

Tensor methods for strongly convex strongly concave saddle point problems and strongly monotone variational inequalities

Petr Ostroukhov¹ · Rinat Kamalov² · Pavel Dvurechensky³ · Alexander Gasnikov⁴

Received: date / Accepted: date

Abstract In this paper we propose three p -th order tensor methods for μ -strongly-convex-strongly-concave saddle point problems (SPP). The first method is based on the assumption of p -th order smoothness of the objective and it achieves a convergence rate of $O\left(\left(\frac{L_p R^{p-1}}{\mu}\right)^{\frac{2}{p+1}} \log \frac{\mu R^2}{\varepsilon_G}\right)$, where R is an estimate of the initial distance to the solution, and ε_G is the error in terms of duality gap. Under additional assumptions of first and second order smoothness of the objective we connect the first method with a locally superlinear converging algorithm and develop a second method with the complexity of $O\left(\left(\frac{L_p R^{p-1}}{\mu}\right)^{\frac{2}{p+1}} \log \frac{L_2 R \max\{1, \frac{L_1}{\mu}\}}{\mu} + \log \frac{\log \frac{L_1^3}{2\mu^2 \varepsilon_G}}{\log \frac{L_1 L_2}{\mu^2}}\right)$. The third method is a modified version of the second method, and it solves gradient norm minimization SPP with $\tilde{O}\left(\left(\frac{L_p R^p}{\varepsilon_\nabla}\right)^{\frac{2}{p+1}}\right)$ oracle calls, where ε_∇ is an error in terms of norm of the gradient of the objective. Since we treat SPP as a particular case of variational inequalities, we also propose three methods for strongly

The work of P. Ostroukhov was fulfilled in Sirius (Sochi) in August 2020 and was supported by Andrei M. Raigorodskii Scholarship in Optimization. The research of P. Dvurechensky was partially supported by the Ministry of Science and Higher Education of the Russian Federation (Goszadaniye) 075-00337-20-03, project no. 0714-2020-0005. The research of A. Gasnikov was funded by Math+ AA4-2 Scholarship in Optimization.

¹ Moscow Institute of Physics and Technology, Dolgoprudny, Russia; Institute for Information Transmission Problems RAS, Moscow, Russia E-mail: ostroukhov@phystech.edu ·

² Moscow Institute of Physics and Technology, Dolgoprudny, Russia; V. A. Trapeznikov Institute of Control Sciences of Russian Academy of Sciences, Moscow, Russia E-mail: kamalov.ra@phystech.edu ·

³ Weierstrass Institute for Applied Analysis and Stochastics, Berlin; Institute for Information Transmission Problems RAS, Moscow E-mail: pavel.dvurechensky@wias-berlin.de ·

⁴ Moscow Institute of Physics and Technology, Dolgoprudny, Russia; Weierstrass Institute for Applied Analysis and Stochastics, Berlin; Institute for Information Transmission Problems RAS, Moscow, Russia E-mail: gasnikov@yandex.ru

monotone variational inequalities with the same complexity as the described above.

Keywords Variational inequality · Saddle point problem · High-order smoothness · Tensor methods · Gradient norm minimization

1 Introduction

In this work we focus on two types of saddle point problems (SPP). The first one is the classic minimax problem:

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

where $g : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}$ is a convex over \mathcal{X} and concave over \mathcal{Y} , and the sets \mathcal{X}, \mathcal{Y} are convex. This is a particular case of a more general problem, called monotone variational inequality (MVI). In MVI we have a monotone operator $F : \mathcal{Z} \rightarrow \mathbb{R}^n$ over a convex set $\mathcal{Z} \subset \mathbb{R}^n$ and we need to find

$$z^* \in \mathcal{Z} : \forall z \in \mathcal{Z}, \langle F(z), z^* - z \rangle \leq 0. \quad (2)$$

If we set $\mathcal{Z} = \mathcal{X} \times \mathcal{Y}$ and $F(z) = (\nabla_x g(x, y), -\nabla_y g(x, y))$, then MVI is equivalent to the min-max SPP (1).

The second problem is gradient norm minimization of SPP:

$$\min_{(x, y) \in \mathcal{X} \times \mathcal{Y}} \|\nabla g(x, y)\|_2. \quad (3)$$

For both problems we consider unconstrained case with $\mathcal{X} = \mathbb{R}^n$ and $\mathcal{Y} = \mathbb{R}^m$. Additionally, we assume $g(x, y)$ is μ -strongly convex in $x \in \mathbb{R}^n$ and μ -strongly concave in $y \in \mathbb{R}^m$.

There is a number of papers on numerical methods for SPP (1) in convex-concave setting [13, 18, 20, 27, 28]. One of the most popular among first-order methods for this setting is the Mirror-Prox algorithm [18], which treats saddle-point problems via solving the corresponding MVI. According to [19], this method achieves optimal complexity of $O(1/\varepsilon)$ iterations for first-order methods applied to smooth convex-concave SPP in large dimensions.

Additional assumption of strong convexity and strong concavity lead to better results. The algorithms from [8, 15, 23, 25, 27] achieve iteration complexity of $O(L/\mu \log(1/\varepsilon))$. In [14] the authors proposed an algorithm with complexity $O(L/\sqrt{\mu_x \mu_y} \log^3(1/\varepsilon))$, which matches up to a logarithmic factor the lower bound, obtained in [29]. It worths to mention that $\log^3(1/\varepsilon)$ factor can be improved, namely, it is possible to achieve iteration complexity of $O(L/\sqrt{\mu_x \mu_y} \log(1/\varepsilon))$ (see [5]).

The methods listed above use first-order oracles, and it is known from optimization that tensor methods, which use higher-order derivatives, have faster convergence rate, yet for the price of more expensive iteration. The idea of using derivatives of high order in optimization is not new (see [10]). The most common type of high-order methods use second-order oracles, for example

Newton method [21, 24] and its modifications such as the cubic regularized Newton method [22]. Recently the idea of exploiting oracles beyond the second order started to attract increased attention, especially in convex optimization [1, 3, 4, 6, 7].

However, much less is known on high-order methods for SPP and MVIs. In [17] the authors propose a second-order method based on their Hybrid Proximal Extragradient framework [16]. The resulting complexity is $O(1/\varepsilon^{\frac{2}{3}})$. A recent work [2] shows how to modify Mirror-Prox method using oracles beyond second order and improves complexity to reach duality gap ε to $O(1/\varepsilon^{\frac{2}{p+1}})$ for convex-concave problems with p -th order Lipschitz derivatives. The paper [11] proposes a cubic regularized Newton method for solving SPP, which has global linear and local superlinear convergence rate if $\nabla g(x, y)$ and $\nabla^2 g(x, y)$ are Lipschitz-continuous and $g(x, y)$ is strongly convex in x and strongly concave in y .

In our work we make a next step and propose a Tensor method for strongly monotone variational inequalities and, as a corollary, a Tensor method for saddle point problems with strongly-convex-strongly-concave objective. Standing on the ideas from [2] and [11], our work can be split into three parts.

Firstly, we apply restart technique [26] to the HighOrderMirrorProx Algorithm 1 from [2], which is possible because of strong convexity and strong concavity of the objective. Such a modification improves the algorithm complexity to $O\left(\left(\frac{L_p R^{p-1}}{\mu}\right)^{\frac{2}{p+1}} \log \frac{\mu R^2}{\varepsilon_G}\right)$, where R is an upper bound for the initial distance to the solution $\|(x_1, y_1) - (x^*, y^*)\|_2$ and L_p is the Lipschitz constant of the p -th derivative, and ε_G is the error in terms of duality gap.

Secondly, using an estimate of the area of local superlinear convergence, when the algorithm reaches this area, we switch to the Cubic-Regularized Newton Algorithm 3 from [11] to obtain local superlinear convergence of our algorithm. The total complexity of the final Algorithm 4 becomes

$O\left(\left(\frac{L_p R^{p-1}}{\mu}\right)^{\frac{2}{p+1}} \log \frac{L_2 R \max\{1, \frac{L_1}{\mu}\}}{\mu} + \log \frac{\log \frac{L_1^3}{2\mu^2 \varepsilon_G}}{\log \frac{L_1 L_2}{\mu^2}}\right)$, where L_1 and L_2 are Lipschitz constants for first and second order derivatives respectively. We want to emphasize, that the obtained $\log \log(1/\varepsilon)$ dependency on ε cannot be improved even in convex optimization [12].

Thirdly, we apply framework from [4] to the Algorithm 4 to solve the problem (3) and obtain the Algorithm 5. Its convergence rate is $\tilde{O}\left(\left(\frac{L_p R^p}{\varepsilon_{\nabla}}\right)^{\frac{2}{p+1}}\right)$, where by tilde we mean additional multiplicative log factor, and ε_{∇} is an error in terms of gradient norm of the objective.

Our paper is organized as follows. First of all, in Section 2 we provide necessary notations and assumptions (Section 2.1). Then, we present the new algorithm and obtain its convergence rate in Section 3. Firstly, in Section 3.1 we talk only about restarted algorithm from [2] and get its complexity. Secondly, in Section 3.2 we describe how to connect it to Algorithm 3 from [11] in its quadratic convergence area and get the final Algorithm 4 convergence rate.

Thirdly, in Section 3.3 we focus on how to wrap Algorithm 4 in a framework from [4] and obtain its complexity. Finally, in Section 4 we discuss our results and present some possible directions for future work.

