# A Fully Adaptive Steepest Descent Method [*]

Z.R. Gabidullina

**Abstract**

For solving pseudo-convex global optimization problems, we present a novel fully adaptive steepest descent method (or ASDM) without any hard-to-estimate parameters. For the step-size regulation in an $\varepsilon$-normalized direction, we use the deterministic rules, which were proposed in J. Optim. Theory Appl. (2019, DOI: 10.1007/S10957-019-01585-W). We obtained the optimistic convergence estimates for the generated by ASDM sequence of iteration points. Namely, the sequence of function values of iterates has the advantage of the strict monotonic behaviour and globally converges to the objective function optimum with the sublinear rate. This rate of convergence is now known to be the best for the steepest descent method in the non-convex objectives context. Preliminary computational tests confirm the efficiency of the proposed method and low computational costs for its realization.

**keywords:** pseudoconvex function, steepest descent, normalization of descent direction, adaptive step-size, rate of convergence
**MSC classes:** 90C30, 65K05

## 1   Introduction

As is broadly known, the development of the original variant of the steepest descent method (or, briefly, SDM) was pioneered by Cauchy (1847) for solving systems of homogeneous equations. For solving the mathematical programs in the other settings, a wide spectrum of various kinds of SDM was investigated by researchers over the many years.

---

[*]Kazan Federal University, e-mail: zgabid@mail.ru, zulfiya.gabidullina@kpfu.ru

A very useful systematic survey of the existing literature related to the different variants of SDM can be found, for example, in [1, 2, 3] (see also the references therein). For an up-to-date survey of the topic, it is appropriate to refer to [2].

To avoid some conflicts and ambiguities that can be arisen in what the term "steepest descent method" means in the literature on optimization framework, we will try here to clarify some terminology. A question can now be addressed. What conditions on an optimization method ensure that the name SDM is properly used for its characterizing. The reasons why we have felt the need for such an explanation are twofold. In a wide range of optimization topics, some gradient method is usually called a steepest descent method in the case of using the uniform descent indicated by opposite to the gradient direction. Although tradition merely associates SDM not only with specific descent directions but with selecting the step-size by exact (or so-called perfect) line search. At the same time, the wide spectrum of papers on global optimization utilize the term "steepest descent method" or modifications of SDM regardless of what strategies is used for the step-size selection. In this paper, we apply in ASDM the rules of regulating the step length which are different from the exact (computationally expensive) line search.

Rather than describing all the various versions of SDM, that researchers have been constructed over the years for the achievement of the best results, we will occupy our attention only in some relevant work.

In [4],[5], there were studied the two versions of SDM for functions being twice differentiable on the Euclidean space. The second of them provides global convergence at a rate which is eventually superlinear and possibly quadratic.

In [6], the development, convergence theory and numerical testing some versions of steepest descent algorithm with adaptive step-size was presented. All of the algorithms are computationally efficient. Based on estimates of the Lipschitz constant, there was proved the convergence of the different variants of SDM to a minimizer or to a stationary point of objective function. The algorithms have been tested on real-life artificial neural network applications and the results have been very satisfactory.

To minimize a continuously differentiable quasiconvex function, SDM with Armijo's step-sizes was proposed in [7]. This method generates a sequence of iterates globally converging to a point at which

2

the gradient of the objective function is equal to zero vector.

In [8], there were derived two-point step sizes for SDM by approximating the secant equation. At the cost of storage of an extra iterate and gradient, these algorithms achieve better performance and cheaper computation than the classical SDM. By Barzilai and Borwein, there was proposed the non-monotonic variant of SDM which is superlinearly convergent for convex quadratic setting in two-dimensional space, and has quite well behaviour for the case of high-dimensional tasks.

In [9], there is used an idea that a limited memory approach might be fashioned by using a limited number of eigenvalue estimates. An improvement of characteristics of SDM has been achieved by the introduction of the Barzilai-Borwein choice of step length, and some other related ideas. There is suggested a method which is able to take advantage of the availability of a few additional 'long' vectors of storage to achieve a significant improvement in performance, both for quadratic and non-quadratic objective functions. The sequence of iterates converges to the point for which the gradient of objective function equals the zero vector.

A step-size formulae, which provides for SDM fast convergence and the monotone property, was presented in [10]. An algorithm with the new step-size in even iterations and exact line search in odd iterations is proposed. Numerical results obtained by the new method confirm that the the exact solution may be found within three iterations for the case of two-dimensional problems. For small-scale problems, the new method is very efficient. A modified version of the new method is also presented, where the new technique for choosing the step length is utilized after every two iterations with exact line searches. The modified algorithm is comparable to the Barzilai-Borwein method for large-scale tasks and better for small-scale tasks.

In [11], there is investigated a generalized hybrid steepest descent method and its convergence theory for solving monotone variational inequality over the fixed point set of a mapping which is not necessarily Lipschitz continuous. This method is used for solving the convex minimization problem for a smooth convex function whose gradient is not necessarily Lipschitzian. There is proved that the sequence of iterates converges strongly to a minimizer $x^*$.

Full convergence of the steepest descent method with inexact line searches was proved in [1]. There were considered two of such procedures and proved, for a convex objective function, convergence of the whole sequence to a minimizer without any level set boundedness

assumption and, for one of them, without any Lipschitz condition.

In [2], there was demonstrated how taking into account the spectral properties of the Hessian matrix, for convex quadratic problems, one can provide the improvement of the practical behaviour of SDM. This allows them to obtain computational results comparable with those of the Barzilai and Borwein algorithm, with the further advantage of monotonic behaviour.

