

# AN INEXACT PRIMAL-DUAL SMOOTHING FRAMEWORK FOR LARGE-SCALE NON-BILINEAR SADDLE POINT PROBLEMS\*

LE THI KHANH HIEN<sup>†</sup>, RENBO ZHAO<sup>‡</sup>, AND WILLIAM B. HASKELL<sup>§</sup>

**Abstract.** We develop an inexact primal-dual first-order smoothing framework to solve a class of non-bilinear saddle point problems with primal strong convexity. Compared with existing methods, our framework yields a significant improvement over the primal oracle complexity, while it has competitive dual oracle complexity. In addition, we consider the situation where the primal-dual coupling term has a large number of component functions. To efficiently handle this situation, we develop a randomized version of our smoothing framework, which allows the primal and dual subproblems in each iteration to be solved by randomized algorithms inexactly in expectation. The convergence of this framework is analyzed both in expectation and with high probability. In terms of the primal and dual oracle complexities, this framework significantly improves over its deterministic counterpart. As an important application, we adapt both frameworks for solving convex optimization problems with many functional constraints. To obtain an  $\varepsilon$ -optimal and  $\varepsilon$ -feasible solution, both frameworks achieve the best-known oracle complexities (in terms of their dependence on  $\varepsilon$ ).

**Key words.** Non-bilinear saddle point problems, Inexact primal-dual smoothing, Convex optimization with functional constraints, Stochastic optimization, Large-scale optimization

**1. Introduction.** Let  $(\mathbb{E}_1, \|\cdot\|_{\mathbb{E}_1})$  and  $(\mathbb{E}_2, \|\cdot\|_{\mathbb{E}_2})$  be finite-dimensional real normed spaces, with dual spaces  $(\mathbb{E}_1^*, \|\cdot\|_{\mathbb{E}_1^*})$  and  $(\mathbb{E}_2^*, \|\cdot\|_{\mathbb{E}_2^*})$ , respectively. We consider the following convex-concave saddle point problem (SPP)

$$(1.1) \quad \min_{x \in \mathcal{X}} \max_{\lambda \in \Lambda} \{S(x, \lambda) \triangleq f(x) + g(x) + \Phi(x, \lambda) - h(\lambda)\},$$

where  $\mathcal{X} \subseteq \mathbb{E}_1$  and  $\Lambda \subseteq \mathbb{E}_2$  are nonempty, convex and closed sets,  $\mathcal{X}$  is bounded and the functions  $f, g : \mathbb{E}_1 \rightarrow \mathbb{R} \triangleq \mathbb{R} \cup \{+\infty\}$  and  $h : \mathbb{E}_2 \rightarrow \mathbb{R}$  are convex, closed and proper (CCP) functions. Define  $\text{dom } g \triangleq \{x \in \mathbb{E}_1 : g(x) < +\infty\}$  and  $\text{dom } h$  similarly. We assume that  $\mathcal{X} \subseteq \text{dom } g$  and  $\Lambda \subseteq \text{dom } h$ . (Otherwise, we can take  $\mathcal{X} \cap \text{dom } g$  and  $\Lambda \cap \text{dom } h$  to be the new constraint sets that satisfy the above assumptions on  $\mathcal{X}$  and  $\Lambda$ .) We also assume that both  $g$  and  $h$  admit tractable Bregman proximal projections (BPP) on  $\mathcal{X}$  and  $\Lambda$ , respectively. (See Section 2 for its precise definition.) In addition, we assume that  $f$  is differentiable on an open set  $\mathcal{X}' \supseteq \mathcal{X}$ , and  $\mu$ -strongly convex (s.c.) and  $L$ -smooth on  $\mathcal{X}$  (where  $L \geq \mu > 0$ ), i.e.,

$$(1.2) \quad \frac{\mu}{2} \|x - x'\|_{\mathbb{E}_1}^2 \leq f(x) - f(x') - \langle \nabla f(x'), x - x' \rangle \leq \frac{L}{2} \|x - x'\|_{\mathbb{E}_1}^2, \quad \forall x, x' \in \mathcal{X},$$

where  $\langle \cdot, \cdot \rangle$  denotes the duality pairing between  $\mathbb{E}_1^*$  (resp.  $\mathbb{E}_2^*$ ) and  $\mathbb{E}_1$  (resp.  $\mathbb{E}_2$ ).

We next state our assumptions on the function  $\Phi : \mathbb{E}_1 \times \mathbb{E}_2 \rightarrow [-\infty, +\infty]$ . First, it is convex-concave, i.e., for any  $(x, \lambda) \in \mathbb{E}_1 \times \mathbb{E}_2$ ,  $\Phi(\cdot, \lambda)$  is convex and differentiable on  $\mathbb{E}_1$  and  $\Phi(x, \cdot)$  is concave and differentiable on  $\mathbb{E}_2$ . In addition,  $\Phi$  satisfies the  $(L_{xx}, L_{\lambda x}, L_{\lambda\lambda})$ -smoothness condition (where  $L_{xx}, L_{\lambda x}, L_{\lambda\lambda} \geq 0$ ) on  $\mathcal{X} \times \Lambda$ , i.e., for

---

\* L. T. K. Hien and R. Zhao contribute equally to this work.

**Funding:** This work was funded by A\*STAR Project Number 1421200078 “ABCD: Analyzing Big Corrupted Data”.

<sup>†</sup>Department of Mathematics and Operations Research, University of Mons, Belgium ([thikhanh.hien@umons.ac.be](mailto:thikhanh.hien@umons.ac.be)).

<sup>‡</sup>Operations Research Center, Massachusetts Institute of Technology, USA ([renboz@mit.edu](mailto:renboz@mit.edu)).

<sup>§</sup>Department of Industrial Systems Engineering and Management, National University of Singapore, Singapore ([ischwab@nus.edu.sg](mailto:ischwab@nus.edu.sg)).

any  $x, x' \in \mathcal{X}$  and  $\lambda, \lambda' \in \Lambda$ ,

$$(1.3a) \quad \|\nabla_x \Phi(x, \lambda) - \nabla_x \Phi(x', \lambda)\|_{\mathbb{E}_1^*} \leq L_{xx} \|x - x'\|_{\mathbb{E}_1},$$

$$(1.3b) \quad \|\nabla_x \Phi(x, \lambda) - \nabla_x \Phi(x, \lambda')\|_{\mathbb{E}_1^*} \leq L_{\lambda x} \|\lambda - \lambda'\|_{\mathbb{E}_2},$$

$$(1.3c) \quad \|\nabla_\lambda \Phi(x, \lambda) - \nabla_\lambda \Phi(x', \lambda)\|_{\mathbb{E}_2^*} \leq L_{\lambda x} \|x - x'\|_{\mathbb{E}_1},$$

$$(1.3d) \quad \|\nabla_\lambda \Phi(x, \lambda) - \nabla_\lambda \Phi(x, \lambda')\|_{\mathbb{E}_2^*} \leq L_{\lambda\lambda} \|\lambda - \lambda'\|_{\mathbb{E}_2},$$

where  $x \mapsto \nabla_x \Phi(x, \lambda)$  and  $\lambda \mapsto \nabla_\lambda \Phi(x, \lambda)$  denote the gradients of  $\Phi(\cdot, \lambda)$  and  $\Phi(x, \cdot)$ , respectively. For later use, let us define the (primal) condition number

$$(1.4) \quad \kappa_{\mathcal{X}} \triangleq (L + L_{xx})/\mu.$$

In this work, we assume that the saddle function  $\Phi(\cdot, \cdot)$  in (1.1) has the following finite-sum structure, i.e.,

$$(1.5) \quad \Phi(x, \lambda) \triangleq (1/n) \sum_{i=1}^n \Phi_i(x, \lambda),$$

where for each  $i \in [n] \triangleq \{1, \dots, n\}$  and any  $(x, y) \in \mathbb{E}_1 \times \mathbb{E}_2$ ,  $\Phi_i : \mathbb{E}_1 \times \mathbb{E}_2 \rightarrow [-\infty, +\infty]$  is convex-concave and satisfies the  $(L_{xx}^i, L_{\lambda x}^i, L_{\lambda\lambda}^i)$ -smoothness condition on  $\mathcal{X} \times \Lambda$ . As a result, the smoothness parameters of  $\Phi$  can be bounded as  $L_{xx} \leq (1/n) \sum_{i=1}^n L_{xx}^i$ ,  $L_{\lambda x} \leq (1/n) \sum_{i=1}^n L_{\lambda x}^i$  and  $L_{\lambda\lambda} \leq (1/n) \sum_{i=1}^n L_{\lambda\lambda}^i$ . In addition, we are particularly interested in the setting where the number of component functions (i.e.,  $n$ ) is *large*.

For well-posedness, we assume that for the SPP in (1.1), at least one saddle point  $(x^*, \lambda^*)$  exists, i.e., there exists  $(x^*, \lambda^*) \in \mathcal{X} \times \Lambda$  such that

$$(1.6) \quad S(x^*, \lambda) \leq S(x^*, \lambda^*) \leq S(x, \lambda^*), \quad \forall (x, \lambda) \in \mathcal{X} \times \Lambda.$$

**1.1. Primal and Dual First-order Oracles.** Since we are interested in developing primal-dual first-order methods to solve the SPP in (1.1), where the function  $\Phi(\cdot, \cdot)$  has the finite-sum structure as in (1.5), we set up the primal and dual first-order oracles as follows: Upon receiving  $(x, \lambda, i) \in \mathcal{X} \times \Lambda \times [n]$  (where  $[n] \triangleq \{1, \dots, n\}$ ), the primal oracle  $\mathcal{O}^P$  returns  $\nabla_x \Phi_i(x, \lambda)$  and the dual oracle  $\mathcal{O}^D$  returns  $\nabla_\lambda \Phi_i(x, \lambda)$ . In addition,  $\mathcal{O}^P$  returns  $\nabla f(x)$  upon receiving  $(x, 0)$ . Accordingly, we define the *primal oracle complexity* and *dual oracle complexity* to be the number of oracle calls to  $\mathcal{O}^P$  and  $\mathcal{O}^D$ , respectively.

**1.2. Applications.** The SPP in (1.1), where  $\Phi(\cdot, \cdot)$  has the finite-sum structure as in (1.5), has applications in numerous fields, including game theory, image processing, machine learning and statistics. In Section 5, we will focus on one of the most important ones, i.e., convex optimization problems with functional constraints (and their stochastic extension). Apart from this, we also illustrate another important application below. For more applications, we refer readers to [12, 9, 3, 11].

*Maximum Margin Clustering* [30]. Let  $\mathcal{D}$  denote a set of  $m$  objects, which consists of two unknown disjoint subsets, denoted by  $\mathcal{D}_1$  and  $\mathcal{D}_2$ , respectively (i.e.,  $\mathcal{D}_1 \cap \mathcal{D}_2 = \emptyset$  and  $\mathcal{D}_1 \cup \mathcal{D}_2 = \mathcal{D}$ ). Given  $n$  noisy samples  $\{\mathcal{S}_i\}_{i=1}^n$  of  $\mathcal{D}$  (where each  $\mathcal{S}_i$  consists of  $m$  data points), we wish to find a label kernel matrix  $M \in \mathbb{S}_+^m$  that can assign each object in  $\mathcal{D}$  to  $\mathcal{D}_1$  or  $\mathcal{D}_2$  (where  $\mathbb{S}_+^m$  denotes the positive semi-definite cone of dimension  $m$ ). To do so, we first compute a kernel matrix  $K_i \in \mathbb{S}_+^m$  from each  $\mathcal{S}_i$  ( $i \in [n]$ ). Define  $e \triangleq (1, 1, \dots, 1) \in \mathbb{R}^m$ , we then solve the following SPP:

$$\min_{M \in \mathcal{M}} \max_{\lambda_i \in \Lambda, \forall i \in [n]} (1/n) \sum_{i=1}^n - \langle \lambda_i \lambda_i^T, K_i \circ M + \alpha I \rangle + 2\lambda_i^T e,$$

TABLE 1  
*Comparison of primal and dual oracle complexities with existing methods.*

Algorithms	Primal Oracle Comp.	Dual Oracle Comp.
PDHG-type [11]	$O(n/\varepsilon)$	$O(n/\varepsilon)$
Mirror-Prox [20]	$O(n/\varepsilon)$	$O(n/\varepsilon)$
Det. IPDS (Algo. 1)	$\tilde{O}(n\sqrt{\kappa_{\mathcal{X}}}/\varepsilon)$	$\tilde{O}(n/\varepsilon)$
Rand. IPDS (Algo. 2) <sup>1</sup>	$\tilde{O}((n + \sqrt{n\kappa_{\mathcal{X}}})/\sqrt{\varepsilon})$	$\tilde{O}(n/\sqrt{\varepsilon} + \sqrt{n}/\varepsilon)$

<sup>1</sup> Both primal and dual oracle complexities of Rand. IPDS correspond to obtaining *expected* duality gap (cf. Theorem 4.2).

where  $\Lambda \triangleq \{\lambda \in \mathbb{R}^m : 0 \leq \lambda_i \leq C < +\infty, \forall i \in [m]\}$ ,  $\mathcal{M} \triangleq \{M \in \mathbb{S}_+^m : \text{diag}(M) = e, |e^T M| \leq \ell < +\infty\}$  and  $\alpha > 0$ . In addition,  $\circ$  denotes the Hadamard product and  $\langle \cdot, \cdot \rangle$  denotes the Frobenius inner product. We can easily see that (1.2) is a convex-concave SPP with finite-sum structure. Moreover, both constraint sets  $\mathcal{M}$  and  $\Lambda^n$  are nonempty, closed and convex. In addition, to achieve high clustering accuracy, the required number of samples  $n$  may potentially be large.

**1.3. Related Works.** We focus on reviewing the works on solving (convex-concave) *non-bilinear non-smooth* SPPs, where the primal-dual coupling term  $\Phi(\cdot, \cdot)$  is not bilinear. (Note that  $\Phi(\cdot, \cdot)$  is bilinear if  $\Phi(x, y) = \langle Ax, y \rangle$ , where  $A : \mathbb{E}_1 \rightarrow \mathbb{E}_2^*$  is a (bounded) linear operator.) In the past few years, bilinear SPPs have been extensively studied, with many efficient algorithms proposed, e.g., Nesterov smoothing [22] and primal-dual hybrid gradient (PDHG) [7]. For details, see [8, 32].

For non-bilinear non-smooth SPPs without composite structure, i.e., the saddle function  $S(\cdot, \cdot)$  cannot be decomposed into a smooth term  $\Phi(\cdot, \cdot)$  and non-smooth functions  $g$  and  $h$  with tractable BPPs, the existing algorithms are mainly based on primal-dual subgradient (e.g., [19, 23]). To reach an  $\varepsilon$ -duality gap (defined in (2.13)), these works achieve the optimal oracle complexity  $O(1/\varepsilon^2)$ .

However, if the non-smooth SPPs do possess composite structure, e.g., the one in (1.1), then the algorithms for smooth SPPs can be extended to this case and their oracle complexities can be much better than  $O(1/\varepsilon^2)$ . For example, the Mirror-Prox method [20] was extended in [12] to solve (1.1) with an oracle complexity of  $O(1/\varepsilon)$ . As another example, Nesterov smoothing was extended in [14] to solve a special case of (1.1), where  $L_{\lambda\lambda} = 0$  (i.e.,  $\Phi(x, \cdot)$  is linear). Recently, the PDHG method, originally developed for the bilinear SPPs, was extended in [11] to solve (1.1). In addition, the stochastic (and accelerated) versions of the aforementioned methods have also been developed in [13, 9, 32] to tackle the situation where only stochastic first-order oracles are available. (See Section 4.5 for details.) However, most of these methods focus on the case where  $\mu = 0$ , so the favorable condition  $\mu > 0$  may not be exploited (to further improve the oracle complexity). The only two works that have utilized this condition are [12, 11], which are based on the Mirror-Prox and PDHG methods, respectively. However, both works only consider the special case where  $L_{\lambda\lambda} = 0$ . Therefore, in this work, we aim to answer the following question:

*Can we develop an algorithmic framework that works for both  $\mu > 0$  and  $L_{\lambda\lambda} \geq 0$ , yet with improved oracle complexity over the existing methods?*

**1.4. Main Contributions.** We make the following three main contributions. First, we develop a novel (deterministic) inexact primal-dual smoothing (IPDS)

framework (i.e., Algorithm 1) for solving the non-bilinear SPP in (1.1) with primal strong convexity (i.e.,  $\mu > 0$ ). To the best of our knowledge, this is the first inexact smoothing framework developed for such a problem. To reach an  $\varepsilon$ -duality gap (defined in (2.13)), the primal and dual oracle complexities are  $\tilde{O}(n\sqrt{\kappa_{\mathcal{X}}/\varepsilon})$  and  $\tilde{O}(n/\varepsilon)$ , respectively (where  $\tilde{O}(\cdot)$  hides the  $\log n$  and  $\log(1/\varepsilon)$  factors). Compared with existing works (cf. Table 1), the primal oracle complexity of our framework is significantly better, while the dual oracle complexity is competitive. In addition, in contrast to the methods in [12, 11], which can *only* improve the primal oracle complexity  $O(n/\varepsilon)$  (to e.g.,  $O(n\sqrt{\kappa_{\mathcal{X}}/\varepsilon})$ ) when  $L_{\lambda\lambda} = 0$ , our framework applies to any value  $L_{\lambda\lambda} \geq 0$ .

Second, we develop a randomized version of our IPDS framework (i.e., Algorithm 2), by allowing each sub-problem to be solved inexactly *in expectation*. This framework is particularly useful in the regime where  $n$  is large in (1.5). Indeed, to reach an  $\varepsilon$ -expected duality gap, the primal and dual oracle complexities are  $\tilde{O}((n + \sqrt{n\kappa_{\mathcal{X}}})/\sqrt{\varepsilon})$  and  $\tilde{O}(n/\sqrt{\varepsilon} + \sqrt{n}/\varepsilon)$ , respectively, which significantly improve over those of Algorithm 1. In addition, we show that Algorithm 2 also converges *with high probability*. We believe that the techniques used in our stochastic analysis are of independent interest, as they can be applied to other inexact frameworks (e.g., the inexact augmented Lagrangian method (ALM) [29, 25, 31]).

