

---

# Accelerated Primal-Dual Gradient Method for Smooth and Convex-Concave Saddle-Point Problems with Bilinear Coupling

---

**Dmitry Kovalev**  
KAUST\*  
dakovalev1@gmail.com

**Alexander Gasnikov**  
MIPT† ISP RAS‡ HSE§  
gasnikov@yandex.ru

**Peter Richtárik**  
KAUST  
richtarik@gmail.com

## Abstract

In this paper we study the convex-concave saddle-point problem  $\min_x \max_y f(x) + y^T \mathbf{A}x - g(y)$ , where  $f(x)$  and  $g(y)$  are smooth and convex functions. We propose an Accelerated Primal-Dual Gradient Method (APDG) for solving this problem, achieving (i) an optimal linear convergence rate in the strongly-convex-strongly-concave regime, matching the lower complexity bound (Zhang et al., 2021), and (ii) an accelerated linear convergence rate in the case when only one of the functions  $f(x)$  and  $g(y)$  is strongly convex or even none of them are. Finally, we obtain a linearly convergent algorithm for the general smooth and convex-concave saddle point problem  $\min_x \max_y F(x, y)$  without the requirement of strong convexity or strong concavity.

## 1 Introduction

In this paper we revisit the well studied smooth convex-concave saddle point problem with a bilinear coupling function, which takes the form

$$\min_{x \in \mathbb{R}^{d_x}} \max_{y \in \mathbb{R}^{d_y}} F(x, y) = f(x) + y^T \mathbf{A}x - g(y), \quad (1)$$

where  $f(x) : \mathbb{R}^{d_x} \rightarrow \mathbb{R}$  and  $g(y) : \mathbb{R}^{d_y} \rightarrow \mathbb{R}$  are smooth and convex functions, and  $\mathbf{A} \in \mathbb{R}^{d_y \times d_x}$  is a coupling matrix.

Problem (1) has a large number of application, some of which we now briefly introduce.

### 1.1 Empirical risk minimization

A classical application is the regularized empirical risk minimization (ERM) with linear predictors, which is a classical supervised learning problem. Given a data matrix  $\mathbf{A} = [a_1, \dots, a_n]^T \in \mathbb{R}^{n \times d}$ , where  $a_i \in \mathbb{R}^d$  is the feature vector of the  $i$ -th data entry, our goal is to find a solution of

$$\min_x f(x) + \ell(\mathbf{A}x), \quad (2)$$

where  $f(x) : \mathbb{R}^d \rightarrow \mathbb{R}$  is a convex regularizer,  $\ell(y) : \mathbb{R}^n \rightarrow \mathbb{R}$  is a convex loss function, and  $x \in \mathbb{R}^d$  is a linear predictor. Alternatively, one can solve the following equivalent saddle-point reformulation

---

\*King Abdullah University of Science and Technology, Thuwal, Saudi Arabia

†Moscow Institute of Physics and Technology, Dolgoprudny, Russia

‡Institute for System Programming RAS, Research Center for Trusted Artificial Intelligence, Moscow, Russia

§National Research University Higher School of Economics, Moscow, Russia

of problem (2):

$$\min_x \max_y f(x) + y^\top \mathbf{A}x - \ell^*(y). \quad (3)$$

The saddle-point reformulation is often preferable. For example, when such a formulation admits a finite sum structure (Zhang and Lin, 2015; Wang and Xiao, 2017), this may reduce the communication complexity in the distributed setting (Xiao et al., 2019), and one may also better exploit the underlying sparsity structure (Lei et al., 2017).

## 1.2 Reinforcement learning

In reinforcement learning (RL) we are given a sequence  $\{(s_t, a_t, r_t, s_{t+1})\}_{t=1}^n$  generated by a policy  $\pi$ , where  $s_t$  is the state at time step  $t$ ,  $a_t$  is the action taken at time step  $t$  by policy  $\pi$  and  $r_t$  is the reward after taking action  $a_t$ . A key step in many RL algorithms is to estimate the value function of a given policy  $\pi$ , which is defined as

$$V^\pi(s) = \mathbb{E} \left[ \sum_{t=0}^{\infty} \gamma^t r_t \mid s_0 = s, \pi \right], \quad (4)$$

where  $\gamma \in (0, 1)$  is a discount factor. A common approach to this problem is to use a linear approximation  $V^\pi(s) = \phi(s)^\top x$ , where  $\phi(s)$  is a feature vector of a state  $s$ . The model parameter  $x$  is often estimated by minimizing the mean squared projected Bellman error

$$\min_x \|\mathbf{A}x - b\|_{\mathbf{C}^{-1}}^2, \quad (5)$$

where  $\mathbf{C} = \sum_{t=1}^n \phi(s_t)\phi(s_t)^\top$ ,  $b = \sum_{t=1}^n r_t \phi(s_t)$  and  $\mathbf{A} = \mathbf{C} - \gamma \sum_{t=1}^n \phi(s_t)\phi(s_{t+1})^\top$ . One can observe that it is hard to apply gradient-based methods to problem (5) because this would require one to compute an inverse of the matrix  $\mathbf{C}$ . In order to tackle this issue, one can solve an equivalent saddle-point reformulation proposed by Du et al. (2017) instead. This reformulation is given by

$$\min_x \max_y -2y^\top \mathbf{A}x - \|y\|_{\mathbf{C}}^2 + 2b^\top y, \quad (6)$$

and is an instance of problem (1). Solving this reformulation with gradient methods does not require matrix inversion.

## 1.3 Minimization under affine constraints

Next, consider the problem of convex minimization under affine constraints,

$$\min_{\mathbf{A}x=b} f(x), \quad (7)$$

where  $b \in \text{range} \mathbf{A}$ . This problem covers a wide range of applications, including inverse problems in imaging (Chambolle and Pock, 2016), sketched learning-type applications (Keriven et al., 2018), network flow optimization (Zargham et al., 2013) and optimal transport (Peyré et al., 2019).

Another important application of problem (7) is decentralized distributed optimization (Kovalev et al., 2020; Scaman et al., 2017; Li et al., 2020; Nedic et al., 2017; Arjevani et al., 2020; Ye et al., 2020). In this setting, the distributed minimization problem is often reformulated as

$$\min_{\sqrt{\mathbf{W}}(x_1, \dots, x_n)^\top = 0} \left[ f(x_1, \dots, x_n) = \sum_{i=1}^n f_i(x_i) \right], \quad (8)$$

where  $f_i(x_i)$  is a function stored locally by a computing node  $i \in \{1, \dots, n\}$  and  $\mathbf{W} \in \mathbb{R}^{n \times n}$  is the Laplacian matrix of a graph representing the communication network. The constraint enforces consensus among the nodes:  $x_1 = \dots = x_n$ .

One can observe that problem (7) is equivalent to the saddle-point formulation

$$\min_x \max_y f(x) + y^\top \mathbf{A}x - y^\top b, \quad (9)$$

which is another instance of problem (1). State-of-the-art methods often focus on this formulation instead of directly solving (7). In particular, Salim et al. (2021) and Kovalev et al. (2020) obtained optimal algorithms for solving (7) and (8) using this saddle-point approach.

Table 1: Comparison of method (APGD, Algorithm 1) with existing state-of-the-art algorithms for solving problem (1) in the 5 different cases described in section 5.

<b>Strongly-convex-strongly-concave case (section 5.1)</b>	
Algorithm 1	$\mathcal{O}\left(\max\left\{\sqrt{\frac{L_x}{\mu_x}}, \sqrt{\frac{L_y}{\mu_y}}, \frac{L_{xy}}{\sqrt{\mu_x\mu_y}}\right\}\log\frac{1}{\epsilon}\right)$
Lower bound Zhang et al. (2021b)	$\mathcal{O}\left(\max\left\{\sqrt{\frac{L_x}{\mu_x}}, \sqrt{\frac{L_y}{\mu_y}}, \frac{L_{xy}}{\sqrt{\mu_x\mu_y}}\right\}\log\frac{1}{\epsilon}\right)$
DIPPA Xie et al. (2021)	$\tilde{\mathcal{O}}\left(\max\left\{\sqrt[4]{\frac{L_x^2 L_y}{\mu_x^2 \mu_y}}, \sqrt[4]{\frac{L_x L_y^2}{\mu_x \mu_y^2}}, \frac{L_{xy}}{\sqrt{\mu_x \mu_y}}\right\}\log\frac{1}{\epsilon}\right)$
Proximal Best Response Wang and Li (2020)	$\tilde{\mathcal{O}}\left(\max\left\{\sqrt{\frac{L_x}{\mu_x}}, \sqrt{\frac{L_y}{\mu_y}}, \sqrt{\frac{L_{xy}L}{\mu_x\mu_y}}\right\}\log\frac{1}{\epsilon}\right)$
<b>Affinely constrained minimization case (section 5.2)</b>	
Algorithm 1	$\mathcal{O}\left(\frac{L_{xy}}{\mu_{xy}}\sqrt{\frac{L_x}{\mu_x}}\log\frac{1}{\epsilon}\right)$
Lower bound Salim et al. (2021)	$\mathcal{O}\left(\frac{L_{xy}}{\mu_{xy}}\sqrt{\frac{L_x}{\mu_x}}\log\frac{1}{\epsilon}\right)$
OPAPC Kovalev et al. (2020)	$\mathcal{O}\left(\frac{L_{xy}}{\mu_{xy}}\sqrt{\frac{L_x}{\mu_x}}\log\frac{1}{\epsilon}\right)$
<b>Strongly-convex-concave case (section 5.3)</b>	
Algorithm 1	$\mathcal{O}\left(\max\left\{\sqrt{\frac{L_x L_y}{\mu_{xy}}}, \frac{L_{xy}}{\mu_{xy}}\sqrt{\frac{L_x}{\mu_x}}, \frac{L_{xy}^2}{\mu_{xy}^2}\right\}\log\frac{1}{\epsilon}\right)$
Lower bound	N/A
Alt-GDA Zhang et al. (2021a)	$\mathcal{O}\left(\max\left\{\frac{L^2}{\mu_{xy}^2}, \frac{L}{\mu_x}\right\}\log\frac{1}{\epsilon}\right)$
<b>Bilinear case (section 5.4)</b>	
Algorithm 1	$\mathcal{O}\left(\frac{L_{xy}^2}{\mu_{xy}^2}\log\frac{1}{\epsilon}\right)$
Lower bound Ibrahim et al. (2020)	$\mathcal{O}\left(\frac{L_{xy}}{\mu_{xy}}\log\frac{1}{\epsilon}\right)$
Azizian et al. (2020)	$\mathcal{O}\left(\frac{L_{xy}}{\mu_{xy}}\log\frac{1}{\epsilon}\right)$
<b>Convex-concave case (section 5.5)</b>	
Algorithm 1	$\mathcal{O}\left(\max\left\{\sqrt{\frac{L_x L_y L_{xy}}{\mu_{xy}^2}}, \frac{L_{xy}^2}{\mu_{xy}^2}\right\}\log\frac{1}{\epsilon}\right)$
Lower bound	N/A

## 1.4 Bilinear min-max problems

Unconstrained bilinear saddle-point problems of the form

$$\min_{x \in \mathbb{R}^{d_x}} \max_{y \in \mathbb{R}^{d_y}} a^\top x + y^\top \mathbf{A}x - b^\top y \quad (10)$$

are another special case of problem (1), one where both  $f(x)$  and  $g(y)$  are linear functions. While such problems do not usually play an important role in practice, they are often a good testing ground for theoretical purposes (Gidel et al., 2019; Azizian et al., 2020; Zhang et al., 2021a; Mokhtari et al., 2020; Daskalakis et al., 2018; Liang and Stokes, 2019).

## 2 Literature Review and Contributions

In this work we are interested in algorithms able to solve problem (1) with a linear iteration complexity. That is, we are interested in methods that can provably find an  $\epsilon$ -accurate solution of problem (1) in a number of iterations proportional to  $\log\frac{1}{\epsilon}$  (see Definitions 2 and 3). This is typically achieved when

functions  $f(x)$  and  $g(x)$  are assumed to be strongly convex (see Definition 1). An example of this is the celebrated extragradient method of Korpelevich (1976).

Recent work has shown that linear iteration complexity can be achieved also in the less restrictive case when only one of the functions  $f(x)$  and  $g(x)$  is strongly convex. This was first shown by Du and Hu (2019), and later improved on by Zhang et al. (2021a).

*However, and this is the starting point of our research, to the best of our knowledge, there are no algorithms with linear iteration complexity in the case when neither  $f(x)$  nor  $g(x)$  is strongly convex.*

## 2.1 Acceleration

Loosely speaking, we say that an algorithm is *non-accelerated* if its iteration complexity is proportional to at least the first power of the condition numbers associated with the problem, such as  $L_x/\mu_x$  and  $L_y/\mu_y$ , where  $L_x$  and  $L_y$  are smoothness constants, and  $\mu_x$  and  $\mu_y$  are strong convexity constants (see Assumption 1 and Assumption 2). In contrast, the iteration complexity of an *accelerated* algorithm is proportional to the square root of such condition numbers, e.g.,  $\sqrt{L_x/\mu_x}$  and  $\sqrt{L_y/\mu_y}$ .

There were several recent attempts to design accelerated algorithms for solving problem (1) (Xie et al., 2021; Wang and Li, 2020; Alkousa et al., 2020). These attempts rely on *stacking multiple algorithms on top of each other*, and result in complicated methods. For example, Lin et al. (2020) use a non-accelerated algorithm as a sub-routine for the inexact accelerated proximal-point method. This approach allows them to obtain accelerated algorithms for solving problem (1) in a straightforward and tractable way. However, this approach has significant drawbacks: the algorithms obtained this way have (i) additional logarithmic factors in their iteration complexity, and (ii) a complex nested structure with the requirement to manually set inner loop sizes, which is a byproduct of the design process based on combining multiple algorithms. This drawback limits the performance of the resulting algorithms in theory, and requires additional fine tuning in practice.

A philosophically different approach to designing such algorithms—one that we adopt in this work—is to attempt to provide a *direct* acceleration of a suitable algorithm for solving problem (1), similarly to what Nesterov (1983) did for convex minimization problems. While this technically more demanding, algorithms obtained this way typically don't have the aforementioned drawbacks. Hence, we follow the latter approach in this work.

## 2.2 Main contributions

In this work we propose an Accelerated Primal-Dual Gradient Method (APDG; Algorithm 1) for solving problem (1) and provide a theoretical analysis of its convergence properties (Theorem 1). In particular, we prove the following results.

- (i) When both functions  $f(x)$  and  $g(y)$  are strongly convex, Algorithm 1 achieves the optimal linear convergence rate, matching the lower bound obtained by Zhang et al. (2021b). To the best of our knowledge, Algorithm 1 is the first optimal algorithm in this regime.
- (ii) We establish linear convergence of Algorithm 1 in the case when *only one* of the functions  $f(x)$  or  $g(y)$  is strongly convex, and  $\mathbf{A}$  is a full row or full column rank matrix, respectively. This improves upon the results provided by Du and Hu (2019); Zhang et al. (2021a).
- (iii) We establish linear convergence of the Algorithm 1 in the case when *neither* of the functions  $f(x)$  nor  $g(y)$  is strongly convex, and the matrix  $\mathbf{A}$  is square and full rank. To the best of our knowledge, Algorithm 1 is the first algorithm achieving linear convergence in this setting.

Table 1 provides a brief comparison of the complexity of Algorithm 1 (Theorem 1) with the current state of the art. Please refer to section 5 for a detailed discussion of this result and comparison with related work.

### 2.3 General min-max problem and additional contributions

In our work we also consider the saddle-point problem

$$\min_{x \in \mathbb{R}^{d_x}} \max_{y \in \mathbb{R}^{d_y}} F(x, y), \quad (11)$$

where  $F(x, y): \mathbb{R}^{d_x} \times \mathbb{R}^{d_y} \rightarrow \mathbb{R}$  is a smooth function, which is convex in  $x$  and concave in  $y$ . One can observe that the main problem (1) is a special case of this more general problem (11).

As an additional contribution, we propose a Gradient Descent-Ascent Method with Extrapolation (GDAE) for solving the general convex-concave saddle-point problem (11), and provide a theoretical analysis of its convergence properties.

- (i) When the function  $F(x, y)$  is strongly convex in  $x$  and strongly concave in  $y$ , GDAE achieves a linear convergence rate, which recovers the convergence result of Cohen et al. (2020).
- (ii) Under certain assumptions on the way the variables  $x$  and  $y$  are coupled by the function  $F(x, y)$ , we establish linear convergence of GDAE in the case when the function  $F(x, y)$  is strongly-convex-concave, convex-strongly-concave, or even just convex-concave. To the best of our knowledge, GDAE is the first algorithm achieving linear convergence under such assumptions.