2 Preliminaries

We use $z \in \mathbb{R}^n \times \mathbb{R}^m$ to denote the pair (x, y) , $\nabla^p g(z)[h_1, \dots, h_p]$, $p \geq 1$ to denote directional derivative of g at z along directions $h_i \in \mathbb{R}^n \times \mathbb{R}^m$, $i = 1, \dots, p$. The norm of the p -th order derivative is defined as

$$\|\nabla^p g(z)\|_2 := \max_{h_1, \dots, h_p \in \mathbb{R}^n \times \mathbb{R}^m} \{|\nabla^p g(z)[h_1, \dots, h_p]| : \|h_i\|_2 \leq 1, i = 1, \dots, p\}$$

or equivalently

$$\|\nabla^p g(z)\|_2 := \max_{h \in \mathbb{R}^n \times \mathbb{R}^m} \{|\nabla^p g(z)[h]^p| : \|h\|_2 \leq 1\}.$$

Here we denote $\nabla^p g(z)[h, \dots, h]$ as $\nabla^p g(z)[h]^p$. Also here and below $\|\cdot\|_2$ is a Euclidean norm for vectors.

Taylor approximation of some function f at point z up to the order of p we denote by

$$\Phi_{z,p}^f(\hat{z}) := \sum_{i=0}^p \frac{1}{i!} \nabla^i f(z)[\hat{z} - z]^i.$$

For ease of notation, the Taylor approximation of the objective g we denote by $\Phi_{(x,y),p}(\hat{x}, \hat{y}) \equiv \Phi_{z,p}(\hat{z}) \equiv \Phi_{z,p}^g(\hat{z})$.

By $D : \mathcal{Z} \times \mathcal{Z} \rightarrow \mathbb{R}^n$ we denote Bregman divergence induced by a function $d : \mathcal{Z} \rightarrow \mathbb{R}$, which is continuously-differentiable and 1-strongly convex. The definition of Bregman divergence is

$$D(z_1, z_2) := d(z_1) - d(z_2) - \langle \nabla d(z_2), z_1 - z_2 \rangle.$$

In our paper we use half of squared Euclidean distance as Bregman divergence

$$D(z_1, z_2) = \frac{1}{2} \|z_1 - z_2\|_2^2. \quad (4)$$

During the analysis of convergence of our approach for gradient norm minimization (3) we will need the regularized Taylor approximation of objective g :

$$\Omega_{(x,y),p,L_p}(\hat{x}, \hat{y}) := \Phi_{(x,y),p}(\hat{x}, \hat{y}) + \frac{L_p(\sqrt{2})^{p-1}}{(p+1)!} \|\hat{x} - x\|_2^{p+1} - \frac{L_p(\sqrt{2})^{p-1}}{(p+1)!} \|\hat{y} - y\|_2^{p+1}.$$

Its min-max point we denote by

$$T_{p,L_p}^g(x, y) \in \text{Arg min}_{\tilde{x} \in \mathbb{R}^n} \max_{\tilde{y} \in \mathbb{R}^m} \{\Omega_{(x,y),p,L_p}(\tilde{x}, \tilde{y})\}.$$

As we mentioned earlier, in this paper we consider two types of SPP: classic minimax problem (1) and gradient norm minimization (3). We need to introduce the definitions of approximate solutions of these problems. We use different indices in error notations for these problems to avoid ambiguity.

Firstly, the problem (1) is usually solved in terms of the duality gap

$$G_{\mathcal{X} \times \mathcal{Y}}(x, y) := \max_{y' \in \mathcal{Y}} g(x, y') - \min_{x' \in \mathcal{X}} g(x', y). \quad (5)$$

Since in our case $\mathcal{X} = \mathbb{R}^n$ and $\mathcal{Y} = \mathbb{R}^m$, we drop the notations of these sets from index of the duality gap and denote duality gap just as $G(x, y)$. Then, we define ε_G -approximate solution of (1):

$$\tilde{x}^* \in \mathbb{R}^n, \tilde{y}^* \in \mathbb{R}^m \Rightarrow G(\tilde{x}^*, \tilde{y}^*) \leq \varepsilon_G. \quad (6)$$

Secondly, for the problem (3) we don't need any additional functionals, and ε_{∇} -approximate solution of (3) is of the form

$$\tilde{x}^* \in \mathbb{R}^n, \tilde{y}^* \in \mathbb{R}^m \Rightarrow \|\nabla g(\tilde{x}^*, \tilde{y}^*)\|_2 \leq \varepsilon_{\nabla}. \quad (7)$$

2.1 Assumptions

We assume objective g is strongly convex, strongly concave and p -times differentiable.

Assumption 1 $g(x, y)$ is μ -strongly convex in x and μ -strongly concave in y .

Recall that the definition of strong convexity and strong concavity is as follows.

Definition 1 $g : \mathbb{R}^n \times \mathbb{R}^m \rightarrow \mathbb{R}$ is called μ -strongly convex and μ -strongly concave if

$$\forall x_1, x_2 \in \mathbb{R}^n, y \in \mathbb{R}^m \Rightarrow \langle \nabla_x g(x_1, y) - \nabla_x g(x_2, y), x_1 - x_2 \rangle \geq \mu \|x_1 - x_2\|_2^2, \quad (8)$$

$$\forall y_1, y_2 \in \mathbb{R}^m, x \in \mathbb{R}^n \Rightarrow \langle -\nabla_y g(x, y_1) + \nabla_y g(x, y_2), y_1 - y_2 \rangle \geq \mu \|y_1 - y_2\|_2^2. \quad (9)$$

Before showing the connection between problem (1) and MVI (2) we need the definition of strong monotonicity.

Definition 2 $F : \mathcal{Z} \rightarrow \mathbb{R}^n$ is strongly monotone if

$$\langle F(z_1) - F(z_2), z_1 - z_2 \rangle \geq \mu \|z_1 - z_2\|_2^2. \quad (10)$$

Denote $z = \begin{pmatrix} x \\ y \end{pmatrix}$, and operator $F : \mathbb{R}^n \times \mathbb{R}^m \rightarrow \mathbb{R}^n \times \mathbb{R}^m$:

$$F(z) = F(x, y) := \begin{pmatrix} \nabla_x g(x, y) \\ -\nabla_y g(x, y) \end{pmatrix}. \quad (11)$$

According to these definitions, the min-max problem (1) can be tackled via solving the MVI problem (2) with the specific operator F given in (11). In our work we use the following assumptions.

Assumption 2 $F(z)$ satisfies first order Lipschitz condition:

$$\begin{aligned} & \|F(z_1) - F(z_2)\|_2 \leq L_1 \|z_1 - z_2\|_2 \\ \Leftrightarrow & \|\nabla g(z_1) - \nabla g(z_2)\|_2 \leq L_1 \|z_1 - z_2\|_2. \end{aligned} \quad (12)$$

Assumption 3 $F(z)$ satisfies second order Lipschitz condition:

$$\begin{aligned} & \|\nabla F(z_1) - \nabla F(z_2)\|_2 \leq L_2 \|z_1 - z_2\|_2 \\ \Leftrightarrow & \|\nabla^2 g(z_1) - \nabla^2 g(z_2)\|_2 \leq L_2 \|z_1 - z_2\|_2. \end{aligned} \quad (13)$$

Assumption 4 $F(z)$ satisfies p -th order Lipschitz condition (p -smooth):

$$\begin{aligned} & \|\nabla^{p-1} F(z_1) - \nabla^{p-1} F(z_2)\|_2 \leq L_p \|z_1 - z_2\|_2 \\ \Leftrightarrow & \|\nabla^p g(z_1) - \nabla^p g(z_2)\|_2 \leq L_p \|z_1 - z_2\|_2. \end{aligned} \quad (14)$$

We should note, that, to be consistent with [2], we define p -th order smoothness (Lipschitzness) of F as a property of $(p-1)$ -th derivative of F , and, therefore, as a property of p -th derivative of g .

3 Main results

Firstly, in this section we propose the algorithm for finding ε_G -approximate solution to problem (6), where $g: \mathbb{R}^n \times \mathbb{R}^m \rightarrow \mathbb{R}$ is p -smooth and μ -strongly-convex-concave (assumptions 4 and 1), which allows to achieve iteration complexity of $O\left(\left(\frac{L_p R^{p-1}}{\mu}\right)^{\frac{2}{p+1}} \log \frac{\mu R^2}{\varepsilon_G}\right)$, where $R \geq \|z_1 - z^*\|_2$. This algorithm is a restarted modification of Algorithm 1.

Secondly, we develop the algorithm for tackling the same problem, where g is first, second and p -th order Lipschitz and μ -strongly-convex-concave function (all assumptions 1, 2, 3, 4). It involves the idea of exploiting previous algorithm and then switching to the Algorithm 3 in its quadratic convergence area. Thus, we obtain the Algorithm 4, that allows to achieve iteration complexity of $O\left(\left(\frac{L_p R^{p-1}}{\mu}\right)^{\frac{2}{p+1}} \log \frac{L_2 R \max\{1, \frac{L_1}{\mu}\}}{\mu} + \log \frac{\log \frac{L_1^3}{2\mu^2 \varepsilon_G}}{\log \frac{L_1 L_2}{\mu^2}}\right)$.

Thirdly, we propose the algorithm to find ε_∇ -approximate solution to problem (7), where all the assumptions 1, 2, 3, 4 hold. To achieve this we use the Algorithm 4, which we mentioned earlier, inside the framework from [4]. Final complexity of such algorithm in terms of norm of the gradient is $\tilde{O}\left(\left(\frac{L_p R^p}{\varepsilon_\nabla}\right)^{\frac{2}{p+1}}\right)$, where by tilde we mean additional multiplicative log factor.