Finally, we apply both of our aforementioned frameworks (i.e., Algorithms 1 and 2) to convex optimization problems with (potentially many) functional constraints (more precisely, to their associated Lagrangian problems). To do so, we manage to overcome two challenges: the unboundedness of the constraint set  $\Lambda$  and the dependence of the smoothness parameter  $L_{xx}$  on the dual variable  $\lambda$ . (For details, see Section 5.3.) To obtain an  $\varepsilon$ -optimal and  $\varepsilon$ -feasible solution (cf. (5.6)), both Algorithms 1 and 2 achieve the state-of-the-art (primal) oracle complexities  $\tilde{O}(1/\sqrt{\varepsilon})$ . (Note that the dual oracle complexity is  $O(1)$  in this case; see Section 5.1.) Compared to other first-order methods with similar oracle complexities but *specifically* designed for constrained convex optimization problems, our frameworks enjoy much wider applicability.

**2. Preliminaries.** We first define the primal and dual functions and the duality gap associated with the saddle function  $S(\cdot, \cdot)$ , their smoothed versions, and the inexact solutions of the optimization problems appearing in these definitions. We then prove some smoothness properties of the formerly defined quantities.

**2.1. Definitions.** First, we define the primal function  $\psi^{\text{P}} : \mathbb{E}_1 \rightarrow \overline{\mathbb{R}}$  and the dual function  $\psi^{\text{D}} : \mathbb{E}_2 \rightarrow \overline{\mathbb{R}}$  associated with  $S(\cdot, \cdot)$  as

$$(2.1) \quad \psi^{\text{P}}(x) \triangleq \sup_{\lambda \in \Lambda} S(x, \lambda) = f(x) + g(x) + \widehat{\psi}^{\text{P}}(x), \quad \forall x \in \mathbb{E}_1,$$

$$(2.2) \quad \psi^{\text{D}}(\lambda) \triangleq \inf_{x \in \mathcal{X}} S(x, \lambda) = \widehat{\psi}^{\text{D}}(\lambda) - h(\lambda), \quad \forall \lambda \in \mathbb{E}_2,$$

where

$$(2.3) \quad \widehat{\psi}^{\text{P}}(x) \triangleq \sup_{\lambda \in \Lambda} \{\widehat{S}^{\text{D}}(x, \lambda) \triangleq \Phi(x, \lambda) - h(\lambda)\}, \quad \forall x \in \mathbb{E}_1,$$

$$(2.4) \quad \widehat{\psi}^{\text{D}}(\lambda) \triangleq \inf_{x \in \mathcal{X}} \{\widehat{S}^{\text{P}}(x, \lambda) \triangleq f(x) + g(x) + \Phi(x, \lambda)\}, \quad \forall \lambda \in \mathbb{E}_2.$$

Let  $\omega : \mathbb{E}_2 \rightarrow \overline{\mathbb{R}}$  be a CCP function that is 1-s.c. and continuous on  $\Lambda$  and essentially smooth, i.e.,  $\omega$  is continuously differentiable on  $\text{int dom } \omega \neq \emptyset$ , and  $\|\nabla \omega(\lambda_k)\|_* \rightarrow +\infty$  if  $\lambda_k \rightarrow \lambda \in \text{bd dom } \omega$  [4]. (Note that for any set  $\mathcal{K}$ ,  $\text{int } \mathcal{K}$  and  $\text{bd } \mathcal{K}$  denote the interior and boundary of  $\mathcal{K}$ , respectively.) In addition, for any  $\alpha > 0$  and  $v \in \mathbb{E}_2^*$ , the following

problem (which is defined by the triple  $(\omega, h, \Lambda)$ )

$$(2.5) \quad \min_{\lambda \in \Lambda} h(\lambda) + \langle v, \lambda \rangle + \alpha^{-1} \omega(\lambda)$$

has a (unique) *easily computable* solution in  $\Lambda^\circ \triangleq \Lambda \cap \text{int dom } \omega$ . (Note that (2.5) is equivalent to BPP; see [2] for details and examples.) We call  $\omega$  a *distance generating function* (DGF) w.r.t.  $(h, \Lambda)$ . Additionally, we say that  $h$  has a tractable BPP on  $\Lambda$  if and only if such an  $\omega$  exists. Since we also assume that  $g$  has a tractable BPP on  $\mathcal{X}$  (cf. Section 1), there exists a DGF  $\bar{\omega} : \mathbb{E}_1 \rightarrow \bar{\mathbb{R}}$  w.r.t.  $(g, \mathcal{X})$ . Similar to  $\Lambda^\circ$ , we define  $\mathcal{X}^\circ \triangleq \mathcal{X} \cap \text{int dom } \bar{\omega}$ . The assumption that (2.5) can be solved easily is typical in the first-order methods for composite optimization on normed spaces (see e.g., [28]). Since we will employ such methods to solve the sub-problems in our framework (cf. Sections 3.2 and 4.1), we also make this assumption throughout the whole work.

Based on  $\omega$ , we define the dual-regularized saddle function

$$(2.6) \quad S_\rho(x, \lambda) \triangleq S(x, \lambda) - \rho \omega(\lambda),$$

where  $\rho > 0$  is the smoothing parameter. Accordingly, the primal function  $\psi_\rho^P : \mathbb{E}_1 \rightarrow \bar{\mathbb{R}}$  associated with  $S_\rho(\cdot, \cdot)$  is

$$(2.7) \quad \psi_\rho^P(x) \triangleq \sup_{\lambda \in \Lambda} S_\rho(x, \lambda) = f(x) + g(x) + \widehat{\psi}_\rho^P(x), \quad \forall x \in \mathbb{E}_1,$$

where

$$(2.8) \quad \widehat{\psi}_\rho^P(x) \triangleq \sup_{\lambda \in \Lambda} \{ \widehat{S}_\rho^D(x, \lambda) \triangleq \Phi(x, \lambda) - h(\lambda) - \rho \omega(\lambda) \}, \quad \forall x \in \mathbb{E}_1.$$

Next, we introduce the optimal solutions of the optimization problems in (2.4) and (2.8). Since  $f$  is  $\mu$ -s.c. on  $\mathcal{X}$ , the minimization problem in (2.2) has the unique solution

$$(2.9) \quad x^*(\lambda) \triangleq \arg \min_{x \in \mathcal{X}} \widehat{S}^P(x, \lambda), \quad \forall \lambda \in \mathbb{E}_2.$$

In addition, for any  $\lambda \in \mathbb{E}_2$ , we call  $\tilde{x}_\gamma(\lambda) \in \mathcal{X}$  an  $\gamma$ -inexact solution (where  $\gamma \geq 0$ ) if

$$(2.10) \quad \widehat{S}^P(\tilde{x}_\gamma(\lambda), \lambda) - \widehat{\psi}^D(\lambda) = \widehat{S}^P(\tilde{x}_\gamma(\lambda), \lambda) - \widehat{S}^P(x^*(\lambda), \lambda) \leq \gamma.$$

Similar to (2.4), since  $\widehat{S}_\rho^D(x, \cdot)$  is  $\rho$ -strongly concave on  $\Lambda$ , the maximization problem in (2.8) has the unique solution

$$(2.11) \quad \lambda_\rho^*(x) \triangleq \arg \max_{\lambda \in \Lambda} \widehat{S}_\rho^D(x, \lambda), \quad \forall x \in \mathbb{E}_1.$$

Additionally, for any  $x \in \mathbb{E}_1$ , we call  $\tilde{\lambda}_{\rho, \eta}(x) \in \Lambda$  an  $\eta$ -inexact solution (where  $\eta \geq 0$ ) if

$$(2.12) \quad \widehat{\psi}_\rho^P(x) - \widehat{S}_\rho^D(x, \tilde{\lambda}_{\rho, \eta}(x)) = \widehat{S}_\rho^D(x, \lambda_\rho^*(x)) - \widehat{S}_\rho^D(x, \tilde{\lambda}_{\rho, \eta}(x)) \leq \eta.$$

Finally, we define the duality gap  $\Delta : \mathbb{E}_1 \times \mathbb{E}_2 \rightarrow \bar{\mathbb{R}}$  associated with  $S(\cdot, \cdot)$  as

$$(2.13) \quad \Delta(x, \lambda) \triangleq \psi^P(x) - \psi^D(\lambda), \quad \forall (x, \lambda) \in \mathbb{E}_1 \times \mathbb{E}_2.$$

Clearly,  $(x^*, \lambda^*) \in \mathcal{X} \times \Lambda$  is a saddle point of (1.1) if and only if  $\Delta(x^*, \lambda^*) = 0$ . As a consequence, for any  $\varepsilon > 0$ , we call  $(\bar{x}, \bar{\lambda}) \in \mathcal{X} \times \Lambda$  an  $\varepsilon$ -saddle point of (1.1) if  $\Delta(\bar{x}, \bar{\lambda}) \leq \varepsilon$ . Additionally, we define the  $(\rho)$ -smoothed duality gap  $\Delta_\rho : \mathbb{E}_1 \times \mathbb{E}_2 \rightarrow \bar{\mathbb{R}}$  as

$$(2.14) \quad \Delta_\rho(x, \lambda) \triangleq \psi_\rho^P(x) - \psi^D(\lambda), \quad \forall (x, \lambda) \in \mathbb{E}_1 \times \mathbb{E}_2.$$

**2.2. Smoothness Properties.** We first show that the optimal solution  $x^*(\cdot)$  in (2.9) is Lipschitz on  $\mathbb{E}_2$  and the function  $\widehat{\psi}^D$  in (2.4) is smooth on  $\mathbb{E}_2$ .

PROPOSITION 2.1. *The function  $\widehat{\psi}^D$  is differentiable on  $\mathbb{E}_2$  and for any  $\lambda \in \mathbb{E}_2$ ,  $\nabla \widehat{\psi}^D(\lambda) = \nabla_\lambda \widehat{S}^P(x^*(\lambda), \lambda) = \nabla_\lambda \Phi(x^*(\lambda), \lambda)$ . In addition,  $x^*(\cdot)$  is  $(2L_{\lambda x}/\mu)$ -Lipschitz on  $\mathbb{E}_2$  and  $\nabla \widehat{\psi}^D$  is  $L_D$ -Lipschitz on  $\mathbb{E}_2$ , where*

$$(2.15) \quad L_D \triangleq L_{\lambda\lambda} + 2L_{\lambda x}^2/\mu.$$

*Proof.* Since for any  $(x, \lambda) \in \mathbb{E}_1 \times \mathbb{E}_2$ ,  $\mathcal{X}$  is compact,  $\widehat{S}^P(\cdot, \lambda)$  is  $\mu$ -s.c. on  $\mathcal{X}$ , and  $\widehat{S}^P(x, \cdot)$  is differentiable on  $\mathbb{E}_2$ , we can invoke Danskin's Theorem [5, Proposition B.25] to conclude that  $\widehat{\psi}^D$  is differentiable on  $\mathbb{E}_2$  and  $\nabla \widehat{\psi}^D(\lambda) = \nabla_\lambda \widehat{S}^P(x^*(\lambda), \lambda)$ . To show Lipschitz continuity of  $x^*(\cdot)$ , note that for any  $\lambda_1, \lambda_2 \in \mathbb{E}_2$ , since  $\widehat{S}^P(\cdot, \lambda_1)$  is  $\mu$ -s.c.,

$$(2.16) \quad \|x^*(\lambda_2) - x^*(\lambda_1)\|^2 \leq \frac{2}{\mu} (\widehat{S}^P(x^*(\lambda_2), \lambda_1) - \widehat{S}^P(x^*(\lambda_1), \lambda_1)).$$

On the other hand,

$$(2.17) \quad \begin{aligned} & \widehat{S}^P(x^*(\lambda_2), \lambda_1) - \widehat{S}^P(x^*(\lambda_1), \lambda_1) \\ & \stackrel{(a)}{\leq} (\widehat{S}^P(x^*(\lambda_2), \lambda_1) - \widehat{S}^P(x^*(\lambda_1), \lambda_1)) - (\widehat{S}^P(x^*(\lambda_2), \lambda_2) - \widehat{S}^P(x^*(\lambda_1), \lambda_2))) \\ & \stackrel{(b)}{=} (\Phi(x^*(\lambda_2), \lambda_1) - \Phi(x^*(\lambda_1), \lambda_1)) - (\Phi(x^*(\lambda_2), \lambda_2) - \Phi(x^*(\lambda_1), \lambda_2))) \\ & = \int_0^1 \langle \nabla_\lambda \Phi(x^*(\lambda_2), \lambda_2 + t(\lambda_1 - \lambda_2)) \\ & \quad - \nabla_\lambda \Phi(x^*(\lambda_1), \lambda_2 + t(\lambda_1 - \lambda_2)), \lambda_1 - \lambda_2 \rangle dt \\ & \stackrel{(c)}{\leq} \|\nabla_\lambda \Phi(x^*(\lambda_2), \lambda_2 + t(\lambda_1 - \lambda_2)) - \nabla_\lambda \Phi(x^*(\lambda_1), \lambda_2 + t(\lambda_1 - \lambda_2))\|_* \|\lambda_1 - \lambda_2\| \\ & \stackrel{(d)}{\leq} L_{\lambda x} \|x^*(\lambda_2) - x^*(\lambda_1)\| \|\lambda_1 - \lambda_2\|, \end{aligned}$$

where in (a) we use the fact that  $\widehat{S}^P(x^*(\lambda_2), \lambda_2) - \widehat{S}^P(x^*(\lambda_1), \lambda_2) \leq 0$  since  $x^*(\lambda_2)$  is the minimizer of  $\widehat{S}^P(\cdot, \lambda_2)$  on  $\mathcal{X}$ , in (b) we use the definition of  $\widehat{S}^P(\cdot, \cdot)$  in (2.9), in (c) we use the definition of the dual norm  $\|\cdot\|_*$  and in (d) we use the Lipschitz continuity of  $\nabla_\lambda \Phi(\cdot, \lambda)$  in (1.3c). Therefore, by combining (2.16) and (2.17), we have

$$(2.18) \quad \|x^*(\lambda_1) - x^*(\lambda_2)\| \leq \frac{2L_{\lambda x}}{\mu} \|\lambda_1 - \lambda_2\|.$$

As a result,

$$\begin{aligned} \|\nabla \widehat{\psi}^D(\lambda_1) - \nabla \widehat{\psi}^D(\lambda_2)\|_* &= \|\nabla_\lambda \Phi(x^*(\lambda_1), \lambda_1) - \nabla_\lambda \Phi(x^*(\lambda_2), \lambda_2)\|_* \\ &\leq \|\nabla_\lambda \Phi(x^*(\lambda_1), \lambda_1) - \nabla_\lambda \Phi(x^*(\lambda_2), \lambda_1)\|_* \\ &\quad + \|\nabla_\lambda \Phi(x^*(\lambda_2), \lambda_1) - \nabla_\lambda \Phi(x^*(\lambda_2), \lambda_2)\|_* \\ &\leq L_{\lambda x} \|x^*(\lambda_1) - x^*(\lambda_2)\| + L_{\lambda\lambda} \|\lambda_1 - \lambda_2\| \\ &\leq (2L_{\lambda x}^2/\mu + L_{\lambda\lambda}) \|\lambda_1 - \lambda_2\|, \end{aligned}$$

where in the last inequality we use (2.18).  $\square$

By a symmetric argument, we can also conclude that  $\widehat{\psi}_\rho^P$  is differentiable on  $\mathbb{E}_1$  and  $\nabla \widehat{\psi}_\rho^P$  is  $(L_{xx} + 2L_{\lambda x}^2/\rho)$ -Lipschitz on  $\mathbb{E}_1$ . For this reason, we will call  $\widehat{\psi}_\rho^P$  the  $(\rho)$ -smoothed primal function and consequently,  $\Delta_\rho$  the  $(\rho)$ -smoothed duality gap.

Based on Proposition 2.1, we can establish the following results involving  $\tilde{x}_\gamma(\lambda)$ , i.e., the  $\gamma$ -inexact solution of (2.2).

PROPOSITION 2.2. *For any  $\gamma \geq 0$  and  $\lambda, \lambda' \in \mathbb{E}_2$ , we have*

$$(2.19) \quad \|\nabla_\lambda \widehat{S}^P(\tilde{x}_\gamma(\lambda), \lambda) - \nabla \widehat{\psi}^D(\lambda)\|_* \leq L_{\lambda x} \sqrt{2\gamma/\mu},$$

$$(2.20) \quad 0 \leq \widehat{S}^P(\tilde{x}_\gamma(\lambda), \lambda) - \widehat{\psi}^D(\lambda') + \langle \nabla_\lambda \widehat{S}^P(\tilde{x}_\gamma(\lambda), \lambda), \lambda' - \lambda \rangle \leq L_D \|\lambda - \lambda'\|^2 + 2\gamma.$$

*Proof.* First, for any  $\lambda \in \mathbb{E}_2$ , we note that  $\nabla_\lambda \widehat{S}^P(\tilde{x}_\gamma(\lambda), \lambda) = \nabla_\lambda \Phi(\tilde{x}_\gamma(\lambda), \lambda)$  and  $\nabla \widehat{\psi}^D(\lambda) = \nabla_\lambda \Phi(x^*(\lambda), \lambda)$  (cf. Proposition 2.1). Therefore,

$$(2.21) \quad \begin{aligned} \|\nabla_\lambda \Phi(\tilde{x}_\gamma(\lambda), \lambda) - \nabla \widehat{\psi}^D(\lambda)\|_*^2 &= \|\nabla_\lambda \Phi(\tilde{x}_\gamma(\lambda), \lambda) - \nabla_\lambda \Phi(x^*(\lambda), \lambda)\|_*^2 \\ &\stackrel{(a)}{\leq} L_{\lambda x}^2 \|\tilde{x}_\gamma(\lambda) - x^*(\lambda)\|^2 \\ &\stackrel{(b)}{\leq} (2L_{\lambda x}^2/\mu) (\widehat{S}^P(\tilde{x}_\gamma(\lambda), \lambda) - \widehat{S}^P(x^*(\lambda), \lambda)) \\ &\stackrel{(c)}{\leq} 2L_{\lambda x}^2 \gamma / \mu, \end{aligned}$$

where in (a) we use (1.3c), in (b) we use the  $\mu$ -strong convexity of  $\widehat{S}^P(\cdot, \lambda)$  on  $\mathcal{X}$  and in (c) we use the definition of  $\tilde{x}_\gamma(\lambda)$  in (2.10). This proves (2.19).

