

1 **AN ACCELERATED INEXACT PROXIMAL POINT METHOD FOR**
2 **SOLVING NONCONVEX-CONCAVE MIN-MAX PROBLEMS**

3 WEIWEI KONG* AND RENATO D.C. MONTEIRO*

4 **Abstract.** This paper presents smoothing schemes for obtaining approximate stationary points
5 of unconstrained or linearly-constrained composite nonconvex-concave min-max (and hence non-
6 smooth) problems by applying well-known algorithms to composite smooth approximations of the
7 original problems. More specifically, in the unconstrained (resp. constrained) case, approximate
8 stationary points of the original problem are obtained by applying, to its composite smooth approx-
9 imation, an accelerated inexact proximal point (resp. quadratic penalty) method presented in a
10 previous paper by the authors. Iteration complexity bounds for both smoothing schemes are also
11 established. Finally, numerical results are given to demonstrate the efficiency of the unconstrained
12 smoothing scheme.

13 **Key words.** quadratic penalty method, composite nonconvex problem, iteration-complexity,
14 inexact proximal point method, first-order accelerated gradient method, minimax problem.

15 **AMS subject classifications.** 47J22, 90C26, 90C30, 90C47, 90C60, 65K10.

16 **1. Introduction.** The first goal of this paper is to present and study the complex-
17 ity of an accelerated inexact proximal point smoothing (AIPP-S) scheme for find-
18 ing approximate stationary points of the (potentially nonsmooth) min-max compos-
19 ite nonconvex optimization (CNO) problem

20 (1.1)
$$\min_{x \in X} \{\hat{p}(x) := p(x) + h(x)\}$$

21 where h is a proper lower-semicontinuous convex function, X is a nonempty convex
22 set, and p is a max function given by

23 (1.2)
$$p(x) := \max_{y \in Y} \Phi(x, y) \quad \forall x \in X,$$

24 for some nonempty compact convex set Y and function Φ which, for some scalar
25 $m > 0$ and open set $\Omega \supseteq X$, is such that: (i) Φ is continuous on $\Omega \times Y$; (ii) the
26 function $-\Phi(x, \cdot) : Y \mapsto \mathbb{R}$ is lower-semicontinuous and convex for every $x \in X$; and
27 (ii) for every $y \in Y$, the function $\Phi(\cdot, y) + m\|\cdot\|^2/2$ is convex, differentiable, and its
28 gradient is Lipschitz continuous on $X \times Y$. Here, the objective function is the sum of a
29 convex function h and the pointwise supremum of (possibly nonconvex) differentiable
30 functions which is generally a (possibly nonconvex) nonsmooth function.

31 When Y is a singleton, the max term in (1.1) becomes smooth and (1.1) reduces
32 to a smooth CNO problem for which many algorithms have been developed in the
33 literature. In particular, accelerated inexact proximal points (AIPP) methods, i.e.
34 methods which use an accelerated composite gradient variant to approximately solve
35 the generated sequence of prox subproblems, have been developed for it (see, for
36 example, [4, 15]). When Y is not a singleton, (1.1) can no longer be directly solved by
37 an AIPP method due to the nonsmoothness of the max term. The AIPP-S scheme
38 developed in this paper is instead based on a perturbed version of (1.1) in which the
39 max term in (1.1) is replaced by a smooth approximation and the resulting smooth
40 CNO problem is solved by an AIPP method.

*School of Industrial and Systems Engineering, Georgia Institute of Technology, Atlanta, GA,
30332-0205. (E-mails: wkong37@gatech.edu & monteiro@isye.gatech.edu). The works of these
authors were partially supported by ONR Grant N00014-18-1-2077 and NSERC Grant PGSD3-
516700-2018.

41 Throughout our presentation, it is assumed that efficient oracles for evaluating
 42 the quantities $\Phi(x, y)$, $\nabla_x \Phi(x, y)$, and $h(x)$ and for obtaining exact solutions of the
 43 problems

$$44 \quad (1.3) \quad \min_{x \in X} \left\{ \lambda h(x) + \frac{1}{2} \|x - x_0\|^2 \right\}, \quad \max_{y \in Y} \left\{ \lambda \Phi(x_0, y) - \frac{1}{2} \|y - y_0\|^2 \right\}$$

45 for any (x_0, y_0) and $\lambda > 0$, are available. Throughout this paper, the terminology
 46 “oracle call” is used to refer to a collection of the above oracles of size $\mathcal{O}(1)$ where
 47 each of them appears at least once. We refer to the computation of the solution of
 48 the first problem above as a h -resolvent evaluation. In this manner, the computation
 49 of the solution of the second one is a $[-\Phi(x_0, \cdot)]$ -resolvent evaluation.

50 We first develop an AIPP-S scheme that obtains a stationary point based on a
 51 primal-dual formulation of (1.1). More specifically, given a tolerance pair $(\rho_x, \rho_y) \in$
 52 \mathbb{R}_{++}^2 , it is shown that an instance of this scheme obtains a quadruple $(\bar{u}, \bar{v}, \bar{x}, \bar{y})$ such
 53 that

$$54 \quad (1.4) \quad \begin{pmatrix} \bar{u} \\ \bar{v} \end{pmatrix} \in \begin{pmatrix} \nabla_x \Phi(\bar{x}, \bar{y}) \\ 0 \end{pmatrix} + \begin{pmatrix} \partial h(\bar{x}) \\ \partial [-\Phi(\bar{x}, \cdot)](\bar{y}) \end{pmatrix}, \quad \|\bar{u}\| \leq \rho_x, \quad \|\bar{v}\| \leq \rho_y$$

55 in $\mathcal{O}(\rho_x^{-2} \rho_y^{-1/2})$ oracle calls, where $\partial \phi(z)$ is the subdifferential of a convex function ϕ
 56 at a point z (see (1.9) with $\varepsilon = 0$). We then show that another instance of this scheme
 57 can obtain an approximate stationary point based on the directional derivative of \hat{p} .
 58 More specifically, given a tolerance pair $\delta > 0$, it is shown that this instance computes
 59 a point $x \in X$ such that

$$60 \quad (1.5) \quad \exists \hat{x} \in X \text{ s.t. } \inf_{\|d\| \leq 1} \hat{p}'(\hat{x}; d) \geq -\delta, \quad \|\hat{x} - x\| \leq \delta,$$

61 in $\mathcal{O}(\delta^{-3})$ oracle calls, where $\hat{p}'(x; d)$ is the directional derivative of \hat{p} at the point x
 62 along the direction d (see (1.10)).

63 The second goal of this paper is to develop a quadratic penalty AIPP-S (QP-
 64 AIPP-S) scheme to obtain approximate stationary points of a linearly constrained
 65 version of (1.1), namely

$$66 \quad (1.6) \quad \min_{x \in X} \{p(x) + h(x) : \mathcal{A}x = b\}$$

67 where p is as in (1.2), \mathcal{A} is a linear operator, and b is in the range of \mathcal{A} . The scheme is
 68 a penalty-type method which approximately solves a sequence of penalty subproblems
 69 of the form
 70 of the form

$$71 \quad (1.7) \quad \min_{x \in X} \left\{ p(x) + h(x) + \frac{c}{2} \|\mathcal{A}x - b\|^2 \right\}$$

72 for an increasing sequence of positive penalty parameters c . Similar to the approach
 73 used for the first goal of this paper, the method considers a perturbed variant of
 74 (1.7) in which the objective function is replaced by a smooth approximation and
 75 the resulting problem is solved by the quadratic-penalty AIPP (QP-AIPP) method
 76 proposed in [15]. For a given tolerance triple $(\rho_x, \rho_y, \eta) \in \mathbb{R}_{++}^3$, it is shown that the
 77 method computes a quintuple $(\bar{u}, \bar{v}, \bar{x}, \bar{y}, \bar{r})$ satisfying

$$78 \quad (1.8) \quad \begin{pmatrix} \bar{u} \\ \bar{v} \end{pmatrix} \in \begin{pmatrix} \nabla_x \Phi(\bar{x}, \bar{y}) + \mathcal{A}^* \bar{r} \\ 0 \end{pmatrix} + \begin{pmatrix} \partial h(\bar{x}) \\ \partial [-\Phi(\bar{x}, \cdot)](\bar{y}) \end{pmatrix},$$

$$\|\bar{u}\| \leq \rho_x, \quad \|\bar{v}\| \leq \rho_y, \quad \|\mathcal{A}\bar{x} - b\| \leq \eta.$$

79 in $\mathcal{O}(\rho_x^{-2}\rho_y^{-1/2} + \rho_x^{-2}\eta^{-1})$ oracle calls.

80 Finally, it is worth mentioning that all of the above complexities are obtained under the mild assumption that the optimal value in each of the respective optimization problems, namely (1.1) and (1.6) is bounded below. Moreover, it is neither assumed that X be bounded nor that (1.1) or (1.6) has an optimal solution.

84

85 *Related Works.* Since the case when $\Phi(\cdot, \cdot)$ in (1.1) is convex-concave has been well-studied in the literature (see, for example, [1, 11, 13, 21, 22, 23, 27]), we will make no more mention of it here. Instead, we will focus on papers that consider (1.1) where $\Phi(\cdot, y)$ is differentiable and nonconvex for every $y \in Y$ and there are mild conditions on $\Phi(x, \cdot)$ for every $x \in X$.

90 Letting δ_C denote the indicator function of a closed convex set $C \subseteq \mathcal{X}$ (see Subsection 1.1), $\overline{\text{Conv}}(\mathcal{X})$ denote the set of proper lower semicontinuous convex functions on \mathcal{X} , and $\rho := \min\{\rho_x, \rho_y\}$, Tables 1.1 and 1.2 compare the assumptions and iteration complexities obtained in this work with corresponding ones derived in the earlier papers [24, 26] and the subsequent works [17, 25, 30]. It is worth mentioning that the above works consider termination conditions that are slightly different than the ones in this paper. It is shown in Subsection 2.1 that they are actually equivalent to the ones in this paper up to multiplicative constants that are independent of the tolerances, i.e., ρ_x, ρ_y, δ .

Algorithm	Oracle Complexity	Use Cases			
		$D_h = \infty$	$h \equiv 0$	$h \equiv \delta_C$	$h \in \overline{\text{Conv}}(\mathcal{X})$
PGSF [24]	$\mathcal{O}(\rho^{-3})$	✗	✓	✓	✗
Minimax-PPA [17]	$\mathcal{O}(\rho^{-2.5} \log^2(\rho^{-1}))$	✗	✓	✓	✗
FNE Search [25]	$\mathcal{O}(\rho_x^{-2}\rho_y^{-1/2} \log(\rho^{-1}))$	✓	✓	✓	✗
AIPP-S	$\mathcal{O}(\rho_x^{-2}\rho_y^{-1/2})$	✓	✓	✓	✓

TABLE 1.1

Comparison of iteration complexities and assumptions under notions equivalent to (1.4) with $\rho := \min\{\rho_x, \rho_y\}$.

Algorithm	Oracle Complexity	Use Cases			
		$D_h = \infty$	$h \equiv 0$	$h \equiv \delta_C$	$h \in \overline{\text{Conv}}(\mathcal{X})$
PG-SVRG [26]	$\mathcal{O}(\delta^{-6} \log \delta^{-1})$	✗	✓	✓	✓
Minimax-PPA [17]	$\mathcal{O}(\delta^{-3} \log^2(\delta^{-1}))$	✗	✓	✓	✗
Prox-DIAG [30]	$\mathcal{O}(\delta^{-3} \log^2(\delta^{-1}))$	✓	✓	✗	✗
AIPP-S	$\mathcal{O}(\delta^{-3})$	✓	✓	✓	✓

TABLE 1.2

Comparison of iteration complexities and assumptions under notions equivalent to (1.5).

99 To the best of our knowledge, this work is the first one to analyze the complexity of a smoothing scheme for finding approximate stationary points of (1.6).

101

102 *Organization of the paper.* Subsection 1.1 presents notation and some basic definitions that are used in this paper. Subsection 1.2 presents several motivating applications that are of the form in (1.1). Section 2 is divided into two subsections. The first one precisely states the assumptions underlying problem (1.1) and discusses four notions of stationary points. The second one presents a smooth approximation of the

107 function p in (1.1). Section 3 is divided into two subsections. The first one reviews
 108 the AIPP method in [15] and its iteration complexity. The second one presents the
 109 AIPP-S scheme its iteration complexities for finding approximate stationary points
 110 as in (1.4) and (1.5). Section 4 is also divided into two subsections. The first one
 111 reviews the QP-AIPP method in [15] and its iteration complexity. The second one
 112 presents the QP-AIPP-S scheme its iteration complexity for finding points satisfying
 113 (1.8). Section 5 presents some computational results. Section 6 gives some conclud-
 114 ing remarks. Finally, several appendices at the end of this paper contain proofs of
 115 technical results needed in our presentation.

116 **1.1. Notation and basic definitions.** This subsection provides some basic
 117 notation and definitions.

118 The set of real numbers is denoted by \mathbb{R} . The set of non-negative real numbers
 119 and the set of positive real numbers is denoted by \mathbb{R}_+ and \mathbb{R}_{++} respectively. The
 120 set of natural numbers is denoted by \mathbb{N} . For $t > 0$, define $\log_1^+(t) := \max\{1, \log(t)\}$.
 121 Let \mathbb{R}^n denote a real-valued n -dimensional Euclidean space with standard norm $\|\cdot\|$.
 122 Given a linear operator $A : \mathbb{R}^n \mapsto \mathbb{R}^p$, the operator norm of A is denoted by $\|A\| :=$
 123 $\sup\{\|Az\|/\|z\| : z \in \mathbb{R}^n, z \neq 0\}$.

124 The following notation and definitions are for a general complete inner product
 125 space \mathcal{Z} , whose inner product and its associated induced norm are denoted by $\langle \cdot, \cdot \rangle$
 126 and $\|\cdot\|$ respectively. Let $\psi : \mathcal{Z} \mapsto (-\infty, \infty]$ be given. The effective domain of ψ is
 127 denoted as $\text{dom } \psi := \{z \in \mathcal{Z} : \psi(z) < \infty\}$ and ψ is said to be proper if $\text{dom } \psi \neq \emptyset$.
 128 The set of proper, lower semi-continuous, convex functions $\psi : \mathcal{Z} \mapsto (-\infty, \infty]$ is
 129 denoted by $\overline{\text{Conv}}(\mathcal{Z})$. Moreover, for a convex set $Z \subseteq \mathcal{Z}$, we denote $\overline{\text{Conv}}(Z)$ to be
 130 set of functions in $\overline{\text{Conv}}(\mathcal{Z})$ whose effective domain is equal to Z . For $\varepsilon \geq 0$, the
 131 ε -subdifferential of $\psi \in \overline{\text{Conv}}(\mathcal{Z})$ at $z \in \text{dom } \psi$ is denoted by

$$132 \quad (1.9) \quad \partial_\varepsilon \psi(z) := \{w \in \mathbb{R}^n : \psi(z') \geq \psi(z) + \langle w, z' - z \rangle - \varepsilon, \forall z' \in \mathcal{Z}\},$$

133 and we denote $\partial\psi \equiv \partial_0\psi$. The *directional derivative* of ψ at $z \in \mathcal{Z}$ in the direction
 134 $d \in \mathcal{Z}$ is denoted by

$$135 \quad (1.10) \quad \psi'(z; d) := \lim_{t \rightarrow 0} \frac{\psi(z + td) - \psi(z)}{t}.$$

136 It is well-known that if ψ is differentiable at $z \in \text{dom } \psi$, then for a given direction
 137 $d \in \mathcal{Z}$ we have $\psi'(z; d) = \langle \nabla\psi(z), d \rangle$.

138 For a given $Z \subseteq \mathcal{Z}$, the indicator function of Z , denoted by δ_Z , is defined as
 139 $\delta_Z(z) = 0$ if $z \in Z$ and $\delta_Z(z) = \infty$ if $z \notin Z$. Moreover, the closure, interior, and
 140 relative interior of Z are denoted by $\text{cl } Z$, $\text{int } Z$, and $\text{ri } Z$, respectively. The support
 141 function of Z at a point z is denoted by $\sigma_Z(z) := \sup_{z' \in Z} \langle z, z' \rangle$.

142 **1.2. Motivating applications.** This subsection lists motivating applications
 143 that are of the form in (1.1). In Section 5, we examine the performance of our
 144 proposed smoothing scheme on some special instances of these applications.