Algorithm 1 HighOrderMirrorProx [Algorithm 1 in [2]]

-
- 1: **Input** $z_1 \in \mathcal{Z}, p \geq 1, T > 0$.
 - 2: **for** $t = 1$ **to** T **do**
 - 3: Determine γ_t, \hat{z}_t such that:

$$\begin{aligned} \hat{z}_t &= \arg \min_{z \in \mathcal{Z}} \{ \gamma_t \langle \Phi_{z_t, p}^F(\hat{z}_t), z - z_t \rangle + D(z, z_t) \}, \\ \frac{p!}{32L_p \|\hat{z}_t - z_t\|_2^{p-1}} &\leq \gamma_t \leq \frac{p!}{16L_p \|\hat{z}_t - z_t\|_2^{p-1}}, \\ z_{t+1} &= \arg \min_{z \in \mathcal{Z}} \{ \gamma_t F(\hat{z}_t), z - \hat{z}_t \rangle + D(z, z_t) \}. \end{aligned}$$

- 4: Define $\Gamma_T \stackrel{\text{def}}{=} \sum_{t=1}^T \gamma_t$
 - 5: **return** $\bar{z}_T \stackrel{\text{def}}{=} \frac{1}{\Gamma_T} \sum_{t=1}^T \gamma_t \hat{z}_t$.
-

3.1 Restarted HighOrderMirrorProx

As mentioned earlier, in this subsection we provide restarted modification of Algorithm 1. But, initially, we need to give some additional information from [2].

Since our goal is an approximate solution to MVI, we define its ε -approximate solution as

$$z^* \in \mathcal{Z} : \forall z \in \mathcal{Z} \Rightarrow \langle F(z), z^* - z \rangle \leq \varepsilon. \quad (15)$$

At the same time, the bounds of Algorithm 1 is of the form

$$\forall z \in \mathcal{Z} \Rightarrow \frac{1}{\Gamma_T} \sum_{t=1}^T \gamma_t \langle F(z_t), z_t - z \rangle \leq \varepsilon, \quad (16)$$

where points z_t and $\gamma_t > 0$ are produced by the Algorithm 1, and $\Gamma_T = \sum_{t=1}^T \gamma_t$. The following lemma establishes the relation between (15) and (16).

Lemma 1 (Lemma 2.7 from [2]) *Let $F : \mathcal{Z} \rightarrow \mathbb{R}^n$, be monotone, $z_t \in \mathcal{Z}$, $t = 1, \dots, T$, and let $\gamma_t > 0$. Let $\bar{z}_t = \frac{1}{\Gamma_T} \sum_{t=1}^T \gamma_t z_t$. Assume (16) holds. Then \bar{z}_t is an ε -approximate solution to (2).*

MVI problem (2), which is sometimes called "weak MVI", is closely connected to strong MVI problem, where we need to find

$$z^* \in \mathcal{Z} : \forall z \in \mathcal{Z} \Rightarrow \langle F(z^*), z^* - z \rangle \leq 0. \quad (17)$$

If F is continuous and monotone, the problems (2) and (17) are equivalent.

The convergence rate of the Algorithm 1 is stated in the following lemma.

Lemma 2 (Lemma 4.1 from [2]) *Suppose $F : \mathcal{Z} \rightarrow \mathbb{R}^n$ is p -th order Lipschitz and let $\Gamma_T = \sum_{t=1}^T \gamma_t$. Then, the iterates $\{\hat{z}_t\}_{t \in [T]}$, generated by Algorithm 1, satisfy*

$$\forall z \in \mathcal{Z} \Rightarrow \frac{1}{\Gamma_T} \sum_{t=1}^T \langle \gamma_t F(\hat{z}_t), \hat{z}_t - z \rangle \leq \frac{16L_p}{p!} \left(\frac{D(z, z_1)}{T} \right)^{\frac{p+1}{2}}. \quad (18)$$

Algorithm 2 Restarted HighOrderMirrorProx

```

1: Input  $z_1 \in \mathcal{Z}, p \geq 1, 0 < \varepsilon_G < 1, R : R \geq \|z_1 - z^*\|_2$ .
2:  $k = 1$ 
3:  $\tilde{z}_1 = z_1$ 
4: for  $i \in [n]$ , where  $n = \lceil \frac{1}{2} \log \frac{\mu R^2}{\varepsilon_G} \rceil$  do
5:   Set  $R_i = \frac{R}{2^{i-1}}$ 
6:   Set  $T_i = \left\lceil \left( \frac{64L_p R_i^{p-1}}{\mu} \right)^{\frac{2}{p+1}} \right\rceil$ 
7:   Run Algorithm 1 with  $\tilde{z}_i, p, T_i$  as input
8:    $\tilde{z}_{i+1} = \bar{z}_{T_i}$ 
9: return  $\tilde{z}_i$ 

```

Thus, these two lemmas tell us, that if z_t and γ_t are generated by the Algorithm 1, and the right hand side of (18) is smaller than ε , then $\bar{z}_t = \frac{1}{T} \sum_{t=1}^T \gamma_t z_t$ is an ε -solution to regular MVI (15). Hence, it is also a solution to a convex-concave SPP. The natural way to improve the method for convex-concave problem in tighter strongly-convex-strongly-concave setting is to use restarts [26]. As a result, we obtain Algorithm 2.

Theorem 1 Suppose $F : \mathbb{R}^n \times \mathbb{R}^m \rightarrow \mathbb{R}^n \times \mathbb{R}^m$, that is defined in (11), is p -th order Lipschitz and μ -strongly monotone (Assumptions 1 and 4 hold). Denote R such that $R \geq \|z_1 - z^*\|_2$. Then Algorithm 2 complexity is

$$O \left(\left(\frac{L_p R^{p-1}}{\mu} \right)^{\frac{2}{p+1}} \log \frac{\mu R^2}{\varepsilon_G} \right). \quad (19)$$

Proof From (17) and (18) we get the following:

$$\sum_{t=1}^T \gamma_t \langle F(\hat{z}_t) - F(z^*); \hat{z}_t - z^* \rangle \leq \frac{16L_p}{p!} \left(\frac{\|z_1 - z^*\|_2^2}{2T} \right)^{\frac{p+1}{2}}. \quad (20)$$

From this and the fact that $F(x)$ is μ -strongly monotone we have

$$\begin{aligned} \mu \|\bar{z}_T - z^*\|_2^2 &\stackrel{(*)}{\leq} \frac{\mu}{T} \sum_{t=1}^T \gamma_t \|\hat{z}_t - z^*\|_2^2 \stackrel{(10)}{\leq} \frac{1}{T} \sum_{t=1}^T \gamma_t \langle F(\hat{z}_t) - F(z^*); \hat{z}_t - z^* \rangle \\ &\stackrel{(20)}{\leq} \frac{16L_p}{p!} \left(\frac{\|z_1 - z^*\|_2^2}{2T} \right)^{\frac{p+1}{2}}, \end{aligned} \quad (21)$$

where (*) follows from convexity of $\|z\|_2^2$.

Now we restart the method every time the distance to solution decreases at least twice. Let T_i be such that $\|\bar{z}_{T_i} - z^*\|_2 \leq \frac{\|\tilde{z}_i - z^*\|_2}{2}$, where \tilde{z}_i is the point, where we restart our algorithm. Denote $R_1 = R \geq \|\tilde{z}_1 - z^*\|_2$, $R_i =$

$R_1/2^{i-1} \geq \|\tilde{z}_i - z^*\|_2$. Then the number of iterations before $(i+1)$ -th restart is

$$\begin{aligned} \mu \|\tilde{z}_{T_i} - z^*\|_2^2 &\stackrel{(21)}{\leq} \frac{16L_p}{p!} \left(\frac{\|\tilde{z}_i - z^*\|_2^2}{2T_i} \right)^{\frac{p+1}{2}} \leq \frac{16L_p}{p!} \left(\frac{R_i^2}{2T_i} \right)^{\frac{p+1}{2}} \leq \frac{\mu \|\tilde{z}_i - z^*\|_2^2}{4} \leq \frac{\mu R_i^2}{4} \\ \Leftrightarrow T_i &\geq \frac{R_i^2}{2} \left(\frac{64L_p}{p! \mu R_i^2} \right)^{\frac{2}{p+1}} \geq \left(\frac{64L_p R_i^{p-1}}{\mu} \right)^{\frac{2}{p+1}} = \left\lceil \left(\frac{64L_p R_i^{p-1}}{\mu} \right)^{\frac{2}{p+1}} \right\rceil. \end{aligned}$$

Next we need to obtain the number of restarts n , required to achieve the desired accuracy. From (20) we get

$$\begin{aligned} \frac{1}{T_n} \sum_{t=1}^{T_n} \gamma_t \langle F(\hat{z}_t) - F(z^*); \hat{z}_t - z^* \rangle &\leq \frac{16L_p}{p!} \left(\frac{\|\tilde{z}_n - z^*\|_2^2}{2T_n} \right)^{\frac{p+1}{2}} \\ &\leq 16L_p \left(\frac{R_n^2}{\left(\frac{64L_p R_n^{p-1}}{\mu} \right)^{\frac{2}{p+1}}} \right)^{\frac{p+1}{2}} \\ &= \frac{\mu R_n^2}{4} = \frac{\mu R^2}{2^{2n}} \leq \varepsilon_G. \end{aligned}$$

$$\Leftrightarrow n \geq \frac{1}{2} \log \frac{\mu R^2}{\varepsilon_G} = \left\lceil \frac{1}{2} \log \frac{\mu R^2}{\varepsilon_G} \right\rceil.$$

Finally, the total number of iterations is

$$\begin{aligned} N &= \sum_{i=1}^n T_i = \sum_{i=1}^n \left\lceil \left(\frac{64L_p R_i^{p-1}}{\mu} \right)^{\frac{2}{p+1}} \right\rceil \leq \left(\frac{64L_p}{\mu} \right)^{\frac{2}{p+1}} \sum_{i=1}^n R_i^{\frac{2(p-1)}{p+1}} + n \\ &\leq \left(\frac{64L_p R^{p-1}}{\mu} \right)^{\frac{2}{p+1}} n + n \\ &= \left(\frac{64L_p R^{p-1}}{\mu} \right)^{\frac{2}{p+1}} \left\lceil \frac{1}{2} \log \frac{\mu R^2}{\varepsilon_G} \right\rceil + \left\lceil \frac{1}{2} \log \frac{\mu R^2}{\varepsilon_G} \right\rceil \\ &= O \left(\left(\frac{L_p R^{p-1}}{\mu} \right)^{\frac{2}{p+1}} \log \frac{\mu R^2}{\varepsilon_G} \right). \end{aligned}$$

This completes the proof. \square

3.2 Local quadratic convergence

Just like in previous subsection, besides introducing the Algorithm 3 and its convergence rate we need to provide some prerequisite information from [11].

