

## **Stopping rules for accelerated gradient methods with additive noise in gradient**

Artem Vasin<sup>1</sup>, Alexander Gasnikov<sup>1,2,3</sup>, Vladimir Spokoiny<sup>2,3</sup>

submitted: February 11, 2021

<sup>1</sup> Moscow Institute of Physics and Technology  
Institutskiy Pereulok, 9  
Dolgoprudny, Moscow Region  
141701 Russian Federation  
E-Mail: vasin.aa@phystech.edu

<sup>2</sup> Weierstrass Institute  
Mohrenstr. 39  
10117 Berlin  
Germany  
E-Mail: vladimir.spokoiny@wias-berlin.de  
alexander.gasnikov@wias-berlin.de

<sup>3</sup> Institute for Information Transmission Problems of RAS  
Bolshoy Karetny per. 19, build.1  
127051 Moscow  
Russian Federation  
E-Mail: gasnikov@yandex.ru

No. 2812  
Berlin 2021



---

2020 *Mathematics Subject Classification.* 90C30, 90C25, 68Q25.

*Key words and phrases.* Accelerated methods, inexact gradient, stopping rule, inverse problems.

The research of A. Gasnikov was supported by the Ministry of Science and Higher Education of the Russian Federation (Goszadaniye) 075-00337-20-03, project no. 0714-2020-0005. The work of A. Vasin was supported by Andrei M. Raigorodskii Scholarship in Optimization.

Edited by  
Weierstraß-Institut für Angewandte Analysis und Stochastik (WIAS)  
Leibniz-Institut im Forschungsverbund Berlin e. V.  
Mohrenstraße 39  
10117 Berlin  
Germany

Fax: +49 30 20372-303  
E-Mail: [preprint@wias-berlin.de](mailto:preprint@wias-berlin.de)  
World Wide Web: <http://www.wias-berlin.de/>

# Stopping rules for accelerated gradient methods with additive noise in gradient

Artem Vasin, Alexander Gasnikov, Vladimir Spokoiny

## Abstract

In this article, we investigate an accelerated first-order method, namely, the method of similar triangles, which is optimal in the class of convex (strongly convex) problems with a Lipschitz gradient. The paper considers a model of additive noise in a gradient and a Euclidean prox-structure for not necessarily bounded sets. Convergence estimates are obtained in the case of strong convexity and its absence, and a stopping criterion is proposed for not strongly convex problems.

## 1 Introduction

We consider  $L$ -smooth ( $\mu$ -strongly) convex optimization problem ( $\mu \geq 0$ ):

$$\min_{x \in Q} f(x).$$

This means that  $Q$  is convex set, and for all  $x, y \in Q$ :

$$f(x) + h \langle \nabla f(x), y - x \rangle \leq \frac{\mu}{2} \|y - x\|_2^2 + f(y),$$

$$\| \nabla f(y) - \nabla f(x) \|_2 \leq L \|y - x\|_2.$$

In the analysis of the rates of convergence of different first-order methods these relations are typically rewrite as follows [15, 9, 6, 25, 4, 26, 37, 34, 50, 21, 48, 23, 13]

$$\begin{aligned} f(x) + h \langle \nabla f(x), y - x \rangle &\leq \frac{\mu}{2} \|y - x\|_2^2 + f(y) \\ f(x) + h \langle \nabla f(x), y - x \rangle &\leq \frac{L}{2} \|y - x\|_2^2. \end{aligned} \quad (1)$$

Note, that the last relation is a consequence of the previous ones and in general is not equivalent to them [49, 26].

In many applications, especially for gradient-free methods (when estimating the gradient by finite differences [11, 44, 7]) optimization problems in infinite dimensional spaces (such examples arise when solving inverse problems [31, 27]) instead of an access to  $\nabla f(x)$  we have an access to its inexact approximation  $\tilde{\nabla} f(x)$ .

The two most popular conception of inexactness of gradient in practice are [42]: for all  $x \in Q$

$$\| \tilde{\nabla} f(x) - \nabla f(x) \|_2 \leq \delta, \quad (2)$$

$$\| \tilde{\nabla} f(x) - \nabla f(x) \|_2 \leq \alpha \| \nabla f(x) \|_2, \quad \alpha \in [0, 1). \quad (3)$$

For the first conception (2) several results about the accumulation of error can be found in [42, 12, 10, 1], but all these results are still far from to be optimistic in general. The reason was described in [41]. We can explain this reason by very simple example:

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

where  $0 < \mu = \lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_n = L, L \geq 2\mu$ . The solution of this problem is  $x = 0$ . Assume that inexactness takes place only in the first component. That is instead of  $\partial f(x)/\partial x^1 = \mu x^1$  we have an access to  $\partial f(x)/\partial x^1 = \mu x^1 + \delta$ . For simple gradient dynamic

$$x_k = x_{k-1} - \frac{1}{L} \nabla f(x_{k-1}),$$

we can conclude that for all  $k \in \mathbb{N}$

$$x_k^1 = \frac{\delta}{L} \frac{1 - (1 - \mu/L)^k}{(1 - \mu/L)} = \frac{\delta}{8\mu}. \quad (5)$$

Hence<sup>1</sup>

$$f(x_k) - f(x^*) = \frac{\delta^2}{2\mu}.$$

So we have a problem with (5), since  $\mu$  can be too small ( $\mu \ll \varepsilon$  – degenerate regime, where  $\varepsilon$  – desired accuracy in function value) in denominator of the RHS. We may expect even more serious troubles for accelerated gradient methods, since they are more sensitive to the level of noise [16, 26]. The solution of this problem is well known (see, for example, [41, 42, 35]): to propose a stopping rule for the considered algorithm or to use regularization  $\mu \approx \varepsilon$  [26]. Roughly speaking, for non accelerated algorithms in [41, 42] it was proved that if  $\delta \approx \varepsilon^2$ , then it's possible to reach  $\varepsilon$ -accuracy in function value (with almost the same number of iterations as for no noise case  $\delta = 0$ ) by applying computationally convenient stopping rule.

*In this paper we show that it's sufficient to have  $\delta \approx \varepsilon$  both for primal-dual non accelerated and accelerated gradient type methods [37, 26]. Primal-duality of methods is used to build computationally convenient stopping rule in degenerate regime. We emphasize, that the results  $\delta \approx \varepsilon$  has a simple explanation (see section 2) and one might think that it is well known. But to the best of our knowledge the best results for accelerated methods require  $\delta \approx \varepsilon^{3/2}$ . So we consider our observation (that  $\delta \approx \varepsilon$ ) to be an important part of this paper, although it has rather simple explanation.*

The situation with the second criteria (3) is significantly better. For non accelerated algorithms inexactness in this case lead only to the deceleration of convergence  $(1 - \alpha)^{-1}$ -times [42]. This result holds true with the relaxed strong convexity assumption [26] (Polyak–Lojasiewicz condition). For accelerated case to the best of our knowledge this is an open problem to estimate accumulation of an error [26].

*In this paper we show that if  $\alpha \approx \left(\frac{1}{L}\right)^{3/4}$  in  $\mu$ -strongly convex case and (on  $k$ -th iteration)  $\alpha_k \approx \left(\frac{1}{k}\right)^{3/2}$  in degenerate regime we do not have any deceleration. Numerical experiments demonstrate that in general for  $\alpha$  larger than mentioned above thresholds the convergence may slow down a lot up to divergence for considered accelerated method.*

Note, that close results (with the requirement  $\alpha \approx \left(\frac{1}{L}\right)^{5/4}$ ) in the case  $\mu \approx \varepsilon$  were recently obtained by using another techniques in Stochastic Optimization with decision dependent distribution [18] and

<sup>1</sup>This bound corresponds to the worst-case philosophy concerning the choice of considered example for considered class of methods [36, 37, 9, 26]. We expect more interesting results here by considering average-case complexity [46, 40] (spectrum  $\overline{f\lambda_i g}$  average).

Policy Evaluation in Reinforcement Learning via reduction to stochastic Variational Inequality with Markovian noise [33]. In [33, 18] it was assumed that

$$k\| \nabla f(x) - \nabla f(x) \|_2 \leq B\|x - x^*\|_2, \quad \alpha \in [0, 1]. \quad (6)$$

Since  $x^*$  is a solution, from Fermat's principle  $\nabla f(x^*) = 0$ . Therefore,

$$k\| \nabla f(x) \|_2 = k\| \nabla f(x) - \nabla f(x^*) \|_2 \leq L\|x - x^*\|_2.$$

So if (3) holds true then (6) also holds true with  $B = \alpha L$ .

## 2 Ideas behind the results

Important results in gradient error accumulation for first-order methods were developed in the cycle of works of O. Devolder, F. Glineur and Yu. Nesterov 2011–2014 [14, 16, 17, 15]. In these works authors were motivated by (1). The idea is to “relax” (1), assuming inexactness in gradient. So they introduce inexact gradient  $\tilde{\nabla} f(x)$ , satisfying for all  $x, y \in Q$

$$\begin{aligned} f(x) + h\langle \tilde{\nabla} f(x), y - x \rangle &\leq \frac{\mu}{2} \|y - x\|_2^2 + \delta + f(y) \\ f(x) + h\langle \tilde{\nabla} f(x), y - x \rangle &\geq \frac{L}{2} \|y - x\|_2^2 + \delta. \end{aligned} \quad (7)$$

Such a definition allows to develop precise theory for error accumulation for first-order methods.

Namely, it was proved that for non-accelerated gradient methods

$$f(x_k) - f(x^*) = O\left(\min\left\{\frac{LR^2}{k} + \delta, LR^2 \exp\left(\frac{\mu}{L}k\right) + \delta\right\}\right), \quad (8)$$

and for accelerated ones [16, 20]

$$f(x_k) - f(x^*) = O\left(\min\left\{\frac{LR^2}{k^2} + k\delta, LR^2 \exp\left(\sqrt{\frac{\mu}{L}}\frac{k}{2}\right) + \sqrt{\frac{L}{\mu}}\delta\right\}\right), \quad (9)$$

where  $R = \|x_{start} - x^*\|_2$  – the distance between starting point and the solution  $x^*$ . If  $x^*$  is not unique we take such  $x^*$  that is the closest to  $x_{start}$ . Both of these bounds are unimprovable [16, 17]. See also [15, 22, 32] for “intermediate” situations between accelerated and non-accelerated methods.

Following to [17] we may reduce conception (2) to (7) by putting

$$\delta = \delta_{(7)} = \frac{\delta_{(2)}^2}{2L} + \frac{\delta_{(2)}^2}{\mu}, \quad \frac{\delta_{(2)}^2}{\mu} \quad (10)$$

and changing 2-times constant  $\mu, L$ . The key observations here are

$$h\langle \tilde{\nabla} f(x) - \nabla f(x), y - x \rangle \leq \frac{1}{2L} k\| \nabla f(x) - \nabla f(x) \|_2^2 + \frac{L}{2} \|y - x\|_2^2,$$

$$h\langle \tilde{\nabla} f(x) - \nabla f(x), y - x \rangle \geq \frac{1}{\mu} k\| \nabla f(x) - \nabla f(x) \|_2^2 - \frac{\mu}{4} \|y - x\|_2^2.$$

So, when  $\mu > 0$  for non-accelerated methods this result is almost the same as we've obtained by considering example (4). To reach  $f(x_k) - f(x^*) = \varepsilon$  when  $\mu \propto \varepsilon$  we should put  $\delta_{(2)} \propto \varepsilon$  that is good and rather expected. Unfortunately, for accelerated methods from this approach we will have  $\delta_{(2)} \propto \varepsilon^{3/2}$ . That is far from what we've declared in section 1. To improve this it's worth to propose more detailed conception rather than (7).

In the following works [16, 19, 20, 47, 48] the conception (7) was further developed

$$f(x) + h\Gamma f(x), y \quad xi + \frac{\mu}{2}ky \quad xk_2^2 \quad \delta_1ky \quad xk_2 \quad f(y) \\ f(x) + h\Gamma f(x), y \quad xi + \frac{L}{2}ky \quad xk_2^2 + \delta_2. \quad (11)$$

In this case (8) and (9) take a form for non-accelerated gradient methods

$$f(x_k) - f(x^*) \\ = O \left( \min \left\{ \frac{LR^2}{k} + R\delta_1 + \delta_2, LR^2 \exp \left( \frac{\mu}{L}k \right) + R\delta_1 + \delta_2 \right\} \right), \quad (12)$$

and for accelerated ones [16, 20]

$$f(x_k) - f(x^*) \\ = O \left( \min \left\{ \frac{LR^2}{k^2} + R\delta_1 + k\delta_2, LR^2 \exp \left( \sqrt{\frac{\mu}{L}}\frac{k}{2} \right) + R\delta_1 + \sqrt{\frac{L}{\mu}}\delta_2 \right\} \right), \quad (13)$$

where  $R$  is the maximal distance between generated points and the solution.

Thus from (12), (13) we may conclude that if  $R$  is bounded,<sup>3</sup> then by choosing

$$\delta_1 = \delta_{(2)}, \delta_2 = \frac{\delta_{(2)}^2}{2L},$$

we will have the desired result: it is possible to reach  $f(x_k) - f(x^*) = \varepsilon$  with  $\delta_{(2)} \propto \varepsilon$ .

But in general situation there is a problem in the assumption "if  $R$  is bounded". As we may see from example (4) in general degenerate regime only such bound

$$R \propto R + \frac{\delta_{(2)}}{\mu} \propto R + \frac{\delta_{(2)}}{\varepsilon}$$

takes place [26]. This dependence spoils the result. The growth of  $R$  we observe in different experiments. In the paper below we investigate this problem. In particular, we propose an alternative approach to regularization<sup>4</sup> that is based on "early stopping"<sup>5</sup> of considered iterative procedure by developing proper stopping rule.

<sup>2</sup>If  $\mu \propto \varepsilon$ , we can regularize the problem and guarantee the required condition [26]. Another advantage of strong convexity is possibility to use the norm of inexact gradient for the stopping criteria [26], like in [41]. But regularization requires some prior knowledge about the size of the solution [26]. Since we typically don't have such information the procedure becomes more difficult via applying the restarts [25, 26].

<sup>3</sup>In many situations this is true. For example, when  $Q$  is bounded, when  $\mu \propto \varepsilon$ .

<sup>4</sup>By using regularization we can guarantee  $\mu \propto \varepsilon$  and therefore with  $\delta_{(2)} \propto \varepsilon$  we have the desired  $R \propto R$ .

<sup>5</sup>This terminology is popular also in Machine Learning community, where "early stopping" is used also as alternative to regularization to prevent overfitting [29].

