

Sequential minimax optimization methods for bilevel optimization with strongly convex lower-level objective function

Zhaosong Lu ^{*} Sanyou Mei ^{*} Jin Zhang [†]

November 16, 2023

Abstract

In this paper we study a class of unconstrained and constrained bilevel optimization problems with strongly convex lower-level objective function. Specifically, we propose sequential minimax optimization methods for solving them, whose subproblems are nonconvex-strongly-concave unconstrained minimax problems and suitably solved by a first-order method developed by the authors that leverages the strong concavity structure. Under suitable assumptions, an *operation complexity* of $\mathcal{O}(\varepsilon^{-3} \log \varepsilon^{-1})$ and $\mathcal{O}(\varepsilon^{-6} \log \varepsilon^{-1})$, measured by their fundamental operations, is established for the proposed methods for finding an ε -KKT solution of the unconstrained and constrained bilevel optimization problems respectively, which improves the previously best known operation complexity by a factor of $1/\varepsilon$.

Keywords: bilevel optimization, minimax optimization, first-order methods, operation complexity

Mathematics Subject Classification: 90C26, 90C30, 90C47, 90C99, 65K05

1 Introduction

Bilevel optimization is a two-level hierarchical optimization, which is typically in the form of

$$\begin{aligned} f^* = \min \quad & f(x, y) \\ \text{s.t.} \quad & y \in \underset{z}{\text{Argmin}} \{ \tilde{f}(x, z) | \tilde{g}(x, z) \leq 0 \}.^1 \end{aligned} \quad (1)$$

Bilevel optimization has widely been used in many areas, including adversarial training [44, 45, 57], continual learning [35], hyperparameter tuning [3, 15, 48], image reconstruction [8], meta-learning [4, 24, 51], neural architecture search [13, 33], reinforcement learning [19, 27], and Stackelberg games [59]. More applications about it can be found in [2, 7, 9, 10, 11, 54] and the references therein. Theoretical properties including optimality conditions of (1) have been extensively studied in the literature (e.g., see [11, 12, 42, 58, 62]).

Numerous methods have been developed for solving some special cases of (1). For example, constraint-based methods [18, 53], deterministic gradient-based methods [14, 15, 16, 20, 43, 50, 51], and stochastic gradient-based methods [6, 17, 21, 22, 25, 26, 28, 29, 61] were proposed for solving (1) with $\tilde{g} \equiv 0$, f , \tilde{f} being smooth, and \tilde{f} being *strongly convex* with respect to y . For a similar case as this but with \tilde{f} being *convex* with respect to y , a zeroth-order method was recently proposed in [5], and also numerical methods were developed in [30, 56, 32] by solving (1) as a single or sequential smooth constrained optimization problems. Besides, when all the functions in (1) are smooth and \tilde{f} , \tilde{g} are *convex* with respect to y , gradient-type methods were proposed by solving a mathematical program with equilibrium constraints resulting from replacing the lower-level optimization problem of (1) by its first-order optimality conditions (e.g., see [1, 41, 49]). Recently, difference-of-convex (DC) algorithms were developed in [63] for solving (1) with f being a DC function, and \tilde{f} , \tilde{g} being convex functions. In addition, penalty type of methods

^{*}Department of Industrial and Systems Engineering, University of Minnesota, USA (email: zhaosong@umn.edu, mei00035@umn.edu). The work of these authors was partially supported by NSF Award IIS-2211491.

[†]Department of Mathematics, SUSTech International Center for Mathematics, Southern University of Science and Technology, Shenzhen 518055, China (email: zhangj9@sustech.edu.cn). The work of this author was partially supported by National Natural Science Foundation of China Grant 12271161.

¹For ease of reading, throughout this paper the tilde symbol is particularly used for the functions related to the lower-level optimization problem. Besides, “Argmin” denotes the set of optimal solutions of the associated problem.

were proposed in [23, 38, 52] for solving (1). Notably, the paper [38] demonstrates for the first time that bilevel optimization can be approximately solved as minimax optimization. Specifically, it reformulates bilevel optimization as minimax optimization using a novel double penalty scheme and proposes a first-order penalty method with *complexity guarantees* for bilevel optimization via solving a single minimax problem. Subsequently, the work [39] proposed a first-order method for solving (1) via solving a sequence of minimax problems. Lately, a practically efficient multi-stage gradient descent and ascent algorithm (GDA) was developed in [60] for (1) with $\tilde{g} \equiv 0$, f being convex and Lipschitz continuous, and \tilde{f} being strongly convex and Lipschitz smooth via solving the aforementioned minimax reformulation of (1). More discussion on algorithmic development for bilevel optimization can be found in [2, 7, 11, 34, 55, 58] and the references therein.

In this paper, we consider problem (1) that has at least one optimal solution and satisfies the following assumptions.

Assumption 1. (i) $f(x, y) = f_1(x, y) + f_2(x)$ and $\tilde{f}(x, y) = \tilde{f}_1(x, y) + \tilde{f}_2(y)$ are continuous on $\mathcal{X} \times \mathcal{Y}$, where $f_2 : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{+\infty\}$ and $\tilde{f}_2 : \mathbb{R}^m \rightarrow \mathbb{R} \cup \{+\infty\}$ are proper closed convex functions, $\tilde{f}_1(x, \cdot)$ is σ -strongly-convex for any given $x \in \mathcal{X}$, and f_1, \tilde{f}_1 are respectively $L_{\nabla f_1}$ - and $L_{\nabla \tilde{f}_1}$ -smooth on $\mathcal{X} \times \mathcal{Y}$ with $\mathcal{X} := \text{dom } f_2$ and $\mathcal{Y} := \text{dom } \tilde{f}_2$.

(ii) The proximal operator associated with f_2 and \tilde{f}_2 can be exactly evaluated.

(iii) $\tilde{g} : \mathbb{R}^n \times \mathbb{R}^m \rightarrow \mathbb{R}^l$ is $L_{\tilde{g}}$ -Lipschitz continuous and $L_{\nabla \tilde{g}}$ -smooth on $\mathcal{X} \times \mathcal{Y}$, and $\tilde{g}_i(x, \cdot)$ is convex for all $x \in \mathcal{X}$ and $i = 1, 2, \dots, l$.

(iv) The sets \mathcal{X} and \mathcal{Y} (namely, $\text{dom } f_2$ and $\text{dom } \tilde{f}_2$) are compact.

Due to the sophisticated structure specified in Assumption 1, existing methods are generally not applicable to problem (1) except the methods in [38, 39]. While the latter methods are applicable to problem (1), they do not exploit the strong convexity structure of $\tilde{f}_1(x, \cdot)$. Consequently, they may not be the most efficient methods for solving (1).

In this paper, we propose sequential minimax optimization (SMO) methods for solving problem (1) and its special case with $\tilde{g} \equiv 0$. Our approaches follow a similar framework as [38, Algorithm 2] and [39, Algorithm 1] respectively, but we enhance them by leveraging the strong concavity of $\tilde{f}_2(x, \cdot)$. As a result, our methods achieve a substantially improved operation complexity compared to [38, Algorithm 2] and [39, Algorithm 1].

The main contributions of this paper are summarized below.

- We propose an SMO method for solving an unconstrained bilevel optimization problem with strongly convex lower-level objective function. Under suitable assumptions, we show that this method achieves an operation complexity of $\mathcal{O}(\varepsilon^{-3} \log \varepsilon^{-1})$, measured by its fundamental operations, for finding an ε -KKT solution of the problem, which improves the operation complexity of [38, Algorithm 2] by a factor of $1/\varepsilon$.
- We propose an SMO method for solving a constrained bilevel optimization problem with strongly convex lower-level objective function. Under suitable assumptions, we show that this method achieves an operation complexity of $\mathcal{O}(\varepsilon^{-6} \log \varepsilon^{-1})$, measured by its fundamental operations, for finding an ε -KKT solution of the problem, which improves the operation complexity of [38, Algorithm 4] and [39, Algorithm 1] by a factor of $1/\varepsilon$.

The rest of this paper is organized as follows. In Subsection 1.1, we introduce some notation and terminology. In Section 2, we propose an SMO method for solving an unconstrained bilevel optimization problem with strongly convex lower-level objective function and study its complexity. In Section 3, we propose an SMO method for solving a constrained bilevel optimization problem with strongly convex lower-level objective function and study its complexity. Finally, we provide the proof of the main results in Section 4.

1.1 Notation and terminology

The following notation will be used throughout this paper. Let \mathbb{R}^n denote the Euclidean space of dimension n and \mathbb{R}_+^n denote the nonnegative orthant in \mathbb{R}^n . The standard inner product and Euclidean norm are denoted by $\langle \cdot, \cdot \rangle$ and $\|\cdot\|$, respectively. For any $v \in \mathbb{R}^n$, let v_+ denote the nonnegative part of v , that is, $(v_+)_i = \max\{v_i, 0\}$ for all i . For any two vectors u and v , $(u; v)$ denotes the vector resulting

from stacking v under u . Given a point x and a closed set S in \mathbb{R}^n , let $\text{dist}(x, S) = \min_{x' \in S} \|x' - x\|$ and \mathcal{I}_S denote the indicator function associated with S .

A function or mapping ϕ is said to be L_ϕ -Lipschitz continuous on a set S if $\|\phi(x) - \phi(x')\| \leq L_\phi \|x - x'\|$ for all $x, x' \in S$. In addition, it is said to be $L_{\nabla\phi}$ -smooth on S if $\|\nabla\phi(x) - \nabla\phi(x')\| \leq L_{\nabla\phi} \|x - x'\|$ for all $x, x' \in S$. A function is said to be σ -strongly-convex if it is strongly convex with modulus $\sigma > 0$. For a closed convex function $p : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{\infty\}$,² the proximal operator associated with p is denoted by prox_p , that is,

$$\text{prox}_p(x) = \arg \min_{x' \in \mathbb{R}^n} \left\{ \frac{1}{2} \|x' - x\|^2 + p(x') \right\} \quad \forall x \in \mathbb{R}^n.$$

Given that evaluation of $\text{prox}_{\gamma p}(x)$ is often as cheap as $\text{prox}_p(x)$, we count the evaluation of $\text{prox}_{\gamma p}(x)$ as one evaluation of proximal operator of p for any $\gamma > 0$ and $x \in \mathbb{R}^n$.

Finally, we introduce two types of approximate solutions for a general minimax problem

$$\Psi^* = \min_x \max_y \Psi(x, y), \quad (2)$$

where $\Psi(\cdot, y) : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{\infty\}$ is a lower semicontinuous function, $\Psi(x, \cdot) : \mathbb{R}^m \rightarrow \mathbb{R} \cup \{-\infty\}$ is an upper semicontinuous function, and Ψ^* is finite.

Definition 1. A point (x, y) is said to be a primal-dual stationary point of the minimax problem (2) if

$$0 \in \partial_x \Psi(x, y), \quad 0 \in \partial_y \Psi(x, y).$$

In addition, for any $\epsilon > 0$, a point (x_ϵ, y_ϵ) is said to be an ϵ -primal-dual stationary point of the minimax problem (2) if

$$\text{dist}(0, \partial_x \Psi(x_\epsilon, y_\epsilon)) \leq \epsilon, \quad \text{dist}(0, \partial_y \Psi(x_\epsilon, y_\epsilon)) \leq \epsilon.$$

2 Unconstrained bilevel optimization with strongly convex lower-level objective function

In this section, we consider problem (1) with $\tilde{g} \equiv 0$, which is an unconstrained bilevel optimization problem with strongly convex lower-level objective function given by

$$\begin{aligned} f^* &= \min_{x, y} f(x, y) \\ \text{s.t. } & y \in \underset{z}{\text{Argmin}} \tilde{f}(x, z). \end{aligned} \quad (3)$$

Assume that problem (3) has at least one optimal solution. In addition, f and \tilde{f} satisfy Assumption 1.

A first-order penalty method was recently proposed in [38, Algorithm 2] for solving an unconstrained minimax problem in the form of (3) in which $\tilde{f}_1(x, \cdot)$ is however merely convex for any given $x \in \mathcal{X}$. Specifically, it applies a first-order method to the minimax problem

$$\min_{x, y} \max_z \left\{ P_\rho(x, y, z) := f(x, y) + \rho(\tilde{f}(x, y) - \tilde{f}(x, z)) \right\} \quad (4)$$

for a suitably chosen $\rho > 0$. While this method is applicable to problem (3), it does not exploit the strong convexity structure of $\tilde{f}_1(x, \cdot)$. Besides, it solves a single minimax problem with a large ρ , which is typically an ill-conditioned problem. Consequently, this method may not be the most efficient method for solving (3).

In what follows, we propose a sequential minimax optimization (SMO) method for solving problem (3). Our approach solves a sequence of nonconvex-strongly-concave minimax problems in the form of (4) using the recently developed first-order method [37, Algorithm 1] that leverages the strong convexity of $\tilde{f}_1(x, \cdot)$. As a result, our method achieves a substantially improved operation complexity compared to [38, Algorithm 2].

²For convenience, ∞ stands for $+\infty$.

Algorithm 1 A sequential minimax optimization method for problem (3)

Input: $\varepsilon, \tau \in (0, 1]$, $\epsilon_0 \in (\tau\varepsilon, 1]$, $\epsilon_k = \epsilon_0\tau^k$, $\rho_k = \epsilon_k^{-1}$, $(x^0, y^0) \in \mathcal{X} \times \mathcal{Y}$ with $\tilde{f}(x^0, y^0) \leq \min_y \tilde{f}(x^0, y) + \tau\varepsilon$.

- 1: **for** $k = 0, 1, \dots$ **do**
- 2: Call Algorithm 5 in Appendix B with $\epsilon \leftarrow \epsilon_k$, $\tilde{\epsilon}_0 \leftarrow \epsilon_k/2$, $\tilde{x}^0 \leftarrow (x^k, y^0)$, $\tilde{y}^0 \leftarrow y^0$, $\sigma_y \leftarrow \sigma\rho_k$ and $L_{\nabla h} \leftarrow L_{\nabla f_1} + 2\epsilon_k^{-1}L_{\nabla \tilde{f}_1}$ to find an ϵ_k -primal-dual stationary point $(x^{k+1}, y^{k+1}, z^{k+1})$ of problem (4) with $\rho = \rho_k$.
- 3: If $\epsilon_k \leq \varepsilon$, terminate the algorithm and output (x^{k+1}, y^{k+1}) .
- 4: **end for**

Remark 1. The initial point (x^0, y^0) of Algorithm 1 can be found by an additional procedure. Indeed, one can first choose any $x^0 \in \mathcal{X}$ and then apply an accelerated proximal gradient method [46] to the problem $\min_y \tilde{f}(x^0, y)$ for finding $y^0 \in \mathcal{Y}$ such that $\tilde{f}(x^0, y^0) \leq \min_y \tilde{f}(x^0, y) + \tau\varepsilon$.

To characterize the approximate solution found by Algorithm 1, we review a terminology called an ε -KKT solution of problem (3), which was introduced in [38, Definition 3].

Definition 2. For any $\varepsilon > 0$, (x, y) is said to be an ε -KKT solution of problem (3) if there exists $(z, \rho) \in \mathbb{R}^m \times \mathbb{R}_+$ such that

$$\begin{aligned} \text{dist}\left(0, \partial f(x, y) + \rho \partial \tilde{f}(x, y) - (\rho \nabla_x \tilde{f}(x, z); 0)\right) &\leq \varepsilon, \quad \text{dist}(0, \rho \partial_z \tilde{f}(x, z)) \leq \varepsilon, \\ \tilde{f}(x, y) - \min_{z'} \tilde{f}(x, z') &\leq \varepsilon. \end{aligned}$$

We next study operation complexity of Algorithm 1. To proceed, we define

$$D_{\mathbf{x}} := \max\{\|u - v\| \mid u, v \in \mathcal{X}\}, \quad D_{\mathbf{y}} := \max\{\|u - v\| \mid u, v \in \mathcal{Y}\}, \quad (5)$$

$$\tilde{f}_{\text{hi}} := \max\{\tilde{f}(x, y) \mid (x, y) \in \mathcal{X} \times \mathcal{Y}\}, \quad \tilde{f}_{\text{low}} := \min\{\tilde{f}(x, y) \mid (x, y) \in \mathcal{X} \times \mathcal{Y}\}, \quad (6)$$

$$f_{\text{low}} := \min\{f(x, y) \mid (x, y) \in \mathcal{X} \times \mathcal{Y}\}, \quad (7)$$

$$\widehat{K} := \lceil (\log \varepsilon - \log \epsilon_0) / \log \tau \rceil_+, \quad \widehat{\mathbb{K}} := \{0, 1, \dots, \widehat{K} + 1\}, \quad \widehat{\mathbb{K}} - 1 = \{k - 1 \mid k \in \widehat{\mathbb{K}}\}. \quad (8)$$

By Assumption 1, one can observe that $D_{\mathbf{x}}$, $D_{\mathbf{y}}$, \tilde{f}_{hi} , \tilde{f}_{low} and f_{low} are finite.

We are now ready to present a theorem regarding *operation complexity* of Algorithm 1, measured by the amount of evaluations of ∇f_1 , $\nabla \tilde{f}_1$ and proximal operator of f_2 and \tilde{f}_2 , for finding an $\mathcal{O}(\varepsilon)$ -KKT solution of (3), whose proof is deferred to Subsection 4.1.

