

One-Point Gradient-Free Methods for Smooth and Non-Smooth Saddle-Point Problems*

Aleksandr Beznosikov^{1,2}, Vasilii Novitskii¹, and Alexander Gasnikov^{1,2}

¹ Moscow Institute of Physics and Technology, Dolgoprudny, Russia

² Higher School of Economics, Russia

Abstract. In this paper, we analyze gradient-free methods with one-point feedback for stochastic saddle point problems $\min_x \max_y \varphi(x, y)$. For non-smooth and smooth cases, we present an analysis in a general geometric setup with arbitrary Bregman divergence. For problems with higher order smoothness, the analysis is carried out only in the Euclidean case. The estimates we have obtained repeat the best currently known estimates of gradient-free methods with one-point feedback for problems of imagining a convex or strongly convex function. The paper uses three main approaches to recovering the gradient through finite differences: standard with a random direction, as well as its modifications with kernels and residual feedback. We also provide experiments to compare these approaches for the matrix game.

Keywords: saddle-point problem · zeroth order method · one-point feedback · stochastic optimization.

1 Introduction

This paper is devoted to solving the saddle-point problem:

$$\min_{x \in \mathcal{X}} \max_{y \in \mathcal{Y}} \varphi(x, y). \quad (1)$$

It has many practical applications. These are the already well-known and classic matrix game and Nash equilibrium, as well as modern machine learning problems: Generative Adversarial Networks (GANs) [11] and Reinforcement Learning (RL) [12]. We assume that only zeroth-order information about the function is available, i.e. only its values, not a gradient, hessian, etc. This concept is called a Black-Box and arises in optimization [13], adversarial training [7], RL [9]. To make the problem statement more complex, but close to practice, it is natural to assume that we have access inexact values of function $\varphi(x, y, \xi)$, for example, with some random noise ξ . But even with the help of such an oracle, it is

*The research of A. Beznosikov and A. Gasnikov in Algorithm 1, Theorems 1-3 was supported by Russian Science Foundation (project No. 21-71-30005). The research of V. Novitskii in Algorithms 2, Theorems 4-7 was partially supported by Andrei Raigorodskii scholarship.

possible to recover some estimate of the gradient of a function in terms of finite differences.

Let us highlight two main approaches to such gradient estimates. The first approach is more well researched in the literature and is called a two-point feedback:

$$\frac{n}{2\tau}(\varphi(x + \tau\mathbf{e}_x, y + \tau\mathbf{e}_y, \xi) - \varphi(x - \tau\mathbf{e}_x, y - \tau\mathbf{e}_y, \xi)) \begin{pmatrix} \mathbf{e}_x \\ -\mathbf{e}_y \end{pmatrix}.$$

An important feature of this approach is that it is assumed that we were able to obtain the values of the function in points $(x + \tau\mathbf{e}_x, y + \tau\mathbf{e}_y)$ and $(x - \tau\mathbf{e}_x, y - \tau\mathbf{e}_y)$ with the same realization of the noise ξ . From the point of view of theoretical analysis, such an assumption is strong and gives good guarantees of convergence [8,16,13]. But from a practical point of view, this is a very idealistic assumption. Therefore, it is proposed to consider the concept of one-point feedback (which this paper is about):

$$\frac{n}{2\tau}(\varphi(x + \tau\mathbf{e}_x, y + \tau\mathbf{e}_y, \xi^+) - \varphi(x - \tau\mathbf{e}_x, y - \tau\mathbf{e}_y, \xi^-)) \begin{pmatrix} \mathbf{e}_x \\ -\mathbf{e}_y \end{pmatrix}.$$

In general $\xi^+ \neq \xi^-$. As far as we know, the use of methods with one-point approximation for saddle-point problems has not been studied at all in the literature. This is the main goal of our work.

1.1 Related works

Since the use of one-point feedback for saddle-point problems is new in the literature, we present related papers in two categories: two-point gradient-free methods for saddle-point problems, and one-point methods for minimization problems. Partially the results of these works are transferred to Table 1.

Two-point for saddle-point problems. Here, we first highlight work for non-smooth saddle-point problems [5], as well as work for smooth ones [15]. Note that in these papers an optimal estimate was obtained in the non-smooth case, and in the smooth case only for a special class of "firmly smooth" saddle-point problems. Also note the work devoted to coordinated methods for matrix games [6], which is also close to our topic.

One-point for minimization problems. First of all, we present works that analyze functions with higher order smoothness: [2,1,14]. These works are united by the technique of special random kernels, which allow you to use the smoothness of higher orders. Note that there is an error in work [2], therefore Table 1 shows the corrected result (according to the note from [1]). The special case of higher order smoothness is also interesting – the ordinary smoothness, it is also analyzed in [2,1,14], in addition we note the papers [10,17]. A nonsmooth analysis is presented in [10,17]. Note that in paper [10], not only the Euclidean setup is analyzed, but also the general case with an arbitrary Bregman divergence, which gives additional advantages in the estimates of the convergence (see Table 1).

1.2 Our contribution

In the nonsmooth case, we consider convex-concave and strongly-convex-strongly-concave problems with bounded $\nabla_x \varphi(x, y)$, $\nabla_y \varphi(x, y)$ on the optimization set. Our algorithm is modofocation of Mirror Descent with arbitrary Bregman divergence. The estimates we obtained coincide with the estimates for convex optimization with one-pointed feedback [10,17]. Using the correct geometry helps to reduce the contribution of the problem dimension to the final convergence estimate. In particular, in the entropy setting, convergence depends on the dimension of the problem linearly (see Table 1 for more details in convex-concave case and Table 2 – in strongly-convex-strongly-concave).

In the smooth case we obtained the estimates of the convergence rate with arbitrary Bregman divergence for convex-concave case and in Euclidean setup for strongly-convex-strongly-concave case. These estimates also coincide with the estimates for convex optimization with one-point feedback [10].

To the best of our knowledge this is the first time when exploiting higher-order smoothness helps to improve performance in saddle-point problems in both strongly-convex-strongly-concave and convex-concave cases. The results also coincide with the estimates for minimization [14,1].

In Tables 1 and 2 one can find a comparison of the oracle complexity of known results with zeroth-order methods for saddle-point problems in related works. Factor q depends on geometric setup of our problem and gives a benefit when we work in the Hölder, but non-Euclidean case (use non-Euclidean prox), i.e. $\|\cdot\| = \|\cdot\|_p$ and $p \in [1; 2]$, then $\|\cdot\|_* = \|\cdot\|_q$, where $1/p + 1/q = 1$. Then q takes values from 2 to ∞ , in particular, in the Euclidean case $q = 2$, but when the optimization set is a simplex, $q = \infty$. In higher-order smooth case we consider functions satisfying so called generalized Hölder condition with parameter $\beta > 2$ (see inequality (20) below). Note that it is prefer to use higher-order smooth methods rather than smooth methods only if $\beta > 3$.

2 Preliminaries

To begin with, we introduce some notation and definitions that we use in the work.

2.1 Notation

We use $\langle x, y \rangle \stackrel{\text{def}}{=} \sum_{i=1}^n x_i y_i$ to denote inner product of $x, y \in \mathbb{R}^n$ where x_i is the i -th component of x in the standard basis in \mathbb{R}^n . Then it induces ℓ_2 -norm in \mathbb{R}^n in the following way $\|x\|_2 \stackrel{\text{def}}{=} \sqrt{\langle x, x \rangle}$. We define ℓ_p -norms as $\|x\|_p \stackrel{\text{def}}{=} (\sum_{i=1}^n |x_i|^p)^{1/p}$ for $p \in (1, \infty)$ and for $p = \infty$ we use $\|x\|_\infty \stackrel{\text{def}}{=} \max_{1 \leq i \leq n} |x_i|$. The dual norm $\|\cdot\|_q$ for the norm $\|\cdot\|_p$ is denoted in the following way: $\|y\|_q \stackrel{\text{def}}{=} \max \{ \langle x, y \rangle \mid \|x\|_p \leq 1 \}$. Operator $\mathbb{E}[\cdot]$ is full mathematical expectation and operator $\mathbb{E}_\xi[\cdot]$ express conditional mathematical expectation.

Case	Oracle	Prob.	Complexity	Reference
non-smooth	two-point	SP	$\mathcal{O}\left(n^{\frac{2}{q}} \cdot \varepsilon^{-2}\right)$	[5]
	one-point	Min	$\mathcal{O}\left(n^{1+\frac{2}{q}} \cdot \varepsilon^{-4}\right)$	[10]
		SP	$\mathcal{O}\left(n^{1+\frac{2}{q}} \cdot \varepsilon^{-4}\right)$	this paper
smooth	two-point	SP	$\mathcal{O}\left(\left[n^{\frac{2}{q}} \text{ or } n\right] \cdot \varepsilon^{-2}\right)$	[15]
	one-point	Min	$\tilde{\mathcal{O}}\left(n^2 \cdot \varepsilon^{-3}\right)$	[10]
		SP	$\tilde{\mathcal{O}}\left(n^2 \cdot \varepsilon^{-3}\right)$	this paper
higher order smooth	one-point	Min	$\tilde{\mathcal{O}}\left(n^{2+\frac{2}{\beta-1}} \cdot \varepsilon^{-2-\frac{2}{\beta-1}}\right)$	[14,1]
		SP	$\tilde{\mathcal{O}}\left(n^{2+\frac{2}{\beta-1}} \cdot \varepsilon^{-2-\frac{2}{\beta-1}}\right)$	this paper

Table 1. Comparison of oracle complexity of one-point/two-point 0th-order methods for non-smooth/smooth **convex** minimization (Min) and **convex-concave** saddle-point (SP) problems under different assumptions. ε means the accuracy of the solution, n – dimension of the problem, $q = 2$ for the Euclidean case and $q = \infty$ for setup of $\|\cdot\|_1$ -norm.

Case	Oracle	Prob.	Complexity	Reference
non-smooth	one-point	Min	$\tilde{\mathcal{O}}\left(n^2 \cdot \varepsilon^{-3}\right)$	[10]
		SP	$\tilde{\mathcal{O}}\left(n^2 \cdot \varepsilon^{-3}\right)$	this paper
smooth	two-point	SP	$\mathcal{O}\left(n \cdot \varepsilon^{-1}\right)$	[15]
	one-point	Min	$\tilde{\mathcal{O}}\left(n^2 \cdot \varepsilon^{-2}\right)$	[10]
		SP	$\tilde{\mathcal{O}}\left(n^2 \cdot \varepsilon^{-2}\right)$	this paper
higher order smooth	one-point	Min	$\tilde{\mathcal{O}}\left(n^{2+\frac{1}{\beta-1}} \cdot \varepsilon^{-\frac{\beta}{\beta-1}}\right)$	[14,1]
		SP	$\tilde{\mathcal{O}}\left(n^{2+\frac{1}{\beta-1}} \cdot \varepsilon^{-\frac{\beta}{\beta-1}}\right)$	this paper

Table 2. Comparison of oracle complexity of one-point/two-point 0th-order methods for non-smooth/smooth **strongly-convex** minimization (Min) and **strongly-convex-strongly-concave** saddle-point (SP) problems under different assumptions.

Definition 1 (μ -strong convexity). Function $f(x)$ is μ -strongly convex w.r.t. $\|\cdot\|$ -norm on $\mathcal{X} \subseteq \mathbb{R}^n$ when it is continuously differentiable and there is a constant $\mu > 0$ such that the following inequality holds:

$$f(y) \geq f(x) + \langle \nabla f(x), y - x \rangle + \frac{\mu}{2} \|y - x\|^2, \quad \forall x, y \in \mathcal{X}.$$

Definition 2 (Prox-function). Function $d(z) : \mathcal{Z} \rightarrow \mathbb{R}$ is called prox-function if $d(z)$ is 1-strongly convex w.r.t. $\|\cdot\|$ -norm and differentiable on \mathcal{Z} .

Definition 3 (Bregman divergence). Let $d(z) : \mathcal{Z} \rightarrow \mathbb{R}$ is prox-function. For any two points $z, w \in \mathcal{Z}$ we define Bregman divergence $V_z(w)$ associated with $d(z)$ as follows:

$$V_z(w) = d(z) - d(w) - \langle \nabla d(w), z - w \rangle.$$

We denote the Bregman-diameter $\Omega_{\mathcal{Z}}$ of \mathcal{Z} w.r.t. $V_{z_1}(z_2)$ as $\Omega_{\mathcal{Z}} \stackrel{\text{def}}{=} \max\{\sqrt{2V_{z_1}(z_2)} \mid z_1, z_2 \in \mathcal{Z}\}$.

Definition 4 (Prox-operator). Let $V_z(w)$ Bregman divergence. For all $x \in \mathcal{Z}$ define prox-operator of ξ :

$$\text{prox}_x(\xi) = \arg \min_{y \in \mathcal{Z}} (V_x(y) + \langle \xi, y \rangle).$$

Now we are ready to formally describe the problem statement, as well as the necessary assumptions.

2.2 Settings and assumptions

As mentioned earlier, we consider the saddle-point problem (1), where $\varphi(\cdot, y)$ is convex function defined on compact convex set $\mathcal{X} \subset \mathbb{R}^{n_x}$, $\varphi(x, \cdot)$ is concave function defined on compact convex set $\mathcal{Y} \subset \mathbb{R}^{n_y}$. For convenience, we denote $\mathcal{Z} = \mathcal{X} \times \mathcal{Y}$ and then $z \in \mathcal{Z}$ means $z \stackrel{\text{def}}{=} (x, y)$, where $x \in \mathcal{X}$, $y \in \mathcal{Y}$. When we use $\varphi(z)$, we mean $\varphi(z) = \varphi(x, y)$.

Assumption 1 (Diameter of \mathcal{Z}) Let the compact set \mathcal{Z} have diameter Ω .

Assumption 2 (M -Lipschitz continuity) Function $\varphi(z)$ is M -Lipschitz continuous in certain neighbourhood of \mathcal{Z} with $M > 0$ w.r.t. norm $\|\cdot\|_2$ when

$$|\varphi(z) - \varphi(z')| \leq M \|z - z'\|_2, \quad \forall z, z' \in \mathcal{Z}.$$

One can prove that for all $z \in \mathcal{Z}$ we have

$$\|\tilde{\nabla} \varphi(z)\|_2 \leq M. \tag{2}$$

Assumption 3 (μ -strong convexity–strong concavity) Function $\varphi(z)$ is μ -strongly-convex-strongly-concave in \mathcal{Z} with $\mu > 0$ w.r.t. norm $\|\cdot\|_2$ when $\varphi(\cdot, y)$ is μ -strongly-convex for all y and $\varphi(x, \cdot)$ is μ -strongly-concave for all x w.r.t. $\|\cdot\|_2$.

Hereinafter, by $\tilde{\nabla} \varphi(z)$ we mean a block vector consisting of two vectors $\nabla_x \varphi(x, y)$ and $-\nabla_y \varphi(x, y)$. Recall that we do not have access to oracles $\nabla_x \varphi(x, y)$ or $\nabla_y \varphi(x, y)$. We only can use an inexact stochastic zeroth-order oracle $\tilde{\varphi}(x, y, \xi, \delta)$ at each iteration. Our model corresponds to the case when the oracle gives an inexact noisy function value. We have stochastic unbiased noise, depending on

the random variable ξ and biased deterministic noise δ . One can write it the following way:

$$\tilde{\varphi}(x, y, \xi) = \varphi(x, y) + \xi + \delta(x, y). \quad (3)$$

Note that δ depends on point (x, y) , and ξ is generated randomly regardless of this point.