In [12], there are presented results describing the properties of the gradient norm for the SDM applied to quadratic objective functions. There are also made some general observations that apply to nonlinear problems, relating the gradient norm, the objective function value, and the path generated by the iterates.

By the method from [3], there is guaranteed the well definedness of the generated sequence. Under mild assumptions on the multicriteria function, there was justified that each accumulation point (if they exist) satisfies first-order necessary conditions for Pareto optimality. Under assumptions of quasi-convexity of the multicriteria function and non-negativity of the Riemannian manifold curvature, full convergence of the sequence to a Pareto critical was proved.

The main contributions in this paper are as follows. We propose a novel fully adaptive steepest descent method with the step length regulation for solving pseudo-convex unconstrained optimization tasks. This relaxation algorithm allows one to generate the sequence of iterates $\{x_k\}, \ k = 0, 1, \ldots$ such that the sequence of its function values $\{f(x_k)\}, \ k = 0, 1, \ldots$ has the advantage of the strict monotonic behaviour and converges globally to the objective function optimum with the following rate $O(1/k)$. This rate is traditionally called the sublinear one. To the best of our knowledge, a convergence rate of $O(1/k)$ is now known to be the best for SDM in a non-convex objectives context.

The sublinear convergence rate takes place under the following conditions relating the original problem: 1) an objective function $f(x)$ is pseudo-convex on some convex set $D \subseteq \mathbb{R}^n$ (the set $D$ may coincide, for instance, with the Lebesgue set (corresponding to a starting point of the iterates sequence) of the objective function or with the whole Euclidean space and etc), 2) the function $f(x)$ is required to be satisfied to so-called Condition A introduced in [13]. We note that this condition will be defined explicitly below in Section 2 (see Definition 2.2).

Here, we underline that, for the execution of ASDM, there is no need to make use any priori information regarding the auxiliary con-

stant defined in Condition A. With respect to this fact, the presented version of ASDM compares favorably with the other variants of SDM considered above. The fulfillment of Condition A allows us further to adaptively regulate the step-size without appling any complicated line search techniques. Indeed, there is no need to use the one-dimensional exact minimization of the objective function in the selected descent direction. We apply the two deterministic rules of the step length adjustment. These strategies guarantee to determine the step-size by utilizing the finite procedures of diminishing an initial value of a certain parameter. The latter is a user-selected parameter which is diminished until a moment when the condition applied for the step length regulation becomes fulfilled. We notice that the step-size computing strategies provide the strict relaxation of the objective function at each iteration. Note that the concept of the objective function relaxation allows to interpret of the relaxation properties of minimization methods. Due to this interpretation, we can to evaluate how the objective function value is decreased at each iteration. The event of this value being diminished on the positive magnitude corresponds to the property of the strict relaxation which is described by the following inequality: $f(x_k) > f(x_{k+1})$, $k = 0, 1, \ldots$

The fully adaptive character of the presented variant of SDM is established namely by combining the simultaneous control of adapting an $\varepsilon-$normalization parameter of the descent direction as well as the step-size regulation in tandem. We establish the finiteness of all the procedures for the step-size regulation as well as the adaptation of the $\varepsilon-$normalization parameter.

Due to all the mentioned properties of ASDM, it seems that the results of the paper may be potentially useful in various applied domains covering their theoretical as well as practical aspects (in pseudo-convex programming, variational inequality problem solving, and many others (in particular, data classification techniques and neural networks simulation)).

The rest of this paper is organized as follows. In Section 2 we present some preliminaries which are necessary for our convergence rate analysis of ASDM. In Section 3, there is formulated ASDM and justified its convergence rate. In Section 4, there are drawn some conclusions.

# 2 Definitions and Preliminaries

In this paper we aim to to explore the following problem:

$$\min_{x \in \mathbb{R}^n} f(x), \tag{1}$$

where $f(x) : \mathbb{R}^n \to \mathbb{R}^1$ is a continuously differentiable pseudo-convex function satisfying the so-called Condition $A$ (introduced in [13]) on a convex set $D \subseteq \mathbb{R}^n$. To solve this problem, we propose a new efficient algorithm, which has the estimates of the rate of its convergence and allows one to adaptively regulate both the parameter of an $\varepsilon-$normalization of a descent direction and a step length.

We begin with some notations:

$$\nabla f(x) = \left( \frac{\partial f(x)}{\partial \xi_1}, \frac{\partial f(x)}{\partial \xi_2}, \ldots, \frac{\partial f(x)}{\partial \xi_n} \right)$$

is the gradient of the function $f(x)$ at the point $x = (\xi_1, \xi_2, \ldots, \xi n)$, $x_0$ stands for a starting point of the iterative consequence $\{x_k\}$, $k \in \mathbb{N}$ generated by minimizing the objective function.

Let $\| \cdot \|$ stand for the Euclidean norm of a vector in $\mathbb{R}^n$, $\langle \cdot, \cdot \rangle$ stand for the usual inner product, $f^* := \min_{x \in \mathbb{R}^n} f(x)$, $X^* := \{x \in \mathbb{R}^n : f(x) = f^*\}$, $\mathbb{N} = \{0, 1, \ldots\}$, $\mathbf{0}$ be a zero vector of $\mathbb{R}^n$, and $p_k^*$ correspond to a projection of the iterative point $x_k$ onto the set $X^*$, $k \in \mathbb{N}$. In the literature on optimization, $p_k^*$, $k \in \mathbb{N}$ are sometimes called accumulation points.