Theorem 1 (Complexity of Algorithm 1). Suppose that Assumption 1 holds. Let f^* , f , \tilde{f} , $D_{\mathbf{x}}$, $D_{\mathbf{y}}$, \tilde{f}_{hi} , \tilde{f}_{low} , f_{low} and \widehat{K} be defined in (3), (5), (6), (7) and (8), σ , $L_{\nabla f_1}$ and $L_{\nabla \tilde{f}_1}$ be given in Assumption 1, ε , τ , $\{\rho_k\}$, x^0 and y^0 be given in Algorithm 1, and

$$\widehat{L} = L_{\nabla f_1} + 2L_{\nabla \tilde{f}_1}, \quad \hat{\alpha} = \min\left\{1, \sqrt{8\sigma/\widehat{L}}\right\}, \quad (9)$$

$$\hat{\delta} = (2 + \hat{\alpha}^{-1})(D_{\mathbf{x}}^2 + D_{\mathbf{y}}^2)\widehat{L} + \max\left\{2\sigma, \widehat{L}/4\right\} D_{\mathbf{y}}^2, \quad (10)$$

$$\widehat{M} = \frac{16 \max\left\{\frac{1}{4L_{\nabla \tilde{f}_1}}, \frac{2}{\hat{\alpha}L_{\nabla \tilde{f}_1}}\right\}}{\left[9\widehat{L}^2 / \min\{2L_{\nabla \tilde{f}_1}, \sigma\} + 3\widehat{L}\right]^{-2}} \left(\hat{\delta} + 2\hat{\alpha}^{-1}(f^* - f_{\text{low}} + \tilde{f}_{\text{hi}} - \tilde{f}_{\text{low}} + \widehat{L}(D_{\mathbf{x}}^2 + D_{\mathbf{y}}^2))\right), \quad (11)$$

$$\widehat{T} = \left[16(1 + f(x^0, y^0) - f_{\text{low}})\widehat{L} + 8\sigma^{-2}\widehat{L}^2 + 7\right]_+, \quad (12)$$

$$\begin{aligned} \widehat{N} &= 3397 \max\left\{2, \sqrt{\widehat{L}/(2\sigma)}\right\} \widehat{T} \epsilon_0^{-3} (1 - \tau^3)^{-1} \\ &\quad \times (\tau\varepsilon/\epsilon_0)^{-3} \left(14\widehat{K} \log(1/\tau) + 14 \log(1/\epsilon_0) + 2(\log \widehat{M})_+ + 2 + 2 \log(2\widehat{T})\right). \end{aligned} \quad (13)$$

Then Algorithm 1 outputs an approximate solution $(x^{\hat{K}+1}, y^{\hat{K}+1})$ of problem (3) satisfying

$$\text{dist}\left(0, \partial f(x^{\hat{K}+1}, y^{\hat{K}+1}) + \rho_{\hat{K}} \partial \tilde{f}(x^{\hat{K}+1}, y^{\hat{K}+1}) - (\rho_{\hat{K}} \nabla_x \tilde{f}(x^{\hat{K}+1}, z^{\hat{K}+1}); 0)\right) \leq \varepsilon, \quad (14)$$

$$\text{dist}\left(0, \rho_{\hat{K}} \partial \tilde{f}(x^{\hat{K}+1}, z^{\hat{K}+1})\right) \leq \varepsilon, \quad (15)$$

$$\tilde{f}(x^{\hat{K}+1}, y^{\hat{K}+1}) \leq \min_z \tilde{f}(x^{\hat{K}+1}, z) + \varepsilon \left(1 + f(x^0, y^0) - f_{\text{low}} + \varepsilon^2 (L_{\nabla f_1}^{-1} + \sigma^{-2} \hat{L})/2\right), \quad (16)$$

after at most \hat{N} evaluations of ∇f_1 , $\nabla \tilde{f}_1$ and proximal operator of f_2 and \tilde{f}_2 , respectively.

Remark 2. (i) One can observe from Theorem 1 that $\hat{N} = \mathcal{O}(\varepsilon^{-3} \log \varepsilon^{-1})$. As a result, Algorithm 1 enjoys an operation complexity of $\mathcal{O}(\varepsilon^{-3} \log \varepsilon^{-1})$, measured by the amount of evaluations of ∇f_1 , $\nabla \tilde{f}_1$ and proximal operator of f_2 and \tilde{f}_2 , for finding an $\mathcal{O}(\varepsilon)$ -KKT solution $(x^{\hat{K}+1}, y^{\hat{K}+1})$ of (3) satisfying

$$\begin{aligned} \text{dist}\left(0, \partial f(x^{\hat{K}+1}, y^{\hat{K}+1}) + \rho_{\hat{K}} \partial \tilde{f}(x^{\hat{K}+1}, y^{\hat{K}+1}) - (\rho_{\hat{K}} \nabla_x \tilde{f}(x^{\hat{K}+1}, z^{\hat{K}+1}); 0)\right) &\leq \varepsilon, \\ \text{dist}\left(0, \rho_{\hat{K}} \partial \tilde{f}(x^{\hat{K}+1}, z^{\hat{K}+1})\right) &\leq \varepsilon, \quad \tilde{f}(x^{\hat{K}+1}, y^{\hat{K}+1}) - \min_z \tilde{f}(x^{\hat{K}+1}, z) = \mathcal{O}(\varepsilon), \end{aligned}$$

where $z^{\hat{K}+1}$ is given in Algorithm 1 and $\rho_{\hat{K}} = \epsilon_0^{-1} \tau^{-\hat{K}}$.

(ii) It shall be mentioned that an $\mathcal{O}(\varepsilon)$ -KKT solution of (3) can be found by [38, Algorithm 2] with an operation complexity of $\mathcal{O}(\varepsilon^{-4} \log \varepsilon^{-1})$ (see [38, Theorem 2]). As a result, the operation complexity of Algorithm 1 improves that of [38, Algorithm 2] by a factor of $1/\epsilon$.

3 Constrained bilevel optimization with strongly convex lower-level objective function

In this section, we consider problem (1), which is a constrained bilevel optimization problem with strongly convex lower-level objective function satisfying Assumption 1. Recall from Assumption 1 that $\mathcal{X} = \text{dom } f_2$ and $\mathcal{Y} = \text{dom } \tilde{f}_2$. We now make some additional assumptions for problem (1).

Assumption 2. (i) f and \tilde{f} are L_f - and $L_{\tilde{f}}$ -Lipschitz continuous on $\mathcal{X} \times \mathcal{Y}$, respectively.

(ii) $\tilde{g}_i(x, \cdot)$ is convex and there exists $\hat{z}_x \in \mathcal{Y}$ for each $x \in \mathcal{X}$ such that $\tilde{g}_i(x, \hat{z}_x) < 0$ for all $i = 1, 2, \dots, l$ and $G := \min\{-\tilde{g}_i(x, \hat{z}_x) | x \in \mathcal{X}, i = 1, \dots, l\} > 0$.³

First-order penalty type of methods were recently proposed in [38, Algorithm 4] and [39, Algorithm 1] for solving constrained bilevel optimization in the form of (1) in which $f_1(x, \cdot)$ is however merely convex for any given $x \in \mathcal{X}$. Specifically, [38, Algorithm 4] solves constrained bilevel optimization by applying a first-order method to the minimax problem

$$\min_{x, y} \max_z P_{\rho, \mu}(x, y, z)$$

for suitably chosen $\rho, \mu > 0$, where

$$\begin{aligned} P_{\rho, \mu}(x, y, z) &= (f_1(x, y) + \rho \tilde{f}_1(x, y) + \rho \mu \|\tilde{g}(x, y)_+\|^2 - \rho \tilde{f}_1(x, z) - \rho \mu \|\tilde{g}(x, z)_+\|^2) \\ &\quad + (f_2(x) + \rho \tilde{f}_2(y) - \rho \tilde{f}_2(z)). \end{aligned}$$

In addition, [39, Algorithm 1] solves constrained bilevel optimization by applying a first-order method to a sequence of minimax problems in the form of

$$\min_{x, y} \max_z \mathcal{L}(x, y, z, \lambda; \rho, \mu), \quad (17)$$

³The latter part of this assumption can be weakened to the one that the pointwise Slater's condition holds for the lower-level part of (1), that is, there exists $\hat{z}_x \in \mathcal{Y}$ such that $\tilde{g}(x, \hat{z}_x) < 0$ for each $x \in \mathcal{X}$. Indeed, if $G > 0$, Assumption 2(iii) clearly holds. Otherwise, one can solve the perturbed counterpart of (1) with $\tilde{g}(x, z)$ being replaced by $\tilde{g}(x, z) - \epsilon$ for some suitable $\epsilon > 0$ instead, which satisfies Assumption 2(iii).

where \mathcal{L} is defined as

$$\mathcal{L}(x, y, z, \lambda; \rho, \mu) = f(x, y) + \rho \tilde{f}(x, y) + \frac{1}{2\mu} \|\lambda + \mu \tilde{g}(x, y)\|_+^2 - \rho \tilde{f}(x, z) - \frac{1}{2\mu} \|\lambda + \mu \tilde{g}(x, z)\|_+^2. \quad (18)$$

While these methods are applicable to problem (3), they do not exploit the strong convexity structure of $\tilde{f}_1(x, \cdot)$. Consequently, these methods may not be the most efficient methods for solving (3).

In what follows, we propose a sequential minimax optimization (SMO) method for solving problem (1). Our approach follows a similar framework as [39, Algorithm 3]. Specifically, at each iteration, our method finds an approximate primal-dual stationary point of a minimax problem in the form of (17). Yet, given that this problem is a nonconvex-strongly-concave unconstrained minimax problem, we solve it using the recently developed first-order method [37, Algorithm 1] that leverages the strong convexity of $\tilde{f}_1(x, \cdot)$. As a result, our method achieves a substantially improved operation complexity compared to [38, Algorithm 4] and [39, Algorithm 1].

To present our method, we define

$$\tilde{\mathcal{L}}(x, z, \lambda; \rho, \mu) := \tilde{f}(x, z) + \frac{1}{2\rho\mu} \|\lambda + \mu \tilde{g}(x, z)\|_+^2. \quad (19)$$

We are now ready to present our method for solving problem (1).

Algorithm 2 A sequential minimax optimization method for problem (1)

Input: $\varepsilon, \tau \in (0, 1)$, $\epsilon_0 \in (\tau\varepsilon, 1]$, $x^0 \in \mathcal{X}$, $z^0 \in \mathcal{Y}$, $\epsilon_k = \epsilon_0 \tau^k$, $\rho_k = \epsilon_k^{-1}$, $\mu_k = \epsilon_k^{-3}$, $\eta_k = \epsilon_k$ and $\lambda^0 \in \mathbb{R}_+^l$.

1: **for** $k = 0, 1 \dots$ **do**

2: Call Algorithm 3 (see Appendix A) with $\Psi(\cdot) \leftarrow \tilde{\mathcal{L}}(x^k, \cdot, \lambda^k; \rho_k, \mu_k)$, $\tilde{\epsilon} \leftarrow \eta_k / D_{\mathbf{y}}$, $\sigma_\phi \leftarrow \sigma$, $L_{\nabla\phi} \leftarrow \tilde{L}_k$, $\tilde{x}^0 \leftarrow y^k$ to find an approximate solution y_{init}^k of $\min_z \tilde{L}(x^k, z, \lambda^k; \rho_k, \mu_k)$ such that

$$\tilde{\mathcal{L}}(x^k, y_{\text{init}}^k, \lambda^k; \rho_k, \mu_k) - \min_z \tilde{\mathcal{L}}(x^k, z, \lambda^k; \rho_k, \mu_k) \leq \eta_k, \quad (20)$$

where $\tilde{\mathcal{L}}$ is given in (19) and

$$\tilde{L}_k = L_{\nabla\tilde{f}_1} + \rho_k^{-1}(\mu_k L_{\tilde{g}}^2 + \mu_k \tilde{g}_{\text{hi}} L_{\nabla\tilde{g}} + \|\lambda^k\| L_{\nabla\tilde{g}}). \quad (21)$$

3: Call Algorithm 5 (see Appendix B) with $\epsilon \leftarrow \epsilon_k$, $\tilde{\epsilon}_0 \leftarrow \epsilon_k / 2$, $\tilde{x}^0 \leftarrow (x^k, y_{\text{init}}^k)$, $\tilde{y}^0 \leftarrow z^k$, $\sigma_y \leftarrow \sigma \rho_k$, and $L_{\nabla h} \leftarrow L_k$ to find an ϵ_k -primal-dual stationary point $(x^{k+1}, y^{k+1}, z^{k+1})$ of

$$\min_{x, y} \max_z \mathcal{L}(x, y, z, \lambda^k; \rho_k, \mu_k), \quad (22)$$

where

$$L_k = L_{\nabla f_1} + 2\rho_k L_{\nabla\tilde{f}_1} + 2\mu_k L_{\tilde{g}}^2 + 2\mu_k \tilde{g}_{\text{hi}} L_{\nabla\tilde{g}} + 2\|\lambda^k\| L_{\nabla\tilde{g}}. \quad (23)$$

4: Set $\lambda^{k+1} = [\lambda^k + \mu_k \tilde{g}(x^{k+1}, z^{k+1})]_+$.

5: If $\epsilon_k \leq \varepsilon$, terminate the algorithm and output (x^{k+1}, y^{k+1}) .

6: **end for**

Remark 3. (i) Notice that $\tilde{L}(x^k, y, \lambda^k; \rho_k, \mu_k) = \phi(y) + \tilde{f}_2(y)$ with $\phi(y) = \tilde{f}_1(x^k, y) + \|\lambda^k + \mu_k \tilde{g}(x^k, y)\|_+^2 / (2\rho_k \mu_k)$.

By Assumptions 1 and 2, Theorem 3 (see Appendix A), and (90), it is not hard to observe that y_{init}^k satisfying (20) can be successfully found in step 2 of Algorithm 2 by applying Algorithm 3 to the problem $\min_z \tilde{L}(x^k, z, \lambda^k; \rho_k, \mu_k)$.

(ii) In view of Theorem 4 (see Appendix B), one can see that an ϵ_k -primal-dual stationary point of (22) can be successfully found in step 3 of Algorithm 2 by applying Algorithm 5 to problem (22). Consequently, Algorithm 2 is well-defined.

We next study *iteration and operation complexity* of Algorithm 2. To characterize the approximate solution found by Algorithm 2, we first review a terminology called an ε -KKT solution of problem (1), which was introduced in [38, Definition 4].

Definition 3. For any $\varepsilon > 0$, (x, y) is said to be an ε -KKT solution of problem (1) if there exists $(z, \rho, \lambda_{\mathbf{y}}, \lambda_{\mathbf{z}}) \in \mathbb{R}^m \times \mathbb{R}_+ \times \mathbb{R}_+^l \times \mathbb{R}_+^l$ such that

$$\begin{aligned} \text{dist} \left(0, \partial f(x, y) + \rho \partial \tilde{f}(x, y) - \rho (\nabla_x \tilde{f}(x, z) + \nabla_x \tilde{g}(x, z) \lambda_{\mathbf{z}}; 0) + \nabla \tilde{g}(x, y) \lambda_{\mathbf{y}} \right) &\leq \varepsilon, \\ \text{dist} \left(0, \rho (\partial_z \tilde{f}(x, z) + \nabla_z \tilde{g}(x, z) \lambda_{\mathbf{z}}) \right) &\leq \varepsilon, \\ \|\tilde{g}(x, z)\|_+ &\leq \varepsilon, \quad |\langle \lambda_{\mathbf{z}}, \tilde{g}(x, z) \rangle| \leq \varepsilon, \\ |\tilde{f}(x, y) - \tilde{f}^*(x)| &\leq \varepsilon, \quad \|\tilde{g}(x, y)\|_+ \leq \varepsilon, \quad |\langle \lambda_{\mathbf{y}}, \tilde{g}(x, y) \rangle| \leq \varepsilon, \end{aligned}$$

where

$$\tilde{f}^*(x) := \min\{\tilde{f}(x, z) | \tilde{g}(x, z) \leq 0\}. \quad (24)$$

To proceed, recall that $\mathcal{X} = \text{dom } f_2$ and $\mathcal{Y} = \text{dom } \tilde{f}_2$. We define

$$\tilde{f}_{\text{hi}}^* := \sup\{\tilde{f}^*(x) | x \in \mathcal{X}\}, \quad (25)$$

$$f_{\text{hi}} := \max\{f(x, y) | (x, y) \in \mathcal{X} \times \mathcal{Y}\}, \quad \tilde{g}_{\text{hi}} := \max\{\|\tilde{g}(x, y)\| | (x, y) \in \mathcal{X} \times \mathcal{Y}\}, \quad (26)$$

$$K := \lceil (\log \varepsilon - \log \epsilon_0) / \log \tau \rceil_+, \quad \mathbb{K} := \{0, 1, \dots, K+1\}, \quad \mathbb{K} - 1 = \{k-1 | k \in \mathbb{K}\}. \quad (27)$$

It then follows from Assumption 1(iii) that

$$\|\nabla \tilde{g}(x, y)\| \leq L_{\tilde{g}} \quad \forall (x, y) \in \mathcal{X} \times \mathcal{Y}. \quad (28)$$

In addition, by Assumption 1 and the compactness of \mathcal{X} and \mathcal{Y} , one can observe that f_{hi} and \tilde{g}_{hi} are finite. Besides, \tilde{f}_{hi}^* is also finite (see Lemma 3).

We are now ready to present an *iteration and operation complexity* of Algorithm 2, measured by the amount of evaluations of ∇f_1 , $\nabla \tilde{f}_1$, $\nabla \tilde{g}$ and proximal operator of f_2 and \tilde{f}_2 , for finding an $\mathcal{O}(\varepsilon)$ -KKT solution of (1), whose proof is deferred to Subsection 4.2.