Assumption 4 (Noise restrictions) *Stochastic noise ξ is unbiased with bounded variance, δ is bounded, i.e. there exists $\Delta, \sigma > 0$ such that*

$$\mathbb{E}\xi = 0, \quad \mathbb{E}[\xi^2] \leq \sigma^2, \quad |\delta| \leq \Delta. \quad (4)$$

3 Theoretical results

Since we do not have access to $\nabla_x \varphi(x, y)$ or $\nabla_y \varphi(x, y)$, it is proposed to replace them with finite differences. We present two variants: using a random euclidean direction [16,10] in non-smooth case and a kernel approximation [1,14] in smooth. These two concepts will be discussed in more detail later in the respective sections. As mentioned earlier, we work with one-point feedback. We use Mirror Descent as the basic algorithm, but with approximations instead of gradient.

3.1 Non-smooth case

Random euclidean direction. For $\mathbf{e} \in \mathcal{RS}_2^n(1)$ (a random vector uniformly distributed on the Euclidean unit sphere) and some constant τ let $\tilde{\varphi}(z + \tau \mathbf{e}, \xi) \stackrel{\text{def}}{=} \tilde{\varphi}(x + \tau \mathbf{e}_x, y + \tau \mathbf{e}_y, \xi)$, where \mathbf{e}_x is the first part of \mathbf{e} size of dimension n_x , and \mathbf{e}_y is the second part of dimension n_y . Then define estimation of the gradient through the difference of functions:

$$g(z, \mathbf{e}, \tau, \xi^\pm) = \frac{n(\tilde{\varphi}(z + \tau \mathbf{e}, \xi^+) - \tilde{\varphi}(z - \tau \mathbf{e}, \xi^-))}{2\tau} \begin{pmatrix} \mathbf{e}_x \\ -\mathbf{e}_y \end{pmatrix}, \quad (5)$$

where $n = n_x + n_y$. It is important that ξ^+ and ξ^- are different variables – this corresponds to the one-point concept. Next, we present Algorithm 1 – a modification of Mirror Descent with (5). Note that any Bregman divergence can be used in the prox operator. This allows us to take into account the geometric setup of the problem. \mathbf{e}_k and ξ_k^\pm are generated independently of the previous iterations and of each other. Here $\bar{z}_N = \frac{1}{N+1} \sum_{i=0}^N z_i$. Below we give technical facts about (5). Note that we do not provide proofs in the main part of the paper, they are all in the Appendix.

Algorithm 1 zoopMD

Input: z_0, N, γ, τ .
for $k = 0, 1, 2, \dots, N$ **do**
 $z_{k+1} = \text{prox}_{z_k}(\gamma_k \cdot g(z_k, \mathbf{e}_k, \tau, \xi_k^\pm))$.
end for
Output: \bar{z}_N .

account the geometric setup of the problem. \mathbf{e}_k and ξ_k^\pm are generated independently of the previous iterations and of each other. Here $\bar{z}_N = \frac{1}{N+1} \sum_{i=0}^N z_i$. Below we give technical facts about (5). Note that we do not provide proofs in the main part of the paper, they are all in the Appendix.

Lemma 1 (see Lemma 2 from [4] or Lemma 1 from [5]). For $g(z, \mathbf{e}, \tau, \xi^\pm)$ defined in (5) under Assumptions 2 and 4 the following inequality holds:

$$\mathbb{E} [\|g(z, \mathbf{e}, \tau, \xi^\pm)\|_q^2] \leq 3a_q^2 \left(3nM^2 + \frac{n^2(\sigma^2 + \Delta^2)}{\tau^2} \right), \quad (6)$$

where a_q^2 is determined by $\mathbb{E}[\|e\|_q^2] \leq \sqrt{\mathbb{E}[\|e\|_q^4]} \leq a_q^2$ and the following statement is true

$$a_q^2 = \min\{2q - 1, 32 \log n - 8\} n^{\frac{2}{q}-1}, \quad \forall n \geq 3. \quad (7)$$

Next we define an important object for further theoretical discussion – a smoothed version of the function φ (see [13,16]).

Definition 5. Function $\hat{\varphi}(z)$ defines on set \mathcal{Z} satisfies:

$$\hat{\varphi}(z) = \mathbb{E}_{\mathbf{e}} [\varphi(z + \tau \mathbf{e})]. \quad (8)$$

To define smoothed version correctly it is important that the function φ is specified not only on an admissible set \mathcal{Z} , but in a certain neighborhood of it. This is due to the fact that for any point z belonging to the set, the point $z + \tau \mathbf{e}$ can be outside it.

Lemma 2 (see Lemma 8 from [16]). Let $\varphi(z)$ is μ -strongly-convex-strongly-concave (convex-concave with $\mu = 0$) and \mathbf{e} be from $\mathcal{RS}_2^n(1)$. Then function $\hat{\varphi}(z)$ is μ -strongly-convex-strongly-concave and under Assumption 2 satisfies:

$$\sup_{z \in \mathcal{Z}} |\hat{\varphi}(z) - \varphi(z)| \leq \tau M. \quad (9)$$

Lemma 3 (see Lemma 10 from [16] and Lemma 2 from [4]). Under Assumption 4 it holds that

$$\tilde{\nabla} \hat{\varphi}(z) = \mathbb{E}_{\mathbf{e}} \left[\frac{n(\varphi(z + \tau \mathbf{e}) - \varphi(z - \tau \mathbf{e}))}{2\tau} \begin{pmatrix} \mathbf{e}_x \\ -\mathbf{e}_y \end{pmatrix} \right], \quad (10)$$

$$\|\mathbb{E}_{\mathbf{e}, \xi} [g(z, \mathbf{e}, \tau, \xi^\pm)] - \tilde{\nabla} \hat{\varphi}(z)\|_q \leq \frac{\Delta n a_q}{\tau}. \quad (11)$$

Now we are ready to present the main results of this section. Let begin with **convex-concave** case (Assumption 3 with $\mu = 0$)

Theorem 1. Let problem (1) with function $\varphi(x, y)$ be solved using Algorithm 1 with the oracle (5). Assume, that the set \mathcal{Z} , the convex-concave function $\varphi(x, y)$ and its inexact modification $\tilde{\varphi}(x, y)$ satisfy Assumptions 1, 2, 4. Denote by N the number of iterations and $\gamma_k = \gamma = \text{const}$. Then the rate of convergence is given by the following expression:

$$\mathbb{E} [\varepsilon_{sad}(\bar{z}_N)] \leq \frac{3\Omega^2}{2\gamma(N+1)} + \frac{3\gamma M_{all}^2}{2} + \frac{\Delta \Omega n a_q}{\tau} + 2\tau M.$$

Ω is a diameter of \mathcal{Z} , $M_{all}^2 = 3 \left(3nM^2 + \frac{n^2(\sigma^2 + \Delta^2)}{\tau^2} \right) a_q^2$ and

$$\varepsilon_{sad}(\bar{z}_N) = \max_{y' \in \mathcal{Y}} \varphi(\bar{x}_N, y') - \min_{x' \in \mathcal{X}} \varphi(x', \bar{y}_N). \quad (12)$$

Let analyze the results:

Corollary 1. *Under the assumptions of the Theorem 1 let ε be accuracy of the solution of the problem (1) obtained using Algorithm 1. Assume that*

$$\gamma = \Theta \left(\frac{\Omega}{n^{\frac{1}{4} + \frac{1}{2q}} MN^{\frac{3}{4}}} \right), \quad \tau = \Theta \left(\frac{\sigma}{M} \cdot \frac{n^{\frac{1}{4} + \frac{1}{2q}}}{N^{\frac{1}{4}}} \right), \quad \Delta = \mathcal{O} \left(\frac{\varepsilon\tau}{\Omega n a_q} \right), \quad (13)$$

then the number of iterations to find ε -solution

$$N = \mathcal{O} \left(\frac{n^{1 + \frac{2}{q}}}{\varepsilon^4} [C^4(n, q) M^4 \Omega^4 + \sigma^4] \right),$$

or with

$$\gamma = \Theta \left(\frac{\Omega}{n^{\frac{1}{q}} MN^{\frac{3}{4}}} \right), \quad \tau = \Theta \left(\frac{\sigma}{M} \cdot \frac{n^{\frac{1}{2}}}{N^{\frac{1}{4}}} \right), \quad \Delta = \mathcal{O} \left(\frac{\varepsilon\tau}{\Omega n a_q} \right),$$

$$N = \mathcal{O} \left(\frac{n^{\frac{4}{q}} C^4(n, q) M^4 \Omega^4 + \frac{n^2}{\varepsilon^4} \sigma^4}{\varepsilon^4} \right),$$

where $C(n, q) \stackrel{def}{=} \min\{2q - 1, 32 \log n - 8\}$.

Analyse separately cases with $p = 1$ and $p = 2$.

$p, (1 \leq p \leq 2)$	$q, (2 \leq q \leq \infty)$	N , Number of iterations
$p = 2$	$q = 2$	$\mathcal{O}(n^2 \varepsilon^{-4})$
$p = 1$	$q = \infty$	$\mathcal{O}(n \log^4 n \cdot \varepsilon^{-4})$

Table 3. Summary of convergence estimation for non-smooth case: $p = 2$ and $p = 1$.

Next we consider μ -strongly-convex-strongly-concave. Here we work with $V_z(w) = \frac{1}{2} \|z - w\|_2^2$.

Theorem 2. *Let problem (1) with function $\varphi(x, y)$ be solved using Algorithm 1 with $V_z(w) = \frac{1}{2} \|z - w\|_2^2$ and the oracle (5). Assume, that the set \mathcal{Z} , the function $\varphi(x, y)$ and its inexact modification $\tilde{\varphi}(x, y)$ satisfy Assumptions 1, 2,*

3, 4. Denote by N the number of iterations and $\gamma_k = \frac{1}{\mu k}$. Then the rate of convergence is given by the following expression:

$$\mathbb{E} [\varphi(\bar{x}_N, y^*) - \varphi(x^*, \bar{y}_N)] \leq \frac{M_{all}^2 \log(N+1)}{2\mu(N+1)} + \frac{\Delta n \Omega}{\tau} + 2\tau M$$

Ω is a diameter of \mathcal{Z} , $M_{all}^2 = 3 \left(3nM^2 + \frac{n^2(\sigma^2 + \Delta^2)}{\tau^2} \right)$.

From here one can get

Corollary 2. Under the assumptions of the Theorem 2 let ε be accuracy of the solution of the problem (1) obtained using Algorithm 1. Assume that

$$\tau = \Theta \left(\sqrt[3]{\frac{\sigma^2}{\mu M}} \cdot \sqrt[3]{\frac{n^2}{N}} \right), \quad \Delta = \mathcal{O} \left(\frac{\varepsilon \tau}{\Omega n} \right),$$

then the number of iterations to find ε -solution

$$N = \tilde{\mathcal{O}} \left(\frac{nM^2}{\mu\varepsilon} + \frac{M^2 n^2 \sigma^2}{\mu\varepsilon^3} \right).$$

Random euclidean direction with residual feedback. In this part of the work we use the technique from [17]. In more detail, in Algorithm 1 we replace $g(z_k, \mathbf{e}_k, \tau, \xi_k^\pm)$ with

$$\begin{aligned} & \tilde{g}(z_k, z_{k-1}, \mathbf{e}_k, \mathbf{e}_{k-1}, \xi_k, \xi_{k-1}) \\ &= \frac{n(\tilde{\varphi}(z_k + \tau \mathbf{e}_k, \xi_k) - \tilde{\varphi}(z_{k-1} + \tau \mathbf{e}_{k-1}, \xi_{k-1}))}{\tau} \begin{pmatrix} (\mathbf{e}_k)_x \\ -(\mathbf{e}_k)_y \end{pmatrix}. \end{aligned} \quad (14)$$

The main advantage of this technique is that it requires only one call to the oracle per iteration.

We consider only convex-concave case in the Euclidean setup, i.e. $V_z(w) = \frac{1}{2} \|z - w\|_2^2$. Let us carry out reasoning similar to the analysis of Theorem 1.

Lemma 4. For $\tilde{g}_k \stackrel{def}{=} \tilde{g}(z_k, z_{k-1}, \mathbf{e}_k, \mathbf{e}_{k-1}, \xi_k, \xi_{k-1})$ defined in (14) under Assumptions 2 and 4 the following inequalities holds:

$$\mathbb{E} [\|\tilde{g}_k\|_2^2] \leq \alpha^k \mathbb{E} [\|\tilde{g}_0\|_2^2] + \left(\frac{12n^2(\sigma^2 + \Delta^2)}{\tau^2} + 12n^2 M^2 \right) \frac{1}{1 - \alpha}, \quad (15)$$

where $\alpha = \frac{6\gamma^2 n^2 M^2}{\tau^2} < 1$.

Lemma 5. Under Assumption 4 it holds that

$$\tilde{\nabla} \hat{\varphi}(z_k) = \mathbb{E}_{\mathbf{e}_k} \left[\frac{n(\varphi(z_k + \tau \mathbf{e}_k) - \varphi(z_k + \tau \mathbf{e}_{k-1}))}{\tau} \begin{pmatrix} (\mathbf{e}_k)_x \\ -(\mathbf{e}_k)_y \end{pmatrix} \right], \quad (16)$$

$$\|\mathbb{E}_{\mathbf{e}_k} [\tilde{g}_k] - \tilde{\nabla} \hat{\varphi}(z_k)\|_2 \leq \frac{\Delta n}{\tau}. \quad (17)$$

Theorem 3. *Let problem (1) with function $\varphi(x, y)$ be solved using Algorithm 1 with $V_z(w) = \frac{1}{2}\|z - w\|_2^2$ and the oracle (14). Assume, that the set \mathcal{Z} , the convex-concave function $\varphi(x, y)$ and its inexact modification $\tilde{\varphi}(x, y)$ satisfy Assumptions 1, 2, 4. Denote by N the number of iterations and $\gamma_k = \gamma = \text{const}$. Then the rate of convergence is given by the following expression:*

$$\begin{aligned} \mathbb{E}[\varepsilon_{sad}(\bar{z}_N)] &\leq \frac{3\Omega^2}{2\gamma(N+1)} + \frac{3\gamma}{2(N+1)(1-\alpha)} \mathbb{E}[\|\tilde{g}_0\|_2^2] \\ &\quad + \frac{3\gamma}{2(1-\alpha)} \left(\frac{12n^2(\sigma^2 + \Delta^2)}{\tau^2} + 12n^2M^2 \right) + 2\tau M + \frac{\Delta\Omega n}{\tau}. \end{aligned}$$

Ω is a diameter of \mathcal{Z} , $\alpha = \frac{6\gamma^2 n^2 M^2}{\tau^2} < 1$.