To the extent of our knowledge, the class of smooth pseudo-convex functions was pioneered by Mangasarian in [14]. The above-mentioned class represents a generalization of the family of all continuously differentiable convex functions.

**Definition 2.1** (pseudo-convexity) A function $f(x)$, which is given and continuously differentiable on an open and convex set $G$ from $\mathbb{R}^n$, is called pseudo-convex, if there is fulfilled the following implication:

$$\langle \nabla f(x), y - x \rangle \geq 0 \Rightarrow f(x) \leq f(y), \ \forall x, y \in G,$$

or equivalently,

$$f(y) < f(x) \Rightarrow \langle \nabla f(x), y - x \rangle < 0, \ \forall x, y \in G.$$

In the case of pseudo-convex functions, the necessary and sufficient

conditions of optimality are established in the following theorem.

**Theorem 2.1** *(basic first-order conditions for optimality)* ([14], p.282)
*For the point $x^* \in G$ to furnish the minimum of $f(x)$ over $G$, it is necessary and sufficient for all $x \in G$ to hold*

$$\langle \nabla f(x^*), x - x^* \rangle \geq 0.$$

**Definition 2.2** (Condition $A$) We say that a continuous function $f(x)$ satisfies Condition $A$ on the convex set $D \subseteq \mathbb{R}^n$ if there exist a non-negative symmetric function $\tau(x, y)$ and $\mu > 0$ such that

$$f(\alpha x + (1 - \alpha)y) \geq \alpha f(x) + (1 - \alpha)f(y) - \alpha(1 - \alpha)\mu\tau(x, y),$$
$$\forall x, y \in D, \ \alpha \in [0, 1].$$

For $x, y \in D \subseteq \mathbb{R}^n$, we say that some function $\tau(x, y)$ is symmetric if $\tau(x, y) = \tau(y, x), \ \tau(x, x) = 0$. Condition $A$ characterizes a sufficiently wide class of functions $A(\mu, \tau(x, y))$. It was demonstrated in [13],[15, 16] that the functions class $A(\mu, \|x - y\|^2)$, in particular, is broader than $C^{1,1}(D)$ - the commonly known class of functions whose gradient vectors have Lipschitzian property on the convex set $D \subseteq \mathbb{R}^n$. By the way, we notice that namely the Lipschitz condition for gradients of functions being minimized has been determined as the necessary assumption in justifying the theoretical estimates of the convergence rate for many modern smooth optimization methods. In [16], there were presented some examples of functions that satisfy Condition A. For functions from $A(\mu, \tau(x, y))$, there also were explored their main properties and criteria for membership in the studied class. Furthermore, for a smooth function satisfying Condition A on a convex set D, there was proved in [16] the following remarkable differential inequality:

$$f(x) - f(y) \geq \langle \nabla f(x), x - y \rangle - \mu\tau(x, y). \tag{2}$$

**Theorem 2.2** *(relation between two classes of functions)* ([17], p.1082)
*If $D$ is convex subset of $\mathbb{R}^n$, $f(x) \in C^{1,1}(D)$, then $f(x)$ satisfies Condition A on $D$ with coefficient $\mu = L/2$ and function $\tau(x, y) = \|x - y\|^2$, where $L$ is a Lipschitz constant for the gradient of $f(x)$.*

**Definition 2.3** ($\varepsilon-$normalized descent direction) For functions from the class $A(\mu, \|x - y\|^v), \ v \geq 2$, a vector $s \neq \mathbf{0}$ is called an $\varepsilon-$normalized

7

descent direction ($\varepsilon > 0$) of the function $f$ at the point $x \in D \subseteq \mathbb{R}^n$ if the following inequality is fulfilled:

$$\langle \nabla f(x), s \rangle + \varepsilon \|s\|^v \leq 0.$$

**Lemma 2.1** *($\varepsilon-$normalization) If some descent direction $s$ is not $\varepsilon-$normalized, then the vector constructed in such a way that $\bar{s} = \dfrac{ts}{\varepsilon \|s\|^v}$ is $\varepsilon-$normalized under the condition $0 < t \leq |\langle \nabla f(x), s \rangle|$.*

**Proof.** By construction, we obviously obtain the following relation:

$$\langle \nabla f(x), \bar{s} \rangle + \varepsilon \|\bar{s}\|^v = \frac{t}{\varepsilon \|s\|^v} \langle \nabla f(x), s \rangle + \frac{t^v}{\varepsilon^{v-1} \|s\|^{v(v-1)}} =$$

$$= \frac{t}{\varepsilon \|s\|^v} \left[ \langle \nabla f(x), s \rangle + t \left[ \frac{t}{\varepsilon \|s\|^v} \right]^{v-2} \right] \leq 0,$$

because $\left[ \dfrac{t}{\varepsilon \|s\|^v} \right]^{v-2} \leq 1.$ $\square$

Under the condition $v = 2$, if for some fixed point $x \in \mathbb{R}^n$ the vectors $z - x$ are $\varepsilon-$normalized descent directions at the point $x$, then it is easy to observe that all the points $z \in \mathbb{R}^n$ belong to the $n-$dimensional ball of radius $R = \|\nabla f(x)\|/2\varepsilon$ with center at the point $u = x - \nabla f(x)/2\varepsilon$.