Algorithm 3 CRN-SPP [Algorithm 1 in [11]]

```

1: Input  $z_0, \varepsilon, \bar{\gamma} > 0, \rho, \alpha \in (0, 1), g$  satisfies Assumptions 1, 2 and 3.
2: while  $m(z_k) > \varepsilon$  do
3:    $\gamma_k = \bar{\gamma}$ 
4:   while True do
5:     Solve the subproblem  $(\tilde{x}_{k+1}, \tilde{y}_{k+1}) = \arg \min_x \max_y g_k(x, y; \gamma_k)$ 
6:     if  $\gamma_k(\|\tilde{x}_{k+1} - x_k\| + \|\tilde{y}_{k+1} - y_k\|) > \mu$  then
7:        $\gamma_k = \rho\gamma_k$ 
8:     else
9:       break
10:     $d_k = (\tilde{x}_{k+1} - x_k; \tilde{y}_{k+1} - y_k)$ 
11:    if  $m(z_k + \alpha d_k) < m(z_k + d_k)$  then
12:       $z_{k+1} = z_k + \alpha d_k$ 
13:    else if  $m(z_k + \alpha d_k) \geq m(z_k + d_k)$  then
14:       $z_{k+1} = z_k + d_k$ 
15:     $k = k + 1$ 
16: return  $z_k$ 

```

Because of strong convexity and strong concavity of $g(x, y)$ a unique solution z^* to a SPP (1) exists, and $F(z^*) = 0$. Thus, we can use the following merit function from [11] during analysis of Algorithm 3 complexity.

$$m(z) := \frac{1}{2} \|F(z)\|_2^2 = \frac{1}{2} (\|\nabla_x g(x, y)\|_2^2 + \|\nabla_y g(x, y)\|_2^2). \quad (22)$$

Algorithm 3 solves additional saddle point subproblem on each step, that we denote as

$$\begin{aligned} & \min_{x \in \mathbb{R}^n} \max_{y \in \mathbb{R}^m} g_k(x, y, \gamma_k) := \\ & g(z_k) + \langle \nabla g(z_k), z - z_k \rangle + \frac{1}{2} \nabla^2 g(z_k) [z - z_k]^2 + \frac{\gamma_k}{3} \|x - x_k\|_2^3 - \frac{\gamma_k}{3} \|y - y_k\|_2^3, \end{aligned}$$

where γ_k is some constant.

This proposition provides the relation between the merit function $m(z)$ and the duality gap under assumptions 1 and 2.

Proposition 1 (Proposition 2.5 from [11]) *Let assumptions 1 and 2 hold. For problem (1) and any point $z = (x, y)$ the duality gap (5) and the merit function (22) satisfy the following inequalities*

$$\frac{\mu}{L_1^2} m(z) \leq G(x, y) \leq \frac{L_1}{\mu^2} m(z). \quad (23)$$

The next theorem proves local quadratic convergence of the Algorithm 3, and it is based on Theorem 3.6 from [11].

Theorem 2 (Theorem 3.6 from [11]) *Suppose $F : \mathcal{Z} \rightarrow \mathbb{R}^n$ is μ -strongly monotone, first and second order Lipschitz operator (assumptions 1, 2 and 3 hold). Let $\{z_k\}$ be generated by Algorithm 3 with $\bar{\gamma} = \frac{L_2 \mu^2}{2L^2}$, $\xi = \max \left\{ 1, \frac{L_1}{\mu} \right\}$ and*

$$z_0 : \|z_0 - z^*\|_2 \leq \frac{\mu}{L_2 \xi}. \quad (24)$$

Algorithm 4 Restarted HighOrderMirrorProx with local quadratic convergence

1: **Input** $z_1 \in \mathcal{Z}, p \geq 1, 0 < \varepsilon_G < 1, R : R \geq \|z_1 - z^*\|_2, \rho \in (0, 1), \alpha \in (0, 1)$.
2: $\tilde{z}_1 = z_1$
3: **for** $i \in [n]$, where $n = \lceil \log \frac{L_2 R \xi}{\mu} + 1 \rceil$ **do**
4: Set $R_i = \frac{R}{2^{i-1}}$
5: Set $T_i = \left\lfloor \frac{R_i^2}{2} \left(\frac{64L_p}{p\mu R_i} \right)^{\frac{2}{p+1}} \right\rfloor$
6: Run Algorithm 1 with \tilde{z}_i, p, T_i as input
7: $\tilde{z}_{i+1} = \tilde{z}_{T_i}$
8: Run Algorithm 3 with $\tilde{z}_{i+1}, \tilde{\varepsilon} = \frac{\mu^2 \varepsilon_G}{L}, \tilde{\gamma} = \frac{L_2 \mu^2}{2L_1^2}, \rho, \alpha, g$ as input
9: **return** z_k

Then

$$\forall k \geq 0 \quad \|z_{k+1} - z^*\|_2 \leq \frac{L_2 \xi}{\mu} \|z_k - z^*\|_2^2, \quad (25)$$

Proof Here we provide only the modified part of its proof. The rest of it can be found in [11].

If $z_{k+1} = \tilde{z}_{k+1} = z_k + d_k$, then

$$\|z_{k+1} - z^*\|_2 = \|\tilde{z}_{k+1} - z^*\|_2 \leq \frac{L_2}{\mu} \|z^k - z^*\|_2^2 \leq \frac{L_2 \xi}{\mu} \|z_k - z^*\|_2^2.$$

Else if $z_{k+1} = \hat{z}_{k+1} = z_k + \alpha d_k$, then

$$\|z_{k+1} - z^*\|_2 = \|\hat{z}_{k+1} - z^*\|_2 \leq \frac{L_1 L_2}{\mu^2} \|z^k - z^*\|_2^2 \leq \frac{L_2 \xi}{\mu} \|z_k - z^*\|_2^2.$$

Hence, we get (25).

Now we need to find the area, where (25) works:

$$\begin{aligned} \exists c : \forall k \geq 0 : \|z_k - z^*\|_2 \leq c &\Rightarrow \|z_{k+1} - z^*\|_2 \leq \frac{L_2 \xi}{\mu} \|z_k - z^*\|_2^2 \\ &\Leftrightarrow \|z_{k+1} - z^*\|_2 \leq \frac{L_2 \xi}{\mu} \|z_k - z^*\|_2 \leq \frac{L_2 \xi c^2}{\mu} = c \\ &\Leftrightarrow c = \frac{\mu}{L_2 \xi}. \end{aligned}$$

Thus, we get (24). \square

Our idea is to use Algorithm 2 until it reaches the area (24) and then switch to Algorithm 3. Algorithm 4 provides the pseudocode of this idea. From Proposition 1, our Theorem 1 and Theorem 2, we obtain the complexity of Algorithm 4.

Theorem 3 Suppose $F : \mathbb{R}^n \times \mathbb{R}^m \rightarrow \mathbb{R}^n \times \mathbb{R}^m$, that is defined in (11), is μ -strongly monotone, first, second and p -th order Lipschitz operator (all assumptions 1, 2, 3, 4 hold). Denote $R : R \geq \|z_1 - z^*\|_2$ and $\xi = \max \left\{ 1, \frac{L_1}{\mu} \right\}$. Then the complexity of Algorithm 4 is

$$O \left(\left(\frac{L_p R^{p-1}}{\mu} \right)^{\frac{2}{p+1}} \log \frac{L_2 \xi R}{\mu} + \log \frac{\log \frac{L_1^3}{2\mu^2 \varepsilon_G}}{\log \frac{L_1 L_2}{\mu^2}} \right). \quad (26)$$

Proof First of all, we need to find the number of restarts n of Algorithm 2 to reach the area of local quadratic convergence of Algorithm 3 from (24): $\|\tilde{z}_n - z^*\|_2 \leq \frac{\mu}{L_2 \xi}$. We can choose such n , that

$$\|\tilde{z}_n - z^*\|_2 \leq R_n \leq \frac{\mu}{L_2 \xi}.$$

Therefore, the number of restarts is

$$\frac{R}{2^{n-1}} \leq \frac{\mu}{L_2 \xi} \Leftrightarrow n = \left\lceil \log \frac{L_2 R \xi}{\mu} + 1 \right\rceil.$$

