Algorithm 1 (which will be introduced later) only has access to the zero-order oracle, which has the following definition.

Definition 1.3 (Zero-order oracle). The zero-order oracle \tilde{f}_p returns only the objective function value $f(x_k)$ at the requested point x_k with stochastic noise ξ_p :

$$\tilde{f}_p(x_k) = f(x_k) + \xi_p$$

where $p \in \{1, 2\}$ and we suppose that the following assumptions on stochastic noise hold

- $\xi_1 \neq \xi_2$ such that $\mathbb{E}[\xi_1^2] \leq \Delta^2$ and $\mathbb{E}[\xi_2^2] \leq \Delta^2$, where $\Delta \geq 0$ is level noise;
- the random variables ξ_1 and ξ_2 are independent from $\mathbf{e} \in S^d(1)$ is a random vector uniformly distributed on the Euclidean unit sphere, and r is a random value uniformly distributed on the interval.

We impose constraints on the Kernel function.

Assumption 1.4 (Kernel function). Let function $K: [-1, 1] \rightarrow \mathbb{R}$ satisfying:

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

155 Definition 1.3 is common among gradient-free works 156 (Lobanov, 2023). In particular, a zero-order oracle will 157 produce the exact function value when the noise level is 158 0. We would also like to point out that we relaxed the re-159 striction on stochastic noise by not assuming a zero mean. 160 We only need the assumption that random variables ξ_1 and 161 ξ_2 are independent from e and r. Assumption 1.4 is often 162 found in papers using the gradient approximation – Kernel 163 approximation. An example of such a function is weighted 164 sums of Lejandre polynomial (Bach & Perchet, 2016). **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 in \mathbb{R}^d as $||x|| := \sqrt{\langle x, x \rangle}$. We use the notation $B^d(r) := \{x \in \mathbb{R}^d : ||x|| \le r\}$ to denote Euclidean ball, $S^d(r) := \{x \in \mathbb{R}^d : ||x|| = r\}$ to denote Euclidean sphere. Operator $\mathbb{E}[\cdot]$ denotes full expectation.

1.2. Related Works and Our Contributions

In Table 1, we provide an overview of the convergence results of the most related works, in particular we provide estimates on the iteration complexity. Research studying the problem (1) with a zero-order oracle (see Definition 1.3), assuming that the function f has increased smoothness $(\beta > 2)$, see Assumption 1.1) comes from (Polyak & Tsybakov, 1990). After 20-30 years, this problem has received widespread attention. However, as we can see, all previous works "fought" (improved/considered) exclusively for oracle complexity (which matches the iteration complexity), without paying attention to other optimality criteria of the gradient-free algorithm. In this paper, we ask another question: Is estimation on iteration complexity unimprovable? And as we can see from Table 1 or Theorem 3.1, we significantly improve the iteration complexity without worsening the oracle complexity, and also provide the best estimates among those we have seen on Δ .

More specifically, our contributions are the following:

- We provide a detailed explanation of the technique for creating a gradient-free algorithm that takes advantage of the increased smoothness of the function via Kernel approximation.
- We generalize existing convergence results for accelerated stochastic gradient descent to the case where the gradient oracle is biased, thereby demonstrating how bias accumulates in the convergence of the algorithm. This result may be of independent interest.
- We close the question regarding the iteration complex-

- ity search by providing an improved estimate (see Ta-ble 1) that is, we provide an optimal estimate.
- 167 168 • We find the maximum noise level Δ at which the al-169 gorithm will still achieve the desired accuracy ε (see 170 Table 1 and Theorem 3.1). Moreover, we show that 171 if overbatching is done, the positive effect on the er-172 ror floor is preserved in a strongly convex problem 173 formulation.

175

176

177

178

179

 We show the importance of considering the maximum noise level Δ as a third optimality criterion along with standard two using an example of a numerical experiment of interest for ML (logistic regression problem).

180 Paper Organization This paper has the following struc-181 ture. In Section 2, we present a first-order algorithm on 182 the basis of which a novel gradient-free algorithm will be 183 created. And in Section 3 we provide the main result of this 184 paper, namely the convergence results of the novel acceler-185 ated gradient-free optimization algorithm. In Section 4, we 186 provide experiments. While Section 5 concludes this paper. 187 The missing proofs of the paper are presented in Appendix. 188

¹⁸⁹ 2. Search for First-Order Algorithm as a Base

190 As mentioned earlier, the basic idea of creating a gradient-191 free algorithm is to take advantage of first-order algorithms. 192 That is, in this subsection, we find the first-order algorithm 193 on which we will base to create a novel gradient-free algorithm by replacing the true gradient with a gradient approxi-195 mation. Since gradient approximations use randomization 196 on the sphere $e(e.g., l_1, l_2$ randomization, or Kernel approx-197 imation), it is important to look for a first-order algorithm 198 that solves a stochastic optimization problem (due to the 199 artificial stochasticity of e). Furthermore, since the gradient 200 approximation from a zero-order oracle concept has a bias, 201 it is also important to find a first-order algorithm that will 202 use a biased gradient oracle. Using these criteria, we formulate an optimization problem to find the most appropriate 204 first-order algorithm. 205

206 2.1. Statement Problem

208

209

210

211

212 213 214 Due to the presence of artificial stochasticity in the gradient approximation, we reformulate the original optimization problem as follows:

$$f^* = \min_{x \in Q \subseteq \mathbb{R}^d} \left\{ f(x) := \mathbb{E} \left[f(x, \omega) \right] \right\}.$$
(2)

We assume that the function satisfies the *L*-smoothness assumption, since it is a basic assumption in papers on firstorder optimization algorithms.

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

it holds $\forall x, y \in \mathbb{R}^d$

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

Next, we define a biased gradient oracle that uses a firstorder algorithm.

Definition 2.2 (Biased Gradient Oracle). A map $\mathbf{g} : \mathbb{R}^d \times \mathcal{D} \to \mathbb{R}^d$ s.t.

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

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

We assume that the bias and gradient noise are bounded.

Assumption 2.3 (Bounded bias). There exists constant $\delta \ge 0$ such that $\forall x \in \mathbb{R}^d$ the following inequality holds

$$\|\mathbf{b}(x)\| = \|\mathbb{E}\left[\mathbf{g}(x,\omega)\right] - \nabla f(x)\| \le \delta.$$
(3)

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

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

Assumption 2.3 is standard for analysis, bounding bias. Assumption 2.4 is a more general condition for strong growth due to the presence of σ^2 .

2.2. First-Order Algorithm as a Base

Now that the problem is formally defined (see Subsection 2.1), we can find an appropriate first-order algorithm. Since one of the main goals of this research is to improve iteration complexity, we have to look for a accelerated batched first-order optimization algorithm. And the most appropriate optimization algorithm which has the following update rule:

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

has the following convergence rate presented in Lemma 2.5.

Lemma 2.5 ((Vaswani et al., 2019), Theorem 1). Let the function f satisfy Assumption 1.2 and 2.1, and the gradient oracle (see Definition 2.2 with $\delta = 0$) satisfy Assumptions 2.3 and 2.4, then with $\tilde{\rho} = \max\{1, \rho\}$ and with the chosen parameters $\gamma_k, a_{k+1}, \alpha_k, \eta$ the Accelerated Stochastic Gradient Descent has the following convergence rate:

$$F_{N} \leq \left(1 - \sqrt{\frac{\mu}{\tilde{\rho}^{2}L}}\right)^{N} \left[f(x_{0}) - f^{*} + \frac{\mu}{2} \|x_{0} - x^{*}\|^{2}\right] \\ + \frac{\sigma^{2}}{\sqrt{\tilde{\rho}^{2}\mu L}},$$

where $F_N = \mathbb{E}[f(x_N)] - f^*$.

As can be seen from Lemma 2.5, that the presented First Order Accelerated Algorithm is not appropriate for creating a gradient-free algorithm, since this algorithm uses an unbiased gradient oracle, and also does not use the batching technique. Therefore, we are ready to present one of the significant results of this work, namely generalizing the results of Lemma 2.5 to the case with an biased gradient oracle and also adding batching (where *B* is a batch size).



Figure 2. Bias influence on the algorithm convergence

237

238

239

240

241

242

243 244 245

246

247

248

249

250

251

252

253

254

255

256 257

258

259

261

262

263

264

265

266

268

269

270

271

272

273

274

Theorem 2.6. Let the function f satisfy Assumption 1.2 and 2.1, and the gradient oracle (see Definition 2.2) satisfy Assumptions 2.3 and 2.4, then with $\tilde{\rho}_B = \max\{1, \frac{\rho}{B}\}$ and with the chosen parameters $\gamma_k = \frac{1}{\sqrt{2\mu\eta\rho}}, \ \beta_k = 1 - \frac{\mu\eta}{2\rho}, \ b_{k+1} = \frac{\sqrt{2\mu}}{(1 - \sqrt{\frac{\mu\eta}{2\rho}})^{(k+1)/2}}, \ a_{k+1} = \frac{1}{(1 - \sqrt{\frac{\mu\eta}{2\rho}})^{(k+1)/2}}, \ \alpha_k = \frac{\gamma_k \beta_k b_{k+1}^2 \eta}{\gamma_k \beta_k b_{k+1}^2 \eta + 2a_k^2}, \ \eta \leq \frac{1}{2\rho L}$ Accelerated Stochastic Gradient Descent with batching has the following convergence rate $(F_N = \mathbb{E}[f(x_N)] - f^*)$:

$$F_{N} \leq \left(1 - \sqrt{\frac{\mu}{\tilde{\rho}_{B}^{2}L}}\right)^{N} \left[f(x_{0}) - f^{*} + \frac{\mu}{2} \|x_{0} - x^{*}\|^{2}\right] \\ + \frac{\sigma^{2}}{\sqrt{\tilde{\rho}_{B}^{2}\mu LB^{2}}} + \left(1 - \sqrt{\frac{\mu}{\tilde{\rho}_{B}^{2}L}}\right)^{N} \tilde{R}\delta + \frac{\delta^{2}}{\sqrt{4\mu L}},$$

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

As can be seen from Theorem 2.6, this result is very similar to the result of Lemma 2.5, moreover, they will be the same if we take $\delta = 0$ and B = 1. It is also worth noting that the third summand does not affect convergence much (the noise does not accumulate due to the decreasing sequence), so we will not consider it in the future for simplicity. Finally, it is worth noting that the Algorithm presented in (Vaswani et al., 2019) can converge as closely as possible to the problem solution (see the red line in Figure 2), while the Algorithm using the biased gradient oracle can only converge to the error floor (see the blue line in Figure 2). This is explained by the last summand from Theorem 2.6. However, convergence to the error floor opens questions about how this asymptote can be controlled. And as shown in (Gasnikov et al., 2024a), the convergence of gradient-free algorithms to the asymptote depends directly on the noise level: the more noise, the better the algorithm can achieve the error

floor. This fact is another clear motivation for finding the maximum noise level. For a detailed proof of Theorem 2.6, see the supplementary materials (Appendix B).