We next prove (2.20). First, for any  $\lambda, \lambda' \in \mathbb{E}_2$ ,

$$(2.22) \quad \widehat{\psi}^D(\lambda') \stackrel{(a)}{\leq} \widehat{S}^P(\tilde{x}_\gamma(\lambda), \lambda') \stackrel{(b)}{\leq} \widehat{S}^P(\tilde{x}_\gamma(\lambda), \lambda) + \langle \nabla_\lambda \widehat{S}^P(\tilde{x}_\gamma(\lambda), \lambda), \lambda' - \lambda \rangle,$$

where (a) follows from (2.4) and (b) follows from the concavity of  $\widehat{S}^P(\tilde{x}_\gamma(\lambda), \cdot)$  on  $\mathbb{E}_2$ . This proves the left-hand side (LHS) of (2.20). To show the right-hand side (RHS), we note that  $\widehat{\psi}^D$  is concave and  $L_D$ -smooth on  $\mathbb{E}_2$  (cf. Proposition 2.1). Thus we can invoke the descent lemma [5], such that for all  $\lambda, \lambda' \in \mathbb{E}_2$ ,

$$(2.23) \quad \begin{aligned} \widehat{\psi}^D(\lambda') &\geq \widehat{\psi}^D(\lambda) + \langle \nabla \widehat{\psi}^D(\lambda), \lambda' - \lambda \rangle - (L_D/2) \|\lambda - \lambda'\|^2 \\ &\geq \widehat{S}^P(\tilde{x}_\gamma(\lambda), \lambda) - \gamma + \langle \nabla_\lambda \widehat{S}^P(\tilde{x}_\gamma(\lambda), \lambda), \lambda' - \lambda \rangle \\ &\quad + \langle \nabla \widehat{\psi}^D(\lambda) - \nabla_\lambda \widehat{S}^P(\tilde{x}_\gamma(\lambda), \lambda), \lambda' - \lambda \rangle - (L_D/2) \|\lambda - \lambda'\|^2 \\ &\geq \widehat{S}^P(\tilde{x}_\gamma(\lambda), \lambda) - \gamma + \langle \nabla_\lambda \widehat{S}^P(\tilde{x}_\gamma(\lambda), \lambda), \lambda' - \lambda \rangle \\ &\quad - \|\nabla \widehat{\psi}^D(\lambda) - \nabla_\lambda \widehat{S}^P(\tilde{x}_\gamma(\lambda), \lambda)\|_* \|\lambda - \lambda'\| - (L_D/2) \|\lambda - \lambda'\|^2 \\ &\stackrel{(a)}{\geq} \widehat{S}^P(\tilde{x}_\gamma(\lambda), \lambda) - \gamma + \langle \nabla_\lambda \widehat{S}^P(\tilde{x}_\gamma(\lambda), \lambda), \lambda' - \lambda \rangle \\ &\quad - L_{\lambda x} \sqrt{2\gamma/\mu} \|\lambda - \lambda'\| - (L_D/2) \|\lambda - \lambda'\|^2, \\ &\stackrel{(b)}{\geq} \widehat{S}^P(\tilde{x}_\gamma(\lambda), \lambda) - 2\gamma + \langle \nabla_\lambda \widehat{S}^P(\tilde{x}_\gamma(\lambda), \lambda), \lambda' - \lambda \rangle - L_D \|\lambda - \lambda'\|^2, \end{aligned}$$

where (a) follows from (2.19) and (b) follows from the AM-GM inequality, i.e.,

$$(2.24) \quad L_{\lambda x} \sqrt{2\gamma/\mu} \|\lambda - \lambda'\| \leq (L_{\lambda x}^2/\mu) \|\lambda - \lambda'\|^2 + \gamma \leq (L_D/2) \|\lambda - \lambda'\|^2 + \gamma.$$

We then rearrange (2.23) to obtain the RHS of (2.20).  $\square$

**3. Deterministic IPDS Framework.** We develop the IPDS framework based on the idea of *smoothed duality gap reduction*. First, we make an important assumption, followed by a few remarks.

ASSUMPTION 3.1. *The set  $\Lambda$  is bounded.*

**Algorithm 1** Deterministic inexact primal-dual smoothing framework

**Input:** Initial smoothing parameter  $\rho_0 > 0$ , nonnegative error sequences  $\{\eta_k\}_{k \in \mathbb{Z}_+}$  and  $\{\gamma_k\}_{k \in \mathbb{Z}_+}$ , interpolation sequence  $\{\tau_k\}_{k \in \mathbb{Z}_+} \subseteq (0, 1)$  and deterministic first-order algorithms  $\mathbf{N}_1$  and  $\mathbf{N}_2$ .

**Initialize:**  $x^0 \in \mathcal{X}$ ,  $\lambda^0 \in \Lambda$  and  $k = 0$ .

**Repeat** (until some convergence criterion is met)

1. Use  $\mathbf{N}_1$  to find  $\tilde{\lambda}_{\rho_k, \eta_k}(x^k) \in \Lambda$  such that

$$(3.3) \quad \widehat{\psi}_{\rho_k}^{\mathbf{P}}(x^k) - \widehat{S}_{\rho_k}^{\mathbf{D}}(x^k, \tilde{\lambda}_{\rho_k, \eta_k}(x^k)) \leq \eta_k.$$

2. Set  $\hat{\lambda}^k = \tau_k \lambda^k + (1 - \tau_k) \tilde{\lambda}_{\rho_k, \eta_k}(x^k)$ .
3. Use  $\mathbf{N}_2$  to find  $\tilde{x}_{\gamma_k}(\hat{\lambda}^k) \in X$  such that

$$(3.4) \quad \widehat{S}^{\mathbf{P}}(\tilde{x}_{\gamma_k}(\hat{\lambda}^k), \hat{\lambda}^k) - \widehat{\psi}^{\mathbf{D}}(\hat{\lambda}^k) \leq \gamma_k.$$

4. Set  $x^{k+1} = \tau_k x^k + (1 - \tau_k) \tilde{x}_{\gamma_k}(\hat{\lambda}^k)$ .
5. Set  $\rho_{k+1} = \tau_k \rho_k$ .
6. Use  $\mathbf{N}_1$  to find  $\tilde{\lambda}_{\rho_{k+1}, \eta_k}(x^{k+1}) \in \Lambda$  such that

$$(3.5) \quad \widehat{\psi}_{\rho_{k+1}}^{\mathbf{P}}(x^{k+1}) - \widehat{S}_{\rho_{k+1}}^{\mathbf{D}}(x^{k+1}, \tilde{\lambda}_{\rho_{k+1}, \eta_k}(x^{k+1})) \leq \eta_k.$$

7. Set  $\lambda^{k+1} = \tau_k \lambda^k + (1 - \tau_k) \tilde{\lambda}_{\rho_{k+1}, \eta_k}(x^{k+1})$ .
8. Set  $k = k + 1$ .

**Output:**  $(x^{\text{out}}, \lambda^{\text{out}}) \triangleq (x^k, \lambda^k)$ .

To see the implication of this assumption, for any  $(x, \lambda) \in \mathcal{X} \times \Lambda$ , we may bound

$$(3.1) \quad \begin{aligned} |\Delta(x, \lambda) - \Delta_\rho(x, \lambda)| &= |\sup_{\lambda \in \Lambda} S(x, \lambda) - \sup_{\lambda \in \Lambda} S_\rho(x, \lambda)| \\ &\leq \sup_{\lambda \in \Lambda} |S(x, \lambda) - S_\rho(x, \lambda)| \\ &= \rho \sup_{\lambda \in \Lambda} |\omega(\lambda)| \\ &< +\infty, \end{aligned}$$

where the last inequality follows from Assumption 3.1, the closedness of  $\Lambda$  and the continuity of  $\omega$  on  $\Lambda$ . For convenience, define

$$(3.2) \quad B_{\omega, \Lambda} \triangleq \sup_{\lambda \in \Lambda} |\omega(\lambda)|.$$

Since  $\omega$  is also 1-s.c. on  $\Lambda$ , we conclude that  $B_{\omega, \Lambda} < +\infty$  if and only if  $\Lambda$  is bounded.

*Remark 3.2.* We provide a few remarks about Assumption 3.1 (which is equivalent to  $B_{\omega, \Lambda} < +\infty$ ). First, it is important in connecting the smoothed duality gap  $\Delta_\rho$  to the duality gap  $\Delta$ . Indeed, in our analysis, we will first analyze the convergence rate of the smoothed duality gap, and then show that the same rate holds for the (original) duality gap if  $B_{\omega, \Lambda} < +\infty$ . Note that we *do not* need this assumption to derive any convergence results regarding the smoothed duality gap. Finally, we note that Assumption 3.1 also appears in many other algorithms for solving SPPs, e.g., Mirror-Prox [20], HPE-type [14] and PDHG [32], although for different reasons. Indeed, it is typical in the works where the duality gap is used as the convergence criterion, and is not specific to our work.

**3.1. Framework Description.** The framework is presented in Algorithm 1. We now describe the main ideas behind this framework. From (3.1), we observe that  $\Delta(x, \lambda) \leq \Delta_\rho(x, \lambda) + \rho B_{\omega, \Lambda}$ . Therefore, if there exists a primal-dual pair  $(x, \lambda) \in \mathcal{X} \times \Lambda$  such that the smoothed duality gap  $\Delta_\rho(x, \lambda)$  is small, then with a small smoothing parameter  $\rho$ , the duality gap  $\Delta(x, \lambda)$  will also be small. This leads us to develop a framework that “sufficiently” reduces both the smoothed duality gap and smoothing parameter in each iteration. Indeed, in step 5 of Algorithm 1, the smoothing parameter  $\rho_k$  is reduced by a factor of  $(1 - \tau_k) \in (0, 1)$ . The same factor is also used to interpolate the primal and dual iterates (cf. steps 2, 4 and 7). As a key difference with the smoothing frameworks for bilinear SPPs (e.g., [22, 21]), our framework does not require the sub-problems (2.4) and (2.8) to be solved exactly. Instead, we only need the inexact solutions satisfying certain accuracy criteria (involving parameters  $\gamma_k$  and  $\eta_k$ ; cf. (3.3), (3.4) and (3.5)). In principle, such solutions can be computed via any first-order method. (For the implementation details, we refer readers to Section 3.2.) From the description above, we see that the success of our framework hinges upon the proper choices of  $\tau_k$ ,  $\gamma_k$  and  $\eta_k$ , which ensure the reduction of the smoothed duality gap  $\Delta_{\rho_k}(x^k, \lambda^k)$  in each iteration, and simultaneously decrease  $\rho_k$  (cf. Section 3.3).

**3.2. Solving sub-problems.** In Algorithm 1, it is important for us to solve the sub-problems in steps 1, 3 and 6 inexactly in an efficient manner. Since both of the optimization problems in (2.4) and (2.8) have composite forms, it is natural for us to employ the optimal first-order algorithm in [24] to solve them. More specifically, we choose  $\mathbf{N}_1$  and  $\mathbf{N}_2$  to be the accelerated proximal gradient method (denoted as APG) in [24, Equation (4.9)].

We now briefly review the convergence rate of this method. Let  $\mathcal{U}$  be a nonempty, convex and compact set (in  $\mathbb{E}_1$  or  $\mathbb{E}_2$ ). In addition, let  $\phi_1$  and  $\phi_2$  be CCP functions such that  $\phi_1$  is  $L'$ -smooth on  $\mathcal{U}$ , and  $\phi_2$  is  $\mu'$ -s.c. on  $\mathcal{U}$  (where  $\mu' > 0$ ) and admits an easily computable Bregman projection on  $\mathcal{U}$  (with DGF  $\pi$ ; cf. Section 2.1). Consider

$$(3.6) \quad \min_{u \in \mathcal{U}} \{ \Psi(u) \triangleq \phi_1(u) + \phi_2(u) \},$$

whose (unique) solution is denoted by  $u^* \in \mathcal{U}$ . Define  $\kappa' \triangleq L'/\mu'$ . From [24, Theorem 6], if we use APG to solve (3.6), then for any starting point  $u^0 \in \mathcal{U}^o \triangleq \mathcal{U} \cap \text{int dom } \pi$ ,

$$(3.7) \quad \begin{aligned} \Psi(u^N) - \Psi(u^*) &\leq (L'/4)(1 + 1/\sqrt{2\kappa})^{-2(N-1)} \|u^0 - u^*\|^2 \\ &\leq (L'/4)(1 + 1/\sqrt{2\kappa})^{-2(N-1)} D_{\mathcal{U}}^2, \quad \forall N \in \mathbb{N}, \end{aligned}$$

where  $u^N$  denotes the  $N$ -th iterate of APG (where  $N \in \mathbb{N}$ ) and

$$(3.8) \quad D_{\mathcal{U}} \triangleq \max_{u, u' \in \mathcal{U}} \|u - u'\| < +\infty$$

denotes the diameter of the constraint set  $\mathcal{U}$ . In other words, to obtain an  $\varepsilon$ -optimality gap, the number of iterations of APG that we need is

$$(3.9) \quad \left\lceil \sqrt{\frac{\kappa'}{2}} \log \left( \frac{L' D_{\mathcal{U}}^2}{4\varepsilon} \right) \right\rceil + 1 = O \left( \sqrt{\kappa'} \log \left( \frac{L'}{\varepsilon} \right) \right).$$

Note that (3.9) also holds for when  $\phi_1$ , instead of  $\phi_2$ , is  $\mu$ -s.c. on  $\mathcal{U}$ ; see [24, Section 5.1].

Therefore, based on (3.9), to find an  $\eta$ -inexact solution of (2.8) (cf. (2.12)), the number of dual oracle calls (cf. Section 1.1) made by  $\mathbf{N}_1$  is

$$(3.10) \quad C_{\mathbf{N}_1} \triangleq n \left\{ \left\lceil \sqrt{\frac{L_{\lambda\lambda}}{2\rho}} \log \left( \frac{L_{\lambda\lambda} D_{\Lambda}^2}{4\eta} \right) \right\rceil + 1 \right\} = O \left( n \sqrt{\frac{L_{\lambda\lambda}}{\rho}} \log \left( \frac{L_{\lambda\lambda}}{\eta} \right) \right),$$

and to find a  $\gamma$ -inexact solution of (2.4) (cf. (2.10)), the number of primal oracle calls made by  $\mathbf{N}_2$  is

$$(3.11) \quad C_{\mathbf{N}_2} \triangleq (n+1) \left\{ \left\lceil \sqrt{\frac{L+L_{xx}}{2\mu}} \log \left( \frac{(L+L_{xx})D_{\mathcal{X}}^2}{4\gamma} \right) \right\rceil + 1 \right\} \\ = O(n\sqrt{\kappa_{\mathcal{X}} \log((L+L_{xx})/\gamma)}).$$

Note that in (3.10) and (3.11),  $D_{\Lambda}$  and  $D_{\mathcal{X}}$  denote the diameters of the sets  $\Lambda$  and  $\mathcal{X}$ , respectively (defined similarly to  $D_{\mathcal{U}}$  as in (3.8)), and are both *finite* (cf. Assumption 3.1). Therefore, we do *not* need to know the optimal solution or the optimal objective value of the problem in (2.8) (resp. (2.4)) in order to find  $\tilde{\lambda}_{\rho,\eta}(x)$  (resp.  $\tilde{x}_{\gamma}(\lambda)$ ).

**3.3. Convergence Analysis.** As the first step of our analysis, we prove that in each iteration  $k$ , if the smoothing parameter  $\rho_k$  is chosen to be sufficiently large, then the smoothed duality gap is reduced.