Next we analyze the results:

Corollary 3. *Under the assumptions of the Theorem 3 let ε be accuracy of the solution of the problem (1) obtained using Algorithm 1 with (14). Assume that*

$$\gamma = \left(\frac{\Omega\tau}{6nMN^{\frac{1}{2}}} \right), \quad \tau = \Theta \left(\frac{\sigma}{M} \cdot \frac{n^{\frac{1}{2}}}{N^{\frac{1}{4}}} \right), \quad \Delta = \mathcal{O} \left(\frac{\varepsilon\tau}{\Omega n} \right),$$

then the number of iterations to find ε -solution

$$N = \mathcal{O} \left(\frac{n^2}{\varepsilon^4} [M^4\Omega^4 + \sigma^4] \right).$$

3.2 Smooth case

Assumption 5 (Gradient's Lipschitz continuity) *The gradient $\nabla\varphi(z)$ of the function φ is L -Lipschitz continuous in certain neighbourhood of \mathcal{Z} with $L > 0$ w.r.t. norm $\|\cdot\|_2$ when*

$$|\nabla\varphi(z) - \nabla\varphi(z')| \leq L\|z - z'\|_2, \quad \forall z, z' \in \mathcal{Z}.$$

Lemma 6 (see Lemma A.3 from [1]). *Let $\varphi(z)$ be convex-concave (or μ -strongly-convex-strongly-concave) and \mathbf{e} be from $\mathcal{RS}_2^n(1)$. Then function $\hat{\varphi}(z)$ is convex-concave (μ -strongly-convex-strongly-concave) too and under Assumption 5 satisfies:*

$$\sup_{z \in \mathcal{Z}} |\hat{\varphi}(z) - \varphi(z)| \leq \frac{L\tau^2}{2}. \quad (18)$$

Theorem 4. *Let problem (1) with function $\varphi(x, y)$ be solved using Algorithm 1 with the oracle (5). Assume, that the set \mathcal{Z} , the convex-concave function $\varphi(x, y)$ and its inexact modification $\tilde{\varphi}(x, y)$ satisfy Assumptions 1, 4, 5. Denote by N the number of iterations and $\gamma_k = \gamma = \text{const}$. Then the rate of convergence is given by the following expression:*

$$\mathbb{E}[\varepsilon_{sad}(\bar{z}_N)] \leq \frac{3\Omega^2}{2\gamma(N+1)} + \frac{3\gamma M_{all}^2}{2} + \frac{\Delta\Omega n a_q}{\tau} + L\tau^2.$$

Ω is a diameter of \mathcal{Z} , $M_{\text{all}}^2 = 3 \left(3nM^2 + \frac{n^2(\sigma^2 + \Delta^2)}{\tau^2} \right) a_q^2$.

Let's analyze the results:

Corollary 4. *Under the assumptions of the Theorem 4 let ε be accuracy of the solution of the problem (1) obtained using Algorithm 1. Assume that*

$$\gamma = \Theta \left(\frac{\Omega}{n^{\frac{1}{3} + \frac{2}{3q}} MN^{\frac{2}{3}}} \right), \quad \tau = \Theta \left(\frac{\sigma}{M} \cdot \frac{n^{\frac{1}{6} + \frac{1}{3q}}}{N^{\frac{1}{6}}} \right), \quad \Delta = \mathcal{O} \left(\frac{\varepsilon\tau}{\Omega n a_q} \right), \quad (19)$$

then the number of iterations to find ε -solution

$$N = \mathcal{O} \left(\frac{n^{1 + \frac{2}{q}}}{\varepsilon^3} \left[M^3 \Omega^3 + \frac{L^3 \sigma^3}{M^3} \right] \right).$$

Theorem 5. *Let problem (1) with function $\varphi(x, y)$ be solved using Algorithm 1 with $V_z(w) = \frac{1}{2} \|z - w\|_2^2$ and the oracle (5). Assume, that the set \mathcal{Z} , the function $\varphi(x, y)$ and its inexact modification $\tilde{\varphi}(x, y)$ satisfy Assumptions 1, 3, 4, 5. Denote by N the number of iterations and $\gamma_k = \frac{1}{\mu k}$. Then the rate of convergence is given by the following expression:*

$$\mathbb{E} [\varphi(\bar{x}_N, y^*) - \varphi(x^*, \bar{y}_N)] \leq \frac{M_{\text{all}}^2 \log(N+1)}{2\mu(N+1)} + \frac{\Delta n \Omega}{\tau} + L\tau^2.$$

Ω is a diameter of \mathcal{Z} , $M_{\text{all}}^2 = 3 \left(3nM^2 + \frac{n^2(\sigma^2 + \Delta^2)}{\tau^2} \right)$.

Let's analyze the results:

Corollary 5. *Under the assumptions of the Theorem 5 let ε be accuracy of the solution of the problem (1) obtained using Algorithm 1. Assume that*

$$\tau = \Theta \left(\sqrt[4]{\frac{\sigma^2}{\mu L}} \cdot \frac{n^{\frac{1}{2}}}{N^{\frac{1}{4}}} \right), \quad \Delta = \mathcal{O} \left(\frac{\varepsilon\tau}{\Omega n} \right),$$

then the number of iterations to find ε -solution

$$N = \tilde{\mathcal{O}} \left(\frac{nM^2}{\mu\varepsilon} + \frac{Ln^2\sigma^2}{\mu\varepsilon^2} \right).$$

3.3 Higher-order smooth case

In this paragraph we study higher-order smooth functions φ functions satisfying so called generalized Hölder condition with parameter $\beta > 2$ (see inequality (20) below).

Higher order smoothness Let l denote maximal integer number strictly less than β . Let $\mathcal{F}_\beta(L_\beta)$ denote the set of all functions $\varphi : \mathbb{R}^n \rightarrow \mathbb{R}$ which are differentiable l times and for all $z, z_0 \in U_{\varepsilon_0}(\mathcal{Z})$ satisfy Hölder condition:

$$\left| \varphi(z) - \sum_{0 \leq |m| \leq l} \frac{1}{m!} D^m \varphi(z_0) (z - z_0)^m \right| \leq L_\beta \|z - z_0\|^\beta, \quad (20)$$

where $L_\beta > 0$, the sum is over multi-index $m = (m_1, \dots, m_n) \in \mathbb{N}^n$, we use the notation $m! = m_1! \dots m_n!$, $|m| = m_1 + \dots + m_n$ and we defined

$$D^m \varphi(z_0) z^m = \frac{\partial^{|m|} \varphi(z_0)}{\partial^{m_1} z_1 \dots \partial^{m_n} z_n} z_1^{m_1} \dots z_n^{m_n}, \quad \forall z = (z_1, \dots, z_n) \in \mathbb{R}^n.$$

Let $\mathcal{F}_{\mu, \beta}(L_\beta)$ denote the set of μ -strongly-convex-strongly-concave functions $\varphi \in \mathcal{F}_\beta(L_\beta)$.

To use the higher-order smoothness we propose smoothing kernel though this is not the only way. We propose to use Algorithm 2 which uses the kernel smoothing technique. In fact the Algorithm 2 arises from Algorithm 1 in the Euclidean setting ($V_z(w) = \frac{1}{2} \|z - w\|_2^2$).

Algorithm 2 Zero-order Stochastic Projected Gradient

Requires: Kernel $K : [-1, 1] \rightarrow \mathbb{R}$, step size $\gamma_k > 0$, parameters τ_k .

Initialization: Generate scalars r_1, \dots, r_N uniformly on $[-1, 1]$ and vectors e_1, \dots, e_N uniformly on the Euclidean unit sphere $S_n = \{e \in \mathbb{R}^n : \|e\| = 1\}$.

for $k = 1, \dots, N$ **do**

1. $\tilde{\varphi}_k^+ := \varphi(z_k + \tau_k r_k e_k) + \xi_k^+$, $\tilde{\varphi}_k^- := \varphi(z_k - \tau_k r_k e_k) + \xi_k^-$

2. Define $\tilde{g}_k := \frac{n}{2\tau_k} (\tilde{\varphi}_k^+ - \tilde{\varphi}_k^-) \begin{pmatrix} (\mathbf{e}_k)_x \\ -(\mathbf{e}_k)_y \end{pmatrix} K(r_k)$

3. Update $z_{k+1} := \Pi_Q(z_k - \gamma_k \tilde{g}_k)$

end for

Output: $\{z_k\}_{k=1}^N$.

To use the higher-order smoothness we propose we need to introduce additional noise assumption:

Assumption 6 For all $k = 1, 2, \dots, N$ it holds that

1. $\mathbb{E}[\xi_k^{+2}] \leq \sigma^2$ and $\mathbb{E}[\xi_k^{-2}] \leq \sigma^2$ where $\sigma \geq 0$;
2. the random variables ξ_k^+ and ξ_k^- are independent from e_k and r_k , the random variables e_k and r_k are independent.

In other words we assume that $\delta(x, y)$ in (3) is equal to zero. We do not assume here neither zero-mean of ξ_k^+ and ξ_k^- nor i.i.d of $\{\xi_k^+\}_{k=1}^N$ and $\{\xi_k^-\}_{k=1}^N$ as item 2 from Assumption 6 allows to avoid that.

Kernel For gradient estimator \tilde{g}_k we use the kernel

$$K : [-1, 1] \rightarrow \mathbb{R},$$

satisfying

$$\mathbb{E}[K(r)] = 0, \mathbb{E}[rK(r)] = 1, \mathbb{E}[r^j K(r)] = 0, j = 2, \dots, l, \mathbb{E}[|r|^\beta |K(r)|] \leq \infty, \quad (21)$$

where r is a uniformly distributed on $[-1, 1]$ random variable. This helps us to get better bounds on the gradient bias $\|\tilde{g}_k - \nabla f(x_k)\|$ (see Theorem 6 for details). The examples of possible kernels are presented in Appendix E.

For Theorem 6 and Theorem 7 we need to introduce the constants

$$\kappa_\beta = \int |u|^\beta |K(u)| du \quad (22)$$

and

$$\kappa = \int K^2(u) du. \quad (23)$$

It is proved in [2] that κ_β and κ do not depend on n , they depend only on β :

$$\kappa_\beta \leq 2\sqrt{2}(\beta - 1), \quad (24)$$

$$\kappa \leq \sqrt{3}\beta^{3/2}. \quad (25)$$

Theorem 6. *Let $\varphi \in \mathcal{F}_{\mu, \beta}(L)$ with $\mu, L > 0$ and $\beta > 2$. Let Assumption 6 hold and let \mathcal{Z} be a convex compact subset of \mathbb{R}^n . Let φ be M -Lipschitz on the Euclidean τ_1 -neighborhood of \mathcal{Z} (see τ_k below).*

Then the rate of convergence is given by Algorithm 2 with parameters

$$\tau_k = \left(\frac{3\kappa\sigma^2 n}{2(\beta - 1)(\kappa_\beta L)^2} \right)^{\frac{1}{2\beta}} k^{-\frac{1}{2\beta}}, \quad \alpha_k = \frac{2}{\mu k}, \quad k = 1, \dots, N$$

satisfies

$$\begin{aligned} \mathbb{E}[\varphi(\bar{x}_N, y^*) - \varphi(x^*, \bar{y}_N)] &\leq \max_{y \in \mathcal{Y}} \mathbb{E}[\varphi(\bar{x}_N, y)] - \min_{x \in \mathcal{X}} \mathbb{E}[\varphi(x, \bar{y}_N)] \\ &\leq \frac{1}{\mu} \left(n^{2-\frac{1}{\beta}} \frac{A_1}{N^{\frac{\beta-1}{\beta}}} + A_2 \frac{n(1 + \ln N)}{N} \right), \end{aligned}$$

where $\bar{z}_N = \frac{1}{N} \sum_{k=1}^N z_k$, $A_1 = 3\beta(\kappa\sigma^2)^{\frac{\beta-1}{\beta}} (\kappa_\beta L)^{\frac{2}{\beta}}$, $A_2 = 9\kappa G^2$, κ_β and κ are constants depending only on β , see (22) and (23).

We emphasize that the usage of kernel smoothing technique, measure concentration inequalities and the assumption that ξ_k is independent from e_k or r_k (Assumption 6) lead to the results better than the state-of-the-art ones for $\beta > 2$. The last assumption also allows us not to assume neither zero-mean of ξ_k^+ and ξ_k^- nor i.i.d of $\{\xi_k^+\}_{k=1}^N$ and $\{\xi_k^-\}_{k=1}^N$.

Theorem 7. Let $\varphi \in \mathcal{F}_\beta(L)$ with $L > 0$ and $\beta > 2$. Let Assumption 6 hold and let \mathcal{Z} be a convex compact subset of \mathbb{R}^n . Let φ be M -Lipschitz on the Euclidean τ_1 -neighborhood of \mathcal{Z} (τ_k is parameter from Theorem 6 for the regularized function $\varphi_\mu(z)$ whose description is given below). Let \bar{z}_N denote $\frac{1}{N} \sum_{k=1}^N z_k$.

Let's define $N(\varepsilon)$:

$$N(\varepsilon) = \max \left\{ \left(R\sqrt{2A_1} \right)^{\frac{2\beta}{\beta-1}} \frac{n^{2+\frac{1}{\beta-1}}}{\varepsilon^{2+\frac{2}{\beta-1}}}, \left(R\sqrt{2c'A_2} \right)^{2(1+\rho)} \frac{n^{1+\rho}}{\varepsilon^{2(1+\rho)}} \right\},$$

where $A_1 = 3\beta(\kappa\sigma^2)^{\frac{\beta-1}{\beta}} (\kappa_\beta L)^{\frac{2}{\beta}}$, $A_2 = 9\kappa G^2$ – constants from Theorem 6, $\rho > 0$ – arbitrarily small positive number, c' – constant which depends on ρ .

Then the rate of convergence is given by the following expression:

$$\mathbb{E} [\varphi(\bar{x}_N, y^*) - \varphi(x^*, \bar{y}_N)] \leq \max_{y \in \mathcal{Y}} \mathbb{E} [\varphi(\bar{x}_N, y)] - \min_{x \in \mathcal{X}} \mathbb{E} [\varphi(x, \bar{y}_N)] \leq \varepsilon \quad (26)$$

after $N(\varepsilon)$ steps of Algorithm 2 with settings from Theorem 6 for the regularized function: $\varphi_\mu(z) := \varphi(z) + \frac{\mu}{2} \|x - x_0\|^2 - \frac{\mu}{2} \|y - y_0\|^2$, where $\mu \leq \frac{\varepsilon}{R^2}$, $R = \|z_0 - z^*\|$, $z_0 \in \mathcal{Z}$ – arbitrary point.

4 Experiments

In our experiments we consider the classical bilinear problem on a probability simplex:

$$\min_{x \in \Delta_n} \max_{y \in \Delta_k} [y^T C x], \quad (27)$$

This problem has many different applications and interpretations, one of the main ones is a matrix game (see Part 5 in [3]), i.e. the element c_{ij} of the matrix are interpreted as a winning, provided that player X has chosen the i th strategy and player Y has chosen the j th strategy, the task of one of the players is to maximize the gain, and the opponent's task – to minimize.

The step of our algorithms can be written as follows (see [5]):

$$[x_{k+1}]_i = \frac{[x_k]_i \exp(-\gamma_k [g_x]_i)}{\sum_{j=1}^n [x_k]_j \exp(-\gamma_k [g_x]_j)}, \quad [y_{k+1}]_i = \frac{[y_k]_i \exp(\gamma_k [g_y]_i)}{\sum_{j=1}^n [y_k]_j \exp(\gamma_k [g_y]_j)},$$

where under g_x, g_y we mean parts of g which are responsible for x and for y . Note that we do not present a generalization of Algorithm 2 in an arbitrary Bregman setup, but we want to check in practice.

We take matrix 50×50 . All elements of the matrix are generated from the uniform distribution from 0 to 1. Next, we select one row of the matrix and

generate its elements from the uniform from 5 to 10. Finally, we take one element from this row and generate it uniformly from 1 to 5. Finally, the matrix is normalized. Further, with each call of the function value $y^T Cx$ we add stochastic noise with constant variance (which is on average 5% or 10% of the function value).

The main goal of our experiments is to compare three gradient-free approaches: Algorithm 1 with (5) and (14) approximations, as well as Algorithm 2. We also added a first order method for comparison. Parameters γ and τ are selected with the help of grid-search so that the convergence is the fastest, but stable. See Figure 1 for results.

