# Analogues of Switching Subgradient Schemes for Relatively Lipschitz-Continuous Convex Programming Problems

Alexander A. Titov[1], Fedor S. Stonyakin[1,2],
Mohammad S. Alkousa[1], Seydamet S. Ablaev[2], and
Alexander V. Gasnikov[1]

[1]Moscow Institute of Physics and Technology, Moscow, Russia,
email: a.a.titov@phystech.edu, gasnikov@yandex.ru, mohammad.math84@gmail.com.
[2]V. I. Vernadsky Crimean Federal University, Simferopol, Russia,
email: fedyor@mail.ru.

May 7, 2021

## Abstract

Recently some specific classes of non-smooth and non-Lipschitz convex optimization problems were selected by Yu. Nesterov along with H. Lu. We consider convex programming problems with similar smoothness conditions for the objective function and functional constraints. We introduce a new concept of an inexact model and propose some analogues of switching subgradient schemes for convex programming problems for the relatively Lipschitz-continuous objective function and functional constraints. Some class of online convex optimization problems is considered. The proposed methods are optimal in the class of optimization problems with relatively Lipschitz-continuous objective and functional constraints.

**Keywords**: Convex Programming Problem, Switching Subgradient Scheme, Relative Lipschitz-Continuity, Inexact Model, Stochastic Mirror Descent, Online Optimization Problem.

1

# Introduction

Different relaxations of the classical smoothness conditions for functions are interesting for a large number of modern applied optimization problems. In particular, in [2] there were proposed conditions of the relative smoothness of the objective function, which mean the replacement of the classic Lipschitz condition by the following weak version

$$f(y) \leq f(x) + \langle \nabla f(x), y - x \rangle + LV_d(y, x), \qquad (0.1)$$

for any $x$ and $y$ from the domain of the objective function $f$ and some $L > 0$, $V_d(y, x)$ is an analogue of the distance between the points $x$ and $y$ (often called the *Bregman divergence*). Such a distance is widely used in various fields of science, particularly in optimization. Usually, the *Bregman divergence* is defined on the base of the auxiliary 1-strongly convex and continuously-differentiable function $d : Q \subset \mathbb{R}^n \to \mathbb{R}$ (*distance generating function*) as follows

$$V_d(y, x) = d(y) - d(x) - \langle \nabla d(x), y - x \rangle \quad \forall x, y \in Q, \qquad (0.2)$$

where $Q$ is a convex closed set, $\langle \cdot, \cdot \rangle$ is a scalar product in $\mathbb{R}^n$. In particular, for the standard Euclidean norm $\| \cdot \|_2$ and the Euclidean distance in $\mathbb{R}^n$, we can assume that $V_d(y, x) = d(y - x) = \frac{1}{2}\|y - x\|_2^2$ for arbitrary $x, y \in Q$. However, in many applications, it often becomes necessary to use non-Euclidean norms. Moreover, the considered condition of relative smoothness in [2, 18] implies only the convexity (but not strong convexity) of the distance generating function $d$. As shown in [18], the concept of relative smoothness makes it possible to apply a variant of the gradient method to some problems which were previously solved only by using interior-point methods. In particular, we talk about the well-known problem of the construction of an optimal ellipsoid which covers a given set of points. This problem is important in the field of statistics and data analysis.

A similar approach to the Lipschitz property and non-smooth problems was proposed in [19] (see also [27]). This approach is based on an analogue of the Lipschitz condition for the objective function $f : Q \to \mathbb{R}$ with Lipschitz constant $M_f > 0$, which involves replacing the boundedness of the norm of the subgradient, i.e. $\|\nabla f(x)\|_* \leq M_f$, with the so-called *relative Lipschitz condition*

$$\|\nabla f(x)\|_* \leqslant \frac{M_f \sqrt{2V_d(y, x)}}{\|y - x\|} \quad \forall x, y \in Q, \ y \neq x,$$

where $\|\cdot\|_*$ denotes the conjugate norm, see Section 1. below. Moreover, the distance generating function $d$ is not necessarily strongly convex. In [19] there were proposed deterministic and stochastic Mirror Descent algorithms for optimization problems with convex relatively Lipschitz-continuous objective functionals. Note that the applications of the relative Lipschitz-continuity to the well-known classical support vector machine (SVM) problem and to the problem of minimizing the maximum of convex quadratic functions (intersection of $n$ ellipsoids problem) were discussed in [19].

In this paper we propose a new concept of an inexact model for objective functional and functional constraint. More precisely, we introduce some analogues of the concepts of an inexact oracle [9] and an inexact model [32] for objective functionals. However, unlike [9, 32], we do not generalize the smoothness condition. We relax the Lipschitz condition and consider a recently proposed generalization of *relative Lipschitz-continuity* [19, 27]. We propose some optimal Mirror Descent methods, in different settings of Relatively Lipschitz-continuous convex optimization problems.

The Mirror Descent method originated in the works of A. Nemirovski and D. Yudin more than 30 years ago [24, 25] and was later analyzed in [5]. It can be considered as the non-Euclidean extension of subgradient methods. The method was used in many applications [22, 23, 31]. Standard subgradient methods employ the Euclidean distance function with a suitable step-size in the projection step. The Mirror Descent extends the standard projected subgradient methods by employing a nonlinear distance function with an optimal step-size in the nonlinear projection step [21]. The Mirror Descent method not only generalizes the standard subgradient descent method, but also achieves a better convergence rate and it is applicable to optimization problems in Banach spaces, while the subgradient descent is not [10]. Also, in some works [4, 25] there was proposed an extension of the Mirror Descent method for constrained problems.

Also, in recent years, online convex optimization (OCO) has become a leading online learning framework, due to its powerful modeling capability for a lot of problems from diverse domains. OCO plays a key role in solving problems where statistical information is being updated [15, 16]. There are a lot of examples of such problems: Internet network, consumer data sets or financial market, machine learning applications, such as adaptive routing in networks, dictionary learning, classification and regression (see [33] and references therein). In recent years, methods for solving online opti-

mization problems have been actively developed, in both deterministic and stochastic settings [7, 14, 17, 28]. Among them one can mention the Mirror Descent method for the deterministic setting of the problem [29, 30] and for the stochastic case [1, 12, 34], which allows to solve problems for an arbitrary distance function.

This paper is devoted to Mirror Descent methods for convex programming problems with a relatively Lipschitz-continuous objective function and functional constraints. It consists of an introduction and 6 main sections. In Section 1. we consider the problem statement and define the concept of an inexact $(\delta, \phi, V)$–model for the objective function. Also we propose some modifications of the Mirror Descent method for the concept of Model Generality. Section 2. is devoted to some special cases of problems with the properties of relative Lipschitz continuity, here we propose two versions of the Mirror Descent method in order to solve the problems under consideration. In Sections 3. and 4. we consider the stochastic and online (OCO) setting of the optimization problem respectively. In Section 5. one can find numerical experiments which demonstrate the efficiency of the proposed methods.

The contribution of the paper can be summarized as follows:

- Continuing the development of Yurii Nesterov's ideas in the direction of the relative smoothness and non-smoothness [27] there was introduced the concept of an inexact $(\delta, \phi, V)$–model of the objective function. For the proposed model we proposed some variants of the well-known Mirror Descent method, which provides an $(\varepsilon + \delta)$–solution of the optimization problem, where $\varepsilon$ is the controlled accuracy. There was considered the applicability of the proposed method to the case of the stochastic setting of the considered optimization problem.

- We also considered a special case of the relative Lipschitz condition for objective function. The proposed Mirror Descent algorithm was specified for the case of such functions. Furthermore, there was introduced one more modification of the algorithm with another approach to the step selection. There was also considered the possibility of applying the proposed methods to the case of several functional constraints.