Now we explain how to reduce relative inexactness (3) to (7) and to apply (9) when  $\mu \approx \varepsilon$ . Since  $f(x)$  has Lipschitz gradient from (3), (7) we may derive that after  $k$  iterations (where  $k$  is greater than  $\sqrt{L/\mu}$  on a logarithmic factor  $\log(LR^2/\varepsilon)$ , where  $\varepsilon$  – accuracy in function value)

$$\begin{aligned}
 f(x_k) - f(x) &\stackrel{(9),(10)}{\leq} \frac{\varepsilon}{2} + \sqrt{\frac{L}{\mu}} \frac{\delta_{(2)}^2}{\mu} + \sqrt{\frac{L}{\mu}} \frac{\delta_{(2)}^2}{\mu} \\
 &\stackrel{(3),(7)}{\leq} \sqrt{\frac{L}{\mu}} \frac{\alpha^2 \max_{t=1,\dots;k} k R f(x_t) k_2^2}{\mu} + \sqrt{\frac{L}{\mu}} \frac{2L\alpha^2 \max_{t=1,\dots;k} (f(x_k) - f(x))}{\mu} \\
 &\leq \sqrt{\frac{L}{\mu}} \frac{4L\alpha^2 (f(x_0) - f(x))}{\mu}.
 \end{aligned} \tag{14}$$

To guarantee that (restart condition)

$$f(x_k) - f(x) \leq \frac{1}{2} (f(x_0) - f(x))$$

we should have  $\alpha \leq \left(\frac{1}{L}\right)^{5/4}$ . Then we restart the method. After  $\log(L/\varepsilon)$  restarts we can guarantee the desired  $\varepsilon$ -accuracy in function value. In degenerate case the calculations are more tricky, but the idea remains the same with the replacing  $\sqrt{L/\mu}$  to  $k$  (see (9)) that lead to  $\alpha_k \leq \left(\frac{1}{k}\right)^{5/2}$ . More accurate analysis in the subsequent part of the paper allows to improve these bounds:

$$\alpha \leq \left(\frac{1}{L}\right)^{3/4}, \alpha_k \leq \left(\frac{1}{k}\right)^{3/2}.$$

Below we'll concentrate only on accelerated method and choose the method with one projection (Similar Triangles Method (STM)), see [28, 10, 30, 47, 23] and reference there in. We decided to choose this method because: 1) it's primal-dual [28]; 2) has a nice theory of how to bound  $\bar{R}$  in no noise regime [28, 37] ( $\bar{R} \approx R$ ) and noise one [30]; 3) and has previously been intensively investigated, see [23] and references there in.

### 3 Some motivation for inexact gradients

In this section we describe only two directions where inexact gradient play an important role. We emphasise that although the results below are not new, the way they are presented is of some value in our opinion and can be useful for specialist in these directions.

#### 3.1 Gradient-free methods

In this section we consider convex optimization problem:

$$\min_{x \in \mathbb{R}^n} f(x).$$

In some applications we do not have an access to gradient  $\nabla f(x)$  of target function, but can calculate the value of  $f(x)$  with accuracy  $\delta_f$  [11]:

$$|f(x) - f(x)| \leq \delta_f.$$

In this case there exist different conceptions for full gradient estimation (see [7] and references therein). For example (below we assuming that  $f$  has  $L_p$ -Lipschitz  $p$ -order derivatives in 2-norm),

**( $p$ -order finite-differences)**

$$\nabla_i f(x) = \frac{f(x + he_i) - f(x - he_i)}{2h} \text{ for } p = 2,$$

where  $e_i$  is  $i$ -th ort. Here we have

$$\delta = O\left(L_p h^p + \frac{\delta_f}{h}\right)$$

in the conception (2), see [7]. Optimal choice of  $h$  guarantees  $\delta = O\left(\frac{\delta_f}{n^{p+1}}\right)$ . From section 1 we know that it is possible to solve the problem with accuracy (in function value)  $\varepsilon = \delta$ . Hence,

$$\delta_f = O\left(\frac{\varepsilon}{n}\right)^{\frac{p+1}{p}}.$$

Unfortunately, such simple idea does not give tight lower bound in the class of algorithm that has sample complexity  $\text{Poly}(n, \frac{1}{\varepsilon})$  [44] (obtained for  $p = 0$ , that is only Lipschitz-continuity of  $f$  required):

$$\delta_f = \max\left\{\frac{\varepsilon^2}{n}, \frac{\varepsilon}{n}\right\}. \quad (15)$$

Note, that instead of finite-difference approximation approach in some applications we can use kernel approach [43, 3]. The interest to this alternative has grown last time [2, 39].

**(Gaussian Smoothed Gradients)**

$$\nabla f(x) = \frac{1}{h} \mathbb{E} f(x + he) e,$$

where  $e \sim N(0, I_n)$  is standard normal random vector. Here we have

$$\delta = O\left(n^{p-2} L_p h^p + \frac{\delta_f}{h}\right)$$

in the conception (2), see [38, 7]. Optimal choice of  $h$  guarantees  $\delta = O\left(\frac{\delta_f}{n^{p+1}}\right)$ . Hence,

$$\delta_f = O\left(\frac{\varepsilon}{n}\right)^{\frac{p+1}{p}}.$$

That is also does not match the lower bound. Moreover, here (and in the approach below) we have additional difficulty: how to estimate  $f(x)$ . We can do it only roughly, for example, by using Monte Carlo approach [7]. This is a payment for the better quality of approximation!

<sup>6</sup>Note, that the approach describe above required that function values should be available not only in  $Q$ , but also in some (depends on approach we used) vicinity of  $Q$ . This problem can be solved in a two different ways. The first one is "margins inward approach" [8]. The second one is "continuation"  $f$  to  $\mathbb{R}^n$  with preserving of convexity and Lipschitz continuity [44]:  $f_{new}(x) := f(\text{proj}_Q(x)) + \alpha \min_{y \in Q} \|x - y\|_2$ .



**(Sphere Smoothed Gradients)**

$$r f(x) = \frac{n}{h} \mathbb{E} f(x + he)e,$$

where  $e$  is random vector with uniform distribution in a unit sphere (with center at 0) in  $\mathbb{R}^n$ . Here we have

$$\delta = O\left(L_\rho h^p + \frac{n\delta_f}{h}\right)$$

in the conception (2), see [7]. Optimal choice of  $h$  guarantees  $\delta \sim (n\delta_f)^{\frac{p}{p+1}}$ . Hence,

$$\delta_f \sim \frac{\varepsilon^{\frac{p+1}{p}}}{n}.$$

That is also does not match the lower bound. One can consider that the last two approach are almost the same, but below we describe more accurate result concerning Sphere smoothing. We do not know how to obtain such a result for Gaussian smoothing. The results is as follows [16, 44]: For Sphere smoothed gradient in conception (7) we have

$$\delta \leq 2L_0 h + \frac{\rho_{-} n \delta_f R}{h}, \quad (16)$$

where  $L_0$  is Lipschitz constant of  $f$  and  $L = \min\left\{L_1, \frac{7L_0^2}{h}\right\}$  in (7), when  $p = 1$  and  $L = \frac{7L_0^2}{h}$ , when  $p = 0$ . The bound (16) is more accurate than the previous ones, since it corresponds to the first part of the lower bound (15). Indeed, by choosing properly  $h$  in (16) we obtain  $\varepsilon \sim \delta \sim n^{1-4} \delta_f^{1-2}$ . Hence,

$$\delta_f \sim \frac{\varepsilon^2}{n}.$$

The rest part ( $\delta_f \sim \frac{\varepsilon^2}{n}$ ) of lower bound (15) is also tight, see [5].

The last calculations (see (16)) additionally confirm that the conception of inexactness and algorithms we use and develop in section 2 are also tight (optimal) enough. Otherwise, it'd be hardly possible to reach lower bound by using gradient-free methods reduction to gradient ones and proposed analysis of an error accumulation for gradient-type methods.

**3.2 Inverse problems**

Another rather big direction of research where gradients are typically available only approximately is optimization in a Hilbert spaces [51]. Such optimization problems arise, in particular, in inverse problems theory [31].

We start with the reminder of how to calculate a derivative in general Hilbert space. Let

$$J(q) := J(q, u(q)),$$

where  $u(q)$  is determine as unique solution of

$$G(q, u) = 0.$$

Assume that  $G_q(q, u)$  is invertible, then

$$G_q(q, u) + G_u(q, u) r u(q) = 0,$$

hence

$$r u(q) = - [G_u(q, u)]^{-1} G_q(q, u).$$

Therefore

$$r J(q) := J_q(q, u) + J_u(q, u) r u(q) = J_q(q, u) - J_u(q, u) [G_u(q, u)]^{-1} G_q(q, u).$$

The same result could be obtained by considering Lagrange functional

$$L(q, u; \psi) = J(q, u(q)) + \langle \psi, G(q, u) \rangle$$

with

$$L_u(q, u; \psi) = 0, G_q(q, u) = 0$$

and

$$r J(q) = L_q(q, u; \psi).$$

Indeed, by simple calculations we can relate these two approaches, where

$$\psi(q, u) = - [G_u(q, u)^T]^{-1} J_u(q, u)^T.$$

Now we demonstrate this technique on inverse problem for elliptic initial-boundary value problem.

Let  $u$  be the solution of the following problem (P)

$$\begin{aligned} u_{xx} + u_{yy} &= 0, \quad x, y \in (0, 1), \\ u(1, y) &= q(y), \quad y \in (0, 1), \\ u_x(0, y) &= 0, \quad y \in (0, 1), \\ u(x, 0) &= u(x, 1) = 0, \quad x \in (0, 1). \end{aligned}$$

The first two relations

$$\begin{aligned} u_{xx} - u_{yy} &= 0, \quad x, y \in (0, 1), \\ q(y) - u(1, y) &= 0, \quad y \in (0, 1), \end{aligned}$$

we denote as  $G(q, u) = G^-(q, u) = 0$  and the last two ones as  $u \in Q$ .

Assume that we want to estimate  $q(y) \in L_2(0, 1)$  by observing  $b(y) = u(0, y) \in L_2(0, 1)$ , where  $u(x, y) \in L_2((0, 1) \times (0, 1))$  is the (unique) solution of (P) [31]. This is an inverse problem. We can reduce this problem to optimization one [31]:

$$\min_q J(q) := \min_{u: G(q; u)=0; u \in Q} J(q, u) := J(u) = \int_0^1 \int_0^1 u(0, y) - b(y) f^2 dy. \quad (17)$$

We can solve (17) numerically. This problem is convex quadratic optimization problem. We can directly apply Lagrange multipliers principle to (17), see [51]:

$$\begin{aligned} L(q, u; \psi := (\psi(x, y), \lambda(y))) &= J(u) + \langle \psi, G^-(q, u) \rangle = \int_0^1 \int_0^1 u(0, y) - b(y) f^2 dy \\ &+ \int_0^1 \int_0^1 (u_{xx} + u_{yy}) \psi(x, y) dx dy + \int_0^1 (q(y) - u(1, y)) \lambda(y) dy. \end{aligned}$$

To obtain conjugate problem for  $\psi$  we should vary  $L(q, u; \psi)$  on  $\delta u$  satisfying  $u \in Q$ :

$$\delta_u L(q, u; \psi) = 2 \int_0^1 (u(0, y) - b(y)) \delta u(0, y) dy + \int_0^1 \int_0^1 (\delta u_{xx} + \delta u_{yy}) \psi(x, y) dx dy - \int_0^1 \delta u(1, y) \lambda(y) dy, \quad (18)$$

where

$$\begin{aligned} \delta u_x(0, y) &= 0, \quad y \in (0, 1), \\ \delta u(x, 0) &= \delta u(x, 1) = 0, \quad x \in (0, 1). \end{aligned}$$

Using integration by part, from (18) we can derive

$$\begin{aligned} \delta_u L(q, u; \psi) &= \int_0^1 (2(u(0, y) - b(y)) - \psi_x(0, y)) \delta u(0, y) dy \\ &\quad - \int_0^1 \psi(1, y) \delta u_x(1, y) dy - \int_0^1 \psi(x, 1) \delta u_y(x, 1) dx + \int_0^1 \psi(x, 0) \delta u_y(x, 0) dy + \\ &\quad - \int_0^1 \int_0^1 (\psi_{xx} + \psi_{yy}) \delta u(x, y) dx dy + \int_0^1 (\psi_x(1, y) - \lambda(y)) \delta u(1, y) dy. \end{aligned}$$

Consider corresponding conjugate problem (D)

$$\begin{aligned} \psi_{xx} + \psi_{yy} &= 0, \quad x, y \in (0, 1), \\ \psi_x(0, y) &= 2(u(0, y) - b(y)), \quad y \in (0, 1), \\ \psi(1, y) &= 0, \quad y \in (0, 1), \\ \psi(x, 0) &= \psi(x, 1) = 0, \quad x \in (0, 1) \end{aligned}$$

and additional relation between Lagrange multipliers

$$\lambda(y) = \psi_x(1, y), \quad y \in (0, 1). \quad (19)$$

These relations appear since  $\delta_u L(q, u; \psi) = 0$  and  $\delta u(0, y), \delta u_x(1, y), \delta u(1, y) \in L_2(0, 1); \delta u_y(x, 1), \delta u_y(x, 0) \in L_2(0, 1); \delta u(x, y) \in L_2((0, 1) \times (0, 1))$  are arbitrary.

Since [45]

$$J(q) = \min_{u: (q; u) \in (P)} J(u) = \min_{u: G(q; u) = 0; u \in Q} J(u) = \min_{u \in Q} \max_{\psi \in (D)} L(q, u; \psi),$$

from the Demyanov–Danskin's formula [45]<sup>7</sup>

$$r J(q) = r_q \min_{u \in Q} \max_{\psi \in (D)} L(q, u; \psi) = L_q(q, u(q); \psi(q)),$$

where  $u(q)$  is the solution of (P) and  $\psi(q)$  is the solution of (D) where

$$\psi_x(0, y) = 2(u(0, y) - b(y)), \quad y \in (0, 1)$$

<sup>7</sup>The same result in more simple situation (without additional constraint  $u \in Q$ ) we consider at the beginning of this section. We don't apply Demyanov–Danskin's formula and use inverse function theorem.

and  $u(0, y)$  depends on  $q(y)$  via (P) and, at the same time, the pair  $(u(q), \psi(q))$  is the solution of

$$\min_{u \in Q} \max_{\psi \in (D)} L(q, u; \psi)$$

saddle-point problem. Since  $\delta L(q, u; \psi) = 0$  entails  $G(q, u) = 0$  that is form (P) if we add  $u \in Q$  and  $\delta_u L(q, u; \psi) = 0$ , when  $u \in Q$  entails (D) as we've shown above.

Note also that

$$L_q(q, u(q); \psi(q))(y) = \lambda(y), \quad y \in (0, 1).$$

Hence, due to (19)

$$r J(q)(y) = \psi_x(1, y), \quad y \in (0, 1)$$

So we reduce  $r J(q)(y)$  calculation to the solution of two correct initial-boundary value problem for elliptic equation in a square (P) and (D) [31].