Theorem 2 (iteration and operation complexity of Algorithm 2). Suppose that Assumptions 1 and 2 hold. Let $\{(x_k, y_k, z_k, \lambda^k)\}_{k \in \mathbb{K}}$ be generated by Algorithm 2, $D_{\mathbf{x}}$, $D_{\mathbf{y}}$, f_{low} , \tilde{f}_{low} , f^* , \tilde{f}_{hi}^* , f_{hi} , \tilde{g}_{hi} and K be defined in (5), (6), (7), (1), (25), (26) and (27), σ , L_f , $L_{\tilde{f}}$, $L_{\nabla f_1}$, $L_{\nabla \tilde{f}_1}$, $L_{\tilde{g}}$, $L_{\nabla \tilde{g}}$ and G be given in Assumptions 1 and 2, and ε , ϵ_0 , τ , μ_K , ρ_K and λ_0 be given in Algorithm 2. Let

$$\vartheta = \frac{1}{2} \|\lambda^0\|^2 + \frac{\tilde{f}_{\text{hi}}^* - \tilde{f}_{\text{low}}}{1 - \tau^4} + \frac{D_{\mathbf{y}} \epsilon_0}{1 - \tau^3}, \quad (29)$$

$$L = L_{\nabla f_1} + 2L_{\nabla \tilde{f}_1} + 2L_{\tilde{g}}^2 + 2\tilde{g}_{\text{hi}} L_{\nabla \tilde{g}} + 2\sqrt{2\vartheta} L_{\nabla \tilde{g}}, \quad \tilde{L} = L_{\nabla \tilde{f}_1} + L_{\tilde{g}}^2 + \tilde{g}_{\text{hi}} L_{\nabla \tilde{g}} + \sqrt{2\vartheta} L_{\nabla \tilde{g}}, \quad (30)$$

$$\alpha = \min\{1, \sqrt{8\sigma/L}\}, \quad \delta = (2 + \alpha^{-1})(D_{\mathbf{x}}^2 + D_{\mathbf{y}}^2)L + \max\{1/D_{\mathbf{y}}, L/4\}D_{\mathbf{y}}^2, \quad (31)$$

$$\begin{aligned} M &= 16 \max\{1/(4L_{\tilde{g}}^2), 2/(\alpha L_{\tilde{g}}^2)\} [9L^2 / \min\{2L_{\tilde{g}}^2, \sigma\} + 3L]^2 \\ &\quad \times \left(\delta + 2\alpha^{-1} \left(f^* - f_{\text{low}} + \tilde{f}_{\text{hi}}^* - \tilde{f}_{\text{low}} + L_{\tilde{f}} D_{\mathbf{y}} + 3\vartheta + \tilde{g}_{\text{hi}}^2 + L(D_{\mathbf{x}}^2 + D_{\mathbf{y}}^2) \right) \right), \end{aligned} \quad (32)$$

$$T = \left\lceil 16(f_{\text{hi}} - f_{\text{low}} + 1)L + 8(1 + \sigma^{-2}L^2) \right\rceil_+, \quad (33)$$

$$\lambda_{\mathbf{y}}^{K+1} = [\lambda_{\mathbf{y}}^K + \mu_K \tilde{g}(x^{K+1}, y^{K+1})]_+, \quad \lambda_{\mathbf{z}}^{K+1} = \rho_K^{-1} [\lambda_{\mathbf{z}}^K + \mu_K \tilde{g}(x^{K+1}, z^{K+1})]_+. \quad (34)$$

Suppose that $\varepsilon^{-2} - 8\tau^{-3}G^{-2}\vartheta \geq 0$. Then the following statements hold.

- (i) Algorithm 2 terminates after $K+1$ outer iterations and outputs an approximate point (x^{K+1}, y^{K+1}) of problem (1) satisfying

$$\begin{aligned} \text{dist} \left(0, \partial f(x^{K+1}, y^{K+1}) + \rho_K \partial \tilde{f}(x^{K+1}, y^{K+1}) - \rho_K (\nabla_x \tilde{f}(x^{K+1}, z^{K+1}) + \nabla_x \tilde{g}(x^{K+1}, z^{K+1}) \lambda_{\mathbf{z}}^{K+1}; 0) \right. \\ \left. + \nabla \tilde{g}(x^{K+1}, y^{K+1}) \lambda_{\mathbf{y}}^{K+1} \right) &\leq \varepsilon, \end{aligned} \quad (35)$$

$$\text{dist} \left(0, \rho_K (\partial_z \tilde{f}(x^{K+1}, z^{K+1}) + \nabla_z \tilde{g}(x^{K+1}, z^{K+1}) \lambda_{\mathbf{z}}^{K+1}) \right) \leq \varepsilon, \quad (36)$$

$$\|\tilde{g}(x^{K+1}, z^{K+1})\|_+ \leq 2\varepsilon^2 G^{-1}(\epsilon_0 + L_{\tilde{f}}) D_{\mathbf{y}}, \quad (37)$$

$$|\langle \lambda_{\mathbf{z}}^{K+1}, \tilde{g}(x^{K+1}, z^{K+1}) \rangle| \leq 2\varepsilon^2 G^{-1}(\epsilon_0 + L_{\tilde{f}}) D_{\mathbf{y}} \max\{\|\lambda^0\|, 2G^{-1}(\epsilon_0 + L_{\tilde{f}}) D_{\mathbf{y}}\}, \quad (38)$$

$$\|[\tilde{g}(x^{K+1}, y^{K+1})]_+\| \leq 2\varepsilon^2 G^{-1}(\epsilon_0 + L_f + L_{\tilde{f}})D_{\mathbf{y}}, \quad (39)$$

$$|\langle \lambda_{\mathbf{y}}^{K+1}, \tilde{g}(x^{K+1}, z^{K+1}) \rangle| \leq 2\varepsilon G^{-1}(\epsilon_0 + L_f + L_{\tilde{f}})D_{\mathbf{y}} \max\{\|\lambda^0\|, 2G^{-1}(\epsilon_0 + L_f + L_{\tilde{f}})D_{\mathbf{y}}\}, \quad (40)$$

$$|\tilde{f}(x^{K+1}, y^{K+1}) - \tilde{f}^*(x^{K+1})| \leq \max\left\{2\varepsilon^2 G^{-2}L_{\tilde{f}}(\epsilon_0 + L_f + L_{\tilde{f}})D_{\mathbf{y}}^2, \varepsilon^3 \max\{\|\lambda^0\|, 2G^{-1}(\epsilon_0 + L_{\tilde{f}})D_{\mathbf{y}}\}/2\right. \\ \left. + \varepsilon \left(f_{\text{hi}} - f_{\text{low}} + 1 + L_{\tilde{g}}^{-2}/4 + \sigma^{-2}L/2\right)\right\}. \quad (41)$$

(ii) The total number of evaluations of ∇f_1 , $\nabla \tilde{f}_1$, $\nabla \tilde{g}$ and proximal operator of f_2 and \tilde{f}_2 performed in Algorithm 2 is N , respectively, where

$$N = 3397 \max\left\{2, \sqrt{L/(2\sigma)}\right\} T(1 - \tau^6)^{-1} \\ \times (\tau\varepsilon)^{-6} (38K \log(1/\tau) + 38 \log(1/\epsilon_0) + 2(\log M)_+ + 2 + 2 \log(2T)) \\ + 2(\tau\varepsilon)^{-1}(1 - \tau) \left\lceil \sqrt{\frac{\tilde{L}}{\sigma}} + 1 \right\rceil \max\left\{1, \left\lceil 2 \log(2\tilde{L}D_{\mathbf{y}}^2) + 6K \log(1/\tau) - 6 \log \epsilon_0 \right\rceil\right\} + K. \quad (42)$$

Remark 4. (i) One can observe from Theorem 2 that Algorithm 2 enjoys an iteration complexity of $\mathcal{O}(\log \varepsilon^{-1})$ and an operation complexity of $\mathcal{O}(\varepsilon^{-6} \log \varepsilon^{-1})$, measured by the amount of evaluations of ∇f_1 , $\nabla \tilde{f}_1$, $\nabla \tilde{g}$ and proximal operator of f_2 and \tilde{f}_2 , for finding an $\mathcal{O}(\varepsilon)$ -KKT solution (x^{K+1}, y^{K+1}) of (1) such that

$$\text{dist}\left(0, \partial f(x^{K+1}, y^{K+1}) + \rho_K \partial \tilde{f}(x^{K+1}, y^{K+1}) + \nabla \tilde{g}(x^{K+1}, y^{K+1})\lambda_{\mathbf{y}}^{K+1} \right. \\ \left. - \rho_K(\nabla_x \tilde{f}(x^{K+1}, z^{K+1}) + \nabla_x \tilde{g}(x^{K+1}, z^{K+1})\tilde{\lambda}_{\mathbf{z}}^{K+1}; 0)\right) \leq \varepsilon, \\ \text{dist}\left(0, \rho_K(\partial_z \tilde{f}(x^{K+1}, z^{K+1}) + \nabla_z \tilde{g}(x^{K+1}, z^{K+1})\lambda_{\mathbf{z}}^{K+1})\right) \leq \varepsilon, \\ \|[\tilde{g}(x^{K+1}, z^{K+1})]_+\| = \mathcal{O}(\varepsilon^2), \quad |\langle \lambda_{\mathbf{z}}^{K+1}, \tilde{g}(x^{K+1}, z^{K+1}) \rangle| = \mathcal{O}(\varepsilon^2), \\ \|[\tilde{g}(x^{K+1}, y^{K+1})]_+\| = \mathcal{O}(\varepsilon^2), \quad |\langle \lambda_{\mathbf{y}}^{K+1}, \tilde{g}(x^{K+1}, y^{K+1}) \rangle| = \mathcal{O}(\varepsilon), \\ |\tilde{f}(x^{K+1}, y^{K+1}) - \tilde{f}^*(x^{K+1})| = \mathcal{O}(\varepsilon),$$

where \tilde{f}^* is defined in (24), $\rho_K = (\epsilon_0 \tau^K)^{-1}$, and $\lambda_{\mathbf{y}}^{K+1}, \lambda_{\mathbf{z}}^{K+1} \in \mathbb{R}_+^l$ are given in (34).

(ii) It shall be mentioned that an $\mathcal{O}(\varepsilon)$ -KKT solution of (1) can be found by [38, Algorithm 4] and [39, Algorithm 1] with an operation complexity of $\mathcal{O}(\varepsilon^{-4} \log \varepsilon^{-1})$ (see [38, Theorem 4] and [39, Theorem 1]). As a result, the operation complexity of Algorithm 2 improves that of [38, Algorithm 4] and [39, Algorithm 1] by a factor of $1/\varepsilon$.

4 Proof of the main results

In this section we provide a proof of our main results presented in Sections 2 and 3, which are particularly Theorems 1 and 2.

4.1 Proof of the main results in Section 2

In this subsection we prove Theorem 1. We first establish two lemmas below.

Lemma 1. Suppose that Assumption 1 holds. Let f^* , f , \tilde{f} , $D_{\mathbf{x}}$, $D_{\mathbf{y}}$, \tilde{f}_{hi} , \tilde{f}_{low} and f_{low} be defined in (3), (5), (6) and (7), σ , $L_{\nabla f_1}$ and $L_{\nabla \tilde{f}_1}$ be given in Assumption 1, ϵ_k , ρ_k , x^0 , y^0 and z_ϵ be given in

Algorithm 1, and

$$\widehat{L}_k = L_{\nabla f_1} + 2\rho_k L_{\nabla \tilde{f}_1}, \quad \hat{\alpha}_k = \min \left\{ 1, \sqrt{8\sigma\rho_k/\widehat{L}_k} \right\}, \quad (43)$$

$$\hat{\delta}_k = (2 + \hat{\alpha}_k^{-1})(D_{\mathbf{x}}^2 + D_{\mathbf{y}}^2)\widehat{L}_k + \max \left\{ 2\sigma\rho_k, \hat{\alpha}_k\widehat{L}_k/4 \right\} D_{\mathbf{y}}^2, \quad (44)$$

$$\widehat{M}_k = \frac{16 \max \left\{ \frac{1}{2\widehat{L}_k}, \min \left\{ \frac{1}{2\sigma\rho_k}, \frac{4}{\hat{\alpha}_k\widehat{L}_k} \right\} \right\} \left[\hat{\delta}_k + 2\hat{\alpha}_k^{-1}(f^* - f_{\text{low}} + \rho_k(\tilde{f}_{\text{hi}} - \tilde{f}_{\text{low}}) + \widehat{L}_k(D_{\mathbf{x}}^2 + D_{\mathbf{y}}^2)) \right]}{\left[9\widehat{L}_k^2 / \min\{\widehat{L}_k, \sigma\rho_k\} + 3\widehat{L}_k \right]^{-2} \epsilon_k^2}, \quad (45)$$

$$\widehat{T}_k = \left\lceil 16(1 + f(x^0, y^0) - f_{\text{low}})\widehat{L}_k\epsilon_k^{-2} + 8\sigma^{-2}\rho_k^{-2}\widehat{L}_k^2 + 7 \right\rceil, \quad (46)$$

$$\widehat{N}_k = 3397 \max \left\{ 2, \sqrt{\widehat{L}_k/(2\sigma\rho_k)} \right\} \left((\widehat{T}_k + 1)(\log \widehat{C}_k)_+ + \widehat{T}_k + 1 + 2\widehat{T}_k \log(\widehat{T}_k + 1) \right). \quad (47)$$

Then for all $0 \leq k \in \mathbb{K} - 1$, an ϵ_k -primal-dual stationary point $(x^{k+1}, y^{k+1}, z^{k+1})$ of problem (4) is successfully found at step 2 of Algorithm 1 that satisfies

$$P_{\rho_k}(x^{k+1}, y^{k+1}, z^{k+1}) \leq 1 + f(x^0, y^0) + \epsilon_k^2(\widehat{L}_k^{-1} + \sigma^{-2}\rho_k^{-2}\widehat{L}_k)/2. \quad (48)$$

Moreover, the total number of evaluations of ∇f_1 , $\nabla \tilde{f}_1$ and proximal operator of f_2 and \tilde{f}_2 performed at step 2 in iteration k of Algorithm 1 is no more than \widehat{N}_k .

Proof. Observe that problem (4) with $\rho = \rho_k$ can be viewed as

$$\min_{x,y} \max_z \{ P_{\rho_k}(x, y, z) = h(x, y, z) + p(x, y) - q(z) \},$$

where $h(x, y, z) = f_1(x, y) + \rho_k \tilde{f}_1(x, y) - \rho_k \tilde{f}_1(x, z)$, $p(x, y) = f_2(x) + \rho_k \tilde{f}_2(y)$, and $q(z) = \rho_k \tilde{f}_2(z)$. Hence, problem (4) is in the form of (91) with $H = P_{\rho_k}$. By Assumption 1 and $\rho_k = \epsilon_k^{-1}$, one can see that $h(x, y, \cdot)$ is $\sigma\rho_k$ -strongly-convex on \mathcal{Y} , and h is \widehat{L}_k -smooth on its domain, where \widehat{L}_k is given in (43). Also, notice from Algorithm 1 that $\tilde{\epsilon}_0 = \epsilon_k/2$. Consequently, Algorithm 5 can be suitably applied to problem (4) with $\rho = \rho_k$ for finding an ϵ_k -stationary point $(x^{k+1}, y^{k+1}, z^{k+1})$ of it.

In addition, notice from Algorithm 1 that $\tilde{f}(x^0, y^0) \leq \min_y \tilde{f}(x^0, y) + \tau\epsilon$ and $\tau\epsilon \leq \epsilon_k$ for all $0 \leq k \in \mathbb{K} - 1$. Using this, (4), $\tau\epsilon \leq \epsilon_k$ and $\rho_k = \epsilon_k^{-1}$, we obtain

$$\max_z P_{\rho_k}(x^0, y^0, z) = f(x^0, y^0) + \rho_k(\tilde{f}(x^0, y^0) - \min_z \tilde{f}(x^0, z)) \leq f(x^0, y^0) + \rho_k\epsilon_k = f(x^0, y^0) + 1. \quad (49)$$

By this and (96) with $H = P_{\rho_k}$, $\epsilon = \epsilon_k$, $\tilde{\epsilon}_0 = \epsilon_k/2$, $\tilde{x}^0 = (x^0, y^0)$, $D_q = D_{\mathbf{y}}$, $\sigma_y = \sigma\rho_k$ and $L_{\nabla h} = \widehat{L}_k$, one has

$$\begin{aligned} P_{\rho_k}(x^{k+1}, y^{k+1}, z^{k+1}) &\leq \max_z P_{\rho_k}(x^0, y^0, z) + \epsilon_k^2(\widehat{L}_k^{-1} + \sigma^{-2}\rho_k^{-2}\widehat{L}_k)/2 \\ &\stackrel{(49)}{\leq} 1 + f(x^0, y^0) + \epsilon_k^2(\widehat{L}_k^{-1} + \sigma^{-2}\rho_k^{-2}\widehat{L}_k)/2. \end{aligned}$$

Hence, (48) holds as desired.

