

---

# STOPPING RULES FOR GRADIENT METHODS FOR NON-CONVEX PROBLEMS WITH ADDITIVE NOISE IN GRADIENT

---

A PREPRINT

**Fedor Stonyakin**

Moscow Institute of Physics and Technology  
Moscow, Russia  
V. I. Vernadsky Crimean Federal University  
Simferopol, Russia  
fedyor@mail.ru

**Ilya Kuruzov**

Moscow Institute of Physics and Technology  
Moscow, Russia  
kuruzov.ia@phystech.edu

**Boris Polyak**

Moscow Institute of Physics and Technology  
Moscow, Russia  
Institute for Control Sciences  
Moscow, Russia  
boris@ipu.ru

December 13, 2022

## ABSTRACT

We study the gradient method under the assumption that an additively inexact gradient is available for, generally speaking, non-convex problems. The non-convexity of the objective function, as well as the use of an inexactness specified gradient at iterations, can lead to various problems. For example, the trajectory of the gradient method may be far enough away from the starting point. On the other hand, the unbounded removal of the trajectory of the gradient method in the presence of noise can lead to the removal of the trajectory of the method from the desired exact solution. The results of investigating the behavior of the trajectory of the gradient method are obtained under the assumption of the inexactness of the gradient and the condition of gradient dominance. It is well known that such a condition is valid for many important non-convex problems. Moreover, it leads to good complexity guarantees for the gradient method. A rule of early stopping of the gradient method is proposed. Firstly, it guarantees achieving an acceptable quality of the exit point of the method in terms of the function. Secondly, the stopping rule ensures a fairly moderate distance of this point from the chosen initial position. In addition to the gradient method with a constant step, its variant with adaptive step size is also investigated in detail, which makes it possible to apply the developed technique in the case of an unknown Lipschitz constant for the gradient. Some computational experiments have been carried out which demonstrate effectiveness of the proposed stopping rule for the investigated gradient methods.

**Keywords:** Non-convex Optimization · Polyak-Łojasiewicz Condition · Inexact Gradient · Stopping Rule · Adaptive Method

**Mathematics Subject Classification (2000):** 49M37 · 90C25 · 65K05

## 1. Introduction

Gradient methods are relatively simple, and they require a low iteration cost as well as a small amount of memory, which explains their popularity. In data analysis, non-convex problems often arise under the standard assumption that the gradient of the objective function  $f$  is Lipschitz-continuous with some constant  $L > 0$  (or in other words, the function  $f$  is  $L$ -smooth):

$$\|\nabla f(x) - \nabla f(y)\| \leq L\|x - y\| \quad \forall x, y \in \mathbb{R}^n, \quad (1)$$

where  $\|\cdot\|$  (here and everywhere in the paper) denotes the Euclidean norm. For these problems, by applying the gradient-type methods, the generated sub-sequence of points, generated by applying the gradient-type methods, converges to the zero value of the  $\|\nabla f(x)\|$ . For this fact we have the following known result (see, for example, [6, 14]).

**Theorem 1.1.** *Let  $f$  be an  $L$ -smooth function. Let us consider the gradient method*

$$x_{k+1} = x_k - \frac{1}{L}\nabla f(x_k) \quad (2)$$

for the following optimization problem:

$$\min_{x \in \mathbb{R}^n} f(x). \quad (3)$$

Then the following inequality holds:

$$\min_{k=0, \dots, N-1} \|\nabla f(x_k)\| \leq \sqrt{\frac{2L(f(x_0) - f(x_*))}{N}}, \quad (4)$$

where  $x_0$  is a starting point of the method and  $x_*$  is one of the exact solutions of the problem (3).

Let  $f$  be an  $L$ -smooth function and its gradient satisfy the Polyak-Łojasiewicz condition (for brevity, we will write PL-condition) for some constant  $\mu > 0$  [13] (see also the recent papers [8, 1], and the references therein):

$$f(x) - f^* \leq \frac{1}{2\mu} \|\nabla f(x)\|^2 \quad \forall x \in \mathbb{R}^n, \quad (5)$$

where  $f^* = f(x_*)$  is the value of the function  $f$  at one of the exact solutions  $x_*$  of the optimization problem under consideration. Then the Gradient Descent Method converges at the rate of a geometric progression

$$f(x_N) - f^* \leq \left(1 - \frac{\mu}{L}\right)^N (f(x_0) - f^*) \leq \exp\left(-\frac{\mu}{L}N\right) (f(x_0) - f^*), \quad (6)$$

$$\|x_* - x_0\| \leq \frac{\sqrt{2L(f(x_0) - f^*)}}{\mu}. \quad (7)$$

From [8] it is known that the PL-condition (5) implies the following so-called quadratic growth condition:

$$f(x) - f^* \geq \frac{\mu}{2} \inf_{x_*} \|x - x_*\|^2 \quad \forall x \in \mathbb{R}^n,$$

whence one can obtain that (6) means that the Gradient Descent method also converges in argument at the rate of a geometric progression

$$\inf_{x_*} \|x_N - x_*\|^2 \leq \frac{2}{\mu} \exp\left(-\frac{\mu}{L}N\right) (f(x_0) - f^*).$$

It is worth noting that the gradient dominance condition (5) is certainly holds for a strongly convex objective function  $f$ . However, there are known examples, where PL-condition holds, but one cannot be sure even that  $f$  is convex (see, for example, [12]). So from [6], we can consider the problem of finding some solution to a system of nonlinear equations  $g(x) = 0$  (written in a vector form), where  $g : \mathbb{R}^n \rightarrow \mathbb{R}^m$ ,  $m \leq n$  the problem of finding some solution to this system.

Let us introduce the Jacobian matrix  $J(x) = \frac{\partial g(x)}{\partial x} = \left\| \frac{\partial g_i(x)}{\partial x_j} \right\|_{i,j=1}^{m,n}$  of the mapping  $g$  and assume that there exists  $\mu > 0$  such that for all  $x \in \mathbb{R}^n$  the Jacobian matrix is uniformly non-singular, i.e.  $\lambda_{\min} \left( J(x) [J(x)]^\top \right) \geq \mu$ . In this case, the function  $f(x) = \|g(x)\|^2$  satisfies condition (5) for an arbitrary  $x_*$  such that  $f(x_*) = 0$ , i.e.  $g(x_*) = 0$  [10]. We would like to mention separately the review [1], which describes in detail a deep learning-motivated example of a non-linear equation-related minimization problem with over-parametrization for a non-convex smooth function with PL-condition.

### 1.1. The formulation of the problem

In this paper, we consider the problem of minimizing the function  $f$  which satisfies PL-condition (5) and has  $L$ -Lipschitz continuous gradient with some constant  $L > 0$

$$\|\nabla f(x) - \nabla f(y)\| \leq L\|x - y\| \quad \forall x, y \in \mathbb{R}^n. \quad (8)$$

We suppose that the method has access not to the exact, but to the approximate value of the gradient  $\tilde{\nabla}f(x)$  at any requested point  $x$ , which means the following

$$\nabla f(x) = \tilde{\nabla}f(x) + v(x), \quad \text{and} \quad \|v(x)\| \leq \Delta \quad (9)$$

for some fixed  $\Delta > 0$ . Then (5) means that

$$f(x) - f^* \leq \frac{1}{\mu} (\|\tilde{\nabla}f(x)\|^2 + \Delta^2) \quad \forall x \in \mathbb{R}^n. \quad (10)$$

Therefore,  $\|\tilde{\nabla}f(x)\|^2 + \Delta^2 \geq \mu(f(x) - f^*)$ , where

$$\|\tilde{\nabla}f(x)\|^2 \geq \mu(f(x) - f^*) - \Delta^2 \quad \forall x \in \mathbb{R}^n. \quad (11)$$

It is worth noting that the issue of studying the influence of gradient errors on the estimates of the convergence rate of the first-order methods attracted many researchers (see, for example, [14, 4, 3, 2, 16]). However, we will focus on the distinguished class of non-convex problems. The non-convexity of the objective function of the problem, as well as the use of an inexactness of the specified gradient at iterations, can lead to various problems. In particular, in the absence of any early stopping rules, divergence of the gradient method trajectory from the starting point can be quite a large. It is problematic when the initial point of the method already has some appropriate properties. On the other hand, the unlimited divergence of the trajectory of the Gradient Descent method in the stochastic setting can lead to a larger distance from the desired exact solution. Let us describe some situations of this type.

As a simple example of a non-strongly convex function that satisfies the gradient dominance condition, we consider

$$f(x) = \langle Ax, x \rangle, \quad (12)$$

where  $A = \text{diag}(L, \mu, 0)$  is a 3-order diagonal matrix with exactly two positive entries  $L > \mu > 0$ . If for the problem of minimizing the function (12) we assume that there is a gradient error  $v(x) = (0, 0, \Delta)$  for  $\Delta > 0$ , then for  $x_0 = (0, 0, 0)$ ,  $h_k > 0$  and  $x_{k+1} = x_k - h_k \tilde{\nabla}f(x_k)$ , we have  $\lim_{k \rightarrow \infty} \|x_{k+1}\|_2 = \infty$ .

Further, we can consider the Rosenbrock function of two variables  $x = (x^{(1)}, x^{(2)})$ :

$$f(x) = 100 \left( x^{(2)} - (x^{(1)})^2 \right)^2 + \left( 1 - x^{(1)} \right)^2.$$

Let our method starts from  $x_0 = (1, 1) = x_*$ . Then at each step of the gradient method, the error of the gradient  $v(x_k)$  is such that  $x_k^{(2)} = (x_k^{(1)})^2$  and without stopping rule the trajectory can go very far from the exact solution  $x_*$ . Similarly, the trajectory of the gradient method can be unbounded for the objective function of two variables  $f(x) = (x^{(2)} - (x^{(1)})^2)^2$ .

The purpose of this paper is to study the estimate of the distance  $\|x_N - x_0\|$  for points  $x_N$  produced by the Gradient Descent method and to propose an early stopping rule that guarantees some compromise, such as a significant divergence of the trajectory from the chosen starting point of the method. Note that the early stopping rules in iterative procedures are being actively studied for various types of problems. Apparently, for the first time, the ideology of early stopping of iterations was proposed in [5]. This paper is devoted to a technique for the approximate solution of ill-posed or ill-conditioned problems arising during regularization (in the mentioned work, the authors considered the problem of solving a linear equation). In this case, an early stop is aimed at overcoming the problem of the potential accumulation of errors in the regularization of the original problem. The topic of our paper is related to well-known approaches related to the early termination of first-order methods in the case of using inexact information about the gradient at iterations (see [14], Ch. 6, paragraph 1, and also, for example, the recent preprint [16]). However, the results known to us for convex (not strongly convex) problems differ from those obtained in this note. The main difference is that usually either the achievement of the worst level in function is guaranteed (compared with the comment after theorem 2, section 1, chapter 6 of [14]) or estimates such as  $\|x_N - x_*\| \leq \|x_0 - x_*\|$  without examining  $\|x_N - x_0\|$ .

Here  $\{x_k\}_{k \in \mathbb{N}}$  is the sequence generated by the method,  $x_*$  is the exact solution of the minimization problem closest to the starting point of the method  $x_0$ .

In this paper, we obtained Theorem 2.2 devoted to the Gradient Descent method with a constant step-size with a sufficiently small value of the inexact gradient. It indicates the level of accuracy with respect to the function that can be guaranteed after the proposed early stopping rule is fulfilled. It is important to note that this result can be applied to any  $L$ -smooth non-convex problem. Further, using PL-condition, this result is refined in Theorem 2.3, which describes the estimate of a sufficient number of iterations to achieve the desired quality of the output point  $\hat{x}$  by the function  $f(\hat{x}) - f^* = O\left(\frac{\Delta^2}{\mu}\right)$ . Moreover, it contains an estimate (26) of the distance from  $\hat{x}$  to the starting point  $x_0$ . The obtained results are compared with the well-known distance estimate [13] from the starting point  $x_0$  to the nearest exact solution  $x_*$  (see remark 2.5).

However, the method with a constant step-size imposes the need to efficiently estimate the Lipschitz constant of the gradient of the objective function, which can be problematic in practice. Moreover, many real problems lead to functions that have not an Lipschitz continuous gradient, and a condition such as (1) holds for such functions only locally on some subset  $Q \subset \mathbb{R}^n$ . Therefore, we propose variations of Theorems 2.2 and 2.3 for the Gradient Descent method with an adaptively selected step-size. This makes it possible to apply an analog of Theorem 2.2 with the early stopping rule (31) to an arbitrary non-convex problem without additional conditions. If PL-condition is guaranteed, then the execution of (31) automatically guarantees the achievement of an acceptable quality level for the solution of the problem of minimizing  $f$  by the function.