- We considered an online optimization problem and proposed the modification of the Mirror Descent algorithm for such a case. Moreover, there were conducted some numerical experiments which demonstrate the effectiveness of the proposed methods.

# 1. Inexact Model for Relative Non-Smooth Functionals and Mirror Descent Algorithm

Let $(E, \| \cdot \|)$ be a normed finite-dimensional vector space and $E^*$ be the conjugate space of $E$ with the norm:

$$\|y\|_* = \max_x \{\langle y, x \rangle, \|x\| \leq 1\},$$

where $\langle y, x \rangle$ is the value of the continuous linear functional $y$ at $x \in E$.

Let $Q \subset E$ be a (simple) closed convex set. Consider two subdifferentiable functions $f, g : Q \to \mathbb{R}$. In this paper we consider the following optimization problem

$$f(x) \to \min_{x \in Q, \, g(x) \leq 0}. \tag{1.3}$$

Let $d : Q \to \mathbb{R}$ be any convex (not necessarily strongly-convex) differentiable function, we will call it the *reference function*. Suppose we have a constant $\Theta_0 > 0$, such that $d(x^*) \leq \Theta_0^2$, where $x^*$ is a solution of (1.3). Note that if there is a set, $X_* \subset Q$, of optimal points for the problem (1.3), we may assume that

$$\min_{x^* \in X_*} d(x^*) \leq \Theta_0^2.$$

Let us introduce some generalization of the concept of Relative Lipschitz continuity [27]. Consider one more auxiliary function $\phi : \mathbb{R} \to \mathbb{R}$, which is strictly increasing and $\phi(0) = 0$. Clearly, due to the strict monotonicity of $\phi(\cdot)$, there exists the inverse function $\phi^{-1}(\cdot)$.

**Definition 1.1.** *Let $\delta > 0$. We say that $f$ and $g$ admit the $(\delta, \phi, V)$–model at the point $y \in Q$ if*

$$f(x) + \psi_f(y, x) \leq f(y), \quad -\psi_f(y, x) \leq \phi_f^{-1}(V_d(y, x)) + \delta \tag{1.4}$$

$$g(x) + \psi_g(y, x) \leq g(y), \quad -\psi_g(y, x) \leq \phi_g^{-1}(V_d(y, x)) + \delta, \tag{1.5}$$

*where $\psi_f(y, x)$ and $\psi_g(y, x)$ are convex functions on $y$ and $\psi_f(x, x) = \psi_g(x, x) = 0$ for all $x \in Q$.*

Let $h > 0$. For problems with a $(\delta, \phi, V)$–model, the proximal mapping operator (Mirror Descent step) is defined as follows

$$Mirr_h(x, \psi) = \arg\min_{y \in Q} \left\{ \psi(y, x) + \frac{1}{h} V_d(y, x) \right\}.$$

The following lemma describes the main property of this operator.

**Lemma 1.1** (Main Lemma). *Let $f$ be a convex function, which satisfies (1.4), $h > 0$ and $x^+ = hMirr_h(x, \psi_f)$. Then for any $y \in Q$*

$$h(f(x) - f(y)) \leq -h\psi_f(y, x) \leq \phi_f^*(h) + V_d(y, x) - V_d(y, x^+) + h\delta,$$

*where $\phi_f^*$ is the conjugate function of $\phi_f$.*

*Proof.* From the definition of $x^+$

$$x^+ = hMirr_h(x, \psi_f) = \arg\min_{y \in Q} \{h\psi_f(y, x) + V_d(y, x)\},$$

for any $y \in Q$, we have

$$h\psi_f(y, x) - h\psi_f(x^+, x) + \langle \nabla d(x^+) - \nabla d(x), y - x^+ \rangle \geq 0.$$

Further, $h(f(x) - f(y)) \leq -h\psi_f(y, x) \leq$

$$\leq -h\psi_f(x^+, x) + \langle \nabla d(x^+) - \nabla d(x), y - x^+ \rangle$$
$$= -h\psi_f(x^+, x) + V_d(y, x) - V_d(y, x^+) - V_d(x^+, x) + h\delta$$
$$\leq h\phi_f^{-1}(V_d(x^+, x)) + V_d(y, x) - V_d(y, x^+) - V_d(x^+, x) + h\delta$$
$$\leq \phi_f^*(h) + \phi_f(\phi_f^{-1}(V_d(x^+, x))) + V_d(y, x) - V_d(y, x^+) - V_d(x^+, x) + h\delta$$
$$= \phi_f^*(h) + V_d(x^+, x) + V_d(y, x) - V_d(y, x^+) - V_d(x^+, x) + h\delta$$
$$= \phi_f^*(h) + V_d(y, x) - V_d(y, x^+) + h\delta.$$

$\square$

For problem (1.3) with an inexact $(\delta, \phi, V)$–model, we consider a Mirror Descent algorithm, listed as Algorithm 1 below. For this proposed algorithm, we will call step $k$ productive if $g(x^k) \leq \varepsilon$, and non-productive if the reverse inequality $g(x^k) > \varepsilon$ holds. Let $I$ and $|I|$ denote the set of indexes of productive steps and their number, respectively. Similarly, we use the notation $J$ and $|J|$ for non-productive steps.

Let $x^*$ denote the exact solution of the problem (1.3). The next theorem provides the complexity and quality of the proposed Algorithm 1.

**Theorem 1.1** (Modified MDA for Model Generality). *Let $f$ and $g$ be convex functionals, which satisfy (1.4), (1.5) respectively and $\varepsilon > 0, \delta > 0$ be fixed positive numbers. Assume that $\Theta_0 > 0$ is a known constant such that $d(x^*) \leq \Theta_0^2$. Then, after the stopping of Algorithm 1, the following inequalities hold:*

$$f(\widehat{x}) - f(x^*) \leq \varepsilon + \delta \quad and \quad g(\widehat{x}) \leq \varepsilon + \delta.$$

**Algorithm 1** Modified MDA for $(\delta, \phi, V)$–model.

**Require:** $\varepsilon > 0, \delta > 0,\ h^f > 0, h^g > 0, \Theta_0 : d(x^*) \le \Theta_0^2.$

1: $x^0 = \arg\min_{x \in Q} d(x).$

2: $I =: \emptyset$ and $J =: \emptyset$

3: $N \leftarrow 0$

4: **repeat**

5:    **if** $g\left(x^N\right) \le \varepsilon + \delta$ **then**

6:       $x^{N+1} = Mirr_{h^f}\left(x^N, \psi_f\right),$    "productive step"

7:       $N \to I$

8:    **else**

9:       $x^{N+1} = Mirr_{h^g}\left(x^N, \psi_g\right),$    "non-productive step"

10:       $N \to J$

11:    **end if**

12:    $N \leftarrow N + 1$

13: **until** $\Theta_0^2 \le \varepsilon\left(|J|h^g + |I|h^f\right) - |J|\phi_g^*(h^g) - |I|\phi_f^*(h^f).$

**Ensure:** $\widehat{x} := \frac{1}{|I|}\sum_{k \in I} x^k.$

---

*Proof.* By Lemma 1.1, we have for all $k \in I$ and $y \in Q$

$$h^f\left(f(x^k) - f(y)\right) \le \phi_f^*(h^f) + V_d(y, x^k) - V_d(y, x^{k+1}) + h^f\delta, \qquad (1.6)$$

Similarly, for all $k \in J$ and $y \in Q$

$$h^g\left(g(x^k) - g(y)\right) \le \phi_g^*(h^g) + V_d(y, x^k) - V_d(y, x^{k+1}) + h^g\delta, \qquad (1.7)$$