Please refer to the Appendix for a detailed description of these results and related work.

### 3 Basic Definitions and Assumptions

We start by formalizing the notions of smoothness and strong convexity of a function.

**Definition 1.** *Function  $h(z): \mathbb{R}^d \rightarrow \mathbb{R}$  is  $L$ -smooth and  $\mu$ -strongly convex for  $L \geq \mu \geq 0$ , if for all  $z_1, z_2 \in \mathbb{R}^d$  the following inequality holds:*

$$\frac{\mu}{2} \|z_1 - z_2\|^2 \leq D_h(z_1, z_2) \leq \frac{L}{2} \|z_1 - z_2\|^2. \quad (12)$$

Above,  $D_h(z_1, z_2) = h(z_1) - h(z_2) - \langle \nabla h(z_2), z_1 - z_2 \rangle$  is the Bregman divergence associated with the function  $h(z)$ .

We are now ready to state the main assumptions that we impose on problem (1). We start with Assumptions 1 and 2 that formalize the strong-convexity and smoothness properties of functions  $f(x)$  and  $g(y)$ .

**Assumption 1.** *Function  $f(x)$  is  $L_x$ -smooth and  $\mu_x$ -strongly convex for  $L_x \geq \mu_x \geq 0$ .*

**Assumption 2.** *Function  $g(y)$  is  $L_y$ -smooth and  $\mu_y$ -strongly convex for  $L_y \geq \mu_y \geq 0$ .*

Note, that  $\mu_x$  and  $\mu_y$  are allowed to be zero. That is, both  $f(x)$  and  $g(y)$  are allowed to be non-strongly convex.

The following assumption formalizes the spectral properties of matrix  $\mathbf{A}$ .

**Assumption 3.** *There exist constants  $L_{xy} > \mu_{xy}, \mu_{yx} \geq 0$  such that*

$$\begin{aligned} \mu_{xy}^2 &\leq \begin{cases} \lambda_{\min}^+(\mathbf{A}\mathbf{A}^\top) & \nabla g(y) \in \text{range}\mathbf{A} \text{ for all } y \in \mathbb{R}^{d_y} \\ \lambda_{\min}(\mathbf{A}\mathbf{A}^\top) & \text{otherwise} \end{cases} \\ \mu_{yx}^2 &\leq \begin{cases} \lambda_{\min}^+(\mathbf{A}^\top\mathbf{A}) & \nabla f(x) \in \text{range}\mathbf{A}^\top \text{ for all } x \in \mathbb{R}^{d_x} \\ \lambda_{\min}(\mathbf{A}^\top\mathbf{A}) & \text{otherwise} \end{cases} \\ L_{xy}^2 &\geq \lambda_{\max}(\mathbf{A}^\top\mathbf{A}) = \lambda_{\max}(\mathbf{A}\mathbf{A}^\top), \end{aligned}$$

where  $\lambda_{\min}(\cdot)$ ,  $\lambda_{\min}^+(\cdot)$  and  $\lambda_{\max}(\cdot)$  denote the smallest, smallest positive and largest eigenvalue of a matrix, respectively, and  $\text{range}$  denotes the range space of a matrix.

By  $\mathcal{S} \subset \mathbb{R}^{d_x} \times \mathbb{R}^{d_y}$  we denote the solution set of problem (1). Note that  $(x^*, y^*) \in \mathcal{S}$  if and only if  $(x^*, y^*)$  satisfies the first-order optimality conditions

$$\begin{cases} \nabla_x F(x^*, y^*) = \nabla f(x^*) + \mathbf{A}^\top y^* = 0, \\ \nabla_y F(x^*, y^*) = -\nabla g(y^*) + \mathbf{A}x^* = 0. \end{cases} \quad (13)$$

Our main goal is to propose an algorithm for finding a solution to problem (1). Numerical iterative algorithms typically find an approximate solution of a given problem. We formalize this through the following definition.

**Definition 2.** Let the solution set  $\mathcal{S}$  be nonempty. We call a pair of vectors  $(x, y) \in \mathbb{R}^{d_x} \times \mathbb{R}^{d_y}$  an  $\epsilon$ -accurate solution of problem (1) for a given accuracy  $\epsilon > 0$  if it satisfies

$$\min_{(x^*, y^*) \in \mathcal{S}} \max \{ \|x - x^*\|^2, \|y - y^*\|^2 \} \leq \epsilon. \quad (14)$$

We also want to propose an *efficient* algorithm for solving problem (1). That is, we want to propose an algorithm with the the lowest possible *iteration complexity*, which we define next.

**Definition 3.** The iteration complexity of an algorithm for solving problem (1) is the number of iterations the algorithm requires to find an  $\epsilon$ -accurate solution of this problem. At each iteration the algorithm is allowed to perform  $\mathcal{O}(1)$  computations of the gradients  $\nabla f(x)$  and  $\nabla g(y)$  and matrix-vector multiplications with matrices  $\mathbf{A}$  and  $\mathbf{A}^\top$ .

## 4 Accelerated Primal-Dual Gradient Method

---

### Algorithm 1 APDG: Accelerated Primal-Dual Gradient Method

---

- 1: **Input:**  $x^0 \in \text{range} \mathbf{A}^\top, y^0 \in \text{range} \mathbf{A}, \eta_x, \eta_y, \alpha_x, \alpha_y, \beta_x, \beta_y > 0, \tau_x, \tau_y, \sigma_x, \sigma_y \in (0, 1], \theta \in (0, 1)$
  - 2:  $x_f^0 = x^0$
  - 3:  $y_f^0 = y^{-1} = y^0$
  - 4: **for**  $k = 0, 1, 2, \dots$  **do**
  - 5:    $y_m^k = y^k + \theta(y^k - y^{k-1})$
  - 6:    $x_g^k = \tau_x x^k + (1 - \tau_x)x_f^k$
  - 7:    $y_g^k = \tau_y y^k + (1 - \tau_y)y_f^k$
  - 8:    $x^{k+1} = x^k + \eta_x \alpha_x (x_g^k - x^k) - \eta_x \beta_x \mathbf{A}^\top (\mathbf{A}x^k - \nabla g(y_g^k)) - \eta_x (\nabla f(x_g^k) + \mathbf{A}^\top y_m^k)$
  - 9:    $y^{k+1} = y^k + \eta_y \alpha_y (y_g^k - y^k) - \eta_y \beta_y \mathbf{A} (\mathbf{A}^\top y^k + \nabla f(x_g^k)) - \eta_y (\nabla g(y_g^k) - \mathbf{A}x^{k+1})$
  - 10:    $x_f^{k+1} = x_g^k + \sigma_x (x^{k+1} - x^k)$
  - 11:    $y_f^{k+1} = y_g^k + \sigma_y (y^{k+1} - y^k)$
  - 12: **end for**
- 

In this section we present the Accelerated Primal-Dual Gradient Method (APDG; Algorithm 1) for solving problem (1). First, we prove an outline of the key ideas used in the development of this algorithm.

### 4.1 Algorithm development strategy

First, we observe that problem (1) is equivalent to the problem of finding a zero of a sum of two monotone operators,  $G_1, G_2: \mathbb{R}^{d_x} \times \mathbb{R}^{d_y} \rightarrow \mathbb{R}^{d_x} \times \mathbb{R}^{d_y}$ , defined as

$$G_1: (x, y) \mapsto (\nabla f(x), \nabla g(y)), \quad (15)$$

$$G_2: (x, y) \mapsto (\mathbf{A}^\top y, -\mathbf{A}x). \quad (16)$$

Indeed,  $G_1(x^*, y^*) + G_2(x^*, y^*) = 0$  is just another way to write the optimality conditions (13).

**The Forward Backward algorithm.** A natural way to tackle this problem is via *Forward Backward algorithm* (Bauschke and Combettes, 2011), the iterates of which have the form

$$(x^{k+1}, y^{k+1}) = J_{G_2}((x^k, y^k) - G_1(x^k, y^k)), \quad (17)$$

where the operator  $J_{G_2}$  is the inverse of the operator  $I + G_2$ , and  $I$  is the identity operator. Note that  $J_{G_2}$  can be written as  $J_{G_2}: (x, y) \mapsto (x^+, y^+)$ , where  $(x^+, y^+) \in \mathbb{R}^{d_x} \times \mathbb{R}^{d_y}$  is a solution of the linear system

$$\begin{cases} x^+ = x - \mathbf{A}^\top y^+ \\ y^+ = y + \mathbf{A}x^+ \end{cases}. \quad (18)$$

**Linear extrapolation step.** Next, notice that the computation of operator  $J_{G_2}$  requires solving the linear system (18). This is expensive<sup>5</sup> and has to be done at each iteration of the Forward Backward algorithm. Let us instead consider the related problem

$$\begin{cases} x^+ = x - \mathbf{A}^\top y_m \\ y^+ = y + \mathbf{A}x^+ \end{cases}, \quad (19)$$

where  $y_m \in \mathbb{R}^{d_y}$  is a newly introduced variable. It's easy to observe that (19) is equivalent to (18) when  $y_m = y^+$ . Next, notice that choosing  $y_m = y$  makes (19) easy to solve. However, it turns out that the convergence analysis of an algorithm with this approximation may be challenging (Zhang et al., 2021a), especially if we want to combine it with other techniques, such as acceleration. Our key idea is to propose a better alternative: the *linear extrapolation step*

$$y_m = y + \theta(y - y^-), \quad (20)$$

where  $y^- \in \mathbb{R}^{d_y}$  corresponds to  $y$  obtained from the previous iteration of the Forward Backward algorithm, and  $\theta \in (0, 1]$  is an extrapolation parameter. The linear extrapolation step was introduced by Chambolle and Pock (2011) in the analysis of the Primal-Dual Hybrid Gradient algorithm<sup>6</sup>.

**Nesterov acceleration.** Next, we note that operator  $G_1$  is equal to the gradient of the (potential) function  $(x, y) \mapsto f(x) + g(y)$  function. This function is smooth and convex due to Assumptions 1 and 2. This allows us to incorporate the *Nesterov acceleration* mechanism in the Forward Backward algorithm. Nesterov acceleration is known to be a powerful tool which allows to improve convergence properties of gradient methods (Nesterov, 1983, 2003).

## 4.2 Convergence of the algorithm

We are now ready to study the convergence properties of Algorithm 1. We are interested in the case when the following condition holds:

$$\min \{ \max \{ \mu_x, \mu_{yx} \}, \max \{ \mu_y, \mu_{xy} \} \} > 0. \quad (21)$$

In this case one can show that the solution set  $\mathcal{S}$  of problem (1) is nonempty. Moreover, *strong duality* holds in this case, as captured by the following lemma.

**Lemma 1.** *Let Assumptions 1, 2 and 3 and condition (21) hold. Let  $p$  be the optimal value of the primal problem*

$$p = \min_{x \in \mathbb{R}^{d_x}} [P(x) = f(x) + g^*(\mathbf{A}x)], \quad (22)$$

*and let  $d$  be the optimal value of the dual problem*

$$d = \max_{y \in \mathbb{R}^{d_y}} [D(y) = -g(y) - f^*(-\mathbf{A}^\top y)]. \quad (23)$$

*Then  $p = d$  is finite and  $(x^*, y^*) \in \mathcal{S}$  if and only if  $x^*$  is a solution of the primal problem (22) and  $y^*$  is a solution of the dual problem (23).*

Under the aforementioned conditions, Algorithm 1 achieves linear convergence. That is, the iteration complexity is proportional to  $\log \frac{1}{\epsilon}$ .

<sup>5</sup>The solution of (18) can be written in a closed form and requires to compute an inverse matrix  $(\mathbf{I} + \mathbf{A}^\top \mathbf{A})^{-1}$  or  $(\mathbf{I} + \mathbf{A} \mathbf{A}^\top)^{-1}$ , where  $\mathbf{I}$  is the identity matrix of an appropriate size.

<sup>6</sup>However, the Primal-Dual Hybrid Gradient algorithm is not applicable in our case since it requires to compute the proximal operator of  $f(x)$  and  $g(y)$  at each iteration. Moreover, Chambolle and Pock (2011) established linear convergence of this algorithm in the strongly-convex-strongly-concave setting only.

**Theorem 1.** *Let Assumptions 1, 2 and 3 and condition (21) hold. Then there exist parameters of Algorithm 1 such that its iteration complexity for finding an  $\epsilon$ -accurate solution of problem (1) is*

$$\mathcal{O}\left(\min\{T_a, T_b, T_c, T_d\} \log \frac{C}{\epsilon}\right), \quad (24)$$

where  $T_a, T_b, T_c, T_d$  are defined as

$$\begin{aligned} T_a &= \max\left\{\sqrt{\frac{L_x}{\mu_x}}, \sqrt{\frac{L_y}{\mu_y}}, \frac{L_{xy}}{\sqrt{\mu_x \mu_y}}\right\}, & T_b &= \max\left\{\frac{\sqrt{L_x L_y}}{\mu_{xy}}, \frac{L_{xy}}{\mu_{xy}} \sqrt{\frac{L_x}{\mu_x}}, \frac{L_{xy}^2}{\mu_{xy}^2}\right\}, \\ T_c &= \max\left\{\frac{\sqrt{L_x L_y}}{\mu_{yx}}, \frac{L_{xy}}{\mu_{yx}} \sqrt{\frac{L_y}{\mu_y}}, \frac{L_{xy}^2}{\mu_{yx}^2}\right\}, & T_d &= \max\left\{\frac{\sqrt{L_x L_y} L_{xy}}{\mu_{xy} \mu_{yx}}, \frac{L_{xy}^2}{\mu_{yx}^2}, \frac{L_{xy}^2}{\mu_{xy}^2}\right\}, \end{aligned}$$

and  $C > 0$  is some constant, which does not depend on  $\epsilon$ , but possibly depends on  $L_x, \mu_x, L_y, \mu_y, L_{xy}, \mu_{xy}, \mu_{yx}$ .

## 5 Discussion of Theorem 1 and Related Work

In this section we comment on the iteration complexity result for Algorithm 1 provided in Theorem 1. We consider important and illustrative special cases of this complexity result and draw connections with the existing results in the literature.

### 5.1 Strongly convex and strongly concave case

In this case  $\mu_x, \mu_y > 0$ . We can always assume  $\mu_{xy} = \mu_{yx} = 0$  in Assumption 3. Then, Algorithm 1 has iteration complexity given by

$$\mathcal{O}\left(\max\left\{\sqrt{\frac{L_x}{\mu_x}}, \sqrt{\frac{L_y}{\mu_y}}, \frac{L_{xy}}{\sqrt{\mu_x \mu_y}}\right\} \log \frac{1}{\epsilon}\right). \quad (25)$$

This improves the current state-of-the-art results

$$\tilde{\mathcal{O}}\left(\max\left\{\sqrt[4]{\frac{L_x^2 L_y}{\mu_x^2 \mu_y}}, \sqrt[4]{\frac{L_x L_y^2}{\mu_x \mu_y^2}}, \frac{L_{xy}}{\sqrt{\mu_x \mu_y}}\right\} \log \frac{1}{\epsilon}\right) \quad (26)$$

due to Xie et al. (2021), and

$$\tilde{\mathcal{O}}\left(\max\left\{\sqrt{\frac{L_x}{\mu_x}}, \sqrt{\frac{L_y}{\mu_y}}, \sqrt{\frac{L_{xy} L}{\mu_x \mu_y}}\right\} \log \frac{1}{\epsilon}\right), \quad (27)$$

due to Wang and Li (2020), where  $\tilde{\mathcal{O}}(\cdot)$  hides additional logarithmic factors, and  $L = \max\{L_x, L_y, L_{xy}\}$ . Moreover, our result (25) matches the lower complexity bound provided by Zhang et al. (2021b). Hence, *Algorithm 1 is optimal in this regime.*

Apart from our work, algorithms that achieve optimal complexity (25) were developed in three independent works by Thekumparampil et al. (2022); Jin et al. (2022); Du et al. (2022). However, to the best of our knowledge these works were published or appeared on arXiv in 2022, while our work appeared on arXiv in 2021. Hence, Algorithm 1 is the first algorithm which achieves the lower complexity bound (25) for smooth and strongly-convex-strongly-concave saddle-point problems with bilinear coupling.