LEMMA 3.3. *In Algorithm 1, for any  $k \in \mathbb{Z}_+ \triangleq \mathbb{N} \cup \{0\}$ , if  $\rho_{k+1} \geq 4(1-\tau_k)^2 L_D$ , then*

$$(3.12) \quad \Delta_{\rho_{k+1}}(x^{k+1}, \lambda^{k+1}) \leq \tau_k \Delta_{\rho_k}(x^k, \lambda^k) + 2\gamma_k + 2\eta_k.$$

*Proof.* Fix any  $k \in \mathbb{Z}_+$ . Since  $\widehat{S}_{\rho_k}^{\mathbf{D}}(x^k, \cdot)$  is  $\rho_k$ -strongly concave on  $\Lambda$ ,

$$(3.13) \quad \widehat{\psi}_{\rho_k}^{\mathbf{P}}(x^k) - \widehat{S}_{\rho_k}^{\mathbf{D}}(x^k, \lambda) = \widehat{S}_{\rho_k}^{\mathbf{D}}(x^k, \lambda_{\rho_k}^*(x^k)) - \widehat{S}_{\rho_k}^{\mathbf{D}}(x^k, \lambda) \\ \geq \frac{\rho_k}{2} \|\lambda - \lambda_{\rho_k}^*(x^k)\|^2, \quad \forall \lambda \in \Lambda.$$

As a result, for all  $\lambda \in \Lambda$ , we have

$$(3.14) \quad \Delta_{\rho_k}(x^k, \lambda^k) \\ = \psi_{\rho_k}^{\mathbf{P}}(x^k) - \psi^{\mathbf{D}}(\lambda^k) \\ = f(x^k) + g(x^k) + \widehat{\psi}_{\rho_k}^{\mathbf{P}}(x^k) - \psi^{\mathbf{D}}(\lambda^k) \\ \stackrel{(a)}{\geq} f(x^k) + g(x^k) + \widehat{S}_{\rho_k}^{\mathbf{D}}(x^k, \lambda) + \frac{\rho_k}{2} \|\lambda - \lambda_{\rho_k}^*(x^k)\|^2 - \psi^{\mathbf{D}}(\lambda^k) \\ = S(x^k, \lambda) - \rho_k \omega(\lambda) + \frac{\rho_k}{2} \|\lambda - \lambda_{\rho_k}^*(x^k)\|^2 - (\widehat{\psi}^{\mathbf{D}}(\lambda^k) - h(\lambda^k)) \\ = \widehat{S}^{\mathbf{P}}(x^k, \lambda) - h(\lambda) - \rho_k \omega(\lambda) + \frac{\rho_k}{2} \|\lambda - \lambda_{\rho_k}^*(x^k)\|^2 - \widehat{\psi}^{\mathbf{D}}(\lambda^k) + h(\lambda^k),$$

where in (a) we use (3.13). Define  $z_k(\lambda) \triangleq \tau_k \lambda^k + (1-\tau_k)\lambda$ . We then multiply both sides of (3.14) by  $\tau_k > 0$ , and obtain

$$(3.15) \quad \tau_k \Delta_{\rho_k}(x^k, \lambda^k) \\ \stackrel{(a)}{\geq} \tau_k \widehat{S}^{\mathbf{P}}(x^k, \lambda) - \tau_k h(\lambda) - \rho_{k+1} \omega(\lambda) \\ + \frac{\rho_{k+1}}{2} \|\lambda - \lambda_{\rho_k}^*(x^k)\|^2 - \tau_k \widehat{\psi}^{\mathbf{D}}(\lambda^k) + \tau_k h(\lambda^k) \\ = \tau_k \widehat{S}^{\mathbf{P}}(x^k, \lambda) + (1-\tau_k) \widehat{S}^{\mathbf{P}}(\tilde{x}_{\gamma_k}(\hat{\lambda}^k), \lambda) - \tau_k h(\lambda) - \rho_{k+1} \omega(\lambda) \\ + \frac{\rho_{k+1}}{2} \|\lambda - \lambda_{\rho_k}^*(x^k)\|^2 - \tau_k \widehat{\psi}^{\mathbf{D}}(\lambda^k) - (1-\tau_k) \widehat{S}^{\mathbf{P}}(\tilde{x}_{\gamma_k}(\hat{\lambda}^k), \lambda) + \tau_k h(\lambda^k) \\ \stackrel{(b)}{\geq} \widehat{S}^{\mathbf{P}}(x^{k+1}, \lambda) - \tau_k h(\lambda) - \rho_{k+1} \omega(\lambda) + \frac{\rho_{k+1}}{2} \|\lambda - \lambda_{\rho_k}^*(x^k)\|^2$$

$$\begin{aligned}
 & -\tau_k(\widehat{S}^{\text{P}}(\tilde{x}_{\gamma_k}(\hat{\lambda}^k), \hat{\lambda}^k) + \langle \nabla_{\lambda} \widehat{S}^{\text{P}}(\tilde{x}_{\gamma_k}(\hat{\lambda}^k), \hat{\lambda}^k), \lambda^k - \hat{\lambda}^k \rangle) \\
 & - (1 - \tau_k)(\widehat{S}^{\text{P}}(\tilde{x}_{\gamma_k}(\hat{\lambda}^k), \hat{\lambda}^k) + \langle \nabla_{\lambda} \widehat{S}^{\text{P}}(\tilde{x}_{\gamma_k}(\hat{\lambda}^k), \hat{\lambda}^k), \lambda - \hat{\lambda}^k \rangle) + \tau_k h(\lambda^k) \\
 \stackrel{\text{(c)}}{=} & S_{\rho_{k+1}}(x^{k+1}, \lambda) + (1 - \tau_k)h(\lambda) - \widehat{S}^{\text{P}}(\tilde{x}_{\gamma_k}(\hat{\lambda}^k), \hat{\lambda}^k) \\
 & + \frac{\rho_{k+1}}{2} \|\lambda - \lambda_{\rho_k}^*(x^k)\|^2 - \langle \nabla_{\lambda} \widehat{S}^{\text{P}}(\tilde{x}_{\gamma_k}(\hat{\lambda}^k), \hat{\lambda}^k), z_k(\lambda) - \hat{\lambda}^k \rangle + \tau_k h(\lambda^k) \\
 \stackrel{\text{(d)}}{\geq} & S_{\rho_{k+1}}(x^{k+1}, \lambda) + \frac{\rho_{k+1}}{2} \|\lambda - \lambda_{\rho_k}^*(x^k)\|^2 \\
 & - (\widehat{\psi}^{\text{D}}(z_k(\lambda)) + L_{\text{D}} \|\hat{\lambda}^k - z_k(\lambda)\|^2 + 2\gamma_k) + h(z_k(\lambda)),
 \end{aligned}$$

where in (a) we use  $\tau_k \rho_k = \rho_{k+1}$ , in (b) we use the convexity of  $\widehat{S}^{\text{P}}(\cdot, \lambda)$ , the LHS of (2.20) and the concavity of  $\widehat{S}^{\text{P}}(\tilde{x}_{\gamma_k}(\hat{\lambda}^k), \cdot)$ , in (c) we use the definition of  $\widehat{S}^{\text{P}}$  and  $S_{\rho_{k+1}}$  in (2.4) and (2.6), respectively) and in (d) we use the RHS of (2.20) and the convexity of  $h$ . Note that if we take  $\lambda = \tilde{\lambda}_{\rho_{k+1}, \eta_k}(x^{k+1})$ , then  $z_k(\lambda) = \lambda^{k+1}$  by step 7. In addition, from steps 2 and 7, we have

$$(3.16) \quad \hat{\lambda}^k - \lambda^{k+1} = (1 - \tau_k)(\tilde{\lambda}_{\rho_k, \eta_k}(x^k) - \tilde{\lambda}_{\rho_{k+1}, \eta_k}(x^{k+1})),$$

and from (3.3) and (3.13), we have

$$(3.17) \quad \frac{\rho_k}{2} \|\tilde{\lambda}_{\rho_k, \eta_k}(x^k) - \lambda_{\rho_k}^*(x^k)\|^2 \leq \widehat{\psi}_{\rho_k}^{\text{P}}(x^k) - \widehat{S}_{\rho_k}^{\text{D}}(x^k, \tilde{\lambda}_{\rho_k, \eta_k}(x^k)) \leq \eta_k.$$

This observation leads us to bound  $\|\tilde{\lambda}_{\rho_{k+1}, \eta_k}(x^{k+1}) - \lambda_{\rho_k}^*(x^k)\|^2$  as

$$\begin{aligned}
 (3.18) \quad & \|\tilde{\lambda}_{\rho_{k+1}, \eta_k}(x^{k+1}) - \lambda_{\rho_k}^*(x^k)\|^2 \\
 & \stackrel{\text{(a)}}{\geq} \frac{1}{2} \|\tilde{\lambda}_{\rho_{k+1}, \eta_k}(x^{k+1}) - \tilde{\lambda}_{\rho_k, \eta_k}(x^k)\|^2 - \|\tilde{\lambda}_{\rho_k, \eta_k}(x^k) - \lambda_{\rho_k}^*(x^k)\|^2 \\
 & \stackrel{\text{(b)}}{\geq} \frac{1}{2} \|\tilde{\lambda}_{\rho_{k+1}, \eta_k}(x^{k+1}) - \tilde{\lambda}_{\rho_k, \eta_k}(x^k)\|^2 - \frac{2\eta_k}{\rho_k} \\
 & \stackrel{\text{(c)}}{=} \frac{\|\lambda^{k+1} - \hat{\lambda}^k\|^2}{2(1 - \tau_k)^2} - \frac{2\eta_k}{\rho_k},
 \end{aligned}$$

where in (a) we use  $\|a + b\|^2 \leq 2(\|a\|^2 + \|b\|^2)$ , in (b) we use (3.17) and in (c) we use (3.16). We then substitute  $\lambda = \tilde{\lambda}_{\rho_{k+1}, \eta_k}(x^{k+1})$  and (3.18) into (3.15), and obtain

$$\begin{aligned}
 \tau_k \Delta_{\rho_k}(x^k, \lambda^k) & \geq S_{\rho_{k+1}}(x^{k+1}, \tilde{\lambda}_{\rho_{k+1}, \eta_k}(x^{k+1})) + \frac{\rho_{k+1}}{2} \left( \frac{\|\lambda^{k+1} - \hat{\lambda}^k\|^2}{2(1 - \tau_k)^2} - \frac{2\eta_k}{\rho_k} \right) \\
 & - (\widehat{\psi}^{\text{D}}(\lambda^{k+1}) + L_{\text{D}} \|\hat{\lambda}^k - \lambda^{k+1}\|^2 + 2\gamma_k) \\
 & \stackrel{\text{(a)}}{\geq} \widehat{\psi}_{\rho_{k+1}}^{\text{P}}(x^{k+1}) - (1 + \tau_k)\eta_k + \left( \frac{\rho_{k+1}}{4(1 - \tau_k)^2} - L_{\text{D}} \right) \|\lambda^{k+1} - \hat{\lambda}^k\|^2 \\
 & - \widehat{\psi}^{\text{D}}(\lambda^{k+1}) - 2\gamma_k \\
 & \stackrel{\text{(b)}}{\geq} \Delta_{\rho_{k+1}}(x^{k+1}, \lambda^{k+1}) - 2\eta_k - 2\gamma_k,
 \end{aligned}$$

where in (a) we use (3.5) and in (b) we use  $\tau_k \in (0, 1)$  and  $\rho_{k+1} \geq 4(1 - \tau_k)^2 L_{\text{D}}$ . We hence complete the proof.  $\square$

In Lemma 3.3, we notice that if  $\Delta_{\rho_k}(x^k, \lambda^k) > 2(\gamma_k + \eta_k)/(1 - \tau_k)$ , then the smoothed duality gap will be reduced, i.e.,  $\Delta_{\rho_{k+1}}(x^{k+1}, \lambda^{k+1}) < \Delta_{\rho_k}(x^k, \lambda^k)$ . Indeed, from our choices of  $\tau_k$ ,  $\gamma_k$  and  $\eta_k$  in Theorem 3.5 (see below), the reduction holds as

long as  $\Delta_{\rho_k}(x^k, \lambda^k) > \varepsilon/2$ . This corroborates our description in Section 3.1.

Before proving our main convergence results, let us state a result about linear recursion, whose proof simply follows from induction.

LEMMA 3.4. *Let  $\{\alpha_k\}_{k \in \mathbb{Z}_+}$ ,  $\{\beta_k\}_{k \in \mathbb{Z}_+}$  and  $\{a_k\}_{k \in \mathbb{Z}_+}$  be real sequences. If for all  $k \in \mathbb{Z}_+$ ,  $\alpha_k \geq 0$  and*

$$(3.19) \quad a_{k+1} \leq \alpha_k a_k + \beta_k,$$

then for all  $K \in \mathbb{N}$ ,

$$(3.20) \quad a_K \leq \left( \prod_{k=0}^{K-1} \alpha_k \right) a_0 + \sum_{k=1}^K \left( \prod_{j=k}^{K-1} \alpha_j \right) \beta_{k-1},$$

where we define the empty product  $\prod_{j=K}^{K-1} \alpha_j \triangleq 1$ .

Based on Lemmas 3.3 and 3.4, our main results follow immediately.

THEOREM 3.5. *In Algorithm 1, if we choose  $\rho_0 = 8L_D$  and for any  $k \in \mathbb{Z}_+$ ,*

$$(3.21) \quad \tau_k = \frac{k+1}{k+3}, \quad \gamma_k = \frac{\varepsilon}{4(k+3)} \quad \text{and} \quad \eta_k = \frac{\varepsilon}{4(k+3)},$$

then for any starting point  $(x^0, \lambda^0) \in \mathcal{X} \times \Lambda$  and  $K \in \mathbb{N}$ ,

$$(3.22) \quad \Delta_{\rho_K}(x^K, \lambda^K) \leq B'_\Delta(K, \varepsilon) \triangleq \frac{2\Delta_{\rho_0}(x^0, \lambda^0)}{(K+1)(K+2)} + \frac{\varepsilon}{2}.$$

Furthermore, if Assumption 3.1 holds, then

$$(3.23) \quad \Delta(x^K, \lambda^K) \leq B_\Delta(K, \varepsilon) \triangleq \frac{32L_D B_{\omega, \Lambda} + 2\Delta(x^0, \lambda^0)}{(K+1)(K+2)} + \frac{\varepsilon}{2}.$$

*Proof.* Based on the choice of  $\rho_0$  and  $\{\tau_k\}_{k \in \mathbb{Z}_+}$ , for any  $K \in \mathbb{N}$ , we have

$$(3.24) \quad \prod_{k=0}^{K-1} \tau_k = \frac{2}{(K+1)(K+2)} \implies \rho_K = \rho_0 \prod_{k=0}^{K-1} \tau_k = \frac{16L_D}{(K+1)(K+2)}.$$

Therefore, we can easily verify that the condition  $\rho_K \geq 4(1 - \tau_{K-1})^2 L_D$  in Lemma 3.3 is satisfied. Consequently, by the recursion in (3.12) and Lemma 3.4, we have

$$(3.25) \quad \begin{aligned} \Delta_{\rho_K}(x^K, \lambda^K) &\leq \Delta_{\rho_0}(x^0, \lambda^0) \prod_{k=0}^{K-1} \tau_k + \sum_{k=1}^K 2(\gamma_{k-1} + \eta_{k-1}) \prod_{j=k}^{K-1} \tau_j \\ &= \frac{2\Delta_{\rho_0}(x^0, \lambda^0)}{(K+1)(K+2)} + \sum_{k=1}^K \frac{\varepsilon}{k+2} \cdot \frac{(k+1)(k+2)}{(K+1)(K+2)} \\ &= \frac{2\Delta_{\rho_0}(x^0, \lambda^0)}{(K+1)(K+2)} + \frac{\varepsilon}{2} \frac{K(K+3)}{(K+1)(K+2)}. \end{aligned}$$

We then obtain (3.22) by noticing that  $K(K+3) \leq (K+1)(K+2)$ . Based on (3.22), to derive (3.23), we simply use (3.1) and (3.24).  $\square$

*Remark 3.6.* From (3.23), since  $L_D$  only depends on  $L_{\lambda x}$  and  $L_{\lambda \lambda}$  (cf. (2.15)), we note that the convergence of the duality gap in Algorithm 1 only requires the Lipschitz continuity of  $\nabla_{\lambda} \Phi(\cdot, \lambda)$  and  $\nabla_{\lambda} \Phi(x, \cdot)$  (cf. (1.3c) and (1.3d)), but not the Lipschitz continuity of  $\nabla f$ ,  $\nabla_x \Phi(\cdot, \lambda)$  and  $\nabla_x \Phi(x, \cdot)$  (cf. (1.2), (1.3a) and (1.3b)). However, the latter smoothness conditions are needed in order to use the optimal first-order algorithm, as introduced in Section 3.2, to solve the sub-problems in (3.4). By doing so, in Algorithm 1, we can achieve the overall primal oracle complexity  $\tilde{O}(1/\sqrt{\varepsilon})$ . For details, we refer readers to Section 3.4.

Note that Theorem 3.5 indicates that in Algorithm 1, to achieve an  $\varepsilon$ -duality gap, the number of iterations we need is

$$(3.26) \quad K_{\det} \triangleq \left\lceil \frac{2\sqrt{16L_D B_{\omega, \Lambda} + \Delta(x^0, \lambda^0)}}{\sqrt{\varepsilon}} \right\rceil + 1 = O\left(\sqrt{\frac{L_D}{\varepsilon}}\right).$$

Furthermore, by the definition of  $L_D$  in (2.15), we have  $K_{\det} = O(\sqrt{L_{\lambda\lambda}/\varepsilon} + L_{\lambda x}/\sqrt{\mu\varepsilon})$ .

**3.4. Oracle Complexity.** Based on the results in Sections 3.2 and 3.3 (specifically, (3.10), (3.11) and (3.26)), we may analyze the primal and dual oracle complexities needed in Algorithm 1 to achieve an  $\varepsilon$ -duality gap (i.e.,  $\Delta(x^{\text{out}}, \lambda^{\text{out}}) \leq \varepsilon$ ).

**THEOREM 3.7.** *Let Assumption 3.1 hold. In Algorithm 1, for any starting point  $(x^0, \lambda^0) \in \mathcal{X} \times \Lambda$ , let  $C_{\det}^P$  and  $C_{\det}^D$  denote the primal and dual oracle complexities (cf. Section 1.1) to achieve an  $\varepsilon$ -duality gap, respectively. Then we have*

$$(3.27) \quad C_{\det}^P = O\left(n\sqrt{\kappa_{\mathcal{X}}L_D/\varepsilon} \log((L + L_{xx})L_D/\varepsilon)\right),$$

$$(3.28) \quad C_{\det}^D = O\left(n(\sqrt{L_{\lambda\lambda}L_D/\varepsilon}) \log(L_{\lambda\lambda}L_D/\varepsilon)\right).$$

*Proof.* Since  $\gamma_k = \varepsilon/(4(k+3)) = O(\varepsilon/k)$  (cf. (3.21)), based on (3.11), we have

$$(3.29) \quad \begin{aligned} C_{\det}^P &= \sum_{k=1}^{K_{\det}} O\left(n\sqrt{\kappa_{\mathcal{X}}} \log((L + L_{xx})k/\varepsilon)\right) \\ &= O\left(n\sqrt{\kappa_{\mathcal{X}}}\left(K_{\det} \log(\kappa_{\mathcal{X}}) + \log(K_{\det}!)\right)\right) \\ &\stackrel{(a)}{=} O\left(n\sqrt{\kappa_{\mathcal{X}}}\sqrt{L_D/\varepsilon}\left(\log((L + L_{xx})/\varepsilon) + \log(L_D/\varepsilon)\right)\right) \\ &= O\left(n\sqrt{\kappa_{\mathcal{X}}L_D/\varepsilon} \log((L + L_{xx})L_D/\varepsilon)\right), \end{aligned}$$

where in (a) we use the fact that  $\log(K!) = \Theta(K \log K)$  for any  $K \in \mathbb{N}$  and (3.26).