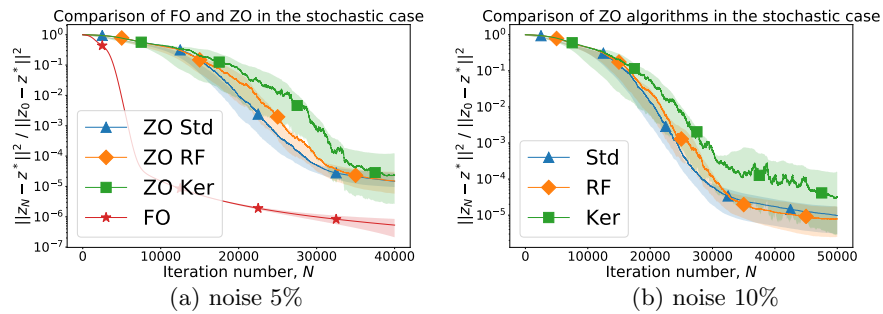


Fig. 1. Algorithm 1 with (5) (ZO Std) and (14) (ZO RF) approximations, Algorithm 2 (ZO Ker) and Mirror Descent (FO) applied to solve saddle-problem (27) with noise level: (a) 5%, (b) 10%.

Based on the results of the experiments, we note that the gradient-free methods converge more slowly than the first-order method – which is predictable. The convergence of zeroth-order methods is approximately the same, the only thing that can be noted is that the method with a kernel is subject to larger fluctuations.

References

1. Akhavan, A., Pontil, M., Tsybakov, A.B.: Exploiting higher order smoothness in derivative-free optimization and continuous bandits. arXiv preprint arXiv:2006.07862 (2020)
2. Bach, F., Perchet, V.: Highly-smooth zero-th order online optimization. In: Conference on Learning Theory. pp. 257–283. PMLR (2016)
3. Ben-Tal, A., Nemirovski, A.: Lectures on Modern Convex Optimization: Analysis, Algorithms, and Engineering Applications (2019)
4. Beznosikov, A., Gorbunov, E., Gasnikov, A.: Derivative-free method for composite optimization with applications to decentralized distributed optimization. arXiv preprint arXiv:1911.10645 (2019)

5. Beznosikov, A., Sadiev, A., Gasnikov, A.: Gradient-free methods for saddle-point problem. arXiv preprint arXiv:2005.05913 (2020)
6. Carmon, Y., Jin, Y., Sidford, A., Tian, K.: Coordinate methods for matrix games. arXiv preprint arXiv:2009.08447 (2020)
7. Chen, P.Y., Zhang, H., Sharma, Y., Yi, J., Hsieh, C.J.: Zoo. Proceedings of the 10th ACM Workshop on Artificial Intelligence and Security - AISec '17 (2017). <https://doi.org/10.1145/3128572.3140448>, <http://dx.doi.org/10.1145/3128572.3140448>
8. Duchi, J.C., Jordan, M.I., Wainwright, M.J., Wibisono, A.: Optimal rates for zero-order convex optimization: the power of two function evaluations. arXiv preprint arXiv:1312.2139 (2013)
9. Fazel, M., Ge, R., Kakade, S., Mesbahi, M.: Global convergence of policy gradient methods for the linear quadratic regulator. In: International Conference on Machine Learning. pp. 1467–1476. PMLR (2018)
10. Gasnikov, A.V., Krymova, E.A., Lagunovskaya, A.A., Usmanova, I.N., Fedorenko, F.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)
11. Goodfellow, I.: Nips 2016 tutorial: Generative adversarial networks. arXiv preprint arXiv:1701.00160 (2016)
12. Jin, Y., Sidford, A.: Efficiently solving MDPs with stochastic mirror descent. In: III, H.D., Singh, A. (eds.) Proceedings of the 37th International Conference on Machine Learning. Proceedings of Machine Learning Research, vol. 119, pp. 4890–4900. PMLR (13–18 Jul 2020)
13. Nesterov, Y., Spokoiny, V.G.: Random gradient-free minimization of convex functions. Foundations of Computational Mathematics **17**(2), 527–566 (2017)
14. Novitskii, V., Gasnikov, A.: Improved exploiting higher order smoothness in derivative-free optimization and continuous bandit. arXiv preprint arXiv:2101.03821 (2021)
15. Sadiev, A., Beznosikov, A., Dvurechensky, P., Gasnikov, A.: Zeroth-order algorithms for smooth saddle-point problems. arXiv preprint arXiv:2009.09908 (2020)
16. Shamir, O.: An optimal algorithm for bandit and zero-order convex optimization with two-point feedback. Journal of Machine Learning Research **18**(52), 1–11 (2017)
17. Zhang, Y., Zhou, Y., Ji, K., Zavlanos, M.M.: Improving the convergence rate of one-point zeroth-order optimization using residual feedback. arXiv preprint arXiv:2006.10820 (2020)

A General facts

Lemma 7 (see inequality 5.3.18 from [3]). *Let $d(z) : \mathcal{Z} \rightarrow \mathbb{R}$ is prox-function and $V_z(w)$ define Bregman divergence associated with $d(z)$. The following equation holds for $x, y, u \in X$:*

$$\langle \nabla d(x) - \nabla d(y), u - x \rangle = V_y(u) - V_x(u) - V_y(x). \quad (28)$$

Lemma 8 (Fact 5.3.2 from [3]). *Given norm $\|\cdot\|$ on space \mathcal{Z} and prox-function $d(z)$, let $z \in \mathcal{Z}$, $w \in \mathbb{R}^n$ and $z_+ = \text{prox}_z(w)$. Then for all $u \in \mathcal{Z}$*

$$\langle w, z_+ - u \rangle \leq V_z(u) - V_{z_+}(u) - V_z(z_+). \quad (29)$$

Lemma 9. For arbitrary integer $n \geq 1$ and arbitrary set of positive numbers a_1, \dots, a_n we have

$$\left(\sum_{i=1}^m a_i \right)^2 \leq m \sum_{i=1}^m a_i^2. \quad (30)$$

Lemma 10 (Lemma 9 from [16]). For any function g which is M -Lipschitz with respect to the ℓ_2 -norm, it holds that if e is uniformly distributed on the Euclidean unit sphere, then

$$\sqrt{\mathbb{E}[(g(e) - \mathbb{E}g(e))^4]} \leq \frac{3M^2}{n}.$$

B Proofs for Section 3.1

Lemma 11. For $g(z, \mathbf{e}, \tau, \xi^\pm)$ defined in (5) under Assumptions 2 and 4 the following inequality holds:

$$\mathbb{E} [\|g(z, \mathbf{e}, \tau, \xi^\pm)\|_q^2] \leq 3a_q^2 \left(3nM^2 + \frac{n^2(\sigma^2 + \Delta^2)}{\tau^2} \right),$$

where a_q^2 is determined by $\mathbb{E}[\|e\|_q^2] \leq \sqrt{\mathbb{E}[\|e\|_q^4]} \leq a_q^2$ and the following statement is true

$$a_q^2 = \min\{2q - 1, 32 \log n - 8\} n^{\frac{2}{q}-1}, \quad \forall n \geq 3.$$

Proof. Using a simple fact (30), we obtain the following inequalities:

$$\begin{aligned} \mathbb{E} [\|g(z, \mathbf{e}, \tau, \xi^\pm)\|_q^2] &= \mathbb{E} \left[\left\| \frac{n}{2\tau} (\tilde{\varphi}(z + \tau\mathbf{e}, \xi^+) - \tilde{\varphi}(z - \tau\mathbf{e}, \xi^-)) \mathbf{e} \right\|_q^2 \right] \\ &= \mathbb{E} \left[\left\| \frac{n}{2\tau} (\varphi(z + \tau\mathbf{e}) + \xi^+ + \delta(z + \tau\mathbf{e}) - \varphi(z - \tau\mathbf{e}) - \xi^- - \delta(z - \tau\mathbf{e})) \mathbf{e} \right\|_q^2 \right] \\ &\leq \frac{3n^2}{4\tau^2} \mathbb{E} [\|(\varphi(z + \tau\mathbf{e}) - \varphi(z - \tau\mathbf{e})) \mathbf{e}\|_q^2] + \frac{3n^2}{4\tau^2} \mathbb{E} [\|(\xi^+ - \xi^-) \mathbf{e}\|_q^2] \\ &\quad + \frac{3n^2}{4\tau^2} \mathbb{E} [\|(\delta(z + \tau\mathbf{e}) - \delta(z - \tau\mathbf{e})) \mathbf{e}\|_q^2] \\ &\leq \frac{3n^2}{4\tau^2} \mathbb{E} [(\varphi(z + \tau\mathbf{e}, \xi) - \varphi(z - \tau\mathbf{e}, \xi))^2 \|\mathbf{e}\|_q^2] + \frac{3n^2}{2\tau^2} \mathbb{E} [((\xi^+)^2 + (\xi^-)^2) \|\mathbf{e}\|_q^2] \\ &\quad + \frac{3n^2}{2\tau^2} \mathbb{E} [((\delta(z + \tau\mathbf{e}))^2 + (\delta(z - \tau\mathbf{e}))^2) \|\mathbf{e}\|_q^2]. \end{aligned}$$

By independence of ξ^\pm and \mathbf{e} , we have

$$\begin{aligned}
\mathbb{E} [\|g(z, \mathbf{e}, \tau, \xi^\pm)\|_q^2] &\leq \frac{3n^2}{4\tau^2} \mathbb{E}_\xi \left[\mathbb{E}_\mathbf{e} \left[(\varphi(z + \tau\mathbf{e}) - \alpha - \varphi(z - \tau\mathbf{e}) + \alpha)^2 \|\mathbf{e}\|_q^2 \right] \right] \\
&\quad + \frac{3n^2}{2\tau^2} \mathbb{E}_\mathbf{e} \left[\mathbb{E}_\xi \left[((\xi^+)^2 + (\xi^-)^2) \|\mathbf{e}\|_q^2 \right] \right] \\
&\quad + \frac{3n^2}{2\tau^2} \mathbb{E} \left[((\delta(z + \tau\mathbf{e}))^2 + (\delta(z - \tau\mathbf{e}))^2) \|\mathbf{e}\|_q^2 \right] \\
&\leq \frac{3n^2}{2\tau^2} \mathbb{E}_\xi \left[\mathbb{E}_\mathbf{e} \left[((\varphi(z + \tau\mathbf{e}) - \alpha)^2 + (\varphi(z - \tau\mathbf{e}) - \alpha)^2) \|\mathbf{e}\|_q^2 \right] \right] \\
&\quad + \frac{3n^2}{2\tau^2} \mathbb{E}_\mathbf{e} \left[\mathbb{E}_\xi \left[((\xi^+)^2 + (\xi^-)^2) \|\mathbf{e}\|_q^2 \right] \right] \\
&\quad + \frac{3n^2}{2\tau^2} \mathbb{E} \left[((\delta(z + \tau\mathbf{e}))^2 + (\delta(z - \tau\mathbf{e}))^2) \|\mathbf{e}\|_q^2 \right].
\end{aligned}$$

Taking into account the symmetric distribution of \mathbf{e} and Cauchy–Schwarz inequality:

$$\begin{aligned}
\mathbb{E} [\|g(z, \mathbf{e}, \tau, \xi^\pm)\|_q^2] &\leq \frac{3n^2}{\tau^2} \mathbb{E}_\xi \left[\mathbb{E}_\mathbf{e} \left[(\varphi(z + \tau\mathbf{e}) - \alpha)^2 \|\mathbf{e}\|_q^2 \right] \right] + \frac{3n^2}{2\tau^2} \mathbb{E}_\mathbf{e} \left[\mathbb{E}_\xi \left[((\xi^+)^2 + (\xi^-)^2) \|\mathbf{e}\|_q^2 \right] \right] \\
&\quad + \frac{3n^2}{2\tau^2} \mathbb{E} \left[((\delta(z + \tau\mathbf{e}))^2 + (\delta(z - \tau\mathbf{e}))^2) \|\mathbf{e}\|_q^2 \right] \\
&\leq \frac{3n^2}{\tau^2} \mathbb{E}_\xi \left[\sqrt{\mathbb{E}_\mathbf{e} \left[(\varphi(z + \tau\mathbf{e}, \xi) - \alpha)^4 \right]} \sqrt{\mathbb{E}_\mathbf{e} \left[\|\mathbf{e}\|_q^4 \right]} \right] \\
&\quad + \frac{3n^2}{2\tau^2} \mathbb{E}_\mathbf{e} \left[\mathbb{E}_\xi \left[((\xi^+)^2 + (\xi^-)^2) \|\mathbf{e}\|_q^2 \right] \right] \\
&\quad + \frac{3n^2}{2\tau^2} \mathbb{E} \left[((\delta(z + \tau\mathbf{e}))^2 + (\delta(z - \tau\mathbf{e}))^2) \|\mathbf{e}\|_q^2 \right] \\
&\leq \frac{3n^2 a_q^2}{\tau^2} \mathbb{E}_\xi \left[\sqrt{\mathbb{E}_\mathbf{e} \left[(\varphi(z + \tau\mathbf{e}, \xi) - \alpha)^4 \right]} \right] + \frac{3n^2 a_q^2 (\sigma^2 + \Delta^2)}{\tau^2}.
\end{aligned}$$

In the last inequalities we use (4) and (7). Substituting $\alpha = \mathbb{E}[\varphi(z + \tau\mathbf{e})]$, applying Lemma 10 with the fact that $\varphi(z + \tau\mathbf{e})$ is τM -Lipschitz w.r.t. \mathbf{e} in terms of the $\|\cdot\|_2$ -norm we get

$$\mathbb{E} [\|g(z, \mathbf{e}, \tau, \xi^\pm)\|_q^2] \leq 3a_q^2 \left(3nM^2 + \frac{n^2(\sigma^2 + \Delta^2)}{\tau^2} \right).$$

□

Lemma 12. *Let $\varphi(z)$ is μ -strongly-convex-strongly-concave (convex-concave with $\mu = 0$) and \mathbf{e} be from $\mathcal{RS}_2^n(1)$. Then function $\hat{\varphi}(z)$ is μ -strongly-convex-strongly-concave and under Assumption 2 satisfies:*

$$\sup_{z \in \mathcal{Z}} |\hat{\varphi}(z) - \varphi(z)| \leq \tau M.$$

Proof. Using definition (8) of $\hat{\varphi}$:

$$|\hat{\varphi}(z) - \varphi(z)| = |\mathbb{E}_{\mathbf{e}}[\varphi(z + \tau\mathbf{e})] - \varphi(z)| = |\mathbb{E}_{\mathbf{e}}[\varphi(z + \tau\mathbf{e}) - \varphi(z)]|.$$

Since $\varphi(z)$ is M -Lipschitz, we get

$$|\mathbb{E}_{\mathbf{e}}[\varphi(z + \tau\mathbf{e}) - \varphi(z)]| \leq |\mathbb{E}_{\mathbf{e}}[M\|\tau\mathbf{e}\|_2]| \leq M\tau.$$

□

Lemma 13. *Under Assumption 4 it holds that*

$$\begin{aligned} \tilde{\nabla}\hat{\varphi}(z) &= \mathbb{E}_{\mathbf{e}} \left[\frac{n(\varphi(z + \tau\mathbf{e}) - \varphi(z - \tau\mathbf{e}))}{2\tau} \begin{pmatrix} \mathbf{e}_x \\ -\mathbf{e}_y \end{pmatrix} \right], \\ \|\mathbb{E}_{\mathbf{e},\xi}[g(z, \mathbf{e}, \tau, \xi^\pm)] - \tilde{\nabla}\hat{\varphi}(z)\|_q &\leq \frac{\Delta na_q}{\tau}. \end{aligned} \quad (31)$$

Proof. The proof of (31) is given in [16] and follows from the Stokes' theorem. Then

$$\mathbb{E}_{\mathbf{e},\xi}[g(z, \mathbf{e}, \tau, \xi^\pm)] - \tilde{\nabla}\hat{\varphi}(z) = \mathbb{E}_{\mathbf{e}} \left[\frac{n(\delta(z + \tau\mathbf{e}) - \delta(z - \tau\mathbf{e}))}{2\tau} \begin{pmatrix} \mathbf{e}_x \\ -\mathbf{e}_y \end{pmatrix} \right].$$

Using inequalities (4) and definition of a_q completes the proof.