3. Zero-Order Accelerated Batched SGD

Now that we have a proper first-order algorithm, we can move on to creating a novel gradient-free algorithm. To do this, we need to use the gradient approximation instead of the gradient oracle. In this work, we are going to use exactly the Kernel approximation because it takes into account the advantages of increased smoothness of the function, and which has the following

$$\mathbf{g}(x,\mathbf{e}) = d\frac{\tilde{f}_1(x+hr\mathbf{e}) - \tilde{f}_2(x-hr\mathbf{e})}{2h}K(r)\mathbf{e}, \quad (5)$$

where h > 0 is a smoothing parameter, $\mathbf{e} \in S^d(1)$ is a random vector uniformly distributed on the Euclidean unit sphere, r is a random value uniformly distributed on the interval $r \in [0, 1], K : [-1, 1] \rightarrow \mathbb{R}$ is a Kernel function. Then a novel gradient-free method aimed at solving the original problem (1) is presented in Algorithm 1. The missing hyperparameters are given in the Theorem 2.6.

Algorithm 1 Zero-Order Accelerated Batched Stochastic Gradient Descent (ZO-ABSGD)

Input: iteration number N, batch size B, Kernel $K : [-1,1] \to \mathbb{R}$, step size η , smoothing parameter h, $x_0 = y_0 = z_0 \in \mathbb{R}^d$, $a_0 = 1$, $\rho = 4d\kappa$. for k = 0 to N - 1 do 1. Sample vectors $\mathbf{e}_1, \mathbf{e}_2..., \mathbf{e}_B \in S^d(1)$ and scalars $r_1, r_2, ..., r_B \in [-1, 1]$ independently

2. Calculate
$$\mathbf{g}_k = \frac{1}{B} \sum_{i=1}^{B} \mathbf{g}(x_k, \mathbf{e}_i)$$
 via (5)

3.
$$y_k \leftarrow \alpha_k z_k + (1 - \alpha_k) x_k$$

4. $x_{k+1} \leftarrow y_k - \eta \mathbf{g}_k$

5.
$$z_{k+1} \leftarrow \beta_k z_k + (1 - \beta_k) y_k - \gamma_k \eta \mathbf{g}_k$$

end for

```
Return: x_N
```

Now, in order to obtain an estimate of the convergence rate of Algorithm 1, we need to evaluate the bias as well as the second moment of the gradient approximation (5). Let's start with the bias of the gradient approximation:

Bias of gradient approximation Using the variational representation of the Euclidean norm, and definition of gradient approximation (5) we can write:

$$\left\| \frac{d}{2h} \mathbb{E} \left[\left(\tilde{f}_1(x + hr\mathbf{e}) - \tilde{f}_2(x - hr\mathbf{e}) \right) K(r) \mathbf{e} \right] - \nabla f(x) \right\|$$

$$\stackrel{\text{(1)}}{=} \left\| \frac{d}{h} \mathbb{E} \left[f(x + hr\mathbf{e}) K(r) \mathbf{e} \right] - \nabla f(x) \right\|$$

$$\stackrel{\text{(2)}}{=} \left\| \mathbb{E} \left[\nabla f(x + hr\mathbf{u}) rK(r) \right] - \nabla f(x) \right\|$$

$$275 = \sup_{z \in S_2^d(1)} \mathbb{E} \left[\| (\nabla_z f(x + hr \mathbf{u}) - \nabla_z f(x)) r K(r) \| \right]$$

$$276 \qquad \overset{\mathfrak{g},\mathfrak{g}}{\leq} \kappa_\beta h^{\beta-1} \frac{L}{(l-1)!} \mathbb{E} \left[\| u \|^{\beta-1} \right]$$

$$278 \qquad \leq \kappa_\beta h^{\beta-1} \frac{L}{(l-1)!} \frac{d}{d+\beta-1}$$

$$281 \qquad \leq \kappa_\beta L h^{\beta-1}, \qquad (6)$$

283

284

285

286

287

288

289

293

31

31

320

321

322

323

324

where $u \in B^d(1)$; ① = the equality is obtained from the fact, namely, distribution of e is symmetric' 2 = the equality is obtained from a version of Stokes' theorem (Zorich & Paniagua, 2016); ③ = Taylor expansion (see Appendix for more detail); () =assumption that $|R(hr\mathbf{u})| \leq$ $\frac{L}{(l-1)!} \|hr\mathbf{u}\|^{\beta-1} = \frac{L}{(l-1)!} |r|^{\beta-1} h^{\beta-1} \|\mathbf{u}\|^{\beta-1}.$

290 Now we find an estimate of the second moment of the gra-291 dient approximation (5). 292

Bounding second moment of gradient approximation By definition gradient approximation (5) and Wirtinger-Poincare inequality we have

$$\begin{aligned}
 & \mathbb{E}\left[\left\|\mathbf{g}(x_{k},\mathbf{e})\right\|^{2}\right] \\
 &= \frac{d^{2}}{4h^{2}}\mathbb{E}\left[\left\|\left(\tilde{f}(x_{k}+hr\mathbf{e})-\tilde{f}(x_{k}-hr\mathbf{e})\right)K(r)\mathbf{e}\right\|^{2}\right] \\
 &\leq \frac{\kappa d^{2}}{2h^{2}}\left(\mathbb{E}\left[\left(f(x_{k}+hr\mathbf{e})-f(x_{k}-hr\mathbf{e})\right)^{2}\right]+2\Delta^{2}\right) \\
 &\leq \frac{\kappa d^{2}}{2h^{2}}\left(\mathbb{E}\left[\left\|\nabla f(x_{k}+hr\mathbf{e})+\nabla f(x_{k}-hr\mathbf{e})\right\|^{2}\right] \\
 &\leq \frac{\kappa d}{2}\mathbb{E}\left[\left\|\nabla f(x_{k}+hr\mathbf{e})+\nabla f(x_{k}-hr\mathbf{e})\right\|^{2}\right] \\
 &+\frac{\kappa d^{2}\Delta^{2}}{h^{2}} \\
 &\leq \kappa d\mathbb{E}\left[\left\|\nabla f(x_{k}+hr\mathbf{e})\pm\|\nabla f(x_{k})\right\|\right\|^{2}\right] \\
 &+\kappa d\mathbb{E}\left[\left\|\nabla f(x_{k}-hr\mathbf{e})\pm\|\nabla f(x_{k})\right\|\right\|^{2}\right] \\
 &+\frac{\kappa d^{2}\Delta^{2}}{h^{2}} \\
 &\leq 4d\kappa \left\|\nabla f(x_{k})\right\|^{2}+4d\kappa L^{2}h^{2}+\frac{\kappa d^{2}\Delta^{2}}{h^{2}}.$$
(7)

314
$$\rho$$

315 σ^2
316 Now substituting into Theorem 2.6 instead of
317 $\delta \to \kappa_\beta Lh^{\beta-1}$ from (6), $\rho \to 4d\kappa$ from (7) and
318 $\sigma^2 \to 4d\kappa L^2h^2 + \frac{\kappa d^2\Delta^2}{h^2}$ from (7), we obtain convergence
for the neural analysis free method (see Algorithm 1) with

ice for the novel gradient-free method (see Algorithm 1) with $\rho_B = \max\{1, \frac{4d\kappa}{B}\}$:

$$F_N \le \underbrace{\left(1 - \sqrt{\frac{\mu}{\rho_B^2 L}}\right)^N \left[f(x_0) - f^* + \frac{\mu}{2} \|x_0 - x^*\|^2\right]}_{\mathbb{T}}$$

$$\begin{array}{ccc} 325 & & & & \\ 326 & & & \\ 327 & & & \\ 328 & & & \\ 329 & & & \\ \end{array} + \underbrace{\frac{4d\kappa L^2 h^2}{\sqrt{\rho_B^2 \mu L B^2}}}_{\textcircled{2}} + \underbrace{\frac{\kappa d^2 \Delta^2}{h^2 \sqrt{\rho_B^2 \mu L B^2}}}_{\textcircled{3}} + \underbrace{\frac{\kappa_\beta^2 L^2 h^{2(\beta-1)}}{\sqrt{4\mu L}}}_{\textcircled{9}} \end{array}$$

We are now ready to present the main result of this paper.

Theorem 3.1. Let the function f satisfy Assumptions 1.1 and 1.2, and let the Kernel approximation with zero-order oracle (see Definition 1.3) satisfy Assumptions 1.4 and 2.3–2.4, then the novel Zero-Order Accelerated Batched Stochastic Gradient Descent (see Algorithm 1) converges to the desired accuracy ε at the following parameters

- 1. Case: B = 1: with smoothing parameter $h \leq 1$ $\varepsilon^{1/2}\mu^{1/4}$, after $N = \mathcal{O}\left(\sqrt{\frac{d^2L}{\mu}}\log\frac{1}{\varepsilon}\right)$ successive iterations, $T = N \cdot B = \mathcal{O}\left(\sqrt{\frac{d^2L}{\mu}}\log\frac{1}{\varepsilon}\right)$ oracle calls and at $\Delta \lesssim \frac{\varepsilon \mu^{1/2}}{\sqrt{d}}$ maximum noise level.
- 2. Case: $1 < B < 4d\kappa$: with parameter $h \leq \varepsilon^{1/2} \mu^{1/4}$, after $N = \mathcal{O}\left(\sqrt{\frac{d^2L}{B^2\mu}}\log\frac{1}{\varepsilon}\right)$ successive iterations, $T = N \cdot B = \mathcal{O}\left(\sqrt{\frac{d^2L}{\mu}}\log \frac{1}{\varepsilon}\right)$ oracle calls and at $\Delta \lesssim \frac{\varepsilon \mu^{1/2}}{\sqrt{4}}$ maximum noise level.
- 3. Case: $B = 4d\kappa$: with smoothing parameter $h \lesssim$ $\varepsilon^{1/2}\mu^{1/4}$, after $N = \mathcal{O}\left(\sqrt{\frac{L}{\mu}}\log{\frac{1}{\varepsilon}}\right)$ successive iterations, $T = N \cdot B = \mathcal{O}\left(\sqrt{\frac{d^2L}{\mu}}\log \frac{1}{\varepsilon}\right)$ oracle calls and at $\Delta \lesssim \frac{\varepsilon \mu^{1/2}}{\sqrt{d}}$ maximum noise level.
- 4. Case: $B > 4d\kappa$: with parameter $h \lesssim (\varepsilon \sqrt{\mu})^{\frac{1}{2(\beta-1)}}$, after $N = \mathcal{O}\left(\sqrt{\frac{L}{\mu}}\log{\frac{1}{\varepsilon}}\right)$ successive iterations, T = $N \cdot B = \max\{\tilde{\mathcal{O}}\left(\sqrt{\frac{d^2L}{\mu}}\right), \tilde{\mathcal{O}}\left(\frac{d^2\Delta^2}{(\varepsilon\mu)^{\frac{\beta}{\beta-1}}}\right)\}$ oracle calls and at $\Delta \lesssim \frac{(\varepsilon \sqrt{\mu})^{\frac{\beta}{2(\beta-1)}}}{d} B^{1/2}$ maximum noise level.