Taking summation, in each side of (1.6) and (1.7), over productive and non-productive steps, we get

$$\sum_{k \in I} h^f\left(f(x^k) - f(x^*)\right) + \sum_{k \in J} h^g\left(g(x^k) - g(x^*)\right) \le$$

$$\le \sum_{k \in I} \phi_f^*(h^f) + \sum_{k \in J} \phi_g^*(h^g) + \sum_k \left(V_d(x^*, x^k) - V_d(x^*, x^{k+1})\right) + \sum_{k \in I} h^f\delta + \sum_{k \in J} h^g\delta \le$$

$$\sum_{k \in I} \phi_f^*(h^f) + \sum_{k \in J} \phi_g^*(h^g) + \Theta_0^2 + \sum_{k \in I} h^f\delta + \sum_{k \in J} h^g\delta.$$

Since for any $k \in J$, $g(x^k) - g(x^*) > \varepsilon + \delta$, we have

$$\sum_{k \in I} h^f\left(f(\widehat{x}) - f(x^*)\right) \le \sum_{k \in I} \phi_f^*(h^f) + \sum_{k \in J} \phi_g^*(h^g) + \Theta_0^2 - \varepsilon\sum_{k \in J} h^g + \sum_{k \in I} h^f\delta =$$

7

$$= |I|\left(\phi_f^*(h^f) + \delta h^f\right) + |J|\phi_g^*(h^g) - |J|h^g\varepsilon + \Theta_0^2 \le \varepsilon|I|h^f + \delta|I|h^f.$$

So, for $\widehat{x} := \frac{1}{|I|}\sum_{k\in I} x^k$, after the stopping criterion of Algorithm 1 is satisfied, the following inequalities hold

$$f(\widehat{x}) - f(x^*) \le \varepsilon + \delta \quad \text{and} \quad g(\widehat{x}) \le \varepsilon + \delta.$$

$\square$

## 2. The Case of Relatively Lipschitz-Continuous Functionals

Suppose hereinafter that the objective function $f$ and the constraint $g$ satisfy the so-called relative Lipschitz condition, with constants $M_f > 0$ and $M_g > 0$, i.e. the functions $\phi_f^{-1}$ and $\phi_g^{-1}$ from (1.4) and (1.5) are modified as follows:

$$\phi_f^{-1}\left(V_d(y,x)\right) = M_f\sqrt{2V_d(y,x)}, \tag{2.8}$$

$$\phi_g^{-1}\left(V_d(y,x)\right) = M_g\sqrt{2V_d(y,x)} \tag{2.9}$$

Note that the functions $f,g$ must still satisfy the left inequalities in (1.4),(1.5):

$$f(x) + \psi_f(y,x) \le f(y), \quad -\psi_f(y,x) \le M_f\sqrt{2V_d(y,x)} + \delta \tag{2.10}$$

$$g(x) + \psi_g(y,x) \le g(y), \quad -\psi_g(y,x) \le M_g\sqrt{2V_d(y,x)} + \delta, \tag{2.11}$$

For this particular case we say that $f$ and $g$ admit the $(\delta, M_f, V)$– and $(\delta, M_g, V)$–model at each point $x \in Q$ respectively. The following Remark 2 provides the explicit form of $\phi_f, \phi_g$ and their conjugate functions $\phi_f^*, \phi_g^*$.

**Remark 2.1.** Let $M_f > 0$ and $M_g > 0$. Then functions $\phi_f$ and $\phi_g$ which correspond to (2.8) and (2.9) are defined as follows:

$$\phi_f(t) = \frac{t^2}{2M_f^2}, \quad \phi_g(t) = \frac{t^2}{2M_g^2}.$$

Their conjugate functions have the following form:

$$\phi_f^*(y) = \frac{y^2 M_f^2}{2}, \tag{2.12}$$

$$\phi_g^*(y) = \frac{y^2 M_g^2}{2}. \tag{2.13}$$

For the case of relatively Lipschitz-continuous objective function and constraint, we consider a modification of Algorithm 1, the modified algorithm is listed as Algorithm 2, below. The difference between Algorithms 1 and 2 is represented in the control of productivity and the stopping criterion.

---

**Algorithm 2** Mirror Descent for Relatively Lipschitz-continuous functions, version 1

---
:

**Require:** $\varepsilon > 0, \delta > 0, M_f > 0, M_g > 0, \Theta_0 : d(x^*) \le \Theta_0^2$
1: $x^0 = \arg\min_{x \in Q} d(x)$.
2: $I =: \emptyset$
3: $N \leftarrow 0$
4: **repeat**
5:     **if** $g\left(x^N\right) \le M_g \varepsilon + \delta$ **then**
6:        $h^f = \frac{\varepsilon}{M_f}$,
7:        $x^{N+1} = Mirr_{h^f}\left(x^N, \psi_f\right)$,     "productive step"
8:        $N \to I$
9:     **else**
10:        $h^g = \frac{\varepsilon}{M_g}$,
11:        $x^{N+1} = Mirr_{h^g}\left(x^N, \psi_g\right)$,     "non-productive step"
12:     **end if**
13:     $N \leftarrow N + 1$
14: **until** $N \ge \frac{2\Theta_0^2}{\varepsilon^2}$.
**Ensure:** $\widehat{x} := \frac{1}{|I|} \sum\limits_{k \in I} x^k$.

---

For the proposed Algorithm 2, we have the following theorem, which provides an estimate of its complexity and the quality of the solution of the problem.

**Theorem 2.1.** *Let $f$ and $g$ be convex functions, which satisfy (2.10) and (2.11) for $M_f > 0$ and $M_g > 0$.*

*Let $\varepsilon > 0, \delta > 0$ be fixed positive numbers. Assume that $\Theta_0 > 0$ is a known constant such that $d(x^*) \le \Theta_0^2$. Then, after the stopping of Algorithm 2, the following inequalities hold:*

$$f(\widehat{x}) - f(x^*) \le M_f \varepsilon + \delta \quad and \quad g(\widehat{x}) \le M_g \varepsilon + \delta.$$

*Proof.* By Lemma 1.1, we have

$$\sum_{k\in I} h^f \left( f(x^k) - f(x^*) \right) + \sum_{k\in J} h^g \left( g(x^k) - g(x^*) \right) \le \sum_{k\in I} \phi_f^*(h^f) + \sum_{k\in J} \phi_g^*(h^g) +$$
$$+ \Theta_0^2 + \sum_{k\in I} h^f \delta + \sum_{k\in J} h^g \delta$$

Since for any $k \in J$, $g(x^k) - g(x^*) > M_g \varepsilon + \delta$ we have

$$\sum_{k\in I} h^f \left( f(\widehat{x}) - f(x^*) \right) \le \sum_{k\in I} \phi_f^*(h^f) + \sum_{k\in J} \phi_g^*(h^g) + \Theta_0^2 - M_g \varepsilon \sum_{k\in J} h^g + \sum_{k\in I} h^f \delta$$
$$= |I|(\phi_f^*(h^f) + \delta h^f) + |J|\phi_g^*(h^g) - |J|\varepsilon^2 + \Theta_0^2.$$

Taking into account the explicit form of the conjugate functions (2.12), (2.13) one can get:

$$h^f \left( f(\widehat{x}) - f(x^*) \right) \le |I| \left( \frac{M_f^2 h^{f^2}}{2} + \delta h^f \right) + |J| \frac{M_g^2 h^{g^2}}{2} - |J|\varepsilon^2 + \Theta_0^2$$
$$= |I| \left( \frac{\varepsilon^2}{2} + \delta h^f \right) + |J| \frac{\varepsilon^2}{2} - |J|\varepsilon^2 + \Theta_0^2$$
$$\le M_f \varepsilon |I| h^f + \delta |I| h^f,$$

supposing that the stopping criterion is satisfied.