Set

$$\zeta = \begin{cases} (\varepsilon \cdot \mu^{-1})^{1/(v-1)}, & \text{if } \varepsilon < \mu, \\ 1, & \text{if } \varepsilon \geq \mu. \end{cases}$$

We further study very useful properties of $\varepsilon-$normalized descent directions. These properties provide a strict relaxation of the objective function in gradient methods.

**Lemma 2.2** *(main properties of $\varepsilon-$normalized descent directions) ([17], p.1083) Let $s$ be some $\varepsilon-$normalized descent direction for the function $f$ at the point $x$ where $v \geq 2$, $f(x) \in A(\mu, \|x - y\|^v)$, then for all $\beta \in ]0, 1[$ there exists a constant $\hat{\lambda} = \hat{\lambda}(\beta) > 0$ ($\hat{\lambda} = (1 - \beta)^{1/(v-1)} \zeta$) such that for all $\lambda \in \left]0, \hat{\lambda}\right]$ it holds*

$$f(x) - f(x + \lambda s) \geq -\lambda \beta \cdot \langle \nabla f(x), s \rangle, \tag{3}$$

$$f(x) - f(x + \lambda s) \geq \lambda \beta \cdot \varepsilon \|s\|^v. \tag{4}$$

8

For establishing the convergence of the adaptive algorithm, which will be proposed below, there are needed the inequalities (3)–(4). In particular, they imply a strict relaxation of the objective function. Namely, it holds

$$f(x + \lambda s) < f(x), \ \forall \lambda \in (0, \ (1 - \beta)^{1/v-1}\zeta], \ \beta \in (0, 1).$$

Based on Lemma 2.2, there can be described the rules of computing the step-size satisfying (3)–(4). Suppose that $s$ is some $\varepsilon-$normalized direction of descent for $f$ at the point $x$. Additionally, let the following conditions be fulfilled:   $\beta \in \ ]0, 1[, \ \eta = (1 - \beta)^{1/v-1}, \ \hat{i} = 1, \ J(\hat{i}) = \{\hat{i}, \hat{i} + 1, \hat{i} + 2, \ldots\}$. Next we need to find $i^*$ - the least index $i \in J(\hat{i})$ for which there holds the following inequality:

$$f(x) - f(x + \eta^i s) \geq -\eta^i \beta \cdot \langle \nabla f(x), s \rangle, \tag{5}$$

or the more weak inequality:

$$f(x) - f(x + \eta^i s) \geq \eta^i \beta \cdot \varepsilon \|s\|^v. \tag{6}$$

We further set $\lambda = \eta^{i^*}$. In what follows, we have in view that there is utilized Rule 1 (or Rule 2) when we follow the first (or the second) of the above-mentioned strategies for calculating the step-size. The step length determined in accordance with these strategies satisfies (3) or (4), respectively.

Further, there should be attentionally explored the case when $s$ is the $\varepsilon-$normalized direction of descent, but it is not $\mu-$normalized (this case is possible only for $\varepsilon < \mu$). Under the assumption that $0 < \varepsilon < \mu$, for the event of choosing $\lambda$ according to Rule 1 or Rule 2, we demonstrate that the step length is bounded from below. This evidently yields that the described procedures of diminishing the step-size are finite.

**Lemma 2.3** *(exact lower estimate of the step length)([17], p.1084) If*
*(b) $f(x) \in A(\mu, \|x - y\|^v), \ v \geq 2$,*
*(c) $0 < \varepsilon < \mu, \ \beta \in \ ]0, 1[$,*
*(d) $s-$ is an $\varepsilon-$normalized descent direction of the function $f$ at the point $x$, but it is not $\mu-$normalized,*
*(e) $i^*$ is the smallest index $i = 1, 2, \ldots$, for which there is fulfilled the condition of Rule 1 or Rule 2, $\lambda = \eta^{i^*}$;*

*Then the following estimate holds:*

$$\lambda > \left(\varepsilon\mu^{-1} \cdot (1-\beta)^2\right)^{1/(v-1)} > 0.$$

**Remark 2.1** (exact lower estimate of the constant $\mu$) Due to Lemma 2.3, there comes immediately the following estimate:

$$\mu > \varepsilon \cdot (1-\beta)^2 \lambda^{1-v}. \tag{7}$$

Later, the estimate (7) will be utilized in ASDM for adaptive regulation of the parameter for the $\varepsilon-$normalization of the descent direction.

# 3  Adaptive Steepest Descent Algorithm and its Convergence

This section is aimed at providing the principles of selecting the $\varepsilon-$normalization parameter for the descent direction. We note that these principles are universal and may be utilized for developing the various adaptive algorithms with normalized descent directions (see, for instance, [17]).

The method convergence for the fixed parameter $\varepsilon$ (in the case of an arbitrary ratio of the parameter $\varepsilon$ and the value of $\mu$ in Condition A) follows from the convergence of the adaptive variant of SDM. We notice that generally saying the constant $\mu$ is unknown beforehand. In practice, the selection of the $\varepsilon$ values close to the $\mu$ value is therefore decisive for the algorithm convergence. If one chooses the too small parameter $\varepsilon$, then, according to Rule 1 and Rule 2, this may imply significant diminishing the step-size. In the case of selecting the unjustifiably large value of $\varepsilon$, the convergence of the adaptive algorithm may be slowed down. Consequently, it is expedient to evaluate the parameter $\varepsilon$ in the process of executing the algorithm. The inequalities (5)–(7) allow us to make an adjustment to the value $\varepsilon$ by increasing it if the previous choice was unsuccessful. Now, we specify further details of a procedure for pointwise adapting the parameter $\varepsilon$ during the iterative implementation of the algorithm.