As can be seen from Theorem 3.1, Algorithm 1 indeed improves the iteration complexity compared to previous works (see Table 1), reaching the optimal estimate in a class of algorithms based on first-order algorithms at batch size $B = 4d\kappa$. However, if we consider the case $B \in [1, 4d\kappa]$, then when the batch size increases from 1, the algorithm improves the convergence rate (without changing the oracle complexity), but achieves the same error floor. This is not very good, because the asymptote does not depend on either the batch size or the smoothness order of the function. However, if we take the batch size larger than $B > 4d\kappa$, we will significantly improve the maximal noise level by worsening the oracle complexity. That is, in the overbatching condition, the error floor depends on both the batch size and the smoothness order, which can play a critical role in real life. For a detailed proof, see Appendix D.

Remark 3.2 (Convex case.). It is not difficult to show that the results of Theorem 3.1 generalize to the convex case (see Assumption 1.2 with $\mu = 0$), preserving the same dependence on *B*, namely in the case $B \in [1; 4d\kappa]$ and $h \lesssim \varepsilon^{3/4}$ we have the following convergence estimates for Algorithm 1:

$$N = \mathcal{O}\left(\sqrt{\frac{d^2 L R^2}{B^2 \varepsilon}}\right); \qquad T = \mathcal{O}\left(\sqrt{\frac{d^2 L R^2}{\varepsilon}}\right)$$

338 and

334

335

337

339

340 341 342

343

345

346 347

353

367

368 369

$$\Delta \lesssim \frac{\varepsilon^{3/2}}{\sqrt{d}}$$

We can also observe that the optimal estimate of iteration complexity in the convex setup is achieved when $B = 4d\kappa$. Moreover, the maximum noise level behaves in a similar way:

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

and

$$\Delta \lesssim \frac{\varepsilon^{\frac{3\beta+1}{4(\beta-1)}}}{d} B^{1/2}.$$

It can be seen that if we take $\mu \sim \varepsilon$, the oracle complexity is the same in the worst case, and the maximum noise level is inferior depending on the order of smoothness compared to the strongly convex set (which is surprising).

362 *Remark* 3.3 (Deterministic adversarial noise). It should be 363 noted that when considering deterministic adversarial noise 364 $(|\tilde{\xi}(x)| \leq \Delta)$ in a zero-order oracle instead of stochastic 365 (see ξ_p with $p = \{1, 2\}$ in Definition 1.3), Theorem 3.1 will 366 preserve the results except for the maximum noise level:

$$\Delta \lesssim \frac{(\varepsilon \sqrt{\mu})^{\frac{\beta}{2(\beta-1)}}}{d} B^{1/2} \to \Delta \lesssim \frac{(\varepsilon \sqrt{\mu})^{\frac{\beta}{2(\beta-1)}}}{d}.$$

This can be explained by the fact that deterministic noise is more adversarial because it accumulates not only in the second moment of the gradient approximation, but also in the bias! The results in the convex case will change similarly.

Remark 3.4 (High probability deviations bound). Given that Algorithm 1 in strongly convex setting demonstrates a linear convergence rate and employs a randomization (see e.g. $\mathbf{e} \in S^d(1)$), we can derive exact estimates of high deviation probabilities using Markov's inequality (Anikin et al., 2017):

$$\mathcal{P}\left(f(x_{N_{(\varepsilon\theta)}}) - f^* \ge \varepsilon\right) \le \theta \frac{\mathbb{E}\left[f(x_{N_{(\varepsilon\theta)}})\right] - f^*}{\varepsilon\theta}$$

Remark 3.5 (Non-convex setup (PL)). It should be noted that our algorithm will have global convergence for a subclass of non-convex functions that satisfy the Polyak—Lojasiewicz (PL) condition (see, Karimi et al., 2016). It is not hard to see that the results will have a similar dependence on the batch size:

$$N = \tilde{\mathcal{O}}\left(\frac{d}{B}\tilde{\mu}^{-1}\right); \qquad T = \tilde{\mathcal{O}}\left(d\tilde{\mu}^{-1}\right)$$

and

$$\Delta \lesssim \frac{\varepsilon \tilde{\mu}}{\sqrt{d}},$$

where $\tilde{\mu}$ from PL Assumption (see, Karimi et al., 2016). We can also observe that the optimal estimate of iteration complexity in the convex setup is achieved when $B = 4d\kappa$. Also, the maximum noise level behaves similarly:

$$\begin{split} N &= \tilde{\mathcal{O}}\left(\tilde{\mu}^{-1}\right);\\ T &= \max\left[\tilde{\mathcal{O}}\left(d\tilde{\mu}^{-1}\right), \tilde{\mathcal{O}}\left(\frac{d^{2}\Delta^{2}}{\varepsilon^{\frac{\beta}{\beta-1}}\tilde{\mu}^{\frac{2\beta-1}{\beta-1}}}\right)\right] \end{split}$$

and

$$\Delta \lesssim \frac{(\varepsilon \tilde{\mu})^{\frac{\beta}{2(\beta-1)}}}{d} B^{1/2}$$

Similarly to the cases discussed above, when considering deterministic adversarial noise, the dependence on the batch size will disappear in the estimation of the maximum noise level. The transition to High probability deviations bounds is also valid. And if we compare with the estimates of Theorem 3.1, provided $\mu \sim \varepsilon$ from the strong convexity condition, and $\tilde{\mu} \sim \varepsilon$ from the PL condition, then the iteration complexity is the same, but the oracle complexity in the PL case is inferior to the strongly convex case. This can be explained by the fact that the PL condition covers a subclass of non-convex functions.

4. Numerical Experiments

In this section, we show the importance of considering the maximum noise level Δ as a third optimality criterion along with the standard two. We consider a problem of interest in machine learning, namely the logistic regression problem:

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

Here we can understand $\log(1 + \exp(-y_i \cdot (Ax)_i)) = f_i(x)$ as the loss at the *i*-th data point, $x \in \mathbb{R}^d$ as a vector of parameters (or weights), $y \in \{-1, 1\}^M$ as a vector of labels, and $A \in \mathbb{R}^{M \times d}$ as a matrix of instances. For our experiments we use data from the LIBSVM library (Chang & Lin, 2011),

 $\leq \theta$.

namely the a9a data. In the gradient approximation (5), we choose as the kernel function K(r) the Legendre polynomials, for which it is shown in (Bach & Perchet, 2016) that the parameters κ and κ_{β} depend only on the smoothness order β . We have the following values for different β :

390

395

396

397

398

399

400

401

402

403

404

405

406

407

408

409 410

411

412

413

414

415

416

417

418

419

420

421

422

$$K(r) = \frac{15r}{4}(5 - 7r^2) \qquad \text{for } \beta = 3, 4;$$

$$K(r) = \frac{195r}{16}(99r^4 - 126r^2 + 35) \quad \text{ for } \beta = 5, 6.$$

To show the effectiveness of our Algorithm 1 (ZO-ABSGD) we compare with SOTA accelerated gradient-free algorithms, namely ZO-VARAG from (Chen et al., 2020), ARDFDS from (Gorbunov et al., 2022). We also compare our Algorithm 1 with RDFDS from (Gorbunov et al., 2022) to demonstrate the superiority of the accelerated algorithm over the unaccelerated ones, which are all previous works (see Table 1).



Figure 3. Comparison of SOTA gradient-free algorithms convergence. Here we optimize f(x) with the parameters: d = 123(problem dimensionality), B = 1000 (batch size), $\Delta = 10^{-5}$ (noise level), $\eta = 10^{-4}$ (step size), $h = 10^{-4}$ (smoothing parameter). In all experiments, the hyperparameters of the algorithms are tuned.

423 Figure 3 shows both standard results, such as the superiority 424 of accelerated methods over unaccelerated methods, and the 425 outperformance, the robustness of our algorithm. It is not 426 hard to see that the ZO-VARAG algorithm outperforms the 427 convergence rate on the first iterations, but converges to an 428 error floor thereafter. This effect (convergence to the asymp-429 tote) can be explained by the fact that in (Chen et al., 2020) 430 an accelerated ZO-VARAG algorithm was proposed, which 431 is not robust to adversarial noise. Regarding the RDFDS 432 and ARDFDS algorithms, as the Figure shows they are also 433 robust to adversarial stochastic noise like our algorithm. 434 The robust convergence of the algorithms from (Gorbunov 435 et al., 2022) can be explained by the fact that in (Gorbunov 436 et al., 2022) algorithms were proposed that are robust to 437 deterministic adversarial noise (DAN). As we know DAN is 438 more antagonistic than stochastic adversarial noise because 439

it accumulates not only in the variance but also in the bias of the gradient approximation. Despite this, ZO-ABSGD has better convergence compared to ARDFDS because the proposed 1 takes advantage of increased smoothness ($\beta = 3$), unlike its counterpart. Thus, this Figure 3 demonstrates not only the advantage of our algorithm, but also the importance in the design and analysis of algorithms robust to adversarial noise!

5. Conclusion

In this paper, we proposed a novel accelerated gradient-free algorithm to solve the black-box optimization problem under the assumption of increased smoothness and strong convexity of the objective function. By choosing a first-order accelerated algorithm and generalizing it to the Batched algorithm with a biased gradient oracle, we were able to improve the iteration complexity, reaching optimal estimates. Moreover, we have shown the importance of considering the maximum noise level as a third optimality criterion in a numerical experiment of interest in machine learning. And finally, we believe that this work offers a new perspective on black-box optimization and opens avenues for future research.

Impact Statement

This paper presents work whose goal is to advance the field of Machine Learning. There are many potential societal consequences of our work, none which we feel must be specifically highlighted here.