We next show that at most \widehat{N}_k evaluations of ∇f_1 , $\nabla \tilde{f}_1$, and proximal operator of f_2 and \tilde{f}_2 are respectively performed in step 2 of Algorithm 1 at iteration k . Let (x^*, y^*) be an optimal solution of (3). It then follows that $f(x^*, y^*) = f^*$ and $\tilde{f}(x^*, y^*) = \min_z \tilde{f}(x^*, z)$. By these and the definition of P_{ρ} in (4), one has

$$\max_z P_{\rho_k}(x^*, y^*, z) = f(x^*, y^*) + \rho_k(\tilde{f}(x^*, y^*) - \min_z \tilde{f}(x^*, z)) = f(x^*, y^*) = f^*,$$

which implies that

$$\min_{x,y} \max_z P_{\rho_k}(x, y, z) \leq \max_z P_{\rho_k}(x^*, y^*, z) = f^*. \quad (50)$$

Also, by (4), (6) and (7), one has

$$\min_{x,y} \max_z P_{\rho_k}(x, y, z) \stackrel{(4)}{=} \min_{x,y} \{ f(x, y) + \rho_k(\tilde{f}(x, y) - \min_z \tilde{f}(x, z)) \} \geq \min_{(x,y) \in \mathcal{X} \times \mathcal{Y}} f(x, y) \stackrel{(7)}{=} f_{\text{low}}, \quad (51)$$

$$\min_{(x,y,z) \in \mathcal{X} \times \mathcal{Y} \times \mathcal{Y}} P_{\rho_k}(x, y, z) \stackrel{(4)}{=} \min_{(x,y,z) \in \mathcal{X} \times \mathcal{Y} \times \mathcal{Y}} \{ f(x, y) + \rho_k(\tilde{f}(x, y) - \tilde{f}(x, z)) \} \stackrel{(6)(7)}{\geq} f_{\text{low}} + \rho_k(\tilde{f}_{\text{low}} - \tilde{f}_{\text{hi}}). \quad (52)$$

For ease of the rest proof, let

$$H = P_{\rho_k}, \quad H^* = \min_{x,y} \max_z P_{\rho_k}(x, y, z), \quad H_{\text{low}} = \min\{P_{\rho_k}(x, y, z) | (x, y, z) \in \mathcal{X} \times \mathcal{Y} \times \mathcal{Z}\}. \quad (53)$$

In view of these, (49), (50), (51) and (52), we obtain that

$$\max_z H(x^0, y^0, z) \stackrel{(49)}{\leq} f(x^0, y^0) + 1, \quad f_{\text{low}} \stackrel{(51)}{\leq} H^* \stackrel{(50)}{\leq} f^*, \quad H_{\text{low}} \stackrel{(52)}{\geq} f_{\text{low}} + \rho_k(\tilde{f}_{\text{low}} - \tilde{f}_{\text{hi}}).$$

Using these and Theorem 4 with $\epsilon = \epsilon_k$, $\tilde{x}^0 = (x^0, y^0)$, $D_p = \sqrt{D_{\mathbf{x}}^2 + D_{\mathbf{y}}^2}$, $D_q = D_{\mathbf{y}}$, $\tilde{\epsilon}_0 = \epsilon_k/2$, $\sigma_y = \sigma\rho_k$, $L_{\nabla h} = \hat{L}_k$, $\tilde{\alpha} = \hat{\alpha}$, $\tilde{\delta} = \hat{\delta}$, and H, H^*, H_{low} given in (53), we can conclude that step 2 of Algorithm 1 at iteration k performs at most \hat{N}_k evaluations of ∇f_1 , $\nabla \tilde{f}_1$ and proximal operator of f_2 and \tilde{f}_2 respectively for finding an ϵ_k -primal-dual stationary point $(x^{k+1}, y^{k+1}, z^{k+1})$ of problem (4) satisfying (48). \square

Lemma 2. Suppose that Assumption 1 holds and $(x^{k+1}, y^{k+1}, z^{k+1})$ is generated by Algorithm 1 for some $0 \leq k \in \hat{\mathbb{K}}$. Let P_{ρ} , $D_{\mathbf{y}}$, f_{low} , \tilde{f} and ρ_k be given in (4), (5) and (7), respectively. Then we have

$$\begin{aligned} \text{dist}\left(0, \partial f(x^{k+1}, y^{k+1}) + \rho_k \partial \tilde{f}(x^{k+1}, y^{k+1}) - (\rho_k \nabla_x \tilde{f}(x^{k+1}, z^{k+1}); 0)\right) &\leq \epsilon_k, \\ \text{dist}\left(0, \rho_k \partial \tilde{f}(x^{k+1}, z^{k+1})\right) &\leq \epsilon_k, \\ \tilde{f}(x^{k+1}, y^{k+1}) &\leq \min_z \tilde{f}(x^{k+1}, z) + \rho_k^{-1}(\max_z P_{\rho_k}(x^{k+1}, y^{k+1}, z) - f_{\text{low}}). \end{aligned}$$

Proof. Since $(x^{k+1}, y^{k+1}, z^{k+1})$ is an ϵ_k -stationary point of (4), it follows from Definition 1 that

$$\text{dist}\left(0, \partial_{x,y} P_{\rho_k}(x^{k+1}, y^{k+1}, z^{k+1})\right) \leq \epsilon_k, \quad \text{dist}\left(0, \partial_z P_{\rho_k}(x^{k+1}, y^{k+1}, z^{k+1})\right) \leq \epsilon_k.$$

Using these and the definition of P_{ρ} in (4), we have

$$\begin{aligned} \text{dist}\left(0, \partial f(x^{k+1}, y^{k+1}) + \rho_k \partial \tilde{f}(x^{k+1}, y^{k+1}) - (\rho_k \nabla_x \tilde{f}(x^{k+1}, z^{k+1}); 0)\right) &\leq \epsilon_k, \\ \text{dist}\left(0, \rho_k \partial \tilde{f}(x^{k+1}, z^{k+1})\right) &\leq \epsilon_k. \end{aligned}$$

In addition, by (4), we have

$$f(x^{k+1}, y^{k+1}) + \rho_k(\tilde{f}(x^{k+1}, y^{k+1}) - \min_z \tilde{f}(x^{k+1}, z)) = \max_z P_{\rho_k}(x^{k+1}, y^{k+1}, z),$$

which along with (7) implies that

$$\tilde{f}(x^{k+1}, y^{k+1}) - \min_z \tilde{f}(x^{k+1}, z) \leq \rho_k^{-1}(\max_z P_{\rho_k}(x^{k+1}, y^{k+1}, z) - f_{\text{low}}).$$

This completes the proof of this lemma. \square

We are now ready to prove Theorem 1.

Proof of Theorem 1. Observe from the definition of \hat{K} in (8) and $\epsilon_k = \epsilon_0 \tau^k$ that \hat{K} is the smallest nonnegative integer such that $\epsilon_{\hat{K}} \leq \varepsilon$. Hence, Algorithm 1 terminates and outputs $(x^{\hat{K}+1}, y^{\hat{K}+1})$ after $\hat{K} + 1$ outer iterations. In addition, by (9), (43), one has

$$L_{\nabla f_1} \stackrel{(43)}{\leq} \hat{L}_k, \quad 2\rho_k L_{\nabla \tilde{f}_1} \stackrel{(43)}{\leq} \hat{L}_k \stackrel{(9)}{\leq} \rho_k \hat{L}. \quad (54)$$

It follows from this, $\epsilon_{\hat{K}} = \rho_{\hat{K}}^{-1} \leq \varepsilon$, (48) and Lemma 2 that (14), (15) and (16) hold.

Let \hat{K} and \hat{N} be given in (8) and (13). Recall from Lemma 1 that the number of evaluations of ∇f_1 , $\nabla \tilde{f}_1$, proximal operator of f_2 and \tilde{f}_2 performed by Algorithm 5 at iteration k of Algorithm 1 is at most \hat{N}_k , where \hat{N}_k is given in (47). It then follows that the total number of evaluations of ∇f_1 , $\nabla \tilde{f}_1$, and proximal operator of f_2 and \tilde{f}_2 performed in Algorithm 2 is no more than $\sum_{k=0}^{\hat{K}} \hat{N}_k$, respectively. As a result, to prove the last statement of this theorem, it suffices to show that $\sum_{k=0}^{\hat{K}} \hat{N}_k \leq \hat{N}$.

To this end, using $\rho_k \geq 1 \geq \epsilon_k$, (9), (10), (11), (12), (43), (44), (45), (46) and (54), we obtain that

$$1 \geq \hat{\alpha}_k \geq \min \left\{ 1, \sqrt{8\sigma\rho_k/(\rho_k\hat{L})} \right\} \geq \hat{\alpha}, \quad (55)$$

$$\hat{\delta}_k \leq (2 + \hat{\alpha}^{-1})(D_{\mathbf{x}}^2 + D_{\mathbf{y}}^2)\rho_k\hat{L} + \max\{2\sigma\rho_k, \rho_k\hat{L}/4\}D_{\mathbf{y}}^2 \leq \rho_k\hat{\delta}, \quad (56)$$

$$\begin{aligned} \widehat{M}_k &\leq \frac{16 \max \left\{ 1/(4\rho_k L_{\nabla \tilde{f}_1}), 2/(\hat{\alpha}\rho_k L_{\nabla \tilde{f}_1}) \right\}}{\left[9\rho_k^2 \hat{L}^2 / \min\{2\rho_k \hat{L}_{\nabla \tilde{f}_1}, \sigma\rho_k\} + 3\rho_k \hat{L} \right]^{-2} \epsilon_k^2} \\ &\quad \times \left(\rho_k \hat{\delta} + 2\hat{\alpha}^{-1} \left(f^* - f_{\text{low}} + \rho_k(\tilde{f}_{\text{hi}} - \tilde{f}_{\text{low}}) + \rho_k \hat{L}(D_{\mathbf{x}}^2 + D_{\mathbf{y}}^2) \right) \right) \end{aligned} \quad (57)$$

$$\begin{aligned} &\leq \frac{16\rho_k^{-1} \max \left\{ 1/(4L_{\nabla \tilde{f}_1}), 2/(\hat{\alpha}L_{\nabla \tilde{f}_1}) \right\}}{\rho_k^{-2} \left[9\hat{L}^2 / \min\{2\hat{L}_{\nabla \tilde{f}_1}, \sigma\} + 3\hat{L} \right]^{-2} \epsilon_k^2} \rho_k \left(\hat{\delta} + 2\hat{\alpha}^{-1} \left(f^* - f_{\text{low}} + \tilde{f}_{\text{hi}} - \tilde{f}_{\text{low}} + \hat{L}(D_{\mathbf{x}}^2 + D_{\mathbf{y}}^2) \right) \right) \\ &= \epsilon_k^{-2} \rho_k^2 \widehat{M}, \end{aligned} \quad (58)$$

$$\widehat{T}_k \leq \left[16(1 + f(x^0, y^0) - f_{\text{low}}) \epsilon_k^{-2} \rho_k \hat{L} + 8\sigma^{-2} \hat{L}^2 + 7 \right]_{+} \leq \epsilon_k^{-2} \rho_k \widehat{T}, \quad (59)$$

where (55) follows from (9), (43) and (54); (56) is due to (10), (44), (55) and $\mu_k \geq 1 \geq \epsilon_k$; (57) is due to (45), (54), (55), (56) and $\epsilon_k \in (0, 1]$; (58) follows from $\rho_k \geq 1 \geq \epsilon_k$ and (11); and (59) is due to (12) and (54). By the above inequalities, (47), (54), $T > 1$ and $\rho_k \geq 1 \geq \epsilon_k$, one has

$$\begin{aligned} \sum_{k=0}^{\widehat{K}} \widehat{N}_k &\leq \sum_{k=0}^{\widehat{K}} 3397 \max \left\{ 2, \sqrt{\rho_k \hat{L}/(2\sigma\rho_k)} \right\} \\ &\quad \times \left((\epsilon_k^{-2} \rho_k \widehat{T} + 1)(\log(\epsilon_k^{-2} \rho_k^2 \widehat{M}))_+ + \epsilon_k^{-2} \rho_k \widehat{T} + 1 + 2\epsilon_k^{-2} \rho_k \widehat{T} \log(\epsilon_k^{-2} \rho_k \widehat{T} + 1) \right) \\ &\leq \sum_{k=0}^{\widehat{K}} 3397 \max \left\{ 2, \sqrt{\hat{L}/(2\sigma)} \right\} \epsilon_k^{-2} \rho_k \left((\widehat{T} + 1)(\log(\epsilon_k^{-2} \rho_k^2 \widehat{M}))_+ + \widehat{T} + 1 + 2\widehat{T} \log(\epsilon_k^{-2} \rho_k \widehat{T} + 1) \right) \\ &\leq \sum_{k=0}^{\widehat{K}} 3397 \max \left\{ 2, \sqrt{\hat{L}/(2\sigma)} \right\} \epsilon_k^{-2} \rho_k \widehat{T} \left(2(\log(\epsilon_k^{-2} \rho_k^2 \widehat{M}))_+ + 2 + 2\log(2\epsilon_k^{-2} \rho_k \widehat{T}) \right) \\ &\leq \sum_{k=0}^{\widehat{K}} 3397 \max \left\{ 2, \sqrt{\hat{L}/(2\sigma)} \right\} \widehat{T} \epsilon_k^{-2} \rho_k \left(6\log \rho_k - 8\log \epsilon_k + (\log \widehat{M})_+ + 2 + 2\log(2\widehat{T}) \right), \end{aligned} \quad (60)$$

where the first inequality follows from $\epsilon_k \in (0, 1]$, (47), (54), (58) and (59), and the second and third inequalities are due to the fact that $\rho_k \geq 1 \geq \epsilon_k$ and $T > 1$. By the definition of \widehat{K} in (8), one has $\tau^{\widehat{K}} \geq \tau\epsilon/\epsilon_0$. Also, notice from Algorithm 2 that $\rho_k = \epsilon_k^{-1} = (\epsilon_0 \tau^k)^{-1}$. It then follows from these and (60) that

$$\begin{aligned} \sum_{k=0}^{\widehat{K}} \widehat{N}_k &\leq \sum_{k=0}^{\widehat{K}} 3397 \max \left\{ 2, \sqrt{\hat{L}/(2\sigma)} \right\} \widehat{T} \times \epsilon_k^{-3} \left(14\log(1/\epsilon_k) + 2(\log \widehat{M})_+ + 2 + 2\log(2\widehat{T}) \right) \\ &= 3397 \max \left\{ 2, \sqrt{\hat{L}/(2\sigma)} \right\} \widehat{T} \sum_{k=0}^{\widehat{K}} \epsilon_0^{-3} \tau^{-3k} \left(14k \log(1/\tau) + 14\log(1/\epsilon_0) + 2(\log \widehat{M})_+ + 2 + 2\log(2\widehat{T}) \right) \\ &\leq 3397 \max \left\{ 2, \sqrt{\hat{L}/(2\sigma)} \right\} \widehat{T} \sum_{k=0}^{\widehat{K}} \epsilon_0^{-3} \tau^{-3k} \left(14\widehat{K} \log(1/\tau) + 14\log(1/\epsilon_0) + 2(\log \widehat{M})_+ + 2 + 2\log(2\widehat{T}) \right) \\ &\leq 3397 \max \left\{ 2, \sqrt{\hat{L}/(2\sigma)} \right\} \widehat{T} \epsilon_0^{-3} \tau^{-3\widehat{K}} (1 - \tau^3)^{-1} \left(14\widehat{K} \log(1/\tau) + 14\log(1/\epsilon_0) + 2(\log \widehat{M})_+ + 2 + 2\log(2\widehat{T}) \right) \\ &\leq 3397 \max \left\{ 2, \sqrt{\hat{L}/(2\sigma)} \right\} \widehat{T} \epsilon_0^{-3} (1 - \tau^3)^{-1} \\ &\quad \times (\tau\epsilon/\epsilon_0)^{-3} \left(14\widehat{K} \log(1/\tau) + 14\log(1/\epsilon_0) + 2(\log \widehat{M})_+ + 2 + 2\log(2\widehat{T}) \right) \stackrel{(13)}{\leq} \widehat{N}, \end{aligned}$$

where the third inequality is due to $\sum_{k=0}^{\hat{K}} \tau^{-3k} \leq \tau^{-3\hat{K}}/(1 - \tau^3)$, the fourth inequality follows from $\tau^{\hat{K}} \geq \tau\varepsilon/\epsilon_0$, and the last inequality is due to (13). \square

4.2 Proof of the main results in Section 3

In this subsection we provide a proof of our main result presented in Section 3, which is particularly Theorem 2. Before proceeding, one can observe from (19) and (24) that

$$\min_z \tilde{\mathcal{L}}(x, z, \lambda; \rho, \mu) \leq \tilde{f}^*(x) + \frac{\|\lambda\|^2}{2\rho\mu} \quad \forall x \in \mathcal{X}, \lambda \in \mathbb{R}_+^l, \rho, \mu > 0, \quad (61)$$

which will be frequently used later.

We next establish several technical lemmas that will be used to prove Theorem 2 subsequently.

Lemma 3. *Suppose that Assumptions 1 and 2 hold. Let \tilde{f}^* , \tilde{f}_{hi}^* , $D_{\mathbf{y}}$, $L_{\tilde{f}}$ and G be given in (5), (24), (25), and Assumptions 1 and 2, respectively. Then the following statements hold.*

- (i) $\lambda^* \geq 0$ and $\|\lambda^*\| \leq G^{-1}L_{\tilde{f}}D_{\mathbf{y}}$ for all $\lambda^* \in \Lambda^*(x)$ and $x \in \mathcal{X}$, where $\Lambda^*(x)$ denotes the set of optimal Lagrangian multipliers of problem (24) for any $x \in \mathcal{X}$.
- (ii) The function \tilde{f}^* is Lipschitz continuous on \mathcal{X} and \tilde{f}_{hi}^* is finite.

