

1 **COMPLEXITY-OPTIMAL AND PARAMETER-FREE FIRST-ORDER**  
2 **METHODS FOR FINDING STATIONARY POINTS OF COMPOSITE**  
3 **OPTIMIZATION PROBLEMS\***

4 WEIWEI KONG<sup>†</sup>

5 **Abstract.** This paper develops and analyzes an accelerated proximal descent method for finding  
6 stationary points of nonconvex composite optimization problems. The objective function is of the  
7 form  $f + h$  where  $h$  is a proper closed convex function,  $f$  is a differentiable function on the domain  
8 of  $h$ , and  $\nabla f$  is Lipschitz continuous on the domain of  $h$ . The main advantage of this method is  
9 that it is “parameter-free” in the sense that it does not require knowledge of the Lipschitz constant  
10 of  $\nabla f$  or of any global topological properties of  $f$ . It is shown that the proposed method can  
11 obtain an  $\varepsilon$ -approximate stationary point with iteration complexity bounds that are optimal, up  
12 to logarithmic terms over  $\varepsilon$ , in both the convex and nonconvex settings. Some discussion is also  
13 given about how the proposed method can be leveraged in other existing optimization frameworks,  
14 such as min-max smoothing and penalty frameworks for constrained programming, to create more  
15 specialized parameter-free methods. Finally, numerical experiments are presented to support the  
16 practical viability of the method.

17 **Key words.** nonconvex composite optimization, first-order accelerated gradient method, iteration  
18 complexity, inexact proximal point method, parameter-free, adaptive, optimal complexity

19 **AMS subject classifications.** 47J22, 65K10, 90C25, 90C26, 90C30, 90C60

20 **1. Introduction.** Consider the nonsmooth composite optimization problem

21 (1.1) 
$$\phi_* = \min_{z \in \mathbb{R}^n} \{\phi(z) := f(z) + h(z)\}$$

22 where  $h : \mathbb{R}^n \mapsto (\infty, \infty]$  is a proper closed convex function,  $f$  is a (possibly noncon-  
23 vex) continuously differentiable function on an open set containing the domain of  $h$   
24 (denoted as  $\text{dom } h$ ), and  $\nabla f$  is Lipschitz continuous. It is well known that the above  
25 assumption on  $f$  implies the existence of positive scalars  $m$  and  $M$  such that

26 (1.2) 
$$-\frac{m}{2}\|x - x'\|^2 \leq f(x) - f(x') - \langle \nabla f(x'), x - x' \rangle \leq \frac{M}{2}\|x - x'\|^2$$

27 for every  $x, x' \in \text{dom } h$ . The quantity  $(m, M)$  is often called a *curvature pair* of  $\phi$  (see,  
28 for example, [24, 25]), and the first inequality of (1.2) is often called *weak-convexity*  
29 when  $m > 0$  (see, for example, [8, 9]).

30 Recently, there has been a surge of interest in developing efficient algorithms  
31 for finding  $\varepsilon$ -stationary points of (1.1), which consist of a pair  $(\bar{z}, \bar{v}) \in \text{dom } h \times \mathbb{R}^n$   
32 satisfying

33 (1.3) 
$$\bar{v} \in \nabla f(\bar{z}) + \partial h(\bar{z}), \quad \|\bar{v}\| \leq \varepsilon.$$

34 While complexity-optimal algorithms exist for the case where both  $m$  and  $M$  are  
35 known, a *parameter-free* algorithm — one without knowledge of  $(m, M)$  — with op-  
36 timal iteration complexity remains elusive.

---

\***Funding:** The work of this author has been supported by the Exascale Computing Project (17-SC-20-SC), a collaborative effort of the U.S. Department of Energy Office of Science and the National Nuclear Security Administration.

**Versions:** v1.0 (May 25, 2022), v2.0 (February 11, 2023), v3.0 (September 17, 2023), v4.0 (February 9, 2024)

<sup>†</sup>Work done at the Computer Science and Mathematics Division, Oak Ridge National Laboratory, Oak Ridge, TN, 37830. Email: [wwkong92@gmail.com](mailto:wwkong92@gmail.com)

37 Our goal in this paper is to develop, analyze, and extend a parameter-free *accel-*  
 38 *erated proximal descent* (PF.APD) algorithm that obtains, up-to-logarithmic terms,  
 39 optimal iteration complexities regardless of the convexity of  $f$ . Roughly speaking,  
 40 PF.APD generates a sequence of iterates  $\{(z_k, m_k)\} \subseteq \text{dom } h \times \mathbb{R}_{++}$  which satisfies

$$41 \quad (1.4) \quad z_{k+1} \approx \underset{z \in \text{dom } h}{\text{argmin}} \left\{ \frac{\phi(z)}{2m_{k+1}} + \frac{1}{2} \|z - z_k\|^2 \right\}, \quad \phi(z_{k+1}) \leq \phi(z_k).$$

42 for every  $k \geq 0$ . Notice that the first expression in (1.4) is an inexact *proximal* point  
 43 update with stepsize  $1/(2m_{k+1})$ , while the inequality in (1.4) implies  $\{\phi(z_k)\}$  is a  
 44 *descent* sequence. More precisely, the  $(k+1)$ -th iteration of PF.APD is as follows:

45 **Iteration  $k+1$ :**

- (i) Given  $\hat{m} \in \mathbb{R}_{++}$ , find a *proximal descent* point  $z_{k+1} \in \text{dom } h$  in which there exists  $\hat{u} \in \mathbb{R}^n$  satisfying

$$(1.5) \quad \hat{u} \in \nabla f(z_{k+1}) + \partial(h + \hat{m}\|\cdot - z_k\|^2)(z_{k+1}),$$

$$(1.6) \quad \|\hat{u} + \hat{m}(z_k - z_{k+1})\|^2 \leq 2\theta\hat{m}[\phi(z_k) - \phi(z_{k+1})],$$

$$(1.7) \quad \|\hat{u}\|^2 \leq \hat{m}^2 \|z_{k+1} - z_k\|^2,$$

for some  $\theta > 0$ .

- (ii) If a key inequality fails during the execution of step (i), change  $\hat{m}$  and try step (i) again. Else, set  $m_{k+1} = \hat{m}$ .

46 To find  $z_{k+1}$  in step (i) in the above outline, PF.APD specifically applies a  
 47 parameter-free *accelerated* composite gradient (PF.ACG) algorithm to the subprob-  
 48 lem  $\min_{z \in \text{dom } h} \{\phi(z)/(2\hat{m}) + \|z - z_k\|^2/2\}$  until a finite set of key descent inequalities  
 49 holds. During the execution of PF.ACG, several inequalities are also checked to en-  
 50 sure its convergence (specifically the ones in (3.5)), and execution is halted if at least  
 51 one of these inequalities does not hold. These inequalities are always guaranteed to  
 52 hold when  $\hat{m} \geq m$  but may fail to hold when  $\hat{m} < m$ .

53 It is worth mentioning that the main difficulties preventing the extension of exist-  
 54 ing complexity-optimal methods to parameter-free ones is their dependence on *global*  
 55 topological conditions that strongly depend on the knowledge of  $(m, M)$ , e.g., (1.2),  
 56 convexity of  $f$ , or knowledge of the Lipschitz modulus of  $\nabla f$ . Hence, one of the novel-  
 57 ties of PF.APD is its ability to relax these conditions to a finite set of *local* topological  
 58 conditions that only depend on the generated sequence of iterates.

59 **1.1. Literature Review.** To keep our notation concise, we will make use of

$$60 \quad (1.8) \quad \Delta_0 := \phi(z_0) - \inf_{z \in \mathbb{R}^n} \phi(z), \quad d_0 := \inf_{z_* \in \mathbb{R}^n} \left\{ \|z_0 - z_*\| : \phi(z_*) = \inf_{z \in \mathbb{R}^n} \phi(z) \right\},$$

61 with the assumption that  $\Delta_0 < \infty$  but  $d_0$  may be infinite. Furthermore, we break  
 62 our discussion between the convex and nonconvex settings and between two types of  
 63 methods:

64 *I.* Algorithms that find  $\hat{z} \in \text{dom } h$  satisfying  $\phi(\hat{z}) - \inf_{z \in \mathbb{R}^n} \phi(z) \leq \varepsilon$ ;

65 *II.* Algorithms that find  $\bar{z} \in \text{dom } h$  satisfying  $\text{dist}(0, \nabla f(\bar{z}) + \partial h(\bar{z})) \leq \varepsilon$ .

66 It is worth mentioning that complexity-optimal *type-I* methods are not necessarily  
 67 complexity-optimal *type-II* methods, as noted in [34].

68 **Convex Setting.** For this discussion, we assume  $\phi$  to be convex. Paper [32]  
 69 presents the first complexity-optimal *type-I* methods, under the assumption that  
 70  $\max\{m, M\}$  is known. Papers [14, 15, 35, 38] (resp. paper [39]) present parameter-free  
 71 complexity-optimal *type-I* methods for the case of  $h \equiv 0$  (resp.  $h$  being the indicator  
 72 of a closed convex set). Paper [1] extends the method in [39] to another parameter-free  
 73 complexity-optimal *type-I* method for general convex functions  $h$ .

74 The regularized accelerated method described in [34] is one of the earliest nearly-  
 75 optimal (up to logarithmic terms) *type-II* methods for the case of  $h \equiv 0$ . However, its  
 76 complexity is obtained under the strong assumption that: (i)  $\max\{m, M\}$  is known,  
 77 (ii) that there exists  $z_* \in \text{dom } h$  such that  $\phi(z_*) = \inf_{z \in \mathbb{R}^n} \phi(z)$ , (iii) and that a lower  
 78 bound for  $d_0$  is known. Whether a parameter-free complexity-optimal *type-II* method  
 79 exists in the convex setting is still unknown.

80 **Nonconvex Setting.** For this discussion, we assume  $\phi$  to be nonconvex. One  
 81 of the most well-known parameter-free *type-II* algorithms is the proximal gradient  
 82 descent (PGD) method with backtracking line search. In [35], it was shown that  
 83 this method has a  $\mathcal{O}(\varepsilon^{-2})$  *type-II* complexity bound when  $f$  is weakly-convex and a  
 84 suboptimal  $\mathcal{O}(\varepsilon^{-1})$  *type-II* bound when  $f$  is convex.

85 One of the earliest accelerated *type-II* methods is found in [12] under the assump-  
 86 tion that  $\text{dom } h$  is bounded. Following this, paper [13] presented a parameter-free  
 87 extension of the method in [12] that handles Hölder continuous gradients of  $f$ . In  
 88 a separate line of research, [25] presented a *type-II* accelerated method whose main  
 89 steps are variants of the (accelerated) FISTA algorithm in [5] and assumes  $\text{dom } h$  is  
 90 bounded. A variant of this method, with improved iteration complexity bounds in  
 91 the convex setting, was examined in [43]. It is worth noting that some of the methods  
 92 in [12, 13, 25, 43] have optimal *type-I* bounds when  $f$  is convex but all the methods  
 93 have suboptimal *type-II* bounds even when  $f$  is convex.

94 Motivated by the developments in [12], other papers, e.g., [6, 10, 23, 40], developed  
 95 similar accelerated methods under different assumptions on  $f$  and  $h$ . Recently, [18]  
 96 proposed a parameter-dependent accelerated inexact proximal point (AIPP) method  
 97 that has an optimal iteration complexity bound of  $\mathcal{O}(\sqrt{Mm}\Delta_0/\varepsilon^2)$  when  $f$  is weakly  
 98 convex but has no advantage when  $f$  is convex. The work in [19] proposed an adap-  
 99 tive version of AIPP where  $(m, M)$  were estimated locally, but a lower bound for  
 100  $\max\{m, M\}$  was still required. A version of [18] in which the outer proximal point  
 101 scheme is replaced with an accelerated one was examined in [24], in which a moder-  
 102 ately worse iteration-complexity bound was established.

103 **Tangentially Related Works.** The developments in [17, 18, 21] strongly influ-  
 104 enced and motivated the technical developments of both PF.ACG and PF.APD. Since  
 105 PF.APD shares strong similarities with AIPP in [18], we mention one of the former's  
 106 technical improvements on the latter. To begin, note that AIPP is a double-loop  
 107 method that repeatedly calls an ACG-type method on a sequence of prox subprob-  
 108 lems to generate a sequence of *outer* iterates  $\{(z_k, v_k, \varepsilon_k)\}$  (at the end of each ACG  
 109 call) satisfying

$$110 \quad (1.9) \quad v_k \in \partial_{\varepsilon_k} \left( \frac{\phi}{2m} + \frac{1}{2} \|\cdot - z_{k-1}\|^2 \right) (z_k), \quad \|v_k\|^2 + 2\varepsilon_k \leq \sigma^2 \|v_k + z_{k-1} - z_k\|^2,$$

111 where  $\sigma \in (0, 1)$  and  $\partial_{\varepsilon} \psi(x) := \{u \in \mathbb{R}^n : \psi(z') \geq \psi(z) + \langle u, z' - z \rangle - \varepsilon, \quad \forall z' \in \mathbb{R}^n\}$ .  
 112 An expensive refinement procedure, whose effectiveness strongly depends on (1.9)  
 113 and knowledge of  $\max\{m, M\}$ , is then applied to each  $(z_k, v_k, \varepsilon_k)$  to obtain  $(\bar{z}, \bar{v})$

114 satisfying the inclusion in (1.3). In contrast, the iterates generated at every *inner*  
 115 iteration of PF.APD always satisfy the inclusion in (1.3), for a different choice of  
 116  $\bar{v}$  (see Lemma 3.3), and, consequently, the termination of PF.APD can be checked  
 117 at *every* one of its inner iterations *without* the need for an expensive refinement  
 118 procedure. It is worth mentioning those relative prox-stationarity criteria, such as  
 119 (1.7) and (1.9), were previously analyzed in [42] and, more recently, in [2, 26, 28–31].

120 We now make a brief comparison between PF.APD and two adaptive proximal  
 121 methods in the literature. First, compared to the redistributed prox-bundle (RPB)  
 122 method in [16], both PF.APD and RPB are double-loop methods consisting of (i) outer  
 123 (or “serious”) iterations that consider prox-subproblems of the form in (1.4) and some  
 124  $\lambda > 0$  and (ii) inner (or “null”) iterations that consider composite subproblems of the  
 125 form  $\min_{y \in \mathbb{R}^n} \{\Phi_{j,k}(y) + h(y)\}$  for the  $k$ -th subproblem and  $j$ -th iteration, until there  
 126 is a sufficient decrease in  $\phi(z_k)$ . However, PF.APD chooses  $\Phi_{j,k}$  to be a quadratic  
 127 approximation of  $\Phi_k$  centered on a specially chosen point (see the update of  $y_{k+1}$   
 128 in Algorithm 3.1), while RPB chooses  $\Phi_{j,k}$  to be the maximum of a different set of  
 129 quadratic approximations, which is generally more difficult to minimize. Moreover,  
 130 PF.APD uses values of  $\nabla f(\cdot)$  and elements of  $\partial h(\cdot)$  in its construction of  $\Phi_{j,k}$  whereas  
 131 RPB uses elements of the limiting subdifferential of  $\phi$ .

132 Second, compared to the Catalyst Acceleration Framework (CAF) in [40], both  
 133 PF.APD and CAF consider inexactly solving proximal subproblems as in (1.4) us-  
 134 ing ACG subroutine and subproblem termination conditions similar to (2.3)–(2.4).  
 135 However, CAF obtains the inequality in (1.4) by inexactly solving a second prox-  
 136 subproblem (with a different prox center) and applying an extra interpolation step.  
 137 As a consequence, CAF requires nearly double the work of PF.APD. Moreover, the  
 138 line search strategy (analogous to Algorithm 3.1 and Algorithm 3.3) employed by CAF  
 139 in [40, Algorithm 3] is static in that it prescribes a large number of ACG iterations,  
 140 whereas the line search strategy in PF.APD is dynamic in that it checks a finite set  
 141 of simple inequalities at each ACG iteration.

142 **1.2. Contributions.** Throughout, we refer to the two types of algorithms de-  
 143 scribed in the previous subsection. Given a starting point  $z_0 \in \text{dom } h$  and a tolerance  
 144  $\varepsilon > 0$ , it is shown that PF.APD has the following nice properties:

- 145 (i) for any  $\hat{m} > 0$ , it always obtains a pair  $(\bar{z}, \bar{v}) \in \text{dom } h \times \mathbb{R}^n$  satisfying (1.3);
  - 146 (ii) if  $f$  is nonconvex, then it stops in  $\tilde{O}(\sqrt{mM\Delta_0}/\varepsilon^2)$  resolvent evaluations<sup>1</sup>;
  - 147 (iii) if  $f$  is convex, then it stops in  $\tilde{O}(\sqrt{M} \min\{\sqrt{\Delta_0}/\varepsilon, d_0/\sqrt{\varepsilon}\})$  resolvent evaluations;
- 148 Both of the above complexity bounds are optimal (up to logarithmic terms in) in terms  
 149 of  $\Delta_0$ ,  $M$ ,  $m$ , and  $\varepsilon$  (although suboptimal by a factor of  $\sqrt{d_0}$  in the convex case).  
 150 Moreover, it appears to be the first time that a *type-II* parameter-free method has  
 151 obtained such bounds<sup>2</sup>. Improved iteration complexity bounds are also obtained when  
 152  $d_0$  is known. Also, all of the above results are obtained under the mild assumption  
 153 that the optimal value in (1.1) is finite and does not assume the boundedness of  $\text{dom } h$   
 154 (cf. [25, 43]) nor that an optimal solution of (1.1) exists.