References

- Akhavan, A., Pontil, M., and Tsybakov, A. Exploiting higher order smoothness in derivative-free optimization and continuous bandits. *Advances in Neural Information Processing Systems*, 33:9017–9027, 2020.
- Akhavan, A., Chzhen, E., Pontil, M., and Tsybakov, A. A gradient estimator via 11-randomization for online zeroorder optimization with two point feedback. *Advances in Neural Information Processing Systems*, 35:7685–7696, 2022.
- Akhavan, A., Chzhen, E., Pontil, M., and Tsybakov, A. B. Gradient-free optimization of highly smooth functions: improved analysis and a new algorithm. *arXiv preprint arXiv:2306.02159*, 2023.
- Alashqar, B., Gasnikov, A., Dvinskikh, D., and Lobanov, A. Gradient-free federated learning methods with 1 1 and 1 2-randomization for non-smooth convex stochastic optimization problems. *Computational Mathematics and Mathematical Physics*, 63(9):1600–1653, 2023.

- Anikin, A. S., Gasnikov, A. V., Dvurechensky, P., Tyurin, A.,
 and Chernov, A. V. Dual approaches to the minimization
 of strongly convex functionals with a simple structure
 under affine constraints. *Computational Mathematics*and Mathematical Physics, 57:1262–1276, 2017.
- Bach, F. and Perchet, V. Highly-smooth zero-th order online
 optimization. In *Conference on Learning Theory*, pp.
 257–283. PMLR, 2016.

449

456

460

- Bogolubsky, L., Dvurechenskii, P., Gasnikov, A., Gusev,
 G., Nesterov, Y., Raigorodskii, A. M., Tikhonov, A.,
 and Zhukovskii, M. Learning supervised pagerank with
 gradient-based and gradient-free optimization methods. *Advances in neural information processing systems*, 29,
 2016.
- Chang, C.-C. and Lin, C.-J. Libsvm: a library for support
 vector machines. ACM transactions on intelligent systems
 and technology (TIST), 2(3):1–27, 2011.
- Chen, Y., Orvieto, A., and Lucchi, A. An acceler-461 ated DFO algorithm for finite-sum convex functions. 462 In III, H. D. and Singh, A. (eds.), Proceedings of 463 the 37th International Conference on Machine Learn-464 ing, volume 119 of Proceedings of Machine Learn-465 ing Research, pp. 1681-1690. PMLR, 13-18 Jul 2020. 466 URL https://proceedings.mlr.press/v119/ 467 chen20r.html. 468
- 470 Conn, A. R., Scheinberg, K., and Vicente, L. N. Introduction
 471 to derivative-free optimization. SIAM, 2009.
- Dvinskikh, D., Tominin, V., Tominin, I., and 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.
- Gasnikov, A., Novitskii, A., Novitskii, V., Abdukhakimov,
 F., Kamzolov, D., Beznosikov, A., Takac, M., Dvurechensky, P., and 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.
- Gasnikov, A., Dvinskikh, D., Dvurechensky, P., Gorbunov,
 E., Beznosikov, A., and Lobanov, A. Randomized
 gradient-free methods in convex optimization. In *Encyclopedia of Optimization*, pp. 1–15. Springer, 2023.
- 490
 491
 492
 492
 493
 494
 494
 494
 Gasnikov, A., Lobanov, A., and Bashirov, N. The "overbatching" effect? yes, or how to improve error floor in black-box optimization problems. *arXiv preprint arXiv*, 2024a.

- Gasnikov, A., Lobanov, A., and Stonyakin, F. Highly smooth zeroth-order methods for solving optimization problems under the pl condition. *Computational Mathematics and Mathematical Physics*, 64(4):739–770, 2024b.
- Gorbunov, E., Dvurechensky, P., and Gasnikov, A. An accelerated method for derivative-free smooth stochastic convex optimization. *arXiv preprint arXiv:1802.09022*, 2018.
- Gorbunov, E., Dvurechensky, P., and Gasnikov, A. An accelerated method for derivative-free smooth stochastic convex optimization. *SIAM Journal on Optimization*, 32 (2):1210–1238, 2022. doi: 10.1137/19M1259225. URL https://doi.org/10.1137/19M1259225.
- Karimi, H., Nutini, J., and Schmidt, M. Linear convergence of gradient and proximal-gradient methods under the polyak-łojasiewicz condition. In Machine Learning and Knowledge Discovery in Databases: European Conference, ECML PKDD 2016, Riva del Garda, Italy, September 19-23, 2016, Proceedings, Part I 16, pp. 795– 811. Springer, 2016.
- Kimiaei, M. and Neumaier, A. Efficient unconstrained black box optimization. *Mathematical Programming Computation*, 14(2):365–414, 2022.
- Kornilov, N., Shamir, O., Lobanov, A., Dvinskikh, D., Gasnikov, A., Shibaev, I., Gorbunov, E., and Horváth, S. Accelerated zeroth-order method for non-smooth stochastic convex optimization problem with infinite variance. *Advances in Neural Information Processing Systems*, 36, 2024.
- Lobanov, A. Stochastic adversarial noise in the "black box" optimization problem. In *International Conference on Optimization and Applications*, pp. 60–71. Springer, 2023.
- Lobanov, A. and Gasnikov, A. Accelerated zero-order sgd method for solving the black box optimization problem under "overparametrization" condition. In *International Conference on Optimization and Applications*, pp. 72–83. Springer, 2023.
- Lobanov, A., Anikin, A., Gasnikov, A., Gornov, A., and Chukanov, S. Zero-order stochastic conditional gradient sliding method for non-smooth convex optimization. In *International Conference on Mathematical Optimization Theory and Operations Research*, pp. 92–106. Springer, 2023a.
- Lobanov, A., Veprikov, A., Konin, G., Beznosikov, A., Gasnikov, A., and Kovalev, D. Non-smooth setting of stochastic decentralized convex optimization problem over timevarying graphs. *Computational Management Science*, 20 (1):48, 2023b.

495 496 497 498 400	Lobanov, A., Gasnikov, A., and Krasnov, A. Acceler- ation exists! optimization problems when oracle can only compare objective function values. <i>arXiv preprint</i> <i>arXiv:2402.09014</i> , 2024.
500 501 502	Nesterov, Y. Introductory lectures on convex optimization: A basic course, volume 87. Springer Science & Business Media, 2003.
503 504 505 506	Novitskii, V. and Gasnikov, A. Improved exploiting higher order smoothness in derivative-free optimization and con- tinuous bandit. <i>arXiv preprint arXiv:2101.03821</i> , 2021.
507 508 509	Polyak, B. T. and Tsybakov, A. B. Optimal order of accuracy of search algorithms in stochastic optimization. <i>Problemy</i> <i>Peredachi Informatsii</i> , 26(2):45–53, 1990.
510 511 512 513 514	Rios, L. M. and Sahinidis, N. V. Derivative-free optimiza- tion: a review of algorithms and comparison of software implementations. <i>Journal of Global Optimization</i> , 56(3): 1247–1293, 2013.
515 516 517	Stich, S. U. Unified optimal analysis of the (stochastic) gradient method. <i>arXiv preprint arXiv:1907.04232</i> , 2019.
518 519 520 521 522	Vaswani, S., Bach, F., and Schmidt, M. Fast and faster convergence of sgd for over-parameterized models and an accelerated perceptron. In <i>The 22nd international</i> <i>conference on artificial intelligence and statistics</i> , pp. 1195–1204. PMLR, 2019.
525 524 525 526 527 528 529	Zorich, V. A. and Paniagua, O. <i>Mathematical analysis II</i> , volume 220. Springer, 2016.
530 531 532 533 534 525	
535 536 537 538 539	
540 541 542 543 544	
545 546 547 548 549	

APPENDIX Maximum Noise Level as Third Optimality Criterion in Black-Box Optimization Problem

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

$$\|a_1 + \dots + a_n\|^2 \le n \|a_1\|^2 + \dots + n \|a_n\|^2.$$
(8)

A.2. Fenchel-Young inequality

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

$$\langle a, b \rangle \le \frac{\|a\|^2}{2\lambda} + \frac{\lambda \|b\|^2}{2}.$$
(9)

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.

$$\|\nabla f(x) - \nabla f(y)\| \le L \|x - y\|, \quad \forall x, y \in \mathbb{R}^d.$$
(10)

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

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

and if f is additionally convex, then

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

A.4. Wirtinger-Poincare inequality

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

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

A.5. Taylor expansion

Using the Taylor expansion we have

$$\nabla_z f(x + hr\mathbf{u}) = \nabla_z f(x) + \sum_{1 \le |n| \le l-1} \frac{(rh)^{|n|}}{n!} D^{(n)} \nabla_z f(x) \mathbf{u}^n + R(hr\mathbf{u}), \tag{12}$$

where by assumption

$$|R(hr\mathbf{u})| \le \frac{L}{(l-1)!} \|hr\mathbf{u}\|^{\beta-1} = \frac{L}{(l-1)!} |r|^{\beta-1} h^{\beta-1} \|\mathbf{u}\|^{\beta-1}.$$
(13)

A.6. Kernel property

If e is uniformly distributed on $S_2^d(1)$ we have $\mathbb{E}[\mathbf{e}\mathbf{e}^T] = (1/d)I_{d\times d}$, 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 $2 \le |n| \le l$ we have

$$\mathbb{E}\left[\frac{d}{h}\left(\langle \nabla f(x), hr\mathbf{e} \rangle + \sum_{2 \le |n| \le l} \frac{(rh)^{|n|}}{n!} D^{(n)} f(x) \mathbf{e}^n\right) K(r) \mathbf{e}\right] = \nabla f(x).$$
(14)

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, (Bach & Perchet, 2016)) that κ_{β} and κ do not depend on d, they depend only on β , such that for $\beta > 1$:

$$\kappa_{\beta} \le 2\sqrt{2}(\beta - 1),\tag{15}$$

$$\kappa \le 3\beta^3. \tag{16}$$

B. Missing proof of Theorem 2.6

In this Section we demonstrate a missing proof of Theorem 2.6, namely a generalization of Lemma 2.5 to the case with a biased gradient oracle (see Definition 2.2). Therefore, our reasoning is based on the proof of Lemma 2.5 (Vaswani et al., 2019).