The last section of the paper is devoted to numerical experiments which explain the purpose of using stopping rule (21) for some specific examples of object functions in problems: logistic regression, Rosenbrock and Nesterov-Skokov functions, quadratic function.

## 2. The Proposed Approach and Main Theoretical Results

### 2.1. Variant of the Gradient Descent method with a constant step-size

We assume that the values of the parameters  $L > 0$  and  $\Delta > 0$  are known. Also, the Gradient Descent method of the following form

$$x_{k+1} = x_k - \frac{1}{L} \tilde{\nabla} f(x_k) \quad (13)$$

can be applied to solve the minimization problem of the function  $f$ . In view of (8) for the method (13), we get

$$\begin{aligned} f(x_{k+1}) &\leq f(x_k) + \langle \nabla f(x_k), x_{k+1} - x_k \rangle + \frac{L}{2} \|x_{k+1} - x_k\|^2 \\ &= f(x_k) - \frac{1}{L} \langle \nabla f(x_k), \tilde{\nabla} f(x_k) \rangle + \frac{1}{2L} \|\tilde{\nabla} f(x_k)\|^2 \\ &= f(x_k) + \frac{1}{2L} \left( \|\nabla f(x_k)\|^2 - 2 \langle \nabla f(x_k), \tilde{\nabla} f(x_k) \rangle + \|\tilde{\nabla} f(x_k)\|^2 \right) - \frac{\|\nabla f(x_k)\|^2}{2L} \\ &= f(x_k) + \frac{1}{2L} \|\nabla f(x_k) - \tilde{\nabla} f(x_k)\|^2 - \frac{\|\nabla f(x_k)\|^2}{2L} \\ &\leq f(x_k) + \frac{\Delta^2}{2L} - \frac{1}{2L} \|\nabla f(x_k)\|^2, \end{aligned}$$

i.e.

$$f(x_{k+1}) - f(x_k) \leq \frac{\Delta^2}{2L} - \frac{1}{2L} \|\nabla f(x_k)\|^2. \quad (14)$$

Summing up inequalities (14) over  $k = \overline{0, N-1}$  leads us to an estimate

$$\min_{k=0, \dots, N-1} \|\nabla f(x_k)\| \leq \sqrt{\Delta^2 + \frac{2L(f(x_0) - f(x_*))}{N}} \leq \Delta + \sqrt{\frac{2L(f(x_0) - f(x_*))}{N}}. \quad (15)$$

Note that, in contrast to (4), the estimate (15) points to the potential divergence of the Gradient Descent method in the case of an additively inexact gradient. Specific examples of such situations were described above.

Taking into account (5), we get

$$f(x_{k+1}) - f(x_k) \leq \frac{\Delta^2}{2L} - \frac{2\mu(f(x_k) - f^*)}{2L} = -\frac{\mu}{L}(f(x_k) - f^*) + \frac{\Delta^2}{2L},$$

thus

$$\begin{aligned}
f(x_{k+1}) - f^* &\leq \left(1 - \frac{\mu}{L}\right) (f(x_k) - f^*) + \frac{\Delta^2}{2L} \\
&\leq \left(1 - \frac{\mu}{L}\right)^{k+1} (f(x_0) - f^*) + \frac{\Delta^2}{2L} \left(1 + 1 - \frac{\mu}{L} + \dots + \left(1 - \frac{\mu}{L}\right)^k\right) \\
&< \left(1 - \frac{\mu}{L}\right)^{k+1} (f(x_0) - f^*) + \frac{\Delta^2}{2\mu},
\end{aligned}$$

i.e.

$$f(x_{k+1}) - f^* \leq \left(1 - \frac{\mu}{L}\right)^{k+1} (f(x_0) - f^*) + \frac{\Delta^2}{2\mu}. \quad (16)$$

**Remark 2.1.** *It is important to note that bounds (15) and (16) cannot be improved for the Gradient Descent method with an additively inexact gradient in the general case. For example, the lower estimates of accuracy with respect to the function  $O\left(\frac{\Delta^2}{2\mu}\right)$  are known even on the class of strongly convex functions (see, for example, section 2.11.1 of the manual [17], as well as references therein). In this regard, we consider the following example:*

$$\min_{x \in \mathbb{R}^n} f(x) := \frac{1}{2} \sum_{i=1}^n \lambda_i \left(x^{(i)}\right)^2, \quad (17)$$

where  $0 \leq \mu = \lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_n = L$ , such that  $L \geq 2\mu$ . The exact solution of problem (17) is  $x_* = 0 \in \mathbb{R}^n$ . Suppose that an inexact gradient is available at the current point of the feasible area. Besides, the error is only in the calculation of the first component of the gradient. That is, instead of  $\partial f(x)/\partial x^{(1)} = \mu x^{(1)}$ , we have only  $\tilde{\partial} f(x)/\partial x^{(1)} = \mu x^{(1)} - \Delta$ , for some  $\Delta > 0$ . Then for the simplest Gradient Descent method (13) one can obtain that for  $x_0^{(1)} \geq 0$  and sufficiently large  $k \in \mathbb{N}$  ( $k \gg L/\mu$ ) the following inequality holds:

$$x_k^{(1)} \geq \frac{\Delta}{L} \frac{1 - (1 - \mu/L)^k}{1 - (1 - \mu/L)} \simeq \frac{\Delta}{\mu}. \quad (18)$$

Therefore,  $f(x_k) - f(x_*) \gtrsim \frac{\Delta^2}{2\mu}$ .

Further, in view of

$$\|\nabla f(x_k)\|^2 \geq \frac{\|\tilde{\nabla} f(x_k)\|^2}{2} - \Delta^2,$$

from (14), the inexact gradient satisfies the following inequality:

$$f(x_{k+1}) - f(x_k) \leq \frac{\Delta^2}{2L} - \frac{1}{2L} \left( \frac{\|\tilde{\nabla} f(x_k)\|^2}{2} - \Delta^2 \right),$$

whence we have

$$f(x_{k+1}) - f(x_k) \leq \frac{\Delta^2}{L} - \frac{1}{4L} \|\tilde{\nabla} f(x_k)\|^2. \quad (19)$$

Inequality (19) shows that if the value  $\|\tilde{\nabla} f(x_k)\|$  is sufficiently large, it can be guaranteed that  $f(x_{k+1}) < f(x_k)$ . Thus, for any  $C > 2$ , an alternative arises: either the inequality  $\|\tilde{\nabla} f(x_k)\| \leq C\Delta$  holds, or

$$f(x_{k+1}) - f(x_k) < -\frac{\Delta^2}{L} \left( \frac{C^2}{4} - 1 \right).$$

In the first case, the inequality  $\|\tilde{\nabla} f(x_k)\| \leq C\Delta$  guarantees the achievement of an acceptable quality of the output point  $x_k$  with respect to the function due to PL-condition. In the second case, we can guarantee the decreasing with respect to the function for  $C > 2$ .

So, it is possible to get  $x_k$  such that the value of  $f(x_k)$  is close enough to the minimum  $f^*$ . For definiteness, let us choose  $C = \sqrt{6}$  (to get a ‘‘convenient’’ coefficient) and consider 2 scenarios:

1.  $\|\tilde{\nabla} f(x_k)\| > \Delta\sqrt{6}$ , whence, taking (19) into account, we obtain the inequality

$$f(x_{k+1}) - f(x_k) < -\frac{\Delta^2}{2L}. \quad (20)$$

2.

$$\|\tilde{\nabla}f(x_k)\| \leq \Delta\sqrt{6}, \quad (21)$$

whence, in view of (10) we have

$$f(x_k) - f^* \leq \frac{7\Delta^2}{\mu}. \quad (22)$$

Let us consider estimate (22) acceptable for the function level and agree to terminate process (13) if (21) is satisfied.

Let us investigate an alternative situation in which for any  $k = 0, 1, \dots, N - 1$ , it is true that  $\|\tilde{\nabla}f(x_k)\| > \Delta\sqrt{6}$  and (20) holds, where

$$f(x_0) - f(x_N) = \sum_{k=0}^{N-1} (f(x_k) - f(x_{k+1})) > \frac{N\Delta^2}{2L},$$

i.e.  $N < \frac{2L}{\Delta^2}(f(x_0) - f^*)$ , which indicates the end of the process. Thus, we have the following result.

**Theorem 2.2.** *Let stopping criterion (21) be satisfied for the first time at the  $N$ -th iteration of the Gradient Descent method (13). Then the output point  $\hat{x} = x_N$  is guaranteed to satisfy the inequality*

$$f(\hat{x}) - f^* \leq \frac{7\Delta^2}{\mu}.$$

In this case, the following estimate for the number of iterations before stopping criterion is valid

$$N < \frac{2L}{\Delta^2}(f(x_0) - f^*). \quad (23)$$

It is clear that for a small value of the parameter  $\Delta > 0$ , the right-hand side of inequality (23) leads to a significantly overestimated number of iterations. At the same time, the conducted computational experiments (see Section 4, below) showed no increase in the number of iterations with a significant decrease in  $\Delta > 0$  due to the proposed early stopping rule (21).

However, in the case of a known  $\mu$ , the estimate for the number of steps  $N$  can be improved if the quality in (22) is assumed to be sufficient. Using inequality (19) we get  $\frac{1}{4L}\|\tilde{\nabla}f(x_k)\|^2 \leq \frac{\Delta^2}{L} + f(x_k) - f(x_{k+1})$ , and due to  $\tilde{\nabla}f(x_k) = L(x_k - x_{k+1})$ , we have the following estimation for every  $k \geq 0$ :

$$\begin{aligned} \|x_{k+1} - x_k\|^2 &\leq \frac{4\Delta^2}{L^2} + \frac{4(f(x_k) - f(x_{k+1}))}{L} \\ &\leq \frac{4\Delta^2}{L^2} + \frac{4(f(x_k) - f^*)}{L} \\ &\leq \frac{4\Delta^2}{L^2} + \frac{4\Delta^2}{\mu L} + \frac{4}{L} \left(1 - \frac{\mu}{L}\right)^k (f(x_0) - f^*). \end{aligned}$$

Whence one can obtain the final estimation:

$$\|x_{k+1} - x_k\| \leq 2\Delta\sqrt{\frac{1}{L^2} + \frac{1}{\mu L}} + 2\left(1 - \frac{\mu}{L}\right)^{\frac{k}{2}} \sqrt{\frac{f(x_0) - f^*}{L}}.$$

Next, summing the inequalities above for  $k = 0 \dots N - 1$ , we have

$$\|x_N - x_0\| \leq \sum_{k=0}^{N-1} \|x_{k+1} - x_k\| \leq 2N\Delta\sqrt{\frac{1}{L^2} + \frac{1}{\mu L}} + 2\sum_{k=0}^{N-1} \left(1 - \frac{\mu}{L}\right)^{\frac{k}{2}} \sqrt{\frac{f(x_0) - f^*}{L}}. \quad (24)$$

If at some step (21) is satisfied, then the required accuracy by function (22) will be achieved. Therefore, we estimate  $N$  in an alternative situation ((21) does not hold for all  $k = 0, 1, \dots, N - 1$ ). We use inequality (16) and impose the requirement that the level of approximation with respect to the function  $f(x_N) - f^* \leq \frac{7\Delta^2}{\mu}$ . In view of (16), it suffices to require that

$$\left(1 - \frac{\mu}{L}\right)^N (f(x_0) - f^*) \leq \frac{6\Delta^2}{\mu},$$

or

$$\left(1 - \frac{\mu}{L}\right)^N \leq e^{-\frac{\mu N}{L}} \leq \frac{6\Delta^2}{\mu(f(x_0) - f^*)},$$

where  $N \leq \left\lceil \frac{L}{\mu} \ln \frac{\mu(f(x_0) - f^*)}{6\Delta^2} \right\rceil$ . In this case, (24) takes the following form:

$$\|x_N - x_0\| \leq \frac{2\Delta}{\mu} \sqrt{1 + \frac{L}{\mu} \left\lceil \ln \frac{\mu(f(x_0) - f^*)}{6\Delta^2} \right\rceil} + \frac{4\sqrt{L(f(x_0) - f^*)}}{\mu}.$$

**Theorem 2.3.** *Let one of the following alternatives hold:*

1. *The Gradient Descent method (13) works  $N_*$  steps where  $N_*$  is such that*

$$N_* = \left\lceil \frac{L}{\mu} \ln \frac{\mu(f(x_0) - f^*)}{6\Delta^2} \right\rceil. \quad (25)$$