□

Theorem 1. Let problem (1) with function $\varphi(x, y)$ be solved using Algorithm 1 with the oracle (5). Assume, that the set \mathcal{Z} , the convex-concave function $\varphi(x, y)$ and its inexact modification $\tilde{\varphi}(x, y)$ satisfy Assumptions 1, 2, 4. Denote by N the number of iterations and $\gamma_k = \gamma = \text{const}$. Then the rate of convergence is given by the following expression:

$$\mathbb{E}[\varepsilon_{sad}(\bar{z}_N)] \leq \frac{3\Omega^2}{2\gamma(N+1)} + \frac{3\gamma M_{all}^2}{2} + \frac{\Delta\Omega na_q}{\tau} + 2\tau M.$$

Ω is a diameter of \mathcal{Z} , $M_{all}^2 = 3\left(3nM^2 + \frac{n^2(\sigma^2 + \Delta^2)}{\tau^2}\right) a_q^2$ and

$$\varepsilon_{sad}(\bar{z}_N) = \max_{y' \in \mathcal{Y}} \varphi(\bar{x}_N, y') - \min_{x' \in \mathcal{X}} \varphi(x', \bar{y}_N).$$

Proof. We divided the proof into three steps.

Step 1. Let $g_k \stackrel{\text{def}}{=} \gamma g(z_k, \mathbf{e}_k, \tau, \xi_k^\pm)$. By the step of Algorithm 1, $z_{k+1} = \text{prox}_{z_k}(g_k)$. Taking into account (29), we get that for all $u \in \mathcal{Z}$

$$\langle g_k, z_{k+1} - u \rangle = \langle g_k, z_{k+1} - z_k + z_k - u \rangle \leq V_{z_k}(u) - V_{z_{k+1}}(u) - V_{z_k}(z_{k+1}).$$

By simple transformations:

$$\begin{aligned} \langle g_k, z_k - u \rangle &\leq \langle g_k, z_k - z_{k+1} \rangle + V_{z_k}(u) - V_{z_{k+1}}(u) - V_{z_k}(z_{k+1}) \\ &\leq \langle g_k, z_k - z_{k+1} \rangle + V_{z_k}(u) - V_{z_{k+1}}(u) - \frac{1}{2} \|z_{k+1} - z_k\|_p^2. \end{aligned}$$

In last inequality we use the property of the Bregman divergence: $V_x(y) \geq \frac{1}{2} \|x - y\|_p^2$. Using Hölder's inequality and the fact: $ab - b^2/2 \leq a^2/2$, we have

$$\begin{aligned} \langle g_k, z_k - u \rangle &\leq \|g_k\|_q \|z_k - z_{k+1}\|_p + V_{z_k}(u) - V_{z_{k+1}}(u) - \frac{1}{2} \|z_{k+1} - z_k\|_p^2 \\ &\leq V_{z_k}(u) - V_{z_{k+1}}(u) + \frac{1}{2} \|g_k\|_q^2. \end{aligned} \quad (32)$$

Summing (32) over all k from 0 to N and by the definitions of g_k and Ω (diameter of \mathcal{Z}): $\forall u \in \mathcal{Z}$

$$\gamma \sum_{k=0}^N \langle g(z_k, \mathbf{e}_k, \tau, \xi_k^\pm), z_k - u \rangle \leq \frac{\Omega^2}{2} + \frac{\gamma^2}{2} \sum_{k=0}^N \|g(z_k, \mathbf{e}_k, \tau, \xi_k^\pm)\|_q^2. \quad (33)$$

Let $\Delta_k \stackrel{\text{def}}{=} g(z_k, \mathbf{e}_k, \tau, \xi_k^\pm) - \tilde{\nabla} \hat{\varphi}(z_k)$ and $D(u) \stackrel{\text{def}}{=} \sum_{k=0}^N \gamma \langle \Delta_k, u - z_k \rangle$. Substituting the definition of $D(u)$ in (33), we have for all $u \in \mathcal{Z}$

$$\gamma \sum_{k=0}^N \langle \tilde{\nabla} \hat{\varphi}(z_k), z_k - u \rangle \leq \frac{\Omega^2}{2} + \frac{\gamma^2}{2} \sum_{k=0}^N \|g(z_k, \mathbf{e}_k, \tau, \xi_k^\pm)\|_q^2 + D(u). \quad (34)$$

By $\tilde{\nabla} \hat{\varphi}(z)$ we mean a block vector consisting of two vectors $\nabla_x \hat{\varphi}(x, y)$ and $-\nabla_y \hat{\varphi}(x, y)$.

Step 2. We consider a relationship between functions $\hat{\varphi}(z)$ and $\varphi(z)$. Combining (12) and (9) we get

$$\varepsilon_{sad}(\bar{z}_N) \leq \max_{y' \in \mathcal{Y}} \hat{\varphi}(\bar{x}_N, y') - \min_{x' \in \mathcal{X}} \hat{\varphi}(x', \bar{y}_N) + 2\tau M.$$

Then, by the definition of \bar{x}_N and \bar{y}_N (see (12)), Jensen's inequality and convexity-concavity of $\hat{\varphi}$:

$$\begin{aligned} \varepsilon_{sad}(\bar{z}_N) &\leq \max_{y' \in \mathcal{Y}} \hat{\varphi} \left(\frac{1}{N+1} \left(\sum_{k=0}^N x_k \right), y' \right) - \min_{x' \in \mathcal{X}} \hat{\varphi} \left(x', \frac{1}{N+1} \left(\sum_{k=0}^N y_k \right) \right) \\ &\quad + 2\tau M \\ &\leq \max_{y' \in \mathcal{Y}} \frac{1}{N+1} \sum_{k=0}^N \hat{\varphi}(x_k, y') - \min_{x' \in \mathcal{X}} \frac{1}{N+1} \sum_{k=0}^N \hat{\varphi}(x', y_k) + 2\tau M. \end{aligned}$$

Given the fact of linear independence of x' and y' :

$$\varepsilon_{sad}(\bar{z}_N) \leq \max_{(x', y') \in \mathcal{Z}} \frac{1}{N+1} \sum_{k=0}^N (\hat{\varphi}(x_k, y') - \hat{\varphi}(x', y_k)) + 2\tau M.$$

Using convexity and concavity of the function $\hat{\varphi}$:

$$\begin{aligned}
 \varepsilon_{sad}(\bar{z}_N) &\leq \max_{(x', y') \in \mathcal{Z}} \frac{1}{N+1} \sum_{k=1}^N (\hat{\varphi}(x_k, y') - \hat{\varphi}(x', y_k)) + 2\tau M \\
 &= \max_{(x', y') \in \mathcal{Z}} \frac{1}{N+1} \sum_{k=1}^N (\hat{\varphi}(x_k, y') - \hat{\varphi}(x_k, y_k) + \hat{\varphi}(x_k, y_k) - \hat{\varphi}(x', y_k)) \\
 &\quad + 2\tau M \\
 &\leq \max_{(x', y') \in \mathcal{Z}} \frac{1}{N+1} \sum_{k=1}^N (\langle \nabla_y \hat{\varphi}(x_k, y_k), y' - y_k \rangle + \langle \nabla_x \hat{\varphi}(x_k, y_k), x_k - x' \rangle) \\
 &\quad + 2\tau M \\
 &\leq \max_{u \in \mathcal{Z}} \frac{1}{N+1} \sum_{k=0}^N \langle \tilde{\nabla} \hat{\varphi}(z_k), z_k - u \rangle + 2\tau M. \tag{35}
 \end{aligned}$$

Step 3. Combining expressions (34), (35), (6) and taking full mathematical expectation, we get

$$\mathbb{E} [\varepsilon_{sad}(\bar{z}_N)] \leq \frac{\Omega^2}{2\gamma(N+1)} + \frac{\gamma M_{all}^2}{2} + \frac{1}{\gamma(N+1)} \mathbb{E} \left[\max_{u \in \mathcal{Z}} D(u) \right] + 2\tau M. \tag{36}$$

Let's estimate $D(u)$. For this we prove the following lemma:

Lemma 14 (see Lemma 5.3.2 from [3]).

$$\mathbb{E} \left[\max_{u \in \mathcal{Z}} D(u) \right] \leq \Omega^2 + \frac{\gamma(N+1)\Delta\Omega n a_q}{\tau} + \gamma^2 M_{all}^2 (N+1), \tag{37}$$

where $M_{all}^2 \stackrel{\text{def}}{=} 3 \left(cnM^2 + \frac{n^2(\sigma^2 + \Delta^2)}{\tau^2} \right) a_q^2$ is from Lemma 1.

Proof. Let define sequence v : $v_1 \stackrel{\text{def}}{=} z_1$, $v_{k+1} \stackrel{\text{def}}{=} \text{prox}_{v_k}(-\rho\gamma\Delta_k)$ for some $\rho > 0$:

$$\begin{aligned}
 D(u) &= \gamma \sum_{k=0}^N \langle -\Delta_k, z_k - u \rangle \\
 &= \gamma \sum_{k=0}^N \langle -\Delta_k, z_k - v_k \rangle + \gamma \sum_{k=0}^N \langle -\Delta_k, v_k - u \rangle. \tag{38}
 \end{aligned}$$

By the definition of v and an optimal condition for the prox-operator, we have for all $u \in \mathcal{Z}$

$$\langle -\gamma\rho\Delta_k - \nabla d(v_{k+1}) + \nabla d(v_{k+1}), u - v_{k+1} \rangle \geq 0.$$

Rewriting this inequality, we get

$$\langle -\gamma\rho\Delta_k, v_k - u \rangle \leq \langle -\gamma\rho\Delta_k, v_k - v_{k+1} \rangle + \langle \nabla d(v_{k+1}) - \nabla d(v_k), u - v_{k+1} \rangle.$$

Using (28):

$$\langle -\gamma\rho\Delta_k, v_k - u \rangle \leq \langle -\gamma\rho\Delta_k, v_k - v_{k+1} \rangle + V_{v_k}(u) - V_{v_{k+1}}(u) - V_{v_k}(v_{k+1}).$$

Bearing in mind the Bregman divergence property $2V_x(y) \geq \|x - y\|_p^2$:

$$\langle -\gamma\rho\Delta_k, v_k - u \rangle \leq \langle -\gamma\rho\Delta_k, v_k - v_{k+1} \rangle + V_{v_k}(u) - V_{v_{k+1}}(u) - \frac{1}{2}\|v_{k+1} - v_k\|_p^2.$$

Using the definition of the conjugate norm:

$$\begin{aligned} \langle -\gamma\rho\Delta_k, v_k - u \rangle &\leq \|\gamma\rho\Delta_k\|_q \cdot \|v_k - v_{k+1}\|_p + V_{v_k}(u) - V_{v_{k+1}}(u) - \frac{1}{2}\|v_{k+1} - v_k\|_p^2 \\ &\leq \frac{\rho^2\gamma^2}{2}\|\Delta_k\|_q^2 + V_{v_k}(u) - V_{v_{k+1}}(u). \end{aligned}$$

Summing over k from 0 to N :

$$\sum_{k=0}^N \gamma\rho\langle -\Delta_k, v_k - u \rangle \leq V_{v_1}(u) - V_{v_{N+1}}(u) + \frac{\rho^2\gamma^2}{2} \sum_{k=0}^N \|\Delta_k\|_q^2.$$

Notice that $V_x(y) \geq 0$ and $V_{v_1}(u) \leq \Omega^2/2$:

$$\sum_{k=0}^N \gamma\langle -\Delta_k, v_k - u \rangle \leq \frac{\Omega^2}{2\rho} + \frac{\rho\gamma^2}{2} \sum_{k=0}^N \|\Delta_k\|_q^2. \quad (39)$$

Substituting (39) into (38):

$$D(u) \leq \sum_{k=0}^N \gamma\langle \Delta_k, v_k - z_k \rangle + \frac{\Omega^2}{2\rho} + \frac{\rho\gamma^2}{2} \sum_{k=0}^N \|\Delta_k\|_q^2.$$

The right side is independent of u , then

$$\max_{u \in \mathcal{Z}} D(u) \leq \sum_{k=0}^N \gamma\langle \Delta_k, v_k - z_k \rangle + \frac{\Omega^2}{2\rho} + \frac{\rho\gamma^2}{2} \sum_{k=0}^N \|\Delta_k\|_q^2. \quad (40)$$

Taking the full expectation:

$$\mathbb{E} \left[\max_{u \in \mathcal{Z}} D(u) \right] \leq \mathbb{E} \left[\sum_{k=1}^N \gamma\langle \Delta_k, v_k - z_k \rangle \right] + \frac{\Omega^2}{2\rho} + \frac{\rho\gamma^2}{2} \mathbb{E} \left[\sum_{k=0}^N \|\Delta_k\|_q^2 \right].$$

Using the independence of $\mathbf{e}_1, \dots, \mathbf{e}_N, \xi_1^\pm, \dots, \xi_N^\pm$, we have

$$\mathbb{E} \left[\max_{u \in \mathcal{Z}} D(u) \right] \leq \mathbb{E} \left[\sum_{k=0}^N \gamma \mathbb{E}_{\mathbf{e}_k, \xi_k} [\langle \Delta_k, v_k - z_k \rangle] \right] + \frac{\Omega^2}{2\rho} + \frac{\rho\gamma^2}{2} \mathbb{E} \left[\sum_{k=0}^N \|\Delta_k\|_q^2 \right].$$

Note that $v_k - z_k$ does not depend on \mathbf{e}_k, ξ_k . Then

$$\mathbb{E} \left[\max_{u \in \mathcal{Z}} D(u) \right] \leq \mathbb{E} \left[\sum_{k=0}^N \gamma \langle \mathbb{E}_{\mathbf{e}_k, \xi_k} [\Delta_k], v_k - z_k \rangle \right] + \frac{\Omega^2}{2\rho} + \frac{\rho\gamma^2}{2} \mathbb{E} \left[\sum_{k=0}^N \|\Delta_k\|_q^2 \right].$$

By (11) and definition of diameter Ω we get

$$\mathbb{E} \left[\max_{u \in \mathcal{Z}} D(u) \right] \leq \frac{\Delta\Omega na_q}{\tau} \sum_{k=0}^N \gamma + \frac{\Omega^2}{2\rho} + \frac{\rho\gamma^2}{2} \sum_{k=0}^N \mathbb{E} [\|\Delta_k\|_q^2].$$

To prove the lemma, it remains to estimate $\mathbb{E} [\|\Delta_k\|_q^2]$:

$$\begin{aligned} \mathbb{E} [\|\Delta_k\|_q^2] &\leq \mathbb{E} \left[\|g(z_k, \mathbf{e}_k, \tau, \xi_k^\pm, \delta_k^\pm) - \tilde{\nabla} \hat{\varphi}(z_k)\|_q^2 \right] \\ &\leq 2\mathbb{E} [\|g(z_k, \mathbf{e}_k, \tau, \xi_k^\pm, \delta_k^\pm)\|_q^2] + 2\mathbb{E} [\|\tilde{\nabla} \hat{\varphi}(z_k)\|_q^2] \\ &\leq 2\mathbb{E} [\|g(z_k, \mathbf{e}_k, \tau, \xi_k^\pm, \delta_k^\pm)\|_q^2] + 2\mathbb{E} \left[\left\| \frac{n(\varphi(z + \tau\mathbf{e}) - \varphi(z - \tau\mathbf{e}))}{2\tau} \mathbf{e} \right\|_q^2 \right]. \end{aligned}$$