Next we switch to Algorithm 3 and we need to obtain its number of iterations until convergence. Denote by ε' the accuracy of solution in terms of the merit function (22). Owing to first order Lipschitzness of $F(z)$ and the fact that $F(z^*) = 0$, we can get

$$\varepsilon' = m(z_k) = \frac{1}{2} \|F(z_k)\|_2^2 = \frac{1}{2} \|F(z_k) - F(z^*)\|_2^2 \leq \frac{L_1^2}{2} \|z_k - z^*\|_2^2. \quad (27)$$

Now we establish a connection between the solution in terms of merit function $m(z)$ and the duality gap $G(x, y)$. From (27) and (23) we get the following:

$$\begin{aligned} \varepsilon_G = G(x, y) &= \max_{y' \in \mathbb{R}^n} f(x, y') - \min_{x' \in \mathbb{R}^n} f(x', y) \leq \frac{L_1}{\mu^2} m(z_k) = \frac{L_1}{\mu^2} \varepsilon' \\ &\Leftrightarrow \frac{\mu^2 \varepsilon_G}{L_1} \leq \varepsilon'. \end{aligned} \quad (28)$$

Then, from (25), (24), (27) and (28) we can obtain the needed number of iterations k

$$\begin{aligned} &\frac{\mu^2 \varepsilon_G}{L_1} \stackrel{(27),(28)}{\leq} \frac{L_1^2}{2} \|z_k - z^*\|_2^2 \\ &\stackrel{(25)}{\leq} \frac{L_1^2}{2} \left(\frac{L_1 L_2}{\mu^2} \|z_{k-1} - z^*\|_2^2 \right)^2 \leq \frac{L_1^2}{2} \left(\frac{L_1 L_2}{\mu^2} \left(\frac{L_1 L_2}{\mu^2} \|z_{k-2} - z^*\|_2^2 \right)^2 \right)^2 \leq \dots \\ &\leq \frac{L_1^2}{2} \left(\frac{L_1 L_2}{\mu^2} \right)^{2^{k-1}-2} \|z_1 - z^*\|_2^{2^k} \stackrel{(24)}{\leq} \frac{L_1^2}{2} \left(\frac{L_1 L_2}{\mu^2} \right)^{2^{k-1}-2} \left(\frac{\mu^2}{L_1 L_2} \right)^{2^k} \\ &\Leftrightarrow \frac{2\mu^2 \varepsilon_G}{L_1^3} \leq \left(\frac{\mu^2}{L_1 L_2} \right)^{2^{k-1}+2} \Leftrightarrow \log \frac{2\mu^2 \varepsilon_G}{L_1^3} \leq (2^{k-1} + 2) \log \frac{\mu^2}{L_1 L_2} \end{aligned}$$

Since $\log(\mu^2/L_1L_2) < 0$,

$$\log \frac{2\mu^2\varepsilon_G}{L_1^3} \leq 2^{k-1} \log \frac{\mu^2}{L_1L_2} \Leftrightarrow k = \left\lceil \log \frac{\log \frac{L_1^3}{2\mu^2\varepsilon_G}}{\log \frac{L_1L_2}{\mu^2}} \right\rceil + 1.$$

Finally, the total number of iterations of Algorithm 4 is

$$\begin{aligned} N &= \sum_{i=1}^n T_i + k \\ &\leq \left(\frac{64L_pR^{p-1}}{\mu} \right)^{\frac{2}{p+1}} \left[\log \frac{L_2\xi R}{\mu} + 1 \right] + \left[\log \frac{L_2\xi R}{\mu} + 1 \right] + \left[\log \frac{\log \frac{L_1^3}{2\mu^2\varepsilon_G}}{\log \frac{L_1L_2}{\mu^2}} \right] + 1 \\ &= O \left(\left(\frac{L_pR^{p-1}}{\mu} \right)^{\frac{2}{p+1}} \log \frac{L_2\xi R}{\mu} + \log \frac{\log \frac{L_1^3}{2\mu^2\varepsilon_G}}{\log \frac{L_1L_2}{\mu^2}} \right) \end{aligned}$$

□

3.3 Gradient norm minimization

In this subsection we apply the framework from [4] to Algorithm 4, introduce Algorithm 5 for problem (7) and analyze its complexity in terms of the norm of the gradient $\|\nabla g(x, y)\|_2$.

Firstly, we need to introduce some technical lemmas.

Lemma 3 *If $g(x, y)$ is p -Lipchitz (14), then its partial p -th order derivatives are also Lipschitz.*

$$\forall \hat{x}, x \in \mathbb{R}^n, \hat{y}, y \in \mathbb{R}^m \Rightarrow \|\nabla_{x^i y^{p-i}}^p g(\hat{x}, \hat{y}) - \nabla_{x^i y^{p-i}}^p g(x, y)\|_2 \leq L_p \|\hat{z} - z\|_2. \quad (29)$$

Proof Here we provide proof only for $\nabla_{x\dots x}^p$. For other partial derivatives the proof is analogous.

From definition of $\|\cdot\|_2$

$$\begin{aligned} \|\nabla_{x\dots x}^p g(\hat{x}, \hat{y}) - \nabla_{x\dots x}^p g(x, y)\|_2 &= \max_{\|s\|_2 \leq 1} |(\nabla_{x\dots x}^p g(\hat{x}, \hat{y}) - \nabla_{x\dots x}^p g(x, y))[s]^p| \\ &= \max_{\|s\|_2 \leq 1} \left| (\nabla^p g(\hat{x}, \hat{y}) - \nabla^p g(x, y)) \left[\begin{pmatrix} s \\ 0 \end{pmatrix} \right]^p \right| \\ &\leq \max_{\|h\|_2 \leq 1} |(\nabla^p g(\hat{x}, \hat{y}) - \nabla^p g(x, y))[h]^p| \\ &= \|\nabla^p g(\hat{x}, \hat{y}) - \nabla^p g(x, y)\|_2 \leq L_p \|\hat{z} - z\|_2. \end{aligned}$$

□

Lemma 4 *Let $\nabla_{x\dots x}^p g(x, y)$ be Lipschitz (29). Then*

$$\forall n \in [p] \Rightarrow \|\nabla_{x\dots x}^{p-n} g(\hat{z}) - \nabla_{x\dots x}^{p-n} \Phi_{(x,y),p}(\hat{z})\|_2 \leq \frac{L_p(\sqrt{2})^n}{(n+1)!} \|\hat{z} - z\|_2^{n+1}. \quad (30)$$

Proof We prove this by induction.

The base of induction $n = 1$ follows from the definition of Taylor approximation. Denote $f(z) = \nabla_{x\dots x}^{p-1} g(z)$.

$$\begin{aligned} & \|\nabla_{x\dots x}^{p-1} g(\hat{z}) - \nabla_{x\dots x}^{p-1} \Phi_{(x,y),p}(\hat{z})\|_2 \\ = & \|\nabla_{x\dots x}^{p-1} g(\hat{z}) - \nabla_{x\dots x}^{p-1} g(z) - \nabla_{x\dots xx}^p g(z)[\hat{x} - x] - \nabla_{x\dots xy}^p g(z)[\hat{y} - y]\|_2 \\ & = \|f(\hat{z}) - f(z) - \nabla f(z)[\hat{z} - z]\|_2 \\ & = \left\| \int_0^1 \langle \nabla f(z + \tau(\hat{z} - z)) - \nabla f(z); \hat{z} - z \rangle d\tau \right\|_2 \\ & \leq \int_0^1 \left\| \begin{pmatrix} \nabla_{x\dots xx}^p g(z + \tau(\hat{z} - z)) \\ \nabla_{x\dots xy}^p g(z + \tau(\hat{z} - z)) \end{pmatrix} - \begin{pmatrix} \nabla_{x\dots xx}^p g(z) \\ \nabla_{x\dots xy}^p g(z) \end{pmatrix} \right\|_2 \|\hat{z} - z\|_2 d\tau \\ = & \int_0^1 \sqrt{\|\nabla_{x\dots xx}^p g(z + \tau(\hat{z} - z)) - \nabla_{x\dots xx}^p g(z)\|_2^2 + \|\nabla_{x\dots xy}^p g(z + \tau(\hat{z} - z)) - \nabla_{x\dots xy}^p g(z)\|_2^2} \\ & \quad \cdot \|\hat{z} - z\|_2 d\tau \\ & \stackrel{(29)}{\leq} \sqrt{2} L_p \|\hat{z} - z\|_2^2 \int_0^1 \tau d\tau = \frac{L_p \sqrt{2}}{2} \|\hat{z} - z\|_2^2. \end{aligned}$$

Now assume it holds for $n = p - 1$:

$$\begin{aligned} & \|\nabla_x g(\hat{z}) - \nabla_x \Phi_{(x,y),p}(\hat{z})\|_2 \\ = & \left\| \nabla_x g(\hat{z}) - \nabla_x g(z) - (\nabla_{xx}^2 g(z)[\hat{x} - x] - \nabla_{xy}^2 g(z)[\hat{y} - y]) - \dots - \right. \\ & \quad \left. - \nabla_x \left(\frac{1}{p!} \nabla^p g(z)[\hat{z} - z]^p \right) \right\|_2 \\ & \leq \frac{L_p(\sqrt{2})^{p-1}}{p!} \|\hat{z} - z\|_2^p. \quad (31) \end{aligned}$$