Before proceeding directly to the proof, we recall the update rules of First-order Accelerated SGD from (Vaswani et al., 2019):

$$y_k = \alpha_k z_k + (1 - \alpha_k) x_k; \tag{17}$$

$$x_{k+1} = y_k - \eta \mathbf{g}_k; \tag{18}$$

$$z_{k+1} = \beta_k z_k + (1 - \beta_k) y_k - \gamma_k \eta \mathbf{g}_k, \tag{19}$$

where we choose the parameters γ_k , α_k , β_k , a_k , b_k such that the following equations are satisfied:

$$\gamma_k = \frac{1}{2\rho} \cdot \left[1 + \frac{\beta_k (1 - \alpha_k)}{\alpha_k} \right]; \tag{20}$$

$$\alpha_k = \frac{\gamma_k \beta_k b_{k+1}^2 \eta}{\gamma_k \beta_k b_{k+1}^2 \eta + 2a_k^2};\tag{21}$$

$$\beta_k > 1 - \gamma_k \mu n; \tag{22}$$

$$\sum_{k=1}^{k} \sum_{k=1}^{k} \sqrt{n\rho} b_{k+1}; \tag{23}$$

$$a_{k+1} = \gamma_k \sqrt{\eta} \rho b_{k+1}; \tag{23}$$
$$b_{k+1} \le \frac{b_k}{\sqrt{2}}. \tag{24}$$

$$a_{k+1} \le \frac{b_k}{\sqrt{\beta_k}}.$$
(24)

Now, we're ready to move on to the proof itself. Let $r_{k+1} = ||z_{k+1} - x^*||$ and $\mathbf{g}_k = \mathbf{g}(y_k, \omega_k)$ from Definition 2.2, then using equation (19):

$$r_{k+1}^{2} = \|\beta_{k}z_{k} + (1-\beta_{k})y_{k} - x^{*} - \gamma_{k}\eta\mathbf{g}_{k}\|^{2}$$

$$r_{k+1}^{2} = \|\beta_{k}z_{k} + (1-\beta_{k})y_{k} - x^{*}\|^{2} + \gamma_{k}^{2}\eta^{2}\|\mathbf{g}_{k}\|^{2} + 2\gamma_{k}\eta\langle x^{*} - \beta_{k}z_{k} - (1-\beta_{k})y_{k}, \mathbf{g}_{k}\rangle$$

....

Taking expecation wrt to ξ_k ,

$$\mathbb{E}[r_{k+1}^2] = \mathbb{E}[\|\beta_k z_k + (1 - \beta_k)y_k - x^*\|^2] + \gamma_k^2 \eta^2 \mathbb{E} \|\mathbf{g}_k\|^2 \\ + 2\gamma_k \eta \mathbb{E} \left[\langle x^* - \beta_k z_k - (1 - \beta_k)y_k, \mathbf{g}_k \rangle \right]$$

660	$\stackrel{(2.4)}{\leq} \ \beta_k z_k + (1 - \beta_k) y_k - x^*\ ^2 + \gamma_k^2 \eta^2 \rho \ \nabla f(y_k)\ ^2$	
661	+ $2\gamma_{L}n\langle x^{*}-\beta_{L}z_{L}-(1-\beta_{L})y_{L}, \mathbb{E}[\mathbf{g}_{L}]\rangle + \gamma_{L}^{2}n^{2}\sigma^{2}$	
002	$ 2 / \kappa / (\infty) \rangle = 2 / \kappa / (\infty) \rangle 2 / (\kappa / \kappa / $	
003	$= \ \beta_k(z_k - x^*) + (1 - \beta_k)(y_k - x^*)\ ^2 + \gamma_k^2 \eta^2 \rho \ \nabla f(y_k)\ ^2$	
665	$+ 2\gamma_k\eta\langle x^* - \beta_k z_k - (1 - \beta_k)y_k, \mathbb{E}\left[\mathbf{g}_k ight] angle + \gamma_k^2\eta^2\sigma^2$	
003	$< \beta \parallel_{\infty} = m^* \parallel^2 + (1 - \beta) \parallel_{\alpha} = m^* \parallel^2 + n^2 m^2 \rho \parallel_{\nabla} f(\alpha) \parallel^2$	
000	$\leq \rho_k \ z_k - x\ + (1 - \rho_k) \ y_k - x\ + \gamma_k \eta \rho \ \nabla f(y_k)\ $	
667	$+2\gamma_k\eta\langle x^*-eta_k z_k-(1-eta_k)y_k,\mathbb{E}\left[\mathbf{g}_k ight] angle+\gamma_k^2\eta^2\sigma^2$ (By	convexity of $\ \cdot\ ^2$
669	$= \beta_k r_k^2 + (1 - \beta_k) \ y_k - x^*\ ^2 + \gamma_k^2 \eta^2 \rho \ \nabla f(y_k)\ ^2$	
670	$+ 2\gamma_{h}n\langle x^{*} - \beta_{h}z_{h} - (1 - \beta_{h})y_{h}, \mathbb{E}[\mathbf{g}_{h}]\rangle + \gamma_{h}^{2}n^{2}\sigma^{2}$	
671	$= -\frac{1}{2} + \frac{1}{2} + \frac$	
672	$= \beta_k r_k^2 + (1 - \beta_k) \ y_k - x^*\ ^2 + \gamma_k^2 \eta^2 \rho \ \nabla f(y_k)\ ^2$	
673	$+ 2\gamma_k\eta \left< \beta_k(y_k - z_k) + x^* - y_k, \mathbb{E}\left[\mathbf{g}_k\right] \right> + \gamma_k^2 \eta^2 \sigma^2$	
674	$(17) 0 2 (1 0) ^2 ^2 \nabla f(\cdot) ^2$	
675	$= \beta_k r_k^{-} + (1 - \beta_k) \ y_k - x^{-}\ + \gamma_k^{-} \eta^{-} \rho \ \nabla f(y_k)\ $	
676	$+2\gamma_{k}n\left\langle \frac{\beta_{k}(1-\alpha_{k})}{\alpha_{k}}(x_{k}-u_{k})+x^{*}-u_{k}\mathbb{E}\left[\mathbf{g}_{k}\right]\right\rangle +\gamma_{k}^{2}n^{2}\sigma^{2}$	
677	$ \begin{array}{c} 1 & 2 & k & k \\ 1 & 2 $	
678	$= \beta_k r_k^2 + (1 - \beta_k) \ y_k - x^*\ ^2 + \gamma_k^2 \eta^2 \rho \ \nabla f(y_k)\ ^2$	
679	$\begin{bmatrix} \rho & (1 & c \end{pmatrix} \end{bmatrix}$	
680	$+2\gamma_k\eta\left \frac{\rho_k(1-lpha_k)}{2}\langle \mathbb{E}\left[\mathbf{g}_k\right],(x_k-y_k) ight angle+\langle \mathbb{E}\left[\mathbf{g}_k\right],x^*-y_k ight angle$	
681	$\begin{bmatrix} \alpha_k \\ \vdots \end{bmatrix}$	
682	$+ \gamma_k^2 \eta^2 \sigma^2$	
683	$\leq \beta_{1} r^{2} + (1 - \beta_{1}) \ u_{1} - r^{*}\ ^{2} + \gamma^{2} r^{2} a \ \nabla f(u_{1})\ ^{2}$	
684	$ \geq \rho_k r_k + (1 \rho_k) \ g_k x \ + r_k r_k p \ \nabla f(g_k)\ $	
685	$+2\gamma_{k}n\left \frac{\beta_{k}(1-\alpha_{k})}{\beta_{k}(1-\alpha_{k})}\left(f(x_{k})-f(y_{k})\right)+\left\langle\mathbb{E}\left[\mathbf{g}_{k}\right],x^{*}-y_{k}\right\rangle\right +\gamma_{k}^{2}n^{2}\sigma^{2}$	
686	$\left[\begin{array}{ccc} \alpha_k \end{array}\right] = \left[\begin{array}{ccc} \alpha_k \end{array}\right] = \left[\begin{array}{ccc} \alpha_k \end{array}\right]$	
687	$\left[\beta_k(1-\alpha_k)\right]_{\text{TE}}\left[-\sum f(\alpha_k) - \sum f(\alpha_k)\right]$	(D -1,
688	$+ 2\gamma_k \eta \left[\frac{\alpha_k}{\alpha_k} \langle \mathbb{E} \left[\mathbf{g}_k \right] - \nabla J \left(y_k \right), x_k - y_k \rangle \right].$	(By convexity)
689		

By strong convexity,

$$\mathbb{E}[r_{k+1}^{2}] \leq \beta_{k} r_{k}^{2} + (1 - \beta_{k}) \|y_{k} - x^{*}\|^{2} + \gamma_{k}^{2} \eta^{2} \rho \|\nabla f(y_{k})\|^{2}
+ 2\gamma_{k} \eta \left[\frac{\beta_{k}(1 - \alpha_{k})}{\alpha_{k}} \left(f(x_{k}) - f(y_{k}) \right) + f^{*} - f(y_{k}) - \frac{\mu}{2} \|y_{k} - x^{*}\|^{2} \right]
+ 2\gamma_{k} \eta \left[\frac{\beta_{k}(1 - \alpha_{k})}{\alpha_{k}} \left\langle \mathbb{E}\left[\mathbf{g}_{k}\right] - \nabla f(y_{k}), x_{k} - y_{k} \right\rangle + \left\langle \mathbb{E}\left[\mathbf{g}_{k}\right] - \nabla f(y_{k}), x^{*} - y_{k} \right\rangle \right]
+ \gamma_{k}^{2} \eta^{2} \sigma^{2}.$$
(25)

By Lipschitz continuity of the gradient,

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

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

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

Taking expectation wrt ξ_k , we obtain

$$\begin{aligned} & \overset{709}{710} \\ & \overset{710}{711} \\ & \overset{711}{712} \\ & & & \mathbb{E}[f(x_{k+1}) - f(y_k)] \leq -\eta \left\| \nabla f(y_k) \right\|^2 + \frac{L\rho\eta^2}{2} \left\| \nabla f(y_k) \right\|^2 + \frac{L\eta^2\sigma^2}{2} \\ & & & -\eta \left\langle \nabla f(y_k), \mathbb{E}\left[\mathbf{g}_k\right] - \nabla f(y_k) \right\rangle \end{aligned}$$