Using Lemma 1, we have $\mathbb{E} [\|\Delta_k\|_q^2] \leq 4M_{all}^2$, whence

$$\mathbb{E} \left[\max_{u \in \mathcal{Z}} D(u) \right] \leq \frac{\Omega^2}{2\rho} + \frac{\gamma(N+1)\Delta\Omega na_q}{\tau} + 2\rho\gamma^2 M_{all}^2(N+1).$$

Taking $\rho = 1/2$ ends the proof of lemma. □

(36) with this lemma gives

$$\mathbb{E} [\varepsilon_{sad}(\bar{z}_N)] \leq \frac{3\Omega^2}{2\gamma(N+1)} + \frac{3\gamma M_{all}^2}{2} + \frac{\Delta\Omega na_q}{\tau} + 2\tau M.$$

This completes the proof of the theorem. □

Theorem 2. Let problem (1) with function $\varphi(x, y)$ be solved using Algorithm 1 with $V_z(w) = \frac{1}{2}\|z-w\|_2^2$ and the oracle (5). Assume, that the set \mathcal{Z} , the function $\varphi(x, y)$ and its inexact modification $\tilde{\varphi}(x, y)$ satisfy Assumptions 1, 2, 3, 4. Denote by N the number of iterations and $\gamma_k = \frac{1}{\mu k}$. Then the rate of convergence is given by the following expression:

$$\mathbb{E} [\varphi(\bar{x}_N, \bar{y}^*) - \varphi(x^*, \bar{y}_N)] \leq \frac{M_{all}^2 \log(N+1)}{2\mu(N+1)} + \frac{\Delta n \Omega}{\tau} + 2\tau M$$

Proof. We start this proof from substituting definition of g_k and $u = z^*$ in (32):

$$2\gamma_k \langle g(z_k, \mathbf{e}_k, \tau, \xi_k^\pm), z_k - z^* \rangle \leq \|z_k - z^*\|^2 - \|z_{k+1} - z^*\|^2 + \gamma_k^2 \|g(z_k, \mathbf{e}_k, \tau, \xi_k^\pm)\|^2.$$

With small rearrangement

$$\begin{aligned} 2\gamma_k \langle \tilde{\nabla} \hat{\varphi}(z_k), z_k - z^* \rangle &\leq \|z_k - z^*\|^2 - \|z_{k+1} - z^*\|^2 + \gamma_k^2 \|g(z_k, \mathbf{e}_k, \tau, \xi_k^\pm)\|^2 \\ &\quad + 2\gamma_k \langle \tilde{\nabla} \hat{\varphi}(z_k) - g(z_k, \mathbf{e}_k, \tau, \xi_k^\pm), z_k - z^* \rangle. \end{aligned}$$

On the other hand with (9) and Lemma 3 we get

$$\begin{aligned} \varphi(x_k, y^*) - \varphi(x^*, y_k) &= \hat{\varphi}(x_k, y^*) + |\varphi(x_k, y^*) - \hat{\varphi}(x_k, y^*)| \\ &\quad - \hat{\varphi}(x^*, y_k) + |\varphi(x^*, y_k) - \hat{\varphi}(x^*, y_k)| \\ &\leq \hat{\varphi}(x_k, y^*) - \hat{\varphi}(x^*, y_k) + 2\tau M \\ &\leq \hat{\varphi}(x_k, y^*) - \hat{\varphi}(x_k, y_k) + \hat{\varphi}(x_k, y_k) - \hat{\varphi}(x^*, y_k) + 2\tau M \\ &\leq \langle -\nabla_y \hat{\varphi}(x_k, y_k), y_k - y^* \rangle - \frac{\mu}{2} \|y_k - y^*\|^2 \\ &\quad + \langle -\nabla_x \hat{\varphi}(x_k, y_k), x_k - x^* \rangle - \frac{\mu}{2} \|x_k - x^*\|^2 + 2\tau M \\ &= \langle \tilde{\nabla} \hat{\varphi}(z_k), z_k - z^* \rangle - \frac{\mu}{2} \|z_k - z^*\|^2 + 2\tau M. \end{aligned}$$

By connecting we have

$$\begin{aligned} 2\gamma_k (\varphi(x_k, y^*) - \varphi(x^*, y_k)) &\leq (1 - \mu\gamma_k) \|z_k - z^*\|^2 - \|z_{k+1} - z^*\|^2 + \gamma_k^2 \|g(z_k, \mathbf{e}_k, \tau, \xi_k^\pm)\|^2 \\ &\quad + 2\gamma_k \langle \tilde{\nabla} \hat{\varphi}(z_k) - g(z_k, \mathbf{e}_k, \tau, \xi_k^\pm), z_k - z^* \rangle + 4\gamma_k \tau M. \end{aligned}$$

Taking the total expectation and taking into account that $z_k - z^*$ does not depend on \mathbf{e}_k, ξ_k :

$$\begin{aligned} \mathbb{E}[\varphi(x_k, y^*) - \varphi(x^*, y_k)] &\leq \left(\frac{1}{2\gamma_k} - \frac{\mu}{2} \right) \mathbb{E} \|z_k - z^*\|^2 \\ &\quad - \frac{1}{2\gamma_k} \mathbb{E} \|z_{k+1} - z^*\|^2 + \frac{\gamma_k}{2} \mathbb{E} \|g(z_k, \mathbf{e}_k, \tau, \xi_k^\pm)\|^2 \\ &\quad + \mathbb{E} \langle \mathbb{E}_{\mathbf{e}_k, \xi_k} [\tilde{\nabla} \hat{\varphi}(z_k) - g(z_k, \mathbf{e}_k, \tau, \xi_k^\pm)], z_k - z^* \rangle + 2\tau M. \end{aligned}$$

With (6), (11) with $a_q = 1$ (Euclidean case) we get

$$\begin{aligned} \mathbb{E}[\varphi(x_k, y^*) - \varphi(x^*, y_k)] &\leq \left(\frac{1}{2\gamma_k} - \frac{\mu}{2} \right) \mathbb{E} \|z_k - z^*\|^2 - \frac{1}{2\gamma_k} \mathbb{E} \|z_{k+1} - z^*\|^2 \\ &\quad + \frac{\gamma_k M_{all}^2}{2} + \frac{\Delta n \Omega}{\tau} + 2\tau M. \end{aligned}$$

Summing over all k from 0 to N , we have

$$\begin{aligned} \mathbb{E} \left[\sum_{k=0}^N \varphi(x_k, y^*) - \sum_{k=0}^N \varphi(x^*, y_k) \right] &\leq \sum_{k=1}^{N-1} \left(\frac{1}{2\gamma_k} - \frac{1}{2\gamma_{k-1}} - \frac{\mu}{2} \right) \mathbb{E} \|z_k - z^*\|^2 \\ &\quad + \left(\frac{1}{2\gamma_0} - \frac{\mu}{2} \right) \|z_0 - z^*\|^2 + \frac{M_{all}^2}{2} \sum_{k=0}^N \gamma_k \\ &\quad + \frac{\Delta n \Omega (N+1)}{\tau} + 2\tau M (N+1). \end{aligned}$$

With $\gamma_k = \frac{1}{\mu(k+1)}$ we get

$$\mathbb{E} \left[\sum_{k=0}^N \varphi(x_k, y^*) - \sum_{k=0}^N \varphi(x^*, y_k) \right] \leq \frac{M_{all}^2 \log(N+1)}{2\mu} + \frac{\Delta n \Omega (N+1)}{\tau} + 2\tau M (N+1).$$

It remains only to apply Jensen's inequality to the left-hand side:

$$\mathbb{E} [\varphi(\bar{x}_N, y^*) - \varphi(x^*, \bar{y}_N)] \leq \frac{M_{all}^2 \log(N+1)}{2\mu(N+1)} + \frac{\Delta n \Omega}{\tau} + 2\tau M.$$

□

Lemma 15. For $\tilde{g}_k \stackrel{def}{=} \tilde{g}(z_k, z_{k-1}, \mathbf{e}_k, \mathbf{e}_{k-1}, \xi_k, \xi_{k-1})$ defined in (14) under Assumptions 2 and 4 the following inequalities holds:

$$\mathbb{E} [\|\tilde{g}_k\|_2^2] \leq \alpha^k \mathbb{E} [\|\tilde{g}_0\|_2^2] + \left(\frac{12n^2(\sigma^2 + \Delta^2)}{\tau^2} + 12n^2M^2 \right) \frac{1}{1-\alpha},$$

where $\alpha = \frac{6\gamma^2 n^2 M^2}{\tau^2} < 1$.

Proof.

$$\begin{aligned} &\mathbb{E} [\|\tilde{g}(z_k, z_{k-1}, \mathbf{e}_k, \mathbf{e}_{k-1}, \xi_k, \xi_{k-1})\|_2^2] \\ &= \frac{n^2}{\tau^2} \mathbb{E} \left[(\tilde{\varphi}(z_k + \tau \mathbf{e}_k, \xi_k) - \tilde{\varphi}(z_{k-1} + \tau \mathbf{e}_{k-1}, \xi_{k-1}))^2 \right] \\ &= \frac{n^2}{\tau^2} \mathbb{E} \left[(\varphi(z_k + \tau \mathbf{e}_k) + \xi_k + \delta(z_k + \tau \mathbf{e}_k) - \varphi(z_{k-1} + \tau \mathbf{e}_{k-1}) - \xi_{k-1} - \delta(z_{k-1} + \tau \mathbf{e}_{k-1}))^2 \right]. \end{aligned}$$

With a simple fact (30), we get

$$\begin{aligned} \mathbb{E} [\|\tilde{g}_k\|_2^2] &\leq \frac{6n^2}{\tau^2} \mathbb{E} [\xi_k^2 + \delta^2(z_k + \tau \mathbf{e}_k) + \xi_{k-1}^2 + \delta^2(z_{k-1} + \tau \mathbf{e}_{k-1})] \\ &\quad + \frac{6n^2}{\tau^2} \mathbb{E} [(\varphi(z_k + \tau \mathbf{e}_k) - \varphi(z_{k-1} + \tau \mathbf{e}_{k-1}))^2] \\ &\quad + \frac{6n^2}{\tau^2} \mathbb{E} [(\varphi(z_{k-1} + \tau \mathbf{e}_{k-1}) - \varphi(z_{k-1} + \tau \mathbf{e}_k))^2]. \end{aligned}$$

Next we use (2) and (4) and have

$$\begin{aligned}\mathbb{E} [\|\tilde{g}_k\|_2^2] &\leq \frac{12n^2(\sigma^2 + \Delta^2)}{\tau^2} + \frac{6n^2M^2}{\tau^2} \mathbb{E} \|z_k - z_{k-1}\|_2^2 + 6n^2M^2 \mathbb{E} [\|\mathbf{e}_{k-1} - \mathbf{e}_k\|_2^2] \\ &\leq \frac{12n^2(\sigma^2 + \Delta^2)}{\tau^2} + \frac{6n^2M^2}{\tau^2} \mathbb{E} \|z_k - z_{k-1}\|_2^2 + 12n^2M^2.\end{aligned}$$

Considering the step of Algorithm 1 we can rewrite as follows:

$$\mathbb{E} [\|\tilde{g}_k\|_2^2] \leq \frac{12n^2(\sigma^2 + \Delta^2)}{\tau^2} + \frac{6\gamma^2n^2M^2}{\tau^2} \mathbb{E} \|\tilde{g}_{k-1}\|_2^2 + 12n^2M^2.$$

Then we run recursion

$$\mathbb{E} [\|\tilde{g}_k\|_2^2] \leq \left(\frac{6\gamma^2n^2M^2}{\tau^2}\right)^k \mathbb{E} \|\tilde{g}_0\|_2^2 + \left(\frac{12n^2(\sigma^2 + \Delta^2)}{\tau^2} + 12n^2M^2\right) \sum_{i=0}^{k-1} \left(\frac{6\gamma^2n^2M^2}{\tau^2}\right)^i.$$

With $\alpha = \frac{6\gamma^2n^2M^2}{\tau^2} < 1$

$$\mathbb{E} [\|\tilde{g}_k\|_2^2] \leq \alpha^k \mathbb{E} \|\tilde{g}_0\|_2^2 + \left(\frac{12n^2(\sigma^2 + \Delta^2)}{\tau^2} + 12n^2M^2\right) \frac{1}{1 - \alpha}.$$

□

Theorem 3. Let problem (1) with function $\varphi(x, y)$ be solved using Algorithm 1 with $V_z(w) = \frac{1}{2}\|z - w\|_2^2$ and the oracle (14). Assume, that the set \mathcal{Z} , the convex-concave function $\varphi(x, y)$ and its inexact modification $\tilde{\varphi}(x, y)$ satisfy Assumptions 1, 2, 4. Denote by N the number of iterations and $\gamma_k = \gamma = \text{const}$. Then the rate of convergence is given by the following expression:

$$\begin{aligned}\mathbb{E} [\varepsilon_{sad}(\bar{z}_N)] &\leq \frac{3\Omega^2}{2\gamma(N+1)} + \frac{3\gamma}{2(N+1)(1-\alpha)} \mathbb{E} [\|\tilde{g}_0\|_2^2] \\ &\quad + \frac{3\gamma}{2(1-\alpha)} \left(\frac{12n^2(\sigma^2 + \Delta^2)}{\tau^2} + 12n^2M^2\right) + 2\tau M + \frac{\Delta\Omega n}{\tau}.\end{aligned}$$

Ω is a diameter of \mathcal{Z} , $\alpha = \frac{6\gamma^2n^2M^2}{\tau^2} < 1$.

Proof. We begin our proof right away by obtaining the inequality similarly to (36) but by (15), not (6)

$$\begin{aligned}
 \mathbb{E} [\varepsilon_{sad}(\bar{z}_N)] &\leq \frac{\Omega^2}{2\gamma(N+1)} + \frac{\gamma}{2(N+1)} \mathbb{E} [\|\tilde{g}_0\|_2^2] \sum_{k=0}^N \alpha^k \\
 &\quad + \frac{\gamma}{2(1-\alpha)} \left(\frac{12n^2(\sigma^2 + \Delta^2)}{\tau^2} + 12n^2 M^2 \right) \\
 &\quad + \frac{1}{\gamma(N+1)} \mathbb{E} \left[\max_{u \in \mathcal{Z}} \tilde{D}(u) \right] + 2\tau M \\
 &\leq \frac{\Omega^2}{2\gamma(N+1)} + \frac{\gamma}{2(N+1)(1-\alpha)} \mathbb{E} [\|\tilde{g}_0\|_2^2] \\
 &\quad + \frac{\gamma}{2(1-\alpha)} \left(\frac{12n^2(\sigma^2 + \Delta^2)}{\tau^2} + 12n^2 M^2 \right) \\
 &\quad + \frac{1}{\gamma(N+1)} \mathbb{E} \left[\max_{u \in \mathcal{Z}} \tilde{D}(u) \right] + 2\tau M, \tag{41}
 \end{aligned}$$

where $\tilde{D}(u) \stackrel{\text{def}}{=} \sum_{k=0}^N \gamma \langle \tilde{\Delta}_k, u - z_k \rangle$ with $\tilde{\Delta}_k \stackrel{\text{def}}{=} \tilde{g}_k - \tilde{\nabla} \tilde{\varphi}(z_k)$. Let estimate $\tilde{D}(u)$. For this we prove the following lemma:

Lemma 16.

$$\begin{aligned}
 \mathbb{E} \left[\max_{u \in \mathcal{Z}} \tilde{D}(u) \right] &\leq \frac{\gamma(N+1)\Delta\Omega n}{\tau} + \Omega^2 + \frac{\gamma^2}{1-\alpha} \|\tilde{g}_0\|_2^2 \\
 &\quad + \frac{\gamma^2(N+1)}{1-\alpha} \left(\frac{12n^2(\sigma^2 + \Delta^2)}{\tau^2} + 12n^2 M^2 \right). \tag{42}
 \end{aligned}$$