155 For convenience, we compare in Table 1.1 the best iteration complexity bounds  
 156 of some of the parameter-free methods listed in the previous subsection with two in-  
 157 stances of PF.APD. For shorthands, PGD is the adaptive proximal gradient descent  
 158 method in [35], UPF is the UPFAG method in [13], ANCF is the ADAP-NC-FISTA  
 159 method in [25], VRF is the VAR-FISTA method in [43], and APD is as in Algo-

<sup>1</sup>The notation  $\tilde{O}(\cdot)$  ignores any terms that logarithmically depend on the tolerance  $\varepsilon$ .

<sup>2</sup>Compare this to the complexity-optimal methods in [34] and [18] which require knowledge of  $d_0$  and  $(m, M)$ , respectively.

160 **rithm 3.4** in this paper with  $m_0 = 1$ .

Algorithm	$f$ convex	$f$ nonconvex	$D_h < \infty$
PGD [35]	$\mathcal{O}\left(\frac{M^{3/2}d_0}{\varepsilon}\right)$	$\mathcal{O}\left(\frac{M^2\Delta_0}{\varepsilon^2}\right)$	No
UPF [13]	N/A	$\mathcal{O}\left(\frac{M\Delta_0}{\varepsilon^2}\right)$	No
ANCF [25]	$\mathcal{O}\left(\frac{M^{2/3}[\Delta_0^{1/3}+d_0^{2/3}]}{\varepsilon^{2/3}} + \frac{MD_h}{\varepsilon}\right)$	$\mathcal{O}\left(mM^2\left[\frac{mD_h^2+\Delta_0}{\varepsilon^2}\right]\right)$	Yes
VRF [43]	$\mathcal{O}\left(\frac{M^{2/3}[\Delta_0^{1/3}+D_h^{2/3}]}{\varepsilon^{2/3}}\right)$	$\mathcal{O}\left(mM^2D_h^2\left[\frac{1+m^2}{\varepsilon^2}\right]\right)$	Yes
APD	$\tilde{\mathcal{O}}\left(\sqrt{M}\left[\min\left\{\frac{\sqrt{\Delta_0}}{\varepsilon}, \frac{d_0}{\sqrt{\varepsilon}}\right\}\right]\right)$	$\tilde{\mathcal{O}}\left(\frac{\sqrt{mM}\Delta_0}{\varepsilon^2}\right)$	No
<b>Known Lower Bounds</b>	$\Omega\left(\sqrt{M}\left[\min\left\{\frac{\sqrt{\Delta_0}}{\varepsilon}, \sqrt{\frac{d_0}{\varepsilon}}\right\}\right]\right)$	$\Omega\left(\frac{\sqrt{mM}\Delta_0}{\varepsilon^2}\right)$	-

TABLE 1.1

Lower bounds and iteration-complexity bounds of various parameter-free type-II composite optimization algorithms for finding  $\varepsilon$ -stationary points as in (1.3). The scalar  $D_h$  denotes the diameter of  $\text{dom } h$  and it is assumed that  $d_0$ ,  $\Delta_0$ ,  $m$ , and  $M$  are not known but  $M$  is greater than or equal to  $m$  for the listed algorithms. The lower bounds for the convex (resp. nonconvex) case can be found in [7, Theorem 1] (resp. [48, Theorem 4.5]).

161 Notice that the analysis for UPFAG does not include an iteration complexity  
 162 bound for finding stationary points when  $f$  is convex, while ANCF and VRF suffer  
 163 from the requirement that  $\text{dom } h$  must be bounded. Moreover, up until this point,  
 164 PGD was the only parameter-free *type-II* algorithm with an established iteration  
 165 complexity bound for the unbounded case when  $f$  is convex. None of the parameter-  
 166 free methods before this work, in the nonconvex setting, could obtain the optimal  
 167 complexity bound in [18].

168 In addition to the development of PF.APD, some details are given regarding  
 169 how PF.APD could be used in other existing optimization frameworks, including  
 170 min-max smoothing and penalty frameworks for constrained optimization. The main  
 171 advantages of these resulting frameworks are that (i) they are parameter-free and (ii)  
 172 they have improved complexities when  $f$  in (1.1) is convex, without requiring any  
 173 adjustments to their inputs.

174 Finally, numerical experiments are given to support the practical efficiency of  
 175 PF.ADP on some randomly generated problem instances. These experiments specif-  
 176 ically show that PF.APD consistently outperforms several existing parameter-free  
 177 methods in practice.

178 **1.3. Organization.** Section 2 presents background material. Section 3 presents  
 179 PF.ACG, PF.APD, and their iteration complexity bounds. Section 4 gives the proofs  
 180 of several important technical results. Section 5 describes how PF.APD can be used  
 181 in existing optimization frameworks. Section 6 presents some numerical experiments.  
 182 Section 7 gives some concluding remarks. Several technical appendices follow after  
 183 the above sections.

184 **1.4. Notation and Basic Definitions.**  $\mathbb{R}_+$  and  $\mathbb{R}_{++}$  denote the set of nonneg-  
 185 ative and positive real numbers, respectively.  $\mathbb{R}^n$  denotes an  $n$ -dimensional Euclid-  
 186 ean space with inner product and norm denoted by  $\langle \cdot, \cdot \rangle$  and  $\|\cdot\|$ , respectively.  
 187  $\text{dist}(x, X)$  denotes the Euclidean distance of a point  $x$  to a set  $X$ . For any  $t > 0$ ,

188 we denote  $\log^+(t) := \max\{\log t, 1\}$ . For a function  $h : \mathbb{R}^n \rightarrow (-\infty, \infty]$  we denote  
 189  $\text{dom } h := \{x \in \mathbb{R}^n : h(x) < +\infty\}$  to be the domain of  $h$ . Moreover,  $h$  is considered  
 190 proper if  $\text{dom } h \neq \emptyset$ . The set of all lower semi-continuous proper convex functions defined  
 191 in  $\mathbb{R}^n$  is denoted by  $\overline{\text{Conv}}(\mathbb{R}^n)$ . The convex subdifferential of a proper function  
 192  $h : \mathbb{R}^n \rightarrow (-\infty, \infty]$  is given by

$$193 \quad (1.10) \quad \partial h(z) := \{u \in \mathbb{R}^n : h(z') \geq h(z) + \langle u, z' - z \rangle, \quad \forall z' \in \mathbb{R}^n\}$$

194 for every  $z \in \mathbb{R}^n$ . If  $\psi$  is a real-valued function which is differentiable at  $\bar{z} \in \mathbb{R}^n$ , then  
 195 its affine/linear approximation  $\ell_\psi(\cdot, \bar{z})$  at  $\bar{z}$  is given by

$$196 \quad (1.11) \quad \ell_\psi(z; \bar{z}) := \psi(\bar{z}) + \langle \nabla \psi(\bar{z}), z - \bar{z} \rangle \quad \forall z \in \mathbb{R}^n.$$

197 **2. Background.** This section gives some necessary background for presenting  
 198 PF.ACG and PF.APD. More specifically, [Subsection 2.1](#) describes and comments on  
 199 the problem of interest, while [Subsection 2.2](#) presents a general proximal descent  
 200 scheme which serves as a template for PF.APD.

201 **2.1. Problem of Interest.** To reiterate, we are interested in the following com-  
 202 posite optimization problem:

203 **Problem  $\mathcal{CO}$ :** Given  $\varepsilon \in \mathbb{R}_{++}$  and a function  $\phi = f + h$  satisfying:  
[A1](#)  $h \in \overline{\text{Conv}}(\mathbb{R}^n)$  and the resolvent  $(\lambda \partial h + \text{id})^{-1}$  is easy to compute  
 for any  $\lambda > 0$ ,  
[A2](#)  $f$  is continuously differentiable on an open set  $\Omega \supseteq \text{dom } h$ , and  
 $\nabla f$  is  $\mathcal{M}$ -Lipschitz continuous on  $\text{dom } h$  for some  $\mathcal{M} \in \mathbb{R}_{++}$ ,  
[A3](#)  $\phi_* = \inf_{z \in \mathbb{R}^n} \phi(z) > -\infty$ ,  
 find a pair  $(\bar{z}, \bar{v}) \in \text{dom } h \times \mathbb{R}^n$  satisfying [\(1.3\)](#).

204 Of the three above assumptions, only [A1](#) is a necessary condition that is used  
 205 to ensure PF.APD is well-defined. Assumptions [A2](#)–[A3](#), on the other hand, are  
 206 sufficient conditions that are used to show that PF.APD stops in a finite number of  
 207 iterations. It is possible to replace assumption [A2](#) with more general smoothness  
 208 conditions (e.g., Hölder continuity [[13,36](#)]) at the cost of a possibly more complicated  
 209 analysis. It is known<sup>3</sup> that assumption [A2](#) holds if and only if

$$210 \quad (2.1) \quad |f(z) - \ell_f(z; z')| \leq \frac{\mathcal{M}}{2} \|z - z'\|^2, \quad \forall z, z' \in \text{dom } h,$$

211 which implies  $(\mathcal{M}, \mathcal{M})$  is a curvature pair of  $\phi$ .

212 We now comment on criterion [\(1.3\)](#). First, it is related to the directional derivative  
 213 of  $\phi$ :

$$214 \quad \min_{\|d\|=1} \phi'(z; d) = \min_{\|d\|=1} \max_{\zeta \in \partial h(z)} \langle \nabla f(z) + \zeta, d \rangle = \max_{\zeta \in \partial h(z)} \min_{\|d\|=1} \langle \nabla f(z) + \zeta, d \rangle$$

$$215 \quad = - \min_{\zeta \in \partial h(z)} \|\nabla f(z) + \zeta\| = -\text{dist}(0, \nabla f(z) + \partial h(z)).$$

216

<sup>3</sup>The proof of the forward direction is well-known (see, for example, [[4,37](#)]) while the proof of the reverse direction can be found, for example, in [[17](#), Proposition 2.1.55]. For the special case where  $f$  is convex and real-valued, the proof of the reverse direction can be found, for example, in [[3](#), Theorem 18.15] and [[33](#), 2.1.5].

217 Consequently, if  $\bar{z} \in \text{dom } h$  is a local minimum of  $\phi$  then  $\min_{\|d\|=1} \phi'(\bar{z}; d) \geq 0$  and  
 218 the above relation implies that (1.3) holds with  $\varepsilon = 0$ . That is, (1.3) is a necessary  
 219 condition for local optimality of a point  $\bar{z} \in \text{dom } h$ . Second, when  $f$  is convex then  
 220 (1.3) with  $\varepsilon = 0$  implies that  $0 \in \nabla f(\bar{z}) + \partial h(\bar{z}) = \partial \phi(\bar{z})$  and  $\bar{z}$  is a global minimum.  
 221 Given the first comment, (1.3) is equivalent to global optimality of a point  $\bar{z} \in \text{dom } h$   
 222 when  $f$  is convex.

223 **2.2. General Proximal Descent Scheme.** Our interest in this subsection is  
 224 the general proximal descent scheme in Algorithm 2.1, which follows the ideas in  
 225 (1.5)–(1.7). Its iteration scheme serves as a template for the PF.APD presented in  
 226 Subsection 3.2.

---

**Algorithm 2.1** General Proximal Descent Scheme

---

**Data:**  $(f, h)$  as in (A1)–(A3),  $z_0 \in \text{dom } h$ ;

**Parameters:**  $\theta \in \mathbb{R}_+$ ;

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

2: **find**  $(z_{k+1}, u_{k+1}) \in \text{dom } h \times \mathbb{R}^n$  and  $m_{k+1} \in \mathbb{R}_{++}$  satisfying

$$(2.2) \quad u_{k+1} \in \nabla f(z_{k+1}) + 2m_{k+1}(z_{k+1} - z_k) + \partial h(z_{k+1}),$$

$$(2.3) \quad \|u_{k+1} + 2m_{k+1}(z_k - z_{k+1})\|^2 \leq 2\theta m_{k+1} [\phi(z_k) - \phi(z_{k+1})],$$

$$(2.4) \quad \|u_{k+1}\|^2 \leq m_{k+1}^2 \|z_{k+1} - z_k\|^2.$$


---

227 Before presenting the properties of Algorithm 2.1, let us comment on its steps.  
 228 First, (2.2)–(2.4) are analogous to (1.5)–(1.7) because of assumption (A1). Second, if  
 229  $f + m_{k+1}\|\cdot\|^2$  is convex and  $u_{k+1} = 0$  then (2.2) implies that

$$230 \quad z_{k+1} = \underset{z \in \text{dom } h}{\text{argmin}} \left\{ \frac{\phi(z)}{2m_{k+1}} + \frac{1}{2} \|z - z_{k+1}\|^2 \right\},$$

231 which is a proximal point update with stepsize  $1/(2m_{k+1})$ . Third, (2.3) implies that  
 232 Algorithm 2.1 is a descent scheme, i.e.,  $\phi(z_{k+1}) \leq \phi(z_k)$  for  $k \geq 0$ . Hence, in view of  
 233 the second comment, this justifies its qualifier as a “proximal descent” scheme.

234 It is also worth mentioning that (2.3)–(2.4) are similar to conditions in existing  
 235 literature. More specifically, a version of (2.3) can be found in the descent scheme  
 236 of [19], while an inequality similar to (2.4) can be found in the GIPP framework of [18]  
 237 with  $\sigma = 1$ ,  $\tilde{\varepsilon} = 0$ , and  $v_{k+1} = u_{k+1}/m_{k+1}$ . However, the addition of condition (2.2)  
 238 appears to be new.

239 We now present the most important properties of Algorithm 2.1. The first result  
 240 supports the importance of conditions (2.2)–(2.3).

241 **LEMMA 2.1.** *Given  $z_0 \in X$ , let  $\{(z_{k+1}, u_{k+1})\}_{k \geq 0}$  denote a sequence of iterates  
 242 satisfying (2.2)–(2.3). Moreover, let  $\Delta_0$  be as in (1.8), and define*

$$243 \quad v_{k+1} := u_{k+1} + 2m_{k+1}(z_k - z_{k+1}), \quad \Lambda_{k+1} := \sum_{j=0}^k \frac{1}{m_{j+1}}, \quad \forall k \geq 0.$$

244 *Then, for every  $k \geq 0$ ,*

- 245 (a)  $v_{k+1} \in \nabla f(z_{k+1}) + \partial h(z_{k+1})$ ;  
 246 (b)  $\min_{0 \leq j \leq k} \|v_{j+1}\|^2 \leq 2\theta \Delta_0 \Lambda_{k+1}^{-1}$ .

247 *Proof.* (a) This follows immediately from (2.2) and the definition of  $v_{k+1}$ .

248 (b) Summing up both sides of (2.3) from 0 to  $k$ , the definition of  $v_{k+1}$ , and the  
 249 definition of  $\phi_*$ , we have that

$$\begin{aligned}
 250 \quad \Lambda_{k+1} \min_{0 \leq j \leq k} \|v_{j+1}\|^2 &\leq \sum_{j=0}^k \frac{\|v_{j+1}\|^2}{m_{j+1}} \stackrel{(2.3)}{\leq} 2\theta \sum_{j=0}^k [\phi(z_j) - \phi(z_{j+1})] \\
 251 \quad &= 2\theta [\phi(z_0) - \phi(z_{k+1})] \leq 2\theta [\phi(z_0) - \phi_*] = 2\theta \Delta_0. \quad \square
 \end{aligned}$$

253 Notice that Lemma 2.1(b) implies that if  $\lim_{k \rightarrow \infty} \Lambda_{k+1} \rightarrow \infty$  then we have that  
 254  $\lim_{k \rightarrow \infty} \min_{j \leq k} \|v_{j+1}\| \rightarrow 0$ . Moreover, if  $\sup_{k \geq 0} m_{k+1} < \infty$  then for any  $\varepsilon > 0$ , there  
 255 exists some finite  $j \geq 0$  such that  $\|v_{j+1}\| \leq \varepsilon$ .