This result can be also interpreted in a little bit different manner. We introduce a linear operator

$$A : q(y) := u(1, y) \nabla u(0, y).$$

Here  $u(x, y)$  is the solution of problem (P). It was shown in [31] that

$$A : L_2(0, 1) \rightarrow L_2(0, 1).$$

Conjugate operator is [31]

$$A : p(y) := \psi_x(0, y) \nabla \psi_x(1, y), \quad A : L_2(0, 1) \rightarrow L_2(0, 1).$$

Here  $\psi(x, y)$  is the solution of conjugate problem (D). So, by considering

$$J(q)(y) = k A q - b k_2^2,$$

we can write

$$r J(q)(y) = A (2(Aq - b)),$$

that completely corresponds to the same scheme as described above:

1. Based on  $q(y)$  we solve (P) and obtain  $u(0, y) = Aq(y)$  and define  $p(y) = 2(u(0, y) - b(y))$ .
2. Based on  $p(y)$  we solve (D) and calculate  $r J(q)(y) = A p(y) = \psi_x(1, y)$ .

So inexactness in gradient  $r J(q)$  arises since we can solve (P) and (D) only numerically.

The described above technique can be applied to many different inverse problems [31] and optimal control problems [51]. Note that for optimal control problems in practice another strategy widely used. Namely, instead of approximate calculation of gradient, optimization problem replaced by approximate one (for example, by using finite-differences schemes). For this reduced (finite-dimensional) problem the gradient is typically available precisely [24]. Moreover, in [24] the described above Lagrangian approach is based to explain the core of automatic differentiation where the function calculation tree represented as system of explicitly solvable interlocking equations.

## 4 Basic assumptions and problem description

We consider convex optimization problem on a convex (not necessarily bounded) set  $Q \subset \mathbb{R}^n$ :

$$\min_{x \in Q} f(x).$$

Assume that

$$\| \nabla f(x) - \nabla f(x) \|_2 \leq \delta, \quad (20)$$

where  $\nabla f(x)$  oracle gradient value. We consider two cases:  $Q$  is a compact set and  $Q$  is unbounded, for example  $\mathbb{R}^n$ . We define the constant:

$$R = \|x_{start} - x^*\|_2$$

to be the distance between the solution  $x^*$  and starting point  $x_{start}$ , if  $x^*$  is not unique we take such  $x^*$  that is the closest to  $x_{start}$ . We assume that function  $f$  has Lipschitz gradient with constant  $L_f$ :

$$\| \nabla f(x) - \nabla f(y) \|_2 \leq L_f \|x - y\|_2. \quad (21)$$

This implies inequality:

$$\| \nabla f(y) - \nabla f(x) \|_2 \leq L_f \|x - y\|_2 + \delta, \quad (22)$$

We will use following lemma:

**Lemma 4.1** (Fenchel inequality). *Let  $(E, \langle \cdot, \cdot \rangle, \|\cdot\|_E)$  – euclidean space, then  $\forall \lambda \in \mathbb{R}_+, \forall u, v \in E$  the inequality holds:*

$$\langle u, v \rangle \leq \frac{1}{2\lambda} \|u\|_E^2 + \frac{\lambda}{2} \|v\|_E^2.$$

From previous assumptions we can get upper bound with inexact oracle.

**Claim 1.**  $\forall x, y \in Q$ , the following estimate holds:

$$f(y) \leq f(x) + \langle \nabla f(x), y - x \rangle + \frac{L}{2} \|y - x\|_2^2 + \delta_2,$$

where  $L = 2L_f, \delta_2 = \frac{\delta^2}{2L_f}$ .

*Proof.* The proof follows from

$$\begin{aligned} f(y) &\leq f(x) + \langle \nabla f(x), y - x \rangle + \frac{L_f}{2} \|y - x\|_2^2 + \delta \\ &\leq f(x) + \langle \nabla f(x), y - x \rangle + \frac{1}{2L_f} \|\nabla f(x) - \nabla f(x)\|_2^2 + \frac{L_f}{2} \|y - x\|_2^2 + \delta \\ &\leq f(x) + \langle \nabla f(x), y - x \rangle + \frac{L}{2} \|y - x\|_2^2 + \delta_2. \end{aligned}$$

□

We also assume strong convexity of  $f$  with parameter  $\mu$ , however  $\mu$  may equal zero – this corresponds to the ordinary convexity, supposed initially. Further we will use only a consequence of this:

$$f(x) + h\Gamma f(x), y \quad xi + \frac{\mu}{2}kx \quad yk_2^2 \subset f(y). \quad (23)$$

We obtain similar to claim 1 two lower bounds with inexact oracle.

**Claim 2.**  $\forall x, y \in Q$ , the following estimate holds:

$$f(x) + h\Gamma f(x), y \quad xi + \frac{\mu}{2}kx \quad yk_2^2 \quad \delta_1 kx \quad yk_2 \subset f(y),$$

where  $\delta_1 = \delta$ .

*Proof.* Using Cauchy inequality and (23) we obtain:

$$\begin{aligned} f(x) + h\Gamma f(x), y \quad xi + \frac{\mu}{2}kx \quad yk_2^2 \quad \delta_1 kx \quad yk_2 \subset f(x) + \\ + h\Gamma f(x), y \quad xi + \frac{\mu}{2}kx \quad yk_2^2 \quad k\Gamma f(x) \quad \Gamma f(x)k_2 kx \quad yk_2 \subset \\ \subset f(x) + h\Gamma f(x), y \quad xi + \frac{\mu}{2}kx \quad yk_2^2 \\ h\Gamma f(x) \quad \Gamma f(x), y \quad xi = f(x) + h\Gamma f(x), y \quad xi + \frac{\mu}{2}kx \quad yk_2^2 \subset f(y) \quad ) \\ f(x) + h\Gamma f(x), y \quad xi + \frac{\mu}{2}kx \quad yk_2^2 \quad \delta_1 kx \quad yk_2 \subset f(y). \end{aligned}$$

□

**Claim 3.**  $\forall x, y \in Q$ , if in (23)  $\mu \neq 0$ , the following estimate holds,

$$f(x) + h\Gamma f(x), y \quad xi + \frac{\mu}{4}ky \quad xk_2^2 \quad \delta_3 \subset f(y),$$

where  $\delta_3 = \frac{\delta^2}{2}$ .

*Proof.* Trivial calculations bring

$$\begin{aligned} f(x) + h\Gamma f(x), y \quad xi + \frac{\mu}{4}kx \quad yk_2^2 \quad \delta_3 = f(x) + h\Gamma f(x), y \quad xi + \\ + h\Gamma f(x) \quad \Gamma f(x), y \quad xi + \frac{\mu}{4}kx \quad yk_2^2 \quad \delta_3. \end{aligned}$$

Using lemma 1 we obtain:

$$\begin{aligned} f(x) + h\Gamma f(x), y \quad xi + \frac{\mu}{4}kx \quad yk_2^2 \quad \delta_3 \subset f(x) + \\ + h\Gamma f(x), y \quad xi + \frac{\delta^2}{\mu} + \frac{\mu}{4}kx \quad yk_2^2 + \frac{\mu}{4}ky \quad xk_2^2 \quad \delta_3 = \\ = f(x) + h\Gamma f(x), y \quad xi + \frac{\mu}{2}ky \quad xk_2^2 \subset f(y). \end{aligned}$$

□

The last two inequalities give different results in convergence under certain conditions. We will study two models based on statements 2, 3 and we will denote them by the index  $\tau$ , that is denote:

$$\begin{aligned} \mu_1 &= \mu, \\ \mu_2 &= \frac{\mu}{2}. \end{aligned} \quad (24)$$

Further in the text, we will use statements 3 and 2 in the notation corresponding to (24).

## 5 Similar Triangles Method and its properties

In this section we describe an accelerated method we choose to investigate gradient-error accumulation.

---

**Algorithm 1**  $STM(L, \mu, \tau, x_{start}), Q = \mathbb{R}^n$

---

**Input:** Starting point  $x_{start}$ , number of steps  $N$

**Output:**  $x_N$

- 1: **Set**  $x_0 = x_{start}$
- 2: **Set**  $A_0 = \frac{1}{L}, \alpha_0 = \frac{1}{L},$
- 3: **Set**  $\psi_0(x) = \frac{1}{2}kx - x_0k_2^2 + \alpha_0 \left( f(x_0) + h\Gamma f(x_0), x - x_0 + \frac{1}{2}kx - x_0k_2^2 \right),$
- 4: **Set**  $z_0 = \underset{y \in Q}{\operatorname{argmin}} \psi_0(y),$
- 5: **Set**  $x_0 = z_0.$
- 6: **for**  $k = 1, 2 \dots N$  **do**
- 7:  $\alpha_k = \frac{1 + \frac{A_{k-1}}{2L}}{2L} + \sqrt{\frac{1 + \frac{A_{k-1}}{4L^2}}{4L^2} + \frac{A_{k-1}}{1 + \frac{A_{k-1}}{L}}},$
- 8:  $A_k = A_{k-1} + \alpha_k,$
- 9:  $\tilde{x}_k = \frac{A_{k-1}x_{k-1} + k z_{k-1}}{A_k},$
- 10:  $\psi_k(x) = \psi_{k-1}(x) + \alpha_k \left( (f(\tilde{x}_k) + h\Gamma f(\tilde{x}_k), x - \tilde{x}_k + \frac{1}{2}kx - \tilde{x}_k k_2^2) \right),$
- 11:  $z_k = \underset{y \in Q}{\operatorname{argmin}} \psi_k(y),$
- 12:  $x_k = \frac{A_{k-1}x_{k-1} + k z_k}{A_k}.$
- 13: **end for**
- 14: **return**  $x_N$

---

Figure 5 describes the position of the vertices. On the sides, not their lengths are marked, but the relationships in the corresponding sides in the similarity of triangles. In the case  $Q = \mathbb{R}^n$ , we can simplify the step of the algorithm by replacing it with:

$$z_k = z_{k-1} - \frac{\alpha_k}{1 + A_k \mu} \left( \Gamma f(\tilde{x}_k) + \mu (z_{k-1} - \tilde{x}_k) \right).$$

We define constant:

$$R = \max_{0 \leq k \leq N} f(z_k - x_k, k z_k - x_k, k \tilde{x}_k - x_k)g.$$

We will also write down several identities that we will need in the proofs

$$\begin{aligned} A_k(x_k - \tilde{x}_k) &= \alpha_k(z_k - \tilde{x}_k) + A_{k-1}(x_{k-1} - \tilde{x}_k), \\ \frac{1 + \mu A_{k-1}}{2}kz_k - z_{k-1}k_2^2 &= \frac{L}{2}kx_k - \tilde{x}_k k_2^2, \\ A_{k-1}k\tilde{x}_k - x_{k-1}k_2 &= \alpha_k k \tilde{x}_k - z_{k-1}k_2. \end{aligned} \tag{25}$$

Some of the identities can be obtained from geometric considerations, for example, from a figure, others by direct substitution into the definitions of the sequences  $x_k, \tilde{x}_k, z_k$ . Also very important are the estimates for the sequence  $A_k$ .

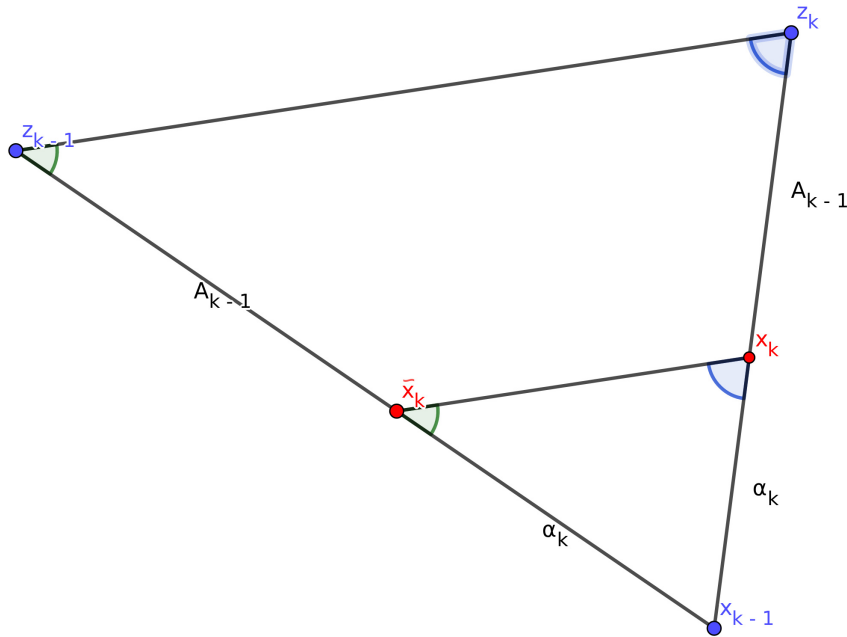


Figure 1: Geometry of Similar Triangles method [28]

**Claim 4.** If  $\mu \notin 0$  and  $\delta k \geq 2 \in \mathbb{N}$  the following inequality holds:

$$A_k > A_{k-1} \lambda_{;L},$$

where

$$\theta_{;L} = \frac{\mu}{L}, \lambda_{;L} = \left( 1 + \frac{1}{2} \theta_{;L} + \frac{1}{2} \sqrt{\theta^2_{;L} + 4\theta_{;L}} \right).$$

*Proof.* Using the definition of the sequences  $A_k$  and solving the quadratic equation, we obtain:

$$\begin{aligned} A_k(1 + \mu A_{k-1}) &= L\alpha_k^2 = L(A_k - A_{k-1})^2 = LA_k^2 - 2LA_kA_{k-1} + LA_{k-1}^2, \\ \therefore A_k^2 - A_k(1 + A_{k-1}(2 + \theta_{;L})) + A_{k-1}^2 &= 0, D = (1 + 2A_{k-1} + \theta_{;L})^2 - 4A_{k-1}^2 \\ A_{k:apex} &= \frac{1}{2} + \left( 1 + \frac{1}{2} \theta_{;L} \right) A_{k-1} \Rightarrow A_k = \frac{1}{2} \left( (1 + (2 + \theta_{;L})) A_{k-1} + \sqrt{D} \right) \\ \sqrt{D} &= \sqrt{1 + (2\theta_{;L} + 4) A_{k-1} + (\theta^2_{;L} + 4\theta_{;L}) A_{k-1}^2} > A_{k-1} \sqrt{\theta^2_{;L} + 4\theta_{;L}} \\ \Rightarrow A_k &> \frac{1}{2} \left( 2 + \theta_{;L} + \sqrt{\theta^2_{;L} + 4\theta_{;L}} \right) A_{k-1} = \\ &= \left( 1 + \frac{1}{2} \theta_{;L} + \frac{1}{2} \sqrt{\theta^2_{;L} + 4\theta_{;L}} \right) A_{k-1} \\ \Rightarrow A_k &> \frac{1}{L} \lambda^k_{;L}, \lambda_{;L} = \left( 1 + \frac{1}{2} \theta_{;L} + \frac{1}{2} \sqrt{\theta^2_{;L} + 4\theta_{;L}} \right). \end{aligned}$$