So, for the output value of the form $\widehat{x} = \frac{1}{|I|} \sum_{k\in I} x^k$, the following inequalities hold:

$$f(\widehat{x}) - f(x^*) \le M_f \varepsilon + \delta \quad \text{and} \quad g(\widehat{x}) \le M_g \varepsilon + \delta.$$

$\square$

Also, for the case of relatively Lipschitz-continuous objective function and constraint, we consider another modification of Algorithm 1, which is listed as the following Algorithm 3. Note that the difference lies in the choice of steps $h^f, h^g$ and the stopping criterion.

By analogy with the proof of Theorem 2 one can obtain the following result concerning the quality of the convergence of the proposed Algorithm 3.

**Theorem 2.2.** *Let $f$ and $g$ be convex functions, which satisfy (2.10) and (2.11) for $M_f > 0$ and $M_g > 0$. Let $\varepsilon > 0, \delta > 0$ be fixed positive numbers. Assume that $\Theta_0 > 0$ is a known constant such that $d(x^*) \le \Theta_0^2$.*

**Algorithm 3** Mirror Descent for Relatively Lipschitz-continuous functions, version 2.

---

**Require:** $\varepsilon > 0, \delta > 0, M_f > 0, M_g > 0, \Theta_0 : d(x^*) \leq \Theta_0^2$.

1: $x^0 = \arg\min_{x \in Q} d(x)$.

2: $I =: \emptyset$ and $J =: \emptyset$

3: $N \leftarrow 0$

4: **repeat**

5:     **if** $g\left(x^N\right) \leq \varepsilon + \delta$ **then**

6:        $h^f = \frac{\varepsilon}{M_f^2}$,

7:        $x^{k+1} = Mirr_{h^f}\left(x^N, \psi_f\right)$,    "productive step"

8:        $N \rightarrow I$

9:     **else**

10:       $h^g = \frac{\varepsilon}{M_g^2}$,

11:       $x^{N+1} = Mirr_{h^g}\left(x^N, \psi_g\right)$,    "non-productive step"

12:       $N \rightarrow J$

13:     **end if**

14:     $N \leftarrow N + 1$

15: **until** $\frac{2\Theta_0^2}{\varepsilon^2} \leq \frac{|I|}{M_f^2} + \frac{|J|}{M_g^2}$.

**Ensure:** $\widehat{x} := \frac{1}{|I|} \sum\limits_{k \in I} x^k$.

---

*Then, after the stopping of Algorithm 3, the following inequalities hold:*

$$f(\widehat{x}) - f(x^*) \leq \varepsilon + \delta \quad and \quad g(\widehat{x}) \leq \varepsilon + \delta.$$

*Moreover, the required number of iterations of Algorithm 3 does not exceed*

$$N = \frac{2M^2 \Theta_0^2}{\varepsilon^2}, \ where \ M = \max\{M_f, M_g\}.$$

**Remark 2.2.** Clearly, Algorithms 2 and 3 are optimal in terms of the lower bounds [25]. More precisely, let us understand hereinafter the optimality of the Mirror Descent methods as the complexity $O(\frac{1}{\varepsilon^2})$ (it is well-kown that this estimate is optimal for Lipschitz-continuous functionals [25]).

**Remark 2.3** (The case of several functional constraints)**.** Let us consider a set of convex functions $f$ and $g_p : Q \rightarrow \mathbb{R}$, $p \in [m] \stackrel{\text{def}}{=} \{1, 2, \ldots, m\}$. We will focus on the following constrained optimization problem

$$\min\{f(x): \ x \in Q \ \text{and} \ g_p(x) \leq 0 \ \text{for all} \ p \in [m]\}. \tag{2.14}$$

It is clear that instead of a set of functionals $\{g_p(\cdot)\}_{p=1}^m$ we can consider one functional constraint $g : Q \to \mathbb{R}$, such that $g(x) = \max_{p \in [m]}\{g_p(x)\}$. Therefore, by this setting, problem (2.14) will be equivalent to problem (1.3).

Assume that for any $p \in [m]$, the functional $g_p$ satisfies the following condition

$$-\psi_{g_p}(y, x) \leq M_{g_p}\sqrt{2V_d(y, x)} + \delta$$

For problem (2.14), we propose a modification of Algorithms 2 and 3 (the modified algorithms are listed as Algorithm 6 and 7 in Appendix A). The idea of the proposed modification allows saving the running time of algorithms due to consideration of not all functional constraints on non-productive steps.

**Remark 2.4** (Composite Optimization Problems [6,20,26])**.** Proposed methods are applicable to the composite optimization problems, specifically

$$\min\{f(x) + r(x) : \quad x \in Q, \ g(x) + \eta(x) \leq 0\},$$

where $r, \eta : Q \to \mathbb{R}$ are so-called simple convex functionals (i. e. the proximal mapping operator $Mirr_h(x, \psi)$ is easily computable). For this case, for any $x, y \in Q$, we have

$$\psi_f(y, x) = \langle \nabla f(x), y - x \rangle + r(y) - r(x)$$

$$\psi_g(y, x) = \langle \nabla g(x), y - x \rangle + \eta(y) - \eta(x).$$

# 3.   Stochastic Mirror Descent Algorithm

Let us, in this section, consider the stochastic setting of the problem (1.3). This means that we can still use the value of the objective function and functional constraints, but instead of their (sub)gradient, we use their stochastic (sub)gradient. Namely, we consider the first-order unbiased oracle that produces $\nabla f(x, \xi)$ and $\nabla g(x, \zeta)$, where $\xi$ and $\zeta$ are random vectors and

$$\mathbb{E}[\nabla f(x, \xi)] = \nabla f(x), \quad \nabla \mathbb{E}[g(x, \zeta)] = \nabla g(x).$$

Assume that for each $x, y \in Q$

$$\langle \nabla f(x, \xi), x - y \rangle \leq M_f\sqrt{2V_d(y, x)} \text{ and } \langle \nabla g(x, \zeta), x - y \rangle \leqslant M_g\sqrt{2V_d(y, x)}, \tag{3.15}$$

where $M_f, M_g > 0$, Let us consider the proximal mapping operator for $f$

$$Mirr_h(x, \xi) = \arg\min_{y \in Q} \left\{ \frac{1}{h} V_d(y, x) + \langle \nabla f(x, \xi), y \rangle \right\},$$

and similarly, we consider the proximal mapping operator for $g$. The following lemma describes the main property of this operator.

**Lemma 3.1.** *Let $f$ be a convex function which satisfies (1.4), $h > 0, \delta > 0$, $\xi$ be a random vector and $\tilde{x} = Mirr_h(x, \xi)$. Then for all $y \in Q$*

$$h(f(x) - f(y)) \le \phi_f^*(h) + V_d(y, x) - V_d(y, \tilde{x}) + h\langle \nabla f(x, \xi) - \nabla f(x), y - x \rangle + h\delta,$$

*where, as earlier, $\phi_f^*(h) = \frac{h^2 M_f^2}{2}$.*

Suppose $\varepsilon > 0$ is a given positive real number. We say that a (random) point $\widehat{x} \in Q$ is an expected $\varepsilon$–solution to the problem (1.3), in the stochastic setting, if

$$\mathbb{E}[f(\widehat{x})] - f(x^*) \le \varepsilon \text{ and } g(\widehat{x}) \le \varepsilon. \tag{3.16}$$

In order to solve the stochastic setting of the considered problem (1.3), we propose the following algorithm.

The following theorem gives information about the efficiency of the algorithm. The proof of this theorem is given in Appendix B.