2. *For some  $N \leq N_*$ , at the  $N$ -th iteration of the method (13), stopping criterion (21) is satisfied for the first time.*

Then for the output point  $\hat{x}$  ( $\hat{x} = x_N$  or  $\hat{x} = x_{N_*}$ ) of the method (13), the following inequalities hold:

$$\begin{aligned} f(\hat{x}) - f^* &\leq \frac{7\Delta^2}{\mu}, \\ \|\hat{x} - x_0\| &\leq \frac{2\Delta}{\mu} \sqrt{1 + \frac{L}{\mu} \left\lceil \ln \frac{\mu(f(x_0) - f^*)}{6\Delta^2} \right\rceil} + \frac{4\sqrt{L(f(x_0) - f^*)}}{\mu}. \end{aligned} \quad (26)$$

**Remark 2.4.** *Since it is often difficult to estimate the value of the parameter  $\mu$  and usually  $f^*$  is not known, the estimate of the number of iterations (25) is difficult to use in practice. If the implementation works only according to stopping rule (21), then we can only confirm an upper bound on the number of iterations of the form (23), but in this case we cannot guarantee (26). However, the estimate (24) remains relevant. Moreover, the estimation of the value  $\|\hat{x} - x_0\|$  can be refined if the value of the parameter  $\mu$  is not available. Indeed, in view of (19) for the Gradient Descent method (13) with a constant step-size it holds that  $\frac{1}{4L} \|\tilde{\nabla} f(x_k)\|^2 \leq \frac{\Delta^2}{L} + f(x_k) - f(x_{k+1})$ . Whence we have  $\|x_{k+1} - x_k\|^2 \leq \frac{4\Delta^2}{L^2} + \frac{4(f(x_k) - f(x_{k+1}))}{L}$ , i.e.  $\|x_{k+1} - x_k\| \leq \frac{2\Delta}{L} + 2\sqrt{\frac{(f(x_k) - f(x_{k+1}))}{L}}$ . Further, after summing the inequalities above over  $k = 0, N-1$ , we have:*

$$\begin{aligned} \|x_0 - x_N\| &\leq \sum_{k=0}^{N-1} \|x_k - x_{k+1}\| \leq \frac{2N\Delta}{L} + 2 \sum_{k=0}^{N-1} \sqrt{\frac{f(x_k) - f(x_{k+1}))}{L}} \\ &\leq \frac{2N\Delta}{L} + \sqrt{N} \sqrt{\sum_{k=0}^{N-1} \frac{f(x_k) - f(x_{k+1}))}{L}} \\ &= \frac{2N\Delta}{L} + 2\sqrt{N} \sqrt{\frac{f(x_0) - f(x_N)}{L}} \\ &\leq \frac{2N\Delta}{L} + 2\sqrt{N} \sqrt{\frac{f(x_0) - f^*}{L}}. \end{aligned}$$

It is clear that for small values of the error  $\Delta > 0$  the following inequality

$$\|x_0 - x_N\| \leq \frac{2N\Delta}{L} + 2\sqrt{N} \sqrt{\frac{f(x_0) - f^*}{L}}$$

may turn out to be worse than (26). Taking into account (23), we get

$$\begin{aligned} \|x_0 - x_N\| &\leq \frac{2\Delta}{L} \cdot \frac{2L(f(x_0) - f^*)}{\Delta^2} + 2\sqrt{\frac{2L(f(x_0) - f^*)}{\Delta^2}} \cdot \frac{f(x_0) - f^*}{L} \\ &= \frac{4 + 2\sqrt{2}}{\Delta} (f(x_0) - f^*). \end{aligned}$$

**Remark 2.5.** By (7) and (26) the quantity  $\|\hat{x} - x_0\|$  can be comparable with  $\|x_* - x_0\|$  for a sufficiently small  $\Delta > 0$ .

**Remark 2.6.** In view of (26), it suffices to require that conditions (5) and (8) are satisfied only in  $R$ -neighborhood of the  $x_0$ , where

$$R = \frac{2\Delta}{\mu} \sqrt{1 + \frac{L}{\mu}} \left[ \ln \frac{\mu(f(x_0) - f^*)}{6\Delta^2} \right] + \frac{4\sqrt{L(f(x_0) - f^*)}}{\mu}.$$

## 2.2. Some Variant of the Gradient Descent Method with an Adaptive Step-Size Policy

In many applied optimization problems, it is difficult to estimate the Lipschitz constant of the gradient of the objective function. For example, the well-known Rosenbrock function and its multidimensional generalizations (for example, the Nesterov-Skokov function [11]) have only a locally Lipschitz gradient. Thus, it is impossible to estimate for them the Lipschitz constant of the gradient without additional restrictions on the domain in which the method operates. Therefore, we present a generalization of the universal gradient method from [9] for working with an inexact gradient of the functions satisfying PL-condition.

For  $L$ -smooth functions we have the following well-known inequality:

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

For the inexact gradient (9) we can get a similar inequality:

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

This inequality contains an exact calculation of the value of the function  $f$  at an arbitrary point from the dom $f$ . For most important applications with an inexact gradient, we do not have an opportunity to make such a calculation. An important example of such problems is some optimization problems in the Hilbert space [15] and, in a particular case, inverse problems [7]. Therefore, further, we will discuss the possibility of using an inexact function value when checking the iteration exit criterion.

Let us assume that we can calculate the inexact value  $\tilde{f}$  of the function  $f$  at any point  $x$ , so that

$$|f(x) - \tilde{f}(x)| \leq \delta. \quad (27)$$

Then we have the following inequality:

$$\tilde{f}(x) \leq \tilde{f}(y) + \langle \tilde{\nabla} f(y), x - y \rangle + L \|x - y\|^2 + \frac{\Delta^2}{2L} + 2\delta, \quad \forall x, y \in \mathbb{R}^n. \quad (28)$$

Further, when  $\mu$  is known, we select the constant  $L$  in such a way that (28) is satisfied for the points from the neighboring iterations (see Algorithm 1).

---

**Algorithm 1.** Adaptive Gradient Descent with Inexact Gradient.

---

**Require:**  $L_{\min} \geq 0, L_0 \geq L_{\min}, \delta \geq 0, \Delta \geq 0$ .

- 1: Set  $k := 0$
- 2: Calculate

$$x_{k+1} = x_k - \frac{1}{2L_k} \tilde{\nabla} f(x_k) \quad (29)$$

- 3: If the following inequality holds:

$$\tilde{f}(x_{k+1}) \leq \tilde{f}(x_k) + \langle \tilde{\nabla} f(x_k), x_{k+1} - x_k \rangle + L_k \|x_{k+1} - x_k\|^2 + \frac{\Delta^2}{2L_k} + 2\delta, \quad (30)$$

then  $k := k + 1, L_k := \max\left(\frac{L_{k-1}}{2}, L_{\min}\right)$  and go to Step 2. Otherwise,  $L_k := 2L_k$  and go to Step 3.

- 4: **return**  $x_k$
- 

Similarly to the approach of the method with a constant step-size proposed above, in the case of a sufficiently small inexact gradient

$$\|\tilde{\nabla} f(x_k)\| \leq 2\Delta \quad (31)$$



we agree to interrupt Algorithm 1. In this case, according to (10) we can guarantee that  $f(x_k) - f^* \leq \frac{5\Delta^2}{\mu}$ .

An alternative case, where condition (31) is not satisfied, can be investigated similarly to the constant step-size case in Section 2.1. A detailed proof is given in Appendix A.

The theoretical results about the operation of Algorithm 1 are presented in the following theorem.

**Theorem 2.7.** *Suppose  $f(x)$  satisfies PL-condition (5) and conditions (27),  $\Delta^2 \geq 16L\delta$  hold. Let the parameter  $L_{\min}$  in Algorithm 1 be such that  $L_{\min} \geq \frac{\mu}{4}$  and one of the following alternatives holds:*

1. Algorithm 1 works  $N_*$  steps where  $N_*$  is such that

$$N_* = \left\lceil \frac{8L}{\mu} \log \frac{\mu(f(x_0) - f^*)}{\Delta^2} \right\rceil. \quad (32)$$

2. For some  $N \leq N_*$ , at the  $N$ -th iteration of Algorithm 1, stopping criterion (31) is satisfied for the first time.

Then for the output point  $\hat{x}$  ( $\hat{x} = x_N$  or  $\hat{x} = x_{N_*}$ ) of Algorithm 1, we have the following inequalities

$$\begin{aligned} f(\hat{x}) - f^* &\leq \frac{5\Delta^2}{\mu}, \\ \|\hat{x} - x_0\| &\leq 8\frac{\Delta}{\mu} \sqrt{\frac{1}{2}\gamma^2 + 4\gamma\frac{L}{\mu} \log \frac{\mu(f(x_0) - f^*)}{\Delta^2}} + 16\frac{\sqrt{\gamma L(f(x_0) - f^*)}}{\mu}, \end{aligned} \quad (33)$$

where  $\gamma = \frac{L}{L_{\min}}$ . Also, the total number of calls to the subroutine for calculating inexact values of the objective function and step (29) is not more than  $2N + \log \frac{2L}{L_0}$ .

As we can see, estimate (33) from Theorem 2.7 for the Gradient Descent with an adaptive step-size differs significantly from estimate (26) from Theorem 2.3 for the method with a constant step-size, namely, by the presence of the factor  $\gamma$ . In the worst case, the ratio of these two estimates can be  $O\left(\frac{L}{\mu}\right)$ . However, as it will be shown in experiments, the distances  $\|\hat{x} - x_0\|$  for the methods differ insignificantly. In addition, note, that Algorithm 1 uses subroutines for finding the inexact value of the objective function more often than the gradient method with a constant step. But the number of calls to these subroutines in adaptive Algorithm 1 is not more than  $2N + \log \frac{2L}{L_0}$ . This means that the "cost" of an iteration of the adaptive algorithm is on average comparable to about two iterations of the non-adaptive method (13). At the same time, the accuracy achieved by the proposed methods is also approximately equal.

**Remark 2.8.** *Note that condition (31) is satisfied for any  $L_k \geq L$ . By construction, we obtain that  $L_k \leq 2L$ . In the estimates above, the quantity  $2L$  estimates the maximum value of the parameter  $L_k$ . The estimates above remain valid if  $L$  is replaced by  $\frac{1}{2} \max_{j \leq k} L_j$  and  $\gamma$  by  $\frac{\max_{j \leq k} L_j}{2 \min_{j \leq k} L_j}$ . Similarly, we can replace the algorithm parameter  $L_{\min}$  with  $\min_{j \leq k} L_j$ .*

**Remark 2.9.** *Note that the estimate for the number of iterations (32) in Theorem 2.7 indicates the finiteness of the process, but it is strongly overestimated. In practice, the following relation is a more interesting:*

$$N_* = \left\lceil \frac{4\hat{L}}{\mu} \log \frac{\mu(f(x_0) - f^*)}{\Delta^2} \right\rceil,$$

where  $\hat{L} = \frac{\mu}{4} \frac{1}{1 - \left(\prod_{j=0}^{N_*-1} \left(1 - \frac{\mu}{4L_j}\right)\right)^{\frac{1}{N_*}}}$  is a parameter depending on the fitted parameters  $L_j$  in Algorithm 1.

**Remark 2.10.** *Also note that we can relax the requirement  $L_{\min} \geq \frac{\mu}{4}$  to  $L_{\min} > 0$ . In this case, the estimate for the distance from the starting point to the point  $x_N$  at the  $N$ -th iteration (see the proof of the expression (41)) has the following form*

$$\|x_N - x_0\| \leq N\Delta \sqrt{\frac{1}{2L_{\min}^2} + \frac{4}{\mu L_{\min}}} + 16\sqrt{\frac{L}{L_{\min}}} \frac{\sqrt{L(f(x_0) - f^*)}}{\mu}.$$

But we can no longer use estimate (33) from Theorem 2.7. In this case, it is possible to evaluate the sufficient number of iterations of Algorithm 1, assuming that the stopping condition  $\|\tilde{\nabla} f(x_k)\| \leq 2\Delta$  is not satisfied. Further, we obtain an estimate (see the proof of (37)) for  $\Delta^2 > 16L\delta$ :

$$N < \frac{2L}{\Delta^2 - 16L\delta} (f(x_0) - f^*).$$

For experimental comparison of the Gradient Descent methods with constant and adaptive steps, a single stopping criterion must be chosen. If we consider criterion (21) for the adaptive Algorithm 1 instead of (31), then the results of Theorem 2.7 about the number of iterations (32) and the estimate of the distance from  $x_0$  to  $\hat{x}$  will remain valid. Thus, criterion (21) makes it possible to achieve the same theoretical guarantees. Therefore, further in the experimental comparison of the variants of the gradient method (Algorithms 1 and (13)), we will use stopping criterion (21).