Proof. Let's start with (40). All other steps are done in the same way.

$$\max_{u \in \mathcal{Z}} \tilde{D}(u) \leq \sum_{k=0}^N \gamma \langle \tilde{\Delta}_k, v_k - z_k \rangle + \frac{\Omega^2}{2\rho} + \frac{\rho\gamma^2}{2} \sum_{k=0}^N \|\tilde{\Delta}_k\|_2^2.$$

Taking the full expectation:

$$\mathbb{E} \left[\max_{u \in \mathcal{Z}} \tilde{D}(u) \right] \leq \mathbb{E} \left[\sum_{k=1}^N \gamma \langle \tilde{\Delta}_k, v_k - z_k \rangle \right] + \frac{\Omega^2}{2\rho} + \frac{\rho\gamma^2}{2} \mathbb{E} \left[\sum_{k=0}^N \|\tilde{\Delta}_k\|_2^2 \right].$$

Using the independence of $\mathbf{e}_1, \dots, \mathbf{e}_N, \xi_1^\pm, \dots, \xi_N^\pm$, we have

$$\begin{aligned}
 \mathbb{E} \left[\max_{u \in \mathcal{Z}} D(u) \right] &\leq \mathbb{E} \left[\sum_{k=0}^N \gamma \mathbb{E}_{\xi_k} \left[\langle \tilde{\Delta}_k, v_k - z_k \rangle \right] \right] + \mathbb{E} \left[\sum_{k=0}^N \gamma \mathbb{E}_{\mathbf{e}_k} \left[\langle \Delta_k - \tilde{\Delta}_k, v_k - z_k \rangle \right] \right] \\
 &\quad + \frac{\Omega^2}{2\rho} + \frac{\rho\gamma^2}{2} \mathbb{E} \left[\sum_{k=0}^N \|\Delta_k\|_2^2 \right].
 \end{aligned}$$

Note that $v_k - z_k$ does not depend on \mathbf{e}_k , ξ_k and $\mathbb{E}_{\xi_k} \tilde{\Delta}_k = 0$. Then

$$\begin{aligned} \mathbb{E} \left[\max_{u \in \mathcal{Z}} \tilde{D}(u) \right] &\leq \mathbb{E} \left[\sum_{k=0}^N \gamma \langle \mathbb{E}_{\mathbf{e}_k} [\tilde{\Delta}_k], v_k - z_k \rangle \right] + \frac{\Omega^2}{2\rho} + \frac{\rho\gamma^2}{2} \mathbb{E} \left[\sum_{k=0}^N \|\tilde{\Delta}_k\|_2^2 \right] \\ &= \mathbb{E} \left[\sum_{k=0}^N \gamma \langle \mathbb{E}_{\mathbf{e}_k} [\tilde{\Delta}_k], v_k - z_k \rangle \right] + \frac{\Omega^2}{2\rho} + \frac{\rho\gamma^2}{2} \mathbb{E} \left[\sum_{k=0}^N \|\tilde{\Delta}_k\|_2^2 \right]. \end{aligned}$$

By (11) and definition of diameter Ω we get

$$\mathbb{E} \left[\max_{u \in \mathcal{Z}} \tilde{D}(u) \right] \leq \frac{\Delta\Omega n}{\tau} \sum_{k=0}^N \gamma + \frac{\Omega^2}{2\rho} + \frac{\rho\gamma^2}{2} \sum_{k=0}^N \mathbb{E} \left[\|\tilde{\Delta}_k\|_2^2 \right].$$

To prove the lemma, it remains to estimate $\mathbb{E} \left[\|\tilde{\Delta}_k\|_2^2 \right]$:

$$\mathbb{E} \left[\|\tilde{\Delta}_k\|_2^2 \right] \leq 2\mathbb{E} \left[\|\check{g}_k\|_2^2 \right] + 2\mathbb{E} \left[\left\| \frac{n(\varphi(z + \tau\mathbf{e}) - \varphi(z - \tau\mathbf{e}))}{2\tau} \mathbf{e} \right\|_2^2 \right].$$

Using Lemma 4, we have

$$\begin{aligned} \mathbb{E} \left[\max_{u \in \mathcal{Z}} \tilde{D}(u) \right] &\leq \frac{\Delta\Omega n}{\tau} \sum_{k=0}^N \gamma + \frac{\Omega^2}{2\rho} + 2\rho\gamma^2 \sum_{k=0}^N \alpha^k \mathbb{E} \left[\|\check{g}_0\|_2^2 \right] \\ &\quad + 2\rho\gamma^2(N+1) \left(\frac{12n^2(\sigma^2 + \Delta^2)}{\tau^2} + 12n^2M^2 \right) \frac{1}{1-\alpha}. \end{aligned}$$

Taking $\rho = 1/2$ ends the proof of lemma. □

(41) with this lemma gives

$$\begin{aligned} \mathbb{E} [\varepsilon_{sad}(\bar{z}_N)] &\leq \frac{3\Omega^2}{2\gamma(N+1)} + \frac{3\gamma}{2(N+1)(1-\alpha)} \mathbb{E} \left[\|\check{g}_0\|_2^2 \right] \\ &\quad + \frac{3\gamma}{2(1-\alpha)} \left(\frac{12n^2(\sigma^2 + \Delta^2)}{\tau^2} + 12n^2M^2 \right) + 2\tau M + \frac{\Delta\Omega n}{\tau}. \end{aligned}$$

This completes the proof of the theorem. □

C Proofs for Section 3.2

The proofs of the Theorems 4 and 5 copy the proofs of the Theorems 1 and 2 except for the usage Lemma 6 instead of Lemma 2. So the term $2M\tau$ in Theorems 1 and 2 is replaced by the term $L\tau^2$ in Theorems 4 and 5.

D Proofs for Section 3.3

Theorem 6. Let $\varphi \in \mathcal{F}_{\mu,\beta}(L)$ with $\mu, L > 0$ and $\beta > 2$. Let Assumption 6 hold and let \mathcal{Z} be a convex compact subset of \mathbb{R}^n . Let φ be M -Lipschitz on the Euclidean τ_1 -neighborhood of \mathcal{Z} (see τ_k below).

Then the rate of convergence is given by Algorithm 2 with parameters

$$\tau_k = \left(\frac{3\kappa\sigma^2 n}{2(\beta-1)(\kappa_\beta L)^2} \right)^{\frac{1}{2\beta}} k^{-\frac{1}{2\beta}}, \quad \alpha_k = \frac{2}{\mu k}, \quad k = 1, \dots, N$$

satisfies

$$\begin{aligned} \mathbb{E}[\varphi(\bar{x}_N, y^*) - \varphi(x^*, \bar{y}_N)] &\leq \max_{y \in \mathcal{Y}} \mathbb{E}[\varphi(\bar{x}_N, y)] - \min_{x \in \mathcal{X}} \mathbb{E}[\varphi(x, \bar{y}_N)] \\ &\leq \frac{1}{\mu} \left(n^{2-\frac{1}{\beta}} \frac{A_1}{N^{\frac{\beta-1}{\beta}}} + A_2 \frac{n(1+\ln N)}{N} \right), \end{aligned}$$

where $\bar{z}_N = \frac{1}{N} \sum_{k=1}^N z_k$, $A_1 = 3\beta(\kappa\sigma^2)^{\frac{\beta-1}{\beta}} (\kappa_\beta L)^{\frac{2}{\beta}}$, $A_2 = 9\kappa G^2$, κ_β and κ are constants depending only on β , see (22) and (23).

Proof. Step 1. Fix an arbitrary $z \in \mathcal{Z}$. As z_{k+1} is the Euclidean projection we have $\|z_{k+1} - z\|^2 \leq \|z_k - \gamma_k \tilde{g}_k - z\|^2$ which is equivalent to

$$\langle \tilde{g}_k, z_k - z \rangle \leq \frac{\|z_k - z\|^2 - \|z_{k+1} - z\|^2}{2\gamma_k} + \frac{\gamma_k}{2} \|\tilde{g}_k\|^2. \quad (43)$$

Using the strong convexity-concavity and combining x and y parts of the argument z together we have

$$\begin{aligned} \varphi(x_k, y) - \varphi(x, y_k) &= \varphi(x_k, y) - \varphi(x_k, y_k) + \varphi(x_k, y_k) - \varphi(x, y_k) \\ &\leq \langle -\nabla_y \varphi(x_k, y_k), y_k - y \rangle - \frac{\mu}{2} \|y_k - y\|^2 \\ &\quad + \langle -\nabla_x \varphi(x_k, y_k), x_k - x \rangle - \frac{\mu}{2} \|x_k - x\|^2 \\ &= \langle \tilde{\nabla} \varphi(z_k), z_k - z \rangle - \frac{\mu}{2} \|z_k - z\|^2. \end{aligned} \quad (44)$$

Combining the last two inequations we obtain

$$\begin{aligned} \varphi(x_k, y) - \varphi(x, y_k) &\leq \langle \tilde{\nabla} \varphi(z_k) - \tilde{g}_k, z_k - z \rangle + \frac{\|z_k - z\|^2 - \|z_{k+1} - z\|^2}{2\gamma_k} \\ &\quad + \frac{\gamma_k}{2} \|\tilde{g}_k\|^2 - \frac{\mu}{2} \|z_k - z\|^2. \end{aligned} \quad (45)$$

Taking conditional expectation given z_k with respect to r_k , ξ_k^+ and ξ_k^- we obtain

$$\begin{aligned} \varphi(x_k, y) - \varphi(x, y_k) &\leq \langle \tilde{\nabla} \varphi(z_k) - \mathbb{E}[\tilde{g}_k | z_k], z_k - z \rangle + \frac{\gamma_k}{2} \mathbb{E}[\|\tilde{g}_k\|^2 | z_k] \\ &\quad + \frac{\|z_k - z\|^2 - \mathbb{E}[\|z_{k+1} - z\|^2 | z_k]}{2\gamma_k} - \frac{\mu}{2} \|z_k - z\|^2. \end{aligned} \quad (46)$$

Step 2 (Bounding bias term). Our aim is to bound the first term in (46), namely $\langle \tilde{\nabla} \varphi(z_k) - \mathbb{E}[\tilde{g}_k | z_k], z_k - z \rangle$. Using the Taylor expansion we have

$$\begin{aligned} \varphi(z_k + \tau_k r_k e_k) &= \varphi(z_k) + \langle \nabla \varphi(z_k), \tau_k r_k e_k \rangle \\ &+ \sum_{2 \leq |m| \leq l} \frac{(\tau_k r_k)^{|m|}}{m!} D^{(m)} \varphi(z_k) e_k^m + R(\tau_k r_k e_k), \end{aligned} \quad (47)$$

where by assumption $|R(\tau_k r_k e_k)| \leq L \|\tau_k r_k e_k\|^\beta = L(\tau_k \cdot |r_k|)^\beta$. Thus,

$$\begin{aligned} \tilde{g}_k &= \langle \nabla \varphi(x_k), \tau_k r_k e_k \rangle + \sum_{2 \leq |m| \leq l, |m| \text{ odd}} \frac{(\tau_k r_k)^{|m|}}{m!} D^{(m)} \varphi(z_k) e_k^m \\ &+ \frac{1}{2} R(\tau_k r_k e_k) - \frac{1}{2} R(-\tau_k r_k e_k) + \xi_k^+ - \xi_k^- \frac{n}{\tau_k} K(r_k) \begin{pmatrix} (\mathbf{e}_k)_x \\ -(\mathbf{e}_k)_y \end{pmatrix}. \end{aligned} \quad (48)$$

Using the properties of the smoothing kernel K , independence of \mathbf{e}_k and r_k (Assumption 6) and the fact that $\mathbb{E}[e_k e_k^T] = \frac{1}{n} \mathbb{I}_{n \times n}$ we obtain

$$\mathbb{E}_{e_k, r_k} \left[\langle \nabla \varphi(z_k), \tau_k r_k e_k \rangle \frac{n}{\tau_k} K(r_k) \begin{pmatrix} (\mathbf{e}_k)_x \\ -(\mathbf{e}_k)_y \end{pmatrix} \middle| z_k \right] = \tilde{\nabla} \varphi(z_k). \quad (49)$$

Using the fact that $\mathbb{E}[r_k^{|m|} K(r_k)] = 0$ if $2 \leq |m| \leq l$ or $|m| = 0$ and Assumption 6 we have

$$\mathbb{E} \left[\left(\sum_{2 \leq |m| \leq l, |m| \text{ odd}} \frac{(\tau_k r_k)^{|m|}}{m!} D^{(m)} \varphi(z_k) e_k^m + \xi_k^+ - \xi_k^- \right) \frac{n}{\tau_k} K(r_k) \begin{pmatrix} (\mathbf{e}_k)_x \\ -(\mathbf{e}_k)_y \end{pmatrix} \middle| x_k \right] = 0. \quad (50)$$

Substituting (48), (49) and (50) in the first term in (46) and using the definition of κ_β (see (22)) we obtain

$$\begin{aligned} & \left| \langle \tilde{\nabla} \varphi(z_k) - \mathbb{E}[\tilde{g}_k | z_k], z_k - z \rangle \right| = \\ &= \left| \mathbb{E} \left[\left(\frac{1}{2} R(\tau_k r_k e_k) - \frac{1}{2} R(-\tau_k r_k e_k) \right) \frac{n}{\tau_k} K(r_k) \left\langle \begin{pmatrix} (\mathbf{e}_k)_x \\ -(\mathbf{e}_k)_y \end{pmatrix}, z_k - z \right\rangle \middle| z_k \right] \right| \\ & \leq L \tau_k^{\beta-1} \cdot \mathbb{E}_{r_k} [|r_k|^\beta K(r_k)] \cdot n |\mathbb{E}_{e_k} [\langle \mathbf{e}_k, z_k - z \rangle | z_k]| \\ & \leq \kappa_\beta L \sqrt{n} \tau_k^{\beta-1} \|z_k - z\|, \end{aligned} \quad (51)$$

where in the last two inequalities the symmetry of Euclidean sphere and the fact from concentration measure theory that $|\mathbb{E}_e [\langle e, s \rangle]|^2 \leq \mathbb{E}_e [\langle e, s \rangle^2] = \frac{\|s\|^2}{n}$ were

used . Applying the inequality $ab \leq 1/2(a^2 + b^2)$ to the last expression in (51) we finally get

$$\left| \langle \tilde{\nabla} \varphi(z_k) - \mathbb{E}[\tilde{g}_k | z_k], z_k - z \rangle \right| \leq \frac{(\kappa_\beta L)^2}{\mu} n \tau_k^{2(\beta-1)} + \frac{\mu}{4} \|z_k - z\|^2. \quad (52)$$

Step 3 (Bounding second moment of gradient estimator). Our aim is to estimate $\mathbb{E}[\|\tilde{g}_k\|^2 | z_k]$ which is the second term in (46). The expectation here is with respect to r_k, ξ_k^+ and ξ_k^- . To lighten the presentation and without loss of generality we drop the lower script k in all quantities.