And consider $n = p$

$$\begin{aligned}
& |g(\hat{z}) - \Phi_{(x,y),p}(\hat{z})| \\
&= |g(\hat{z}) - g(z) - \nabla_x g(z)[\hat{x} - x] - \nabla_y g(z)[\hat{y} - y] - \dots - \frac{1}{p!} \nabla^p g(z)[\hat{z} - z]^p| \\
&\leq \int_0^1 \left\| \begin{pmatrix} \nabla_x g(z + \tau(\hat{z} - z)) \\ \nabla_y g(z + \tau(\hat{z} - z)) \end{pmatrix} - \begin{pmatrix} \nabla_x g(z) \\ \nabla_y g(z) \end{pmatrix} - \right. \\
&\quad \left. - \tau \begin{pmatrix} \nabla_{xx}^2 g(z)[\hat{x} - x] + \nabla_{xy}^2 g(z)[\hat{y} - y] \\ \nabla_{yx}^2 g(z)[\hat{x} - x] + \nabla_{yy}^2 g(z)[\hat{y} - y] \end{pmatrix} - \dots - \right. \\
&\quad \left. - \frac{\tau^{p-1}}{p!} \begin{pmatrix} \nabla_x(\nabla^p g(z)[\hat{z} - z]^p) \\ \nabla_y(\nabla^p g(z)[\hat{z} - z]^p) \end{pmatrix} \right\|_2 \|\hat{z} - z\|_2 d\tau \\
&= \int_0^1 \left(\|\nabla_x g(z + \tau(\hat{z} - z)) - \nabla_x g(z) - \right. \\
&\quad \left. - \tau(\nabla_{xx}^2 g(z)[\hat{x} - x] + \nabla_{xy}^2 g(z)[\hat{y} - y]) - \dots - \right. \\
&\quad \left. - \frac{\tau^{p-1}}{p!} \nabla_x(\nabla^p g(z)[\hat{z} - z]^p) \right\|_2^2 + \\
&\quad + \|\nabla_y g(z + \tau(\hat{z} - z)) - \nabla_y g(z) - \\
&\quad - \tau(\nabla_{yx}^2 g(z)[\hat{x} - x] + \nabla_{yy}^2 g(z)[\hat{y} - y]) - \dots - \\
&\quad \left. - \frac{\tau^{p-1}}{p!} \nabla_y(\nabla^p g(z)[\hat{z} - z]^p) \right\|_2^2)^{1/2} \|\hat{z} - z\|_2 d\tau.
\end{aligned}$$

If we denote $\hat{z} = z + \tau(\hat{z} - z)$ in (31), each of two factors under the square root is indeed what we had for $n = p - 1$. Finally,

$$\begin{aligned}
\|\nabla_x g(\hat{z}) - \nabla_x \Phi_{(x,y),p}(\hat{z})\|_2 &\leq \sqrt{2} \frac{L_p(\sqrt{2})^{p-1}}{p!} \|\hat{z} - z\|_2^{p+1} \int_0^1 \tau^p d\tau \\
&= \frac{L_p(\sqrt{2})^p}{(p+1)!} \|\hat{z} - z\|_2^{p+1}.
\end{aligned}$$

For any other partial derivative in (30) the result is the same and can be obtained in a similar way. \square

The next lemma is a modified version of Lemma 5.2 from [9] for SPP.

Lemma 5 (Lemma 5.2 from [9]) *Let $(\tilde{x}, \tilde{y}) = T_{p,M}^g(x, y)$, $p \geq 2$, where $M \geq \sqrt{2}pL_p > \frac{1}{\sqrt{2}}pL_p$ and assumption 4 hold. Then*

$$\|\nabla g(\tilde{x}, \tilde{y})\|_2^{\frac{p+1}{p}} \frac{M^{\frac{3p+1}{2p}}}{2^{\frac{2p^2+p+1}{2p}} p(p+1)!} \leq g(x, \tilde{y}) - g(\tilde{x}, y). \quad (32)$$

Proof

$$\|\nabla g(\tilde{x}, \tilde{y})\|_2^2 = \|\nabla_x g(\tilde{x}, \tilde{y})\|_2^2 + \|\nabla_y g(\tilde{x}, \tilde{y})\|_2^2.$$

Firstly, consider ∇_x :

$$\begin{aligned} \|\nabla_x g(\tilde{x}, \tilde{y})\|_2^2 &= \|\nabla_x g(\tilde{x}, \tilde{y}) - \nabla_x \Phi_{(x,y),p}(\tilde{x}, \tilde{y}) + \nabla_x \Phi_{(x,y),p}(\tilde{x}, \tilde{y}) - \\ &\quad - \nabla_x \Omega_{(x,y),p,M}(\tilde{x}, \tilde{y}) + \nabla_x \Omega_{(x,y),p,M}(\tilde{x}, \tilde{y})\|_2^2 \\ &\leq \left(\|\nabla_x g(\tilde{x}, \tilde{y}) - \nabla_x \Phi_{(x,y),p}(\tilde{x}, \tilde{y})\|_2 + \right. \\ &\quad \left. + \|\nabla_x \Phi_{(x,y),p}(\tilde{x}, \tilde{y}) - \nabla_x \Omega_{(x,y),p,M}(\tilde{x}, \tilde{y})\|_2 + \|\nabla_x \Omega_{(x,y),p,M}(\tilde{x}, \tilde{y})\|_2 \right)^2 \\ &\leq \left(\frac{2^{\frac{p-1}{2}} L_p}{p!} \|\tilde{z} - z\|_2^p + \frac{2^{\frac{p-1}{2}} M}{p!} \|\tilde{x} - x\|_2^p \right)^2 \leq 2^p M^2 \|\tilde{z} - z\|_2^{2p}. \end{aligned}$$

For ∇_y in a similar way we get the same result

$$\|\nabla_y g(\tilde{x}, \tilde{y})\|_2^2 \leq 2^p M^2 \|\tilde{z} - z\|_2^{2p}.$$

Summing these two results, we obtain

$$\|\nabla g(\tilde{x}, \tilde{y})\|_2^2 \leq 2^{p+1} M (\|\tilde{x} - x\|_2^2 + \|\tilde{y} - y\|_2^2)^p. \quad (33)$$

Secondly, consider point (\tilde{x}, y) . From (30) it is obvious that

$$|g(\tilde{x}, y) - \Phi_{(x,y),p}(\tilde{x}, y)| \leq \frac{L_p(\sqrt{2})^p}{(p+1)!} \|(\tilde{x}, y) - (x, y)\|_2^{p+1} = \frac{L_p(\sqrt{2})^p}{(p+1)!} \|\tilde{x} - x\|_2^{p+1}.$$

From this fact we get

$$\begin{aligned} g(\tilde{x}, y) &\leq \Phi_{(x,y),p}(\tilde{x}, y) + \frac{L_p(\sqrt{2})^p}{(p+1)!} \|\tilde{x} - x\|_2^{p+1} \\ &= \Phi_{(x,y),p}(\tilde{x}, y) + \frac{L_p(\sqrt{2})^{p-1}}{(p+1)!} \|\tilde{x} - x\|_2^{p+1} - \\ &\quad - \left(\frac{M(\sqrt{2})^{p-1}}{(p+1)!} \|\tilde{x} - x\|_2^{p+1} - \frac{L_p(\sqrt{2})^p}{(p+1)!} \|\tilde{x} - x\|_2^{p+1} \right) \\ &= \Omega_{(x,y),p,M}(\tilde{x}, y) - (M - L_p \sqrt{2}) \frac{(\sqrt{2})^{p-1} \|\tilde{x} - x\|_2^{p+1}}{(p+1)!} \\ &\leq \Omega_{(x,y),p,M}(\tilde{x}, \tilde{y}) - (M - L_p \sqrt{2}) \frac{(\sqrt{2})^{p-1} \|\tilde{x} - x\|_2^{p+1}}{(p+1)!}. \end{aligned}$$

Since $M \geq \sqrt{2} p L_p \Leftrightarrow -L_p \sqrt{2} \geq -\frac{M}{p}$, we have

$$\Omega_{(x,y),p,M}(\tilde{x}, \tilde{y}) - g(\tilde{x}, y) \geq \frac{M(p-1)(\sqrt{2})^{p-1} \|\tilde{x} - x\|_2^{p+1}}{p(p+1)!} \geq \frac{M \|\tilde{x} - x\|_2^{p+1}}{p(p+1)!}. \quad (34)$$

Now consider the point (x, \tilde{y}) . In a similar way we can get the following result:

$$g(x, \tilde{y}) - \Omega_{(x,y),p,M}(x, \tilde{y}) \geq \frac{M \|\tilde{y} - y\|_2^{p+1}}{p(p+1)!}. \quad (35)$$

From the sum of (34) and (35) we obtain

$$g(x, \tilde{y}) - g(\tilde{x}, y) \geq \frac{M}{p(p+1)!} \left(\|\tilde{x} - x\|_2^{p+1} + \|\tilde{y} - y\|_2^{p+1} \right). \quad (36)$$

Finally, we need to connect (33) and (36). From Hölder's inequality we can get

$$\left(\sum_{i=1}^n x_i^p \right)^{\frac{1}{p}} \leq n^{\frac{q-p}{qp}} \left(\sum_{i=1}^n x_i^q \right)^{\frac{1}{q}},$$

where $q, p \in \mathbb{N}$, $q > p \geq 1$. Now, from (33) it follows that

$$\left(\frac{\|\nabla g(\tilde{x}, \tilde{y})\|_2^2}{2^{p+1}M} \right)^{\frac{1}{2p}} \leq (\|\tilde{x} - x\|_2^2 + \|\tilde{y} - y\|_2^2)^{\frac{1}{2}}.$$

And, from (36) we can get

$$\left(\frac{p(p+1)!(g(x, \tilde{y}) - g(\tilde{x}, y))}{M} \right)^{\frac{1}{p+1}} \geq \left(\|\tilde{x} - x\|_2^{p+1} + \|\tilde{y} - y\|_2^{p+1} \right)^{\frac{1}{p+1}}.$$

Since $p \geq 2$, we obtain the final result

$$\|\nabla g(\tilde{x}, \tilde{y})\|_2^{\frac{p+1}{p}} \frac{M^{\frac{3p+1}{2p}}}{2^{\frac{2p^2+p+1}{2p}} p(p+1)!} \leq g(x, \tilde{y}) - g(\tilde{x}, y).$$

□

Now we have all the needed information to estimate the final convergence rate of the Algorithm 5 for gradient norm minimization.