### 3. Numerical Experiments

#### 3.1. The quadratic form minimization problem

In this section, we compare the number of iterations of method (13), required to stop according to criterion (21), and the estimate for the number of iterations (25) to achieve estimate (22). To obtain the theoretical estimate for the number of iterations (25) we need the values of the constants  $L$  and  $\mu$ . Therefore, as the first example, we consider a quadratic function for which these constants are easy to calculate.

As an inexact gradient, we use an exact gradient with random noise (9). In our experiments, we consider the following types of the additive inexactness  $v(x)$  in (9):

- **Random.** Randomly generated from a uniform distribution, i.e.  $v(x) \sim \mathcal{U}(S_1^n(0))$ , where  $S_1^n(0)$  is the  $n$  dimensional sphere with radius 1 at the center 0.
- **Antigradient.**  $v(x) = -\frac{\nabla f(x)}{\|\nabla f(x)\|}$ .
- **Constant.**  $v(x) = v \in \mathbb{R}^n$ , such that  $\|v\| = 1$ .

Let us start with a simple example that allows us to estimate the parameters  $L$  and  $\mu$ . As shown in [13], the function  $f(x) = \frac{1}{2}\langle x, Ax \rangle$  satisfies PL-condition if the operator  $A$  is non-negative definite and its spectrum is separated from zero. In such a case,  $\mu$  is the smallest nonzero eigenvalue of the matrix  $A$ . At the same time, the Lipschitz constant of the gradient is the largest eigenvalue of the matrix  $A$ . Thus, we consider the following problem of quadratic programming:

$$\min_{x \in \mathbb{R}^n} \sum_{j=k+1}^n d_j x_j^2, \tag{34}$$

where  $k$  is the number of zero eigenvalues of the matrix  $A$ , and  $d_j$  are some positive constants. Thus, we have a quadratic form with a non-negative definite diagonal matrix. In this case, we can explicitly find the constants  $\mu = \min_{j=k+1, n} (d_j), L = \max_{j=k+1, n} (d_j)$ .

In the conducted experiments, we take  $L = 1$  and change  $\mu$  from 0 to 1. The parameters  $d_j$  will be taken uniformly random from the interval  $[\mu, L]$ . We take the dimension  $n = 100$  and  $k = 10$  of zero eigenvalues. Let us compare the required number of iterations to achieve condition (21) and the estimate of  $N_*$  from Theorem 2.3. As an inexactness, we will take **Random** noise  $v(x)$ . The results for problem (34) are presented in Table 1.

$\mu$	$\Delta$	$N$	$N_*$	$\mu$	$\Delta$	$N$	$N_*$
0.01	$10^{-7}$	1528	3817	0.1	$10^{-7}$	169	406
	$10^{-4}$	841	2436		$10^{-4}$	104	267
	$10^{-1}$	155	1054		$10^{-1}$	40	129
0.9	$10^{-7}$	10	48	0.99	$10^{-7}$	6	44
	$10^{-4}$	8	33		$10^{-4}$	5	30
	$10^{-1}$	5	17		$10^{-1}$	3	16

**Table 1.** Comparison of the iteration number  $N$  to achieve condition (21) and the estimate  $N_*$  from Theorem 2.3.

In Table 1 we can see that in all cases  $N < N_*$ . It means that stopping condition (21) is reached earlier than the theoretical estimate of the number of iterations  $N_*$  justified using PL-condition (see Theorem 2.3). Also, we can note that the method converges much faster than the stated estimate for large values of  $\mu$ . At the same time, for small values of  $\mu$ , the value of  $N_*$  exceeds  $N$  by at most 2.5 times. For the other types of the noise of the gradient, a similar picture is observed.

Now, we compare the results of the Gradient Descent with a constant step-size (13) and the proposed Gradient Descent method with an adaptive step-size (Algorithm 1) when using stopping criterion (21). In Tables 2 and 3, there are presented the results of the experiments for the quadratic function (34). The experiments were carried out for the uniformly distributed noise  $v(x)$  on the sphere. In these experiments, the inexactness  $\delta = 16\Delta^2$  in the function was taken. Note that in this case, the correlation of inexactness satisfies the condition of Theorem 2.7.

$\mu$	$\Delta$	Constant		Adaptive $L$	
		Iters	Time, ms	Iters	Time, ms
0.01	$10^{-7}$	1525	139.02	668	243.45
	$10^{-4}$	837	76.72	421	151.83
	$10^{-1}$	158	14.88	75	24.67
0.1	$10^{-7}$	169	15.84	72	26.40
	$10^{-4}$	104	10.00	46	16.39
	$10^{-1}$	42	4.40	19	6.62
0.99	$10^{-7}$	6	1.01	6	2.16
	$10^{-4}$	5	0.92	5	1.75
	$10^{-1}$	3	0.55	3	1.12

**Table 2.** Comparison of the running time of the algorithms and number of iterations to achieve the accuracy  $\|\tilde{\nabla} f(x)\| \leq \sqrt{6}\Delta$  for the quadratic problem.

From Table 2, we can see that the adaptive method is inferior in real time to the Gradient Descent for all parameters  $\mu$  and  $\Delta$ . However, it needs a smaller number of iterations for big values of  $\mu$ .

### 3.2. The problem of minimizing the logistic regression function

Now let us check the work of the proposed stopping criterion in the case when it is rather difficult to estimate the constant  $\mu$  of the function which satisfies PL-condition. In this case, we will not be able to use estimate (25). This situation has been discussed in Remark 2.4. The detailed experiments presented in B.2.

However, we note that, as shown by the previous experiment, condition (21) can be achieved in a significantly smaller number of steps compared to the theoretical estimate of the number of iterations  $N_*$  from Theorem 2.3. We will consider the following optimization problem associated with logistic regression  $f(x) = \frac{1}{m} \sum_{i=1}^m \log(1 + \exp(-y_i \langle w_i, x \rangle))$ , where  $y = (y_1, \dots, y_m)^\top \in [-1, 1]^m$  is the feasible variable vector,  $W = [w_1 \dots w_m] \in \mathbb{R}^{n \times m}$  is the feature matrix, where the vector  $w_i \in \mathbb{R}^n$  is from the same space as the optimized weight vector  $w$ .

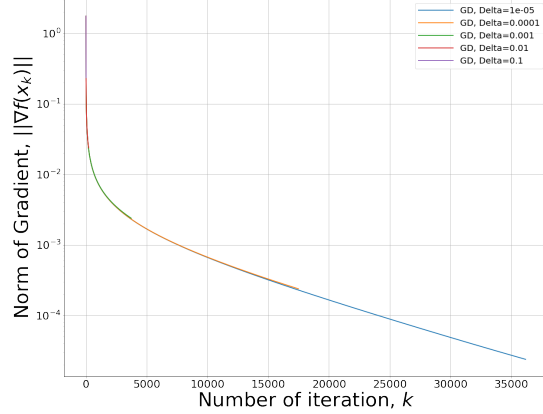
Note that this problem may not have a finite solution in the general case. So we will create such an artificial data set that there is a finite vector  $x^*$  minimizing the given function. The details of data generation is presented in Appendix B.2.

In the conducted experiments, we chose  $n = 200, m = 700$  and  $k = 10 < \min(n, \frac{m}{2})$ . We consider in this section the case of constant inexactness. From Fig. 1 it can be seen that the trajectories of the method are not the same. Moreover, adding inexactness slows down the convergence. On the other hand, the trajectories have become more similar compared to the case of the inexactness directed along the minus of the gradient (see Appendix B.2).

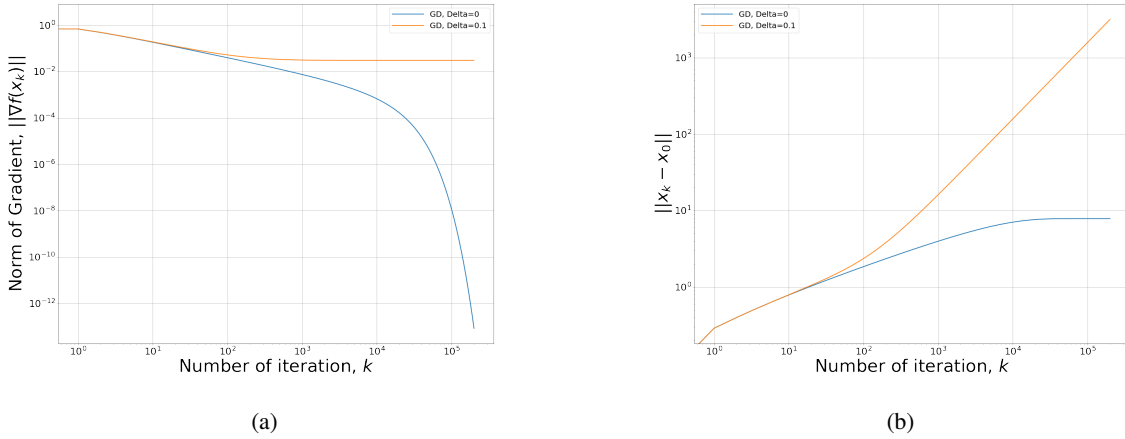
However, in this case, in Fig. 2(b) it can be seen that without using the stopping criterion, the distance  $\|x_k - x_0\|$  grows rather quickly. Thus, in the case of randomly generated gradient noise after  $10^5$  iterations of Gradient Descent method (13) the distance was 1.25 times larger compared to the result without noise in the gradient. At the same time, in the case of a constant gradient specification error, these values differ by more than two orders of magnitude (see Appendix B.2).

### 3.3. Some experiments with the Rosenbrock-type function

In this subsection, we describe results of our investigation of the behavior of the proposed adaptive Algorithm 1 for some non-convex problems. The details are presented in Appendix B.3 and B.4. Firstly, we considered the well-known two-dimensional Rosenbrock function  $f(x_1, x_2) = 100(x_2 - (x_1)^2)^2 + (x_1 - 1)^2$ .



**Fig. 1.** The rate of convergence of the gradient method in the gradient norm for different values of  $\Delta$  for the problem of minimizing logistic regression using stopping criterion (21) for the constant  $v$ .



**Fig. 2.** The results of the gradient method with respect to the norm of the gradient without using the stopping criterion for  $\Delta = 0.1$  for the problem of logistic regression minimization for the constant inaccuracy  $\Delta v$ . (a) The convergence rate with respect to the norm of the gradient; (b) the distance from the starting point to  $x_k$ .

This function is not convex, and it satisfies the Lipschitz condition for the gradient only locally. Indeed, if we consider the line  $x_2 = 0$ , then we get  $f(x_1, 0) = 100x_1^4 + (x_1 - 1)^2$ . The gradient of this function does not satisfy the Lipschitz condition. On the other hand, the Rosenbrock function satisfies locally PL-condition.

In the conducted experiments, we will vary the value of the parameter  $\Delta$  and take  $\delta = \Delta^2$ . In Table 4 in Appendix B.3, we show the results for different types of noise. As previously, from the results presented in Table 4, we can see that the number of required iterations increases with decreasing  $\Delta$  (which also tightens the stopping condition). Moreover, it increases logarithmically, which coincides with the results of Theorem 2.7. We can also note that the resulting distance from the starting point  $x_0$  to the last point does not exceed the distance from the starting point  $x_0$  to the nearest optimal one  $x_* = (1, 1)$  everywhere. In addition, for all considered types of the gradient error (noise), a comparable convergence rate is observed according to the number of iterations until stopping criterion (31) is satisfied, and to the running time for the corresponding values of  $\Delta$ .

Further, let us consider a system of nonlinear equations  $g(x) = 0$ , where  $g_1 = \frac{1}{2}(x_1 - 1)$ ,  $g_i = x_i - 2x_{i-1}^2 + 1$ ,  $i = \overline{2, n}$ . The problem of solving this system is equivalent to minimizing the following Nesterov-Skokov function (see [11])

$$f(x) = \frac{1}{4}(1 - x_1)^2 + \sum_{i=1}^{n-1} (x_{i+1} - 2x_i^2 + 1)^2. \quad (35)$$

This function is analogous to the Rosenbrock function. It is also non-convex and satisfies the Lipschitz gradient condition only locally. Also, function (35) has a global minimum at the point  $(1, 1 \dots 1, 1)^\top$  and an optimal value  $f^* = 0$ . Moreover, this function locally satisfies PL-condition (see the proof in Appendix D).