Proof. Its proof is similar to that of [39, Lemma 1] and thus omitted. \square

Lemma 4. *Suppose that Assumptions 1 and 2 hold. Let \mathbb{K} and ϑ be defined in (27) and (29), μ_k and ρ_k be given in Algorithm 2, and $\{\lambda^k\}_{k \in \mathbb{K}}$ be generated by Algorithm 2. Then we have*

$$\|\lambda^k\|^2 \leq 2\rho_k\mu_k\vartheta \quad \forall 0 \leq k \in \mathbb{K} - 1. \quad (62)$$

Proof. Its proof is similar to that of [39, Lemma 2] and thus omitted. \square

Lemma 5. *Suppose that Assumptions 1 and 2 hold. Let $D_{\mathbf{y}}$, \mathbb{K} and ϑ be defined in (5), (27) and (29), L_f , $L_{\tilde{f}}$ and G be given in Assumptions 1 and 2, and ϵ_0 , ρ_k and μ_k be given in Algorithm 2. Suppose that $(x^{k+1}, y^{k+1}, z^{k+1}, \lambda^{k+1})$ is generated by Algorithm 2 for some $0 \leq k \in \mathbb{K} - 1$ with*

$$\rho_k^{-1}\mu_k \geq 8G^{-2}\vartheta.$$

Then we have

$$\begin{aligned} \|\tilde{g}(x^{k+1}, z^{k+1})\|_+ &\leq \mu_k^{-1}\|\lambda^{k+1}\| \leq 2\mu_k^{-1}G^{-1}(\epsilon_0 + \rho_k L_{\tilde{f}})D_{\mathbf{y}}, \\ \|\tilde{g}(x^{k+1}, y^{k+1})\|_+ &\leq \mu_k^{-1}\|\lambda^k + \mu_k \tilde{g}(x^{k+1}, y^{k+1})\|_+ \leq 2\mu_k^{-1}G^{-1}(\epsilon_0 + L_f + \rho_k L_{\tilde{f}})D_{\mathbf{y}}. \end{aligned}$$

Proof. Its proof is similar to that of [39, Lemma 3] and thus omitted. \square

Lemma 6. *Suppose that Assumptions 1 and 2 hold. Let $D_{\mathbf{y}}$, \mathbb{K} and ϑ be defined in (5), (27) and (29), L_f , $L_{\tilde{f}}$ and G be given in Assumptions 1 and 2, and ϵ_0 , τ , ρ_k and μ_k be given in Algorithm 2. Suppose that $(x^{k+1}, y^{k+1}, z^{k+1}, \lambda^{k+1})$ is generated by Algorithm 2 for some $0 \leq k \in \mathbb{K} - 1$ with*

$$\rho_k^{-1}\mu_k \geq 8\tau^{-2}G^{-2}\vartheta. \quad (63)$$

Let

$$\lambda_{\mathbf{y}}^{k+1} = [\lambda^k + \mu_k \tilde{g}(x^{k+1}, y^{k+1})]_+, \quad \lambda_{\mathbf{z}}^{k+1} = \rho_k^{-1}[\lambda^k + \mu_k \tilde{g}(x^{k+1}, z^{k+1})]_+.$$

Then we have

$$\begin{aligned} &\text{dist}\left(0, \partial f(x^{k+1}, y^{k+1}) + \rho_k \partial \tilde{f}(x^{k+1}, y^{k+1}) - \rho_k (\nabla_x \tilde{f}(x^{k+1}, z^{k+1}) + \nabla_x \tilde{g}(x^{k+1}, z^{k+1})\lambda_{\mathbf{z}}^{k+1}; 0) \right. \\ &\quad \left. + \nabla \tilde{g}(x^{k+1}, y^{k+1})\lambda_{\mathbf{y}}^{k+1}\right) \leq \epsilon_k, \\ &\text{dist}\left(0, \rho_k (\partial_z \tilde{f}(x^{k+1}, z^{k+1}) + \nabla_z \tilde{g}(x^{k+1}, z^{k+1})\lambda_{\mathbf{z}}^{k+1})\right) \leq \epsilon_k, \\ &\|\tilde{g}(x^{k+1}, z^{k+1})\|_+ \leq 2\mu_k^{-1}G^{-1}(\epsilon_0 + \rho_k L_{\tilde{f}})D_{\mathbf{y}}, \\ &|\langle \lambda_{\mathbf{z}}^{k+1}, \tilde{g}(x^{k+1}, z^{k+1}) \rangle| \leq 2\rho_k^{-1}\mu_k^{-1}G^{-1}(\epsilon_0 + \rho_k L_{\tilde{f}})D_{\mathbf{y}} \max\{\|\lambda^0\|, 2G^{-1}(\epsilon_0 + \rho_k L_{\tilde{f}})D_{\mathbf{y}}\}, \\ &\|\tilde{g}(x^{k+1}, y^{k+1})\|_+ \leq 2\mu_k^{-1}G^{-1}(\epsilon_0 + L_f + \rho_k L_{\tilde{f}})D_{\mathbf{y}}, \\ &|\langle \lambda_{\mathbf{y}}^{k+1}, \tilde{g}(x^{k+1}, y^{k+1}) \rangle| \leq 2\mu_k^{-1}G^{-1}(\epsilon_0 + L_f + \rho_k L_{\tilde{f}})D_{\mathbf{y}} \max\{\|\lambda^0\|, 2G^{-1}(\epsilon_0 + L_f + \rho_k L_{\tilde{f}})D_{\mathbf{y}}\}. \end{aligned} \quad (64)$$

Proof. Its proof is similar to that of [39, Lemma 4] and thus omitted. \square

Lemma 7. Suppose that Assumptions 1 and 2 hold. Let $D_{\mathbf{x}}, D_{\mathbf{y}}, \tilde{f}_{\text{low}}, f_{\text{low}}, f^*, L_k, \tilde{f}_{\text{hi}}, f_{\text{hi}}, \tilde{g}_{\text{hi}}, \mathbb{K}$ and ϑ be defined in (5), (6), (7), (1), (23), (25), (26), (27) and (29), $\sigma, L_{\tilde{f}}$ be given in Assumption 1, $\epsilon_k, \rho_k, \mu_k$ and η_k be given in Algorithm 2, and

$$\alpha_k = \min \left\{ 1, \sqrt{8\sigma\rho_k/L_k} \right\}, \quad (65)$$

$$\delta_k = (2 + \alpha_k^{-1})(D_{\mathbf{x}}^2 + D_{\mathbf{y}}^2)L_k + \max \{2\sigma\rho_k, \alpha_k L_k/4\} D_{\mathbf{y}}^2, \quad (66)$$

$$M_k = \frac{16 \max \{1/(2L_k), \min \{1/(2\sigma\rho_k), 4/(\alpha_k L_k)\}\}}{[9L_k^2/\min\{L_k, \sigma\rho_k\} + 3L_k]^{-2} \epsilon_k^2} \times \left(\delta_k + 2\alpha_k^{-1} \left(f^* - f_{\text{low}} + \rho_k(\tilde{f}_{\text{hi}}^* - \tilde{f}_{\text{low}}) + \rho_k L_{\tilde{f}} D_{\mathbf{y}} + 3\rho_k \vartheta + \mu_k \tilde{g}_{\text{hi}}^2 + L_k(D_{\mathbf{x}}^2 + D_{\mathbf{y}}^2) \right) \right), \quad (67)$$

$$T_k = \lceil 16(f_{\text{hi}} - f_{\text{low}} + \rho_k \eta_k) L_k \epsilon_k^{-2} + 8\sigma^{-2} \rho_k^{-2} L_k^2 + 7 \rceil_+, \quad (68)$$

$$N_k = 3397 \max \left\{ 2, \sqrt{L_k/(2\sigma\rho_k)} \right\} ((T_k + 1)(\log M_k)_+ + T_k + 1 + 2T_k \log(T_k + 1)). \quad (69)$$

Then for all $0 \leq k \in \mathbb{K} - 1$, an ϵ_k -primal-dual stationary point $(x^{k+1}, y^{k+1}, z^{k+1})$ of problem (22) is successfully found at step 3 of Algorithm 2 that satisfies

$$\max_z \mathcal{L}(x^{k+1}, y^{k+1}, z, \lambda^k; \rho_k, \mu_k) \leq f_{\text{hi}} + \rho_k \eta_k + \frac{1}{2\epsilon_k^2} (L_k^{-1} + \sigma^{-2} \rho_k^{-2} L_k). \quad (70)$$

Moreover, the total number of evaluations of $\nabla f_1, \nabla \tilde{f}_1, \nabla \tilde{g}$ and proximal operator of f_2 and \tilde{f}_2 performed at step 3 in iteration k of Algorithm 2 is no more than N_k , respectively.

Proof. Observe from (18) and Assumption 1 that problem (22) can be viewed as

$$\min_{x,y} \max_z \{h(x, y, z) + p(x, y) - q(z)\}$$

with

$$h(x, y, z) = f_1(x, y) + \rho_k \tilde{f}_1(x, y) + \frac{1}{2\mu_k} \|\lambda^k + \mu_k \tilde{g}(x, y)\|_+^2 - \rho_k \tilde{f}_1(x, z) - \frac{1}{2\mu_k} \|\lambda^k + \mu_k \tilde{g}(x, z)\|_+^2, \\ p(x, y) = f_2(x) + \rho_k \tilde{f}_2(y), \quad q(z) = \rho_k \tilde{f}_2(z).$$

By (28) and Assumption 1, it can be verified that $\|\lambda^k + \mu_k \tilde{g}(x, y)\|_+^2/(2\mu_k)$ and $\|\lambda^k + \mu_k \tilde{g}(x, z)\|_+^2/(2\mu_k)$ are both $(\mu_k L_{\tilde{g}}^2 + \mu_k \tilde{g}_{\text{hi}} L_{\nabla \tilde{g}} + \|\lambda^k\| L_{\nabla \tilde{g}})$ -smooth on $\mathcal{X} \times \mathcal{Y}$. Using this and the fact that f_1 and \tilde{f}_1 are respectively $L_{\nabla f_1}$ - and $L_{\nabla \tilde{f}_1}$ -smooth on $\mathcal{X} \times \mathcal{Y}$, we can see that $h(x, y, \cdot)$ is $\sigma\rho_k$ -strongly-convex on \mathcal{Y} , and $h(x, y, z)$ is L_k -smooth on $\mathcal{X} \times \mathcal{Y} \times \mathcal{Y}$ for all $0 \leq k \in \mathbb{K} - 1$, where L_k is given in (23). Consequently, it follows from Theorem 4 (see Appendix B) that an ϵ_k -primal-dual stationary point $(x^{k+1}, y^{k+1}, z^{k+1})$ of problem (22) is successfully found by Algorithm 5 at step 3 of Algorithm 2.

In addition, by (7), (18), (19), one has

$$\min_{x,y} \max_z \mathcal{L}(x, y, z, \lambda^k; \rho_k, \mu_k) \stackrel{(18)(19)}{=} \min_{x,y} \left\{ f(x, y) + \rho_k \tilde{L}(x, y, \lambda^k; \rho_k, \mu_k) - \min_z \rho_k \tilde{L}(x, z, \lambda^k; \rho_k, \mu_k) \right\} \\ \geq \min_{(x,y) \in \mathcal{X} \times \mathcal{Y}} f(x, y) \stackrel{(7)}{=} f_{\text{low}}. \quad (71)$$

Let (x^*, y^*) be an optimal solution of (1). It then follows that $f(x^*, y^*) = f^*, \tilde{f}(x^*, y^*) = \tilde{f}^*(x^*)$ and $\tilde{g}(x^*, y^*) \leq 0$, where f^* and \tilde{f}^* are defined in (1) and (24), respectively. Using these, (6), (18), (19), (25) and (62), we obtain that

$$\min_{x,y} \max_z \mathcal{L}(x, y, z, \lambda^k; \rho_k, \mu_k) \leq \max_z \mathcal{L}(x^*, y^*, z, \lambda^k; \rho_k, \mu_k) \\ \stackrel{(18)(19)}{=} f(x^*, y^*) + \rho_k \tilde{f}(x^*, y^*) + \frac{1}{2\mu_k} \|\lambda^k + \mu_k \tilde{g}(x^*, y^*)\|_+^2 - \min_z \rho_k \tilde{L}(x^*, z, \lambda^k; \rho_k, \mu_k) \\ \leq f^* + \rho_k \tilde{f}^*(x^*) + \frac{1}{2\mu_k} \|\lambda^k\|^2 - \min_z \left\{ \rho_k \tilde{f}(x^*, z) + \frac{1}{2\mu_k} \|\lambda^k + \mu_k \tilde{g}(x^*, z)\|_+^2 \right\} \\ \stackrel{(6)(25)}{\leq} f^* + \rho_k(\tilde{f}_{\text{hi}}^* - \tilde{f}_{\text{low}}) + \frac{1}{2\mu_k} \|\lambda^k\|^2 \stackrel{(62)}{\leq} f^* + \rho_k(\tilde{f}_{\text{hi}}^* - \tilde{f}_{\text{low}}) + \rho_k \vartheta, \quad (72)$$

where the second inequality is due to $\tilde{f}(x^*, y^*) = \tilde{f}^*(x^*)$, $\tilde{g}(x^*, y^*) \leq 0$ and (19). Also, by (5), (6), (18), (26) and (62), one has

$$\begin{aligned}
& \min_{(x,y,z) \in \mathcal{X} \times \mathcal{Y} \times \mathcal{Y}} \mathcal{L}(x, y, z, \lambda^k; \rho_k, \mu_k) \\
& \stackrel{(18)}{\geq} \min_{(x,y,z) \in \mathcal{X} \times \mathcal{Y} \times \mathcal{Y}} \left\{ f(x, y) + \rho_k(\tilde{f}(x, y) - \tilde{f}(x, z)) - \frac{1}{2\mu_k} \|[\lambda^k + \mu_k \tilde{g}(x, z)]_+\|^2 \right\} \\
& \geq \min_{(x,y,z) \in \mathcal{X} \times \mathcal{Y} \times \mathcal{Y}} \left\{ f(x, y) - \rho_k L_{\tilde{f}} \|y - z\| - \frac{1}{2\mu_k} (\|\lambda^k\| + \mu_k \|[\tilde{g}(x, z)]_+\|)^2 \right\} \\
& \geq \min_{(x,y,z) \in \mathcal{X} \times \mathcal{Y} \times \mathcal{Y}} \left\{ f(x, y) - \rho_k L_{\tilde{f}} \|y - z\| - \frac{1}{\mu_k} \|\lambda^k\|^2 - \mu_k \|[\tilde{g}(x, z)]_+\|^2 \right\} \\
& \geq f_{\text{low}} - \rho_k L_{\tilde{f}} D_{\mathbf{y}} - 2\rho_k \vartheta - \mu_k \tilde{g}_{\text{hi}}^2,
\end{aligned} \tag{73}$$

where the second inequality is due to $\lambda^k \in \mathbb{R}_+^l$ and $L_{\tilde{f}}$ -Lipschitz continuity of \tilde{f} (see Assumption 1(i)), and the last inequality is due to (5), (6), (26) and (62). Notice from step 2 of Algorithm 2 that y_{init}^k is an approximate solution of $\min_z \tilde{L}(x^k, z, \lambda^k; \rho_k, \mu_k)$ satisfying (20). It then follows from (18), (19), (20) and (26) that

$$\begin{aligned}
\max_z \mathcal{L}(x^k, y_{\text{init}}^k, z, \lambda^k; \rho_k, \mu_k) & \stackrel{(18)(19)}{=} f(x^k, y_{\text{init}}^k) + \rho_k \left(\tilde{\mathcal{L}}(x^k, y_{\text{init}}^k, \lambda^k; \rho_k, \mu_k) - \min_z \tilde{\mathcal{L}}(x^k, z, \lambda^k; \rho_k, \mu_k) \right) \\
& \stackrel{(20)}{\leq} f(x^k, y_{\text{init}}^k) + \rho_k \eta_k \stackrel{(26)}{\leq} f_{\text{hi}} + \rho_k \eta_k.
\end{aligned} \tag{74}$$

To complete the rest of the proof, let

$$H(x, y, z) = \mathcal{L}(x, y, z, \lambda^k; \rho_k, \mu_k), \quad H^* = \min_{x,y} \max_z \mathcal{L}(x, y, z, \lambda^k; \rho_k, \mu_k), \tag{75}$$

$$H_{\text{low}} = \min \{ \mathcal{L}(x, y, z, \lambda^k; \rho_k, \mu_k) | (x, y, z) \in \mathcal{X} \times \mathcal{Y} \times \mathcal{Y} \}. \tag{76}$$

In view of these, (71), (72), (73) and (74), we obtain that

$$\begin{aligned}
\max_z H(x^k, y_{\text{init}}^k, z) & \stackrel{(74)}{\leq} f_{\text{hi}} + \rho_k \eta_k, \\
f_{\text{low}} & \stackrel{(71)}{\leq} H^* \stackrel{(72)}{\leq} f^* + \rho_k (\tilde{f}_{\text{hi}}^* - \tilde{f}_{\text{low}}) + \rho_k \vartheta, \\
H_{\text{low}} & \stackrel{(73)}{\geq} f_{\text{low}} - \rho_k L_{\tilde{f}} D_{\mathbf{y}} - 2\rho_k \vartheta - \mu_k \tilde{g}_{\text{hi}}^2.
\end{aligned}$$