### 5.2 Affinely-constrained minimization case

In this case  $\mu_x > 0$  and  $\mu_y = 0$ . Firstly, we consider the case when  $L_y = 0$ , i.e.,  $g(y)$  is a linear function. Then, problem (1) is equivalent to the smooth and strongly-convex affinely-constrained minimization problem (7). Algorithm 1 enjoys the linear convergence rate

$$\mathcal{O}\left(\max\left\{\frac{L_{xy}}{\mu_{xy}} \sqrt{\frac{L_x}{\mu_x}}, \frac{L_{xy}^2}{\mu_{xy}^2}\right\} \log \frac{1}{\epsilon}\right), \quad (28)$$



where  $\mu_{xy} = \lambda_{\min}^+(\mathbf{A}\mathbf{A}^\top) > 0$  due to Assumption 3. This result recovers the complexity of the APAPC algorithm (Kovalev et al., 2020). It is possible to incorporate the Chebyshev acceleration mechanism (Arioli and Scott, 2014) into Algorithm 1 for solving problem (7) to obtain the improved complexity

$$\mathcal{O}\left(\frac{L_{xy}}{\mu_{xy}}\sqrt{\frac{L_x}{\mu_x}}\log\frac{1}{\epsilon}\right). \quad (29)$$

This matches the complexity of the OPAPC algorithm of Kovalev et al. (2020); Salim et al. (2021), which was shown to be optimal (Salim et al., 2021; Scaman et al., 2017).

### 5.3 Strongly convex and concave case

We also allow  $L_y > 0$ , i.e., function  $g(y)$  is a general, not necessarily linear, smooth and convex function. It is often possible that  $\mu_{xy} > 0$  due to Assumption 3; for instance, when  $\mathbf{A}$  is a full row rank matrix. Then, Algorithm 1 enjoys the following linear iteration complexity:

$$\mathcal{O}\left(\max\left\{\frac{\sqrt{L_x L_y}}{\mu_{xy}}, \frac{L_{xy}}{\mu_{xy}}\sqrt{\frac{L_x}{\mu_x}}, \frac{L_{xy}^2}{\mu_{xy}^2}\right\}\log\frac{1}{\epsilon}\right). \quad (30)$$

This case was previously studied by Du and Hu (2019); Du et al. (2017); Zhang et al. (2021a). Du and Hu (2019) provided an analysis for an algorithm called Sim-GDA, and established its iteration complexity

$$\mathcal{O}\left(\max\left\{\frac{L_x^3 L_y L_{xy}^2}{\mu_x^2 \mu_{xy}^4}, \frac{L_x^3 L_{xy}^4}{\mu_x^3 \mu_{xy}^4}\right\}\log\frac{1}{\epsilon}\right). \quad (31)$$

This result is substantially worse than our complexity (30); possibly due to a suboptimal analysis. Subsequently, Zhang et al. (2021a) provided an improved analysis for the Sim-GDA algorithm, obtaining the complexity

$$\mathcal{O}\left(\max\left\{\frac{L^3}{\mu_x \mu_{xy}^2}, \frac{L^2}{\mu_x^2}\right\}\log\frac{1}{\epsilon}\right). \quad (32)$$

They also studied the Alt-GDA method, obtaining the complexity

$$\mathcal{O}\left(\max\left\{\frac{L^2}{\mu_{xy}^2}, \frac{L}{\mu_x}\right\}\log\frac{1}{\epsilon}\right), \quad (33)$$

where  $L = \max\{L_x, L_y, L_{xy}\}$ . However, these results are local, i.e., they are valid only if the initial iterates of these algorithms are close enough to the solution of problem (1). Moreover, these results are still worse than our rate (30) because Sim-GDA and Alt-GDA do not utilize the Nesterov acceleration mechanism, while our Algorithm 1 does.

### 5.4 Bilinear case

In this case  $\mu_x = \mu_y = L_x = L_y = 0$ . That is, functions  $f(x)$  and  $g(y)$  are linear. Then, problem (1) turns into the bilinear min-max problem (10), and  $\mu_{xy}^2 = \mu_{yx}^2 = \lambda_{\min}^+(\mathbf{A}^\top \mathbf{A}) > 0$  due to Assumption 3. The iteration complexity of Algorithm 1 becomes

$$\mathcal{O}\left(\frac{L_{xy}^2}{\mu_{xy}^2}\log\frac{1}{\epsilon}\right). \quad (34)$$

This recovers the results of Daskalakis et al. (2018); Liang and Stokes (2019); Gidel et al. (2018, 2019); Mishchenko et al. (2020); Mokhtari et al. (2020) for the bilinear min-max problem (10). However, this result is worse than the complexity lower bound

$$\mathcal{O}\left(\frac{L_{xy}}{\mu_{xy}}\log\frac{1}{\epsilon}\right), \quad (35)$$

obtained in the work of Ibrahim et al. (2020), which was reached by Azizian et al. (2020); Du et al. (2022)<sup>7</sup>.

<sup>7</sup>We provide these results for completeness. The result of Azizian et al. (2020) is better than our result (34) for Algorithm 1 because they specifically focus on solving the bilinear min-max problem (10), while Algorithm 1 aims to solve the much more general convex-concave saddle-point problem (1).

## 5.5 Convex-concave case

In this case  $\mu_y = \mu_x = 0$ . It is often possible that  $\mu_{xy} = \mu_{yx} > 0$  due to Assumption 3, for example, when  $\mathbf{A}$  is a square and full rank matrix. Then, the iteration complexity of Algorithm 1 becomes

$$\mathcal{O} \left( \max \left\{ \frac{\sqrt{L_x L_y L_{xy}}}{\mu_{xy}^2}, \frac{L_{xy}^2}{\mu_{xy}^2} \right\} \log \frac{1}{\epsilon} \right), \quad (36)$$

which is still linear. This complexity result generalizes the result (34) for bilinear min-max problems as it allows for general, not necessarily linear, convex and smooth functions  $f(x)$  and  $g(x)$ . To the best of our knowledge, Algorithm 1 is the first algorithm which can achieve linear convergence for smooth and non-strongly convex non-strongly concave min-max problems with bilinear coupling.

## Acknowledgements

The work of Alexander Gasnikov was supported by a grant for research centers in the field of artificial intelligence, provided by the Analytical Center for the Government of the Russian Federation in accordance with the subsidy agreement (agreement identifier 000000D730321P5Q0002) and the agreement with the Ivannikov Institute for System Programming of the Russian Academy of Sciences dated November 2, 2021 No. 70-2021-00142.

## References

- Alkousa, M., Gasnikov, A., Dvinskikh, D., Kovalev, D., and Stonyakin, F. (2020). Accelerated methods for saddle-point problem. *Computational Mathematics and Mathematical Physics*, 60(11):1787–1809.
- Arioli, M. and Scott, J. (2014). Chebyshev acceleration of iterative refinement. *Numerical Algorithms*, 66(3):591–608.
- Arjevani, Y., Bruna, J., Can, B., Gürbüzbalaban, M., Jegelka, S., and Lin, H. (2020). Ideal: Inexact decentralized accelerated augmented lagrangian method. *arXiv preprint arXiv:2006.06733*.
- Azizian, W., Scieur, D., Mitliagkas, I., Lacoste-Julien, S., and Gidel, G. (2020). Accelerating smooth games by manipulating spectral shapes. In *International Conference on Artificial Intelligence and Statistics*, pages 1705–1715. PMLR.
- Bauschke, H. H. and Combettes, P. L. (2011). *Convex Analysis and Monotone Operator Theory in Hilbert Spaces*. Springer.
- Chambolle, A. and Pock, T. (2011). A first-order primal-dual algorithm for convex problems with applications to imaging. *Journal of mathematical imaging and vision*, 40(1):120–145.
- Chambolle, A. and Pock, T. (2016). An introduction to continuous optimization for imaging. *Acta Numerica*, 25:161–319.
- Cohen, M. B., Sidford, A., and Tian, K. (2020). Relative lipschitzness in extragradient methods and a direct recipe for acceleration. *arXiv preprint arXiv:2011.06572*.
- Daskalakis, C., Ilyas, A., Syrgkanis, V., and Zeng, H. (2018). Training gans with optimism. In *International Conference on Learning Representations (ICLR 2018)*.
- Du, S. S., Chen, J., Li, L., Xiao, L., and Zhou, D. (2017). Stochastic variance reduction methods for policy evaluation. In *International Conference on Machine Learning*, pages 1049–1058. PMLR.
- Du, S. S., Gidel, G., Jordan, M. I., and Li, C. J. (2022). Optimal extragradient-based bilinearly-coupled saddle-point optimization. *arXiv preprint arXiv:2206.08573*.
- Du, S. S. and Hu, W. (2019). Linear convergence of the primal-dual gradient method for convex-concave saddle point problems without strong convexity. In *The 22nd International Conference on Artificial Intelligence and Statistics*, pages 196–205. PMLR.

- Gidel, G., Berard, H., Vignoud, G., Vincent, P., and Lacoste-Julien, S. (2018). A variational inequality perspective on generative adversarial networks. *arXiv preprint arXiv:1802.10551*.
- Gidel, G., Hemmat, R. A., Pezeshki, M., Le Priol, R., Huang, G., Lacoste-Julien, S., and Mitliagkas, I. (2019). Negative momentum for improved game dynamics. In *The 22nd International Conference on Artificial Intelligence and Statistics*, pages 1802–1811. PMLR.
- Ibrahim, A., Azizian, W., Gidel, G., and Mitliagkas, I. (2020). Linear lower bounds and conditioning of differentiable games. In *International Conference on Machine Learning*, pages 4583–4593. PMLR.
- Jin, Y., Sidford, A., and Tian, K. (2022). Sharper rates for separable minimax and finite sum optimization via primal-dual extragradient methods. *arXiv preprint arXiv:2202.04640*.
- Keriven, N., Bourrier, A., Gribonval, R., and Pérez, P. (2018). Sketching for large-scale learning of mixture models. *Information and Inference: A Journal of the IMA*, 7(3):447–508.
- Korpelevich, G. M. (1976). The extragradient method for finding saddle points and other problems. *Matecon*, 12:747–756.
- Kovalev, D., Salim, A., and Richtárik, P. (2020). Optimal and practical algorithms for smooth and strongly convex decentralized optimization. *Advances in Neural Information Processing Systems*, 33.
- Lei, Q., Yen, I. E.-H., Wu, C.-y., Dhillon, I. S., and Ravikumar, P. (2017). Doubly greedy primal-dual coordinate descent for sparse empirical risk minimization. In *International Conference on Machine Learning*, pages 2034–2042. PMLR.
- Li, H., Lin, Z., and Fang, Y. (2020). Optimal accelerated variance reduced extra and diging for strongly convex and smooth decentralized optimization. *arXiv e-prints*, pages arXiv–2009.
- Liang, T. and Stokes, J. (2019). Interaction matters: A note on non-asymptotic local convergence of generative adversarial networks. In *The 22nd International Conference on Artificial Intelligence and Statistics*, pages 907–915. PMLR.
- Lin, T., Jin, C., and Jordan, M. I. (2020). Near-optimal algorithms for minimax optimization. In *Conference on Learning Theory*, pages 2738–2779. PMLR.
- Mishchenko, K., Kovalev, D., Shulgin, E., Richtárik, P., and Malitsky, Y. (2020). Revisiting stochastic extragradient. In *International Conference on Artificial Intelligence and Statistics*, pages 4573–4582. PMLR.
- Mokhtari, A., Ozdaglar, A., and Pattathil, S. (2020). 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*, pages 1497–1507. PMLR.
- Nedic, A., Olshevsky, A., and Shi, W. (2017). Achieving geometric convergence for distributed optimization over time-varying graphs. *SIAM Journal on Optimization*, 27(4):2597–2633.
- Nesterov, Y. (2003). *Introductory lectures on convex optimization: A basic course*, volume 87. Springer Science & Business Media.
- Nesterov, Y. and Scramali, L. (2006). Solving strongly monotone variational and quasi-variational inequalities.
- Nesterov, Y. E. (1983). A method for solving the convex programming problem with convergence rate  $o(1/k^2)$ . In *Dokl. akad. nauk Sssr*, volume 269, pages 543–547.
- Peyré, G., Cuturi, M., et al. (2019). Computational optimal transport: With applications to data science. *Foundations and Trends® in Machine Learning*, 11(5-6):355–607.
- Salim, A., Condat, L., Kovalev, D., and Richtárik, P. (2021). An optimal algorithm for strongly convex minimization under affine constraints. *arXiv preprint arXiv:2102.11079*.

- Scaman, K., Bach, F., Bubeck, S., Lee, Y. T., and Massoulié, L. (2017). Optimal algorithms for smooth and strongly convex distributed optimization in networks. In *international conference on machine learning*, pages 3027–3036. PMLR.
- Thekumparampil, K. K., He, N., and Oh, S. (2022). Lifted primal-dual method for bilinearly coupled smooth minimax optimization. In *International Conference on Artificial Intelligence and Statistics*, pages 4281–4308. PMLR.
- Wang, J. and Xiao, L. (2017). Exploiting strong convexity from data with primal-dual first-order algorithms. In *International Conference on Machine Learning*, pages 3694–3702. PMLR.
- Wang, Y. and Li, J. (2020). Improved algorithms for convex-concave minimax optimization. *arXiv preprint arXiv:2006.06359*.
- Xiao, L., Yu, A. W., Lin, Q., and Chen, W. (2019). Dscovr: Randomized primal-dual block coordinate algorithms for asynchronous distributed optimization. *The Journal of Machine Learning Research*, 20(1):1634–1691.
- Xie, G., Han, Y., and Zhang, Z. (2021). Dippa: An improved method for bilinear saddle point problems. *arXiv preprint arXiv:2103.08270*.
- Ye, H., Luo, L., Zhou, Z., and Zhang, T. (2020). Multi-consensus decentralized accelerated gradient descent. *arXiv preprint arXiv:2005.00797*.
- Zargham, M., Ribeiro, A., Ozdaglar, A., and Jadbabaie, A. (2013). Accelerated dual descent for network flow optimization. *IEEE Transactions on Automatic Control*, 59(4):905–920.
- Zhang, G., Wang, Y., Lessard, L., and Grosse, R. (2021a). Don’t fix what ain’t broke: Near-optimal local convergence of alternating gradient descent-ascent for minimax optimization. *arXiv preprint arXiv:2102.09468*.
- Zhang, J., Hong, M., and Zhang, S. (2021b). On lower iteration complexity bounds for the convex concave saddle point problems. *Mathematical Programming*, pages 1–35.
- Zhang, Y. and Lin, X. (2015). Stochastic primal-dual coordinate method for regularized empirical risk minimization. In *International Conference on Machine Learning*, pages 353–361. PMLR.

## Checklist

1. For all authors...
  - (a) Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope? [Yes]
  - (b) Did you describe the limitations of your work? [No]
  - (c) Did you discuss any potential negative societal impacts of your work? [No] This is a theoretical work with no foreseeable negative societal impact.
  - (d) Have you read the ethics review guidelines and ensured that your paper conforms to them? [Yes]
2. If you are including theoretical results...
  - (a) Did you state the full set of assumptions of all theoretical results? [Yes] Assumptions 1 to 3.
  - (b) Did you include complete proofs of all theoretical results? [Yes] see Appendix.
3. If you ran experiments...
  - (a) Did you include the code, data, and instructions needed to reproduce the main experimental results (either in the supplemental material or as a URL)? [N/A]
  - (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were chosen)? [N/A]
  - (c) Did you report error bars (e.g., with respect to the random seed after running experiments multiple times)? [N/A]
  - (d) Did you include the total amount of compute and the type of resources used (e.g., type of GPUs, internal cluster, or cloud provider)? [N/A]
4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
  - (a) If your work uses existing assets, did you cite the creators? [N/A]
  - (b) Did you mention the license of the assets? [N/A]
  - (c) Did you include any new assets either in the supplemental material or as a URL? [N/A]
  - (d) Did you discuss whether and how consent was obtained from people whose data you're using/curating? [N/A]
  - (e) Did you discuss whether the data you are using/curating contains personally identifiable information or offensive content? [N/A]
5. If you used crowdsourcing or conducted research with human subjects...
  - (a) Did you include the full text of instructions given to participants and screenshots, if applicable? [N/A]
  - (b) Did you describe any potential participant risks, with links to Institutional Review Board (IRB) approvals, if applicable? [N/A]
  - (c) Did you include the estimated hourly wage paid to participants and the total amount spent on participant compensation? [N/A]

# Appendix

In Appendix A we present the new GDAE algorithm, in Appendix B we provide a proof of Lemma 1, in Appendix C we provide a proof of Theorem 1, and in Appendix D we provide a proof of Theorem 2.