As it was seen from the results of the previous experiments, our proposed stopping criterion (31) of Algorithm 1 can work equally well for all considered types of noise in the gradient. In the current experiments for the Nesterov-Skokov function, we used the random noise of the gradient which is uniformly distributed on the sphere. For the experiments, the starting point is  $(-1, 1, \dots, 1, 1)^\top$  and therefore  $\|x_0 - x_*\| = 2$ . We will vary the value of the inexactness  $\Delta$  and the dimension of the problem  $n$ .

Table 5 in Appendix B.4 shows the results of the adaptive gradient method 1 for the Nesterov-Skokov function (35). Firstly, we see that as the dimension of  $n$  increases, the difference between the required time to solve the problem for different  $\Delta$  grows significantly. Secondly, for different  $n$  with the same  $\Delta$ , the method converges to a solution with significantly different accuracy. We can also note that  $\|x_N - x_0\|$  exceeds  $\|x_0 - x_*\|$  by at most 2 times. Moreover, significant upward deviations are observed for the cases when numerous iterations are made ( $n = 5$  and  $\Delta = 10^{-4}, 10^{-3}$ ). It can also be noted that even for sufficiently small values of the norm of the gradient, the accuracy by the function turns out to be quite low (which is typical for the Nesterov-Skokov function).

## 4. Conclusion

This paper studies stopping criteria for the gradient method with an inexact gradient. The authors focus on the case of non-convex functions. The paper presents a stopping criterion that finds a compromise between the accuracy of the obtained point and the distance to the starting point. Moreover, it is shown that the method moves away from the starting point to a distance comparable to the distance to the nearest solution if the function satisfies PL-condition.

Besides, the paper considers the cases of a constant and adaptive step size in the gradient methods. For both cases, we present theoretical analysis and the number of iterations required to approach the stopping criterion or to find the point with the required quality.

In addition, the paper contains numerical experiments demonstrating the work of the stopping criterion. In particular, there are experiments on a quadratic function (convex, but not strongly convex), demonstrating the stopping criterion to be approached faster than the theoretical estimation of the iteration number  $N_*$ . Also, we present experiments on the problem of logistic regression where the objective function is convex and meets PL-condition only locally. The proposed stopping criterion on this function stops the growth of the distance  $\|x_k - x_0\|$ . Moreover, we present experiments on non-convex functions: the Rosenbrock function and its multidimensional generalization, which is the Nesterov-Skokov function. The first function demonstrates that our stopping criterion works with general types of inexactness. The second function demonstrates that even a small inexactness can lead to quite a high value of function and this value cannot be improved. Also, we demonstrate that for some noises, the gradient method can move away quite far on the Nesterov-Skokov function without a stopping criterion.

## Acknowledgements

The research by F. Stonyakin in Section 2.1 was supported by the strategic academic leadership program "Priority 2030" (Agreement 075-02-2021-1316, 30.09.2021). The research by B. Polyak in Section 1 was supported by the Russian Science Foundation (project No. 21-71-30005).

## References

- [1] Belkin, M.: Fit without fear: remarkable mathematical phenomena of deep learning through the prism of interpolation. *Acta Numerica* **30**, 203–248 (2021)
- [2] d’Aspremont, A.: Smooth optimization with approximate gradient. *SIAM Journal on Optimization* **19**(3), 1171–1183 (2008)
- [3] Devolder, O.: Exactness, inexactness and stochasticity in first-order methods for large-scale convex optimization. Ph.D. thesis, ICTEAM and CORE, Université Catholique de Louvain (2013)
- [4] Devolder, O., Glineur, F., Nesterov, Y.: First-order methods of smooth convex optimization with inexact oracle. *Mathematical Programming* **146**(1), 37–75 (2014)
- [5] Emelin, I.V., Krasnosel’skii, M.A.: The stoppage rule in iterative procedures of solving ill-posed problems. *Autom. Remote Control* **39**, 1783–1787 (1979). Translation from *Avtom. Telemekh.* 1978, No. 12, 59-63 (in Russian)
- [6] Gasnikov, A.V.: In: *Modern Numerical Optimization Methods: The Universal Gradient Descent Method*. Moscow, MCCME (2021). (in Russian)
- [7] Kabanikhin, S.I.: Inverse and ill-posed problems. In: *Inverse and Ill-posed Problems*. deGruyter (2011)
- [8] Karimi, H., Nutini, J., Schmidt, M.: Linear convergence of gradient and proximal-gradient methods under the polyak-Lojasiewicz condition. In: *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*, pp. 795–811. Springer (2016)
- [9] Nesterov, Y.: Universal gradient methods for convex optimization problems. *Mathematical Programming* **152**(1), 381–404 (2015)
- [10] Nesterov, Y., Polyak, B.T.: Cubic regularization of newton method and its global performance. *Mathematical Programming* **108**(1), 177–205 (2006)
- [11] Nesterov, Y., Skokov, V.: On the issue of testing unconstrained optimization algorithms. In: *Numerical methods of mathematical programming*, pp. 77–91. Moscow (1980). (in Russian)
- [12] Polyak, B., Tremba, A.: New versions of newton method: step-size choice, convergence domain and under-determined equations. *Optimization Methods and Software* **35**(6), 1272–1303 (2020)
- [13] Polyak, B.T.: Gradient methods for minimizing functionals. *Comput. Math. Math. Phys.* **3**(4), 864–878 (1963). (in Russian)
- [14] Polyak, B.T.: *Introduction to optimization*. optimization software. Inc., New York **1**, 32 (1987)
- [15] Vasilyev, F.: *Optimization methods*. Moscow: FP (2002). (in Russian)
- [16] Vasin, A., Gasnikov, A., Spokoiny, V.: Stopping rules for accelerated gradient methods with additive noise in gradient. *arXiv preprint arXiv:2102.02921* (2021)
- [17] Vorontsova, E., Hildbrand, R., Gasnikov, A., Stonyakin, F.: *Convex optimization* (2021). (in Russian)

## A. The proof of Theorem 2.7

In the case, if at some  $k$ -th iteration of Algorithm 1 the condition (31) is satisfied, then according to the PL-condition we obtain that

$$f(x_k) - f^* \leq \frac{5\Delta^2}{\mu}.$$

Let us study the estimation of the quality of the output point of Algorithm 1 under conditions when (31) is not satisfied. Note that for each iteration  $k \geq 1$  condition (30) is satisfied, i.e.

$$\tilde{f}(x_{k+1}) \leq \tilde{f}(x_k) + \langle \tilde{\nabla} f(x_k), x_{k+1} - x_k \rangle + L_k \|x_{k+1} - x_k\|^2 + \frac{\Delta^2}{2L_k} + 2\delta.$$

Then by using condition (27), we get

$$f(x_{k+1}) \leq f(x_k) + \langle \tilde{\nabla} f(x_k), x_{k+1} - x_k \rangle + L_k \|x_{k+1} - x_k\|^2 + \frac{\Delta^2}{2L_k} + 4\delta.$$

Moreover, taking into account (29), we have the following inequality

$$f(x_{k+1}) - f(x_k) \leq \frac{\Delta^2}{2L_k} - \frac{1}{4L_k} \|\tilde{\nabla} f(x_k)\|^2 + 2\delta. \quad (36)$$

This condition, together with the relation for the inexactness  $8L_k\delta \leq \Delta^2$ , tells us that if  $\|\tilde{\nabla} f(x_k)\| \geq C\Delta$  for  $C > \sqrt{3}$ , then  $f(x_{k+1}) < f(x_k)$  is guaranteed and the method converges to the minimum. For definiteness, we take  $C = 2$ . Then if  $\|\tilde{\nabla} f(x_k)\| > 2\Delta$  and (36) holds for every  $k = 0, 1, \dots, N-1$ , then

$$f(x_0) - f(x_N) = \sum_{k=0}^{N-1} (f(x_k) - f(x_{k+1})) > \frac{\Delta^2}{2} \sum_{k=0}^{N-1} \frac{1}{L_k} - 2\delta N \geq \frac{N\Delta^2}{4L} - 4\delta N,$$

i.e. at  $16L\delta < \Delta^2$ , we have

$$N < \frac{2L}{\Delta^2 - 16L\delta} (f(x_0) - f^*), \quad (37)$$

which indicates the end of the process.

On the other hand, from estimate (36) one can get an estimate for the function residual at the  $k$ -th iteration. Using the PL-condition, we get

$$f(x_{k+1}) - f(x_k) \leq \frac{\Delta^2}{2L_k} - \frac{\mu}{4L_k} (f(x_k) - f^*) + 4\delta,$$

whence

$$\begin{aligned} f(x_{k+1}) - f^* &\leq \left(1 - \frac{\mu}{4L_k}\right) (f(x_k) - f^*) + \frac{\Delta^2}{2L_k} + 2\delta \\ &\leq \prod_{j=0}^k \left(1 - \frac{\mu}{4L_j}\right) (f(x_0) - f^*) + \frac{\Delta^2}{2} \left(\frac{1}{L_k} + \sum_{j=0}^{k-1} \frac{1}{L_j} \prod_{i=j+1}^k \left(1 - \frac{\mu}{4L_i}\right)\right) + \\ &\quad + 4\delta \left(\sum_{j=0}^{k-1} \prod_{i=j+1}^k \left(1 - \frac{\mu}{4L_i}\right)\right). \end{aligned}$$

Let us estimate the second term. We denote by  $S_k = \left(\frac{1}{L_k} + \sum_{j=0}^{k-1} \frac{1}{L_j} \prod_{i=1}^{k-j} \left(1 - \frac{\mu}{4L_i}\right)\right)$ . Note that  $S_k$  with  $k \geq 1$  satisfies the recursive formula  $S_k = \frac{1}{L_k} + S_{k-1} \left(1 - \frac{\mu}{4L_k}\right)$ . Let us consider two cases. In the first case, if  $\mu S_{k-1} \geq 4$ , then  $S_k = \frac{4 - \mu S_{k-1}}{4L_k} + S_{k-1} \leq S_{k-1}$ . In the second case, for  $\mu S_{k-1} < 4$  we get that  $S_{k-1} < \frac{4}{\mu}$  and  $S_k < \frac{1}{L_k} + \frac{4}{\mu} \left(1 - \frac{\mu}{4L_k}\right) = \frac{4}{\mu}$ . Thus  $S_k \leq \max\left(S_{k-1}, \frac{4}{\mu}\right)$ . By sequentially expanding, we obtain the estimate

$$S_k \leq \max\left(S_0, \frac{4}{\mu}\right) = \max\left(\frac{1}{L_0}, \frac{4}{\mu}\right) \leq \frac{4}{\mu}.$$

We also estimate the third term using  $L_j \leq 2L$  as

$$\sum_{j=0}^{k-1} \prod_{i=j+1}^k \left(1 - \frac{\mu}{4L_i}\right) \leq \sum_{j=0}^{k-1} \left(1 - \frac{\mu}{8L}\right)^j \leq \frac{8L}{\mu}.$$

Thus, we obtain the estimate

$$f(x_{k+1}) - f^* \leq \prod_{j=0}^k \left(1 - \frac{\mu}{4L_j}\right) (f(x_0) - f^*) + \frac{2\Delta^2}{\mu} + \frac{32L}{\mu} \delta. \quad (38)$$

The estimate of (38) depends on how the algorithm works. Using the inequality  $L_k \leq 2L$ , we obtain a result about the convergence which is determined only by the parameters of the function

$$f(x_{k+1}) - f^* \leq \left(1 - \frac{\mu}{8L}\right)^{k+1} (f(x_0) - f^*) + \frac{2\Delta^2}{\mu} + \frac{32L}{\mu} \delta. \quad (39)$$

Next, we use the relation for the inexactness  $16L\delta \leq \Delta^2$  and estimate (39):

$$f(x_{k+1}) - f^* \leq \left(1 - \frac{\mu}{8L}\right)^{k+1} (f(x_0) - f^*) + \frac{4\Delta^2}{\mu}. \quad (40)$$

Let us estimate the distance  $\|x_0 - x_N\|$  in the same way as it was done for the gradient descent method with a constant step. Using inequalities (36), (40) and taking into account that  $\tilde{\nabla} f(x_k) = 2L_k(x_{k+1} - x_k)$ , we get the following estimate

$$\begin{aligned} \|x_{k+1} - x_k\|^2 &= \frac{\|\tilde{\nabla} f(x_k)\|^2}{4L_k^2} \leq \frac{\Delta^2}{2L_k^2} + \frac{f(x_{k+1}) - f(x_k)}{L_k} \\ &\leq \frac{\Delta^2}{2L_k^2} + \frac{f(x_{k+1}) - f^*}{L_k} \\ &\leq \frac{\Delta^2}{2L_k^2} + \frac{4\Delta^2}{\mu L_k} + \frac{1}{L_k} \left(1 - \frac{\mu}{8L}\right)^k (f(x_0) - f(x^*)) \\ &\leq \frac{4\Delta^2}{2L_{\min}^2} + \frac{\Delta^2}{\mu L_{\min}} + \frac{1}{L_{\min}} \left(1 - \frac{\mu}{8L}\right)^k (f(x_0) - f(x^*)). \end{aligned}$$