**Theorem 3.1.** *Let $f$ and $g$ be convex functions, which satisfy (1.4) and (1.5). Let $\varepsilon > 0, \delta > 0$ be fixed positive numbers. Then, after the stopping of Algorithm 4, the following inequalities hold:*

$$\mathbb{E}[f(\widehat{x})] - f(x^*) \le \varepsilon + \delta \quad and \quad g(\widehat{x}) \le \varepsilon + \delta.$$

**Remark 3.1.** It should be noted how the optimality of the proposed method can be understood. With the special assumptions (2.10) – (2.11) and choice of $h^f, h^g$, the complexity of the algorithm is $O(\frac{1}{\varepsilon^2})$, which is optimal in such class of problems.

## 4.  Online Optimization Problem

In this section we consider the online setting of the optimization problem (1.3). Namely

$$\frac{1}{N} \sum_{i=1}^{N} f_i(x) \to \min_{x \in Q, g(x) \le 0}, \tag{4.17}$$

**Algorithm 4** Modified Mirror Descent for the stochastic setting.

---

**Require:** $\varepsilon > 0, \delta > 0, h^f > 0, h^g > 0, \Theta_0 : d(x^*) \le \Theta_0^2$.

1: $x^0 = \arg\min_{x \in Q} d(x)$.

2: $I =: \emptyset$ and $J =: \emptyset$

3: $N \leftarrow 0$

4: **repeat**

5:     **if** $g\left(x^N\right) \le \varepsilon + \delta$ **then**

6:        $x^{N+1} = Mirr_{h^f}\left(x^N, \xi^N, \psi_f\right),$    "productive step"

7:        $N \to I$

8:     **else**

9:        $x^{N+1} = Mirr_{h^g}\left(x^N, \zeta^N, \psi_g\right),$    "non-productive step"

10:       $N \to J$

11:     **end if**

12:     $N \leftarrow N + 1$

13: **until** $\Theta_0^2 \le \varepsilon\left(|J|h^g + |I|h^f\right) - |J|\phi_g^*(h^g) - |I|\phi_f^*(h^f)$.

**Ensure:** $\widehat{x} := \frac{1}{|I|}\sum\limits_{k \in I} x^k$.

---

under the assumption that all $f_i : Q \to \mathbb{R}$ $(i = 1, \ldots, N)$ and g satisfy (2.10) and (2.11) with constants $M_i > 0, i = 1, \ldots, N$ and $M_g > 0$.

In order to solve problem (4.17), we propose an algorithm (listed as Algorithm 5 below). This algorithm produces $N$ productive steps and in each step, the (sub)gradient of exactly one functional of the objectives is calculated. As a result of this algorithm, we get a sequence $\{x^k\}_{k \in I}$ (on productive steps), which can be considered as a solution to problem (4.17) with accuracy $\kappa$ (see (4.18)).

Assume that $M = \max\{M_i, M_g\}, h^f = h^g = h = \frac{\varepsilon}{M}$.

For Algorithm 5, we have the following result.

**Theorem 4.1.** *Suppose all $f_i : Q \to \mathbb{R}$ $(i = 1, \ldots, N)$ and g satisfy (2.10) and (2.11) with constants $M_i > 0, i = 1, \ldots, N$ and $M_g > 0$, Algorithm 5 works exactly $N$ productive steps. Then after the stopping of this Algorithm, the following inequality holds*

$$\frac{1}{N}\sum_{i=1}^{N} f_i(x^k) - \min_{x \in Q} \frac{1}{N}\sum_{i=1}^{N} f_i(x) \le \kappa,$$

*moreover, when the regret is non-negative, there will be no more than $O(N)$ non-productive steps.*

14

**Algorithm 5** Modified Mirror Descent for the online setting.

**Require:** $\varepsilon > 0, \delta > 0, M > 0, N, \Theta_0 : d(x^*) \leq \Theta_0^2$.

1: $x^0 = \arg\min_{x \in Q} d(x)$.

2: $i := 1, k := 0$

3: set $h = \frac{\varepsilon}{M^2}$

4: **repeat**

5:     **if** $g\left(x^k\right) \leq \varepsilon + \delta$ **then**

6:         $x^{k+1} = Mirr_h\left(x^k, \psi_{f_i}\right),$     "productive step"

7:         $i = i + 1,$

8:         $k = k + 1,$

9:     **else**

10:       $x^{k+1} = Mirr_h\left(x^k, \psi_g\right),$     "non-productive step"

11:       $k = k + 1,$

12:     **end if**

13: **until** $i = N + 1$.

14: Guaranteed accuracy:

$$\kappa = \frac{|J|}{N}\left(-\frac{\varepsilon}{2}\right) + \left(\frac{\varepsilon}{2} + \delta\right) + \frac{M^2\Theta_0^2}{N\varepsilon}. \tag{4.18}$$

The proof of this theorem is given in Appendix C. In particular, note that the proposed method is optimal [15]: if for some $C > 0$, $\kappa \sim \varepsilon \sim \delta = \frac{C}{\sqrt{N}}$, then $|J| \sim O(N)$.

# 5. Numerical Experiments

To show the practical performance of the proposed Algorithms 2, 3 and their modified versions, Algorithm 6 and 7, in the case of many functional constraints, a series of numerical experiments were performed[1], for the well-known *Fermat-Torricelli-Steiner* problem, but with some non-smooth functional constraints.

For a given set $\{P_k = (p_{1k}, p_{2k}, \ldots, p_{nk}); \ k \in [r]\}$ of $r$ points, in $n$-dimensional Euclidean space $\mathbb{R}^n$, we need to solve the considered optimization problem

---

[1]All experiments were implemented in Python 3.4, on a computer fitted with Intel(R) Core(TM) i7-8550U CPU @ 1.80GHz, 1992 Mhz, 4 Core(s), 8 Logical Processor(s). RAM of the computer is 8 GB.

(1.3), where the objective function $f$ is given by

$$f(x) := \frac{1}{r} \sum_{k=1}^{r} \sqrt{(x_1 - p_{1k})^2 + \ldots + (x_n - p_{nk})^2} = \frac{1}{r} \sum_{k=1}^{r} \|x - P_k\|_2 . \quad (5.19)$$

The functional constraint has the following form

$$g(x) = \max_{i \in [m]} \{g_i(x) = \alpha_{i1}x_1 + \alpha_{i2}x_2 + \ldots + \alpha_{in}x_n\}. \quad (5.20)$$

The coefficients $\alpha_{i1}, \alpha_{i2}, \ldots, \alpha_{in}$, for all $i \in [m]$, in (5.20) and the coordinates of the points $P_k$, for all $k \in [r]$, are drawn from the normal (Gaussian) distribution with the location of the mode equaling 1 and the scale parameter equaling 2.

We choose the standard Euclidean norm and the Euclidean distance function in $\mathbb{R}^n$, $\delta = 0$, starting point $x^0 = \left( \frac{1}{\sqrt{n}}, \ldots, \frac{1}{\sqrt{n}} \right) \in \mathbb{R}^n$ and $Q$ is the unit ball in $\mathbb{R}^n$.

We run Algorithms 2, 3, 6 and 7, for $m = 200, n = 500, r = 100$ and different values of $\varepsilon \in \{\frac{1}{2^i} : i = 1, 2, 3, 4, 5\}$. The results of the work of these algorithms are represented in Table 1 below. These results demonstrate the comparison of the number of iterations (Iter.), the running time (in seconds) of each algorithm and the qualities of the solution, produced by these algorithms with respect to the objective function $f$ and the functional constraint $g$, where we calculate the values of these functions at the output $x^{\text{out}} := \widehat{x}$ of the algorithms. We set $f^{\text{best}} := f(x^{\text{out}})$ and $g^{\text{out}} := g(x^{\text{out}})$.