Maximum Noise Level as Third Optimality Criterion

$$\begin{aligned}
& \begin{bmatrix} f(x_{k+1}) - f(y_k) \end{bmatrix} \stackrel{(9)}{\leq} \left[-\frac{\eta}{2} + \frac{L\rho\eta^2}{2} \right] \|\nabla f(y_k)\|^2 + \frac{L\eta^2\sigma^2}{2} \\
& + \frac{\eta}{2} \|\mathbb{E}[\mathbf{g}_k] - \nabla f(y_k)\|^2 . \\
& + \frac{\eta}{2} \|\mathbb{E}[\mathbf{g}_k] - \nabla f(y_k)\|^2 . \\
& \end{bmatrix} \\
& \begin{bmatrix} f(x_{k+1}) - f(y_k) \end{bmatrix} \leq \left(\frac{-\eta}{4} \right) \|\nabla f(y_k)\|^2 + \frac{L\eta^2\sigma^2}{2} + \frac{\eta}{2} \|\mathbb{E}[\mathbf{g}_k] \| \\
& \end{bmatrix} \\
& \mathbb{E}[f(x_{k+1}) - f(y_k)] \leq \left(\frac{-\eta}{4} \right) \|\nabla f(y_k)\|^2 + \frac{L\eta^2\sigma^2}{2} + \frac{\eta}{2} \|\mathbb{E}[\mathbf{g}_k] \| \\
& \mathbb{E}[\mathbf{g}_k] \| \\
& \mathbb{E}[f(x_{k+1}) - f(y_k)] \leq \left(\frac{-\eta}{4} \right) \|\nabla f(y_k)\|^2 + \frac{L\eta^2\sigma^2}{2} + \frac{\eta}{2} \|\mathbb{E}[\mathbf{g}_k] \| \\
& \mathbb{E}[\mathbf{g}_k] \| \\
& \mathbb{E}[f(x_{k+1}) - f(y_k)] \leq \left(\frac{-\eta}{4} \right) \|\nabla f(y_k)\|^2 + \frac{L\eta^2\sigma^2}{2} + \frac{\eta}{2} \|\mathbb{E}[\mathbf{g}_k] \| \\
& \mathbb{E}[\mathbf{g}_k] \| \\
& \mathbb{E}[f(x_{k+1}) - f(y_k)] \leq \left(\frac{-\eta}{4} \right) \|\nabla f(y_k)\|^2 + \frac{L\eta^2\sigma^2}{2} + \frac{\eta}{2} \|\mathbb{E}[\mathbf{g}_k] \| \\
& \mathbb{E}[\mathbf{g}_k] \| \\
& \mathbb{E}[f(x_{k+1}) - f(y_k)] \leq \left(\frac{-\eta}{4} \right) \|\nabla f(y_k)\|^2 + \frac{L\eta^2\sigma^2}{2} + \frac{\eta}{2} \|\mathbb{E}[\mathbf{g}_k] \| \\
& \mathbb{E}[f(x_k) + \eta^2 + \eta^2$$

725

 If $\eta \leq \frac{1}{2\rho L}$,

$$\mathbb{E}[f(x_{k+1}) - f(y_k)] \le \left(\frac{-\eta}{4}\right) \|\nabla f(y_k)\|^2 + \frac{L\eta^2 \sigma^2}{2} + \frac{\eta}{2} \|\mathbb{E}[\mathbf{g}_k] - \nabla f(y_k)\|^2 \\ \|\nabla f(y_k)\|^2 \le \left(\frac{4}{\eta}\right) \mathbb{E}[f(y_k) - f(x_{k+1})] + 2L\eta\sigma^2 + 2 \|\mathbb{E}[\mathbf{g}_k] - \nabla f(y_k)\|^2.$$
(26)

From equations (25) and (26), we get

$$\mathbb{E}[r_{k+1}^2] \leq \beta_k r_k^2 - 4\gamma_k^2 \eta \rho \mathbb{E} f(x_{k+1}) + 2\gamma_k \eta f^* + \left[2\gamma_k \eta \cdot \frac{\beta_k (1 - \alpha_k)}{\alpha_k}\right] f(x_k) + \left[2\gamma_k \eta \cdot \frac{\beta_k (1 - \alpha_k)}{\alpha_k}\right] \langle \mathbb{E}[\mathbf{g}_k] - \nabla f(u_k) | x_k - u_k \rangle$$

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

Multiplying by b_{k+1}^2 , 770 $b_{k+1}^2 \mathbb{E}[r_{k+1}^2] \le b_{k+1}^2 \beta_k r_k^2 - 4b_{k+1}^2 \gamma_k^2 \eta \rho \mathbb{E} f(x_{k+1}) + 2b_{k+1}^2 \gamma_k \eta f^*$ 772 773 $+ \left[2b_{k+1}^2 \gamma_k \eta \cdot \frac{\beta_k(1-\alpha_k)}{\alpha_k}\right] f(x_k)$ 774 775 + $\left[2b_{k+1}^2\gamma_k\eta\cdot\frac{\beta_k(1-\alpha_k)}{\alpha_k}\right]\langle \mathbb{E}\left[\mathbf{g}_k\right]-\nabla f(y_k), x_k-y_k\rangle$ $+2b_{k+1}^2\gamma_k\eta \langle \mathbb{E}\left[\mathbf{g}_k\right] - \nabla f(u_k), x^* - u_k \rangle$ + $2b_{k+1}^2 \gamma_k^2 \eta^2 \sigma^2 + 2b_{k+1}^2 \gamma_k^2 \eta^2 \rho \|\mathbb{E}[\mathbf{g}_k] - \nabla f(y_k)\|^2$. 779 780 781 Since $b_{k+1}^2 \beta_k \leq b_k^2, b_{k+1}^2 \gamma_k^2 \eta \rho = a_{k+1}^2, \frac{\gamma_k \eta \beta_k (1-\alpha_k)}{\alpha_k} = \frac{2a_k^2}{b_{k+1}^2}$ 782 783 $b_{k+1}^2 \mathbb{E}[r_{k+1}^2] \le b_k^2 r_k^2 - 4a_{k+1}^2 \mathbb{E} f(x_{k+1}) + 2b_{k+1}^2 \gamma_k \eta f^* + 4a_k^2 f(x_k)$ 784 785 $+4a_k^2 \langle \mathbb{E}[\mathbf{g}_k] - \nabla f(y_k), x_k - y_k \rangle$ 786 $+2b_{k+1}^2\gamma_k\eta \langle \mathbb{E}[\mathbf{g}_k] - \nabla f(y_k), x^* - y_k \rangle$ 787 + $\frac{2a_{k+1}^2\sigma^2\eta}{1}$ + $2a_{k+1}^2\eta \|\mathbb{E}[\mathbf{g}_k] - \nabla f(y_k)\|^2$ 788 790 $= b_k^2 r_k^2 - 4a_{k+1}^2 \left[\mathbb{E} f(x_{k+1}) - f^* \right] + 4a_k^2 \left[f(x_k) - f^* \right]$ $+2\left[b_{k+1}^{2}\gamma_{k}\eta-2a_{k+1}^{2}+2a_{k}^{2}\right]f^{*}$ 792 $+4a_{k}^{2}\langle \mathbb{E}\left[\mathbf{g}_{k}\right]-\nabla f(y_{k}), x_{k}-y_{k}\rangle$ $+2b_{k+1}^2\gamma_k\eta \langle \mathbb{E}[\mathbf{g}_k] - \nabla f(y_k), x^* - y_k \rangle$ + $\frac{2a_{k+1}^2\sigma^2\eta}{2}$ + $2a_{k+1}^2\eta \|\mathbb{E}[\mathbf{g}_k] - \nabla f(y_k)\|^2$. 796 797 Since $[b_{k+1}^2 \gamma_k \eta - a_{k+1}^2 + a_k^2] = 0$, 799 800 $b_{k+1}^2 \mathbb{E}[r_{k+1}^2] \le b_k^2 r_k^2 - 4a_{k+1}^2 [\mathbb{E}f(x_{k+1}) - f^*] + 4a_k^2 [f(x_k) - f^*]$ 801 $+4a_{L}^{2}\langle \mathbb{E}[\mathbf{g}_{k}]-\nabla f(u_{k}), x_{k}-x^{*}\rangle$ 802 803 $+4a_{k+1}^2 \langle \mathbb{E}[\mathbf{g}_k] - \nabla f(y_k), x^* - y_k \rangle$ $+\frac{2a_{k+1}^{2}\sigma^{2}\eta}{\sigma}+2a_{k+1}^{2}\eta\left\|\mathbb{E}\left[\mathbf{g}_{k}\right]-\nabla f(y_{k})\right\|^{2}.$ 805 806 807 Denoting $\mathbb{E} f(x_{k+1}) - f^*$ as Φ_{k+1} , we obtain 808 809 $4a_{k+1}^2\Phi_{k+1} - 4a_k^2\Phi_k \stackrel{(2.3)}{\leq} b_k^2r_k^2 - b_{k+1}^2\mathbb{E}[r_{k+1}^2]$ 810 811 $+4a_k^2\delta\tilde{R}-4a_{k+1}^2\delta\tilde{R}$ 812 $+\frac{2a_{k+1}^2\sigma^2\eta}{\sigma^2}+2a_{k+1}^2\eta\delta^2,$ 813 814 815 where $\tilde{R} = \max_k \{ \|x_k - x^*\|, \|y_k - x^*\| \}.$ 816 817 By summing over k we obtain: 818 819 $4\sum_{k=0}^{N-1} \left[a_{k+1}^2 \Phi_{k+1} - a_k^2 \Phi_k\right] \le \sum_{k=0}^{N-1} \left[b_k^2 r_k^2 - b_{k+1}^2 \mathbb{E}[r_{k+1}^2]\right]$ 820 821 822 $+4\sum_{k=1}^{N-1} \left[a_k^2\delta\tilde{R} - a_{k+1}^2\delta\tilde{R}\right]$ 823

+
$$\sum_{k=0}^{N-1} \left[\frac{2a_{k+1}^2 \sigma^2 \eta}{\rho} \right]$$
 + $2 \sum_{k=0}^{N-1} \left[a_{k+1}^2 \eta \delta^2 \right]$.