For the iterate $x_k$ of the adaptive algorithm, let $\varepsilon_k > 0$ be the value of a parameter for an $\varepsilon-$normalization of descent direction. Let $s_k$ be an $\varepsilon-$normalized descent direction of the function $f$ at $x_k$; the iterative step-size $\lambda_k$ is chosen in accordance with one of the Rules 1–2. Suppose that $i_k$ is the least index $i \in J(\hat{i})$, for which there is fulfilled

the condition of selecting the iteration step-size for $x = x_k$, $\varepsilon = \varepsilon_k$, $s = s_k$. According to Lemma 2.2, if $i_k = \hat{i}$, then for implementing the next - $(k+1)-$th iteration it is expedient to remain unchanged the value of the normalization parameter, i.e. to put $\varepsilon_{k+1} = \varepsilon_k$. During the process of dropping the step-size, let the verified condition (5) (or condition (6)) be fulfilled for $i_k > \hat{i}$. Then, according to (7), there should be increased the current value of the parameter for the $\varepsilon-$normalization of the descent direction, for example, as follows: $\varepsilon_{k+1} = \varepsilon_k \cdot \zeta_k$. Regardless of what rule is chosen for computing the step-size, here we have

$$\zeta_k = (1 - \beta)^{1-i_k}. \tag{8}$$

The fact that $\mu$ is finite implies that after the finite number of increases, the value of the parameter $\varepsilon$ can exceed $\mu$ and cease to vary. Suppose that $j > 0$ is the number of iteration, on which it is fulfilled $\varepsilon_j \geq \mu$. One then has $\varepsilon_k \geq \varepsilon_j \geq \mu$, $\forall k \geq j$. In this event, after some iteration ($j \geq 0$ the adaptive algorithm starts to be implemented with the fixed constant for the $\varepsilon-$normalization of a descent direction. We emphasize that beginning with the $j-$th iteration, the step-size becomes unchanged: $\lambda_k = \eta$, $\forall k \geq j$. At that time, for computing the iterative step length, one needs only one calculation of the objective function value at the point $x_k + \eta s_k$ (for verifying the fulfillment of the condition applied for selecting the step-size).

**Algorithm**.
Step 0. Initialization. Select $x_0 \in \mathbb{R}^n$, $\beta \in ]0, 1[$, $\varepsilon_0 > 0$. Set the iteration counter $k$ to 0.
Step 1. For the objective function $f(x)$, calculate the gradient vector at $x_k$. Verify the optimality criterion:
if $\nabla f(x_k) = \mathbf{0}$, then terminate the algorithm implementation (since $x_k$ is a solution of the problem (1)). Otherwise, set $p_k = -\nabla f(x_k)$,

$$s_k = \begin{cases} p_k, & \text{if } \langle \nabla f(x_k), p_k \rangle + \varepsilon_k \|p_k\|^v \leq 0, \\ \dfrac{t_k p_k}{\varepsilon_k \|p_k\|^v} = \dfrac{p_k}{\varepsilon_k \|p_k\|^{v-2}}, & \text{othewise.} \end{cases}$$

Here $t_k = |\langle \nabla f(x_k), p_k \rangle| = \|\nabla f(x_k)\|^2$.
Step 2. Let $i_k$ be the least index $i \in J(\hat{i})$ for which there holds the condition from Rule 1 or Rule 2 when $x = x_k$, $s = s_k$, $\varepsilon = \varepsilon_k$. Set

11

$\lambda_k = \eta^{i_k}$.

Step 3. Compute the next iterate $x_{k+1} = x_k + \lambda_k s_k$.

Step 4. Update $\varepsilon_{k+1} = \zeta_k \varepsilon_k$. Set $k = k + 1$ and go to Step 1.

Clearly, there should be applied the same rule for choosing the step-size at each iteration point.

**Remark 3.1** (verification of the descent direction) Let $\bar{s}_k = \|p_k\|$. The vector $s_k$ generated by the algorithm is the $\varepsilon-$normalized descent direction. This is obvious when $s_k = \bar{s}_k$. In the case of $s_k = \dfrac{t_k \bar{s}_k}{\varepsilon_k \|\bar{s}_k\|^v}$, Lemma 2.1 confirms that $s_k$ is also $\varepsilon-$ normalized. In addition, it clearly holds $\|s_k\| \le \|\bar{s}_k\|$. Really,

$$
\|s_k\| = \begin{cases} \|\bar{s}_k\|, & \text{if } \langle \nabla f(x_k), \bar{s}_k \rangle + \varepsilon_k \|\bar{s}_k\|^v \le 0, \\ \dfrac{t_k}{\varepsilon_k \|\bar{s}_k\|^{v-1}} < \|\bar{s}_k\|, & \text{otherwise,} \end{cases}
$$

because $\dfrac{t_k}{\varepsilon_k \|\bar{s}_k\|^v} = \dfrac{-\langle \nabla f(x_k), \bar{s}_k \rangle}{\varepsilon_k \|\bar{s}_k\|^v} < 1$.

Suppose that $\{x_k\}$ is the sequence generated by ASDM. To explore the rate convergence of ASDM for the pseudo-convex setting, it is necessary to establish a criterion of global optimality for a solution. Usually, we do not know beforehand the optimum value of the function being minimized. Therefore, it is crucial to work out directly the stopping criterion of ASDM for the pseudo-convex setting. We note that the set $D$ in the formulation of the next theorem may coincide, for instance, with the Lebesgue set of the function $f(x)$ at the point $x_0 \in \mathbb{R}^n$ or with the whole space.