Using these and Theorem 4 (see Appendix B) with $\tilde{x}^0 = (x^k, y_{\text{init}}^k)$, $\epsilon = \epsilon_k$, $\tilde{\epsilon}_0 = \epsilon_k/2$, $\sigma_y = \sigma \rho_k$, $L_{\nabla h} = L_k$, $\tilde{\alpha} = \alpha_k$, $\tilde{\delta} = \delta_k$, $D_p = \sqrt{D_{\mathbf{x}}^2 + D_{\mathbf{y}}^2}$, $D_q = D_{\mathbf{y}}$, and H, H^*, H_{low} given in (75) and (76), we can conclude that the ϵ_k -primal-dual stationary point $(x^{k+1}, y^{k+1}, z^{k+1})$ of problem (22) found at step 3 of Algorithm 2 satisfies (70). Moreover, the total number of evaluations of ∇f_1 , $\nabla \tilde{f}_1$, $\nabla \tilde{g}$ and proximal operator of f_2 and \tilde{f}_2 performed by Algorithm 5 at step 3 of Algorithm 2 is no more than N_k , respectively. \square

Lemma 8. *Suppose that Assumptions 1 and 2 hold. Let $D_{\mathbf{y}}, f_{\text{low}}, \tilde{f}^*, L_k, f_{\text{hi}}$ and \mathbb{K} be defined in (5), (7), (24), (23) and (27), $\sigma, L_f, L_{\tilde{f}}$ and G be given in Assumptions 1 and 2, and $\epsilon_k, \rho_k, \mu_k, \eta_k$ and λ^0 be given in Algorithm 2. Suppose that $(x^{k+1}, y^{k+1}, \lambda^{k+1})$ is generated by Algorithm 2 for some $0 \leq k \in \mathbb{K} - 1$ satisfying (63). Then we have*

$$\begin{aligned}
|\tilde{f}(x^{k+1}, y^{k+1}) - \tilde{f}^*(x^{k+1})| & \leq \max \left\{ 2\mu_k^{-1} G^{-2} L_{\tilde{f}} (\epsilon_0 + L_f + \rho_k L_{\tilde{f}}) D_{\mathbf{y}}^2, \right. \\
& \quad \rho_k^{-1} \mu_k^{-1} \max\{\|\lambda^0\|, 2G^{-1}(\epsilon_0 + \rho_k L_{\tilde{f}}) D_{\mathbf{y}}\} / 2 \\
& \quad \left. + \rho_k^{-1} \left(f_{\text{hi}} - f_{\text{low}} + \rho_k \eta_k + \frac{1}{2\epsilon_k^2} (L_k^{-1} + \sigma^{-2} \rho_k^{-2} L_k) \right) \right\}.
\end{aligned}$$

Proof. We first claim that

$$\|\lambda^k\| \leq \max\{\|\lambda^0\|, 2G^{-1}(\epsilon_0 + \rho_k L_{\tilde{f}}) D_{\mathbf{y}}\}. \tag{77}$$

Indeed, (77) clearly holds if $k = 0$. We now assume that $k > 0$. Notice from Algorithm 2 that $\mu_{k-1} = \tau^3 \mu_k$ and $\rho_{k-1} = \tau \rho_k$, which along with (63) imply that $\rho_{k-1}^{-1} \mu_{k-1} \geq 8G^{-2}\vartheta$. By this and Lemma 4 with k replaced by $k-1$, one can conclude that $\|\lambda^k\| \leq 2G^{-1}(\epsilon_0 + \rho_{k-1}L_{\tilde{f}})D_{\mathbf{y}}$. This together with $\rho_{k-1} < \rho_k$ implies that (77) holds as desired.

Using this, (7), (18), (19), (61) and (77), we have

$$\begin{aligned}
& \max_z \mathcal{L}(x^{k+1}, y^{k+1}, z, \lambda^k; \rho_k, \mu_k) \\
& \stackrel{(18)(19)}{=} f(x^{k+1}, y^{k+1}) + \rho_k \tilde{f}(x^{k+1}, y^{k+1}) + \frac{1}{2\mu_k} \|[\lambda^k + \mu_k \tilde{g}(x^{k+1}, y^{k+1})]_+\|^2 - \min_z \rho_k \tilde{\mathcal{L}}(x^{k+1}, z, \lambda^k; \rho_k, \mu_k) \\
& \geq f(x^{k+1}, y^{k+1}) + \rho_k \tilde{f}(x^{k+1}, y^{k+1}) - \min_z \rho_k \tilde{\mathcal{L}}(x^{k+1}, z, \lambda^k; \rho_k, \mu_k) \\
& \stackrel{(7)(61)}{\geq} f_{\text{low}} + \rho_k (\tilde{f}(x^{k+1}, y^{k+1}) - \tilde{f}^*(x^{k+1})) - \frac{1}{2\mu_k} \|\lambda^k\|^2 \\
& \stackrel{(77)}{\geq} f_{\text{low}} + \rho_k (\tilde{f}(x^{k+1}, y^{k+1}) - \tilde{f}^*(x^{k+1})) - \mu_k^{-1} \max\{\|\lambda^0\|, 2G^{-1}(\epsilon_0 + \rho_k L_{\tilde{f}})D_{\mathbf{y}}\}/2.
\end{aligned}$$

This together with (70) implies that

$$\begin{aligned}
\tilde{f}(x^{k+1}, y^{k+1}) - \tilde{f}^*(x^{k+1}) & \leq \rho_k^{-1} \left(f_{\text{hi}} - f_{\text{low}} + \rho_k \eta_k + \frac{1}{2\epsilon_k^2} (L_k^{-1} + \sigma^{-2} \rho_k^{-2} L_k) \right) \\
& \quad + \rho_k^{-1} \mu_k^{-1} \max\{\|\lambda^0\|, 2G^{-1}(\epsilon_0 + \rho_k L_{\tilde{f}})D_{\mathbf{y}}\}/2.
\end{aligned} \tag{78}$$

On the other hand, let $\lambda^* \in \mathbb{R}_+^l$ be an optimal Lagrangian multiplier of problem (24) with $x = x^{k+1}$. It then follows from Lemma 3(i) that $\|\lambda^*\| \leq G^{-1}L_{\tilde{f}}D_{\mathbf{y}}$. Using these, Lemma 2, (24) and (64), we have

$$\begin{aligned}
\tilde{f}^*(x^{k+1}) & = \min_y \left\{ \tilde{f}(x^{k+1}, y) + \langle \lambda^*, \tilde{g}(x^{k+1}, y) \rangle \right\} \leq \tilde{f}(x^{k+1}, y^{k+1}) + \langle \lambda^*, \tilde{g}(x^{k+1}, y^{k+1}) \rangle \\
& \leq \tilde{f}(x^{k+1}, y^{k+1}) + \|\lambda^*\| \|[\tilde{g}(x^{k+1}, y^{k+1})]_+\| \leq \tilde{f}(x^{k+1}, y^{k+1}) + 2\mu_k^{-1} G^{-2} L_{\tilde{f}} (\epsilon_0 + L_f + \rho_k L_{\tilde{f}}) D_{\mathbf{y}}^2.
\end{aligned}$$

The conclusion of this lemma then follows from this and (78). \square

Lemma 9. Suppose that Assumptions 1 and 2 hold. Let $D_{\mathbf{y}}$, \tilde{L}_k and \mathbb{K} be defined in (5), (21) and (27), σ be given in Assumption 1, η_k be given in Algorithm 2, and

$$\tilde{N}_k = 2 \left\lceil \sqrt{\tilde{L}_k \sigma^{-1}} \right\rceil \max \left\{ 1, \left\lceil 2 \log(2\eta_k^{-1} \tilde{L}_k D_{\mathbf{y}}^2) \right\rceil \right\} + 1. \tag{79}$$

Then for all $0 \leq k \in \mathbb{K} - 1$, y_{init}^k satisfying (20) is found at step 3 of Algorithm 2 by Algorithm 3 in no more than \tilde{N}_k evaluations of $\nabla \tilde{f}_1$, $\nabla \tilde{g}$ and the proximal operator of \tilde{f}_2 , respectively.

Proof. Notice from (19) and Algorithm 2 that y_{init}^k satisfying (20) is found by Algorithm 3 applied to the problem

$$\min_y \left\{ \tilde{L}(x^k, y, \lambda^k; \rho_k, \mu_k) = \phi(y) + P(y) \right\},$$

where $\phi(y) = \tilde{f}_1(x^k, y) + \|[\lambda^k + \mu_k \tilde{g}(x^k, y)]_+\|^2 / (2\rho_k \mu_k)$ and $P(y) = \tilde{f}_2(y)$. By Assumption 1 and (28), one can see that ϕ is σ -strongly-convex and \tilde{L}_k -smooth on $\text{dom } P$, where \tilde{L}_k is given in (21). It then follows from this and Theorem 3 (see Appendix A) and (90) with $\tilde{\epsilon} = \eta_k / D_{\mathbf{y}}$, $D_P = D_{\mathbf{y}}$, $\sigma_\phi = \sigma$ and $L_{\nabla \phi} = \tilde{L}_k$ that Algorithm 3 finds y_{init}^k satisfying (20) in no more than \tilde{T}_k iterations, where

$$\tilde{T}_k = \left\lceil \sqrt{\tilde{L}_k \sigma^{-1}} \right\rceil \max \left\{ 1, \left\lceil 2 \log(2\eta_k^{-1} \tilde{L}_k D_{\mathbf{y}}^2) \right\rceil \right\}.$$

Notice that the first step of Algorithm 3 requires one evaluation of $\nabla \phi$ and the proximal operator of P , respectively, and each iteration of Algorithm 3 requires two evaluation of $\nabla \phi$ and the proximal operator of P , respectively. Hence, the conclusion of this lemma holds. \square

We are now ready to prove Theorem 2.

Proof of Theorem 2. (i) Observe from the definition of K in (27) and $\epsilon_k = \epsilon_0 \tau^k$ that K is the smallest nonnegative integer such that $\epsilon_K \leq \varepsilon$. Hence, Algorithm 2 terminates and outputs (x^{K+1}, y^{K+1}) after $K + 1$ outer iterations. Also, one can see from Algorithm 2 that

$$\rho_K = \epsilon_K^{-1}, \quad \mu_K = \epsilon_K^{-3}, \quad \eta_K = \epsilon_K. \quad (80)$$

Moreover, notice from the assumption of Theorem 2 that $\varepsilon^{-2} - 8\tau^{-2}G^{-2}\vartheta \geq 0$. It then follows from this and (80) that

$$\rho_K^{-1}\mu_K = \epsilon_K^{-2} \geq \varepsilon^{-2} \geq 8\tau^{-2}G^{-2}\vartheta,$$

which implies that (63) holds for $k = K$. In addition, by (23), (30), (62) and $\mu_k \geq \rho_k \geq 1$, one has that for all $0 \leq k \in \mathbb{K} - 1$,

$$\begin{aligned} 2\mu_k L_{\tilde{g}}^2 &\leq L_k \stackrel{(23)}{=} L_{\nabla f_1} + 2\rho_k L_{\nabla \tilde{f}_1} + 2\mu_k L_{\tilde{g}}^2 + 2\mu_k \tilde{g}_{\text{hi}} L_{\nabla \tilde{g}} + 2\|\lambda^k\| L_{\nabla \tilde{g}} \\ &\stackrel{(62)}{\leq} L_{\nabla f_1} + 2\rho_k L_{\nabla \tilde{f}_1} + 2\mu_k L_{\tilde{g}}^2 + 2\mu_k \tilde{g}_{\text{hi}} L_{\nabla \tilde{g}} + 2\sqrt{2\rho_k \mu_k \vartheta} L_{\nabla \tilde{g}} \leq \mu_k L. \end{aligned} \quad (81)$$

It then follows from $\epsilon_K \leq \varepsilon$, (80) and Lemmas 6 and 8 that (35)-(41) hold.

(ii) Let K and N be given in (27) and (42). Recall from Lemmas 7 and 9 that the number of evaluations of ∇f_1 , $\nabla \tilde{f}_1$, $\nabla \tilde{g}$, proximal operator of f_2 and \tilde{f}_2 performed by Algorithms 3 and 5 at iteration k of Algorithm 2 is at most $N_k + \tilde{N}_k$, where N_k and \tilde{N}_k are given in (69) and (79), respectively. By this and statement (i) of this theorem, one can observe that the total number of evaluations of ∇f_1 , $\nabla \tilde{f}_1$, $\nabla \tilde{g}$ and proximal operator of f_2 and \tilde{f}_2 performed in Algorithm 2 is no more than $\sum_{k=0}^K (N_k + \tilde{N}_k)$, respectively. As a result, to prove statement (ii) of this theorem, it suffices to show that $\sum_{k=0}^K (N_k + \tilde{N}_k) \leq N$.

To this end, using $\mu_k \geq \rho_k \geq 1 \geq \epsilon_k$, (31), (32), (33), (65), (66), (67), (68) and (81), we obtain that

$$1 \geq \alpha_k \geq \min \left\{ 1, \sqrt{8\sigma\rho_k/(\mu_k L)} \right\} \geq \rho_k^{1/2} \mu_k^{-1/2} \alpha, \quad (82)$$

$$\delta_k \leq (2 + \rho_k^{-1/2} \mu_k^{1/2} \alpha^{-1})(D_{\mathbf{x}}^2 + D_{\mathbf{y}}^2) \mu_k L + \max\{2\sigma\rho_k, \mu_k L/4\} D_{\mathbf{y}}^2 \leq \rho_k^{-1/2} \mu_k^{3/2} \delta, \quad (83)$$

$$\begin{aligned} M_k &\leq \frac{16 \max \left\{ 1/(4\mu_k L_{\tilde{g}}^2), 2/(\rho_k^{1/2} \mu_k^{-1/2} \alpha \mu_k L_{\tilde{g}}^2) \right\}}{\left[9\mu_k^2 L^2 / \min\{2\mu_k L_{\tilde{g}}^2, \sigma\rho_k\} + 3\mu_k L \right]^{-2} \epsilon_k^2} \times \left(\rho_k^{-1/2} \mu_k^{3/2} \delta \right. \\ &\quad \left. + 2\rho_k^{-1/2} \mu_k^{1/2} \alpha^{-1} \left(f^* - f_{\text{low}} + \rho_k(\tilde{f}_{\text{hi}}^* - \tilde{f}_{\text{low}}) + \rho_k L_{\tilde{f}} D_{\mathbf{y}} + 3\rho_k \vartheta + \mu_k \tilde{g}_{\text{hi}}^2 + \mu_k L(D_{\mathbf{x}}^2 + D_{\mathbf{y}}^2) \right) \right) \end{aligned} \quad (84)$$

$$\begin{aligned} &\leq \frac{16\rho_k^{-1/2} \mu_k^{-1/2} \max \left\{ 1/(4L_{\tilde{g}}^2), 2/(\alpha L_{\tilde{g}}^2) \right\}}{\rho_k^2 \mu_k^{-4} \left[9L^2 / \min\{2L_{\tilde{g}}^2, \sigma\} + 3L \right]^{-2} \epsilon_k^2} \times \rho_k^{-1/2} \mu_k^{3/2} \\ &\quad \times \left(\delta + 2\alpha^{-1} \left(f^* - f_{\text{low}} + \tilde{f}_{\text{hi}}^* - \tilde{f}_{\text{low}} + L_{\tilde{f}} D_{\mathbf{y}} + 3\vartheta + \tilde{g}_{\text{hi}}^2 + L(D_{\mathbf{x}}^2 + D_{\mathbf{y}}^2) \right) \right) = \epsilon_k^{-2} \rho_k^{-3} \mu_k^5 M, \end{aligned} \quad (85)$$

$$T_k \leq \left\lceil 16(f_{\text{hi}} - f_{\text{low}} + \rho_k \eta_k) \epsilon_k^{-2} \mu_k L + 8\sigma^{-2} \rho_k^{-2} \mu_k^2 L^2 + 7 \right\rceil_+ \leq \epsilon_k^{-2} \mu_k T, \quad (86)$$

where (82) follows from (31), (65) and (81); (83) is due to (31), (66), (82) and $\mu_k \geq 1 \geq \epsilon_k$; (84) is due to (67), (81), (82), (83) and $\epsilon_k \in (0, 1]$; (85) follows from $\mu_k \geq \rho_k \geq 1 \geq \epsilon_k$ and (32); and (86) is due to (33), (81) and the fact that $\epsilon_k \in (0, 1]$ and $\rho_k \eta_k = 1$. By the above inequalities, (69), (81), $T > 1$ and $\mu_k \geq 1 \geq \epsilon_k$, one has

$$\begin{aligned} \sum_{k=0}^K N_k &\leq \sum_{k=0}^K 3397 \max \left\{ 2, \sqrt{\mu_k L / (2\sigma\rho_k)} \right\} \\ &\quad \times ((\epsilon_k^{-2} \mu_k T + 1)(\log(\epsilon_k^{-2} \rho_k^{-3} \mu_k^5 M))_+ + \epsilon_k^{-2} \mu_k T + 1 + 2\epsilon_k^{-2} \mu_k T \log(\epsilon_k^{-2} \mu_k T + 1)) \\ &\leq \sum_{k=0}^K 3397 \max \left\{ 2, \sqrt{L / (2\sigma)} \right\} \rho_k^{-1/2} \mu_k^{1/2} \times \epsilon_k^{-2} \mu_k ((T + 1)(\log(\epsilon_k^{-2} \rho_k^{-3} \mu_k^5 M))_+ + T + 1 + 2T \log(\epsilon_k^{-2} \mu_k T + 1)) \end{aligned}$$

$$\begin{aligned}
&\leq \sum_{k=0}^K 3397 \max \left\{ 2, \sqrt{L/(2\sigma)} \right\} \epsilon_k^{-2} \rho_k^{-1/2} \mu_k^{3/2} T \left(2(\log(\epsilon_k^{-2} \rho_k^{-3} \mu_k^5 M))_+ + 2 + 2 \log(2\epsilon_k^{-2} \mu_k T) \right) \\
&\leq \sum_{k=0}^K 3397 \max \left\{ 2, \sqrt{L/(2\sigma)} \right\} T \epsilon_k^{-2} \rho_k^{-1/2} \mu_k^{3/2} (12 \log \mu_k - 6 \log \rho_k - 8 \log \epsilon_k + 2(\log M)_+ + 2 + 2 \log(2T)),
\end{aligned} \tag{87}$$