## A A Novel Gradient Method for General Convex-Concave Saddle-Point Problems

---

### Algorithm 2 GDAE: Gradient Descent-Ascent with Extrapolation

---

**Input:**  $x^0 \in \mathbb{R}^{d_x}$ ,  $y^0 \in \mathbb{R}^{d_y}$ ,  $\eta_x, \eta_y > 0$ ,  $\theta \in (0, 1)$   
 $x^{-1} = x^0$   
 $y^{-1} = y^0$   
**for**  $k = 0, 1, 2, \dots$  **do**  
 $x^{k+1} = x^k - \eta_x \nabla_x F(x^k, y^k) - \eta_x \theta (\nabla_x F(x^{k-1}, y^k) - \nabla_x F(x^{k-1}, y^{k-1}))$   
 $y^{k+1} = y^k + \eta_y \nabla_y F(x^{k+1}, y^k)$   
**end for**

---

In this section we present a new method—Gradient Descent-Ascent Method with Extrapolation (GDAE; Algorithm 2)—for solving problem (11).

### A.1 Assumptions and definitions

First, we state the main assumptions that we impose on problem (11).

**Assumption 4.** *Function  $F(x, y)$  is  $L_x$ -smooth and  $\mu_x$ -strongly convex in  $x$  and  $L_y$ -smooth and  $\mu_y$ -strongly concave in  $y$ , where  $L_x \geq \mu_x \geq 0$ ,  $L_y \geq \mu_y \geq 0$ .*

Assumption 4 generalizes the smoothness and strong convexity Assumptions 1 and 2 imposed on problem (1).

**Assumption 5.** *There exists a constant  $L_{xy} > 0$  such that for all  $x, x_1, x_2 \in \mathbb{R}^{d_x}$  and  $y, y_1, y_2 \in \mathbb{R}^{d_y}$ , the following inequalities hold:*

$$\begin{aligned} \|\nabla_x F(x, y_1) - \nabla_x F(x, y_2)\| &\leq L_{xy} \|y_1 - y_2\|, \\ \|\nabla_y F(x_1, y) - \nabla_y F(x_2, y)\| &\leq L_{xy} \|x_1 - x_2\|. \end{aligned} \quad (37)$$

**Assumption 6.** *There exist constants  $\mu_{xy}, \mu_{yx} \geq 0$  such that for all  $x, x_1, x_2 \in \mathbb{R}^{d_x}$  and  $y, y_1, y_2 \in \mathbb{R}^{d_y}$ , the following inequalities hold:*

$$\begin{aligned} \|\nabla_x F(x, y_1) - \nabla_x F(x, y_2)\| &\geq \mu_{xy} \|y_1 - y_2\|, \\ \|\nabla_y F(x_1, y) - \nabla_y F(x_2, y)\| &\geq \mu_{yx} \|x_1 - x_2\|. \end{aligned} \quad (38)$$

Assumptions 5 and 6 combined form a generalized version of Assumption 3 for problem (11). Indeed, if one assumes that (37) and (38) hold for problem (1), then the following inequalities hold

$$\begin{aligned} \mu_{xy}^2 &\leq \lambda_{\min}(\mathbf{A}\mathbf{A}^\top) \leq L_{xy}^2, \\ \mu_{yx}^2 &\leq \lambda_{\min}(\mathbf{A}^\top \mathbf{A}) \leq L_{xy}^2, \end{aligned} \quad (39)$$

which can be seen as a simplified version of Assumption 3.

Next, we recall several basic definitions. Similarly to section 3, by  $\mathcal{S} \subset \mathbb{R}^{d_x} \times \mathbb{R}^{d_y}$  we denote the solution set of problem (11). Note that  $(x^*, y^*) \in \mathcal{S}$  if and only if  $(x^*, y^*)$  satisfies the optimality conditions

$$\begin{cases} \nabla_x F(x^*, y^*) = 0, \\ \nabla_y F(x^*, y^*) = 0. \end{cases} \quad (40)$$

We also use notions of iteration complexity for achieving an  $\epsilon$ -accurate solution analogous to Definitions 2 and 3, respectively.

## A.2 Algorithm development

We now present the main ingredients and intuition behind the development of our method.

**Implicit gradient descent-ascent.** First, we recall the iterations of the Forward-Backward algorithm (17), which can be written in the form

$$\begin{cases} x^+ = x - \eta_x \nabla f(x) - \eta_x \mathbf{A}^\top y^+ \\ y^+ = y - \eta_y \nabla g(y) + \eta_y \mathbf{A} x^+ \end{cases}, \quad (41)$$

where  $\eta_x, \eta_y > 0$  are stepsizes. Iterations (41) can also be written in terms of the gradients  $\nabla_x F(x, y)$  and  $\nabla_y F(x, y)$ ,

$$\begin{cases} x^+ = x - \eta_x \nabla_x F(x, y^+) \\ y^+ = y + \eta_y \nabla_y F(x^+, y) \end{cases}, \quad (42)$$

which makes the method applicable to the general problem (11).

Iterations (42) were the foundation for the development of Algorithm 1 for solving problem (1), with strong convergence properties established by Theorem 1. Hence, we expect that this approach would work for solving the more general problem (11). However, (42) is an implicit algorithm and can't be applied in its current state.

**Gradient extrapolation.** In analogy to the development of Algorithm 1, we want to find a good approximation of the implicit iterations (42). A naive solution would be using the approximation

$$\begin{cases} x^+ = x - \eta_x \nabla_x F(x, y_m) \\ y^+ = y + \eta_y \nabla_y F(x^+, y) \end{cases}, \quad (43)$$

where  $y_m \approx y^+$ . Similarly to Section 4.1, we could use  $y_m = y$ , which would lead to the Alt-GDA algorithm (Zhang et al., 2021a), or  $y_m = y + \theta(y - y^-)$ , which is a linear extrapolation step (Chambolle and Pock, 2011).

The linear extrapolation step with  $\theta = 1$  is based on the ‘‘assumption’’ that  $y^+ \approx y_m = y + (y - y^-)$ , or equivalently,  $y^+ - y \approx y - y^-$ . We can use a similar intuition for the gradients  $\nabla_x F(x, y)$  rather than the iterates  $y$ . In particular, we ‘‘assume’’ that

$$\nabla_x F(x, y^+) - \nabla_x F(x, y) \approx \nabla_x F(x^-, y) - \nabla_x F(x^-, y^-),$$

or equivalently,

$$\begin{cases} \nabla_x F(x, y^+) \approx \Delta_x \\ \Delta_x = \nabla_x F(x, y) + (\nabla_x F(x^-, y) - \nabla_x F(x^-, y^-)) \end{cases}.$$

This intuition leads to the following novel update rule, which we call *gradient extrapolation step*:

$$\begin{cases} \Delta_x = \nabla_x F(x, y) + \theta(\nabla_x F(x^-, y) - \nabla_x F(x^-, y^-)) \\ x^+ = x - \eta_x \Delta_x \end{cases}.$$

Above,  $\theta \in (0, 1]$  is the extrapolation parameter. We use this gradient extrapolation step together with the update rule for  $y$  from (42) in the design of our Algorithm 2.

## A.3 Convergence of Algorithm 2 and related work

We now present Theorem 2, which establishes linear convergence rate for Algorithm 2 under Assumptions 4, 5 and 6.

**Theorem 2.** *Let Assumptions 4, 5 and 6 and condition (21) hold. Then there exist parameters of Algorithm 2 such that the iteration complexity for finding an  $\epsilon$ -accurate solution of problem (11) is*

$$\mathcal{O}\left(\min\{T_a, T_b, T_c, T_d\} \log \frac{C}{\epsilon}\right), \quad (44)$$

where  $T_a, T_b, T_c, T_d$  are defined as

$$\begin{aligned} T_a &= \max \left\{ \frac{L_x}{\mu_x}, \frac{L_y}{\mu_y}, \frac{L_{xy}}{\sqrt{\mu_x \mu_y}} \right\}, & T_b &= \max \left\{ \frac{L_x}{\mu_x}, \frac{L_x L_y}{\mu_{xy}^2}, \frac{L_{xy}^2}{\mu_{xy}^2} \right\}, \\ T_c &= \max \left\{ \frac{L_y}{\mu_y}, \frac{L_x L_y}{\mu_{yx}^2}, \frac{L_{xy}^2}{\mu_{yx}^2} \right\}, & T_d &= \max \left\{ \frac{L_x L_y}{\mu_{xy}^2}, \frac{L_x L_y}{\mu_{yx}^2}, \frac{L_{xy}^2}{\mu_{xy}^2}, \frac{L_{xy}^2}{\mu_{yx}^2} \right\}, \end{aligned}$$

and  $C > 0$  is some constant, which does not depend on  $\epsilon$ , but possibly depends on  $L_x, \mu_x, L_y, \mu_y, L_{xy}, \mu_{xy}, \mu_{yx}$ .

Consider the case when  $\mu_x, \mu_y > 0$ . In this case the iteration complexity of Algorithm 2 becomes

$$\mathcal{O} \left( \max \left\{ \frac{L_x}{\mu_x}, \frac{L_y}{\mu_y}, \frac{L_{xy}}{\sqrt{\mu_x \mu_y}} \right\} \log \frac{1}{\epsilon} \right). \quad (45)$$

This recovers the result of Cohen et al. (2020). Moreover, when  $\mu_x = \mu_y$ , this result recovers the complexity of solving problem (11) by a number of known algorithms, including the extragradient method (Korpelevich, 1976), optimistic gradient method (Daskalakis et al., 2018; Gidel et al., 2018), and the dual extrapolation method (Nesterov and Scramali, 2006).

Finally, consider then opposite case when at least one of the constants  $\mu_x$  and  $\mu_y$  is zero. To the best of our knowledge, there are no algorithms that can achieve a linear convergence. However, Algorithm 2 can still achieve linear iteration complexity provided that condition (21) is satisfied.

## B Proof of Lemma 1

**Part 1.** Let us first show that primal problem (22) has at least a single solution  $x^* \in \mathbb{R}^{d_x}$ .

Condition (21) implies that  $\max\{\mu_x, \mu_{yx}\} > 0$ . If  $\mu_x > 0$  then function  $P(x)$  is obviously strongly convex and primal problem indeed has a solution. Consider the opposite case  $\mu_x = 0$ . Then  $\mu_{yx} > 0$  due to condition (21).

Assumption 3 and  $\mu_{yx} > 0$  imply that  $\nabla f(x) \in \text{range} \mathbf{A}^\top$  for all  $x \in \mathbb{R}^{d_x}$ . Hence,

$$f(x + x') = f(x) \text{ for all } x \in \mathbb{R}^{d_x}, x' \in \ker \mathbf{A}. \quad (46)$$

Using the definition of  $P(x)$  we get

$$\begin{aligned} P(x + x') &= f(x + x') + g^*(\mathbf{A}(x + x')) \\ &= f(x) + g^*(\mathbf{A}x) \\ &= P(x) \end{aligned}$$

for all  $x \in \mathbb{R}^{d_x}, x' \in \ker \mathbf{A}$ . From this one can conclude that

$$\min_{x \in \mathbb{R}^{d_x}} P(x) = \min_{x \in x^0 + \text{range} \mathbf{A}^\top} P(x).$$

for any vector  $x^0 \in \mathbb{R}^{d_x}$ . From the definition of  $P(x)$  it follows that  $P(x)$  is  $\mu_{yx}$ -strongly convex on any affine space  $x^0 + \text{range} \mathbf{A}^\top$  for arbitrary  $x^0 \in \mathbb{R}^{d_x}$ . Hence, problem  $\min_{x \in x^0 + \text{range} \mathbf{A}^\top} P(x)$  has a unique solution and primal problem  $\min_{x \in \mathbb{R}^{d_x}} P(x)$  has at least a single solution  $x^*$ .

**Part 2.** Let us show that there exists  $y^* \in \mathbb{R}^{d_y}$  such that  $(x^*, y^*) \in \mathcal{S}$ , i.e.,  $(x^*, y^*)$  satisfy optimality conditions (13).

Let us show that  $-\nabla f(x^*) \in \mathbf{A}^\top \partial g^*(\mathbf{A}x^*)$ . We use condition (21) which implies  $\max\{\mu_y, \mu_{xy}\} > 0$ . If  $\mu_y > 0$ , then function  $g^*(y)$  is smooth and our statement is trivial. Consider the opposite case  $\mu_y = 0$ . Then  $\mu_{xy} > 0$  due to condition (21).

Assumption 3 and  $\mu_{xy} > 0$  imply that  $\nabla g(y) \in \text{range} \mathbf{A}$  for all  $y \in \mathbb{R}^{d_y}$ . Hence,  $\text{dom} g^*(\cdot) \subset \text{range} \mathbf{A}$ . Let  $h(x) = g^*(\mathbf{A}x)$ . From standard theory it follows that  $-\nabla f(x^*) \in \partial h(x^*)$  or

$$h(x) \geq h(x^*) - \langle \nabla f(x^*), x - x^* \rangle \text{ for all } x \in \mathbb{R}^{d_x},$$



From this one can conclude that

$$\langle \nabla f(x^*), x - x^* \rangle \geq 0 \text{ for all } x \in x^* + \ker \mathbf{A},$$

which implies  $\nabla f(x^*) \in (\ker \mathbf{A})^\perp = \text{range} \mathbf{A}^\top$ . Hence, there exists vector  $y^* \in \mathbb{R}^{d_y}$  such that  $-\nabla f(x^*) = \mathbf{A}^\top y^*$ . Now, we can write

$$h(x) \geq h(x^*) + \langle \mathbf{A}^\top y^*, x - x^* \rangle \text{ for all } x \in \mathbb{R}^{d_x},$$

which is equivalent to

$$g^*(\mathbf{A}x) \geq g^*(\mathbf{A}x^*) + \langle y^*, \mathbf{A}x - \mathbf{A}x^* \rangle \text{ for all } x \in \mathbb{R}^{d_x}.$$

The latter can be written as

$$g^*(y) \geq g^*(\mathbf{A}x^*) + \langle y^*, y - \mathbf{A}x^* \rangle \text{ for all } y \in \text{range} \mathbf{A}.$$

But  $\text{dom } g^*(\cdot) \subset \text{range} \mathbf{A}$ , which means that  $g^*(y) = +\infty$  for all  $y \notin \text{range} \mathbf{A}$ . This implies

$$g^*(y) \geq g^*(\mathbf{A}x^*) + \langle y^*, y - \mathbf{A}x^* \rangle \text{ for all } y \in \mathbb{R}^{d_y},$$

which is a definition of  $y^* \in \partial g^*(\mathbf{A}x^*)$ . An equivalent for this is  $\nabla g(y^*) = \mathbf{A}x^*$ , which together with  $-\nabla f(x^*) = \mathbf{A}^\top y^*$  form optimality condition (13).