**Theorem 3.1** (constructive measure of optimality for ASDM) Let $f(x)$ be a continuously differentiable pseudo-convex function on some convex set $D \subseteq \mathbb{R}^n$, $X^* \subset D$. Then, for the fulfillment of the equality $f(x_k) = f^*$, $x_k \in D$, $k = 1, 2, \ldots$, it is sufficient to hold

$$\nabla f(x_k) = \mathbf{0}. \tag{9}$$

**Proof.** The equality (9) obviously yields that $\langle \nabla f(x_k), x - x_k \rangle = 0, \forall x \in D$. By definition of a pseudo-convex function, we then have $f(x_k) \le f(x)$, $\forall x \in D$. Since $X^* \subset D$, this is what we want to prove. $\square$

Further, we assume that $x_k \notin X^*$, $\forall k = 0, 1, \ldots$ With the purpose of investigating the convergence of optimization methods in the pseudo-convex setting, there is usually described in the literature on optimization an auxiliary numeric sequence $\{\theta_k\}$ as follows:

$$\theta_k > 0, \, 0 < \theta_k \cdot (f(x_k) - f(x^*)) \leq \langle \nabla f(x_k), x_k - x^* \rangle, \, x^* \in X^*, k \in \mathbb{N}. \tag{10}$$

By the definition of pseudo-convex functions, there must exist such values $\theta_k$. For instance, in the case of a continuously differentiable convex function, it holds $\theta_k = 1, k = 0, 1, \ldots$ The estimates of elements of the sequence $\{\theta_k\}$ were studied, for instance, in [13, 16].

Before formulating the theorem on the convergence of the sequence $\{x_k\}$ constructed by ASDM to a solution of problem (1), there should be reminded the following well-known fact related to convergence of some special numeric sequence.

**Lemma 3.1** *(sublinear rate of convergence for numeric sequences) ([18], p.102) If a numeric sequence $\{a_k\}$ is such that*

$$a_k \geq 0, \, a_k - a_{k+1} \geq q \cdot a_k^2, \, k = 1, \, 2, \, \ldots,$$

*where $q$ is some positive constant, then the following estimate holds:*

$$a_k \sim O(1/k),$$

*i.e. there will be found a constant*
$$q_1 > 0 \quad \text{such that} \quad 0 \leq a_k \leq q_1 \cdot k^{-1}, \, k = 1, 2, \ldots$$

The next auxiliary lemma is needed for the purpose of proving the convergence theorem. This lemma establishes the upper bound for the adapted values of the normalization parameter for each iteration.

**Lemma 3.2** *(boundedness of adapted values of the normalization parameter) ([17], p.1088) If*
*(b1) $f(x) \in A(\mu, \|x - y\|^v)$, $v \geq 2$,*
*(c1) $s_k$ is the $\varepsilon_k-$normalized descent direction,*
*(d1) $i_k$ is the least index $i \in J(\hat{i})$ for which there holds one of the conditions (5) or (6) under the assumptions that $s = s_k$, $x = x_k$, $\varepsilon = \varepsilon_k$, $\lambda_k = \eta^{i_k}$,*
*(e1) $\{x_k\}$ is some iterative sequence constructed by the rule:*
$$x_{k+1} = x_k + \lambda_k s_k, \, k \in \mathbb{N},$$

*(f1)* $\varepsilon_0 > 0$, $\varepsilon_{k+1} = \varepsilon_k \cdot (1-\beta)^{1-i_k}$, $k \in \mathbb{N}$.

*Then it is fulfilled* $\varepsilon_k \leq \bar{\varepsilon}$, $\forall k \in \mathbb{N}$, *where* $\bar{\varepsilon} = max\left\{\varepsilon_0, \dfrac{\mu}{1-\beta}\right\} > 0$.

The next two universal theorems provide the possibility of evaluating the expected decrease of the objective function value when one uses the $\varepsilon_k$−normalized direction of descent and step-size chosen according to one of the above-described rules (Rule 1 or Rule 2). Here we speak of these next two theorems as universal ones in the sense that their formulation and proof do not depend on concrete algorithms. Indeed, there are important only the facts that all of the descent directions are normalized and the step-size is regulated by Rule 1 or Rule 2. For the proof details of these theorems, we refer the interested reader to [17].

**Theorem 3.2** *(estimate of the magnitude of decreasing the objective function value when the step length is selected according to Rule 1)([17], p.1089) If*
*(b2) the conditions (b1), (c1), (e1) and (f1) of Lemma 3.2 are fulfilled,*
*(c2) $i_k$ is the smallest index $i \in J(\hat{i})$ for which there is fulfilled the condition (5) with $x = x_k$, $s = s_k$, $\varepsilon = \varepsilon_k$, $\lambda_k = \eta^{i_k}$, $\eta = (1-\beta)^{1/(v-1)}$, $\beta \in ]0, 1[$.*
   *Then there will be found a constant $\bar{C} > 0$ such that for all $k \in \mathbb{N}$ the following relation holds:*

$$f(x_k) - f(x_{k+1}) \geq -\bar{C} \cdot \langle \nabla f(x_k), s_k \rangle \geq$$
$$-\bar{C} \cdot (\langle \nabla f(x_k), s_k \rangle + \varepsilon_k \|s_k\|^v), \quad (11)$$