256 The next result shows that if  $m_{k+1}$  is bounded relative to the global topology of  
 257  $f$ , and conditions (2.2)–(2.4) hold, then a more refined bound of  $\min_{j \leq k} \|v_{j+1}\|$  can  
 258 be obtained. To keep the notation concise, we make use of the following quantity:

$$259 \quad (2.5) \quad R_\tau(\hat{z}) := \inf_{z \in \mathbb{R}^n} \left\{ R_\tau(z, \hat{z}) := \frac{\phi(z) - \phi_*}{\tau} + \frac{1}{2} \|z - \hat{z}\|^2 \right\}.$$

260 It is easy to see that  $R_\tau(z')$  is the Moreau envelope of  $\phi/\tau$  at  $z'$  shifted by  $-\phi_*/\tau$ .

261 LEMMA 2.2. *Given  $z_0 \in X$ , let  $\{(v_{j+1}, z_{j+1}, \Lambda_{j+1})\}_{j \geq 0}$  be as in Lemma 2.1 and  
 262  $k \geq 0$  be fixed. Moreover, suppose (2.4) holds and that there exists  $\tilde{m} > 0$  such that  
 263  $f + \tilde{m}\|\cdot\|^2/2$  is convex. If  $\min_{0 \leq j \leq k} m_{j+1} \geq \tilde{m}$  and  $\max_{0 \leq j \leq k} m_{j+1} \leq (1 + \nu)\tilde{m}$  for  
 264 some  $\nu > 0$ , then*

$$265 \quad (2.6) \quad \phi(z_{k+1}) + \frac{m_{k+1}}{2} \|z_{k+1} - z_k\|^2 \leq \inf_{z \in \mathbb{R}^n} \left\{ \phi(z) + \frac{\nu \tilde{m}}{2} \|z - z_k\|^2 \right\},$$

266 and if  $k \geq 1$  then it holds that

$$267 \quad (2.7) \quad \min_{1 \leq j \leq k} \|v_{j+1}\|^2 \leq 2\theta \nu \tilde{m} \left[ \frac{R_{\nu \tilde{m}}(z_0)}{\Lambda_{k+1} - m_1^{-1}} \right].$$

268 *Proof.* Using the assumption that  $m_{k+1} \geq \tilde{m}$  and (2.2), we have that  $f(\cdot) +$   
 269  $m_{k+1}\|\cdot - z_k\|^2$  is  $\tilde{m}$ -strongly convex and, hence,

$$\begin{aligned}
 270 \quad u_{k+1} &\in \nabla f(z_{k+1}) + 2m_{k+1}(z_{k+1} - z_k) + \partial h(z_{k+1}) \\
 271 \quad &= \nabla f(z_{k+1}) - \tilde{m}(z_{k+1} - z_{k+1}) + 2m_{k+1}(z_{k+1} - z_k) + \partial h(z_{k+1}) \\
 272 \quad (2.8) \quad &= \partial \left( \phi - \frac{\tilde{m}}{2} \|\cdot - z_{k+1}\|^2 + m_{k+1} \|\cdot - z_k\|^2 \right) (z_{k+1}). \\
 273
 \end{aligned}$$

274 Using (2.8), (2.4), and the bound  $\langle a, b \rangle \geq -\|a\|^2/(2m_{k+1}) - m_{k+1}\|b\|^2/2$  for any  
 275  $a, b \in \mathbb{R}^n$ , it holds for any  $z \in \mathbb{R}^n$  that

$$\begin{aligned}
 276 \quad &\phi(z) + m_{k+1} \|z - z_k\|^2 \\
 277 \quad &\stackrel{(2.8)}{\geq} \phi(z_{k+1}) + m_{k+1} \|z_k - z_{k+1}\|^2 + \frac{\tilde{m}}{2} \|z - z_{k+1}\|^2 + \langle u_{k+1}, z - z_{k+1} \rangle \\
 278 \quad &\geq \phi(z_{k+1}) + m_{k+1} \|z_k - z_{k+1}\|^2 - \frac{1}{2m_{k+1}} \|u_{k+1}\|^2 + \frac{\tilde{m} - m_{k+1}}{2} \|z - z_{k+1}\|^2 \\
 279 \quad &\stackrel{(2.4)}{\geq} \phi(z_{k+1}) + \frac{m_{k+1}}{2} \|z_k - z_{k+1}\|^2 + \frac{\tilde{m} - m_{k+1}}{2} \|z - z_{k+1}\|^2. \\
 280
 \end{aligned}$$



281 Re-arranging terms and using the assumption  $m_{k+1} \leq (1 + \nu)\tilde{m}$ , we then have that

$$282 \quad \phi(z_{k+1}) + \frac{m_{k+1}}{2} \|z_k - z_{k+1}\|^2 \leq \phi(z) + \frac{\nu\tilde{m}}{2} \|z - z_k\|^2,$$

283 which implies (2.6) as  $z \in \mathbb{R}^n$  was arbitrary. To show (2.7), we use (2.6) at  $k = 1$ ,  
284 (2.3), and the definition of  $v_{k+1}$  to conclude that

$$\begin{aligned} 285 \quad R_{\nu\tilde{m}}(z_0) &= \inf_{z \in \mathbb{R}^n} \left\{ \frac{\phi(z) - \phi_*}{\nu\tilde{m}} + \frac{1}{2} \|z - z_0\|^2 \right\} \geq \frac{\phi(z_1) - \phi_*}{\nu\tilde{m}} + \frac{m_1}{2\nu\tilde{m}} \|z_1 - z_0\|^2 \\ 286 \quad &\stackrel{(2.6)}{\geq} \frac{\phi(z_1) - \phi(z_{k+1})}{\nu\tilde{m}} = \frac{\sum_{j=1}^k [\phi(z_j) - \phi(z_{j+1})]}{\nu\tilde{m}} \stackrel{(2.3)}{\geq} \frac{1}{2\theta\nu\tilde{m}} \sum_{j=1}^k \frac{\|v_{j+1}\|^2}{m_{j+1}} \\ 287 \quad &\geq \frac{\sum_{j=1}^k m_{j+1}^{-1}}{2\theta\nu\tilde{m}} \left( \inf_{1 \leq j \leq k} \|v_{j+1}\|^2 \right) = \frac{\Lambda_{k+1} - m_1^{-1}}{2\theta\nu\tilde{m}} \inf_{1 \leq j \leq k} \|v_{j+1}\|^2. \quad \square \end{aligned}$$

289 Similar to the previous lemma, the above result also implies that if  $\lim_{k \rightarrow \infty} \Lambda_{k+1} \rightarrow \infty$   
290 then we have  $\lim_{k \rightarrow \infty} \min_{j \leq k} \|v_{j+1}\| \rightarrow 0$ . However, it is more general in the sense  
291 that the rate of convergence depends on  $R_{\nu\tilde{m}}(z_0)$  instead of  $\Delta_0$ , and the former can  
292 be bounded as

$$293 \quad (2.9) \quad R_{\nu\tilde{m}}(z_0) \leq \min \{R_{\nu\tilde{m}}(z_0, z_0), R_{\nu\tilde{m}}(z_*, z_0)\} \leq \min \left\{ \frac{\Delta_0}{\nu\tilde{m}}, \frac{d_0^2}{2} \right\},$$

294 where  $z_*$  is any optimal solution of (1.1) that is the closest to  $z_0$  and  $(\Delta_0, d_0)$  are as  
295 in (1.8). This fact will be important when we establish an iteration complexity bound  
296 for PF.APD in the convex setting.

297 **3. Parameter-Free Algorithms.** This section presents PF.ACG, PF.APD,  
298 and their iteration complexity bounds. More specifically, [Subsection 3.1](#) presents  
299 PF.ACG, while [Subsection 3.2](#) presents PF.APD.

300 It is also worth recalling that PF.APD is an implementation of the general de-  
301 scent scheme of the previous section that repeatedly calls PF.ACG to obtain a single  
302 iteration of the scheme mentioned above.

303 **3.1. PF.ACG Algorithm.** Broadly speaking, PF.ACG is a modification of  
304 the well-known FISTA [5, 11] algorithm for minimizing  $\mu$ -strongly convex composite  
305 functions. Specifically, both PF.ACG and FISTA consider the composite optimization  
306 problem

$$307 \quad (3.1) \quad \min_{x \in \mathbb{R}^n} \{ \psi(x) := \psi^s(x) + \psi^n(x) \}$$

308 where  $(\psi^s, \psi^n)$  satisfies the following assumptions:

309 **(B1)**  $\psi^n \in \text{Conv}(\mathbb{R}^n)$  and the resolvent  $(\lambda\partial\psi^n + \text{id})^{-1}$  is easy to compute for any  
310  $\lambda > 0$ ,

311 **(B2)**  $\psi^s$  is continuously differentiable on an open set  $\Omega \supseteq \text{dom } \psi^n$ , and  $\nabla\psi^s$  is  
312  $L_*$ -Lipschitz continuous on  $\text{dom } \psi^n$  for some  $L_* \in \mathbb{R}_{++}$ .

313 Similar to (2.1), note that **(B2)** implies

$$314 \quad (3.2) \quad |\psi^s(x) - \ell_{\psi^s}(x; x')| \leq \frac{L_*}{2} \|x - x'\|^2 \quad \forall x, x' \in \text{dom } \psi^n.$$

315 PF.ACG differs from FISTA in that it adds two stopping conditions that help  
316 implement a single iteration of [Algorithm 2.1](#). Specifically, for a given function pair

317  $(f, h)$  satisfying [\(A1\)](#)–[\(A2\)](#), hyperparameters  $(\sigma, \theta, \mu) \in \mathbb{R}_{++}^3$ , and an initial point  
 318  $\hat{z} \in \text{dom } h$ , if PF.ACG is invoked with

$$319 \quad (3.3) \quad \psi^s(\cdot) = \frac{f(\cdot)}{2\hat{m}} + \frac{1}{2}\|\cdot - \hat{z}\|^2, \quad \psi^n(\cdot) = \frac{h(\cdot)}{2\hat{m}},$$

320 for some  $\hat{m} > 0$ , then either (i) PF.ACG has found a pair  $(y, u)$  satisfying conditions  
 321 [\(2.2\)](#)–[\(2.4\)](#) with  $(z_{k+1}, u_{k+1}, m_{k+1}, z_k) = (y, u, m, \hat{z})$ , or (ii) some local  $\mu$ -strong con-  
 322 vexity condition has failed, and the estimate of  $\mu$  or the function pair  $(\psi^s, \psi^n)$  has to  
 323 be changed.

324 We now present the details of PF.ACG and its key properties. To help our  
 325 discussion, we first give the complete pseudocode of PF.ACG through [Algorithm 3.1](#)  
 326 and [Algorithm 3.2](#). More specifically, [Algorithm 3.1](#) presents the accelerated gradient  
 327 FISTA update and (Lipschitz constant) line search strategy used in PF.ACG, while  
 328 [Algorithm 3.2](#) describes the other steps of PF.ACG and how [Algorithm 3.1](#) is invoked.

---

**Algorithm 3.1** Line Search and Accelerated Gradient Step Subroutine

---

**Data:**  $(\psi^s, \psi^n)$  as in [\(B1\)](#)–[\(B2\)](#),  $(\hat{y}, \hat{x}) \in \text{dom } \psi^n \times \mathbb{R}^n$ ,  $\hat{A} \geq 0$ ,  $\mu \in \mathbb{R}_{++}$ ,  $\hat{L} \in [\mu, \infty)$ ;

**Hyper-parameters:**  $\beta \in (1, \infty)$ ;

**Outputs:**  $(A, \tilde{x}, y, x, L) \in \mathbb{R}_+ \times \mathbb{R}^n \times \text{dom } \psi^n \times \mathbb{R}^n \times \mathbb{R}_+$  and function  $q$ ;

1:  $\psi \leftarrow \psi^s + \psi^n$

2: **for**  $\ell \leftarrow 0, 1, \dots$  **do**

3:  $L \leftarrow \hat{L}\beta^\ell$

▷ **Step 1: Accelerated gradient step.**

4:  $\xi \leftarrow 1 + \mu\hat{A}$  and find  $\hat{a}$  satisfying  $\hat{a}^2 = \hat{\xi}(\hat{a} + \hat{A})/L$

5:  $A \leftarrow \hat{A} + \hat{a}$

6:  $\tilde{x} \leftarrow \frac{\hat{A}}{A}\hat{y} + \frac{\hat{a}}{A}\hat{x}$

7:  $y \leftarrow \text{argmin}_{z \in \mathbb{R}^n} \left\{ \ell_{\psi^s}(z; \tilde{x}) + \psi^n(z) + \frac{L+\mu}{2}\|z - \tilde{x}\|^2 \right\}$

8:  $x \leftarrow \hat{x} + \frac{\hat{a}}{1+A\mu} [L(y - \tilde{x}) + \mu(y - \hat{x})]$

▷ **Step 2: Descent condition check.**

9: **if** the inequality

$$(3.4) \quad \psi^s(y) - \ell_{\psi^s}(y; \tilde{x}) \leq \frac{L}{2}\|y - \tilde{x}\|^2$$

holds, then **return**  $(A, \tilde{x}, y, x, L)$

---

329 We next present some key properties about [Algorithm 3.2](#) and its iterates. As  
 330 their proof is mostly technical, we moved it to [Subsection 4.1](#).

331 **LEMMA 3.1.** *For every  $j \geq 0$ ,*

332 (a)  $A_{j+1} \geq (1/L_0) \prod_{i=1}^j [1 + \sqrt{\mu/(2L_i)}]$  and

$$333 \quad (3.7) \quad L_j \leq L_{j+1} \leq \bar{L} := \max\{1, \alpha L_*\}.$$

334 (b)  $r_{j+1} \in \nabla\psi^s(y_{j+1}) + \partial\psi^n(y_{j+1})$ ;

335 (c) if  $\psi^s$  is  $\mu$ -strongly convex, then [\(3.5\)](#) holds;

336 (d) if [\(3.5\)](#) holds and

$$337 \quad (3.8) \quad A_{j+1} \geq \frac{16\bar{L}^2}{\mu} \max\left\{\frac{1}{\sigma^2}, \frac{4\theta}{\theta-2}\right\} =: \mathcal{A}_{\mu, \bar{L}}(\sigma, \theta)$$

338 then [\(3.6\)](#) holds.

---

**Algorithm 3.2** Parameter-Free Accelerated Composite Gradient (PF.ACG) Algorithm

---

**Data:**  $(\psi^s, \psi^n)$  as in (B1)–(B2),  $y_0 \in \text{dom } \psi^n$ ,  $\mu \in \mathbb{R}_{++}$ ,  $L_0 \in [\mu, \infty)$ ;

**Hyper-parameters:**  $\sigma \in \mathbb{R}_{++}$ ,  $\theta \in (2, \infty)$ ,  $\beta \in (1, \infty)$ ;

**Outputs:**  $(y_{j+1}, u_{j+1}, L_{j+1}) \in \text{dom } \psi^n \times \mathbb{R}^n \times \mathbb{R}_{++}$ ;

1:  $(x_0, A_0) \leftarrow (y_0, 0)$

2:  $\psi(\cdot) \leftarrow \psi^s(\cdot) + \psi^n(\cdot)$

3: **for**  $j \leftarrow 0, 1, \dots$  **do**

    ▷ **Step 1: Line search for  $L_{j+1}$  and accelerated gradient step.**

4: **call** **Algorithm 3.1** with data  $(\psi^s, \psi^n)$ ,  $(\hat{y}, \hat{x}) \equiv (y_j, x_j)$ ,  $\hat{A} \equiv A_j$ ,  $\hat{\xi} \equiv \xi_j$ ,  $\mu$ ,  
     $\hat{L} \equiv L_j$  and hyper-parameter  $\beta$  to obtain  $(A_{j+1}, \tilde{x}_j, y_{j+1}, x_{j+1}, L_{j+1})$

    ▷ **Step 2: "Bad" termination check.**

5:  $r_{j+1} \leftarrow \nabla \psi^s(y_{j+1}) - \nabla \psi^s(\tilde{x}_j) + (L_{j+1} + \mu)(\tilde{x}_j - y_{j+1})$

6: **if** the inequalities

$$(3.5) \quad \begin{aligned} \mu A_{j+1} \|y_{j+1} - \tilde{x}_j\|^2 &\leq \|y_{j+1} - y_0\|^2, \\ \psi(y_0) &\geq \psi(y_{j+1}) + \langle r_{j+1}, y_0 - y_{j+1} \rangle, \end{aligned}$$

    do not hold, then **return**  $(y_{j+1}, r_{j+1}, L_{j+1})$

    ▷ **Step 3: "Good" termination check.**

7: **if** the inequalities

$$(3.6) \quad \begin{aligned} \|r_{j+1}\|^2 &\leq \sigma^2 \|y_{j+1} - y_0\|^2, \\ \|r_{j+1} + y_0 - y_{j+1}\|^2 &\leq \theta [\psi(y_0) - \psi(y_{j+1}) + \frac{1}{2} \|y_{j+1} - y_0\|^2], \end{aligned}$$

    hold, then **return**  $(y_{j+1}, r_{j+1}, L_{j+1})$

---

339 We now give a complexity bound for **Algorithm 3.2** and a condition for guaran-  
340 teeing its successful termination.

341 **PROPOSITION 3.2.** *The following properties hold about **Algorithm 3.2**:*

342 (a) it stops in

$$(3.9) \quad \left[ 1 + 2\sqrt{\frac{2\bar{L}}{\mu}} \log^{1+} \{ \bar{L} \mathcal{A}_{\mu, \bar{L}}(\sigma, \theta) \} \right],$$

344 where  $\bar{L}$  and  $\mathcal{A}_{\mu, \bar{L}}$  are as in (3.7) and (3.8), respectively.

345 (b) if  $\psi^s$  is  $\mu$ -strongly convex, then it always terminates in its Step 3 with a triple  
346  $(y_{j+1}, u_{j+1}, L_{j+1})$  satisfying (3.6) and  $L_0 \leq L_{j+1} \leq \bar{L}$ .

347 *Proof.* (a) Let  $J + 1$  denote the quantity in (3.9) and suppose **Algorithm 3.2** has  
348 not terminated at the end of iteration  $J + 1$ . Moreover, denote  $\mathcal{A} := \mathcal{A}_{\mu, \bar{L}}(\sigma, \theta)$ . Using  
349 **Lemma 3.1**(a), we first have

$$(3.10) \quad A_{J+1} \geq \frac{1}{L_0} \prod_{i=1}^J \left( 1 + \sqrt{\frac{\mu}{2L_i}} \right) \geq \frac{1}{\bar{L}} \left( 1 + \sqrt{\frac{\mu}{2\bar{L}}} \right)^J$$

351 Using the above bound, the fact that  $J \geq 2\sqrt{2\bar{L}/\mu} \log(\bar{L}\mathcal{A})$  from the definition in  
352 (3.9), the bound  $\mu \leq \bar{L}$ , and the fact that  $\log(1+t) \geq t/2$  on  $t \in [0, 1]$ , it holds that

$$353 \quad \log(\bar{L}\mathcal{A}) \leq \frac{J}{2} \sqrt{\frac{\mu}{2\bar{L}}} \leq J \log \left( 1 + \sqrt{\frac{\mu}{2\bar{L}}} \right) \stackrel{(3.10)}{\leq} \log(\bar{L}A_{J+1})$$

354 which implies  $A_{J+1} \geq \mathcal{A}$ . Hence, it follows from **Lemma 3.1**(d) that (3.6) holds. In  
355 view of Step 3 of **Algorithm 3.2** this implies that termination has to have occurred at or

356 before iteration  $J + 1$ , which contradicts our initial assumption. Thus, [Algorithm 3.2](#)  
 357 must have terminated by iteration  $J + 1$ .

358 (b) This follows immediately from part (a) and [Lemma 3.1\(c\)](#).  $\square$

359 The last result of this subsection shows how to invoke [Algorithm 3.2](#) so that its  
 360 successful termination implements a single iteration of [Algorithm 2.1](#).