**Part 3.** We showed that there exists a pair of vectors  $(x^*, y^*) \in \mathbb{R}^{d_x} \times \mathbb{R}^{d_y}$  which is a saddle point of the function  $F(x, y)$  in problem (1). Hence, strong duality holds and proof of the rest of Lemma 1 is trivial.  $\square$

## C Proof of Theorem 1

**Lemma 2.** *There exists a solution  $(x^*, y^*) \in \mathcal{S}$  of the problem (1) such that for all  $k = 0, 1, 2, \dots$  the iterates of Algorithm 1 satisfy*

$$\begin{aligned}\|\mathbf{A}(x^k - x^*)\| &\geq \mu_{yx}\|x^k - x^*\|, \\ \|\mathbf{A}^\top(y^k - y^*)\| &\geq \mu_{xy}\|y^k - y^*\|.\end{aligned}\tag{47}$$

*Proof.* The proof of this lemma is a trivial extension of the derivations from the proof of Lemma 1.  $\square$

**Lemma 3.** *Let  $\tau_x$  be defined as*

$$\tau_x = (\sigma_x^{-1} + 1/2)^{-1}.\tag{48}$$

*Let  $\alpha_x$  be defined as*

$$\alpha_x = \mu_x.\tag{49}$$

*Let  $\beta_x$  be defined as*

$$\beta_x = \min \left\{ \frac{1}{2L_y}, \frac{1}{2\eta_x L_{xy}^2} \right\}.\tag{50}$$

*Then, the following inequality holds:*

$$\begin{aligned}\frac{1}{\eta_x}\|x^{k+1} - x^*\|^2 &\leq \left( \frac{1}{\eta_x} - \mu_x - \beta_x \mu_{yx}^2 \right) \|x^k - x^*\|^2 + \left( \mu_x + L_x \sigma_x - \frac{1}{2\eta_x} \right) \|x^{k+1} - x^k\|^2 \\ &\quad + D_g(y_g^k, y^*) - D_f(x_g^k, x^*) - \frac{2}{\sigma_x} D_f(x_f^{k+1}, x^*) + \left( \frac{2}{\sigma_x} - 1 \right) D_f(x_f^k, x^*) \\ &\quad - 2\langle \mathbf{A}^\top(y_m^k - y^*), x^{k+1} - x^* \rangle.\end{aligned}\tag{51}$$

*Proof.* Using Line 8 of the Algorithm 1 we get

$$\begin{aligned}\frac{1}{\eta_x}\|x^{k+1} - x^*\|^2 &= \frac{1}{\eta_x}\|x^k - x^*\|^2 + \frac{2}{\eta_x}\langle x^{k+1} - x^k, x^{k+1} - x^* \rangle - \frac{1}{\eta_x}\|x^{k+1} - x^k\|^2 \\ &= \frac{1}{\eta_x}\|x^k - x^*\|^2 + 2\alpha_x \langle x_g^k - x^k, x^{k+1} - x^* \rangle - 2\beta_x \langle \mathbf{A}^\top(\mathbf{A}x^k - \nabla g(y_g^k)), x^{k+1} - x^* \rangle \\ &\quad - 2\langle \nabla f(x_g^k) + \mathbf{A}^\top y_m^k, x^{k+1} - x^* \rangle - \frac{1}{\eta_x}\|x^{k+1} - x^k\|^2.\end{aligned}$$

Using the parallelogram rule we get

$$\begin{aligned}\frac{1}{\eta_x}\|x^{k+1} - x^*\|^2 &= \frac{1}{\eta_x}\|x^k - x^*\|^2 + \alpha_x (\|x_g^k - x^*\|^2 - \|x_g^k - x^{k+1}\|^2 - \|x^k - x^*\|^2 + \|x^{k+1} - x^k\|^2) \\ &\quad - 2\beta_x \langle \mathbf{A}x^k - \nabla g(y_g^k), \mathbf{A}(x^{k+1} - x^*) \rangle - 2\langle \nabla f(x_g^k) + \mathbf{A}^\top y_m^k, x^{k+1} - x^* \rangle - \frac{1}{\eta_x}\|x^{k+1} - x^k\|^2.\end{aligned}$$

Using the optimality condition  $\nabla g(y^*) = \mathbf{A}x^*$ , which follows from (13), and the parallelogram rule we get

$$\begin{aligned}\frac{1}{\eta_x}\|x^{k+1} - x^*\|^2 &= \frac{1}{\eta_x}\|x^k - x^*\|^2 + \alpha_x (\|x_g^k - x^*\|^2 - \|x_g^k - x^{k+1}\|^2 - \|x^k - x^*\|^2 + \|x^{k+1} - x^k\|^2) \\ &\quad + \beta_x (\|\mathbf{A}(x^{k+1} - x^k)\|^2 - \|\mathbf{A}(x^k - x^*)\|^2 + \|\nabla g(y_g^k) - \nabla g(y^*)\|^2 - \|\nabla g(y_g^k) - \mathbf{A}(x^{k+1})\|^2) \\ &\quad - 2\langle \nabla f(x_g^k) + \mathbf{A}^\top y_m^k, x^{k+1} - x^* \rangle - \frac{1}{\eta_x}\|x^{k+1} - x^k\|^2.\end{aligned}$$

Using Assumption 3, equation 47 and  $L_y$ -smoothness of  $g$  we get

$$\begin{aligned}
\frac{1}{\eta_x} \|x^{k+1} - x^*\|^2 &\leq \frac{1}{\eta_x} \|x^k - x^*\|^2 + \alpha_x \|x_g^k - x^*\|^2 - \alpha_x \|x^k - x^*\|^2 + \alpha_x \|x^{k+1} - x^k\|^2 \\
&\quad + \beta_x L_{xy}^2 \|x^{k+1} - x^k\|^2 - \beta_x \mu_{yx}^2 \|x^k - x^*\|^2 + 2\beta_x L_y D_g(y_g^k, y^*) \\
&\quad - 2\langle \nabla f(x_g^k) + \mathbf{A}^\top y_m^k, x^{k+1} - x^* \rangle - \frac{1}{\eta_x} \|x^{k+1} - x^k\|^2 \\
&= \left( \frac{1}{\eta_x} - \alpha_x - \beta_x \mu_{yx}^2 \right) \|x^k - x^*\|^2 + \left( \beta_x L_{xy}^2 + \alpha_x - \frac{1}{\eta_x} \right) \|x^{k+1} - x^k\|^2 \\
&\quad + 2\beta_x L_y D_g(y_g^k, y^*) + \alpha_x \|x_g^k - x^*\|^2 - 2\langle \nabla f(x_g^k) + \mathbf{A}^\top y_m^k, x^{k+1} - x^* \rangle.
\end{aligned}$$

Using the optimality condition  $\nabla f(x^*) + \mathbf{A}^\top y^* = 0$ , which follows from (13), we get

$$\begin{aligned}
\frac{1}{\eta_x} \|x^{k+1} - x^*\|^2 &\leq \left( \frac{1}{\eta_x} - \alpha_x - \beta_x \mu_{yx}^2 \right) \|x^k - x^*\|^2 + \left( \beta_x L_{xy}^2 + \alpha_x - \frac{1}{\eta_x} \right) \|x^{k+1} - x^k\|^2 + 2\beta_x L_y D_g(y_g^k, y^*) \\
&\quad + \alpha_x \|x_g^k - x^*\|^2 - 2\langle \nabla f(x_g^k) - \nabla f(x^*), x^{k+1} - x^* \rangle - 2\langle \mathbf{A}^\top (y_m^k - y^*), x^{k+1} - x^* \rangle \\
&= \left( \frac{1}{\eta_x} - \alpha_x - \beta_x \mu_{yx}^2 \right) \|x^k - x^*\|^2 + \left( \beta_x L_{xy}^2 + \alpha_x - \frac{1}{\eta_x} \right) \|x^{k+1} - x^k\|^2 \\
&\quad + 2\beta_x L_y D_g(y_g^k, y^*) + \alpha_x \|x_g^k - x^*\|^2 - 2\langle \nabla f(x_g^k) - \nabla f(x^*), x^{k+1} - x^k + x^k - x_g^k + x_g^k - x^* \rangle \\
&\quad - 2\langle \mathbf{A}^\top (y_m^k - y^*), x^{k+1} - x^* \rangle
\end{aligned}$$

Using  $\mu_y$ -strong convexity of  $f$  and Lines 6 and 10 of the Algorithm 1 we get

$$\begin{aligned}
\frac{1}{\eta_x} \|x^{k+1} - x^*\|^2 &\leq \left( \frac{1}{\eta_x} - \alpha_x - \beta_x \mu_{yx}^2 \right) \|x^k - x^*\|^2 + \left( \beta_x L_{xy}^2 + \alpha_x - \frac{1}{\eta_x} \right) \|x^{k+1} - x^k\|^2 + 2\beta_x L_y D_g(y_g^k, y^*) \\
&\quad + \alpha_x \|x_g^k - x^*\|^2 - \frac{2}{\sigma_x} \langle \nabla f(x_g^k) - \nabla f(x^*), x_f^{k+1} - x_g^k \rangle + \frac{2(1-\tau_x)}{\tau_x} \langle \nabla f(x_g^k) - \nabla f(x^*), x_f^k - x_g^k \rangle \\
&\quad - 2D_f(x_g^k, x^*) - \mu_x \|x_g^k - x^*\|^2 - 2\langle \mathbf{A}^\top (y_m^k - y^*), x^{k+1} - x^* \rangle \\
&= \left( \frac{1}{\eta_x} - \alpha_x - \beta_x \mu_{yx}^2 \right) \|x^k - x^*\|^2 + \left( \beta_x L_{xy}^2 + \alpha_x - \frac{1}{\eta_x} \right) \|x^{k+1} - x^k\|^2 + (\alpha_x - \mu_x) \|x_g^k - x^*\|^2 \\
&\quad + 2\beta_x L_y D_g(y_g^k, y^*) - 2D_f(x_g^k, x^*) - \frac{2}{\sigma_x} \langle \nabla f(x_g^k) - \nabla f(x^*), x_f^{k+1} - x_g^k \rangle \\
&\quad + \frac{2(1-\tau_x)}{\tau_x} \langle \nabla f(x_g^k) - \nabla f(x^*), x_f^k - x_g^k \rangle - 2\langle \mathbf{A}^\top (y_m^k - y^*), x^{k+1} - x^* \rangle.
\end{aligned}$$

Using convexity of  $D_f(x, x^*)$  with respect to  $x$ , which follows from the convexity of  $f$ , we get

$$\begin{aligned}
\frac{1}{\eta_x} \|x^{k+1} - x^*\|^2 &\leq \left( \frac{1}{\eta_x} - \alpha_x - \beta_x \mu_{yx}^2 \right) \|x^k - x^*\|^2 + \left( \beta_x L_{xy}^2 + \alpha_x - \frac{1}{\eta_x} \right) \|x^{k+1} - x^k\|^2 + (\alpha_x - \mu_x) \|x_g^k - x^*\|^2 \\
&\quad + 2\beta_x L_y D_g(y_g^k, y^*) - 2D_f(x_g^k, x^*) - \frac{2}{\sigma_x} \langle \nabla f(x_g^k) - \nabla f(x^*), x_f^{k+1} - x_g^k \rangle \\
&\quad + \frac{2(1-\tau_x)}{\tau_x} (D_f(x_f^k, x^*) - D_f(x_g^k, x^*)) - 2\langle \mathbf{A}^\top (y_m^k - y^*), x^{k+1} - x^* \rangle.
\end{aligned}$$

Using  $L_x$ -smoothness of  $D_f(x, x^*)$  with respect to  $x$ , which follows from the  $L_x$ -smoothness of  $f$ , we get

$$\begin{aligned}
\frac{1}{\eta_x} \|x^{k+1} - x^*\|^2 &\leq \left( \frac{1}{\eta_x} - \alpha_x - \beta_x \mu_{yx}^2 \right) \|x^k - x^*\|^2 + \left( \beta_x L_{xy}^2 + \alpha_x - \frac{1}{\eta_x} \right) \|x^{k+1} - x^k\|^2 + (\alpha_x - \mu_x) \|x_g^k - x^*\|^2 \\
&\quad + 2\beta_x L_y D_g(y_g^k, y^*) - 2D_f(x_g^k, x^*) - \frac{2}{\sigma_x} \left( D_f(x_f^{k+1}, x^*) - D_f(x_g^k, x^*) - \frac{L_x}{2} \|x_f^{k+1} - x_g^k\|^2 \right) \\
&\quad + \frac{2(1-\tau_x)}{\tau_x} (D_f(x_f^k, x^*) - D_f(x_g^k, x^*)) - 2\langle \mathbf{A}^\top (y_m^k - y^*), x^{k+1} - x^* \rangle.
\end{aligned}$$

Using Line 10 of the Algorithm 1 we get

$$\begin{aligned}
\frac{1}{\eta_x} \|x^{k+1} - x^*\|^2 &\leq \left( \frac{1}{\eta_x} - \alpha_x - \beta_x \mu_{yx}^2 \right) \|x^k - x^*\|^2 + \left( \beta_x L_{xy}^2 + \alpha_x - \frac{1}{\eta_x} \right) \|x^{k+1} - x^k\|^2 + (\alpha_x - \mu_x) \|x_g^k - x^*\|^2 \\
&\quad + 2\beta_x L_y D_g(y_g^k, y^*) - 2D_f(x_g^k, x^*) - \frac{2}{\sigma_x} \left( D_f(x_f^{k+1}, x^*) - D_f(x_g^k, x^*) - \frac{L_x \sigma_x^2}{2} \|x^{k+1} - x^k\|^2 \right) \\
&\quad + \frac{2(1 - \tau_x)}{\tau_x} (D_f(x_f^k, x^*) - D_f(x_g^k, x^*)) - 2\langle \mathbf{A}^\top (y_m^k - y^*), x^{k+1} - x^* \rangle \\
&= \left( \frac{1}{\eta_x} - \alpha_x - \beta_x \mu_{yx}^2 \right) \|x^k - x^*\|^2 + \left( \beta_x L_{xy}^2 + \alpha_x + L_x \sigma_x - \frac{1}{\eta_x} \right) \|x^{k+1} - x^k\|^2 \\
&\quad + (\alpha_x - \mu_x) \|x_g^k - x^*\|^2 + 2\beta_x L_y D_g(y_g^k, y^*) + \left( \frac{2}{\sigma_x} - \frac{2}{\tau_x} \right) D_f(x_g^k, x^*) - \frac{2}{\sigma_x} D_f(x_f^{k+1}, x^*) \\
&\quad + \left( \frac{2}{\tau_x} - 2 \right) D_f(x_f^k, x^*) - 2\langle \mathbf{A}^\top (y_m^k - y^*), x^{k+1} - x^* \rangle.
\end{aligned}$$

Using the definition of  $\tau_x$ ,  $\alpha_x$  and  $\beta_x$  we get

$$\begin{aligned}
\frac{1}{\eta_x} \|x^{k+1} - x^*\|^2 &\leq \left( \frac{1}{\eta_x} - \mu_x - \beta_x \mu_{yx}^2 \right) \|x^k - x^*\|^2 + \left( \mu_x + L_x \sigma_x - \frac{1}{2\eta_x} \right) \|x^{k+1} - x^k\|^2 \\
&\quad + D_g(y_g^k, y^*) - D_f(x_g^k, x^*) - \frac{2}{\sigma_x} D_f(x_f^{k+1}, x^*) + \left( \frac{2}{\sigma_x} - 1 \right) D_f(x_f^k, x^*) \\
&\quad - 2\langle \mathbf{A}^\top (y_m^k - y^*), x^{k+1} - x^* \rangle.
\end{aligned}$$

□

**Lemma 4.** Let  $\tau_y$  be defined as

$$\tau_y = (\sigma_y^{-1} + 1/2)^{-1}. \quad (52)$$

Let  $\alpha_y$  be defined as

$$\alpha_y = \mu_y. \quad (53)$$

Let  $\beta_y$  be defined as

$$\beta_y = \min \left\{ \frac{1}{2L_x}, \frac{1}{2\eta_y L_{xy}^2} \right\}. \quad (54)$$

Then, the following inequality holds:

$$\begin{aligned}
\frac{1}{\eta_y} \|y^{k+1} - y^*\|^2 &\leq \left( \frac{1}{\eta_y} - \mu_y - \beta_y \mu_{xy}^2 \right) \|y^k - y^*\|^2 + \left( \mu_y + L_y \sigma_y - \frac{1}{2\eta_y} \right) \|y^{k+1} - y^k\|^2 \\
&\quad + D_f(x_g^k, x^*) - D_g(y_g^k, y^*) - \frac{2}{\sigma_y} D_g(y_f^{k+1}, y^*) + \left( \frac{2}{\sigma_y} - 1 \right) D_g(y_f^k, y^*) \\
&\quad + 2\langle \mathbf{A} (x^{k+1} - x^*), y^{k+1} - y^* \rangle.
\end{aligned} \quad (55)$$