829 Let's substitute $a_{k+1}^2 = b_{k+1}^2 \gamma_k^2 \eta \rho$:

$$4b_N^2 \gamma_{N-1}^2 \eta \rho \Phi_N \leq 4a_0^2 \Phi_0 + b_0^2 r_0^2 - b_N^2 \mathbb{E}\left[r_N^2\right] + 4a_0^2 \delta \tilde{R} - 4a_N^2 \delta \tilde{R} + \sum_{k=0}^{N-1} \left[\frac{2a_{k+1}^2 \sigma^2 \eta}{\rho}\right] + 2\sum_{k=0}^{N-1} \left[a_{k+1}^2 \eta \delta^2\right].$$

837 Divide the left and right parts by $4\rho\eta$:

$$b_N^2 \gamma_{N-1}^2 \Phi_N \le \frac{a_0^2}{\rho \eta} \Phi_0 + \frac{b_0^2 r_0^2}{4\rho \eta} + \frac{a_0^2 \tilde{R}}{\rho \eta} \delta + \frac{\eta \sigma^2}{2\rho} \sum_{k=0}^{N-1} \left[b_{k+1}^2 \gamma_k^2 \right] + \frac{\eta}{2} \delta^2 \sum_{k=0}^{N-1} \left[b_{k+1}^2 \gamma_k^2 \right].$$

Next, we show that according to (20)-(24) the following relation is correct:

$$\gamma_k^2 - \gamma_k \left[\frac{1}{2\rho} - \mu \eta \gamma_{k-1}^2 \right] = \gamma_{k-1}^2$$

847 Namely, 848

 $\begin{array}{cc} 867\\ 868 \end{array} \quad \text{If } \gamma_k = C \text{, then} \end{array}$

$$\begin{array}{l}
869\\
870\\
870\\
871\\
872\\
873\\
874\\
875\\
876\\
877\\
878\\
879
\end{array}
\qquad \gamma_k = \frac{1}{\sqrt{2\mu\eta\rho}} \\
\beta_k = 1 - \sqrt{\frac{\mu\eta}{2\rho}} \\
\beta_{k+1} = \frac{b_0}{\left(1 - \sqrt{\frac{\mu\eta}{2\rho}}\right)^{(k+1)/2}} \\
b_{k+1} = \frac{b_0}{\left(1 - \sqrt{\frac{\mu\eta}{2\rho}}\right)^{(k+1)/2}} \\
b_{k+1} = \frac{b_0}{\left(1 - \sqrt{\frac{\mu\eta}{2\rho}}\right)^{(k+1)/2}} \\
\frac{b_0}{\left(1 - \sqrt{\frac{\mu\eta}{2\rho}}\right)^{(k+1)/2}} \\
879
\end{array}$$

If $b_0 = \sqrt{2\mu}$, $a_{k+1} = \frac{1}{\left(1 - \sqrt{\frac{\mu\eta}{2\rho}}\right)^{(k+1)/2}}.$ The above equation implies that $a_0 = 1$. Now the above relations allow us to obtain the following inequality: $\frac{2\mu}{\left(1-\sqrt{\frac{\mu\eta}{2\alpha}}\right)^N}\frac{1}{2\mu\eta\rho}\Phi_N \le \frac{1}{\rho\eta}\Phi_0 + \frac{2\mu r_0^2}{4\rho\eta} + \frac{R}{\rho\eta}\delta$ $+ \frac{\sigma^2}{\rho^2} \sum_{k=0}^{N-1} \left| \frac{1}{\left(1 - \sqrt{\frac{\mu\eta}{2\sigma}}\right)^{(k+1)}} \right|$ $+ \frac{1}{2\rho} \delta^2 \sum_{k=0}^{N-1} \left| \frac{1}{\left(1 - \sqrt{\frac{\mu\eta}{2\pi}}\right)^{(k+1)}} \right|;$ $\frac{1}{\left(1-\sqrt{\frac{\mu\eta}{2a}}\right)^N}\Phi_N \le \Phi_0 + \frac{\mu}{2}r_0^2 + \tilde{R}\delta$ $+ \frac{\sigma^2 \eta}{\rho} \sum_{k=0}^{N-1} \left| \frac{1}{\left(1 - \sqrt{\frac{\mu \eta}{2\sigma}}\right)^{(k+1)}} \right|$ $+ \frac{\eta}{2} \delta^2 \sum_{k=0}^{N-1} \left| \frac{1}{\left(1 - \sqrt{\frac{\mu\eta}{2}}\right)^{(k+1)}} \right|;$ $\frac{1}{\left(1-\sqrt{\frac{\mu\eta}{2\alpha}}\right)^N}\Phi_N \le \Phi_0 + \frac{\mu}{2}r_0^2 + \tilde{R}\delta$ $+ \frac{\sigma^2 \sqrt{2\eta}}{\sqrt{\rho \mu}} \cdot \frac{1}{\left(1 - \sqrt{\frac{\mu \eta}{2\rho}}\right)^N}$ $+ \frac{\sqrt{\eta\rho}}{\sqrt{2\mu}} \delta^2 \cdot \frac{1}{\left(1 - \sqrt{\frac{\mu\eta}{2\rho}}\right)^{(k+1)}};$ $\mathbb{E}\left[f(x_N)\right] - f^* \le \left(1 - \sqrt{\frac{\mu\eta}{2\rho}}\right)^N \left[f(x_0) - f^* + \frac{\mu}{2}r_0^2\right]$ $+\left(1-\sqrt{\frac{\mu\eta}{2\mu}}\right)^{N}\tilde{R}\delta+\frac{\sigma^{2}\sqrt{2\eta}}{\sqrt{2\mu}}+\frac{\sqrt{\eta\rho}}{\sqrt{2\mu}}\delta^{2};$ $\mathbb{E}\left[f(x_N)\right] - f^* \le \left(1 - \sqrt{\frac{\mu}{4\rho^2 L}}\right)^N \left[f(x_0) - f^* + \frac{\mu}{2} \|x_0 - x^*\|^2\right]$ $+\left(1-\sqrt{\frac{\mu}{4\rho^2 L}}\right)^N \tilde{R}\delta + \frac{\sigma^2}{\sqrt{\rho^2 \mu L}} + \frac{1}{\sqrt{4\mu L}}\delta^2.$

By adding batching, given that $\tilde{\rho}_B = \max\{1, \frac{\rho}{B}\}$ and $\sigma_B^2 = \frac{\sigma^2}{B}$ we have the convergence rate for accelerated batched SGD

935 with biased gradient oracle and parameter $\eta \lesssim \frac{1}{2\rho_B L}$:

$$\mathbb{E}\left[f(x_N)\right] - f^* \le \left(1 - \sqrt{\frac{\mu}{4\tilde{\rho}_B^2 L}}\right)^N \left[f(x_0) - f^* + \frac{\mu}{2} \|x_0 - x^*\|^2\right] \\ + \left(1 - \sqrt{\frac{\mu}{4\tilde{\rho}_B^2 L}}\right)^N \tilde{R}\delta + \frac{\sigma_B^2}{\sqrt{\tilde{\rho}_B^2 \mu L}} + \frac{1}{\sqrt{4\mu L}}\delta^2.$$

C. Properties of the Kernel Approximation

945 In this Section, we extend the explanations for obtaining the bias and second moment estimates of the gradient approximation.

946 Using the variational representation of the Euclidean norm, and definition of gradient approximation (5) we can write:

$$\begin{aligned} \|\mathbf{b}(x_k)\| &= \|\mathbb{E}\left[\mathbf{g}(x_k, \mathbf{e})\right] - \nabla f(x_k)\| \\ &= \left\|\frac{d}{2h}\mathbb{E}\left[\left(\tilde{f}(x_k + hr\mathbf{e}) - \tilde{f}(x_k - hr\mathbf{e})\right)K(r)\mathbf{e}\right] - \nabla f(x_k)\right\| \\ &\stackrel{\textcircled{\tiny 0}}{=} \left\|\frac{d}{h}\mathbb{E}\left[f(x_k + hr\mathbf{e})K(r)\mathbf{e}\right] - \nabla f(x_k)\right\| \\ &\stackrel{\textcircled{\tiny 0}}{=} \|\mathbb{E}\left[\nabla f(x_k + hr\mathbf{u})rK(r)\right] - \nabla f(x_k)\| \\ &= \sup_{z \in S_2^d(1)} \mathbb{E}\left[\left(\nabla_z f(x_k + hr\mathbf{u}) - \nabla_z f(x_k)\right)rK(r)\right] \\ \stackrel{(12),(13)}{\leq} \kappa_\beta h^{\beta-1}\frac{L}{(l-1)!}\mathbb{E}\left[\|u\|^{\beta-1}\right] \\ &\leq \kappa_\beta h^{\beta-1}\frac{L}{(l-1)!}\frac{d}{d+\beta-1} \\ &\lesssim \kappa_\beta Lh^{\beta-1}, \end{aligned}$$

where $u \in B^d(1)$, (1) = the equality is obtained from the fact, namely, distribution of e is symmetric, (2) = the equality is obtained from a version of Stokes' theorem (Zorich & Paniagua, 2016).

966 By definition gradient approximation (5) and Wirtinger-Poincare inequality (11) we have

$$\mathbb{E}\left[\left\|\mathbf{g}(x_{k},\mathbf{e})\right\|^{2}\right] = \frac{d^{2}}{4h^{2}} \mathbb{E}\left[\left\|\left(\tilde{f}(x_{k}+hr\mathbf{e})-\tilde{f}(x_{k}-hr\mathbf{e})\right)K(r)\mathbf{e}\right\|^{2}\right]\right]$$

$$= \frac{d^{2}}{4h^{2}} \mathbb{E}\left[\left(f(x_{k}+hr\mathbf{e})-f(x_{k}-hr\mathbf{e})+(\xi_{1}-\xi_{2})\right)\right)^{2}K^{2}(r)\right]$$

$$\stackrel{(8)}{\leq} \frac{\kappa d^{2}}{2h^{2}} \left(\mathbb{E}\left[\left(f(x_{k}+hr\mathbf{e})-f(x_{k}-hr\mathbf{e})\right)^{2}\right]+2\Delta^{2}\right)\right]$$

$$\stackrel{(11)}{\leq} \frac{\kappa d^{2}}{2h^{2}} \left(\frac{h^{2}}{d} \mathbb{E}\left[\left\|\nabla f(x_{k}+hr\mathbf{e})+\nabla f(x_{k}-hr\mathbf{e})\right\|^{2}\right]+2\Delta^{2}\right)$$

$$= \frac{\kappa d^{2}}{2h^{2}} \left(\frac{h^{2}}{d} \mathbb{E}\left[\left\|\nabla f(x_{k}+hr\mathbf{e})+\nabla f(x_{k}-hr\mathbf{e})\pm2\nabla f(x_{k})\right\|^{2}\right]+2\Delta^{2}\right)$$

$$\stackrel{(10)}{\leq} \underbrace{4d\kappa}_{\rho} \left\|\nabla f(x_{k})\right\|^{2} + \underbrace{4d\kappa L^{2}h^{2}}_{\sigma^{2}} + \underbrace{\frac{\kappa d^{2}\Delta^{2}}{h^{2}}}_{\sigma^{2}}.$$

983 D. Missing proof of Theorem 3.1

985 Let us consider case B = 1, then we have the following convergence rate:

$$\mathbb{E}\left[f(x_N)\right] - f^* \le \underbrace{\left(1 - \sqrt{\frac{\mu}{(4d\kappa)^2 L}}\right)^N \left[f(x_0) - f^* + \frac{\mu}{2} \|x_0 - x^*\|^2\right]}_{(1)} + \underbrace{\frac{4d\kappa L^2 h^2}{\sqrt{(4d\kappa)^2 \mu L}}}_{(2)}$$

$$+\underbrace{\frac{\kappa d^2\Delta^2}{h^2\sqrt{(4d\kappa)^2\mu L}}}_{\textcircled{\tiny \textcircled{0}}}+\underbrace{\frac{\kappa_\beta^2L^2h^{2(\beta-1)}}{\sqrt{4\mu L}}}_{\textcircled{\tiny \textcircled{0}}}.$$