361 **LEMMA 3.3.** *Suppose [Algorithm 3.2](#) is called with  $(\psi^s, \psi^n)$  as in (3.3) for some*  
 362  *$m > 0$  and  $\hat{z} \in \text{dom } \psi^n$ ,  $\sigma = 1/4$ , and  $y_0 = \hat{z}$ . If the call terminates in Step 3*  
 363 *with an output triple  $(y_{j+1}, r_{j+1}, L_{j+1})$ , then the quadruple  $(z_{k+1}, u_{k+1}, m_{k+1}, z_k) =$*   
 364  *$(y_{j+1}, 2mr_{j+1}, m, \hat{z})$  satisfies (2.2)–(2.4).*

365 *Proof.* Using [Lemma 3.1\(b\)](#), it holds that

$$366 \quad u_{k+1} = 2mr_{j+1} \in 2m [\nabla \psi^s(y_{j+1}) + \partial \psi^n(y_{j+1})]$$

$$367 \quad = \nabla f(z_{k+1}) + 2m_{k+1}(z_{k+1} - z_k) + \partial \psi^n(z_{k+1}).$$

369 which is exactly (2.2). Now, using the first inequality in (3.6), the choice of  $\sigma = 1/4$ ,  
 370 and the fact that  $y_0 = \hat{z} = z_k$ , we have

$$371 \quad \|u_{k+1}\|^2 = 4m^2 \|r_{j+1}\|^2 \stackrel{(3.6)}{\leq} m^2 \|y_{j+1} - y_0\|^2 = m_{k+1}^2 \|z_{k+1} - z_k\|^2,$$

372 which is exactly (2.4). Finally, the second condition of (3.6), the relation  $\psi(\cdot) =$   
 373  $\phi(\cdot)/(2m_{k+1}) + \|\cdot - y_0\|^2/2$ , and the fact that  $y_0 = \hat{z} = z_k$  imply

$$374 \quad \|u_{k+1} + 2m_{k+1}(z_k - z_{k+1})\|^2 = 4m_{k+1}^2 \|r_{j+1} + y_{j+1} - y_0\|^2$$

$$375 \quad \stackrel{(3.6)}{\leq} 4\theta m_{k+1}^2 \left[ \psi(y_0) - \psi(y_{j+1}) + \frac{1}{2} \|y_{j+1} - y_0\|^2 \right] = 2\theta m_{k+1} [\phi(z_k) - \phi(z_{k+1})],$$

377 which is exactly (2.3). Combing all previous inequalities yields the desired conclusion.  $\square$

378 Some remarks are in order. We first remark on [Algorithm 3.1](#):

- 379 1. In view of (3.2), the number of iterations in its  $(j + 1)$ -th call stops is bounded
- 380 above by  $1 + \log_\beta(L_{j+1}/L_j)$ .
- 381 2. The update for  $y$  is equivalent to

$$382 \quad y = \underset{z \in \text{dom } \psi^n}{\text{argmin}} \left\{ \frac{\psi^n(z)}{L + \mu} + \frac{1}{2} \left\| z - \left( \tilde{x} - \frac{\nabla \psi^s(\tilde{x})}{L + \mu} \right) \right\|^2 \right\}$$

383 which is a single call to the prox-oracle of  $\psi^n/(L + \mu)$ .

- 384 3. The descent condition (3.4) is well-known in existing literature for adaptive
- 385 FISTA-type methods (see, for example, [41, Subsection 4.3]).

386 We now remark on [Algorithm 3.2](#) and its associated results:

- 387 4. It is shown in [Lemma 3.1](#) that (i)  $r_{j+1}$  is a stationarity residual for the iterate
- 388  $y_{j+1}$  and (ii)  $\{L_j\}_{j \geq 0}$  forms a nondecreasing sequence of nonnegative scalars.
- 389 5. Step 1 is generally where most of the computation is done, wherein (possibly)
- 390 multiple accelerated gradient steps are performed using [Algorithm 3.1](#). It is also
- 391 the only step that requires evaluating the prox oracle for  $\psi^n$ .
- 392 6. It is shown in [Proposition 3.2\(b\)](#) that both inequalities in Step 2 hold when  $\psi^s$  is
- 393  $\mu$ -strongly convex. The first (resp. second) inequality of (3.5) is used to ensure
- 394 that the first (resp. second) inequality of (3.6) holds when enough iterations are
- 395 performed. See the analysis in [Subsection 4.1](#) for more details.

- 396 7. Condition (3.6) is chosen so that Algorithm 3.2 implements a single step of  
 397 Algorithm 2.1 if it stops in Step 3 and it is given the right inputs (see Lemma 3.3).  
 398 8. Suppose Algorithm 3.2 terminates in  $J$  iterations. Then, the number iterations  
 399 of Algorithm 3.1 taken by Algorithm 3.2 is

$$400 \quad \sum_{j=0}^{J-1} \left[ 1 + \log_{\beta} \frac{L_{j+1}}{L_j} \right] = J + \log_{\beta} \frac{L_J}{L_0} \leq J + \log_{\beta} \frac{\bar{L}}{L_0}.$$

401 Thus, on average (up to a  $(1/J) \log_{\beta}(\bar{L}/L_0)$  additive term) Algorithm 3.2 uses  
 402 only one accelerated gradient step or two function and prox oracle calls. It is  
 403 worth mentioning that Nesterov's universal fast gradient method [36, Section  
 404 4] uses on average (up to a  $(1/J) \log_{\beta}(\bar{L}/L_0)$  additive term) four function/prox  
 405 oracle calls per invocation.

406 **3.2. PF.APD Algorithm.** Broadly speaking, PF.APD is a *double-loop* method  
 407 consisting of *outer iterations* and (possibly) several *inner iterations* per outer iteration.  
 408 More specifically, the  $(k+1)$ -th outer iteration of PF.APD repeatedly applies  
 409 Algorithm 3.2 to the proximal subproblem

$$410 \quad z_{k+1} \approx \operatorname{argmin}_{z \in \operatorname{dom} h} \left\{ \frac{\phi(z)}{2\hat{m}} + \frac{1}{2} \|z - z_k\|^2 \right\},$$

411 for increasing values of  $\hat{m} > 0$ , where  $z_k$  is an approximate solution to the  $k$ -th  
 412 subproblem. On the other hand, the inner iterations refer to the iterations performed  
 413 by Algorithm 3.2.

414 We now present the details of PF.APD and its key properties. To help our  
 415 discussion, we first give the complete pseudocode of PF.APD through Algorithm 3.1  
 416 and Algorithm 3.4. More specifically, Algorithm 3.1 presents the (lower curvature)  
 417 line search strategy used in PF.APD, while Algorithm 3.4 describes the other steps  
 418 of PF.APD and how Algorithm 3.3 is invoked.

419 We next present three important properties about Algorithm 3.4 and its iterates.  
 420 As its proof is mostly technical, we move it to Subsection 4.2. Moreover, to ensure  
 421 that the resulting properties account for the possible asymmetry in (1.2), we make  
 422 use of the scalars

$$423 \quad (3.13) \quad \begin{aligned} m_* &:= \operatorname{argmin}_{z, z' \in \operatorname{dom} h, t \geq 0} \left\{ t : f(z) - \ell_f(z; z') \geq -\frac{t}{2} \|z - z'\|^2 \right\}, \\ M_* &:= \operatorname{argmin}_{z, z' \in \operatorname{dom} h, t \geq 0} \left\{ t : f(z) - \ell_f(z; z') \leq \frac{t}{2} \|z - z'\|^2 \right\}, \end{aligned}$$

424 which are the values of a curvature pair of  $f$ .

425 **PROPOSITION 3.4.** *Define the scalars*

$$426 \quad (3.14) \quad \begin{aligned} \bar{m} &:= \max\{m_0, (\alpha + \beta)m_*\}, \quad \bar{M} := \beta [\max\{M_0, M_*\} + 2\bar{m}], \\ \bar{\mathcal{L}}_0 &:= \frac{\bar{M}}{2m_0} + 1, \quad P_0 := \log^{1+} \left\{ \bar{\mathcal{L}}_0 \mathcal{A}_{\frac{1}{2}, \bar{\mathcal{L}}_0} \left( \frac{1}{4}, \theta \right) \right\}, \end{aligned}$$

427 where  $(m_*, M_*)$  and  $\mathcal{A}_{\mu, \bar{\mathcal{L}}}(\cdot, \cdot)$  are as in (3.13) and (3.8), respectively. Then, for every  
 428  $k \geq 0$ , the following statements hold about Algorithm 3.4 and its iterates:

---

**Algorithm 3.3** Line Search and Proximal Descent Step

---

**Data:**  $(\psi^s, \psi^n, f, h)$  as in (3.3),  $\hat{z} \in \text{dom } h$ ,  $\hat{m} \in \mathbb{R}_{++}$ ,  $\hat{M} \in [m, \infty)$ ;

**Hyper-parameters:**  $\theta \in (2, \infty)$ ,  $\alpha \in (1, \infty)$ ,  $\beta \in (1, \infty)$ ;

**Outputs:**  $(z, u, m, M) \in \text{dom } h \times \mathbb{R}^n$ ;

- 1:  $M \leftarrow \hat{M}$
  - 2:  $\phi(\cdot) \leftarrow f(\cdot) + h(\cdot)$
  - 3: **for**  $\ell \leftarrow 0, 1, \dots$  **do**
  - 4:    $m \leftarrow \hat{m}\alpha^\ell$   
    ▷ **Step 1:**  $(\ell + 1)^{\text{th}}$  proximal subproblem.
  - 5:   **call** **Algorithm 3.2** with data  $(\psi^s, \psi^n)$ ,  $y_0 \equiv \hat{z}$ ,  $\mu \equiv 1/2$ ,  
     $L_0 \equiv M/(2m) + 1$ , and hyper-parameters  $\sigma \equiv 1/4$ ,  $\theta$ ,  $\beta$ , to obtain an  
    output tuple  $(z, r, L)$
  - 6:    $u \leftarrow 2mr$
  - 7:    $M \leftarrow 2m(L - 1)$   
    ▷ **Step 2:** Proximal descent check.
  - 8:   **if** the inequalities  
(3.11)                    $\|u + 2m(z - \hat{z})\|^2 \leq 2\theta m [\phi(\hat{z}) - \phi(z)],$   
                               $\|u\|^2 \leq m^2 \|z - \hat{z}\|^2,$   
    hold, then **return**  $(z, u, m, M)$
- 

---

**Algorithm 3.4** Parameter-Free Accelerated Proximal Descent (PF.APD) Algorithm

---

**Data:**  $(f, h)$  as in (A1)–(A3),  $z_0 \in \text{dom } h$ ,  $m_0 \in \mathbb{R}_{++}$ ,  $M_0 \in [m_0, \infty)$ ,  $\varepsilon \in \mathbb{R}_{++}$ ;

**Hyper-parameters:**  $\theta \in (2, \infty)$ ,  $\alpha \in (1, \infty)$ ,  $\beta \in (1, \infty)$ ;

**Outputs:**  $(z_{k+1}, v_{k+1}) \in \text{dom } h \times \mathbb{R}^n$ ;

- 1: **for**  $k \leftarrow 0, 1, \dots$  **do**  
    ▷ **Step 1:** Line search for  $m_{k+1}$  and proximal descent step.
  - 2:    $\hat{m} \leftarrow \begin{cases} m_k/\alpha, & \text{if } k \geq 1 \text{ and } m_k < \dots < m_0, \\ m_k, & \text{otherwise} \end{cases}$
  - 3:   **call** **Algorithm 3.3** with data  
(3.12)                    $\psi^s(\cdot) = \frac{f(\cdot)}{2\hat{m}} + \frac{1}{2} \|\cdot - z_k\|^2, \quad \psi^n(\cdot) = \frac{h(\cdot)}{2\hat{m}},$   
     $(f, h)$ ,  $\hat{z} \equiv z_k$ ,  $\hat{m} \equiv \hat{m}$ ,  $\hat{M} \equiv M_k$ , and hyper-parameters  $\theta$ ,  $\alpha$ ,  $\beta$   
    to obtain  $(z_{k+1}, u_{k+1}, m_{k+1}, M_{k+1})$   
    ▷ **Step 2:** Stationarity termination check.
  - 4:    $v_{k+1} \leftarrow 2m_{k+1}(u_{k+1} + z_k - z_{k+1})$
  - 5:   **if**  $\|v_{k+1}\| \leq \varepsilon$  **then**
  - 6:     **return**  $(z_{k+1}, v_{k+1})$
- 

- 429 (a)  $M_k \leq M_{k+1} \leq \bar{M} < \infty$  and  $\{1/m_k\}$  is bitonic<sup>4</sup> and bounded below by  $1/\bar{m}$ ;  
430 (b) its  $(k + 1)$ -th outer iteration performs at most  $T_{k+1}$  inner iterations, where

431 (3.15)                    $T_{k+1} \leq 20 \left( 1 + \log_\alpha \frac{m_{k+1}}{m_k} + \frac{1}{\sqrt{\alpha} - 1} \sqrt{\frac{\bar{M}}{2m_k}} \right) P_0;$

---

<sup>4</sup>A sequence  $\{a_k\}_{k=0}^n$  is *bitonic* if there exists  $0 \leq j \leq n$  such that  $a_0 \leq \dots \leq a_j \geq \dots \geq a_n$ . Note that monotone sequences are bitonic as well.

432 (c) it performs a finite number of outer iterations  $K(\varepsilon)$ , where

$$433 \quad (3.16) \quad K(\varepsilon) \leq 1 + \sum_{k=0}^{K(\varepsilon)-2} \frac{\bar{m}}{m_{k+1}} < 1 + \frac{2\theta\Delta_0\bar{m}}{\varepsilon^2};$$

434 (d) if, in addition,  $f$  is convex, then  $m_k = \alpha^{-k}m_0$  for every  $k \geq 0$  and  $K(\varepsilon)$  in  
435 (3.16) also satisfies

$$436 \quad (3.17) \quad K(\varepsilon) \leq 2 + \log_\alpha \left[ \alpha^4 + \frac{2\theta m_0^2 R_{m_0}(z_0)}{\varepsilon^2} \right],$$

437 where  $R_\tau(\cdot)$  is as in (2.5);

438 (e)  $v_{k+1} \in \nabla f(z_{k+1}) + \partial h(z_{k+1})$  and its final iterate  $(\bar{z}, \bar{v}) = (z_{k+1}, v_{k+1})$  solves  
439 Problem  $\mathcal{CO}$ .

440 We are now ready to give some important iteration complexity bounds on [Algo-](#)  
441 [rithm 3.4](#).

442 **THEOREM 3.5.** Define  $Q_0 := 20P_0 [1 + \log_\alpha(\bar{m}/m_0)]$ , where  $\bar{m}$  and  $P_0$  are as in  
443 (3.14), respectively. Then, [Algorithm 3.4](#) stops and outputs a pair  $(\bar{z}, \bar{v}) = (z_{k+1}, v_{k+1})$   
444 solving Problem  $\mathcal{CO}$  in  $\bar{T}$  inner iterations, where

$$445 \quad (3.18) \quad \bar{T} \leq Q_0 + \frac{20P_0}{\sqrt{\alpha}-1} \sqrt{\bar{M} \left[ 1 + \frac{2\theta\Delta_0\bar{m}}{\varepsilon^2} \right] \left[ \frac{1}{m_0} + \frac{2\theta\Delta_0}{\varepsilon^2} \right]},$$

446 and  $\Delta_0$  is as in (1.8). Moreover, if  $f$  is convex, then

$$447 \quad (3.19) \quad \bar{T} \leq Q_0 + \frac{20P_0\alpha}{(\sqrt{\alpha}-1)^2} \left[ \frac{\alpha^2}{\sqrt{m_0}} + \frac{\sqrt{\theta \min\{2\Delta_0, m_0 d_0^2\}}}{\varepsilon} \right],$$

448 where  $d_0$  is as in (1.8).