□



**Corollary 5.1.**

$$\lambda_{;L} > \left(1 + \frac{1}{2}\sqrt{\theta_{;L}}\right)^2 = \left(1 + \sqrt{\theta_{;L}} + \frac{1}{4}\theta_{;L}\right),$$

$$\left(1 + \frac{1}{2}\sqrt{\theta_{;L}}\right)^2 > e^{\frac{1}{2}\theta_{;L}}.$$

**Claim 5.** If  $\mu \in (0, \theta_{;L}] \subset \mathbb{N}$  the following inequality holds:

$$\frac{\sum_{j=0}^k A_j}{A_k} \leq 1 + \sqrt{\frac{L}{\mu}}.$$

*Proof.* According to the previous designations:

$$\lambda_{;L} = \left(1 + \frac{1}{2}\theta_{;L} + \frac{1}{2}\sqrt{\theta_{;L}^2 + 4\theta_{;L}}\right), \theta_{;L} = \frac{\mu}{L}.$$

Using previous claim we can reduce this amount exponentially.

$$\frac{\sum_{j=0}^k A_j}{A_k} \leq \sum_{j=0}^k \lambda_{;L}^{-j} = \frac{\lambda_{;L}^{k+1} - 1}{\lambda_{;L}^k (\lambda_{;L} - 1)} \leq \frac{\lambda_{;L} - 1}{\lambda_{;L} - 1} \leq 1 + \sqrt{\frac{L}{\mu}}.$$

□

**Claim 6.** If  $\mu = 0$  then:

$$A_k > \frac{(k+1)^2}{4L}.$$

*Proof.* If  $\mu = 0$  then  $A_k = L\alpha_k^2$  and solving quadratic equation we get:

$$\alpha_k = \frac{1 + \sqrt{1 + 4L^2\alpha_{k-1}^2}}{2L}.$$

Then by induction it is easy to get that:

$$\alpha_k > \frac{k+1}{2L} \Rightarrow A_k = L\alpha_k^2 > \frac{(k+1)^2}{2L}.$$

□

**Claim 7.** If  $\mu = 0$  we have:

$$\frac{\sum_{j=0}^k A_j}{A_k} \leq k.$$

*Proof.* The proof follows from the simple calculations:

$$\frac{\sum_{j=0}^k A_j}{A_k} \leq \frac{\sum_{j=0}^k a_j^2}{a_k^2} \leq \frac{k\alpha_k^2}{\alpha_k^2} = \frac{k\alpha_k^2}{\left(\frac{1}{2L} + \sqrt{\frac{1}{4L^2} + \alpha_k^2}\right)^2} \leq k.$$

□

**Lemma 5.2.**  $\delta k > 1$  the following inequality holds:

$$\psi_k(z_k) > \psi_{k-1}(z_{k-1}) + \frac{1 + \mu A_{k-1}}{2} k z_k - z_{k-1} k_2^2 + \alpha_k \left( f(x_k) + h\Gamma f(x_k), z_k - x_{k-1} + \frac{\mu}{2} k z_k - x_k k_2^2 \right).$$

*Proof.* From the definition of the  $\psi_{k-1}$  function, it has a minimum at the point  $z_{k-1}$ , then:

$$\begin{aligned} h\Gamma \psi_{k-1}(z_{k-1}), z_k - z_{k-1} &> 0, \quad \Gamma \psi_{k-1}(z_{k-1}) = (z_{k-1} - x_0) + \\ &+ \sum_{j=0}^{k-1} \alpha_j \left( \Gamma f(x_j) + \mu (z_{k-1} - x_j) \right) \\ \Rightarrow \psi_k(z_k) &= \psi_{k-1}(z_{k-1}) + \alpha_k \left( f(x_k) + h\Gamma f(x_k), z_k - x_{k-1} + \frac{\mu}{2} k z_k - x_k k_2^2 \right) = \\ &= \frac{1}{2} k z_k - x_0 k_2^2 + \sum_{j=0}^{k-1} \alpha_j \left( f(x_j) + h\Gamma f(x_j), z_k - x_j + \frac{\mu}{2} k z_k - x_j k_2^2 \right) + \\ &+ \alpha_k \left( f(x_k) + h\Gamma f(x_k), z_k - x_{k-1} + \frac{\mu}{2} k z_k - x_k k_2^2 \right). \end{aligned}$$

From (25) and the above we obtain:

$$\begin{aligned} \psi_k(z_k) &> \frac{1}{2} k z_k - x_0 k_2^2 + h z_k - x_0, z_k - z_{k-1} + \frac{1}{2} k z_k - z_k k_2^2 + \\ &+ \sum_{j=0}^{k-1} \alpha_j \left( f(x_j) + h\Gamma f(x_j), z_k - x_j + \frac{\mu}{2} k z_k - x_j k_2^2 \right) + \\ &+ \alpha_k \left( f(x_k) + h\Gamma f(x_k), z_k - x_{k-1} + \frac{\mu}{2} k z_k - x_k k_2^2 \right) = \\ &= \sum_{j=0}^{k-1} \alpha_j \left( h\Gamma f(x_j) + \mu (z_{k-1} - x_j), z_{k-1} - z_{k-1} \right) + \\ &+ \sum_{j=0}^{k-1} \alpha_j \left( f(x_j) + h\Gamma f(x_j), z_k - x_j + \frac{\mu}{2} k z_k - x_j k_2^2 \right) + \\ &+ \alpha_k \left( f(x_k) + h\Gamma f(x_k), z_k - x_{k-1} + \frac{\mu}{2} k z_k - x_k k_2^2 \right) + \frac{1}{2} k z_k - x_0 k_2^2 + \frac{1}{2} k z_k - z_k k_2^2. \end{aligned}$$

Using the linearity of the dot product, we split the sum by two and apply to:

$$\mu h z_k - x_j, z_{k-1} - z_{k-1}.$$

Equality from (25), and finally get:

$$\begin{aligned}
\psi_k(z_k) &> \frac{1}{2}kz_{k-1} \quad x_0k_2^2 + \frac{1 + \mu A_{k-1}}{2}kz_{k-1} \quad z_kk_2^2 + \\
&+ \sum_{j=0}^{k-1} \alpha_j \left( f(x_j) + h\Gamma f(x_j), z_{k-1} \quad x_ji + \frac{\mu}{2}kz_{k-1} \quad x_jk_2^2 \right) + \\
&+ \alpha_k \left( f(x_k) + h\Gamma f(x_k), z_k \quad x_ki + \frac{\mu}{2}kz_k \quad x_kk_2^2 \right) = \\
&= \psi_{k-1}(z_{k-1}) + \frac{1 + \mu A_{k-1}}{2}kz_k \quad z_{k-1}k_2^2 + \\
&+ \alpha_k \left( f(x_k) + h\Gamma f(x_k), z_k \quad x_ki + \frac{\mu}{2}kz_k \quad x_kk_2^2 \right).
\end{aligned}$$

□

### Remark 1.

In the case  $\mu = 0$ , we obtain a corollary from the strongly convexity of functions  $\psi_k$  and their definition, that is:

$$\begin{aligned}
\psi_k(z_k) &= \psi_{k-1}(z_{k-1}) + \alpha_k \left( f(x_k) + h\Gamma f(x_k), z_k \quad x_ki \right) \\
\psi_k(z_k) &> \psi_{k-1}(z_{k-1}) + \frac{1}{2}kz_k \quad z_{k-1}k_2^2 + \alpha_k \left( f(x_k) + h\Gamma f(x_k), z_k \quad x_ki \right).
\end{aligned}$$

## 6 Main results

Here we will describe some results based on the previously presented lemmas and statements.

### 6.1 Additive noise and main theorems.

**Theorem 6.1.** *For  $k \in \mathbb{N}$  the following inequality holds:*

$$A_k f(x_k) \leq \psi_k(z_k) + \delta_2 \sum_{j=0}^k A_j + 2R\delta_1 A_k.$$

*Proof.* Base,  $k = 0$ :

$$\begin{aligned}
f(x_0) &\leq f(x_0) + h\Gamma f(x_0), x_0 \quad x_0i + \frac{L}{2}kx_0 \quad x_0k_2^2 + \delta_2 \leq \\
&\leq L\psi_0(z_0) + \frac{L\mu}{2}kz_0 \quad x_0k_2^2 + \delta_2 \leq L\psi_0(z_0) + \delta_2.
\end{aligned}$$

Induction step:

$$\begin{aligned}
A_k f(x_k) &\leq A_k \delta_1 kx_{k-1} \quad x_kk_2 \leq \\
&\leq A_k \left( f(x_k) + h\Gamma f(x_k), x_k \quad x_ki + \frac{L}{2}kx_k \quad x_kk_2^2 + \delta_2 \right) \leq A_k \delta_1 kx_{k-1} \quad x_kk_2.
\end{aligned}$$

Using equations (25) we obtain:

$$\begin{aligned}
& A_k f(x_k) \leq A_{k-1} \delta_1 k x_{k-1} - x_k k_2 \leq \\
& \leq A_{k-1} \left( f(x_k) + h \Gamma f(x_k), x_{k-1} - x_k i \right) + \alpha_k \left( f(x_k) + h \Gamma f(x_k), z_k - x_k i \right) + \\
& + \frac{(1 + \mu_1 A_{k-1})}{2} k z_k - z_{k-1} k_2^2 + A_k \delta_2 - A_{k-1} \delta_1 k x_{k-1} - x_k k_2 \leq \\
& \leq A_{k-1} f(x_{k-1}) + \alpha_k (f(x_k) + h \Gamma f(x_k), z_k - x_k i) + \\
& + \frac{1 + \mu_1 A_{k-1}}{2} k z_k - z_{k-1} k_2^2 + A_k \delta_2.
\end{aligned}$$

Using the induction hypothesis, we obtain:

$$\begin{aligned}
A_k f(x_k) & \leq A_{k-1} \delta_1 k x_{k-1} - x_k k_2 \leq \psi_{k-1}(z_{k-1}) + \delta_2 \sum_{j=0}^{k-1} A_j + 2R \delta_1 A_{k-1} + \\
& + \frac{1 + \mu_1 A_{k-1}}{2} k z_k - z_{k-1} k_2^2 + \alpha_k \left( f(x_k) + h \Gamma f(x_k), z_k - x_k i \right) + A_k \delta_2.
\end{aligned}$$

Using lemma 5.2 we can get:

$$\begin{aligned}
A_k f(x_k) & \leq A_{k-1} \delta_1 k x_{k-1} - x_k k_2 + \psi_k(z_k) + \delta_2 \sum_{j=0}^k A_j + 2R \delta_1 A_{k-1} = \\
& = \psi_k(z_k) + \delta_2 \sum_{j=0}^k A_j + 2R \delta_1 A_{k-1} + \alpha_k k x_k - z_{k-1} k_2 \leq \\
& \leq \psi_k(z_k) + \delta_2 \sum_{j=0}^k A_j + 2R \delta_1 A_{k-1} + \alpha_k (k z_{k-1} - x_{k-2} + k x_k - x_{k-2}) \delta_1 \leq \\
& \leq \psi_k(z_k) + \delta_2 \sum_{j=0}^k A_j + 2R \delta_1 A_{k-1} + 2\alpha_k R \\
& ) \quad A_k f(x_k) \leq \psi(z_k) + \delta_2 \sum_{j=0}^k A_j + 2R \delta_1 A_k.
\end{aligned}$$

□

## Remark 2.

We should note that this inequality is true both in the case of  $\mu \neq 0$  and in the case of  $\mu = 0$ .

**Theorem 6.2.** *If  $\mu \neq 0$   $\forall k \in \mathbb{N}$  the following inequality holds:*

$$A_k f(x_k) \leq \psi_k(z_k) + \delta_2 \sum_{j=0}^k A_j + \delta_3 \sum_{j=0}^{k-1} A_j.$$

*The proof repeats verbatim theorem 6.1, except for claim 2, replaced by claim 3.*

**Theorem 6.3.** *If  $\delta_1 = \delta_2 = \delta_3 = 0$  then  $R = R$ .*

*Proof.* Using theorem 6.1 we get  $A_k f(x_k) \in \psi_k(z_k)$  then:

$$\begin{aligned} \frac{1}{2}kz_k \quad x \quad k_2^2 &= \frac{1}{2}kz_k \quad x \quad k_2^2 + A_k f(x_k) \quad A_k f(x_k) \in \\ &\in \psi_k(z_k) + \frac{1}{2}kz_k \quad x \quad k_2^2 \quad A_k f(x_k) \in \sum_{j=0}^k \alpha_j (f(x_j) + h\Gamma f(x_j), x \quad x_{k+1} + \\ &+ \frac{\mu}{2}kx \quad x_k k_2^2) + \frac{1}{2}kx \quad x_0 k_2^2 \in A_k f(x_k) \quad A_k f(x_k) + \frac{1}{2}kx_0 \quad x \quad k_2^2 = \frac{1}{2}R^2. \end{aligned}$$

We now prove bound for the sequence  $x_k$ , similarly for  $x_k$ . We prove by induction, so assume fairness for  $k-1$ , base is obvious.

$$\begin{aligned} kx_k \quad x \quad k_2^2 &= k \frac{A_{k-1}}{A_k} (x_{k-1} \quad x) + \frac{\alpha_k}{A_k} (z_k \quad x) k_2^2 \in \\ &\in \frac{A_{k-1}}{A_k} kx_{k-1} \quad x \quad k_2^2 + \frac{\alpha_k}{A_k} kz_k \quad x \quad k_2^2 = R^2. \end{aligned}$$

□

**Theorem 6.4** (convergence in function). *Both inequalities take place with  $\mu \notin 0$*

$$\begin{aligned} f(x_N) \quad f(x) &\in LR^2 \exp\left(\frac{1}{2}\sqrt{\frac{\mu}{L}}N\right) + \left(1 + \sqrt{\frac{L}{\mu_1}}\right) \delta_2 + 3R\delta_1, \\ f(x_N) \quad f(x) &\in LR^2 \exp\left(\frac{1}{2}\sqrt{\frac{\mu}{2L}}N\right) + \left(1 + \sqrt{\frac{L}{\mu_2}}\right) \delta_2 + \left(1 + \sqrt{\frac{L}{\mu_2}}\right) \delta_3. \end{aligned}$$

*Proof.* Using, all of the above is easy to show what is required, the proof of both convergence is the same with the replacement of theorem 6.1 by theorem 6.2 and replacement claim 2 by claim 3, therefore, we present only the proof of the first inequality.