Similarly, we can analyze the dual oracle complexity for solving (3.3). Since  $\rho_k = O(L_D/k^2)$  (cf. (3.24)) and  $\eta_k = O(\varepsilon/k)$  (cf. (3.21)), based on (3.10), we have

$$(3.30) \quad \begin{aligned} C_{\det,1}^D &= \sum_{k=1}^{K_{\det}} O\left(n\sqrt{L_{\lambda\lambda}k^2/L_D} \log(L_{\lambda\lambda}k/\varepsilon)\right) \\ &= O\left(n\sqrt{L_{\lambda\lambda}/L_D}\left(\log(L_{\lambda\lambda}/\varepsilon)\sum_{k=1}^{K_{\det}} k + \sum_{k=1}^{K_{\det}} k \log k\right)\right) \\ &\stackrel{(a)}{=} O\left(n\sqrt{L_{\lambda\lambda}/L_D}(L_D/\varepsilon)\left(\log(L_{\lambda\lambda}/\varepsilon) + \log(L_D/\varepsilon)\right)\right) \\ &= O\left(n\sqrt{L_{\lambda\lambda}L_D/\varepsilon} \log(L_{\lambda\lambda}L_D/\varepsilon)\right), \end{aligned}$$

where in (a) we use  $\sum_{k=1}^K k = \Theta(K^2)$ ,  $\sum_{k=1}^K k \log k = \Theta(K^2 \log K)$  and (3.26). We can also repeat this analysis to conclude that the dual oracle complexity for solving (3.5), i.e.,  $C_{\det,2}^D$  has the same order as  $C_{\det,1}^D$ . Since  $C_{\det}^D = C_{\det,1}^D + C_{\det,2}^D$ , we complete the proof.  $\square$

**4. Randomized IPDS Framework.** When the saddle function  $\Phi(\cdot, \cdot)$  has a large number of components (i.e.,  $n$  is large), we propose to find the inexact solutions in steps 1, 3 and 6 of Algorithm 1 using randomized first-order methods. This is because randomized first-order methods, especially those incorporating the variance-reduction techniques (e.g., [26, 15]), enjoy superior oracle complexities compared to their deterministic counterparts, for solving finite-sum convex composite problems. Based on this idea, we develop our randomized IPDS framework, which is shown in

**Algorithm 2** Randomized primal-dual smoothed gap reduction framework

**Input:** Initial smoothing parameter  $\rho_0 > 0$ , nonnegative error sequences  $\{\eta_k\}_{k \in \mathbb{Z}_+}$  and  $\{\gamma_k\}_{k \in \mathbb{Z}_+}$ , interpolation sequence  $\{\tau_k\}_{k \in \mathbb{Z}_+} \subseteq (0, 1)$  and randomized first-order algorithms  $M_1$  and  $M_2$ .

**Initialize:**  $x^0 \in \mathcal{X}$ ,  $\lambda^0 \in \Lambda$  and  $k = 0$ .

**Repeat** (until some convergence criterion is met)

1. Use  $M_1$  to find  $\tilde{\lambda}_{\rho_k, \eta_k}(x^k) \in \Lambda$  such that

$$(4.1) \quad \mathbb{E}[\psi_{\rho_k}^P(x^k) - S_{\rho_k}(x^k, \tilde{\lambda}_{\rho_k, \eta_k}(x^k)) \mid \mathcal{F}_{k,0}] \leq \eta_k \quad \text{a.s.}$$

2. Set  $\hat{\lambda}^k = \tau_k \lambda^k + (1 - \tau_k) \tilde{\lambda}_{\rho_k, \eta_k}(x^k)$ .
3. Use  $M_2$  to find  $\tilde{x}_{\gamma_k}(\hat{\lambda}^k) \in X$  such that

$$(4.2) \quad \mathbb{E}[S(\tilde{x}_{\gamma_k}(\hat{\lambda}^k), \hat{\lambda}^k) - \psi^D(\hat{\lambda}^k) \mid \mathcal{F}_{k,1}] \leq \gamma_k \quad \text{a.s.}$$

4. Set  $x^{k+1} = \tau_k x^k + (1 - \tau_k) \tilde{x}_{\gamma_k}(\hat{\lambda}^k)$ .
5. Set  $\rho_{k+1} = \tau_k \rho_k$ .
6. Use  $M_1$  to find  $\tilde{\lambda}_{\rho_{k+1}, \eta_k}(x^{k+1}) \in \Lambda$  such that

$$(4.3) \quad \mathbb{E}[\psi_{\rho_{k+1}}^P(x^{k+1}) - S_{\rho_{k+1}}(x^{k+1}, \tilde{\lambda}_{\rho_{k+1}, \eta_k}(x^{k+1})) \mid \mathcal{F}_{k,2}] \leq \eta_k \quad \text{a.s.}$$

7. Set  $\lambda^{k+1} = \tau_k \lambda^k + (1 - \tau_k) \tilde{\lambda}_{\rho_{k+1}, \eta_k}(x^{k+1})$ .
8. Set  $k = k + 1$ .

**Output:**  $(x^k, \lambda^k)$ .

**Algorithm 2.**

Note that at each iteration  $k$ , in steps 1, 3 and 6 of Algorithm 2, the inexact solutions that we aim to find are functions of some *stochastic* iterates, i.e.,  $x^k$ ,  $\hat{\lambda}^k$  and  $x^{k+1}$ . Therefore, to analyze such inexact solutions, we need to properly condition on the past information. To this end, let us denote the probability space for all the stochastic processes in Algorithm 2 by  $(\Omega, \mathcal{B}, \text{Pr})$  (where  $\mathcal{B}$  denotes the Borel  $\sigma$ -field on  $\Omega$ ) and define a filtration  $\bigcup_{k \in \mathbb{Z}_+} \{\mathcal{F}_{k,i}\}_{i=0}^2$ , where  $\mathcal{F}_{0,0} \triangleq \{\emptyset, \Omega\}$  and for any  $k \in \mathbb{Z}_+$ ,

$$(4.4) \quad \mathcal{F}_{k,1} \triangleq \sigma\{\mathcal{F}_{k,0} \cup \sigma\{\tilde{\lambda}_{\rho_k, \eta_k}(x^k)\}\},$$

$$(4.5) \quad \mathcal{F}_{k,2} \triangleq \sigma\{\mathcal{F}_{k,1} \cup \sigma\{\tilde{x}_{\gamma_k}(\hat{\lambda}^k)\}\},$$

$$(4.6) \quad \mathcal{F}_{k+1,0} \triangleq \sigma\{\mathcal{F}_{k,2} \cup \sigma\{\tilde{\lambda}_{\rho_{k+1}, \eta_k}(x^{k+1})\}\}.$$

Here we overload the notation  $\sigma\{\cdot\}$  to represent the  $\sigma$ -field generated by either a family of (Borel-measurable) sets or a random variable. From this definition, we clearly have

$$(4.7) \quad \mathcal{F}_{k,0} \subseteq \mathcal{F}_{k,1} \subseteq \mathcal{F}_{k,2} \subseteq \mathcal{F}_{k+1,0}, \quad \forall k \in \mathbb{Z}_+.$$

In addition, we have  $x^0, \lambda^0 \in \mathcal{F}_{0,0}$  and for any  $k \in \mathbb{Z}_+$ ,

$$(4.8) \quad \tilde{\lambda}_{\rho_k, \eta_k}(x^k), \hat{\lambda}^k \in \mathcal{F}_{k,1}, \quad \tilde{x}_{\gamma_k}(\hat{\lambda}^k), x^{k+1} \in \mathcal{F}_{k,2}, \quad \tilde{\lambda}_{\rho_{k+1}, \eta_k}(x^{k+1}), \lambda^{k+1} \in \mathcal{F}_{k+1,0},$$

where for any function  $\xi$  and  $\sigma$ -field  $\mathcal{F}$ ,  $\xi \in \mathcal{F}$  denotes that  $\xi$  is measurable w.r.t.  $\mathcal{F}$ .

**4.1. Solving sub-problems.** Similar to the deterministic case (cf. Section 3.2), we choose both  $M_1$  and  $M_2$  to be the optimal first-order randomized method in [15, Algorithm 3]. We denote this method as RPD since it is based on the randomized

primal-dual gradient. Consider the optimization problem in (3.6), where the smooth function  $\phi_1$  has a finite-sum structure, i.e.,

$$(4.9) \quad \phi_1(u) = (1/m) \sum_{i=1}^m \varphi_i(u)$$

and each  $\varphi_i$  is convex and  $L'_i$ -smooth on  $\mathcal{U}$  (so that  $L' = (1/m) \sum_{i=1}^m L'_i$ ). From the convergence results in [15, Corollary 1], for any starting point  $u^0 \in \mathcal{U}^o$ , to have  $\mathbb{E}[\Psi(\tilde{u}^N) - \Psi(u^*)] \leq \varepsilon$  (where  $\tilde{u}^N$  denotes the  $N$ -th iterate of RPD), it suffices to let

$$(4.10) \quad \begin{aligned} N &= 2(m + \sqrt{8m\kappa'}) \log \left( 2(L'/\sqrt{\mu'} + \sqrt{\mu'})^2 (m + \sqrt{8m\kappa'}) D_{\varpi}(u^*, u^0) / \varepsilon \right) \\ &\leq 2(m + \sqrt{8m\kappa'}) \log \left( 2(L'/\sqrt{\mu'} + \sqrt{\mu'})^2 (m + \sqrt{8m\kappa'}) R_{\varpi, \mathcal{U}}(u^0) / \varepsilon \right) \\ &= O((m + \sqrt{m\kappa'}) \log(L'\kappa'(m + \sqrt{m\kappa'}) R_{\varpi, \mathcal{U}}(u^0) / \varepsilon)), \end{aligned}$$

where  $\varpi$  denotes the DGF w.r.t.  $(\phi_2, \mathcal{U})$ ,  $D_{\varpi}(u, u') \triangleq \varpi(u) - \varpi(u') - \langle \nabla \varpi(u'), u - u' \rangle$  (for any  $u \in \mathcal{U}$  and  $u' \in \mathcal{U}^o$ ) denotes the Bregman distance induced by  $\varpi$  on  $\mathcal{U}$ , and

$$(4.11) \quad R_{\varpi, \mathcal{U}}(u^0) \triangleq \sup_{u \in \mathcal{U}} D_{\varpi}(u, u^0) < +\infty,$$

since  $\varpi$  is continuous on the compact set  $\mathcal{U}$ . Therefore, similar to the (deterministic) APG algorithm, we do not need to know  $u^*$  or  $\Psi(u^*)$  to use RPD for solving (3.6).

In the context of Algorithm 2, it follows that we should fix  $\bar{x} \in \mathcal{X}^o$  and use it as the starting point to solve (4.2) in each iteration  $k$ . Similarly, we can solve (4.1) and (4.3) in each iteration  $k$  by using a fixed  $\bar{\lambda} \in \Lambda^o$  as the starting point. By doing so, both  $R_{\bar{\varpi}, \mathcal{X}}(\bar{x})$  and  $R_{\omega, \Lambda}(\bar{\lambda})$  are finite constants that are independent of  $k$ .

Based on (4.10), for *any*  $x \in \mathcal{X}$ , if we denote  $C_{M_1}$  as the number of dual oracle calls of  $M_1$  to find an  $\eta$ -inexact solution of (2.8) in expectation, i.e.,  $\tilde{\lambda}_{\rho, \eta}(x)$  such that  $\mathbb{E}[\widehat{\psi}_{\rho}^P(\lambda) - \widehat{S}_{\rho}^D(x, \tilde{\lambda}_{\rho, \eta}(x))] \leq \eta$ , then

$$(4.12) \quad C_{M_1} = O \left( \left( n + \sqrt{nL_{\lambda\lambda}/\rho} \right) \log \left( L_{\lambda\lambda} (n + \sqrt{nL_{\lambda\lambda}/\rho}) / (\rho\eta) \right) \right),$$

where  $\bar{\lambda}$  is any point in  $\Lambda^o$ . Similarly, for *any*  $\lambda \in \Lambda$ , if we denote  $C_{M_2}$  as the number of primal oracle calls of  $M_2$  to find a  $\gamma$ -inexact solution of (2.4), i.e.,  $\tilde{x}_{\gamma}(\lambda)$  such that  $\mathbb{E}[\widehat{S}^P(\tilde{x}_{\gamma}(\lambda), \lambda) - \widehat{\psi}^D(\lambda)] \leq \gamma$ , then

$$(4.13) \quad C_{M_2} = O \left( (n + \sqrt{n\kappa_{\mathcal{X}}}) \log \left( (L + L_{xx})(n + \sqrt{n\kappa_{\mathcal{X}}}) / (\mu\gamma) \right) \right).$$

**4.2. Convergence Analysis.** We analyze the convergence rate of Algorithm 2 in expectation. For convergence results w.h.p., we refer readers to Section 4.4.

LEMMA 4.1. *In Algorithm 2, for any  $k \in \mathbb{Z}_+$ , if  $\rho_{k+1} \geq 4(1 - \tau_k)^2 L_D$ , then*

$$(4.14) \quad \mathbb{E}[\Delta_{\rho_{k+1}}(x^{k+1}, \lambda^{k+1}) | \mathcal{F}_{k,0}] \leq \tau_k \Delta_{\rho_k}(x^k, \lambda^k) + 2\gamma_k + 2\eta_k \quad a.s.$$

*Proof.* To prove this lemma, one only needs to properly incorporate the inexact criteria in (4.1), (4.2) and (4.3) (which involve conditional expectations) into the proof of Lemma 3.3. The key steps are: i) taking conditional expectation over the steps in the proof of Lemma 3.3 by using the measurability results in (4.8) and ii) applying the tower property of conditional expectation by using the nested relation in (4.7).

Specifically, at the  $k$ -th iteration, we first modify the proof of Proposition 2.1 and show that

$$(4.15) \quad \begin{aligned} \mathbb{E}[\widehat{S}^P(\tilde{x}_{\gamma}(\hat{\lambda}^k), \hat{\lambda}^k) - \widehat{\psi}^D(\lambda^{k+1}) + \langle \nabla_{\lambda} \widehat{S}^P(\tilde{x}_{\gamma}(\hat{\lambda}^k), \hat{\lambda}^k), \lambda^{k+1} - \hat{\lambda}^k \rangle | \mathcal{F}_{k,1}] \\ \leq L_D \mathbb{E}[\|\hat{\lambda}^k - \lambda^{k+1}\|^2 | \mathcal{F}_{k,1}] + 2\gamma_k. \end{aligned}$$

(For notational brevity, we omit ‘a.s.’ here and for all the inequalities below.) Furthermore, since  $\mathcal{F}_{k,0} \subseteq \mathcal{F}_{k,1}$ , if we take  $\mathbb{E}[\cdot | \mathcal{F}_{k,0}]$  over (4.15), then we have

$$(4.16) \quad \mathbb{E}[\widehat{S}^{\text{P}}(\tilde{x}_\gamma(\hat{\lambda}^k), \hat{\lambda}^k) - \widehat{\psi}^{\text{D}}(\lambda^{k+1}) + \langle \nabla_\lambda \widehat{S}^{\text{P}}(\tilde{x}_\gamma(\hat{\lambda}^k), \hat{\lambda}^k), \lambda^{k+1} - \hat{\lambda}^k \rangle | \mathcal{F}_{k,0}] \\ \leq L_{\text{D}} \mathbb{E}[\|\hat{\lambda}^k - \lambda^{k+1}\|^2 | \mathcal{F}_{k,0}] + 2\gamma_k.$$

In addition, from (3.18), we have

$$(4.17) \quad \mathbb{E}[\|\tilde{\lambda}_{\rho_{k+1}, \eta_k}(x^{k+1}) - \lambda_{\rho_k}^*(x^k)\|^2 | \mathcal{F}_{k,0}] \geq \frac{\mathbb{E}[\|\lambda^{k+1} - \hat{\lambda}^k\|^2 | \mathcal{F}_{k,0}]}{2(1 - \tau_k)^2} - \frac{2\eta_k}{\rho_k}.$$

Now, we can take  $\mathbb{E}[\cdot | \mathcal{F}_{k,0}]$  over Equation (c) in (3.15), and use (4.16), (4.17) and the fact that  $x^k, \lambda^k \in \mathcal{F}_{k,0}$  to obtain

$$(4.18) \quad \tau_k \Delta_{\rho_k}(x^k, \lambda^k) \geq \mathbb{E}[S_{\rho_{k+1}}(x^{k+1}, \tilde{\lambda}_{\rho_{k+1}, \eta_k}(x^{k+1})) + \frac{\rho_{k+1} \|\lambda^{k+1} - \hat{\lambda}^k\|^2}{4(1 - \tau_k)^2} - \tau_k \eta_k \\ - (\psi^{\text{D}}(\lambda^{k+1}) + L_{\text{D}} \|\hat{\lambda}^k - \lambda^{k+1}\|^2 + 2\gamma_k) | \mathcal{F}_{k,0}].$$

Again, since  $\mathcal{F}_{k,0} \subseteq \mathcal{F}_{k,2}$ , if we take  $\mathbb{E}[\cdot | \mathcal{F}_{k,0}]$  over (4.3), then we have

$$(4.19) \quad \mathbb{E}[\psi_{\rho_{k+1}}^{\text{P}}(x^{k+1}) - S_{\rho_{k+1}}(x^{k+1}, \tilde{\lambda}_{\rho_{k+1}, \eta_k}(x^{k+1})) | \mathcal{F}_{k,0}] \leq \eta_k.$$

We then substitute (4.19) into (4.18), and use the condition  $\rho_{k+1} \geq 4(1 - \tau_k)^2 L_{\text{D}}$  to conclude that

$$(4.20) \quad \tau_k \Delta_{\rho_k}(x^k, \lambda^k) \geq \mathbb{E}[\Delta_{\rho_{k+1}}(x^{k+1}, \lambda^{k+1}) | \mathcal{F}_{k,0}] - 2\eta_k - 2\gamma_k. \quad \square$$

Based on Lemma 4.1, we can derive the convergence rate of Algorithm 2 in expectation. The proof directly follows that of Theorem 3.5 and the tower property of conditional expectation, hence it is omitted.

**THEOREM 4.2.** *In Algorithm 2, if we choose the input parameters  $\rho_0, \{\tau_k\}_{k \in \mathbb{Z}_+}, \{\gamma_k\}_{k \in \mathbb{Z}_+}$  and  $\{\eta_k\}_{k \in \mathbb{Z}_+}$  in the same way as in Theorem 3.5, then for any starting point  $(x^0, \lambda^0) \in \mathcal{X} \times \Lambda$  and  $K \in \mathbb{N}$ ,  $\mathbb{E}[\Delta_{\rho_K}(x^K, \lambda^K)] \leq B'_\Delta(K, \varepsilon)$  (defined in (3.22)). Moreover, if Assumption 3.1 holds, then  $\mathbb{E}[\Delta(x^K, \lambda^K)] \leq B_\Delta(K, \varepsilon)$  (defined in (3.23)).*

Denote  $K_{\text{stoc}}$  as the number of iterations needed to achieve an  $\varepsilon$ -expected duality gap in Algorithm 2. Based on Theorem 4.2, we have that  $K_{\text{stoc}} = K_{\text{det}} = O(\sqrt{L_{\text{D}}/\varepsilon})$ .