449 *Proof.* The fact that [Algorithm 3.4](#) stops in a finite number of inner iterations  
450 with a pair solving Problem  $\mathcal{CO}$  is immediate from [Proposition 3.4](#). Furthermore, the  
451 previous proposition also implies that the total number of inner iterations in a single  
452 call of [Algorithm 3.4](#) is at most

$$453 \quad \sum_{k=0}^{K(\varepsilon)-1} T_{k+1} \leq 20P_0 \sum_{k=0}^{K(\varepsilon)-1} \left( 1 + \log_\alpha \frac{m_{k+1}}{m_k} + \frac{1}{\sqrt{\alpha}-1} \sqrt{\frac{\bar{M}}{2m_k}} \right)$$

$$454 \quad \leq 20P_0 \left( 1 + \log_\alpha \frac{m_{K(\varepsilon)+1}}{m_0} + \frac{\sqrt{\bar{M}}}{\sqrt{\alpha}-1} \sum_{k=0}^{K(\varepsilon)-1} \frac{1}{\sqrt{m_k}} \right)$$

$$455 \quad (3.20) \quad \leq Q_0 + \frac{20P_0\sqrt{\bar{M}}}{\sqrt{\alpha}-1} \sum_{k=0}^{K(\varepsilon)-1} \frac{1}{\sqrt{m_k}},$$

$$456$$

457 where  $T_{k+1}$  and  $K(\varepsilon)$  are as in (3.15) and (3.16), respectively. Let us now bound the  
458 sum  $\sum_{k=0}^{K(\varepsilon)-1} m_k^{-1/2}$ . Using [Proposition 3.4\(c\)](#) and the fact that  $\|z\|_1 \leq \sqrt{n}\|z\|_2$  for  
459 any  $z \in \mathbb{R}^n$ , we first have

$$460 \quad \sum_{k=0}^{K(\varepsilon)-1} \frac{1}{\sqrt{m_k}} \leq \left[ K(\varepsilon) \sum_{k=0}^{K(\varepsilon)-1} \frac{1}{m_k} \right]^{1/2} \leq \sqrt{\left( 1 + \frac{2\theta\Delta_0\bar{m}}{\varepsilon^2} \right) \left( \frac{1}{m_0} + \frac{2\theta\Delta_0}{\varepsilon^2} \right)}.$$

461 Using (3.20) and the above bound yields (3.18).

462 Now, let  $\mathcal{R}_0 := R_{m_0}(z_0)$  and suppose  $f$  is convex. Using Proposition 3.4(d),  
 463 (2.9) with  $\tilde{m} = m_0\alpha^{-K(\varepsilon)}$  and  $\nu = \tilde{m}/m_0$ , and the inequality  $\sqrt{a+b} \leq \sqrt{a} + \sqrt{b}$  for  
 464  $a, b \in \mathbb{R}$ , we have

$$\begin{aligned}
 465 \quad & \sum_{k=0}^{K(\varepsilon)-1} \frac{1}{\sqrt{m_k}} = \sum_{k=0}^{K(\varepsilon)-1} \sqrt{\frac{\alpha^k}{m_0}} \leq \frac{\alpha^{K(\varepsilon)/2}}{\sqrt{m_0}(\sqrt{\alpha}-1)} \\
 466 \quad & \stackrel{(d)}{\leq} \frac{\alpha}{\sqrt{m_0}(\sqrt{\alpha}-1)} \sqrt{\alpha^4 + \frac{2\theta m_0^2 R_{m_0}(z_0)}{\varepsilon^2}} \leq \frac{\alpha}{(\sqrt{\alpha}-1)\varepsilon} \left[ \frac{\alpha^2}{\sqrt{m_0}} + \frac{\sqrt{2\theta m_0 \mathcal{R}_0}}{\varepsilon} \right] \\
 467 \quad & \stackrel{(2.9)}{\leq} \frac{\alpha}{(\sqrt{\alpha}-1)\varepsilon} \left[ \frac{\alpha^2}{\sqrt{m_0}} + \frac{\sqrt{\theta \min\{2\Delta_0, m_0 d_0^2\}}}{\varepsilon} \right]. \\
 468
 \end{aligned}$$

469 Combining (3.20) and the above bound yields (3.19).  $\square$

470 Some remarks are in order. We first remark on Algorithm 3.3:

- 471 1. In view of assumption  $\langle A2 \rangle$  and Proposition 3.2, the number of iterations in its  
 472  $k$ -th call is bounded above by  $1 + \log_\alpha(m_{k+1}/m_k)$ .
- 473 2. The checks in its Step 2 correspond to (2.3) and (2.4), respectively.
- 474 3. If the  $\ell$ -th call to Algorithm 3.2 ends with a “bad termination”, i.e., Step 2 in  
 475 Algorithm 3.2, then (3.11) does not hold, the estimate  $m$  is increased by a factor  
 476 of  $\alpha$ , and the algorithm proceeds to the  $(\ell + 1)$ -th iteration.

477 We now remark on Algorithm 3.4 and its associated results:

- 478 4. It is shown in Proposition 3.4 that (i)  $v_{j+1}$  is a stationarity residual for the  
 479 iterate  $z_{j+1}$  and (ii)  $\{M_k\}_{k \geq 0}$  and  $\{m_k\}_{k \geq 0}$  are nondecreasing and nonnegative.
- 480 5.  $Q_0$  in (3.18)–(3.19) bounds the total number of inner iterations performed by un-  
 481 successful calls to Algorithm 3.2, i.e., those that stop in Step 2 of Algorithm 3.2.
- 482 6. While  $m_0$  and  $M_0$  are free parameters, a good initial value<sup>5</sup> for them is an  
 483 estimate of the local Lipschitz constant  $\tilde{L}_0$  of  $\nabla f$  at  $z_0$ . Similar to the approach  
 484 in [32], one can estimate  $\tilde{L}_0$  by sampling some  $\hat{z} \in \text{dom } h$  with  $\hat{z} \neq z_0$  and  
 485 choosing  $\tilde{L}_0 = \|\nabla f(z_0) - \nabla f(\hat{z})\|/\|z_0 - \hat{z}\|$ .

486 Before ending the section, we discuss how different choices of  $m_0$  affect the complex-  
 487 ities in (3.18) and (3.19) when  $m_* \leq M_*$ :

- 488 7. In the general case, choosing  $m_0 = 1$  implies that the bound in (3.18) (resp.  
 489 (3.19)) is  $\mathcal{O}(\sqrt{M_* m_*} \Delta_0 / \varepsilon^2)$  (resp.  $\mathcal{O}(\sqrt{M_*} \Delta_0 / \varepsilon)$ ) which matches the complexity  
 490 of the AIPP in [18] and is optimal<sup>6</sup> for finding stationary points of (1) in the  
 491 weakly-convex (resp. convex) setting in terms of  $m_*$ ,  $M_*$ ,  $\Delta_0$ , and  $\varepsilon$ .
- 492 8. If  $d_0$  is known, then choosing  $m_0 = \varepsilon/d_0$  implies (3.19) is  $\tilde{\mathcal{O}}(\sqrt{M_* d_0} / \sqrt{\varepsilon})$  which  
 493 is optimal<sup>7</sup>, up to logarithmic terms, for finding stationary points of (1) in the  
 494 convex setting in terms of  $M_*$ ,  $d_0$ , and  $\varepsilon$ .

495 **4. Technical Proofs.** This section gives the proofs of several technical results in  
 496 Section 3. More specifically, it presents the proofs of Lemma 3.1 and Proposition 3.4.

497 **4.1. Proof of Lemma 3.1.** To avoid repetition, we let

$$498 \quad (4.1) \quad \{(A_j, \tilde{x}_j, y_j, x_j, L_j)\}_{j \geq 0}$$

<sup>5</sup>This is motivated by the fact that  $m_0$  and  $M_0$  are bounded by the Lipschitz constant of  $\nabla f$ .

<sup>6</sup>See [48, Theorem 4.7].

<sup>7</sup>See [37, Section 2.2.2] or [7, Theorem 1].



499 denote the sequence of iterates generated by a single call to [Algorithm 3.2](#) and define

$$a_i := A_{i+1} - A_i, \quad \xi_i := 1 + \mu A_i,$$

500

$$\tilde{q}_{i+1}(\cdot) := \ell_{\psi^s}(\cdot; \tilde{x}_i) + \psi^n(\cdot) + \frac{\mu}{2} \|\cdot - \tilde{x}_i\|^2,$$

$$q_{i+1}(\cdot) := \tilde{q}_{i+1}(y_{i+1}) + L_{i+1} \langle \tilde{x}_i - y_{i+1}, \cdot - y_{i+1} \rangle + \frac{\mu}{2} \|\cdot - y_{i+1}\|^2,$$

501 for every  $i \geq 0$ . Recall also that each iterate in [\(4.1\)](#) is obtained in a finite number of  
502 iterations of [Algorithm 3.1](#) in view of [\(3.2\)](#) and [\(3.4\)](#).

503 We first present some basic technical properties about  $\tilde{q}$  and  $q$ .

504 **LEMMA 4.1.** *If  $\psi^s$  is  $\mu$ -strongly convex, then, for every  $j \geq 0$ ,*

505 (a)  $\tilde{q}_{j+1}(y_{j+1}) = q_{j+1}(y_{j+1})$  and  $\tilde{q}_{j+1}(\cdot) \leq q_{j+1}(\cdot) \leq \psi(\cdot)$ ;

506 (b)  $y_{j+1} = \min_{x \in \mathbb{R}^n} \{q_{j+1}(x) + L_{j+1} \|x - \tilde{x}_{j+1}\|^2/2\}$ ;

507 (c)  $x_{j+1} = \operatorname{argmin}_{x \in \mathbb{R}^n} \{a_j q_{j+1}(x) + \xi_{j+1} \|x - x_j\|^2/2\}$ .

508 *Proof.* (a) See [\[17, Lemma B.0.1\]](#).

509 (b) Let  $\Psi(\cdot) = q_{j+1}(\cdot) + L_{j+1} \|\cdot - \tilde{x}_{j+1}\|^2/2$ . It follows from the definition of  $q_{j+1}$   
510 that  $\nabla \Psi(y_{j+1}) = 0$  and, hence,  $y_{j+1}$  satisfies the optimality condition of the given  
511 inclusion.

512 (c) Using the definition of  $q_{j+1}$ , the given optimality condition of  $x_{j+1}$  holds if  
513 and only if

$$514 \quad x_{j+1} = x_j - \frac{a_j \nabla q_{j+1}(x_j)}{\xi_{j+1}} = x_j + \frac{a_j [L(y_{j+1} - \tilde{x}_j) + \mu(y_{j+1} - x_j)]}{1 + \mu A_{j+1}}$$

515 which is equivalent to the update for  $x_{j+1}$  in [Algorithm 3.2](#) (given by [Algorithm 3.1](#)). $\square$

516 The next result presents an important technical bound on the residual  $\|y_{j+1} - \tilde{x}_j\|^2$ .

517 **LEMMA 4.2.** *If  $\psi^s$  is  $\mu$ -strongly convex, then, for every  $j \geq 0$  and  $y \in \mathbb{R}^n$ ,*

$$518 \quad (4.2) \quad \frac{\mu A_{j+1}}{2} \|y_{j+1} - \tilde{x}_j\|^2 + A_{j+1} \psi(y_{j+1}) + \frac{\xi_{j+1}}{2} \|y - x_{j+1}\|^2$$

$$519 \quad \leq A_j q_{j+1}(y_j) + a_j q_{j+1}(y) + \frac{\xi_j}{2} \|y - x_j\|^2.$$

521 *Proof.* Let  $y \in \mathbb{R}^n$  be fixed. We first derive two auxiliary technical inequalities.  
522 For the first one, we use the fact that  $a_j q_{j+1} + \xi_j \|\cdot - x_j\|^2/2$  is  $\xi_{j+1}$ -strongly convex,  
523 the definition of  $\xi_{j+1}$ , and the optimality of  $x_{j+1}$  in [Lemma 4.1\(c\)](#) to obtain

$$524 \quad (4.3) \quad a_j q_{j+1}(y) + \frac{\xi_j}{2} \|y - x_j\|^2 - \frac{\xi_{j+1}}{2} \|y - x_{j+1}\|^2 \geq a_j q_{j+1}(x_{j+1}) + \frac{\xi_j}{2} \|x_{j+1} - x_j\|^2.$$

525 For the second one, let  $r_{j+1} := (A_j y_j + a_j x_{j+1})/A_{j+1}$ . Using the convexity of  $q_{j+1}$ ,  
526 the updates in [Algorithm 3.1](#) and [Algorithm 3.2](#), and [Lemma 4.1\(a\)](#)–(b), we obtain

$$527 \quad A_j q_{j+1}(y_j) + a_j q_{j+1}(x_{j+1}) + \frac{\xi_j}{2} \|x_{j+1} - x_j\|^2$$

$$528 \quad \geq A_{j+1} \left[ q_{j+1}(r_{j+1}) + \frac{\xi_j}{2a_j^2} \left\| r_{j+1} - \frac{A_j y_j + a_j x_j}{A_{j+1}} \right\|^2 \right]$$

$$529 \quad = A_{j+1} \left[ q_{j+1}(r_{j+1}) + \frac{L_{j+1}}{2} \|r_{j+1} - \tilde{x}_j\|^2 \right] \geq A_{j+1} \min_{x \in \mathbb{R}^n} \left\{ q_{j+1}(x) + \frac{L_{j+1}}{2} \|x - \tilde{x}_j\|^2 \right\}$$

(4.4)

$$530 \quad \stackrel{\text{Lemma 4.1(a)-(b)}}{=} A_{j+1} \left[ \tilde{q}_{j+1}(y_{j+1}) + \frac{L_{j+1}}{2} \|y_{j+1} - \tilde{x}_j\|^2 \right]$$

531

532 Combining (4.3), (4.4), and (3.4) with  $L = L_{j+1}$ , we conclude that

$$\begin{aligned}
533 \quad & A_j q_{j+1}(y_j) + a_j q_{j+1}(y) + \frac{\xi_j}{2} \|y - x_j\|^2 - \frac{\xi_{j+1}}{2} \|y - x_{j+1}\|^2 \\
534 \quad & \stackrel{(4.3)}{\geq} A_j q_{j+1}(y_j) + a_j q_{j+1}(x_{j+1}) + \frac{\xi_j}{2} \|x_{j+1} - x_j\|^2 \\
535 \quad & \stackrel{(4.4)}{\geq} \tilde{q}_{j+1}(y_{j+1}) + \frac{L_{j+1}}{2} \|y_{j+1} - \tilde{x}_j\|^2 \stackrel{(3.4)}{\geq} \psi(y_{j+1}) + \frac{\mu}{2} \|y_{j+1} - \tilde{x}_j\|^2. \quad \square
\end{aligned}$$

537 The following result further refines the previous bound on  $\|y_{j+1} - \tilde{x}_j\|^2$ .

538 LEMMA 4.3. *If  $\psi^s$  is  $\mu$ -strongly convex, then, for every  $j \geq 0$ ,*

$$539 \quad (4.5) \quad \mu A_{j+1} \|y_{j+1} - \tilde{x}_j\|^2 \leq \|y_{j+1} - y_0\|^2 - \xi_{j+1} \|y_{j+1} - x_{j+1}\|^2.$$

540 *Proof.* Let  $j \geq 0$  be fixed and suppose  $\psi^s$  is  $\mu$ -strongly convex. Moreover, define

$$541 \quad \Psi_i := A_i [\psi(y_i) - \psi(y_j)] + \frac{\xi_i}{2} \|y_j - x_i\|^2 \quad \forall i \geq 0.$$

542 Using Lemma 4.2 with  $y = y_j$ , Lemma 4.1(a), the fact that  $a_j = A_{j+1} - A_j$ , and the  
543 definition of  $\Psi_i$  above, we have that for every  $i \geq 0$ ,

$$\begin{aligned}
544 \quad & \frac{\mu A_{i+1}}{2} \|y_{i+1} - \tilde{x}_i\|^2 \\
545 \quad & \stackrel{(4.2)}{\leq} A_i q_{i+1}(y_i) + a_i q_{i+1}(y_j) + \frac{\xi_i}{2} \|y_i - x_i\|^2 - \Psi_{i+1} - A_{i+1} \psi(y_j) \\
546 \quad & \stackrel{\text{Lemma 4.1(a)}}{\leq} A_i \psi(y_i) + a_i \psi(y_j) + \frac{\xi_i}{2} \|y_i - x_i\|^2 - \Psi_{i+1} - A_{i+1} \psi(y_j) \\
547 \quad & = \Psi_i - \Psi_{i+1}.
\end{aligned}$$

549 Summing the above inequality from  $i = 0$  to  $j$  and using the fact that  $A_{i+1} \geq 0$  for  
550 every  $i$  and  $(x_0, A_0, \xi_0) = (y_0, 0, 1)$ , we conclude that

$$\begin{aligned}
551 \quad & \frac{\mu A_{j+1}}{2} \|y_{j+1} - \tilde{x}_j\|^2 \leq \sum_{i=0}^j \frac{\mu A_{i+1}}{2} \|y_{i+1} - \tilde{x}_i\|^2 \leq \Psi_0 - \Psi_{j+1} \\
552 \quad & = \frac{\xi_0}{2} \|y_j - x_0\|^2 - \frac{\xi_j}{2} \|y_{j+1} - x_{j+1}\|^2 = \frac{1}{2} \|y_j - y_0\|^2 - \frac{\xi_j}{2} \|y_{j+1} - x_{j+1}\|^2. \quad \square
\end{aligned}$$

554 We are now ready to prove Lemma 3.1.

555 *Proof of Lemma 3.1.* (a) See [17, Lemma B.0.2] for the bound on  $A_{j+1}$ . The  
556 bound on  $L_j$  follows from how Algorithm 3.1 is called in Algorithm 3.2, the update  
557 rule for  $L$  in Algorithm 3.1, and (3.2) which follows from assumption (B2).

558 (b) Using the optimality of  $y_{j+1}$  given by Algorithm 3.1 and Algorithm 3.2 and  
559 the definition of  $r_{j+1}$ , it follows that

$$560 \quad 0 \in \nabla \psi^s(\tilde{x}_j) + \partial \psi^n(y_{j+1}) + (L_{j+1} + \mu)(y_{j+1} - \tilde{x}_j) = \nabla \psi^s(y_{j+1}) + \partial \psi^n(y_{j+1}) - r_{j+1}.$$

561 (c) The first bound in (3.4) is an immediate consequence of Lemma 4.3. For the  
562 second bound in (3.4), note that part (b) and the assumption that  $\psi^s$  implies that  
563  $r_{j+1} \in \partial \psi(y_{j+1})$ . The conclusion now follows from the previous inclusion and the  
564 definition of the subdifferential.