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

$$\left(1 - \sqrt{\frac{\mu}{(4d\kappa)^2 L}}\right)^N \left[f(x_0) - f^* + \frac{\mu}{2} \left\|x_0 - x^*\right\|^2\right] \le \varepsilon \quad \Rightarrow \quad \left[N = \tilde{\mathcal{O}}\left(\sqrt{\frac{d^2 L}{\mu}}\right).\right]$$

1001 **From terms**
$$(2)$$
, (4) we find the smoothing parameter h :

$$\begin{aligned} &\textcircled{2}: \quad \frac{4d\kappa L^2 h^2}{\sqrt{(4d\kappa)^2 \mu L}} \leq \varepsilon \quad \Rightarrow \quad h^2 \lesssim \varepsilon \sqrt{\mu} \quad \Rightarrow \quad \boxed{h \lesssim (\varepsilon \sqrt{\mu})^{1/2};} \\ &\textcircled{4}: \quad \frac{\kappa_{\beta}^2 L^2 h^{2(\beta-1)}}{\sqrt{4\mu L}} \leq \varepsilon \quad \Rightarrow \quad h^{2(\beta-1)} \lesssim \varepsilon \sqrt{\mu} \quad \Rightarrow \quad h \lesssim (\varepsilon \sqrt{\mu})^{\frac{1}{2(\beta-1)}}. \end{aligned}$$

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

$$\frac{\kappa d^2 \Delta^2}{h^2 \sqrt{(4d\kappa)^2 \mu L}} \leq \varepsilon \quad \Rightarrow \quad \Delta^2 \lesssim \frac{\varepsilon \sqrt{\mu} h^2}{d} \quad \Rightarrow \quad \Delta \lesssim \frac{\varepsilon \sqrt{\mu}}{\sqrt{d}}.$$

The oracle complexity in this case is obtained as follows:

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

Consider now the case $1 < B < 4d\kappa$, then we have the convergence rate:

$$\mathbb{E}[f(x_N)] - f^* \leq \underbrace{\left(1 - \sqrt{\frac{\mu B^2}{(4d\kappa)^2 L}}\right)^N \left[f(x_0) - f^* + \frac{\mu}{2} \|x_0 - x^*\|^2\right]}_{(2)} + \underbrace{\frac{4d\kappa L^2 h^2}{\sqrt{(4d\kappa)^2 \mu L}}}_{(3)} + \underbrace{\frac{\kappa d^2 \Delta^2}{h^2 \sqrt{(4d\kappa)^2 \mu L}}}_{(3)} + \underbrace{\frac{\kappa^2 \mu^2 L^2 h^{2(\beta-1)}}{\sqrt{4\mu L}}}_{(4)}.$$

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

$$\left(1 - \sqrt{\frac{B^2 \mu}{(4d\kappa)^2 L}}\right)^N \left[f(x_0) - f^* + \frac{\mu}{2} \|x_0 - x^*\|^2\right] \le \varepsilon \quad \Rightarrow \quad \boxed{N = \tilde{\mathcal{O}}\left(\sqrt{\frac{d^2 L}{B^2 \mu}}\right).}$$

From terms (2), (4) we find the smoothing parameter h:

1039
1040
1041
1042
(2):
$$\frac{4d\kappa L^2 h^2}{\sqrt{(4d\kappa)^2 \mu L}} \le \varepsilon \Rightarrow h^2 \lesssim \varepsilon \sqrt{\mu} \Rightarrow h \lesssim (\varepsilon \sqrt{\mu})^{1/2};$$

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

$$\frac{\kappa d^2 \Delta^2}{h^2 \sqrt{(4d\kappa)^2 \mu L}} \leq \varepsilon \quad \Rightarrow \quad \Delta^2 \lesssim \frac{\varepsilon \sqrt{\mu} h^2}{d} \quad \Rightarrow \quad \overline{\Delta \lesssim \frac{\varepsilon \sqrt{\mu}}{\sqrt{d}}}.$$

050 The oracle complexity in this case is obtained as follows:

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

1056 Now let us move to the case where $B = 4d\kappa$, then we have convergence rate:

$$\mathbb{E}\left[f(x_{N})\right] - f^{*} \leq \underbrace{\left(1 - \sqrt{\frac{\mu}{L}}\right)^{N} \left[f(x_{0}) - f^{*} + \frac{\mu}{2} \|x_{0} - x^{*}\|^{2}\right]}_{\mathbb{Q}} + \underbrace{\frac{d\Delta^{2}}{\frac{h^{2}\sqrt{\mu L}}{2}}}_{\mathbb{Q}} + \underbrace{\frac{d\Delta^{2}}{\frac{h^{2}\sqrt{\mu L}}{2}} + \underbrace{\frac{\kappa_{\beta}^{2}L^{2}h^{2(\beta-1)}}{\sqrt{4\mu L}}}_{\mathbb{Q}}}_{\mathbb{Q}}.$$

¹⁰⁶⁶ From term (1), we find iteration number N required for Algorithm 1 to achieve ε -accuracy:

$$\left(1 - \sqrt{\frac{\mu}{L}}\right)^{N} \left[f(x_{0}) - f^{*} + \frac{\mu}{2} \|x_{0} - x^{*}\|^{2}\right] \leq \varepsilon \quad \Rightarrow \quad \left[N = \tilde{\mathcal{O}}\left(\sqrt{\frac{L}{\mu}}\right).\right]$$

¹⁰⁷² **From terms (2)**, **(4)** we find the smoothing parameter h:

$$\begin{split} & @: \quad \frac{L^2 h^2}{\sqrt{\mu L}} \leq \varepsilon \quad \Rightarrow \quad h^2 \lesssim \varepsilon \sqrt{\mu} \quad \Rightarrow \quad \boxed{h \lesssim (\varepsilon \sqrt{\mu})^{1/2};} \\ & @: \quad \frac{\kappa_{\beta}^2 L^2 h^{2(\beta-1)}}{\sqrt{4\mu L}} \leq \varepsilon \quad \Rightarrow \quad h^{2(\beta-1)} \lesssim \varepsilon \sqrt{\mu} \quad \Rightarrow \quad h \lesssim (\varepsilon \sqrt{\mu})^{\frac{1}{2(\beta-1)}}. \end{split}$$

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

$$\frac{d\Delta^2}{h^2\sqrt{\mu L}} \le \varepsilon \quad \Rightarrow \quad \Delta^2 \lesssim \frac{\varepsilon\sqrt{\mu}h^2}{d} \quad \Rightarrow \quad \Delta \lesssim \frac{\varepsilon\sqrt{\mu}}{\sqrt{d}}.$$

1085 The oracle complexity in this case is obtained as follows:

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

¹⁰⁹¹ Finally, consider the case when $B > 4d\kappa$, then we have convergence rate:

1092
1093
1094
1095
1096
1097
1098
1099

$$\mathbb{E}\left[f(x_N)\right] - f^* \leq \underbrace{\left(1 - \sqrt{\frac{\mu}{L}}\right)^N \left[f(x_0) - f^* + \frac{\mu}{2} \|x_0 - x^*\|^2\right]}_{\textcircled{0}} + \underbrace{\frac{4d\kappa L^2 h^2}{\sqrt{\mu L B^2}}}_{\textcircled{0}} + \underbrace{\frac{\kappa d^2 \Delta^2}{h^2 \sqrt{\mu L B^2}}}_{\textcircled{0}} + \underbrace{\frac{\kappa d^2 \Delta^2}{\sqrt{\mu L B^2}}}_{\underbrace{0}} + \underbrace{\frac{\kappa$$

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

 $\left(1 - \sqrt{\frac{\mu}{L}}\right)^{N} \left[f(x_{0}) - f^{*} + \frac{\mu}{2} \|x_{0} - x^{*}\|^{2}\right] \leq \varepsilon \quad \Rightarrow \quad \left|N = \tilde{\mathcal{O}}\left(\sqrt{\frac{L}{\mu}}\right).\right.$

From terms (2), (4) we find the smoothing parameter h:

$$\begin{aligned} &\textcircled{2}: \quad \frac{4d\kappa L^2 h^2}{\sqrt{\mu L B^2}} \leq \varepsilon \quad \Rightarrow \quad h^2 \lesssim \frac{\varepsilon \sqrt{\mu}}{d} B \quad \Rightarrow \quad h \lesssim \sqrt{\frac{\varepsilon \sqrt{\mu} B}{d}}; \\ &\textcircled{4}: \quad \frac{\kappa_{\beta}^2 L^2 h^{2(\beta-1)}}{\sqrt{4\mu L}} \leq \varepsilon \quad \Rightarrow \quad h^{2(\beta-1)} \lesssim \varepsilon \sqrt{\mu} \quad \Rightarrow \quad \boxed{h \lesssim (\varepsilon \sqrt{\mu})^{\frac{1}{2(\beta-1)}}.} \end{aligned}$$

From term (3), we find the maximum noise level Δ (via batch size B) at which Algorithm 1 can still achieve ε accuracy:

$$\frac{\kappa d^2 \Delta^2}{h^2 \sqrt{\mu L B^2}} \leq \varepsilon \quad \Rightarrow \quad \Delta^2 \lesssim \frac{(\varepsilon \sqrt{\mu})^{1 + \frac{1}{\beta - 1}B}}{d^2} \quad \Rightarrow \quad \left[\Delta \lesssim \frac{(\varepsilon \sqrt{\mu})^{\frac{\beta}{2(\beta - 1)}} B^{1/2}}{d} \right].$$

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

$$\frac{\kappa d^2 \Delta^2}{h^2 \sqrt{\mu L B^2}} \leq \varepsilon \quad \Rightarrow \quad B \gtrsim \frac{\kappa d^2 \Delta^2}{(\varepsilon \sqrt{\mu})^{1 + \frac{1}{\beta - 1}}} \quad \Rightarrow \quad B = \mathcal{O}\left(\frac{d^2 \Delta^2}{(\varepsilon \sqrt{\mu})^{\frac{\beta}{\beta - 1}}}\right).$$

Then the oracle complexity $T = N \cdot B$ in this case has the following form:

$T = \max\left\{\tilde{\mathcal{O}}\right.$	$\left(\sqrt{\frac{d^2L}{\mu}}\right),\tilde{\mathcal{O}}$	$\left(\frac{d^2\Delta^2}{(\varepsilon\mu)^{\frac{\beta}{\beta-1}}}\right)\right\}.$