$$\begin{aligned} A_N f(x_N) &\in \psi_N(z_N) + \delta_2 \sum_{j=0}^N A_j + 2R\delta_1 A_N \in \frac{1}{2}kx \quad x_0 k_2^2 + \\ &+ \delta_2 \sum_{j=0}^N A_j + 2R\delta_1 A_N + \sum_{j=0}^N \alpha_k (f(x_j) + h\Gamma f(x_j), x \quad x_{j+1} + \frac{\mu_1}{2}kx \quad x_j k_2^2) \in \\ &\in \delta_2 \sum_{j=0}^N A_j + 2R\delta_1 A_N + \sum_{j=0}^N \alpha_k (R\delta_1 + f(x)) + \frac{1}{2}R^2 = \\ &= \delta_2 \sum_{j=0}^N A_j + 3R\delta_1 A_N + A_N f(x) + \frac{1}{2}R^2 \quad ( ) \\ ( ) \quad f(x_N) \quad f(x) &\in LR^2 \exp\left(\frac{1}{2}\sqrt{\frac{\mu}{L}}N\right) + \left(1 + \sqrt{\frac{L}{\mu_1}}\right) \delta_2 + 3R\delta_1. \end{aligned}$$

□

**Remark 3.**

If  $\mu = 0$  we can get analogue of the first convergence, repeating the proof using claims 6, 7

$$f(x_N) - f(x) \leq \frac{4LR^2}{N^2} + 3R\delta_1 + N\delta_2.$$

**Remark 4.**

Suppose  $\mu = 0$ , then consider the auxiliary problem:

$$\begin{aligned} f(x) &= f(x) + \frac{\mu}{2}kx - x_0k_2 \quad \min_{x \in Q}, \\ \Gamma f(x) &= \Gamma f(x) + \mu(x - x_0), \\ \Gamma f(x) &= \Gamma f(x) + \mu(x - x_0), \\ k\Gamma f(x) - \Gamma f(x)k_2 &= k\Gamma f(x) - \Gamma f(x)k_2 \leq \delta. \end{aligned}$$

The resulting function will satisfy the condition that the gradient is Lipschitz, that is  $\exists x, y \in Q$ :

$$\begin{aligned} k\Gamma f(x) - \Gamma f(y)k_2 &= k(\Gamma f(x) - \Gamma f(y)) + \mu(x - y)k_2 \leq \\ &\leq k(\Gamma f(x) - \Gamma f(y))k_2 + \mu kx - yk_2 \leq \\ &\leq L_f kx - yk_2 + \mu kx - yk_2 \leq (L_f + \mu)kx - yk_2. \end{aligned}$$

We will assume, that  $\mu < 1$ . That is, we can let  $L = 2(L_f + 1) = L + 2 > L_f$ . The resulting function will already be strongly convex, which means that the second model is applicable to it  $\tau = 2$ . Using theorem 6.4 we can get the following inequality:

$$\begin{aligned} x &= \operatorname{argmin}_{x \in Q} f(x), \\ R &= kx - x_0k_2, \\ f(x_k) - f(x) &\leq \frac{LR^2}{2\lambda_{\frac{k}{2}, 2L}} + \left(1 + \sqrt{\frac{2L+4}{\mu}}\right) (\delta_2 + \delta_3) \\ f(x_k) - f(x) &\leq LR^2 \exp\left(\frac{1}{2}\sqrt{\frac{\mu}{2(L+2)k}}\right) + \left(1 + \sqrt{\frac{2L+4}{\mu}}\right) \left(\frac{1}{L} + \frac{1}{\mu}\right) \delta^2, \\ f(x) &\leq f(x) + \frac{\mu}{2}R^2. \end{aligned}$$

Then we can get convergence rate for not regularized function:

$$\begin{aligned} f(x_k) - f(x) &\leq f(x_k) - f(x) \leq f(x_k) - f(x) + \frac{\mu}{2}R^2 \leq \\ &\leq LR^2 \exp\left(\frac{1}{2}\sqrt{\frac{\mu}{2(L+2)k}}\right) + \\ &\quad + \left(1 + \sqrt{\frac{2L+4}{\mu}}\right) \left(\frac{1}{L} + \frac{1}{\mu}\right) \delta^2 + \frac{\mu}{2}R^2. \end{aligned}$$

Using strong convexity of the function  $f$  we get:

$$\begin{aligned} f(x) + \frac{\mu}{2}R^2 &\leq f(x) + \frac{\mu}{2}R^2 = f(x) \leq f(x) = f(x) + \frac{\mu}{2}R^2 \\ R &\leq R. \end{aligned}$$

Finally we get convergence:

$$f(x_k) - f(x^*) \leq LR^2 \exp\left(\frac{1}{2} \sqrt{\frac{\mu}{2(L+2)}} k\right) + \left(1 + \sqrt{\frac{2L+4}{\mu}}\right) \left(\frac{1}{L} + \frac{1}{\mu}\right) \delta^2 + \frac{\mu}{2} R^2.$$

We choose value for parameter  $\mu$  in the remark 9.

**Remark 5.**

If we consider the problem in the first model  $\tau = 1$ , the case  $\mu = 0$  and assume that  $kx_k \in R$ . Then we choose a starting point for the *STM* algorithm in a ball of radius  $R$ , specifically put  $x_{start} = 0$ .

$$R = kx_k \quad x_0 k_2 \in R.$$

Let us formulate a stopping rule for the this model ( $\delta\zeta > 0$ ).

$$f(x_k) - f(x^*) \leq k\delta_2 + R\delta_1 + \delta_1 \sum_{j=1}^k \frac{\alpha_j}{A_k} kx_j \quad z_{j-1} k_2 + \zeta.$$

**Lemma 6.5** (Bound for  $R$ ). *Before the stopping criterion is satisfied, the following inequality holds:*

$$R \leq R.$$

*Proof.* Note, that from  $kz_k \in R$  we get  $kx_k \in R$ ,  $kx_k \in R$  similarly to theorem 6.3. But it's worth noting that to estimate  $kx_k \in R$ , only inequalities are required for all  $j \leq k-1$ .

$$\begin{aligned} kx_k - x_k k_2 &= k \frac{A_{k-1}}{A_k} (x_{k-1} - x) + \frac{\alpha_k}{A_k} (z_{k-1} - x) k_2 \leq \\ &\leq \frac{A_{k-1}}{A_k} kx_{k-1} - x_k k_2 + \frac{\alpha_k}{A_k} kz_{k-1} - x_k k_2 \leq R. \end{aligned}$$

An analysis of the proof of theorem 6.2 gives a stronger convergence:

$$A_k f(x_k) \leq \psi_k(z_k) + \delta_2 \sum_{j=0}^k A_j + \delta_1 \sum_{j=1}^k \alpha_j kx_j \quad z_{j-1} k_2^2.$$

Then, using the convexity of the function  $\psi_k$  we get:

$$\begin{aligned}
& A_k f(x_k) + \frac{1}{2} k z_k \quad x \quad k_2 \leq \frac{1}{2} k z_k \quad x \quad k_2 + \psi_k(z_k) + \delta_2 \sum_{j=0}^k A_j + \\
& \quad + \delta_1 \sum_{j=1}^k \alpha_j k x_j \quad z_j \quad 1 k_2 \leq \psi_k(x) + \delta_2 \sum_{j=0}^k A_j + \\
& \quad + \delta_1 \sum_{j=1}^k \alpha_j k x_j \quad z_j \quad 1 k_2 \leq \frac{1}{2} R^2 + A_k f(x) + \delta_2 \sum_{j=0}^k A_j + \\
& + \delta_1 \sum_{j=1}^k \alpha_j k x_j \quad z_j \quad 1 k_2 + \delta_1 \sum_{j=0}^k \alpha_j k z_k \quad x \quad k_2 \quad \frac{1}{2} (R^2 \quad k z_k \quad x \quad k_2) > \\
& > A_k \left( (f(x_k) \quad f(x)) \quad \left( k \delta_2 + \delta_1 \sum_{j=1}^k \frac{\alpha_j}{A_k} k x_j \quad z_j \quad 1 k_2 + R \delta_1 + \zeta \right) \right) > 0.
\end{aligned}$$

□

Therefore, when the stopping criterion is met, we will receive the estimate:

$$f(x_k) \quad f(x) \leq k \delta_2 + \delta_1 R + \delta_1 \sum_{j=0}^k \alpha_j k x_j \quad z_j \quad 1 k_2 + \zeta.$$

From remark 3 we get an estimate of the number of iterations:

$$N_{stop} > 2 \sqrt{\frac{LR^2}{\zeta}}.$$

$$\begin{aligned}
f(x_N) \quad f(x) & \leq \frac{4LR^2}{N^2} + N \delta_2 + R \delta_1 + \delta_1 \sum_{j=1}^N \frac{\alpha_j}{A_N} k x_j \quad z_j \quad 1 k_2 \leq \\
& \leq N \delta_2 + R \delta_1 + \delta_1 \sum_{j=1}^N k x_j \quad z_j \quad 1 k_2 + \zeta \quad \frac{4LR^2}{N^2} \leq \zeta, \quad N^2 > \frac{4LR^2}{\zeta}.
\end{aligned}$$

Summing up, we obtain the following theorem:

**Theorem 6.6.** For model  $\tau = 1$  with  $\mu = 0$ , using stopping rule:

$$f(x_N) \quad f(x) \leq N \delta_2 + R \delta_1 + \delta_1 \sum_{j=1}^N \frac{\alpha_j}{A_N} k x_j \quad z_j \quad 1 k_2 + \zeta.$$

We can guarantee, that:

$$R \leq R.$$

And the criterion is reached after:

$$N_{stop} = \left\lceil 2 \sqrt{\frac{LR^2}{\zeta}} \right\rceil + 1.$$



## 6.2 Relative noise.

Recall that we denote:

$$L = 2L_f.$$

Where  $L_f$  – Lipschitz constant of  $r f$ . From theorem 6.2 and similar to theorem 6.4 reasoning we get:

$$\begin{aligned} f(x_k) - f(x) &\leq \frac{R^2}{A_k} + \frac{3}{2} \sum_{j=0}^{k-1} \frac{A_j \alpha^2 k r f(x_k) k_2^2}{A_k \mu} + \frac{3\alpha^2 k r f(x_k) k_2^2}{2\mu}, \\ k &= f(x_k) - f(x), \\ k &\leq \frac{R^2}{A_k} + \frac{3}{2} \sum_{j=0}^{k-1} \frac{A_j \alpha^2 k r f(x_j) k_2^2}{A_k \mu} + \frac{3\alpha^2 k r f(x_k) k_2^2}{2\mu}. \end{aligned}$$

We define:

$$\begin{aligned} \theta &= \frac{3L\alpha^2}{2\mu(1 - \frac{3L}{2})}, \\ \lambda &= \frac{R^2}{1 - \frac{3L}{2}}. \end{aligned}$$

From inequality:

$$k r f(x_k) k_2^2 \leq L (f(x_k) - f(x)).$$

We obtain:

$$k \leq \frac{\lambda}{A_k} + \theta \sum_{j=0}^{k-1} \frac{A_j}{A_k} \cdot j.$$

In these designations by induction we can obtain:

**Claim 8.**

$$k \leq \frac{(1 + \theta)^{k-1}}{A_k} \lambda + \theta \frac{A_0 (1 + \theta)^{k-1}}{A_k} \cdot 0.$$

*Proof.* Base,  $k = 1$  is obvious. Induction step:

$$\begin{aligned} k &\leq \frac{\lambda}{A_k} + \theta \sum_{j=0}^{k-1} \frac{A_j}{A_k} \cdot j \leq \\ &\leq \frac{\lambda}{A_k} + \sum_{j=0}^{k-1} \left( \frac{A_j (1 + \theta)^{j-1}}{A_k} \lambda + \theta \frac{A_0 (1 + \theta)^{j-1}}{A_k} \cdot 0 \right) + \frac{A_0}{A_k} \cdot 0 \leq \\ &\leq \frac{\lambda}{A_k} + \sum_{j=0}^{k-1} \left( \frac{\lambda (1 + \theta)^j}{A_k} + \theta \frac{A_0 (1 + \theta)^j}{A_k} \cdot 0 \right) + \frac{A_0}{A_k} \cdot 0 = \\ &= \frac{(1 + \theta)^{k-1}}{A_k} \lambda + \theta \frac{A_0 (1 + \theta)^{k-1}}{A_k} \cdot 0. \end{aligned}$$

□

That is we can formulate the following inequality:

$$f(x_k) - f(x) \leq \frac{\lambda(1+\theta)^k}{A_k} + \theta \frac{A_0(1+\theta)^k}{A_k} (f(x_0) - f(x)).$$

Using corollary 5.1 we can estimate:

$$A_k > \left(1 + \sqrt{\frac{\mu}{2L}}\right)^k A_0. \quad (26)$$

We will choose an alpha such that:

$$\frac{1+\theta}{A_k} \leq \frac{1}{1 + \frac{1}{3^{1/2}} \sqrt{\frac{\mu}{L}}}.$$

Using (26) and definition of  $\theta$  we obtain, that we should choose  $\alpha$  from:

$$\begin{aligned} \alpha &\leq \sqrt{\frac{1}{\frac{1}{2^{1/2}} \left(\frac{\mu}{L}\right)^{3/2} + 2L}} \\ \alpha &= O\left(\left(\frac{\mu}{L}\right)^{3/4}\right) \end{aligned} \quad (27)$$

From simple inequality:

$$1 + \frac{1}{3^{1/2}} \sqrt{\frac{\mu}{L}} > \exp\left(\frac{1}{6^{1/2}} \sqrt{\frac{\mu}{L}}\right).$$

We get the following theorem:

**Theorem 6.7.** *If in the model described in (3) in the strongly convex case we can choose  $\alpha$  according to (27) we obtain:*

$$f(x_k) - f(x) \leq \left(\frac{LR^2}{1-\alpha^2} + \frac{3L\alpha^2}{2\mu(1-\alpha^2)} (f(x_0) - f(x))\right) \exp\left(\frac{1}{6^{1/2}} \sqrt{\frac{\mu}{L}}\right).$$

**Corollary 6.8.** *Under the conditions of the theorem, we obtain convergence in the argument:*

$$kx_k - x \leq R^2 \left(\frac{2L}{\mu(1-\alpha^2)} + \frac{3L^2\alpha^2}{4\mu^2(1-\alpha^2)}\right) \exp\left(\frac{1}{6^{1/2}} \sqrt{\frac{\mu}{L}}\right).$$

*Proof.* This is a direct consequence of the inequalities:

$$\begin{aligned} f(x_k) - f(x) &\leq \frac{L}{4} kx_k - x \leq k_2^2 \\ f(x_k) - f(x) &> \frac{\mu}{2} kx_k - x \leq k_2^2. \end{aligned}$$

□

## 7 Conclusions and observations

### Remark 6.

Using Theorem 6.3 and assume, that  $Q$  – compact set we can denote  $R$  as  $\text{diam}(Q)$  instead of  $kx_0 - x_{k_2}$ , then we can also bound  $R \subset R$  and this will simplify bounds in theorem 6.4.