*where* $\bar{C} = \min\{C_1, C_2\}$, $C_1 = \beta(1-\beta)^{1/(v-1)}$,
$$C_2 = \left((1-\beta)^2 \varepsilon_0/\mu\right)^{1/(v-1)} \beta.$$

**Theorem 3.3** *(estimate of the magnitude of decreasing the objective function value when the step length is selected according to Rule 2)([17], p.1089) Let*
*(b3) the conditions (b1), (c1), (e1) and (f1) of Lemma 3.2 be fulfilled,*
*(c3) the values of the iterative step-size $\lambda_k$, $\forall k \in \mathbb{N}$ be determined using (6). Then there will be found a constant $\bar{C} > 0$ such that for all*

14

$k \in \mathbb{N}$ *the inequality (11) holds with the following constants:*

$$\bar{C} = \min\{C_1, C_2\}, \; C_1 = \beta(1 - \beta)^{1/(v-1)},$$
$$C_2 = \left(\varepsilon_0 \mu^{-1}(1 - \beta)^2\right)^{1/(v-1)} \left(1 + \varepsilon_0 \mu^{-1}(1 - \beta)\beta^{-1}\right)^{-1}.$$

**Theorem 3.4** *(sublinear rate of convergence of ASDM)  If*
*(b4)  $f(x)$ is a continuously differentiable pseudo-convex function on the convex set  $D \subseteq \mathbb{R}^n$ ($X^* \subset D$) satisfying Condition A with some constant  $\mu$  and a function  $\tau(x, y) = \|x - y\|^v, \; v \geq 2$,*
*(c4)  a numeric sequence  $\{\theta_k\}$, which is defined by (10), satisfies the condition:  $\exists \theta > 0$  such that  $\theta_k \geq \theta, \; \forall k$,*
*(d4)  there exists a constant  $\gamma > 0$  such that  $\|\nabla f(x)\| \leq \gamma < \infty$, $\forall x \in D$,*
*(e4)  the Lebesgue set of the function  $f(x)$  at the point  $x_0 \in \mathbb{R}^n$, which is denoted by  $M(f, x_0) := \{x \in \mathbb{R}^n : f(x) \leq f(x_0)\}$, is bounded,*
*(g4)  a step-size  $\lambda_k, \; k \in \mathbb{N}$  is chosen according to one of the rules (Rule 1 or Rule 2).*
*Then the sequence  $\{x_k\}, \; k \in \mathbb{N}$  converges weakly , i.e.*

$$f(x_k) - f^* \sim O(1/k),$$

*or equivalently, there exists a constant  $C_3 > 0$  such that it holds*

$$f(x_k) - f^* \leq C_3 \cdot k^{-1}.$$

**Proof.** Clearly, $f(p_k^*) = f^*$, $f(x_k) > f(p_k^*)$, $\forall k \in \mathbb{N}$. By definition of pseudo-convex functions, it is fulfilled $\langle \nabla f(x_k), p_k^* - x_k \rangle < 0$. According to the assertions of Theorems 3.2–3.3, regardless of whether Rule 1 will be selected for calculating the step-size or Rule 2, there may be found a constant $\bar{C} > 0$ such that the inequality (11) holds for all $k \in \mathbb{N}$. Choose the subset of indices $\mathbb{N}_1 \subset \mathbb{N}$ such that $s_k = p_k$, $k \in \mathbb{N}_1$. One then has $s_k = \dfrac{p_k}{\varepsilon_k \|p_k\|^{v-2}}$ $\forall k \in \mathbb{N}_2 = \mathbb{N} \backslash \mathbb{N}_1$. For all $k \in \mathbb{N}_1$, we have the estimate:

$$f(x_k) - f(x_{k+1}) \geq -\bar{C} \cdot \langle \nabla f(x_k), s_k \rangle = \bar{C} \cdot \|\nabla f(x_k)\|^2.$$

From Lemma 3.2, due to (11), for all $k \in \mathbb{N}_2$, it follows the relation

$$f(x_k) - f(x_{k+1}) \geq -\bar{C} \cdot \langle \nabla f(x_k), s_k \rangle =$$
$$\frac{-\bar{C}}{\varepsilon_k \|p_k\|^{v-2}} \cdot \langle \nabla f(x_k), p_k \rangle \geq \frac{\bar{C}}{\bar{\varepsilon}\gamma^{v-2}} \|\nabla f(x_k)\|^2.$$

Thus, for all $k \in \mathbb{N}$, one has arrived at the inequality

$$f(x_k) - f(x_{k+1}) \geq \tilde{C} \cdot \|\nabla f(x_k)\|^2, \tag{12}$$

where $\tilde{C} = \bar{C} \min \left\{ 1, \dfrac{1}{\bar{\varepsilon} \gamma^{\upsilon-2}} \right\}$. By virtue of the condition (e4),

$$diam\, M(f, x_0) = \sup\{\|x - y\|, \, x, \, y \in M(f, x_0)\} = \bar{\eta} < +\infty.$$

Therefore, we immediately have the following estimate

$$\|\nabla f(x_k)\|^2 \geq \|\nabla f(x_k)\|^2 \|x_k - p_k^*\|^2 \cdot \frac{1}{\bar{\eta}^2} \geq \frac{1}{\bar{\eta}^2} \cdot \langle \nabla f(x_k), x_k - p_k^* \rangle^2.$$