After summing these inequalities over  $k$  from 0 to  $N - 1$ , we get the following result

$$\|x_{k+1} - x_k\| \leq \Delta \sqrt{\frac{1}{2L_{\min}^2} + \frac{4}{\mu L_{\min}}} + \left(1 - \frac{\mu}{8L}\right)^{\frac{k}{2}} \sqrt{\frac{f(x_0) - f(x^*)}{L_{\min}}}.$$

Thus,

$$\|x_N - x_0\| \leq \sum_{k=0}^{N-1} \|x_{k+1} - x_k\| \leq N\Delta \sqrt{\frac{1}{2L_{\min}^2} + \frac{4}{\mu L_{\min}}} + 16\sqrt{\frac{L}{L_{\min}}} \frac{\sqrt{L(f(x_0) - f^*)}}{\mu}. \quad (41)$$

We can estimate the number of iterations of Algorithm 1 from the inequality (40) as follows

$$N \leq \left\lceil \frac{8L}{\mu} \log \frac{\mu(f(x_0) - f^*)}{2\Delta^2} \right\rceil.$$

Let us introduce the notation  $\gamma = \frac{L}{L_{\min}}$ . Then we can estimate the final estimation of the distance from the starting point  $x_0$  to the current one  $x_N$  as follows

$$\|x_N - x_0\| \leq 8\frac{\Delta}{\mu} \sqrt{\frac{1}{2}\gamma^2 + 4\gamma} \frac{L}{\mu} \log \frac{\mu(f(x_0) - f^*)}{2\Delta^2} + 16\frac{\sqrt{\gamma L(f(x_0) - f^*)}}{\mu}.$$

Note that the factor  $\gamma = \frac{L}{L_{\min}}$  appeared in this estimate, which depends on an unknown constant and the parameter of the algorithm  $L_{\min}$ . In the worst case, we have  $\gamma = \frac{4L}{\mu}$ .



Let us estimate the number of additional calculations of the value of the function and operations of the form (29) of the adaptive gradient method 1 in comparison with the gradient descent with a constant step-size (13). Let  $i_k$  computations of the inexact gradient be made at the  $k$ -th step. Then  $2^{i_1-1} = \frac{2L_k}{L_{k-1}}$ . Then we note that

$$\prod_{k=1}^N 2^{i_k-1} = \prod_{k=1}^N \frac{2L_k}{L_{k-1}} = 2^N \frac{L_N}{L_0}.$$

As mentioned above, it is true that  $L_N \leq 2L$ . Then the total number of additional function evaluations and steps of the form (29)  $I(N) = \sum_{i=1}^N i_k$  is estimated from above as follows

$$2^{I(N)-2N} \leq \frac{2L}{L_0}.$$

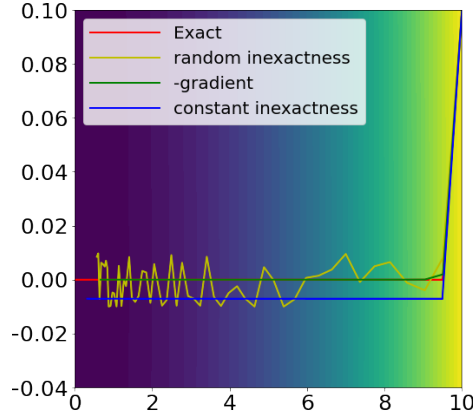
Therefore,

$$I(N) \leq 2N + \log \frac{2L}{L_0}.$$

## B. Numerical Experiments

### B.1. The quadratic form minimization problem

In this section, we present additional experiments for quadratic problem presented in Section 3.1.



**Fig. 3.** Convergence of the gradient method with stopping criterion (21) in the two-dimensional case for different  $v$ : zero (exact), randomly generated at each iteration (random inexactness), co-directional with minus gradient (Antigradient), and constant (constant inexactness).

Now, let us consider a quadratic form in dimension  $n = 2$  with coefficients 0.05 and 1. In this case, we can draw the convergence trajectory of the method for different types of noise (see Fig. 3). We chose the starting point  $x_0 = (10, 0.1)^\top$  and noise level  $\Delta = 10^{-2}$ . In the case of constant inexactness, the vector  $v = (1, 1)^\top$  was taken.

We can see that, if the inexact gradient is collinear to the exact one, then the trajectories coincide almost completely for the gradient method with exact and inexact gradients. At the same time, we do not see changes in the coordinate  $x_k$ , on which the function almost does not depend. For the constant noise, we observe a constant displacement of the trajectory from the path of the inexact gradient descent method.

Now we compare the results of the gradient descent with a constant step-size (13) and the proposed gradient descent method with an adaptive step-size (Algorithm 1) when we use stopping criterion (21). In Tables 2 and 3, there were presented the results of the experiments for the quadratic function (34). The experiments were carried out for the uniformly distributed noise  $v(x)$  on the sphere. In these experiments, the inexactness  $\delta = 16\Delta^2$  in the function was taken. Note that in this case, the correlation of inexactness satisfies the condition of Theorem 2.7.

From Table 2, we can see that the adaptive method is inferior in real time to the gradient descent for all parameters  $\mu$  and  $\Delta$ . However, it needs a smaller number of iterations for small  $\mu$ .

$\mu$	$\Delta$	Constant			Adaptive $L$		
		$\ x_N - x_0\ $	$\frac{\ \nabla f(x_N)\ }{\Delta}$	$f(x_N) - f^*$	$\ x_N - x_0\ $	$\frac{\ \nabla f(x_N)\ }{\Delta}$	$f(x_N) - f^*$
0.01	$10^{-7}$	948.7	2.34	$0.25 \cdot 10^{-11}$	948.7	2.01	$0.18 \cdot 10^{-11}$
	$10^{-4}$	948.7	2.37	$0.26 \cdot 10^{-5}$	948.7	2.04	$0.19 \cdot 10^{-5}$
	$10^{-1}$	946.2	2.39	2.50	946.5	2.30	1.93
0.1	$10^{-7}$	948.7	2.17	$0.22 \cdot 10^{-12}$	948.7	1.90	$0.14 \cdot 10^{-12}$
	$10^{-4}$	948.7	1.95	$0.17 \cdot 10^{-6}$	948.7	2.14	$0.19 \cdot 10^{-6}$
	$10^{-1}$	948.3	2.15	0.20	948.4	1.80	0.14
0.9	$10^{-7}$	948.7	0.93	$0.46 \cdot 10^{-14}$	948.7	0.92	$0.45 \cdot 10^{-14}$
	$10^{-4}$	948.7	0.86	$0.39 \cdot 10^{-8}$	948.7	0.95	$0.48 \cdot 10^{-8}$
	$10^{-1}$	948.7	0.99	$0.52 \cdot 10^{-2}$	948.7	0.88	$0.41 \cdot 10^{-2}$
0.99	$10^{-7}$	948.7	0.99	$0.49 \cdot 10^{-14}$	948.7	0.94	$0.44 \cdot 10^{-14}$
	$10^{-4}$	948.7	0.94	$0.44 \cdot 10^{-8}$	948.7	0.93	$0.43 \cdot 10^{-8}$
	$10^{-1}$	948.7	1.02	$0.52 \cdot 10^{-2}$	948.6	0.99	$0.49 \cdot 10^{-2}$

**Table 3.** Comparison of algorithms in terms of the achieved accuracy in terms of the gradient norm and the distance from the start point to the last point. The distance from the starting point to the nearest optimal one is 948.7.

Note that according to the results presented in Table 3, the compared methods lead to the achievement of approximately the same quality of the approximate solution. In this case, the trajectories of the methods are moved away approximately equally from the starting point. In this case, the distance  $\|x_N - x_0\|$  is approximately equal to the distance  $x_0$  to the nearest exact solution  $x_*$ . As we can see, the achieved accuracy is no less than  $\frac{7\Delta^2}{\mu}$ .

### B.2. The problem of minimizing the logistic regression function

Now let us check the work of the proposed stopping criterion in the case when it is rather difficult to estimate the constant  $\mu$  of the function which satisfies the PL-condition. In this case, we will not be able to use estimate (25). This situation has been discussed in Remark 2.4.

However, we note that, as shown by the previous experiment, condition (21) can be achieved in a significantly smaller number of steps compared to the theoretical estimate of the number of iterations  $N_*$  from Theorem 2.3. We will consider the following optimization problem associated with logistic regression:

$$f(x) = \frac{1}{m} \sum_{i=1}^m \log(1 + \exp(-y_i \langle w_i, x \rangle)), \quad (42)$$

where  $y = (y_1, \dots, y_m)^\top \in [-1, 1]^m$  is the feasible variable vector,  $W = [w_1 \dots w_m] \in \mathbb{R}^{n \times m}$  is the feature matrix, where the vector  $w_i \in \mathbb{R}^n$  is from the same space as the optimized weight vector  $w$ .

Note that this problem may not have a finite solution in the general case. So we will create such an artificial data set that there is a finite vector  $x^*$  minimizing the given function. To do this, we generate  $W$  and  $y$  as follows:

1. We construct  $k \leq \min(n, \frac{m}{2})$  orthogonal vectors with the unit norm and combine them into a matrix  $W_B \in \mathbb{R}^{n \times k}$ .
2. Construct a matrix  $\widetilde{W} = W_B V^\top \in \mathbb{R}^{n \times m - 2k}$ , where  $V \in \mathbb{R}^{m - 2k \times k}$  is some random matrix that defines the expansion of the vectors from the matrix  $W$  in the basis  $W_B$ .
3. Construct some vector  $x_0$  and define new vectors  $\widetilde{y} = \text{sign}(\widetilde{W} x_0)$ ,  $y_1 = \text{sign}(W_B x_0)$ .
4. Define the feasible variable vector  $y = [y_1 | -y_1 \widetilde{y}] \in [-1, 1]^m$ .
5. Define the feature matrix  $W = [W_B | W_B \widetilde{W}] \in \mathbb{R}^{m \times n}$ .

**Proposition B.1.** For the function (42) with such data, the following statements hold:

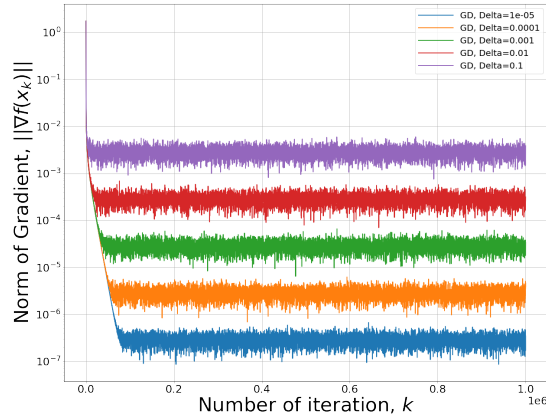
1. The function  $f$  satisfies the PL-condition on any compact set  $K$ ;

2. The function  $f$  has a Lipschitz gradient with the constant  $L = \frac{\lambda_{\max}(W^T W)}{4m}$ ;
3. There is a finite  $x_*$ , where the objective function reaches its minimal value;
4. The set of minimum points  $X_*$  is unbounded if  $k < n$ .

The proof of Proposition B.1 is presented in Appendix C.

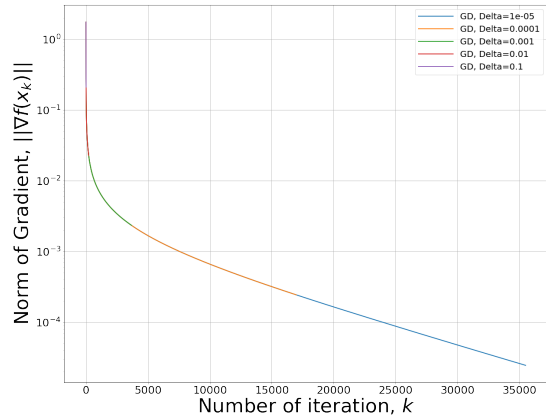
In the conducted experiments, we chose  $n = 200, m = 700$  and  $k = 10 < \min(n, \frac{m}{2})$ . Accordingly, there are 700 objects with 200 features and a feature matrix of rank 10. Thus, the set of the solutions for the minimization problem of the function  $f$  with such data is non-empty and unbounded.