Theorem 4 *Assume the function $g(x, y) : \mathbb{R}^n \times \mathbb{R}^m \rightarrow \mathbb{R}$ is convex by x and concave by y , p times differentiable on \mathbb{R}^n with L_p -Lipschitz p -th derivative. Let \tilde{z} be generated by Algorithm 5. Then*

$$\|\nabla g(\tilde{z})\|_2 \leq \varepsilon_{\nabla},$$

and the total complexity of Algorithm 5 is

$$O \left(\left(\frac{L_p R^p}{\varepsilon_{\nabla}} \right)^{\frac{2}{p+1}} \log \frac{L_2 R^2 \xi}{\varepsilon_{\nabla}} \right),$$

where $\xi = \max \left\{ 1, \frac{4RL_1}{\varepsilon} \right\}$.

Proof Denote $z_{\mu}^* = (x_{\mu}^*, y_{\mu}^*)$ the saddle point of $g_{\mu}(z)$. First of all, since $g_{\mu}(x, y)$ is strongly-convex-strongly-concave function, we can apply restart technique to it every time the distance to its saddle point $\|z - z_{\mu}^*\|_2$ reduces twice. To check this, we consider upper estimate of the distance to the solution of regular

Algorithm 5 Restarted HighOrderMirrorProx with local quadratic convergence for gradient norm minimization

1: **Input** $z_1 \in \mathcal{Z}, p \geq 1, 0 < \varepsilon_\nabla < 1, R : R \geq \|z_1 - z^*\|_2, \rho \in (0, 1), \alpha \in (0, 1)$.

2: **Define:**

$$\begin{aligned} \tilde{z}_1 &= z_1, \quad M = \sqrt{2}pL_p, \quad \mu = \frac{\varepsilon}{4R}, \quad \xi = \max\left\{1, \frac{4RL_1}{\varepsilon_\nabla}\right\}, \\ \varepsilon' &= \frac{M^{\frac{3p+1}{2p}} \varepsilon_\nabla^{\frac{p+1}{p}}}{2^{\frac{2p^2+3p+3}{2p}} p(p+1)!}, \\ g_\mu(x, y) &= g(x, y) + \frac{\mu}{2}(\|x - x_1\|_2^2 - \|y - y_1\|_2^2). \end{aligned}$$

3: **for** $i \in [n]$, where $n = \lceil \log \frac{L_2 R \xi}{\mu} + 1 \rceil$ **do**

4: Set $R_i = \frac{R}{2^{i-1}}$

5: Set $T_i = \left\lceil \left(\frac{64L_p R_i^{p-1}}{p! \mu} \right)^{\frac{2}{p+1}} \right\rceil$

6: Run Algorithm 1 for g_μ with \tilde{z}_i, p, T_i as input

7: $\tilde{z}_{i+1} = \tilde{z}_{T_i}$

8: Run Algorithm 3 with $\tilde{z}_{i+1}, \varepsilon', \bar{\gamma} = \frac{L_2 \mu^2}{2L_1^2}, \rho, \alpha, g_\mu$ as input

9: **Find** $\tilde{z} = T_{p, M}^{g_\mu}(z_k)$

10: **Output** \tilde{z} .

function $R : R \geq \|z^* - z\|_2$ and show, that on each i -th restart $\|z_\mu^* - z_i\|_2 \leq \|z^* - z_i\|_2 \leq R_i$. We prove this by induction.

$$\begin{aligned} g(x_\mu^*, y_1) + \frac{\mu}{2}\|x_\mu^* - x_1\|_2^2 &= g_\mu(x_\mu^*, y_1) \leq g_\mu(x^*, y_1) = g(x^*, y_1) + \frac{\mu}{2}\|x^* - x_1\|_2^2 \\ &\leq g(x_\mu^*, y_1) + \frac{\mu}{2}\|x^* - x_1\|_2^2 \\ \Leftrightarrow \|x_\mu^* - x_1\|_2 &\leq \|x^* - x_1\|_2. \end{aligned}$$

$$\begin{aligned} g(x_1, y_\mu^*) - \frac{\mu}{2}\|y_\mu^* - y_1\|_2^2 &= g_\mu(x_1, y_\mu^*) \geq g_\mu(x_1, y^*) = g(x_1, y^*) - \frac{\mu}{2}\|y^* - y_1\|_2^2 \\ &\geq g(x_1, y_\mu^*) - \frac{\mu}{2}\|y^* - y_1\|_2^2 \\ \Leftrightarrow \|y_\mu^* - y_1\|_2 &\leq \|y^* - y_1\|_2. \end{aligned}$$

This gives us

$$\|z_\mu^* - z_1\|_2 \leq \|z^* - z_1\|_2 \leq R.$$

Now suppose, that $\|z_\mu^* - z_i\|_2 \leq \|z^* - z_i\|_2 \leq R_i = R/2^{i-1}$. Consider $i+1$. From the proof of Theorem 1 and our choice of T_i in Algorithm 5, we know, that

$$\begin{aligned} \mu\|z_{i+1} - z_\mu^*\|_2^2 &= \mu\|\tilde{z}_{T_i} - z_\mu^*\|_2^2 \leq \frac{16L_p}{p!} \left(\frac{R_i^2}{2T_i} \right)^{\frac{p+1}{2}} \leq \mu R_{i+1}^2 \\ \Leftrightarrow \|z_{i+1} - z_\mu^*\|_2 &\leq R_{i+1}. \end{aligned}$$

From Theorem 3 we already know the number of restarts to reach the area of quadratic convergence: $n = \left\lceil \log \frac{L_2 R \xi}{\mu} + 1 \right\rceil$.

Next, we need to show, that Algorithm 5 converges in terms of $\|\nabla g_\mu(z)\|_2$. Let $\tilde{z} = (\tilde{x}, \tilde{y})$ be the output of Algorithm 5. From the definition of g_μ we get

$$\begin{aligned} \|\nabla g(\tilde{x}, \tilde{y})\|_2^2 &= \|\nabla_x g_\mu(\tilde{x}, \tilde{y}) - \mu(\tilde{x} - x_1)\|_2^2 + \|\nabla_y g_\mu(\tilde{x}, \tilde{y}) + \mu(\tilde{y} - y_1)\|_2^2 \\ &\leq (\|\nabla_x g_\mu(\tilde{x}, \tilde{y})\|_2 + \mu\|\tilde{x} - x\|_2)^2 + (\|\nabla_y g_\mu(\tilde{x}, \tilde{y})\|_2 + \mu\|\tilde{y} - y\|_2)^2 \\ &\leq 2(\|\nabla_x g_\mu(\tilde{x}, \tilde{y})\|_2^2 + \|\nabla_y g_\mu(\tilde{x}, \tilde{y})\|_2^2) + 2\mu^2(\|\tilde{x} - x\|_2^2 + \|\tilde{y} - y\|_2^2) \\ &= 2\|\nabla g_\mu(\tilde{x}, \tilde{y})\|_2^2 + 2\mu^2\|\tilde{z} - z_1\|_2^2 \\ &\Leftrightarrow \|\nabla g(\tilde{x}, \tilde{y})\|_2 \leq \sqrt{2\|\nabla g_\mu(\tilde{x}, \tilde{y})\|_2^2 + 2\mu^2\|\tilde{z} - z_1\|_2^2}. \end{aligned}$$

Firstly, we estimate $\|\nabla g_\mu(\tilde{x}, \tilde{y})\|_2$. From (32) we know, that

$$\begin{aligned} \|\nabla g_\mu(\tilde{x}, \tilde{y})\|_2^{\frac{p+1}{p}} \frac{M^{\frac{3p+1}{2p}}}{2^{\frac{2p^2+p+1}{2p}} p(p+1)!} &\stackrel{(32)}{\leq} g_\mu(x, \tilde{y}) - g_\mu(\tilde{x}, y) \\ &\leq \max_{\tilde{y} \in \mathbb{R}^m} g_\mu(x, \tilde{y}) - \min_{\tilde{x} \in \mathbb{R}^n} g_\mu(\tilde{x}, y) = G_\mu(x, y) \leq \varepsilon'. \\ \Leftrightarrow \|\nabla g_\mu(\tilde{x}, \tilde{y})\|_2 &\leq \left(\frac{2^{\frac{2p^2+p+1}{2p}} p(p+1)! \varepsilon'}{M^{\frac{3p+1}{2p}}} \right)^{\frac{p}{p+1}} = \frac{\varepsilon_\nabla}{2}. \end{aligned} \quad (37)$$