In general, from the conducted experiments, we can see that Algorithm 2 and its modified version (Algorithm 6) work faster than Algorithms 3 and its modified version (Algorithm 7). But note that Algorithms 3 and 7 guarantee a better quality of the resulting solution to the considered problem, with respect to the objective function $f$ and the functional constraint (5.20). Also, we can see the efficiency of the modified Algorithm 7, which saves the running time of the algorithm, due to consideration of not all functional constraints on non-productive steps.

# Conclusion

In the paper, there was introduced the concept of an inexact $(\delta, \phi, V)$–model of the objective function. There were considered some modifications of the

Table 1: The results of Algorithms 2, 3, 6 and 7, with $m = 200, n = 500, r = 100$ and different values of $\varepsilon$.

| $1/\varepsilon$ | Algorithm 2 | | | | Algorithm 6 | | | |
|---|---|---|---|---|---|---|---|---|
| | Iter. | Time (sec.) | $f^{\text{best}}$ | $g^{\text{out}}$ | Iter. | Time (sec.) | $f^{\text{best}}$ | $g^{\text{out}}$ |
| 2 | 16 | 5.138 | 22.327427 | 2.210041 | 16 | 4.883 | 22.327427 | 2.210041 |
| 4 | 64 | 20.911 | 22.303430 | 2.016617 | 64 | 20.380 | 22.303430 | 2.016617 |
| 8 | 256 | 84.343 | 22.283362 | 1.858965 | 256 | 79.907 | 22.283362 | 2.015076 |
| 16 | 1024 | 317.991 | 22.274366 | 1.199792 | 1024 | 317.033 | 22.273177 | 1.988190 |
| 32 | 4096 | 1253.717 | 22.272859 | 0.607871 | 4096 | 1145.033 | 22.269038 | 1.858965 |
| | Algorithm 3 | | | | Algorithm 7 | | | |
| 2 | 167 | 9.455 | 22.325994 | 0.417002 | 164 | 7.373 | 22.325604 | 0.391461 |
| 4 | 710 | 39.797 | 22.305980 | 0.204158 | 667 | 29.954 | 22.305654 | 0.188497 |
| 8 | 2910 | 158.763 | 22.289320 | 0.103493 | 2583 | 119.055 | 22.289302 | 0.088221 |
| 16 | 11613 | 626.894 | 22.280893 | 0.051662 | 10155 | 468.649 | 22.280909 | 0.045343 |
| 32 | 46380 | 2511.261 | 22.277439 | 0.026000 | 40149 | 1723.136 | 22.277450 | 0.022639 |

Mirror Descent algorithm, in particular for stochastic and online optimization problems. A significant part of the work was devoted to the research of a special case of relative Lipschitz condition for objective function and functional constraints. The proposed methods are applicable for a wide class of problems because relative Lipschitz-continuity is an essential generalization of the classical Lipschitz-continuity. However, for relatively Lipschitz-continuous problems, we could not propose adaptive methods like [3]. Note that Algorithm 3 and its modified version Algorithm 7 are partially adaptive since the resulting number of iterations is not fixed, due to the stopping criterion, although the step-sizes are fixed.

# References

[1] Alkousa M. S.: On Some Stochastic Mirror Descent Methods for Constrained Online Optimization Problems. Computer Research and Modeling, **11**(2), 205–217 (2019)

[2] Bauschke H. H., Bolte J., Teboulle M.: A Descent Lemma Beyond Lip-

schitz Gradient Continuity: First-Order Methods Revisited and Applications. Mathematics of Operations Research, **42**(2), 330–348 (2017)

[3] Bayandina A., Dvurechensky P., Gasnikov A., Stonyakin F., Titov A.: Mirror descent and convex optimization problems with non-smooth inequality constraints. In: Large-Scale and Distributed Optimization. Springer, Cham, 181–213 (2018)

[4] Beck A., Ben-Tal A., Guttmann-Beck N., Tetruashvili L.: The comirror algorithm for solving nonsmooth constrained convex problems. Operations Research Letters, **38**(6), 493–498 (2010)

[5] Beck A., Teboulle M.: Mirror descent and nonlinear projected subgradient methods for convex optimization. Oper. Res. Lett., **31**(3), 167–175 (2003)

[6] Beck A., Teboulle M.: A Fast Iterative Shrinkage-Thresholding Algorithm for Linear Inverse Problems. SIAM J. Imaging Sci., **2**(1), 183-вЂ"202 (2009)

[7] Belmega E. V., Mertikopoulos P.: Online and Stochastic Optimization beyond Lipschitz Continuity: A Riemannian Approach. Published as a conference paper at ICLR (2020). `https://openreview.net/pdf?id=rkxZyaNtwB`

[8] Boyd S., Vandenberghe L.: Convex Optimization. Cambridge University Press, New York (2004)

[9] Devolder O., Glineur F., Nesterov Yu.: First-order methods of smooth convex optimization with inexact oracle. Math. Program. **146**(1–2), 37–75 (2014)

[10] Doan T. T., Bose S., Nguyen D. H., Beck C. L.: Convergence of the Iterates in Mirror Descent Methods. IEEE Control Systems Letters, **3**(1), 114–119 (2019)

[11] Gasnikov A. V.: Modern numerical optimization methods. The method of universal gradient descent (2018). (in Russian). `https://arxiv.org/ftp/arxiv/papers/1711/1711.00394.pdf`

[12] Gasnikov A. V., Lagunovskaya A. A., Usmanova I. N., Fedorenko F. A., Krymova E. A.: Stochastic online optimization. Single-point and multi-point non-linear multi-armed bandits. Convex and strongly-convex case. Automation and Remote Control, **78**(2), 224–234 (2017)

[13] Hanzely F., Richtarik P., Xiao L.: Accelerated Bregman proximal gradient methods for relatively smooth convex optimization. (2018) https://arxiv.org/pdf/1808.03045.pdf

[14] Hao Y., Neely M. J., Xiaohan W.: Online Convex Optimization with Stochastic Constraints. Published in NIPS, 1427–1437 (2017)

[15] Hazan E., Kale S.: Beyond the regret minimization barrier: Optimal algorithms for stochastic strongly-convex optimization. JMLR. **15**, 2489–2512 (2014)

[16] Hazan E.: Introduction to online convex optimization. Foundations and Trends in Optimization, **2**(3–4), 157–325 (2015)

[17] Jenatton R., Huang J., Archambeau C.: Adaptive Algorithms for Online Convex Optimization with Long-term Constraints. Proceedings of The 33rd International Conference on Machine Learning, PMLR 48, 402–411 (2016)

[18] Lu, H., Freund, R. M., Nesterov, Y. (2018). Relatively smooth convex optimization by first-order methods, and applications. SIAM Journal on Optimization, **28**(1), 333-354.

[19] Lu H.: Relative Continuity for Non-Lipschitz Nonsmooth Convex Optimization Using Stochastic (or Deterministic) Mirror Descent. Informs Journal on Optimization **1**(4), 288–303 (2019).

[20] Lu H., Freund R. M., Nesterov Yu.: Relatively smooth convex optimization by first-order methods, and applications. SIAM J. Optim. **28**(1), 333вЪ‰-354 (2018)

[21] Luong D. V. N., Parpas P., Rueckert D., Rustem B.: A Weighted Mirror Descent Algorithm for Nonsmooth Convex Optimization Problem. J Optim Theory Appl. **170**(3), 900–915 (2016)

[22] Nazin A. V., Miller B. M.: Mirror Descent Algorithm for Homogeneous Finite Controlled Markov Chains with Unknown Mean Losses. Proceedings of the 18th World Congress The International Federation of Automatic Control Milano (Italy) August 28 – September 2 (2011)

[23] Nazin A., Anulova S., Tremba A.: Application of the Mirror Descent Method to Minimize Average Losses Coming by a Poisson Flow. European Control Conference (ECC) June 24–27 (2014)

[24] Nemirovskii A.: Efficient methods for large-scale convex optimization problems. Ekonomika i Matematicheskie Metody, (1979). (in Russian)

[25] Nemirovsky A., Yudin D.: Problem Complexity and Method Efficiency in Optimization. J. Wiley & Sons, New York (1983).