145 **1.2.1. Maximum of a finite number of nonconvex functions.** Given a
 146 family of functions $\{f_i\}_{i=1}^k$ that are continuously differentiable everywhere with Lip-
 147 schitz continuous gradients and a closed convex set $C \subseteq \mathbb{R}^n$. The problem of interest
 148 is the minimization of $\max_{1 \leq i \leq k} f_i$ over the set C , i.e.,

$$149 \quad \min_{x \in C} \max_{1 \leq i \leq k} f_i(x),$$

150 which is clearly an instance of (1.1) where $Y = \{y \in \mathbb{R}_+^k : \sum_{i=1}^k y_i = 1\}$, $\Phi(x, y) =$
 151 $\sum_{i=1}^k y_i f_i(x)$, and $h(x) = \delta_C(x)$.

152 **1.2.2. Robust regression.** Given a set of observations $\sigma := \{\sigma_i\}_{i=1}^n$ and a
 153 compact convex set $\Theta \in \mathbb{R}^k$, let $\{\ell_\theta(\cdot|\sigma)\}_{\theta \in \Theta}$ be a family of nonconvex loss functions
 154 in which: (i) $\ell_\theta(x|\sigma)$ is concave in θ for every $x \in \mathbb{R}^n$; and (ii) $\ell_\theta(x|\sigma)$ is continuously
 155 differentiable in x with Lipschitz continuous gradient for every $\theta \in \Theta$. The problem
 156 of interest is to minimize the worst-case loss in Θ , i.e.,

$$157 \quad \min_{x \in \mathbb{R}^n} \max_{\theta \in \Theta} \ell_\theta(x|\sigma),$$

158 which is clearly an instance of (1.1), where $Y = \Theta$, $\Phi(x, y) = \ell_y(x|\sigma)$, and $h(x) = 0$.

159 **1.2.3. Min-max games with an adversary.** Let $\{\mathcal{U}_j(x_1, \dots, x_k, y)\}_{j=1}^k$ be a
 160 set of utility functions in which: (i) \mathcal{U}_j is nonconvex and continuously differentiable
 161 in its first k arguments, but concave in its last argument; (ii) $\nabla_{x_i} \mathcal{U}_j(x_1, \dots, x_k, y)$ is
 162 Lipschitz continuous for every $1 \leq i \leq k$. Given input constraint sets $\{B_i\}_{i=1}^k$ and
 163 B_y , the problem of interest is to maximize the total utility of the players (indices 1
 164 to k) given that the adversary (index $k+1$) seeks to maximize his own utility, i.e.,

$$165 \quad \min_{x_1, \dots, x_k} \max_y \left\{ - \sum_{i=1}^k \mathcal{U}_j(x_1, \dots, x_k, y) : x_i \in B_i, i = 0, \dots, k \right\},$$

167 which is clearly an instance of (1.1) where $x = (x_1, \dots, x_k)$, $Y = B_y$, $\Phi(x, y) =$
 168 $-\sum_{i=1}^k \mathcal{U}_j(x_1, \dots, x_k, y)$, and $h(x) = \delta_{B_1 \times \dots \times B_k}(x)$.

169 **2. Preliminaries.** This section presents some preliminary material and is di-
 170 vided into two subsections. The first one precisely describes the assumptions and
 171 various notions of stationary points for problem (1.1) and briefly compares two ap-
 172 proaches for obtaining them. The second one presents a smooth approximation of the
 173 max function p in (1.1) and some of its properties.

174 **2.1. Assumptions and notions of stationary points.** This subsection de-
 175 scribes the assumptions and four notions of stationary points for for problem (1.1).
 176 It is worth mentioning that the complexities of the smoothing scheme of Section 3
 177 are presented with respect to two of these notions. In order to understand how these
 178 results can be translated to the other two alternative notions, which have been used
 179 in a few papers dealing with problem (1.1), we present a few results discussing some
 180 useful relations between all these notions.

181 Throughout our presentation, we let \mathcal{X} and \mathcal{Y} be finite dimensional inner product
 182 spaces. We also make the following assumptions on problem (1.1):

- 183 (A0) $X \subset \mathcal{X}$ and $Y \subset \mathcal{Y}$ are nonempty convex sets, and Y is also compact;
- 184 (A1) there exists an open set $\Omega \supseteq X$ such that $\Phi(\cdot, \cdot)$ is finite and continuous on
 185 $\Omega \times Y$; moreover, $\nabla_x \Phi(x, y)$ exists and is continuous at every $(x, y) \in \Omega \times Y$;
- 186 (A2) $h \in \text{Conv}(X)$ and $-\Phi(x, \cdot) \in \text{Conv}(Y)$ for every $x \in \Omega$;
- 187 (A3) there exist scalars $(L_x, L_y) \in \mathbb{R}_{++}^2$, and $m \in (0, L_x]$ such that

$$188 \quad (2.1) \quad \Phi(x, y) - [\Phi(x', y) + \langle \nabla_x \Phi(x', y), x - x' \rangle] \geq -\frac{m}{2} \|x - x'\|^2,$$

$$189 \quad (2.2) \quad \|\nabla_x \Phi(x, y) - \nabla_x \Phi(x', y')\| \leq L_x \|x - x'\| + L_y \|y - y'\|,$$

191 for every $x, x' \in X$ and $y, y' \in Y$;

192 (A4) $\hat{p}_* := \inf_{x \in X} \hat{p}(x)$ is finite, where \hat{p} is as in (1.1);

193 We make three remarks about the above assumptions. First, it is well-known that
 194 condition (2.2) implies that

$$195 \quad (2.3) \quad \Phi(x', y) - [\Phi(x, y) + \langle \nabla_x \Phi(x, y), x' - x \rangle] \leq \frac{L_x}{2} \|x' - x\|^2,$$

196 for every $(x', x, y) \in X \times X \times Y$. Second, functions satisfying the lower curvature
 197 condition in (2.1) are often referred to as weakly convex functions (see, for example,
 198 [5, 6, 7, 8]). Third, the aforementioned weak convexity condition implies that, for any
 199 $y \in Y$, the function $\Phi(\cdot, y) + m\|\cdot\|^2/2$ is convex, and hence $p + m\|\cdot\|^2/2$ is as well.
 200 Note that while \hat{p} is generally nonconvex and nonsmooth, it has the nice property
 201 that $\hat{p} + m\|\cdot\|^2/2$ is convex.

202 We now discuss two stationarity conditions of (1.1) under assumptions (A0)–(A3).
 203 First, denoting

$$204 \quad (2.4) \quad \hat{\Phi}(x, y) := \Phi(x, y) + h(x) \quad \forall (x, y) \in X \times Y,$$

205 it is well-known that (1.1) is related to the saddle-point problem which consists of
 206 finding a pair $(x^*, y^*) \in X \times Y$ such that

$$207 \quad (2.5) \quad \hat{\Phi}(x^*, y) \leq \hat{\Phi}(x^*, y^*) \leq \hat{\Phi}(x, y^*),$$

208 for every $(x, y) \in X \times Y$. More specifically, (x^*, y^*) satisfies (2.5) if and only if x^*
 209 is an optimal solution of (1.1), y^* is an optimal solution of the dual of (1.1), and
 210 there is no duality gap between the two problems. Using the composite structure
 211 described above for $\hat{\Phi}$, it can be shown that a necessary condition for (2.5) to hold is
 212 that (x^*, y^*) satisfy the stationarity condition

$$213 \quad (2.6) \quad \begin{pmatrix} 0 \\ 0 \end{pmatrix} \in \begin{pmatrix} \nabla_x \Phi(x^*, y^*) \\ 0 \end{pmatrix} + \begin{pmatrix} \partial h(x^*) \\ \partial [-\Phi(x^*, \cdot)](y^*) \end{pmatrix}.$$

214 When $m = 0$, the above condition also becomes sufficient for (2.5) to hold. Second,
 215 it can be shown that $p'(x^*; d)$ is well-defined for every $d \in \mathcal{X}$ and that a necessary
 216 condition for $x^* \in X$ to be a local minimum of (1.1) is that it satisfies the stationarity
 217 condition

$$218 \quad (2.7) \quad \inf_{\|d\| \leq 1} \hat{p}'(x^*; d) \geq 0.$$

219 When $m = 0$, the above condition also becomes sufficient for x^* to be a global
 220 minimum of (1.1). Moreover, in view of Lemma 19 in Appendix D with $(\bar{u}, \bar{v}, \bar{x}, \bar{y}) =$
 221 $(0, 0, x^*, y^*)$, it follows that x^* satisfies (2.7) if and only if there exists $y^* \in Y$ such
 222 that (x^*, y^*) satisfies (2.6).

223 Note that finding points that satisfy (2.6) or (2.7) exactly is generally a difficult
 224 task. Hence, in this section and the next one, we only consider approximate versions
 225 of (2.6) or (2.7), which are (1.4) and (1.5), respectively. For ease of future reference,
 226 we say that:

- 227 (i) a quadruple $(\bar{u}, \bar{v}, \bar{x}, \bar{y})$ is a (ρ_x, ρ_y) -**primal-dual stationary point** of (1.1)
 228 if it satisfies (1.4);
- 229 (ii) a point \hat{x} is a δ -**directional stationary point** of (1.1) if it satisfies the first
 230 inequality in (1.5).

231 It is worth mentioning that (1.5) is generally hard to verify for a given point $x \in X$.
 232 This is primarily because the definition requires us to check an infinite number of
 233 directional derivatives for a (potentially) nonsmooth function at points \hat{x} near \bar{x} . In
 234 contrast, the definition of an approximate primal-dual stationary point is generally
 235 easier to verify because the quantities $\|\bar{u}\|$ and $\|\bar{v}\|$ can be measured directly, and the
 236 inclusions in (1.4) are easy to verify when the prox oracles for h and $\Phi(x, \cdot)$, for every
 237 $x \in X$, are readily available.

238 The next result, whose proof is given in Appendix D, shows that a (ρ_x, ρ_y) -primal-
 239 dual stationary point, for small enough ρ_x and ρ_y , yields a point x satisfying (1.5).
 240 Its statement makes use of the diameter of Y defined as

$$241 \quad (2.8) \quad D_y := \sup_{y, y' \in Y} \|y - y'\|.$$

242

243 **PROPOSITION 1.** *If the quadruple $(\bar{u}, \bar{v}, \bar{x}, \bar{y})$ is a (ρ_x, ρ_y) -primal-dual stationary*
 244 *point of (1.1), then there exists a point $\hat{x} \in X$ such that*

$$245 \quad \inf_{\|d\| \leq 1} \hat{p}'(\hat{x}; d) \geq -\rho_x - 2\sqrt{2mD_y\rho_y}, \quad \|\bar{x} - \hat{x}\| \leq \sqrt{\frac{2D_y\rho_y}{m}}.$$

246 The iteration complexities in this paper (see Section 3) are stated with respect to
 247 the two notions of stationary points (1.4) and (1.5). However, it is worth discussing
 248 below two other notions of stationary points that are common in the literature as well
 249 as some results that relate all four notions.

250 Given $(\lambda, \varepsilon) \in \mathbb{R}_{++}^2$, a point x is said to be a (λ, ε) -prox stationary point of (1.1)
 251 if the function $\hat{p} + \|\cdot\|^2/(2\lambda)$ is strongly convex and

$$252 \quad (2.9) \quad \frac{1}{\lambda}\|x - x_\lambda\| \leq \varepsilon, \quad x_\lambda = \operatorname{argmin}_{u \in \mathcal{X}} \left\{ \hat{P}_\lambda(u) := \hat{p}(u) + \frac{1}{2\lambda}\|u - x\|^2 \right\}.$$

253 The above notion is considered, for example, in [17, 26, 30]. The result below, whose
 254 proof is given in Appendix D, shows how it is related to (1.5).

255 **PROPOSITION 2.** *For any given $\lambda \in (0, 1/m)$, the following statements hold:*

256 (a) *for any $\varepsilon > 0$, if $x \in X$ satisfies (1.5) and*

$$257 \quad (2.10) \quad 0 < \delta \leq \frac{\lambda^3 \varepsilon}{\lambda^2 + 2(1 - \lambda m)(1 + \lambda)},$$

258 *then x is a (λ, ε) -prox stationary point;*

259 (b) *for any $\delta > 0$, if $x \in X$ is a (λ, ε) -prox stationary point for some $\varepsilon \leq$
 260 $\delta \cdot \min\{1, 1/\lambda\}$, then x satisfies (1.5) with $\hat{x} = x_\lambda$, where x_λ is as in (2.9).*

261 Note that for a fixed $\lambda \in (0, 1/m)$ such that $\max\{\lambda^{-1}, (1 - \lambda m)^{-1}\} = \mathcal{O}(1)$, the
 262 largest δ in part (a) is $\mathcal{O}(\varepsilon)$. Similarly, for part (b), if $\lambda^{-1} = \mathcal{O}(1)$ then largest ε in
 263 part (b) is $\mathcal{O}(\delta)$. Combining these two observations, it follows that (2.9) and (1.5)
 264 are equivalent (up to a multiplicative factor) under the assumption that $\delta = \Theta(\varepsilon)$.

265 Given $(\rho_x, \rho_y) \in \mathbb{R}_{++}^2$, a pair (\bar{x}, \bar{y}) is said to be a (ρ_x, ρ_y) -first-order Nash equi-
 266 librium point of (1.1) if

$$267 \quad (2.11) \quad \inf_{\|d_x\| \leq 1} \mathcal{S}'_{\bar{y}}(\bar{x}; d_x) \geq -\rho_x, \quad \sup_{\|d_y\| \leq 1} \mathcal{S}'_{\bar{x}}(\bar{y}; d_y) \leq \rho_y,$$

268 where $\mathcal{S}_{\bar{y}} := \Phi(\cdot, \bar{y}) + h(\cdot)$ and $\mathcal{S}_{\bar{x}} := \Phi(\bar{x}, \cdot)$. The above notion is considered, for
 269 example, in [17, 24, 25]. The next result, whose proof is given in Appendix D, shows
 270 that (2.11) is equivalent to (1.4).

271 PROPOSITION 3. A pair (\bar{x}, \bar{y}) is a (ρ_x, ρ_y) -first-order Nash equilibrium point if
 272 and only if there exists $(\bar{u}, \bar{v}) \in \mathcal{X} \times \mathcal{Y}$ such that $(\bar{u}, \bar{v}, \bar{x}, \bar{y})$ satisfies (1.4).

273 We now end this subsection by briefly discussing some approaches for finding
 274 approximate stationary points of (1.1). One approach is to apply a proximal descent
 275 type method directly to problem (1.1), but this would lead to subproblems with
 276 nonsmooth convex composite functions. A second approach is based on first applying
 277 a smoothing method to (1.1) and then using a prox-convexifying descent method such
 278 as the one in [15] to solve the perturbed unconstrained smooth problem. An advantage
 279 of the second approach, which is the one pursued in this paper, is that it generates
 280 subproblems with smooth convex composite objective functions. The next subsection
 281 describes one possible way to smooth the (generally) nonsmooth function p in (1.1).

282 **2.2. Smooth approximation.** This subsection presents a smooth approxima-
 283 tion of the function p in (1.1).

284 For every $\xi > 0$, consider the smoothed function p_ξ defined by

$$285 \quad (2.12) \quad p_\xi(x) := \max_{y \in Y} \left\{ \Phi_\xi(x, y) := \Phi(x, y) - \frac{1}{2\xi} \|y - y_0\|^2 \right\} \quad \forall x \in X,$$

287 for some $y_0 \in Y$. The following proposition presents the key properties of p_ξ and its
 288 related quantities.