Secondly, we estimate $\mu\|\tilde{z} - z_1\|_2$. By definition of R we know, that

$$\|z^* - z_1\|_2 \leq R.$$

And since \tilde{z} is closer to solution than z_1 , we have

$$\|\tilde{z} - z^*\|_2 \leq \|z^* - z_1\|_2 \leq R.$$

From these facts and triangle inequality we get

$$\mu\|\tilde{z} - z_1\|_2 \leq \mu(\|\tilde{z} - z^*\|_2 + \|z^* - z_1\|_2) \leq 2R\mu = \frac{\varepsilon_\nabla}{2}. \quad (38)$$

Thus, from (37) and (38) we obtain

$$\|\nabla g_\mu(\tilde{x}, \tilde{y})\|_2 \leq \sqrt{2\varepsilon_\nabla^2/4 + 2\varepsilon_\nabla^2/4} = \varepsilon_\nabla.$$

Finally, we need to estimate complexity of the Algorithm 5.

$$\begin{aligned} N &= \sum_{i=1}^n T_i + k \leq \left(\frac{64L_p}{p!\mu} \right)^{\frac{2}{p+1}} \sum_{i=1}^n R_i^{\frac{2(p-1)}{p+1}} + n + k \\ &\leq \left(\frac{64L_p R^{p-1}}{p!\mu} \right)^{\frac{2}{p+1}} \cdot n + n + k \\ &= O \left(\left(\frac{L_p R^p}{\varepsilon_\nabla} \right)^{\frac{2}{p+1}} \log \frac{L_2 R^2 \xi}{\varepsilon_\nabla} \right), \end{aligned}$$

where $\xi = \max\left\{1, \frac{4RL_1}{\varepsilon_\nabla}\right\}$. Here k is the number of iterations of Algorithm 3 inside Algorithm 5. We dropped it due to its $\log \log$ dependence on ε_∇ . \square

4 Discussion

In this work we propose three methods for p -th order tensor methods for strongly-convex-strongly-concave SPP. Two of these methods tackle classical minimax SPP (1) and MVI (2) problems, and the third method aims at gradient norm minimization of SPP (3).

The methods for minimax problem are based on the ideas, developed in the works [2] and [11]. In [2] the authors use p -th order oracle to construct an algorithm for MVI problems with monotone operator. As a corollary, this algorithm allows to solve SPP with convex-concave objective. Because of strong convexity and strong concavity of our problem, we can apply a restart technique to the method from [2] and get better algorithm complexity. To further improve local convergence rate we switch to the algorithm from [11] in the area of its quadratic convergence. This way we get rid of the multiplicative logarithmic factor and get additive $\log \log$ factor in the final complexity estimate and get locally quadratic convergence.

The method for gradient norm minimization relies on the works [9] and [4]. From [9] we take the result, that connects norm of the gradient of the objective with objective residual, and slightly modify it for SPP. This step allows us to use the framework from [4] and use our optimal algorithm for minimax SPP for gradient norm minimization.

In spite of all the improvements, we should remind about many additional assumptions about the problem, which reduces number of real problems, that can suit to it.

One of possible directions for further research are the more general Hölder conditions instead of Lipschitz conditions and uniformly convex case. Additionally, the author in [2] provided implementation details of the Algorithm 1 only for $p = 2$. Therefore, the questions about its realization for $p > 2$ are still opened.

References

1. Bullins, B.: Fast minimization of structured convex quartics. arXiv preprint arXiv:1812.10349 (2018)
2. Bullins, B., Lai, K.A.: Higher-order methods for convex-concave min-max optimization and monotone variational inequalities. arXiv preprint arXiv:2007.04528 (2020)
3. Bullins, B., Peng, R.: Higher-order accelerated methods for faster non-smooth optimization. arXiv preprint arXiv:1906.01621 (2019)
4. Dvurechensky, P., Gasnikov, A., Ostroukhov, P., Uribe, C.A., Ivanova, A.: Near-optimal tensor methods for minimizing the gradient norm of convex function. arXiv preprint arXiv:1912.03381 (2019)

5. Gasnikov, A., Dvinskikh, D., Dvurechensky, P., Kamzolov, D., Pasechnykh, D., Matykhin, V., Tupitsa, N., Chernov, A.: Accelerated meta-algorithm for convex optimization. *Computational Mathematics and Mathematical Physics* **61**(1) (2020)
6. Gasnikov, A., Dvurechensky, P., Gorbunov, E., Vorontsova, E., Selikhanovych, D., Uribe, C.A.: Optimal tensor methods in smooth convex and uniformly convex optimization. In: A. Beygelzimer, D. Hsu (eds.) *Proceedings of the Thirty-Second Conference on Learning Theory, Proceedings of Machine Learning Research*, vol. 99, pp. 1374–1391. PMLR, Phoenix, USA (2019). URL <http://proceedings.mlr.press/v99/gasnikov19a.html>. ArXiv:1809.00382
7. Gasnikov, A., Dvurechensky, P., Gorbunov, E., Vorontsova, E., Selikhanovych, D., Uribe, C.A., Jiang, B., Wang, H., Zhang, S., Bubeck, S., Jiang, Q., Lee, Y.T., Li, Y., Sidford, A.: Near optimal methods for minimizing convex functions with lipschitz p -th derivatives. In: A. Beygelzimer, D. Hsu (eds.) *Proceedings of the Thirty-Second Conference on Learning Theory, Proceedings of Machine Learning Research*, vol. 99, pp. 1392–1393. PMLR, Phoenix, USA (2019). URL <http://proceedings.mlr.press/v99/gasnikov19b.html>
8. Gidel, G., Berard, H., Vignoud, G., Vincent, P., Lacoste-Julien, S.: A variational inequality perspective on generative adversarial networks. arXiv preprint arXiv:1802.10551 (2018)
9. Grapiglia, G.N., Nesterov, Y.: Tensor methods for finding approximate stationary points of convex functions. arXiv preprint arXiv:1907.07053 (2019)
10. Hoffmann, K.H., Kornstaedt, H.J.: Higher-order necessary conditions in abstract mathematical programming. *Journal of Optimization Theory and Applications* **26**(4), 533–568 (1978). DOI 10.1007/BF00933151. URL <https://doi.org/10.1007/BF00933151>
11. Huang, K., Zhang, J., Zhang, S.: Cubic regularized newton method for saddle point models: a global and local convergence analysis. arXiv preprint arXiv:2008.09919 (2020)
12. Kornowski, G., Shamir, O.: High-order oracle complexity of smooth and strongly convex optimization. arXiv preprint arXiv:2010.06642 (2020)
13. Korpelevich, G.: The extragradient method for finding saddle points and other problems. *Ekonomika i Matematicheskie Metody* **12**, 747–756 (1976)
14. Lin, T., Jin, C., Jordan, M., et al.: Near-optimal algorithms for minimax optimization. arXiv preprint arXiv:2002.02417 (2020)
15. Mokhtari, A., Ozdaglar, A., Pattathil, S.: A unified analysis of extra-gradient and optimistic gradient methods for saddle point problems: Proximal point approach. In: *International Conference on Artificial Intelligence and Statistics*, pp. 1497–1507. PMLR (2020)
16. Monteiro, R.D., Svaiter, B.F.: On the complexity of the hybrid proximal extragradient method for the iterates and the ergodic mean. *SIAM Journal on Optimization* **20**(6), 2755–2787 (2010)
17. Monteiro, R.D., Svaiter, B.F.: Iteration-complexity of a newton proximal extragradient method for monotone variational inequalities and inclusion problems. *SIAM Journal on Optimization* **22**(3), 914–935 (2012)
18. Nemirovski, A.: Prox-method with rate of convergence $o(1/t)$ for variational inequalities with lipschitz continuous monotone operators and smooth convex-concave saddle point problems. *SIAM Journal on Optimization* **15**(1), 229–251 (2004)
19. Nemirovsky, A., Yudin, D.: *Problem Complexity and Method Efficiency in Optimization*. J. Wiley & Sons, New York (1983)
20. Nesterov, Y.: Dual extrapolation and its applications to solving variational inequalities and related problems. *Mathematical Programming* **109**(2-3), 319–344 (2007). First appeared in 2003 as CORE discussion paper 2003/68
21. Nesterov, Y., Nemirovskii, A.: *Interior-point polynomial algorithms in convex programming*. SIAM (1994)
22. Nesterov, Y., Polyak, B.: Cubic regularization of newton method and its global performance. *Mathematical Programming* **108**(1), 177–205 (2006). DOI 10.1007/s10107-006-0706-8. URL <http://dx.doi.org/10.1007/s10107-006-0706-8>
23. Nesterov, Y., Scramali, L.: Solving strongly monotone variational and quasi-variational inequalities. Available at SSRN 970903 (2006)
24. Nocedal, J., Wright, S.: *Numerical optimization*. Springer Science & Business Media (2006)

-
25. Rockafellar, R.T.: Monotone operators and the proximal point algorithm. *SIAM journal on control and optimization* **14**(5), 877–898 (1976)
 26. Stonyakin, F., Gasnikov, A., Dvurechensky, P., Alkousa, M., Titov, A.: Generalized Mirror Prox for monotone variational inequalities: Universality and inexact oracle. arXiv:1806.05140 (2018)
 27. Tseng, P.: On linear convergence of iterative methods for the variational inequality problem. *Journal of Computational and Applied Mathematics* **60**(1-2), 237–252 (1995)
 28. Tseng, P.: On accelerated proximal gradient methods for convex-concave optimization. Tech. rep., MIT (2008). URL <http://www.mit.edu/~dimitrib/PTseng/papers/apgm.pdf>
 29. Zhang, J., Hong, M., Zhang, S.: On lower iteration complexity bounds for the saddle point problems. arXiv preprint arXiv:1912.07481 (2019)