Consequently, for all $k \in \mathbb{N}$ from $(10), (12)$, and $(c4)$ it follows

$$f(x_k) - f(x_{k+1}) \geq \frac{\tilde{C}}{\bar{\eta}^2} \cdot \langle \nabla f(x_k), x_k - p_k^* \rangle^2 = C_3 \cdot (f(x_k) - f(p_k^*))^2, \tag{13}$$

where $C_3 = \theta^2 \dfrac{\tilde{C}}{\bar{\eta}^2}$. Due to Lemma 3.1, the latter implies that the sequence $\{x_k\}$, $k \in \mathbb{N}$ is weakly convergent to a solution of (1), since there holds the following estimate for the convergence rate:

$$f(x_k) - f^* \leq C_3^{-1} k^{-1}. \qquad \square$$

**Remark 3.2** *There is no a need to strictly require the fulfillment of the condition (b4) of Theorem 3.4 in the whole space $\mathbb{R}^n$. For instance, it is sufficient to have that $D = M(f, x_0)$, where $x_0$ is a starting point for ASDM. Besides, there may simply be chosen some convex set for which it holds $X^* \subset D$.*

Preliminary computational tests confirm the efficiency of the proposed method and the strict monotone property of the used step-size rules. These tests show that the results of minimization depends on the user-selected parameters of ASDM such as $\beta$ and $\varepsilon$. We also observe that at each iteration (after the first one) there is needed only one function and gradient evaluation. Our preliminary experiments has demonstrated the ability of ASDM to lead to the minimum neighborhood at low computational costs.

# 4 Conclusions

We proposed a fully adaptive variant of the steepest descent method. There are used some novel rules for the calculation of the step length in which the iteration step-size is regulated additionally by an adaptation of the $\varepsilon-$normalization parameter for the descent direction. The finiteness of all the procedures of adaptive controlling both the parameter of an $\varepsilon-$normalization of a descent direction and a step length was established. For the problem of unconstrained minimizing a smooth pseudo-convex function, we justified the sublinear rate of the convergence for the adaptive variant of the steepest descent method.

One of the motivating ideas was that of using in the future the adaptive steepest descent method to solve the problems of sets separation (by minimizing the error function) and related problems of data mining (in particular, neural network classification of data).

# References

[1] Burachik, R., Graña Drummond, L.M., Iusem, A.N., Svaiter,B.F.: Full convergence of the steepest descent method with inexact line searches. Optimization 32(2), 137-146 (1995)

[2] De Asmundis, R., Di Serafino, D., Riccio F., Toraldo, G.: On spectral properties of steepest descent methods. IMA Journal of Numerical Analysis 33, 1416-1435 (2013) doi:10.1093/imanum/drs056

[3] Bento, G.C., Ferreira, O.P., Oliveira, P.R.: Unconstrained steepest descent method for multicriteria optimization on riemannian manifolds. Journal of Optimization Theory and Applications 154(1), 88-107 (2012) doi:10.1007/s10957-011-9984-2

[4] Goldstein, A.A.: On steepest descent. SIAM J. on Control, A 3(1), 147-151 (1965)

[5] Goldstein, A.A., Price, J.F.: An effective algorithm for minimization. Numer. Math. 10, 184-189 (1967) doi:10.1007/BF02162162

[6] Vrahatis, M.N., Androulakis,G.S., Lambrinos,J.N., and Magoulas,G.D.: A class of gradient unconstrained minimization algorithms with adaptive stepsize. J. Comput. and Appl. Math. 114, 367-386 (2000)

[7] Kiwiel, K.C., Murty, K.: Convergence of the steepest descent method for minimizing quasiconvex functions Journal of Optimization Theory and Applications 89(1), 221-226 (1996)

[8] Barzilai, J., Borwein J.M.: Two point step size gradient methods. IMA J. Numer. Anal. 8,141-148 (1988)

[9] Fletcher, R.: A limited memory steepest descent method. Math. Program. 135, 413-436 (2012) doi:10.1007/s10107-011-0479-6

[10] Ya-xiang Yuan: A new stepsize for the steepest descent method. Journal of Computational Mathematics 24(2), 149-156 (2006)

[11] Sahu, D.R., Yao, J.C.: A generalized hybrid steepest descent method and applications. Nonlinear Var. Anal. 1(1), 111-126 (2017)

[12] Nocedal, J., Sartenaer, A., Zhu, C.: On the behavior of the gradient norm in the steepest descent method. Computational Optimization and Applications 22, 5-35 (2002) doi: 10.1023/A:1014897230089

[13] Gabidullina, Z.R.: Relaxation methods with step regulation for solving constrained optimization problems. Journal of Mathematical Sciences 73(5), 538-543 (1995)

[14] Mangasarian, O.L.: Pseudo-convex functions. J.Soc. Industr.and Appl. Math.- Ser.A Control 3, 281-290 (1965)

[15] Gabidullina, Z.R.: Convergence of the constrained gradient method for a class of nonconvex functions. Journal of Soviet Mathematics 50(5), 1803-1809 (1990)

[16] Gabidullina, Z.R.: Adaptive methods with step length regulation for solving pseudo-convex programming problems. Dissertation for the Degree of Candidate of Science in Physics and Mathematics. Kazan (1994)

[17] Gabidullina, Z.R.: Adaptive Conditional Gradient Method. Journal of Optimization Theory and Applications 183(3), 1077-1098 (2019)

[18] Vasil'ev, F.P.: Numerical Methods for Solving Extremum Problems, Nauka, Moscow (1980)