where the first inequality follows from $\epsilon_k \in (0, 1]$, (69), (81), (85) and (86), and the second and third inequalities are due to the fact that $\mu_k \geq 1 \geq \epsilon_k$ and $T > 1$. By the definition of K in (27), one has $\tau^K \geq \tau\epsilon/\epsilon_0$. Also, notice from Algorithm 2 that $\rho_k = \epsilon_k^{-1} = (\epsilon_0 \tau^k)^{-1}$ and $\mu_k = \epsilon_k^{-3} = (\epsilon_0 \tau^k)^{-3}$. It then follows from these and (87) that

$$\begin{aligned}
\sum_{k=0}^K N_k &\leq \sum_{k=0}^K 3397 \max \left\{ 2, \sqrt{L/(2\sigma)} \right\} T \\
&\quad \times \epsilon_k^{-6} (38 \log(1/\epsilon_k) + 2(\log M)_+ + 2 + 2 \log(2T)) \\
&= 3397 \max \left\{ 2, \sqrt{L/(2\sigma)} \right\} T \sum_{k=0}^K \epsilon_0^{-6} \tau^{-6k} (38k \log(1/\tau) + 38 \log(1/\epsilon_0) + 2(\log M)_+ + 2 + 2 \log(2T)) \\
&\leq 3397 \max \left\{ 2, \sqrt{L/(2\sigma)} \right\} T \sum_{k=0}^K \epsilon_0^{-6} \tau^{-6k} (38K \log(1/\tau) + 38 \log(1/\epsilon_0) + 2(\log M)_+ + 2 + 2 \log(2T)) \\
&\leq 3397 \max \left\{ 2, \sqrt{L/(2\sigma)} \right\} T \epsilon_0^{-6} \tau^{-6K} (1 - \tau^6)^{-1} (38K \log(1/\tau) + 38 \log(1/\epsilon_0) + 2(\log M)_+ + 2 + 2 \log(2T)) \\
&\leq 3397 \max \left\{ 2, \sqrt{L/(2\sigma)} \right\} T \epsilon_0^{-6} (1 - \tau^6)^{-1} \\
&\quad \times (\tau\epsilon/\epsilon_0)^{-6} (38K \log(1/\tau) + 38 \log(1/\epsilon_0) + 2(\log M)_+ + 2 + 2 \log(2T)),
\end{aligned} \tag{88}$$

where the second last inequality is due to $\sum_{k=0}^K \tau^{-6k} \leq \tau^{-6K}/(1 - \tau^6)$, and the last inequality follows from $\tau^K \geq \tau\epsilon/\epsilon_0$.

In addition, observe from (21), (30), (62) and $\rho_k^{-1} \mu_k \geq 1$, one has that for all $0 \leq k \in \mathbb{K} - 1$,

$$\tilde{L}_k = L_{\nabla \tilde{f}_1} + \rho_k^{-1} (\mu_k L_{\tilde{g}}^2 + \mu_k \tilde{g}_{\text{hi}} L_{\nabla \tilde{g}} + \|\lambda^k\| L_{\nabla \tilde{g}}) \leq L_{\nabla \tilde{f}_1} + \rho_k^{-1} (\mu_k L_{\tilde{g}}^2 + \mu_k \tilde{g}_{\text{hi}} L_{\nabla \tilde{g}} + \sqrt{2\rho_k \mu_k} \vartheta L_{\nabla \tilde{g}}) \leq \rho_k^{-1} \mu_k \tilde{L}.$$

Using this, (79), $\epsilon_k = \epsilon_0 \tau^k$, $\rho_k = \epsilon_k^{-1}$, $\mu_k = \epsilon_k^{-3}$ and $\eta_k = \epsilon_k$, we have

$$\begin{aligned}
\sum_{k=1}^K \tilde{N}_k &\leq \sum_{k=1}^K \left(2 \left\lceil \sqrt{\frac{\rho_k^{-1} \mu_k \tilde{L}}{\sigma}} \right\rceil \max \left\{ 1, \left\lceil 2 \log \left(2\eta_k^{-1} \rho_k^{-1} \mu_k \tilde{L} D_{\mathbf{y}}^2 \right) \right\rceil \right\} + 1 \right) \\
&= \sum_{k=1}^K 2 \left\lceil (\epsilon_0 \tau^k)^{-1} \sqrt{\frac{\tilde{L}}{\sigma}} \right\rceil \max \left\{ 1, \left\lceil 2 \log(2\tilde{L} D_{\mathbf{y}}^2) + 6k \log(1/\tau) - 6 \log \epsilon_0 \right\rceil \right\} + K \\
&\leq \sum_{k=1}^K 2(\epsilon_0 \tau^k)^{-1} \left\lceil \sqrt{\frac{\tilde{L}}{\sigma}} + 1 \right\rceil \max \left\{ 1, \left\lceil 2 \log(2\tilde{L} D_{\mathbf{y}}^2) + 6K \log(1/\tau) - 6 \log \epsilon_0 \right\rceil \right\} + K \\
&\leq 2\epsilon_0^{-1} \tau^{-K} (1 - \tau) \left\lceil \sqrt{\frac{\tilde{L}}{\sigma}} + 1 \right\rceil \max \left\{ 1, \left\lceil 2 \log(2\tilde{L} D_{\mathbf{y}}^2) + 6K \log(1/\tau) - 6 \log \epsilon_0 \right\rceil \right\} + K \\
&\leq 2(\tau\epsilon)^{-1} (1 - \tau) \left\lceil \sqrt{\frac{\tilde{L}}{\sigma}} + 1 \right\rceil \max \left\{ 1, \left\lceil 2 \log(2\tilde{L} D_{\mathbf{y}}^2) + 6K \log(1/\tau) - 6 \log \epsilon_0 \right\rceil \right\} + K
\end{aligned}$$

where the second last inequality is due to $\sum_{k=0}^K \tau^{-k} \leq \tau^{-K}/(1 - \tau)$, and the last inequality follows from $\tau^K \geq \tau\epsilon/\epsilon_0$. This together with (42) and (88) implies that $\sum_{k=1}^K (N_k + \tilde{N}_k) \leq N$. Hence, statement (ii) of Theorem 2 holds. \square

References

- [1] G. B. Allende and G. Still. Solving bilevel programs with the KKT-approach. *Mathematical programming*, 138(1):309–332, 2013.
- [2] J. F. Bard. *Practical bilevel optimization: algorithms and applications*, volume 30. Springer Science & Business Media, 2013.
- [3] K. P. Bennett, G. Kunapuli, J. Hu, and J.-S. Pang. Bilevel optimization and machine learning. In *IEEE World Congress on Computational Intelligence*, pages 25–47. Springer, 2008.
- [4] L. Bertinetto, J. F. Henriques, P. Torr, and A. Vedaldi. Meta-learning with differentiable closed-form solvers. In *International Conference on Learning Representations*, 2018.
- [5] L. Chen, J. Xu, and J. Zhang. On bilevel optimization without lower-level strong convexity. *arXiv preprint arXiv:2301.00712*, 2023.
- [6] T. Chen, Y. Sun, Q. Xiao, and W. Yin. A single-timescale method for stochastic bilevel optimization. In *International Conference on Artificial Intelligence and Statistics*, pages 2466–2488, 2022.
- [7] B. Colson, P. Marcotte, and G. Savard. An overview of bilevel optimization. *Annals of operations research*, 153(1):235–256, 2007.
- [8] C. Crockett, J. A. Fessler, et al. Bilevel methods for image reconstruction. *Foundations and Trends® in Signal Processing*, 15(2-3):121–289, 2022.
- [9] S. Dempe. *Foundations of bilevel programming*. Springer Science & Business Media, 2002.
- [10] S. Dempe, V. Kalashnikov, G. A. Pérez-Valdés, and N. Kalashnykova. Bilevel programming problems. *Energy Systems. Springer, Berlin*, 10:978–3, 2015.
- [11] S. Dempe and A. Zemkoho. Bilevel optimization. In *Springer optimization and its applications. Vol. 161*. Springer, 2020.
- [12] S. Dempe and A. B. Zemkoho. The bilevel programming problem: reformulations, constraint qualifications and optimality conditions. *Mathematical Programming*, 138(1):447–473, 2013.
- [13] M. Feurer and F. Hutter. Hyperparameter optimization. In *Automated machine learning*, pages 3–33. Springer, Cham, 2019.
- [14] L. Franceschi, M. Donini, P. Frasconi, and M. Pontil. Forward and reverse gradient-based hyperparameter optimization. In *International Conference on Machine Learning*, pages 1165–1173, 2017.
- [15] L. Franceschi, P. Frasconi, S. Salzo, R. Grazzi, and M. Pontil. Bilevel programming for hyperparameter optimization and meta-learning. In *International Conference on Machine Learning*, pages 1568–1577, 2018.
- [16] R. Grazzi, L. Franceschi, M. Pontil, and S. Salzo. On the iteration complexity of hypergradient computation. In *International Conference on Machine Learning*, pages 3748–3758, 2020.
- [17] Z. Guo, Q. Hu, L. Zhang, and T. Yang. Randomized stochastic variance-reduced methods for multi-task stochastic bilevel optimization. *arXiv preprint arXiv:2105.02266*, 2021.
- [18] P. Hansen, B. Jaumard, and G. Savard. New branch-and-bound rules for linear bilevel programming. *SIAM Journal on scientific and Statistical Computing*, 13(5):1194–1217, 1992.
- [19] M. Hong, H.-T. Wai, Z. Wang, and Z. Yang. A two-timescale stochastic algorithm framework for bilevel optimization: Complexity analysis and application to actor-critic. *SIAM Journal on Optimization*, 33(1):147–180, 2023.
- [20] X. Hu, N. Xiao, X. Liu, and K.-C. Toh. An improved unconstrained approach for bilevel optimization. *arXiv preprint arXiv:2208.00732*, 2022.
- [21] F. Huang and H. Huang. Biadam: Fast adaptive bilevel optimization methods. *arXiv preprint arXiv:2106.11396*, 2021.

- [22] M. Huang, K. Ji, S. Ma, and L. Lai. Efficiently escaping saddle points in bilevel optimization. *arXiv preprint arXiv:2202.03684*, 2022.
- [23] Y. Ishizuka and E. Aiyoshi. Double penalty method for bilevel optimization problems. *Annals of Operations Research*, 34(1):73–88, 1992.
- [24] K. Ji, J. D. Lee, Y. Liang, and H. V. Poor. Convergence of meta-learning with task-specific adaptation over partial parameters. *Advances in Neural Information Processing Systems*, 33:11490–11500, 2020.
- [25] K. Ji, J. Yang, and Y. Liang. Bilevel optimization: Convergence analysis and enhanced design. In *International conference on machine learning*, pages 4882–4892, 2021.
- [26] P. Khanduri, S. Zeng, M. Hong, H.-T. Wai, Z. Wang, and Z. Yang. A near-optimal algorithm for stochastic bilevel optimization via double-momentum. *Advances in neural information processing systems*, 34:30271–30283, 2021.
- [27] V. Konda and J. Tsitsiklis. Actor-critic algorithms. *Advances in neural information processing systems*, 12, 1999.
- [28] J. Kwon, D. Kwon, S. Wright, and R. Nowak. A fully first-order method for stochastic bilevel optimization. *arXiv preprint arXiv:2301.10945*, 2023.
- [29] J. Li, B. Gu, and H. Huang. A fully single loop algorithm for bilevel optimization without hessian inverse. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 36, pages 7426–7434, 2022.
- [30] Y. Li, G.-H. Lin, J. Zhang, and X. Zhu. A novel approach for bilevel programs based on Wolfe duality. *arXiv preprint arXiv:2302.06838*, 2023.
- [31] Q. Lin, Z. Lu, and L. Xiao. An accelerated randomized proximal coordinate gradient method and its application to regularized empirical risk minimization. *SIAM Journal on Optimization*, 25(4):2244–2273, 2015.
- [32] B. Liu, M. Ye, S. Wright, P. Stone, and Q. Liu. Bome! bilevel optimization made easy: A simple first-order approach. *Advances in Neural Information Processing Systems*, 35:17248–17262, 2022.
- [33] H. Liu, K. Simonyan, and Y. Yang. Darts: Differentiable architecture search. In *International Conference on Learning Representations*, 2018.
- [34] R. Liu, J. Gao, J. Zhang, D. Meng, and Z. Lin. Investigating bi-level optimization for learning and vision from a unified perspective: A survey and beyond. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2021.
- [35] D. Lopez-Paz and M. Ranzato. Gradient episodic memory for continual learning. *Advances in neural information processing systems*, 30, 2017.
- [36] Z. Lu and S. Mei. A first-order augmented Lagrangian method for constrained minimax optimization. *arXiv preprint arXiv:2301.02060*, 2023.
- [37] Z. Lu and S. Mei. A first-order method for nonconvex-strongly-concave constrained minimax optimization. *Preprint*, 2023. <https://zhaosong-lu.github.io/ResearchPapers/strongly-cvx-minimax.pdf>.
- [38] Z. Lu and S. Mei. First-order penalty methods for bilevel optimization. *arXiv preprint arXiv:2301.01716*, 2023.
- [39] Z. Lu and S. Mei. Solving bilevel optimization via sequential minimax optimization. *Preprint*, 2023.
- [40] Z. Lu and Z. Zhou. Iteration complexity of first-order augmented Lagrangian methods for convex conic programming. *SIAM journal on optimization*, 33(2):1159–1190, 2023.
- [41] Z.-Q. Luo, J.-S. Pang, and D. Ralph. *Mathematical programs with equilibrium constraints*. Cambridge University Press, 1996.

- [42] X. Ma, W. Yao, J. J. Ye, and J. Zhang. Combined approach with second-order optimality conditions for bilevel programming problems. 2023. To appear in *Journal of Convex Analysis*.
- [43] D. Maclaurin, D. Duvenaud, and R. Adams. Gradient-based hyperparameter optimization through reversible learning. In *International conference on machine learning*, pages 2113–2122, 2015.
- [44] A. Madry, A. Makelov, L. Schmidt, D. Tsipras, and A. Vladu. Towards deep learning models resistant to adversarial attacks. In *International Conference on Learning Representations*, 2018.
- [45] J. A. Mirrlees. The theory of moral hazard and unobservable behaviour: Part I. *The Review of Economic Studies*, 66(1):3–21, 1999.
- [46] Y. Nesterov. Gradient methods for minimizing composite functions. *Mathematical programming*, 140(1):125–161, 2013.
- [47] Y. E. Nesterov. *Introductory Lectures on Convex Optimization: a basic course*. Kluwer Academic Publishers, Massachusetts, 2004.
- [48] T. Okuno, A. Takeda, A. Kawana, and M. Watanabe. On ℓ_p -hyperparameter learning via bilevel nonsmooth optimization. *Journal of Machine Learning Research*, 22(1):11093–11139, 2021.
- [49] J. Outrata, M. Kocvara, and J. Zowe. *Nonsmooth approach to optimization problems with equilibrium constraints: theory, applications and numerical results*, volume 28. Springer Science & Business Media, 2013.
- [50] F. Pedregosa. Hyperparameter optimization with approximate gradient. In *International conference on machine learning*, pages 737–746, 2016.
- [51] A. Rajeswaran, C. Finn, S. M. Kakade, and S. Levine. Meta-learning with implicit gradients. *Advances in neural information processing systems*, 32, 2019.
- [52] H. Shen and T. Chen. On penalty-based bilevel gradient descent method. *arXiv preprint arXiv:2302.05185*, 2023.
- [53] C. Shi, J. Lu, and G. Zhang. An extended Kuhn–Tucker approach for linear bilevel programming. *Applied Mathematics and Computation*, 162(1):51–63, 2005.
- [54] K. Shimizu, Y. Ishizuka, and J. F. Bard. *Nondifferentiable and two-level mathematical programming*. Springer Science & Business Media, 2012.
- [55] A. Sinha, P. Malo, and K. Deb. A review on bilevel optimization: from classical to evolutionary approaches and applications. *IEEE Transactions on Evolutionary Computation*, 22(2):276–295, 2017.
- [56] D. Sow, K. Ji, Z. Guan, and Y. Liang. A primal-dual approach to bilevel optimization with multiple inner minima. *arXiv preprint arXiv:2203.01123*, 2022.
- [57] C. Szegedy, W. Zaremba, I. Sutskever, J. Bruna, D. Erhan, I. Goodfellow, and R. Fergus. Intriguing properties of neural networks. *arXiv preprint arXiv:1312.6199*, 2013.
- [58] L. N. Vicente and P. H. Calamai. Bilevel and multilevel programming: A bibliography review. *Journal of Global optimization*, 5(3):291–306, 1994.
- [59] H. Von Stackelberg. *Market structure and equilibrium*. Springer Science & Business Media, 2010.
- [60] X. Wang, R. Pan, R. Pi, and T. Zhang. Effective bilevel optimization via minimax reformulation. *arXiv preprint arXiv:2305.13153*, 2023.
- [61] J. Yang, K. Ji, and Y. Liang. Provably faster algorithms for bilevel optimization. *Advances in Neural Information Processing Systems*, 34:13670–13682, 2021.
- [62] J. J. Ye. Constraint qualifications and optimality conditions in bilevel optimization. In *Bilevel Optimization*, pages 227–251. Springer, 2020.
- [63] J. J. Ye, X. Yuan, S. Zeng, and J. Zhang. Difference of convex algorithms for bilevel programs with applications in hyperparameter selection. *Mathematical Programming*, pages 1–34, 2022.