Let us apply the gradient descent method (13) for various inexactness  $\Delta$  in the gradient with stopping criterion (21). As an inexact gradient, we will use a gradient with random noise **Random**.



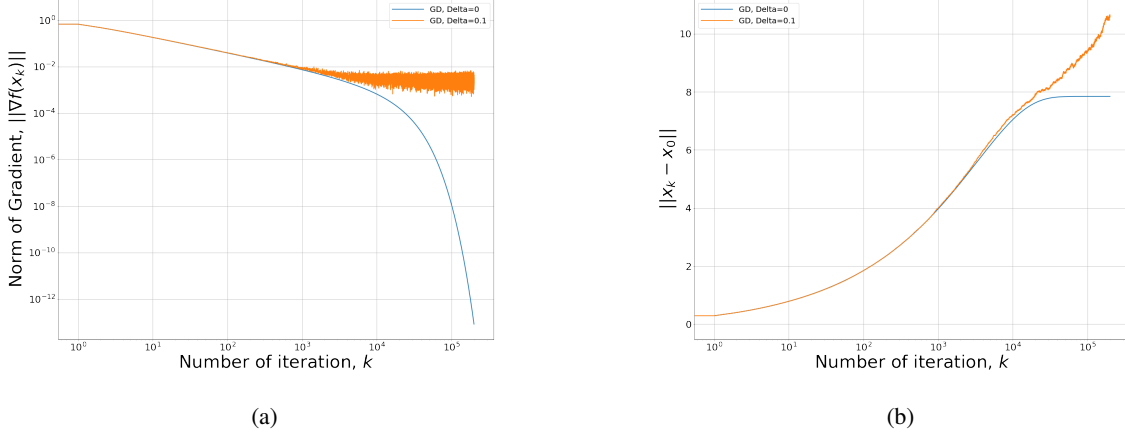
**Fig. 4.** The convergence rate of the gradient method with respect to the norm of the gradient for different values of the inexactness  $\Delta$  for the problem of logistic regression minimization at the first  $N = 10^5$  iterations without using the stopping criterion.

In Fig. 4 there was shown the plot of the convergence of the gradient method with different levels of the noise  $\Delta$  without using the stopping criterion proposed in this article. It can be seen that the method reaches points with a gradient norm of order  $\Delta$ , but it cannot converge closer.



**Fig. 5.** The convergence rate of the gradient method with respect to the norm of the gradient for different values of the inexactness  $\Delta$  for the problem of logistic regression minimization with the use of stopping criterion (21).

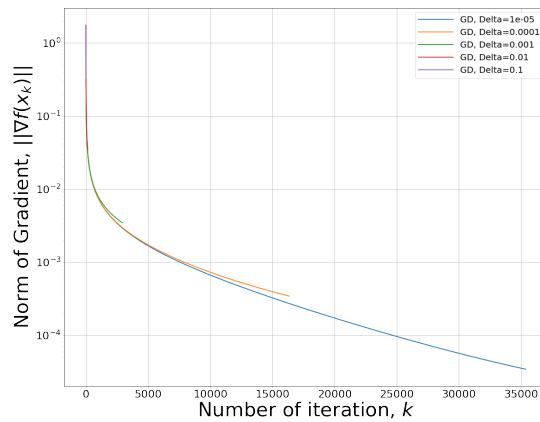
At the same time, in Fig. 5 there was shown the plot of the convergence of the gradient method with the use of stopping criterion (21). As we can see, the method stops when it reaches the accuracy  $\|\tilde{\nabla}f(x_k)\| \sim \Delta$ . We also note, in this example, that the trajectories of the methods practically coincide until the corresponding accuracy is achieved.



**Fig. 6.** The results of the gradient method with respect to the norm of the gradient without using the stopping criterion for  $\Delta = 0.1$  for the problem of logistic regression minimization for the inexactness  $\Delta v(x)$ . (a) The convergence rate with respect to the norm of the gradient; (b) the distance from the starting point to  $x_k$ .

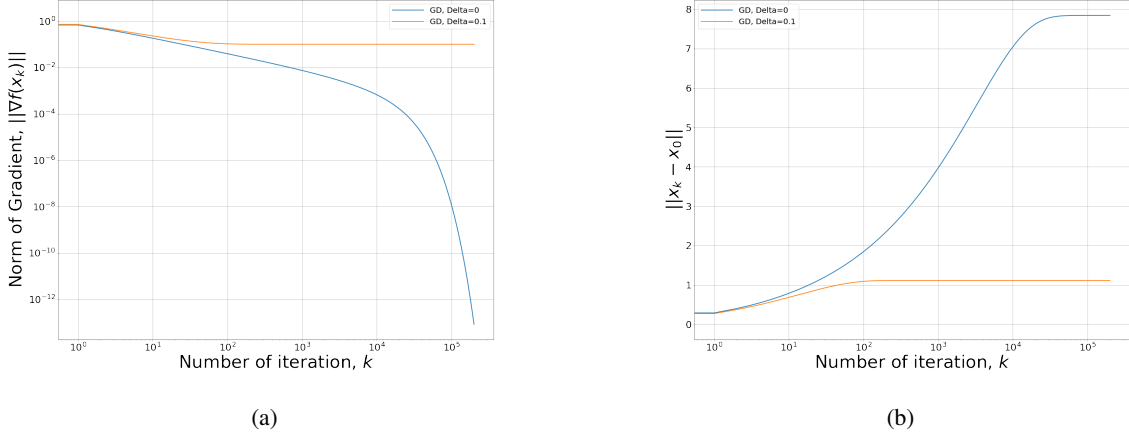
In Fig. 6(a) the plot of the convergence of the gradient method with the noise level  $\Delta = 0.1$  without using the stopping is shown. We can see that, after reaching  $\|\tilde{\nabla}f(x_k)\| \sim \Delta = 0.1$ , gradient method (13) slows down significantly compared to the noise-free method (2). Also note that if the gradient method works without using the stopping criterion, the distance from the starting point grows uncontrollably. Moreover, this growth exceeds the distance increase for the method without noise.

Next, we consider other types of problems in which the additive inexactness of the gradient proposed in section 3.1 occurs. The results for the inexact case **Antigradient** are shown in Fig. 7 and 8. In Fig. 7, we can see that the trajectories of the method begin to noticeably differ from the previous case. The inexact gradient method converges more slowly as the value of the parameter  $\Delta$  increases.



**Fig. 7.** The convergence rate of the gradient method with respect to the norm of the gradient for different values of the inexactness  $\Delta$  for the problem of logistic regression minimization using stopping criterion (21) for  $v = -\frac{\Delta}{\|\nabla f(x)\|} \nabla f(x)$ .

However, on the graph 8(b), we can see that when methods use such an inexact gradient, the error does not accumulate. Indeed, the method stabilizes near a point such that  $\|\nabla f(x)\| \approx 0.1 = \Delta$ , as it can be seen from Fig. 8(a). Moreover, the method stops at the distance of  $\|x_k - x_0\| \approx 1$  and does not move further.



**Fig. 8.** The results of the gradient method with respect to the norm of the gradient without using the stopping criterion for  $\Delta = 0.1$  for the problem of logistic regression minimization for  $v = -\frac{\Delta}{\|\nabla f(x)\|} \nabla f(x)$ . (a) The convergence rate with respect to the norm of the gradient; (b) the distance from the starting point to  $x_k$ .

Thus, in the case of  $v = -\frac{\Delta}{\|\nabla f(x)\|} \nabla f(x)$ , we can see that there is no problem of too large growth of the distance from the starting point to the resulting one.

The situation is essentially different for the inexactness in form of constant vector  $v$ , which is the same at all iterations. From Fig. 1 it can be seen that the trajectories of the method are also not the same. Moreover, adding inexactness slows down the convergence somewhat. On the other hand, the trajectories have become more similar compared to the case of the inexactness directed along the antigradient.

### B.3. Some experiments with the Rosenbrock function

In this subsection, we investigate the behavior of the proposed adaptive Algorithm 1 for the well-known two-dimensional Rosenbrock function

$$f(x_1, x_2) = 100(x_2 - (x_1)^2)^2 + (x_1 - 1)^2.$$

This function is not convex, and it satisfies the Lipschitz condition for the gradient only locally. Indeed, if we consider the line  $x_2 = 0$ , then we get  $f(x_1, 0) = 100x_1^4 + (x_1 - 1)^2$ . The gradient of this function does not satisfy the Lipschitz condition. On the other hand, the Rosenbrock function satisfies locally the PL-condition. Indeed, let us consider the system of nonlinear equations  $g_1(x_1, x_2) = 10(-x_2 + (x_1)^2)$ ,  $g_2(x_1, x_2) = x_1 - 1$ . The Jacobian of this system is

$$J = \begin{bmatrix} 20x_1 & -10 \\ 1 & 0 \end{bmatrix},$$

and consequently,

$$JJ^T = \begin{bmatrix} 400x_1^2 + 100 & 20x_1 \\ 20x_1 & 1 \end{bmatrix} \succ 0.$$

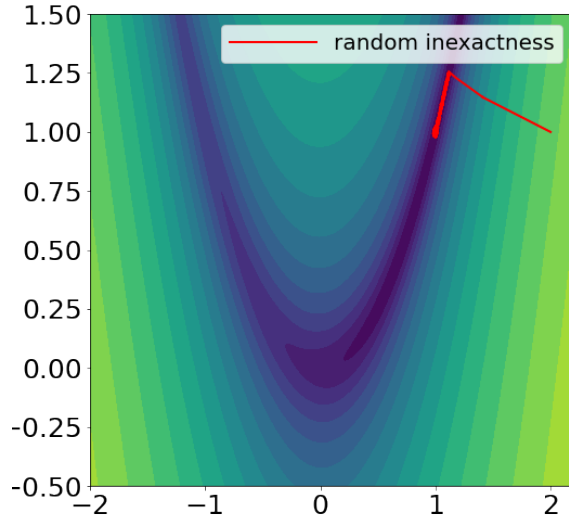
Thus, for any compact set there exists some constant  $\mu$ , such that  $JJ^T \succeq \mu I$ . Then, according to the results given in the introduction [10], the Rosenbrock function satisfies the PL-condition with the constant  $\mu$  on the corresponding compact set.

In the conducted experiments, we vary the value of the parameter  $\Delta$  and take  $\delta = \Delta^2$ . The starting point for all parameters is  $x_1 = 1, x_2 = 2$ . The distance from the initial point to the optimal point  $\mathbf{1} = (1, 1)$  is 1. In Table 4, we show the results for different types of noise. The vector  $v = (1, 0)^T$  was taken as a constant noise. In this experiment (and also in the next subsection), we will use stopping criterion (31).

Inexactness	$\Delta$	Iters	Time, ms	$\ x_N - x_0\ $	$\frac{\ \nabla f(x_N)\ }{\Delta}$	$f(x_N) - f_*$
Random	$10^{-4}$	7266	2273.26	0.999	2.70	$0.89 \cdot 10^{-7}$
	$10^{-3}$	5412	2594.61	0.994	2.90	$1.00 \cdot 10^{-5}$
	$10^{-2}$	3690	3163.32	0.940	2.61	$0.89 \cdot 10^{-3}$
Antigradient	$10^{-4}$	7188	2615.61	0.999	2.99	$0.11 \cdot 10^{-6}$
	$10^{-3}$	5493	2490.56	0.993	3.00	$0.11 \cdot 10^{-4}$
	$10^{-2}$	3536	3031.05	0.931	2.99	$0.12 \cdot 10^{-2}$
Constant	$10^{-4}$	7491	2301.32	1.000	1.54	$0.27 \cdot 10^{-7}$
	$10^{-3}$	5697	2490.32	0.997	1.87	$0.24 \cdot 10^{-5}$
	$10^{-2}$	3965	3485.89	0.965	1.93	$0.30 \cdot 10^{-3}$

**Table 4.** The results of the adaptive gradient descent for the 2D Rosenbrock function using stopping criterion (31).

As previously, from the results presented in Table 4, we can see that the number of required iterations increases with decreasing  $\Delta$  (which also tightens the stopping condition). Moreover, it increases logarithmically, which coincides with the results of Theorem 2.7. We can also note that the resulting distance from the starting point  $x_0$  to the last point does not exceed the distance from the starting point  $x_0$  to the nearest optimal one  $x_* = (1, 1)$  everywhere. In addition, for all considered types of the gradient error (noise), a comparable convergence rate is observed according to the number of iterations until stopping criterion (31) is satisfied, and to the running time for the corresponding values of  $\Delta$ .