*Proof.* The proof is similar to the proof of the previous lemma. □

**Lemma 5.** Let  $\eta_x$  be defined as

$$\eta_x = \min \left\{ \frac{1}{4(\mu_x + L_x \sigma_x)}, \frac{\delta}{4L_{xy}} \right\}, \quad (56)$$

and let  $\eta_y$  be defined as

$$\eta_y = \min \left\{ \frac{1}{4(\mu_y + L_y \sigma_y)}, \frac{1}{4L_{xy} \delta} \right\}, \quad (57)$$

where  $\delta > 0$  is a parameter. Let  $\theta$  be defined as

$$\theta = \theta(\delta, \sigma_x, \sigma_y) = 1 - \max \{ \rho_a(\delta, \sigma_x, \sigma_y), \rho_b(\delta, \sigma_x, \sigma_y), \rho_c(\delta, \sigma_x, \sigma_y), \rho_d(\delta, \sigma_x, \sigma_y) \}, \quad (58)$$

where  $\rho_a(\delta, \sigma_x, \sigma_y)$ ,  $\rho_b(\delta, \sigma_x, \sigma_y)$ ,  $\rho_c(\delta, \sigma_x, \sigma_y)$ ,  $\rho_d(\delta, \sigma_x, \sigma_y)$  are defined as

$$\rho_a(\delta, \sigma_x, \sigma_y) = \left[ \max \left\{ \frac{4(\mu_x + L_x \sigma_x)}{\mu_x}, \frac{2}{\sigma_x}, \frac{4(\mu_y + L_y \sigma_y)}{\mu_y}, \frac{2}{\sigma_y}, \frac{4L_{xy}}{\mu_x \delta}, \frac{4L_{xy} \delta}{\mu_y} \right\} \right]^{-1}, \quad (59)$$

$$\rho_b(\delta, \sigma_x, \sigma_y) = \left[ \max \left\{ \frac{4(\mu_x + L_x \sigma_x)}{\mu_x}, \frac{2}{\sigma_x}, \frac{8L_x(\mu_y + L_y \sigma_y)}{\mu_{xy}^2}, \frac{2}{\sigma_y}, \frac{2L_{xy}^2}{\mu_{xy}^2}, \frac{8L_x L_{xy} \delta}{\mu_{xy}^2}, \frac{4L_{xy}}{\mu_x \delta} \right\} \right]^{-1}, \quad (60)$$

$$\rho_c(\delta, \sigma_x, \sigma_y) = \left[ \max \left\{ \frac{4(\mu_y + L_y \sigma_y)}{\mu_y}, \frac{2}{\sigma_y}, \frac{8L_y(\mu_x + L_x \sigma_x)}{\mu_{yx}^2}, \frac{2}{\sigma_x}, \frac{2L_{xy}^2}{\mu_{yx}^2}, \frac{8L_y L_{xy} \delta}{\mu_{yx}^2}, \frac{4L_{xy} \delta}{\mu_y} \right\} \right]^{-1}, \quad (61)$$

$$\rho_d(\delta, \sigma_x, \sigma_y) = \left[ \max \left\{ \frac{8L_y(\mu_x + L_x \sigma_x)}{\mu_{yx}^2}, \frac{2}{\sigma_x}, \frac{8L_x(\mu_y + L_y \sigma_y)}{\mu_{xy}^2}, \frac{2}{\sigma_y}, \frac{8L_y L_{xy}}{\delta \mu_{yx}^2}, \frac{8L_x L_{xy} \delta}{\mu_{xy}^2}, \frac{2L_{xy}^2}{\mu_{yx}^2}, \frac{2L_{xy}^2}{\mu_{xy}^2} \right\} \right]^{-1}. \quad (62)$$

Let  $\Psi^k$  be the following Lyapunov function:

$$\begin{aligned} \Psi^k &= \frac{1}{\eta_x} \|x^k - x^*\|^2 + \frac{1}{\eta_y} \|y^k - y^*\|^2 + \frac{2}{\sigma_x} D_f(x_f^k, x^*) + \frac{2}{\sigma_y} D_g(y_f^k, y^*) \\ &\quad + \frac{1}{4\eta_y} \|y^k - y^{k-1}\|^2 - 2\langle y^k - y^{k-1}, \mathbf{A}(x^k - x^*) \rangle. \end{aligned} \quad (63)$$

Then, the following inequalities hold

$$\Psi^k \geq \frac{3}{4\eta_x} \|x^k - x^*\|^2 + \frac{1}{\eta_y} \|y^k - y^*\|^2, \quad (64)$$

$$\Psi^{k+1} \leq \theta \Psi^k. \quad (65)$$

*Proof.* After adding up (51) and (55) we get

$$\begin{aligned} (\text{LHS}) &\leq \left( \frac{1}{\eta_x} - \mu_x - \beta_x \mu_{yx}^2 \right) \|x^k - x^*\|^2 + \left( \frac{1}{\eta_y} - \mu_y - \beta_y \mu_{xy}^2 \right) \|y^k - y^*\|^2 \\ &\quad + \left( \mu_x + L_x \sigma_x - \frac{1}{2\eta_x} \right) \|x^{k+1} - x^k\|^2 + \left( \mu_y + L_y \sigma_y - \frac{1}{2\eta_y} \right) \|y^{k+1} - y^k\|^2 \\ &\quad + \left( \frac{2}{\sigma_x} - 1 \right) D_f(x_f^k, x^*) + \left( \frac{2}{\sigma_y} - 1 \right) D_g(y_f^k, y^*) + 2\langle y^{k+1} - y^k, \mathbf{A}(x^{k+1} - x^*) \rangle. \end{aligned}$$

where (LHS) is given as

$$(\text{LHS}) = \frac{1}{\eta_x} \|x^{k+1} - x^*\|^2 + \frac{1}{\eta_y} \|y^{k+1} - y^*\|^2 + \frac{2}{\sigma_x} D_f(x_f^{k+1}, x^*) + \frac{2}{\sigma_y} D_g(y_f^{k+1}, y^*).$$

Using Line 5 of the Algorithm 1 and Assumption 3 we get

$$\begin{aligned} (\text{LHS}) &\leq \left( \frac{1}{\eta_x} - \mu_x - \beta_x \mu_{yx}^2 \right) \|x^k - x^*\|^2 + \left( \frac{1}{\eta_y} - \mu_y - \beta_y \mu_{xy}^2 \right) \|y^k - y^*\|^2 \\ &\quad + \left( \mu_x + L_x \sigma_x - \frac{1}{2\eta_x} \right) \|x^{k+1} - x^k\|^2 + \left( \mu_y + L_y \sigma_y - \frac{1}{2\eta_y} \right) \|y^{k+1} - y^k\|^2 \\ &\quad + \left( \frac{2}{\sigma_x} - 1 \right) D_f(x_f^k, x^*) + \left( \frac{2}{\sigma_y} - 1 \right) D_g(y_f^k, y^*) \\ &\quad + 2\langle y^{k+1} - y^k, \mathbf{A}(x^{k+1} - x^*) \rangle - 2\theta \langle y^k - y^{k-1}, \mathbf{A}(x^{k+1} - x^*) \rangle \\ &\leq \left( \frac{1}{\eta_x} - \mu_x - \beta_x \mu_{yx}^2 \right) \|x^k - x^*\|^2 + \left( \frac{1}{\eta_y} - \mu_y - \beta_y \mu_{xy}^2 \right) \|y^k - y^*\|^2 \end{aligned}$$

$$\begin{aligned}
& + \left( \mu_x + L_x \sigma_x - \frac{1}{2\eta_x} \right) \|x^{k+1} - x^k\|^2 + \left( \mu_y + L_y \sigma_y - \frac{1}{2\eta_y} \right) \|y^{k+1} - y^k\|^2 \\
& + \left( \frac{2}{\sigma_x} - 1 \right) D_f(x_f^k, x^*) + \left( \frac{2}{\sigma_y} - 1 \right) D_g(y_f^k, y^*) \\
& + 2\langle y^{k+1} - y^k, \mathbf{A}(x^{k+1} - x^*) \rangle - 2\theta \langle y^k - y^{k-1}, \mathbf{A}(x^k - x^*) \rangle + 2\theta L_{xy} \|y^k - y^{k-1}\| \|x^{k+1} - x^k\|.
\end{aligned}$$

Using the definition of  $\eta_x$  and  $\eta_y$  and the fact that  $\theta < 1$  we get

$$\begin{aligned}
(\text{LHS}) & \leq \left( \frac{1}{\eta_x} - \mu_x - \beta_x \mu_{yx}^2 \right) \|x^k - x^*\|^2 + \left( \frac{1}{\eta_y} - \mu_y - \beta_y \mu_{xy}^2 \right) \|y^k - y^*\|^2 \\
& - \frac{1}{4\eta_x} \|x^{k+1} - x^k\|^2 - \frac{1}{4\eta_y} \|y^{k+1} - y^k\|^2 + \left( \frac{2}{\sigma_x} - 1 \right) D_f(x_f^k, x^*) + \left( \frac{2}{\sigma_y} - 1 \right) D_g(y_f^k, y^*) \\
& + 2\langle y^{k+1} - y^k, \mathbf{A}(x^{k+1} - x^*) \rangle - 2\theta \langle y^k - y^{k-1}, \mathbf{A}(x^k - x^*) \rangle + \frac{\theta}{2\sqrt{\eta_x \eta_y}} \|y^k - y^{k-1}\| \|x^{k+1} - x^k\| \\
& \leq \left( \frac{1}{\eta_x} - \mu_x - \beta_x \mu_{yx}^2 \right) \|x^k - x^*\|^2 + \left( \frac{1}{\eta_y} - \mu_y - \beta_y \mu_{xy}^2 \right) \|y^k - y^*\|^2 \\
& - \frac{1}{4\eta_x} \|x^{k+1} - x^k\|^2 - \frac{1}{4\eta_y} \|y^{k+1} - y^k\|^2 + \left( \frac{2}{\sigma_x} - 1 \right) D_f(x_f^k, x^*) + \left( \frac{2}{\sigma_y} - 1 \right) D_g(y_f^k, y^*) \\
& + 2\langle y^{k+1} - y^k, \mathbf{A}(x^{k+1} - x^*) \rangle - 2\theta \langle y^k - y^{k-1}, \mathbf{A}(x^k - x^*) \rangle + \frac{\theta}{4\eta_x} \|x^{k+1} - x^k\|^2 + \frac{\theta}{4\eta_y} \|y^k - y^{k-1}\|^2 \\
& \leq \left( \frac{1}{\eta_x} - \mu_x - \beta_x \mu_{yx}^2 \right) \|x^k - x^*\|^2 + \left( \frac{1}{\eta_y} - \mu_y - \beta_y \mu_{xy}^2 \right) \|y^k - y^*\|^2 \\
& + \frac{\theta}{4\eta_y} \|y^k - y^{k-1}\|^2 - \frac{1}{4\eta_y} \|y^{k+1} - y^k\|^2 + \left( \frac{2}{\sigma_x} - 1 \right) D_f(x_f^k, x^*) + \left( \frac{2}{\sigma_y} - 1 \right) D_g(y_f^k, y^*) \\
& + 2\langle y^{k+1} - y^k, \mathbf{A}(x^{k+1} - x^*) \rangle - 2\theta \langle y^k - y^{k-1}, \mathbf{A}(x^k - x^*) \rangle.
\end{aligned}$$

Using the definition of  $\beta_x$  and  $\beta_y$  we get

$$\begin{aligned}
(\text{LHS}) & \leq \left( 1 - \eta_x \mu_x - \min \left\{ \frac{\eta_x \mu_{yx}^2}{2L_y}, \frac{\mu_{yx}^2}{2L_{xy}^2} \right\} \right) \frac{1}{\eta_x} \|x^k - x^*\|^2 + \left( 1 - \eta_y \mu_y - \min \left\{ \frac{\eta_y \mu_{xy}^2}{2L_x}, \frac{\mu_{xy}^2}{2L_{xy}^2} \right\} \right) \frac{1}{\eta_y} \|y^k - y^*\|^2 \\
& + \frac{\theta}{4\eta_y} \|y^k - y^{k-1}\|^2 - \frac{1}{4\eta_y} \|y^{k+1} - y^k\|^2 + \left( \frac{2}{\sigma_x} - 1 \right) D_f(x_f^k, x^*) + \left( \frac{2}{\sigma_y} - 1 \right) D_g(y_f^k, y^*) \\
& + 2\langle y^{k+1} - y^k, \mathbf{A}(x^{k+1} - x^*) \rangle - 2\theta \langle y^k - y^{k-1}, \mathbf{A}(x^k - x^*) \rangle \\
& \leq \left( 1 - \max \left\{ \eta_x \mu_x, \min \left\{ \frac{\eta_x \mu_{yx}^2}{2L_y}, \frac{\mu_{yx}^2}{2L_{xy}^2} \right\} \right\} \right) \frac{1}{\eta_x} \|x^k - x^*\|^2 \\
& + \left( 1 - \max \left\{ \eta_y \mu_y, \min \left\{ \frac{\eta_y \mu_{xy}^2}{2L_x}, \frac{\mu_{xy}^2}{2L_{xy}^2} \right\} \right\} \right) \frac{1}{\eta_y} \|y^k - y^*\|^2 \\
& + \frac{\theta}{4\eta_y} \|y^k - y^{k-1}\|^2 - \frac{1}{4\eta_y} \|y^{k+1} - y^k\|^2 + \left( \frac{2}{\sigma_x} - 1 \right) D_f(x_f^k, x^*) + \left( \frac{2}{\sigma_y} - 1 \right) D_g(y_f^k, y^*) \\
& + 2\langle y^{k+1} - y^k, \mathbf{A}(x^{k+1} - x^*) \rangle - 2\theta \langle y^k - y^{k-1}, \mathbf{A}(x^k - x^*) \rangle.
\end{aligned}$$

Using the definition of  $\theta$  we get

$$\begin{aligned}
(\text{LHS}) & \leq \theta \left( \frac{1}{\eta_x} \|x^k - x^*\|^2 + \frac{1}{\eta_y} \|y^k - y^*\|^2 + \frac{1}{4\eta_y} \|y^k - y^{k-1}\|^2 - 2\langle y^k - y^{k-1}, \mathbf{A}(x^k - x^*) \rangle \right) \\
& + \theta \left( \frac{2}{\sigma_x} D_f(x_f^k, x^*) + \frac{2}{\sigma_y} D_g(y_f^k, y^*) \right) - \frac{1}{4\eta_y} \|y^{k+1} - y^k\|^2 + 2\langle y^{k+1} - y^k, \mathbf{A}(x^{k+1} - x^*) \rangle.
\end{aligned}$$

After rearranging and using the definition of  $\Psi^k$  we get

$$\Psi^{k+1} \leq \theta \Psi^k.$$

Finally, using the definition of  $\Psi^k$ ,  $\eta_x$  and  $\eta_y$  we get

$$\begin{aligned}
\Psi^k &\geq \frac{1}{\eta_x} \|x^k - x^*\|^2 + \frac{1}{\eta_y} \|y^k - y^*\|^2 + \frac{1}{4\eta_y} \|y^k - y^{k-1}\|^2 - 2\langle y^k - y^{k-1}, \mathbf{A}(x^k - x^*) \rangle \\
&\geq \frac{1}{\eta_x} \|x^k - x^*\|^2 + \frac{1}{\eta_y} \|y^k - y^*\|^2 + \frac{1}{4\eta_y} \|y^k - y^{k-1}\|^2 - 2L_{xy} \|y^k - y^{k-1}\| \|x^k - x^*\| \\
&\geq \frac{1}{\eta_x} \|x^k - x^*\|^2 + \frac{1}{\eta_y} \|y^k - y^*\|^2 + \frac{1}{4\eta_y} \|y^k - y^{k-1}\|^2 - \frac{1}{2\sqrt{\eta_x \eta_y}} \|y^k - y^{k-1}\| \|x^k - x^*\| \\
&\geq \frac{1}{\eta_x} \|x^k - x^*\|^2 + \frac{1}{\eta_y} \|y^k - y^*\|^2 + \frac{1}{4\eta_y} \|y^k - y^{k-1}\|^2 - \frac{1}{4\eta_x} \|x^k - x^*\|^2 - \frac{1}{4\eta_y} \|y^k - y^{k-1}\|^2 \\
&= \frac{3}{4\eta_x} \|x^k - x^*\|^2 + \frac{1}{\eta_y} \|y^k - y^*\|^2.
\end{aligned}$$

□

*Proof of Theorem 1.* From (64) and (65) we can conclude that

$$\frac{3}{4\eta_x} \|x^k - x^*\|^2 + \frac{1}{\eta_y} \|y^k - y^*\|^2 \leq \theta^k \Psi^0.$$

This implies the following inequality

$$\max \{ \|x^k - x^*\|^2, \|y^k - x^*\|^2 \} \leq \theta^k \Psi^0 \max \{ 4\eta_x/3, \eta_y \}.$$

Hence, we can conclude that

$$\max \{ \|x^k - x^*\|^2, \|y^k - x^*\|^2 \} \leq \epsilon,$$

as long as the number of iterations  $k$  satisfies

$$k \geq \frac{1}{1 - \theta} \log \frac{C}{\epsilon},$$

where  $C = \Psi^0 \max \{ 4\eta_x/3, \eta_y \}$ , which does not depend on  $\epsilon$ . From (58) we obtain

$$\frac{1}{1 - \theta} = \min \left\{ \frac{1}{\rho_a(\delta, \sigma_x, \sigma_y)}, \frac{1}{\rho_b(\delta, \sigma_x, \sigma_y)}, \frac{1}{\rho_c(\delta, \sigma_x, \sigma_y)}, \frac{1}{\rho_d(\delta, \sigma_x, \sigma_y)} \right\}.$$