565 (d) Suppose  $A_{j+1} \geq \mathcal{A}_{\mu, \bar{L}} := \mathcal{A}_{\mu, \bar{L}}(\sigma, \theta)$  and (3.4) holds. We separate this proof  
 566 into two parts. We first prove the bound in (3.4). Using the definitions of  $r_{j+1}$   
 567 and  $\bar{L}$ , part (c), the fact that  $\mu \leq L_0 \leq L_{j+1}$ , assumption  $\langle \mathbf{B2} \rangle$ , and the relation  
 568  $(a+b)^2 \leq 2a^2 + 2b^2$  for  $a, b \in \mathbb{R}$ , we have that

$$\begin{aligned}
 569 \quad \|r_{j+1}\|^2 &= \|\nabla\psi^s(y_{j+1}) - \nabla\psi^s(\tilde{x}_j) + (L_{j+1} + \mu)(\tilde{x}_j - y_{j+1})\|^2 \\
 570 &\leq 2\|\nabla\psi^s(y_{j+1}) - \nabla\psi^s(\tilde{x}_j)\|^2 + 2(L_{j+1} + \mu)^2\|\tilde{x}_j - y_{j+1}\|^2 \\
 571 &\leq 2[L_*^2 + (L_{j+1} + \mu)^2]\|\tilde{x}_j - y_{j+1}\| \leq 16\bar{L}^2\|\tilde{x}_j - y_{j+1}\|^2 \\
 572 &\stackrel{(4.5)}{\leq} \frac{16\bar{L}}{\mu A_{j+1}}\|y_{j+1} - y_0\|^2. \\
 573
 \end{aligned}$$

574 It follows from the above bound and the definition of  $\mathcal{A}_{\mu, \bar{L}}$  that

$$575 \quad \|r_{j+1}\|^2 \leq \frac{16\bar{L}}{\mu A_{j+1}}\|y_{j+1} - y_0\|^2 \leq \frac{16\bar{L}^2}{\mu \mathcal{A}_{\mu, \bar{L}}}\|y_{j+1} - y_0\|^2 \leq \sigma^2\|y_{j+1} - y_0\|^2$$

576 and, hence, the first condition of (3.6) holds.

577 To show the second condition of (3.6), let  $\gamma := \sqrt{(2-\theta)/\theta}$ . Using the fact that  
 578  $\gamma \in (0, 1)$ , (3.4),  $\mu \leq L_{j+1}$ , and the bound

$$579 \quad \|a+b\|^2 \leq (1+\gamma)\|a\|^2 + (1+\gamma^{-1})\|b\|^2 \quad \forall a, b \in \mathbb{R}^n,$$

580 we then have that

$$\begin{aligned}
 581 \quad \|r_{j+1}\|^2 &\stackrel{(3.4)}{\leq} \frac{L^2}{\mu A_{j+1}}\|y_{j+1} - y_0\|^2 \leq \frac{4(\mu + L_{j+1})^2}{\mu A_{j+1}}\|y_{j+1} - y_0\|^2 \\
 582 &\leq \frac{16\bar{L}^2}{\mu \mathcal{A}_{\mu, \bar{L}}(\sigma, \theta)}\|y_{j+1} - y_0\|^2 \leq \frac{\gamma^2}{4}\|y_{j+1} - y_0\|^2 \stackrel{\gamma \in (0,1)}{\leq} \left(\frac{\gamma}{1+\gamma}\right)^2 \|y_{j+1} - y_0\|^2 \\
 583 &\leq \left(\frac{\gamma}{1+\gamma}\right)^2 (1+\gamma)\|r_{j+1} + y_{j+1} - y_0\|^2 + \left(\frac{\gamma}{1+\gamma}\right)^2 \left(1 + \frac{1}{\gamma}\right) \|r_{j+1}\|^2 \\
 584 &= \frac{\gamma^2}{1+\gamma}\|r_{j+1} + y_{j+1} - y_0\|^2 + \frac{\gamma}{1+\gamma}\|r_{j+1}\|^2, \\
 585
 \end{aligned}$$

586 which implies  $\|r_{j+1}\|^2 \leq \gamma^2\|r_{j+1} + y_{j+1} - y_0\|^2$ . It then follows from the second bound  
 587 in (3.4) and the previous inequality that

$$\begin{aligned}
 588 \quad 2[\psi(y_0) - \psi(y_{j+1})] &\stackrel{(3.5)}{\geq} 2\langle r_{j+1}, y_0 - y_{j+1} \rangle \\
 589 &= \|r_{j+1} + y_0 - y_{j+1}\|^2 - \|r_{j+1}\|^2 - \|y_0 - y_{j+1}\|^2 \\
 590 &\geq (1-\gamma^2)\|r_{j+1} + y_0 - y_{j+1}\|^2 - \|y_0 - y_{j+1}\|^2 \\
 591 &= \frac{2}{\theta}\|r_{j+1} + y_0 - y_{j+1}\|^2 - \|y_0 - y_{j+1}\|^2. \quad \square \\
 592
 \end{aligned}$$

## 593 4.2. Proof of Proposition 3.4.

594 *Proof of Proposition 3.4.* (a) Note that the  $k$ -th successful call of Algorithm 3.2  
 595 is such that its input  $\psi^s$  has the curvature pair

$$596 \quad (4.6) \quad (L_{k+1}^-, L_{k+1}^+) := \left( \max \left\{ 0, \frac{m_*}{2m_{k+1}} - 1 \right\}, \frac{M_*}{2m_{k+1}} + 1 \right).$$

597 Hence, it follows from Step 1 of [Algorithm 3.3](#), [Proposition 3.2\(b\)](#) with  $\mu = 1/2$ , and  
598 the definition of  $\bar{m}$  imply that the last call of [Algorithm 3.2](#) at the  $k$ -th iteration of  
599 [Algorithm 3.4](#) obtains  $m_{k+1}$  being at most  $\alpha m_k \leq \bar{m}$ . Consequently,  $\{1/m_k\}$  (resp.  
600  $\{m_k\}$ ) is bounded below by  $1/\bar{m}$  (resp. bounded above by  $\bar{m}$ ). The fact that  $\{1/m_j\}$   
601 is bitonic follows from the the definition of  $\hat{m}$  in Step 1 of [Algorithm 3.4](#), the call to  
602 [Algorithm 3.3](#) in of [Algorithm 3.4](#), and the fact that in [Algorithm 3.4](#) the returned  
603 scalar  $m$  in is always lower bounded by the input  $\hat{m}$ . To show the bound on  $M_k$ , note  
604 that the curvature pair of  $\psi^s$  in (4.6) implies that  $\nabla\psi^s$  is  $L_*$ -Lipschitz continuous  
605 where  $L_* = \max\{L_{k+1}^-, L_{k+1}^+\}$ . It then follows from the upper previous bound on  
606  $m_{k+1}$  and [Lemma 3.1\(a\)](#) that

$$607 \quad \frac{M_k}{2m_{k+1}} + 1 \leq \frac{M_{k+1}}{2m_{k+1}} + 1 \leq \beta \left[ \frac{\max\{M_0, M_*\}}{2m_{k+1}} + 1 \right]$$

$$608 \quad \leq \frac{\beta [\max\{M_0, M_*\} + 2\bar{m}]}{2m_{k+1}} = \frac{\bar{M}}{2m_{k+1}},$$

$$609$$

610 which immediately implies  $M_{k+1} \geq M_k$  and  $M_{k+1} \leq \bar{M}$ .

611 (b) Let an outer iteration index  $k \geq 1$  be fixed and define

$$612 \quad \mathcal{L}_\ell := \frac{\bar{M}}{2m_k \alpha^\ell} + 1, \quad \mathcal{I}_\ell := \left[ 1 + 4\sqrt{\mathcal{L}_\ell P_0} \right], \quad \bar{\ell} := 1 + \log_\alpha(m_{k+1}/m_k),$$

613 where  $P_0$  is as in (3.14). Using [Proposition 3.2\(a\)](#) with  $(\mu, \sigma) = (1/2, 1/4)$ , part (a),  
614 the fact that  $P_0 \geq 1$ , and assumptions [\(A1\)](#)–[\(A2\)](#), it follows that the number of inner  
615 iterations performed by [Algorithm 3.4](#) at outer iteration  $k$  is bounded above by

$$616 \quad \sum_{\ell=0}^{\bar{\ell}} \mathcal{I}_\ell \leq 2 \sum_{\ell=0}^{\bar{\ell}} \left( 1 + 4\sqrt{\mathcal{L}_\ell P_0} \right) \leq 2 \sum_{\ell=0}^{\bar{\ell}} \left( 1 + 4 \left[ \sqrt{\frac{\bar{M}}{2m_k \alpha^\ell}} + 1 \right] P_0 \right)$$

$$617 \quad \leq 10P_0 \sum_{\ell=0}^{\bar{\ell}} \left( \sqrt{\frac{\bar{M}}{2m_k \alpha^\ell}} + 1 \right) = 10 \left[ \bar{\ell} + \sqrt{\frac{\bar{M}}{2m_k}} \sum_{\ell=0}^{\bar{\ell}} \alpha^{-\ell/2} \right] P_0$$

$$618 \quad = 10 \left[ \bar{\ell} + \sqrt{\frac{\bar{M}}{2m_k}} \left( \frac{1 - \alpha^{-\bar{\ell}/2}}{\sqrt{\alpha} - 1} \right) \right] P_0 \leq 10 \left[ \bar{\ell} + \frac{1}{\sqrt{\alpha} - 1} \sqrt{\frac{\bar{M}}{2m_k}} \right] P_0$$

$$619 \quad \leq 20 \left[ 1 + \log_\alpha \frac{m_{k+1}}{m_k} + \frac{1}{\sqrt{\alpha} - 1} \sqrt{\frac{\bar{M}}{2m_k}} \right] P_0.$$

$$620$$

621 (c) In view of [Proposition 3.4\(a\)](#), let  $\bar{K}$  be an index satisfying

$$622 \quad \frac{\bar{K} - 1}{\bar{m}} \leq \sum_{k=0}^{\bar{K}-2} \frac{1}{m_{k+1}} < \frac{2\theta\Delta_0}{\varepsilon^2} \leq \sum_{k=0}^{\bar{K}-1} \frac{1}{m_{k+1}}.$$

623 Using [Lemma 3.3](#), the choice of inputs to [Algorithm 3.2](#), and [Lemma 2.1\(b\)](#), and the  
624 last of the above inequalities, we have that

$$625 \quad \inf_{0 \leq k \leq \bar{K}-1} \|v_{j+1}\|^2 \leq \frac{2\theta\Delta_0}{\sum_{k=0}^{\bar{K}-1} m_{k+1}^{-1}} \leq \varepsilon^2.$$

626 Hence, because of the termination condition in Step 2 of [Algorithm 3.4](#), it follows that  
 627 the number of outer iterations  $K(\varepsilon)$  is at most  $\bar{K}$ . Using the fact that  $m_{k+1} > 0$  for  
 628 every  $k \geq 0$ , the bounds in [\(3.16\)](#) immediately follow.

629 (d) Since  $f$  is convex,  $\psi^s$  in [\(3.12\)](#) is  $(1/2)$ -strongly convex at every (outer) it-  
 630 eration of [Algorithm 3.4](#). Consequently, using [Proposition 3.2\(b\)](#) with  $\mu = 1/2$ , the  
 631 inputs and outputs given to [Algorithm 3.2](#) by [Algorithm 3.3](#), and the definition of  
 632  $\psi^s$ , it follows that every call to [Algorithm 3.3](#) by [Algorithm 3.4](#) stops at (line search)  
 633 iteration  $\ell = 0$ , i.e., the conditions in [\(3.11\)](#) are satisfied when they are first checked.  
 634 Using the update rule in Step 1 of [Algorithm 3.4](#) and the previous conclusion, we have  
 635 that  $m_{k+1} = m_k/\alpha$  for every  $k \geq 0$ . Inductively, it then follows that  $m_k = \alpha^{-k}m_0$   
 636 for every  $k \geq 0$ . We now prove the claimed complexity bound. In view of the fact  
 637 that  $\{1/m_k\}$  is bounded below from part (a), let  $\bar{K}$  be the smallest index such that  
 638  $\bar{K} \geq 2$  and

$$639 \quad (4.7) \quad \sum_{k=1}^{\bar{K}-1} \frac{1}{m_{k+1}} \leq \frac{\alpha^2}{m_0} + \frac{2\theta m_0 R_{m_0}(z_0)}{\varepsilon^2} \leq \sum_{k=1}^{\bar{K}} \frac{1}{m_{k+1}}$$

640 Using the fact that  $\{m_k\}$  is nonincreasing, [\(2.7\)](#) with  $\tilde{m} = m_0\alpha^{-\bar{K}}$  and  $\nu = m_0/\tilde{m}$ ,  
 641 the identity  $m_k = \alpha^{-k}m_0$ , and the same type of arguments as in part (c), we have  
 642 that

$$643 \quad \min_{1 \leq k \leq \bar{K}-1} \|v_{k+1}\|^2 \leq \frac{2\theta\nu\tilde{m}R_{\nu\tilde{m}}(z_0)}{\sum_{k=1}^{\bar{K}-1} m_{k+1}^{-1}} = \frac{2\theta m_0 R_{m_0}(z_0)}{-\alpha^2 m_0^{-1} + \sum_{k=1}^{\bar{K}} m_{k+1}^{-1}} \leq \varepsilon^2,$$

644 and, hence, the number of outer iterations  $K(\varepsilon)$  is bounded above by  $\bar{K}$ . It now  
 645 remains to show that  $\bar{K}$  is bounded above by the expression on the right-hand side  
 646 of [\(3.17\)](#). Using the identity  $m_k = \alpha^{-k}m_0$  and the right-hand side of [\(4.7\)](#), we have

$$647 \quad \frac{\alpha^2}{m_0} + \frac{2\theta m_0 R_{m_0}(z_0)}{\varepsilon^2} \geq \sum_{k=1}^{\bar{K}-1} \frac{1}{m_{k+1}} = \frac{\alpha^2}{m_0} \sum_{k=0}^{\bar{K}-2} \alpha^k \geq \frac{\alpha^{\bar{K}-1}}{m_0(\alpha-1)} \geq \frac{\alpha^{\bar{K}-2}}{m_0}.$$

648 Applying the function  $\log_\alpha(\cdot)$  to both sides of the above inequality and re-arranging  
 649 terms yields the desired bound on  $\bar{K}$ .

650 (e) Using the definition of  $v_{k+1}$  and [Lemma 3.3](#) with  $\psi^s$  as in [\(3.12\)](#), we have

$$\begin{aligned} 651 \quad v_{k+1} &\in 2m_{k+1} [\nabla\psi^s(z_{k+1}) + \partial\psi^s(z_{k+1})] + 2m_{k+1}(z_k - z_{k+1}) \\ 652 \quad &= 2m_{k+1} \left[ \frac{\nabla f(z_{k+1})}{2m_{k+1}} + (z_{k+1} - z_k) + \frac{\partial h(z_{k+1})}{m_{k+1}} \right] + 2m_{k+1}(z_k - z_{k+1}) \\ 653 \quad &= \nabla f(z_{k+1}) + \partial h(z_{k+1}). \end{aligned}$$

655 The fact that the last iterate solves Problem  $\mathcal{CO}$  follows from the above inclusion and  
 656 the termination condition in Step 2 of [Algorithm 3.4](#).  $\square$

657 **5. Applications.** This section describes a few possible applications of [Algo-](#)  
 658 [rithm 3.4](#) in more general optimization frameworks.

659 *Min-Max Smoothing.* In [\[22\]](#), a smoothing framework was proposed for finding  
 660  $\varepsilon$ -stationary points of the nonconvex-concave min-max problem

$$661 \quad (5.1) \quad \min_{x \in \mathbb{R}^n} \max_{y \in \mathbb{R}^l} [\phi(x, y) + h(x)]$$

662 where  $h$  is as in assumption [\(A1\)](#),  $\phi(\cdot, y)$  is  $m_x$ -weakly convex and differentiable,  
 663  $-\phi(x, \cdot)$  is proper closed convex, and  $\nabla_x \phi(\cdot, \cdot)$  is Lipschitz continuous.

664 The framework considers finding an  $\varepsilon$ -stationary point of  $h$  plus a smooth approx-  
 665 imation  $\hat{p}$  of  $\max_{y \in Y} \phi(\cdot, y)$ . Choosing a special smoothing constant such that the cur-  
 666 vature pair  $(\hat{m}, \hat{M})$  of  $\hat{p}$  satisfies  $\hat{m} = m_x$  and  $\hat{M} = \Theta(\varepsilon^{-1} D_y)$  (resp.  $\hat{M} = \Theta(D_y^2 \varepsilon^{-2})$ ),  
 667 where  $D_y$  is diameter of  $\text{dom}(-\phi(x, \cdot))$ , it was shown that an  $\varepsilon$ -stationary point of  $\hat{p}$   
 668 yields an  $\varepsilon$ -primal-dual (resp. directional) stationary point of [\(5.1\)](#).

669 If we use PF.APD with  $m_0 = \varepsilon$  to obtain an  $\varepsilon$ -stationary point of  $\hat{p}$  as above,  
 670 then an  $\varepsilon$ -primal-dual (resp. directional) stationary point of [\(5.1\)](#) is obtained in  
 671  $\tilde{\mathcal{O}}(\varepsilon^{-2.5})$  (resp.  $\tilde{\mathcal{O}}(\varepsilon^{-3})$ ) inner iterations, and this matches, up to logarithmic terms,  
 672 the complexity bounds for the smoothing method in [\[22\]](#). Moreover, when  $\phi(\cdot, y)$   
 673 is convex, the above complexity is  $\tilde{\mathcal{O}}(\varepsilon^{-1})$  (resp.  $\tilde{\mathcal{O}}(\varepsilon^{-1.5})$ ), and this appears to be  
 674 the first parameter-free approach that could be used for min-max optimization. This  
 675 approach also has the strong advantage that it does not need to know  $D_y$ .

676 *Penalty Method.* In [\[19\]](#), a penalty method is proposed for finding  $\varepsilon$ -KKT points  
 677 of the linearly-constrained nonconvex optimization problem

$$678 \quad (5.2) \quad \min_{x \in \mathbb{R}^n} \{ \phi(x) := f(x) + h(x) : Ax = b \}$$

679 where  $(f, h)$  are as in [\(A1\)](#)–[\(A3\)](#). It was shown that if the penalty method uses an  
 680 algorithm  $\mathcal{A}$  that needs  $\mathcal{O}(T_{m, M}(\varepsilon))$  iterations to obtain an  $\varepsilon$ -stationary point of  $\phi$ ,  
 681 then the total number of inner iterations of the penalty method (for finding an  $\varepsilon$ -KKT  
 682 point) is  $\tilde{\mathcal{O}}(T_{m, \varepsilon^{-2}}(\varepsilon))$ .