**Fig. 9.** The trajectory of the gradient method on 2D Rosenbrock function with Random Inexactness in gradient

Note that in this example, the gradient method without a stopping criterion approaches some level for the function and next iterations are meaningless. So, in Fig. 9 we can see that the method stopped to improve the function value after some iterations.

#### B.4. Some experiments with the Nesterov-Skovok function

Let us consider a system of nonlinear equations  $g(x) = 0$ , where  $g_1 = \frac{1}{2}(x_1 - 1)$ ,  $g_i = x_i - 2x_{i-1}^2 + 1$ ,  $i = \overline{2, n}$ . The problem of solving this system is equivalent to minimizing the following Nesterov-Skovok function (see [11]):

$$f(x) = \frac{1}{4}(1 - x_1)^2 + \sum_{i=1}^{n-1} (x_{i+1} - 2x_i^2 + 1)^2. \tag{43}$$

This function is analogous to the Rosenbrock function. It is also non-convex and satisfies the Lipschitz gradient condition only locally. Also, function (43) has a global minimum at the point  $(1, 1 \dots 1, 1)^\top$  and an optimal value  $f^* = 0$ .

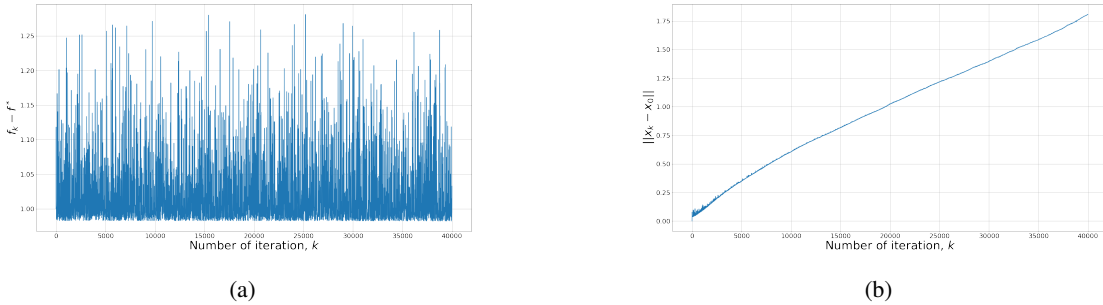
Let  $J$  be the Jacobian of the function  $g$ . Then note that  $JJ^\top$  is a tridiagonal matrix, and one can easily verify that all its minors are positive for any  $x$  (see the proof in appendix D). Whence it follows that for any compact set there exists some constant  $c$  such that  $JJ^\top \succeq cI$ . Thus, this function locally satisfies the PL-condition.

As it was seen from the results of the previous experiments, our proposed stopping criterion (31) of Algorithm 1 can work equally well for all considered types of noise in the gradient. In the current experiments, for the Nesterov-Skovok function, we used the random noise of the gradient which is uniformly distributed on the sphere. For the experiments, the starting point is  $(-1, 1, \dots 1, 1)^\top$  and therefore  $\|x_0 - x_*\| = 2$ . We vary the value of the inexactness  $\Delta$  and the dimension of the problem  $n$ .

$n$	$\Delta$	Iters	Time, ms	$\ x_N - x_0\ $	$\frac{\ \nabla f(x_N)\ }{\Delta}$	$f(x_N) - f^*$
3	$10^{-4}$	14097	230.58	1.996	2.86	$0.20 \cdot 10^{-4}$
	$10^{-3}$	2477	247.64	2.155	2.93	$0.11 \cdot 10^{-2}$
	$10^{-2}$	606	383.63	2.650	2.19	$0.87 \cdot 10^{-2}$
5	$10^{-4}$	73028	275.03	2.930	2.93	$0.30 \cdot 10^{-3}$
	$10^{-3}$	15765	292.39	3.312	2.65	$0.49 \cdot 10^{-2}$
	$10^{-2}$	6	200.87	0.036	1.45	0.98
7	$10^{-4}$	2898	316.51	0.049	2.69	0.98
	$10^{-3}$	103	164.23	0.036	2.07	0.98
	$10^{-2}$	17	104.77	0.036	1.42	0.98

**Table 5.** The results of the adaptive gradient descent for the Nesterov-Skovok function with the use of stopping criterion (31).

Table 5 shows the results of the adaptive gradient method 1 for the Nesterov-Skovok function (43). Firstly, we see that as the dimension of  $n$  increases, the difference between the required time to solve the problem for different  $\Delta$  grows significantly. Secondly, for different  $n$  with the same  $\Delta$ , the method converges to a solution with significantly different accuracy. So for  $\Delta = 10^{-4}$  the accuracy for  $n = 7$  and  $n = 3$  differs by more than 100 times. This is explained by the decrease in the constant  $\mu$  in the PL-condition as  $n$  grows. We can also note that  $\|x_N - x_0\|$  exceeds  $\|x_0 - x_*\|$  by at most 2 times. Moreover, significant upward deviations are observed for the cases when numerous iterations are made ( $n = 5$  and  $\Delta = 10^{-4}, 10^{-3}$ ). It can also be noted that even for sufficiently small values of the norm of the gradient, the accuracy by the function turns out to be quite low (which is typical for the Nesterov-Skovok function). Thus, for  $n = 7$  and  $\Delta = 10^{-4}$  we get a point with  $\|\nabla f(\hat{x})\| \approx 10^{-4}$ , but  $f(\hat{x}) - f^* \approx 0.98$ .



**Fig. 10.** The convergence of the gradient method for the Nesterov-Skovok function in the 7 dimensional space with the Antigradient inexactness in the gradient: (a) the function value  $f(x_k) - f^*$ ; (b) the distance  $\|x_k - x_0\|$ .

In Fig. 10 we can see the convergence rate of the gradient method with the Antigradient inexactness. First, we can see in Fig. 10(a) that the method approaches level 0.98 quite quickly and does not improve after that. On the other hand, from Fig. 10(b) one can see that the method without the stopping criterion moves away from the starting point quite far.

### C. Proof of Proposition B.1

Further,  $x_* \in X_*$ . Condition 1 in Proposition B.1 holds to any logistic regression function. Let us estimate the Lipschitz constant  $L$  of the gradient of function (42). It is known that  $L = \max_{x \in \mathbb{R}^n} \lambda_{\max}(\nabla^2 f(x))$ . On the other hand,  $\nabla^2 f(x) = \nabla^2 g(Ax) = A^\top H_g(z) \Big|_{z=Ax} A$ , where  $g(z) = \frac{1}{m} \sum_{i=1}^m \log(1 + \exp(z_i))$ , and  $A = [-y_1 w_1 \cdots -y_m w_m]^\top \in \mathbb{R}^{m \times n}$  and the matrix  $H_g(z)$  is the Hessian of the function  $g$  at the point  $z$ . Note that  $H_g(z)$  is a diagonal matrix with entries  $\frac{1}{m} \frac{e^{z_i}}{(e^{z_i} + 1)^2} \leq \frac{1}{4m}$ . Thus, we have the following estimate from above:

$$L \leq \|A\|_2^2 \max_z \|H_g(z)\|_2 = \frac{\|A\|_2^2}{4m}.$$

So, statement 2 in Proposition B.1 holds.

Further, let us introduce new notations. Let  $E_1$  be a subspace given by the basis  $W_B$  and  $E_2$  be a subspace orthogonal to  $E_1$ . Note that if  $k < n$ , then the dimension of  $E_2$  is at least 1. Then if there exists a minimum point  $x_*$ , then at any point from the set  $x_* + E_2 \subseteq X^*$ , the objective function takes the minimal value. Therefore, the set of the solutions is unbounded.

Now let us prove statement 3 in Proposition B.1. Note that the created matrix  $W$  has a rank  $k \leq n$ . Accordingly, all vectors  $w_i$  belong to the  $k$ -dimensional subspace  $E_1$ , given by the basis  $W_B$ . Also, for any vector  $\tilde{x} \in E_2$  from the subspace orthogonal to  $E_1$  and for any vector  $x \in \mathbb{R}^n$ , it is true that  $f(x + \tilde{x}) = f(x)$ . Thus,  $f(E_1) = f(\mathbb{R}^n)$ . Note that function (42) is bounded from below by 0, and hence  $f^* = \inf_x f(x) \geq 0 > -\infty$ . Thus, we consider a sequence  $\{x_j\}_j \in E_1$  such that  $f(x_j) \rightarrow f^*$ . Let us transform the sum, taking into account that the first  $2k$  vectors  $w_j$  are rows of the matrices  $W_B$  and  $-W_B$

$$\begin{aligned} \sum_{i=1}^m \log(1 + \exp(-y_i \langle w_i, x_j \rangle)) &\geq \sum_{i=1}^{2k} \log(1 + \exp(-y_i \langle w_i, x_j \rangle)) \\ &= \sum_{i=1}^k \log\left(\left(1 + e^{-y_i \langle w_i, x_j \rangle}\right) \left(1 + e^{y_i \langle w_i, x_j \rangle}\right)\right) \\ &= \sum_{i=1}^k \log(2 + 2\text{ch}(y_i \langle w_i, x_j \rangle)). \end{aligned}$$

From the fact that  $f^*$  is finite and from the constructed lower bound, it follows that  $|\langle a_i, x_j \rangle| \leq C, \forall j$  for some constant  $C > 0$ , i.e.  $\|W_B x_j\| \leq kC, \forall j$ . On the subspace  $E_1$ , the matrix  $W_B$  defines an invertible operator. Hence,  $\|\cdot\|_{W_B^\top W_B}$  is a norm on the subspace  $E_1$ . Therefore, in view of the equivalence of norms in a finite-dimensional space, we have that  $\|x_j\| \leq C_1 \forall j$  for some constant  $C_1 > 0$  depending only on the constants  $C, k$  and the parameters of the matrix  $W_B$ . Thus, any sequence of elements of the space in which the sequence of values of the function converges to  $f^*$  is bounded. This means that from this bounded sequence, we can extract a convergent subsequence  $\{x_{j_l}\}_l$ . The limit of this subsequence is the desired point  $x_*$ , which is a finite vector. So, statement 3 in Proposition B.1 holds.

As mentioned before, for any vector  $\tilde{x} \in E_2$  and for any vector  $x \in \mathbb{R}^n$ , the equality  $f(x + \tilde{x}) = f(x)$  holds. By constructing  $W_B$ , the dimension of  $E_2$  is at least 1. So, from this and statement 3, we have that statement 4 in Proposition B.1 holds.

### D. The Nesterov-Skovok function

Let us consider the following known Nesterov-Skovok function [11]

$$f(x) = \frac{1}{4} (1 - x_1)^2 + \sum_{i=1}^{n-1} (x_{i+1} - 2x_i^2 + 1)^2 = \sum_{i=1}^n g_i^2(x), \quad (44)$$

where  $g_1 = \frac{1}{2}(x_1 - 1), g_i = x_i - 2x_{i-1}^2 + 1, i = \overline{2, n}$ . Then the Jacobian of the system  $g(x) = 0$  is a two-diagonal matrix. On the main diagonal  $J_{11} = \frac{1}{2}$  and  $J_{ii} = 1$ , on the side diagonal  $J_{i,i-1} = -2x_{i-1}$ . Then the matrix  $JJ^\top$  is a



tridiagonal symmetric matrix of the form

$$\begin{bmatrix} \frac{1}{4} & -2x_1 & 0 & 0 & \dots & 0 & 0 \\ -2x_1 & 16x_1^2 + 1 & -4x_2 & 0 & \dots & 0 & 0 \\ 0 & -4x_2 & 16x_2^2 + 1 & -4x_3 & \dots & 0 & 0 \\ \dots & \dots & \dots & \dots & \dots & \dots & \dots \\ 0 & 0 & 0 & 0 & \dots & 16x_{n-2}^2 + 1 & -4x_{n-1} \\ 0 & 0 & 0 & 0 & \dots & -4x_{n-1} & 16x_{n-1}^2 + 1 \end{bmatrix}.$$

Then the first principal minor is  $f_1 = \frac{1}{4}$ , the second is  $f_2 = \frac{1}{4}$ . The recursive formula for a tridiagonal matrix is  $f_k = (16x_{k-1}^2 + 1)f_{k-1} - 16x_{k-1}^2 f_{k-2}$ . Then we can prove that  $f_j \geq f_{j-1}$  for all  $j > 1$ . Thus, the matrix is strictly positive definite for any  $x$ . Then for any compact set  $W$ , one can choose a constant  $c$  which limits from below all the eigenvalues of the matrix  $JJ^\top$ , which means that  $JJ^\top \succeq cI$ . Then, according to the results [10] given in the introduction, the function satisfies the PL-condition with a constant  $c > 0$  on the corresponding compact set.