We can now try to approximately optimize parameters  $\delta > 0$  and  $\sigma_x, \sigma_y \in (0, 1]$  to obtain the smallest possible values of  $\rho_a(\delta, \sigma_x, \sigma_y)^{-1}$ ,  $\rho_b(\delta, \sigma_x, \sigma_y)^{-1}$ ,  $\rho_c(\delta, \sigma_x, \sigma_y)^{-1}$ ,  $\rho_d(\delta, \sigma_x, \sigma_y)^{-1}$ . This can be done in a closed form and the result is the following:

$$\begin{aligned}
\frac{1}{\rho_a} &\leq 4 + 4 \max \left\{ \sqrt{\frac{L_x}{\mu_x}}, \sqrt{\frac{L_y}{\mu_y}}, \frac{L_{xy}}{\sqrt{\mu_x \mu_y}} \right\} \text{ for } \delta = \sqrt{\frac{\mu_y}{\mu_x}}, \sigma_x = \sqrt{\frac{\mu_x}{2L_x}}, \sigma_y = \sqrt{\frac{\mu_x}{2L_x}}, \\
\frac{1}{\rho_b} &\leq 4 + 8 \max \left\{ \frac{\sqrt{L_x L_y}}{\mu_{xy}}, \frac{L_{xy}}{\mu_{xy}} \sqrt{\frac{L_x}{\mu_x}}, \frac{L_{xy}^2}{\mu_{xy}^2} \right\} \text{ for } \delta = \sqrt{\frac{\mu_{xy}^2}{2\mu_x L_x}}, \sigma_x = \sqrt{\frac{\mu_x}{2L_x}}, \sigma_y = \min \left\{ 1, \sqrt{\frac{\mu_{xy}^2}{4L_x L_y}} \right\}, \\
\frac{1}{\rho_c} &\leq 4 + 8 \max \left\{ \frac{\sqrt{L_x L_y}}{\mu_{yx}}, \frac{L_{xy}}{\mu_{yx}} \sqrt{\frac{L_y}{\mu_y}}, \frac{L_{xy}^2}{\mu_{yx}^2} \right\} \text{ for } \delta = \sqrt{\frac{2\mu_y L_y}{\mu_{yx}^2}}, \sigma_x = \min \left\{ 1, \sqrt{\frac{\mu_{yx}^2}{4L_x L_y}} \right\}, \sigma_y = \sqrt{\frac{\mu_y}{2L_y}}, \\
\frac{1}{\rho_d} &\leq 2 + 8 \max \left\{ \frac{\sqrt{L_x L_y L_{xy}}}{\mu_{xy} \mu_{yx}}, \frac{L_{xy}^2}{\mu_{yx}^2}, \frac{L_{xy}^2}{\mu_{xy}^2} \right\} \text{ for } \delta = \frac{\mu_{xy}}{\mu_{yx}} \sqrt{\frac{L_y}{L_x}}, \sigma_x = \min \left\{ 1, \sqrt{\frac{\mu_{yx}^2}{4L_x L_y}} \right\}, \sigma_y = \min \left\{ 1, \sqrt{\frac{\mu_{xy}^2}{4L_x L_y}} \right\}.
\end{aligned}$$

Note, that we set  $\mu_y = 0$  in the bound for  $\rho_b^{-1}$ ,  $\mu_x = 0$  in the bound for  $\rho_c^{-1}$  and  $\mu_x = \mu_y = 0$  in the bound for  $\rho_d^{-1}$ . This is a valid move, because any convex function is 0-strongly convex by the definition of strong convexity. □

## D Proof of Theorem 2

**Lemma 6.** *Problem (11) has a unique solution  $(x^*, y^*)$ .*

*Proof.* Consider operator  $T: \mathbb{R}^{d_x} \times \mathbb{R}^{d_y} \rightarrow \mathbb{R}^{d_x} \times \mathbb{R}^{d_y}$  defined as  $T: (x, y) \mapsto (x - t_x \nabla_x F(x, y), y - t_y \nabla_y F(x, y))$  for some fixed  $t_x, t_y > 0$ . It is obvious that  $(x, y)$  is a fixed point of operator  $T$  if and only if  $(x, y)$  is a solution to problem (11). If one can show that this operator is contractive, then it has a unique fixed point due to Banach fixed-point theorem. The proof of the fact that  $T$  is contractive is similar to the proof of the rest of Theorem 2.  $\square$

**Lemma 7.** *Let  $\eta_x$  be defined as*

$$\eta_x = \min \left\{ \frac{1}{8L_x}, \frac{\delta}{4L_{xy}} \right\}, \quad (66)$$

*and let  $\eta_y$  be defined as*

$$\eta_y = \min \left\{ \frac{1}{8L_y}, \frac{1}{4\delta L_{xy}} \right\}, \quad (67)$$

*where  $\delta > 0$  is a parameter. Let  $\theta$  be defined as*

$$\theta = \theta(\delta) = 1 - \max \{ \rho_a(\delta), \rho_b(\delta), \rho_c(\delta), \rho_d(\delta) \}, \quad (68)$$

*where  $\rho_a(\delta), \rho_b(\delta), \rho_c(\delta), \rho_d(\delta)$  are defined as*

$$\frac{1}{\rho_a(\delta)} = \max \left\{ \frac{8L_x}{\mu_x}, \frac{8L_y}{\mu_y}, \frac{4L_{xy}}{\delta\mu_x}, \frac{4L_{xy}\delta}{\mu_y} \right\}, \quad (69)$$

$$\frac{1}{\rho_b(\delta)} = \max \left\{ \frac{8L_x}{\mu_x}, \frac{512L_xL_y}{\mu_{xy}^2}, \frac{4L_{xy}}{\delta\mu_x}, \frac{256L_xL_{xy}\delta}{\mu_{xy}^2}, \frac{256L_yL_{xy}}{\mu_{xy}^2\delta}, \frac{128L_{xy}^2}{\mu_{xy}^2} \right\}, \quad (70)$$

$$\frac{1}{\rho_c(\delta)} = \max \left\{ \frac{8L_y}{\mu_y}, \frac{512L_xL_y}{\mu_{yx}^2}, \frac{4L_{xy}\delta}{\mu_y}, \frac{256L_xL_{xy}\delta}{\mu_{yx}^2}, \frac{256L_yL_{xy}}{\mu_{yx}^2\delta}, \frac{128L_{xy}^2}{\mu_{yx}^2} \right\}, \quad (71)$$

$$\frac{1}{\rho_d(\delta)} = \max \left\{ \frac{512L_xL_y}{\min\{\mu_{xy}^2, \mu_{yx}^2\}}, \frac{256L_xL_{xy}\delta}{\min\{\mu_{xy}^2, \mu_{yx}^2\}}, \frac{256L_yL_{xy}}{\min\{\mu_{xy}^2, \mu_{yx}^2\}\delta}, \frac{128L_{xy}^2}{\min\{\mu_{xy}^2, \mu_{yx}^2\}} \right\}. \quad (72)$$

*Let  $\Psi^k$  be the following Lyapunov function:*

$$\Psi^k = \frac{1}{\eta_x} \|x^k - x^*\|^2 + \frac{1}{\eta_y} \|y^k - y^*\|^2 - 2\langle \nabla_x F(x^{k-1}, y^k) - \nabla_x F(x^{k-1}, y^{k-1}), x^k - x^* \rangle + \frac{5}{16\eta_y} \|y^k - y^{k-1}\|^2. \quad (73)$$

*Then, the following inequalities hold*

$$\Psi^k \geq \frac{3}{4\eta_x} \|x^k - x^*\|^2 + \frac{1}{\eta_y} \|y^k - y^*\|^2, \quad (74)$$

$$\Psi^{k+1} \leq \theta \Psi^k. \quad (75)$$

*Proof.* Using Line 5 of the Algorithm 2 we get.

$$\begin{aligned} \frac{1}{\eta_x} \|x^{k+1} - x^*\|^2 &= \frac{1}{\eta_x} \|x^k - x^*\|^2 + \frac{2}{\eta_x} \langle x^{k+1} - x^k, x^{k+1} - x^* \rangle - \frac{1}{\eta_x} \|x^{k+1} - x^k\|^2 \\ &= \frac{1}{\eta_x} \|x^k - x^*\|^2 - \frac{1}{\eta_x} \|x^{k+1} - x^k\|^2 \\ &\quad - 2\langle \nabla_x F(x^k, y^k) + \theta(\nabla_x F(x^{k-1}, y^k) - \nabla_x F(x^{k-1}, y^{k-1})), x^{k+1} - x^* \rangle \\ &= \frac{1}{\eta_x} \|x^k - x^*\|^2 - \frac{1}{\eta_x} \|x^{k+1} - x^k\|^2 - 2\langle \nabla_x F(x^k, y^{k+1}), x^{k+1} - x^k + x^k - x^* \rangle \\ &\quad + 2\langle \nabla_x F(x^k, y^{k+1}) - \nabla_x F(x^k, y^k), x^{k+1} - x^* \rangle \\ &\quad - 2\theta \langle \nabla_x F(x^{k-1}, y^k) - \nabla_x F(x^{k-1}, y^{k-1}), x^{k+1} - x^* \rangle. \end{aligned}$$



Using the Assumption 4 we get

$$\begin{aligned} \frac{1}{\eta_x} \|x^{k+1} - x^*\|^2 &\leq \left(\frac{1}{\eta_x} - \mu_x\right) \|x^k - x^*\|^2 + \left(L_x - \frac{1}{\eta_x}\right) \|x^{k+1} - x^k\|^2 - 2(F(x^{k+1}, y^{k+1}) - F(x^*, y^{k+1})) \\ &\quad + 2\langle \nabla_x F(x^k, y^{k+1}) - \nabla_x F(x^k, y^k), x^{k+1} - x^* \rangle \\ &\quad - 2\theta \langle \nabla_x F(x^{k-1}, y^k) - \nabla_x F(x^{k-1}, y^{k-1}), x^{k+1} - x^* \rangle. \end{aligned}$$

Using the Assumption 5 we get

$$\begin{aligned} \frac{1}{\eta_x} \|x^{k+1} - x^*\|^2 &\leq \left(\frac{1}{\eta_x} - \mu_x\right) \|x^k - x^*\|^2 + \left(L_x - \frac{1}{\eta_x}\right) \|x^{k+1} - x^k\|^2 - 2(F(x^{k+1}, y^{k+1}) - F(x^*, y^{k+1})) \\ &\quad + 2\langle \nabla_x F(x^k, y^{k+1}) - \nabla_x F(x^k, y^k), x^{k+1} - x^* \rangle - 2\theta \langle \nabla_x F(x^{k-1}, y^k) - \nabla_x F(x^{k-1}, y^{k-1}), x^k - x^* \rangle \\ &\quad + 2L_{xy}\theta \|x^{k+1} - x^k\| \|y^k - y^{k-1}\|. \end{aligned}$$

Similarly, we can obtain the following upper-bound on  $\frac{1}{\eta_y} \|y^{k+1} - y^*\|^2$ :

$$\frac{1}{\eta_y} \|y^{k+1} - y^*\|^2 \leq \left(\frac{1}{\eta_y} - \mu_y\right) \|y^k - y^*\|^2 + \left(L_y - \frac{1}{\eta_y}\right) \|y^{k+1} - y^k\|^2 + 2(F(x^{k+1}, y^{k+1}) - F(x^{k+1}, y^*)).$$

Summing up the upper-bounds on  $\frac{1}{\eta_x} \|x^{k+1} - x^*\|^2$  and  $\frac{1}{\eta_y} \|y^{k+1} - y^*\|^2$  gives

$$\begin{aligned} \text{(LHS)} &\leq \left(\frac{1}{\eta_x} - \mu_x\right) \|x^k - x^*\|^2 + \left(L_x - \frac{1}{\eta_x}\right) \|x^{k+1} - x^k\|^2 \\ &\quad + \left(\frac{1}{\eta_y} - \mu_y\right) \|y^k - y^*\|^2 + \left(L_y - \frac{1}{\eta_y}\right) \|y^{k+1} - y^k\|^2 \\ &\quad + 2L_{xy}\theta \|x^{k+1} - x^k\| \|y^k - y^{k-1}\| - 2\theta \langle \nabla_x F(x^{k-1}, y^k) - \nabla_x F(x^{k-1}, y^{k-1}), x^k - x^* \rangle \\ &\quad + 2(F(x^*, y^{k+1}) - F(x^{k+1}, y^*)), \end{aligned}$$

where (LHS) is defined as

$$\text{(LHS)} = \frac{1}{\eta_x} \|x^{k+1} - x^*\|^2 + \frac{1}{\eta_y} \|y^{k+1} - y^*\|^2 - 2\langle \nabla_x F(x^k, y^{k+1}) - \nabla_x F(x^k, y^k), x^{k+1} - x^* \rangle.$$

The Assumption 4 states, that function  $F(x, y)$  is  $L_x$ -smooth in  $x$  and  $L_y$ -smooth in  $y$ . Hence, using the optimality conditions (40) we get

$$\begin{aligned} \text{(LHS)} &\leq \left(\frac{1}{\eta_x} - \mu_x\right) \|x^k - x^*\|^2 + \left(L_x - \frac{1}{\eta_x}\right) \|x^{k+1} - x^k\|^2 \\ &\quad + \left(\frac{1}{\eta_y} - \mu_y\right) \|y^k - y^*\|^2 + \left(L_y - \frac{1}{\eta_y}\right) \|y^{k+1} - y^k\|^2 \\ &\quad + 2L_{xy}\theta \|x^{k+1} - x^k\| \|y^k - y^{k-1}\| - 2\theta \langle \nabla_x F(x^{k-1}, y^k) - \nabla_x F(x^{k-1}, y^{k-1}), x^k - x^* \rangle \\ &\quad - 2(F(x^{k+1}, y^*) - F(x^*, y^*)) - 2(F(x^*, y^*) - F(x^*, y^{k+1})) \\ &\leq \left(\frac{1}{\eta_x} - \mu_x\right) \|x^k - x^*\|^2 + \left(L_x - \frac{1}{\eta_x}\right) \|x^{k+1} - x^k\|^2 \\ &\quad + \left(\frac{1}{\eta_y} - \mu_y\right) \|y^k - y^*\|^2 + \left(L_y - \frac{1}{\eta_y}\right) \|y^{k+1} - y^k\|^2 \\ &\quad + 2L_{xy}\theta \|x^{k+1} - x^k\| \|y^k - y^{k-1}\| - 2\theta \langle \nabla_x F(x^{k-1}, y^k) - \nabla_x F(x^{k-1}, y^{k-1}), x^k - x^* \rangle \\ &\quad - \frac{\delta_x}{L_x} \|\nabla_x F(x^{k+1}, y^*)\|^2 - \frac{\delta_y}{L_y} \|\nabla_y F(x^*, y^{k+1})\|^2, \end{aligned}$$

where  $\delta_x, \delta_y \in (0, 1]$  are some parameters, that will be defined later. Using the Assumption 6 we get

$$\begin{aligned}
(\text{LHS}) &\leq \left(\frac{1}{\eta_x} - \mu_x\right) \|x^k - x^*\|^2 + \left(L_x - \frac{1}{\eta_x}\right) \|x^{k+1} - x^k\|^2 \\
&\quad + \left(\frac{1}{\eta_y} - \mu_y\right) \|y^k - y^*\|^2 + \left(L_y - \frac{1}{\eta_y}\right) \|y^{k+1} - y^k\|^2 \\
&\quad + 2L_{xy}\theta\|x^{k+1} - x^k\|\|y^k - y^{k-1}\| - 2\theta\langle \nabla_x F(x^{k-1}, y^k) - \nabla_x F(x^{k-1}, y^{k-1}), x^k - x^* \rangle \\
&\quad - \frac{\delta_x}{2L_x} \|\nabla_x F(x^{k+1}, y^*) - \nabla_x F(x^{k+1}, y^k)\|^2 + \frac{\delta_x}{L_x} \|\nabla_x F(x^{k+1}, y^k)\|^2 \\
&\quad - \frac{\delta_y}{2L_y} \|\nabla_y F(x^*, y^{k+1}) - \nabla_y F(x^k, y^{k+1})\|^2 + \frac{\delta_y}{L_y} \|\nabla_y F(x^k, y^{k+1})\|^2 \\
&\leq \left(\frac{1}{\eta_x} - \mu_x - \frac{\delta_y \mu_{yx}^2}{2L_y}\right) \|x^k - x^*\|^2 + \left(L_x - \frac{1}{\eta_x}\right) \|x^{k+1} - x^k\|^2 \\
&\quad + \left(\frac{1}{\eta_y} - \mu_y - \frac{\delta_x \mu_{xy}^2}{2L_x}\right) \|y^k - y^*\|^2 + \left(L_y - \frac{1}{\eta_y}\right) \|y^{k+1} - y^k\|^2 \\
&\quad + 2L_{xy}\theta\|x^{k+1} - x^k\|\|y^k - y^{k-1}\| - 2\theta\langle \nabla_x F(x^{k-1}, y^k) - \nabla_x F(x^{k-1}, y^{k-1}), x^k - x^* \rangle \\
&\quad + \frac{\delta_x}{L_x} \|\nabla_x F(x^{k+1}, y^k)\|^2 + \frac{\delta_y}{L_y} \|\nabla_y F(x^k, y^{k+1})\|^2
\end{aligned}$$