289 PROPOSITION 4. Let $\xi > 0$ be given and assume that the function Φ satisfies
 290 conditions (A0)–(A3). Let $p_\xi(\cdot)$ and $\Phi_\xi(\cdot, \cdot)$ be as defined in (2.12) and define

$$291 \quad (2.13) \quad Q_\xi := \xi L_y + \sqrt{\xi(L_x + m)}, \quad L_\xi := L_y Q_\xi + L_x \leq \left(L_y \sqrt{\xi} + \sqrt{L_x} \right)^2,$$

$$y_\xi(x) := \operatorname{argmax}_{y' \in Y} \Phi_\xi(x, y'),$$

292 for every $x \in X$. Then, the following properties hold:

- 293 (a) $y_\xi(\cdot)$ is Q_ξ -Lipschitz continuous on X ;
- 294 (b) $p_\xi(\cdot)$ is continuously differentiable on X and $\nabla p_\xi(x) = \nabla_x \Phi(x, y_\xi(x))$ for
 295 every $x \in X$;
- 296 (c) $\nabla p_\xi(\cdot)$ is L_ξ -Lipschitz continuous on X ;
- 297 (d) for every $x, x' \in X$, we have

$$298 \quad (2.14) \quad p_\xi(x) - [p_\xi(x') + \langle \nabla p_\xi(x'), x - x' \rangle] \geq -\frac{m}{2} \|x - x'\|^2;$$

299 *Proof.* Note that the inequality in (2.13) follows from (a), the fact that $m \leq L_x$,
 300 and the bound

$$301 \quad L_\xi = L_y \left[\xi L_y + \sqrt{\xi(L_x + m)} \right] + L_x \leq \xi L_y^2 + 2\sqrt{\xi L_x} + L_x = \left(L_y \sqrt{\xi} + \sqrt{L_x} \right)^2.$$

302 The other conclusions of (a)–(c) follow from Lemma 13 and Proposition 14 in Appen-
 303 dix B with $(\Psi, q, y) = (\Phi_\xi, p_\xi, y_\xi)$. We now show that the conclusion of (d) is true.
 304 Indeed, if we consider (2.1) at $(y, x') = (y_\xi(x'), x')$, the definition of Φ_ξ , and use the
 305 definition of ∇p_ξ in (b), then

$$306 \quad -\frac{m}{2} \|x - x'\|^2 \leq \Phi(x', y_\xi(x)) - [\Phi(x, y_\xi(x)) + \langle \nabla_x \Phi(x, y_\xi(x)), x' - x \rangle]$$

$$307 \quad = \Phi_\xi(x', y_\xi(x)) - [p_\xi(x) + \langle \nabla p_\xi(x), x' - x \rangle] \leq p_\xi(x') - [p_\xi(x) + \langle \nabla p_\xi(x), x' - x \rangle],$$

309 where the last inequality follows from the optimality of y . □

310 We now make two remarks about the above properties. First, the Lipschitz con-
 311 stants of y_ξ and ∇p_ξ depend on the value of ξ while the weak convexity constant m in
 312 (2.14) does not. Second, as $\xi \rightarrow \infty$, it holds that $p_\xi \rightarrow p$ pointwise and $Q_\xi, L_\xi \rightarrow \infty$.
 313 These remarks are made more precise in the next result.

314 **LEMMA 5.** *For every $\xi > 0$, it holds that $-\infty < p(x) - D_y^2/(2\xi) \leq p_\xi(x) \leq p(x)$
 315 for every $x \in X$, where D_y is as in (2.8).*

316 *Proof.* The fact that $p(x) > -\infty$ follows immediately from assumption (A4). To
 317 show the other bounds, observe that for every $y_0 \in Y$, we have

$$318 \quad \Phi(x, y) + h(x) \geq \Phi(x, y) - \frac{1}{2\xi} \|y - y_0\|^2 + h(x) \geq \Phi(x, y) - \frac{D_y^2}{2\xi} + h(x)$$

319 for every $(x, y) \in X \times Y$. Taking the supremum of the bounds over $y \in Y$ and using
 320 the definitions of p and p_ξ yields the remaining bounds. \square

321 **3. Unconstrained min-max optimization.** This section presents our pro-
 322 posed AIPP-S scheme for solving the min-max CNO problem (1.1) and is divided
 323 into two subsections. The first one reviews an AIPP method for solving smooth CNO
 324 problems. The second one presents the AIPP-S scheme and its iteration complexity
 325 for finding stationary points as in (1.4) and (1.5).

326 Before proceeding, we briefly outline the idea of the AIPP-S scheme. The main
 327 idea is to apply the AIPP method described in the next subsection to the smooth
 328 CNO problem

$$329 \quad (3.1) \quad \min_{x \in X} \{\hat{p}_\xi(x) := p_\xi(x) + h(x)\},$$

330 where p_ξ is as in (2.12) and ξ is a positive scalar that will depend on the tolerances
 331 in (1.4) and (1.5). The above smoothing approximation scheme is similar to the one
 332 used in [23]; the approximation function p_ξ used in both schemes is smooth, but the
 333 one here is nonconvex while the one in [23] is convex. Moreover, while [23] uses an
 334 ACG variant to approximately solve (3.1), the AIPP-S scheme uses the AIPP method
 335 discussed below for this purpose.

336 **3.1. AIPP method for smooth CNO problems.** This subsection describes
 337 the AIPP method studied in [15], and its corresponding iteration complexity result,
 338 for solving a class of smooth CNO problems.

339 We first describe the problem that the AIPP method is intended to solve. Let \mathcal{X}
 340 be a finite-dimensional inner product and consider the smooth CNO problem

$$341 \quad (3.2) \quad \phi_* := \inf_{x \in \mathcal{X}} [\phi(x) := f(x) + h(x)]$$

342 where $h : \mathcal{X} \mapsto (-\infty, \infty]$ and function f satisfy the following assumptions:

- 343 (P1) $h \in \overline{\text{Conv}}(\mathcal{X})$ and f is differentiable on $\text{dom } h$;
 344 (P2) for some $M \geq m > 0$, the function f satisfies

$$345 \quad (3.3) \quad -\frac{m}{2} \|x' - x\|^2 \leq f(x') - [f(x) + \langle \nabla f(x), x' - x \rangle],$$

$$346 \quad (3.4) \quad \|\nabla f(x') - \nabla f(x)\| \leq M \|x' - x\|,$$

348 for any $x, x' \in \text{dom } h$;

- 349 (P3) ϕ_* defined in (3.2) is finite.

350 We now make four remarks about the above assumptions. First, it is well-known
 351 that a necessary condition for $x^* \in \text{dom } h$ to be a local minimum of (3.2) is that x^*
 352 is a stationary point of ϕ , i.e. $0 \in \nabla f(x^*) + \partial h(x^*)$. Second, it is well-known that
 353 (3.4) implies that (3.3) holds for any $m \in [-M, M]$. Third, it is easy to see from
 354 Proposition 4 that p_ξ in (2.12) satisfies assumption (P2) with $(M, f) = (L_\xi, p_\xi)$ where
 355 L_ξ is as in (2.13). Fourth, it is also easy to see that the function p_ξ in (2.12) satisfies
 356 assumption (P3) with $\phi_* = \inf_{x \in X} \hat{p}_\xi(x)$ in view of assumption (A4) and Lemma 5.

357 For the purpose of discussing future complexity results, we consider the following
 358 notion of an approximate stationary point of (3.2): given a tolerance $\bar{\rho} > 0$, a pair
 359 $(\bar{x}, \bar{u}) \in \text{dom } h \times \mathcal{X}$ is said to be a $\bar{\rho}$ -approximate stationary point of (3.2) if

$$360 \quad (3.5) \quad \bar{u} \in \nabla f(\bar{x}) + \partial h(\bar{x}), \quad \|\bar{u}\| \leq \bar{\rho}.$$

361 We now state the AIPP method for finding a pair (\bar{x}, \bar{u}) satisfying (3.5).

362 AIPP method

363 **Input:** a function pair (f, h) , a scalar pair $(m, M) \in \mathbb{R}_{++}^2$ satisfying (P2), scalars
 364 $\lambda \in (0, 1/(2m)]$ and $\sigma \in (0, 1)$, an initial point $x_0 \in \text{dom } h$, and a tolerance $\bar{\rho} > 0$;

365 **Output:** a pair $(\bar{x}, \bar{u}) \in \text{dom } h \times \mathcal{X}$ satisfying (3.5);

- 366 (0) set $k = 1$ and define $\hat{\rho} := \bar{\rho}/4$, $\hat{\varepsilon} := \bar{\rho}^2/[32(M + \lambda^{-1})]$, and $M_\lambda := M + \lambda^{-1}$;
 367 (1) call the accelerated composite gradient (ACG) method in Appendix A with
 368 inputs $z_0 = x_{k-1}$, $(\mu, L) = (1/2, \lambda M + 1/2)$, $\psi_s = \lambda f + \|\cdot - x_{k-1}\|^2/4$, and
 369 $\psi_n = \lambda h + \|\cdot - x_{k-1}\|^2/4$ in order to obtain a triple $(x, u, \varepsilon) \in \mathcal{X} \times \mathcal{X} \times \mathbb{R}_+$
 370 satisfying

$$371 \quad (3.6) \quad u \in \partial_\varepsilon \left(\lambda \phi + \frac{1}{2} \|\cdot - x_{k-1}\|^2 \right) (x), \quad \|u\|^2 + 2\varepsilon \leq \sigma \|x_{k-1} - x + u\|^2;$$

- 372 (2) if $\|x_{k-1} - x + u\| \leq \lambda \hat{\rho}/5$, then go to (3); otherwise set $(x_k, \tilde{u}_k, \tilde{\varepsilon}_k) = (x, u, \varepsilon)$,
 373 increment $k = k + 1$ and go to (1);
 374 (3) restart the previous call to the ACG method in step 1 to find a triple $(\tilde{x}, \tilde{u}, \tilde{\varepsilon})$
 375 such that $\tilde{\varepsilon} \leq \hat{\varepsilon} \lambda$ and $(x, u, \varepsilon) = (\tilde{x}, \tilde{u}, \tilde{\varepsilon})$ satisfies (3.6);
 376 (4) compute

$$377 \quad (3.7) \quad \bar{x} := \operatorname{argmin}_{x' \in \mathcal{X}} \left\{ \langle \nabla f(x), x' - x \rangle + h(x') + \frac{M_\lambda}{2} \|x' - x\|^2 \right\},$$

$$378 \quad (3.8) \quad \bar{u} := M_\lambda(x - \bar{x}) + \nabla f(\bar{x}) - \nabla f(x),$$

380 where M_λ is as in step 0, and output the pair (\bar{x}, \bar{u}) .

381 We now make four remarks about the above AIPP method. First, at the k^{th}
 382 iteration of the method, its step 1 invokes an ACG method, whose description is given
 383 in Appendix A, to approximately solve the strongly convex proximal subproblem
 384

$$385 \quad (3.9) \quad \min_{x \in \mathcal{X}} \left\{ \lambda \phi(x) + \frac{1}{2} \|x - x_{k-1}\|^2 \right\}$$

386 according to (3.6). Second, Lemma 12 shows that every ACG iterate (z, u, ε) satisfies
 387 the inclusion in (3.6), and hence, only the inequality in (3.6) needs to be verified.
 388 Third, note that (3.4) implies that the gradient of the function ψ_s defined in step 1 of

389 the AIPP method is $(\lambda M + 1/2)$ -Lipschitz continuous. As a consequence, Lemma 12
 390 with $L = \lambda M + 1/2$ implies that the triple (z, u, ε) in step 1 of any iteration of the
 391 AIPP method can be obtained in $\mathcal{O}(\sqrt{[\lambda M + 1]/\sigma})$ ACG iterations.

392 Note that the above method differs slightly from the one presented in [15] in that
 393 it adds step 4 in order to directly output a $\bar{\rho}$ -approximate stationary point as in (3.5).
 394 The justification for the latter claim follows from [15, Lemma 12], [15, Theorem 13],
 395 and [15, Corollary 14], which also imply the following complexity result.

396 **PROPOSITION 6.** *The AIPP method terminates with a $\bar{\rho}$ -approximate stationary*
 397 *point of (3.2) in*

$$398 \quad (3.10) \quad \mathcal{O} \left(\sqrt{\lambda M + 1} \left[\frac{R(\phi; \lambda)}{\sqrt{\sigma}(1 - \sigma)^2 \lambda^2 \bar{\rho}^2} + \log_1^+(\lambda M) \right] \right)$$

399 ACG iterations, where

$$400 \quad (3.11) \quad R(\phi; \lambda) = \inf_{x'} \left\{ \frac{1}{2} \|x_0 - x'\|^2 + \lambda [\phi(x') - \phi_*] \right\}.$$

401 Note that scaling $R(\phi; \lambda)$ by $1/\lambda$ and then shifting by ϕ_* results in the λ -Moreau
 402 envelope¹ of ϕ . Moreover, $R(\phi; \lambda)$ admits the upper bound

$$403 \quad (3.12) \quad R(\phi; \lambda) \leq \min \left\{ \frac{1}{2} d_0^2, \lambda [\phi(x_0) - \phi_*] \right\}$$

404 where $d_0 := \inf \{ \|x_0 - x_*\| : x_* \text{ is an optimal solution of (3.2)} \}$.

405 **3.2. AIPP-S scheme for min-max CNO problems.** We are now ready to
 406 state the AIPP-S scheme for finding approximate stationary points of the uncon-
 407 strained min-max CNO problem (1.1).

408 It is stated in a incomplete manner in the sense that it does not specify how the
 409 parameter ξ and the tolerance ρ used in its step 2 are chosen. Two invocations of
 410 this method, with different choices of ξ and ρ , are considered in Propositions 8 and
 411 9, which describe the iteration complexities for finding approximate stationary points
 412 as in (1.4) and (1.5), respectively.

413

AIPP-S scheme

414 **Input:** a triple $(m, L_x, L_y) \in \mathbb{R}_{++}^3$ satisfying (A3), a smoothing constant $\xi > 0$, an
 415 initial point $(x_0, y_0) \in X \times Y$, and a tolerance $\rho > 0$;

416 **Output:** a pair $(x, u) \in X \times \mathcal{X}$;

417 (0) set L_ξ as in (2.13), $\sigma = 1/2$, $\lambda = 1/(4m)$, and define p_ξ as in (2.12);

418 (1) apply the AIPP method with inputs (m, L_ξ) , (p_ξ, h) , λ , σ , x_0 , and ρ to obtain
 419 a pair (x, u) satisfying

$$420 \quad (3.13) \quad u \in \nabla p_\xi(x) + \partial h(x), \quad \|u\| \leq \rho;$$

421 (2) output the pair (x, u) .

422

423 We now give four remarks about the above method. First, the AIPP method
 424 invoked in step 2 terminates due to [15, Theorem 13] and the third and fourth remarks

¹See [28, Chapter 1.G] for an exact definition.

425 following assumptions (P1)–(P3). Second, since the AIPP-S scheme is a one-pass
 426 method (as opposed to an iterative method), the complexity of the AIPP-S scheme is
 427 essentially that of the AIPP method. Third, similar to the smoothing scheme of [23]
 428 which assumes $m = 0$, the AIPP-S scheme is also a smoothing scheme for the case in
 429 which $m > 0$. On the other hand, in contrast to the algorithm of [23] which uses an
 430 ACG variant, AIPP-S invokes the AIPP method to solve (3.1) due to its nonconvexity.
 431 Finally, while the AIPP method in step 2 is called with $(\sigma, \lambda) = (1/2, 1/(4m))$, it
 432 can also be called with any $\sigma \in (0, 1)$ and $\lambda \in (0, 1/(2m))$ to establish the desired
 433 termination of the AIPP-S scheme.

434 For the remainder of this subsection, our goal will be to show that a careful
 435 selection of the parameter ξ and the tolerance ρ will allow the AIPP-S method to
 436 generate approximate stationary points as in (1.5) and (1.4).

437 Before proceeding, we first present a bound on the quantity $R(\hat{p}_\xi; \lambda)$ in terms of
 438 the data in problem (1.1). Its importance derives from the fact that the AIPP method
 439 applied to the smoothed problem (3.1) yields the bound (3.10) with $\phi = \hat{p}_\xi$.

440 LEMMA 7. *For every $\xi > 0$ and $\lambda \geq 0$, it holds that*

$$441 \quad (3.14) \quad R(\hat{p}_\xi; \lambda) \leq R(\hat{p}; \lambda) + \frac{\lambda D_y^2}{2\xi},$$

442 where $R(\cdot, \cdot)$ and D_y are as in (3.11) and (2.8), respectively.

443 *Proof.* Using Lemma 5 and the definitions of \hat{p} and \hat{p}_ξ , it holds that

$$444 \quad (3.15) \quad \hat{p}_\xi(x) - \inf_{x'} \hat{p}_\xi(x') \leq \hat{p}(x) - \inf_{x'} \hat{p}(x') + \frac{D_y^2}{2\xi}, \quad \forall x \in X.$$

445 Multiplying the above expression by $(1 - \sigma)\lambda$ and adding the quantity $\|x_0 - x\|^2/2$
 446 yields the inequality

$$447 \quad \frac{1}{2}\|x_0 - x\|^2 + (1 - \sigma)\lambda \left[\hat{p}_\xi(x) - \inf_{x'} \hat{p}_\xi(x') \right] \\ 448 \quad (3.16) \quad \leq \frac{1}{2}\|x_0 - x\|^2 + (1 - \sigma)\lambda \left[\hat{p}(x) - \inf_{\bar{x}} \hat{p}(x') \right] + (1 - \sigma) \frac{\lambda D_y^2}{2\xi} \quad \forall x \in X, \\ 449$$

450 Taking the infimum of the above expression, and using the definition of $R(\cdot; \cdot)$ in
 451 (3.11) yields the desired conclusion. \square

452 We now show how the AIPP-S scheme generates a (ρ_x, ρ_y) -primal-dual stationary
 453 point, i.e. one satisfying (1.4). Recall the definition of “oracle call” in the paragraph
 454 containing (1.3).