A An optimal first-order method for unconstrained strongly convex optimization problems

In this part we review an optimal first-order method for solving a class of strongly convex optimization problems in the form of

$$\Psi^* = \min_x \{\Psi(x) := \phi(x) + P(x)\} \quad (89)$$

where $P : \mathbb{R}^m \rightarrow (-\infty, \infty]$ is a closed convex function, $\phi : \mathbb{R}^m \rightarrow (-\infty, \infty]$ is a σ_ϕ -strongly-convex function, and $\nabla\phi$ is $L_{\nabla\phi}$ -Lipschitz continuous on $\text{dom } P$. In addition, we assume that $\text{dom } P$ is compact and let $D_P := \max_{x,y \in \text{dom } P} \|x - y\|$.

We next present an optimal first-order method [40, Algorithm 4] for solving problem (89), which is a slight variant of Nesterov's optimal first-order methods [31, 47].

Algorithm 3 An optimal first-order method for problem (89)

Input: $\tilde{\epsilon} > 0$ and $\tilde{x}^0 \in \text{dom } P$.

1: Compute

$$x^0 = \text{prox}_{P/L_{\nabla\phi}} \left(\tilde{x}^0 - \frac{1}{L_{\nabla\phi}} \nabla\phi(\tilde{x}^0) \right).$$

2: Set $z^0 = x^0$ and $\alpha = \sqrt{\sigma_\phi/L_{\nabla\phi}}$.

3: **for** $k = 0, 1, \dots$ **do**

4: Set $y^k = (x^k + \alpha z^k)/(1 + \alpha)$.

5: Compute z^{k+1} as

$$z^{k+1} = \arg \min_z \left\{ \ell(z; y^k) + \frac{\alpha L_{\nabla\phi}}{2} \|z - \alpha y^k - (1 - \alpha) z^k\|^2 \right\},$$

where

$$\ell(x; y) := \phi(y) + \langle \nabla\phi(y), x - y \rangle + P(x).$$

6: Set $x^{k+1} = (1 - \alpha)x^k + \alpha z^{k+1}$.

7: Compute

$$\tilde{x}^{k+1} = \text{prox}_{P/L_{\nabla\phi}} \left(x^{k+1} - \frac{1}{L_{\nabla\phi}} \nabla\phi(x^{k+1}) \right).$$

8: Terminate the algorithm and output \tilde{x}^{k+1} if

$$\|\tilde{x}^{k+1} - x^k\| \leq \frac{\tilde{\epsilon}}{2L_{\nabla\phi}}.$$

9: **end for**

The following result provides an *iteration-complexity* of Algorithm 3 for finding an approximate optimal solution of problem (89), which was established in [40, Proposition 4].

Theorem 3. *Let $\{\tilde{x}^k\}$ be the sequence generated by Algorithm 3. Then for any given $\tilde{\epsilon} > 0$, an approximate solution \tilde{x}^{k+1} of problem (89) satisfying $\text{dist}(0, \partial\Psi(\tilde{x}^{k+1})) \leq 2L_{\nabla\phi}\|\tilde{x}^{k+1} - x^{k+1}\| \leq \tilde{\epsilon}$ is generated by running Algorithm 3 for at most \tilde{T} iterations, where*

$$\tilde{T} = \left\lceil \sqrt{\frac{L_{\nabla\phi}}{\sigma_\phi}} \max \left\{ 1, \left\lceil 2 \log \frac{2L_{\nabla\phi}D_P}{\tilde{\epsilon}} \right\rceil \right\} \right\rceil.$$

Remark 5. *By the convexity of Ψ , $D_P = \max_{x,y \in \text{dom } P} \|x - y\|$, and Theorem 3, it is not hard to show that the output \tilde{x}^{k+1} of Algorithm 3 satisfies*

$$\Psi(\tilde{x}^{k+1}) - \Psi^* \leq \text{dist}(0, \partial\Psi(\tilde{x}^{k+1}))D_P \leq \tilde{\epsilon}D_P. \quad (90)$$

B A first-order method for nonconvex-strongly-concave minimax problem

In this part we present a first-order method proposed in [36, Algorithm 2] for finding an ϵ -primal-dual stationary point of the nonconvex-concave minimax problem

$$H^* = \min_x \max_y \{H(x, y) := h(x, y) + p(x) - q(y)\}, \quad (91)$$

which has at least one optimal solution and satisfies the following assumptions.

Assumption 3. (i) $p : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{\infty\}$ and $q : \mathbb{R}^m \rightarrow \mathbb{R} \cup \{\infty\}$ are proper convex functions and continuous on $\text{dom } p$ and $\text{dom } q$, respectively, and moreover, $\text{dom } p$ and $\text{dom } q$ are compact.

(ii) The proximal operator associated with p and q can be exactly evaluated.

(iii) h is $L_{\nabla h}$ -smooth on $\text{dom } p \times \text{dom } q$, and moreover, $h(x, \cdot)$ is σ_y -strongly-concave for any $x \in \text{dom } p$.

For ease of presentation, we define

$$D_p = \max\{\|u - v\| \mid u, v \in \text{dom } p\}, \quad D_q = \max\{\|u - v\| \mid u, v \in \text{dom } q\}, \quad (92)$$

$$H_{\text{low}} = \min\{H(x, y) \mid (x, y) \in \text{dom } p \times \text{dom } q\}. \quad (93)$$

Given an iterate (x^k, y^k) , the first-order method [36, Algorithm 2] finds the next iterate (x^{k+1}, y^{k+1}) by applying a modified optimal first-order method [36, Algorithm 1] to the strongly-convex-strongly-concave minimax problem

$$\min_x \max_y \{h_k(x, y) = h(x, y) + L_{\nabla h} \|x - x^k\|^2\}. \quad (94)$$

For ease reference, we next present the modified optimal first-order method [36, Algorithm 1] in Algorithm 4 below for solving the strongly-convex-strongly-concave minimax problem

$$\min_x \max_y \{\bar{h}(x, y) + p(x) - q(y)\}, \quad (95)$$

where $\bar{h}(x, y)$ is $\bar{\sigma}_x$ -strongly-convex- $\bar{\sigma}_y$ -strongly-concave and $L_{\nabla \bar{h}}$ -smooth on $\text{dom } p \times \text{dom } q$ for some $\bar{\sigma}_x, \bar{\sigma}_y > 0$. In Algorithm 4, the functions \hat{h} , a_x^k and a_y^k are defined as follows:

$$\begin{aligned} \hat{h}(x, y) &= \bar{h}(x, y) - \bar{\sigma}_x \|x\|^2/2 + \bar{\sigma}_y \|y\|^2/2, \\ a_x^k(x, y) &= \nabla_x \hat{h}(x, y) + \bar{\sigma}_x (x - \bar{\sigma}_x^{-1} z_g^k)/2, \quad a_y^k(x, y) = -\nabla_y \hat{h}(x, y) + \bar{\sigma}_y y + \bar{\sigma}_x (y - y_g^k)/8, \end{aligned}$$

where y_g^k and z_g^k are generated at iteration k of Algorithm 4 below.

Algorithm 4 A modified optimal first-order method for problem (95)

Input: $\tau > 0$, $\bar{z}^0 = z_f^0 \in -\bar{\sigma}_x \text{dom } p$,⁴ $\bar{y}^0 = y_f^0 \in \text{dom } q$, $(z^0, y^0) = (\bar{z}^0, \bar{y}^0)$, $\bar{\alpha} = \min \left\{ 1, \sqrt{8\bar{\sigma}_y/\bar{\sigma}_x} \right\}$, $\eta_z = \bar{\sigma}_x/2$, $\eta_y = \min \{1/(2\bar{\sigma}_y), 4/(\bar{\alpha}\bar{\sigma}_x)\}$, $\beta_t = 2/(t+3)$, $\zeta = (2\sqrt{5}(1+8L_{\nabla\bar{h}}/\bar{\sigma}_x))^{-1}$, $\gamma_x = \gamma_y = 8\bar{\sigma}_x^{-1}$, and $\hat{\zeta} = \min\{\bar{\sigma}_x, \bar{\sigma}_y\}/L_{\nabla\bar{h}}^2$.

- 1: **for** $k = 0, 1, 2, \dots$ **do**
- 2: $(z_g^k, y_g^k) = \bar{\alpha}(z^k, y^k) + (1 - \bar{\alpha})(z_f^k, y_f^k)$.
- 3: $(x^{k,-1}, y^{k,-1}) = (-\bar{\sigma}_x^{-1}z_g^k, y_g^k)$.
- 4: $x^{k,0} = \text{prox}_{\zeta\gamma_x p}(x^{k,-1} - \zeta\gamma_x a_x^k(x^{k,-1}, y^{k,-1}))$.
- 5: $y^{k,0} = \text{prox}_{\zeta\gamma_y q}(y^{k,-1} - \zeta\gamma_y a_y^k(x^{k,-1}, y^{k,-1}))$.
- 6: $b_x^{k,0} = \frac{1}{\zeta\gamma_x}(x^{k,-1} - \zeta\gamma_x a_x^k(x^{k,-1}, y^{k,-1}) - x^{k,0})$.
- 7: $b_y^{k,0} = \frac{1}{\zeta\gamma_y}(y^{k,-1} - \zeta\gamma_y a_y^k(x^{k,-1}, y^{k,-1}) - y^{k,0})$.
- 8: $t = 0$.
- 9: **while**
- 10: $\gamma_x \|a_x^k(x^{k,t}, y^{k,t}) + b_x^{k,t}\|^2 + \gamma_y \|a_y^k(x^{k,t}, y^{k,t}) + b_y^{k,t}\|^2 > \gamma_x^{-1} \|x^{k,t} - x^{k,-1}\|^2 + \gamma_y^{-1} \|y^{k,t} - y^{k,-1}\|^2$
- 11: **do**
- 12: $x^{k,t+1/2} = x^{k,t} + \beta_t(x^{k,0} - x^{k,t}) - \zeta\gamma_x(a_x^k(x^{k,t}, y^{k,t}) + b_x^{k,t})$.
- 13: $y^{k,t+1/2} = y^{k,t} + \beta_t(y^{k,0} - y^{k,t}) - \zeta\gamma_y(a_y^k(x^{k,t}, y^{k,t}) + b_y^{k,t})$.
- 14: $x^{k,t+1} = \text{prox}_{\zeta\gamma_x p}(x^{k,t} + \beta_t(x^{k,0} - x^{k,t}) - \zeta\gamma_x a_x^k(x^{k,t+1/2}, y^{k,t+1/2}))$.
- 15: $y^{k,t+1} = \text{prox}_{\zeta\gamma_y q}(y^{k,t} + \beta_t(y^{k,0} - y^{k,t}) - \zeta\gamma_y a_y^k(x^{k,t+1/2}, y^{k,t+1/2}))$.
- 16: $b_x^{k,t+1} = \frac{1}{\zeta\gamma_x}(x^{k,t} + \beta_t(x^{k,0} - x^{k,t}) - \zeta\gamma_x a_x^k(x^{k,t+1/2}, y^{k,t+1/2}) - x^{k,t+1})$.
- 17: $b_y^{k,t+1} = \frac{1}{\zeta\gamma_y}(y^{k,t} + \beta_t(y^{k,0} - y^{k,t}) - \zeta\gamma_y a_y^k(x^{k,t+1/2}, y^{k,t+1/2}) - y^{k,t+1})$.
- 18: $t \leftarrow t + 1$.
- 19: **end while**
- 20: $(x_f^{k+1}, y_f^{k+1}) = (x^{k,t}, y^{k,t})$.
- 21: $(z_f^{k+1}, w_f^{k+1}) = (\nabla_x \hat{h}(x_f^{k+1}, y_f^{k+1}) + b_x^{k,t}, -\nabla_y \hat{h}(x_f^{k+1}, y_f^{k+1}) + b_y^{k,t})$.
- 22: $z^{k+1} = z^k + \eta_z \bar{\sigma}_x^{-1}(z_f^{k+1} - z^k) - \eta_z(x_f^{k+1} + \bar{\sigma}_x^{-1}z_f^{k+1})$.
- 23: $y^{k+1} = y^k + \eta_y \bar{\sigma}_y(y_f^{k+1} - y^k) - \eta_y(w_f^{k+1} + \bar{\sigma}_y y_f^{k+1})$.
- 24: $x^{k+1} = -\bar{\sigma}_x^{-1}z^{k+1}$.
- 25: $\hat{x}^{k+1} = \text{prox}_{\hat{\zeta}p}(x^{k+1} - \hat{\zeta}\nabla_x \bar{h}(x^{k+1}, y^{k+1}))$.
- 26: $\hat{y}^{k+1} = \text{prox}_{\hat{\zeta}q}(y^{k+1} + \hat{\zeta}\nabla_y \bar{h}(x^{k+1}, y^{k+1}))$.
- 27: Terminate the algorithm and output $(\hat{x}^{k+1}, \hat{y}^{k+1})$ if

$$\|\hat{\zeta}^{-1}(x^{k+1} - \hat{x}^{k+1}, \hat{y}^{k+1} - y^{k+1}) - (\nabla \bar{h}(x^{k+1}, y^{k+1}) - \nabla \bar{h}(\hat{x}^{k+1}, \hat{y}^{k+1}))\| \leq \tau.$$

26: **end for**

We are now ready to present the first-order method [36, Algorithm 2] for finding an ϵ -primal-dual stationary point of (91) in Algorithm 5 below.

Algorithm 5 A first-order method for problem (91)

Input: $\epsilon > 0$, $\tilde{\epsilon}_0 \in (0, \epsilon/2]$, $(\tilde{x}^0, \tilde{y}^0) \in \text{dom } p \times \text{dom } q$, $(x^0, y^0) = (\tilde{x}^0, \tilde{y}^0)$, and $\tilde{\epsilon}_k = \tilde{\epsilon}_0/(k+1)$.

- 1: **for** $k = 0, 1, 2, \dots$ **do**
- 2: Call Algorithm 4 with $\bar{h} \leftarrow h_k$, $\tau \leftarrow \epsilon_k$, $\bar{\sigma}_x \leftarrow L_{\nabla h}$, $\bar{\sigma}_y \leftarrow \sigma_y$, $L_{\nabla \bar{h}} \leftarrow 3L_{\nabla h}$, $\bar{z}^0 = z_f^0 \leftarrow -\bar{\sigma}_x x^k$, $\bar{y}^0 = y_f^0 \leftarrow y^k$, and denote its output by (x^{k+1}, y^{k+1}) , where h_k is given in (94).
- 3: Terminate the algorithm and output $(x_\epsilon, y_\epsilon) = (x^{k+1}, y^{k+1})$ if

$$\|x^{k+1} - x^k\| \leq \epsilon/(4L_{\nabla h}).$$

4: **end for**

The following theorem presents the iteration complexity of Algorithm 5, whose proof is given in [36, Theorem 2].

⁴For convenience, $-\bar{\sigma}_x \text{dom } p$ stands for the set $\{-\bar{\sigma}_x u | u \in \text{dom } p\}$.

Theorem 4 (Complexity of Algorithm 5). *Suppose that Assumption 3 holds. Let H^* , H , D_p , D_q , and H_{low} be defined in (91), (92) and (93), $L_{\nabla h}$ be given in Assumption 3, ϵ , ϵ_0 and x^0 be given in Algorithm 5, and*

$$\begin{aligned}
\tilde{\alpha} &= \min \left\{ 1, \sqrt{8\sigma_y/L_{\nabla h}} \right\}, \\
\tilde{\delta} &= (2 + \tilde{\alpha}^{-1})L_{\nabla h}D_p^2 + \max \{ 2\sigma_y, \tilde{\alpha}L_{\nabla h}/4 \} D_q^2, \\
\tilde{K} &= \left\lceil 16(\max_y H(\tilde{x}^0, y) - H^*)L_{\nabla h}\epsilon^{-2} + 32\tilde{\epsilon}_0^2(1 + \sigma_y^{-2}L_{\nabla h}^2)\epsilon^{-2} - 1 \right\rceil_+, \\
\tilde{N} &= 3397 \max \left\{ 2, \sqrt{L_{\nabla h}/(2\sigma_y)} \right\} \\
&\quad \times \left((\tilde{K} + 1) \left(\log \frac{4 \max \left\{ \frac{1}{2L_{\nabla h}}, \min \left\{ \frac{1}{2\sigma_y}, \frac{4}{\tilde{\alpha}L_{\nabla h}} \right\} \right\} \left(\tilde{\delta} + 2\tilde{\alpha}^{-1}(H^* - H_{\text{low}} + L_{\nabla h}D_p^2) \right)}{[9L_{\nabla h}^2/\min\{L_{\nabla h}, \sigma_y\} + 3L_{\nabla h}]^{-2} \tilde{\epsilon}_0^2} \right) \right)_+ \\
&\quad + \tilde{K} + 1 + 2\tilde{K} \log(\tilde{K} + 1) \right).
\end{aligned}$$

Then Algorithm 5 terminates and outputs an ϵ -primal-dual stationary point (x_ϵ, y_ϵ) of (91) in at most $\tilde{K} + 1$ outer iterations that satisfies

$$\max_y H(x_\epsilon, y) \leq \max_y H(\tilde{x}^0, y) + 2\tilde{\epsilon}_0^2 (L_{\nabla h}^{-1} + \sigma_y^{-2}L_{\nabla h}). \quad (96)$$

Moreover, the total number of evaluations of ∇h and proximal operator of p and q performed in Algorithm 5 is no more than \tilde{N} , respectively.