Using Lines 5 and 6 of the Algorithm 2 and the Lipschitzness property of  $\nabla_x F(x, y)$  and  $\nabla_y F(x, y)$  we get

$$\begin{aligned}
(\text{LHS}) &\leq \left(\frac{1}{\eta_x} - \mu_x - \frac{\delta_y \mu_{yx}^2}{2L_y}\right) \|x^k - x^*\|^2 + \left(L_x - \frac{1}{\eta_x}\right) \|x^{k+1} - x^k\|^2 \\
&\quad + \left(\frac{1}{\eta_y} - \mu_y - \frac{\delta_x \mu_{xy}^2}{2L_x}\right) \|y^k - y^*\|^2 + \left(L_y - \frac{1}{\eta_y}\right) \|y^{k+1} - y^k\|^2 \\
&\quad + 2L_{xy}\theta\|x^{k+1} - x^k\|\|y^k - y^{k-1}\| - 2\theta\langle \nabla_x F(x^{k-1}, y^k) - \nabla_x F(x^{k-1}, y^{k-1}), x^k - x^* \rangle \\
&\quad + \frac{2\delta_x}{L_x} \|\nabla_x F(x^{k+1}, y^k) - \nabla_x F(x^k, y^k) - \theta(\nabla_x F(x^{k-1}, y^k) - \nabla_x F(x^{k-1}, y^{k-1}))\|^2 + \frac{2\delta_x}{L_x \eta_x^2} \|x^{k+1} - x^k\|^2 \\
&\quad + \frac{2\delta_y}{L_y} \|\nabla_y F(x^k, y^{k+1}) - \nabla_y F(x^{k+1}, y^k)\|^2 + \frac{2\delta_y}{L_y \eta_y^2} \|y^{k+1} - y^k\|^2 \\
&\leq \left(\frac{1}{\eta_x} - \mu_x - \frac{\delta_y \mu_{yx}^2}{2L_y}\right) \|x^k - x^*\|^2 + \left(L_x - \frac{1}{\eta_x}\right) \|x^{k+1} - x^k\|^2 \\
&\quad + \left(\frac{1}{\eta_y} - \mu_y - \frac{\delta_x \mu_{xy}^2}{2L_x}\right) \|y^k - y^*\|^2 + \left(L_y - \frac{1}{\eta_y}\right) \|y^{k+1} - y^k\|^2 \\
&\quad + 2L_{xy}\theta\|x^{k+1} - x^k\|\|y^k - y^{k-1}\| - 2\theta\langle \nabla_x F(x^{k-1}, y^k) - \nabla_x F(x^{k-1}, y^{k-1}), x^k - x^* \rangle \\
&\quad + 4\delta_x L_x \|x^{k+1} - x^k\|^2 + \frac{4\delta_x L_{xy}^2 \theta^2}{L_x} \|y^k - y^{k-1}\|^2 + \frac{2\delta_x}{L_x \eta_x^2} \|x^{k+1} - x^k\|^2 \\
&\quad + 4\delta_y L_y \|y^{k+1} - y^k\|^2 + \frac{4\delta_y L_{xy}^2}{L_y} \|x^{k+1} - x^k\|^2 + \frac{2\delta_y}{L_y \eta_y^2} \|y^{k+1} - y^k\|^2.
\end{aligned}$$

Now, we set  $\delta_x = \min\{1, c_x \eta_x L_x\}$ ,  $\delta_y = \min\{1, c_y \eta_y L_y\}$ , where  $c_x, c_y > 0$  will be defined later, and obtain

$$\begin{aligned}
(\text{LHS}) &\leq \left(\frac{1}{\eta_x} - \mu_x - \frac{\delta_y \mu_{yx}^2}{2L_y}\right) \|x^k - x^*\|^2 + \left(L_x - \frac{1}{\eta_x}\right) \|x^{k+1} - x^k\|^2 \\
&\quad + \left(\frac{1}{\eta_y} - \mu_y - \frac{\delta_x \mu_{xy}^2}{2L_x}\right) \|y^k - y^*\|^2 + \left(L_y - \frac{1}{\eta_y}\right) \|y^{k+1} - y^k\|^2
\end{aligned}$$

$$\begin{aligned}
& + 2L_{xy}\theta\|x^{k+1} - x^k\|\|y^k - y^{k-1}\| - 2\theta\langle\nabla_x F(x^{k-1}, y^k) - \nabla_x F(x^{k-1}, y^{k-1}), x^k - x^*\rangle \\
& + 4c_x\eta_x L_x^2\|x^{k+1} - x^k\|^2 + 4c_x\eta_x L_{xy}^2\theta^2\|y^k - y^{k-1}\|^2 + \frac{2c_x}{\eta_x}\|x^{k+1} - x^k\|^2 \\
& + 4c_y\eta_y L_y^2\|y^{k+1} - y^k\|^2 + 4c_y\eta_y L_{xy}^2\|x^{k+1} - x^k\|^2 + \frac{2c_y}{\eta_y}\|y^{k+1} - y^k\|^2.
\end{aligned}$$

Using the definition of  $\eta_x$  and  $\eta_y$  we get

$$\begin{aligned}
(\text{LHS}) & \leq \left(\frac{1}{\eta_x} - \mu_x - \frac{\delta_y\mu_{yx}^2}{2L_y}\right)\|x^k - x^*\|^2 + \left(L_x - \frac{1}{\eta_x}\right)\|x^{k+1} - x^k\|^2 \\
& + \left(\frac{1}{\eta_y} - \mu_y - \frac{\delta_x\mu_{xy}^2}{2L_x}\right)\|y^k - y^*\|^2 + \left(L_y - \frac{1}{\eta_y}\right)\|y^{k+1} - y^k\|^2 \\
& - 2\theta\langle\nabla_x F(x^{k-1}, y^k) - \nabla_x F(x^{k-1}, y^{k-1}), x^k - x^*\rangle \\
& + 4c_x\eta_x L_x^2\|x^{k+1} - x^k\|^2 + \frac{(c_x + 1)\theta^2}{4\eta_y}\|y^k - y^{k-1}\|^2 + \frac{2c_x}{\eta_x}\|x^{k+1} - x^k\|^2 \\
& + 4c_y\eta_y L_y^2\|y^{k+1} - y^k\|^2 + \frac{c_y + 1}{4\eta_x}\|x^{k+1} - x^k\|^2 + \frac{2c_y}{\eta_y}\|y^{k+1} - y^k\|^2.
\end{aligned}$$

Now, we choose  $c_x = c_y = \frac{1}{4}$  and get

$$\begin{aligned}
(\text{LHS}) & \leq \left(\frac{1}{\eta_x} - \mu_x - \frac{\delta_y\mu_{yx}^2}{2L_y}\right)\|x^k - x^*\|^2 + \left(L_x - \frac{1}{\eta_x}\right)\|x^{k+1} - x^k\|^2 \\
& + \left(\frac{1}{\eta_y} - \mu_y - \frac{\delta_x\mu_{xy}^2}{2L_x}\right)\|y^k - y^*\|^2 + \left(L_y - \frac{1}{\eta_y}\right)\|y^{k+1} - y^k\|^2 \\
& - 2\theta\langle\nabla_x F(x^{k-1}, y^k) - \nabla_x F(x^{k-1}, y^{k-1}), x^k - x^*\rangle \\
& + \eta_x L_x^2\|x^{k+1} - x^k\|^2 + \frac{5\theta^2}{16\eta_y}\|y^k - y^{k-1}\|^2 + \frac{1}{2\eta_x}\|x^{k+1} - x^k\|^2 \\
& + \eta_y L_y^2\|y^{k+1} - y^k\|^2 + \frac{5}{16\eta_x}\|x^{k+1} - x^k\|^2 + \frac{1}{2\eta_y}\|y^{k+1} - y^k\|^2 \\
& = \left(\frac{1}{\eta_x} - \mu_x - \frac{\delta_y\mu_{yx}^2}{2L_y}\right)\|x^k - x^*\|^2 + \frac{\eta_x L_x + \eta_x^2 L_x^2 - 3/16}{\eta_x}\|x^{k+1} - x^k\|^2 \\
& + \left(\frac{1}{\eta_y} - \mu_y - \frac{\delta_x\mu_{xy}^2}{2L_x}\right)\|y^k - y^*\|^2 + \frac{\eta_y L_y + \eta_y^2 L_y^2 - 3/16}{\eta_y}\|y^{k+1} - y^k\|^2 \\
& - 2\theta\langle\nabla_x F(x^{k-1}, y^k) - \nabla_x F(x^{k-1}, y^{k-1}), x^k - x^*\rangle + \frac{5\theta^2}{16\eta_y}\|y^k - y^{k-1}\|^2 - \frac{5}{16\eta_y}\|y^{k+1} - y^k\|^2.
\end{aligned}$$

Using the definition of  $\eta_x$  and  $\eta_y$  we get

$$\begin{aligned}
(\text{LHS}) & \leq \left(\frac{1}{\eta_x} - \mu_x - \frac{\delta_y\mu_{yx}^2}{2L_y}\right)\|x^k - x^*\|^2 + \left(\frac{1}{\eta_y} - \mu_y - \frac{\delta_x\mu_{xy}^2}{2L_x}\right)\|y^k - y^*\|^2 \\
& - 2\theta\langle\nabla_x F(x^{k-1}, y^k) - \nabla_x F(x^{k-1}, y^{k-1}), x^k - x^*\rangle + \frac{5\theta^2}{16\eta_y}\|y^k - y^{k-1}\|^2 - \frac{5}{16\eta_y}\|y^{k+1} - y^k\|^2.
\end{aligned}$$

Using the definition of  $\delta_x$  and  $\delta_y$  we get

$$\begin{aligned}
(\text{LHS}) & \leq \left(1 - \max\left\{\eta_x\mu_x, \min\left\{\frac{\eta_x\mu_{yx}^2}{2L_y}, \frac{\eta_x\eta_y\mu_{yx}^2}{8}\right\}\right\}\right)\frac{1}{\eta_x}\|x^k - x^*\|^2 \\
& + \left(1 - \max\left\{\eta_y\mu_y, \min\left\{\frac{\eta_y\mu_{xy}^2}{2L_x}, \frac{\eta_y\eta_x\mu_{xy}^2}{8}\right\}\right\}\right)\frac{1}{\eta_y}\|y^k - y^*\|^2 \\
& - 2\theta\langle\nabla_x F(x^{k-1}, y^k) - \nabla_x F(x^{k-1}, y^{k-1}), x^k - x^*\rangle + \frac{5\theta^2}{16\eta_y}\|y^k - y^{k-1}\|^2 - \frac{5}{16\eta_y}\|y^{k+1} - y^k\|^2.
\end{aligned}$$

Using the definition of  $\eta_x, \eta_y$  and  $\theta$  we get

$$\begin{aligned} \text{(LHS)} &\leq \frac{\theta}{\eta_x} \|x^k - x^*\|^2 + \frac{\theta}{\eta_y} \|y^k - y^*\|^2 - 2\theta \langle \nabla_x F(x^{k-1}, y^k) - \nabla_x F(x^{k-1}, y^{k-1}), x^k - x^* \rangle + \frac{5\theta}{16\eta_y} \|y^k - y^{k-1}\|^2 \\ &\quad - \frac{5}{16\eta_y} \|y^{k+1} - y^k\|^2. \end{aligned}$$

After rearranging and using the definition of  $\Psi^k$  we get

$$\Psi^{k+1} \leq \theta \Psi^k.$$

Finally, using the definition of  $\Psi^k, \eta_x$  and  $\eta_y$  we get

$$\begin{aligned} \Psi^k &= \frac{1}{\eta_x} \|x^k - x^*\|^2 + \frac{1}{\eta_y} \|y^k - y^*\|^2 - 2\langle \nabla_x F(x^{k-1}, y^k) - \nabla_x F(x^{k-1}, y^{k-1}), x^k - x^* \rangle + \frac{5}{16\eta_y} \|y^k - y^{k-1}\|^2 \\ &\geq \frac{1}{\eta_x} \|x^k - x^*\|^2 + \frac{1}{\eta_y} \|y^k - y^*\|^2 - 2L_{xy} \|y^k - y^{k-1}\| \|x^k - x^*\| + \frac{5}{16\eta_y} \|y^k - y^{k-1}\|^2 \\ &\geq \frac{1}{\eta_x} \|x^k - x^*\|^2 + \frac{1}{\eta_y} \|y^k - y^*\|^2 - \frac{1}{4\eta_x} \|x^k - x^*\|^2 - \frac{1}{4\eta_y} \|y^k - y^{k-1}\|^2 + \frac{5}{16\eta_y} \|y^k - y^{k-1}\|^2 \\ &\geq \frac{3}{4\eta_x} \|x^k - x^*\|^2 + \frac{1}{\eta_y} \|y^k - y^*\|^2. \end{aligned}$$

□

*Proof of Theorem 2.* From (74) and (75) we can conclude that

$$\frac{3}{4\eta_x} \|x^k - x^*\|^2 + \frac{1}{\eta_y} \|y^k - y^*\|^2 \leq \theta^k \Psi^0.$$

This implies the following inequality

$$\max \{ \|x^k - x^*\|^2, \|y^k - y^*\|^2 \} \leq \theta^k \Psi^0 \max \{ 4\eta_x/3, \eta_y \}.$$

Hence, we can conclude that

$$\max \{ \|x^k - x^*\|^2, \|y^k - y^*\|^2 \} \leq \epsilon,$$

as long as the number of iterations  $k$  satisfies

$$k \geq \frac{1}{1-\theta} \log \frac{C}{\epsilon},$$

where  $C = \Psi^0 \max \{ 4\eta_x/3, \eta_y \}$ , which does not depend on  $\epsilon$ . From (68) we obtain

$$\frac{1}{1-\theta} = \min \left\{ \frac{1}{\rho_a(\delta)}, \frac{1}{\rho_b(\delta)}, \frac{1}{\rho_c(\delta)}, \frac{1}{\rho_d(\delta)} \right\}.$$

Now, we find the parameter  $\delta$  to obtain the following upper bounds on  $\rho_a(\delta), \rho_b(\delta), \rho_c(\delta), \rho_d(\delta)$ :

$$\frac{1}{\rho_a} = \max \left\{ \frac{8L_x}{\mu_x}, \frac{8L_y}{\mu_y}, \frac{4L_{xy}}{\sqrt{\mu_x \mu_y}} \right\} \text{ for } \delta = \sqrt{\frac{\mu_y}{\mu_x}}, \quad (76)$$

$$\frac{1}{\rho_b} = \max \left\{ \frac{8L_x}{\mu_x}, \frac{512L_x L_y}{\mu_{xy}^2}, \frac{128L_{xy}^2}{\mu_{xy}^2} \right\} \text{ for } \delta = \max \left\{ \frac{\mu_{xy}}{8\sqrt{\mu_x L_x}}, \sqrt{\frac{L_y}{L_x}} \right\}, \quad (77)$$

$$\frac{1}{\rho_c} = \max \left\{ \frac{8L_y}{\mu_y}, \frac{512L_x L_y}{\mu_{yx}^2}, \frac{128L_{xy}^2}{\mu_{yx}^2} \right\} \text{ for } \delta = \min \left\{ \frac{8\sqrt{\mu_y L_y}}{\mu_{yx}}, \sqrt{\frac{L_y}{L_x}} \right\}, \quad (78)$$

$$\frac{1}{\rho_d} = \max \left\{ \frac{512L_x L_y}{\mu_{xy}^2}, \frac{512L_x L_y}{\mu_{yx}^2}, \frac{128L_{xy}^2}{\mu_{xy}^2}, \frac{128L_{xy}^2}{\mu_{yx}^2} \right\} \text{ for } \delta = \sqrt{\frac{L_y}{L_x}}. \quad (79)$$

□