455 PROPOSITION 8. *For a given tolerance pair $(\rho_x, \rho_y) \in \mathbb{R}_{++}^2$, let (x, u) be the pair
 456 output by the AIPP-S scheme with input parameter ξ and tolerance ρ satisfying $\xi \geq$
 457 D_y/ρ_y and $\rho = \rho_x$. Moreover, define*

$$458 \quad (3.17) \quad (\bar{u}, \bar{v}) := \left(u, \frac{y_0 - y_\xi(x)}{\xi} \right), \quad (\bar{x}, \bar{y}) := (x, y_\xi(x)),$$

459 where y_ξ is as in (2.13). Then, the following statements hold:

460 (a) the AIPP-S scheme performs

$$461 \quad (3.18) \quad \mathcal{O} \left(\Omega_\xi \left[\frac{m^2 R(\hat{p}; 1/(4m))}{\rho_x^2} + \frac{m D_y^2}{\xi \rho_x^2} + \log_1^+(\Omega_\xi) \right] \right)$$

462 oracle calls, where $R(\cdot; \cdot)$ and D_y are as in (3.11) and (2.8), respectively, and

$$463 \quad (3.19) \quad \Omega_\xi := 1 + \frac{\sqrt{\xi}L_y + \sqrt{L_x}}{\sqrt{m}};$$

464 (b) the quadruple $(\bar{u}, \bar{v}, \bar{x}, \bar{y})$ is a (ρ_x, ρ_y) -primal-dual stationary point of (1.1).

465 *Proof.* (a) Using the inequality in (2.13), it holds that

$$466 \quad (3.20) \quad \sqrt{\frac{L_\xi}{4m}} + 1 \leq 1 + \sqrt{\frac{L_\xi}{4m}} \leq 1 + \frac{\sqrt{\xi}L_y + \sqrt{L_x}}{2\sqrt{m}} = \Theta(\Omega_\xi).$$

468 Moreover, using Proposition 6 with $(\phi, M) = (\hat{p}_\xi, L_\xi)$, Lemma 7, and bound (3.20),
 469 it follows that the number of ACG iterations performed by the AIPP-S scheme is on
 470 the order given by (3.18). Since step 1 of the AIPP invokes once the ACG variant in
 471 Appendix A with a pair (ψ_s, ψ_n) of the form

$$472 \quad \psi_s = \lambda p_\xi + \frac{1}{4} \|\cdot - \bar{z}\|^2, \quad \psi_n = \lambda h + \frac{1}{4} \|\cdot - \bar{z}\|^2$$

473 for some \bar{z} and each iteration of this ACG variant performs $\mathcal{O}(1)$ gradient evaluations
 474 of ψ_s , $\mathcal{O}(1)$ function evaluations of ψ_s and ψ_n , and $\mathcal{O}(1)$ ψ_n -resolvent evaluations, it
 475 follows from Proposition 4(b) and the definition of an ‘‘oracle call’’ in the paragraph
 476 containing (1.3) that each one of the above ACG iterations requires $\mathcal{O}(1)$ oracle calls.
 477 Statement (a) now follows from the above two conclusions.

478 (b) It follows from the definitions of p_ξ , tolerance ρ , and (\bar{y}, \bar{u}) in (2.12), the choice
 479 of ξ and ρ , and (3.17), respectively, Proposition 4(b), and the inclusion in (3.13) that
 480 $\|\bar{u}\| \leq \rho_x$ and

$$481 \quad \bar{u} \in \nabla p_\xi(\bar{x}) + \partial h(\bar{x}) = \nabla_x \Phi(\bar{x}, y_\xi(\bar{x})) + \partial h(\bar{x}) = \nabla_x \Phi(\bar{x}, \bar{y}) + \partial h(\bar{x}).$$

482 Hence, we conclude that the top inclusion and the upper bound on $\|\bar{u}\|$ in (1.4) hold.
 483 Next, the optimality condition of $\bar{y} = y_\xi(\bar{x})$ as a solution to (2.12) and the definition
 484 of \bar{v} in in (2.12) give

$$485 \quad (3.21) \quad 0 \in \partial [-\Phi(\bar{x}, \cdot)](\bar{y}) + \frac{\bar{y} - y_0}{\xi} = \partial [-\Phi(\bar{x}, \cdot)](\bar{y}) - \bar{v}$$

486 Moreover, the definition of ξ implies that $\|\bar{v}\| = \|\bar{y} - y_0\|/\xi \leq D_y/(D_y/\rho_y) = \rho_y$.
 487 Hence, combining (3.21) and the previous identity, we conclude that the bottom
 488 inclusion and the upper bound on $\|\bar{v}\|$ in (1.4) hold. \square

489 We now make three remarks about Proposition 8. First, recall that $R(\hat{p}; 1/(4m))$
 490 in the complexity (3.18) can be majorized by the rightmost quantity in (3.12) with
 491 $(\phi, \lambda) = (\hat{p}, 1/(4m))$. Second, under the assumption that $\xi = D_y/\rho_y$, the complexity
 492 of AIPP-S scheme reduces to

$$493 \quad (3.22) \quad \mathcal{O} \left(m^{3/2} \cdot R(\hat{p}; 1/(4m)) \cdot \left[\frac{L_x^{1/2}}{\rho_x^2} + \frac{L_y D_y^{1/2}}{\rho_x^2 \rho_y^{1/2}} \right] \right)$$

494 under the reasonable assumption that the $\mathcal{O}(\rho_x^{-2} + \rho_x^{-2} \rho_y^{-1/2})$ term in (3.18) dominates
 495 the other terms. Third, recall from the last remark following the previous proposition

496 that when Y is a singleton, (1.1) becomes a special instance of (3.2) and the AIPP-
 497 S scheme becomes equivalent to the AIPP method of Subsection 3.1. It similarly
 498 follows that the complexity in (3.22) reduces to $\mathcal{O}(\rho_x^{-2})$ and, in view of this remark,
 499 the $\mathcal{O}(\rho_x^{-2}\rho_y^{-1/2})$ term in (3.22) is attributed to the (possible) nonsmoothness in (1.1).

500 We next show how the AIPP-S scheme generates a point that is *near* a δ -
 501 directional stationary point, i.e., one satisfying (1.5). Recall the definition of “oracle
 502 call” in the paragraph containing (1.3).

503 **PROPOSITION 9.** *Let a tolerance pair $\delta > 0$ be given and consider the AIPP-S*
 504 *scheme with input parameter ξ and tolerance ρ satisfying $\xi \geq D_y/\tau$ and $\rho = \delta/2$ for*
 505 *some $\tau \leq \min\{m\delta^2/2D_y, \delta^2/32mD_y\}$. Then, the following statements hold:*

506 (a) *the AIPP-S scheme performs*

$$507 \quad (3.23) \quad \mathcal{O}\left(\Omega_\xi \left[\frac{R(\hat{p}; \lambda)}{\lambda^2 \delta^2} + \frac{D_y^2}{\lambda \xi \delta^2} + \log_1^+(\Omega_\xi) \right]\right)$$

508 *oracle calls where Ω_ξ , $R(\cdot; \cdot)$, and D_y are as in (3.19), (3.11), and (2.8),*
 509 *respectively;*

510 (b) *the first argument x in the pair output by the AIPP-S scheme satisfies (1.5).*

511 *Proof.* (a) Using Proposition 8 with $(\rho_x, \rho_y) = (\delta/2, \tau)$ and our assumption on τ
 512 it follows that the AIPP-S stops in a number of ACG iterations bounded above by
 513 (3.23).

514 (b) Let (x, u) be the $\bar{\rho}$ -approximate stationary point of (3.1) generated by the
 515 AIPP-S scheme (see step 2) under the given assumption on ξ and $\bar{\rho}$. Defining (\bar{v}, \bar{y})
 516 as in (3.17), it follows from Proposition 8 with $(\rho_x, \rho_y) = (\delta/2, \tau)$ that (u, \bar{v}, x, \bar{y}) is
 517 a $(\delta/2, \tau)$ -primal-dual stationary point of (1.1). As a consequence, it follows from
 518 Proposition 1 with $(\rho_x, \rho_y) = (\delta/2, \tau)$ that there exists a point \hat{x} satisfying

$$519 \quad (3.24) \quad \|\hat{x} - x\| \leq \sqrt{\frac{2D_y\tau}{m}}, \quad \inf_{\|d\| \leq 1} \hat{p}'(\hat{x}; d) \geq -\frac{\delta}{2} - 2\sqrt{2mD_y\tau}.$$

521 Combining the above bounds with our assumption on τ yields the desired conclusion
 522 in view of (1.5). \square

523 We now give four remarks about the above result. First, recall that $R(\hat{p}; 1/(4m))$
 524 in the complexity (3.23) can be majorized by the rightmost quantity in (3.12) with
 525 $(\phi, \lambda) = (\hat{p}, 1/(4m))$. Second, Proposition 9(b) states that, while x not a stationary
 526 point itself, it is near a δ -directional stationary point \hat{x} . Third, under the assumption
 527 that the bounds on ξ and τ in Proposition 9 hold at equality, the complexity of the
 528 AIPP-S scheme reduces to

$$529 \quad (3.25) \quad \mathcal{O}\left(m^{3/2} \cdot R(\hat{p}; 1/(4m)) \cdot \left[\frac{L_x^{1/2}}{\delta^2} + \frac{L_y D_y}{\delta^3} \right]\right)$$

530 under the reasonable assumption that the $\mathcal{O}(\delta^{-2} + \delta^{-3})$ term in (3.23) dominates
 531 the other $\mathcal{O}(\delta^{-1})$ terms. Fourth, when Y is a singleton, it is easy to see that (1.1)
 532 becomes a special instance of (3.2), the AIPP-S scheme becomes equivalent to the
 533 AIPP method of Subsection 3.1, and the complexity in (3.25) reduces to $\mathcal{O}(\delta^{-2})$. In
 534 view of the last remark, the $\mathcal{O}(\delta^{-3})$ term in (3.25) is attributed to the (possible)
 535 nonsmoothness in (1.1).

536 **4. Linearly-constrained min-max optimization.** This section presents our
537 proposed QP-AIPP-S scheme for solving the linearly constrained min-max CNO prob-
538 lem (1.6), and it is divided into two subsections. The first one reviews a QP-AIPP
539 method for solving smooth linearly-constrained CNO problems. The second one
540 presents the QP-AIPP-S scheme and its iteration complexity for finding stationary
541 points as in (1.8). Throughout our presentation, we let \mathcal{X}, \mathcal{Y} , and \mathcal{U} be finite dimen-
542 sional inner product spaces.

543 Before proceeding, let us give the precise assumptions underlying the problem of
544 interest and discuss the relevant notion of stationarity. For problem (1.6) suppose
545 that assumptions (A0)–(A3) hold and that the linear operator $\mathcal{A} : \mathcal{X} \mapsto \mathcal{U}$ and vector
546 $b \in \mathcal{U}$ satisfy:

547 (A5) $\mathcal{A} \neq 0$ and $\mathcal{F} := \{x \in X : \mathcal{A}x = b\} \neq \emptyset$;

548 (A6) there exists $\hat{c} \geq 0$ such that $\inf_{x \in X} \{\hat{p}(x) + \hat{c}\|\mathcal{A}x - b\|^2/2\} > -\infty$.

549 Note that (A4) in Subsection 2.1 is replaced by (A6) which is required by the QP-
550 AIPP method of the next subsection.

551 Analogous to the first remark following (2.6), it is known that if (x^*, y^*) satisfies
552 (2.5) for every $(x, y) \in \mathcal{F} \times Y$ and $\hat{\Phi}$ as in (2.4), then there exists a multiplier $r^* \in \mathcal{U}$
553 such that

$$554 \quad (4.1) \quad \begin{pmatrix} 0 \\ 0 \end{pmatrix} \in \begin{pmatrix} \nabla_x \Phi(x^*, y^*) + A^* r^* \\ 0 \end{pmatrix} + \begin{pmatrix} \partial h(x^*) \\ \partial [-\Phi(x^*, \cdot)](y^*) \end{pmatrix},$$

555 holds. Hence, in view of the third remark in the paragraph following (2.7), we only
556 consider the approximate version of (4.1) which is (1.8).

557 We now briefly outline the idea of the QP-AIPP-S scheme. The main idea is to
558 apply the QP-AIPP method described in the next subsection to the smooth linearly-
559 constrained CNO problem

$$560 \quad (4.2) \quad \min_{x \in X} \{p_\xi(x) + h(x) : \mathcal{A}x = b\},$$

561 where p_ξ is as in (1.2) and ξ is a positive scalar that will depend on the tolerances
562 in (1.8). This idea is similar to the one in Section 3 in that it applies an accelerated
563 solver to a perturbed version of the problem of interest.

564 **4.1. QP-AIPP method for constrained smooth CNO problems.** This
565 subsection describes the QP-AIPP method studied in [15], and its corresponding
566 iteration complexity, for solving linearly-constrained smooth CNO problems.

567 We begin by describing the problem that the QP-AIPP method intends to solve.
568 Consider the linearly-constrained smooth CNO problem

$$569 \quad (4.3) \quad \hat{\phi}_* := \inf_{x \in \mathcal{X}} \{\phi(x) := f(x) + h(x) : \mathcal{A}x = b\}$$

570 where $h : \mathcal{X} \mapsto (-\infty, \infty]$ and a function f satisfy assumptions (P1)–(P3), the operator
571 $\mathcal{A} : \mathcal{X} \mapsto \mathcal{U}$ is linear, $b \in \mathcal{U}$, and the following additional assumptions hold:

572 (Q1) $\mathcal{A} \neq 0$ and $\mathcal{F} := \{x \in \text{dom } h : \mathcal{A}x = b\} \neq \emptyset$;

573 (Q2) there exists $\hat{c} \geq 0$ such that $\hat{\phi}_{\hat{c}} > -\infty$ where

$$574 \quad (4.4) \quad \hat{\phi}_c := \inf_{x \in \mathcal{X}} \left\{ \phi_c(x) := \phi(x) + \frac{c}{2} \|\mathcal{A}x - b\|^2 \right\}, \quad \forall c \geq 0.$$

575 We now give some remarks about the above assumptions. First, similar to problem
576 (3.2), it is well-known that a necessary condition for $x^* \in \text{dom } h$ to be a local minimum

577 of (4.3) is that x^* satisfies $0 \in \nabla f(x^*) + \partial h(x^*) + \mathcal{A}^* r^*$ for some $r^* \in \mathcal{U}$. Second,
578 it is straightforward to verify that (p, h, \mathcal{A}, b) in (1.6) satisfy (Q1)–(Q2) in view of
579 assumptions (A5)–(A6). Third, since every feasible solution of (4.3) is also a feasible
580 solution of (4.4), it follows from assumptions (Q2) that $\hat{\phi}_* \geq \hat{\phi}_{\hat{c}} > -\infty$. Fourth, if
581 $\inf_{x \in \mathcal{X}} \phi(x) > -\infty$ (e.g., $\text{dom } h$ is compact) then (Q2) holds with $\hat{c} = 0$.