### Remark 7.

With the same assumption  $\mu \notin 0$  we obtain a comparison of the two convergences in the Theorem 6.4. Recall that:

$$\delta_1 = \delta, \quad \delta_2 = \frac{\delta^2}{L}, \quad \delta_3 = \frac{\delta^2}{\mu}.$$

So if

$$\delta < \frac{3R}{\frac{1+\sqrt{L}}{L} + \frac{\sqrt{L}(\rho_{\frac{1}{2}} - 1)}{L}}.$$

Then the accumulation of noise in the model corresponding to  $\tau = 2$ , that described in (3) is less than in model  $\tau = 1$ , described in (2).

### Remark 8.

If we use model  $\tau = 2$ , described in theorem 6.4 one can set the desired accuracy of the solution.

$$f(x_N) - f(x^*) \subset \varepsilon.$$

Then we get from theorem 6.4 that:

$$f(x_N) - f(x^*) \subset LR^2 \exp\left(\frac{1}{2}\sqrt{\frac{\mu}{2L}}N\right) + \left(1 + \sqrt{\frac{L}{\mu_2}}\right)\delta_2 + \left(1 + \sqrt{\frac{L}{\mu_2}}\right)\delta_3,$$

$$f(x_N) - f(x^*) \subset LR^2 \exp\left(\frac{1}{2}\sqrt{\frac{\mu}{2L}}N\right) + \left(\frac{L + \mu}{\sqrt{\mu^3 L}}(\rho_{\frac{1}{2}} - 1)\right)\delta^2.$$

That is we can get estimates for  $\delta$  value and number of steps  $N$ :

$$\left(\frac{L + \mu}{\sqrt{\mu^3 L}}(\rho_{\frac{1}{2}} - 1)\right)\delta^2 \subset \frac{\varepsilon}{2},$$

$$\delta \subset \rho_{\frac{1}{2}} \frac{\varepsilon}{2} \sqrt{\frac{\rho_{\frac{1}{2}} - 1}{2}} \sqrt{\frac{L + \mu}{\sqrt{\mu^3 L}}},$$

$$\delta = O\left(\frac{\rho_{\frac{1}{2}} \varepsilon (L + \mu)^{\frac{1}{2}}}{(\mu^3 L)^{\frac{1}{4}}}\right);$$

$$LR^2 \exp\left(\frac{1}{2}\sqrt{\frac{\mu}{2L}}N\right) \subset \frac{\varepsilon}{2},$$

$$N > 2\sqrt{\frac{2L}{\mu}} (\ln 2LR^2 + \ln \varepsilon^{-1}),$$

$$N = O\left(\sqrt{\frac{L}{\mu}} \ln \frac{LR^2}{\varepsilon}\right).$$

**Remark 9.**

Using remark 4 and previous remark 8, we can found similar bounds. Remind that:

$$f(x_N) - f(x) \leq LR^2 \exp\left(\frac{1}{2}\sqrt{\frac{\mu}{2(L+2)}}N\right) + \left(1 + \sqrt{\frac{2L+4}{\mu}}\right) \left(\frac{1}{L} + \frac{1}{\mu}\right) \delta^2 + \frac{\mu}{2}R^2.$$

However we should value of the parameter  $\mu$ . We will let:

$$\mu = \frac{2}{3} \frac{\varepsilon}{R^2}.$$

Using inequality:

$$\delta^2 \left(1 + \sqrt{\frac{2L+4}{\mu}}\right) \left(\frac{\mu+L}{\mu L}\right) \leq \frac{\varepsilon}{3}.$$

And the selected value of the parameter  $\mu$  we get required value of error  $\delta$ :

$$\delta \leq \left(\frac{2}{243}\right)^{\frac{1}{4}} \frac{1}{\sqrt{1 + \frac{2L+4}{\mu}}} R^{\frac{3}{2}\varepsilon^{\frac{5}{4}}},$$

$$\delta = O\left(L^{\frac{1}{4}} R^{\frac{3}{2}\varepsilon^{\frac{5}{4}}}\right).$$

Similarly, get an estimate of the number of steps:

$$LR^2 \exp\left(\frac{1}{2}\sqrt{\frac{\mu}{2(L+2)}}N\right) \leq \frac{\varepsilon}{3},$$

$$N > \frac{\rho}{12L+24R} \ln 2LR^2 + 2 \frac{\rho}{2L+4} \frac{1}{\varepsilon} \ln \frac{1}{\varepsilon},$$

$$N = O\left(\sqrt{\frac{L}{\varepsilon}} \ln \frac{LR^2}{\varepsilon}\right).$$

**Remark 10.**

Using remark 5 and theorem 6.6 we can apply it to problem:

$$Ax = b,$$

$$A \in \text{GL}_n(\mathbb{R}).$$

Solving such a problem is equivalent to solving the convex optimization problem:

$$f(x) = \frac{1}{2} \|Ax - b\|_2^2 \rightarrow \min,$$

$$r f(x) = A^T (Ax - b).$$

We will assume similarly the estimate of the norm  $x$  :

$$\|x\|_2 \leq R.$$

Let the original problem be solved with an  $\varepsilon_1$  accuracy in the sense:

$$\begin{aligned} \|kAx - bk\|_2 &\leq \varepsilon_1, \\ f(x) - f(x^*) &= \frac{1}{2}kAx - bk\|_2^2 \leq \varepsilon, \\ \varepsilon &= \frac{1}{2}\varepsilon_1^2. \end{aligned}$$

When the algorithm stops, we get the convergence:

$$\begin{aligned} |f(x_{N_{stop}}) - f(x^*)| &\leq N\delta_2 + 3\delta_1 R, \\ N_{stop} &= \left\lceil 2\sqrt{\frac{LR^2}{\zeta}} \right\rceil + 1. \end{aligned}$$

Then we choose  $\delta, \zeta$  from the following conditions:

$$\begin{cases} \zeta \leq \frac{1}{3}, \\ \delta \leq \left(\frac{1}{6\sqrt{3R}}\right)\varepsilon^{\frac{3}{4}}, \\ \delta \leq \varepsilon\frac{1}{9R}. \end{cases}$$

For example, we can let:

$$\begin{aligned} \delta &= C_{R;R;L}\varepsilon, \\ C_{R;R;L} &= \min \left\{ \frac{L^{\frac{1}{4}}}{6\sqrt{3R}}, \frac{1}{9R} \right\}. \end{aligned}$$

Then the number of steps required is expressed as:

$$N'' = \left\lceil 2\sqrt{\frac{3LR^2}{\varepsilon}} \right\rceil + 1.$$

Accordingly, the estimate required for solving the problem of linear equations:

$$N''_1 = \left\lceil 2\frac{\sqrt{3LR^2}}{\varepsilon_1} \right\rceil + 1.$$

### Remark 11.

The work considered a model of additive noise in equation (20), similar to [41], that is we can consider that:

$$\begin{aligned} r f(x) &= r f(x) + r_x, \\ \|r_x\|_2 &\leq \delta. \end{aligned}$$

Similarly to this work, a stopping criterion was proposed for the *STM* algorithm, as was proposed for gradient descent.

$$x_{k+1} = x_k - \frac{1}{L} r f(x_k).$$

Note that in the same noise model, the convergence estimate in both considered cases will be:

$$\begin{aligned}
 j_N &= \operatorname{argmin}_{1 \leq k \leq N} f(x_k), \\
 y_N &= x_{j_N}, \\
 f(y_N) - f(x) &= O\left(\frac{LR^2}{N} + \frac{\delta^2}{L} + R\delta\right), \\
 f(y_N) - f(x) &= O\left(LR^2 \exp\left(\frac{\mu}{L}N\right) + \frac{\delta^2}{L} + R\delta\right), \\
 f(y_N) - f(x) &= O\left(LR^2 \exp\left(\frac{\mu}{2L}\right) + \frac{\delta^2}{L} + \frac{\delta^2}{\mu}\right), \\
 R &= \max_{k \in N} \|x_k - x\|_2.
 \end{aligned}$$

Despite the fact that in the work [17], a slightly different model was considered, namely  $(\delta, L)$  and  $(\delta, L, \mu)$  oracle (equation 3.1 Definition 1 in [17]) similar orders of convergence were obtained, that is theorem 6.4 and relevant remark 3. Namely, function satisfies the  $(\delta, L, \mu)$  model at point  $x \in Q$  means, that exists functions  $f(x)$  and  $\psi(x, y)$ , such that:

$$\begin{aligned}
 & \forall y \in Q \\
 & \frac{\mu}{2} \|x - y\|_2^2 \leq f(x) - f(y) - \psi(x, y) \leq \frac{L}{2} \|x - y\|_2^2 + \delta.
 \end{aligned}$$

Similarly to papers [47], [17], the results also hold in the case of an unbounded set  $Q$  (result in [47] is on the page 26, obtained for fast adaptive gradient method page 13). Stopping criteria are also formulated, which give an estimate on  $R$  for a non-compact  $Q$ , remind that:

$$R = \max_{0 \leq k \leq N} \|x_k - x\|_2, \|x_k - x\|_2, \|z_k - x\|_2.$$

We also note that a similar models of  $(\delta, \cdot, L)$  and  $(\delta, \cdot, L, \mu)$  oracle was considered in the work [48]. Moreover, the function satisfies  $(\delta, \cdot, L, \mu)$ -model

$$\begin{aligned}
 f(y) &\leq f + \psi(y, x) + \|x - y\|_2 + \delta + LV(y, x), \\
 f + \psi(x, x) + \mu V(y, x) &\leq f(x), \\
 f(x) - \delta &\leq f(x) \leq f(x), \\
 \psi(x, x) &= 0.
 \end{aligned}$$

Here  $V(x, y)$  – Bregman divergence. At the same time, an adaptive analogue of STM was considered. As well as similar estimations for a  $\delta$  and number of steps  $N$ , following [17] (page 24, remarks 11 – 14), namely there are remarks 8, 9, 10. Also considered an example of using regularization to obtain convergence in the model  $\tau = 1$ , for the case  $\mu = 0$ .

## 8 Numerical experiments

For testing *STM* for degenerate problems, the function described in [37] on page 69, that is:

$$f(x) = \frac{L}{8} \left( x_1^2 + \sum_{j=0}^{k-1} (x_j - x_{j+1})^2 + x_k^2 \right) - \frac{L}{4} x_1,$$

$$x = \left( 1 - \frac{1}{k+1}, \dots, 1 - \frac{k}{k+1}, 0, \dots, 0 \right)^T,$$

$1 \leq k \leq \dim x.$

These two plots reflect the convergence of the method at the first 50 000 and 10 000 iterations, respectively, at different  $\delta$ .

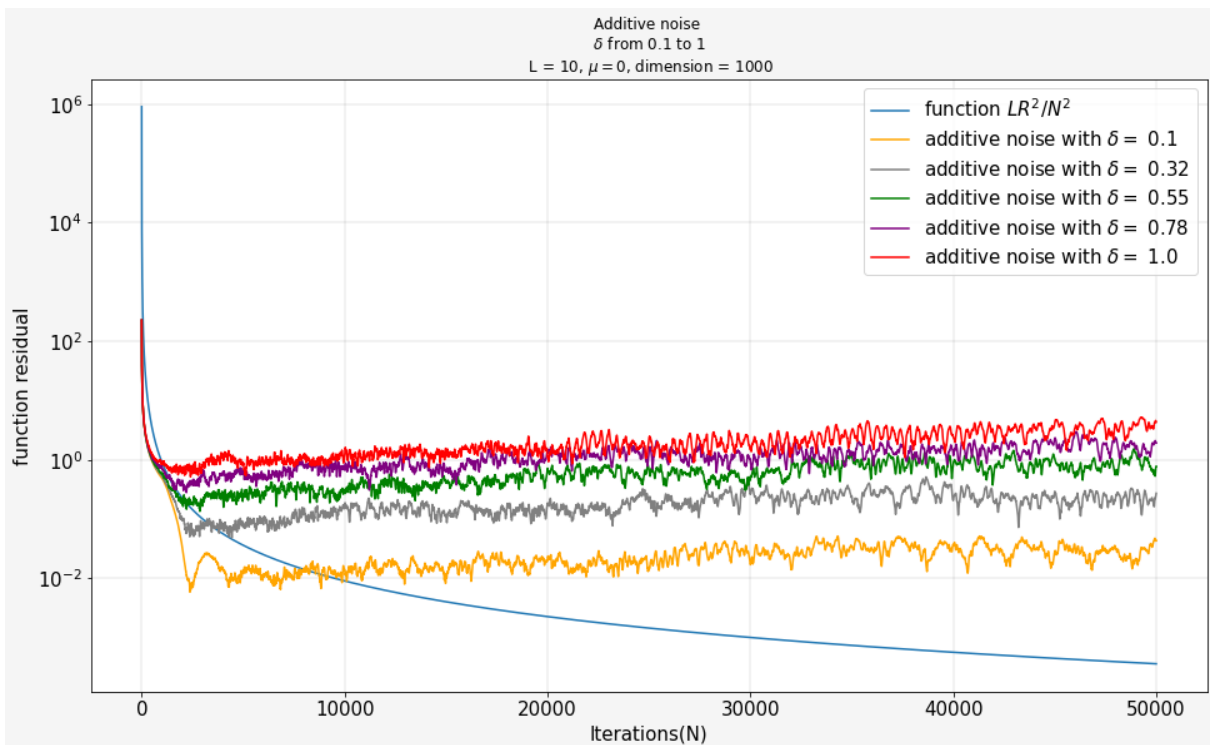


Figure 2: First test – first 50 000 steps.