**4.3. Oracle Complexity.** We analyze the primal and dual oracle complexities of Algorithm 2 to achieve an  $\varepsilon$ -expected duality gap, i.e.,  $\mathbb{E}[\Delta(x^{\text{out}}, \lambda^{\text{out}})] \leq \varepsilon$ .

**THEOREM 4.3.** *Let Assumption 3.1 hold. In Algorithm 2, for any starting point  $(x^0, \lambda^0) \in \mathcal{X} \times \Lambda$ , denote  $C_{\text{stoc}}^{\text{P}}$  and  $C_{\text{stoc}}^{\text{D}}$  as the primal and dual oracle complexities to achieve an  $\varepsilon$ -expected duality gap, respectively. Then we have*

$$(4.21) \quad C_{\text{stoc}}^{\text{P}} = O\left(\left(n + \sqrt{n\kappa_{\mathcal{X}}}\right) \sqrt{\frac{L_{\text{D}}}{\varepsilon}} \log\left(\frac{\kappa_{\mathcal{X}} L_{\text{D}} (n + \sqrt{n\kappa_{\mathcal{X}}})}{\varepsilon}\right)\right),$$

$$(4.22) \quad C_{\text{stoc}}^{\text{D}} = O\left(\left(n \sqrt{\frac{L_{\text{D}}}{\varepsilon}} + \frac{\sqrt{nL_{\lambda\lambda}L_{\text{D}}}}{\varepsilon}\right) \log\left(\frac{L_{\lambda\lambda} (n + \sqrt{nL_{\lambda\lambda}/L_{\text{D}}})}{\varepsilon}\right)\right).$$

*Proof.* The proof follows the same line of argument as that of Theorem 3.7, hence we only outline the important steps. Based on the choice of  $\gamma_k$  in (3.21) and the complexity of  $M_2$  in (4.13), we have

$$C_{\text{stoc}}^{\text{P}} = \sum_{k=1}^{K_{\text{det}}} O\left((n + \sqrt{n\kappa_{\mathcal{X}}}) \log\left((L + L_{xx})(n + \sqrt{n\kappa_{\mathcal{X}}})k/(\mu\varepsilon)\right)\right)$$

$$\begin{aligned}
 &= O\left((n + \sqrt{n\kappa_{\mathcal{X}}})\left(K_{\text{stoc}} \log\left((L + L_{xx})(n + \sqrt{n\kappa_{\mathcal{X}}})/(\mu\varepsilon)\right) + \log(K_{\text{stoc}}!)\right)\right) \\
 &= O\left((n + \sqrt{n\kappa_{\mathcal{X}}})\sqrt{L_{\text{D}}}/\varepsilon\left(\log\left((L + L_{xx})(n + \sqrt{n\kappa_{\mathcal{X}}})/(\mu\varepsilon)\right) + \log(L_{\text{D}}/\varepsilon)\right)\right) \\
 &= O\left((n + \sqrt{n\kappa_{\mathcal{X}}})\sqrt{L_{\text{D}}}/\varepsilon \log\left((L + L_{xx})L_{\text{D}}(n + \sqrt{n\kappa_{\mathcal{X}}})/(\mu\varepsilon)\right)\right).
 \end{aligned}$$

Using the same reasoning as in the proof of Theorem 3.7, the dual oracle complexities for solving both (4.1) and (4.3) have the same order, so it suffices to only analyze the complexity for solving (4.1). Specifically, based on (4.12), we have

$$\begin{aligned}
 C_{\text{stoc},1}^{\text{D}} &= \sum_{k=1}^{K_{\text{stoc}}} O\left((n + \sqrt{nL_{\lambda\lambda}k^2/L_{\text{D}}}) \log(L_{\lambda\lambda}(n + \sqrt{nL_{\lambda\lambda}k^2/L_{\text{D}}})k/(L_{\text{D}}\varepsilon))\right) \\
 &\stackrel{(a)}{=} \sum_{k=1}^{K_{\text{stoc}}} O\left((n + k\sqrt{nL_{\lambda\lambda}/L_{\text{D}}})\left(\log(L_{\lambda\lambda}(n + \sqrt{nL_{\lambda\lambda}/L_{\text{D}}})/(L_{\text{D}}\varepsilon)) + \log k\right)\right) \\
 &= O\left((K_{\text{stoc}}n + \sqrt{nL_{\lambda\lambda}/L_{\text{D}}}\sum_{k=1}^{K_{\text{stoc}}} k) \log(L_{\lambda\lambda}(n + \sqrt{nL_{\lambda\lambda}/L_{\text{D}}})/(L_{\text{D}}\varepsilon))\right. \\
 &\quad \left.+ (n\sum_{k=1}^{K_{\text{stoc}}} \log k + \sqrt{L_{\lambda\lambda}/L_{\text{D}}}\sum_{k=1}^{K_{\text{stoc}}} k \log k)\right) \\
 &= O\left((n\sqrt{L_{\text{D}}}/\varepsilon + \sqrt{nL_{\lambda\lambda}L_{\text{D}}}/\varepsilon) \log(L_{\lambda\lambda}(n + \sqrt{nL_{\lambda\lambda}/L_{\text{D}}})/(L_{\text{D}}\varepsilon))\right. \\
 &\quad \left.+ \sqrt{L_{\text{D}}}/\varepsilon \log(L_{\text{D}}/\varepsilon)(n + \sqrt{nL_{\lambda\lambda}/\varepsilon})\right) \\
 &= O\left((n\sqrt{L_{\text{D}}}/\varepsilon + \sqrt{nL_{\lambda\lambda}L_{\text{D}}}/\varepsilon) \log(L_{\lambda\lambda}(n + \sqrt{nL_{\lambda\lambda}/L_{\text{D}}})/\varepsilon)\right),
 \end{aligned}$$

where (a) holds since  $n \leq kn$ . We obtain (4.22) by noting that  $C_{\text{stoc}}^{\text{D}} = \Theta(C_{\text{stoc},1}^{\text{D}})$ .  $\square$

If we compare the results in Theorem 4.3 with those in Theorem 3.7, in terms of the dependence of the primal oracle complexity on  $n$ ,  $\kappa_{\mathcal{X}}$  and  $\varepsilon$ , the randomized framework (i.e., Algorithm 2) indeed yields improvement upon the deterministic one (i.e., Algorithm 1), from  $\tilde{O}(n\sqrt{\kappa_{\mathcal{X}}}/\varepsilon)$  to  $\tilde{O}((n + \sqrt{n\kappa_{\mathcal{X}}})/\sqrt{\varepsilon})$  (where recall that  $\tilde{O}(\cdot)$  omits the the log-factors in  $n$ ,  $\kappa_{\mathcal{X}}$  and  $\varepsilon$ ). Similarly, the dual oracle complexity has also been improved from  $\tilde{O}(n/\varepsilon)$  to  $\tilde{O}(n/\sqrt{\varepsilon} + \sqrt{n}/\varepsilon)$ .

**4.4. Convergence with High Probability.** Apart from convergence in expectation, given an error probability  $\delta \in (0, 1)$ , we can modify the inexact criteria in Algorithm 2 (i.e., (4.1), (4.2) and (4.3)) to obtain an  $\varepsilon$ -duality gap w.p. at least  $1 - \delta$ , i.e.,  $\Pr\{\Delta(x^{\text{out}}, \lambda^{\text{out}}) \leq \varepsilon\} \geq 1 - \delta$ .

**THEOREM 4.4.** *Let Assumption 3.1 hold and  $\varepsilon > 0$  and  $\delta \in (0, 1)$  be given. In Algorithm 2, choose the input parameters  $\rho_0$ ,  $\{\tau_k\}_{k \in \mathbb{Z}_+}$ ,  $\{\gamma_k\}_{k \in \mathbb{Z}_+}$  and  $\{\eta_k\}_{k \in \mathbb{Z}_+}$  in the same way as in Theorem 3.5, fix the total number of iterations  $K \in \mathbb{N}$  and modify the inexact criteria (4.1), (4.2) and (4.3) to*

$$(4.23) \quad \mathbb{E}[\psi_{\rho_k}^{\text{P}}(x^k) - S_{\rho_k}(x^k, \tilde{\lambda}_{\rho_k, \eta_k}(x^k)) \mid \mathcal{F}_{k,0}] \leq \eta_k \delta / (3K) \quad a.s.,$$

$$(4.24) \quad \mathbb{E}[S(\tilde{x}_{\gamma_k}(\hat{\lambda}^k), \hat{\lambda}^k) - \psi^{\text{D}}(\hat{\lambda}^k) \mid \mathcal{F}_{k,1}] \leq \gamma_k \delta / (3K) \quad a.s.,$$

$$(4.25) \quad \mathbb{E}[\psi_{\rho_{k+1}}^{\text{P}}(x^{k+1}) - S_{\rho_{k+1}}(x^{k+1}, \tilde{\lambda}_{\rho_{k+1}, \eta_k}(x^{k+1})) \mid \mathcal{F}_{k,2}] \leq \eta_k \delta / (3K) \quad a.s.,$$

respectively. If we set  $K = K'_{\text{det}} \triangleq 2 \left\lceil \sqrt{\max\{\Delta_{\rho_0}(x^0, \lambda^0), 0\}}/\varepsilon \right\rceil + 1$ , then

$$(4.26) \quad \Pr\{\Delta_{\rho_K}(x^K, \lambda^K) \leq \varepsilon\} \geq 1 - \delta.$$

Furthermore, if we set  $K = K_{\text{det}}$  as in (3.26), then

$$(4.27) \quad \Pr\{\Delta(x^K, \lambda^K) \leq \varepsilon\} \geq 1 - \delta.$$

*Proof.* First, let us define the events  $\mathcal{A}_{0,0} \triangleq \Omega$ , and for any  $k \in \mathbb{Z}_+$ ,

$$\begin{aligned}\mathcal{A}_{k,1} &\triangleq \{\psi_{\rho_k}^{\text{P}}(x^k) - S_{\rho_k}(x^k, \tilde{\lambda}_{\rho_k, \eta_k}(x^k)) \leq \eta_k\}, \\ \mathcal{A}_{k,2} &\triangleq \{S(\tilde{x}_{\gamma_k}(\hat{\lambda}^k), \hat{\lambda}^k) - \psi^{\text{D}}(\hat{\lambda}^k) \leq \gamma_k\}, \\ \mathcal{A}_{k+1,0} &\triangleq \{\psi_{\rho_{k+1}}^{\text{P}}(x^{k+1}) - S_{\rho_{k+1}}(x^{k+1}, \tilde{\lambda}_{\rho_{k+1}, \eta_k}(x^{k+1})) \leq \eta_k\}.\end{aligned}$$

Also, for any measurable event  $\mathcal{A}$ , denote its complement as  $\mathcal{A}^c$  and its indicator function as  $\mathbb{I}_{\mathcal{A}}$ , i.e.,  $\mathbb{I}_{\mathcal{A}}(z) = 1$  if  $z \in \mathcal{A}$  and 0 otherwise.

Fix any  $k \in \{0, \dots, K-1\}$ . From Markov's inequality and (4.23), we have

$$(4.28) \quad \Pr\{\mathcal{A}_{k,1}^c \mid \mathcal{F}_{k,0}\} \leq \mathbb{E}[\psi_{\rho_k}^{\text{P}}(x^k) - S_{\rho_k}(x^k, \tilde{\lambda}_{\rho_k, \eta_k}(x^k)) \mid \mathcal{F}_{k,0}] / \eta_k \leq \delta / (3K) \quad \text{a.s.}$$

Since  $\bigcup_{i=0}^{k-1} \{\mathcal{A}_{i,1}, \mathcal{A}_{i,2}, \mathcal{A}_{i+1,0}\} \subseteq \mathcal{F}_{k,0}$ , we have that

$$(4.29) \quad \mathcal{C}_{k,0} \in \mathcal{F}_{k,0}, \quad \text{where} \quad \mathcal{C}_{k,0} \triangleq \bigcap_{i=0}^{k-1} (\mathcal{A}_{i,1} \cap \mathcal{A}_{i,2} \cap \mathcal{A}_{i+1,0}).$$

(When  $k = 0$ , define  $\mathcal{C}_{0,0} \triangleq \mathcal{A}_{0,0}$ .) In addition, note that  $\Pr\{\mathcal{C}_{k,0}\} > 0$ , since

$$(4.30) \quad \Pr\{\mathcal{C}_{k,0}^c\} = \Pr\left\{\bigcup_{i=0}^{k-1} (\mathcal{A}_{k-1,1}^c \cup \mathcal{A}_{k-1,2}^c \cup \mathcal{A}_{k,0}^c)\right\} \leq (3k)\delta / (3K) \leq \delta < 1.$$

We then take conditional expectation  $\mathbb{E}[\cdot \mid \mathcal{C}_{k,0}]$  in (4.28) to obtain

$$(4.31) \quad \mathbb{E}[\Pr\{\mathcal{A}_{k,1} \mid \mathcal{F}_{k,0}\} \mid \mathcal{C}_{k,0}] \geq 1 - \delta / (3K).$$

On the other hand,

$$\begin{aligned}\mathbb{E}[\Pr\{\mathcal{A}_{k,1} \mid \mathcal{F}_{k,0}\} \mid \mathcal{C}_{k,0}] &= \mathbb{E}[\Pr\{\mathcal{A}_{k,1} \mid \mathcal{F}_{k,0}\} \mathbb{I}_{\mathcal{C}_{k,0}}] / \mathbb{P}(\mathcal{C}_{k,0}) \\ &\stackrel{(a)}{=} \mathbb{E}[\mathbb{I}_{\mathcal{A}_{k,1}} \mathbb{I}_{\mathcal{C}_{k,0}}] / \mathbb{P}(\mathcal{C}_{k,0}) = \Pr\{\mathcal{A}_{k,1} \mid \mathcal{C}_{k,0}\},\end{aligned}$$

where (a) follows since  $\mathcal{C}_{k,0} \in \mathcal{F}_{k,0}$ . Therefore, we have

$$(4.32) \quad \Pr\{\mathcal{A}_{k,1} \mid \mathcal{C}_{k,0}\} \geq 1 - \delta / (3K).$$

Similarly, if we define  $\mathcal{C}_{k,1} \triangleq \mathcal{C}_{k,0} \cap \mathcal{A}_{k,1}$  and  $\mathcal{C}_{k,2} \triangleq \mathcal{C}_{k,1} \cap \mathcal{A}_{k,2}$ , then we also have

$$(4.33) \quad \Pr\{\mathcal{A}_{k,2} \mid \mathcal{C}_{k,1}\} \geq 1 - \delta / (3K), \quad \Pr\{\mathcal{A}_{k+1,0} \mid \mathcal{C}_{k,2}\} \geq 1 - \delta / (3K).$$

From Theorem 3.5, we know that if  $K = K'_{\text{det}}$  and the event  $\bigcap_{k=0}^{K-1} (\mathcal{A}_{k,1} \cap \mathcal{A}_{k,2} \cap \mathcal{A}_{k+1,0})$  occurs, then  $\Delta_{\rho_K}(x^K, \lambda^K) \leq \varepsilon$ . Therefore,

$$\begin{aligned}\Pr\{\Delta_{\rho_K}(x^K, \lambda^K) \leq \varepsilon\} &\geq \Pr\left\{\bigcap_{k=0}^{K-1} (\mathcal{A}_{k,1} \cap \mathcal{A}_{k,2} \cap \mathcal{A}_{k+1,0})\right\} \\ &= \prod_{k=0}^{K-1} \Pr\{\mathcal{A}_{k+1,0} \mid \mathcal{C}_{k,2}\} \Pr\{\mathcal{A}_{k,2} \mid \mathcal{C}_{k,1}\} \Pr\{\mathcal{A}_{k,1} \mid \mathcal{C}_{k,0}\} \\ &\stackrel{(a)}{\geq} (1 - \delta / (3K))^{3K} \stackrel{(b)}{\geq} 1 - \delta,\end{aligned}$$

where (a) follows from (4.32) and (4.33) and (b) follows from Bernoulli's inequality. By the same reasoning, we can also show that if  $K = K_{\text{det}}$ , then (4.27) holds.  $\square$

Based on the inexact criteria in Theorem 4.4, we can also derive the primal and dual oracle complexities of obtaining an  $\varepsilon$ -duality gap w.p. at least  $1 - \delta$ . The derivation is essentially the same as that of Theorem 4.3, hence it is omitted.

**THEOREM 4.5.** *Let Assumption 3.1 hold and  $\varepsilon > 0$  and  $\delta \in (0, 1)$  be given. In Algorithm 2, modify the inexact criteria (4.1), (4.2) and (4.3) in the same way as in Theorem 4.4. For any starting point  $(x^0, \lambda^0) \in \mathcal{X} \times \Lambda$ , denote  $C_{\text{hp}}^{\text{P}}$  and  $C_{\text{hp}}^{\text{D}}$  as the primal and dual oracle complexities to achieve an  $\varepsilon$ -duality gap w.p. at least  $1 - \delta$ .*

Then we have

$$(4.34) \quad C_{\text{hp}}^{\text{P}} = O\left((n + \sqrt{n\kappa_{\mathcal{X}}})\sqrt{L_{\text{D}}/\varepsilon} \log(\kappa_{\mathcal{X}}L_{\text{D}}(n + \sqrt{n\kappa_{\mathcal{X}}})/(\varepsilon\delta))\right),$$

$$(4.35) \quad C_{\text{hp}}^{\text{D}} = O\left((n\sqrt{L_{\text{D}}/\varepsilon} + \sqrt{nL_{\lambda\lambda}L_{\text{D}}/\varepsilon}) \log\left(L_{\lambda\lambda}L_{\text{D}}(n + \sqrt{nL_{\lambda\lambda}/L_{\text{D}}})/(\varepsilon\delta)\right)\right).$$

**4.5. Interface with Stochastic Approximation (SA).** We can also use our randomized framework (i.e., Algorithm 2) to solve the stochastic version of (1.1), where we only have access to stochastic gradients of  $f$ ,  $\Phi(\cdot, \lambda)$  and  $\Phi(x, \cdot)$ , i.e., unbiased estimators of  $\nabla f$ ,  $\nabla_x \Phi(\cdot, \lambda)$  and  $\nabla_\lambda \Phi(x, \cdot)$  with bounded second moments (and potentially other distributional assumptions; for details, we refer readers to [32]). Indeed, in each iteration  $k$ , the sub-problems in steps 1, 3 and 6 can be solved by the optimal stochastic first-order algorithms for stochastic strongly-convex composite problems, e.g., those in [10]. Specifically, in [10], the optimality gaps of the proposed algorithms were analyzed both in expectation (cf. (4.1)) and w.h.p. (cf. (4.28)). These two convergence results can interface with our convergence analyses in Sections 4.2 and 4.4, respectively. Thus, our convergence results for Algorithm 2, both in expectation (cf. Theorem 4.2) and w.h.p. (cf. Theorem 4.4), still hold in this case.