582 Our interest in this subsection is in finding an approximate stationary point of
583 (4.3) in the following sense: given a tolerance pair $(\bar{\rho}, \bar{\eta}) \in \mathbb{R}_{++}^2$, a triple $(\bar{x}, \bar{u}, \bar{r}) \in$
584 $\text{dom } h \times \mathcal{X} \times \mathcal{U}$ is said to be a $(\bar{\rho}, \bar{\eta})$ -approximate stationary point of (4.3) if

$$585 \quad (4.5) \quad \bar{u} \in \nabla f(\bar{x}) + \partial h(\bar{x}) + \mathcal{A}^* \bar{r}, \quad \|\bar{u}\| \leq \bar{\rho}, \quad \|\mathcal{A}\bar{x} - b\| \leq \bar{\eta}.$$

586 We now state the QP-AIPP method for finding $(\bar{x}, \bar{u}, \bar{r})$ satisfying (4.5).
587

QP-AIPP method

588 **Input:** a function pair (f, h) , a scalar pair $(m, M) \in \mathbb{R}_{++}^2$ satisfying (3.3), scalars
589 $\lambda \in (0, 1/(2m)]$ and $\sigma \in (0, 1)$, a scalar \hat{c} satisfying assumption (Q2), an initial point
590 $x_0 \in \text{dom } h$, and a tolerance pair $(\bar{\rho}, \bar{\eta}) \in \mathbb{R}_{++}^2$;

591 **Output:** a triple $(\bar{x}, \bar{u}, \bar{r}) \in \text{dom } h \times \mathcal{X} \times \mathcal{U}$ satisfying (4.5);

- 592 (0) set $c = \hat{c} + M/\|\mathcal{A}\|^2$;
593 (1) define the quantities

$$594 \quad (4.6) \quad M_c := M + c\|\mathcal{A}\|^2, \quad f_c := f + \frac{c}{2}\|\mathcal{A}(\cdot) - b\|^2, \quad \phi_c = f_c + h,$$

- 595 and apply the AIPP method with inputs (m, M_c) , (f_c, h) , λ , σ , x_0 , and $\bar{\rho}$ to
596 obtain a $\bar{\rho}$ -approximate stationary point (\bar{x}, \bar{u}) of (3.2) with $f = f_c$;
597 (2) if $\|\mathcal{A}\bar{x} - b\| > \bar{\eta}$ then set $c = 2c$ and go to (1); otherwise, set $\bar{r} = c(\mathcal{A}\bar{x} - b)$
598 and output the triple $(\bar{x}, \bar{u}, \bar{r})$.
599

600 We now give two remarks about the above method. First, it straightforward to
601 see that QP-AIPP method terminates due to the results in [15, Section 4]. Second,
602 in view of Proposition 6 with $(\phi, M) = (\phi_c, M_c)$, it is easy to see that the number of
603 ACG iterations executed in step 1 at any iteration of the method is

$$604 \quad (4.7) \quad \mathcal{O} \left(\sqrt{\lambda M_c + 1} \left[\frac{R(\phi_c; \lambda)}{\sqrt{\sigma}(1-\sigma)^2 \lambda^2 \bar{\rho}^2} + \log_1^+ (\lambda M_c) \right] \right)$$

605 and that the pair (\bar{x}, \bar{u}) computed in step 1 satisfies the inclusion and the first in-
606 equality in (4.5).

607 We now focus on the iteration complexity of the QP-AIPP method. Before pro-
608 ceeding, we first define the useful quantity

$$609 \quad (4.8) \quad R_c(\phi; \lambda) := \inf_{x'} \left\{ \frac{1}{2} \|x_0 - x'\|^2 + \lambda \left[\phi(x') - \hat{\phi}_c \right] : x' \in \mathcal{F} \right\},$$

610 for every $c \geq \hat{c}$, where ϕ_c is as defined in (4.4). The quantity in (4.8) plays an analogous
611 role as (3.11) in (3.10) and, similar to the discussion following Proposition 6, it is a
612 scaled and shifted λ -Moreau envelope of $\phi + \delta_{\mathcal{F}}$. Moreover, due to [15, Lemma 16], it
613 also admits the upper bound

$$614 \quad (4.9) \quad R_c(\phi; \lambda) \leq R_{\hat{c}}(\phi; \lambda) \leq \min \left\{ \frac{1}{2} \hat{d}_0^2, \lambda \left[\hat{\phi}_* - \hat{\phi}_{\hat{c}} \right] \right\}$$

615 where $\hat{\phi}_*$ is as defined in (4.3) and

$$616 \quad \hat{d}_0 := \inf \{ \|x_0 - x_*\| : x_* \text{ is an optimal solution of (4.3)} \}.$$

617 We now state the iteration complexity of the QP-AIPP method, whose proof may
618 be adapted from [15, Lemma 12] and [15, Theorem 18].

619 **PROPOSITION 10.** *Let a constant \hat{c} as in assumption (Q2), scalar $\sigma \in (0, 1)$,
620 curvature pair $(m, M) \in \mathbb{R}_{++}^2$, and a tolerance pair $(\bar{\rho}, \bar{\eta}) \in \mathbb{R}_+^2$ be given. Moreover,
621 define*

$$622 \quad (4.10) \quad T_{\bar{\eta}} := \frac{2R_{\hat{c}}(\phi; \lambda)}{\bar{\eta}^2(1-\sigma)\lambda} + \hat{c}, \quad \Theta_{\bar{\eta}} := M + T_{\bar{\eta}}\|\mathcal{A}\|^2.$$

623 Then, the QP-AIPP method outputs a triple $(\bar{x}, \bar{u}, \bar{r})$ satisfying (4.5) in

$$624 \quad (4.11) \quad \mathcal{O} \left(\sqrt{\lambda\Theta_{\bar{\eta}} + 1} \left[\frac{R_{\hat{c}}(\phi; \lambda)}{\sqrt{\sigma}(1-\sigma)^2\lambda^2\bar{\rho}^2} + \log_1^+(\lambda\Theta_{\bar{\eta}}) \right] \right)$$

625 ACG iterations.

626 **4.2. QP-AIPP-S scheme for constrained min-max CNO problems.** We
627 are now ready to state the QP-AIPP smoothing scheme for finding an approximate
628 primal-dual stationary point of the linearly-constrained min-max CNO problem (1.6).
629

QP-AIPP-S scheme

630 **Input:** a triple $(m, L_x, L_y) \in \mathbb{R}_{++}^2$ satisfying assumption (A3), a scalar \hat{c} satisfying
631 assumption (A6), a smoothing constant $\xi \geq D_y/\rho_y$, an initial point $(x_0, y_0) \in X \times Y$,
632 and a tolerance triple $(\rho_x, \rho_y, \eta) \in \mathbb{R}_{++}^3$;

633 **Output:** a triple $(\bar{u}, \bar{v}, \bar{x}, \bar{y}, \bar{r})$ satisfying (1.8);

- 634 (0) set L_ξ as in (2.13), $\sigma = 1/2$, $\lambda = 1/(4m)$, and define p_ξ as in (2.12);
635 (1) apply the QP-AIPP method of Subsection 4.1 with inputs (m, L_ξ) , (p_ξ, h) , λ ,
636 σ , \hat{c} , x_0 , and (ρ_x, η) to obtain a triple $(\bar{u}, \bar{x}, \bar{r})$ satisfying

$$637 \quad (4.12) \quad \bar{u} \in \nabla p_\xi(\bar{x}) + \partial h(\bar{x}) + A^*\bar{r}, \quad \|\bar{u}\| \leq \rho_x, \quad \|\mathcal{A}\bar{x} - b\| \leq \eta.$$

- 638 (2) define (\bar{v}, \bar{y}) as in (3.17) and output the quintuple $(\bar{u}, \bar{v}, \bar{x}, \bar{y}, \bar{r})$.
-

639 Some remarks about the above method are in order. First, the QP-AIPP method
640 invoked in step 1 terminates due to the remarks following assumptions (Q1)–(Q2)
641 and the results in Subsection 4.1. Second, since the QP-AIPP-S scheme is a one-pass
642 algorithm (as opposed to an iterative algorithm), the complexity of the QP-AIPP-
643 S scheme is essentially that of the QP-AIPP method. Finally, while the QP-AIPP
644 method in step 2 is called with $(\sigma, \lambda) = (1/2, 1/(4m))$, it can also be called with any
645 $\sigma \in (0, 1)$ and $\lambda \in (0, 1/(2m))$ to establish the desired termination of the QP-AIPP-S
646 scheme.
647

648 We now show how the QP-AIPP-S scheme generates a point $(\bar{u}, \bar{v}, \bar{x}, \bar{y}, \bar{r})$ satisfy-
649 ing (1.8). Recall the definition of “oracle call” in the paragraph containing (1.3).

650 **PROPOSITION 11.** *Let a tolerance triple $(\rho_x, \rho_x, \eta) \in \mathbb{R}_{++}^3$ be given and let the
651 quadruple $(\bar{u}, \bar{v}, \bar{x}, \bar{y}, \bar{r})$ be the output obtained by the QP-AIPP-S scheme. Then the
652 following properties hold:*

653 (a) the QP-AIPP-S scheme terminates in

$$654 \quad (4.13) \quad \mathcal{O} \left(\Omega_{\xi, \eta} \left[\frac{m^2 R_{\hat{c}}(\hat{p}; 1/(4m))}{\rho_x^2} + \frac{m D_y^2}{\xi \rho_x^2} + \log_1^+ (\Omega_{\xi, \eta}) \right] \right)$$

655 oracle calls, where

$$656 \quad (4.14) \quad \Omega_{\xi, \eta} := \Omega_{\xi} + \left(R_{\hat{c}}(\hat{p}; 1/(4m)) + \frac{D_y^2}{m \xi} \right)^{1/2} \frac{\|\mathcal{A}\|}{\eta}$$

657 and Ω_{ξ} , $R(\cdot; \cdot)$, and D_y are as in (3.19), (3.11), and (2.8), respectively;

658 (b) the quintuple $(\bar{u}, \bar{v}, \bar{x}, \bar{y}, \bar{r})$ satisfies (1.8).

659 *Proof.* (a) Let Θ_{η} be as in (4.10) with $M = L_{\xi}$. Using the same arguments as
660 in Lemma 7, it is easy to see that $R_{\hat{c}}(\hat{p}_{\xi}; 1/(4m)) \leq R_{\hat{c}}(\hat{p}; 1/(4m)) + D_y^2/(8m\xi)$, and
661 hence, using (3.20), we have

$$662 \quad \sqrt{\frac{\Theta_{\eta}}{4m} + 1} \leq 1 + \sqrt{\frac{L_{\xi}}{4m}} + \sqrt{\frac{4R_{\hat{c}}(\hat{p}_{\xi}; 1/(4m))\|\mathcal{A}\|^2}{\eta^2}}$$

$$663 \quad (4.15) \quad \leq 1 + \frac{\sqrt{\xi}L_y + \sqrt{L_x}}{2\sqrt{m}} + 2 \left(R_{\hat{c}}(\hat{p}; 1/(4m)) + \frac{D_y^2}{8m\xi} \right)^{1/2} \frac{\|\mathcal{A}\|}{\eta} = \Theta(\Omega_{\xi, \eta}).$$

664

665 The complexity in (4.13) now follows from the above bound and Proposition 10 with
666 $(\phi, f, M) = (p, p_{\xi}, L_{\xi})$.

667 (b) The top inclusion and bounds involving $\|\bar{u}\|$ and $\|\mathcal{A}\bar{x} - b\|$ in (1.8) follow from
668 Proposition 4(b), the definition of \bar{y} in step 2 of the algorithm, and Proposition 10
669 with $f = p_{\xi}$. The bottom inclusion and bound involving $\|\bar{v}\|$ follow from similar
670 arguments given for Proposition 8(b). \square

671 We now make three remarks about the above complexity bound. First, recall that
672 $R_{\hat{c}}(p; 1/(4m))$ in the complexity (11) can be majorized by the rightmost quantity in
673 (4.9) with $\lambda = 1/(4m)$. Second, under the assumption that $\xi = D_y/\rho_y$, the complexity
674 of the QP-AIPP-S scheme reduces to

$$675 \quad (4.16) \quad \mathcal{O} \left(m^{3/2} \cdot R_{\hat{c}}(\hat{p}; 1/(4m)) \cdot \left[\frac{L_x^{1/2}}{\rho_x^2} + \frac{L_y D_y^{1/2}}{\rho_y^{1/2} \rho_x^2} + \frac{m^{1/2} \|\mathcal{A}\| R_{\hat{c}}^{1/2}(p; 1/(4m))}{\eta \rho_x^2} \right] \right),$$

676 under the reasonable assumption that the $\mathcal{O}(\rho_x^{-2} + \eta^{-1} \rho_x^{-2} + \rho_y^{-1/2} \rho_x^{-2})$ term in (4.13)
677 dominates the other terms. Third, when Y is a singleton, it is easy to see that (1.6)
678 becomes a special instance of the linearly-constrained smooth CNO problem (4.3),
679 the QP-AIPP-S of this subsection becomes equivalent to the QP-AIPP method of
680 Subsection 4.1, and the complexity in (4.16) reduces to $\mathcal{O}(\eta^{-1} \rho_x^{-2})$. In view of the last
681 remark, the $\mathcal{O}(\rho_x^{-2} \rho_y^{-1/2})$ term in (4.16) is attributed to the (possible) nonsmoothness
682 in (1.6).

683 Let us now conclude this section with a remark about the penalty subproblem

$$684 \quad (4.17) \quad \min_{x \in X} \left\{ p_{\xi}(x) + h(x) + \frac{c}{2} \|\mathcal{A}x - b\|^2 \right\},$$

685 which is what the AIPP method considers every time it is called in the QP-AIPP-S
686 scheme (see step 1). First, observe that (1.6) can be equivalently reformulated as

$$687 \quad (4.18) \quad \min_{x \in X} \max_{y \in Y, r \in \mathcal{U}} [\Psi(x, y, r) := \Phi(x, y) + h(x) + \langle r, \mathcal{A}x - b \rangle].$$

688 Second, it is straightforward to verify that problem (4.17) is equivalent to

$$689 \quad (4.19) \quad \min_{x \in X} \{\hat{p}_{c,\xi}(x) := p_{c,\xi}(x) + h(x)\},$$

690 where the function $p_{c,\xi} : X \mapsto \mathbb{R}$ is given by $p_{c,\xi}(x) := \max_{y \in Y, r \in \mathcal{U}} \{\Psi(x, y, r) -$
691 $\|r\|^2/(2c) - \|y - y_0\|^2/(2\xi)\}$, for every $x \in X$, and Ψ as in (4.18). As a consequence,
692 problem (4.19) is similar to (3.1) in that a smooth approximate is used in place of the
693 nonsmooth component of the underlying saddle function Ψ .

694 On the other hand, observe that we cannot directly apply the smoothing scheme
695 developed in Subsection 3.2 to (4.19) as the set \mathcal{U} is generally unbounded. One
696 approach that avoids this problem is to invoke the AIPP method of Subsection 3.1 to
697 solve a sequence subproblems of the form in (4.19) for increasing values of c . However,
698 in view of the equivalence of (4.17) and (4.19), this is exactly the approach taken by
699 the QP-AIPP-S scheme of this section.

700 **5. Numerical experiments.** This section presents numerical results that illus-
701 trate the computational efficiency of the our proposed smoothing scheme. It contains
702 three subsections. Each subsection presents computational results for a specific un-
703 constrained nonconvex min-max optimization problem class.

704 Each unconstrained problem considered in this section is of the form in (1.1) and
705 is such that the computation of the function y_ξ in (2.13) is easy. Moreover, for a
706 given initial point $x_0 \in X$, three algorithms are run for each problem instance until a
707 quadruple $(\bar{u}, \bar{v}, \bar{x}, \bar{y})$ satisfying the inclusion of (1.4) and

$$708 \quad (5.1) \quad \frac{\|\bar{u}\|}{\|\nabla p_\xi(z_0)\| + 1} \leq \rho_x, \quad \|\bar{v}\| \leq \rho_y,$$

709 is obtained, where $\xi = D_y/\rho_y$.

710 We now describe the three nonconvex-concave min-max methods that are being
711 compared in this section, namely: (i) the R-AIPP-S method; (ii) the accelerated
712 gradient smoothing (AG-S) scheme; and (iii) the projected gradient step framework
713 (PGSF). Both the AG-S and R-AIPP-S schemes are modifications of the AIPP-S
714 scheme which, instead of using the AIPP method in its step 1, use the AG method
715 of [10] and R-AIPP method of [16], respectively. The PGSF is a simplified variant
716 of Algorithm 2 of [24, Subsection 4.1] which explicitly evaluates the argmax function
717 $\alpha^*(\cdot)$ in [24, Section 4] instead of applying an ACG variant to estimate its evaluation.

718 Regarding the penalty solvers, the AG method is implemented as described in
719 Algorithm 2 of [10] while the R-AIPP method follows the implementation described
720 in [14, Section 5.3].

721 Note that, like the AIPP method, the R-AIPP similarly: (i) invokes at each of its
722 (outer) iterations an ACG method to inexactly solve the proximal subproblem (3.9);
723 and (ii) outputs a $\bar{\rho}$ -approximate stationary point of (3.2). However, the R-AIPP
724 method is more computationally efficient due to three key practical improvements
725 over the AIPP method, namely: (i) it allows the stepsize λ to be significantly larger
726 than the $1/(2m)$ upper bound in the AIPP method using adaptive estimates of m ;
727 (ii) it uses a weaker ACG termination criterion compared to the one in (3.6); and (iii)
728 it does not prespecify the minimum number of ACG iterations as the AIPP method
729 does in its step 1.