We have

$$\begin{aligned} \|\tilde{g}\|^2 &= \frac{n^2}{4\tau^2} \left\| (\varphi(z + \tau re) - \varphi(z - \tau re) + \xi^+ - \xi^-) K(r) \begin{pmatrix} \mathbf{e}_x \\ -\mathbf{e}_y \end{pmatrix} \right\|^2 \\ &= \frac{n^2}{4\tau^2} ((\varphi(z + \tau re) - \varphi(z - \tau re) + \xi^+ - \xi^-))^2 K^2(r). \end{aligned} \quad (53)$$

Using the inequality $(a + b + c)^2 \leq 3(a^2 + b^2 + c^2)$ and Assumption 6 we get

$$\mathbb{E}[\|\tilde{g}\|^2 | z] \leq \frac{3n^2}{4\tau^2} (\mathbb{E}[(\varphi(z + \tau re) - \varphi(z - \tau re))^2 K^2(r) | z] + 2\kappa\sigma^2). \quad (54)$$

Using the symmetry of Euclidean unit sphere and the inequality $(a + b)^2 \leq 2(a^2 + b^2)$ we obtain

$$\begin{aligned} \mathbb{E}[(\varphi(z + e) - \varphi(z - e))^2 | z] &= \mathbb{E}_e[(\varphi(z + e) - \varphi(z - e))^2] \\ &\leq \mathbb{E}_e[((\varphi(z + e) - \mathbb{E}_e[\varphi(z + e)]) - (\varphi(z - e) - \mathbb{E}_e[\varphi(z - e)]))^2] \\ &\leq 2\mathbb{E}_e[(\varphi(z + e) - \mathbb{E}_e[\varphi(z + e)])^2] + 2\mathbb{E}_e[(\varphi(z - e) - \mathbb{E}_e[\varphi(z - e)])^2] \\ &\leq 2\sqrt{\mathbb{E}_e[(\varphi(z + e) - \mathbb{E}_e[\varphi(z + e)])^4]} + 2\sqrt{\mathbb{E}_e[(\varphi(z - e) - \mathbb{E}_e[\varphi(z - e)])^4]} \\ &\leq \frac{12M^2}{n}, \end{aligned} \quad (55)$$

where in the last inequality Lemma 10 was used, so we have

$$\mathbb{E}[(\varphi(z + \tau re) - \varphi(z - \tau re))^2 | z] \leq \frac{12(\tau r)^2 M^2}{n} \leq \frac{12\tau^2 M^2}{n}. \quad (56)$$

By substituting (56) into (54), using independence of e and r and returning the lower script k we finally get

$$\mathbb{E}[\|\tilde{g}_k\|^2 | z_k] \leq \kappa \left(9nM^2 + \frac{3(n\sigma)^2}{2\tau_k^2} \right). \quad (57)$$

Step 4. Let ρ_k^2 denote full expectation $\mathbb{E}[\|z_k - z\|^2]$. Substituting (52) and (57) into (46), taking full expectation we obtain

$$\begin{aligned} \mathbb{E}[\varphi(x_k, y) - \varphi(x, y_k)] &\leq \frac{(\kappa_\beta L)^2}{\mu} n \tau_k^{2(\beta-1)} + \frac{\gamma_k}{2} \kappa \left(9nM^2 + \frac{3(n\sigma)^2}{2\tau_k^2} \right) \\ &\quad + \frac{\rho_k^2 - \rho_{k+1}^2}{2\gamma_k} - \left(\frac{\mu}{2} - \frac{\mu}{4} \right) \rho_k^2. \end{aligned} \quad (58)$$

Using the convexity-concavity of φ and (58) we have

$$\begin{aligned} \mathbb{E}[\varphi(\bar{x}_N, y) - \varphi(x, \bar{y}_N)] &\leq \frac{1}{N} \sum_{k=1}^N \varphi(x_k, y) - \frac{1}{N} \sum_{k=1}^N \varphi(x, y_k) \\ &\leq \frac{1}{N} \sum_{k=1}^N \left(\frac{(\kappa_\beta L)^2}{\mu} n \tau_k^{2(\beta-1)} + \frac{\gamma_k}{2} \kappa \left(9nM^2 + \frac{3(n\sigma)^2}{2\tau_k^2} \right) \right) \\ &\quad + \frac{1}{N} \sum_{k=1}^N \left(\frac{\rho_k^2 - \rho_{k+1}^2}{2\gamma_k} - \frac{\mu}{4} \rho_k^2 \right). \end{aligned} \quad (59)$$

Let $\rho_{N+1}^2 = 0$. Then setting $\gamma_k = \frac{2}{\gamma k}$ yields

$$\begin{aligned} \sum_{k=1}^N \left(\frac{\rho_k^2 - \rho_{k+1}^2}{2\gamma_k} - \frac{\mu}{4} \rho_k^2 \right) &\leq \rho_1^2 \left(\frac{1}{2\gamma_1} - \frac{\mu}{4} \right) + \sum_{k=2}^{N+1} \rho_k^2 \left(\frac{1}{2\gamma_k} - \frac{1}{2\gamma_{k-1}} - \frac{\mu}{4} \right) \\ &= \rho_1^2 \left(\frac{\mu}{4} - \frac{\mu}{4} \right) + \sum_{k=2}^{N+1} \rho_k^2 \left(\frac{\mu}{4} - \frac{\mu}{4} \right) = 0. \end{aligned} \quad (60)$$

Substituting (60) into (58) with $\gamma_k = \frac{2}{\mu k}$ we obtain

$$\begin{aligned} \mathbb{E}[\varphi(\bar{x}_N, y) - \varphi(x, \bar{y}_N)] &\leq \frac{1}{\mu N} \sum_{k=1}^N \left((\kappa_\beta L)^2 n \tau_k^{2(\beta-1)} + \kappa \left(9nM^2 + \frac{3(n\sigma)^2}{2\tau_k^2} \right) \frac{1}{k} \right) \\ &= \frac{1}{\mu N} \sum_{k=1}^N \left(\left[n \cdot (\kappa_\beta L)^2 \tau_k^{2(\beta-1)} + n^2 \cdot \frac{3\kappa\sigma^2}{2k\tau_k^2} \right] + \frac{9\kappa n M^2}{k} \right). \end{aligned} \quad (61)$$

If $\sigma > 0$ then $\tau_k = \left(\frac{3\kappa\sigma^2 n}{2(\beta-1)(\kappa_\beta L)^2} \right)^{\frac{1}{2\beta}} k^{-\frac{1}{2\beta}}$ is the minimizer of square brackets. Plugging this τ_k in (61) and using two inequalities: for the expression in square brackets $\sum_{k=1}^N k^{-1+1/\beta} \leq \beta N^{1/\beta}$ (if $\beta > 2$) and for the term after square

brackets $\sum_{k=1}^N \frac{1}{k} \leq 1 + \ln N$ we get

$$\mathbb{E}[\varphi(\bar{x}_N, y) - \varphi(x, \bar{y}_N)] \leq \frac{1}{\mu} \left(n^{2-\frac{1}{\beta}} \frac{A_1}{N^{\frac{\beta-1}{\beta}}} + A_2 \frac{n(1 + \ln N)}{N} \right).$$

with A_1 and A_2 from the formulation of Theorem 6.

Taking the minimum over x and the maximum over y we finally obtain

$$\begin{aligned} \mathbb{E}[\varphi(\bar{x}_N, y^*) - \varphi(x^*, \bar{y}_N)] &\leq \max_{y \in \mathcal{Y}} \mathbb{E}[\varphi(\bar{x}_N, y)] - \min_{x \in \mathcal{X}} \mathbb{E}[\varphi(x, \bar{y}_N)] \\ &\leq \frac{1}{\mu} \left(n^{2-\frac{1}{\beta}} \frac{A_1}{N^{\frac{\beta-1}{\beta}}} + A_2 \frac{n(1 + \ln N)}{N} \right). \end{aligned}$$

□

Theorem 7. Let $\varphi \in \mathcal{F}_\beta(L)$ with $L > 0$ and $\beta > 2$. Let Assumption 6 hold and let \mathcal{Z} be a convex compact subset of \mathbb{R}^n . Let φ be M -Lipschitz on the Euclidean τ_1 -neighborhood of \mathcal{Z} (τ_k is parameter from Theorem 6 for the regularized function $\varphi_\mu(z)$ whose description is given below). Let \bar{z}_N denote $\frac{1}{N} \sum_{k=1}^N z_k$.

Let's define $N(\varepsilon)$:

$$N(\varepsilon) = \max \left\{ \left(R\sqrt{2A_1} \right)^{\frac{2\beta}{\beta-1}} \frac{n^{2+\frac{1}{\beta-1}}}{\varepsilon^{2+\frac{2}{\beta-1}}}, \left(R\sqrt{2c'A_2} \right)^{2(1+\rho)} \frac{n^{1+\rho}}{\varepsilon^{2(1+\rho)}} \right\},$$

where $A_1 = 3\beta(\kappa\sigma^2)^{\frac{\beta-1}{\beta}} (\kappa_\beta L)^{\frac{2}{\beta}}$, $A_2 = 9\kappa G^2$ – constants from Theorem 6, $\rho > 0$ – arbitrarily small positive number, c' – constant which depends on ρ .

Then the rate of convergence is given by the following expression:

$$\mathbb{E}[\varphi(\bar{x}_N, y^*) - \varphi(x^*, \bar{y}_N)] \leq \max_{y \in \mathcal{Y}} \mathbb{E}[\varphi(\bar{x}_N, y)] - \min_{x \in \mathcal{X}} \mathbb{E}[\varphi(x, \bar{y}_N)] \leq \varepsilon \quad (62)$$

after $N(\varepsilon)$ steps of Algorithm 2 with settings from Theorem 6 for the regularized function: $\varphi_\mu(z) := \varphi(z) + \frac{\mu}{2} \|x - x_0\|^2 - \frac{\mu}{2} \|y - y_0\|^2$, where $\mu \leq \frac{\varepsilon}{R^2}$, $R = \|z_0 - z^*\|$, $z_0 \in \mathcal{Z}$ – arbitrary point.

Proof. Step 1. Let $z^* = (x^*, y^*)$ and $z_\mu^* = (x_\mu^*, y_\mu^*)$ denote the solutions of the saddle-point problems for functions $\varphi(z)$ and $\varphi_\mu(z)$ respectively. Setting $\mu = \frac{\varepsilon}{R^2}$ and using the inequality $\varphi_\mu(\bar{x}_N, y^*) - \varphi_\mu(x^*, \bar{y}_N) \leq \varphi_\mu(\bar{x}_N, y_\mu^*) - \varphi_\mu(x_\mu^*, \bar{y}_N)$ we obtain

$$\begin{aligned}
& \mathbb{E}[\varphi(\bar{x}_N, y^*)] - \mathbb{E}[\varphi(x^*, \bar{y}_N)] \leq \max_{y \in \mathcal{Y}} \mathbb{E}[\varphi(\bar{x}_N, y)] - \min_{x \in \mathcal{X}} \mathbb{E}[\varphi(x, \bar{y}_N)] \\
& = \max_{x \in \mathcal{X}, y \in \mathcal{Y}} \mathbb{E} \left[\varphi_\mu(\bar{x}_N, y) - \varphi_\mu(x, \bar{y}_N) - \frac{\mu x_N^2}{2} + \frac{\mu y^2}{2} + \frac{\mu x^2}{2} - \frac{\mu y_N^2}{2} \right] \\
& \leq \max_{x \in \mathcal{X}, y \in \mathcal{Y}} \mathbb{E} \left[\varphi_\mu(\bar{x}_N, y) - \varphi_\mu(x, \bar{y}_N) + \frac{\mu z^2}{2} \right] \\
& \leq \max_{x \in \mathcal{X}, y \in \mathcal{Y}} \mathbb{E}[\varphi_\mu(\bar{x}_N, y) - \varphi_\mu(x, \bar{y}_N)] + \frac{\varepsilon}{2} \\
& = \max_{y \in \mathcal{Y}} \mathbb{E}[\varphi_\mu(\bar{x}_N, y)] - \min_{x \in \mathcal{X}} \mathbb{E}[\varphi_\mu(x, \bar{y}_N)] + \frac{\varepsilon}{2}
\end{aligned} \tag{63}$$

Step 2. Now we apply Theorem 6 for $\varphi_\mu(z)$ until function error is not greater than $\frac{\varepsilon}{2}$:

$$\max_{y \in \mathcal{Y}} \mathbb{E}[\varphi_\mu(\bar{x}_N, y)] - \min_{x \in \mathcal{X}} \mathbb{E}[\varphi_\mu(x, \bar{y}_N)] \leq \frac{1}{\mu} \left(n^{2-\frac{1}{\beta}} \frac{A_1}{N^{\frac{\beta-1}{\beta}}} + A_2 \frac{n(1+\ln N)}{N} \right) \leq \frac{\varepsilon}{2}. \tag{64}$$

Using that $\mu = \frac{\varepsilon}{R^2}$ the inequality (64) is done if

$$\max \left\{ n^{2-\frac{1}{\beta}} \frac{A_1}{N^{\frac{\beta-1}{\beta}}}, A_2 \frac{n(1+\ln N)}{N} \right\} \leq \frac{\mu \varepsilon}{2} = \frac{\varepsilon^2}{2R^2}. \tag{65}$$

It is true that $1 + \ln N \leq c' N^{\frac{\rho}{\rho+1}}$ for some $c' > 0$. So the inequality (65) holds if

$$N \geq \max \left\{ \left(R\sqrt{2A_1} \right)^{\frac{2\beta}{\beta-1}} \frac{n^{2+\frac{1}{\beta-1}}}{\varepsilon^{2+\frac{2}{\beta-1}}}, \left(R\sqrt{2c'A_2} \right)^{2(1+\rho)} \frac{n^{1+\rho}}{\varepsilon^{2(1+\rho)}} \right\}. \tag{66}$$

The inequalities (63) and (64) yield (62). □

E Kernel examples

A weighted sum of Legendre polynomials is an example of such kernels:

$$K_\beta(r) := \sum_{m=0}^{l(\beta)} p'_m(0) p_m(r), \tag{67}$$

where $l(\beta)$ is maximal integer number strictly less than β and $p_m(r) = \sqrt{2m+1} L_m(r)$, $L_m(u)$ is Legendre polynom. We have

$$\mathbb{E}[p_m p_{m'}] = \delta(m - m').$$

As $\{p_m(r)\}_{m=0}^j$ is a basis for polynomials of degree less than or equal to j we can represent $u^j := \sum_{m=0}^j b_m p_m(r)$ for some integers $\{b_m\}_{m=0}^j$ (they depend on j).

Let's calculate the expectation

$$\mathbb{E}[r^j K_\beta(r)] = \sum_{m=0}^j b_m p'_m(0) = (r^j)'|_{r=0} = \delta(j-1),$$

here $\delta(0) = 1$ and $\delta(x) = 0$ if $x \neq 0$. We proved that the presented $K_\beta(r)$ satisfies (21). We have the following kernels for different betas (see Figure 2):

$$\begin{aligned} K_\beta(r) &= 3r, & \beta &\in [2, 3], \\ K_\beta(r) &= \frac{15r}{4}(5 - 7r^2), & \beta &\in (3, 5], \\ K_\beta(r) &= \frac{105r}{64}(99r^4 - 126r^2 + 35), & \beta &\in (5, 7]. \end{aligned}$$

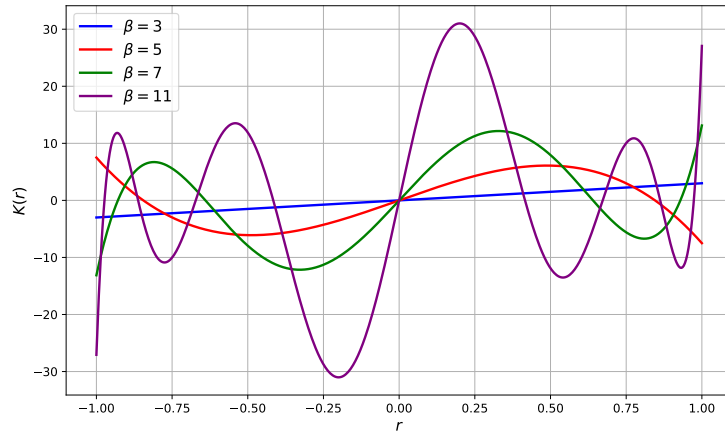


Fig. 2. Examples of kernels from (67)