**5. Convex Optimization with Functional Constraints.** In this section, we apply our IPDS frameworks (i.e., Algorithms 1 and 2) to the Lagrangian (saddle point) problems associated with the constrained convex problems.

**5.1. Problem Setup.** We consider

$$(5.1) \quad \min_{x \in \mathcal{X}} f(x) + r(x) \quad \text{s. t.} \quad g_i(x) \leq 0, \forall i \in [n],$$

where  $\mathcal{X} \neq \emptyset$  is convex and compact,  $f$  is  $\mu$ -s.c. and  $L$ -smooth on  $\mathcal{X}$ ,  $r$  is CCP and admits a tractable BPP on  $\mathcal{X}$  (with DGF  $\bar{\omega}$ ; cf. Section 2.1), and for each  $i \in [n]$ ,  $g_i$  is convex and  $\alpha_i$ -smooth on  $\mathcal{X}$  (where  $\alpha_i \geq 0$ ). We assume that there exists  $\bar{x} \in \mathcal{X}^\circ$  (recall that  $\mathcal{X}^\circ = \mathcal{X} \cap \text{int dom } \bar{\omega}$ ) such that  $g_i(\bar{x}) < 0$ , for any  $i \in [n]$ , so Slater's condition holds for (5.1). Under these conditions, (5.1) has the unique primal optimal solution  $x^* \in \mathcal{X}$ , a dual optimal solution  $\lambda^* \in \mathbb{R}_+^n$  (where  $\mathbb{R}_+ \triangleq [0, +\infty)$ ) and zero duality gap. Moreover, any such  $(x^*, \lambda^*)$  is a saddle point of the Lagrangian problem associated with (5.1):

$$(5.2) \quad \min_{x \in \mathcal{X}} \max_{\lambda \in \mathbb{R}_+^n} \left\{ S(x, \lambda) = f(x) + r(x) + (1/n) \sum_{i=1}^n n \lambda_i g_i(x) \right\},$$

where  $\lambda_i$  denotes the  $i$ -th entry of  $\lambda$ . In addition, any saddle point  $(x^*, \lambda^*)$  of (5.2) is a primal-dual optimal solution pair for (5.1) with zero duality gap [6, Section 5.4]. This establishes the *equivalence* of solving (5.1) and (5.2).

Indeed, we observe that (5.2) has the same form as the SPP in (1.1). Specifically, if we set  $g = r$ ,  $h \equiv 0$ ,  $\Lambda = \mathbb{R}_+^n$  and  $\Phi_i(x, \lambda) = n \lambda_i g_i(x)$  in (1.1), then we recover (5.2). As a result,  $L_{xx}^i(\lambda) = n \lambda_i \alpha_i$ ,  $L_{\lambda\lambda}^i = 0$  and

$$(5.3) \quad L_{\lambda x}^i = n M_i, \quad \text{where} \quad M_i \triangleq \alpha_i D_{\mathcal{X}} + \inf_{x \in \mathcal{X}} \|\nabla g_i(x)\|_*.$$

In (5.3), we recall that  $D_{\mathcal{X}} < +\infty$  denotes the diameter of the set  $\mathcal{X}$ . (To obtain (5.3), we note that  $M_i \leq \sup_{x \in \mathcal{X}} \|\nabla g_i(x)\|_*$ . Then, by the  $\alpha_i$ -smoothness of  $g_i$ , we have  $\|\nabla g_i(x)\|_* \leq \alpha_i \|x - x'\| + \|\nabla g_i(x')\|_* \leq \alpha_i D_{\mathcal{X}} + \|\nabla g_i(x')\|_*$ , for any  $x, x' \in \mathcal{X}$ .) Thus,

$$(5.4) \quad L_{xx}(\lambda) = \sum_{i=1}^n \lambda_i \alpha_i, \quad L_{\lambda x} = M \triangleq \sum_{i=1}^n M_i, \quad L_{\text{D}} = 2M^2/\mu.$$

In addition, since  $\Phi(x, \cdot)$  is linear, we can choose  $\mathbb{E}_2 = (\mathbb{R}^n, \|\cdot\|_2)$ , where  $\|\cdot\|_2$  denotes the Euclidean norm. Subsequently, the problem in (2.8) now has closed-form

solution, i.e.,

$$(5.5) \quad ([g_i(x)]_+/\rho)_{i=1}^n = \arg \max_{\lambda \in \mathbb{R}_+^n} \sum_{i=1}^n \lambda_i g_i(x) - (\rho/2) \|\lambda\|_2^2,$$

where  $[\cdot]_+ \triangleq \max\{0, \cdot\}$  and we choose the DGF w.r.t.  $(0, \mathbb{R}_+^n)$  as  $\omega(\cdot) = (1/2) \|\cdot\|_2^2$ .

However, note that two of the assumptions that we make about (1.1) fail to hold for (5.2). First, in Assumption 3.1, we assume that  $\Lambda$  is bounded in (1.1), but it is unbounded in (5.2). Second, we assume that  $L_{xx}$  is a constant (w.r.t.  $\lambda$ ) in (1.1), but it depends (linearly) on  $\lambda$  in (5.2). That said, both of these challenges can be overcome. Before we present the details below, we first provide some intuition. For the first challenge, recall from Remark 3.2 that the boundedness of  $\Lambda$  is needed for two purposes, i.e., solving the problem in (2.8) inexactly and bounding the duality gap  $\Delta$  via its smoothed counterpart  $\Delta_\rho$ . In the case of (5.2), from (5.5), we see that the problem in (2.8) can be solved *exactly*. In addition, in Section 5.2, we will use a convergence criterion different from the duality gap  $\Delta$ . For these reasons,  $\Lambda$  need not be bounded for (5.2). For the second challenge, we propose to properly bound the growth of  $L_{xx}(\lambda)$  in each iteration of our frameworks via bounding  $\|\lambda\|_\infty$ , i.e., the  $\ell_\infty$ -norm of  $\lambda$ .

**5.2. Convergence Analysis.** For the constrained problem in (5.1), instead of the duality gap, it is more common to use the (primal) optimality gap and constraint violation together as the convergence criterion (e.g., [31]). Specifically, for any  $\varepsilon > 0$ ,  $\bar{x} \in \mathcal{X}$  is an  $\varepsilon$ -optimal and  $\varepsilon$ -feasible solution of (5.1) if

$$(5.6) \quad f(\bar{x}) - f(x^*) \leq \varepsilon, \quad \text{and} \quad [g_i(\bar{x})]_+ \leq \varepsilon, \quad \forall i \in [n].$$

Note that this is a primal convergence criterion. However, if we apply Algorithm 1 or 2 to (5.2), the established convergence results (in Theorems 3.5, 4.2 and 4.4) are all in terms of the smoothed duality gap  $\Delta_\rho$ . Thus, we need to relate  $\Delta_\rho$  to the criteria in (5.6). Indeed, in the following lemma, we will show that if there exists  $\bar{\lambda} \in \mathbb{R}_+^n$  such that both  $\Delta_\rho(\bar{x}, \bar{\lambda})$  and  $\rho$  are sufficiently small, then  $\bar{x} \in \mathcal{X}$  satisfies (5.6).

**LEMMA 5.1.** *Let  $(x^*, \lambda^*) \in \mathcal{X} \times \mathbb{R}_+^n$  be a saddle point of (5.2), so in particular  $x^*$  is the optimal solution of (5.1). For any  $\rho > 0$  and  $\epsilon \geq 0$ , if there exist  $\bar{x} \in \mathcal{X}$  and  $\bar{\lambda} \in \mathbb{R}_+^n$  that satisfy  $\Delta_\rho(\bar{x}, \bar{\lambda}) \leq \epsilon$ , then*

$$(5.7) \quad f(\bar{x}) - f(x^*) \leq \epsilon, \quad [g_i(\bar{x})]_+ \leq V_i(\epsilon, \rho) \triangleq (\lambda_i^* + \|\lambda^*\|_2) \rho + \sqrt{2\epsilon\rho}, \quad \forall i \in [n].$$

*Proof.* Since  $\lambda^* \in \mathbb{R}_+^n$  is an optimal solution of the dual problem  $\max_{\lambda \in \mathbb{R}_+^n} \psi^D(\lambda)$ , we have that  $\psi^D(\bar{\lambda}) \leq \psi^D(\lambda^*)$ . This implies  $\Delta_\rho(\bar{x}, \lambda^*) \leq \Delta_\rho(\bar{x}, \bar{\lambda}) \leq \epsilon$ . From the definition of  $\Delta_\rho$  in (2.14), we have

$$(5.8) \quad \epsilon \geq \Delta_\rho(\bar{x}, \lambda^*) \geq S(\bar{x}, \lambda^*) - (\rho/2) \|\lambda^*\|_2^2 - S(x^*, \lambda^*), \quad \forall x \in \mathcal{X}, \quad \forall \lambda \in \mathbb{R}_+^n.$$

We then choose  $x = x^*$  and  $\lambda = 0$  in (5.8) to obtain

$$(5.9) \quad \epsilon \geq S(\bar{x}, 0) - S(x^*, \lambda^*) = f(\bar{x}) - (f(x^*) + \sum_{i=1}^n \lambda_i^* g_i(x^*)) \geq f(\bar{x}) - f(x^*),$$

where the last step follows from  $\lambda_i^* \geq 0$  and  $g_i(x^*) \leq 0$ , for any  $i \in [n]$ .

Now fix any  $\theta > 0$  and  $i \in [n]$ . Let  $e_i \in \mathbb{R}^n$  denote the  $i$ -th standard basis vector, i.e.,  $(e_i)_i = 1$  and  $(e_i)_j = 0$  for any  $j \in [n] \setminus \{i\}$ . In (5.8), if we choose  $x = x^*$  and  $\lambda = \lambda^* + \theta_i e_i$ , where  $\theta_i = \theta$  if  $g_i(\bar{x}) > 0$  and 0 otherwise, then

$$\begin{aligned} \epsilon &\geq S(\bar{x}, \lambda^*) + \theta_i g_i(\bar{x}) - (\rho/2) \|\lambda^* + \theta_i e_i\|_2^2 - S(x^*, \lambda^*) \\ &\geq \theta_i g_i(\bar{x}) - (\rho/2) \|\lambda^* + \theta_i e_i\|_2^2 \geq \theta [g_i(\bar{x})]_+ - (\rho/2) \|\lambda^* + \theta e_i\|_2^2, \end{aligned}$$

where in the last step we use  $\lambda^* \geq 0$  and  $\theta \geq \theta_i \geq 0$ . After rearranging, we have

$$(5.10) \quad \begin{aligned} [g_i(\bar{x})]_+ &\leq \rho\lambda_i^* + \rho\theta/2 + (\rho\|\lambda^*\|_2^2 + 2\epsilon)/(2\theta) \\ &\stackrel{(a)}{\leq} \rho\lambda_i^* + \sqrt{\rho(\rho\|\lambda^*\|_2^2 + 2\epsilon)} \stackrel{(b)}{\leq} (\lambda_i^* + \|\lambda^*\|_2)\rho + \sqrt{2\epsilon\rho}, \end{aligned}$$

where in (a) we take the infimum over  $\theta > 0$  and in (b) we use  $\sqrt{a+b} \leq \sqrt{a} + \sqrt{b}$ , for any  $a, b \geq 0$ .  $\square$

Using similar arguments, we can derive a stochastic version of Lemma 5.1.

LEMMA 5.2. *Let  $(x^*, \lambda^*) \in \mathcal{X} \times \mathbb{R}_+^n$  be a saddle point of (5.2) and  $(\bar{x}, \bar{\lambda})$  be a (stochastic) primal-dual pair such that  $(\bar{x}, \bar{\lambda}) \in \mathcal{X} \times \mathbb{R}_+^n$  a.s. For any  $\rho > 0$  and  $\epsilon \geq 0$ , if  $(\bar{x}, \bar{\lambda})$  satisfies  $\mathbb{E}[\Delta_\rho(\bar{x}, \bar{\lambda})] \leq \epsilon$ , then*

$$\mathbb{E}[f(\bar{x})] - f(x^*) \leq \epsilon, \quad \mathbb{E}[[g_i(\bar{x})]_+] \leq V_i(\epsilon, \rho), \quad \forall i \in [n],$$

where  $V_i(\epsilon, \rho)$  is defined in (5.7). For any  $\delta \in (0, 1)$ , if we have  $\Pr\{\Delta_\rho(\bar{x}, \bar{\lambda}) \leq \epsilon\} \geq 1 - \delta$  (rather than  $\mathbb{E}[\Delta_\rho(\bar{x}, \bar{\lambda})] \leq \epsilon$ ), then

$$\Pr\{f(\bar{x}) - f(x^*) \leq \epsilon\} \geq 1 - \delta, \quad \Pr\{[g_i(\bar{x})]_+ \leq V_i(\epsilon, \rho)\} \geq 1 - \delta, \quad \forall i \in [n].$$

Based on Lemma 5.1 and the convergence results of Algorithm 1 in terms of the smoothed duality gap (cf. Theorem 3.5), we can easily derive the following results.

THEOREM 5.3. *Let  $(x^*, \lambda^*) \in \mathcal{X} \times \mathbb{R}_+^n$  be a saddle point of (5.2) and  $\epsilon > 0$  be given. If we apply Algorithm 1 to solving (5.2), with the input parameters  $\rho_0, \{\tau_k\}_{k \in \mathbb{Z}_+}, \{\gamma_k\}_{k \in \mathbb{Z}_+}$  chosen in the same way as in Theorem 3.5 and  $\eta_k = 0$  for any  $k \in \mathbb{Z}_+$ , then for any starting point  $(x^0, \lambda^0) \in \mathcal{X} \times \Lambda$  and  $K \in \mathbb{N}$ , we have*

$$(5.11) \quad f(x^K) - f(x^*) \leq W_f(K, \epsilon) \triangleq \frac{2[\Delta_{\rho_0}(x^0, \lambda^0)]_+}{(K+1)(K+2)} + \frac{\epsilon}{2},$$

$$(5.12) \quad [g_i(x^K)]_+ \leq W_{g_i}(K, \epsilon) \triangleq \frac{16(\lambda_i^* + \|\lambda^*\|_2)L_D + 8\sqrt{L_D[\Delta_{\rho_0}(x^0, \lambda^0)]_+}}{(K+1)(K+2)} + \frac{4\sqrt{L_D\epsilon}}{K+1}, \quad \forall i \in [m].$$

*Proof.* In Lemma 5.1, let us take  $\bar{x} = x^K, \bar{\lambda} = \lambda^K, \rho = \rho_K$  and  $\epsilon = [B'_\Delta(K, \epsilon)]_+$  (where  $B'_\Delta(K, \epsilon)$  is defined in (3.22)). Using that  $[a+b]_+ \leq [a]_+ + [b]_+$  (for any  $a, b \in \mathbb{R}$ ), we have  $[B'_\Delta(K, \epsilon)]_+ \leq W_f(K, \epsilon)$ . Thus we obtain (5.11). Using the analytic expression of  $\rho_K$  in (3.24), and  $\epsilon \leq W_f(K, \epsilon)$ , we also obtain (5.12).  $\square$

Similarly, based on Lemma 5.2 and the convergence results of Algorithm 2 in Theorems 4.2 and 4.4, we can show the following results using the same reasoning that leads to Theorem 5.3.

THEOREM 5.4. *Let  $(x^*, \lambda^*) \in \mathcal{X} \times \mathbb{R}_+^n$  be a saddle point of (5.2) and  $\epsilon > 0$  be given. Let us apply Algorithm 2 to solving (5.2), with the input parameters  $\rho_0, \{\tau_k\}_{k \in \mathbb{Z}_+}, \{\gamma_k\}_{k \in \mathbb{Z}_+}$  and  $\{\eta_k\}_{k \in \mathbb{Z}_+}$  chosen in the same way as in Theorem 5.3, and the starting point chosen to be any  $(x^0, \lambda^0) \in \mathcal{X} \times \Lambda$ . Then for any  $K \in \mathbb{N}$ ,*

$$(5.13) \quad \mathbb{E}[f(x^K)] - f(x^*) \leq W_f(K, \epsilon),$$

$$(5.14) \quad \mathbb{E}[[g_i(x^K)]_+] \leq W_{g_i}(K, \epsilon), \quad \forall i \in [m].$$

In addition, for any  $\delta \in (0, 1)$ , if we choose  $K = K'_{\det}$  and modify the inexact crite-

ria (4.1), (4.2) and (4.3) in the same way as in Theorem 4.4, then

$$(5.15) \quad \Pr\{f(x^K) - f(x^*) \leq W_f(K, \varepsilon)\} \geq 1 - \delta,$$

$$(5.16) \quad \Pr\{[g_i(x^K)]_+ \leq W_{g_i}(K, \varepsilon)\} \geq 1 - \delta, \quad \forall i \in [m].$$

From Theorems 5.3 and 5.4, we see that for Algorithm 1 to find an  $\varepsilon$ -optimal and  $\varepsilon$ -feasible solution of (5.1), or for Algorithm 2 to find such a solution in expectation (i.e., a solution that satisfies both (5.13) and (5.14)), the number of iterations needed is the same, which is denoted by  $K_{\text{cons}}$ . Furthermore, we have

$$(5.17) \quad K_{\text{cons}} = O(\sqrt{L_D/\varepsilon}) = O(M/\sqrt{\varepsilon\mu}).$$

**5.3. Oracle Complexity.** From Section 5.1 (specifically, (5.5)), we notice that in Algorithms 1 and 2, the dual sub-problems can be solved exactly. Therefore, we focus on analyzing their primal oracle complexities, where the sub-routines  $\mathbf{N}_2$  and  $\mathbf{M}_2$  remain the same as the ones in Sections 3.2 and 4.1, respectively. Note that based on our oracle model in Section 1.1, in the case of (5.2), the primal oracle  $\mathcal{O}^P$  returns  $\nabla f(x)$  with input  $(x, 0)$  and  $\lambda_i \nabla g_i(x)$  with input  $(x, \lambda, i)$ .