683 If we use the PF.APD with  $m_0 = \varepsilon$  as algorithm  $\mathcal{A}$  above, then an  $\varepsilon$ -KKT point  
 684 of [\(5.2\)](#) is obtained in  $\tilde{\mathcal{O}}(\varepsilon^{-3})$  inner iterations which matches the complexity bound  
 685 for the particular penalty method in [\[19\]](#) (which uses the AIPP in [\[18\]](#) for algorithm  
 686  $\mathcal{A}$ ). Moreover, when  $f$  is convex, the above complexity is  $\tilde{\mathcal{O}}(\varepsilon^{-1.5})$ . Like in the  
 687 above discussion for min-max smoothing, this appears to be the first parameter-free  
 688 approach used for linearly-constrained composite optimization.

689 **6. Numerical Experiments.** This section presents experiments that demon-  
 690 strate the numerical efficiency of PF.APD. Comments about the results are given in  
 691 [Subsection 6.4](#).

692 We first describe the benchmark algorithms, the implementation of APD, and  
 693 the computing environment. The benchmark algorithms are instances of PGD, AIPP,  
 694 ANCF, and UPF described in [Section 1](#) and [Table 1.1](#). Specifically, AIPP uses  $\sigma =$   
 695  $1/4$ , ANCF uses  $\theta = 1.25$ , and UPF uses  $\gamma_1 = \gamma_2 = 0.4$ ,  $\gamma_3 = 1$ ,  $\beta_0 = 1$ , and  $\hat{\lambda}_0 = 1$ .  
 696 Moreover, UPF uses  $\hat{\lambda}_k$  for the initial estimate of  $\hat{\lambda}_{k+1}$  for  $k \geq 1$  and AIPP stops its  
 697 call of ACG when the condition  $\|u_j\|^2 + 2\eta_j \leq \sigma \|x_0 - x_j + u_j\|^2$  holds (inside of ACG)  
 698 instead of prescribing a fixed number of ACG iterations. The implementations for  
 699 ANCF and UPF were generously provided by the respective authors of [\[25\]](#) and [\[13\]](#),  
 700 while the author implemented AIPP and PGD.<sup>8</sup> Note that we did not consider the  
 701 VAR-FISTA method in [\[43\]](#) because: (i) its steps were similar to ANCF and (ii) we  
 702 already had a readily available and optimized code for the ANCF method.

703 The implementation of PF.APD, abbreviated as APD, is as in [Algorithm 3.4](#) with  
 704  $\alpha = \beta = 2$ ,  $\hat{m} = m_k$  for every  $k \geq 1$ , and the following additional updates at the  
 705 beginning of every call to [Algorithm 3.2](#) and the  $(k + 1)$ <sup>th</sup> iteration of [Algorithm 3.4](#),

<sup>8</sup>See [https://github.com/wwkong/nc\\_opt/tree/master/tests/papers/apd](https://github.com/wwkong/nc_opt/tree/master/tests/papers/apd) for the source code of the experiments.

706 respectively:

$$707 \quad (6.1) \quad L_0 \leftarrow \frac{L_0}{1 + \beta/2}, \quad m_{k+1} \leftarrow \max \left\{ m_0, \frac{m_{k+1}}{1 + \alpha/2} \right\}.$$

708 This is done to allow a possible decrease in both of the curvature estimates. While  
 709 we do not show convergence of this modified PF.APD, we believe that convergence  
 710 can be established using similar techniques as in [35]. It is worth mentioning that  
 711 the modification in (6.1) substantially improves upon the numerical performance of  
 712 PF.APD compared to the version given in Algorithm 3.4.

713 All experiments were run in MATLAB 2023a under a 64-bit Windows 11 machine  
 714 with an Intel Core i7-10700K processor and 16 GB of RAM. All algorithms except  
 715 AIPP use an initial curvature estimate of  $(m_0, M_0) = (1, 1)$ , and each algorithm stops  
 716 when it finds a pair  $(\bar{z}, \bar{v})$  solving Problem  $\mathcal{CO}$  for some  $\varepsilon > 0$ . A time limit of 1200  
 717 (resp. 2400) seconds was prescribed for the problems in Subsection 6.1 and 6.3 (resp.  
 718 Subsection 6.2). We also set an (innermost) iteration limit of 500000 (resp. 10000)  
 719 for Subsection 6.2 (resp. Subsection 6.3).

720 **6.1. Quadratic Semidefinite Programming.** The problem of interest is the  
 721 400-variable nonconvex quadratic semidefinite programming (QSDP) problem

$$722 \quad (6.2) \quad \min_{Z \in \mathbb{R}^{35 \times 35}} -\frac{\eta_1}{2} \|DB(Z)\|_2^2 + \frac{\eta_2}{2} \|\mathcal{A}(Z) - b\|_2^2,$$

$$723 \quad \text{s.t. } \text{tr}(Z) = 1, \quad Z \in \mathcal{S}_+^{35},$$

725 where  $\mathcal{S}_+^n$  is the  $n$ -dimensional positive semidefinite cone,  $\text{tr}(Z)$  is the trace of a matrix,  
 726  $b \in \mathbb{R}^{10}$ ,  $D \in \mathbb{R}^{10 \times 10}$  is a diagonal matrix with nonzero entries randomly generated  
 727 from  $\{1, \dots, 1000\}$ ,  $(\eta_1, \eta_2) \in \mathbb{R}_{++}^2$  are chosen to yield a particular curvature pair, and  
 728  $\mathcal{A}, \mathcal{B} : \mathcal{S}_+^{20} \mapsto \mathbb{R}^{10}$  are linear operators defined by

$$729 \quad [\mathcal{A}(Z)]_j = A_j \bullet Z, \quad [\mathcal{B}(Z)]_j = B_j \bullet Z$$

730 for matrices  $\{A_j\}_{j=1}^{10}, \{B_j\}_{j=1}^{10} \subseteq \mathbb{R}^{20 \times 20}$ . Moreover, the entries in these matrices and  
 731  $b$  were sampled from the uniform distribution on  $[0, 1]$ .

732 To build the decomposition in (1.1), we set  $f$  equal to the objective function of  
 733 (6.2),  $h$  equal to the indicator function of the constraint set of (6.2). The starting  
 734 point was set to  $z_0 = I_{20}/20$ , where  $I_{20}$  is an identity matrix, and the tolerance was  
 735 set to  $\varepsilon = 10^{-6}(1 + \|\nabla f(z_0)\|_2)$ .

$m, M$	# of Function Evaluations				# of Gradient Evaluations				Runtime (seconds)			
	UPF	ANCF	AIPP	APD	UPF	ANCF	AIPP	APD	UPF	ANCF	AIPP	APD
$10^2, 10^4$	6.5E4	2.1E4	7.1E4	<b>1.1E3</b>	1.3E4	1.6E4	6.7E4	<b>2.1E3</b>	9.2E1	2.7E1	1.1E2	<b>3.2E0</b>
$10^2, 10^5$	1.9E5	4.4E4	4.1E5	<b>3.3E3</b>	3.8E4	3.3E4	3.9E5	<b>6.7E3</b>	2.6E2	5.8E1	6.5E2	<b>9.9E0</b>
$10^2, 10^6$	3.0E5	5.9E4	7.6E5	<b>7.1E3</b>	6.1E4	4.4E4	7.0E5	<b>1.4E4</b>	4.3E2	7.9E1	1.2E3	<b>2.1E1</b>
$10^3, 10^7$	3.0E5	5.9E4	7.6E5	<b>1.0E4</b>	6.1E4	4.4E4	6.9E5	<b>2.0E4</b>	4.3E2	8.1E1	1.2E3	<b>3.0E1</b>
$10^2, 10^7$	3.3E5	6.6E4	2.6E5	<b>1.2E4</b>	6.5E4	5.0E4	1.3E5	<b>2.4E4</b>	4.5E2	8.6E1	2.5E2	<b>3.4E1</b>
$10^4, 10^7$	5.8E5	1.4E5	8.8E4	<b>2.0E4</b>	1.2E5	1.1E5	4.4E4	<b>4.1E4</b>	7.9E2	1.9E2	8.3E1	<b>5.8E1</b>

TABLE 6.1

*Unique function evaluations, unique gradient evaluations, and runtimes in the QSDP experiments for different curvature pairs  $(m, M)$ . The bolded numbers indicate the best algorithm in terms of the number of evaluations (less is better) and runtime (less is better). Entries marked with “-” are those that did not terminate within the prescribed time limit.*

736 Table 6.1 reports the number of unique function evaluations, unique gradient eval-  
 737 uations, and runtime (in seconds) for different curvature pairs  $(m, M)$ , and Figure 6.1

738 plots the minimum norm of the normalized stationarity residual  $\|\bar{v}\|$  over iteration  
 739 count for each algorithm and curvature pairs  $(m, M) = (10^2, 10^4)$ ,  $(10^2, 10^5)$ , and  
 $(10^2, 10^6)$ .

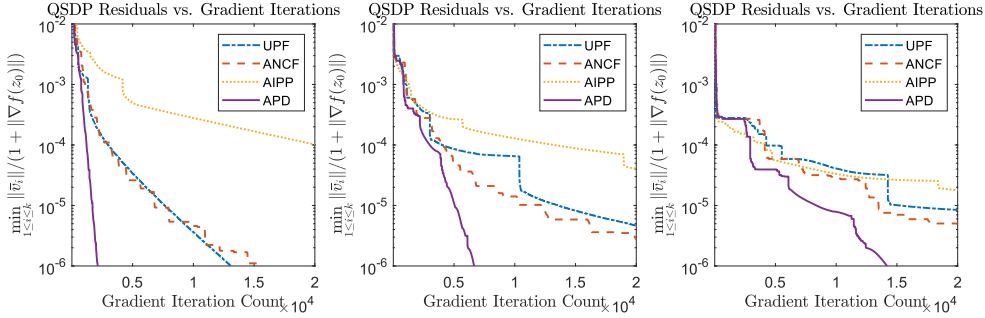


FIGURE 6.1. Plots of the minimum norm of the normalized stationarity residual  $\|\bar{v}\|$  over iteration count in the QSDP experiments. The curvature pairs for the plots are  $(10^2, 10^4)$ ,  $(10^2, 10^5)$ , and  $(10^2, 10^6)$  from left-to-right.

740

741 **6.2. Sparse Vector Recovery.** The problem of interest is the penalized sparse  
 742 vector recovery (SVR) problem [45]

743 (6.3) 
$$\min_{z \in \mathbb{R}^n} \frac{1}{2} \|Az - b\|_2^2 + \frac{\tau}{2} \|z\|_2^2 + \text{LPL}_{\gamma, \delta}(\|z\|_2)$$

744 where  $\tau = 10^{-2}$ ,  $A \in \mathbb{R}^{\ell \times p}$  with  $\ell \geq p$ ,  $b = A\tilde{u}$  where  $u$  is a random vector  
 745 whose entries are sampled uniformly from  $[0, 1]$ , for  $(\gamma, \delta) = (10, 10^{-1})$ , the func-  
 746 tion  $\text{LPL}_{\gamma, \delta}(z) = \gamma[1 - \exp(-z/\delta)]$  is the concave Laplace penalty function [44] at  $z$ .  
 747 The goal of this problem is to find a sparse vector  $\hat{z}$  such that  $A\hat{z}$  is close to  $b$ .

748 Each matrix  $A$  is built from a recommender dataset where each entry corresponds  
 749 to a user-item rating. Specifically, the datasets were taken from the well-known Jester,  
 750 MovieLens 100K, and FilmTrust datasets and the musical instruments and patio,  
 751 lawn, and garden products Amazon Review datasets published by the University of  
 752 California San Diego. The dimensions  $(\ell, p)$  of each matrix generated by the pre-  
 753 vious datasets were  $(24938, 100)$ ,  $(9724, 610)$ ,  $(2071, 1508)$ ,  $(1429, 900)$ ,  $(1686, 962)$ ,  
 754 respectively.

755 To put (6.3) into the form of (1.1), we use the decomposition given in [46] where  
 756  $h$  is a multiple of the 1-norm and  $f$  is the function in (6.3) minus  $h$ . The starting  
 757 point  $z_0$  was set to be a vector whose entries are all equal to  $p$ , and the tolerance  
 758 was set to  $\varepsilon = 10^{-10}(1 + \|\nabla f(z_0)\|_2)$ . Following the analysis in [46], AIPP uses the  
 759 curvature pair  $(m, M) = (2\gamma/\delta^2, \tau + \sigma_{\max}^2(A))$ , where  $\sigma_{\max}(A)$  is the largest singular  
 760 value of  $A$ .

761 Table 6.2 reports the unique function evaluations, unique gradient evaluations,  
 762 and runtime (in seconds) for the different datasets mentioned above, and Figure 6.2  
 763 plots the minimum norm of the normalized stationarity residual  $\|\bar{v}\|$  over the gradient  
 764 count for each algorithm and the first, second, and fourth row of Table 6.2.

765 **6.3. Low-Rank Matrix Completion.** The problem of interest is the penalized  
 766 nonconvex low-rank matrix completion (LRMC) problem [45, 46]



$\ell, p$	# of Function Evaluations				# of Gradient Evaluations				Runtime (seconds)			
	UPF	ANCF	AIPP	APD	UPF	ANCF	AIPP	APD	UPF	ANCF	AIPP	APD
1429, 900	6.9E3	3.5E3	1.5E4	<b>3.7E2</b>	8.0E2	2.6E3	1.1E4	<b>7.3E2</b>	2.7E0	1.2E0	1.1E1	<b>3.4E-1</b>
1686, 962	2.9E4	1.1E4	7.7E4	<b>2.6E3</b>	4.9E3	8.2E3	5.8E4	<b>3.8E3</b>	1.3E1	4.1E0	6.0E1	<b>2.4E0</b>
9724, 610	3.9E4	4.3E4	6.2E4	<b>3.2E3</b>	6.3E3	3.2E4	3.3E4	<b>6.2E3</b>	3.6E1	3.5E1	8.4E1	<b>6.0E0</b>
24938, 100	5.7E5	2.4E5	9.8E5	<b>2.5E4</b>	1.1E5	1.8E5	5.0E5	<b>4.8E4</b>	1.7E2	5.0E1	4.3E2	<b>1.6E1</b>
2071, 1508	-	2.9E5	-	<b>2.8E4</b>	-	2.2E5	-	<b>5.5E4</b>	-	1.3E3	-	<b>2.6E2</b>

TABLE 6.2

Unique function evaluations, unique gradient evaluations, and runtimes in the SVR experiments for different datasets and their dimensions  $(\ell, p)$ . The bolded numbers indicate the best algorithm in terms of the number of evaluations (less is better) and runtime (less is better). Entries marked with “-” are those that did not terminate within the prescribed time or iteration limit.

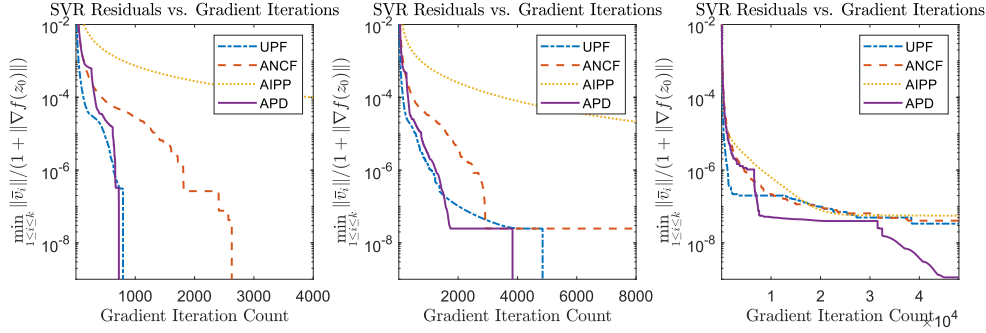


FIGURE 6.2. Plots of the minimum norm of the normalized stationarity residual  $\|\bar{v}\|$  over iteration count in the SVR experiments. The dimensions and upper curvature  $(\ell, p)$  for the plots are  $(1429, 900)$ ,  $(1686, 962)$ , and  $(24938, 100)$  from left-to-right.

$$(6.4) \quad \min_{Z \in \mathbb{R}^{\ell \times p}} \frac{1}{2} \|\Pi_{\Omega}(Z) - \Pi_{\Omega}(X)\|_F^2 + \frac{\tau}{2} \|Z\|_F^2 + (\text{MCP}_{\gamma, \delta} \circ \sigma)(Z),$$

where  $\tau = 10^{-7}$ ,  $X \in \mathbb{R}^{\ell \times p}$  is a reference image,  $\sigma : \mathbb{R}^{\ell \times p} \mapsto \mathbb{R}^{\min\{\ell, p\}}$  maps a matrix to its vector of singular values, for  $(\gamma, \delta) = (450, 10^{-4})$  the function  $\text{MCP}_{\gamma, \delta}(z)$  is the minimax concave penalty (MCP) function [47] at  $z$  (which takes value  $\gamma z - z^2/(2\delta)$  if  $z \leq \gamma\delta$  and  $\gamma^2\delta/2$  otherwise), and, for a given corrupted image  $\Omega$ , the function  $\Pi_{\Omega} : \mathbb{R}^{\ell \times p} \mapsto \mathbb{R}^{\ell \times p}$  is the projection operator that zeros out entries of its input where the corresponding entry in  $\Omega$  is zero. The goal of this problem is to fill in the zero entries of a corrupted image  $\Omega$  of  $X$  so that the resulting image  $\hat{Z}$  is close to  $X$ .