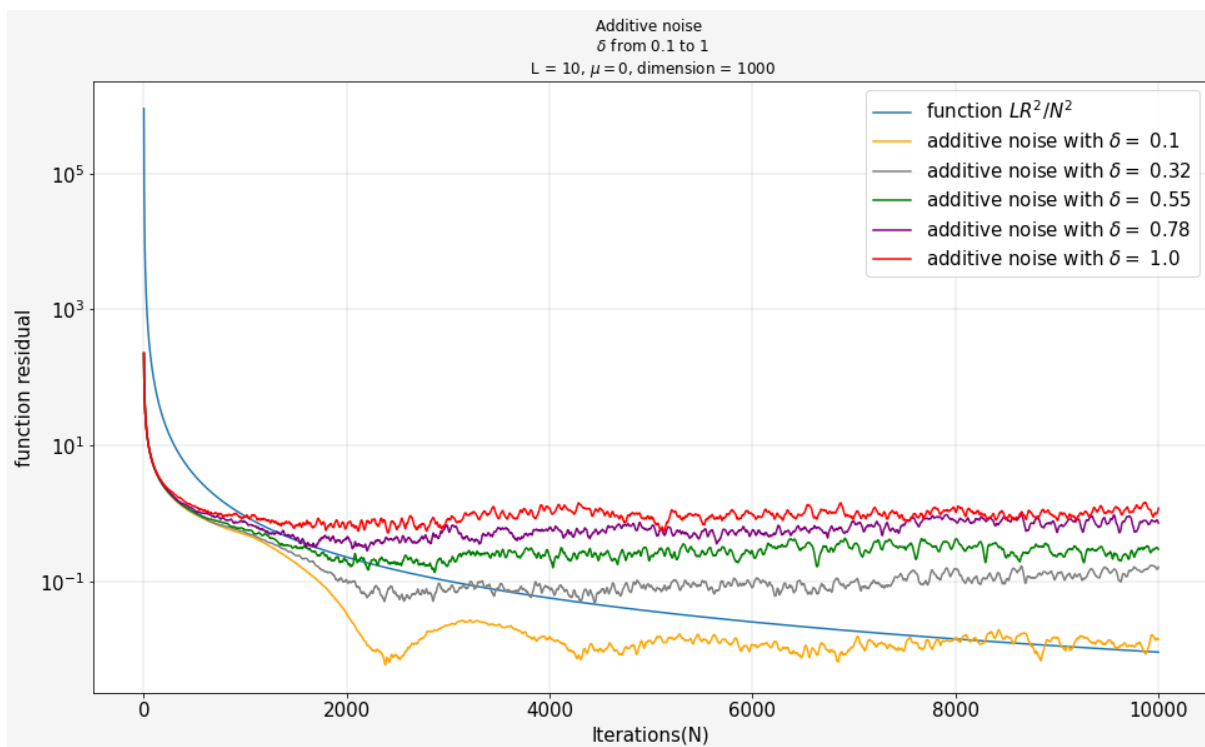


Figure 3: First test – first 10 000 steps.

Let's also consider a drawing with two types of noise.

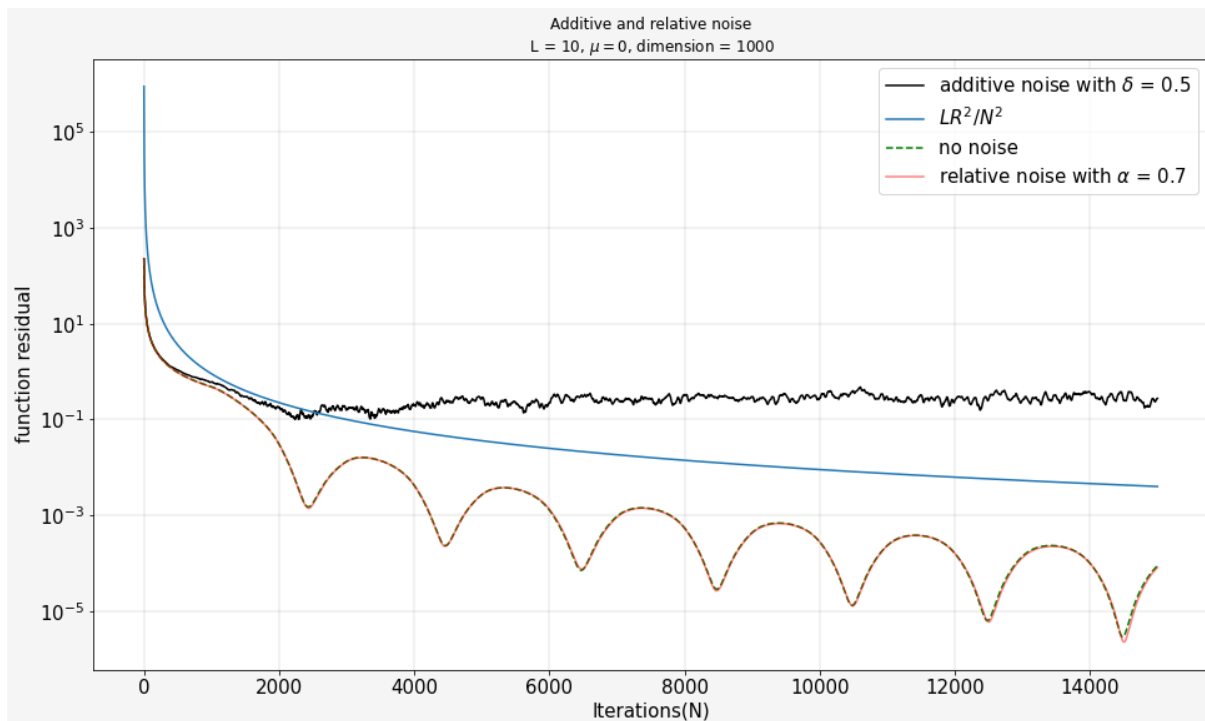


Figure 4: Second test – relative and additive types of noises comparison.



To compare the convergence of a degenerate problem with different  $\alpha$  parameters in the case of relative noise, consider the following graph.

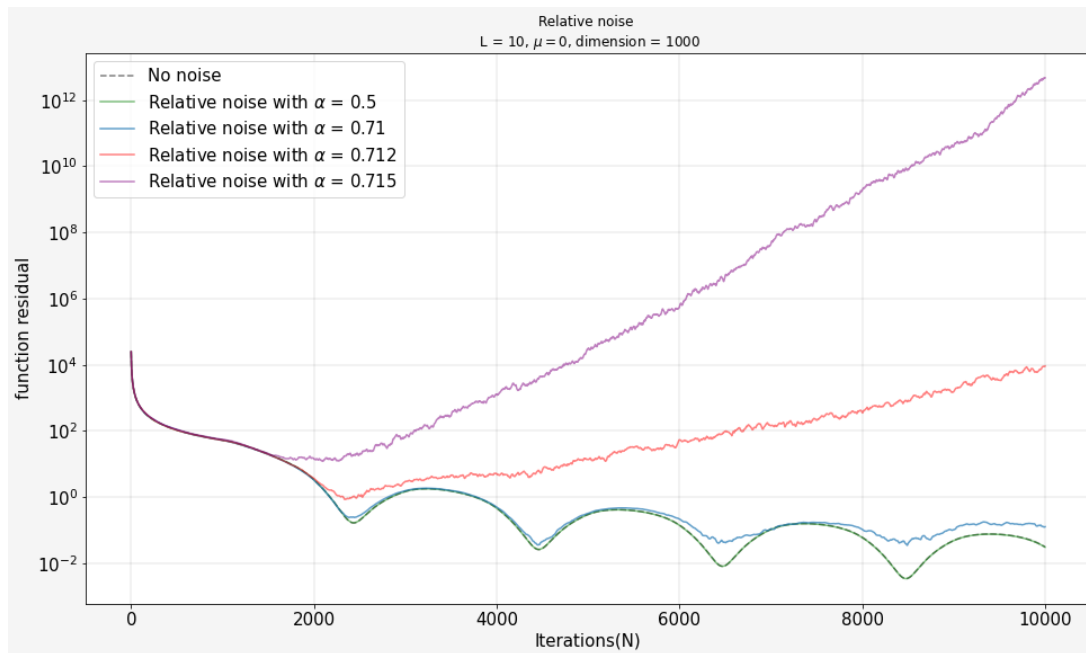


Figure 5: Third test – relative noise with different values of  $\alpha$  for  $\mu = 0$ .

The last figure shows that for  $\alpha \leq 0.71$  the convergence of the method does not deteriorate, but we can assume the existence of such a threshold value  $\alpha = 0.71$ , that at large of  $\alpha$  values the method diverges.

Also for testing on strongly convex functions, an analogue of the finite-dimensional Nesterov function was used from [37] on page 78, that is:

$$f(x) = \frac{\mu(\chi - 1)}{8} \left( x_1^2 + \sum_{j=1}^{n-1} (x_j - x_{j+1})^2 - 2x_1 \right) + \frac{\mu}{2} \|x\|_2^2,$$

$$\chi = \frac{L}{\mu},$$

$$r f(x) = \left( \frac{\mu(\chi - 1)}{4} A + \mu E \right) x - \frac{\mu(\chi - 1)}{4} e_1,$$

$$e_1 = (1, 0, \dots, 0)^T,$$

where  $E$  – identity operator, A is the matrix defined as:

$$\begin{pmatrix} 2 & 1 & 0 & \dots & \dots & 0 \\ 1 & 2 & 1 & 0 & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \ddots & \vdots \\ 0 & \dots & 0 & 1 & 2 & 1 \\ 0 & \dots & \dots & 0 & 1 & 2 \end{pmatrix}.$$

Then minimum  $f, x$ , can be found from systems of linear equations.

Let us consider the graphs of the residuals for different parameters of the delta additive noise.

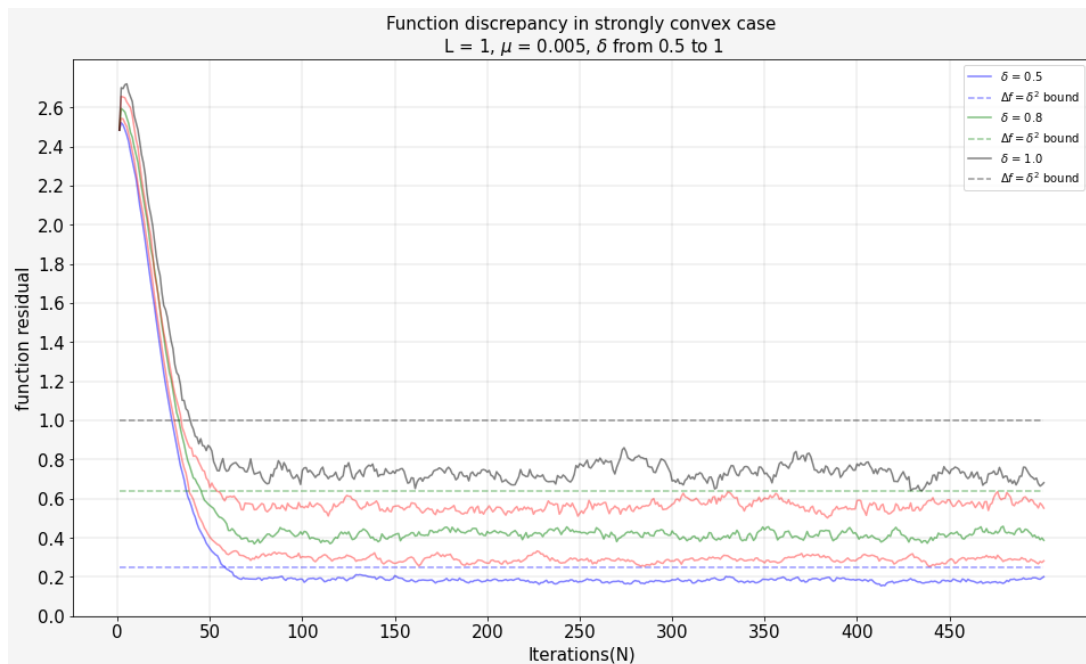


Figure 6: Fourth test –  $\delta \in \{0.5, 0.6, 0.7, 0.8, 0.9, 1.0\}$ .

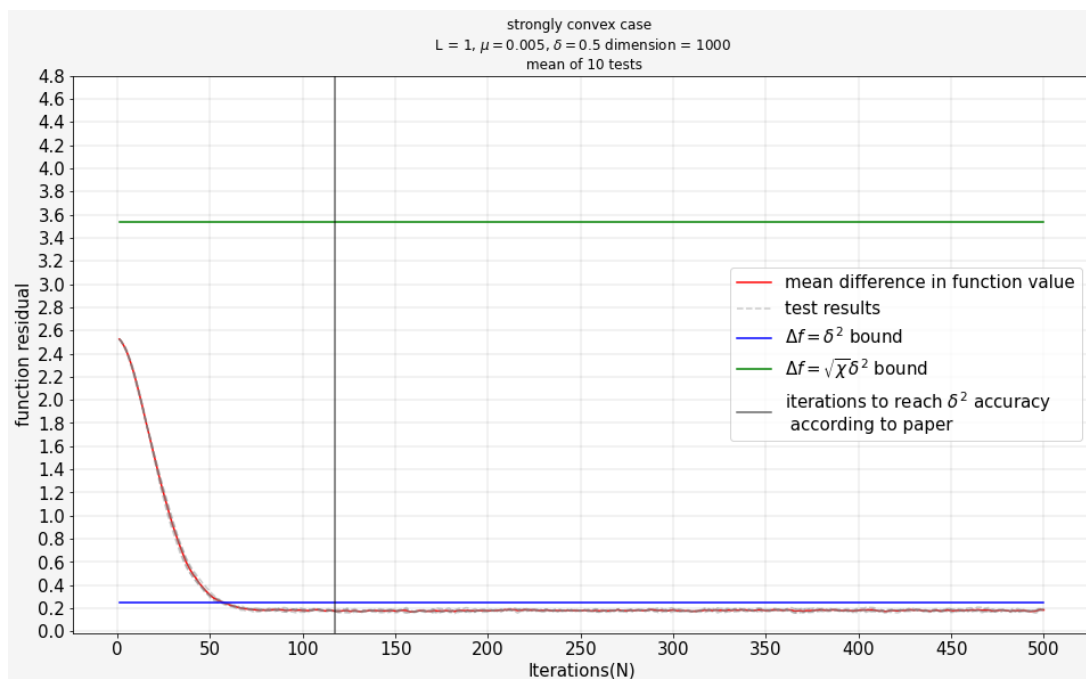


Figure 7: Fifth test – mean of 30 tests, level of approximation and required number of steps.

The last plot confirms theorem 6.4 and remark 8. Similarly to the degenerate case, consider the behavior of the method for different parameters  $\alpha$ .