Compared to the complexity analyses in Sections 3.4, 4.3 and 4.4, the challenge here is that  $L_{xx}(\lambda)$  now depends on  $\lambda$ . This implies that in Algorithm 1 or 2, as  $\hat{\lambda}^k$  changes over iterations,  $L_{xx}(\hat{\lambda}^k)$  also changes. Although this does not affect the iteration complexity of Algorithm 1 or 2 (since  $L_D$  depends only on  $L_{\lambda x}$ ), it does affect the oracle complexity of solving the primal sub-problem at each iteration  $k$ . To overcome this challenge, we propose to bound  $\|\hat{\lambda}^k\|_\infty$  for each  $k \in \mathbb{N}$  (in either the deterministic or stochastic sense).

LEMMA 5.5. *In Algorithm 1, if we choose the input parameters  $\rho_0$ ,  $\{\tau_k\}_{k \in \mathbb{Z}_+}$ ,  $\{\gamma_k\}_{k \in \mathbb{Z}_+}$  and  $\{\eta_k\}_{k \in \mathbb{Z}_+}$  in the same way as in Theorem 5.3, then for any  $k \in \mathbb{N}$ ,*

$$(5.18) \quad \|\hat{\lambda}^k\|_\infty = O(1 + k\sqrt{\varepsilon\mu}/M).$$

*Proof.* From (5.5), we have that  $\tilde{\lambda}_{\rho_k, 0}(x^k) = ([g_i(x^k)]_+ / \rho_k)_{i=1}^n$ . By the bound on  $[g_i(x^k)]_+$  in (5.12) and the expression of  $\rho_k$  in (3.24), we have

$$(5.19) \quad \|\tilde{\lambda}_{\rho_k, 0}(x^k)\|_\infty = O(1 + k\sqrt{\varepsilon\mu}/M).$$

By step 7 and the choice of  $\{\tau_k\}_{k \in \mathbb{Z}_+}$  in Theorem 5.3, we have

$$(5.20) \quad \|\lambda^k\|_\infty \leq \frac{k}{k+2} \|\lambda^{k-1}\|_\infty + \frac{2}{k+2} \|\tilde{\lambda}_{\rho_k, 0}(x^k)\|_\infty.$$

Based on (5.20), we use Lemma 3.4 to conclude that

$$(5.21) \quad \|\lambda^K\|_\infty \leq \frac{2}{(K+1)(K+2)} \left( \|\lambda^0\|_\infty + \sum_{k=0}^{K-1} (k+1) \|\tilde{\lambda}_{\rho_k, 0}(x^k)\|_\infty \right) \\ \stackrel{(a)}{=} O(1 + K\sqrt{\varepsilon\mu}/M),$$

where in (a) we use (5.19). Finally, by step 2, we have

$$(5.22) \quad \|\hat{\lambda}^k\|_\infty \leq \tau_k \|\lambda^k\|_\infty + (1 - \tau_k) \|\tilde{\lambda}_{\rho_k, 0}(x^k)\|_\infty.$$

We then substitute (5.19) and (5.21) into (5.22), and obtain (5.18).  $\square$

Based on Lemma 5.5, we can bound  $L_{xx}(\hat{\lambda}^k)$  via

$$(5.23) \quad L_{xx}(\hat{\lambda}^k) \leq \alpha \|\hat{\lambda}^k\|_\infty = O(\alpha + k\alpha\sqrt{\varepsilon\mu}/M), \quad \text{where } \alpha \triangleq \sum_{i=1}^n \alpha_i.$$

Based on this bound, the oracle complexity of  $\mathbf{N}_2$  for (approximately) solving the sub-problem in (2.4) (cf. (3.11)), and the iteration complexity of Algorithm 1, i.e.,  $K_{\text{cons}}$  in (5.17), we can derive the following result.

**THEOREM 5.6.** *In Algorithm 1, for any starting point  $(x^0, \lambda^0) \in \mathcal{X} \times \mathbb{R}_+$ , denote  $\bar{C}_{\text{det}}$  as the oracle complexity to obtain an  $\varepsilon$ -optimal and  $\varepsilon$ -feasible solution. Then*

$$(5.24) \quad \bar{C}_{\text{det}} = O\left(\frac{nM}{\sqrt{\mu\varepsilon}} \sqrt{(L+\alpha)/\mu} \log\left(\frac{L+\alpha}{\varepsilon}\right)\right).$$

*Proof.* Similar to the analysis in Section 3.4, we have

$$\begin{aligned} \bar{C}_{\text{det}} &\stackrel{(a)}{=} O\left(n \sum_{k=1}^{K_{\text{cons}}} \sqrt{(L+\alpha)/\mu + k\alpha\sqrt{\varepsilon}/(M\sqrt{\mu})} \log(k((L+\alpha)/\varepsilon + k\alpha\sqrt{\mu}/(M\sqrt{\varepsilon})))\right) \\ &\stackrel{(b)}{=} O\left(n \sum_{k=1}^{K_{\text{cons}}} (\sqrt{(L+\alpha)/\mu} + \sqrt{k\alpha/M}(\varepsilon/\mu)^{1/4}) (\log k + \log((L+\alpha)\mu/(M\varepsilon)))\right) \\ &\stackrel{(c)}{=} O\left(n(\sqrt{(L+\alpha)/\mu} K_{\text{cons}} (\log K_{\text{cons}} + \log((L+\alpha)\mu/(M\varepsilon))) \right. \\ &\quad \left. + \sqrt{\alpha/M}(\varepsilon/\mu)^{1/4} K_{\text{cons}}^{3/2} (\log K_{\text{cons}} + \log((L+\alpha)\mu/(M\varepsilon)))\right) \\ &= O\left(n(M\sqrt{L+\alpha}/(\mu\sqrt{\varepsilon}) + M\sqrt{\alpha}/(\mu\sqrt{\varepsilon})) \log((L+\alpha)/\varepsilon)\right), \end{aligned}$$

where in (a) we use  $\gamma_k = \Theta(\varepsilon/k)$ , in (b) we use  $\alpha/\sqrt{\varepsilon} = O((L+\alpha)/\varepsilon)$  and in (c) we use  $\sum_{k=1}^K k^\nu \log k = \Theta(K^{\nu+1} \log K)$ , for any  $\nu \geq 0$ . Finally, by noting that  $\alpha \leq L+\alpha$ , we obtain (5.24).  $\square$

Based on the oracle complexity of  $\mathbf{M}_2$  for (approximately) solving the sub-problem in (2.4) (cf. (4.13)), by using the same arguments as in Theorem 5.6, we can also derive the following oracle complexity for Algorithm 2.

**THEOREM 5.7.** *In Algorithm 2, for any starting point  $(x^0, \lambda^0) \in \mathcal{X} \times \mathbb{R}_+$ , denote  $\bar{C}_{\text{stoc}}$  as the oracle complexity to obtain an  $\varepsilon$ -optimal and  $\varepsilon$ -feasible solution in expectation. Then we have*

$$(5.25) \quad \bar{C}_{\text{stoc}} = O\left(\frac{\sqrt{n}M}{\sqrt{\mu\varepsilon}} \left(\sqrt{n} + \sqrt{(L+\alpha)/\mu}\right) \log\left(\frac{nM(L+\alpha)}{\mu\varepsilon}\right)\right).$$

Let us compare the complexity results in Theorems 5.6 and 5.7. If we interpret the factor  $\kappa_{\text{cons}} \triangleq (L+\alpha)/\mu$  as the ‘‘condition number’’ of the constrained problem in (5.1), and recall that  $K_{\text{cons}} = O(M/\sqrt{\mu\varepsilon})$  (cf. (5.17)), then the oracle complexity of Algorithm 1, i.e.,  $\tilde{O}(nK_{\text{cons}}\sqrt{\kappa_{\text{cons}}})$ , has been reduced to  $\tilde{O}(\sqrt{n}K_{\text{cons}}(\sqrt{n} + \sqrt{\kappa_{\text{cons}}}))$  in Algorithm 2. This is indeed consistent with our observation in Section 4.3, which concerns the oracle complexities of Algorithms 1 and 2 for the SPP in (1.1).

**5.4. Related Works.** Although numerous first-order methods have been proposed for solving the constrained problem in (5.1) when  $f$  is non-strongly convex (i.e.,  $\mu = 0$ ), it appears that there exist few methods that tackle the case where  $f$  is strongly convex (i.e.,  $\mu > 0$ ). Among these, the best-known oracle complexity is  $\tilde{O}(\varepsilon^{-1/2})$ , which has been achieved by three methods, which include the inexact ALM [31], the inexact dual gradient method [18] and the level-set method [17].<sup>1</sup> From Theorems 5.6 and 5.7, we observe that this complexity is also achieved by Algorithms 1 and 2.

<sup>1</sup>Note that the method in [17] requires to start with a strictly feasible solution of (5.1). However, obtaining such a solution ‘‘may be as costly as computing an optimal solution’’ [1]. In contrast, this is not required in any of the other methods discussed in this section (including ours).

Although all of these methods achieve the same complexity, there are two important features that distinguish our methods from the rest. First, our randomized framework (i.e., Algorithm 2) can effectively handle the case where the number of constraints  $n$  is extremely large (cf. Theorem 5.7). Second, both of our frameworks (i.e., Algorithms 1 and 2) are developed for solving the *general* SPP in (1.1), not only the Lagrangian problem in (5.2). Therefore, they have much wider applicability compared to the other methods.

### 5.5. Stochastic Convex Optimization with Expectation Constraints.

We consider a problem related to (5.1), where  $f, g_1, \dots, g_n$  are given in terms of expectation (see e.g., [16]). Specifically, let  $\xi, \zeta_1, \dots, \zeta_n$  be random vectors with sample spaces denoted by  $\Xi, \mathcal{Z}_1, \dots, \mathcal{Z}_n$ , respectively, and  $F : \mathbb{E}_1 \times \Xi \rightarrow \mathbb{R}$  and  $\{G_i : \mathbb{E}_1 \times \mathcal{Z}_i \rightarrow \mathbb{R}\}_{i=1}^n$  be chosen such that

$$(5.26) \quad f(x) \triangleq \mathbb{E}_\xi[F(x, \xi)], \quad g_i(x) \triangleq \mathbb{E}_{\zeta_i}[G_i(x, \zeta_i)], \quad \forall i \in [n],$$

and  $f, g_1, \dots, g_n$  satisfy the structural assumptions stated in Section 5.1. Since the expectations in (5.26) cannot be evaluated with high accuracy in general, from the stochastic programming literature (e.g., [27]), there are two standard approaches to tackle this difficulty, i.e., SA and sample average approximation (SAA). The SA approach involves employing certain mechanisms to generate unbiased estimators of  $\nabla f$  or  $\nabla g_i$  at any  $x \in \mathcal{X}$ , whereas the SAA approach involves generating i.i.d. samples of  $\xi, \zeta_1, \dots, \zeta_n$  (denoted as  $\{\xi^j\}_{j=1}^{m_0}, \{\zeta_1^j\}_{j=1}^{m_1}, \dots, \{\zeta_n^j\}_{j=1}^{m_n}$ , respectively) to approximate the expectations in (5.26), i.e.,

$$(5.27) \quad \begin{aligned} \min_{x \in \mathcal{X}} \{ \hat{f}(x) \triangleq (1/m_0) \sum_{j=1}^{m_0} F(x, \xi_j) \}, \\ \text{s. t. } \{ \hat{g}_i(x) \triangleq (1/m_i) \sum_{j=1}^{m_i} G_i(x, \zeta_i^j) \} \leq 0, \quad \forall i \in [n]. \end{aligned}$$

Note that to achieve high approximation accuracy, the sample sizes  $\{m_i\}_{i=0}^n$  are typically very large. We illustrate that Algorithm 2 can solve (5.26), no matter whether SA or SAA is used. As in Section 1, we still apply Algorithm 2 to the Lagrangian problem associated with (5.26). The ability of Algorithm 2 to interface with SA can be readily seen from Section 4.5. For the case of SAA, we see that (5.27) has a Lagrangian problem that fits into the finite-sum structure in (1.5), which again can be efficiently solved by Algorithm 2 (cf. Theorem 5.7).

**Acknowledgments.** The authors would like to express sincere gratitude to Robert M. Freund and Yangyang Xu for helpful discussions and constructive feedback during the preparation of this manuscript.

### REFERENCES

- [1] Aleksandr Y. Aravkin, James V. Burke, Dmitriy Drusvyatskiy, Michael P. Friedlander, and Scott Roy. Level-set methods for convex optimization. arXiv:1602.01506, 2016.
- [2] A. Auslender and M. Teboulle. Interior gradient and proximal methods for convex and conic optimization. *SIAM J. Optim.*, 16(3):697–725, 2006.
- [3] P. Balamurugan and Francis Bach. Stochastic variance reduction methods for saddle-point problems. In *Proc. NIPS*, 2016.
- [4] Heinz H. Bauschke, Jonathan M. Borwein, and Patrick L. Combettes. Essential smoothness, essential strict convexity, and legendre functions in Banach spaces. *Commun. Contemp. Math.*, 3(4):615–647, 2001.
- [5] Dimitri P. Bertsekas. *Nonlinear Programming*. Athena Scitific, 1999.
- [6] Stephen Boyd and Lieven Vandenberghe. *Convex Optimization*. Cambridge University Press, 2004.

- [7] Antonin Chambolle and Thomas Pock. On the ergodic convergence rates of a first-order primal-dual algorithm. *Math. Program.*, 159(1):253–287, 2016.
- [8] Yunmei Chen, Guanghui Lan, and Yuyuan Ouyang. Optimal primal-dual methods for a class of saddle point problems. *SIAM J. Optim.*, 24(4):1779–1814, 2014.
- [9] Yunmei Chen, Guanghui Lan, and Yuyuan Ouyang. Accelerated schemes for a class of variational inequalities. *Math. Program.*, 165(1):113–149, 2017.
- [10] Saeed Ghadimi and Guanghui Lan. Optimal stochastic approximation algorithms for strongly convex stochastic composite optimization, ii: Shrinking procedures and optimal algorithms. *SIAM J. Optim.*, 23(4):2061–2089, 2013.
- [11] Erfan Yazdandoost Hamedani and Necdet Serhat Aybat. A primal-dual algorithm for general convex-concave saddle point problems. arXiv:1803.01401, 2018.
- [12] Anatoli Juditsky and Arkadi Nemirovski. *First-Order Methods for Nonsmooth Convex Large-Scale Optimization, II: Utilizing Problems Structure*, pages 149–184. MIT Press, 2012.
- [13] Anatoli Juditsky, Arkadi Nemirovski, and Claire Tauvel. Solving variational inequalities with stochastic mirror-prox algorithm. *Stoch. Syst.*, 1(1):17–58, 2011.
- [14] O. Kolossoski and R.D.C. Monteiro. An accelerated non-euclidean hybrid proximal extragradient-type algorithm for convexconcave saddle-point problems. *Optim. Methods Softw.*, 32(6):1244–1272, 2017.
- [15] Guanghui Lan and Yi Zhou. An optimal randomized incremental gradient method. *Math. Program.*, 171(1):167–215, 2018.
- [16] Guanghui Lan and Zhiqiang Zhou. Algorithms for stochastic optimization with expectation constraints. arXiv:1604.03887, 2016.
- [17] Q. Lin, S. Nadarajah, and N. Soheili. A level-set method for convex optimization with a feasible solution path. *SIAM J. Optim.*, 28(4):3290–3311, 2018.
- [18] I. Necoara and V. Nedelcu. Rate analysis of inexact dual first-order methods application to dual decomposition. *IEEE Trans. Autom. Control*, 59(5):1232–1243, 2014.
- [19] A. Nedić and A. Ozdaglar. Subgradient methods for saddle-point problems. *J. Optim. Theory Appl.*, 142(1):205–228, 2009.
- [20] Arkadi Nemirovski. Prox-method with rate of convergence  $O(1/t)$  for variational inequalities with lipschitz continuous monotone operators and smooth convex-concave saddle point problems. *SIAM J. Optim.*, 15(1):229–251, 2005.
- [21] Yu. Nesterov. Excessive gap technique in nonsmooth convex minimization. *SIAM J. Optim.*, 16(1):235–249, 2005.
- [22] Yu. Nesterov. Smooth minimization of non-smooth functions. *Math. Program.*, 103(1):127–152, 2005.
- [23] Yu. Nesterov. Primal-dual subgradient methods for convex problems. *Math. Program.*, 120(1):221–259, 2009.
- [24] Yu. Nesterov. Gradient methods for minimizing composite functions. *Math. Program.*, 140(1):125–161, 2013.
- [25] Andrei Patrascu, Ion Necoara, and Quoc Tran-Dinh. Adaptive inexact fast augmented lagrangian methods for constrained convex optimization. *Optim. Lett.*, 11(3):609–626, 2017.
- [26] Shai Shalev-Shwartz and Tong Zhang. Accelerated proximal stochastic dual coordinate ascent for regularized loss minimization. *Math. Program.*, 155(1):105–145, 2016.
- [27] Alexander Shapiro, Darinka Dentcheva, and Andrzej Ruszczyński. *Lectures on Stochastic Programming Modeling and Theory*. SIAM and MPS, 2009.
- [28] Marc Teboulle. A simplified view of first order methods for optimization. *Math. Program.*, 170(1):67–96, 2018.
- [29] Q. Tran-Dinh, I. Necoara, and M. Diehl. Fast inexact decomposition algorithms for large-scale separable convex optimization. *Optim.*, 65(2):325–356, 2016.
- [30] Linli Xu, James Neufeld, Bryce Larson, and Dale Schuurmans. Maximum margin clustering. In *Proc. NIPS*, pages 1537–1544, 2005.
- [31] Yangyang Xu. Iteration complexity of inexact augmented lagrangian methods for constrained convex programming. arXiv:1711.05812, 2017.
- [32] Renbo Zhao. Optimal algorithms for stochastic three-composite convex-concave saddle point problems. arXiv:1903.01687, 2019.