730 We next state some additional details about the numerical experiments. First,
731 each algorithm is run with a time limit of 4000 seconds. Second, the bold numbers in
732 each of the computational tables in this section highlight the algorithm that performed
733

734 the most efficiently in terms of iteration count or total runtime. Moreover, each of
735 tables contain a column labeled $\hat{p}_\xi(\bar{x})$ that contains the smallest obtained value of the
736 smoothed function in (3.1), across all of the tested algorithms. Third, the description
737 of y_ξ and choice of the constants m , L_x , and L_y for each of the considered optimization
738 problems can be found in [14, Appendix I]. Fourth, y_0 is chosen to be 0 for all of
739 the experiments. Finally, all algorithms described at the beginning of this section are
740 implemented in MATLAB 2019a and are run on Linux 64-bit machines each containing
741 Xeon E5520 processors and at least 8 GB of memory.

742 Before proceeding, it is worth mentioning that the code for generating the results
743 of this section is available online².

744 **5.1. Maximum of a finite number of nonconvex quadratic forms.** This
745 subsection presents computational results for a minmax quadratic vector problem,
746 which is based on a similar problem in [16].

747 We first describe the problem. Given a dimension triple $(n, l, k) \in \mathbb{N}^3$, a set of
748 parameters $\{(\alpha_i, \beta_i)\}_{i=1}^k \subseteq \mathbb{R}_{++}^2$, a set of vectors $\{d_i\}_{i=1}^k \subseteq \mathbb{R}^l$, a set of diagonal
749 matrices $\{D_i\}_{i=1}^k \subseteq \mathbb{R}^{n \times n}$, and matrices $\{C_i\}_{i=1}^k \subseteq \mathbb{R}^{l \times n}$ and $\{B_i\}_{i=1}^k \subseteq \mathbb{R}^{n \times n}$, the
750 problem of interest is the quadratic vector minmax (QVM) problem

$$751 \quad \min_{x \in \mathbb{R}^n} \max_{y \in \mathbb{R}^k} \left\{ \delta_{\Delta^n}(x) + \sum_{i=1}^k y_i g_i(x) : y \in \Delta^k \right\},$$

752 where, for every index $1 \leq i \leq k$, integer $p \in \mathbb{N}$, and $x \in \mathbb{R}^n$, we define $g_i(x) :=$
753 $\alpha_i \|C_i x - d_i\|^2/2 - \beta_i \|D_i B_i x\|^2/2$ and $\Delta^p := \{z \in \mathbb{R}_+^p : \sum_{i=1}^p z_i = 1, z \geq 0\}$.

754 We now describe the experiment parameters for the instances considered. First,
755 the dimensions are set to be $(n, l, k) = (200, 10, 5)$ and only 5.0% of the entries of the
756 submatrices B_i and C_i are nonzero. Second, the entries of B_i , C_i , and d_i (resp., D_i)
757 are generated by sampling from the uniform distribution $\mathcal{U}[0, 1]$ (resp., $\mathcal{U}[1, 1000]$).
758 Third, the initial starting point is $z_0 = I_n/n$, where I_n is the n -dimensional identity
759 matrix. Fourth, with respect to the termination criterion, the inputs, for every $(x, y) \in$
760 $\mathbb{R}^n \times \mathbb{R}^k$, are $\Phi(x, y) = \sum_{i=1}^k y_i g_i(x)$, $h(x) = \delta_{\Delta^n}(x)$, $\rho_x = 10^{-2}$, $\rho_y = 10^{-1}$, and
761 $Y = \Delta^k$. Finally, each problem instance considered is based on a specific curvature
762 pair (m, M) satisfying $m \leq M$, for which each scalar pair $(\alpha_i, \beta_i) \in \mathbb{R}_{++}^2$ is selected
763 so that $M = \lambda_{\max}(\nabla^2 g_i)$ and $-m = \lambda_{\min}(\nabla^2 g_i)$.

764 We now present the results in Table 5.1.

M	m	$\hat{p}_\xi(\bar{x})$	Iteration Count			Runtime		
			R-AIPP-S	AG-S	PGSF	R-AIPP-S	AG	PGSF
10^0	10^0	2.85E-01	23	294	1591	0.66	5.72	22.60
10^1	10^0	2.88E+00	86	1371	14815	1.37	25.96	209.62
10^2	10^0	2.85E+01	217	6270	150493	3.35	118.32	2122.93
10^3	10^0	2.85E+02	1417	28989	-	21.58	546.25	4000.00*

TABLE 5.1
Iteration counts and runtimes for QVM problems.

765 **5.2. Truncated robust regression.** This subsection presents computational
766 results for the robust regression problem in [26].

²See the examples in `./examples/minmax/` from the GitHub repository
https://github.com/wwkong/nc_opt/.

767 It is worth mentioning that [26] also presents a min-max algorithm for obtaining
768 a stationary point as in (5.1). However, its iteration complexity, which is $\mathcal{O}(\rho^{-6})$
769 when $\rho = \rho_x = \rho_y$, is significantly worse than the other algorithms considered in this
770 section and, hence, we choose not to include this algorithm in our benchmarks.

771 We now describe the problem. Given a dimension pair $(n, k) \in \mathbb{N}^2$, a set of n data
772 points $\{(a_j, b_j)\}_{j=1}^n \subseteq \mathbb{R}^k \times \{1, -1\}$ and a parameter $\alpha > 0$, the problem of interest is
773 the truncated robust regression (TRR) problem

$$774 \quad \min_{x \in \mathbb{R}^k} \max_{y \in \mathbb{R}^n} \left\{ \sum_{j=1}^n y_j (\phi_\alpha \circ \ell_j)(x) : y \in \Delta^n \right\}$$

775 where Δ^n is as in Subsection 5.1 with $p = n$, $\phi_\alpha(t) := \alpha \log(1 + t/\alpha)$, and $\ell_j(x) :=$
776 $\log(1 + e^{-b_j(a_j, x)})$, for every $(\alpha, t, x) \in \mathbb{R}_{++} \times \mathbb{R}_{++} \times \mathbb{R}^k$,

777 We now describe the experiment parameters for the instances considered. First,
778 α is set to 10 and the data points $\{(a_i, b_i)\}$ are taken from different datasets in
779 the LIBSVM library³ for which each problem instance is based off of (see the “data
780 name” column in the table below, which corresponds to a particular LIBSVM dataset).
781 Second, the initial starting point is $z_0 = 0$. Third, with respect to the termination
782 criterion, the inputs, for every $(x, y) \in \mathbb{R}^k \times \mathbb{R}^n$, are $\Phi(x, y) = \sum_{j=1}^n y_j (\phi_\alpha \circ \ell_j)(x)$,
783 $h(x) = 0$, $\rho_x = 10^{-5}$, $\rho_y = 10^{-3}$, and $Y = \Delta^n$.

784 We now present the results in Table 5.2.

data name	$\hat{p}_\xi(\bar{x})$	Iteration Count			Runtime		
		R-AIPP-S	AG-S	PGSF	R-AIPP-S	AG	PGSF
heart	6.70E-01	425	1747	6409	6.37	15.54	32.76
diabetes	6.70E-01	852	1642	3718	8.61	24.12	52.77
ionosphere	6.70E-01	1197	8328	54481	8.26	63.82	320.72
sonar	6.70E-01	45350	96209	-	461.52	580.37	4000.00*
breast-cancer	1.11E-03	46097	-	-	476.59	4000.00*	4000.00*

TABLE 5.2
Iteration counts and runtimes for TRR problems

785 **5.3. Power control in the presence of a jammer.** This subsection presents
786 computational results for the power control problem in [18].

787 It is worth mentioning that [18] also presents a min-max algorithm for obtaining
788 stationary points for the aforementioned problem. However, its termination criterion
789 and notion of stationarity are significantly different than what is being considered in
790 this paper and, hence, we choose not to include the algorithm of [18] in our bench-
791 marks.

792 We now describe the problem. Given a dimension pair $(N, K) \in \mathbb{N}^2$, a pair of
793 parameters $(\sigma, R) \in \mathbb{R}_{++}^2$, a 3D tensor $\mathcal{A} \in \mathbb{R}_+^{K \times K \times N}$, and a matrix $B \in \mathbb{R}_+^{K \times N}$, the
794 problem of interest is the power control (PC) problem

$$795 \quad \min_{X \in \mathbb{R}^{K \times N}} \max_{y \in \mathbb{R}^N} \left\{ \sum_{k=1}^K \sum_{n=1}^N f_{k,n}(X, y) : 0 \leq X \leq R, 0 \leq y \leq \frac{N}{2} \right\},$$

³See <https://www.csie.ntu.edu.tw/~cjlin/libsvmtools/datasets/binary.html>.

796 where, for every $(X, y) \in \mathbb{R}^{K \times N} \times \mathbb{R}^N$,

$$797 \quad f_{k,n}(X, y) := -\log \left(1 + \frac{\mathcal{A}_{k,k,n} X_{k,n}}{\sigma^2 + B_{k,n} y_n + \sum_{j=1, j \neq k}^K \mathcal{A}_{j,k,n} X_{j,n}} \right).$$

798 We now describe the experiment parameters for the instances considered. First,
 799 the scalar parameters are set to be $(\sigma, R) = (1/\sqrt{2}, K^{1/K})$ and the quantities \mathcal{A} and
 800 B are set to be the squared moduli of the entries of two Gaussian sampled complex-
 801 valued matrices $\mathcal{H} \in \mathbb{C}^{K \times K \times N}$ and $P \in \mathbb{C}^{K \times N}$. More precisely, the entries of \mathcal{H}
 802 and P are sampled from the standard complex Gaussian distribution $\mathcal{CN}(0, 1)$ with
 803 $\mathcal{A}_{j,k,n} = |\mathcal{H}_{j,k,n}|^2$ and $B_{k,n} = |P_{k,n}|^2$ for every (j, k, n) . Second, the initial starting
 804 point is $z_0 = 0$. Third, with respect to the termination criterion, the inputs, are
 805 $\Phi(X, y) = \sum_{k=1}^K \sum_{n=1}^N f_{k,n}(X, y)$, $h(X) = \delta_{Q_R^{K \times N}}(X)$, $\rho_x = 10^{-1}$, $\rho_y = 10^{-1}$, and
 806 $Y = Q_{N/2}^{N \times 1}$, for every $(X, y) \in \mathbb{R}^{K \times N} \times \mathbb{R}^N$ and $(U, V) \in \mathbb{N}^2$, where $Q_T^{U \times V} := \{z \in$
 807 $\mathbb{R}^{p \times q} : 0 \leq z \leq T\}$ for every $T > 0$. Fourth, each problem instance considered is
 808 based on a specific dimension pair (N, K) .

809 We now present the results in Table 5.3.

N	K	$\hat{p}_\xi(\bar{x})$	Iteration Count			Runtime		
			R-AIPP-S	AG-S	PGSF	R-AIPP-S	AG	PGSF
5	5	-3.64E+00	37	322832	-	0.96	2371.27	4000.00*
10	10	-2.82E+00	54	33399	-	0.75	293.60	4000.00*
25	25	-4.52E+00	183	-	-	9.44	4000.00*	4000.00*
50	50	-4.58E+00	566	-	-	40.89	4000.00*	4000.00*

TABLE 5.3
Iteration counts and runtimes for PC problems.

810 **6. Concluding Remarks.** This section makes some concluding remarks.

811 We first make a final remark about the AIPP-S smoothing scheme. Recall that
 812 the main idea of AIPP-S is to call the AIPP method to obtain a pair satisfying (3.13),
 813 or equivalently⁴,

$$814 \quad (6.1) \quad \inf_{\|d\| \leq 1} (\hat{p}_\xi)'(x; d) \geq -\rho.$$

815 Moreover, using Proposition 8 with $(\rho_x, \rho_y) = (\rho, D_y/\xi)$, it straightforward to see
 816 that that the number of oracle calls, in terms of (ξ, ρ) , is $\mathcal{O}(\rho^{-2}\xi^{1/2})$, which reduces
 817 to $\mathcal{O}(\rho^{-2.5})$ if ξ is chosen so as to satisfy $\xi = \Theta(\rho^{-1})$. The latter complexity bound
 818 improves upon the one obtained for an algorithm in [24] which obtains a point x
 819 satisfying (6.1) with $\xi = \Theta(\rho^{-1})$ in $\mathcal{O}(\rho^{-3})$ oracle calls.

820 We now discuss some possible extensions of this paper. First, it is worth investi-
 821 gating whether complexity results for the AIPP-S method can be derived for the case
 822 where Y is unbounded. Second, it is worth investigating if the notions of stationary
 823 points in Subsection 2.1 are related to first-order stationary points⁵ of the related
 824 mathematical program with equilibrium constraints:

$$825 \quad \min_{(x,y) \in X \times Y} \{\Phi(x, y) + h(y) : 0 \in \partial[-\Phi(\cdot, y)](x)\}.$$

⁴See Lemma 15 with $f = p_\xi$.

⁵See, for example, [19, Chapter 3].

826 Finally, it remains to be seen if a similar prox-type smoothing scheme can be developed
827 for the case in which assumption (A2) is relaxed to the condition that there exists
828 $m_y > 0$ such that $-\Phi(x, \cdot)$ is m_y -weakly convex for every $x \in X$.

829 **Appendix A.** This appendix contains a description and a result about an ACG
830 variant used in the analysis of [15].

831 Part of the input of the ACG variant, which is described below, consists of a pair
832 of functions (ψ_s, ψ_n) satisfying:

- 833 (i) $\psi_n \in \text{Conv}(\mathcal{Z})$ is μ -strongly convex for some $\mu \geq 0$;
834 (ii) ψ_s is a convex differentiable function on $\text{dom } \psi_n$ whose gradient is L -Lipschitz
835 continuous for some $L > 0$.

837 ACG method

838 **Input:** a scalar pair $(\mu, L) \in \mathbb{R}_{++}^2$, a function pair (ψ_n, ψ_s) , and an initial point
839 $z_0 \in \text{dom } \psi_n$;

- 841 (0) set $y_0 = z_0$, $A_0 = 0$, $\Gamma_0 \equiv 0$ and $j = 0$;
842 (1) compute

$$843 \quad A_{j+1} = A_j + \frac{\mu A_j + 1 + \sqrt{(\mu A_j + 1)^2 + 4L(\mu A_j + 1)A_j}}{2L},$$

$$844 \quad \tilde{z}_j = \frac{A_j}{A_{j+1}} z_j + \frac{A_{j+1} - A_j}{A_{j+1}} y_j,$$

$$845 \quad \Gamma_{j+1}(y) = \frac{A_j}{A_{j+1}} \Gamma_j(y) + \frac{A_{j+1} - A_j}{A_{j+1}} [\psi_s(\tilde{z}_j) + \langle \nabla \psi_s(\tilde{z}_j), y - \tilde{z}_j \rangle] \quad \forall y,$$

$$846 \quad y_{j+1} = \underset{y}{\text{argmin}} \left\{ \Gamma_{j+1}(y) + \psi_n(y) + \frac{1}{2A_{j+1}} \|y - y_0\|^2 \right\},$$

$$847 \quad z_{j+1} = \frac{A_j}{A_{j+1}} z_j + \frac{A_{j+1} - A_j}{A_{j+1}} y_{j+1};$$

- 849 (2) compute

$$850 \quad u_{j+1} = \frac{y_0 - y_{j+1}}{A_{j+1}},$$

$$851 \quad \varepsilon_{j+1} = \psi(z_{j+1}) - \Gamma_{j+1}(y_{j+1}) - \psi_n(y_{j+1}) - \langle u_{j+1}, z_{j+1} - y_{j+1} \rangle;$$

- 853 (3) increment $j = j + 1$ and go to (1).
-

854 We now discuss some implementation details of the ACG method. First, a single
855 iteration requires the evaluation of two distinct types of oracles, namely: (i) the eval-
856 uation of the functions $\psi_n, \psi_s, \nabla \psi_s$ at any point in $\text{dom } \psi_n$; and (ii) the computa-
857 tion of the exact solution of subproblems of the form $\min_y \{ \psi_n(y) + \|y - a\|^2 / (2\alpha) \}$
858 for any $a \in \mathcal{Z}$ and $\alpha > 0$. In particular, the latter is needed in the computation of y_{j+1} .
859 Second, because Γ_{j+1} is affine, an efficient way to store it is in terms of a normal
860 vector and a scalar intercept that is updated recursively at every iteration. Indeed, if
861 $\Gamma_j = \alpha_j + \langle \cdot, \beta_j \rangle$ for some $(\alpha_j, \beta_j) \in \mathbb{R} \times \mathcal{Z}$, then step 1 of the ACG method implies
862 that $\Gamma_{j+1} = \alpha_{j+1} + \langle \cdot, \beta_{j+1} \rangle$ where

$$864 \quad \alpha_{j+1} := \frac{A_j}{A_{j+1}} \alpha_j + \frac{A_{j+1} - A_j}{A_{j+1}} [\psi_s(\tilde{z}_j) - \langle \nabla \psi_j(\tilde{z}_j), \tilde{z}_j \rangle],$$

23

$$\beta_{j+1} := \frac{A_j}{A_{j+1}}\beta_j + \frac{A_{j+1} - A_j}{A_{j+1}} [\nabla\psi_s(\tilde{z}_j)].$$

The following result, whose proof is given in [15, Lemma 9], is used to establish the iteration complexity of obtaining the triple (z, u, ε) in step 1 of the AIPP method of Subsection 3.1.