[26] Nesterov Yu.: Gradient methods for minimizing composite functions. Mathematical Programming **140**, 125-вЂ"161 (2013)

[27] Nesterov Yu.: Relative Smoothness: New Paradigm in Convex Optimization. Conference report, EUSIPCO-2019, A Coruna, Spain, September 4, 2019. `http://eusipco2019.org/wp-content/uploads/2019/10/Relative-Smoothness-New-Pa`

[28] Orabona F.: A Modern Introduction to Online Learning. (2020) `https://arxiv.org/pdf/1912.13213.pdf`

[29] Orabona F., Crammer K., Cesa-Bianchi N.: A generalized online mirror descent with applications to classification and regression. Mach Learn 99, 411вЂ"-435 (2015).

[30] Titov A. A., Stonyakin F. S., Gasnikov A. V., Alkousa M. S.: Mirror Descent and Constrained Online Optimization Problems. Optimization and Applications, 9th International Conference OPTIMA-2018 (Petrovac, Montenegro, October 1–5, 2018). Revised Selected Papers. Communications in Computer and Information Science, **974**, 64–78 (2019)

[31] Tremba A., Nazin A.: Extension of a saddle point mirror descent algorithm with application to robust PageRank. 52nd IEEE Conference on Decision and Control December 10–13 (2013)

[32] Tyurin A. I., Gasnikov A. V.: Fast gradient descent method for convex optimization problems with an oracle that generates a $(\delta, L)$–model of a function in a requested point. Computational Mathematics and Mathematical Physics. **59**(7), 1137–1150 (2019)

[33] Yuan J., Lamperski A.: Online convex optimization for cumulative constraints. Published in NIPS, 6140–6149 (2018)

[34] Yunwen L., Ding-Xuan Z.: Convergence of online mirror descent. Appl. Comput. Harmon. Anal. **48**(1), 343–373 (2020)

[35] Zhou Z., Mertikopoulos P., Bambos N., Boyd S., Glynn P.: On the convergence of mirror descent beyond stochastic convex programming. (2018) `https://arxiv.org/pdf/1706.05681.pdf`

# Appendix A. Modified Algorithms for Problems with Several Functional Constraints

---

**Algorithm 6** Modified MDA for Relatively Lipschitz-continuous functions, version 2, several functional constraints. (The modification of Algorithm 2)

---

**Require:** $\varepsilon > 0, \delta > 0, M_f > 0, M_g > 0, \Theta_0 : d(x^*) \leq \Theta_0^2$.

1: $x^0 = \arg\min_{x \in Q} d(x)$.

2: $I =: \emptyset$

3: $N \leftarrow 0$

4: **repeat**

5:    **if** $g(x^N) \leq M_g\varepsilon + \delta$ **then**

6:        $h^f = \frac{\varepsilon}{M_f}$,

7:        $x^{N+1} = Mirr_{h^f}(x^N, \psi_f)$,        "productive step"

8:        $N \rightarrow I$

9:    **else**

10:        // $g_{p(N)}(x^N) > M_g\varepsilon + \delta$ for some $p(N) \in [m]$

11:        $h^{g_{p(N)}} = \frac{\varepsilon}{M_{g_{p(N)}}}$    // $M_{g_{p(N)}}$ is the Lipschitz constant of the constraint $g_{p(N)}$.

12:        $x^{N+1} = Mirr_{h^{g_{p(N)}}}(x^N, \psi_{g_{p(N)}})$,    "non-productive step"

13:    **end if**

14:    $N \leftarrow N + 1$

15: **until** $N \geq \frac{2\Theta_0^2}{\varepsilon^2}$.

**Ensure:** $\widehat{x} := \frac{1}{|I|} \sum_{k \in I} x^k$.

---

For the proposed modified Algorithm 6, the following result holds.

**Theorem 5.1.** *Let $f$ and $g$ be convex functions, which satisfy (2.10) and (2.11) for $M_f > 0$ and $M_g > 0$. Let $\varepsilon > 0, \delta > 0$ be fixed positive numbers. Assume that $\Theta_0 > 0$ is a known constant such that $d(x^*) \leq \Theta_0^2$. Then, after the stopping of Algorithm 6, the following inequalities hold*

$$f(\widehat{x}) - f(x^*) \leq M_g\varepsilon + \delta \quad and \quad g_{p(k)}(\widehat{x}) \leq M_g\varepsilon + \delta,$$

*where, by $g_{p(k)}$ we mean any constraint such that the inequality $g_{p(k)}(x^k) > M_g\varepsilon + \delta$ holds.*

Similarly, for the proposed modified Algorithm 7, we have the following result.

**Algorithm 7** Modified MDA for Relatively Lipschitz-continuous functions, version 2, several functional constraints. (The modification of Algorithm 3)

---

**Require:** $\varepsilon > 0, \delta > 0, M_f > 0, \Theta_0 : d(x^*) \leq \Theta_0^2$.

1: $x^0 = \arg\min_{x \in Q} d(x)$.

2: $I =: \emptyset$ and $J =: \emptyset$

3: $N \leftarrow 0$

4: **repeat**

5:     **if** $g(x^N) \leq \varepsilon + \delta$ **then**

6:        $h^f = \frac{\varepsilon}{M_f^2}$,

7:        $x^{N+1} = Mirr_{h^f}(x^N, \psi_f)$,        "productive step"

8:        $N \to I$

9:     **else**

10:        // $g_{p(N)}(x^N) > \varepsilon + \delta$ for some $p(N) \in [m]$

11:        $h^{g_{p(N)}} = \frac{\varepsilon}{M_{g_{p(N)}}^2}$    // $M_{g_{p(N)}}$ is the Lipschitz constant of the constraint $g_{p(N)}$.

12:        $x^{N+1} = Mirr_{h^{g_{p(N)}}}(x^N, \psi_{g_{p(N)}})$,    "non-productive step"

13:        $N \to J$

14:     **end if**

15:     $N \leftarrow N + 1$

16: **until** $\frac{2\Theta_0^2}{\varepsilon^2} \leq \frac{|I|}{M_f^2} + \sum_{k \in J} \frac{1}{M_{g_{p(k)}}^2}$.

**Ensure:** $\widehat{x} := \frac{1}{|I|} \sum_{k \in I} x^k$.

---

**Theorem 5.2.** *Let $f$ and $g$ be convex functions, which satisfy (2.10) and (2.11) for $M_f > 0$ and $M_g > 0$. Let $\varepsilon > 0, \delta > 0$ be fixed positive numbers. Assume that $\Theta_0 > 0$ is a known constant such that $d(x^*) \leq \Theta_0^2$.*

*Then, after the stopping of Algorithm 7, the following inequalities hold*

$$f(\widehat{x}) - f(x^*) \leq \varepsilon + \delta \quad and \quad g_{p(k)}(\widehat{x}) \leq \varepsilon + \delta.$$

*Moreover, if $g(x) = \max_{p \in [m]}\{g_p(x)\}$ satisfies (2.13), then the required number of iterations of Algorithm 7 does not exceed*

$$N = \frac{2M^2\Theta_0^2}{\varepsilon^2}, \quad where \ M = \max\{M_f, M_g\}.$$

# Appendix B. The proof of Theorem 3.1.

Denote
$$\gamma_k = \begin{cases} \langle \nabla f(x^k, \xi^k) - \nabla f(x^k), x^* - x^k \rangle & \text{if } k \in I, \\ \langle \nabla g(x^k, \zeta^k) - \nabla g(x), x^* - x^k \rangle & \text{if } k \in J. \end{cases} \tag{5.21}$$