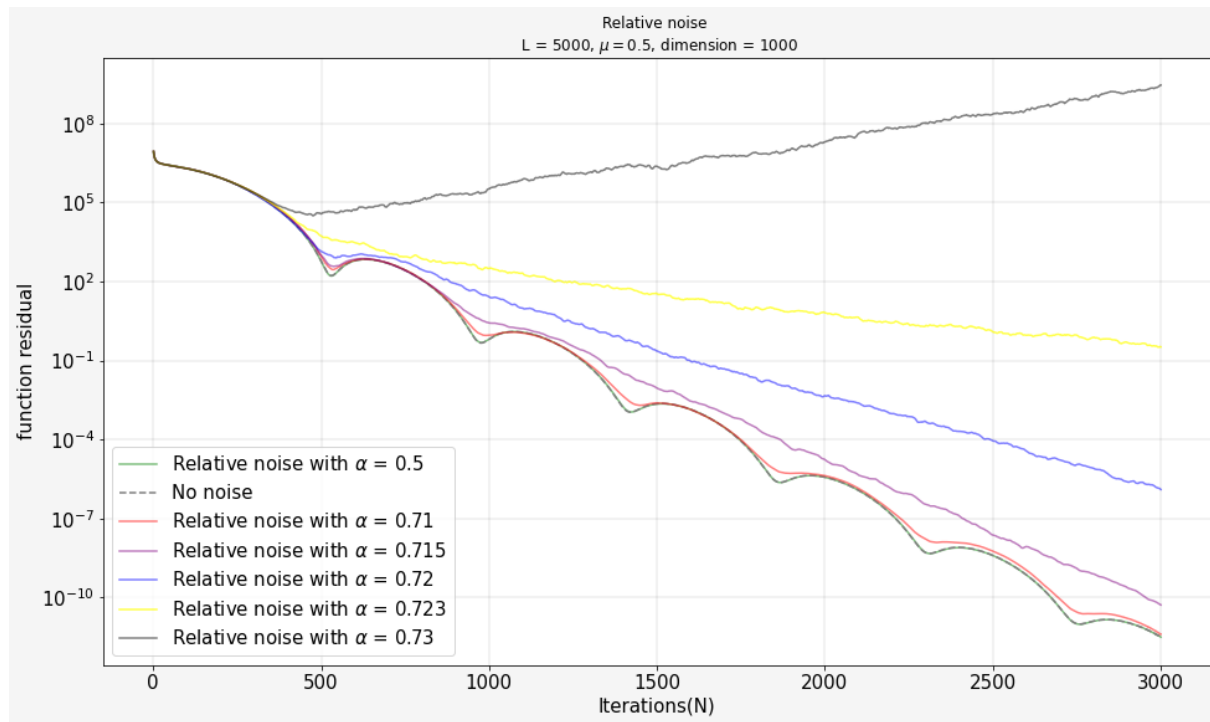


Figure 8: Sixth test – relative noise with different values of  $\alpha$  for  $\mu > 0$ .

Note that in the strongly convex case, we obtain a property similar to the degenerate case: for  $\alpha$  values less than a certain threshold value  $\alpha$ , from the figure we can assume a value of 0.71, the convergence of the method does not deteriorate, and for large  $\alpha$  values, the method diverges.

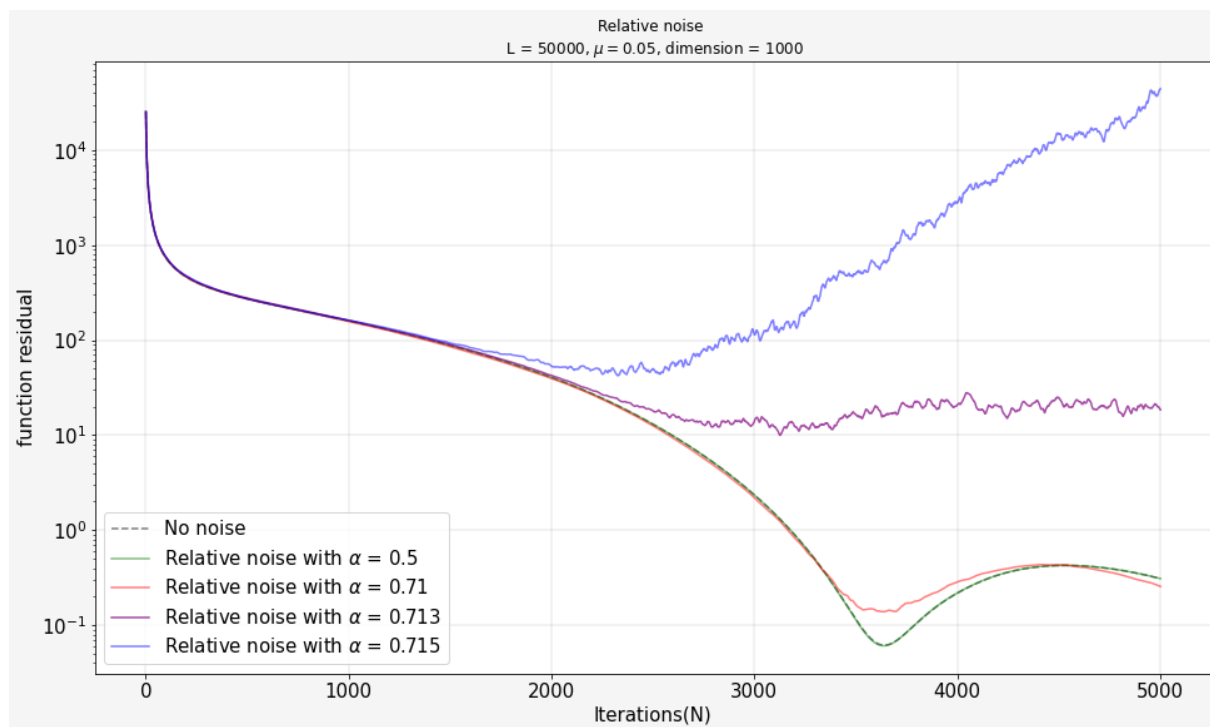


Figure 9: Seventh test – relative noise with different values of  $\alpha$  for other  $L$  and  $\mu$ .

Figure 9 shows, that The figure shows that changing the parameters  $L$  and  $\mu$ , the value of the assumed threshold  $\alpha$  does not change much. We also note that such threshold values turned out to be approximately equal for the degenerate and strongly convex problem.

## 9 Acknowledgment

The authors are grateful to Eduard Gorbunov for useful discussions.

## References

- [1] Ahmad Ajalloeian and Sebastian U Stich. Analysis of sgd with biased gradient estimators. *arXiv preprint arXiv:2008.00051*, 2020.
- [2] Arya Akhavan, Massimiliano Pontil, and Alexandre B Tsybakov. Exploiting higher order smoothness in derivative-free optimization and continuous bandits. *arXiv preprint arXiv:2006.07862*, 2020.
- [3] Francis Bach and Vianney Perchet. Highly-smooth zero-th order online optimization. In Vitaly Feldman, Alexander Rakhlin, and Ohad Shamir, editors, *29th Annual Conference on Learning Theory*, volume 49 of *Proceedings of Machine Learning Research*, pages 257–283, Columbia University, New York, New York, USA, 23–26 Jun 2016. PMLR.
- [4] Amir Beck. *First-order methods in optimization*. SIAM, 2017.
- [5] Alexandre Belloni, Tengyuan Liang, Hariharan Narayanan, and Alexander Rakhlin. Escaping the local minima via simulated annealing: Optimization of approximately convex functions. In Peter Grünwald, Elad Hazan, and Satyen Kale, editors, *Proceedings of The 28th Conference on Learning Theory*, volume 40 of *Proceedings of Machine Learning Research*, pages 240–265, Paris, France, 03–06 Jul 2015. PMLR.
- [6] Aaron Ben-Tal and Arkadi Nemirovski. *Lectures on Modern Convex Optimization (Lecture Notes)*. Personal web-page of A. Nemirovski, 2015.
- [7] Albert S Berahas, Liyuan Cao, Krzysztof Choromanski, and Katya Scheinberg. A theoretical and empirical comparison of gradient approximations in derivative-free optimization. *arXiv preprint arXiv:1905.01332*, 2019.
- [8] Aleksandr Beznosikov, Abdurakhmon Sadiev, and Alexander Gasnikov. Gradient-free methods with inexact oracle for convex-concave stochastic saddle-point problem. In *International Conference on Mathematical Optimization Theory and Operations Research*, pages 105–119. Springer, 2020.
- [9] Sébastien Bubeck. Convex optimization: Algorithms and complexity. *arXiv preprint arXiv:1405.4980*, 2014.
- [10] Michael B Cohen, Jelena Diakonikolas, and Lorenzo Orecchia. On acceleration with noise-corrupted gradients. *arXiv preprint arXiv:1805.12591*, 2018.
- [11] A. Conn, K. Scheinberg, and L. Vicente. *Introduction to Derivative-Free Optimization*. Society for Industrial and Applied Mathematics, 2009.
- [12] Alexandre d’Aspremont. Smooth optimization with approximate gradient. *SIAM Journal on Optimization*, 19(3):1171–1183, 2008.
- [13] Alexandre d’Aspremont, Damien Scieur, and Adrien Taylor. Acceleration methods. *arXiv preprint arXiv:2001.09545*, 2021.
- [14] Olivier Devolder. Stochastic first order methods in smooth convex optimization. *CORE Discussion Paper 2011/70*, 2011.
- [15] Olivier Devolder. *Exactness, inexactness and stochasticity in first-order methods for large-scale convex optimization*. PhD thesis, ICTEAM and CORE, Université Catholique de Louvain, 2013.

- [16] Olivier Devolder, François Glineur, and Yurii Nesterov. First-order methods of smooth convex optimization with inexact oracle. *Mathematical Programming*, 146(1):37–75, 2014.
- [17] Olivier Devolder, François Glineur, Yurii Nesterov, et al. First-order methods with inexact oracle: the strongly convex case. *CORE Discussion Papers*, 2013016:47, 2013.
- [18] Dmitriy Drusvyatskiy and Lin Xiao. Stochastic optimization with decision-dependent distributions. *arXiv preprint arXiv:2011.11173*, 2020.
- [19] Darina Dvinskikh and Alexander Gasnikov. Decentralized and parallelized primal and dual accelerated methods for stochastic convex programming problems. *Journal of Inverse and Ill-posed Problems*, 2021.
- [20] Darina M Dvinskikh, Aleksandr Igorevich Turin, Alexander Vladimirovich Gasnikov, and Sergey Sergeevich Omelchenko. Accelerated and non accelerated stochastic gradient descent in model generality. *Matematicheskie Zametki*, 108(4):515–528, 2020.
- [21] Pavel Dvurechensky. Numerical methods in large-scale optimization: inexact oracle and primal-dual analysis. *HSE. Habilitation*, 2020.
- [22] Pavel Dvurechensky and Alexander Gasnikov. Stochastic intermediate gradient method for convex problems with stochastic inexact oracle. *Journal of Optimization Theory and Applications*, 171(1):121–145, 2016.
- [23] Pavel Dvurechensky, Mathias Staudigl, and Shimrit Shtern. First-order methods for convex optimization. *arXiv preprint arXiv:2101.00935*, 2021.
- [24] Yu G Evtushenko. Optimization and fast automatic differentiation. *Computing Center of RAS, Moscow*, 2013.
- [25] A. V. Gasnikov, E. V. Gasnikova, Yu. E. Nesterov, and A. V. Chernov. Efficient numerical methods for entropy-linear programming problems. *Computational Mathematics and Mathematical Physics*, 56(4):514–524, 2016.
- [26] Alexander Gasnikov. Universal gradient descent. *arXiv preprint arXiv:1711.00394*, 2017.
- [27] Alexander Gasnikov, Sergey Kabanikhin, Ahmed Mohammed, and Maxim Shishlenin. Convex optimization in hilbert space with applications to inverse problems. *arXiv preprint arXiv:1703.00267*, 2017.
- [28] Alexander Vladimirovich Gasnikov and Yu E Nesterov. Universal method for stochastic composite optimization problems. *Computational Mathematics and Mathematical Physics*, 58(1):48–64, 2018.
- [29] Ian Goodfellow, Yoshua Bengio, Aaron Courville, and Yoshua Bengio. *Deep learning*, volume 1. MIT press Cambridge, 2016.
- [30] Eduard Gorbunov, Darina Dvinskikh, and Alexander Gasnikov. Optimal decentralized distributed algorithms for stochastic convex optimization. *arXiv preprint arXiv:1911.07363*, 2019.
- [31] Sergey I Kabanikhin. *Inverse and ill-posed problems: theory and applications*, volume 55. Walter De Gruyter, 2011.

- [32] Dmitry Kamzolov, Pavel Dvurechensky, and Alexander V Gasnikov. Universal intermediate gradient method for convex problems with inexact oracle. *Optimization Methods and Software*, pages 1–28, 2020.
- [33] Georgios Kotsalis, Guanghui Lan, and Tianjiao Li. Simple and optimal methods for stochastic variational inequalities, ii: Markovian noise and policy evaluation in reinforcement learning. *arXiv preprint arXiv:2011.08434*, 2020.
- [34] Guanghui Lan. *First-order and Stochastic Optimization Methods for Machine Learning*. Springer, 2020.
- [35] Arkadi S Nemirovski. Regularizing properties of the conjugate gradient method for ill-posed problems. *Zhurnal Vychislitel'noi Matematiki i Matematicheskoi Fiziki*, 26(3):332–347, 1986.
- [36] A.S. Nemirovsky and D.B. Yudin. *Problem Complexity and Method Efficiency in Optimization*. J. Wiley & Sons, New York, 1983.
- [37] Yurii Nesterov. *Lectures on convex optimization*, volume 137. Springer, 2018.
- [38] Yurii Nesterov and Vladimir Spokoiny. Random gradient-free minimization of convex functions. *Found. Comput. Math.*, 17(2):527–566, April 2017. First appeared in 2011 as CORE discussion paper 2011/16.
- [39] Vasillii Novitskii and Alexander Gasnikov. Improved exploiting higher order smoothness in derivative-free optimization and continuous bandit. *arXiv preprint arXiv:2101.03821*, 2021.
- [40] Fabian Pedregosa and Damien Scieur. Average-case acceleration through spectral density estimation. *arXiv preprint arXiv:2002.04756*, 2020.
- [41] BT Poljak. Iterative algorithms for singular minimization problems. In *Nonlinear Programming 4*, pages 147–166. Elsevier, 1981.
- [42] Boris Polyak. *Introduction to Optimization*. New York, Optimization Software, 1987.
- [43] Boris Teodorovich Polyak and Aleksandr Borisovich Tsybakov. Optimal order of accuracy of search algorithms in stochastic optimization. *Problemy Peredachi Informatsii*, 26(2):45–53, 1990.
- [44] Andrej Risteski and Yuanzhi Li. Algorithms and matching lower bounds for approximately-convex optimization. *Advances in Neural Information Processing Systems*, 29:4745–4753, 2016.
- [45] R Tyrrell Rockafellar. *Convex analysis*, volume 36. Princeton university press, 1970.
- [46] Damien Scieur and Fabian Pedregosa. Universal asymptotic optimality of polyak momentum. In *International Conference on Machine Learning*, pages 8565–8572. PMLR, 2020.
- [47] Fedor Stonyakin, Alexander Tyurin, Alexander Gasnikov, Pavel Dvurechensky, Artem Agafonov, Darina Dvinskikh, Dmitry Pasechnyuk, Sergei Artamonov, and Victorya Piskunova. Inexact relative smoothness and strong convexity for optimization and variational inequalities by inexact model. *arXiv preprint arXiv:2001.09013*, 2020.
- [48] Fedor Stonyakin. Adaptive methods for variational inequalities, minimization problems and functional with generalized growth condition. *MIPT. Habilitation*, 2020.

- [49] Adrien B Taylor, Julien M Hendrickx, and François Glineur. Smooth strongly convex interpolation and exact worst-case performance of first-order methods. *Mathematical Programming*, 161(1-2):307–345, 2017.
- [50] Alexander Tyurin. Development of a method for solving structural optimization problems. *HSE. PhD Thesis*, 2020.
- [51] F. Vasilyev. *Optimization Methods*. Moscow, Russia: FP, 2002.