LEMMA 12. Let $\{(A_j, z_j, u_j, \varepsilon_j)\}$ be the sequence generated by the ACG method. Then, for any $\sigma > 0$, the ACG method obtains a triple (z, u, ε) satisfying

$$(A.1) \quad u \in \partial_\varepsilon(\psi_s + \psi_n)(z) \quad \|u\|^2 + 2\varepsilon \leq \sigma\|z_0 - z + u\|^2$$

in at most $\lceil 2\sqrt{2L}(1 + \sqrt{\sigma})/\sqrt{\sigma} \rceil$ iterations.

Appendix B. This appendix contains results about functions that can be described be as the maximum of a family of differentiable functions.

The technical lemma below, which is a special case of [9, Theorem 10.2.1], presents a key property about max functions.

LEMMA 13. Assume that the triple (X, Y, Ψ) satisfies (A0)–(A1) in Subsection 2.1 with $\Phi = \Psi$. Moreover, define

$$(B.1) \quad q(x) := \sup_{y \in Y} \Psi(x, y), \quad Y(x) := \{y \in Y : \Psi(x, y) = q(x)\}, \quad \forall x \in X.$$

Then, for every $(x, d) \in X \times \mathcal{X}$, it holds that

$$q'(x; d) = \max_{y \in Y(x)} \langle \nabla_x \Psi(x; y), d \rangle.$$

Moreover, if $Y(x)$ reduces to a singleton, say $Y(x) = \{y(x)\}$, then q is differentiable at x and $\nabla q(x) = \nabla_x \Psi(x, y(x))$.

Under assumptions (A0)–(A3) in Subsection 2.1, the next result establishes Lipschitz continuity of the gradient of q . It is worth mentioning that it generalizes related results in [2, Theorem 5.26] (which covers the case where Ψ is bilinear) and [20, Proposition 4.1] (which makes the stronger assumption that $\Psi(\cdot, y)$ is convex for every $y \in Y$).

PROPOSITION 14. Assume that the triple (X, Y, Ψ) satisfies (A0)–(A3) in Subsection 2.1 with $\Phi = \Psi$ and that, for some $\mu > 0$, the function $\Psi(x, \cdot)$ is μ -strongly concave on Y for every $x \in X$, and define

$$(B.2) \quad Q_\mu := \frac{L_y}{\mu} + \sqrt{\frac{L_x + m}{\mu}}, \quad L_\mu := L_y Q_\mu + L_x, \quad y(x) := \operatorname{argmax}_{y \in Y} \Psi(x, y)$$

for every $x \in X$. Then, the following properties hold:

- (a) $y(\cdot)$ is Q_μ -Lipschitz continuous on X ;
- (b) $\nabla q(\cdot)$ is L_μ -Lipschitz continuous on X where q is as in (B.1).

Proof. (a) Let $x, \tilde{x} \in X$ be given and denote $(y, \tilde{y}) = (y(x), y(\tilde{x}))$. Define $\alpha(u) := \Psi(u, y) - \Psi(u, \tilde{y})$ for every $u \in X$, and observe that the optimality conditions of y and \tilde{y} imply that $\alpha(x) \geq \mu\|y - \tilde{y}\|^2/2$ and $-\alpha(\tilde{x}) \geq \mu\|y - \tilde{y}\|^2/2$. Using the previous inequalities, (2.1), (2.2), (2.3), and the Cauchy-Schwarz inequality, we conclude that

$$\mu\|y - \tilde{y}\|^2 \leq \alpha(x) - \alpha(\tilde{x}) \leq \langle \nabla_x \Psi(x, y) - \nabla_x \Psi(x, \tilde{y}), x - \tilde{x} \rangle + \frac{L_x + m}{2} \|x - \tilde{x}\|^2$$

$$\begin{aligned}
903 \quad & \leq \|\nabla_x \Psi(x, y) - \nabla_x \Psi(x, \tilde{y})\| \cdot \|x - \tilde{x}\| + \frac{L_x + m}{2} \|x - \tilde{x}\|^2 \\
904 \quad & \leq L_y \|y - \tilde{y}\| \cdot \|x - \tilde{x}\| + \frac{L_x + m}{2} \|x - \tilde{x}\|^2. \\
905 \quad &
\end{aligned}$$

906 Considering the above as a quadratic inequality in $\|\tilde{y} - y\|$ yields the bound

$$\begin{aligned}
907 \quad & \|y - \tilde{y}\| \leq \frac{1}{2\mu} \left[L_y \|x - \tilde{x}\| + \sqrt{L_y^2 \|x - \tilde{x}\|^2 + 4\mu(L_x + m)\|x - \tilde{x}\|^2} \right] \\
908 \quad & \leq \left[\frac{L_y}{\mu} + \sqrt{\frac{L_x + m}{\mu}} \right] \|x - \tilde{x}\| = Q_\mu \|x - \tilde{x}\| \\
909 \quad &
\end{aligned}$$

910 which is the conclusion of (a).

911 (b) Let $x, \tilde{x} \in X$ be given and denote $(y, \tilde{y}) = (y(x), y(\tilde{x}))$. Using part (a),
912 Lemma 13, and (2.2) we have that

$$\begin{aligned}
913 \quad & \|\nabla q(x) - \nabla q(\tilde{x})\| = \|\nabla_x \Psi(x, y) - \nabla_x \Psi(\tilde{x}, \tilde{y})\| \\
914 \quad & \leq \|\nabla_x \Psi(x, y) - \nabla_x \Psi(x, \tilde{y})\| + \|\nabla_x \Psi(x, \tilde{y}) - \nabla_x \Psi(\tilde{x}, \tilde{y})\| \\
915 \quad & \leq L_y \|y - \tilde{y}\| + L_x \|x - \tilde{x}\| \leq (L_y Q_\mu + L_x) \|x - \tilde{x}\| = L_\mu \|x - \tilde{x}\|, \\
916 \quad &
\end{aligned}$$

917 which is the conclusion of (b). \square

918 **Appendix C.** The main goal of this appendix is to prove Propositions 17 and
919 18, which are used in the proofs of Propositions 1, 2, and 3 given in Appendix D.

920 The following well-known result presents an important property about the direc-
921 tional derivative of a composite function $f + h$.

922 LEMMA 15. Let $h : \mathcal{X} \mapsto (-\infty, \infty]$ be a proper convex function and let f be a
923 differentiable function on $\text{dom } h$. Then, for any $x \in \text{dom } h$, it holds that

$$924 \quad (\text{C.1}) \quad \inf_{\|d\| \leq 1} (f + h)'(x; d) = \inf_{\|d\| \leq 1} [\langle \nabla f(x), d \rangle + \sigma_{\partial h(x)}(d)] = - \inf_{u \in \nabla f(x) + \partial h(x)} \|u\|.$$

925 The proof of Lemma 15 can be found for example in [28, Exercise 8.8(c)]. An
926 alternative and more direct proof is given in [14, Lemma F.1.2]. It is also worth
927 mentioning that if we further assumed that $\text{dom } h = \mathcal{X}$, then the above result would
928 follow from [3, Lemma 5.1].

929 The next technical lemma, which can be found in [29, Corollary 3.3], presents a
930 well-known min-max identity.

931 LEMMA 16. Let a convex set $D \subseteq \mathcal{X}$ and compact convex set $Y \subseteq \mathcal{Y}$ be given.
932 Moreover, let $\psi : D \times Y \mapsto \mathbb{R}$ be a function in which $\psi(\cdot, y)$ is convex lower semicon-
933 tinuous for every $y \in Y$ and $\psi(d, \cdot)$ is concave upper semicontinuous for every $d \in D$.
934 Then,

$$935 \quad \inf_{d \in \mathcal{X}} \sup_{y \in \mathcal{Y}} \psi(d, y) = \sup_{y \in \mathcal{Y}} \inf_{d \in \mathcal{X}} \psi(d, y).$$

936 The next result establishes an identity similar to Lemma 15 but for the case where
937 f is a max function.

938 PROPOSITION 17. Assume the quadruple (Ψ, h, X, Y) satisfies assumptions (A0)–
939 (A3) of Subsection 2.1 with $\Phi = \Psi$. Moreover, suppose that $\Psi(\cdot, y)$ is convex for every
940 $y \in Y$, and let q and $Y(\cdot)$ be as in Lemma 13. Then, for every $\bar{x} \in X$, it holds that

$$941 \quad (\text{C.2}) \quad \inf_{\|d\| \leq 1} (q + h)'(\bar{x}; d) = - \inf_{u \in Q(\bar{x})} \|u\|$$

942 where $Q(\bar{x}) := \partial h(\bar{x}) + \bigcup_{y \in Y(\bar{x})}$. Moreover, if $\partial h(\bar{x})$ is nonempty, then the infimum
 943 on the right-hand side of (C.2) is achieved.

944 *Proof.* Let $\bar{x} \in X$ and define

$$945 \quad (\text{C.3}) \quad \psi(d, y) := (\Psi_y + h)'(\bar{x}; d), \quad \forall (d, x, y) \in \mathcal{X} \times \Omega \times Y.$$

946 We claim that ψ in (C.3) satisfies the assumptions on ψ in Lemma 16 with $Y = Y(\bar{x})$
 947 and D given by

$$948 \quad D := \{d \in \mathcal{Z} : \|d\| \leq 1, d \in F_X(\bar{x})\},$$

949 where $F_X(\bar{x}) := \{t(x - \bar{x}) : x \in X, t \geq 0\}$ is the set of feasible directions at \bar{x} .
 950 Before showing this claim, we use it to show that (C.2) holds. First observe that (A1)
 951 and Lemma 13 imply that $q'(\bar{x}; d) = \sup_{y \in Y} \Psi'_y(\bar{x}; d)$ for every $d \in \mathcal{X}$. Using then
 952 Lemma 16 with $Y = Y(\bar{x})$, Lemma 15 with $(f, x) = (\Psi_{\bar{y}}, \bar{x})$ for every $\bar{y} \in Y(\bar{x})$, and
 953 the previous observation, we have that

$$\begin{aligned} 954 \quad & \inf_{\|d\| \leq 1} (q + h)'(\bar{x}; d) = \inf_{d \in D} (q + h)'(\bar{x}; d) = \inf_{d \in D} \sup_{y \in Y(\bar{x})} (\Psi_y + h)'(\bar{x}; d) \\ 955 \quad & = \inf_{d \in D} \sup_{y \in Y(\bar{x})} \psi(d, y) = \sup_{y \in Y(\bar{x})} \inf_{d \in D} \psi(d, y) = \sup_{y \in Y(\bar{x})} \inf_{\|d\| \leq 1} (\Psi_y + h)'(\bar{x}; d) \\ 956 \quad (\text{C.4}) \quad & = \sup_{y \in Y(\bar{x})} \left[- \inf_{u \in \nabla_x \Phi(\bar{x}, y) + \partial h(\bar{x})} \|u\| \right] = \left[- \inf_{u \in Q(\bar{x})} \|u\| \right]. \\ 957 \end{aligned}$$

958 Let us now assume that $\partial h(\bar{x})$ is nonempty, and hence, $Q(\bar{x})$ is nonempty as well. Note
 959 that continuity of the function $\nabla_x \Psi(\bar{x}, \cdot)$ from assumption (A1) and the compactness
 960 of $Y(\bar{x})$ imply that Q is closed. Moreover, since $\|u\| \geq 0$, it holds that any sequence
 961 $\{u_k\}_{k \geq 1}$ where $\lim_{k \rightarrow \infty} \|u_k\| = \inf_{u \in Q(\bar{x})} \|u\|$ is bounded. Combining the previous
 962 two remarks with the Bolzano-Weierstrass Theorem, we conclude that $\inf_{u \in Q(\bar{x})} \|u\| =$
 963 $\min_{u \in Q(\bar{x})} \|u\|$, and hence (C.2) holds.

964 To complete the proof, we now justify the above claim on ψ . First, for any given
 965 $y \in Y(\bar{x})$, it follows from [27, Theorem 23.1] with $f(\cdot) = \Psi_y(\cdot)$ and the definitions of
 966 q and $Y(\bar{x})$ that

$$967 \quad (\text{C.5}) \quad \psi(d, \bar{y}) = \Psi'_{\bar{y}}(\bar{x}; d) = \inf_{t > 0} \frac{\Psi_y(\bar{x} + td) - q(\bar{x})}{t} \quad \forall d \in \mathcal{X}.$$

968 Since assumption (A2) implies that $\Psi(\bar{x}, \cdot)$ is upper semicontinuous and concave on
 969 Y , it follows from (C.5), [27, Theorem 5.5], and [27, Theorem 9.4] that $\psi(d, \cdot)$ is upper
 970 semicontinuous and concave on Y for every $d \in \mathcal{X}$. On the other hand, since $\Psi(\cdot, y)$
 971 is assumed to be lower semicontinuous and convex on X for every $y \in Y$, it follows
 972 from (C.5), the fact that $\bar{x} \in \text{int} \Omega$, and [27, Theorem 23.4], that $\psi(\cdot, y)$ is lower
 973 semicontinuous and convex on \mathcal{X} , and hence $D \subseteq \mathcal{X}$, for every $y \in Y(\bar{x})$. \square

974 The last technical result is a specialization of the one given in [12, Theorem 4.2.1].

975 **PROPOSITION 18.** *Let a proper closed function $\phi : \mathcal{X} \mapsto (-\infty, \infty]$ and assume*
 976 *that $\phi + \|\cdot\|^2/2\lambda$ is μ -strongly convex for some scalars $\mu, \lambda > 0$. If a quadru-*
 977 *ple $(x^-, x, u, \varepsilon) \in \mathcal{X} \times \text{dom} \phi \times \mathcal{X} \times \mathbb{R}_+$ together with λ satisfy the inclusion $u \in$*
 978 *$\partial_\varepsilon (\phi + \|\cdot - x^-\|^2/[2\lambda])(x)$, then the point $\hat{x} \in \text{dom} \phi$ given by*

$$979 \quad (\text{C.6}) \quad \hat{x} := \underset{x'}{\text{argmin}} \left\{ \phi_\lambda(x') := \phi(x') + \frac{1}{2\lambda} \|x' - x^-\|^2 - \langle u, x' \rangle \right\}$$

980 *satisfies*

$$981 \quad (\text{C.7}) \quad \inf_{\|d\| \leq 1} \phi'(\hat{x}; d) \geq -\frac{1}{\lambda} \|x^- - x + \lambda u\| - \sqrt{\frac{2\varepsilon}{\lambda^2 \mu}}, \quad \|\hat{x} - x\| \leq \sqrt{\frac{2\varepsilon}{\mu}}.$$

982 *Proof.* We first observe that the assumed inclusion implies that $\phi_\lambda(x') \geq \phi_\lambda(x) - \varepsilon$
 983 for every $x' \in X$. Using the previous inequality at $x' = \hat{x}$, the optimality of \hat{x} , and
 984 the μ -strong convexity of ϕ_λ , we have that $\mu \|\hat{x} - x\|^2 / 2 \leq \phi_\lambda(x) - \phi_\lambda(\hat{x}) \leq \varepsilon$ from
 985 which we conclude that $\|\hat{x} - x\| \leq \sqrt{2\varepsilon/\mu}$, i.e., the second inequality in (C.7).