By Lemma 3.1, we have for all $k \in I$

$$h^f \left( f(x^k) - f(x^*) \right) \leq \phi_f^*(h) + V_d(x^k, x^*) - V_d(x^{k+1}, x^*) +$$
$$+ h^f \left\langle \nabla f(x^k, \xi^k) - \nabla f(x^k), x^* - x^k \right\rangle + h^f \delta,$$

the same for all $k \in J$, we have

$$h^g \left( g(x^k) - g(x^*) \right) \leq \phi_g^*(h) + V_d(x^k, x^*) - V_d(x^{k+1}, x^*) +$$
$$+ h^g \left\langle \nabla g(x^k, \zeta^k) - \nabla g(x^k), x^* - x^k \right\rangle + h^g \delta.$$

By taking summation, in each side of both previous inequalities, over productive and non-productive steps, we get

$$\sum_{k \in I} h^f \left( f(x^k) - f(x^*) \right) + \sum_{k \in J} h^g \left( g(x^k) - g(x^*) \right) \leq$$

$$\leq \sum_{k \in I} \phi_f^*(h^f) + \sum_{k \in J} \phi_g^*(h^g) + \sum_k \left( V_d(x^*, x^k) - V_d(x^*, x^{k+1}) \right) + \sum_{k \in I} (h^f \delta + \gamma_k) + \sum_{k \in J} (h^g \delta + \gamma_k) \leq$$

$$\sum_{k \in I} \phi_f^*(h^f) + \sum_{k \in J} \phi_g^*(h^g) + \Theta_0^2 + \sum_{k \in I} h^f \delta + \sum_{k \in J} h^g \delta + \sum_{k \in I} h^f \gamma_k + \sum_{k \in J} h^g \gamma_k.$$

For each $k \in J, g(x^k) - g(x^*) > \varepsilon + \delta$ and we have

$$\sum_{k \in I} h^f (f(\widehat{x}) - f(x^*)) \leq \sum_{k \in I} \phi_f^*(h^f) + \sum_{k \in J} \phi_g^*(h^g) + \Theta_0^2 - \varepsilon \sum_{k \in J} h^g + \sum_{k \in I} h^f \delta$$
$$+ \sum_{k \in I} h^f \gamma_k + \sum_{k \in J} h^g \gamma_k = |I| \left( \phi_f^*(h^f) + \delta h^f \right) + |J| \phi_g^*(h^g) - |J| \varepsilon h^g$$
$$+ \Theta_0^2 + \sum_{k \in I} h^f \gamma_k + \sum_{k \in J} h^g \gamma_k \leq \varepsilon |I| h^f + |I| h^f \delta + \sum_{k \in I} h^f \gamma_k + \sum_{k \in J} h^g \gamma_k.$$

Now from the definition of $\widehat{x}$ (the Ensure of Algorithm 4) and by taking the expectation we obtain

$$\mathbb{E}[f(\widehat{x})] - f(x^*) \leq \varepsilon + \delta + \mathbb{E}\left[ \sum_{k \in I} \frac{\gamma_k}{|I|} \right] + \mathbb{E}\left[ \sum_{k \in J} \frac{\gamma_k}{|J|} \right] = \varepsilon + \delta,$$

as well as $g(\widehat{x}) \leq \varepsilon + \delta$.

# Appendix C. The proof of Theorem 4.1.

By Lemma 1.1, we have for all $k \in I$

$$h \left( f_i(x^k) - f_i(y) \right) \leq \phi^*(h) + V_d(y, x^k) - V_d(y, x^{k+1}) + h\delta, \qquad (5.22)$$

the same for all $k \in J$, we have

$$h \left( g(x^k) - g(y) \right) \leq \phi^*(h) + V_d(y, x^k) - V_d(y, x^{k+1}) + h\delta. \qquad (5.23)$$

By taking summation, in each side of (5.22) and (5.23), over productive and non-productive steps, we get

$$\sum_{i=1}^{N} h \left( f_i(x^k) - f_i(x^*) \right) + \sum_{k \in J} h \left( g(x^k) - g(x^*) \right) \leq (N + |J|) \left( \phi^*(h) + h\delta \right) +$$

$$+ \sum_{k} \left( V_d(x^*, x^k) - V_d(x^*, x^{k+1}) \right)$$

$$\leq (N + |J|) \left( \phi^*(h) + h\delta \right) + \Theta_0^2.$$

Then

$$\sum_{i=1}^{N} f_i(x^k) - f_i(x^*) \leq |J| \left( \frac{M^2 \phi^*(h)}{\varepsilon} + \delta \right) + N \left( \frac{M^2 \phi^*(h)}{\varepsilon} + \delta \right) + \frac{M^2 \Theta_0^2}{\varepsilon} - |J|\varepsilon - |J|\delta$$

$$= |J| \left( \frac{M^2 \phi^*(h)}{\varepsilon} - \varepsilon \right) + N \left( \frac{M^2 \phi^*(h)}{\varepsilon} + \delta \right) + \frac{M^2 \Theta_0^2}{\varepsilon},$$

and then we get

$$\frac{1}{N} \sum_{i=1}^{N} f_i(x^k) - \min_{x \in Q} \frac{1}{N} \sum_{i=1}^{N} f_i(x) \leq \frac{|J|}{N} \left( \frac{M^2 \phi^*(h)}{\varepsilon} - \varepsilon \right) + \left( \frac{M^2 \phi^*(h)}{\varepsilon} + \delta \right) + \frac{M^2 \Theta_0^2}{N \varepsilon}.$$

Recall that $\phi^*(h) = \frac{h^2 M^2}{2}$, so

$$\frac{1}{N} \sum_{i=1}^{N} f_i(x^k) - \min_{x \in Q} \frac{1}{N} \sum_{i=1}^{N} f_i(x) \leq \frac{|J|}{N} \left( \frac{h M^2}{2} - \varepsilon \right) + \left( \frac{h M^2}{2} + \delta \right) + \frac{M^2 \Theta_0^2}{N \varepsilon}$$

$$= \frac{|J|}{N} \left( \frac{\varepsilon}{2} - \varepsilon \right) + \left( \frac{\varepsilon}{2} + \delta \right) + \frac{M^2 \Theta_0^2}{N \varepsilon}.$$

and by virtue of (4.18), we get

$$\frac{1}{N} \sum_{i=1}^{N} f_i(x^k) - \min_{x \in Q} \frac{1}{N} \sum_{i=1}^{N} f_i(x) \leq \kappa. \qquad (5.24)$$

Assuming the non-negativity of the regret (i.e. the left side in (5.24)):

$$0 \leq \sum_{i=1}^{N} f_i(x^k) - f_i(x^*) \leq |J| \left( \frac{hM^2}{2} - \varepsilon \right) + N \left( \frac{hM^2}{2} + \delta \right) + \frac{M^2 \Theta_0^2}{\varepsilon}$$

$$= |J| \left( -\frac{\varepsilon}{2} \right) + N \left( \frac{\varepsilon}{2} + \delta \right) + \frac{M^2 \Theta_0^2}{\varepsilon},$$

so

$$|J| \leq N \left( 1 + \frac{2\delta}{\varepsilon} \right) + \frac{2M^2 \Theta_0^2}{\varepsilon^2}.$$

Suppose $\kappa \sim \varepsilon \sim \delta = \frac{C}{\sqrt{N}}$, for some $C > 0$, then we get

$$|J| \sim O(N) = N \left( 3 + \frac{2M^2 \Theta_0^2}{C^2} \right).$$

It means that the considered method is optimal for OCO, according to [15].