To put (6.4) into the form of (1.1), we use the decomposition given in [46] where  $h$  is a multiple of the nuclear norm and  $f$  is the function in (6.4) minus  $h$ . Experiments were run on different reference images  $X$  given in the first row of Figure 6.3 and  $\Omega$  was set to be a corrupted version of  $X$  where we add Gaussian noise with a 100 dB signal-to-noise ratio and remove 30% of the resulting pixels. For illustration, two corrupted images can be found in the first columns of the last two rows in Figure 6.3. The starting point  $Z_0$  was set to be a matrix whose entries were equal to the average of the grayscale value of  $\Omega$ , and the tolerance was set to  $\varepsilon = 10^{-10}(1 + \|\nabla f(Z_0)\|_F)$ . Following the analysis in [46], AIPP uses the curvature pair  $(m, M) = (2/\delta, 1 + \tau)$ .

Table 6.3 presents the relative error<sup>9</sup> of the final candidate image and runtime (in

<sup>9</sup>For a candidate image  $\hat{Z}$ , this quantity is defined as  $\|\hat{Z} - X\|_F$  divided by  $\max_{Z \in \Xi} \|Z - X\|_F$  where  $\Xi$  is the set of all grayscale images. Its value can range from 0.0 (full recovery) to 1.0.

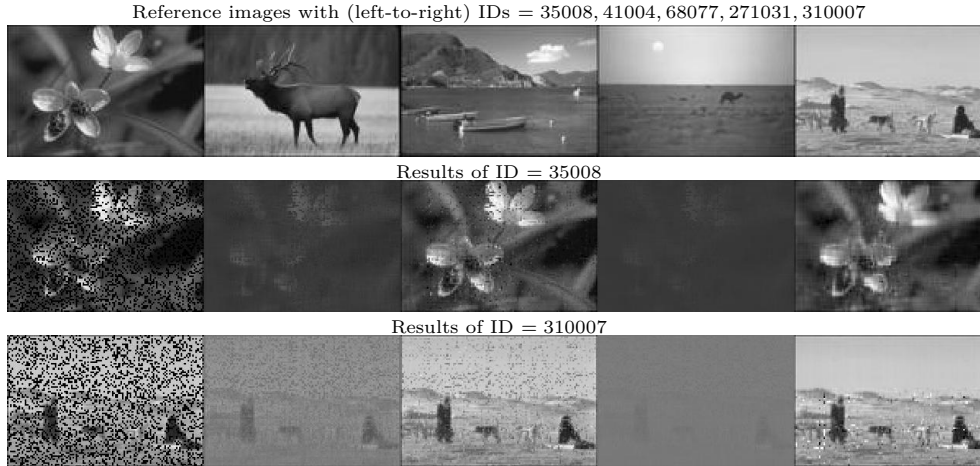


FIGURE 6.3. The first row presents the downscaled ( $80 \times 120$ ) reference images  $X$  taken from the Berkeley Segmentation Dataset, along with their image IDs (in order). The second and third rows present the results of the LPMC experiments for two of the images. Specifically, each of these rows presents (from left to right) the corrupted image  $\Omega$  and the images generated by UPF, ANCF, AIPP, and APD, respectively.

785 seconds) for the different reference images, and the last two rows in Figure 6.3 show the candidate images generated by each method for two of the reference images.

image id	Relative Error				Runtime (seconds)			
	UPF	ANCF	AIPP	APD	UPF	ANCF	AIPP	APD
35008	0.220	0.059	0.241	<b>0.034</b>	104.7	174.5	89.2	<b>44.4</b>
41004	0.259	0.103	0.312	<b>0.072</b>	114.6	175.9	90.6	<b>45.9</b>
68077	0.238	0.075	0.276	<b>0.046</b>	107.6	175.6	89.2	<b>43.0</b>
271031	0.272	0.146	0.363	<b>0.079</b>	117.5	176.8	95.7	<b>48.0</b>
310007	0.265	0.079	0.324	<b>0.048</b>	116.0	186.4	92.7	<b>44.9</b>

TABLE 6.3  
Relative errors and runtimes in the LPMC experiments for different reference images in the LPMC experiments. The bolded numbers indicate the best algorithm in terms of the relative error (less is better) and runtime in seconds (less is better).

786

787 **6.4. Comments about the numerical results.** In Subsection 6.1, APD sub-  
788 stantially outperformed<sup>10</sup> its competitors and its non-adaptive variant AIPP under  
789 the given numerical tolerance  $\epsilon$ . However, Figure 6.1 showed that ANCF was more  
790 comparable to PF.APD when the curvature ratio  $M/m$  was large or a larger (more  
791 lenient) tolerance was given. In Subsection 6.2, APD consistently outperformed its  
792 competitors on all metrics. For the number of gradient evaluations, UPF performed  
793 similarly to APD but was among the worst adaptive methods for function evaluations.  
794 In Subsection 6.3, APD generated higher-quality candidate images compared to its  
795 competitors under a fixed iteration budget. Specifically, it was shown in Figure 6.3  
796 that PF.APD generated images with fewer artifacts, more consistent lighting, and in  
797 a more timely manner.

798 **7. Concluding Remarks.** This paper establishes iteration complexity bounds  
799 for PF.APD that are only optimal, up to logarithmic terms, in terms of  $(M, \Delta_0, \epsilon)$

<sup>10</sup>5-20x (resp. 2-7x) fewer function (resp. gradient) evaluations for ANCF and 27-60x (resp. 2-6x) fewer for UPF.

800 when  $f$  is convex and in terms of  $(m, M, \Delta_0, \varepsilon)$  when  $f$  is weakly-convex. Conse-  
 801 quently, it remains to be seen whether an optimal complexity bound in terms of  $d_0$   
 802 exists for a parameter-free and convexity-unaware method.

803 To alleviate the issues regarding the  $d_0$ -suboptimal complexity of APD (specifi-  
 804 cally, when  $f$  is convex and  $d_0$  is unknown) one could consider running running  $S + 1$   
 805 instances of PF.APD (either in lockstep or in parallel) with different initial estimates  
 806  $m_0 = 1, \varepsilon, \varepsilon/2, \dots, \varepsilon/2^{S-1}$ ; in particular, the whole scheme stops when one of these  
 807 instances stops successfully. The number of resolvent evaluations of this approach is  
 808 at most  $S + 1$  times the minimum of the bound in (3.19) over the different values of  
 809  $m_0$ . Consequently, following the remarks at the end of Section 3, if  $d_0 \leq 2^{S-1}$  then  
 810 one of the  $S + 1$  instances obtains the lower bound in Table 1.1 for the convex case;  
 811 otherwise, the bound for APD in Table 1.1 is obtained. Moreover, if  $S$  is chosen small  
 812 compared to the other terms in (3.19) and  $d_0 \leq 2^{S-1}$ , then the cost is on the same  
 813 order of magnitude as the  $(M, \Delta_0, d_0, \varepsilon)$ -complexity optimal method described at the  
 814 end of Section 3 (which requires knowledge of  $d_0$ ).

815 In addition to the applications in Section 5, it would be interesting to see if  
 816 PF.APD could be leveraged to develop a parameter-free proximal augmented La-  
 817 grangian method, following schemes similar to ones as in [20, 27].

818

#### REFERENCES

- 819 [1] M. AHOOKHOSH AND A. NEUMAIER, *Solving structured nonsmooth convex optimization with*  
 820 *complexity  $\mathcal{O}(\varepsilon^{-1/2})$* , TOP, 26 (2018), pp. 110–145.  
 821 [2] M. M. ALVES, R. D. C. MONTEIRO, AND B. F. SVAITER, *Regularized HPE-type methods for*  
 822 *solving monotone inclusions with improved pointwise iteration-complexity bounds*, SIAM  
 823 J. Optim., 26 (2016), pp. 2730–2743.  
 824 [3] H. H. BAUSCHKE, P. L. COMBETTES, ET AL., *Convex analysis and monotone operator theory*  
 825 *in Hilbert spaces*, vol. 408, Springer, 2011.  
 826 [4] A. BECK, *First-order methods in optimization*, SIAM, 2017.  
 827 [5] A. BECK AND M. TEBULLE, *A fast iterative shrinkage-thresholding algorithm for linear inverse*  
 828 *problems*, SIAM J. Imaging Sci., 2 (2009), pp. 183–202.  
 829 [6] Y. CARMON, J. C. DUCHI, O. HINDER, AND A. SIDFORD, *Accelerated methods for nonconvex*  
 830 *optimization*, SIAM J. Optim., 28 (2018), pp. 1751–1772.  
 831 [7] Y. CARMON, J. C. DUCHI, O. HINDER, AND A. SIDFORD, *Lower bounds for finding stationary*  
 832 *points II: first-order methods*, Math. Program., 185 (2021), pp. 315–355.  
 833 [8] D. DAVIS AND D. DRUSVYATSKIY, *Stochastic model-based minimization of weakly convex func-*  
 834 *tions*, SIAM J. Optim., 29 (2019), pp. 207–239.  
 835 [9] D. DRUSVYATSKIY, *The proximal point method revisited*, arXiv preprint arXiv:1712.06038,  
 836 (2017).  
 837 [10] D. DRUSVYATSKIY AND C. PAQUETTE, *Efficiency of minimizing compositions of convex func-*  
 838 *tions and smooth maps*, Math. Program., 178 (2019), pp. 503–558.  
 839 [11] M. I. FLOREA AND S. A. VOROBYOV, *An accelerated composite gradient method for large-scale*  
 840 *composite objective problems*, IEEE Trans. Signal Process., 67 (2018), pp. 444–459.  
 841 [12] S. GHADIMI AND G. LAN, *Accelerated gradient methods for nonconvex nonlinear and stochastic*  
 842 *programming*, Math. Program., 156 (2016), pp. 59–99.  
 843 [13] S. GHADIMI, G. LAN, AND H. ZHANG, *Generalized uniformly optimal methods for nonlinear*  
 844 *programming*, J. Sci. Comput., 79 (2019), pp. 1854–1881.  
 845 [14] S. GUMINOV, P. DVURECHENSKY, N. TUPITSA, AND A. GASNIKOV, *On a combination of alter-*  
 846 *ating minimization and Nesterov’s momentum*, in Int. Conf. Mach. Learn., PMLR, 2021,  
 847 pp. 3886–3898.  
 848 [15] S. GUMINOV, Y. NESTEROV, P. DVURECHENSKY, AND A. GASNIKOV, *Accelerated primal-dual*  
 849 *gradient descent with linesearch for convex, nonconvex, and nonsmooth optimization prob-*  
 850 *lems*, in Dokl. Math., vol. 99, Springer, 2019, pp. 125–128.  
 851 [16] W. HARE AND C. SAGASTIZÁBAL, *A redistributed proximal bundle method for nonconvex opti-*  
 852 *mization*, SIAM J. Optim., 20 (2010), pp. 2442–2473.  
 853 [17] W. KONG, *Accelerated inexact first-order methods for solving nonconvex composite optimiza-*  
 854 *tion problems*, arXiv preprint arXiv:2104.09685, (2021).

- 855 [18] W. KONG, J. G. MELO, AND R. D. C. MONTEIRO, *Complexity of a quadratic penalty accelerated inexact proximal point method for solving linearly constrained nonconvex composite programs*, SIAM J. Optim., 29 (2019), pp. 2566–2593.
- 856
- 857 [19] W. KONG, J. G. MELO, AND R. D. C. MONTEIRO, *An efficient adaptive accelerated inexact proximal point method for solving linearly constrained nonconvex composite problems*, Comput. Math. Appl., 76 (2020), pp. 305–346.
- 858
- 859
- 860 [20] W. KONG, J. G. MELO, AND R. D. C. MONTEIRO, *Iteration-complexity of a proximal augmented Lagrangian method for solving nonconvex composite optimization problems with nonlinear convex constraints*, arXiv preprint arXiv:2008.07080, (2020).
- 861
- 862 [21] W. KONG, J. G. MELO, AND R. D. C. MONTEIRO, *FISTA and extensions-review and new insights*, arXiv preprint arXiv:2107.01267, (2021).
- 863
- 864 [22] W. KONG AND R. D. C. MONTEIRO, *An accelerated inexact proximal point method for solving nonconvex-concave min-max problems*, SIAM J. Optim., 31 (2021), pp. 2558–2585.
- 865
- 866 [23] H. LI AND Z. LIN, *Accelerated proximal gradient methods for nonconvex programming*, Adv. Neural Inf. Process. Syst., 28 (2015).
- 867
- 868 [24] J. LIANG AND R. D. C. MONTEIRO, *A doubly accelerated inexact proximal point method for nonconvex composite optimization problems*, arXiv preprint arXiv:1811.11378, (2018).
- 869
- 870 [25] J. LIANG, R. D. C. MONTEIRO, AND C.-K. SIM, *A FISTA-type accelerated gradient algorithm for solving smooth nonconvex composite optimization problems*, Comput. Math. Appl., 79 (2021), pp. 649–679.
- 871
- 872 [26] M. MARQUES ALVES, R. D. C. MONTEIRO, AND B. F. SVAITER, *Iteration-complexity of a Rockafellar’s proximal method of multipliers for convex programming based on second-order approximations*, Optimization, 68 (2019), pp. 1521–1550.
- 873
- 874 [27] J. G. MELO, R. D. C. MONTEIRO, AND W. KONG, *Iteration-complexity of an inner accelerated inexact proximal augmented Lagrangian method based on the classical lagrangian function and a full Lagrange multiplier update*, arXiv preprint arXiv:2008.00562, (2020).
- 875
- 876 [28] R. D. C. MONTEIRO, C. ORTIZ, AND B. F. SVAITER, *An adaptive accelerated first-order method for convex optimization*, Comput. Math. Appl., 64 (2016), pp. 31–73.
- 877
- 878 [29] R. D. C. MONTEIRO, M. R. SICRE, AND B. F. SVAITER, *A hybrid proximal extragradient self-concordant primal barrier method for monotone variational inequalities*, SIAM J. Optim., 25 (2015), pp. 1965–1996.
- 879
- 880 [30] R. D. C. MONTEIRO AND B. F. SVAITER, *Convergence rate of inexact proximal point methods with relative error criteria for convex optimization*, Optimization Online preprint, (2010).
- 881
- 882 [31] R. D. C. MONTEIRO AND B. F. SVAITER, *On the complexity of the hybrid proximal extragradient method for the iterates and the ergodic mean*, SIAM J. Optim., 20 (2010), pp. 2755–2787.
- 883
- 884 [32] Y. NESTEROV, *A method for unconstrained convex minimization problem with the rate of convergence  $O(1/k^2)$* , Dokl. Akad. Nauk, 269 (1983), pp. 543–547.
- 885
- 886 [33] Y. NESTEROV, *Introductory lectures on convex optimization: A basic course*, vol. 87, Springer Science & Business Media, 2003.
- 887
- 888 [34] Y. NESTEROV, *How to make the gradients small*, Optim. Math. Optim. Soc. Newsl., (2012), pp. 10–11.
- 889
- 890 [35] Y. NESTEROV, *Gradient methods for minimizing composite functions*, Math. Program., 140 (2013), pp. 125–161.
- 891
- 892 [36] Y. NESTEROV, *Universal gradient methods for convex optimization problems*, Math. Program., 152 (2015), pp. 381–404.
- 893
- 894 [37] Y. NESTEROV, *Lectures on convex optimization*, vol. 137, Springer, 2 ed., 2018.
- 895
- 896 [38] Y. NESTEROV, A. GASNIKOV, S. GUMINOV, AND P. DVURECHENSKY, *Primal–dual accelerated gradient methods with small-dimensional relaxation oracle*, Optim. Methods Softw., (2020), pp. 1–38.
- 897
- 898 [39] A. NEUMAIER, *Osga: a fast subgradient algorithm with optimal complexity*, Math. Program., 158 (2016), pp. 1–21.
- 899
- 900 [40] C. PAQUETTE, H. LIN, D. DRUSVYATSKIY, J. MAIRAL, AND Z. HARCHAOU, *Catalyst acceleration for gradient-based non-convex optimization*, arXiv preprint arXiv:1703.10993, (2017).
- 901
- 902 [41] N. PARIKH, S. BOYD, ET AL., *Proximal algorithms*, Found. Trends Optim., 1 (2014), pp. 127–239.
- 903
- 904 [42] R. T. ROCKAFELLAR, *Augmented Lagrangians and applications of the proximal point algorithm in convex programming*, Math. Oper. Res., 1 (1976), pp. 97–116.
- 905
- 906 [43] C.-K. SIM, *A FISTA-type first order algorithm on composite optimization problems that is adaptable to the convex situation*, arXiv preprint arXiv:2008.09911, (2020).
- 907
- 908 [44] J. TRZASKO AND A. MANDUCA, *Highly undersampled magnetic resonance image reconstruction via homotopic  $\ell_0$ -minimization*, IEEE Trans. Med. Imaging, 28 (2008), pp. 106–121.
- 909
- 910 [45] F. WEN, L. CHU, P. LIU, AND R. C. QIU, *A survey on nonconvex regularization-based sparse*

- 917                    *and low-rank recovery in signal processing, statistics, and machine learning*, IEEE Access,  
918                    6 (2018), pp. 69883–69906.
- 919 [46] Q. YAO AND J. KWOK, *Efficient learning with a family of nonconvex regularizers by redistribut-*  
920 *ing nonconvexity*, in Int. Conf. Mach. Learn., PMLR, 2016, pp. 2645–2654.
- 921 [47] C.-H. ZHANG, *Nearly unbiased variable selection under minimax concave penalty*, Ann. Statist.,  
922                    38 (2010), pp. 894–942.
- 923 [48] D. ZHOU AND Q. GU, *Lower bounds for smooth nonconvex finite-sum optimization*, in Int.  
924                    Conf. Mach. Learn., PMLR, 2019, pp. 7574–7583.