986 To show the other inequality, let $n_\lambda := x^- - x + \lambda u$. Using the definition of ϕ_λ ,
 987 the triangle inequality, and the previous bound on $\|\hat{x} - x\|$, we obtain

$$988 \quad 0 \leq \inf_{\|d\| \leq 1} \phi'_\lambda(\hat{x}; d) = \inf_{\|d\| \leq 1} \phi'(\hat{x}; d) - \frac{1}{\lambda} \langle d, n_\lambda \rangle$$

$$989 \quad (\text{C.8}) \quad \leq \inf_{\|d\| \leq 1} \phi'(\hat{x}; d) + \frac{\|n_\lambda\|}{\lambda} + \frac{\|x - \hat{x}\|}{\lambda} \leq \inf_{\|d\| \leq 1} \phi'(\hat{x}; d) + \frac{\|n_\lambda\|}{\lambda} + \sqrt{\frac{2\varepsilon}{\lambda^2 \mu}},$$

991 which clearly implies the first inequality in (C.7). \square

992 **Appendix D.** This appendix presents the proofs of Propositions 1, 2, and 3.

993 The first technical result shows that an approximate primal-dual stationary point
 994 is equivalent to an approximate directional stationary point of a perturbed version of
 995 problem (1.1).

996 **LEMMA 19.** *Suppose the quadruple (Φ, h, X, Y) satisfies assumptions (A0)–(A3)*
 997 *of Subsection 2.1 and let $(\bar{x}, \bar{u}, \bar{v}) \in X \times \mathcal{X} \times \mathcal{Y}$ be given. Then, there exists $\bar{y} \in Y$*
 998 *such that the quadruple $(\bar{u}, \bar{v}, \bar{x}, \bar{y})$ satisfies the inclusion in (1.4) if and only if*

$$999 \quad (\text{D.1}) \quad \inf_{\|d\| \leq 1} (p_{\bar{u}, \bar{v}} + h)'(\bar{x}; d) \geq 0,$$

1000 where $p_{\bar{u}, \bar{v}} := \max_{y \in Y} [\Phi(x, y) + \langle \bar{v}, y \rangle - \langle \bar{u}, x \rangle]$ for every $x \in \Omega$.

1001 *Proof.* Let $(\bar{x}, \bar{u}, \bar{v}) \in X \times \mathcal{X} \times \mathcal{Y}$ be given, define

$$1002 \quad (\text{D.2}) \quad \Psi(x, y) := \Phi(x, y) + \langle \bar{v}, y \rangle - \langle \bar{u}, x \rangle + m \|x - \bar{x}\|^2 \quad \forall (x, y) \in \Omega \times Y,$$

1004 and let q and $Y(\cdot)$ be as in Lemma 13. It is easy to see that $q = p_{\bar{u}, \bar{v}}$, the function
 1005 Ψ satisfies the assumptions on Ψ in Proposition 17, and \bar{x} satisfies (D.1) if and only
 1006 if $\inf_{\|d\| \leq 1} (q + h)'(\bar{x}; d) \geq 0$. The desired conclusion follows from Proposition 17, the
 1007 previous observation, and the fact that $\bar{y} \in Y(\bar{x})$ if and only if $\bar{v} \in \partial[-\Phi(\bar{x}, \cdot)](\bar{y})$. \square

1008 We are now ready to give the proof of Proposition 1.

1009 *Proof of Proposition 1.* Suppose $(\bar{u}, \bar{v}, \bar{x}, \bar{y})$ is a (ρ_x, ρ_y) -primal-dual stationary
 1010 point of (1.1). Moreover, let Ψ , q , and D_y be as in (D.2), (B.1) and (2.8), respectively,
 1011 and define

$$1012 \quad \hat{q}(x) := q(x) + h(x) \quad \forall x \in X.$$

1013 Using Lemma 19, we first observe that $\inf_{\|d\| \leq 1} \hat{q}(\bar{x}; d) \geq 0$. Since \hat{q} is convex from
 1014 assumption (A3), it follows from the previous bound and Lemma 15 with $(f, h) =$
 1015 $(0, \hat{q})$, that $\min_{u \in \partial \hat{q}(\bar{x})} \|u\| \leq 0$, and hence, $0 \in \partial \hat{q}(\bar{x})$. Moreover, using the Cauchy-
 1016 Schwarz inequality, the second inequality in (1.4), the previous inclusion, and the
 1017 definition of q and Ψ , it follows that for every $x \in \mathcal{X}$,

$$1018 \quad \hat{p}(x) + D_y \rho_y - \langle \bar{u}, x \rangle + m \|x - \bar{x}\|^2 \geq \hat{q}(x) \geq \hat{q}(\bar{x}) \geq \hat{p}(\bar{x}) - D_y \rho_y - \langle \bar{u}, \bar{x} \rangle,$$

1019

1020 and hence that $\bar{u} \in \partial_\varepsilon(\hat{p} + m\|\cdot - \bar{x}\|^2)(\bar{x})$ where $\varepsilon = 2D_y\rho_y$. Using now the first
 1021 inequality in (1.4), Proposition 18 with $(\phi, x, x^-, u) = (\hat{p}, \bar{x}, \bar{x}, \bar{u})$ and also $(\varepsilon, \lambda, \mu) =$
 1022 $(D_y\rho_y, 1/(2m), m)$, we conclude that there exists \hat{x} such that $\|\hat{x} - \bar{x}\| \leq \sqrt{2D_y\rho_y/m}$
 1023 and

$$1024 \quad \inf_{\|d\| \leq 1} \hat{p}'(\hat{x}; d) \geq -\|\bar{u}\| - 2\sqrt{2mD_y\rho_y} \geq -\rho_x - 2\sqrt{2mD_y\rho_y}. \quad \square$$

1025 We next give the proof of Proposition 2.

1026 *Proof of Proposition 2.* (a) We first claim that \hat{P}_λ is α -strongly convex, where
 1027 $\alpha = 1/\lambda - m$. To see this, note that $\Phi(\cdot, y) + m\|\cdot\|^2/2$ is convex for every $y \in Y$
 1028 from assumption (A3). The claim now follows from assumption (A2), the fact that
 1029 the supremum of a collection of convex functions is also convex, and the definition of
 1030 \hat{p} in (1.1).

1031 Suppose the pair (x, δ) satisfies (1.5) and (2.10). If $\hat{x} = x_\lambda$ in (1.5), then clearly
 1032 the second inequality in (1.5), the fact that $\lambda < 1/m$, and (2.10) imply the inequality
 1033 in (2.9), and hence, that x is a (λ, ε) -prox stationary point. Suppose now that $\hat{x} \neq x_\lambda$.
 1034 Using the convexity of \hat{P}_λ , we first have that $\hat{P}'_\lambda(\hat{x}; d) = \inf_{t>0} [\hat{P}_\lambda(\hat{x} + td) - \hat{P}_\lambda(\hat{x})] / t$
 1035 for every $d \in \mathcal{X}$. Denoting $n_\lambda := (x_\lambda - \hat{x})/\|x_\lambda - \hat{x}\|$, using both inequalities in (1.5)
 1036 and the previous identity, we then have that

$$1037 \quad \frac{\hat{P}_\lambda(x_\lambda) - \hat{P}_\lambda(\hat{x})}{\|x_\lambda - \hat{x}\|} \geq \hat{p}'(\hat{x}; n_\lambda) + \left\langle \frac{n_\lambda}{\lambda}, \hat{x} - x \right\rangle \geq -\delta - \frac{\|\hat{x} - x\|}{\lambda} \geq -\delta \left(\frac{1 + \lambda}{\lambda} \right).$$

1039 Using the optimality of x_λ , the α -strong convexity of \hat{P}_λ (see our claim on \hat{p} in the
 1040 first paragraph), and the above bound, we conclude that

$$1041 \quad \frac{1}{2\alpha} \|\hat{x} - x_\lambda\|^2 \leq \hat{P}_\lambda(\hat{x}) - \hat{P}_\lambda(x_\lambda) \leq \delta \left(\frac{1 + \lambda}{\lambda} \right) \|\hat{x} - x_\lambda\|.$$

1042 Thus, $\|\hat{x} - x_\lambda\| \leq 2\alpha\delta(1 + \lambda)/\lambda$. Using the previous bound, the second inequality in
 1043 (1.5), and (2.10) yields

$$1044 \quad \|x - x_\lambda\| \leq \|x - \hat{x}\| + \|\hat{x} - x_\lambda\| \leq \left(1 + 2\alpha \left[\frac{1 + \lambda}{\lambda} \right] \right) \delta \leq \lambda\varepsilon,$$

1045 which implies (2.9), and hence, that x is a (λ, ε) -prox stationary point.

1046 (b) Suppose that the point x is a (λ, ε) -prox stationary point with $\varepsilon \leq \delta \cdot$
 1047 $\min\{1, 1/\lambda\}$. Then the optimality of x_λ , the fact that \hat{P}_λ is convex (see the be-
 1048 ginning of part (a)), the inequality in (2.9), and the Cauchy-Schwarz inequality imply
 1049 that

$$1050 \quad 0 \leq \inf_{\|d\| \leq 1} \left[\hat{p}'(x_\lambda; d) + \frac{1}{\lambda} \langle d, x_\lambda - x \rangle \right] \leq \inf_{\|d\| \leq 1} \hat{p}'(x_\lambda; d) + \varepsilon \leq \inf_{\|d\| \leq 1} \hat{p}'(x_\lambda; d) + \delta,$$

1051 which, together with the fact that $\lambda\varepsilon \leq \delta$, imply that x satisfies (1.5) with $\hat{x} = x_\lambda$. \square

1052 Finally, we give the proof of Proposition 3.

1053 *Proof of Proposition 3.* This follows by using Lemma 15 with $(f, h) = (\Phi(\cdot, \bar{y}), h)$
 1054 and $(f, h) = (0, -\Phi(\bar{x}, \cdot))$. \square

1055

REFERENCES

- 1056 [1] K. ARROW, L. HURWICZ, AND H. UZAWA, *Studies in linear and non-linear programming*, Cam-
1057 bridge Univ. Press, 1958.
- 1058 [2] A. BECK, *First-order methods in optimization*, SIAM, 2017.
- 1059 [3] J. V. BURKE AND J. J. MORÉ, *On the identification of active constraints*, SIAM Journal on
1060 Numerical Analysis, 25 (1988), pp. 1197–1211.
- 1061 [4] Y. CARMON, J. DUCHI, O. HINDER, AND A. SIDFORD, *Accelerated methods for non-convex*
1062 *optimization*, Available on arXiv:1611.00756, (2017).
- 1063 [5] D. DAVIS AND D. DRUSVYATSKIY, *Stochastic model-based minimization of weakly convex func-*
1064 *tions*, SIAM Journal on Optimization, 29 (2019), pp. 207–239.
- 1065 [6] D. DAVIS, D. DRUSVYATSKIY, K. J. MACPHEE, AND C. PAQUETTE, *Subgradient methods for*
1066 *sharp weakly convex functions*, Journal of Optimization Theory and Applications, 179
1067 (2018), pp. 962–982.
- 1068 [7] D. DRUSVYATSKIY AND C. PAQUETTE, *Efficiency of minimizing compositions of convex func-*
1069 *tions and smooth maps*, Mathematical Programming, 178 (2019), pp. 503–558.
- 1070 [8] J. C. DUCHI AND F. RUAN, *Stochastic methods for composite and weakly convex optimization*
1071 *problems*, SIAM Journal on Optimization, 28 (2018), pp. 3229–3259.
- 1072 [9] F. FACCHINEI AND J.-S. PANG, *Finite-dimensional variational inequalities and complementarity*
1073 *problems*, Springer Science & Business Media, 2007.
- 1074 [10] S. GHADIMI AND G. LAN, *Accelerated gradient methods for nonconvex nonlinear and stochastic*
1075 *programming*, Math. Program., 156 (2016), pp. 59–99.
- 1076 [11] Y. HE AND R. D. C. MONTEIRO, *An accelerated HPE-type algorithm for a class of composite*
1077 *convex-concave saddle-point problems*, SIAM J. Optim., 26 (2016), pp. 29–56.
- 1078 [12] J. HIRIART-URRUTY AND C. LEMARECHAL, *Convex Analysis and Minimization Algorithms II*,
1079 Springer, Berlin, 1993.
- 1080 [13] O. KOLOSSOSKI AND R. D. C. MONTEIRO, *An accelerated non-euclidean hybrid proximal*
1081 *extragradient-type algorithm for convex-concave saddle-point problems*, Optim. Methods
1082 Softw., 32 (2017), pp. 1244–1272.
- 1083 [14] W. KONG, *Accelerated inexact first-order methods for solving nonconvex composite optimiza-*
1084 *tion problems*, arXiv preprint arXiv:2104.09685, (2021).
- 1085 [15] W. KONG, J. G. MELO, AND R. D. C. MONTEIRO, *Complexity of a quadratic penalty accel-*
1086 *erated inexact proximal point method for solving linearly constrained nonconvex composite*
1087 *programs*, SIAM Journal on Optimization, 29 (2019), pp. 2566–2593.
- 1088 [16] W. KONG, J. G. MELO, AND R. D. C. MONTEIRO, *An efficient adaptive accelerated inex-*
1089 *act proximal point method for solving linearly constrained nonconvex composite problems*,
1090 Computational Optimization and Applications, 76 (2020), pp. 305–346.
- 1091 [17] T. LIN, C. JIN, AND M. JORDAN, *Near-optimal algorithms for minimax optimization*, arXiv
1092 preprint arXiv:2002.02417, (2020).
- 1093 [18] S. LU, I. TSAKNAKIS, AND M. HONG, *Block alternating optimization for non-convex min-max*
1094 *problems: Algorithms and applications in signal processing and communications*, ICASSP
1095 2019 - 2019 IEEE International Conference on Acoustics, Speech and Signal Processing
1096 (ICASSP), (2019), pp. 4754–4758.
- 1097 [19] Z.-Q. LUO, J.-S. PANG, AND D. RALPH, *Mathematical programs with equilibrium constraints*,
1098 Cambridge University Press, 1996.
- 1099 [20] R. D. C. MONTEIRO AND B. F. SVAITER, *Convergence rate of inexact proximal point meth-*
1100 *ods with relative error criteria for convex optimization*, submitted to SIAM Journal on
1101 Optimization, (2010).
- 1102 [21] A. NEMIROVSKI, *Prox-method with rate of convergence $O(1/t)$ for variational inequalities with*
1103 *lipschitz continuous monotone operators and smooth convex-concave saddle point problems*,
1104 SIAM J. Optim., 15 (2004), pp. 229–251.
- 1105 [22] A. NEMIROVSKI AND D. YUDIN, *Cesari convergence of the gradient method of approximating*
1106 *saddle points of convex-concave functions*, in Dokl. Akad. Nauk, vol. 239, Russian Academy
1107 of Sciences, 1978, pp. 1056–1059.
- 1108 [23] Y. NESTEROV, *Smooth minimization of non-smooth functions*, Math. Program., 103 (2005),
1109 pp. 127–152.
- 1110 [24] M. NOUIEHED, M. SANJABI, T. HUANG, J. LEE, AND M. RAZAVIYAYN, *Solving a class of non-*
1111 *convex min-max games using iterative first order methods*, in Advances in Neural Infor-
1112 mation Processing Systems, 2019, pp. 14905–14916.
- 1113 [25] D. M. OSTROVSKII, A. LOWY, AND M. RAZAVIYAYN, *Efficient search of first-order nash equi-*
1114 *libria in nonconvex-concave smooth min-max problems*, arXiv preprint arXiv:2002.07919,
1115 (2020).
- 1116 [26] H. RAFIQUE, M. LIU, Q. LIN, AND T. YANG, *Non-convex min-max optimization: Provable*
1117 *algorithms and applications in machine learning*, arXiv e-prints, (2018).

- 1118 [27] R. T. ROCKAFELLAR, *Convex Analysis*, Princeton University Press, Princeton, 1970.
1119 [28] R. T. ROCKAFELLAR AND R. J.-B. WETS, *Variational analysis*, vol. 317, Springer Science &
1120 Business Media, 2009.
1121 [29] M. SION, *On general minimax theorems.*, Pacific Journal of mathematics, 8 (1958), pp. 171–176.
1122 [30] K. K. THEKUMPARAMPIL, P. JAIN, P. NETRAPALLI, AND S. OH, *Efficient algorithms for*
1123 *smooth minimax optimization*, in Advances in Neural Information Processing Systems,
1124 2019, pp. 12680–12691.